The spreadsheet version of this table has been made available at https://doi.org/10.5281/zenodo.5646467 (navigating the PDF version below will require a PDF viewer with zoom and search functionality)

Type of Study population and sample size							Torse of	Study population and complexity if	Mathedalom Study			
Number	Authors Auranoff, D and Fernandez-Reves, D and Papadopoulos, M C and Roias, S A and Herbster, M and	Title Identification of diagnostic markers	Journal	Volume Is	sue Pages	Year Location	URL / DOI publicatio	n applicable	design design	Outcome measures, if applicable	Validation type	Main results Key findings that relate to the review question SVM classifier discriminated the proteomic profile of patients with active tuberculosis
	Locamore, A and Tarelli, E and Sheldon, J and Schweek, A and Politol, R and Rayner, C P and Krishna, 1 5	, for tuberculosis by proteomic fingerprinting of serum	Lancet	368	9540 1012-102	1 2006 UK	https://dx.doi.org/10 10169/01-00 6736/0699342-2 article	179 serum samples from patients, 170 serum samples from controls	Case-control study	accuracy, sensitivity, specificity (in k-fold cross-validation + validation set text)	cross-validation + test set	Seek to describe the processing of the processin
	2 Ali, M and Aittokallio, T	Machine learning and feature selection for drug response prediction in precision oncology applications Ovarian Concer Classification Using Serum Proteomic Profiling and Wavelet Features A Comparison of	n Biophysical Reviews	11	1 31-39	2019 Finland	htto:lids.dei.org/10 1007/s12551-018- 0448-2 article	review (not applicable)	Review			prediction performance. Dang et al. 2014. Costello et al. 2015, However, the use of prediction performance Dang et al. 2014. Costello et al. 2015, However, the use of multiple entice perfors from various biological levels. can still improve the prediction multiple entice perfors from various biological levels. can still improve the prediction results
	3. Alqudah, A.M	Machine Learning and Features Selection Algorithms A mass spectrometry-based discovery	Clinical	44	4 165-173	2019	http://dx.doi.org/10 1097/JCE 000000 0000000359 article	262 cancer patients, 191 controls	Case-control study	accuracy, sensitivity, and precision (70% training set, 30% test set split)	training + test set	Recult show that all the presented ML algorithms preferred will five Caption Cacter Results show that all the presented ML right of the presented will five Caption Cacter Coultactions, with Enter that are selected polls from all excellent presents of the present of the present present and present and excellent presents are selected from the present present present present and the present
	Ashton, N J and Nevado-Holgado, A J and Lynham, S and Ward, M and Gupta, V B and Chatterjee, P and Gooze, K and Hone, E and Padrini, S and Bush, A I and Rowe, C C and Wilemagne, V L L and Ames 4 D and Masters, C L and A	and replication of a multi-analyte	Alzheimer's and	13	P1033- 7 P1033	2017	http://dx.doi.org/10 10166.setr 2017.0 meeting 6.1456 abstract	297 participants	Cases only (predicint neocortical amyloid burden)	accuracy (5-fold cross validation, external testing in other cohort)	cross-validation + external cohort validation	t highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP processing, highlights the convergence of pathways involved in coagula- tion, APP proces
	Assawamakin, A and Prueksaaroon, S and Kulawonganunchai, S and Shaw, P J and Varavithya, V and S Ruangrajitpakorn, T and Tongsima, S	bayes classification framework	Biomedical	2013	148014- 148014	2013 Thailand	http://dx.doi.org/10 _1155/2013/14801 4 article	more than 45 microarray and proteomi datasets of different sizes used	Case-control study	AUC, accuracy, sensitivity, specificity (10-fold cross-validation)	cross-validation	Current machine learning approaches are either too complete or perform poorly. See proposed two cells perform confidence of the complete or perform poorly. See proposed two cells perform confidence or perform poorly cells of the complete or performance or performance or performance or performance or performance or proposed proposed the cells performance or performance or proposed production of the cells o
	6 Awedat, K and Abdel-Qader, Land Springstead, J.R.	Prostate cancer recognition based on mass spectrometry sensing data and data fingerprint recovery Molecular classification of AML-MRC	Processing and Control	33	392-399	2017 USA	http://dx.doi.org/10 .1016/j.bispc.2016. 12.003 article	237 blood samples from subjects with different PSA levels	Case-control study (cases with different PSA levels)	accuracy, sensitivity, specificity, PPV, NPV (10-fold CV)	cross-validation	the distribution of the di
	Baser, C and Walter, W and Stengel, A and Hutter, S and Meggendorfer, M and Kern, W and Haferlach 7 C and Haferlach, T	survival MALDI-TOF analysis of blood serum	Blood	134		2019	http://dx.doi.org/10 1182blood-2019- 	619 patients with survival data	Case-control study	accuracy (10-fold CV)	cross-validation	AMA, with mytolopoplasia related changes (AMA, MRC) can be disposed using parties in the control of the disposed using parties in the control of the control
	Barcelo, F and Gomila, R and de Paul, I and Gil, X and Segura, J and Perez-Montana, A and Ilmenez-8 Marco, T and Sampol, A and Portugal, J	proteome can predict the presence o monoclonal gammopathy of undetermined significance	PLoS One	13	e020179 8 e020179	3. 3 2018 Spain	http://dx.doi.org/10 .1371(journal.pone .0201793 article	103 patients clinically diagnosed with MGUS, 108 healthy volunteer donors	Case-control study	accuracy, sensitivity, specificity (20-fold cross-validation)	cross-validation	MAIL TO analys of block seem residence using agreed vector machines can place to the processor of the proces
	9 Barla, A and Jurman, G and Riccadonna, S and Merler, S and Chierici, M and Furfanello, C	Machine learning methods for predictive proteomics A common gene signature across	Brief Bioinform American	9	2 119-128	2008 Italy	https://doi.org/10.1 093/bb/bbn008 article	review (not applicable)	Review			process party in predictive processing, prices and continues seagues of processing and processing prices and prices and pricessing prices and pricessing p
	Baron, D and Ramstein, G and Chesneau, M and Echasseriau, Y and Paller, A and Paul, C and Degauque, N and Hernandez-Fuentes, M and Sanchez-Fueyo, A and Newell, K and Giral, M and 10 Soullilou, J P and Houlgatte, R and Brouard, S	A common gene signature across multiple studies identifies biomarkers and functional regulation in tolerance to renal allografit Improving the Prediction of Survival it	tion	15		2015	http://dx.doi.org/10 .1038/ki.2014.395 article	96 samples with tolerance to renal allograft	Case-control study	accuracy, sensitivity, specificity (6-fold cross-validation + external validation)	cross-validation + test set	A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly in STM or Cases. A gene signature derived from blood cell transcriptional data predicts tolerance to recall allegard correctly
	11 Bashiri, A and Ghazisaeedi, M and Safdari, R and Shahmoradi, L and Ehtesham, H	Cancer Patients by Using Machine Learning Techniques: Experience of Gene Expression Data: A Narrative Review Computational systems medicine- what we can learn from Arnold	Journal of Public Health	46	2 165-172	2017 Iran	https://pubmed.nc bi.nlm.nih.gov/284 51550/ article	review (not applicable)	Case-control study			synamics, application, developing clinical decision reports synamics applications, developing clinical decision reports synamics applications, developing clinical decision report synamics and on these methods for an advantage of the synamics application, developing clinical decision report synamics and the synamics applications, developing clinical decision apport synamic and the submitted for an advantage of the synamics and synamics applications, developing clinical decision apport synamics and make a present and support and individualized transmiss to particular and improve the prognosis and individualized transmiss to particular and improve the prognosis and cransmiss. In a transmission apport synamics and individualized transmiss to particular and improve the prognosis of cransmiss. The authors present tool developments in graph based matchine learning to
	12 Baumbach, J	Schwarzenegger about breast cancer survival	Systems Medicine	2	1 A-29	2019	http://dx.doi.org/10 .108/invsm.2019.2 9005 article	review (not applicable)	Case-control study			DDX general and usually a routal number of samples (cit 2 passens) labels to model computational approaches for network based modeline. we investigated modeline general partner of three severe metables disorders. Prices consequent approaches for network based modeline. we investigated modeline general partner of three severe metables disorders. Prices consequent and dispensio using a sociotion tree paradigm and legistic regression subject (EAR). For the final model building process we assessed the relevance of analysis (EAR) for the final feeb building processes assessed the relevance of
	13 Summartner C and Bohn. C and Summartner D	Modelling of classification rules on metabolic patterns including machine learning and expert knowledge	a J Biomed		2 89.98	2005 Austria	https://doi.org/10.1 016/j.lin/2004.08.0	PAHD, n = 94 cases 1241 randomly	Case, metrol study	accurars, sensitivity, specificity (10-fold cross-validation)	cross,validation	analysis (18.4) for the 18.4 model shalling process we consent the released consentration of control and the control of the co
	15 audingarities, r., amo comin, r., amo camingarities, J. Best, M. G. and Soo, I and Koo, I hard Tannous, J. and Westerman, B. A. and Rustenburg, F. and Schellen, P. and Verschuseren, H. and Post, E. and Koster, J. and Yistra, B. and Amesiane, N. and Dorsman, J. and Smit, F. F. and Verheid, H. M. and Noske, D. P. and Reijneveld, J. C. and Nisson, R. J. A. and Tannous, B. A. and 14. Wesseling, P. and Wurdinger, T.	RNA-Seq of Tumor-Educated Platelet: Enables Blood-Based Pan-Cancer, d Multiclass, and Molecular Pathway Cancer Diagnostics	Cancer Cell	28	5 666-676	Netherlan 2015 s	bite-lists dei org/10 d 10164 coet 2015. 09.018 article	228 patients with localized and metastasized tumors and 55 healthy individuals	,	accuracy, sensitivity, specificity (leave-one-out cross-validation)	cross-validation	We determined the diagnostic potential of TBPs by mBNA sequencing of 283 platelet. A licave-one-out cross-validation support vector machine algorithm (SVM), COCV), samples. We distinguished 282 patients with localized and metastaziaed tumors from trained on mRNA profiles of tumor-educated blood platelets (TBPs), can distinguish substitution of terms with SSRS crussures.
	Bhak, Y and Jeong, H O and Cho, Y S and Jeon, S and Cho, J and Gim, J A and Jeon, Y and Blazyte, A and Park, S G and Kim, H M and Shin, E S and Paik, J W and Lee, H W and Kang, W and Kim, A and Kim 15 Y and Kim, B C and Ham, B J and Bhak, J and Lee, S	Depression and suicide risk prediction n, models using blood-derived multi- omics data	n Translationa I Psychiatry	9	8-8	United 2019 States	http://dx.doi.org/10 1038941398-019 0595-2 article	56 suicide attempters (SAs), 39 patient: with major depressive disorder (MDD), and 87 healthy controls	Case-control study	accuracy, sensitivity, specificity, PPV, NPV (leave-one-out cross-validation)	cross-validation	We developed machine learning models to predict depression and usides risk using blood metaphora and uncongreption and unconformed studies. The summer of the predict depression and usides risk using blood metaphora and uncongreption and usode risk using blood metaphora and uncongreption and uncongreption and usode risk using blood metaphora and uncongreption and uncongreption and usode risk using blood metaphora and uncongreption and u
	16 Bhanot, G and Alexe, G and Venkataraghavan, B and Levine, A J	A robust meta-classification strategy for cancer detection from MS data	Proteomics 2017 Third leee Internationa	6	2 592-604	2006 USA	http://dx.dec.org/10 1502/smic 20050 0192 article	322 serum spectra (63 with normal prostate)	(including cases with different PSA levels)	accuracy, sensitivity, specificity (training data: 215 cases, test data: 107 cases	training + test set	The paper presents a noise analysis and filtering procedure followed by combining the fiversults of several michine learning tools to produce a rebusy prediction of the results of several michine learning tools to produce a rebusy prediction of the spectra based cancer diagnosis (sentificity of 90.31% and a speci-ficity of 98.81%,). The paper presents a noise analysis and filtering procedure followed by combining the results of several machine learning tools to produce a rebust predictor for MS spectra based cancer diagnosis.
			Conference on Research in Computatio nal Intelligence									
	17 Bhattacharjee, S and Singh, Y J and Ray, D Bhorade, S and Bellinger, C and Bernstein, M and Dotton, T and Feller-Kopman, D and Lee, H and	Comparative Performance Analysis of Machine Learning Classifiers on Ovarian Cancer Dataset IMPROVING INDETERMINANT	tion Networks		213-218	2017 USA	https://doi.org/10.1 109/ICRCICN.201 7.8234509 meeting abstract	121 cancer samples, 95 benign	Case-control study	accuracy, sensitivity, specificity (10-fold cross-validation + external test data)	cross-validation + test set	In a comparative evaluation of machine learning methods on mass spectrometry (MS) In a comparative evaluation of machine learning methods on mass spectrometry (MS) for disputing telesping and miligrant formed of contain councer, nearly evaluation of machine learning methods on mass spectrometry (MS) for disputing telesping and miligrant formed of Counties councer, nearly release and professional professio
	senorada, Saind Bellinger, C. and Berinstein, M. and Dotson, I. and Feller-Kogman, D. and Lee, H. and Choi, Y. and Pankratz, D. and Lofaro, L. and Walsh, P. and Huang, J. and Kennedy, G. and Wahildi, M. and 18. Mazzone, P.	PULMONARY NODULE MANAGEMEN WITH THE PERCEPTA GENOMIC SEQUENCING CLASSIFIER		156	A2271- 4 A2272	2019	1016); chest 2019 08 907 meeting abstract	1600 patients training set, 412 samples validation set	Case-control study	accuracy, NPV, PPV (training/test set split)	training + test set	Transcriptome RMA sequencing to classify the probability of miligrancy for larg another patients with configuration to the control of the probability of miligrancy for large another patients with configuration to configuration. The best of responds to application contained STSs - anothering for benefit intermediate pre-set or risk groups. The authors controlled application of the probability of the probabili
	Bochare, A and Gangopadhyay, A and Yesha, Y and Joshi, A and Yesha, Y and Brady, M and Grasso, M 19 A and Rishe, N	Integrating domain knowledge in I supervised machine learning to asses the risk of breast cancer	Medical Engineering s and Informatics	6	2 87-99	2014	http://dx.doi.org/10 .1504IJME12014. 060245 article	1145 cases and 1142 controls	Case-control study	accuracy, sensitivity, specificity (10-fold cross-validation)	cross-validation	the nixed developing power counter for proteocoposal among and expense to the developing power counter for proteocoposal among and expense description. The machine having model proteocoposal among and proteocoposal and feature selection performed better compared to a conventional classification approach.
		Predictive value of targeted										The authors investigated the ability of targeted proteomics to predict presence of highly-risk plaque or advanced recomposed presences in parties with supposed (Cit. The developed machine issuing model that for dispassic, performance and washine investigated the ability of targeted proteomics is predict with supposed (Cit. The developed machine issuing and predict that of dispassic, performance and washine dispassic interpretations (U.M 0.8 of 0.9 of 0.90). Excessing, a dis- vokable dispassic interpretations (U.M 0.8 of 0.90). Excessing, a dispassic interpretations (U.M 0.8 of 0.90), a constant of the washine disreductions (U.M 0.8 of 0.90). The conversal pre- viously disreductions (U.M 0.8 of 0.90). The conversal pre- triations (U.M 0.90). The conversal pre- triati
	Bom, M. and Levin, E and Dissean, E.S. and Disseal, E and Yux Right, C.C. and von Risseam, A.C. and Namaka, Jandhi, M.F., and Leiper, L.A. and Borner, J. Part Toplor, C. and Risseadops, M. and 20 Bajimskers, P. G. and Konnig, W. and Groun, A.K. and Stroen, E.S. G. and Kinaspen, P.	suspected coronary artery disease	EBioMedicin e	39	109-117	2019 Canada	http://ldx.doi.org/10 1016/j.doi.om.201 8.12.033 article	196 patients with suspected coronary artery disease	Case-control study	AUC (10 fold CV + 20% hold out test set)	cross-validation + test set	subset of 34 proteins was predictive for the absence of CAD (JUL - 08 5 105), again outperforming profescion with generally subside characteristics (JUL - 07 1 0 00, g + 0.05). The advant sent by profescion with generally subside characteristics (JUL - 07 1 0 00, g + 0.05). The advant sent by profescion between Closer (E) putent reports to the displace Common, DNA methylation and mithing profescion with general color certain complexity and profescion with general color sent terms of promotion common, DNA methylation and mithing profescion with general color sent terms of promotion complexity synthesis (Microcol Custom Feater on Cgli State methylation extended a complexity synthesis (Microcol Custom Feater on Cgli State methylation extended a complexity synthesis (Microcol Custom Feater on Cgli State methylation extended a complexity synthesis (Microcol Custom Feater on Cgli State methylation extended a complexity synthesis (Microcol Custom Feater on Cgli State methylation extended a complexity synthesis (Microcol Custom Feater on Cgli State methylation extended a complexity synthesis (Microcol Custom Feater on Cgli State methylation sentended and complexity of the complexity of the complexity complexity synthesis (Microcol Custom Feater on Cgli State methylation sentended and complexity profession of the complexity profession of the complexity of the complexity profession o
	21 Bomane, A and Gonçalves, A and Ballister, P J	Pacitaxel Response Can Be Predicted With Interpretable Multi-Variate Classifiers Exploiting DNA-Methylatio and miRNA Data		10		2019 France	http://dx.doi.org/10 3388/figene.2019. 01041 article	1,098 patients	Cases only (treatment response prediction)	t AUC (LOOCV)	cross-validation	considered City Sizes (MLY - 1,99) and a miRNA appreciate based of memory laws of considered City Sizes (MLY - 2,99) and a miRNA appreciate based of memory laws of the mixture of the mix
	Bovelstad, H M and Nygard, S and Storvold, H L and Aldrin, M and Borgan, O and Frigessi, A and 22 Lingsherde, O C	Predicting survival from microarray dataa comparative study Systems genomics of ulcerative colitis	Bioinformati cs	23	6 2080-208	7 2007 Norway	https://doi.org/10.1 093/bioinformatics/ btm305 article	Benchmark comparison across several cancer datasets with > 50 samples per condition	Case-control study	Log rank test p-value (10-fold cross-validation)	cross-validation	combinations of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the simple variable selection methods. For our data sets, rigide pregression has the overall better performance. **Advantage of the approximation is a common from the property of the preference of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the combinations of the gene expression values have much better performance than the combinations of the gene expression values have been performance than the combination of the gene expression values have been performance that the general perfo
	Brooks, J. and Modoo, D. and Sudhakar, P. and Fazekas, D. and Zouffr, A. and Watton, A. and Tremelling, 23 M. and Verstockt, B. and Vermeire, S. and Bender, A. and Carding, S. and Korcsmaros, T.	Combining GWAS and signalling networks for patient stratification an individualised drug targeting in ulcerative colitis	Crohn's and	13	\$006-\$00	7 2019	http://dx.doi.org/10 1093/secro- jeclly/222.000 abstract	377 UC patients	Cases only (clustering using genomic footprint)	"We validated the workflow on a larger cohort of 941 UC patients from the IBD Biobank in Leuven"	external cohort validation	and leaderful Common and differing garbagein; garbragein between them, We thoward that clasters were intelled to gender and god roster of disease, but unrelated to therepoint; opecaling of Britange, With machine learning, we identified a subset of patients from whithin own of the closter, or whom the presence of a regulatory MAMAZ SNP-was a marker for therepoint; opecaling.*

Brown, K K and Choi, Y and Coby, T V and Felherry, K R and Grobburg, S and Intelse, U and Synch, Divard Shellow, M P and Marriese, F J and Pankers, C G and Wolle, P S and Hung, J and Starth, M A and Raphu, G and Romenty, G C	A Prospective validation of a genomic classifier for usual interstitial pneumonia in transbronchial biopsi	American Journal of Respiratory and Critical Care ies Medicine	195		2017	https://www.atsiour nats.org/biolabs/10 1364/siscon- conference_017.1 96.1_Mesting-plass acts_46722 abstract	354 TBB samples	Case-control study	AUC praining / text set split)	training + test set	Address displaces of displaces produced Places (PF) reports the stress of a substitute of the produced places (PF) reports the stress of a substitute of the produced places (PF) reports the produced
25 Cai, Q and Alvarez, J.A and Yang, J. and Yu, T	Network Marker Selection for Untargeted LC-MS Metabolomics D	J Proteome lata Res	16 3	1261-1269	United 2017 States	https://doi.org/10.1 02/lines.herotocres (mo0/86) article	subjects with available high-resolution plasma metabolomics from the Emory-Georgia Tech Predictive Health Initiative Cohort of the Center for Health Discovery and Well Being (N = 371)	Cases only (BMI analysis)	AUC (9-666 CV)	cross-validation	sophis a sequential feature screening procedure and machine feature flowering based or trava- sis select important subservision and identify to septim feature machine; in select important subservision and identifies one that the proposed method has a much light areasority with our commonly used machine machine; approach, for substructions associated with the body mass identified [MIII]. The method interface several underworks with placeble are on functional implications. There are the many price of large counter flaterurin, a well as or interfaced several underworks with placeble are on functional implications. There are the many price of large counter flaterurin, a well as offered several underworks with placeble are on functional implications. There are the many price of large counter, more and and the proposed in the admonstracions (LOCL), squamous of large core (ECCLC) as well as large of large admonstracions (LOCL), squamous of large core (ECCLC) as well as large of large admonstracions (LOCL), squamous of large core (ECCLC) as well as large of large admonstracions (LOCL), squamous of large core (ECCLC) as well as large of large admonstracions (LOCL), squamous of large core (ECCLC) as well as large of large admonstracions (LOCL), squamous of large core (ECCLC) as well as large of large admonstracions (LOCL), squamous of large core (ECCLC) as well as large of large admonstracions (LOCL), squamous of large core (ECCLC) as well as large of large according to the correct state, ECC. (Excellence of procedure (LOCLC), seven as large of large core correct state, the correct state, the Core correct state as the correct state as the correct state as the correct state as the process cor
26 Cai, Z and Ni, D and Zhang, Q and Zhang, I and Nipai, 5 M and Shon, I Casanovo, R and Yurmu, 5 and Simpson, B and film, M and An, Y and Saldana, 5 and Biverse, C and	Classification of lung cancer using ensemble-based feature selection a machine learning methods Blood metabolite markers of	Mol Biosyst	11 3	791-800	2015 China	hitos likisi omri 0.1 0.00/odenb00658c article	More than 100 samples available for the main sample groups LADC and SQCLC in both training and test set	Case-control study	accuracy, precision, recall, F-score (LOCCV)	cross-validation	and middle (Maximum Relanded) were proposed to see the Middle (Maximu
Mocrato, Pand Griswolf, M and Sonetag, D and Waithniet, Land Stainer, K and Innosco, P V and Brikkstorff, G and Alpellund, T and Laurier, L1 and Gudnason, V and Lagido Chajdry, C and 27 Thambisetty, M	preclinical Alzheimer's disease in tw longitudinally followed cohorts of older individuals	Alzheimers	12 7	815-822	2016 Australia	http://dx.doi.org/10 1016/i.julz 2015.1 2.008 https://www.suruk/ assilect.com/node/ 88215/suriste/ machine-learme-	two cohorts of n=93 and n=100 samples	Case-control study	AUC, sensitivity, specificity (6-fold CV)	cross-validation	from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. These findings underscore the importance of large- from two independent calcuts. The importance of large- from two independent calcuts. The importance of large- tion that it is a superior to the importance of large- tion that it is a superior to the importance of large- tion that it is a superior to the importance of large- tion that it is a superior to the importance of large- tion that it is a superior to the importance of large- tion that it is a superior to the importance of large- tion that it is a superior to the importance of large- tion that it is a superior to the importance of large- tion that it is a superior to the importance of large- tion that it is a superior to the importance of large- tion that
28 Chalbonchoe, A and Samarasinghe, S and Kulusiri, D	Machine Learning for Childhood Ac Lymphoblastic Leukaemia Gene Expression Data Analysis: A Review	Bioinformati	5 2	118-133	2010	for childhood, acute: hymohobiasic- teulaamia-gene- expression dish- analysis-a-review article	review (not applicable)				The greats comprehensive review of mulcine larming approaches that have been. We present a comprehensive review of mulcine larming approaches that have been used in place in predictable released [34]. In orders yet data. The emitted have used in place in predictable released [34], in increase yet data residued, been used in four major areas of microarray data analysis; gene selection, distanting, and distantially approaches that the predictable proposed of microarray data analysis; gene selection, distanting, and distantially approaches that the proposed of microarray data analysis; gene selection, distanting, and distantially approaches that the proposed of microarray data analysis; gene selection, distanting, and distantially approaches that the proposed of microarray data analysis; gene selection, distanting, and distantially approaches that the proposed of microarray data analysis; gene selection, distanting, and distantially approaches that the proposed of microarray data analysis; gene selection, distanting, and distantially approaches that the proposed of microarray data analysis; gene selection, distanting, and distantially approaches that the proposed of microarray data analysis; gene selection, distanting, and distantially approaches that the proposed of microarray data analysis; gene selection, distanting, and distantially approaches that the proposed of microarray data analysis; gene selection, distanting, and distantially approaches that the proposed of microarray data analysis; gene selection, distanting, and distantially approaches that the proposed of microarray data analysis; gene selection, distanting, and distantially analysis and dist
29 Chang, Y and Park, H and Yang, H J and Lee, S and Lee, K Y and Kim, T S and Jung, J and Shin, J M	Cancer Drug Response Profile scan (CDRscan): A Deep Learning Model That Predicts Drug Effectiveness fro Cancer Genomic Signature	om Sci Rep	8 1	8857-8857	2018 Australia	http://dx.doi.org/10 .1038h/41598-018- 27214-6 article	787 human cancer cell lines and structural profiles of 244 drugs were considered	Cases only (drug response study)	Risquared, AUC (training/test split)	training + test set	screening savy data encompassing genomic profiles of 732 human cancer cell lines and structural perfiles of 244 degs, the segalest factors to 1,4.78 approved degs and stemtised 14 occology and 22 non-excepting 140 perfiles 140 exception 140 performs
30 Chao, 5 M and Convolly, J and Ng, Y H and Ganecan, I and Bernett, L	Can urinary proteomes be used as a invasive markers for renal involvement in childhood febrile urinary tract infection (UTI)?	Pediatric Nephrology	31 10	1746-1746	2016	http://ktc.doi.org/10 1907/s/00467-016 3466-6 abstract	121 patients (68 males, 53 females)	Case-control study	sensitivity, PPV (10-fold CV)	cross-validation	validation, 62.3% agleistic status as no mail sairning, 55 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning, 65 (12.5% sentishing, 76.6% validation, 62.3% polentic status as on mail sairning
31 Chaudhary, K and Poirios, O B and Lu, L and Garmine, LX	Deep Learning-Based Multi-Omics Integration Robustly Predicts Survivin Liver Cancer		24 6	1248-1259	United 2018 States	https://dci.org/10.1 198/1078-0432-cor- 17-0653 article	360 patients included	Cases only (survival prediction)	"We validated this multi-omics model on five external datasets of various omics types: IB8-IP cohort (n - 230, C-index = 0.75), NCI cohort (n - 212, C-index = 0.87), NCI cohort (n - 212, C-index = 0.87), TabBi-36 cohort (n - 180, C-index = 0.87), TabBi-36 cohort (n - 40, C-index = 0.82)."	external cohort validation	360 HCC patients' data using NNA sequencing (NNA Seq.), mRNNA sequencing (NNA Seq.) and restly-states of the Tribe of Tr
Chong, R.S and Shalegin Rotamiolas, S and Ju, J M and Marietta, E V and Van Dyle, CT and Rijasetaran, J I and Jayaraman, V and Wang, T and Bei, K and Rijasetaran, K E and Krizhna, K and 32 Krizhnamurthy, H K and Murray, J A	Synthetic Neoepitopes of the Transglutaminase-Doamidated Glia Complex as Biomarkers for Diagnos and Monitoring Celiac Disease	ing Gastroenter	156 3	582-591.e1	United 2019 States	http://dec.deci.org/10 10536_pasteo_201 8.10.025 article	serum samples from 90 patients with biopsy-proven Cellac disease and 79 healthy individuals (controls)	Case-control study	AUC, accuracy, senditivity, specificity ("We validated out findings in 82 patients with newly diagnosed CeO and 217 controls.")	external cohort validation	Recreasing papities microarray pattern was used to estimate the articlosy-beinding internsity of each syndroid IT-CDP deptions. In the 201 straining control. The 201 straining control in the 201 straining control. The 201 straining control in the
Chang, WY and Correa, E and Yoshimura, K and Chang, M.C and Dennison, A and Takeda, 5 and 33 Chang, YT	Using probe electrospray ionization mass spectrometry and machine learning for detecting parcreatic cancer with high performance	Journal of Translationa	12 1	171-179	2020 Japan	hilina illenne ncis ni m.h. pordumelarii dessi ^{PMC} 70132211 article		Case-control study	accuracy, sendibing, specificity (1000 independent repetitions of a bootstrape cross-validation process)	cross-validation	amed to validate a unique diagnosis system using Probe Electropary Institution Mass Septemberg 1978 and Marchine Learning Learn de Alignosis of PSC. In the Septemberg 1978 and Marchine Learning Learnin
34 Clark, O and Saffshani, Z and Smirnov, P and Halbe-Kalins, B	Gene isoforms as expression-based biomarkers predictive of drug response in vitro	of Medical Science	187	S348-S348	2018 Canada	https://www.nature .com/articles/s414 87-017-01153-8 article	The data comprised 79,903 experiments for 140 different drugs tested on a panel of up to 778 unique cell lines from 30 tissue types 4,435 patient samples (3,066 patients (340 CRC and 2,759 non-CRC) were	Cases only (drug response prediction is vitro)	n- AUC, accuracy (validation in independent breast cancer data and different pharmacological assay)	external cohort validation	multiple cancer types. We further analyse two independent bravac cancer datasets and find that specific forms of IGEPPRY, MCRITO, RIGH, and MCROS-CO and epigification; passionated with AZESAL (specific), electrons, and packtases, and find that specific for indexed packtases, and find that specific forms of IGEPPR, MCRITO, RIGH, and RIGHD are applicated to a specific forms of IGEPPR, MCRITO, RIGH, and RIGHD are applicated to a specific forms of IGEPPR, MCRITO, RIGH, RIGH, and PROTOCO are a specific form of IGEPPR, MCRITO, RIGHT, RIGH, and PROTOCO are a specific form of IGEPPR, MCRITO, RIGHT, RIGHT, and PROTOCO are a specific form of IGEPPR, MCRITO, RIGHT, RIGHT, and PROTOCO are a specific form of IGEPPR, MCRITO, RIGHT, RIGHT, and PROTOCO are a specific form of IGEPPR, MCRITO, RIGHT, RIGHT, and PROTOCO are a specific form of IGEPPR, MCRITO, RIGHT, RIGHT, and IGEPPR, MCRITO, and PROTOCO are a specific form of IGEPPR, MCRITO, RIGHT, RIGHT, and PROTOCO are a specific form of IGEPPR, MCRITO, RIGHT, RIGHT, and IGEPPR, MCRITO, RIGHT, RIGHT, and IGEPPR and IGENT and IGE
35 Croner, LJ and Kao, A and Benz, R and Blume, J E and Dilton, R and Wilcox, B and Kain, S N	A new blood test for colorectal can in high-risk subjects	cer Clinical Chemistry	63	522-523	2017 Denmark	http://dx.doi.org/10 .11819bs.2020.0 30504 meeting abstract	(340 CRC and 2,759 non-CRC) were randomly assigned to the classifier discovery set. The remaining 1,336 samples (147 CRC and 1,189 non-CRC) were assigned to the validation set)	Case-control study	sensitivity, specificity, PPV, NPC (10-fold CV + training / test set split) "In assembling this review we conducted a broad survey of the different	cross-validation + test set	used to build and test condidate disastfers. The final classifier was a legistic regression using 10 posterior: eight proteins, and genders. In validation, the indeterminate care was 22.78, control to the proteins of the p
36 Cruz, J.A. and Wishart, D.S.	Applications of machine learning in cancer prediction and prognosis	Cancer Informatics 2019 16th Ieee Internationa I Conference on Computatio	2	59-77	2006 Greece	Nitro Tenno redi di mala monatoria cien PECCOTE COST article	review (not applicable)	Case-control study	"Is suitabling the fevere we conducted a finish stury of the distribution of the control of the		"he assembleing this review we conducted a broad survey of the different types of machine learning membrabe bulling such the type of data being integraped and the performance of these methods in career prediction and prognosis. A number of twoici, are noted, including a growing dependence in protein biomodexes, and and a second survey of the second survey of the second survey of the second survey of a heavy relations on "dist" submidges such artificial result infections (s NNNs) interested of more recently developed or more surject interpretable machine learning methods. A number of positioned states also speed to lack an appropriate level of machine learning methods can be used to substitively (15–200) improve the accuracy of predicting cancer succeptibility, recurrence and mortality."
Cuglist, G and Betweensta, S and Guarrers, S and Scientifics, C and Parko, S and Krogh, V and 37 Tumino, R and Vineix, P and Farkell, P and Matsids, G	Improving the prediction of cardiovaccular risk with machine-learning and DNA methylation data	nal Intelligence in Bioinformati cs and Computatio nal Biology -		39-42	2019 USA	bitos ifigiacaniona i cas orgificomenti 1221462 article	584 subjects (292 MI cases and 292 matched controls)	Case-control study	AUC, sensitivity, specificity (nested cross salidation)	cross-validation	"Classically, the cardiovascular risk of individual is evaluated using phenomenological "Classically, the cardiovascular risk of individual is evaluated using phenomenological variables (PP) such as allocal pressure, body years, uncher status, gender, age etc. variables (PP) such as allocal pressure, body years, uncher status, gender, age etc. variables (PP) such as allocal pressure, body years, uncher status, gender, age etc. variables (PP) such as allocal pressure, body years, uncher status, gender, age etc. variables (PP) such as allocal pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, uncher status, gender, age etc. variables (PP) such as a blood pressure, body years, gender, gender, pressure, pressu

38 Das, D and Sto, J and Endowalk, T and Trusts, K	An interpretable machine learning model for diagnosis of Alzheimer's disease	PeerJ	2019	3		2019 Japan	bhuilite doi ora'l 0 7717/mari 0543 article	97 AD subjects + 54 controls	Case-control study	AUC, accuracy, senditivity, specificity (cross-validation + text set)	cross-validation + test set
de Maturana, E L and Alonso, L and Alarcón, P and Martin-Antoniano, I A and Proeds, S and Floros, 39 and Calle, MT Lord Malars, N	E Challenges in the integration of omio and non-omics data	s Genes	10	3		2019 Spain	http://dx.doi.org/10 3309/iganes10030 230 article	review (not applicable)	Review		
40 de Ronde, 11 and Bonder, M3 and Lips, E H and Robenhuis, 5 and Woosle, L F	Breast cancer subtype specific classifiers of response to neoadjuvar chemotherapy do not outperform classifiers trained on all subtypes	PLoS One	9	2	e88551- e88551	2014	http://lik.doi.org/10 1371formal.ones 0088551	374 samples were analyzed	Cases only (treatment response prediction)	t AUC (Insted cross-validation)	cross-validation
	Effect of size and heterogeneity of										
Di Camillo, B and Sanavia, T and Martini, M and Jurman, G and Sambo, F and Barla, A and Squillario 41 M and Furlanello, C and Toffolo, G and Cobelli, C	samples on biomarker discovery:	PLoS One	7	3	e32200- e32200	2012 Italy	http://dx.doc.org/10 1371/journal.pone .0032200 article	3 different datasets (more than 50 samples per group in total)	Case-control study	AUC, sensitivity, specificity (cross validation + Monte Carlo bootstrap resampling)	cross-validation
42 Diaz Cano, 5 and Sutherland, R and Moorhead, J and Blanes, A and Doboon, R	Growth pattern analysis in low gradic clear cell renal cell carcinomas: Prognostic value and biologic significance	Laboratory Investigation	96		226A-226A	2016	http://dx.doi.org/10 1038/labinvast 20 meeting 16.10 abstract	low FG (1-2, 174 cases) vs. high FG (3-4, 139 cases) grade	Subtype comparison	AUC (50-fold cross-validation)	cross-validation
Diggaes, I and Kim, 5 Y and Hs, 22 and Pasketz, D and Wong, M and Reynolds, I and Tom, E and Page, M and Morroe, R and Rossi, I and Livels, V A and Lammar, R B and Ross, R T and Walnit, P 1 43 and Kennedy, G C	MACHINE LEARNING FROM CONCEP TO CLINIC: RELIABLE DETECTION OF BRAF VEODE DNA MUTATIONS IN THYROID MODULES USING HIGH- DIMENSIONAL RNA EXPRESSION DA:	Symposium on Biocomputi	n I		371-382	2015 USA	https://bubmed.nc bi.rkm.nih.gov/255 92597/ article	training (n=181) and independent test (n=535) sets	Case-control study	AUC [ID-fold CV + external test set)	cross-validation + test set
44 Ding, M Q and Chen, Land Cooper, G F and Young; J D and Lu, X	Precision Oncology beyond Targetee Therapy: Combining Omics Data with Machine Learning Matches the Majority of Cancer Cells to Effective Therapeutics	1	16	2	269-278	United 2018 States	http://dx.doi.org/10 .1158/1541- 7786.ms-17-0378 article	transcriptomics data from 727 cell lines was used	Cases only (cancer ce line drug response prediction)	II accuracy, sensitivity, specificity (25-fold cross-validation)	cross-validation
45 Djebburi, A and Labbo, A 46 Dougharry, E R and Hou, a and Bittner, M L	Refining gene signatures: a Bayesian approach Validation of computational method in genomics	cs	ti 10 8	1	410-410 1-19	2009 Canada 2007 USA	http://dx.doi.org/10 1.186/1471.2705- 10.410 article http://dx.doi.org/10 2.174/1486/00/077 80078656 article	the approach was applied to multiple cancer microarray datasets with > 50 samples per group in total review (not applicable)	Case-control study review	AUC, sensitivity, specificity (nested 10-fold CV + external testing). This paper treats the validation issue as it appears in two classes of the paper treats genome valence in personal - districtions and classring. It formulates the problem and reviews valent requisit.	cross-validation + test set
Drouin, A and Gigales, S and Disrape, M and Manchand, M and Tyers, M and Los, V G and Bourpa. 47 A M and Lookette, F and Cochet, I	Predictive computational phenotypi It, and Biomarker discovery using reference-free genome comparisons	ng BMC Genomics	17	1		2016 Canada	http://dec.doi.org/1.0 11880-12884-016- 2888-6 article	17 datasets in which the number of examples ranged from 111 to 556	Antibiotic resistance prediction	error rate (S-fold CV, test evaluation)	cross-validation + test set
48 Drostov, I and Klidt, M and Modlin, I	Graph-theoretic definition of neuroendocrine disease-a tumor specific mathematical toolbox for assessing neoplastic behaviour	Neuroendo rinology	ю 103		45-46	2016	http://dx.doi.org/10 meeting 11590000448725 abstract	Stage 1 (nonulation of 127 exacerbation	Case-control study	AUC, sensitivity, specificity, PPN_LNPV [The model was subdated in two independent sets (Set 1 $[n-115, MRN: n-72]$) Set $2 [n-120, MRN: n-581]$)	training + test set
Elsebaléhi, E and Lee, F and Schendel, E and Haque, A and Kathireason, N and Pathare, T and Syed, f 49 and Al-Ali, R	Large-scale machine learning based on functional networks for biomedic big data with high performance computing platforms	al Journal of Computational Science	11		69-81	2015	https://doi.org/1 0.1016/j.jocs.201 5.09.008 article	cases and 290 non-exacerbation controls) and Stage2 (population of 50 exacerbation cases and 114 non- exacerbation controls)	Case-control study	AUC, sensitivity, specificity (training/test set split)	training + test set
Fan, XJ and Wan, XB and Huang, Y and Gai, H M and Fu, XH and Yang, ZL and Chen, D K and Song, SS X and Wu, P H and Usi, Q and Wang, L and Wang, J P	Epithelial-mesenchymal transition biomarkers and support vector machine guided model in preoperatively predicting regional 5 lymph node metastasis for rectal cancer	Br J Cancer	106	11	1735-1741	2012 China	http://dx.doi.org/10 1038bis-2012.82 article	193 RC patients	Cases only (predicting lymph node metastasis)	; accuracy, constitutey, specificity (training/first set split)	training + test set
S1 Fang, Y and Xu, P and Yang, J and Qin, Y	A quantile regression forest based method to predict drug response an assess prediction reliability	d PLoS One	13	10	e0205155- e0205155	2018 China	http://dx.dei.org/10 .1371/journal.pone .0205155 article	data from 947 cell lines (CCLE dataset)	Cases only (drug response prediction is vitro)	n-Pearson correlation of observed and predicted drug response (out-of-bag validation)	outofbag
Farmakis, D and Kock, T and Mullin, W and Parissis, J and Gogos, B D and Nikolaos, M and Leltakin 52 and Mischak, H and Filippatos, G	Urine proteome analysis in heart failure with reduced specific fractio complicated by the onic kidney complicated by the onic kidney , J diseases feasibility, and clinical and pathogenetic correlates		18	7	822-829	2016 Germany	http://dx.doi.org/10 _1002/sght-544 article	126 individuals, 59 HFFEF patients and 67 controls	Case-control study	AUC, scouracy, sensitivity, specificity (cross-validation + text set)	cross-validation + test set

challenge demand the development and application of free waships a strategies to the integration or closical and applications of any extract pages to the integration or closical and perimenting at these in the integration and applications of the extract pages and perimenting at these integrations are closed as a strategy of the integration or closed and perimenting at these and perimenting and applications of the extract pages can be improve upon the performance of generic production. For the short and are strategies can be improve upon the performance of generic production. For the short are can be understood to the strategy of the perimenting of the

regularization which geniloses for the correlation between the variables selected. Our regularization which penalises for the correlation between the variables selected. Our regularization which penalises for the correlation between the variables selected. Our regularization which penalises for the correlation between the variables selected. Our regularization which penalises for the correlation between the variables selected. Our regularization which penalises for the correlation between the variables selected our selection and is able to exercise the penalises for the penalises f

The manuscript covers several commonly used approaches to evalute classification and clustering methods (e.g. cross-validation and bolstered resubstitution) The manuration covers several commonly used approaches to available constitution and of usering methods by a cross-validation and believe resolutions.

This identification of genome to beneather in a key say beautify in providing diagnosts:

The identification of genome to demand with a serial provided and provided and the control of the control of

Will, searchify yill 27%, NY V (54-60) and 69 V (25-20%). Additionally, multi-variance based makes bearing algorithms understanding designed may be a search of the control of the contro

"We present an interpretable machine learning model for medical dispection claims are all interpretable machine learning model for medical dispection claims are explained to a present many learning model for medical dispection claims are explained to a present many learning model for medical dispection claims are explained to a present many learning model for medical dispection claims are explained to a present many learning model for medical dispection of many intervals, which set them designed in a speciment the dispection with a long rate of 6 ca., compensation of many learning model for medical dispection of many intervals, which set them designed in a speciment the dispection of many intervals, which is settle on the dispection of many intervals, which is settle medical dispection of many intervals, which is settle medical dispection of many intervals, which is settle medical dispection of many intervals. The present many intervals are present to design a settle many learning ALIC = 0.86 + 0.09) *

Also — 6 Ms — 10015*

Only a small number of pallabled studies performed a "rai" integration of mixed and non-miss (Droll) data, manily a predict career extenses. Challeges in Old and integration regions of the number of pallabled studies performed a "rai" integration of mixed and non-miss (Droll) data, manily a predict career extenses. Challeges in Old and integration region for number of pallabled studies performed a "rai" integration of mixed and integration region for number of pallabled studies performed a "rai" integration of mixed and integration region of the number of pallabled studies performed a "rai" integration of mixed and integration region of the number of pallabled studies performed a "rai" integration of data, manily a predict care reactions. Challeges in Old data integration of acts are not existed as the high throughout oneic called an integration of a first and an application of new analysis caregies to data integration of acts and an application of new analysis caregies to distingent to Old data integration of acts and an application of new analysis caregies to distingent to Old data integration of acts and application of new analysis caregies to distingent to Old data integration of acts and application of new analysis caregies to distingent to Old data integration of acts and application of new analysis caregies to distingent to Old data integration of acts and application of new analysis caregies to distingent to Old data integration of acts and application of new analysis caregies and the development and application of new analysis caregies and the integration of acts and the development and application of new analysis caregies and integration of acts and applications of the acts and applications of t

successfully applied to a test set of 25 HFrEF patients and 33 controls, achieving 84% successfully applied to a test set of 25 HFrEF patients and 33 controls, achieving 849 sensitivity and 91% specificity."

												The goal of the druly was to set the hypothesis that specific proteins/peptides are differentially approach in joiners of examine with & without enhancemon as specific differentially approach in joiners of examine with & without enhancemon as part of the druly was to set the hypothesis that specific proteins/peptides are differentially approached in joiners of examine with & without enhancemon as part of proteins in the field, highly cornect conditrations and be reached on both reading and to set to, however good machine is enricing practice dictate the set of robust reading by reading the field are Ling tight enhanced in highly control conditration of the practice of the set of robust are set to be a set of robust reading by reading the field is a Ling tight enhanced in highly control conditration of the practice of the set of robust are set of the study was to set the hypothesis that specific proteins/peptides are the goal of the study was to set the hypothesis that specific proteins/peptides are that applied to the study in the set of the study in the protein of the study in the set of the study is a set of robust are set of the study was to set the hypothesis that specific proteins/peptides are the set of the study was to set the hypothesis that specific proteins/peptides are the set of the study was to set the hypothesis that specific proteins/peptides are the set of the study was to set the hypothesis that specific proteins/peptides are the protein set of the study was to set the hypothesis that specific proteins/peptides are the set of the study was to set the hypothesis that specific proteins/peptides are the set of the study was to set the hypothesis that specific proteins/peptides are the set of the study was to set the hypothesis that specific proteins are the set of the study was to set the hypothesis that specific proteins are the set of the study was to set the hypothesis that specific proteins are the set
Fassbender, A and Waelkens, E and Kyama, C and Bokor, A and Vodolaziaia, A and Verbeeck, N and Van De Piss, R and Ojeda, E and Gevaert, O and Meuleman, C and Peeraer, K and Tomassetti, C and S3 De Moor, B and D'Hooghe, T	Biomarkers in plasma or serum: Pitfalls in data processing A predictive model for survival in n small rell lung ranger (NSCIC) have	d	18	3 19	A-191A 2	2011 India	http://dx.doi.org/10 1177/1933719120 meeting 11183a/867 abstract	254 plasma samples from women with (n=165) & without (n=89) endometriosis	Case-control study	accuracy (data were divided randomly (100 times) into training set (70%) and test set (30%))	training + test set	anding by indinstring of the skill prices in Belduck in large collections an interest of the belduck in the prices of the straing set on the confirmed. The skill prices in Belduck in large collections performs one interest of the straing set
Filmone, N. and Ramos-Gejedo, J. and Oneng, D. and Tuck, D. P. and Shakin, A.R. and Chen, D. and Elber D. and Sung, F.C. and Johnson, B. and Shannon, C. and Pierce-Murray, K. and Gaynor, K. and Dedomest S. C. and Schiller, S. and Aljarapu, S. and Hall, R. and Ayunden, S. and Meng, F. and Brophy, M.T. and Do, N.	on electronic health record (EHR) a co, tumor sequencing data at the Department of Veterans Affairs (VA	nd Journal of Clinical k) Oncology	37		2	2019	https://escoops.cr gistorius.v10.1200/ JCO-2019.37.15_g unoi.109 abstract	356 VA patients newly diagnosed with NSCLC	Case-control study	Precision, recall, and area under the ROC curve (AUC) (5-fold CV)	cross-validation	In ready diagnosed MCC, pierest. Our predictive model for 1-year unread achieves in needy-diagnosed MCC, pierest. Our predictive model for 1-year unread achieves in needy-diagnosed MCC, pierest. Our predictive model for 1-year unread achieves (miles predictive pre
SS Finoschalds; F and Rezasian, I and D'Agnillo, M and Porter, L and Rueda, L and Ngam, A	An Integrative Approach for Identifying Network Biomarkers of Breast Cancer Subtypes Using Genomic, Interactomic, and Transcriptomic Data	J Comput Biol	24	8 751	-766 2	2017 Canada	http://dx.doi.org/10 108/licrots 2017.0 010 article	"We have used the METABRIC data set (Curtis et al., 2012), which contains the copy number values and GE levels of 2000 primary breast tumors with long- term clinical follow-up"	Tumour subtype categorization	AUC, sensitivity, specificity (10-fold CV)	cross-validation	the better prosible prient transment and response to therapy. Concer releave his unless are subserbeed for functionally resided and the "two first concert" to perform functions associated with a temoriganic. We propose a manifeste learning interfaces associated with a temoriganic way propose a manifeste learning interfaces associated with a temoriganic way propose a manifeste learning expensive transment that can be used to identify reventive between the second proposed and the temperature of the second proposed and temperature of the second proposed and the second proposed a
Fong, F and Bar, H Y and Shedden, K and Salya-Cork, K and Oulletts, P and Campagos, F and Meloc S6. A and Molek, S and Shaknoorch, R	Epigenetic profiling of primary CLL reveals novel DNA methylation-base, Lusters and novel mechanisms of leukemogenesis	ed Blood	120	21	2	2012	https://doi.org/10.1 160/absod.v/120.21 3877.5877 abstract	*DNA methylation of over 240 patients with CLL*	Case-control study	AUC (35-fold CV)	cross-validation	world with many 1,1000 merc cause diagnosed every year in the U.St. Vey hoptophescaried trut Def in New 1, 1000 merc cause diagnosed every year in the U.St. Vey hoptophescaried trut Def in New 1, 1000 merc cause diagnosed and very year in the U.St. Vey hoptophescaried trut Def in New 1, 1000 merc cause diagnosed and very year in the U.St. Vey hoptophescaried trut Def in New 1, 1000 merc cause diagnosed and very year in the U.St. Vey hoptophescaried trut Def in New 1, 1000 merc cause diagnosed and very year in the U.St. Vey hoptophescaried truth merc that merc than on the cause of the control of the cause of the ca
57 Gal, O and Auslander, N and Fan, Y and Meerzaman, D	Predicting Complete Remission of Acute Myeloid Leukemia: Machine Learning Applied to Gene Expression	Cancer in Informatics	18		2	2019 USA	http://dx.doi.org/10 _1177/1178995119 835544 article	473 bone marrow specimens from 473 patients	Case-control study	AUC (5-fold CV + test set)	cross-validation + test set	characteristic (PCC) curves. Using the top 7's goes from the Asserter neighbors against (P.CKH) on other C.7's yielded the best area-used refundance (PCC) care against (P.CKH) on other C.7's yielded the best area-used refundance (PCC) care against (PCC) care against (PCC) of the control of the control (PCC) care against (PCC) (P
58 Gamberger, D and Lavrac, N and Zelezny, F and Tolar, J	Induction of comprehensible mode for gene expression datasets by subgroup discovery methodology	ls J Biomed Inform	37	4 269	-284 2	2004 Croatia	http://dx.doi.org/10 .1016/j.jbi.2004.07. 007 article	the approach was applied to multiple cancer microarray datasets with > 50 samples per group in total	Case-control study	sensitivity, specificity, precision (training/test set split)	training + test set	religible to the color part of potentials in other treatments, and offs socialize particular color particular problems, the paper from the flashibility of this special particular problems, the paper forms the flashibility of this special particular part
59 Gaih, Hard Zheng, Z and Yue, Z and Liu, F and Zhou, L and Zhou, X	Evaluation of serum diagnosis of pancreatic cancer by using surface- enhanced laser desorption/ionizati- time-of-flight mass spectrometry	on Int J Mol Med	30	5 10	1-1068 2	2012 China	http://dec.doi.org/10 _98920pmm_2012.1 	serum samples from 132 patients with PCa and 67 healthy controls	Case-control study	sensitivity, specificity (leave one-out cross-validation)	cross-validation	The lating unface withward later desoption/fortable than ordigit must concerning vi(IEID-1764). Support veter manager (MM) analysis of the spectra was used to generate a predictive algorithm based on protein that were maintainly was a predictive algorithm based on protein that were maintainly was a predictive algorithm based on protein that were maintainly was a predictive algorithm based on protein that were maintainly was a predictive algorithm based on protein that were maintainly was a protein that the protein that were maintainly was a protein that the pro
Garl, M.J. and Klose, C. and Surma, M. A. and Tramandez, C. and Malander, O. and Malandez, S. and Gib Borodulin, E. and Haveslinos, A. S. and Salomas, V. and Ronees, E. and Cannotzes, C. V. and Simons, K.	Machine learning of human placeman lipicomes for obesity estimation in large population cohort	a Plos Biology	17	10 25-	25 2	2019 China	http://doi.org/10 1.137 Southerd etion. 30002443 article	Samples of the FMRSX 2012 under newton Egiption. Bipliothic measurements (1,141 according visited individuals) of which 1,061 were used	Case-control study	Required of obesity indicator variables (for repeated 3-0 field Cr)	cross-validation	The control of the co
61 Gong, I Yand For, N S and Huang, V and Boutros, P C	Prediction of early breast cancer patient survival using ensembles of hypoxia signatures	PLoS One	13	e0; 9 e0;	04123- 04123 2	2018 Canada	http://dx.doi.org/10 1971/forumit.com/ 0204123 article	1,564 early breast cancer patients	Cases only (survival prediction)	AUC, accuracy, sendibuty, specificity (10-fold cross-validation + text set)		prepreciating method and the risk dissillations from each were incorporated in a proposation of the control of
£3 Graim, Kand Fried, V and Houbhas, K Eard Stuart, J M	PLATYPUS: A Multiple-View Learnir Predictive Framework for Cancer D Sensitivity Prediction		24	136	-147 2	2019 USA	https://erew.ncbi.nl m.nh.nor/mediati cless/PMC8417892/j.article	At the time of download the Cancer Cell Line Encyclopedia (CCLE) contained genomic, phenotype, clinical, and other annotation data for 1,037 cancer cell lines	Cases only (drug sensitivity prediction)	AUC (cross-validation + external test set)	cross-validation + external cohor validation	We introduce a multi-view multi-relevanting strategy called PATPIS that both the New from multiple data counters that are all early the first generated receives the counter of the counters that are not as counters that are all early into the counters that are not as counters that are not expected as the counters that the counters that are not expected as the counters tha
Grofemund, V and Pradit, PF and Querin, G and Delbot, F and to Chat, G and Pradit-Physis, F and 63 Bods, P	Machine learning in amyotrophic lateral sclerosis: Achievements, pitfalis, and future directions	Frontiers in Neuroscienc e	13		2	2019 France	http://dx.doi.org/10 3389/hvs.2019.0 0135 article	review (not applicable) "The first dataset we used was collected from Genomic of Drug Sensibility in	raviaw			variables, bodogical, and recomminging data. These models also deep raisent with a second control of the contr
64 Guan, N Nand Zhao, Yand Wang, C Candli, J Qand Chen, Xand Pios, X	Anticancer Drug Response Predicti in Cell Lines Using Weighted Graph Regularized Matrix Factorization	Molecular on Therapy - Nucleic Acids	17	16-	-174 2	2019 China	http://dx.doi.org/10 _1016/j.cmtn.2019 _05.017 article	Cancer project (release-5.0, https://www.cancerrxgene.org/downloids), including 652 cancer cell lines, 135 drugs, and 70,676 known response values. The second dataset was collected from the CCLE (https://portals.broadinstitute.org/ccle) which contains 23 drugs and 491 cell lines with 10,870 known responses*		- Pearson correlation coefficient (PCC), root-mean-square error (BMSE), PCCs; and BMSEs averaged over all drugs (10-fold Cr)	cross-validation	"In this work, we presented a rowell method to utilize weighted graph regularized mann factoration (MOMF) for efforms an extraction of the control of the co
65 Guo, 5 and Guo, D and Chen, Land Jiang, Q	A centroid-based gene selection method for microarray data classification	J Theor Biol	400	32	¥1 2	2016 China	http://dx.doi.org/10 :1016/j.pts/2016.03 .034 article	multiple microarray datasets for different cancers with > 50 sample per group in total were used	Case-control study	accuracy + standard deviation (repeated 5-fold CV)	cross-validation	Let depend the control of the contro
66 Guo, Y and Yu, H and Chen, D Q and Zhuo, Y Y	Machine learning distilled metabol biomarkers for early stage renal inj		16	1 10-	10 2	2019 China	http://dx.doi.org/10 .1007/s11308-019- 1824-0 article	"Serum samples were collected from all participants, including S87 CKD patients (CKD1 = 120, CKD2 = 104, CKD3 = 110, CKD4 = 119, CKD5 = 134) and 116 age-matched normal healthy controls."	Case-control study	AUC (10-fold CV)	cross-validation	"Nodiging in the decision of assemble and the second of a second o

67 Mars, H 9	Nonnegative principal component analysis for mass spectral serum profiles and biomarker discovery	BMC Bloinformati cs 11	\$1.51	2010 USA	bits like dok out 10 1180/1471-2005. 11.46.431 article	4 different case/control cancer MS serving profile datasets were analyzed with > 50 samples per group on average	Case-control study	accuracy (LOOCV and 300 trials of 50% holdout cross validations (HOCV))	cross-validation	In this facility we develop a nonequiptive principal component salvalvia algorithm and present and
Hao, Y Y and Duh, Q Y and Kloos, R T and Babiarz, J and Harrell, R M and Traweek, S T and Kim, S Y and g	Identification of Hurthle cell cancers: solving a clinical challenge with genomic sequencing and a trio of machine learning algorithms	Brnc Systems Biology 13	14-14	2019 USA	http://dx.doi.org/10 _11859.12918-019- 06995.g article	318 samples, including 119 Hürthle cell- negative and 199 Hürthle cell-positive samples	Case-control study	AUC, sensitivity, specificity (10-feld nested CV)	cross-validation	Sequencing Classifier (GSC, The accurate algorithmic depotion of this complex biological system analysis first than adapters such as demands or a demand of control of the complex biological system analysis first than adapters such as demands or a distriction performance, specificity among thirthis elders present this is districted, and the second of
d	A computational method to differentiate normal individuals, osteoarthritis and rheumatoid arthritis patients using serum biomarkers	JR Soc Interface 11 9	20140428 20140428	2014 Canada	http://ldx.doi.org/10 1008/hwl/2014.04 28 article	"normal individuals (normal, $n=100$), patients with osteoarthritis (DA, $n=100$), and rhoumatoid arthritis (RA, $n=100$)"	Case-control study	accuracy, sensitively, specificity (training, validation and test set split)	cross-validation + test set	offine-retained between serum amplies of patients with OA, a Supposed RA patient comparation cohorts of monitoral control cont
	Bayesian methods for proteomic biomarker development	EuPA Open Proteomics 9	54-64	2015	http://dx.dei.org/10 .1016(j.eugret.201 5.68.001 article	review (not applicable) "Two datasets were used in this study. The first was the Cancer Genome Atlas (TCGA) methylation brain lower grade	review			robut to own fitting than other approaches, appointing when the number of rampies used for discovery is relatively routful. In the robust per control to severifying than other approaches, exposizily when the number of rampies used for discovery is relatively routful. In the robust per control to Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and demonstrate some of the advertages of using a Bayesian inference and d
71 Hira, 2 M and Gillios, D F p	Identifying significant features in cancer methylation data using gene pathway segmentation	Cancer Informatics 15	189-198	2016 UK	http://dx.doi.org/10 _1127.CN ESIBMO article	gloma (LGG) statuet (http://cncepnos.nhi.gov). In this, there are 170,203 probes and in total there are 170,203 probes and in total are the control of the c	Case-control study	AUC, accuracy (stretified 10-fold CV)	cross-validation	to diagnose whether a patient's cancer will respond to a proposed treatment. Anti-vision profiles and control information which such practications could be distripted and the patient of the patient of the distripted and the patient of the patient o
	Machine learning SNP based prediction for precision medicine	Frontiers in Genetics 10		2019 Australia	htm/lds dei ora/10 3389/lgene 2019. 00267 article	review (not applicable)	review			with large nee population distants with high-quality phenotyping at different stages in the lifectioning, national bearing mode in secretary period of casking included disease risks with high percision. Notable, machine learning prediction that included provide cost effective and practice healthcost with great efficiency." Outdoor included included provide cost effective and practice healthcost with great efficiency." The servival rate of particular control provided cost effective and practice healthcost with great efficiency." The servival rate of particular casks to find which also of particular with great efficiency. The servival rate of particular casks in the find was deep service healthcost with great efficiency." The servival rate of particular is the levent among those as in the service and particular in the levent among those as in the service and particular in the levent among those as in the service and particular and
	Possible detection of paintrastic cancer by plasma protein profiling	Cancer Res 65 2.	10613- : 10622	2005 Japan	http://doi.org/10 .1360008- 8472.xms05.51851 article	71 payroxis cover patients and 71 healthy controls	Case-control study	AUC, sensitivity, specificity (LOOC), enternal test set)	cross-validation + external coho validation	common collid tumors, and darky distriction is one of the most fassible means of improving outcomes. We compared plasma professions between personal contemporary of the contemporary of t
	Genomics analysis of gene expression profiles demonstrates a distinct ARDS signature Machine learning predicts individual	Journal of Respiratory and Critical Care Medicine 197		2018	https://www.atsijour nals.org/soldstart/0 1184/aiscome confuseroe 2018.1 921.1 Meetinghater acts A7533 abstract	training cohort (n = 318, 75%), validation cohort (n = 105, 25%)	Case-control study	sensitivity, specificity (training/test set split)	training + test set	and without ARDS and applied machine learning techniques to a discovery color of an interest of the property
Huang, C and Clayton, E A and Matyunina, L V and McDonald, L D and Benigno, B B and Vannberg, F c 75 and McDonald, J F t	cancer patient responses to therapeutic drugs with high accuracy	Sci Rep 8 1	16444- 16444	2018 USA	.1038941598-018- 34753-5 article	175 cancer patients	Cases only (treatment response prediction)	accuracy, sensitivity, specificity (LOOCV)	cross-validation	from the gene-expression profiles (RNA-seq or microarray) of individual patient tumors. The models were found to predict patient responses with >80% accuracy." tumors. The models were found to predict patient responses with >80% accuracy."
Huang, Y C and Chung. H H and Dutslewicz, E P and Chen, C L and Haish, H Y and Chen, B R and Wang. S 76 M Y and Huu, C C	Predicting Breast Cancer by Paper Spray Ion Mobility Spectrometry Mac Spectrometry and Machine Learning	s Analytical Chemistry 92 2	1653-1657	United 2020 States	http://dx.doi.org/10 1.02/time.menisha m.08/05/86E article	breast core needle biopsies: 29+377 boning, 15442 malignant 7-44 cases, including 18 KD patients, who were tested both prior to receiving at least 3 weeks after IVIQ in receiving at least 3 weeks after IVIQ in receiving at least 3 weeks after IVIQ in week bookerved in the Illumina. HumanMethylatorol 50 Beachlot settle.	Case-control study	accuracy, sensibility, specificity (cross-validation + external validation)	cross-validation + external coho validation	Net we demonstrated that by utilizing paper gary in includes mass spectrometry (FS-MS) capited with mid-dispurations where the model has predictive metabolic and lipidomic profiles of models by extremely (FSMS), predictive metabolic and lipidomic profiles of models beganised and the debt and extremely profiles of models beganised and the debt and extremely profiles of models beganised and the debt and extremely profiles of models beganised and profiles of models beganised and the debt and profiles of models beganised and profiles of models and profiles of mode
P Huang, Y H and Kuo, H C and U, S C and Cal, XY and Uu, S F and Kuo, H C b	HAMP promoter hypomethylation an increased hepcidin levels as biomarkers for Kawasaki disease	d J Mol Cell Cardiol 117	82-87	2018 England	http://dx.doi.org/10 10166.y/moc.2018 02.017 article	for their CpG markers. The remaining cases consisted of another 92 KD patients and 113 controls that were used for validation by pyrosequencing? Blood samples were collected from 64 cases of SLE 30 cases of their market.	Case-control study	AUC, sensitivity, specificity (5-fold CV, external test set)	cross-validation + external coho validation	processoring, We performed a genetic functional study using Lecterars access. A proport vector manning SVM discussification model was doughted to learning. A proport vector manning SVM discussification model was doughted to learning the value and central subjects. We developed a SVM discussification model with a 950 PM and central subjects. We developed a SVM discussification model with a 950 PM and central subjects. We developed a SVM discussification model with a 950 PM and central subjects. We developed a SVM discussification model with a 950 PM and subjects and central subject. We developed a SVM discussification model with a 950 PM and subject and central subject. We developed a SVM discussification model with a 950 PM and subject and central subject. We developed a SVM discussification model with a 950 PM and subject and central subject. We developed a SVM discussification model with a 950 PM and subject and central subject. We developed a SVM discussification model with a 950 PM and subject and subject. We developed a SVM discussification model with a 950 PM and subject and
Huang, Z C and Shi, Y Y and Cai, B and Wang, L L and Wu, Y K and Ying, B W and Feng, W H and Hu, C J s	[Promising diagnostic model for systemic lupus erythematosus using proteomic fingerprint technology]	Sichuan Da Xue Xue Bao Yi Xue Ban 40 3	499.503	2009 Chin-	https://europepmc. org/wrick/med/19 627014	arthritis (RA), 30 cases of Sjogren's syndrome (SS), 25 cases of systemic sclerosis (SSc), as well as 83 healthy controls (segregation SLE from non-SLE)	Caro control ctu4:	sensitivity, specificity (training/test set split)	training + test set	a cluster pattern segregating SE from non-SE with sendotivity of 91% and specificity of 92%. The discriminatory diagnostic pattern correctly identificial SE LA senditivity of 91% and specificity of 92%. The discriminatory diagnostic pattern correctly identificial SE LA senditivity of 95%. The discriminatory diagnostic pattern correctly identificial senditivity of 95% for the blinded text were obtained when comparing SE vs non-SE? —
Notes: B and Nodgrow, K and Malle K and Maller, W and Blassend, M and Mars, O and Hauser, J and Henging Java Democrael, M2 and Soory, D and Staff, D and Former, A and Lenk, C	Combining clinical and genetic variables to predict antidepressant treatment response: A machine	Ban 40 3 European Neuropsych opharmacol		zuus China	http://ikk.doi.org/10/ 1016/j.suronsuro.meeting	countries (segregating SLE from non-SLE)	Cases only (treatment		craining + tost set	In this case, we see explainted models to set the predictive adulty of a combination of Genome Washing with condection promptions (DNA). In this case, we see explainted models to set the predictive ability of a combination of Genome Washing with condecting promptions (DNA). The condection of the promptions of the company of the comp
	learning approach	ogy 27	\$353-\$354	2017	2015.09.010 abstract	430 patients with unipolar depression		Accuracy for remission prediction (10-fold CV)	cross-validation	ecyptobram resulting extractions are a financial resulting and resulting extractions are resulting extractions are resulting extractions and resulting extractions are resulti

ioni, H and Sahoh, M and Salamons, K and Salamonio, K and Salgeou, D and Kasii, H and Adhizawo 80 and Myszawa, K and Taledda, S and Masoupama, K and Yoshimura, K	Lipidome-based rapid diagnosis with machine learning for detection of TGF. British IX. Betta signaling activated area in head. Journal of and neck cancer	10-10 Japan 0723cs article	A total of 740 and 90 mass spectra were obtained from TGF-B-unstimulated and - stimulated 996CC cells, respectively Case control study acuracy (LOOCV)	"We established a rapid diagnostic system based on the combination of probe electrospany institution in easi spectrometry [FE]. Still just an achieve base in great the spectrospan in the combination of probe electrospany institution in easi spectrometry [FE]. Still just an institution is supported by the combination of probe electrospany institution in easi spectrospany. The spectrospan is spectrally probe that is a support easily larger institution in easily spectra, and in the combination of FE, electrospany in electro
hearcir, M.M. and Magna, B.W. and Swershlov, Y. and Coven, M. and Reichelderfor, M. and Pichhard El J. and Sussman, M.R. and Kennedy, G.D.	Noninvasive Detection of Colorectal R, P. Cardinomas Using Serum Protein Biomarkers J Surg Res 246	United Life is 2019 80 article 160-169 2020 States .006 article 2020 80 article	"Blood was drawn from Individuals (n = 231) before coloroscopy or from patients with momentation Circl (n = 50)". Case-control study AUC, sensitivity, specificity (training / test set split) and classes." 19th amplies contain the data measured by MULCS. Less control study AUC, sensitivity, specificity (training / test set split) the data measured by MULCS. Less control study AUC, sensitivity, specificity (training / test set split) when the data measured by MULCS. Less control study AUC, sensitivity, specificity (training / test set split) when the data measured by MULCS.	and an other control research, and purpose in a higher to be in the desired of the size of
\$2. Julia, A and Meller, N	Interpretable per case weighted execution exhibit for cancer SMC execution exhibit for cancer Genomics 17	100-100 ass and 5 100-1001 2016 Germany 2016.2 2016 Germany 2016.2 2016 Germany 2016.2	samples contain mRNA data massands by Not G133 area," TREA dataset. This data set michales 953 samples with the data set michales 953 samples with the data set samples indented as having metastazativel (8 samples). Hence the data set samples concerding for "Tumoria" sour", "affected manifely lymph notes," ("Lagy", and "samples reception". Case only Vol. B. Converting settlification). AAC (praining / text set quick)	"Molecular measurements from carear particular such as game appreciation and DMA methylations can be influenced by several external factor. I manded leaves to take potential bases in the data sinds outcome, this can lead for a problem when they particular particul
Jones, 13 and Wilcon, B E and Bens, R W and Babbar, N and Boragine, G and Burrell, T and Christian and Coose, L1 and Cur, P and Dillon, R and Fatts, S N and Eas, A and Preston, F and Schrickings BS S R and Sanc, H and Smith, W F and You, J and Hilli, W D and Again, D B and Blume, J E	ust, Identified by Multiplex Targeted Mass Colorectal	186 United 1206 discrept 5 120	the present cludy used 274 individual patient blood planns samples, 123 with bloop scorlines decoveral energy and an additional state of 137 age- and gender-matched controls. AUC, sensitively, specificity (cross-validation + external test set)	data as a discovery or life disease case and 60 control case), the size for the favore data as a discovery or life disease case and 60 control case), the size for the favore data as a discovery or life disease cases and 60 control case), the size for the favore data as discovery or life disease cases and 60 control case), the size for the favore cases are control cases, the size for the case of the
Jarmeister, P. and Bockmay, M. and Sangerer, P. and Bockmaye, T. and Trow, D. and Montavon, G. d. Voltered, C. and Armold, A. and Teichmann, D. and Bressem, K. and Schuller, U. and von Laffert, M. 184 Muller, F. R. and Capper, D. and Elasuchen, F.	Machine learning analysis of DNA methylyation profiles distinguishes and primary lung oquamous cell Science arctionate from head and neck Translationa metatatases I Medicine 11 5	http://dx.doi.org/10 1220cc/secologic 1220cc/	Add parients with a biology of primary MISC and a synchronous or metachronous squamous lung tumor Case control study AUC, accuracy (5-fold CV + external test set)	importance, the cost possible in most case; with convent diagnostics. To address this, was performed DAM immelphicum of primary tumous and artisted forese different machine learning methods to distinguish investations with convention of the case in a validation confort of 279 patients with milk of the case in a validation confort of 279 patients with milk of the case in a validation confort of 279 patients with milk of the case in a validation confort of 279 patients with milk of a control of the case in a validation confort of 279 patients with milk of a control of 279 patients with milk of a control of 279 patients with milk of a control of 279 patient with milk of a control of 279 patients with milk of a control of 279 patient with milk of a control of 279 patient with the control of
85 Karrispour Park, A and Epperson, L E and Hunter, L E	A survey of computational tools for downstream analysis of proteomic and other omic datasets Genomics 9	11-11 2015 USA 0000-2 article	review (not applicable) review	dustring and validation, interpretation, and generation of biological information from experimental data. We experimental data set and experimental data set and to experimental data set and experimen
Eban, S R and Mohav, H and Lix, Y and Batchubun, B and Gohl, H and Al Rigal, O and Metablowy. 86 and Cox, B J and Gunderson, E P and Wheeler, M B	The discovery of round predictive blomatiers and only-stage protophysiology for the transformation of the transformation of the stage o	titas litria scripper constitutinis 13 20 20 20 Canada 25022 250 250 250 250 250 250 250 250 2	55 incident case matched to \$5 sun- case control participants Cise-control study AUC, accuracy, sensibility, specificity (45-fold cross-validation)	used a well-characterised prospective color of women with a history of COM programs, and of whom now common histories, were supported to the color of women with a history of COM programs, and of whom now control of the developing, and of whom now control of the developing and the develo
Ibhasial, R D and Cidffi, C E and Califfany, S A and Krasinska, A M and Alazraki, A and Keight-Scott, and Cleeron, R and Cattill-Leon, E and Jones, D P and Perport, B and Caprin, S and Santono, N at 87 AMI, A and Vo., M B	nd Panel for Pediatric Nonalcoholic Fatty Communica	United <u>http://isk.doi.org/10</u> 10 1311-1321 2019 States 11002-Base4.117 article	subjects with NARLD (in = 222) and without NARLD (in = 222) and Case-control study AUC (training set: 2/3 of data, text set 1/3 of data)	which had an area under the recover operating characteristic cone (pURO) of 16 % which had an area under the recover operating characteristic cone (pURO) of 16 % searching of 17%, and quantification model are set of 15% of externing MED cases. A accord discrification model area developed using the homeostast model assessment of inflam restances undertaked for the MED. Similarly, the place professing discrification model area reaction forces, which had an AUROC of 0.5%, constitution of 15%, and quantified and the search of 15% of 15%, and quantified and the search of 15% of 15%, and quantified and the search of 15% of 15%, and quantified and the search of 15% of 15%, and quantified and the search of 15%
Kim, J and Da Rosa, J C and Lee, J and Tomalin, L and Lower, M A and Fize, L and Bernstein, G and 88 Yudez, H and Wolk, R and Knuger, I G and Subrez-Farikar, M	Precision medicine in psoriasis: Machine learning and proteomics join Experiment forces to device a blood-based test al to predict response to tofacilinib or Dermatolog Etanercapt in psoriasis patients y 25	United <u>http://ide.dec.org/10</u> meeting 49-50 2016 States 1111sms 13000 abstract		data datamed using a promitiny extension assert [-] to develop a blood based text to data datamed using a promitiny extension assert [-] to develop a blood based text to prefect response to dutafication of the assertance plan posting assertion. In belatic near application, using only per treatment data, was to be lest performer among the method evaluatant, with however, the least posting and the process of the least performer among the method evaluatant, with however, the least posting and the least posting and the least performer among the method of the least performer among
89 Kim, M and Ch, 1 and Ahn, J	An improved method for prediction of cancer prognosis by network learning. Genes 9	http://dx.doi.org/10 3000/amen/10101 111	"First, we downloaded grow mRM data, CDV data, DM methylation data, SIP data, and clinical state for PASA, BECA, B	prognosis. The proposed method sporties the candidate prognosis are provided by graph harming using the preservation abservation interests (California) and cause of the proposed method sporties are provided interests (California) and cause of the proposed method and provided interests (California) and cause of the proposed method and and cause of the proposed method and and cause of the proposed method and and cause of the proposed method allowed better prediction accuracy than del existing methods."
50 Kim, S and Rhong, J H and See, J and Koo, J Y 50 Kim, S and Rhong, J H and See, J and Koo, J Y 51 Kim, S and Lin, C W and Years, G C	Meta-analytic support vector muchine BioData for integrating multiple onics data. Mining 10 Meta/TSP a meta-analytic top succept again entered for reduct cross-analytic for the control of the control	1 2017 South Korea 1128 article 1 2017 South Korea 1128 article 1 3964-979 2016 South Korea 12001 article 13 3964-979 2016 South Korea 12001 article article 14	boary outcome (i.e., case and control) and bears caree requested or prefiles (CCGA) including mRMs, clay y unable variation (CCM) and degenetic DMA methylation. (http://ccm.org.generian.inl.gov/, 300 (http://ccm.o	"Vis process a make analytic resport vestor make fields 2018 that can accommodate making its models of their commons general accountment of the second of th
con, 5 year daggers, just instance, pard teaming, pard Papers, M and Stroky, N and Ton, 5 and Andrews, 1 and Ton, 5 and Carloy, T and Stroke, M and Wolfers, P is and Markes, 5 D and Carloy, T V and 2 Markes, 1 and Editions, 6 Li and Faghing, G and Remody, G C Markes, 1 and Carbon, 1 and Carbon, 1 And Paghing, G and Remody, G C	Classification of usual interstitial pneumonorial in patients with interstitial pneumonorial machine sanchine sanches	6 473-482 2015 Mexico 2020 13/201405 article	"32 swiged lang Boopies from Ris patient, 55 spaces were identified by the space panel around international personnels, flow as sarvadous, four as respiratory benchmonts, flow as substituted on the second of the s	"Signature parameters" prices in a pregnative (fourth-large desains that determines the prices of th

93 Kin, Y and Biomejer, Tand Zevart, W and Wessell, LFA and Vic. D.J	Genomic data integration by WON- PARAFAC identifies interpretable factors for predicting drug-sensitivity in vivo	y Nat Commun	10	1	5034-5034	Netherland 2019 s	http://de.doi.org/10 d	1815 genes by 935 cell line	Cases only (drug sensitivity prediction)	AUC (10-fold CV)	cross-validation
S4 Kin, Y R and Kin, D and Kin, SY	Prediction of Acquired Taxane Resistance Using a Personalized Pathway-Based Machine Learning Method	Cancer Res Treat	51	2	672-684	Korea 2019 (South)	http://de.doi.org/10 .4143ect.2018.137 article	more than 50 samples per group for most human camer cell line datasets considered	Cases only (drug response prediction in vitro)	AUC (100Cr)	cross-validation
Kingor, T and Kili, S and Alani, Y Y and Yaman, A and Engouse, S B and Stan, M and Turan, S and S Halar, G and Signingh, M S and Bereker, A and Guran, T	Simplifying the interpretation of steroid metabolome data by a machine-learning approach	Hormone Research in Paediatrics	91		128-128	2019	http://dic.doi.org/10 _1169000601808 article	500 healthy controls and 427 treatment- naive children with a disorder of adrenal steroidogenesis	Case-control study	sensibility, specificity (10 fed cross-validation)	cross-validation
Electron, it and Meuronisis it and Meuronisis fall of Chans, It and Walescotts, M. and Yolobia, T. and Honey, M. and Yolobia, T. and Honey, M. and Yolobia, T. and Honey, M. and Kishkawa, D. and Tempore, R. and Essenbirro, N. and Nichkawa, E. and Horrito, A. and Nichkawa, E. and Horrito, A. and Nichkawa, E. and Horrito, M. and Nichkawa, E. and Horrito, A. and Nichkawa, E. and Horrito, M. and Nichkawa, E. and Nichkawa	Genome-wide methylation analysis using the digital restriction enzyme analysis of methylation for tratification of patients with juvenil myelomonocysis leukemia	e Blood	134			2019	http://kix.doi.org/10 .1182/bboos-2010 meeting 127792 abstract	99 children (67 boys and 32 girls) with JMML.	Case-control study	accuracy (braining / foot set split)	training + test set
97 Kong, A and Associat, R	Binary Markov Random Fields and interpretable mass spectra discrimination	Statistical Applications in Genetics and Molecular Biology	16	1	13-30	2017 Germany	http://kk.doi.org/10 1515/sagetb-2010- 0019 article	"A dataset of 238 MALDI colorectal mass spectra and two datasets of 216 and 253 SELDI ovarian mass spectra respectively were used to test our approach."	Case-control study	accuracy (80000)	cross-validation
98 Kravczok, z and Lukoczok, T	The feature selection bias problem in relation to high-dimensional gene data		66		63-71	2016 Poland	http://dx.doi.org/10 1016/i.schrast.201 5.11.001 article	seven microarray datasets with > 50 samples group for multiple datasets were used	Case-control study	accuracy (double LOOCV)	cross-validation
Krittanawong, C and Bomback, A S and Baber, U and Bangalore, S and Messerli, F H and Wilson Tan 99 W H	Future Direction for Using Artificial g, Intelligence to Predict and Manage Hypertension	Curr Hypertens Rep	20	9	75-75	2018 Poland	http://dx.doi.org/10 1007/s11906-018- 0875-x article	review (not applicable)	review		
Foo, C.H.S. and Pavilloi, S. and Loza, M. and Barchaud, F. and Tower, A. and Pavilloi, I. and Rossion, C. and 100 Wilson, S. and Sydamours, Y. and Stork, P. and Chung, E.F. and Adock, 1M and God, Y.	Asthma phenotypes from semi- supenrised machine-learning approach of brochal biopsy and brush transcriptomics in U-biopred	American Journal of Respiratory and Critical Care Medicine	191			2015	http://www.abicor. rails.org/sicked/TO. 1.1164/aprocess. roof-arms.2015.1 21.1.1.5arms.2proc. abb.4.2202. abstract	"Subjects with moderate-to-severe asthma recruited in the U-BIOPRED study underwent fiberoptic bruchscup (or brunchhausp for brunchhalbiopsy (91) and brush (105) samples"	Case-control study	accuracy (cross-validation)	cross-validation
101 Kursa, M B	Robustness of Random Forest-based gene selection methods	BMC Bioinformati cs	i 15		8-8	2014 Poland	http://doi.org/10 .1188/1471.2505. 15:8 article	4 microarray datasets were used, one contained > 50 samples per group	Case-control study	error-sate (training / fest set split)	training + test set
Kowakara, H and Iwabuchi, A and Soya, R and Formato, M and Ishizaki, T and Tsuchida, A and 102 Nagakawa, Y and Estumata, K and Seglemoto, M	Salivary metabolomics for colorectal cancer detection	Annals of Oncology	30		v46-v46	2019	http://dx.doi.org/10 109/Banonoimdz 209-058 article	"231 subjects with CRC, 99 subjects with polyps, and 2272 subjects with healthy controls"	Case-control study	AUC (training / text set split)	training + test set
Lacron Trik, M and Kempousky Harmon, T and Yulle, C and Hedjall, L and Lamanns, S and Trouth, I and Caller, F and Tileron, T and Toure, G and Lacron, M Y and Le (3) Berney Artists, Y	Fuzzy logic selection as a new reliable tool to identify gene signatures in breast cancer - the INNODIAG Study	e Laboratory Investigatio n	93		51A-51A	2013	http://likx.doi.org/10 1038/habrevest 20 meeting 13.14 abstract	7 breast cancer microarray datasets + 151 consecutive invasive breast carcinomas	Case-control study	sensitivity, specificity, error rate (training / text set split)	training + test set
Lai, A and Pance, R and Marjanovic, M and Waller, M and Frantes, E and Eago, D S and Homer, W (104 and Brutzwork, L) and Miller, M H	A gene expression profile test that D distinguishes ovarian from endometrial cancers	Journal of Clinical Oncology	30	15		2012	https://www.ncts.el m.sih.gov/pmcsard meeting clewPMCS326651 abstract	75 metactatic, poorly differentiated or undifferentiated primary FFPE tumor specimens	Differential diagnosis prediction	AUC (braining / text set split)	training + test set

ottain new gene signatures. To validate threat gene signature, we designed probes
for the salected gene on therelooper continuement of the salected gene on therelooper continuement or the salected gene on therelooper continuement or the salected gene on the salected gene on the salected gene on the salected gene on the salected gene of the sal

"We introduce Weighted Orthogonal homogative parallel factor analysis (WOIN-PARAFAC), a data sweepstow method that identifies spars and interpretable factors.

White PARAFACA is data sweepstow method that identifies spars and interpretable factors.

White PARAFACA is data integration method that identifies spars and interpretable factors.

White PARAFACA is data integration method that identifies spars and interpretable factors.

White PARAFACA is data integration method that identifies spars and interpretable factors.

White PARAFACA is data integration method that identifies spars and interpretable factors.

White PARAFACA is data integration method that identifies spars and interpretable factors.

White PARAFACA is data integration method that identifies spars and integration of the factor is based on the integration of the factor is designed in the factor is

prediction model for ATE unity generouslated partnersy and expellutation functions promption and promotion of partnersy and expellutation functions (prodict forminalization and partnersy) (and lossed partners) and transport (prodict forminalization and partners) and partnersy and production for production of production and partners and production for production of production of production and production productions and interpretation of analytical resistant. We have beginned to a understand production of analytical resistant with a featured from production production of analytical resistant. We have beginned to a understand production and production productions. We have beginned and interpretation of analytical resistant, we have been deep and interpretation of analytical resistant. We have breat the performance of this alignethms using our analytical resistant. We have breat the performance of this alignethms using our analytical resistant. We have breat the performance of this alignethms using our analytical resistant. We have breat the performance of this alignethms using our analytical resistant. We have breat the performance of this alignethms using our analytical resistant. We have breat the performance of this alignethms using our analytical resistant. We have treat the performance of this alignethms using our analytical resistant. We have treat the performance of this alignethms using our analytical resistant was a section of the section formers and the analytical resistant of analytical resistant was a section of the performance of this alignethms using our analytical resistant was a section of the section formers and the analytical resistant was a section of the section formers and the analytical resistant of the section formers and the analytical resistant was a section of the section formers and the analytical resistant was a section of the section formers and the analytical resistant was a section of the section of the analytical resistant was a section of the section for

surface disorders of advoration and groundst strondospenses. We have implemented a machine bearing signifer the 6 sites on signifer degree of the support of the strong of the surface of

Ribustrating the innovative All approach for petertial prediction of any stages of hypertension. Addition, we view on enging reason and factor implication of properties and the properties of the control of the control of the control of the control of medicine."

Administ is a temperature disease underfield by inferrit enforcement programs. A finespee of 1D digitative was associated with airway hyper responsioners, allego the phenotyped and making a semi-squared interfines and responsible control of analyses person profiles. Training of a game model for the cluster using an extra structure control objects and the control of the c

Lawton, K.A.and Brown, M.V. and Alexander, D. and U., Z. and Wulff, J.F. and Lawson, R. and Jaffa, M.a. 105 Milburn, M.V. and Rysis, J.A. and Blowsor, R. and Custlowicz, M.E. and Berny, J.D.	disease mimics Integrated machine learning pipelir for aberrant biomarker enrichment	Frontotemp oral Degener	15 5	362-370	2014 England	http://dx.doi.org/10 3109216786912 014308311 article	172 patients recently diagnosed with ALS, 50 healthy controls, and 73 neurological disease mimics The SLE compendium contained 15,497 gene expression measurements with	Case-control study	AUC, sensitivity, specificity (training / test set split)	training + test set		patients with ALS from those with disease minics. Litting all identified biochemicals detected in > 59% of all samples in the metabolomics analysis, samples was classified as ALS or minic with 65% sentitivity and \$1% specificity by ILSSO analysis (ALC of 10%). A subset panel of 32 candidate biomarkers classified these disposis groups with a specificity of 95%/featibility 85% (ALC of all.)? "" "Within a compendium of systemic lupus or sythematous (SLE) patients, we applied the lesseated mythin barraine involved for afformation for the control of the c
106 Le, TT and Blackwood, N O and Taroni, J N and Pu, W and Brotenstein, M K	mAB): characterizing clusters of differentiation within a compendiu of systemic lupus erythematosus patients	m AMIA Annu Symp Proc 2	018	1358-1367	2018 USA	https://www.ncbi.nd m.nh.gov/pmc/strl class/PMC8X71298/ article	observations from healthy control (n=160) samples, treatment-naïve SLE (n=1,290) samples, and SLE samples exposed to various treatments (n =126)	Case-control + treatment response	balanced accuracy (cross-validation + 20% hold-out test set)	cross-validation + test set	hypotheses, i-mAB fostered robust biomarker profiling among interdependent biological fosturer."	profile de novo gene expression futures affecting CDDQ, CDZ and CDDQ gree abternace. Utilities carefully aggregated accordiny data and leveraging a priori hypothesia; i-mell flostered robort biomarker profiling among interdependent biological floatures." "The identification of biomarker signatures in omics molecular profiling is usually performed to predict outcomes in a procinion medicine context, such as polision disclasses acceptablish, disprocis, proposoci, and treatment response. To identify these discusses acceptablish, disprocis, proposoci, and treatment response. To identify these discusses acceptablish.
Lodierce, M. and Vittroer, B. and Martin-Magnistris, M. L. and Scott Boyer, M.P. and Perin, O. and 107 Bergeron, A. and Freder, Y. and Droit, A.	Large-scale automatic feature selection for biomarker discovery in high-dimensional emics data	Frontiers in Genetics	10		2019 Canada	http://doi.org/10 339905pen 2019 00452 article	five microarray datasets were used, including datasets with > 50 samples per group	Case-control study	sourcey (ACC), believed over risk (BBI), Matthew's condition coefficient (MCC), are writer the curve (MCC), smallering, specificity, filed their Squared from (MCC), Correlation Coefficient (CC) (10-bit CV)	cross-validation	collection of amples and their associated characteristics, i.e., the biometeris (p. 2000) are expersion, protein without them, of those publicide along, blottlicker appeals; services are provided and protein a larger surger and protein bearings and protein and blometers for predicting categorists of continuous automost from highly and protein and prote	collection of simples and their associated characteristics, i.e., the biocontains (e.g., to present agreement of the collection of the co
Les, 5 S and Attorood, K and Roder, H and Asmelliath, 5 and Mayer, K and Katolyns, S and Olivera, 1008 and Roder, J and Grigorieva, J and Ohels, L and Iyer, R and Mahalingam, D	hepatocellular carcinoma	Cancer Research	79 13		2019	http://dx.doi.org/10 .1158/1538- 7445 SABCS18- 4530 meeting abstract	156 pts (97 HCC, 59 non-HCC healthy controls)	Case-control study	AUC (training and validation cohort)	external cohort validation	etiologies and Child-Pugh classification []. In independent validation, AUC for the test output prior to thresholding was 0.979, significantly better than AFP AUC 0.915 (P=0.001).*	(P40.001)." "We present a three-gene version of "relative expression analysis" (RXA), a rigorous and outcometic comparison with outlier approaches in a variety of concern trulier.
109 Lin, X and Afsari, B and Marchionni, L and Cope, L and Pamiigiani, G and Naiman, D and Geman, D	The ordering of expression among: few genes can provide simple cance biomarkers and signal BRCA1 mutations		10	256-256	2009 USA	http://dx.doi.org/10 _1186/1471-2106- _10.258 article	118 samples for BRCA1 breast cancers + three datasets used for the ER status cross-study validation, including a dataset with > 50 samples per group	Case-control study	accuracy, sensitivity, specificity (LOOCV, cross-study validation)	cross-validation	cancer and a cross-study validation for predicting [ER status. In the BRCA1 study, RCA vields high accuracy with a simple decision rule in tumors carrying mutations the expression of a "reference gene" falls between the expression of two differentially expressed genes, PPPICB and RNF14." "Ottoosraccoms is a common malignancy with high mortality and poor prognosis due	yields high accuracy with a simple decision rule: in tumors carrying mutations, the expression of a "reference gene" falls between the expression of two differentially expressed genes, PPP1CB and RNF14."
	A Four-Pseudogene Classifier Identified by Machine Learning Sen as a Novel Prognostic Marker for	ves Genes				http://dx.doi.org/10 330/Jeanes 10060		Cases only (survival			to lack of predictive markers. The aim of this study was to identify a prospectic presendagene signature of orteosarcoma by machine learning. A sample of 54 octeosarcoma picetiest? RNA-5aq data with clinical follow-up information was involved in the study. The survival-related pseudogenes were corrected and related signature model was constructed by one orgension analysis (unknike), lacks, and multivariate), Into study 125 survival-related pseudogenes were identified and a four-precedure (IPSI III.5114.114); Bod. Spilic (II.44-04.5); Spilic (II.44-0	"Oblicaciónico il a Common malgianico y leith high mortally sed quel prisposito. Journal officiale de la common de la common de la prisposito de la common del common de la common del la commo
110 Liu, F and Xing, L and Zhang, X and Zhang, X	Survival of Osteosarcoma	(Basel)	10 6		2019 China	414 article	94 osteosarcoma patients A total of 118 samples from the peripheral blood of females, including 47	prediction)	AUC (10-fold CV)	cross-validation	patients, and predicted prognosis with high sensitivity and specificity (AUC: 0.878)."	patients, and predicted prognosis with high sensitivity and specificity (AUC: 0.878)."
111 Uss, Land Liu, Yand Us, C and Zhang, Z and Du, Yand Zhao, H	Analysis of gene expression profile identifies potential biomarkers for atherosclerosis	Mol Med Rep	14 4	3052-3058	2016 China	http://dx.doi.org/10 3892/merr.2016.5 650 article	atherosclerotic and 71 non- atherosclerotic patients, was used for expression profiling.	Case-control study	AUC (5-fold CV)	cross-validation	The present study almost to Mortifly potential biomarker for atherosclerosis via analysis of gone expension optices. I, 19 the Egipathem was used to identify 11 biomarkers, whose receiver operating characteristic convert had an area under convert of 0.90, ledicating that it desired that 11 biomarkers were representation in critical to enable and an extra convertible convertible or the convertible co	
Us, M.C and Jameshid, A and Vener, O and Frields, A.P and Maller, M.C and Cases, G.and Annis, I sac Gross, S and Brodne, J and Millor, M. and Schollenberger, J. and Gustmans, E.N and Fung, E.T and Mackdah, J. and Omard, G.R. and Klies, E.A. and Spigel, O.R. and Hartman, A.R. and Annosotis, A and 112 Solders, M.	Genome-wide cell-free DNA (cfDNA methylation signatures and effect o tissue of origin (TOO) performance	n Clinical	37		2019	https://issocoubs.or gitto/stay10.1200/ ICO.2019.37.15_a meeting uppl 3049 abstract	811 cancer cell methylomes representing 21 tumor types	g Tissue-of-origin prediction	accuracy (training / test set split)	training + test set	a machine learning framework. Improvement was observed across all cancer types and was consistent in early-stage cancer (stage I-III). Respective performances in breast cancer (n = 23) were 87% vs 98%; in lung cancer (n = 22) were 85% vs 88%; in hepatobiliary (n = 10) were 70% vs 90%; and in pancreatic cancer (n = 17) were 94% vs 100%."	a machine learning framework. Improvement was observed across all cancer types and was consistent in early-stage cancer (stage 1-III). Respective performances in breast cancer (n = 23) were 87% vs 96%; in lung cancer (n = 32) were 85% vs 85%; in hepstabilizar (n = 10) were 70% vs 90%; and in pancreatic cancer (n = 17) were 94% vs 100%."
113 Use, WT and Wang, Y and Zhang, J and Ye, F and Hearing, XH and Us, B and Me, Q Y	A novel strategy of integrated microarray analysis identifies CENP CORL and COCO as a cluster of diagnostic biomarkers in lung adenocarcinoma		1 25	43-53	2018 China	https://doi.org/10 10/16/c.comtot.2018 J023643 article	5 different microarray datasets that included 330 samples	Case-control study	accuracy (LOCCV, external feet set)	cross-validation + external cohort validation	is significance (GMGS) and support vestor markine (GMA) analysis programsive to the closed problem of the control of the cont	58 from anxious depression. In this work, our goals are three-fold. First, we test the hypothesis that more clinically homogeneous groups of MDD patients are easier to predict from healthy controls than the entire MDD group using blood metabolomics
Dis, Y and Yish, L and Yang, T and Drinkerburg, W and Peeters, P and Steckler, T and Narayan, V A: 114 Wittenberg, G and Ye, J	Metabolomic biosignature and differentiates melancholic depressi patients from healthy controls	ve BMC Genomics	17	669-669	2016 Belgium	http://dx.doi.org/10 .1186912894-016- 2953-2 article	"the data set consists of 97 healthy control and 90 MDD subjects"	Case-control study	accuracy, sensitivity, specificity (10-fold CV)	cross-validation	data. Second, we develop a novel method for building maximally predictive and robust machine-learning classifiers that retain information on the correlation structure of the metabolomics data to ease biological interpretation. Third, we use	data. Second, we develop a novel method for building maximally predictive and robust machine-larming classifiers that retain information on the correlation structure of the metabolomics data to ease biological interpretation. Third, we use this framework to describe the metabolomics biosignature of melancholic depression." "I this study-visited different eane expension data sets containing 202 cancer, 115
Long, N P and Jung, K H and Yoon, S J and Anh, N H and Night, T O and Kang, Y P and You, H H and M 115 I E and Hong, SS and Boon, S W	Systematic assessment of cervical cancer initiation and progression uncovers genetic panels for deep learning-based early diagnosis and in, proposes novel diagnostic and progressits blomarkers	Oncotarget	8 6S	109436- 109456	2017 Vietnam	http://de.dei.org/10 118/5/Concollerged 228629 article	202 cancer, 115 cervical intraspithelial receptains (CIN), and 105 normal samples	: Case-control study	accuracy, sensitivity, specificity (30-feld CV, external text set)	cross-validation + external cohort validation	integrative systems biology assument in a multi-stage carefuloperois source. Now a training bearding application could have extended beard on the general panel of a separation of the second of the general panel of a separation. (1 — the SES gene deep learning model for the differentiation of constrained model and a materially inflated and constrained of \$47.056 (E)	Integrative systems biology assessment in a multi-stage carcinogenesis manner. Despite of intermity based disposition models were statistished used on the generic purpose of interface, generic of corecal carcinogenesis as well as on the substant undust seatched interface. See a seed of the seed of the terminal production of the seed of seed of the seed of the
116 Long, N P and Nght, T D and Kang, Y P and Anh, N H and Kim, H M and Park, S K and Koon, S W	Toward a standardized strategy of clinical metabolomics for the advancement of precision medicine	Metabolites :	10 2		2020 USA	http://dx.doi.org/10 3300mstabo1002 0051 article	review (not applicable) "The data set GSE44861 comprised 56 CRC tissues and 55 adjacent noncancerous tissues from the United	raview			we will elucidate the potential involvement of machine learning and demonstrate that the need for automated data mining algorithms to improve the quality of future research is undensible. Consequently, we propose a comprehensive metabolomics framework, along with an appropriate checklist refined from current guidelines and	we will elucidate the potential involvement of machine learning and demonstrate that the need for automated data minior algorithms to improve the quality of future
137 Long, N P and Park, Sand Anh, N H and Night, T D and Yoon, S J and Park, J H and Lim, J and Reson, S.	High-Throughput Omics and Statist. Learning Integration for the Discov and Validation of Novel Dispositi W Signatures in Colonectal Cancer	ical sry Int J Mol Sci :	20 2		2019 Vietnam	http://doi.org/10 3395/jms2000059 0 article	States. The data set GSE41258 had 18 GCR and 44 algorized monamerors tissues from the United States between 1992 and 2005. He data are GSE83880 contained 101. CRC tissues and 35 non-neoplastic miscuoli stoses from all patients with stage (ii of CRC from Kones.")	Case-control study	AUC, contribity, specificity (5-times repeated 30-fold CV, test set)	cross-validation + test set	cutting-deep algorithms to introduce convolugations for accurate diagnosis of control states (CPC) and models showed unables of control states (CPC) and models showed unables of cherword mean accuracy 0.598 (standard deviation (D) < 0.001, Morrows produce) 959/EDA < 0.0031, and models and models of 0.001, Morrows (produce) 959/EDA < 0.0031, and models of 0.001, Morrows (produce) 959/EDA < 0.0031, and models of 0.001, Morrows (produce) 959/EDA < 0.0031, and models of 0.001, Morrows (produce) "In this study, see conducted as produced in the same content advancements of the commentation for account cancer. The solution of produce produced in the commentation of paccessing cancer. The solution of produce produced produced in the solution of paccessing cancer. The solution of produced pr	The cash promoted is not all papeads, contribing male distinction transportations and continuing edge paperforms introduce counting regions and continuing edge paperforms introduce counting edge paperforms on the section of colorectal cases (CEC), all models in decederate statistication performances in which it is proposed to be the best of instruction, register the French off, the following were observed many accuracy of 20th (cashed deviation (DS) of 2005), manual societies of continuing contributions of proposed to the policy societies of when inflations in distribution in the policy societies of when inflations in distribution in the concentrations in the legistry societies when indirects with trained some concentrations in the end plays societies when indirects when it was consistent and papeads consistent provides to sensitive recent advancements in the concentrations in each edge papead coloration in contribution of the concentrations of the participation of the contribution of the concentration of the contribution of the c
Long, N P and Yoon, S J and Aoh, N H and Nghi, T D and Lim, D K and Hong, Y J and Hong, S S and 118 Kwon, S W	A systematic review on metabolom based diagnostic biomarker discove and validation in pancreatic cancer	ics- iry Metabolomi cs	14 8	109-109	2018 South Kore	http://dx.doi.org/10 _1007is11306-018- a 1404-2 article	review (not applicable)	review	AUC, sensitivity, specificity (25 discovery studies + different validation strategies across 9 validation studies)		observed in nine studies. The diagnostic area under the curve ranged from 0.68 to 1.00 (sensitivity: 0.43–1.00, specificity: 0.73–1.00). The effects of patients' bio-	observed in nine studies. The diagnostic area under the curve ranged from 0.68 to 1.00 (sensitivity 0.49-1.00, specificity 0.79-1.00). The effects of patients' bio-parameters on metaboleme alterations in a context-dependent manner have not been thoroughly elucidated.*

119 Lopes, C and Tucker, S and Salameh, T and Tucker, C	An unsupervised machine learning method for discovering patient clusters based on genetic signatures	J Biomed Inform	85		30-39	2018 USA	http://dx.doi.org/ _10168_bi-2018.0 004	g 7. article	191 multiple sclerosis patient	Cases only (sub-group stratification)	Rand Index on benchmark clusters (10 fold CV)	cross-validation
Lenuth, Lawf Schiffmunn, S and Schmitz, II and Bruikhonz, R and Lench, F and Fermins, N and 120 Wicker, S and Togodor, Lawf Gerislinger, G and Ultsch, A	Machine-learning based lipid mediator serum concentration patterns allow identification of multiple sciencis patients with high accuracy	Sci Rep	8	1	14884- 14884	2018 Germany	http://dx.doi.org/ 1008/s41998-01 39077-8	g & article	MS patients (n = 102) and healthy subjects (n = 301)	Case-control study	accuracy, sensitivity, specificity (10-fold nested CV)	cross-validation
Is, TP and Nos, XT and Chen, CH and Chang, M C and Lin, HP and Hu, YH and Chiang, Y C and 221 Cheng, W F and Chen, CA	Developing a Prognostic Gene Panel of Epithelial Overlan Cancer Patients by a Machine Learning Model	Cancers	11	2	13-13	2019 Switzerlan	http://likx.dol.org/ 3300/cancars.11 d 20229	o D article	1 different dataset with > 50 samples per group	Cases only (prognosis study)	accuracy, log-rank test p-value (LOCCV)	cross-validation
122 Ms, B and Gong, Y and Meng, F and Yan, G and Song, F 123 Ms, S and Song, J and Magis, A T and Wang, Y and German, D and Price, N D	Identification of a sindeen-gene prognostic bionarize for hug against a machine larning method. Measuring the effect of inter-study workfolling on estimating prediction error	Journal of Cancer	11		1288-1298 e110840- e110840	2020 China United States of 2014 America	http://dx.doi.org/ 7150/ca.34565 http://dx.doi.org/ 13716ournal.po 0110840	o article	TCGA cohort (in = 338) and it (in = 168) "4,470 microurry uniples of 6 long phonotypes from 26 independent sequences of 20 feets, and 20 feets and any other sequences of 20 feets are also sequences of 20 feets and 20 feets are also sequenced as studies".	Cases only (prognosis study) Case-control study	hazard ratio, p-value and C-index (training / feet set sqt)) "Fere we quantify the impact of these combined "study-wfeets" on a disease signature, spedictive performance by comparing two types of validation emblack melaliny-antifered trains validation (Total), which caretast combined study and study of the study of the study and study of the study and st	training + test set e n rcross-validation + external validation
124 Multidaz, N.S and Esimentia, O	Bioinformatics Approaches to Predict Drug Responses from Genomic Sequencing	Methods	1711		277-296	2018 USA	http://dx.dei.org/ 1007/978-1-493 7493-1 34	0 h article	review (not applicable)	review		
125 Mamodhina, P and Vielra, A and Putin, E and Zhavoronkov, A	Applications of Deep Learning in Biomedicine	Mol Pharm	13	5	1445-1454	2016 China	http://dx.doi.org/ .1021lans.molph manual Su0096	o E article	review (not applicable)	review		
Marve, M E and Copper, D and Tones, D T W and Hovestadt, V and von Deimfing, A and Pficter, S M 136 and Benner, A and Zuchnick, M and SII, M	Machine learning workflows to estimate class probabilities for estimate class probabilities for machine control of the control	Nature Protocols	15	2	479-512	2020 Germany	http://dx.doi.org/ 1.008n.41508-0 0251-6	∆ de article	brain tumor 4500, DMA methylation colored of 2,801 tumples with 51 color of 5,000 tumples with 52 with 3.50 tumples per group)	Cases only (sub-group stratification)	accuracy (b × 5-fold neeted cross-validation scheme)	cross-validation
327 Murtinuz, 8 i and Stabenfelds, S.E.	Current trends in biomarker discover and analysis tools for traumatic brain injury		13	1		2019 USA	http://dx.doi.org/ 1189813036-0: 0145-8	o article	review (not applicable) two colorist of a total of 319 subjects: "The first colorit, coming from the main recruitment contest (Culviversity of Salaren), busic composed by 69 hostilty subjects (CREI) and 48 MAPI Datasters	roview		
Measons, M and Trois, J and Agliss, A and Calvanese, G and Tone, P and Caruco, R and Colucci, A 228 and Dallis, M and Federico, A and Balsano, C and Persico, M	Accuracy of metabolomics profiles to non-invasively diagnose NAFLD stage and evolution by mean of machine- learning automated algorithms	s Digestive and Liver Disease	52		e9-e9	2020	http://dx.doi.org/ 1016/j.dd.2019 2.023	meeting abstract	(78 NAFL, 23 NASH, 43 NASH cirrhosis) and the second, coming from the other centers, wis composed by 106 subjects(40 CTR), 34 NAFL, 10 NASH, 18 NASH cirrhosis)*	Case-control study	accuracy (training / text set split)	training + test set

analysis is profromed to identify potential editorologic between the clusters in the context of travers biological branching between the clusters in the proposed method provides the greatest performance out of the methods stead: "The date on increasing performance is performed as produced provides and provides the greatest performance out of the methods stead: "The date on increasing performance is performed as the provides of the greatest performance out of the methods stead in "The provides of the greatest performance out of the methods stead in the provides of the provides and provides a implamentation for relation flowers and configured and, subsepts, subsequent and the configured and analysis of the configured analysis

dependently of the critical particular data (supporting the control particular data of the control particular data (supporting the control particular data) (supporting the control particular data (supporting the control particular data) (supporting the control particular data (supporting the control particular data) (supporting the control particular data (supporting the control particular data) (supporting the control particular data (supporting the control particular data) (supporting t

signifure. These intensive lay genes could achieve a strong power for proposed prediction of LLDs plants in colonic 1918 at 25, p. 4 5.86 of C. Cleake of 2015 at 25 of C. Cleake of 20

129 Metool, K and De No, C and Rahman, R and Ghosh, S and Pul, R	Investigation of model stacking for drug sensitivity prediction	BMC Bioinformati cs	19		71-71	2018 USA	http://doi.org/10 /1.1866/12869-018- 2208-2	article	"we segregate 50 training samples into our vertical and horizontal groups, build involvable predictive model of few 150 or involvable predictive model of few 150 or set of 150 samples, and obtain the prediction MES or candidate models on a set of 50 testing samples. We then add 2 valving samples and resettmate the MSE. We repeat this process used the training set has a solid of 100 samples. The set of 150 samples is testing, and has a loted at 150 samples times with randomly selected training, testing, and has already iteration."	Cases only (drug sensitivity prediction)	normalized AUC (training, testind and validation set)	training + test set
	Performance of a clinical/proteomic panel to predict obstructive peripheral artery disease in patients with and without diabetes mellitus	European Heart Journal	39		117-117	Netherland 2018 s	http://dx.doi.org/10 1093/suphatrijish y584.P732	meeting abstract	"354 patients undergoing peripheral and/or coronary angiography, performance of this diagnostic panel was assessed in patients with (N-94) and without DM (N-260) using Monte Carlo cross validation"	Case-control study	sensitivity, specificity, PPV, NPV (Monte Carlo cross-validation)	cross-validation
McDermott, J E and Wang, J and Mitchell, H and Webb-Robertson, B J and Hafen, R and Ramey, J and	Challenges in biomarker discovery: Combining expert insights with statistical analysis of complex omics data	Expert Opinion on Medical Diagnostics American Journal of Respiratory	7	1	37-51	2013 USA	http://dx.doi.org/10 .1517/17530050.2 012.718329	article	review (not applicable)	review		
132 McGeachie, M and Kelly, R S and Litonjua, A A and Weiss, S T and Lasky-Su, J A	Network of year-3 metabolites indicative of early-life asthma	and Critical Care Medicine	197			2018	http://dx.doi.org/10 3390/mu12051233	meeting abstract	cohort of 411 three-year olds at high-risk for asthm	Case-control study	AUC (five-fold cross-validation)	cross-validation
Melèle, P and Tiscotes, D and Barlow, C and Wee, J and Marinton's, G and Barler, M and Gouley, B and Best, L and Some, S and S	Plasma lipidomic analysis of stable and unttable coronary intery disease	Atheroscler osis Supplement s	11	2	24-24	2010 Australia	http://dx.doi.org/10 .1016/81567- 5688(10)/70103-3	meeting abstract	202 participants (control, n = 60, stable CAD, n = 61, untable CAD, n = 81). Whe used the data from the Geometric Park used the data from the Geometric Carel Care	Case-control study	AUC (multiple cross-validation berations)	cross-validation
Menden, M P and Iorio, F and Garnett, M and McDermott, U and Benes, C H and Ballester, P J and 134 Szez-Rodríguez, J	Machine Learning Prediction of Cancer Cell Sensitivity to Drugs Based on Genomic and Chemical Properties	Plos One	8	4	7-7	United 2013 Kingdom	http://dx.doi.org/10 .1371/journal.gone .0061318	article	more than 15 missing genomic features,	Cases only (drug sensitivity prediction)	R-squared (8-fold cross-validation, hold-out test set)	cross-validation + test set
135 Midorikawa, Y and Troji, 5 and Talayama, T and Aburstani, H	Genomic approach towards personalized anticancer drug therapy	Pharmacoge nomics	13	2	191-199	2012 Japan	http://dx.doi.org/10 2217/squ.11.157	article	review (not applicable)	review		
Mobadersamy, P and Yousell, S and Ampad, M and Gutman, D A and Bamholts-Sloan, I S and 136 Velazquez Vega, JE and Brat, D I and Cooper, L A D	Predicting cancer outcomes from histology and genomics using convolutional networks	Proc Natl Acad Sci U S A	115		E2970- e2979	2018 USA	http://dx.doi.org/10 1073/pnas.17171 39115	article		Cases only (prognosis study)	15 accuracy measurements, including HarrelTs C-Index for measuring concordance between predicted risks and actual survival (Monte Carlo cross-validation)	cross-validation
137 Modlin, Land Kidd, M and Drosdov, Land Bodel, Land Milicaevika, A and Matar, S	Automated finger prick blood genomic diagnosis of neuroendocrine tumors	Neuroendoc rinology	108		132-132	2019	http://dx.doi.org/10 _1159000499965	meeting abstract	whole blood samples from >6,000 NETs and controls	Case-control study	sensitivity, specificity (training/test set split)	training + test set
138 Mohammed, A and Biogert, G and Adamec, J and Helikar, T	Identification of potential tissue- specific caner blames/fern and development of caner brames are genomic classifiers. Plantin-specific caner genes contributes to recurrently perturbed partneys and establish therepositic values/artneys and establish therepositic values/artneys and establish and propositions.	Oncotarget Nat Commun	8	49	85692- 85715 3101-3101	2017 USA 2019 Italy	http://dx.doi.org/10 1883/Docordanat 21192 http://dx.doi.org/10 103884487-019 10886-3	article article	"A total of 2,175 tissue samples, both normal and cancerous, aree collected between the	,	accuracy, sensibility, specificity, precision, F1 score (10 fold cross-validation) log-rank test p-value (cross-validation)	cross-validation
Murugesan, K and Javle, M and Schrock, A B and Ngo, N and Frampton, G M and Alexander, B M and 140 Miller, V A and Beltail-Saab, T and Albacker, L A and Ross, J S and All, S M	cholangiocarcinomas (cHCC-CCA)	Annals of Oncology	30		v256-v257	2019	http://dx.doi.org/10 .1003/armonc/mdz 247.005	article	"1269 HCC, 3965 CCA and 44 cHCC-CCA" (> 50 samples for the main conditions)	Case-control study	error rate, AUC (Random Forest out-of-bag error)	outofbag
Nakamura, M and Bax, H J and Scotto, D and Souri, E A and Sollie, S and Harris, R J and Hammar, N	Immune mediator expression signatures are associated with improved outcome in ovarian carcinoma	Oncolmmun ology	8	6		2019 Sweden	http://dx.doi.org/10 .1080/21624029C2 019.1593811	article	1,656 ovarian carcinoma patient tumors	Cases only (survival prediction)	accuracy, recall, sensivity, Matthew's correlation coefficient, F1 score (5 times 10 fold cross-validation)	cross-validation

"A significant problem in precision medicine is the prediction of drug sentitivity for indebulac career call lines. [...] We update the prediction of drug sentitivity for indebulac career call lines. [...] We update the prediction of drug sentitivity for indebulac career call lines. [...] We update the prediction of drug sentitivity for indebulac career call lines. [...] We update the prediction of drug sentitivity for stabilities of the career call lines. [...] We update the prediction of drug sentitivity for stabilities are supported by the advanced career and inherent takes of career an

Predicting the regione of a specific cancer to a therapy is a major goal in mode oncopy that chould utilizately is not a personaled transment, we developed machine leaving mode to predict the regione of a specific cancer to a therapy is a major goal in mode oncopy that chould utilizately include to a personaled transment, we developed machine leaving mode to predict the response of a specific cancer to a therapy is a major goal in mode oncopy that chould utilizately include a personal transment, we developed machine leaving mode to predict the response of a specific cancer to a therapy is a major goal in mode oncopy that chould utilizately include a personal transment, we developed that the prediction of cancer of lines to day to the control of a personal transment of a control of prediction and cancer all the source of the control of a personal transment on the control of a personal transment on the control of a prediction and cancer are more pre

similarization) latern the visual partners and molecular biomarkers associated with visual partners and molecular biomarkers associated with visual places of the properties of the properties of the properties as increased approach for objective, accurate, and integrated prediction of printer advances.

Authorities approach for objective, accurate, and integrated prediction of printer advances.

Authorities are received for integrated prediction of printer advances.

Authorities are received for integrated prediction of printer advances.

Authorities are received for integrated prediction and biomarker discovery can have been accurate and performance of the printer advances.

Authorities are received for integrated production and biomarker discovery can have conceived printer and performance and performance.

Authorities are required for integrated production and biomarker discovery can have conceived printer and performance.

Authorities are required for integrated performance and performance.

Authorities are required for integrated performance.

**Authorities are required

the identification of potential respondent to 10 FUCKO therapy using income forests agreement.

We developed a compactification of protein that can be be strived to predict and the strip of the product of the strip of the str

342 Nakuriyahul, S	A hybrid gave selection algorithm based on retreaction information for microarray based cancer classification. PLoS On extraourray based cancer classification. PLoS On		e0212333- e0212333	2019 Thailand	http://doi.org/10 1371/formal emp 20212003 article	ten microarray data sets with > 50 samples per group for multiple datasets	Case-control study	accuracy, precision, recall, F-score (nested cross-validation)	cross-validation	This address gave electrics and multi-los learning enthods for concer classification with a control properties of the date in high demonstrating of multi-register, transforcing gave users on a decision against many dates, transforcing gave users on a decision against many and a dec
143 Naorem, LO and Muthalyan, M and Venkatesan, A	Integrated network analysis and Journal of machine learning approach for the identification of key genes of triplenegative breast cancer y		6154-6167	2019 India	http://dx.doi.org/10 1602/sch.27903 article https://da.doi.org/10 org.org/slood/articl	Six microarray data sets consisting of 46 non-TNBC and 405 TNBC samples	3 Case-control study	AUC (training / test set split)	training + test set	Expression formbus, II A naive Sayer based classifier built using the expression profiles of 15 states to jack good junctary sharehold youthly This Common-THEC amplies in the validation test data set with a receiver operating curve of 0.39 to 0.98." We built a personalized prediction model based on clinical and genomic data that
Naths, A and formolol), E S and Barnord, J and Al-loss, K and Padron, E and Madenat, Y F and Kumnanovic, T and Albahadra, N and Steenman, D P and Debrn, A E and Robbs, G J and Garcia- 144 Mannero, G and Liux, A F and Macigneski, I P and Selence, M A	A Personalized Prediction Model to Risk Stratelly Patients with Myelodysplastic Syndromes (MDS) Blood	130		2017	af 15/20undermard 4/2017/90/200156/ A-Paramakad Pradiction-Model to-Reak-Straffy article	"Of 2302 pts, 1471 were included in the training cohort and 831 in the validation cohort"	Cases only (survival prediction)	C-index (training and validation cohort)	external cohort validation	outperformed IPSs and IPSs. in predicting CS and AMX transformation. The new model gives usuaria probabilities a different term grows that are usuaria prosess that are usuaria and ease when charged probabilities at different term and probabilities at different term deposits and are always for a given ease when chigaretics and again was addord. "White treatment with the property-large agents (IMAA) association (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and protong survival as MoSt pasterns (LPA) and destabline (IMC) improves composits and most pasterns (LPA) and destabline (LPA) and destabline (LPA) and destabline destabline and destabline destabline (LPA) and destabline and destabline and destabline and destabline and destabline destabline (LPA) and destabline and dest
Naths, A and Selvers, M A and Sejor, R and Konrolli, R S and Barrourd, J and Al-liss, K and Proproduction, B P and Resto, C M and Secretors, D P and Distrem, A E and Rodox, G J and Garcia- 155 Manners, G and Chart, B L and Missiprovist, I P	Genomic Biomarkers to Predict Response to Hypomethylating Agents in Patients with Myelodyplastic Syndromes (MIDS) Blood	130		2017	https://ushpublicati pris.org/blood/urted er/19/79/pplement \$200.1157/79/8684 Gancanic Biomarkan.to. Practici-Responses to	433 pts with MDS (per 2008 WHO criteria) who received HMAs (230 at our institution (training cohort), and 203 at multiple other academic institutions (validation cohort))	Cases only (treatmen response prediction)	: accuracy, sensitivity, specificity (training and validation cohort)	external cohort validation	costs. We developed an unbiased framework to study the association of several matteriors in predicting response to 1MMs, and possible to 1MMrs of head of the second of th
Nowak, C. and Carlson, R. and Grupen, C. Jand Hyström, F. Hand Alam, M. and Fieldreich, T. R. and Sandtömin, Jan Lerrom Rollg, J. Jan de Lapent, J. and Heelberg, P. O. and Cordeiro, A. C. and Lind, L. and J. Miller, E. and Fall, T. and Armito, J. J. and Armito, J. and J. and J. Armito, J. and J.	Multiplex proteomics for prediction of I major cardiovascular events in type 2 Diabetol diabetes a	0gi 61	\$65-\$65	2018 Brazil	http://dx.doi.org/10 10079:00125-018_meeting 4693-0_abstract	1,211 adults with type 2 diabetes	Case-control study	accuracy with 95%-confidence intervals (training / forst set split)	training + test set	MACE in year 2 disabless. [] Addition of the Sign orders assay to the established risk and cell improved discrimation in the superari solution complet from SEG (1955). (6.23-Ne Sign (19-54, 19-54, 19-55, 19).** (6.23-Ne Sign (19-54, 19-
O'Railly, Paul Onitay, C and Gernon, G and O'Connell, E and Seogles, C and Boyro, S and Serrano, 147 and Sargeall, E	Co-acting game networks prodect L TRAIL responsiveness of tumour cells BMC with high accuracy Genomic	:s 15	1144-1144	2014 Ireland	http://dec.org/10 .1186/167-27664 .156-164 article	Gene expression microarray data for 10 tumor cell lines with known senditivity it have been exactivity to the death ligated opolinie tumor necroid-factor-related apoptosis-inducing ligand (TRAIL)	case-control study	AUC, sensitivity, specificity (training and validation cohort)	external cohort validation	cisclating the area under the receiver operator curve using an independent distanct, table to that the page unal indefendent out of page unal indefendent of page unal indefendent out of page unal indefendent of page unal indefendent out of page unal indefendent out of page unal indefendent of page unal indefendent out of page unal i
148 Oh, J H and Lidan, Y and Gurmani, P and Rosenblet, K P and Gao, J	Protate cancer biomarker discovery Methodologies using high performance mass spectral Programs serum profiling Biomed		33-41	2009 USA	http://kts.doi.org/10 .1016/j.creptb.2000 .04.000 article	Serum samples from 179 prostate cance patients and 74 benign patients	r Case-control study	accuracy, sanothely, specificity, NPY, PPV (DB times 10-feld Cr)	cross-validation	satinguish PAX from beings specimens. [] From the new maker selection algorithm, a paint of glass skillender as except of DAX, is anotherly of ESX, is a specime specimens of DAX, is another process of DAX, is anotherly of ESX, is a specime specimen to be value (PPV) of ESY, and is negative specimens of ESX, is another process of ESX, is another process of ESX, is another of ESX, is also as the process of ESX, is another of ESX, is also as the ESX of ESX of ESX, is also as the ESX of ESX of ESX, is also as the ESX of E
149 Okier, S and Pabhásia, T and Attokallis, T	Genetic variants and their interactions In disease risk prediction-Machine Biolizas Livering and interact pumperties filling	6 1		2013 Finland	http://doi.org/10 1186/1786-0081-6- 2 article	review (not applicable) "The first cohort (EGADODODIOLOL43), breafther tudy cohort) contains fill-Mode data and from CLL purified coils of 316 individuals along with critical data. The cohort was composed of 169 CLL, 27 cohort was composed of 169 CLL, 27 cohort was composed of 169 CLL, 28 cohort for was considered to the cohort for was Mode and five word 1910-point (response SSL) camples. There were 218 CHV modated cases and 64 IGMV immuttace case in 191 mules and 77 females. By tagging at diagnosis, there were 2 Mode.	noview			partons is another energing free for by which to explore how such extended level green, we discretel the siles concepts and egister has been been been been been been been bee
Organiza, A M and Rodrigueto, B A and Verson, N A and López Á, B and Arise, J A D and Versila, N D and 150 Phine, M 1 G and Broston, M M P and López, J L B	on new transcriptional patterns Oncolog			2019 Spain	has tills dutument to 1200 dates 2010.0 0072 article	case, 153 Binet Stage A case, 4 Binet Gag case. The accord color (Edu Case) case. Self color (Edu Case) Self	Case-control study	accuracy, precision, recall, ROC (training and validation cohort)	external cohort validation	"Choicis (prophosyn fri subsensis (CLI)) is the most frequently prophoprofiler strive syndromes in western countries. CLI southtimes (CLI) is the most frequently problems, and treatments is somethy researched from paperties with single or represent of disease properties of disease properties. The contribution of the most frequently problems, and treatment is a final frequently included, and treatment is a final frequently included, and treatment is a final frequently included and treatment is a final frequently included and treatment is a final frequently included and treatment in a final frequently included and treatment is a final frequently included and treatment in a final frequently included and treatment is a final frequently included and treatment included and tre
IST. Orlol, J.D. and Vallejo, E.E. and Estrada, K. and Pena, J.G.T. and Alzhaimers Dix Neutoimaging, Initia	Benchmarking machine learning Broc models for late-onset altheimer's Bioinforr disease prediction from genomic data cs	nati 20 1	17-17	2019 Mexico	http://doi.org/10 1.1889.12898.019 3158-x article	more than 50 samples per group for both discovery and validation cohort	Case-control study	balanced error, accuracy, sensitivity, specificity, AUC (cross-validation, training +validation cohort)	cross-validation + external cohor validation	by the ARInhem's Disease Neurolinging (institute (pCRs)) object. Our experimental by the ARInhem's Disease Neurolinging (institute (pCRs)) object. Our experimental is result demonstrated the discribitation performance of the best models tested violated -72% of area under the ROC curve.*



demonstrally robust, opperforming traditional approaches. However, considerations for application in relicant installation in inclinal missful profiler genian to be evaluated. Considerations from a place and the profiler genian to be evaluated. Considerations for application in relicant installation of profiler genian to be evaluated. Consideration for the profiler genian to be evaluated and the profileration of profileration and profileration of profileration for [FIC97] in metallicity profileration for formation formation for formation for formation for formation for formation formation for formation for formation formation for formation formation for formation f nciderations for annication in clinical metabolic profiling remain to be evaluated

of frast ampic. A machine learning-based polaries identified in a strons colorified or a feature of colorities of 60 pieces interfered in a 12-regularities objects represented in 12-regularities objects represented in 12-regularities objects represented in 12-regularities objects represented in 12-regularities objects o

The raing worderland prevalence of metabolic syndrome (Merit), a cluster of cardiometabolic in defense of type 2 disabeter, relates targely to increasing solventy and extensive by a data to see implication (in the cuttor of predictive of type 2 disabeter, relates targely to increasing solventy and extensive by set of the service of the solventy was to inderify predictive between the disabeter of the solventy was to inderify predictive between the relation tower. The displaces of the solventy was to inderify predictive between the relation tower the displaces of the solventy and to solventy predictive between the relation tower that the relation tower that the relation tower that the relation tower to the relation to the relation tower to the relation to the relation tower to the re integration improved performance and robustness of the prediction (11% misclassification on training set, 8% on validation set)."

considerations for annitration in clinical metabolic profiling remain to be evaluated

inaccurate results. This bias is due to insperfect stratification of analysis in the traverage and the size and the dependency between these stratification errors, i.e. the traverage and the size and the dependency between these stratification errors, i.e. the traverage and the size and the dependency between these stratification errors, i.e. the traverage and the size and the dependency exhibits the size and the dependency exhibits the traverage and the size and

would be profited the commonly used behaviorally law one of our convalidation should not be used to estimate AUC for until diseases."

The same first ingrenization clean was a second convalidation should not be used to estimate AUC for until diseases."

The same first ingrenization clean was a second convalidation should not be used to estimate AUC for until diseases.

The same first ingrenization clean was a second convalidation should not be used to estimate AUC for until diseases.

The same first ingrenization clean was a second convalidation should not be used to estimate AUC for until diseases.

The same first ingrenization clean was a second convalidation should not be used to estimate AUC for until diseases.

The same first ingrenization clean was a second convalidation should not be used to estimate AUC for until diseases.

The same first ingrenization clean was a second convalidation should not be used to estimate AUC for until disease.

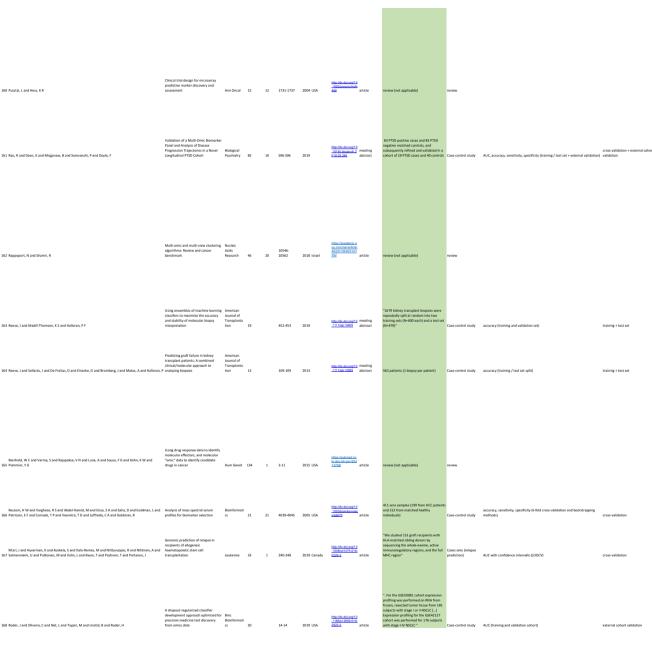
The same first ingrenization clean was a second convalidation should not be used to estimate AUC for until disease.

The same first ingrenization clean was a second convalidation should not be used to estimate AUC for until disease.

The same first ingrenization clean was a second convalidation should not be used to estimate AUC for until disease.

The same first ingrenization of the same

of perspirate blood monomutairs can't and perspirate blood endoughes from healthy controls was also unlifered bloomedure selection. Executing Collection and controls was also unlifered blooded measurement of controls was also unlifered blooded measurement of the controls was also unlifered blooded measurement of controls was also unlifered blooded measurement of the control was also with the control was a



been identified and its predictive accuracy was estimated, the goal of an independent unitation study is red (offered the sensitivity, specificity) and the postion (PV) and negative predictive values (PVP) with greater perceition, and (0) to prove clinical utility of the test. Offerent offered designs may be needed for different clinical students, but there may not be a single best design for any particular indicat scenario. Several designs could just design may be needed (prigor 2, 6) in important question for elegant could predict any could be propredicted primary and (primary 2, 6) important question for the propredictive production of the product of chemotherapy for example)."

This removaries review methodical and distribution laws reviewed to distribution discipling following to discours and validates multiplings profession of responses to the brackers, "Signate associated with PTSD development might energy across multiple feet of the profession and discours and memory and an emprove disposition of the profession of the pro solders is used for external validation [...] We previously found that the multi-omic point results in a mall improvement in diagnostic performance in comparison to individual ringle-omic panels in the initial training and validation orbors (AUC-0.03). The Aucroups (38) is sensitively, 73% sections, 93% sensitively, 74% sections, 93% external validation in the longitudinal cohort suggests that single-omic metabolic panels constituting the multi-omic panel are ingeligitately associated with FTSD status.

and Table." The current convergence of molecular and pharmacological data provides unprecedented apportunities to gain insights into the relationships between the two between the convergence of molecular and pharmacological data provides unprecedented apportunities to gain insights in insight in instance of the relationships between the two between the two precedented apportunities to gain insights convergence of molecular and pharmacological data provides to gain insight insight in instance of the provident and pharmacological data provides to gain insight in instance of the providented apportunities to gain insight in instance of the providented apportunities to gain insight in instance of the providented apportunities of the providented apportunities of the providented apportunities of the providented apportunities. If or careor progression to recognize the paroxyl of potentially influented event both for careor progression to recognize the paroxyl of potentially influented event both for careor progression to recognize the paroxyl of potentially influented event both for careor progression and the register comparative genome. In provident part of the provident progression and provident devent providents and providents and providents are provident to recognize the provident progression and provident developed application of the provident progression and provident developed application of the provident provident providents are provident providents. If the provident provident provident providents are provident providents and provident providents are provident providents. The provident provident providents are provident providents are provident providents. The provident provident providents are provident providents are provident providents. The provident providents are provident providents are provident providents. The provident providents are provident providents are provident providents. The provident providents are provident providents are provident providents. The provident providents are provident providen

underlying multi-omics analysis in general and multi-omics clustering in particular [...] We detected large differences between the p-values derived from the $\chi 2$ [...] We detected large differences between the p-values derived from the y2 approximation control to the "values derived from the permitted to test in the statistical test was used. The differences were specially large due to the manil statistical test was used. The differences were specially large due to the manil semple size, minil tester large is such sometime. As high number of cleanty and due to a low number of events (high number) for the legarnt kest. These p-values are used by ingle and multi-clien methods to assess these preformance, and the legards p-value is often the number of events the value is preformance, and the legards in the value is preformance, and analyses that are belowed on the y2 detection the validity of analyses that are belowed on the y2 detection the validity of analyses that are belowed on the y2 detection the validity of analyses that are belowed on the y2 detection the validity of analyses that are belowed on the y2 detection the validity of analyses that are belowed on the y2 detection the validity of the properties of the y2 detection the validity of the y2 detection the validity of the y2 detection the y2 detection the validity of the y2 detection the validity of the y2 detection the y2 detec

When, we review algorithms for must conice clustering, and discust say juscin in spoking these algorithms. Our review covers method developed specifically for conic and as well as general mid-seem extended specifically for conic and as well as general mid-seem extended specifically for conic and as well as general mid-seem extended specifically for conic and as well as general mid-seem extended specifically for conic and as well as general mid-seem extended specifically for conic and as well as general mid-seem extended specifically and the seem of the specifical sp

"Mass spectrometric profiles of peptides and promise spectra, high dimensionality and unbiasted into the profiles of perticipation of the per

150 Rodelguas CHTS, M E and Prottillo, C and Rodelguas, M, and Durby, P and Mischak, H and Onto, A was a compared to the compa	27, 1 1 d	Sci Rep	8 1	15940- 15940	2018 Germany	into de Administra Competituidade ADMIAS article	Rapid progressors (n = 342) Non-rapid progressors (n = 1140)	Case-control study	AUC, sensitivity, specificity, PPV, RPV (LOCCV + external validation callord)	cross-validation « external cohorte validation	the early identification of those individual most likely to progress will allow deficient false in high its dearly CDD gartents. For CD0723 dissistfers will a passed of 273 unitinary peopless that enables early detection of CDD and prognosis of progression. When have generated united populary detection of CDD and prognosis of progression of progression of CDD73 subclassifies reported to CDD and passed to allow the early people of CDD73 subclassifies reported to CDD and passed to allow the early CDD73 and continued to CDD73	death. Therespectic approaches to lamb gragescalor are limited. Developing both for give a very identification of those indeviduals must have progress and little or progress. and little or progress and little or progress. The CRU273 castion is high risk early CDD patients. The CRU273 castion is a pained of 273 undersp projection that produces only detection of CDD and prognosis of progression. We have generated united capillarly electrophories muss spectrometry based projectionness. CRU273 castionalisms speech for CDD capits on allow the sub- stantification of patients as high risk of CDD pagescrains. [1] invitabilities that section of patients as high risk of CDD pagescrains. [2] invitabilities and considerate patients are simply that of CDD pagescrains. [3] invitabilities and considerate patients of the simply that of CDD pagescrains. [3] invitabilities are considerated pagescrains. The consideration of the co
Quantic Novel, I and Error, Sand Classes, I to de Balant, I, and at Englant Zollman, E and Gosffront, C. and Richard, P and Labolatin, G. and Gosffront, C. and Richard, C. and Gosffront, C. an	C. D. Machine Learning for Better G. Prognostic Stratification and Driver Gene Identification Using Somatic	Oncologist	23 12	1500-1510	2018 France	http://dx.doi.org/10 16549-merchilige L0017-6866 article	97 patients with anequisate objectived originated origi	Cases only (prognosis study)	s error rate, Cindex (cross-validation, external validation)	cross-validation + external cohort- validation	recurrent CVV weets, detected in augustud displaced region and part of production of production of production of the production of the selecting potential genes for for water part of the production of the contract production of the contr	**Tap They could intest analysistic gliomate have variable clinical behavior. We have carefully above that the common 1961-31 salled into its above that the representation of the time of the country of
171 Rychkov, D and Sirota, M and Lin, C	learning to identify novel biomarker for rheumatoid arthritis	s Rheumatolo 8V	70	2206-2206	2018	http://dx.doi.org/10 1002/art.40700 abstract	and 759 healthy controls" (sufficient samples for whole blood)	Case-control study	Cohen's kappa, sensitivity, specificity (5-fold CV)	cross-validation	on the whole blood data, resulting kappa of 0.57 with sensitivity 0.54 and specificity 0.96." The efforts to personalize treatment for patients with breast cancer have led to a focus on the decear characterization of senotypic and phenotypic heterogeneity	on the whole blood data, resulting kappa of 0.57 with sensitivity 0.54 and specificity 0.98." "The efforts to personalize treatment for patients with breast cancer have led to a focus on the deeper characterization of genotypic and phenotypic heterogeneity
177 Seini, G. and Mittal, K. and Rida, P. and Jansson, E. A. M. and Gragineni, K. and Aneja, R. 173 Schreibert, L. and Lustoni, M. and Schmidt, R. and Repolitier, D. and Fuellen, G.	Panoptic view of prognostic models for personalized breast cancer management Tissue-based Alzheimer gene expression markers-comparison of multiple machine learning approximated in the control personal markets are consistent of multiple machine learning approximated investigation of redundancy in manal biomarkets.	Cancers es BMC Bioinformati cs	11 9	266-266	2019 USA 2012 Germany	http://de.doi.org/10 3500/senters/10 97325 article bttp://de.doi.org/10 1150/1471-2405 33-240 article	review (not applicable) PLURI and AD dataset contain > 50 camples per group	review Case-control study	securacy (8-feld CV)	cross-validation	among brace ciscores, 1 This review summirizes the prognostic and predictive mining the provided by commercially variables green persons have last stars and other multivariates or directal -entiric based prognostic predictive models currently variables green persons to be a start of the multivariates or directal -entiric based prognostic predictive models currently varieties and the control of the prognostic pro	among broad cancers. [] This review unmarkers the prognostic and prediction implifys provided by promiserizing viability properties by presented the present of them multi-varies or direct - omes-based proprosticity-inecticities models currently under development, and prediction of the prognostic prognosti prognostic prognosti prognostic prognostic prognostic prognostic
174 Shafi, A and Ngoyen, T and Peysondipour, A and Ngoyen, H and Doghloi, S	A multi-cohort and multi-conics, met analysis flowwards to identify network-based game signatures	a- Frontiers in Genetics	10		United 2019 States	Ministration and 20 Additional 2010 COLDS article	623 simples 233 simples from GBM shareness GB from the battley (son-state) individuals	Cases only (prognosis study)	s Cox p-value (training + validation cohort)	external cohort validation	remains a raisof challenge, in an effort to identify the biomakers that describe the value of the properties of the prop	molecular data of multiple types (e.g., personics, praconosis, proteonics, pro
175 Shain, K.L. and Acharya, C.R. and Smeltaer, S. and Lyenfy, H.K. and Acharya, K.S.	Non-invasive diagnosis of endometricsis: using machine learning instead of the operating room	Fertility and Sterility	112 3	e80-e80	2019	http://dx.doi.org/10 1016/f-fethringert 2 019.07.331 article	"We trained Random Forest classifiers on ten gene-expression based modules, derived from spectral decomposition of the discovery dataset (n 1/4 148) to predict the presence of endometriosis*	Case-control study	AUC, accuracy, NPV, PPV (10-fold CV+ external test set)	cross-validation + external cohort validation	and other pathology, with an accuracy of SMs (prea under ROC 140 ASI; p-value). 6.144-6.5), with an appetive predict berealuse of SMs and optive predictive such services where of SMs and of SMs and of SMs and of SMs and SMs. SMs. We reduced model over-fitting by performing 10-fold cross-validation of our discovery data." 1. **Those than two, blinks of women who, underso oursers for suspented marian.	identified a 288 gene predictor of indometriosis using Bandom Forests that was found to predict the presence of endometriosis, agraditises of the endometrial phase and other pathology, with an accuracy of 89% (pines under ROC L+4 0.84; p-value): 6.14-6.50; with a register predict in evalue of 48% (pines under ROC L+4 0.84; p-value): 6.14-6.50; with a register predict in evalue of 48% and a positive predictive value of 8.1%. We reduced model over-fitting by performing 10 fold cross-validation of our discovered state. ** **Where the tax must harder of various was a supplication of predictions or predictions of the contraction of course discovered that the contraction of the contract
Shork Land Date, YA and Stock Land Nat, I and Jan. Wand Mollow, A G and Average, Hand LaPol. 176: IP and Hoffman, M S and Sellers, T and Moley, T and Molesok, LY and Supplem, B	Measurement of phospholipids may la, improve diagnostic accuracy in ovarian cancer	PLoS One	7 10	e46846- e46846	2012 USA	http://dx.dei.eng/10 1371/format pone d0e6884 article	"a total of 1057 women with suspected ovarian cancer were errolled. [—] Only patients who underweit targety based on clinical supplien of ovarian cancer adaption with DCC, surgical starget with documented (including 23 in whom ECC was confirmed.]. A total of 211 cases and 212 benigns was included in the analysis."	Case-control study	error rate, sensibiley, specificity (5-fold CV)	cross-validation	necipation for on the securior. Our previous multi-suggest phospholipides as posterial believables of our size. In this study, we manuate the assum levels of multiple believables of the security of the description of the security of the security of the security of the security of the Health Security of the security of the security of the security of the security of the Health Security of the security of the security of the security of the security of the security of the security of the security of the security of the security of the security of the ALTS performance is listerately por -early stage cases and those of muchoosis cases with presporter and of probapholipies improving the stepper of early stage security of the security of the security of the security of the "In this stage, we applied units aperformance legislat diventage layer time of flight." "In this stage, we applied units aperformance legislat diventage layer time of flight was sectionary by one fine multiples grade for the security of 25 majors from the partition. A GS samples from health applieds. A MCNS AD description securities may section the property publishes allowed begaters were section may be profited. A named for section were section to profits. A named for section were section to profits and the section of the section of the profits and the section of the section of the profits and the section of the section of the profits and the section of the section of the section of the section of the section of the section of the section of the section of the section of the section of the sectio	all mosplanm of nor thever canner. Our previous results suggest phosphologists as potential of the biomoderist of outside course, in this study, we resulted the summer level of multiple of multiple and the summer level of multiple and the summer level of the summer
Shao, C H and Chen, C L and Lin, J Y and Chen, C J and Fu, S H and Chen, Y T and Chang, Y S and Yu, J 177 and Tsui, K H and Juo, C G and Wu, K P	Metabolite marker discovery for the S detection of bladder cancer by comparative metabolomics	Oncotarget	8 24	38802- 38810	2017 China	http://dx.doi.org/10 .18832/oncotarget. 16993 article	"metabolite profiles of 87 samples from bladder cancer patients and 65 samples from hernia patients"	Case-control study	AUC, accuracy, sensitivity, specificity (5-fold cross-validation + test set)	cross-validation + test set	on the metabolomic profiles and the six marker candidates. The decision tree	metabolomic profiles. [] A machine learning model, decision trees, was built based on the metabolomic profiles and the "or marker candistate." The decision for tree obtained an accuracy of 76.60%, a sensitivity of 71.88%, and a specificity of 86.67% from an independent test."

178 Sharma, A and Rain, B	C-MAGGISSA. Gane selection for cancer classification using multi- objective meta-leveristic and machine learning methods.	Biomed 178	3	219-235	2019 India	http://doi.org/10 10156.seep.2019 J8.5020 article	the proposed machine learning approaches tested on 7 microarray datasets, including a datasets with > 50 manufacts per group. The proposed machine learning	Case-control study	acturies/(DOCY+test-set)	cross-validation = test set	objective operation lyens optimizer (MOSTO) and say assuma algorithm (SSA). There is the finderplination profiles with three bits on election exaulty facts the challenge to materia connecipience and deverty, it also beam algorithm (SSA) materials deverting the contemporary of the contemporary of deverting the contemporary of	other hand, the calculation of MCGNPO requires law computs. Tools off-first hero is used for materializing the calculation of MCGNPO requires law computs. The calculation of MCGNPO requires law computer of the calculation of public plants of the calculation of public plants or application of public plants or pu
179 Shen, Land Tan, E C	logistic regression for cancer classification using microarray data	Comput Biol	2	166-175	2005 Singapore	http://dx.doi.org/10 1109hobb.2005.2 2 article	approachwas tested on 7 microarray datasets, including a datasets with > 50 samples per group	Case-control study	mean error + standard deviation (LODCV, test set)	cross-validation + test set		for comparison. It is shown that our methods have achieved at least equal or better results. They also have the advantage that the output probability can be explicitly given and the regression coefficients are easier to interpret."
German, S Lord Fagas, M and Huang, J and Lin, B and Diggars, J and Torn, E and Huagen, B and 180 Tutlin, B M and Germedy, G	Augmenting pre-operative risk of recurrence stratification in differentiated thyroid carcinoma usin machine learning and high dimensional transcriptional data fron thyroid FNA	g Journal of	: 15		2015	http://imedico.usco.pubs.org/cgi/conta- rutubos.act/33/15- uscypt/004/704/54-04- box/572-6955-4171- 0712- meeting 51mf-00505dd abstract	"81 samples preoperatively collected in previous study and post-curgically claganced as PTC. Each patient was categorised as either ATA fow risk or AT intermediates/high risk using established guidelines for recurrence risk stratification." (< 50 samples per group)	Cases only (risk of recurrence prediction	i) AUC (cross-validation)	cross-validation	prediction of risk for post-operative recurrence. If independently validated in a sufficiently large mounted of patients, cuts horizone of patient care? "The velody used to correct page (1/E) patient care? "The velody used to correct page (1/E) patient care? "The velody used to correct page (1/E) patient care in the velocity of the patient care in the velocity of the velocity	gene pairs. However, its general robustness does not extend to some difficult datasets, such as those involving cancer outcome prediction, which may be due to
181 Sh, P and Roy, S and Zhu, Q F and Roy, M A	Top scoring pairs for feature selection in mutchine learning and applications to cancer outcome prediction			15-15	2011 USA	http://dx.doi.org/10 1380/47/2020 12372 orticle	4 cancer proposits microarray datasets, including data with > 50 samples per group	Case-control study	error rate (LOCCV, best set)	cross-validation + text set	the relationship simple using softenin used by the districtive. We believe that the performance can be extracted by preparing the factors feature selection feature feat	the relatively implies origin planes used by the dissolfer. We believe that the performance can be endured by apposing its inflient behavior selection processors of the performance of
182 Sinners, 1A and Gsi, T	Omnibus risk assessment via accelerated failure time kernel machine modeling	Biometrics 69	4	861-873	United 2013 States	https://www.nchi.el m.mh.gov/penclust daw9MC5890038/ article	"training set of 454 lymph node negative breast cancer patients [] A total of 119 deaths or recurrences were observed"	Case-control study	C-statistic (training + validation data)	cross-validation + test set	interest, it may be unclear in advance which kernel to use for testing and estimation. We propose a robust Omnibus Test that combines information across kernels, and an approach for selecting the best kernel for estimation. The methods are illustrated	interest, it may be under in a demand which learned to use for testing and estimation. We propose a robust combined Set that conditions literature consists information across terms, and an approach for selecting the best learned for estimation. The methods are illustrated with an application in breast cancer."
Son, I and Angell T E and Balains, I and Barth, I Land Blooks, T and Dai, Q and Chousin, R As and Hamil, T Man de Hamil, I and Hamil, U and Exempt, G and Son, S and Stock, T and Uslook, V A and Parta, K H and Repolity, Claud Schaey, M M and Stank, M H and Traverid, S T and Wale, P S and 183 Whiteey, O and York, M and Ladenson, P W	d Clinical validation of the AFIRMA genomic sequencing parathyroid classifier	Thyroid 27		A50-A51	2017	http://dx.doi.org/10 106/inv.2017.29 meeting 046.inv.uchs abstract	476 FNAs-6 parathyroid and 470 thyroid FNAs	Case-control study	senctivity, specificity (training and validation cohort)	cross-validation + external cohon validation	sequencing and machine learning algorithms. [] The final classifier was billionly tested on an independent est set of 15 FAN LEI Bathesia (J. 7. Bethesid (J. 7. Bethesid (J. 7. Bethesid (J. 7. Bethesid (J. 7. Lid solffers had 100% sensitive) (4/4 parathyroid correctly called positive; C 39.3. 1004); and 100% specificity [191/191 thyroid correctly called negative; C 9.9. 1004); "For predicting relevant clinical outcomes, we propose a fittedib statistical machine learning approach that acknowledges and models the interaction between platform- souffice measurements through nordines revent machines and between the control of	**For predicting relevant clinical outcomes, we propose a flexible statistical machine learning approach that acknowledges and models the interaction between platform- seatific measurements fromth profitings where milicarines and province information
184 Shrustow, Sand Wong, W and Manyam, G and Ordona, C and Baladandayathapani, V Samusa, D and Clin, M and Publis, P and Welsteand, S and Black, A and Wongdo Height, A and Hy A and Bob, La Mark, M and A challenge, B and Terrorism, C Sand Sand, M and Scholare, P	Integrating multi-plotform generate data using hierarchical Reyetian relevance vector machines	Eurasip Journal on Bioinformati cs and Systems Biology 201 Alzheimer's	3 1		2013	bits life dat and 10 1180/1602 4153. 2013.2 article	"GBM data have multiple molecular emailmements on over SQS camples the include gene expression, copy number, methylation and microBNA expression."	: Case-control study	mean square prediction error "We nedomly gall the GBM survival data into a traveling data and a text data with 223 (90%) and 25 (10%) patheens, respectively")	training + test set	hap parameters with direct interpretations in terms of the effects of platforms and disc instructions, within a darce pathforms. In a grammater estimation appoints in the control pagas of the contro	disa interactions within and across platforms. The parameter estimation algorithm in our model uses a computationally efficient variational Bayes approach that scales well to large high-throughput disasses. [] We apply our methods of integrating gene/milks expression and microSMA profiles for predicting patient survival times to The Cancer Genome Latti (TCAQ) based globiostors multiforms (SMM) disasses in terms of prediction accuracy, we show that our non-linear and interaction-based integrative method perform better than increal interactions. And on-integrative method on-integrative.
As the first, hear Verille, I are level conditionable, it also facilities in a Linda (it is local productionable). See the Constitution of the Con	approach to diagnose Alzheimer-type dementia in blood: Results from the	and and Dementia: Translationa I Research and Clinical Intervention s 5		933-938	2019 Belgium	http://dx.doi.org/10 10166/ess.2016.1 1801 article	242 copitively normal (CN) people and 115 with AD-type demertia utilizing plasma metabolites	Case-control study	AUC (heated cross-sulfidation, external test set)	cross-validation + external cohorn validation	Need Cons. Violation (NCO, 1_1) On the text data, Disposition of the ACM of OSI. (ISBO 2-19), Killion Device of 18 (18 -6-19) of produced 08 (18 -6	Nested Core x Validation (NOV.) — I On the text deta, D, sendourid the AUC of 285 (IAC-049), Xilliano (NOV.) — I Some has deta, D, sendourid the AUC of 285 (IAC-049), Xilliano (NOV.) — I Some service of ampliado, Para and t tau (Ingelter with age and ingelter of the AUC of 285 (IAC-049), Xilliano (IAC-049
186 Statislacy, A and Affloric, C F and Transactions; I and Hardin, D and Lony, S	A comprehensive evaluation of multicategory classification methods for microarray gene expression cance diagnosis	r Bioinformati cs 21	. 5	631-643	2005 USA	http://dx.doi.org/10 1003bioinformatic wt86033 article	11 datasets spanning 74 diagnostic categories and 41 cancer types and 12 normal tissue types	Case-control study	accuracy, relative classifier information (Design i: rested stratified 30-fold CV coder loop, 3-fold CV inner loop, Sosign 1: rested LOCCV outer loop, 19-fold CV lever floop)	cross-validation	techniques can significantly improve the classification performance of both MC-SVMs and industrial glagatimes. Exemenble classifiers do not generally improve performance of the best non-ensemble models. These results guided the construction of a software system GEMS (Gene Expression Model Selector) that automates high-quality model construction and enforce sound optimization and	techniques can significantly improve the classification performance of both MC-DNM, and other ros-DNM long signifithms. Ensemble classifiers do not generally improve performance of the best non-ensemble models. These results guided the construction of authority seymen DNM (Silver personal hodder listered) that is authorities signify-quality model construction and enforces sound optimization and performance estimation procedures.

								"To develop multi-omic predictors of anticancer therapeutic response we curated data from the CCE, Gr9, and NCIGO databases. The resulting datasets consisted of the gene expression (Affymetrix U133A and Affymetrix			
187 Station, L.C. and Pearl, T. and Chen, Y.W. and Barnholtz-Sloan, J.S.	Computational identification of multi- omic correlates of anticancer therapeutic response	Brnc Genomics	15		8-8	2014 USA	https://bmcgaromi cs.biomedicerbal.c contratisties/10.118 8/1471-2184-15. 87-82 article	U133A plus 2.0), copy number variation (Affymetrix SNP6.0), and mutational status (targeted and whole exome sequencing) of 1299 distinct human	Cases only (drug response study)	precision + standard deviation of precision (10 feld CV, external validation data)	cross-valida
Sweatt, A J and Hodin, H K and Balasubramanian, V and Hsi, A and Blum, L K and Robinson, W H and Haddad, F and Hickey, P M and Condiffe, R and Lawrie, A and Nicolls, M R and Rabinovitch, M and 188 Khatri, P and Zamanian, R T	Phenotypes Using Machine Learning	Circulation Research	124	6	904-919	United 2019 Kingdom	http://dx.doi.org/10 _1161/cicresahe_1 18.313011 article	and 2014, we measured a circulating proteomic paniel of 48 cytokines, chemokines, and factors using multiplex immunoassay."	Case-control study	log-rank test p-value (discovery + validation cohort)	external co
Tabas, E and Longby, A F and Lin, H and Bolm, C R and Sades, S L and Leffouriz, J F and Rurz, D M 189 and Viscole, P and Venction, 1 M and Releas, T G and Pareson, J M and Riss, D M and Long, K T	Molecular characteristics and disease burden metrics determined by next- generation sequencing on circulating tumor DNA correlate with progressio free survival in previously untreated diffuse large 8-cell lymphoma		134			2019	http://dx.doi.org/10 11879boot_2019_ 1228533 abstract	"targeted NGS on plasma samples from 310 previously untreated DLBCL pts enrolled in the GOYA study"	Cases only (predicting progression-free surivial)	: Correlation with progression free survival (training/test set split)	training+t
190 Tan, A C and Gilbert, D	Ensemble machine learning on gene expression data for cancer classification	Appl Bioinformati cs	2	3	S75-83	2003 UK	http://cheseerx.ist. psu.edu/viewdoc/d ownload?doi+10.1. 1.2.9189&rep-rep 1&type=pdf article	Seven microarray datasets were used, including data with > 50 samples group	Case-control study	accuracy, sensitivity, specificity, PPV (10-fold CV)	cross-valida
:331 Tan, A C and Naman, D Q and $N_{\rm NL}$ L and Winslow, R L and German, D	Simple decision rules for classifying human cancers from gone expression profiles	Bioinformati	21	20	3896-3904	2005 USA	http://kk doi.org/10 .1003/bioinformatic arbs631 article	"19 publicly available microarray datasets, with sample sizes ranging from 33 to 327" (more than 50 samples per group for multiple datasets)	Case-control study	accuracy (LOCCV, test set)	cross-valid:
192 Tang, K Land Li, TH and Xiong, W W and Chen, K	Ovarian cancer classification based or dimensionality reduction for SELDI-TOF data	n BMC Bioinformati cs	11		109-109	2010 China	http://dx.doi.org/10 1188/1471-2105 11-100 article	"high-resolution SELDI-TOF ovarian data set for 95 control samples and 121 cancer samples"	Case-control study	accuracy, sensitively, specificity (cross-validation)	cross-valid:
193 Tao, M and Song, T and Du, W and Han, S and Zoo, C and Li, Y and Wang, Y and Yang, Z	Classifying Breast Cancer Subtypes Using Multiple Kernel Learning Based on Omics Data	Genes (Basel)	10	3		2019 China	http://dx.doi.org/10 310/fepares 10030 200 article	"Our dataset contained 606 distinct patient samples of breast cancer, which was divided into five subsyes: 277 furninal A, 0 timinal B, 70 Triple Negative Breast Cancer (TNBC), 11 HER2 (+), and 208 unclear"	Case-control study	accuracy, AUC (19 feld CV)	cross-valida
194 Telbani, A and Afonso, C and Marrier, 5 and Belafi, 5	Omics-Based Strategies in Precision Medicine: Toward a Paradigm Shift in Inhorn Errors of Metabolism Investigations	Int J Mol Sci	17	9		2016 France	http://dx.dei.org/10 _31000/j=ss1700155 S article	review (not applicable)	raview	This is notice, as greater tasks of the set multi-order data analysis transger in a reliant content. The shiftings of orders date shiftings transger in a reliant content. The shiftings of order date shiftings transger in a reliant set of the shiftings of the shiftings of the shiftings approached for indemnity or early order date of the shiftings of the shifting	ŧ.
195 Theofilatos, K and Korfuel, A and Mavroud; S and Comporthwaite, M C and Shpak, M	Discovery of stroke-related blood biomatures from gene augression network models.	BMC Med Genomics	12	1	118-118	2019 Greece	Manufacturan (1) 1388-1590-159 1588-1590-159 1588-1590-159	Slood samples from 82 stroke parlients and 6F controls	Case-control study	accuracy (5-fold cress validation)	cross-validi
156 Tell, T.S. and Donndelinger, F. and Wang, D.	Looking beyond the hyper Applied Al Applied Al and machine learning in translational medicine	Ebiomedicin e	47		607-615	2019 UK	http://doi.org/10 10164.deber 201 268.027 article	review (not applicable)	roview		
Tong, D Land Boccock, D and Coveney, C and Saff, J and Gomez, S G and Querot, S and Rees, R and 197 Bell, G R	A simpler method of preprocessing MALDI-TOF MS data for differential biomarker analysis: Stem cell and melanoma cancer studies	Clinical Proteomics	8	1		2011 UK	http://ike.dei.org/10 .1186/1550-0275-8- 24 article	Melanoma data set: 101 patients analyzed (yielding mass spectral data for 99 samples). Cord blood data set: 158 samples. 70 samples were categorized as containing a "High" number of team cells and the remaining 88 samples with a "Toos" number of stem cells	Case-control study	AUC, accuracy (Monte Carlo cross-validation + external validation set)	cross-valid validation

Three tags stool plantmengement, cludies have a reveal distinction or composed, in progress, and compared on the prediction of dispressions, and compared on the prediction of dispressions, and compared on the prediction of dispressions, and compared on the prediction of the prediction of dispressions, and the prediction of dispressions, and the prediction of the prediction of dispressions, and the prediction of the predict

species in sprigic of action is scheme; total and to general hypothesis about proble interactions among functionally relevant grows, leading to the identification of none informative bismatcher."

"There are general non both rich allowed the special properties of the special pr

Mactine is a Mactine in State of Chen, A and Fulgioti, V and Erizhnan, A 5,070 exorm	J C F learning in schizophrenia L a case-control study using	American ournal of Medical Senetics Fart B- Neuropsych atric Senetics 1	80 2	103-112	2019 Canada	http://dec.doi.org/10 1600/bisima in 1993 8 article	"This study applies ML to WES data from 2,545 individuals with SC2 and 2,545 unaffected individuals".	Case-control study	AUC, accuracy, senditivity, specificity, precision, recall, F1-measure (training test set)	training + test set	"Our hypoth (WES) data of study applier individuals, is supervised to Gradient Boo promising rether receiver algorithm we pathophysio "[To] evaluations are the second to
Troisi, J and Sarno, L and Martinelli, P and Di Carlo, C and Landolfi, A and Scala, G and Rinaldi, M and non-invasiv	iomics-based approach for two diagnosis of P email anomalies c	Metabolomi s :	13 11		2017	http://dx.doi.org/10 1007/s11006-017- 1274-2 article	"Metabolomic profiles have been obtained on serum of 328 mothers (220 controls and 108 cases)"	Case-control study	AUC, accuracy, senditively, specificity, PPV, NPV, F-measure, G-mean (Leave k out cross-validation, external text set)	validation	the second t and optimize randomly di- diagnostic p classified all also, the oth "Fetal malfo intrauterine to build a me
Troisi, J and Sanos, L and Richards, S M and Symes, S J and Adair, D C and Scale, G and Taylor, R S and anomalies: 200 McCowan, L M and Fasanos, A and Martinelli, P and Guide, M way	: The serum metabolomic [Birth Defects Besearch 1	10 9	757-757	2018	http://doi.org/10/meeting 11002/bbg-11955 abstract	Metabolomic profiles were obtained from serum of 654 mothers (320 controls, with a normal fetus and 334 cases with a malformed fetus)	Case-control study	accuracy, sendining, specificity (braining and validation set)	cross-validation + test set	with an inde Zealand SCO controls. Per maternal ser analysis usin 99.4±0.1% (r [1892/1894 a promising "Blood-base spectrum dir
Blood transcription of Stand Hess, I Land Quine, T P and Barve, R and Houng, H and Zhang James, Y and Chang, J spectrum discount of Stand Standow, S S and Sharp, F R and Herts Processe, I and Faranes, S V and Kong, S W and Glast, S 1 samples on	s with and without autism disorder: A combined-	km J Med Senet 8 Beuropsych atr Genet 1	74 3	181-201	2017 Norway	http://dx.del.org/10 1002/eers 8-1251 1 article	"Raw microarray data and clinical mata-data were obtained from seven studies, totaling data distributed and 447 comparation subjects"	Case-control study	AUC, sensitivity, specificity (needed 18 fold CV within 5-times beconstrapped (boots 7) samples + fact sard)	cross-validation + test set	differences in We sought to mega-analytex vivo blood comparison demonstrate perform with our results wand trophic circulating swhich they! "[] this pai 11-regularizz breast cancer in BLASSO, 6 performance breast cancer of 0.7 and 0.7 a
BLASSO I in Societies (Budsta, R and France, L and Vereda, F J and Clarce, M G and Serse, J M model? 202 Unds, D and Anagon, F and Budsta, R and France, L and Vereda, F J and Clarce, M G and serse, J M model?	ntegration of biological is into a regularized linear E	BMC Syst Biol :	12	94-94	2018 Spain	http://dx.doi.org/10 _1186/s12918-018- 0812-8 article	"Our of the 1212 samples, 1013 corresponds to controls for all patients; and 1910 cases for patients; who died from the disease;"	Case-control study	AUC (100 registrose of 10 fold nested CV, with 9 fold CV nested for hyper- parameter tuning)	cross-validation	of 0.7 and 0. approaches, obtained wiffound, BLAS (RI). The Ger given by LAS performed t showed a sig (IFNK) and o related with "The availab knowledge a exploited wiffound or result in more gene sets, as which emplic
	ased outcome prediction at cancer compendia F	LoS One	2 10	e1047- e1047	Netherland 2007 s	http://dx.dxi.org/10 1.137.focumi conte 2007/MZ article	"This compendium contains data from various cancer types and has a total of 1973 arrays; (meet han 50 samples per group for combined datasets)	Case-control study	AUC (stouble-loop cross-validation + external validation)	cross-validation + external cohort validation	modules de the validatic trend in con single breas better comp the module- biology. Fre cycle, E2F re "[] early-st challenging early-stage potentially r promising d was perforn
Wan, N and Weinberg, D and Liu, T Y and Nielhous, K and Arizo, E A and Dirilaba; D and Bannas, A and White, B and Ballyw, M and Berist, M and Ballyw, N and Berist, M and Ballyw, N and Berist, M and Ballyw, B and Berist, B and Ballyw, B and Berist, B and Ballyw, B and Berist, B and Ballyw, B and	learning enables detection of a colorectal cancer by whole- equencing of plasma cell- N E	BMC Cancer :	19 1	832-832	2019 USA	http://dx.doi.org/10 11869112885_019; 6003_8article	N = 545 colorectal cancer and 271 non- cancer controls	Case-control study	AUC, sensitivity, specificity (S-fold CV + confounder-based cross-variablescent)	cross-validation	extracted, as using IchorC and confoun and confoun a colorectal. (/II), we achi 85% (95% CI stage and in institution d accurate ass "Our study a new metabo [] We com
Wang, J and Yan, D and Zhao, A and Hou, X and Zheng, X and Chen, P and Bao, Y and Sa, W and Hu, C conteopora 205 and Zhang, Z L and Sa, W	of potential biomarkers for soise using LC-MS/MS C mic methods I	Osteoporos nt :	30 7	1491-1499	2019 China	http://dx.doi.org/10 _1007x00198-019- 04892-0 article	"Our study recruited 320 participants, including 138 males and 182 postmenopausal females"	Case-control study	AUC, accuracy, sendibility, specificity ("The data sets were randomly split into the recommended ratio of 70% for model training and the other 30% for validation")	training + test set	postmenopa random fore algorithm wi then metabs postmenopa classification the curve (A significantly' "The advanc has made it the tradition
Wang, J and Zuo, Y and Man, Y G and Avital, I and Stoladinovic, A and Liu, M and Yang, X and identification		ournal of Cancer	6 1	54-65	2015 USA	http://dx.doi.org/10 7150/jca.10831 article	review (not applicable)	review			generated by mechanisms massive data pathway and evaluates the discovery use "In this pape diagnostic by network by
Directing of Administration (Administration of Administration of A	diagnostic biomarkers of "s disease by integrating ession data in six brain F	rontiers in Senetics :	10		2019 China	http://dx.doi.org/10 33898/sene 2019. 00157 article	gene expression profiles of 161 samples in six brain regions The TU cohort consists of a total of 89 subjects (40 HCC cases and 49 patients with fiver cirribosis), and the GU cohort comprises of 15 subjects (57 HCC cases	Case-control study	AUC, sensitivity, specificity (LOOCV)	cross-validation	integrated g differential r these coexp discrimination then validate
Ghromis a	e Analysis of Proteomic. E	EEE J Blomed Health Inform	20 S	1225-1231	United 2016 States	http://dei.org/10 _1109/phi-2016_25 74201 article	and 59 patients with liver cirrhosis) (more than 50 samples per group for the	Case-control study	AUC, accuracy, sensitivity, specificity (10-fold cross-validation + test set)	cross-validation + test set	"In this stud metabolites to distinguis to select a p data previou cancer study to separate cancer cases

45 CDPF Court — 1.1 The exemelle model converty classified at I cases and converty classified at I case and converty classified at I cases and converty classified at I case and converty classified at I cases and converty classified at I case and converty classified at I cases and converty classified at I cases and converty classified at I cases and converty classi

And the function of the first contrary during contrary during the first of the contrary during contrary during

ment of high throughout omic technologies during the pact few years descrement of high throughput omit schrodinges during the past for wears delt possible to perform war, omplies assays in an with others time that officional approximation. The regist accumulation and wide availability of omit data of the performance of the very using hippatentialize accritionary (ICC) as an example."

In paper, we gropped as madeline desirable government of the paper, we propose a madeline desirable government of the paper, we propose a madeline desirable government of the paper, we propose a madeline desirable government of the paper was proposed as madeline desirable government of the paper was proposed as madeline desirable government of the paper was proposed as madeline desirable government of the paper was proposed as madeline desirable government of the paper was proposed as madeline from the paper was proposed as madeline from the paper was proposed as the best and paper was proposed as the paper was proposed as the paper was proposed as paper per, we propose a machine-learning-based method of identifying potential. "In this paper, we propose a machine-learning-based method of identifying potential

hypothesis is that machine learning (ML) analysis of whole ecome sequencing (fasts can be used to identify individuals at high risk for inclination (ICC). This (WTS) data can be used to identify individuals at high risk for inclination (ICC). This (WTS) data can be used to identify individuals at high risk for inclination (ICC). This (WTS) data can be used to identify individuals at high risk for inclination (ICC). This (WTS) data can be used to identify in the size of preference of the size of the s

malformations (PM) are structural or functional anomalies that occur during rate development. ___ | Wy performed a characterization of material areas included anomalies. ___ | Wy performed a characterization of material areas included anomalies. ___ | Wy performed a characterization of material areas included reported from the sub-propagate public from the \$550PC (Chort. ___ | The ensemble model correctly classified at case and other public public

as denote the houling of each craimate, showe define profession to my companion of the comp fation data compared to sens-based predictors. We also show that there is a the validation data compared to sens-based predictors. We also show that there is a the CML Muchine tearing models was trained using 16 dot case solidation with the control of the

Wang, N and Cao, Y and Song, W and No, K and Li, T and Wang, J and Xu, B and Si, H Y and No, C J 201 Li, A L	Serum poptide pattern that differentially diagnoses hepatitis illa and virus-related hepatic ellular carcinoma from liver cirrhosis	J Gastroenter ol Hepatol 29	7	1544-1550	2014 China	bhoillis dei ceartú 1111 ligh 12565 article	80 MCC and 67 LC patients "frozen fisses was collected only from 169 patients, of which only 166 contained more than 20% Lumor	Case-control study	AUC, accuracy, sensitivity, specificity (10-fold cross-validation + test set)	cross-validation + test set	(LC) caused by hepatitis is vivia (HRV) (Infection 1_1 LWth a highly systemized psychology and control and with a second to later doubterform formation and must have desirated used an object formation from each fight infection of fight mass spectrometric approach, we investigated seem psychology professer of the NCL and Control	carcinoma (HCL, It is not sufficiently sensitive to differentiate HCL and laver infraosis. (ICL) caused by lepidists lives (HBM) infraosis. In JMM in JMM) optimized poptide extraction and matrix-actised taser decorption/instanton time of flight/time of flight mass spectrometric. JMM in
210 Wang, 53 and U, M C	Impacts of Predictive Genomic Classifier Performance on Subpopulation-Specific Treatment Effects Assessment	Statistics in Biosciences 8	1	129-158	2016	http://dx.doi.org/10 1007/a12561-013- 00022x	cellularity, and gene expression profilin was completed in 33 patients using the Affymentra 138 microarray (8), Among them, 62 NSCLC patients received OBS alone and 71 NSCLC patients received OBS alone and 71 NSCLC patients received ACT." "We accessed the 500K Affymetrix chip genotype data from WTCCC on -1,500 samples from the 1958 British He 1958 British III.	Case-control study	accuracy, sensitivity, specificity, PPV, NPV, permutation p-values (cross-validation + external validation)	cross-validation + external coho validation	based on their expresentative bring parameter sets with varying degree of right in their chascies of the parameter ranging from high yeapons, ordered ranging parameter sets mildly galaxies. We attribute the statistics on the obtain of mining parameter sets the statistics of the statistics of the statistics of the obtain of mining parameter sets the parameter of the statistics of the statistics of the obtained of the statistics of the statis	based on three representative busing parameter sets with varying degree of rigor in their choices of the parameter ranging from highly reprove, moderately rejouse to middly reprove. We articulate the rationalse on the choices of tuning parameter sets, which was been a proposed to miscadisaction of generic bloomater calculation on their assessment of treatment effects in the posterie and register parient subappopations, and also corner parient. Subappopations, and also corner parient.
Wisi, 2 and Whong, E and Gs, If G, and Shang, H and Bradfold, Look Elm, C and Frachiston, E and In C and Gloscow, 1 E and Cheeses, 8 and Stanley, C and Monos, Q and Grant, 5 F and Polydromation 211 and Relevances, H	From disease association to risk ou, associament an optimize deer for s, C genome-wide association studies o type 1 diabetes	n n PLoS Genet S	10	e1000678- e1000678	United States of 2009 America		Cohort, -1,500 camples from the UK Blood Service Confer Group, as well \$\$ -2,000 camples each from the following becase collection: page 1 diabetes (T2D), the 2 diabetes (T2D), the 2 diabetes (T2D), the 2 diabetes (T2D), the content of the 10 conference (T2D), the 2 diabetes (T2D), the content of the 10 conference (T2D), the 2 diabetes (T2D), the 2 diabetes (T2D),	d Case-control study	AUC, accuracy, sensitively, specificity (5-fold cross-validation + test set)	cross-validation + test set	used only a limited number of confirmed susceptibility lock liver was proposed to expositionation describes ensuing approaches with large ensurable of markets may proposed with large ensurable of markets may proposed with a generable of markets may propose with a performance of disease and assessment, the applied all support Vector proposed part from the proposed part of the proposed and assessment of the proposed part of the	us of only a similed number of confirmed assesphibity size. Inverse we propose that up- support the properties of the confirmed assesphibity size in the set properties and up- support values of disease risk assessment. We applied a Support values or support the properties of the set o
White, 8 5 and Rhan, 5 A and Ammad-Ud-Do, M and Podar, 5 and Mason, M I and Tiggion, CE and Dokus, 8 and Mason, M I and Tiggion, CE and Dokus, 8 and Poolsa, K and Tywer, VM and 22x Admittability, 7 and Recommending, 6 and Guizeng, 1	nd Gene expression predicts ox vivo d sencitivity in acute myeloid leuker		13		2018	bibulishi dalami10 11881938- 7483 AU2016 meeting 3883 abstract	"We harmonized two large-scale AML evine studies screened for dury response and profiled transcriptionically—OSE (IDS) AML patient samples and ISO drugs) and FMM (48 AML samples and 480 drugs) and FMM (48 AML samples and 480 drugs)."	Case-control study	correlation, p-value (IS-fold CV)	cross-validation	data asts, we trained a Riige regression model on the ORHUI data size, used the mode to predict regioner in FRMM dataset, are clearly associated for Neurosco consistant on the predict regioner in FRMM dataset, are clearly associated for Neurosco consistant on the prediction of the Control o	dist active, we trained a Ridge regression model on the ORTG data set, used the model of the project regression in Ridmid acts and extraction for Paramot controlled in Paramot paramoters (PRIV - 20%, and paramoters) and
213 Wu, H and Col, L and Li, O and Wang, X and Zhao, S and Zho, F and Zhou, K	Metagenomics Biomarkers Selecte for Prediction of Three Different Diseases in Chinese Population	Biomed Res Int 2018	В	2936257- 2936257	2018 China	http://likt.doi.org/10 _1155/2018/20862 57 article	"microbiome of 806 Chinese individuals (383 controls, 170 with type 2 diabetes, 130 with rheumatoid arthritis, and 123 with liver cirrhosis)"	Case-control study	AUC, F1-score (S-fold CV)	cross-validation	from microbiome and corresponding phenotypes." "To harness the rish information in multi-omics data, we developed GDP (Group lass regularized Deep learning for cancer Prognosis), a computational tool for survival prediction using both clinical and multi-omics data. GDP integrated a deep learning	prediction using both clinical and multi-omics data. GDP integrated a deep learning
214 Xiq, G and Dong, C and Rong, Y and Zhong, I F and Li, M and Wang, K	Group lasso regularized deep learn for cancer prognosis from multi-on and clinical features		3		2019 USA	http://dx.doi.org/10 3300/genest10000 240 article	the used dataset cover more than 50 samples per group	Case-control study			framework and Cox proportional hazard model (CPH) together, and applied group lass regularization incorporate gene-level group prior honeselged in the model training process. We evaluated its performance in both simulated and real data from The Cancer Genome Attiss (TCGA) project.*	
215 Yang, J S and Zhu, Z X and He, S and Ji, Z and Inee	Minimal-redundancy-maximal- relevance flashers selection using different relevance measures for omics data classification	Computatio nal Intelligence in Bioinformati cs and Computatio nal Biology		246-251	2013 USA	britan ilikki angi 10.4 1800-1816 2013 E 2014 II. article	the five selected datasets include multiple datasets with more than 50 samples per group	Case-control study	accuracy (10 nams of 10-fold CV)	cross-validation	search the optimal feature sabed: The opportmental resists on the real-world omics discuss indicates in Microbia PMMM Resizus esistence with Cit more release to detain better (or competitive) classification accuracy than the other two measures. "I, multicates satisfication proteins possion was membodisquical and computational challenges for developing novel and effective statistical approaches. In this paper, introduce a new approach for including integrity discusses trains sourced with interesting and in	sets of genes in which the relative comparison of their expression values leads to class discrimination. For an m-class problem, the classification rule typically depends
216 Yang, 5 and Naiman, D.Q.	Multiclass cancer classification bas on gene expression comparison	Stat Appl id Genet Mol Biol 13	4	477-496	United 2014 States	https://www.ncts.nl m.ch.goo/umstarti class?PMC477527fc_article	one of the considered distanct has not have 50 amples per group ("MILE is a two-stage study where a retrospective stage is generated expression profiles for 2,148 patients and was designed for bommaker discovery. A prospective stag is produced an independent cohort or 1,152 patients and was used for 3,153 patients and was used for 3,154 patients and was used for 3,155 patients and 3,155 patients and 4,155 patients and 5,155 patients and		accuracy (LODCV + test set)	cross-validation + test set	on a small number of mejers sitt, which provide transparent describe householders and dolve for potential belonged interpretations. When that can are approach on sever common jump expension datasets and company in with puppin characteristic constructions. The provided in the company of t	on a small analyse of the gase sets, which prodest transparent decision boundaries and allow for posterior bulgged interpretations. For tiss our approach on severe common gain as expression discusses and compare a with pupping decision. See the comparent with pupping decision and analysis of the comparent and the pupping decision and exercised programs of the desirable cancer to the section of puping decision of the desirable cancer to the section of the desirable cancer to the section of the desirable cancer to the desirable cancer to desirable decisions of the desirable cancer to desirable cancer to the desirable cancer
217 Yang, Yand Huang, Nand Has, Land Rong, W	A clustering-based approach for efficient identification of microRNA combinatorial biomarkers	BMC Genomics 18		210-210	2017 China	http://dx.doi.org/10 11869s12884-017- 2498-5 article	8 unmatched tissues in GSE 40525. And in GSE22220, there are 210 samples from 219 breast cancer patients, including 84 extrogen receptor (ER)-negative tissues, and 135 ER-positive tissues."	Case-control study	accuracy, sensitively, specificity (5-fold cross-validation)	cross-validation	sensitivity and specificity then integle gene blommeters, to order to worde enhancitive search and reduced principal control and the search and reduced principal control and the combination of preparentative duster members are assessed as potential combination of representative duster members are assessed as potential excellent and the search of the	sensibility and specificity than single-gree bloomsfare. In order to work enhancing such and reduced for the control of the c



227 Zou, M and Liu, Z and Zhung, X S and Wang, Y	NCC-AUC: an AUC optimization method to identify multi-biomarker panel for cancer prognosis from genomic and clinical data	Bioinformati cs	31 20	3330-3338	2015 China	http://dx.doi.org/10 1093/baselfernatic wtm/A24 article	Outaset 1: 1981 patients with 328 basal- like tumors, 238 HER2+ tumors, 719 luminal 4, 490 luminal 8 and 200 normal like tumors, Dataset 2: 148 stage IB NSCLC patients	Cases only (prognosis study)	AUC (training + validation data)	training + test set	"In this study, we propose a novel Area Under Curve (AUC) optimization method for multi-biomarker panel identification named basered control dissulter for AUC complication (MCCAC). Or method is noticed by the connection between the control counter for a function of subject. This connection allows as to convert the survival time regression problem to analysis. This connection allows as to convert the survival time regression problem to an allow of the convert of the survival time regression problem to any conscitution problem. Then an optimization model is formulated to referrely manifest AUC and manifest in method or desicted features to construct an aprilication in the amount of control disable time and existent of survival problems. AUC and manifest in minimals the number of selected features to construct in the survival referrely appreciate in the survival referrely. AUC and manifest in minimals the number of selected features to construct a prediction in the survival referrely. AUC and manifest in minimals the number of selected features to construct a prediction in the survival referrely. AUC and manifest and the number of selected features to construct a prediction in the survival research (AUC and manifest problems of the number of selected features to a prediction in the survival research (AUC and manifest and the number of selected features to a prediction in the survival research (AUC and manifest and the number of selected features to a prediction in the survival research (AUC and manifest and the number of selected features to a prediction in the survival research (AUC and manifest and the number of selected features to a prediction in the survival research (AUC and manifest a prediction in
228 Zou, M and Zhang, P J and Chen, L and Tlan, Y P and Wang, Y	Identifying joint biomarker panel fro multiple level dataset by an optimization model	m Biomark Med	10 6	567-575	2016 China	hmolitic dei om/10 2217/8/mm/2015: 0002 article	"IDI colorectal cancer and 95 benign samples" "Whole-blood gene expression profiles were collected from a total of 573 individuals. After preprocessing, the data contained 488 gene profiles (n = 205 PD, n = 205 PD, n	Case-control study	accuracy (LOCCV + test set)	cross-validation + test set	sed holds the gromate to improve disease diagnosis accuracy L_IN contincted IDII and context cancer and the Singing rangels, resemble the melocate concentration by context cancer and the Singing rangels, resemble the melocate concentration by context cancer and the singing range of
Shamir, Ron and Elein, Christine and Amur, David and Vollstedt, Eva Juliane and Bonin, Michael and Usersonic, Manija and Visorg, Evetter C. and Matera, Mas and Potts, Som and Staff, referribel and Cox- comment of the Comment and Comment and Company and reduced, Nation, and William (Some American). Administration of the Comment of th	uol .	Neurology	89 16	1676-1683	2017 Croatia	http://like.doi.org/10 12/12/are/0000000 000004516. article	contained 486 gene profiles (n = 205 PD, n = 233 control), n = 48 other neurodegenerative diseases) that were partitioned into training, validation, and independent test cohorts to identify and validate a gene signature."	Case-control study	AUC (cross-validation + external test set)	cross-validation + external cohor validation	to 64 composition de 27 decembrações de 27 decembra
230 Glaab, E	Using prior knowledge from cellular pathways and molecular networks fe diagnostic specimen classification	or Bioinformati	17 3	440-452	2016 England	http://dx.doi.org/10 160/3/hbhbhs/044 article	review (not applicable) "A total of 792,779 genomic and clinical data points from 3,421 pts were analyzed. The cohort was comprised of five independent datasets: 443 pts from the Beat AM Master Trial Triver et al.	raview			provides an overview of these recent developments and compares pathway- and matterists bear appeared inscriptions appeared in term of their celling in improving model reductions, accuracy and biological interpretability, Different cross: improving model reductions, accuracy and biological interpretability, Different cross: improving model reductions, accuracy and biological interpretability, Different cross: improving model reductions, accuracy and biological interpretability. Different cross: improving model reductions, accuracy and biological interpretability. Different cross: are discussed, and a pressure study is presented an example.*
Shreve, J and Meggendorfer, M and Awoda, H and Mabbergie, S and Walter, W and Hutter, S and Maldhoul, A and Hitton, C and Relation/Ch. N and Ralgaz, Y and Roophal, Y and Adems, Y and Sec C M and Pattle B. J and Razmanov, T and Macingievell, J P and Haferlach, C and Selenes, M A and 231 Hollerfach, T and Razlin, A		Blood	134		2019 USA	http://like.doi.org/10 11828book.2019- 128508 abstract	Nature, 2018), 855 ptc from Cleveland: Clinic, 414 ptc from Munich Leukemia Laboratory (MLL), 1,509 ptc from the German-Austrian Study Group (Papaemmanuil et al, NEIM, 2016), and 200 ptc from The Cancer Genome Attas NEIM, 2013). "According to WHO standards 185/73" "According to WHO standards 185/73"	Subtype categorization	C-index (training * test cohorts)	external cohort validation	"Genomic alterations have a differential impact on OS [overall survival] in each cycapeneter risk group, highlighting the complexity of horoporating these metations for the stratification. A permanishel prediction and order based on chinical parts date an accurately provide survival unique to each individual pit and can significantly desperform full. European testemaked; classifications or any currently available models."
Baer, C and Walter, W and Stengel, A and Hutter, S and Meggendorfer, M and Kern, W and Haferla 232 C and Haferlach, T	Molecular classification of AML-MRC reveals a distinct profile and identifie ch, MRC-like patients with poor overall survival		134		2019 Germany	http://dx.doi.org/10 .11829bood-2019- 128294 article	male; 69 female). The non-MRC cohort (n=573) represents a heterogeneous AML population incl. the WHO defined recurrent cytogenetic abnormalities or t- AML (301 male, 273 female)."	Case-control study	true positive rate, false positive rate (10-fold cross-validation)	cross-validation	"Using partient" history and genetic information instead of morphology allow to identify \$6-99% of AMA. AME as defined in WVO today" "If this result, we represent the most result parties with worst that used deep learning to ball models for cracer prognosis prediction. Deep learning to ball models for cracer prognosis prediction. Deep learning to a deep learning to ball of models for cracer prognosis prediction. Deep learning to ball models for cracer prognosis prediction. Deep learning to ball of models for cracer prognosis prediction. Deep learning to ball of models for cracer prognosis prediction. Deep learning to ball of models for cracer prognosis predictions are also also also also also also also also
233 Zhu, W and Xie, L and Han, J and Guo, X	The application of deep learning in cancer prognosis prediction	Cancers	12 3		2020 China	http://dx.doi.org/10 .3390/cancers120 30903 article	review (not applicable) "The cohort included 189 subjects with CRC 115 with advanced adenoma (AA)	Review			none access prediction when working with large amounts of data. The application of deep learning in carbon prognosis has been shown to be equivalent or better than current approaches, usin as Can PM ² . **Contract approaches, usin as Can P
Angelino, P and Hosseinian Ehrensberger, S and Clarloni, L and Despraz, J and Dotta, G and Perez- 234 Urbs, A and Mongenthaler, S and Distoress; M	immunotranscriptomics signature in blood for early colorectal cancer detection	Annals of Oncology	30	v45-v45	2019 Switzerla	http://dx.doi.org/10 1093/annonc/mdz meeting 1093/annonc/mdz abstract	Cix., 11s with advanced adelenial (AR), 39 with other types of cancer (OC) as well as 218 individuals without any colorectal lesions (CON)." "A total of 2,602 unstimulated saliva samples were collected from 231 subjects with CRC, 99 subjects with	Case-control study	AUC, sensitivity, specificity (independent set validation)	training + test set	"supplier the fraction of the immune system to oxect of critical and exace identification may be able to the system of the syste
Kuwabara, H and Iwabuthi, A and Soya, R and Enomoto, M and Ishizaki, T and Tsuchida, A and 235 Nagakawa, Y and Katsumata, K and Sugimoto, M	Salivary metabolomics for colorectal cancer detection	Annals of Oncology	30	v46-v46	2019 Japan	http://dx.doi.org/10 1093/annonc/mdz 239.058 meeting abstract	subjects with CRC, 99 subjects with polyps, and 2272 subjects with healthy controls"	Case-control study	AUC	training + test set	"Combinations of salivary metabolites show high potential as a screening tool for CRC" "We propose a new approach using interpretable, individualized modeling to predict." "We propose a new approach using interpretable, individualized modeling to predict."
Hilton, C B and Meggendorfer, M and Seleres, M A and Shreve, J and Radakovich, N and Rouphal, and Walter, W and Hutter, 3 and Paction, E and Savona, M B and Gents, AT and Makharjee, S and Solnib, R N and Riggar, and Brick, C M and Konneyli, R S and Ma, B K and Helefisch, C and 200 Mec	Y Geno-clinical model for the diagnosis of bone marrow myeloid neoplasms	s Blood	134		2019 USA	http://dx.doi.org/10 118284604-2010. meeting 128987 abstract	"Of 2471 pts, 1306 had MDS, 223 had ICUS, 107 had CCUS, 478 had CMML, 89 had MDS/MPN, 79 had PV, 90 had ET, and 99 had PMF." "We have collected freshly frozen clinical PDAC tissues (N-46), paratumour	Case-control study	AUC ("The cohort was randomly divided into learner (80%) and validation (20%) cohorts")	training + test set	region developes based on genomic and clinical data without bore marries being data, in a special case of planting as the measurable products. The model also exceptions give with opsocials and supportion that allow for experimental ventorials and supported in a support of the
Liu, X D and Wu, H and Li, Y and Liu, X and Zhang, Z and Yu, L and Qin, Z and Su, Z and Liu, R and He 237 and Dai, M and Liang, Z	Early detection of pancreatic ductal , Q adenocarcinoma using methylation signatures in circulating tumour DNA		30	v261-v262	2019 China	http://dx.dei.org/10 1003/annonc/redz 247.013 meeting abstract	pancreas tissues (N=30), PDAC plasma samples (N=120), chronical pancreatitis plasma samples (N=90), and normal plasma samples (N=100)."	Case-control study	AUC	cross-validation	Tuling marging entrics, colotifies, an identified EAXC operfic Nutrembristion anders, more of white an fractionally entergoid with provisionly response markers. These markers are conditional bounders or non-invasive PAXC correcting. "Advances in schooling in the field of more is investigated provisional provisions "Advances in schooling in the field of more is required the discounty of "Advances in schooling in the field of more is required that discounty of "Advances in schooling in the field of more is required to the speciment of "Advances in schooling in the field of more is required to "Advances in schooling in the field of more is required to "Advances in schooling in the field of more is required to "Advances in schooling in the field of more is required to "Advances in schooling in the field of more is required to "Advances in schooling in the field of more is required to "Advances in the required to "Advances in schooling in the field of more is required to "Advances in schooling in the field of more in the "Advances in schooling in the field of more in the "Advances in the required in the province in "Advances in the required in the required in "Advances in the reduce in the field in the schooling in "Advances in the reduce in the field in the schooling in the "Advances in the reduce in the reduce in "Advances in the reduce in the reduce in the reduce in "Advances in the reduce in the reduce in the reduce in the reduce in "Advances in the reduce in the reduce in "Advances in the redu
238 Adom, D and Rowan, C and Adeniyan, T and Yang, J and Pacesony, S	Biomarkers for Allogeneic HCT Outcomes	Front Immunol	11	673-673	2020 USA	http://dx.doi.org/10 3389/fireru 2020 00923 article	review (not applicable)	Review			markets, from further, are discribed distribution of several and employed and a contract contract of the share of the contract
239 Ahmed, KT and Park, 5 and Jiang, Q and Yes, Y and Heavy, T and Zhang, W	Network-based drug sensitivity prediction	BMC Med Genomics	13	193-193	2020 USA	http://dx.doi.org/10 1.188/s12920-2002- 0829-3 article	"The feature selection methods and prediction models were tested on 144 NSCLC cell lines 8NM-seg gene expression dataset [34]. All the 144 cell lines were screened by the same drugs and the AUC and EDS scores for each drug on each cell line are available in this study."	i Drug response prediction	correlation (70% as the training set, and 30% as the test set)	training + test set	methods, and deep neural network on an NECC cell line distanct, we have made were limited inserted. The third neural design of the second distinct selection method asserted in the second of the second distinct selection method asserted in the second distinct selection method asserted asserted distincts selection method asserted asserted in the second distinct selection method asserted
240 Ahmed, Z and Mohamed, K and Zeeshan, S and Dong, X	Artificial intelligence with multi- functional machine learning platform development for better healthcare and precision medicine	n Database (Oxford)	2020		2020 USA	http://dx.doi.org/10 10936shabase/ba as010 article	review (not applicable)	Review			require addition of suefal enabyle tools, bethodages, databases and approaches (4)) to efficiently approximate relevantity and interpretability of friends, between the public health systems, as well as address enhold and coal tissues related to the public health systems, as well as address enhold and coal tissues related to the public health systems, as well as address enhold and coal tissues related to the priviley and procedure of herafishers and entire date with enhances to the support and procedure of the address enhanced and as address enhanced and enhanced to the enhanced to the priviley application of herafishers and ends as address enhanced and enhanced to the enhanced to
Ambasi, A and Ju, Y E and Lin, L and Olesen, A N and Koch, H and Hedou, J J and Leary, E B and 241 Sempere, V P and Mignot, E and Taheri, S	Proteomic biomarkers of sleep apne.	a Sleep	43	11	2020 Qatar	http://dx.doi.org/10 .1093/aleep/zaas0 88 article	"serum samples from 713 individuals in the Stanford Sleep Cohort"	Case-control study	AUC (10-fold CV + validation)	cross-validation + test set	>=15 vs OAHHC15) trained on SomaScan protein measures alone performed robustly, >=15 vs OAHHC15) trained on SomaScan protein measures alone performed robustly, achieving 76% accuracy in a validation dataset. Multiplies protein assays offer diagnostic potential and provide new insights into the biological basis of skep diagnostic potential and provide new insights into the biological basis of skep
Aslam, M. A. and Xue, C. and Wang, K. and Chen, Y. and Zhang, A. and Cai, W. and Ma, L. and Yang, Y. an 242 Sun, X. and Liu, M. and Pan, Y. and Manir, M. A. and Song, J. and Cui, D.	SVM based classification and prediction system for gastric cancer using dominant features of saliva	Nano Biomedicine and Engineering	12 1	1-13	2020 China	http://dx.doi.org/10 5101/nbe.v12/1.p 1-13 article	"220 saliva samples were collected from the non-cancerous and gastric cancerous persons"	Case-control study	AUC, accuracy, specificity, sensitivity (10-fold cross-validation. "To avoid the overfitting issue, we used an ES approach. We controlled the error of the network during the training phase and stopped the training if the model undergoes the overfitting")	cross-validation	disordered branthing." Taked on the activated results, the signoid kernel has produced the best disselfaction musts with an accessor of 97.12m Ks, specificity of 97.44Ks, and specificity of 97.44Ks
Awada, H and Durmaz, A and Gurmari, C and Kishtagari, A and Meggendorfer, M and Kerr, C M and Kumanovic, T and Durmaz, J and Nagata, Y and Radivopevisth, T and Advazel, A S and Ravanel, F an Carraway, H E and Nasha, A and Hafelefach, C and Sustuthararjal, Y and Sott, J and Visconte, V anv 243 Kantarjian, H M and Kadia, T M and Selexers, M A and Hafelefach, C and Maclojevski, J P	The application of machine learning d improve the subclassification and d prognostication of acute myeloid leukemia	to Blood	136	28-28	2020 USA	http://dx.doi.org/10 1182bbook-2020 139013 meeting abstract	"We collected and analyzed genomic data from a multicenter cohort of 6788 AML patients"	Case-control study	accuracy (cross-validation)	cross-validation	AMA planets is too van for traditional prediction methods. Using never ML methods, however, we were also the decipies and of prognostic subgroups predictive of nurvival, allowing us to move AML into the era of personalized methods, however, we were also the decipies and of prognostic subgroups predictive of nurvival, allowing us to move AML into the era of personalized medicine.*
Bader, J. M. and Geyer, F. E. and Miller, J. B. and Strauss, M.T. and Koch, M. and Leypoldt, F. and Kontrelylessy, P. and Bittery, D. and Schipka, C. G. and Incesoy, E. I. and Peters, O. and Deigendesch, N. 244 and Simon, M. and Jessen, M. N. and Zetterbeg, H. and Mary.	Proteome profiling in cerebrospinal fluid reveals novel biomarkers of Alzheimer's disease	Mol Syst Biol	16 6	e9356- e9356	2020 German	http://doc.doi.org/10 .15252/msb.20199 398 article	"From three independent studies (197 individuals), we characterize differences in proteins by AD status	Case-control chiefe	AUC (k-fold CV, k = 6)	cross-validation	"We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier resides high specificity (ETN) We found that an ensemble method-based classifier res
Tools or annually an			-								

Bergamaschi, A and Su, J and Neg, Y and Colls, F and Ellion, C and Phillips, T and McCarthy, E and McCarthy, E and Scott, A and Libyd, F and Guler, G and Arhworth, A and 245 Quale, S and Lovy, S	Epigenomic detection of multiple cancers in plasma derived cell free DNA	Cancer Research 80)	16	2020 USA	http://dx.doi.org/10 1158/1538: meeting 7-445-0402020-783 abstract	"Cancer and control patient cfDNA cohorts were accrued from multiple site consisting of 48 breast, 55 bug, 32 prostate and 2 pancreatic datasets consisting of 41 and 33 cancer subjects (Set 1 and 2) "We recruited n = 1,002 twins and unrelated healthy adults in the United Kingdom to the PREDICT 1 study and	Case-control study	AUC (S-foldCV)	cross-validation	These findings seggest that Stime Changes in CRNA make non-invalve detections with stage house, procretal, protezte, and large cancers. Frethermore, Stime profiting in CRNA may enable the prediction of finicially relocant features such a strand state of the control of stime of the control of the control of stime of the control of the control of stime of stime in the control of the control of stime of stime in the control of stime of stime of the control of s	
Berry, S E and Valdes, A M and Drew, D A and Asnicar, F and Mazidi, M and Wolf, J and Capdevila, J and Hadjigeorgiou, G and Davies, R and Al Khatb, H and Bonnett, C and Ganesh, S and Bakker, Ean Hart, D and Mangino, M and Merino, J and Unerheegt, Jand Wylart, P and Ordewas, J M and Gardner, 246 D and Delahanty, L M and Chan, A T and Segata, N and Franks, P W and Spector, T D	Human postprandial responses to C food and potential for precision nutrition	Nat Med 28	5 6	964-973	2020 UK	https://www.nature .com/articles/s415 91-920-9934-0 article	Kingdom to the PREDICT 1 study and assessed postprandial metabolic responses in a clinical setting and at home." "In this study, we used the THCA RNA-	Response to food intake prediction	correlation	training + test set	"We developed a machine-learning model that predicted both triglyceride (r = 0.47) and glycemic (r = 0.77) responses to food intake. These findings may be informative for developing personalized diet strategies."	"We developed a machine-learning model that predicted both triglyceride (r = 0.47) and glycemic (r = 0.77) responses to food intake. These findings may be informative for developing personalized diet strategies." "In conduction, 3 6MM-transcripts based SVC prediction model attained considerable
247 Bhalla, S and Kaur, H and Kaur, R and Sharma, S and Raghava, G P S	Expression based biomarkers and models to classify early and late-stag samples of Papillary Thyroid Carcinoma	je PLoS One 19	i 4	e0231629- e0231629	2020 India	http://dx.doi.org/10 1371fournal.pone 0291629 article	"In this study, we used the IHCA RNA- Seq dataset of The Cancer Genome Atlas, consisting of 500 cancer and 58 normal (adjacent non-tumorous) samples"	Case-control study	AUC (cross-validation + external validation)	cross-validation + test set	performance in segregating the early-stage and stee-stage in I/L (Papinary I trylocal Carcinoma) issue samples with F score of 0.75. In addition, precition models base on five protein-coding transcripts categorized tumorous and non-tumorous samples of patients with high F1 score of 0.97." Relatate feature immortance indirects that tumor mutational hurden is the main.	performance in segregating the early-stage and stac-stage in It (Prapitally Hyrydol Carcinoma) itsues amaples with F score of 0.75. In addition, prediction models based on five protein-coding transcripts categorized tumorous and non-tumorous samples of patients with high F1 score of 0.97." "Robitate feature immortance indirecter that tumor mutational hurden is the main
	Machine learning methods to identifysalient genomic predictors of					http://dx.doi.org/10	"paired tumor and normal whole-exome sequencing (WES) data from clinically annotated tumor specimens treated with anti-CTLA4 (n=145, melanoma) and anti-BPL theopole (n=46, MSCLC).				counts did not contribute to predictions across analyses. Novel findings include an	auxiliation of the shifts to predict ICI nations represent union random forest
248 Bigelow, E G and Baras, A and Yarchoan, M and Jaffee, E M	clinical responses toimmune checkpoint inhibitor therapy	Cancer Research 80)	16	2020 USA	7445 AM2020 meeting 5670 abstract	anti-PDI therapies (n=94, NSCLC, melanoma, head and neck, bladder cancer, and others)"	Case-control study	precision, recall ("A random forest classifier was trained and tested on various subsets of the dataset")	training + test set	Unstanding 2, non-seal mid-additionally, analyses suggest the ability to transfer learning on one numor type to identify responders in other disease." Applying elastic nets with cross-validation, we find that 12 metabolites plus patient diabetes status classify NGSY vs NAPL patients with predicted out-of-lample power AUC - 0.281. In addition to discalaropit, and bile acids, predictors include ligid and the complete of the	classifiers, and a semi-quantitative evaluation of the relative contributions of features included the model. Additionally, analyses suggest the ability to transfer learning on one tumor type to identify responders in other diseases." "Applying elastic nets with cross-validation, we find that 12 metabolites plus patient
Brown, E and Karrar, A and Hellings, S and Stepanova, M and Warrack, B and Lam, B and Onorato, J and Felix, S and Aptel, A and Jeffers, T and Rajpur, B and Charles, E and Nader, F and Luo, Y and Reily 249 M and Zhao, L and Thompson, J and Goodman, Z and Younoss, Z	NASH Identification of a Tumor	Hepatology 7	3	S409-S410	2020 USA	http://dx.doi.org/10 1016/90168- meeting 8278/20/31304-0 abstract	"Serum samples were obtained from 10 biopsy-proven NASH (F0-F4) and 50 NAFLD patients without NASH" "a proenostic model was established	Case-control study	AUC, repeated loops of crossvalidation	cross-validation	microsiome-derived metapointes." "As a tumor microenvironment-relevant sene set-based prognostic signature, the	microbiome-derived metaporites." "As a tumor microenvironment-relevant gene set-based prognostic signature, the
250 Cai, W Y and Dong, Z N and Fu, XT and Lin, L Y and Wang, L and Ye, G D and Luo, Q C and Chen, Y C	Microenvironment-relevant Gene se based Prognostic Signature and Related Therapy Targets in Gastric Cancer	Theranostic s 16)	19 8633-8647	2020 China	http://dx.doi.org/10 7150/hno.47938 article	based on gastric cancer gene expression datasets from 1699 patients from five independent cohorts"	Prognostic study	concordance index (C-index; was calculated to determine the discrimination of the nomogram via a bootstrap method with 1000 resamples)	external cohort validation	GPSGC model provides an effective approach to evaluate GC [Gastric Cancer] patient	GPSGC model provides an effective approach to evaluate GC [Gastric Cancer] patient survival outcomes and may prolong overall survival by enabling the selection of individualized argeted therapy." "In this Perspective, we provide a brief overview on the role of gut microbiome in
Cammarota, G and Ianiro, G and Ahern, A and Carbone, C and Temko, A and Claesson, M J and 251 Gasbarrini, A and Tortora, G	Gut microbiome, big data and machine learning to promote precision medicine for cancer	Nat Rev Gastroenter of Hepatol 17	,	10 635-648	2020 Italy	http://dx.doi.org/10 .1038/s41575-020- 0327-3 article	review (not applicable)	Review			concer and focus on the need, rich and limitations of a machine learning-driven op- portunity of the control of	cancer and focus on the need, role and limitations of a machine learning-driven
252 Cascianelli, 5 and Molineris, 1 and Iselia, C and Massereli, M and Medico, E	Machine learning for RNA sequencin based intrinsic subtyping of breast cancer	6- Sci Rep 10) 1	14071- 14071	2020 Italy	http://dx.doi.org/10 _1038s41598-020- 70852-2 article	"A 220-sample training set was extracter randomly from the TCGA dataset, respecting the same 60/40 EN-/ER- proportion of the PAMSO training set. Al the remaining 597 cases were instead included in the TCGA test set." "The ILL1 and ILL2 clinical trials had	Tumor subtype prediction	accuracy (10-fold CV, training/fest set)	cross-validation + external cohon validation		PAMSO algorithm; 2. Define RNA-see-based classification approaches to perform single-sample BC intrinsic subtyping with external-AWCA-based PAMSO or regularized mLR methods. These strategies appeared valuable to favor the use of RNA-seq in BC [Breast Cancer] clinical practice and are worthy of other studies on heterogeneous
	Patient ancestry significantly contributes to molecular					http://dx.doi.org/10	microarray expression data for 1566 female patients of self-described ancestry: AA (n = 216), EA (n = 1118), and NAA (mostly from South America (n = 232); to 3 countries of origin Peru (n = 232); to 3 countries of origin Peru (n = 273); and Guatemala (n = 27)] and 124 male patients of self-described ancestry: AA (n = 14), EA (n = 40).				myeloid cell signatures, logistic regression analysis determined that ancestral background significantly changed 23 of 34 gene signatures. Additionally, the strongest association to gene expression changes was found with autoantibodies, and this also had etiology in ancestry: the AA predisposition to have both RNP and dsDNA	"Although standard therapy affected every gene digitative and significantly increased impold call signatures, logistic regression analysis determined that ancested background significantly changed 23 of 34 giene signatures. Additionally, the strongest association to gene expression changes we stored with autoantibodies, and this also had deology in executy: the AM predisposition to have been RNP and doDNA autoantibodies compared with EA predisposition to have only anti-doDNA. Amachine learning approach was used to determine a gene signature characteristics to
Catalina, M D and Bachall, P and Yeo, A E and Geraci, N S and Petri, M A and Grammer, A C and 253 Upsky, P E	heterogeneity of systemic lupus erythematosus	JCI Insight 5		15	2020 USA	.1172(jci.insight.14 0380 article	93), NAA (n = 17)"	Subtype categorization	accuracy (10-fold CV)	cross-validation	B cell axis in AA SLE patients." "After correction for repeated measures, clustering identified 3 separate metabolic/clinical profiles: 1) right ventricular (RV) dysfunction, arrhythmia and	B cell axis in AA SLE patients." "After correction for repeated measures, clustering identified 3 separate metabolic/clinical profiles: 1) right ventricular (RV) dysfunction, arrhythmia and
Codars, A.M. and Manihoc, C and Ko, J and Bettiglier, T and Weingarten, A and Opotowsky, A and 254 Kuty, S	ARTIFICIAL INTELLIGENCE FACILITATED METABOLOMIC PROFILING IN ADULT CONGENITAL HEART DISEASE (ACHD)	Journal of the American College of Cardiology 75	5	11 552-552	2020 USA	http://dx.dei.org/10 1016/50735 meeting 1097(20)31179-7 abstract	"We analyzed 674 metabolites in 186 serum samples from 101 non-fasting ACHD patients followed regularly at a single institution, including repeated samples at different times."	Subtype categorization	AUC (cross-validation)	cross-validation	dypanea (n=107), 2) complex bisentricular disease with hypoxia and lower deducational level (n=79) and 3) individuals managing well with valuatiar and septal diefects (n=21). Metabolomic data permitted the creation of models associated with prevalent arrhythmia (cross validated AU.O.S.P), seath-reported wearrhythmia (cross validated AU.O.S.P), seath-reported validated seath-reported wearrhythmia (cross validated au.O.S.P), seath-reported validated valid	dyprate [1-07], 2] complex bivertricular disease with hypoxia and lower educational level [1-07] and 3] individual managing well with valuation and septial deflects [1-02]. Metabolomic data permitted the creation of models associated with prevalent arrhythmic forces utilized ALI COS, patient resported elevations associated with prevalent arrhythmic responsibility. ALI COS, patient resported elevations (ALI COS) and 69 dyfurthriction (ALIC COS). In "While ALIM metabos are powerful, they are used il silinative provisions (inclications in that physician interpretation is result for implementation in a real-world setting. We should allow be opposition of the year the silination of the business of the silination of the silination of the silination of the silination of the business of the silination of the silination of the silination of the silination of the silination of the silination of the silination of the silination of the silination of the silination of the
255 Chan, S and Reddy, V and Myers, B and Thibodesus, Q and Brownstone, N and Liao, W	Machine Learning in Dermatology: Current Applications, Opportunities, and Limitations	Dermatolog y and Therapy 10	3	365-386	2020 USA	http://de/dei.org/10 .1007/s/13555-020. 00372-0 article	review (not applicable)	Review			nature of these algorithms. It is also important to make these technologies inclusive or skin of color, Further research in Mt. should be transparent by making algorithms and datasets available to the public for further validation and testing," "Urclogy is a contantly changing sociality with a vide range of therapeutic breakthrought, a huge understanding of the genomic expression profiles for each understanding and the properties of the contantly and the contantly contained to the contained of the contained the contained the contained the contained and the contained the	nature of these algorithms. It is also important to make these technologies inclusive of skin of color. Further research in Michould be transparent by making algorithms and datasets available to the public for further validation and testing." "Urology is a contently changing separably with a wide range of therapeutic breakthrought, a huge understanding of the genomic expression profiles for each value of the public of the profiles of the promise promise proposition profiles for each value of the public of the profiles of the promise profiles for each value of the public of the profiles of the promise profiles for each value of the public of the profiles of the profiles of the profiles of the profiles for each value of the public of the public of the profiles of the profiles of the profiles for each value of the public of the public of the profiles of the profiles of the public of
256 Chenorriago, J and Moreno, C	Precision Medicine, Artificial Intelligence, and Genomic Markers in UrologyDo we need to Tailor our Clinical Practice?	n Urologia Colombiana 25	э з	158-167	2020 Colombia	http://dx.doi.org/10 _1659-0040_ 1714148 article	review (not applicable) "Gene expression, protein expression and copy number variants are used to predict estrogen receptor status (BRCA-ER, N = 381) and breast invasive carcinoms subtypes (BRCA-subtypes, N	Review			clinical practice."	background that comes with them. A critical analysis of these new technologies and pharmacological breakfroughs should be made before considering changing our clinical practice."
Chienis, M. and Buscola, N. and Marcolinis, A. and Francescatto, M. and Zandovo, A. and Trastolla, L. and 257 Appetituals, C. and Auman, G. and Furbandis, C.	I Integrative Network Fusion: A Multi- Omics Approach in Molecular Profile	Frontiers in ng Oncology 10	i.		2020 Italy	http://dx.doi.org/10 31899/org-2020.0 1095 article	carcinoma subtypes (BRCA-subtypes, N. 305), while gene expression, miRNA expression and methylation data is used as predictor layers for acute myeloid leukemia and renal dear cell carcinoma survival (AMIL-OS, N = 157; KIRC-OS, N = 181)*	Disease status, subtype and survival prediction	Matthews Correlation Coefficient (10 × stratified Monte Carlo cross- validation (50% training/validation proportion))	cross-validation + external cohon validation	the proposed approaches are task and/or data dependent, the complexity of numor analysis suggests that network-based approaches are needed [] In this context, it is clear that omici-integration is one of the most promising and demanding challenges of the modern bioinformatics, and hart there is an urgent need to prove the reproducibility, interpretability, and generalization capability of the proposed methods:	analysis suggests that network-based approaches are needed [] In this context, it is clear that omic-integration is one of the most promising and demanding challengs of the modern bioinformatics, and that there is an urgent need to prove the reproducibility, interpretability, and generalization capability of the proposed methods:
Czolk, B and Blueber, J and Sprensen, M and Wilmes, P and Codeans-Morel, F and Slev, P S and 238 Hilger, C and Bindslev-Jensen, C and Olfert, M and Kushn, A	IgE-Mediated Peanut Allergy: Curren and Novel Predictive Biomarkers for Clinical Phenotypes Using Multi-Omi Approaches			594350- 594350	Luxembou 2020 g	http://dx.doi.org/10 33896mmu 2020 504550 article	review (not applicable)	Review			phenotypes." "Perceptidial democratic (RCM) is a funcial infection tunically found in Latin American	powerful multi-omics technologies, togother with integrated data analysis, network- based approaches and unbiased machine learning holds out the prospect of providing clinically useful biomarkers or biomarker signatures being predictive for reaction phenotypes." "Bioteconfoliation property (PCMI) is a fund allegation basicable found in Latin American
de Olisira Lima, E and Noamm, I, C and Morbibba, K N and Esnikawa, C M and Rodrigues, R G M and Dubby, M Z and so Olivino, D N and Dubby, L N Z and so Olivino, D N and Dubby, L N Z and so Shire Riserry, M and 259 Vocentri	Metabolomics and machine learning d approaches combined in pursuit for more accurate paracoccidioidomycosis diagnoses	mSystems 5	3		2020 Brazil	http://dx.doi.org/10 .1128/mSystems.0 0258-20 article	"In total, 343 individuals were included in this study, regardless of age and gender, in two main groups: the test group, consisting of PCM patients (n = 85), and the control group (n = 258)	* Case-control study	accuracy, senditivity, specificity (patients' samples were randomly split into the partition (Pffil) and test partition (Pfiel) in the proportion of 80% and 20%, respectively. Classifice were laranded and visitate in all staps of the method partition (Pfil) and the partition (Pfil) in the proportions of 80% and 20%, respectively).	t cross-validation	countries, repeatally in Brast. The fide infentition of this diseases based on techniques and the register of the fide individual confidence of the diseases based on techniques and the register of the regis	confidence than the routine methods employed today." "Technological advances now enable have scale datasets, including DNA and RNA.
	Integrated multi-omics approaches t	0				http://ibte.dei.org/10					and groups of patients along the genotyper-phenotype continuum of chronic kidney disease (COI). The ability or combine fresh high-dimensional adiataset, in which the number of variables exceeds the number of clinical outcome observations, using computational papinosches such as machine learning, provides an opportunity to re-classify patients in modecularly defined subgroups that better reflect underlying disease mechanisms. Patients with COI are uniquely pointed to benefit from these interatible, multi-cine approaches such as a computation of the kidney beloop. Who don't alive sample.	and groups of patients along the genotypes-phenotype confinuum of chronic kidney disease (CSC) Fall billy to combine them high-dimensional statesets, in which the number of variables exceed the number of clinical patients observations, using compactional approaches task an ematches learning provides an apportunity to re- clarify patients into molecularly skillinds along the state of the control of the control of the clarify patients and control of the control of the control of the control of the clarify patients and control of the control of the control of the patients and control of the control of the control of the patients and control of the control of the control of the patients and control of the patients and control of the patients and control of the patients and control of the patients and patients an
260 Eddy, S and Mariani, L H and Kretzler, M	improve classification of chronic kidney disease Recent advances in precision	Nat Rev Nephrol 16 Curr Opin	5	11 657-668	2020 USA	.1038is41581-020- 0288-5 article	review (not applicable)	Review			used to generate these different types of molecular data are frequently obtained during routine clinical care." "Newly proposed biomarkers offer precise and noninvasive ways to monitor patient's status. Cell-free DNA quantitation is increasingly explored as an indicator of allograft.	used to generate these different types of molecular data are frequently obtained during routine clinical care." "Newly proposed biomarkers offer precise and noninvasive ways to monitor patient's status. Cell-free DNA quantitation is increasingly explored as an indicator of allograft
261 Fu, Sand Zarringar, A	medicine for individualized immunosuppression	Organ Transplant 2! United European	5 4	420-425	2020 USA	1097imot.000000 0000000771 article	review (not applicable) "We used whole metagenome sequencing to analyse composition and	Review				Treaty is policios to liciniare de l'order i fescaria and confidente auto y cominando placente status. Celli-rea Old quantifation in increasingly explored as an indicator of allograft injury and rejection, which can help avoid unneeded biopsies and more frequently monitor graft function." "we show that features of gut microbiome, in combination with already used fecal
Gacesa, R and Vich Vila, A and Collij, V and Imhann, F and Wijmenga, C and Jonkers, D M A E and 262 Kurishikov, A and Fu, J and Zhernakova, A and Weersma, R	Microbiome and fecal biomarkers ca diagnose and classify inflammatory bowel disease		8	166-167	Netherlan 2019 s	msps://orkensibrary .wiky.com/doi/10. d 1177/2050640619 854670 abstract	function of microbiome of fecal samples of 181 IBS patient, 380 IBO cases and 859 healthy controls*	Case-control study	AUC, sensitivity, specificity (independent set validation)	training + test set		
Garcia, 5 and Lauritsen, J and Zhang, Z and Dalgaard, M D and Nielsen, R L and Daugaard, G and 263 Gupta, R	Prediction of nephrotoxicity associated with cisplatinbased chemotherapy in testicular cancer patients	JNCI Cancer Spectrum	4	3	2020 Denmark	https://doi.org/10.1 093/poics/plass03 2 meeting abstract	"Of 433 patients assessed in this study, 26.8% developed nephrotoxicity after bleomycin-etoposide-cisplatin treatment."	Case-control study	AUC ("Training and testing of the algorithm was performed with a 5 outer, 2 inner fold nested cross-validation")	cross-validation + external cohort validation	The this study, we were able to predict patients at risk of developing nephrotoxicity after BE chemotherapy based on clinical and genetic features with a machine learning algorithm (Incilla features selected on the random forest-driven baseline clinical model were known risk factors of renal toxicity and were statistically significant in univariate analysis."	"In this study, we were able to predict parients at risk of developing rephretanisty and the EBP chamedrospy based on clinical and genetic features with a machine learning algorithm. Clinical features selected on the random fovests-driven baseline clinical model were known risk factors of renal tracicity and were statistically significant in univariate analysis."

Gindin, Yand Chuang, J Carid Billin, A and Cannago, M and Hors, R and Chung, C and Myers, R P at 264 Youncesi, Z M and Harrison, S A and Americ, CM and Loomba, R	A random forest classifier based on 30-gene signature distinguishes of patients with bridging fibrosis from those with cirrhosis due to NASH	a Hepatology	72 1	43A-44A	2020 USA	http://dx.doi.org/10 meeting 1602/hep-31578 abstract	"The study included 1,120 adults with advanced fibrosis (F3-F4) due to NASH enrolled in the simturumab and STELLAB trials (discovery cohort, n-994) and the ATLAS trial (validation cohort, n-126)" "We obtained 2.73 TEP expression	Case-control study	AUC (training + validation cohort)	training + test set	gene expression signature that accurately differentiates NASH patients with cirrhosis from those with bridging fibrosis. The functional activities of these genes may suggest	*A machine learning technique applied to hepatic transcriptomic data identified a 30- gene expression signature that accurately differentiates MASH patients with crimosis from those with bridging fibrosis. The functional activities of these genes may suggest novel drivers of fibrosis progression in NASH.*
Growers, C. and Chanks, S. and Thakrul, D. and Partt, H. and Verma, P. and Malik, P. S. and Jayadevo ar 265 Gupts, R. and Alsaja, G. and Sengupts, D. Graham, S. A. and Lea, E. E. and Jeste, D. V. and Van Putters, R. and Twannier, E. W. and Nebeker, C. and 264 Tameshu, Y. read ten, H. C. and Depty, C. A.	based diagnosis of NSCLC Artificial intelligence approaches to neoficting and detectine cognitive	od-BMC Genomics : I Psychiatry Research 2	21 1	744-744	2020 India 2020 USA	https://doi.org/10.1 1880/12864-020- 021672 article https://doi.org/10. 1016/10.0rg/218.12 019.112732 article	profiles spanning oix cancer types: non- maill-call lang career (PISCLT): 59, colovertal cancer (CRC): 44, globilations multiforms (GBM): 40, breast cancer (BRC): 38, pancreatic cancer (PIC): 33, hepatobilary cancer (PIC): 51 and other to the cancer samples, platelets from 54 healthy individuals were also profiled."	Case-control study	AUC (LODO)	cross-validation	"Autificial Intelligence has great potential to advance diagnosis and treatment of patients with neurocognitive disoriers. Multi-feature datasets can improve personalization and predictive ability of machine learning algorithms in healthcare. Development of Explainable Artificial intelligence is warranted to establish trust in models for clinical decision-making."	in determining the existence of cancer, Smills strategies can be developed for inferring the potential center types. In all these case, the goar possile received to be desired, the potential center to be desired, the proposal center to be developed to the case. The case of the case
267 Gumari, A and Sammouda, R and Al-Rukhami, M and AlSsimon, H and B-Zaart, A	Feature selection with ensemble learning for prostate cancer diagnos from microarray gene expression		27 1		Saudi 2021 Arabia	http://dx.doi.org/10 117771460458521 989402 article	"The experiment of this study is conducted on a public dataset of microarray prostate cancer gene expression, consisting of 102 tissue samples (52 prostate tumor and 50 normal tissues)" "The R package TCGA-assemble2 [13] was used for data collection and we obtained 298 samples concluded three	Case-control study	accuracy (10-6dd CV)	cross-validation	data of gene expression. A set of experiments are conducted on a public benchmark dataset using 10-fold cross-validation technique to evaluate the proposed approach. The experimental results revokaled that the proposed approach attains \$5.009% accuracy, which is higher than related work methods on the same dataset."	random committee (RC) ensemble learning to detect prostate cancer from microarray data of gene expersion. A set of experiments are conducted on a public benchmark dataset using 10-feld cross-validation technique to evaluate the proposed approach. The experimental resist revealed that the proposed approach astrain 55.098% accuracy, which is higher than related work methods on the same dataset."
266 Guo, LY and Wu, A H and Wong, YX and Zhang, LP and Choi, H and Llang, X F	Deep learning-based ovarian cancer subtypes identification using multi- omics data	BioData Mining	13 1		2020 China	http://dx.doi.org/10 _1188is13040-020_ 00222-x article	types of omics data: mRNA-seq data (UNC Illumina HISeq, RNA-seq VZ), mRNA-seq data (BCSC Illumina HISeq) and copy number variation (CNV) data (BROAD-MIT Genome wide SNP_6)*	Tumor stratification	silhouette score (external test dazaets)	training + test set	nathways identified based on the classification ultimes have been proved to be	The independent text results in three CEO datasets proved the robustness of our model. The Iterature reviewing those 1y 5 (find) becaminate and 8 (42.19) (\$150 G pathways identified based on the classification ubtypes have been proved to be associated with orderin cancer. "Precision medicine matches each individual to the best treatment in a way that is considered in the control of the state
269 Hajirosodiha, I and Elemento, O	Precision medicine and artificial intelligence: overview and relevance to reproductive medicine	Fertil Steril 1	114 5	908-913	2020 USA	http://dx.doi.org/10 _10166_fertnstert 2 020.00.156article	review (not applicable)	Review			suggestions, and inexactionities, a west as images reclaiming the superior modalities in precision medicine require the sust of artificial intelligence and machine learning. This modern view of practions medicine, adopted early in critaria areas of medicine such as center, has started to impact the field of reproductive medicine. "In summary, we constructed and evaluated AMD [Age-related macular degeneration] prediction models integrating 35MPs and 2 clinical factors. The four models showed AMDC above 0.72 in the training set. Machine learning took have	integration of these modalities in precision medicine require the use of artificial intelligence and machine learning. This modern view of precision medicine, adopted early in certain areas of medicine such as cancer, has started to impact the field of reproductive medicine."
Hao, S and Bai, J and Liu, H and Wang, L and Liu, T and Lin, C and Luo, X and Gao, J and Zhao, J and 270 H and Tang, H	Comparison of machine learning too for the prediction of AMD based on U, genetic, age, and diabetes-related variables in the Chinese population	Regenerativ e Therapy	15	180-186	2020 China	http://dx.doi.org/10 .1016/j.reth.2020.0 9.001 article	"This study included a total of 202 subjects, comprising 82 AMD patients and 120 control subjects." "We finally obtained 488 primary breast	Case-control study	AUC (4-fold CV)	cross-validation	models showed AUNCs above 0.72 in the training set. Machine learning tools have the potential to all in the early diagnost and treatment of patients with AUN. There is still a way to go before the models can be applied in the clinic for AMO prediction, and they should be wildlested as larguer colonist." "We integrated comatic mutations and proviously used data types, including Eup, CIVI, Mothy, and protein, using MRI to predict breast cancer patient survivol. Applying mRMR selected features and MRI classification, we donn dis that the data of the control	the potential to aid in the early disgnosis and treatment of patients with AMD. There is still a way to go before the models can be applied in the clinic for AMD prediction, and they should be validated in a larger cohort. "We integrated somatic mutations and previously used data types, including Exp. CRIAN Methy and question series MRT to predict the post proper existent supplied.
271 Hs, Z and Zhang, 1 and Yuan, X and Zhang, Y	Integrating Somatic Mutations for Breast Cancer Survival Prediction Using Machine Learning Methods Development of circulating free DNA	Frontiers in Genetics	11		2020 China	http://dx.doi.org/10 3389fdgene 2000. 632901 article	tumors together with survival time, and all samples of them included all of the five afforementioned genomic data types. The details of our dataset are illustrated in Table 1." "Candidate markers were developed into a targeted sequencing panel and were	Case-control study	AUC (entire datasets were randomly divided into a learning dataset (80% of the entire dataset) and validation dataset (20%))		integration of somatic mutations enriched the diversity of features and was conductive to the improvement of the prediction models. In all, integrating promising data sources such as commacis mutations and harnessing the powerful feature selection method mBMR and the effective data fusion method MML can increase the prediction accuracy of breast cancer patient survival.*	data sources such as somatic mutations and harnessing the powerful feature selection method mRMR and the effective data fusion method MKL can increase the prediction accuracy of breast cancer patient survival."
Heng, Sand Si, 2 and Li, Jand Yu, Sand Lin, Band Re, 2 and Zhang, Qaard Gen, 2 and Is, Wand Pe 275 Sand Cheng, Lead Ne, Cand Sill, San And Xiao, H Hoshino, A and Kim, H S and Biginez, Lind Gipinez, I E and Coeff, M and Hermondez, 1 and Zhankino Pand Bedringer, Gen Ad Mollan, a Med Hensel, Sand Micha, H and Salmer, Land Bessel, Lind and Lucent, Sand Ci Giammata, A and Other, Are A Makajino, M and Williams, C. and Regoles, I and and Lucent, Sand Ci Giammata, A and Other, Are A Makajino, M and Williams, C. and Regoles, And and Balley, R. And add Islands, I Sand Williams, I Sand Sand Sand Markey, A Mark Garcia-Sandon, and Balley, R. And add Islands, I Sand Williams, I Sand Sand Sandon, And Sandon, And Sandon, And Sandon, And Control, Sandon, Sandon, And Control, Markey, And Control, Mark Garcia-Sandon, And Sandon, And Control, Sand Chang, And Chang,	ng, methylation markers for thyroid nodule diagnostics c. d.	Annals of Oncology :	31	\$1362- \$1362	2020 China	http://dx.doi.org/10 10/166.armon_20 20.10.501 abstract	a Legentus experience (parent arts were validated on plantamo DRA samples (115 PTC, 102 BTN)*	Case-control study	accuracy, sensitivity, specificity (training/fast set)	training + test set+extcohort	thyroid reader based on their maligrancy. They are thus promiting cardidates to develop non-invasive diagnosts. For thyroid causer screening."	"Our study demonstrates that DNA methyllation markers can recountly differentiate through the control of the co
Fand Simpson, A. Land Senger, M. and Maule, C. Mand Simonou, O. Mand Sim, M. and Ollayar, C. M. and Simonou, Simonou, Simonou, Simonou, C. M. and Simonou, C. M. and Simonou, Simonou, Simonou, Simonou, C. M. and Simonou, Sim	M P Y A	Cell 1	182 4	1044- 1061.e18	2020 Japan	http://dx.doi.org/10 1016/i.cell.2020.0 7.000 article	"To confirm that EVPs are ideal diagnostic tools, we analyzed proteomes of TE- (n = 151) and plasma-derived (n = 120) EVPs". "We investigated the two follow-ups of	: Case-control study	sensitivity, specificity (10-fold CV $+$ external test set)	cross-validation + external cohor validation	TEs and plasma, which can classify tumors of unknown primary origin. Thus, EVP proteins can serve as reliable biomarkers for cancer detection and determining cancer type."	particles (ang., including immunoglobulin, revolved most 95% sensitivity/90% specificity in detecting cancer (Finally, we defined specific (FVP proteins in Indetecting cancer, Finally, we defined specific (FVP proteins in TEs and plasma, which can classify tumors of unknown primary origin. Thus, EVP proteins can serve as reliable biomarkers for cancer detection and determining cancer type."
Huang, J and Histh, C and Coric, M and Trill, M and Adam, J and Zuburth, S and Pretin, C and Vizing and Natio, I and Schwerz, MF and Neschen, S and Extendible, G and Subur, K and Lany, M and 224 Schless, F and Engine, C and Adamski, I and Hooke do Angelos, M and Peters, A and Moning Suttley. F	Machine Learning Approaches Revo L. Meetablot Signitures of recident Chronic Edding Josean in Indiabata With Prediabetes and Type 2 Diabet	s	69	12 2756-2765	2020 Germany	benistra aktanittä 2207mioh delle article	the longitudinal cohort KDRA survey 4, conductand in the use of Angeborg. Southern Germany. The first follow up 12–13 years search of Angeborg 12–13 years search of 12–13 years search of 12–13 years search of 12–13 years search of 14–13 years search of 14–13 years search of 14–13 years search of 14–13 years search ye	Case-control study	AUC (three-step feature selection with 100 random repeats of 10 fold cross variations)	cross-validation	"This singulated its day revealed algorithm as committed or displays," and dispressionships (i.e., text. 12 and 1	"This longification study revealed agrificant accumulation of spinigo- and givenerphospholipies (Let EL and PC as CL8) in invideable with postdiothers and T20 up to 6.4 years before their clinical cented of CD. These condition metabolitic values of the control of CD. The condition metabolitic values of the control of CD. The condition of the co
Huang, Z and Johnson, T S and Han, Z and Helm, B and Clao, S and Zhang, C and Sidama, P and Risks 275 M and Yu, C Y and Ching, J and Xiang, S and Zhan, X and Zhang, J and Huang, K	Deep karning based cancer survival III, prosposis from ISM-seq data: approaches and evaluations	BMC Med Genomics :	13	41-41	2020 USA	http://de.del.com/10 1.1886419904-6005 0008-1 article	Squimoux Cell Carcinoma and Carcinoma (CSC); (4) Indea Nack Squimoux Cell Carcinoma (MCNC; §5) (about Nack Squimoux Cell Carcinoma (MCNC; §5) (about Nack Squimoux Cell Carcinoma (MCNC; §5) (about Nack Squimoux Cell Carcinoma (BCNC); (7) Lave Replace(Malta Carcinoma (BCNC); (7) Lave	Cancer survival pragnosis	Crides, p-value of log-rank test (Each dataset was split into training, validation, and testing sets in a proportion of 80, 20, and 20% respectively)	training + test set	superior performances comparing to traditional machine learning models. Among the	three Deep Learning-based models tested, we observed that Cox-nnet, which has the most succinct neural network structure, resulted in better prognosis performances in the measurement of concordance index and p-value of log-rank test. We showed that
Jiang, J and Yao, H and Yang, L and Li, J and Xin, J and Shi, D and Liang, X and Cai, Q and Ren, K and 276 Chen, X and Li, J	acute-on-chronic liver failure Next Generation Sequencing and		70	162A-163A	2019 China	https://easidoube.or nineibrary.wiey.c omidoi/10.1002he p.30940 meeting abstract	n=102), and normal controls (NC n=65) from a prospective multi-center cohort were subjected to a robust and highly streamlined single-run quantification proteomic analysis*	Cancer progression and prognosis prediction	AUC (training + validation cohort)	external cohort validation	disease progression, severity and prognosis."	Related Acute-on-Chronic Liver Failure) disease progression. And the combinatorial predict model provides fundamental information of multiple biomarkers response the disease progression, severity and prognosis."
277 Jovčevska, I	Machine Learning Technologies Are Painting the Epigenetic Portrait of Glioblastoma	Frontiers in Oncology	10		2020 Slovenia	http://dx.doi.org/10 .3389/fonc.2020.0 0798 article	review (not applicable)	Review			"Using a combination of phenotypic, genotypic, and epigenetic parameters in glioblastoma disposition will bring us closer to precision medicine where therapies will be tailored to suit the genetic profile and epigenetic signature of the tumor"	"Using a combination of phenopyic, genotypic, and opigenetic parameters in glioblastoma disposotics will bring us closer to precision medicine where therapies will be tailored to suit the genetic profile and opigenetic signature of the tumor"

Kandimalla, R and Xu, J and Link, A and Matsuyama, T and Yamamura, K and Furker, I and Ustake, H and Hentander-Blac, E and cuzen, J and Buzzarani, E and Trail, 5 and Grave, D and Meliter, 5 I and 28 Bubb, H and Brank, R and N to Hill, 20 and Blagger, F and N. W and Good, A	EpiPanGi-Dx: A cell-free DNA methylation fingerprint for the early Ca detection of gastrointestinal cancers Re Analysis and prediction of	ancer esearch 80		16	2020 USA	http://dx.doi.org/10 1158/1558- 7445 AM/2020- meeting 1084 abstract	"Using this approach, we sequenced 300 plasma specimens from all GI cancers, as the sa age-matched healthy control" "Eight datasets containing a total of 350	Case-control study	AUC (training + validation cohort)	external cohort validation	free DNA methylation biomarkers that offer a robust diagnostic accuracy for the	"Utilizing a novel biomarker discovery approach, we provide first evidence for cell- free DNA methylation biomarkers that office a robust diagnostic accuracy for the identification of specific cancer space, and demonstrate their potential clinical cancer." I have been a specific cancer panel for the early detection of all gestratinetistical cancers." "Prediction models developed based on three genes categorized CCA.
279 Kaur, H and Bhalla, S and Garg, D and Mehta, N and Raghava, G P S	cholangiocarcinoma from Jo	ournal of epatology 73		\$16-\$17	2020 India	http://dx.doi.org/10 1016/S0168 meeting 8278/20/30593-6 abstract	CCA, 133 adjacentnon-tumorous and 90 HCC samples" "The dataset contains gene expression		AUC, accuracy (training + validation set)	training + test set	[Cholangiocarcinoma] with high precision. Thus, they can be further explored for their diagnostic and therapeutic potential for CCA."	Production modes assessed as asset on mere genes caugarized CCA. (Cholangiocarrinoma) with high precision. Thus, they can be further explored for their diagnostic and therapeutic potential for CCA.*
280 Chochnejat, M and Karoosi, K and Banaei Moghoddam, A M and Mossavi Movahedi, A A	Unraveling the molecular heterogeneity in type 2 diabetes: a potential subtype discovery followed by metabolic modeling Gr	MC Med enomics 13	1	119-119	2020 Iran	http://dx.doi.org/10 21200/w.2.20484/ yZ. article	The dataset character gene expression data from participants with glacose tolerance ranging from normal to newly diagnosed T2DM, in which 91 and 63 individuals were healthy and diabetic, respectively."	Case-control study	AUC, accuracy, F1 score, precision, recall (10-fold CV)	cross-validation	in each cluster." "In this study, we tested if the incorporation of network analysis into an ML framework could accurately identify robust drug-response biomarkers using organoid	healthy controls with approximately 50 percent accuracy. Clustering of diabetic patients according to their gene expression patterns and ubsequent more in-depth analysis of each cluster unraveled specific abnormalities leading to insulin resistance in each cluster." "In this study, we tested if the incorporation of network analysis into an ML framework could accurately identify postular draw response insurance unique organization."
281 Kong, J and Lee, H and Kim, D and Han, S K and Ha, D and Shin, K and Kim, S	Network-based machine learning in colorectal and bladder organoid models predicts anti-cancer drug efficacy in patients Co	at ommun 11	1	5485-5485	Republic o 2020 Korea	http://dx.doi.org/10 d 1038945467-020- 19313-8 article	"We downloaded the FPKM-UQ (upper quartile) dataset from TCGA data portal for expression analysis"	Case-control study	"The final predictive performance was measured by comparing the correlation between the observed and predicted drug responses in the test see (FIZ)" "Figil the organized dataset into training (60%), validation (10%), and test (30%) sets", 3-fold CV for training)	training + test set	results for drug-response prediction, which were further tested in external experimental datasets." "For many compounds, even a very small subset of drug-related features is highly predictive of drug sensitivity. Small feature sets selected using prior knowledge are more productive for drugs the sensitivity. The productive sets selected using prior knowledge are more productive for drugs the sensitive productive and outbrooks with a drugs and sensitive productive sensitive se	responses, whereas conventional ML approaches showed less optimal predictive performances: Importantly, our relevok-classed ML model produced interpretable results for drug-response prediction, which wave further tested in external superimental disasses, some a very small subset of drug-related features is highly repredictive of drug semicities, Small relatures sets selected using prior knowledge and
282 Koras, K and Juraeva, D and Kreis, J and Mazur, J and Staub, E and Szcturek, E	Feature selection strategies for drug sensitivity prediction Sc Target analysis of volatile organic	ti Rep 10	1	9377-9377	2020 Poland	http://dx.doi.org/10 _1038941598-020- 65927-9 article	"The total set of samples consisted of 983 cancer cell lines originated from 13 tissue sites	Drug response prediction	Correlation, RMSE (3-fold CV on training data + test set evaluation)	training + test set	wider features sets perform better for drugs affecting general cellular mechanisms. Appropriate feature selection strategies facilitate the development of interpretable models that are indicative for therapy design. "The random forest machine learning algorithm achieved a correct Casilitate patients of 88.5% (area under the curve—AUC 0.94), However, none of the method patients of 48.5% (area under the curve—AUC 0.94), However, none of the method to the curve of t	
Koureas, M and Kirgou, P and Amoutsias, G and Hadjichristodoulou, C and Gourgoulianis, K and 283 Taiskalof, A	compounds in exhaled breath for lung cancer discrimination from other pulmonary diseases and healthy persons M	fetabolites 10	8	1-18	2020 Greece	http://dx.doi.org/10 3390/metabo1008 0317 article	map population sample consisted on 51 patients with confirmed LC, 38 patients with pathological computed tomography (CT) findings not diagnosed with LC, and 53 healthy controls*	Case-control study	AUC (10-fold CV + validation)	cross-validation + test set	used achieved adequate discrimination between LC patients and patients with	used achieved adequate discrimination between LC patients and patients with abnormal computed tomography (CT) findings. Biomarker sets, consisting mainly of the exogenous monoaromatic compounds and 1 and 2 propanol, adequately discriminated LC natients from healthy controls."
284 Lai, YH and Chen, W N and Hou, T C and Lin, C and Tiao, Y and Wu, S	Overall survival prediction of non- small cell lung cancer by integrating microarray and clinical data with deep learning Sc	ti Rep 10	1	4679-4679	2020 Taiwan	http://doi.org/10 _10389-41598-020- 61588-w article	"We separated 256 patients as the training set, 85 patients as the validation set, and 171 patients as the test set"	Case-control study	AUC, accuracy (10-fold CV + validation set)	training + test set	integrative DNN approach is capable of providing the missing information left by the other observed modality. Compared with our microarray DNN, we observed an increase in AUC and accuracy from the integrative DNN.*	integrative DNN approach is capable of providing the missing information left by the other observed modality. Compared with our microarray DNN, we observed an increase in AUC and accuracy from the integrative DNN." "In this large multi-centric prospective clinic-percentic inneiturinal dataset of breast
Lee, S. and Deary, I. O and Ch., I Ha and Si Meglin, A and Dumes, A. and Menveille, G. and Charles, C. are Seopult, S. and Rocosceau, M. and Resce, C. and Thomas, S. and Edand, A. and Citt, P. and Technic, O. and I evry, C. and Martin, A. I. and Everhand, S. and Ganz, P. A. and Partridge, A. H. and Michiels, S. and 285 Deleuze, J. F. and Andre, F. and Vaz-Luis, I.	treatment-induced fatigue by machine learning using genome-wide JN	VCI Cancer pectrum 4	5		2020 USA	http://dx.doi.org/10 1093UNCICSIPK AA000 article	"We accessed germline genome-wide data of 2799 early-stage breast cancer patients from the Cancer Toxicity study (NCT01993498)" "The study was conducted among 452 surgery-resectable patients with lung	Treatment response prediction	AUC (training + test set)	training + test set	the different known dimensions of Tatigoe. Although the ability of our models to identify-clinic angeometic contributors of Tatigoe although the ability of our models, a grup of SNPs and clinical variables was suggested to be associated with the cognitive domain.**	the different known dimensions of Tatigue. Although the ability of our models to identify clinic and genomic contributor of Tatigue differed by tatigue domain, a group of SNPs and clinical variables was suggested to be associated with the cognitive domain."
Li, B and Wang, C and Xu, J and Fang, S and Qiu, F and Su, J and Chu, H and Han-Zhang, H and Mao, 286 and Liu, H and Liu, X and Zhang, W and Zhao, H and Zhang, Z	CI Multiplatform analysis of early-stage Ci cancer signatures in blood Re	linical ancer esearch 26		11	2020 China	11591557- 3265 LigBiop20- A66 abstract	cancer (N=180), colorectal cancer (N=210), liver cancer (N=62), and 290 age_/sex- matched non-cancer controls*	Case-control study	AUC (training + test set)	training + test set	investigation in large-scale clinical studies is ongoing." "In conclusion, we proposed a boosting ensemble predictive framework with the wranner-hosel feature selection alentithe for noelliting antidenressant reatment.	r sequencing as a sensitive ctDNA profiling approach for early cancer detection. Further investigation in large-scale clinical studies is orgoing." "In conclusion, we proposed a boosting ensemble predictive framework with the wranner have feature solerities alengithm for profit into anti-informerscent treatment.
287 Lin, E and Kuo, P H and Liu, Y L and Yu, Y W Y and Yang, A C and Tisal, S J	Microsovicosmost characterization of	als 13 speriment		10 1-12	2020 USA	http://dx.doi.org/10 3390iph13100905 article	"we retained 421 MDD patients for the subsequent analysis" "We performed unsupervised dustering of total 1000 HDC lineatocellular	Case-control study	AUC (repeated 10-fold CV)	cross-validation	selection algorithm may leverage a feasible way to create predictive algorithms for forecasting antidepressant treatment response and remission with clinically meaningful accuracy." "Durands depondent and 2 immuno clusters with different feature. More	response and remission in Taluanese patients with MDD. The present results suggest that care boosting exemble predictive framework with the variety abused parties solicition algorithm may loverage a feasible way to create predictive algorithms for forecasting artisticipscent treatment ereponse and remission with clinically manningful accuracy."
Liu, F and Qin, L and Liao, Z and Song, J and Yuan, C and Liu, Y and Wang, Y and Xu, H and Zhang, Q 288 and Poi, Y and Zhang, H and Pan, Y and Chen, X and Zhang, Z and Zhang, W and Zhang, B	and multi-omics signatures related to H prognosis and immunotherapy an response of hepatocellular carcinoma Or	ematology nd ncology 9	1		2020 China	http://dx.doi.org/10 .1186940164-020- 00165-3 article	carcinoma] samples including discovery and validation group from available public datasets" "The discovery stage involved 160 pairs of ccRCC [clear-cell renal cell carcinoma]	Prognostic subtye stratification	correlation (discovery + validation cohorts)	external cohort validation	importantly, multi-omics signatures, such as MMPP was identified based on three clusters to help us recognize patients with different prognosis and responses to immunotherapy in HCC."	The production of the production of the production of the program of the production of the production of the production of the program of the
289 Uu, P. and Tain, W	Identification of DNA methylation patterns and biomarkers for clear-cell renal cell carrinma by multi-omics data analysis	perj 8			2020 China	http://dx.doi.org/10 _7717/poorj.9654 article	and matched normal tissues for investigation of DNAm and biomarkers as well as 318 cases of coRCC including clinical signatures." "This data describes 20,501 genes in 806 different breast cancer samples. We retained only samples with completin information. After that, 85 TMBC and 466 information. After that, 85 TMBC and 466 information.		AUC (10-fold CV + validation cohort)	cross-validation + external cohor validation	t epigenesis factors, which were the most important causes of advanced cancer and poor clinical properties. The logaum method is a powerful method for feature selection. However, it is unable to use any interactive information. To overcome this drawback, in this paper we first propose Legium-Net regularization to integrate which is the proposed legium which regularization to integrate the proposed legium which is the proposed legium which regularization to integrate the proposed legium which is the proposed legium which regularization to integrate the proposed legium which is the proposed	"The present study provides a comprehensive analysis of cRCE using multi-omics data. These findings is an about could select a snapsic could selectify some movel egigenetic factors, which were the most important causes of advanced cancer and poor critical propriors poor critical propriors." The logarim method is a power bright method for feature selection. However, it is unable to use any prescribe sollogical returner information. To overcome this drawback, in this paper well frest propose largours het engalarization to integrate biological returner. In the space we frest propose largours het engalarization to integrate biological returner. In the space we frest propose largours het engalarization to integrate biological returner. In the space we frest propose largours het engalarization to integrate biological returner.
280 Liu, X Y and Wu, S B and Zeng, W Q and Yuan, Z J and Xu, H B	Combining Genetic Mutation and CI	ti Rep 10 linical fedicine	1	22125- 22125	2020 China	http://doi.org/10 3233/hts-218008 article	non-TNBC were further divided into two groups: training (n= 327; 51 TNBC, 276 non-TNBC) and testing (n= 218; 34 TNBC 184 non-TNBC) sets" "272 samples of the TCGA LUAD cohort were selected according to the overall	Case-control study	AUC (10-fold CV + validation)	cross-validation + test set	regularization logistic regression model (logsum-NL) for gene selection and cancer classification. For a real large dataset, the proposed method has achieved 89.66% (training) and 9.0.02% (testing) AUC performance which are, on average, 5.17% [training] and 4.60% (training) hatter than exhibit the proposed that the property of the performance of the property of the performance of the per	regularization legistic regression model (logism=NLI) for gene selection and cancer resolutions; for an earling selectate, the proposed membed has achieved \$8.56% (training) and 90.02% (testing) ALC performance which are, on average, 5.17% (training) and 4.48% (besting) batter than amilitarizam methods. "The most significant contribution of this article is the integration of the genetic mutation and expression profiles to determine prognostic genes for ULAD patients. If
291 Liu, Y and Liu, F and Hu, X and He, J and Jiang, Y		rsights: ncology 14			2020 China	http://dx.dei.org/10 11779179554620 988290 article	survival and were partitioned into the training set and testing set [75%/25%]* "The TCGA training set contains 226 samples and the test set contains 226 samples. As an external validation set, the GSE17538 data set contains a total of	Prognostic study	AUC (10-fold CV + validation)	cross-validation + test set	needed on the point of the contract pulphone is grained as the property of the popular of determining predict expression and matricles profiles are available, the pipeline of determining DEGs and DMGs in this article can be applied to other types of cancers."	modelcrame separating orders orderes transplanting state of the profiles and available, the pipeline of determining DEGs and DMGs in this article can be applied to other types of cancers."
292 Lu, Y and Wu, S and Cul, C and Yu, M and Wang, S and Yus, Y and Liu, M and Sun, Z	gene signature for estimating sa	ncoTarget and herapy 13		10393- 10408	2020 China	http://dx.doi.org/10 2147/07T 82555 90 article	the 05±17556 details to lontains a total of 244 samples, including 6 mouse samples white among the 238 human samples, 38 samples recorded the survival status of NA, and finally used for follow-up analysis*	Prognostic study	AUC (training + text set + external validation set)	external cohort validation	for predicting the OS [overall survival] of COAD [Colon adenocarcinoma] patients" "We have confirmed a key prediction of the DESNT cancer model by demonstrating	"In this study, the AUC of 9 gene signature screened by multi-omics in the training set and validation set for the years is more than D.S., which is more effective in predicting the proposics of partners." "We have confirmed a key prediction of the DESMT cancer model by demonstrating that the presence of a multi-proportion of the DESMT cancer signature conferrs poorer
Luca, B A and Moulton, V and Ellis, C and Edwards, D R and Campbell, C and Cooper, R A and Clark, 293 and Brewer, D S and Cooper, C S	A novel stratification framework for predicting outcome in patients with prostate cancer Br	r J Cancer 122		10 1467-1476	United 2020 Kingdom	http://dx.doi.org/10 _1038941416-020- 0799-5 article	"There were 1785 samples from primary malignant tissue, and 173 from normal tissue" "a total of 901 TCGA NSCLC samples	Prognostic subtye stratification	Correlation, log-rank p-value (hold-out validation)	training + test set	outcome. The proportion of DESNT signature can be considered a continuous	outcome. The proportion of DESNT signature can be considered a continuous
294 Luo, R and Song, J and Xiao, X and Xie, Z and Zhao, Z and Zhang, W and Miao, S and Tang, Y and Ran		ging Albany NY) 12		14649- 14 14676	2020 China	http://dx.doi.org/10 18830/sping 1095 17 article	were available using the Illumina Infinium HumanMethylation450 platform, including 827 tumor tissues and 74 non-tumor tissues"	Tumor recurrence and immunotherapy response prediction	d AUC (training + external validation)	external cohort validation	lung squamous cell carcinomas (LUSC), a promising DNAm-based risk score model predictive of relapse was constructed and then validated in the other 3 datasets." "When assessing the performance in 30% of test samples after training on 70% of	NSCLC. Base on TCGA NSCLC cohort comprised of lung adenocarcinomas (LUAD) and lung squamous cell carcinomas (LUSC), a promising DNAm-based risk score model predictive of relapse was constructed and then validated in the other 3 datasets." "When assessing the performance in 30% of test samples after training on 70% of
Makarious, M and Iwaki, H and Blauwendraat, C and Leonard, H and Hachemi, S and Kim, J and Van Kuuren-lennen, K and Craig, D and Appelmans, E and Smolensky, L and Bookman, M and Singleton, 295 and Faghri, F and Nalis, M	Biomarker discovery in parkinson's disease using machine learning on A public multi-OMIC datasets: A pilot M study Di	flowement isorders 35		S207-S207	2020 USA	http://dx.doi.org/10 meeting 1002/mda 28288 abstract	"This included 872 samples that had sequenced genomes, clinical data, and "50K normalized transcripts from RNA sequencing"	Case-control study	AUC (30% test samples after training on 70% of samples)	training + test set	samples, multiple modalities implemented in the same predictive model performs best. By incorporating different modalities, we can develop more comprehensive predictive models to better understand the complex disease and identify better biomarkers." "Machine lacquier is a promision approach in the course of more accurate and	samples, multiple modalities implemented in the same predictive model performs best. By incorporating different modalities, we can develop more comprehensive predictive models to better understand the complex disease and identify better biomarkers."
Mantha, S and Dunbar, A and Bolton, K L and Devlin, S and Governhtsyn, D and Donoghue, M and 296 Arcila, M E and Soff, G A	Machine learning for prediction of cancer-associated venous thromboembolism Bi	lood 136		37-37	2020 USA	http://dx.doi.org/10 _1182blood-2000meeting _139579abstract	"12,040 patients were included in the final analysis. There were 855 CAT events during the observation period"	Prognostic study	C-index (cross-validation)	cross-validation	administration. Additional work is needed to identify the optimal algorithm and covariates, including better delineation of which cancer genomic information should be retained." "Multivariate models become a covered worker machines and the LASSO work his	presentables model for prediction of CAT (Cancer-Associated Venous and generalizable models) for prediction of CAT (Cancer-Associated Venous Thombsenbollsnift), the deplication described here, the use of random survival forests: performed without information stood future chamceborary administration. Additional work is needed to identify the optimal algorithm and convinctes, including better dislination of which excere generic information should be retained. **Multivariate models based on support vector machines and the LASSO variable
Mans), M and Palazzo, M and Knott, M E and Beauseroy, P and Yankilevich, P and Giménez, M I and 297 Monge, M E	Coupled Mass-Spectrometry-Based Lipidomics Machine Learning Approach for Early Detection of Clear 1F Cell Renal Cell Carcinoma Re	Proteome es 20	1	841-857	2021 Argentina	http://dx.doi.org/10 :1021facs.jproteo me.0c00963 article	"patients with clear cell renal cell carcinoma (ccRCC) stages I, II, III, and IV (n = 112) and controls (n = 52)" "In a derivation cohort of 636 patients	Case-control study	accuracy (training/test set)	training + test set	selection method yielded two discriminant lipid panels for cRCC detection and early diagnosis. A 16-lipid panel allowed discriminating cRCC patients from controls with 95.7% accuracy in a training set under cross-validation and 77.1% accuracy in an independent text set."	selection method yielded two discriminant lipid panels for cellEC detection and early disapposis. A 16-light panel allowed descriminating cellCD patients from controls with 95.7% accuracy in a training set under cross-validation and 77.1% accuracy in an independent test set."
McCothy, C P and Noumann, JT and Michalburgh, E A and Broken, N E and Goggn, H E and Stormann, N and Schalder, S and Zeller, T and Megaret, C A and Barnes, G and Brynn, B P and 288 Westermann, D and Januarii P., I L		Am Heart ssoc 9		e017221- 16 e017221	2020 USA	http://dx.doi.org/10 _1161/jehn.120.01 7221 article	in a pervation control to the patients referred for coronary angiography, predictors of 270% coronary stenosis were identified from 6 clinical variables and 109 biomarkiers. The final model was first internally validated on a separate cohort (n=275) and then externally validated on a cohort of 241 patients.	s Case-control study	AUC (train + test + external validation)	external cohort validation	the presence of obstructive CAD (Coronary Artery Disease) with high accuracy. The score performs similarly well in the evaluation of acute chest pain in the ED (including patients who had MI neither ruled in nor ruled out) and in outpatients presenting for	"We have derived and externally validated a clinical/proteomic panel that can predict the presence of obstructive CEO with high accuracy. The core performs centrally well in the evaluation of stocks close pain that ED (producing parients with hard his maken and the evaluation of stocks close pain in the ED (producing parients who had his maken with the evaluation of stocks are produced by presenting for evaluation of stocks angles including those with renal injury."

299 Milao, R and Chen, H H and Chong, Q and Xio, LY and Yang, ZY and He, M F and Hao, Z F and Librig, Y	Beyond the limitation of targeted therapy; Improve the application of targeted drugs combining genomic 4 data with machine learning	Pharmacol Res	159	104932- 104932	2020 China	http://dx.dei.org/10 _10186.chm.2020 108952 article	"The GDSC dataset contains 140 drug sensitivity experiments results in 624 cell lines" "Targeted DNA sequencing for more than 500 cancer-associated genes and expense-canhure RNA convenience was	Drug response prediction	AUC, sensitivity, specificity (5-fold CV)	cross-validation	methods combined with genomics data to accurately predict the performance of	"The proposed model of this paper used statistical methods and Machine learning methods combined with genomics data to accurately predict the performance of oncology drugs on cancer cell lines."
Michoda, J and Leibowitz, B and Amar-Farfash, S and Bevin, C and Brecch, A and Kopilivsky, J and Igartus, C and Bell, J S and Beauchamp, K A and White, K and Stumpe, M and Beauchier, N and Tarts 200 T	er, Multimodal prediction of diagnosis t cancers of unknown primary	or Cancer Research	80	16	2020 USA	http://dx.dei.org/10 1158/1658: 7-445 AM2020: meeting 5423 abstract	exome-capture RNA sequencing was carried out in more than 25,000 fresh frozen or paraffin embedded tumor samples, including both primary and metastatic tumors" "We divided 741 ADNI participants with blood microarray data into three groups based on their most recent CDR	Differential diagnosis prediction	accuracy (training/test set)	training + test set		and robustness of machine learning models to predict cancer diagnosis*
301 Miller, J B and Kauwe, J S K	Predicting Clinical Dementia Rating Using Blood RNA Levels	Genes (Basel)	11 6		2020 USA	http://doi.org/10 3300/inparent10000 700 article	assessment: cognitive normal (CDR = 0), mild cognitive impairment (CDR = 0.5),	Differential diagnosis prediction	AUC (10-fold CV)	cross-validation	machine was able to increase predictive accuracy of AD from a 55% baseline to	In individuals using RNA levels from a Blood microarray by taking into account small differences in engression that are individually noneignificant. A support vector machine was able to increase predictive accuracy of AD from a 55% baseline to almost 50%.*
Miller Addins, G and Accessés Morreno, L A and Grove, D and Dwells, R A and Tonelli, A R and Brown 302 M and Allendes, D S and Auceje, F and Rotrott, D M Mongain, D and Focking, M and Healy, C and Sussi, S R and Capney, G and Cannon, M and Zammit,	Secondary Liver Tumors Development of proteomic predictic models for outcomer in the clinical	Communica tions	4 7	1041-1055	2020 USA	http://dx.doi.org/10 1602/hap4.1499 article	= 54), cirrhodis (n = 30), HCC (n = 112), pulmonary hypertension (n = 49), or colorectal cancer liver metastases (n = 51)* "The sample comprised 133 CHR	Differential diagnosis prediction	balanced accuracy (cross-validation)	cross-validation	"The use of machine learning and breath VOCs (Volatile organic compounds) shows promise as an approach to develop intercret, nonlineasive screening tools for chronic liver disease and primary and secondary liver tumors."	
and Netico, B, and McCorn, P and Nocderott, M and Visios, M D and Richer Roster, A Do desire R and Emranter-Vicial, Y and Employee, S and Rulmanne, S and Sachs, G and Note Policy, A and Rutter, B and Partellis, C and De Haan, L and Valmaggio, L and Kempton, M and McGuire, P and 303 Cotter, D	de experiences in adolescence: Machin learning analyses in two nested casecontrol studies Convolutional neural network mode	Schizophren ia Bulletin	46	S238-S239	2020 Ireland	https://www.ncbi.nl m.nh.godoseciari dee/PMC7234235/ abstract	The sample compress 135 cm. participants who were followed clinically for up to 6 years, of whom 49 transitioned to psychosis and 84 did not. "The models were trained and texted on gene expression profiles from combined 10,340 samples of 33 cancer types and 10,340 samples of 33 cancer types and	Case-control study	AUC, PPV, NPV (training + test set)	training + test set	"With external validation, models incorporating proteomic data may contribute to improved prediction of clinical outcomes in individuals at risk of psychosis" "Taken together, we have presented there uniques DN architectures that take high dimension gene expression inputs and perform cancer type prediction while considering their torse out origin, Durmoid solvieved on equivalent 95.7% prediction	improved prediction of clinical outcomes in individuals at risk of psychosis" "Taken together, we have presented three unique CNN architectures that take high dimension seen expression inputs and perform cancer type prediction while
304 Mostavi, M and Chiu, Y C and Huang, Y and Chen, Y	for cancer type prediction based on gene expression Individualized Prediction of Respons to Methotresate Treatment in	BMC Med	13	44-44	2020 USA	11989-12020-020 0677-2 article		Differential diagnosis prediction	accuracy (6x 5-fold CV, 80-20% splitting for training and validation)	cross-validation	accuracy comparing to earlier published studies, however with a drastically simplified CNN construction and with a reduced influence of the tissue origin."	accuracy comparing to earlier published studies, however with a drastically simplified CNN construction and with a reduced influence of the tissue origin.* *Pharmacogenomic biomarkers including gene variants for cancer susceptibility genes
305 Myssoedova, E and Athreya, A and Crowson, C and Weinshilboum, R and Wang, L and Matteson, E	Patients with Rheumatoid Arthritis: Pharmacogenomics-driven Machine	A Arthritis and Rheumatolo 87 Journal of Diabetes	72	4014-4015	2020 USA	http://ktx.doi.org/10 _1602/art.41538 abstract	positive; mean age 54 years; mean baseline Disease Activity Score with 28- joint count [DAS28] 5.65)*	Drug response prediction	AUC (Sx 10-fold CV + external validation)	cross-validation + external cohor validation	(CASC15) and important MTX pathway enzymes (ATIC) combined with baseline DAS28 t score predicted MTX response in patients with early RA more reliably than	(CASCLS) and important MTX pathway enzymes (ATIC) combined with baseline DAS28 score predicted MTX response in patients with early RA more reliably than demographics and baseline DAS28 alone, with replication in an independent cohort*
306 Naz, H and Alnija, S Natha, A and Seleres, M A and Bejar, B and Bash, M I and Other, M and Komrollij, R S and Barnardon and Hitton, C B and Kerr, C M and Steeners, D P and DeZern, A and Roboz, G and Garcia-Manero, or	Deep learning approach for diabete prediction using PIMA Indian datase Genomic biomarkers to predict d, J resistance to hypomethylating agen G in patients with myelodysplastic syndromes using artificial intelligence	Disorders s JCO Precision	19 1	391-403	2020 India 2019 USA	http://dx.doi.org/10 .1007s/40200-020- 00520-5 article http://dx.doi.org/10 .1200PO.19.0011	"a total of 768 instances, from which 268 samples were identified as diabetic and 500 were non-diabetics" "Among 433 patients, 193 (45%) received azacitifine, 176 (40%) received decitabine, and 64 (15%) received HMA alone or in combination"	Case-control study Drug response prediction	accuracy ("splits in an 80/20% ratio into the training and validation set") accuracy Itrainine/hest set1	training + test set training + test set	promising extracted features. Cl. achieves the accuracy of 88.07%which can be used for further development of the automatic prognosis tool." "Genomic biomarkers can identify, with high accuracy, approximately one third of patients with MDS who will not reapon the HMAE. This study highlights the importance of machine learning technologies such as the recommender system absorbtion in translation second-citat for useful clinical formation of the useful clinical control can be according to the control	promising entracted features. DL achieves the accuracy of \$80.07/which can be used for further development promising entracted features. DL achieves the accuracy of \$80.07/which can be used for further developments prognosis tool." "Genomic biomarkers can identify, with high accuracy, approximately prevent bird of patients with MDS who will find respond to HMAD. This study highlights the importance of machine learning technologies such as the recommender system also that the provided of t
Nielsen, R. Land Hellenius, M. and Garcis, S. Land Roager, H. M. and Aytan-Akbag, D. and Hansen, E. B. and Lind, M. Y. and Vogs, J. K. and Disgazed, M. D. and Balt, M. I. and somer, C. B. and Mattapowel, R. A. Warinner, C. A. and Askov, V. and Gold, R. and Kintensen, M. and Friskar, H. and Sontho, M. H. and Christmesen, A. F. and Vestergazed, H. and Hansen, T. and Kintansen, K. and Brist, S. and Petersen, T. N. 30a and Lauritee, L. and Licker, R. and Petersen, T. N. Sonthology, L. and Licker, R. and Advesser, C. and Citter, L. and Licker, R. and Advesser, C. and Citter, L. and Licker, R. and Briston, L. and Licker, R. and Sonthology, L. and Licker, R. and Briston, L. and Licker, R. and Sonthology, L. and Licker, R. and Licker, R. and Licker, R. and Licker, R. and L. and Licker, R. and Licker,	nd Data integration for prediction of	Sci Rep	10 1	20103- 20103	2020 Denmark	http://doi.org/10 .1038/41598-020- 76097-2 article	"Here, we classify weight loss responders (N = 106) and non- responders (N = 97) of overweight non- diabetic middle-aged Danes to two earlier reported dietary trials over 8 weeks"	Treatment response prediction	AUC ("50 shuffle-split fivefold cross-validations was used")	cross-validation	eventually in concert with comprehensive population weight loss strategies. Furthermore, understanding predictive features of weight loss response will drive improved understanding of the interplay between gut microbiota, diet and individual	"By identifying the propentity of study participants likely to experience weight loss, a more effective individual targeting of delay interventions can be Entilizated, eventually in concern with comprehensive population weight loss strategies, eventually in convert with comprehensive population weight loss strategies. Furthermore, understanding predictive features of weight loss repease will drive improved understanding of the interplay between gut microbiota, diet and individual predictipants."
309 Nyamundanda, G and Escon, K and Gulnney, J and Lord, C J and Sadarandam, A	models single omics and phenome data for functional subtyping and personalized cancer medicine	Cancers	12	10 1-14	United 2020 Kingdom	http://lite.doi.org/10 3300/camcars/121 02813 article	"in total, 2043 breast cancer samples were used in this work."	Subgroup stratification	in Silhouette width, cophenetic correlation (external text diatasets)	external cohort validation	"This expert review describes and examines, first, the SVM models employed to forecast breast cancer subtypes using diverse systems science data, including transcriptomics, epigenetics, proteomics, and radiomics, as well as biological	clinical translational potential." "This export review describes and examines, first, the SVM models employed to forecast breast cancer subtypes using diverse systems science data, including transcriptomics, epigenetics, proteomics, and radiomics, as well as biological
310 Ozer, M E and Sarica, P O and Arga, K Y	New Machine Learning Applications Accelerate Personalized Medicine in Breast Cancer: Rise of the Support Vector Machines notibe: Software for huilding	Omics	24 5	241-246	2020 Turkey	http://dx.doi.org/10 1009/cmi 2020.00 01 article	review (not applicable)	Review			performance of the present SVM and other diagnostic and therapeutic prediction models across the data types. We conclude by emphasizing that data integration is a critical bottleneck in systems science, cancer research and development, and health care innovation and that SVM and machine learning approaches offer new solutions	critical bottleneck in systems science, cancer research and development, and health
Pai, S and Weber, P and Isserlin, R and Kaka, H and Hui, S and Shah, M A and Giudice, L and Giugno 311 and Nehr, A K and Baumbach, J and Bader, G D	patient similarity networks binomialRF: interpretable	F1000Res	9	1239-1239	2020 Canada	http://dx.doi.org/10 12698/1000resea csh.26429.2 article	"including 154 Luminal A and 194 tumours of other subtypes" "we conduct a variety of simulations and trials against the Madelon benchmark	Case-control study	AUROC, AUPR, and accuracy (an approximately 70:30 split of samples was used for cross validation)	cross-validation	patient classification from sparse genetic data" "binomialRF extends upon previous methods for identifying interpretable features in	"the netDx Bioconductor package provides a novel workflow for pathway-based patient classification from sparse genetic data" "binomial# extends upon previous methods for identifying interpretable features in RFS and brings them together under a correlated binomial distribution to create an
Rachid Zaim, S and Kenost, C and Berghout, J and Chiu, W and Wilson, L and Zhang, H H and Lussier 312 A	combinatoric efficiency of random r, Y forests to identify biomarker interactions	BMC Bioinformati cs	21 1	374-374	2020 USA	https://doi.org/10.1 188/s12850-020- 03718-9 article	dataset from the University of California – I vine (UCI), and clinical datasets from The Cancer Genome Attas (TGGA) "Clinical and genomic data, including commercially available next-generation sequencing panels, were obtained for	Case-control study	Precision, recall, test error (training and test set)	training + test set	while retaining competitive model selection and classification accuracies." "We developed and externally validated a highly accurate and interpretable model	officiant hypothesis regional mobile "united collaboration consolination of certain and efficiant hypothesis estimal agenthe that identifies bloomarkers" main effects and interactions. Preliminary results in simulations demonstrate computational gains while retaining competitive model selection and castifician accurations." "We developed and externally validated a highly accurate and interpretable model election."
Radsholdt, N. and Maggendorfer. M and Maldounk, I and Salvers, M. A and Stever, J. and Bass. Histor, Cand Boughalk, Yand Mistery, Sand Huster, Sand Madharlers, Sand Ret, Chand Rha, R. R. And Gall, A and Pozzi, S and Gerds, A T and Haferlach, C and Maciejsevski, J P and Haferlach, T and 313 Nazha, A.	A personalized clinical-decision tool improve the diagnostic accuracy of myelodysplastic syndromes	Blood	136	33-35	2020 USA	http://like.doi.org/10 11878bisos5-2000- 159412 abstract	patients (pts) treated at the Cleveland Clinic (CC; 652 pts), Munich Leukemia Laboratory (MLL; 1509 pts), and the University of Pavia in Italy (UP, S36 pts)*	Case-control study	AUC (training + external validation)	external cohort validation	personalized interpretations of its outcome and can aid physicians and hematopathologists in recognising MDS with high accuracy when encountering pts with pancytopenia and with a suspected diagnosis of MDS." "In this article, we compare the usefulness and limitations of traditional statistical methods and ML, when applied to the medical field." Traditional statistical methods	that can distinguish MOS from other mystoid malignancies using clinical and mustional data from a large international colonic. The modifican provide personalisated international colonic management and and physicians and hematopathologists in recognising MOS with high accuracy when recountering pits with panet popera in with a suspected disposic of MOS.* "In this article, we compare the usefulness and limitations of traditional statistical methods and MJ, when applied to the medical field. Traditional assistical methods
314 Rajula, H S R and Verlato, G and Manchia, M and Antonucci, N and Fanos, V	Comparison of Conventional Statistical Methods with Machine Learning in Medicine: Diagnosis, Dru Development, and Treatment	g Medicina- Lithuania	56 9		2020 Italy	http://dx.doi.org/10 3100/marticing/90 90455 article	review (not applicable)	Review			such as in public health. ML could be more suited in highly innovative fields with a huge bulk of data, such as omics, radiodiagnostics, drug development, and	seem to be more useful when the number of cases largely excessed the number of variables under study and a priori inconselegation on the topic under targely is underartial such as in public health. Mt. could be more suited in highly innovative fields with a huge built of data, such as omic, includial generatic, dring development, and personalised treatment. Integration of the use oppreciates should be preferred over a unidirectional choice of either approach." We have developed DMR methylation score for exposure to maternal smoking
Rauschert, S and Melten, PE and Heistala, A and Karhunen, V and Burdge, G and Craig. I M and Godfrey, CM and Litycrop, K-and Mort, T A and Bellin, L1 and Oddy, W H and Pennell, C and Järvelli 315 M R and Sebert, S and H	Machine Learning-Based DNA Methylation Score for Fetal Exposur to Maternal Smoking: Development in, and Validation in Samples Collected from Adolescents and Adults Proteomics and Metabolomics	Environ Health	128 9	97003- 97003	2020 Australia	http://dx.doi.org/10 1209/sheb076 article	"The score was developed and tested in the Raine Study with data from 995 white 17-y-old participants using 10-fold cross-validation"	Case-control study	Sensitivity, specificity (10-fold CV)	cross-validation	during pregnancy, outperforming the three previously developed scores. One possible application of the current rore could be for model adjustment purposes or to assess its association with distal health outcomes where part of the effect can be attributed to maternal smoking. Father, it may provide a biomarker for fetal exposure to maternal smoking."	during pregnancy, outperforming the three previously developed scores. One possible application of the current score could be for model adjustment purposes or to assess its association with distal health outcomes where part of the effect can be attributed to material minoling. Further, it may provide a biomarker for fetal exposure to material smoking."
Ristori, M V and Mortera, S L and Marzano, V and Guerrera, S and Vernocchi, P and laniro, G and 316 Gardini, S and Torre, G and Valeri, G and Vicari, S and Gasbarreii, A and Putignani, L	Approaches towards a Functional Insight onto AUTISM Spectrum Disorders: Phenotype Stratification and Biomarker Discovery	Int J Mol Sci	21	17	2020 Italy	http://dx.doi.org/10 33000jess2117627 4 article	review (not applicable)	Review			ASD and healthy controls, without considering the family and specific characteristics of the pathology, Orther, the sample control sale highly limited. From the point of view of omics data, the biggest limit is that all of the data from the omics are not considered and the data are not integrated with collected clinical data." In this research, we compared three machine learning methods that have been	ASD and healthy controls, without considering the family and specific characteristics of the pathology. Often, the sample cohers it also highly himided. From the point of view of omiss data, the biggest limit is that all of the data from the omiss are not considered and the data are not integrated with collected clinical data." "In this research, we compared three machine learning methods that have been
317 Romero-Rosales, B L and Tamez-Pena, J G and Nicolni, H and Moreno-Trevillo, M G and Trevino, V	Improving predictive models for Alzheimer's disease using GWAS dat by incorporating misclassified sample modeling	25	15 4	e0232103- e0232103	2020 Mexico	http://dx.doi.org/10 137 fijournal.pone 0.032103 article https://www.abstracts .org/abstracturecy waring-goost- lammatkers-for-	"The study has four groups accounting for 5,220 individuals"	Case-control study	AUC (D0 rounds of internal cross-validation (CV) to 80% of the dataset for training and 20% for testing)	cross-validation	wise) and propose the inclusion of markers from misclassified samples to improve overall prediction accuracy. Our results show that the addition of markers from an initial model plus the markers of the model fitted to misclassified samples improves the area under the receiving operative curve by around 5%, reaching "0.84, which is	wise) and propose the inclusion of markers from misclassified samples to improve overall prediction accuracy. Our results show that the addition of markers from an initial model offus the markers of the model fitted to misclassified samples improves
318 Rychlor, D and Neely, J and Sirote, M	Uncovering Novel Biomarkers for Rheumatoid Arthritis from Feature Selection and Machine Learning Approaches on Synovium and Blood Gene Expression Data	Arthritis and Rheumatolo 8Y	72	1503-1504	2020 USA	theoreachies arthritis-form feature-selection and-machine. learning-septrature-and-blood-general meeting blood-general abstract	"The raw data from 13 synovium datasets with 284 samples and 14 blood datasets with 1,885 samples were downloaded and processed"	Case-control study	AUC (training + test set)	training + test set	on public data and validated using multiple independent data sets, coupled with the RAScore may be useful in the early diagnosis and disease and treatment monitoring	"This novel list of biomarkers, identified through a robust feature selection procedure on paids data and unidated using multiple selectands and extra coupled with the complex couples and disease and treatment monitoring of BA."

319 Saorin, A and Di Gregorio, E and Molo, G and Steffen, A and Corona, G	Emerging role of metabolomics in overlan cancer diagnosts	Metabolites 10		10 1-15	2020 Italy	http://dx.deci.org/10 3/90/(matabo 1010 0419 article	review (not applicable)	Review			The most promising circulating signatures of OC (Doralin Cancer) involve metabolites belonging to Spot and A pathways. These metabolite freegrenset find agreement in many facilities, unasing term increases to the Canagosis. However, the contrast apparent of the production apparent to be infinitel because a laid of independent, tage related to the contrast of the c	The most promising circulating algorithms of OC (Dozalna Canors) involve metabolities belonging to lipids and Algorithmys. These metabolities fregerprints find agreement in many facilities, unsafet given investion for OC (algorities, flower), their agreement in many facilities, unsafet given investion for OC (algorities, flower), their validation suddeep promisit bein effective and for OC covering and monitoring. Pathware research should incide better desirguids and legal pathware programment of the pathware research should intellect the state of the pathware research should be proper external volidation in order to them the improve the transferonal research embeddeement of Code Segment values of the pathware should be proper production. In order to repair, continued on the Code Segment values of the pathware should be present production, we can prove that in instead of
320 Schauck, D and Brenner, T and Weigand, M and Uhle, F 321 Schenberg, A V and Boichard, A and Tolgelny, IF and Richard, S B and Kurznock, R	Deep-learning neural networks for accurate diagnosis of sepsis using microarray gene expression data Machine learning model to predict oncologic outcomes for drugs in randomized clinical trials	Intensive Care Medicine Experiment al 7	9	2537-2549	2019 Germany 2020 USA	http://dx.doi.org/10 198640695-019_meeting 0205-yabstract	"septic patients (n=1,354), trauma patients (n=478), and healthy controls (n=38)" A total of 467 progression-free surviva (PFS) and 369 overall survival (OS) data points were used as training sets to buil our ML (random forest) model:	Case-control study I If Treatment response prediction	AUC, sensitivity, specificity (training, validation and test set)	training + test set	learning idiocynoratic features tailored to specific data series, generalized strategies for sample discrimination have developed in the trained artificial neural networks. The combination of artificial neural networks and microarray gene expression data is therefore capable of achieving expsis diagnosis with superior accuracy and thus summont the "rurent filazonchic crone".	learning idiosyncratic features tailored to specific data series, generalized strategies for sample discrimination have devolved in the trained artificial neural networks. The combination of artificial neural networks and microarray gene expression data is therefore capable of artificial neural networks and microarray gene expression data is therefore capable of artificial neural networks and microarray gene expression data is augments the current diagnostic scope."
Senturk, N and Tuncel, G and Kooseglu, S and Dogan, B and Sag, S O and Mocan, G and Temel, S G 322 and Dundar. M and Emoren. M C	Developing evidence based	Gazi Medical Journal 31		P44-P44	2020 Turkey	https://www.ambas g.com/search/sesul to?asbacton-view pcond8id=18278 51078from-aspert article	"268 different BRCA1/2 positive breast cancer patients"			training + test set	*Overall, our developed models will provide the early prediction for BRCA1/BRCA2 related breast cancer cases and will improve to be beneficial for preventive models analyses proposed for today's needs have deepen statement of the proposed of the province of the provin	"Overall, our developed models will provide the early prediction for BRCA1/BRCA2 related breast cancer cases and will improve to be beneficial for preventive medicine and a union around fee tredity among the supplied production."
Singh, M and Singh, S P and Dubey, P K and Rachana, R and Mani, S and Yadav, D and Agarwal, M a 323 Agarwal, S and Agarwal, V and Kaur, H	Advent of Proteomic Tools for	Curr Protein Pept Sci 21	,	10 965-977	2020 India	http://dx.doi.org/10 2174/1389203721 889200515173213 article	review (not applicable) "The pathological stages are known for	Review	source (usuing con en)	THE STATE OF THE S	"the molecular diagnosis of AD incorporates various sophisticated techniques including immuno-sensing, machine learning, nano conjugation-based detections, etc. In the current review description, we have summarized the various diagnostic	"the molecular diagnosis of All incorporates various synthicizand scheinques including immuno-moning, machine learning, nanc consignation based detections, etc. In the current review description, we have summarized the various diagnostic approaches and their relevance is mitigating the long-standing urgency of targeted diagnostic tools for detection of AD"
	Integrative analysis of DNA methylation and gene expression in	Mol Genet Genomics 295				http://dx.dei.org/10 1607/s00438-020-	250 samples (common across both the platforms) with the following distributions: Stage II—167, Stage III—15 Stage III—50, and Stage IV—14. We divided the dataset containing these 25 samples into training (80%) and test (20%) dataset:"		PR.AUC, MCC, Accuracy, Sensitivity and Specificity ("The performance of the		expression to characterize the patterns of DNA methylation in PRCC. Our analysis	"In this study, we performed an integrative analysis of DNA methylation and gene expression to characterize the patterns of DNA methylation in RRCC. Our analysis showed that most policies are hypermethylated in RRCC, and both piper- and hypo-
324 Singh, N P and Vinod, P K Song, X and Yang, X and Narayanan, R and Shankar, V and Ethiraj, S and Wang, X and Duan, N and N	papillary renal cell carcinoma Oral squamous cell carcinoma Ni, diagnosed from saliva metabolic	Proc Natl Acad Sci U S	3	807-824 16167-	2020 India	01804-y article http://dx.doi.org/10 .1073(pnas.20013	"Saliva samples from 373 volunteers, 124 who are healthy, 124 who have premalignant lesions, and 125 who are		models was evaluated on the 20% test dataset") is Accuracy ("20-fold cross-validation was carried out", external validation	training + test set	methylated probes can distinguish normal from cancer samples. The salivary metabolic profile can reflect oral cancer development. Most discovered metabolites in callaw were found to be highly linked to their expression levels within the primary oncological site of oral cavity tissues, demonstrating the potential of	methylated probes can distinguish normal from cancer samples." "The salivary metabolic profile can reflect oral cancer development. Most discovered metabolites in saliva were found to be highly linked to their expression levels within the primary oncological site of oral cavity tissues, demonstrating the potential of
325 Y H and Hu, Q and Zun, R. N. Sesin, A S and Watson, D and Anir, P R and Basus, K and Ullal, Y S and Ghosh, A and Narvekar, Y and Grovey. H and Shabu, D and Praksah, A and Bahura, L and Balakrishnan, V and Ghosh Roy, K and Balgoppalan, S and Alam, A and Parashur, R and Mundleur, N and Christe, J and Macpherson, M D : 36 Kapoor, S and Marcusci, G	profiling Superior therapy response predictic for patients withmyelodysplastic and syndrome (MDS) using cellworks Singula™: Mycare-020-02	A 117 ins Blood 136		9-10	2020 China 2020 USA	95117 article bto://dx.doi.org/10 118284cod-2020- 142214 meeting abstract	OSCC patients" "The performance of Singula" was evaluated in an independent, randomly selected, retrospective cohort of 144 MDS patients."	Therapy response prediction	samples) Accuracy + confidence intervals (training + external cohort)	cross-validation + test set	saiva tor in with molecular diagnosis of USULT. "Collever's Signal" has high accuracy and sensitivity in predicting CR [complete response] for MDS [Myelodoysplasts Syndrome] patient response to physician prescribed therepies. Singula" also has high specificity in identifying patients who are unlikely for respond to physician prescribed therepies and provides alternative extension of the provides alternative source of the provides and provides alternative source of the provides alternative source of the provides and provides alternative source of the provides and provides alternative source of the provides and provides alternative source of the provides alterna	canva for in witro monecular diagnosis of USCL ^{**} "Cellworks Singini" has high accuracy and sensitivity in predicting CR [complete response] for MIDS [Myelodyplastic Syndrome] patient response to physician prescribed therapies. Singiuli" also has high specificity in identifying patients who are unifiliarly to respond to physician prescribed therapies and provides alternative toxicisms of the provides altern
Tabares-Soto, R and Chrosco-Ariss, S and Romero-Cano, V and Buchell, V S and Rodegues-Sotolo, J I 327 and Immero-Yorns, C F	microarray gene expression data PrOTYPE (Predictor of high grade)	ns Peerj Computer Science			2020 Colombia	http://dx.doi.org/10 2717/rearjos.270 article	"This database consists of 174 samples with 12,533 gene expression microarra for 11 different types of cancer."	ys Case-control study	accuracy, confusion matrix (10-fold CV)	cross-validation	"In this work, we show the application of unsupervised and supervised learning approaches of Ma. and IO for the classification of 11 cancer types based on a microarray dataset. We observed that the best average results using the training and validation data are obtained using the raw dataset and the IA algorithm, yielding an accuracy value of 100% (validation set, using the hold out splitting method). [] Additional tests with independent data should be done to discard potential overfitting."	"In this work, we show the application of sunspervised and supervised learning approaches of Ma and Life for destillations of III center types based on a microarray disease. We observed that the best average records using the straining and volations data are obtained using the rese destate and the III application, yielding and solutions data are obtained using the reset destate and the III application, yielding and accuracy value of 100% (volations set, using the hold-out splitting method), [] Additional tests with independent data should be done to discard potential overfitting."
Talhouls, A and George, J and Wang, C and Goode, E and Ramus, S and Doherty, J and Bowtell, D at 128 Anglesis, M	serous Ovarian carcinoma molecula subTIYEJ: The development and validation of a clinical-grade consensus classifier for the molecul dd subtypes of high-grade serous tubo ovarian cancer	r ar Clinical - Cancer Research 26		13	2020 Canada	http://dx.doi.org/10 .1159/1557- .3295.OVCA19- A603 abstract	Adopting two independent approaches we derived and internally validated algorithms for molecular subtype prediction from gene-expression array data in 1,650 tumors.	Differential diagnosis	is accuracy (training + external validation)	external cohort validation	"We validated the Predictor of high-grade-serous Ouation carcinoma molecular subTYPE, or POTYPE, following the Institute of Medicine guidelines for the development of omis-based tests. This simple-base, cost-effective, faily defined, and locked-down finical-grade susay with finitiate molecular subtypes straffication into clinical trial design." "Overall, our analysis downs that utilizing gene expression profiles independent of	"We validated the Predictor of high-grade-serous Ovarian carcinoma molecular subTVPE, or POTTVE, following the institute of Medictine guidelines for the development of omics based less: This simple-to-use, cost-effective, fully defined, and locked-down-filler-gaine saxy will festite medicular subspire stratification into clinical trial design."
329 Taleur, J and Carter, H	Assessing cancer drug responseprediction from gene expression	Cancer Research 80		16	2020 Canada	http://dx.doi.org/10 1158/1538- 7-445_AM2020- 2009 abstract	"Gene expression data for 17,737 gene across 1014 human cancer cell-lines win IC50 concentrations for 251 anti-cance drugs were obtained from the Genomic of Drug Sensitivity in Cancer Project" "A total of four LIAD expression profile (GSE32036, GSE32857, GSE33532, and GSE579037) were retrieved from GEO.	s th r ss Drug response prediction	accuracy (train/fest set)	training + test set	other-omics data for cancer drug response prediction through machine learning frameworks offers modest predictive coghistillice. To increase performance, we suggest augmenting training cise through shared pathway cross-training, optimizing feature encoding to maximize neural network predictive capabilities, and incorporating other-omics data."	other omics data for cancer drug response prediction through machine learning frameworks offers models predictive capitalism. To increase performance, we suggest augmenting training size through shared pathway cross-training, optimizing feature encoding to maximize neural setwork predictive capabilities, and incorporating other -omics data."
330 Tang, B and Wang, Y and Chen, Y and Li, M and Tao, Y	A Novel Early-Stage Lung Adenocardinoma Prognostic Model Based on Feature Selection With Orthogonal Regression	Frontiers in Cell and Developme ntal Biology 8			2020 China	http://dx.doi.org/10 3389/feet 2020 62 0746 article	And the corresponding accession number, platform, and sample information are listed in Table 1. Based on survival information from a total of 479 LUAD samples, the risk score was stratified into high and low groups."	Case-control study	AUC (cross-validation + external validation)	cross-validation + external cohor validation	"In conclusion, the proposed FSOR [feature selection with orthogonal regression] method can deliver better prediction performance for the early-stage prognosis and that the potential is uniprove therapy strategy, but with few predictor consideration and computation burden."	"In conclusion, the proposed FSDR [feature selection with orthogonal regression] method can deliver better prediction performance for the early-stage prognosis and has the potential to improve therapy strategy, but with few predictor consideration and computation business.
331 Tang, W and Cao, Y and Ma, X	Novel prognostic prediction model constructed through machine learn on the basis of methylation-driven genes in kidney renal clear cell carcinoma	Biosci Rep 40	7		2020 China	http://dx.dei.org/10 _1042bar/200160 4 article	"normal samples = 160, tumor samples 325" "To identify the prognostic impact of tumor hypoxia on increased risk for loc	Prognostic study	AUC (10-fold CV, test set)	cross-validation + test set	prognostic prediction model and combined with clinical information to build the tran	"We used the machine learning method to establish a multivariate methylation persponstic prediction model and combined with clinical information to build the trans- omics proposets remogram. [] These results can help in the accurate evaluation of the prognosis of RND posters and provious new cluss and data resources for the further study of the pathogenesis and the development of the disease."
	Hypoxia Methylome Classifier pktM Outperforms Gene Signatures in Identifying PPV-Negative HHSCC Patients at Risk for Locoregional	Internationa					regional recurrence (JR) and all event progression in Human Papilloma Vives DNA negative (HPV-) HMSCC. NAN negative (HPV-) HMSCC. Recurrence (HPV-					
Tawk, B and Widner, U and Schwager, C and Herpoll, E and Tehofor, I and Budach, V and Krauce, And Stuchles, M and Balempas, P and Roedel, C and Grous, A and Zeo, D and Combs, S E and 332 Weichert, W and Bellia, C and Baumann, M and Herold-Mende, C and Debus, J and Abdollahi, A decided to the control of	Failure Post Primary Radiochemotherapy: A German Cancer Consortium Radiation Oncology Group (DKTK-ROG) Multicenter Trial	I Journal of Radiation Oncology Biology Physics 108	3	e552-e553	2020 Germany	http://dx.doi.org/10 1016/i.ioobo.2020 07.1715 article	GESCH (consensus, n = 96) vs those wit intermediate-to low GESOH (n = 146). HDMC was validated in the DKTK-RDG primary cohort." "We performed untargeted most bibliopies on common complex	h Differential diagnosi	is ion correlation (training + validation cohort)	external cohort validation	"A methylation-based classifier of tumor hypoxia is successfully developed and validated to be prognestic for IR [loco-regional recurrence], progression and OS [overall survival] in HPV-HRSCC patients treated with primary RCHT."	*A methylation-based classifier of tumor hypoxia is accessfully developed and validated to be prognostic for IR [loco-regional recurrence], progression and OS [overall sunvival] in HPV-HNSCC patients treated with primary RCHT.*
Tiedt, S and Brandmaier, S and Düring, M and Artati, A and Klein, M and Liebig, T and Holdt, L and 333 Teupser, D and Wang-Sattler, R and Schwedhelm, E and Gieger, C and Dichgans, M	Circulating metabolites differentiate acute ischemic stroke from stroke mimics	Stroke 15	1	77-78	2020 Germany	http://dx.doi.org/10 .117791747493000 meeting 963387 abstract	obtained from patients with ischemic stroke (N = 219) and stroke mimics (N = 138; as defined by absence of a DV positive lesion on MRI)*	VI Case-control study	AUC (training + test set)	training + test set	PORT No. 1 Landson A 400.00	with ischemic stroke (N = 219) and stroke mimics (N = 138; as defined by absence of a
334 Tran, A and Walds, C J and Batt, J and Dos Santos, C C and Hu, P	A machine learning-based clinical to for diagnosing myopathy using mult cohort microarray expression profile	i- J Transl	1	454-454	2020 Canada	http://dx.doi.org/10 _1186/s12967-020- 02630-3 article	"Muscle tissue samples originating fron 1260 patients with muscle weakness."	Differential diagnosi prediction	is AUC (training + test set)	training + test set	clinical tool for muscle disease subtype diagnosis." "Min developed and unlighted an unbiased automated pipeline for transcriptomic	Divey positive lescent in new jordnerning molecular classification tool with the selected "Our results present a well-portnering molecular classification tool with the selected processor of the selection of t
Tran, P. M. H. and Tran, I. K. H. and Nechtman, J. and Doc Santos, B. and Purohit, S. and Satter, K. B. and 335. Dun, B. and Kolhe, R. and Sharma, S. and Bollag, R. and She, J. X.	Comparative analysis of transcriptomic profile, histology, an IDH mutation for classification of gliomas Epithelial-to-mesenchymal transitio is a prognostic marker for patient	d Sci Rep 10 n	1	20651- 20651	2020 USA	http://dx.doi.org/10 1038941598-020- 77777-8 article	"RNAseq and microarray data were obtained for 1032 gliomas from the TCGA and 395 gliomas from REMBRAND" "Pretreatment tumor material from patients of two cohorts, totalling 174	Differential diagnosi and survival predicti	is ion accuracy, log rank test p-value (cross-validation)	cross-validation	subtypes from histology and mutation status. Our analytical pipeline avoids the potential of overhiting a supervised model to michasofiald or michandial samples [] and can be used in establishing gold standard datasets devoid of erroneous and questionable samples for the development of automated tumor classifiers*	dustring, Without any domain knowledge, our discular recapiturated known gloma subpripes from historing and mustion status. Our analytical polipsia avoids the potential of lowerfitting a supervised model to microsciatified or michandided camples [] and can be used in establishing gold standard distances devoid derroneous and questionable samples for the development of automated tumor classifiers."
van der Heijden, M and Essers, P B M and Verhagen, C V M and Willems, S M and Sanders, J and de Roest, R H and Vossen, D M and Leemans, C R and Verheij, M and Brakenholf, R H and van den 336 Brekel, M W M and Vens, C	outcome in advanced stage HNSCC patients treated with chemoradiotherapy Identification of a Transcriptomic	Radiother Oncol 147		186-194	Netherlan 2020 s	d http://dx.doi.org/10 d 10166.redonc.202 0.05.013 article	cisplatin-based chemoradiotherapy treated HPV-negative HNSCC patients, was RNA-sequenced*	Prognostic study	AUC (cross-validation + external validation)	cross-validation + external cohor validation	"EMT in HPV-negative HNSCC co-defines patient outcome after chemoradiotherapy. The generated HNSCC-EMT prediction models can function as strong prognostic biomarkers." "Determining which treatment to provide to men with prostate cancer (PCa) is a	The generated HNSCC-EMT prediction models can function as strong prognostic biomarkers: " Determining which treatment to provide to men with prostate cancer (PCa) is a
Vitzare, B and Leclerco, M and Martin-Magniette, M L and Collins, C and Bergeron, A and Fradet, Y 337 and Droit, A	Prognostic Signature by Machine	all Frontiers in Genetics 11			2020 Canada	http://dx.doi.org/10 .3399/sene.2020. 550804 article	"Gene expression data were extracted from three RNA-Seg datasets cumulating a total of 171 PCa patients"	Prognostic study	balanced error rate ("The resampling strategy was run 200 times with a split of 2/3 for training and 1/3 for text sets")	stratified resampling training and test sets	major challenge for clinicians. [] This study demonstrates the feasibility to regroup different small datasets in one larger to identify a predictive genomic signature that would benefit PCa patients."	

352	Zhuang, H and Chen, Y and Sheng, X and Hong, L and Gao, R and Zh, X	leukemia through a combined multi-	PeerJ	8 e9437		2020 China	http://dx.doi.org/10 7717/peeri 9437 article	"all samples (n = 229) were randomized as test set and training set, respectively"	Prognostic study	AUC, log-rank p-value (training + external test set)	external cohort validation	TCGA databases. We constructed a reliable OS-related 10-gene signature that was	TCGA databases. We constructed a reliable OS-related 10-gene signature that was not dependent on clinical parameters."
35		dimensional -omics data in precision	Oncology 3			2019 USA	http://dx.doi.org/10 1200/PO.19.0001 8 article	review (not applicable)	Review			including both supervised fearning and unsupervised fearning, and their applications to precision oncology, and we discuss future research directions."	"In this review, we survey current knowledge-guided statistical learning methods, including Both supervised learning and unsupervised learning. and their applications to precision oncology, and we discuss future research directions: "Our study examined the expression patterns in AML samples from the GEO and
350) Zhao, T and Khadka, V S and Deng, Y	cross-platform data analyses Knowledge-guided statistical learning	Aging (Albany NY) 1	ž :	14506- 14 14527	2020 USA	http://dx.doi.org/10 18830/sping 1994 98 article	"Lung cancer datasets were obtained from the Gene Expression Omnibus (GEO, n = 287) and The Cancer Genome Atlas (TCGA, n = 216) repositories"	Case-control study	AUC (training + external validation)	external cohort validation	integrative cross-platform data analyses. This data mining and machine learning approach would be an efficient and economical screening method for tumor biomarker discovery"	"We identified 8 incRNAs as potential diagnosits biomarkers for NSCLC through integrative cross-platform data analyses. This data mining and machine learning approach would be an efficient and economical screening method for tumor biomarker discovery"
345	y Zhang, Y H and U, H and Zeng, T and Chen, L and U, Z and Huang, T and Cai, Y D	Identifying Transcriptomic Signatures and Rules for SARS-CoV-2 Infection				2020 China	http://dx.dei.org/10 3399/set 2020 62 7302 article		Case-control study	Matthews Correlation Coefficient (10-fold CV)		distinguishing COVID-19-infected cases from other respiratory patients with or without virus infection, validating the efficacy and accuracy of our prediction. Therefore, the application of machine learning model may efficiently assist in the identification of potential diagnostic biomarkers and candidate drug targets and help establish a standard worlfullow for related analyses in such field.*	distinguishing COVID-19-infected cases from other respiratory patients with or without virus infection, validating the efficacy and accuracy of our prediction. Therefore, the application of machine learning model may efficiently assist in the identification of potential diagnostic biomarkers and candidate drug targets and help establish a standard worlfullow for related analyses in such field."
348	Zhang, Y and Nock, W and Wyon, M and Weber, Z and Adams, E J and Sarah, A and Stockard, S and Talman, D and Singh, J and San, Z and Man, D and Man, P and Lin, M U and Jang, Y Z and Man, D and Wang, P and Sh, L and Hang, W and Sho, Z M and Wardengo, C and Cherole, M and and Londenge, M B and Sardesan, S and S	relapse in triple negative breast	Cancer Research 8) 4		2020 USA	http://dec.org/10 1155/1536. 7445.880319. meeting 94-05-02 abstract	nrTNBC. We compiled primary tumor clinical and multi-omic data, including transcriptome (n=433], copy number alterations (CMAc; n=317), and mutation in 17t cancer-related genes (n=317), then calculated expression and immune signatures*	Prognostic study	AUC (train + text + external validation)	external cohort validation		on timing of relapse. We identify distinct clinical and genomic features that can be
	Thing, L and Ms, F and Qi, A and Liu, L and Zhang, J and Xu, S and Zhong, Q and Chen, Y and Zhang, C Y and Gu, C			i .	19 6656-6659	2020 China	http://doi.org/10 _1009/s00c02399a article	voluntierer and 149 sichemic stroke pacients (including 117 first-lever ischemic stroke patients and 32 recurrent ischemic stroke patients), and the validation set consisting of 30 healthy volunteers and 48 ischemic stroke patients (including 38 firstever ischemic stroke patients) and 10 recurrent ischemic stroke patients! "We identified 45 primary TNBCs from three publicly-available datasets and characterized each as rfMBC, (TMBC, or characterized each as rfMBC, or characterized each as rfMBC, (TMBC, or characterized each as rfMBC, except each each each each each each each each	Case-control study	AUC (training + text set)		"we develop an optimal model to discriminate inhemic stroke patients from healthy persons with 100% sensitivity and 93.18% specificity. This research may facilistic undestanding the roles of starty and metabolistic in stroke occurrence, holding great potential in clinical stroke diagnosis."	persons with 100% sensitivityand 93.18% specificity. This research may facilitate
348	: Zeng, H and Chen, L and Huang, Y and Luo, Y and Ma, X	Integrative Models of Histopathological Image Features and Omics Data Predict Survival in Head and Neck Squamous Cell Carcinoma	Developme			2020 China	http://ldx.doi.org/10_ 33e94ce8_2020_55 3009article	derived from the Cancer Genome Attas: (TCGA) with information of clinical characteristics, genetic mutation, RNA sequencing, protein expression and histopathological images. Patients were randomly assigned into training (n = 108) or validation (n = 108) sets? "the training set consisting of 52 hadron valuntances and 149 ischemic training	Case-control study	AUC (10-fold CV + validation)		significant prognostic biomarkers for overall survival in patients with HNSCC. The integrative models of genomics, transcriptomics, and proteomics along with histopathological image features may more accurately predict survival outcome than single-omics models, which might contribute to the risk stratification and	"The results indicated that histopathological image features had potential as significant prognosits biomarkers for overall survival in patients with HNSCC. The integrate models disponses, transcriptoms, and proteomics along with histopathological image between some one excurately predict annival outcome than single-omics models, which might contribute to the risk stratification and personalized treatment for cancer patients."
345		Responses in Cancer Cell Lines From Cancer Omics and Detection of Drug Effectiveness Related Methylation	Frontiers in Genetics 1	ı		2020 China	http://dx.doi.org/10 3389/fgene 2000 00917 article	"A dataset of 216 HNSCC patients was	Drug response prediction	AUC (5-foldCV + external validation)			
34-	. Yan, Y and Song, D and Zhang, X and Hui, G and Wang, J		Frontiers in Pharmacolo 8y 1	ı		2020 China	http://dx.doi.org/10 .338/fiphur.2020.0 1155 article	that were diagnosed with non-ischemic heart failure at Second Hospital of Jilin University, Changchun, China, From January 2018 to August 2018 [] In addition, 77 matched control subjects without heart failure were used as control."	Case-control study	AUC (boostrap analysis)		In this study, we used the mining strategy to identify the COX-2 and it micro RNAc, which might be used as idenminates for non-ischemic heart failure. Although the milk-469 and milk 129° per pedicated to target the COX-2, their controllation were week. Further studies were need to confirm that their direct correlations, in addition, the sample size was small and analyzed from a single-center, larger studies are needed to confirm that current resulas."	which might be used as biomarkers for non-ischemic heart failure. Although the miR- 4649 and miR-1297 are predicted to target the COX-2, their correlations were week. Further studies were need to confirm that their direct correlations. In addition, the
343	I Xu, K and Liang, X and Justice, A and So-Armah, K and Krystal, J and Sinha, R	DNA methylation biosignature in blood predicts alcohol consumption in Two distinct populations	Neuropsych opharmacol ogy 4		509-510	2019 USA	http://dx.doi.org/10 _1038s41386-019 meeting 0547-9 abstract	distinct populations (Ntotal =1,530)" "In this study, we enrolled 70 patients	Surrogate biomarker study	AUC, correlation (training + external cohort)		epigenetic mechanisms of alcohol consumption. The DNAm signature associated with PEth shows greater utility on prediction hazardous alcohol drinking in comparison with self-report phenotype. These findings suggest that DNA methylation in blood is a	PEth shows greater utility on prediction hazardous alcohol drinking in comparison
342			BMC Genomics 2	ı 1	650-650	2020 China	http://dx.doi.org/10 1185912894-020- 	see Table 1	Differential diagnosis prediction	accuracy (10-fold CV)		benchmark and 6 state-of-the-art feature selection algorithms."	for gene expression data, MCBFS, which performs clustering and feature weighting in a supervised manner. In the algorithm, a multi-scale distance function designed by us
	Xie, G and Wang, X and Wei, R and Wang, J and Zhao, A and Chen, T and Wang, Y and Zhang, H and Xiao, Z and Liu, X and Deng, Y and Wong, L and Rajani, C and Kwee, S and Bian, H and Gao, X and Liu, P and Jiu, W	associated with the presence of advanced liver fibrosis in Chinese patients with chronic hepatitis B viral	BMC Med 1	3 1	144-144	2020 China	http://dx.dei.org/10 _1188612918-020- 01895-w article	cohort 1 (504 HBV associated liver fibrosis patients and 502 normal controls, NC), we selected a panel of 4		AUC (10-fold CV + validation cohort)	cross-validation + external cohort validation	"Our study showed that this 4-metabolite panel has potential usefulness in clinical assessments of CLD progression in patients with chronic hepatitis B virus infection."	assessments of CLD progression in patients with chronic hepatitis B virus infection."
340	Wang, Y and Zhong, J and Li, Z and Jiang, R and Peng, J and Sun, J and Yang, G and Yang, X R and Huang, A and Wang, Y and Jia, Y and Liu, X and Gao, F and Wu, X and Wang, D and Wu, W and Lou, W and Zhou, J and Put.	Development of a novel liquid biopsy test to diagnose and locate gaztrointestinal cancers Serum metabolite profiles are	Journal of Clinical Oncology 3	3 :	15	2020 China	http://dx.doi.org/10 1700/JCO 2020 3 meeting 8.15 sured.1557 abstract	system, we found that cfDNA mutation profiling achieved a sensitivity of 59.6%, 67.2%, and 46.8% for detecting HCC (n = 322), CRC (n = 244) and PC (n = 141) respectively, with a specificity of 95% in	Case-control study	AUC, sensitivity, specificity (10-fold CV, training + validation cohort)		"Plasma cDNA methylome profiling identified effective biomarkers for the detection and tissue-of-origin determination of GI cancers, and outperformed metation-based detection approach. Therefore, a liquid biopsy test capable of detecting and locating GI cancers is fascible and may serve as a valuable tool for early detection and intervention."	and tissue-of-origin determination of GI cancers, and outperformed mutation-based
339	Wang, Y and Wang, Y and Huang, A and Jiang, R and Zheng, J and Li, Z and Peng, J and Sun, J and Liu, C and Pang, G and Yuan, J and Yang, C and Zhou, J and Fan, J		Cancer Research 8	o :	16	2020 China	http://dx.doi.org/10 _1158/1538- meeting 7445.AM2020-782 abstract	"The training cohort consists of 148 hepatocellular carcinoma cases (median age of 63) and 84 healthy individuals (median age of 60)" "Using a one-specified mutation scoring	Case-control study	AUC (10-fold CV + validation cohort)	cross-validation + external cohort		"In conclusion, our results suggest that cancer-derived abnormal methylation pattern of cfDNA provides promising biomarkers for the diagnosis of HCC (hepatocellular carcinoma) with high sensitivity and specificity."
331		An Eight-CircRNA Assessment Model for Predicting Biochemical Recurrence in Prostate Cancer				2020 China	http://dx.doi.org/10 3389/oet 2020 50 9204 article	"The dataset is from the GEO database, using a cohort of 144 gatients in Canada"	Prognostic study	AUC (10-fold CV)		[Biochemical Recurrence] of PCa [Prostate Cancer] patients. We found that the BCR predicting effect may be related to the tumor microenvironment. At the same time, we preliminarily verified the function of circ_14736 and circ_17720 in vitro. Further	predicting effect may be related to the tumor microenvironment. At the same time,