
O uso de conjuntos de dados de expressão gênica na pesquisa de seleção de atributos

Bruno Iochins Grisci

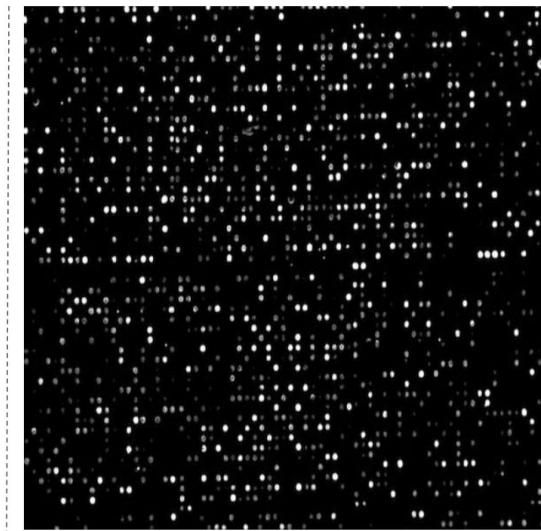
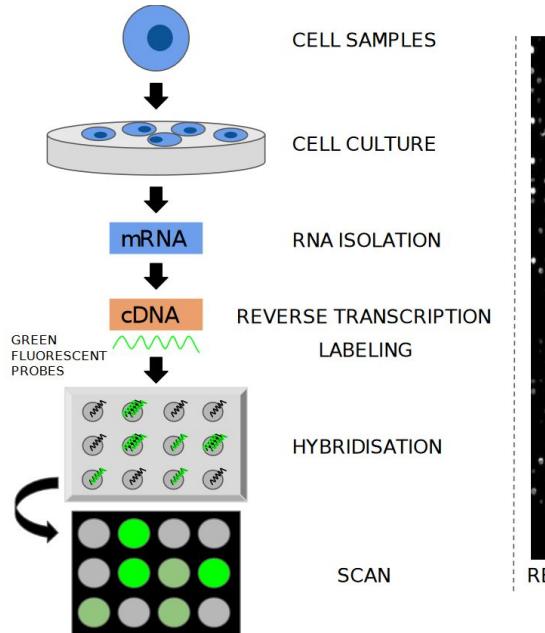
Feature Selection

Generative AI Academy

Part 2: Benchmark datasets

Dados de expressão gênica

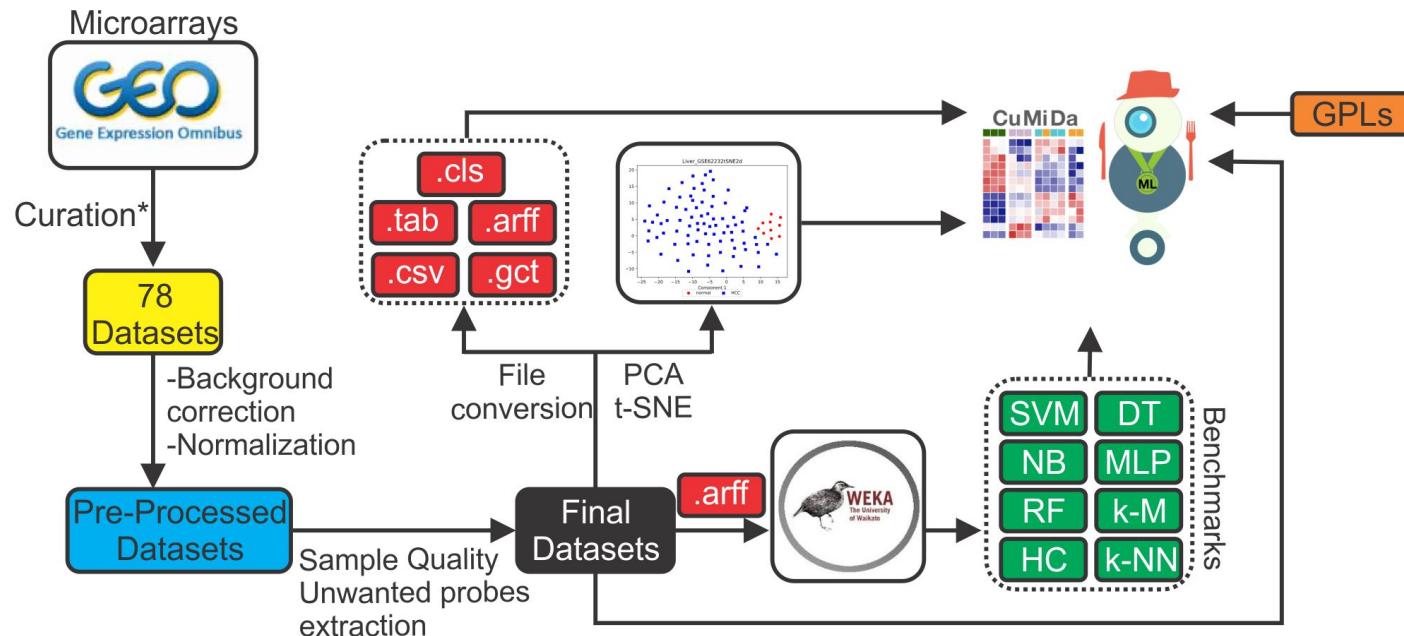
- A estrutura e função de diferentes tipos de células do mesmo organismo são diferentes, enquanto o genoma permanece o mesmo.
- Diferentes genes em cada tipo de célula são expressos.
- Pode ser usado para estudar a diferença na expressão entre tecidos do mesmo tipo sob diferentes condições.
- Novos biomarcadores podem ser potenciais alvos terapêuticos.



Mas na realidade...

	A	B	C	D	E	F	G	H	I	J	K	L
1	Samples	Type	1007_s at	1053_at	117_at	121_at	1255_g_at	1294_at	1316_at	1320_at	1405_i_at	1431_at
2	GSM362958.CEL.gz	HCC	6.80119760710874	4.55318894640395	6.78779002909496	5.43089267353575	3.25022198170724	6.27268846398046	3.4134047464412	3.37490984706688	3.65411613599682	3.80498325183284
3	GSM362959.CEL.gz	HCC	7.58595616226449	4.19354041301431	3.76318259460317	6.00359349944804	3.3093867681399	6.29192656845271	3.75477701364209	3.58760263207689	5.13715938611725	8.6224754628958
4	GSM362960.CEL.gz	HCC	7.8033704814389	4.134075383015	3.43311030956077	5.399505709653068	3.4769439430450601	5.82571255203611	3.50503558852534	3.68733289144417	4.5157535182232	12.681439109932
5	GSM362964.CEL.gz	HCC	6.920844037613297	4.0006514934097	3.75450034120241	5.64529670292514	3.38753047480589	6.47045799349257	3.62924920818033	3.57753407054759	5.19262423621383	11.7594120751621
6	GSM362965.CEL.gz	HCC	6.55648031523744	4.59091049029505	4.06615511575957	6.3445368225612	3.37208126056019	5.43928036335272	3.76221336675644	3.44071407349146	4.9616249065577	10.318552481299
7	GSM362966.CEL.gz	HCC	7.30668391006144	4.36869795997217	3.72489358089459	5.61296709092874	3.2245864145194	6.1120880672554	3.48660154139636	3.3149501822187	4.07110192135605	11.944573918641
8	GSM362970.CEL.gz	HCC	6.51292835592314	4.6456440565343	3.872171884659331	5.6865911726611	3.30420310841847	6.44245508824605	3.37885684430505	3.26577268196747	4.64010462565857	10.1967888008446
9	GSM362971.CEL.gz	HCC	6.82698754213955	4.01936601587141	4.12071959862386	5.7179309504164647	3.21870382202531	6.69032703329845	3.77286897801043	3.51346182154347	5.88813125325086	12.5153547104115
10	GSM362972.CEL.gz	HCC	5.82078827848331	4.1594009295679	3.72900436554028	6.2171154352368	3.18383144008311	6.1305919281839	3.67742123833104	3.401679707856	4.78391990602284	12.6981102815037
11	GSM362976.CEL.gz	HCC	5.9275328265135	3.09591397532946	4.59021756466884	5.54950900963435	3.364777499984	5.07583615012919	3.45456295354757	4.57788427236265	5.6050498565527	4.27799221499875
12	GSM362977.CEL.gz	HCC	5.6285617851269	3.92881959500934	4.8571091301983	6.30642665784121	3.27961837411023	6.54674286974066	3.68233421553639	3.57668134999204	3.7483445874946	3.7483445874946
13	GSM362978.CEL.gz	HCC	6.09841757361098	4.457572870330972	3.99439026148466	5.73292905095275	3.36280682793352	6.97809176114277	3.47418272288103	3.50664240850734	4.20849010621785	4.29009126149271
14	GSM362982.CEL.gz	HCC	6.38762657207279	4.6879538667814	3.85478251355353	6.21822737708569	3.08222302712912	4.23572609261817	3.55217304909753	3.48640332061621	3.76379378225426	7.14459567739736
15	GSM362984.CEL.gz	HCC	6.37522538336142	4.04805125867143	3.67044931930925	5.98933754265255	3.20246625382664	6.4416328687876	3.85477731752766	3.4725149425935	7.17069143839285	11.1891345283344
16	GSM362986.CEL.gz	HCC	6.932664994962496	4.62921242911026	4.21633575129576	5.58179928388985	3.02185457778632	6.20912434026103	3.35290934154755	3.4277749465531	5.46352397101956	12.7311657113786
17	GSM362988.CEL.gz	HCC	7.255358469460416	4.63563029097991	4.32854075845899	5.92269417145194	3.221909044994123	3.5584663738093	3.3640045648158	3.7297853671897	11.7634073358966	11.7634073358966
18	GSM362992.CEL.gz	HCC	6.1699066619275	4.5595903029603	3.44251412571729	6.05767082738335	3.22929038619527	5.9663905720069	3.63740970642988	3.35802913300437	4.32285374565323	8.29601845449646
19	GSM362993.CEL.gz	HCC	8.70212611805967	4.5727589417571091	3.60329349161951	6.50364743723948	3.27761229359352	5.5638834062110307	3.4821352408843	3.68233421553639	3.81659021300882	11.3000047798643
20	GSM362994.CEL.gz	HCC	6.76482163999378	4.2816257303092	3.88421719278073	6.20070231310914	3.19037974490686	4.64989117787010	3.46185272253752	3.31717948115732	3.5482546246707	12.5643275309247
21	GSM363008.CEL.gz	HCC	6.89137236997168	3.89398430136213	4.15418389037624	5.1615195152059	4.61914299664694	5.59398956916429	3.59780910067754	6.67426933341705	11.3200004786843	11.3200004786843
22	GSM363009.CEL.gz	HCC	9.24515459151105	5.52516401623998	4.249273893893	5.32144021626249	3.1495321174594	5.71763061769373	3.51022323584056	4.6211501062122	10.64013498747458	10.64013498747458
23	GSM363010.CEL.gz	HCC	10.2036809385542	4.38750392462001	4.86400382786457	6.04167547322937	3.22782381595483	4.8940545775597	3.54367134935043	3.9351253121344	4.61426213999159	11.6631860481393
24	GSM363011.CEL.gz	HCC	6.78216162520364	5.9337409201189	3.65656705668973	4.07899279401444	5.084037940102139	3.47161243442646	3.784808030512345	3.8121536657709731	5.8278030512345	8.27278030512345
25	GSM363012.CEL.gz	HCC	6.0159272818864	4.94432514265907	4.27235664686957	5.76554771403007	3.15889373220786	5.88695933883493	3.50107439340929	3.6273206180902	4.56650861552184	12.60925170710782
26	GSM363013.CEL.gz	HCC	7.480112535537	4.4882957265875	4.20947617853123	5.16121006448355	3.20182366910594	6.2281552895264	3.67104654835145	3.5287752999281	5.5307799103044	7.3817224553906
27	GSM363015.CEL.gz	HCC	7.73372153457663	5.61893987311495	3.8275012532958	5.05102838028321	3.02710918900597	3.137903456483817	3.4719177008197	3.6201985934873	3.9690470258623	13.3174213352452
28	GSM363016.CEL.gz	HCC	5.868091951266974	4.24639962256086	3.86159690951172	5.83431771922282	3.14179016625245	6.87468444455035	3.68695132989409	3.427561976734	3.92303394246396	9.551368672422325
29	GSM363017.CEL.gz	HCC	6.28898429229956	5.3091460774403119	4.10753938568554	6.0052554226958	3.92690219512167	3.4649327522747	4.10768917344	3.51022323584056	3.76302277656437	11.4588677390233
30	GSM363029.CEL.gz	HCC	6.84321505900151	4.65452880686373	5.4863339138624	4.94066602241705	3.11064895777895	4.5989340954826	3.137062482615	3.5783443808277	3.92602777656437	9.97755055922057
31	GSM363030.CEL.gz	HCC	6.8863897709349	4.19656066484294	4.2024210705086	4.98228087680507	3.05751215465252	5.6891784409978	3.40367769094996	3.30041994003122	3.62365168020821	9.97755055922057
32	GSM363031.CEL.gz	HCC	5.7366416840904	4.3292110700053	3.99689832574993	4.92656687526748	3.1452776472711	6.36206320797895	3.90432661553501	3.53895072957314	4.70835494108223	10.0652927128463
33	GSM363032.CEL.gz	HCC	6.68190339487203	3.65907318228896	3.8776286480877	5.5487926658511	3.12963120334136	6.4337151926667	3.77200741949821	3.67379686009843	6.20879328549831	11.89023268658131
34	GSM363033.CEL.gz	HCC	6.60768163260233	4.171716477834	3.84810776580581	5.01628590431952	3.1570645228569	6.01398600242655	3.67494302626603	3.49381996070804	5.767936336988196	6.21707165280241
35	GSM363034.CEL.gz	HCC	6.96794164857108	4.70058932127659	3.9034592291184	5.33699062327665	3.25378701235011	5.08553463505671	3.736505118154818	3.64212624955036	3.73654204246175	12.2657494635353
36	GSM363035.CEL.gz	HCC	6.26775948784058	3.86720135022778	3.92971488920237	5.18173457424161	4.620401336620377	3.32790719688412	5.81734404457589	3.54364806214208	4.8616909856137	11.9277627380745
37	GSM363036.CEL.gz	HCC	6.03010104997145	4.8486055085024	4.07228434735697	5.3024981394078	3.2793633737947	5.17558789348264	3.30793699307095	3.139545975651585	3.32652874919996	12.35628674919996
38	GSM363037.CEL.gz	HCC	5.96883440092929	5.0727098282114	4.016477647221287	5.1152493091797	3.20668793019304	5.83688576040371	3.50907190473833	3.59986724370337	3.23926953342282	12.1894741115421
39	GSM363038.CEL.gz	HCC	6.2681183983551	5.654278811044	3.9896898325448947	4.92656687526748	3.1452776472711	6.2621528859529	3.5194453586038	3.62498026973503	3.1339307061918	4.6272119271966
40	GSM363039.CEL.gz	HCC	6.15177638832654	4.1807100592337	3.78708038138071	5.73965445035121	3.176549465503121	6.013884658068493	3.7101801402424	3.74184308608491	3.54642617687614	13.342468129499
41	GSM363048.CEL.gz	HCC	6.32087223275949	4.6231477958137	3.84426476947652	5.15226545591636	3.59180393691951	4.670738178334	3.29603236655585	3.56748716292217	3.70368276081107	8.88235247011474
42	GSM363049.CEL.gz	HCC	5.25292538255852	4.716857215331	3.659390623595132	5.36359284244553	3.059765311901	5.18867638033434	3.6614628150378	3.7736917373794	12.26849150415012	12.26849150415012
43	GSM363053.CEL.gz	HCC	6.59501074902191	5.37545270119183	3.6798496694489	5.14456193106467	3.187868453730654	5.85145938206153	3.4076425170495	4.1232396034402	3.83210369886728	9.04127206061433
44	GSM363054.CEL.gz	HCC	6.9774949409328	3.792093170923	4.2345680550824	5.30248234735697	5.28243288475695	4.69057893009893	3.059676018627489	3.0626743990211	5.855680545505	3.5447560612553
45	GSM363056.CEL.gz	HCC	8.14583616094427	5.58254840036409	4.244527755909892	5.39676018627489	3.0262743990211	5.855680545505	3.54230719043621	3.5447560612553	3.38426833831052	3.80498325183284
46	GSM363057.CEL.gz	HCC	6.17618317374308	3.87083017023233	7.23358608907206	5.17933205282205	3.08252431101478	5.715420521510122	3.40253690512998	3.67723814063241	3.54230719043621	11.8019131820272

CuMiDa



Ang et al. (2015)

5.1 Dataset Analysis

Based on Tables 3, 4, and 5, the five most commonly used gene microarray expression datasets in the literatures are leukemia [67], colon [64], prostate [69], Diffuse Large B-Cell Lymphoma (DLBCL) [74], and Small round blue cell tumor (SRBCT) of childhood datasets [76]. The comparison of validation result using highest prediction accuracy based on CV and the number of selected genes in each literature are shown in Table 7, 8, 9, 10, and 11.

Supervised, Unsupervised, and Semi-Supervised Feature Selection: A Review on Gene Selection

Jun Chin Ang, Andri Mirza, Habibollah Haron, and Haza Nuzly Abdull Hamed

Abstract—Recently, feature selection and dimensionality reduction have become fundamental tools for many data mining tasks, especially for processing high-dimensional data such as gene expression microarray data. Gene expression microarray data comprises up to hundreds of thousands of features with relatively small sample size. Because learning algorithms usually do not work well with this kind of data, a challenge to reduce the data dimensionality arises. A huge number of gene selection are applied to select a subset of relevant features for model construction and to seek for better cancer classification performance. This paper presents the basic taxonomy of feature selection, and also reviews the state-of-the-art gene selection methods by grouping the literatures into three categories: supervised, unsupervised, and semi-supervised. The comparison of experimental results on top 5 representative gene expression datasets indicates that the classification accuracy of unsupervised and semi-supervised feature selection is competitive with supervised feature selection.

Index Terms—Feature selection, gene expression, semi-supervised, supervised, unsupervised

1 INTRODUCTION

FEATURE selection is a dimensionality reduction technique that is commonly used in the fields of machine learning, pattern recognition, statistics, and data mining. This technique aims to select a subset of relevant features from the original set of features according to some criteria. Some examples of feature selection techniques include Information Gain, Relief, Chi Squares, Fisher Score, and Lasso. Feature selection usually is used in the domains where the datasets comprise of thousands of features but with relatively small sample size (e.g., gene expression data). Feature selection that is applied to gene expression data is also known as gene selection [1]. Gene selection is necessary as the data usually contains many irrelevant, redundant, and noisy expressions, and also is effective for early tumor detection and cancer discovery as it leads to a more reliable cancer diagnosis or prognosis and better clinical treatment [2].

The gene expression data can either be fully labeled, unlabeled, or partially labeled. This leads to the development of supervised, unsupervised and semi-supervised gene selection to discover the biological patterns and class prediction [3]. Typically, unlabeled data is composed of samples and their features without the presence of labels

process of selecting a feature subset based on some criteria for measuring importance and relevance of the features by using labeled data. Unsupervised feature selection, where there is no prior knowledge about the true functional classes, evaluates feature relevance by exploiting the innate structures of the data, such as data variance, separability, and data distribution. A semi-supervised feature selection integrates a small amount of labeled data into unlabeled data as additional information to improve the performance of an unsupervised feature selection. Even though there are a lot of review papers on gene selection in the literatures [4], [5], [6], [7], [8], but to the best of our knowledge, there is no detailed discussion on the methods and separates them in such three categories as in this paper.

This paper is divided into six sections. Section 2 presents an overview of feature selection and Section 3 gives a review on some gene selection approaches especially that have been proposed over the past five years and further group them in three categories: supervised, unsupervised and semi-supervised approaches. Section 4 describes the challenges inherent in gene selection task. Section 5 discusses the gene selection approaches reviewed in the previous sections and collates the experimental results of gene

1999

5.1 Dataset Analyses

1999

2002

Based on Tables 3, 4, and 5, the five most commonly used gene microarray expression datasets in the literatures are leukemia [67], colon [64], prostate [69], Diffuse Large B-Cell Lymphoma (DLBCL) [74], and Small round blue cell tumor (SRBCT) of childhood datasets [76]. The comparison of validation result using highest prediction accuracy based on CV and the number of selected genes in each literature are shown in Table 7, 8, 9, and 11.

2000 **2001**



Gene expression dataset (Golub et al.)

Data Code (83) Discussion (4) Metadata

▲ 246

New Notebook

Download (1 MiB)



❖ Activity Overview

ACTIVITY STATS

VIEWS

104077

DOWNLOADS

12399

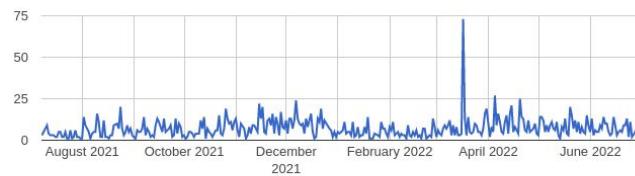
DOWNLOAD PER VIEW RATIO

0.12

TOTAL UNIQUE CONTRIBUTORS

77

Downloads ▾



NOTEBOOKS STATS

NOTEBOOKS

83

NOTEBOOK COMMENTS

104

UPVOTE PER NOTEBOOK RATIO

3.96

NOTEBOOK UPVOTES

329

TOP CONTRIBUTORS



DISCUSSION STATS

TOPICS

4

TOTAL COMMENTS

9

UPVOTE PER POST RATIO

2.56

DISCUSSION UPVOTES

23

<https://www.kaggle.com/datasets/crawford/gene-expression/data>

Bolón-Canedo et al. (2014)

Table 5
Dataset description for binary datasets: s and f are the number of samples and features, respectively.

Dataset	s	f	Distribution	Original Ref.	Year	Where used	Download
B.MD	34	7129	26–74%	[95]	2002	[101]	unknown
Bone Lesion	173	12,625	unknown	[117]	2003	[101]	unknown
Brain	21	12,625	33–67%	[88]	2003	[27,120]	[2]
Brain_Tumor1	60	7129	unknown	unknown	[81]	unknown	unknown
Brain_Tumor2	50	12,625	unknown	unknown	[81]	unknown	unknown
Breast	22	3226	unknown	[58]	2001	[71]	unknown
Breast Cancer	97	24,481	unknown	[118]	2002	[24,23,27]	[8]
Breast-epi	19	24,481	33–44%	[118]	2002	[23,48]	[8]
Breast-train	78	24,481	56–44%	[118]	2002	[22,48]	[8]
BreastTIR	49	7129	49–51%	[123]	2001	[101]	unknown
BR-E849	49	6817	49–51%	[123]	2001	[73]	unknown
C.MD	60	7129	35–65%	[95]	2002	[101]	unknown
Celiac	132	22,185	unknown	[57]	2009	[101]	unknown
CNS/Embryonal-T	60	7129	35–75%	[95]	2002	[124,24,25,23,27,105,125]	[2,3]
Colon	62	2000	35–65%	[14]	1999	[39,101,78,107,80,71,124,11,24,23,27,73,85,48,122,108,31,125,120]	[2,7,8,3]
Colon-epi	202	44,290	unknown	[87]	2000	[101]	unknown
DLBC	77	5470	75–25%	[104]	2002	[39,33,122,125]	[110]
DLBCL	47	4026	49–51%	[13]	2000	[71,24,25,23,27,105]	[8,2]
DLBCL	77	7129	75–25%	[104]	2002	[124,81]	[2,9]
GLU-85	85	22,28	31–65%	[40]	2004	[27]	[5]
Leukemia/ ALL/AML	72	7129	31–65%	[46]	1999	[39,101,107,71,124,11,86,24,25,27,73,85,31,120]	[2,8]
Leukemia/test	34	7129	71–20%	[46]	1999	[119,23,48,108]	[3,2,8]
Leukemia_train	38	7129	59–41%	[46]	1999	[119,23,48,108]	[3,2,8]
Lung	52	918	75–25%	[43]	2001	[101]	unknown
Lung	181	12,533	83–17%	[49]	2002	[24,25,27,48]	[8,2]
Lung	410	2428	34–66%	[103]	2008	[31]	unknown
Lung_test	149	12,533	90–10%	[49]	2002	[23,48]	[8,2]
Lung_train	32	12,533	50–50%	[49]	2002	[23,48]	[8,2]
LUNG181	181	12,600	17–83%	[49]	2002	[73]	unknown
LYM77	77	6817	25–75%	[104]	2002	[73]	[8]
Lymphoma/B- cell	45	4026	49–51%	[13]	2000	[11,108]	[2]
Moffitt colon cancer	122	2619	31–65%	[38]	2005	[31]	unknown
Ovarian	253	15,154	36–64%	[93]	2002	[24,25,23,27,21,107]	[8]
Prostate	102	6033	51–49%	[106]	2002	[78,11]	[2,8]
Prostate	136	12,600	43–57%	[106]	2002	[124,24,25,27]	[8]
Prostate-test	34	12,600	26–74%	[106]	2002	[102,85,23,48]	[8,9]
Prostate-train	102	12,600	49–51%	[106]	2002	[102,86,73,85,81,23,48,105]	[110]
Prostate_Tumor	102	10,059	51–49%	[106]	2002	[39,33,125,120]	unknown
SMK-CAN-187	187	19,983	48–52%	[109]	2007	[27,119,108]	[5]

Information Sciences 282 (2014) 111–135



Contents lists available at ScienceDirect

Information Sciences

journal homepage: www.elsevier.com/locate/ins



A review of microarray datasets and applied feature selection methods

V. Bolón-Canedo ^{a,*}, N. Sánchez-Marcano ^a, A. Alonso-Betanzos ^a, J.M. Benítez ^b, F. Herrera ^{b,c}



CrossMark

Machine learning researchers due
feature selection has been soon
tion, and a huge number of fea
he input dimensionality while
devoted to reviewing the most
d and the microarray databases
interested reader aware of the
the imbalance of the data, their
mental evaluation on the most
methods is presented, bearing
ection method, but to facilitate

Elsevier Inc. All rights reserved.

of research both in bioinfor
cell samples regarding gene
es of tumor. Although there
r of features in the raw data
on task is to separate healthy
re also datasets in which the
are more complicated.

Therefore, microarray data pose a serious challenge for machine learning researchers. Having so many fields relative to co

<https://doi.org/10.1016/j.ins.2014.05.042>

Bolón-Canedo et al. (2014)

4.1. DNA microarray repositories

Although in the initial development of DNA microarray data analysis it was difficult to find datasets to deal with, in recent years there has been a growing number of public microarray data repositories of a wide spectrum of cancer types available for the scientific community. The most famous are listed below:

- *ArrayExpress*, from the European Bioinformatics Institute [1]:
<http://www.ebi.ac.uk/arrayexpress/>

Table 4

Other feature selection methods used on microarray data. Type of evaluation (ranker/subset) and type of data (binary/multiclass).

Method	Original Ref.	Type (r/s)	Data (b/m)
CFS-TGA	[33]	s	m
E1-cp	[25]	s	b
E1-ni	[25]	s	b
E1-ns	[25]	s	b
E2	[25]	s	b
Ensemble RFE	[11]	s	b
EF	[24]	s	m
FAST	[108]	s	m
GADP	[71]	s	m
GC	[78]	r	b
MCF-RFE	[124]	s	b
MFMW	[73]	s	b
R-m-GA	[105]	s	b
SRF	[125]	s	m
SVM-RFE with MRM	[85]	r	b

- *Cancer Program Data Sets*, from the Broad Institute [2]:
<http://www.broadinstitute.org/cgi-bin/cancer/datasets.cgi>
- *Datasets Repository*, from the Bioinformatics Research Group of Universidad Pablo de Olavide [3]:
<http://www.upo.es/eps/bigs/datasets.html>
- *Feature Selection Datasets*, from Arizona State University [5]:
<http://featureselection.asu.edu/datasets.php>
- *Gene Expression Model Selector*, from Vanderbilt University [110]:
<http://www.gems-system.org>
- *Gene Expression Omnibus*, from the National Institutes of Health [6]:
<http://www.ncbi.nlm.nih.gov/geo/>
- *Gene Expression Project*, from Princeton University [7]:
<http://genomics-pubs.princeton.edu/oncology/>
- *Kent Ridge Bio-Medical Dataset Repository*, from the Agency for Sciency, Technology and Research [8]:
<http://datam.izr.a-star.edu.sg/datasets/krbd>
- *Stanford Microarray Database*, from Stanford University [10]:
<http://smd.stanford.edu/>

Problemas com dados

- **Northcutt et al. (2021)**: 3,4% dos rótulos nos 10 datasets mais usados em visão computacional, linguagem natural e áudio contêm erros.
- **Roberts et al. (2021)**: nenhum dos modelos treinados para diagnóstico de COVID-19 a partir de radiografia ou tomografia do tórax pode ser usado na prática clínica por causa de falhas metodológicas, incluindo vieses nos conjuntos de dados.
- **Koch et al. (2021)**: um dos problemas atuais em aprendizado de máquina são estudos de *benchmarking* usando dados para tarefas para as quais eles não foram originalmente desenhados.

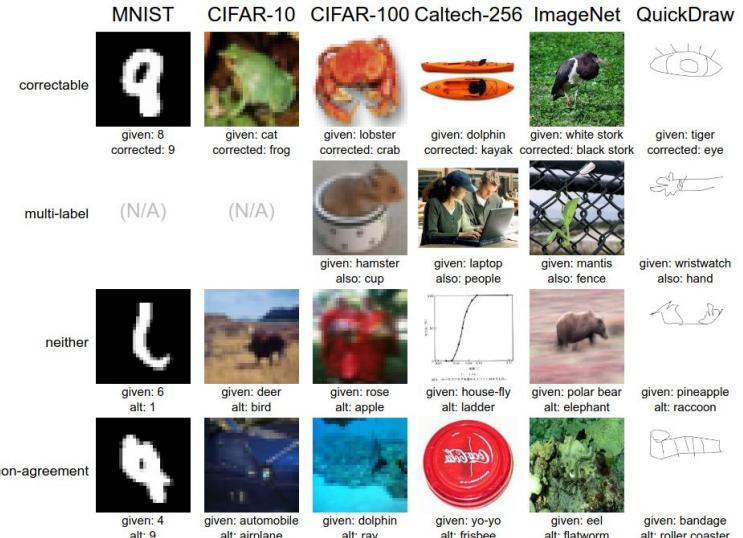
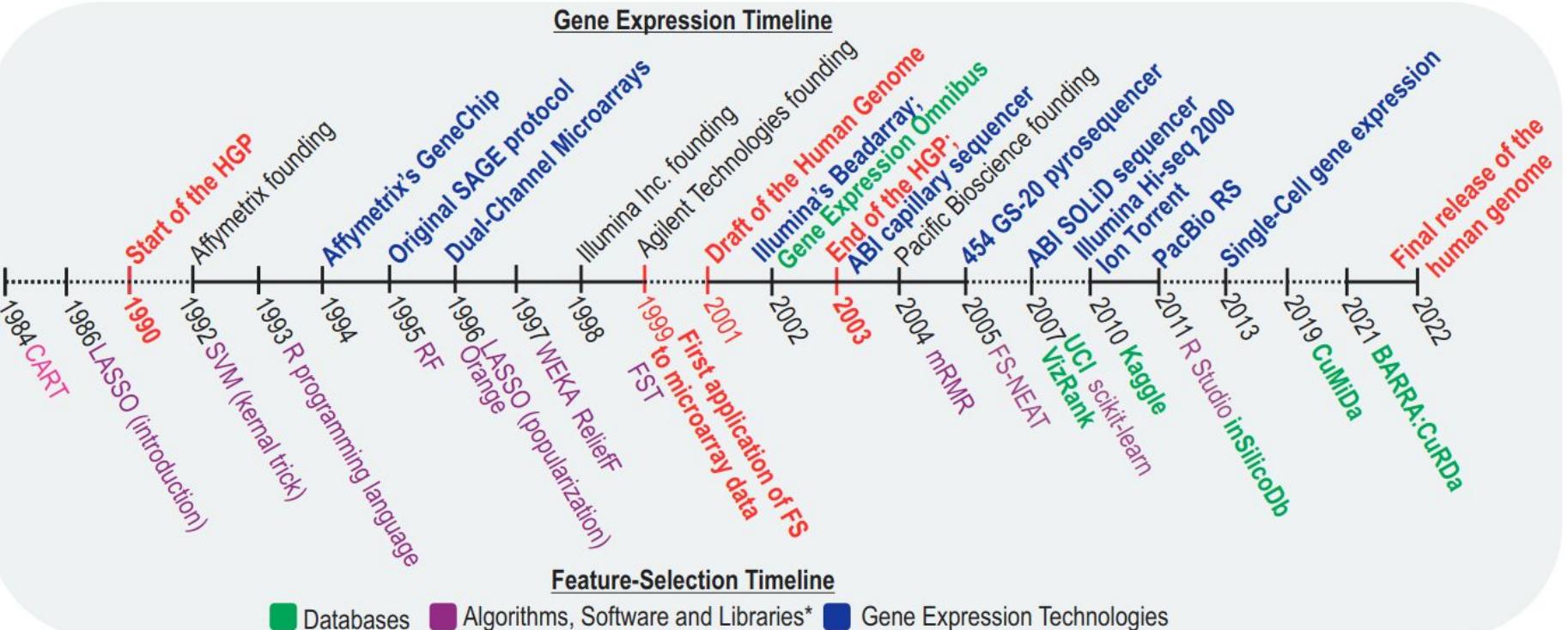


Figure 1: An example label error from each category (Section 4) for image datasets. The figure shows given labels, human-validated corrected labels, also the second label for multi-class data points, and CL-guessed alternatives. A gallery of label errors across all 10 datasets, including text and audio datasets, is available at <https://labelerrors.com>.

[arXiv:2103.14749](https://arxiv.org/abs/2103.14749)





Vamos à revisão

Bases de artigos



	A	B	C	D	E	F	G	H	I	J	K	L	M	
1	Supplementary Material 2: The use of gene expression datasets in feature selection research: 20 years of inherent bias?													
2	Bruno Iochins Grisci, Bruno César Feltes, Joice de Faria Poloni, Pedro Henrique Narloch, and Márcio Dorn													
3	S1-Table PUBLICATIONS													
4	ID	Source	Search1	Search2	Title	Authors	DOI	Year	Publisher	Volume	Issue	Pages	PMID	
5	P0007	WOS	RNA-SEQ	FEATURE SELECTIO	Omic and Electronic Health Record Big Data Analytics for Precision Medicine	Wu, PY; Chei	10.1109/TB	2017	IEEE TRANS	64	2	263-273	277404	
6	P0010	PUBMED	MICROARRA	FEATURE SELECTIO	A BAYESIAN NONPARAMETRIC MIXTURE MODEL FOR SELECTING GENES AND GENE SUB	Zhao Y, Kang	10.1214/14	2014	Ann Appl Stat	8	2	999-1021	259842	
7	P0011	IEEE	MICROARRA	FEATURE SELECTIO	A binary PSO feature selection algorithm for gene expression data	S. Dara; H. B	10.1109/ELC	2014	2014 Intern. NA	NA		1-6	NA	
8	P0012	IEEE	MICROARRA	FEATURE SELECTIO	A Biologically Verified Classification of Microarray Data	R. Mondal; E	10.1109/CIC	2014	2014 Intern. NA	NA		686-690	NA	
9	P0013	SCOPUS	MICROARRA	FEATURE SELECTIO	A biology-driven approach identifies the hypoxia gene signature as a predictor of the outcome	Fardin P, Ba	10.1186/14	2010	Molecular Cell	9	NA	NA	206242	
10	P0015	PUBMED	MICROARRA	FEATURE SELECTIO	A CBR framework with gradient boosting based feature selection for lung cancer subtype	Ramos-Gonz	10.1016/j.cc	2017	Comput Biol	86	NA	98-106	285273	
11	P0016	PUBMED	MICROARRA	FEATURE SELECTIO	A centroid-based gene selection method for microarray data classification.	Guo S, Guo L	10.1016/j.jtbt	2016	J Theor Biol	400	NA	32-41	270567	
12	P0017	PUBMED	MICROARRA	FEATURE SELECTIO	A class imbalance-aware Relief algorithm for the classification of tumors using microarray	He Y, Zhou J,	10.1016/j.cc	2019	Comput Biol	80	NA	121-127	309470	
13	P0018	IEEE	RNA-SEQ	FEATURE SELECTIO	A Class-Information-Based Sparse Component Analysis Method to Identify Differentially Expressed Genes	J. Liu; Y. Xu;	10.1109/TC	2016	IEEE/ACM T	13	2	392-398	NA	
14	P0020	SCOPUS	MICROARRA	FEATURE SELECTIO	A classification model on tumor cancer disease based mutual information and firefly algorithm	Jabbar S.F.	10.21533/p	2019	Periodicals of Computing	7	3	1152-1162	NA	
15	P0021	SCOPUS	MICROARRA	FEATURE SELECTIO	A clustering approach for feature selection in microarray data classification using random forest	Aydadent F	10.3745/JIP	2018	Journal of Infor	14	5	1167-1175	NA	
16	P0023	SCOPUS	MICROARRA	FEATURE SELECTIO	A clustering-based method for gene selection to classify tissue samples in lung cancer	Castellanos-	10.1007/97	2016	Advances in Bioinformatio	477	NA	99-107	NA	
17	P0024	SCOPUS	MICROARRA	FEATURE SELECTIO	A combining dimensionality reduction approach for cancer classification	Han L., Zhou	10.1007/97	2015	Lecture Notes in Bioinform	9426	NA	340-347	NA	
18	P0025	PUBMED	MICROARRA	FEATURE SELECTIO	A comparative analysis of biomarker selection techniques.	Dessi N, Pasquale	10.1155/20	2013	Biomed Res	2013	NA	387673	243249	
19	P0026	IEEE	MICROARRA	FEATURE SELECTIO	A comparative analysis of feature selection algorithms on classification of gene microarray	J. Jeyachidra	10.1109/ICI	2013	2013 Intern. NA	NA		1088-1093	NA	
20	P0027	SCOPUS	MICROARRA	FEATURE SELECTIO	A comparative analysis of feature selection methods for biomarker discovery in study of breast cancer	Zhang X, Jor	10.1007/97	2019	Communicat	1056	CCIS	NA	114-123	NA
21	P0028	PUBMED	MICROARRA	FEATURE SELECTIO	A comparative analysis of swarm intelligence techniques for feature selection in cancer classification	Gunavathi C,	10.1155/20	2014	ScientificWorld	2014	NA	69381	251573	
22	P0029	IEEE	MICROARRA	FEATURE SELECTIO	A Comparative Performance Evaluation of Random Forest Feature Selection on Classification	M. A. Latief;	10.1109/ICI	2019	2019 3rd Int. NA	NA		1-6	NA	
23	P0030	SCOPUS	MICROARRA	FEATURE SELECTIO	A Comparative Performance Evaluation of Supervised Feature Selection Algorithms on Microarray	Arun Kumar	10.1016/j.pnsc	2017	Procedia Comput	115	NA	209-217	NA	
24	P0032	SCOPUS	MICROARRA	FEATURE SELECTIO	A Comparative study and analysis of data mining classifiers for microarray based cancer prediction	Subasree S.,	10.1145/29	2016	ACM Internat	25-26-Aug-2016	NA	NA	NA	NA
25	P0033	SCOPUS	MICROARRA	FEATURE SELECTIO	A comparative study of feature selection and classification methods for gene expression	Abusamra H	10.1016/j.pnsc	2013	Procedia Comput	23	NA	mai./14	NA	
26	P0034	SCOPUS	MICROARRA	FEATURE SELECTIO	A comparative study of feature selection and classification techniques for high-throughput	Alkuhlani A.	10.1007/97	2017	Advances in Bioinformatio	533	NA	793-803	NA	
27	P0035	IEEE	MICROARRA	FEATURE SELECTIO	A Comparative Study of Feature Selection Methods on Genomic Datasets	J. Rahimipour	10.1109/CB	2019	2019 IEEE 3 NA	NA		471-476	NA	
28	P0036	IEEE	MICROARRA	FEATURE SELECTIO	A comparative study of gene selection methods for cancer classification using microarray	M. Babu; K. S.	10.1109/ICI	2016	2016 Second Int. NA	NA		204-211	NA	
29	P0037	SCOPUS	MICROARRA	FEATURE SELECTIO	A Comparative Study of Gene Selection Methods for Microarray Cancer Classification	Alshamali F	10.1007/97	2019	Lecture Note in Bioinform	520	NA	585-595	NA	
30	P0038	SCOPUS	MICROARRA	FEATURE SELECTIO	A Comparative Study of Improvements Filter Methods Bring on Feature Selection Using	Mi Wang Y, Fan	10.1007/97	2014	Lecture Note in Bioinform	8423	LNCS	NA	55-62	NA
31	P0039	PUBMED	MICROARRA	FEATURE SELECTIO	A comparative study of improvements Pre-filter methods bring on feature selection using	Wang Y, Fan	10.1186/20	2014	Health Inf Sci	2	NA	7	258256	
32	P0040	SCOPUS	MICROARRA	FEATURE SELECTIO	A comparative study of multiclass feature selection on RNAseq and microarray data	Zhang S., Wu	10.1504/IIC	2019	International J	12	2	128-142	NA	

<https://doi.org/10.1002/widm.1523>

A	B	C	D	E	F	G				
1	Supplementary Material 2: The use of gene expression datasets in feature selection research: 20 years of inherent bias?									
2	Bruno Iochins Grisci, Bruno César Feltes, Joice de Faria Poloni, Pedro Henrique Narloch, and Márcio Dorn									
3	S2-Table DATASETS									
4	ID	DATA	DATASET	YEAR	METHOD	SPECIES	BACKGROUND	DATASETS	Link Dataset	AKA
665	Marti et al. (2015)			2015	Microarray	<i>Equus ferus</i>	Allergies		https://www.sciencedirect.com/science/article/pii/S0165242715001397?via%3Dihub#se	
202	Arisi et al. (2011)			2011	Microarray	<i>Mus musculus</i>	Alzheimer's Disease		https://content.iospress.com/articles/journal-of-alzheimers-disease/jad101881	
795	GSE11882			2008	Microarray	<i>Homo sapiens</i>	Alzheimer's Disease		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE11882	Berchtold et al.
258	GSE1297			2004	Microarray	<i>Homo sapiens</i>	Alzheimer's Disease		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE1297	Blalock et al.
259	GSE5281			2006	Microarray	<i>Homo sapiens</i>	Alzheimer's Disease		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE5281	
260	Ray et al. (2007)			2007	Microarray	<i>Homo sapiens</i>	Alzheimer's Disease		https://www.nature.com/articles/nm1653#MOESMS	
441	E-TABM-940			2010	Microarray	<i>Homo sapiens</i>	Amyotrophic lateral sclerosis		https://www.ebi.ac.uk/arrayexpress/experiments/E-TABM-940/	
100	Yasvoina et al. (2013)			2013	RNA-Seq	<i>Mus musculus</i>	Amyotrophic lateral sclerosis		https://www.jneurosci.org/content/33/18/7890	
546	GSE83091			2016	Microarray	<i>Homo sapiens</i>	Appendicitis		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE83091	Chawla et al.
350	GSE20129			2010	Microarray	<i>Homo sapiens</i>	Atherosclerosis		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE20129	
379	GSE23746			2010	Microarray	<i>Homo sapiens</i>	Atherosclerosis		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE23746	
156	GSE25507			2010	Microarray	<i>Homo sapiens</i>	Autism		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE25507	
230	E-MTAB-37			2008	Microarray	<i>Homo sapiens</i>	Biology of multiple cancers		https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-37/	
349	GSE2564			2005	Microarray	<i>Homo sapiens</i>	Biology of multiple cancers		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE2564	Getz et al. (2005)
595	GSE32474			2011	Microarray	<i>Homo sapiens</i>	Biology of multiple cancers		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE32474	Reinhold et al.
606	GSE52582			2013	Microarray	<i>Homo sapiens</i>	Biology of multiple cancers		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE52582	Stafford et al.
302	GSE5364			2006	Microarray	<i>Homo sapiens</i>	Biology of multiple cancers		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE5364	
191	GSE68086			2015	RNA-Seq	<i>Homo sapiens</i>	Biology of multiple cancers		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE68086	
166	GSE74251			2015	RNA-Seq	<i>Homo sapiens</i>	Biology of multiple cancers		https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE74251	
44	Pan-Cancer dataset			-1			Biology of multiple cancers			
791	Pfister et al. (2009)			2009	Microarray	<i>Homo sapiens</i>	Biology of multiple cancers		https://mct.aacrjournals.org/content/8/7/1878	
66	Ramaswamy et al. (2001)			2001	Microarray	<i>Homo sapiens</i>	Biology of multiple cancers		https://www.pnas.org/content/98/26/15149.short	
84	Ross et al. (2000)			2000	Microarray	<i>Homo sapiens</i>	Biology of multiple cancers		https://pubmed.ncbi.nlm.nih.gov/10700174/	
187	Staunton et al. (2001)			2001	Microarray	<i>Homo sapiens</i>	Biology of multiple cancers		https://www.pnas.org/content/pnas/98/19/10787.full.pdf	
30	TCGA			-1	RNA-Seq	<i>Homo sapiens</i>	Biology of multiple cancers		http://www.cbioportal.org/study/summary?id=brca_tcga_pub2015	

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	Supplementary Material 2: The use of gene expression datasets in feature selection research: 20 years of inherent bias?												
2	Bruno Iochins Grisci, Bruno César Feltes, Joice de Faria Poloni, Pedro Henrique Narloch, and Márcio Dorn												
3	S3-Table LINKS												
4	Yellow if indirectly cited												
5	PAPER ID	DATASET	DATASET	DATASET	DATASET	DATASET	DATASET	DATASET	DATASET	DATASET	DATASET	DATASET	DATA
6	P0020	Broken Link	Alon (1999)										
7	P0026	Broken Link	Alon (1999)										
8	P0032	Dyrskjot et al. (2)	GSE349_350	GSE3726	Missing	Golub et al. (1 ^c)	Bhattacharjee	Khan et al. (20)	GSE2443				
9	P0038	GSE36961	GSE36964										
10	P0044	Golub et al. (19 ^c)											
11	P0050	GSE15471											
12	P0056	Golub et al. (19 ^c)	Alon (1999)	Broken Link	Pomeroy et al.								
13	P0062	Missing	Golub et al. (1 ^c)	Alon (1999)	Khan et al. (20)								
14	P0068	Golub et al. (19 ^c)	Khan et al. (20)	Armstrong et al.									
15	P0080	GSE54460											
16	P0086	Golub et al. (19 ^c)	Singh et al. (20)	Shipp et al. (20)	Veer et al. (20)	Bhattacharjee	Kong et al. (20)	Alon (1999)	Wu et al. (200)	Boyle et al. (20)	Zhang et al. (20)	Smith et al. (20)	
17	P0092	GSE54460											
18	P0098	Broken Link	Golub et al. (1 ^c)	Khan et al. (20)	Alon (1999)	Singh et al. (20)	Veer et al. (20)						
19	P0104	Broken Link	Golub et al. (1 ^c)	Khan et al. (20)	Shipp et al. (20)	Alon (1999)							
20	P0122	Pollen et al. (20 ^c)											
21	P0128	Golub et al. (19 ^c)	Alon (1999)	Khan et al. (20)	Shipp et al. (20)	Armstrong et al.	Pomeroy et al.						
22	P0134	Golub et al. (19 ^c)	Alon (1999)	Khan et al. (20)	Armstrong et al.	Singh et al. (20)							
23	P0140	Simulation	Authors Data										
24	P0146	Golub et al. (19 ^c)	Pomeroy et al.										
25	P0152	Simulation	TCGA										
26	P0158	Golub et al. (19 ^c)	Pomeroy et al.										
27	P0164	Simulation	Broken Link	Golub et al. (1 ^c)	Armstrong et al.	Pomeroy et al.	Khan et al. (20)						
28	P0170	Alon (1999)	Khan et al. (20)	Singh et al. (20)	Shipp et al. (20)	Armstrong et al.	Bhattacharjee						
29	P0176	Shipp et al. (200)	Golub et al. (1 ^c)	Singh et al. (20)									
30	P0182	Missing	Shipp et al. (20)	Khan et al. (20)	Singh et al. (20)	Bhattacharjee	Pomeroy et al.						

	All	IEEE Xplore	PubMed	WOS	Scopus
Number of publications	1284	375	315	47	547
Mean number of datasets per publication	3.94 ± 3.38	3.76 ± 3.17	3.66 ± 3.47	3.66 ± 3.03	4.26 ± 3.48
Median number of datasets per publication	3	3	2	3	3
Min number of datasets per publication	1	1	1	1	1
Max number of datasets per publication	27	27	25	12	27
Indirect citations ^a	32.3%	39.4%	15.5%	23.4%	37.8%
Broken links ^b	7.7%	8.5%	5.3%	10.6%	8.4%
Missing reference ^c	5.1%	9.6%	0.3%	6.3%	4.7%
Simulated data ^d	6.5%	8.0%	8.5%	6.3%	4.3%
Author's dataset ^e	4.6%	2.1%	9.2%	4.2%	3.8%

^aArticles citing an intermediary study instead of the original source of the data.

^bArticles containing URLs that do not work or are not up-to-date.

^cArticles without proper references for the used datasets.

^dArticles using data from computational simulations and not from biological experiments.

^eArticles using their own datasets.

Citação indireta

4. Datasets

To demonstrate the efficiency of our method, the proposed method is evaluated on five gene expression datasets used in the literature.

(1) Colon dataset: colon dataset [15] is derived from colon cancer patient samples. It consists of the expression levels of 1909 genes of 62 patients among which 40 are colon cancer cases and 22 are normal cases.

(2) Ovarian dataset: ovarian dataset [16] often serves as benchmark for microarray data analysis in most literatures. The dataset provided here includes 91 controls (normal) and 162 ovarian cancers. There are a total of 15,154 genes.

(3) CNS dataset: CNS dataset [15] contains 60 patient samples out of which 39 are normal cases and 21 are cancer cases. There are 7129 genes in the dataset.

(4) Leukemia dataset: leukemia dataset [6] is another dataset widely used in the literature, which is taken as our benchmark dataset. The leukemia dataset contains the expression levels of 7129 genes taken from 72 samples. Labels indicate that there are 47 cancer cases and 25 normal cases.

(5) Breast dataset: breast dataset [17] consists of 97 samples, of which 51 are normal cases and 46 are cancer cases. Large number of features against small sample size is the trademark of the breast cancer dataset. Table 2 shows the summary of these datasets.

5. Guyon I., Weston J., Barnhill S., and Vapnik V., Gene selection for cancer classification using support vector machines, *Machine Learning*. (2002) 46, no. 1–3, 389–422, <https://doi.org/10.1023/A:1012487302797>, 2-s2.0-0036161259.

 [View](#) | [Web of Science®](#)

[Google Scholar](#)

6. Furey T., Cristianini N., Duffy N., Bednarski D. W., Schummer M., and Haussler D., Support vector machine classification and validation of cancer tissue samples using microarray expression data, *Bioinformatics*. (2000) 16, no. 10, 906–914, <https://doi.org/10.1093/bioinformatics/16.10.906>, 2-s2.0-0033636139.

 [View](#) | [CAS](#) | [PubMed](#)

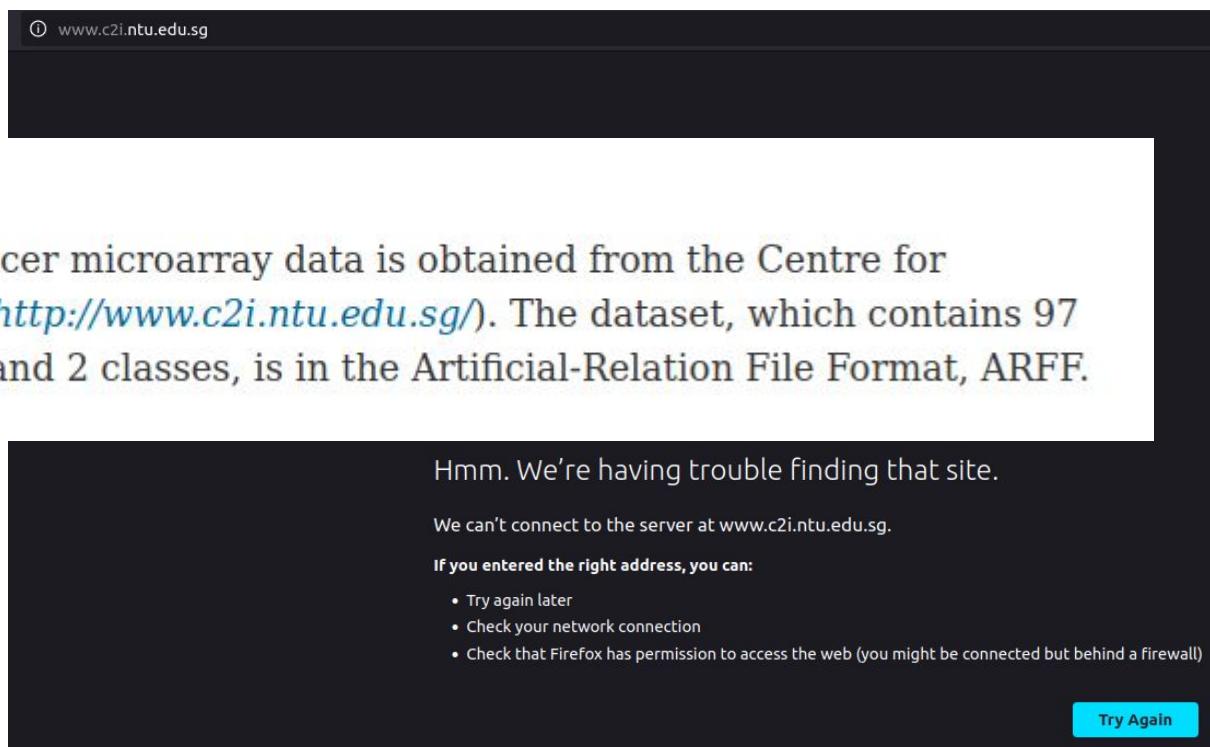
[Web of Science®](#) | [Google Scholar](#)

7. Statnikov A., Aliferis C. F., Tsamardinos I., Hardin D., and Levy S., A comprehensive evaluation of multiclassification methods for microarray gene expression cancer diagnosis, *Bioinformatics*. (2005) 21, no. 5, 621–629, <https://doi.org/10.1093/bioinformatics/bti002>, 2-s2.0-1584436212.

Links quebrados

A. Data Acquisition

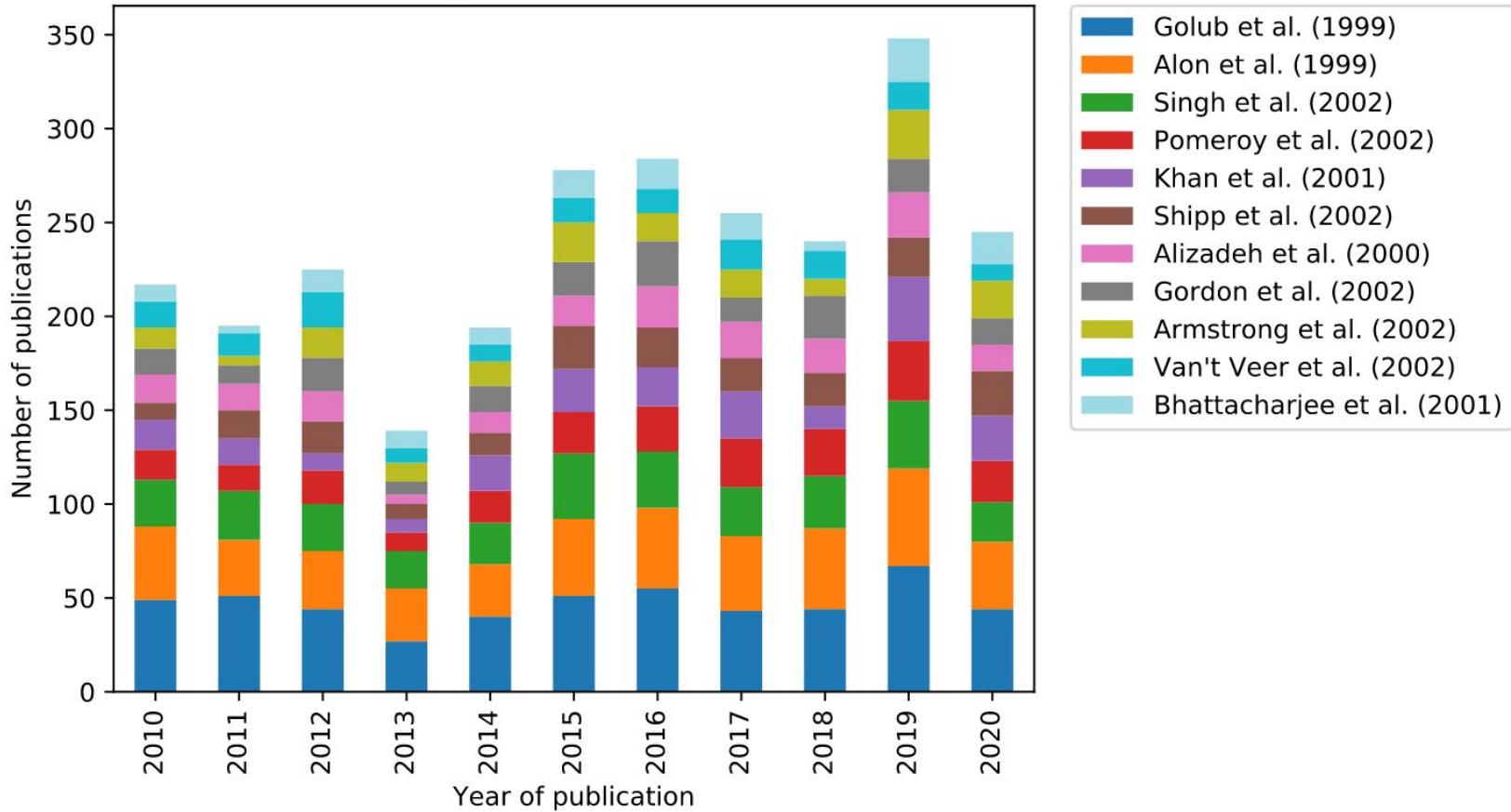
For this study, the breast cancer microarray data is obtained from the Centre for Computational Intelligence (<http://www.c2i.ntu.edu.sg/>). The dataset, which contains 97 instances with 24461 genes and 2 classes, is in the Artificial-Relation File Format, ARFF.

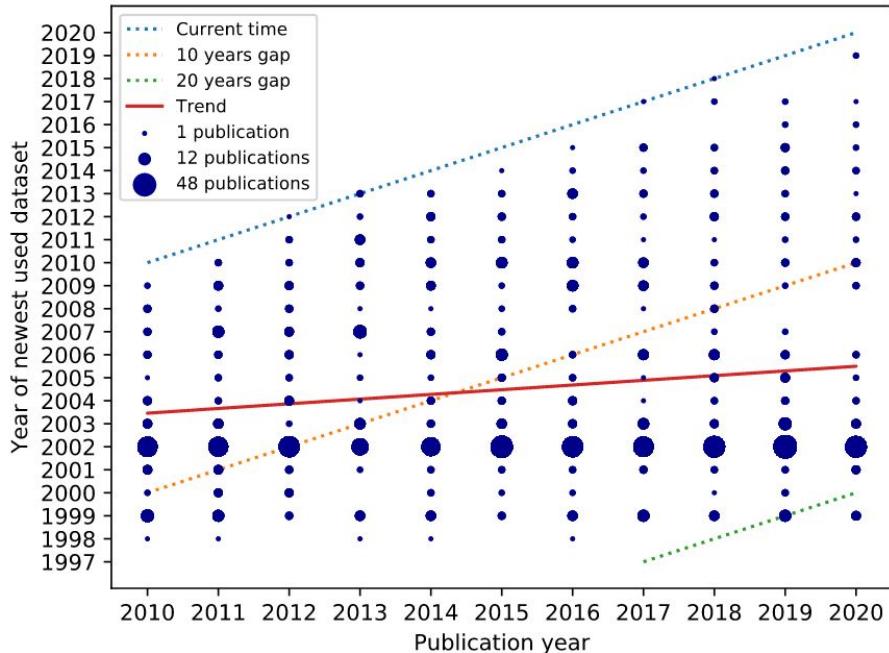


Ranking	Dataset	Year	Samples	Features	Classes	Background	Prevalence in the reviewed publications (%)
1	Golub et al. (1999)	1999	72	7129	2	Leukemia	40.26
2	Alon et al. (1999)	1999	62	2000	2	Colon cancer	32.13
3	Singh et al. (2002)	2002	136	12,600	2	Prostate cancer	22.98
4	Pomeroy et al. (2002)	2002	90	5920	5	Brain cancer	17.67
5	Khan et al. (2001)	2001	83	2309	4	SRBCT	15.94
6	Shipp et al. (2002)	2002	77	7129	2	DLBCL	14.54
7	Alizadeh et al. (2000)	2000	96	4026	9	Lymphoma	13.60
8	Gordon et al. (2002)	2002	181	12,533	2	Lung cancer	13.52
9	Armstrong et al. (2002)	2002	72	11,225	3	Leukemia	12.58
10	Van't Veer et al. (2002)	2002	97	24,481	2	Breast cancer	11.18
11	Bhattacharjee et al. (2001)	2001	203	12,601	5	Lung cancer	10.39

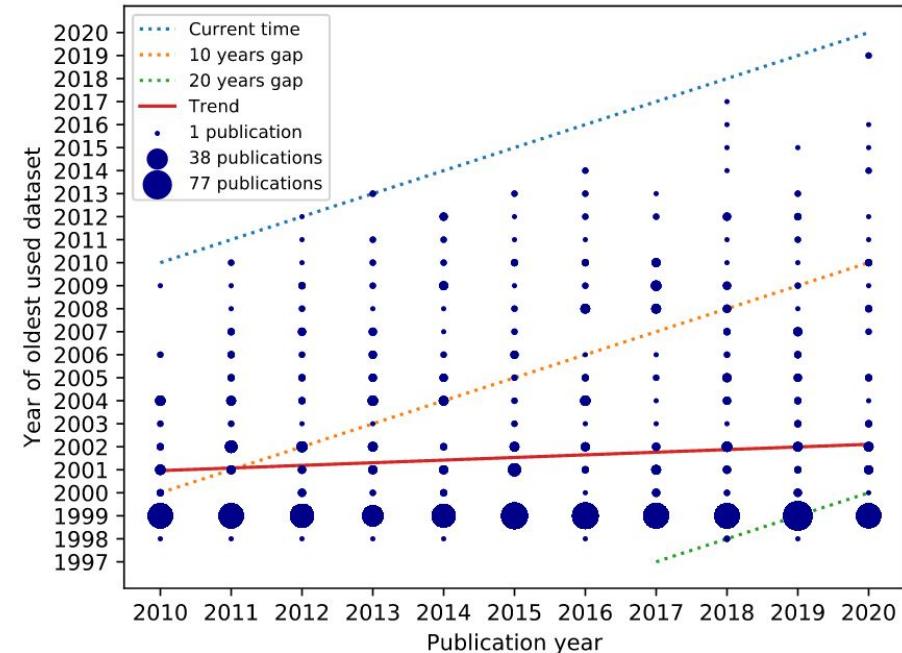
Abbreviations: DLBCL, diffuse large B-cell lymphoma; SRBCT, small round blue cell tumor.

Publications using the eleven most cited datasets

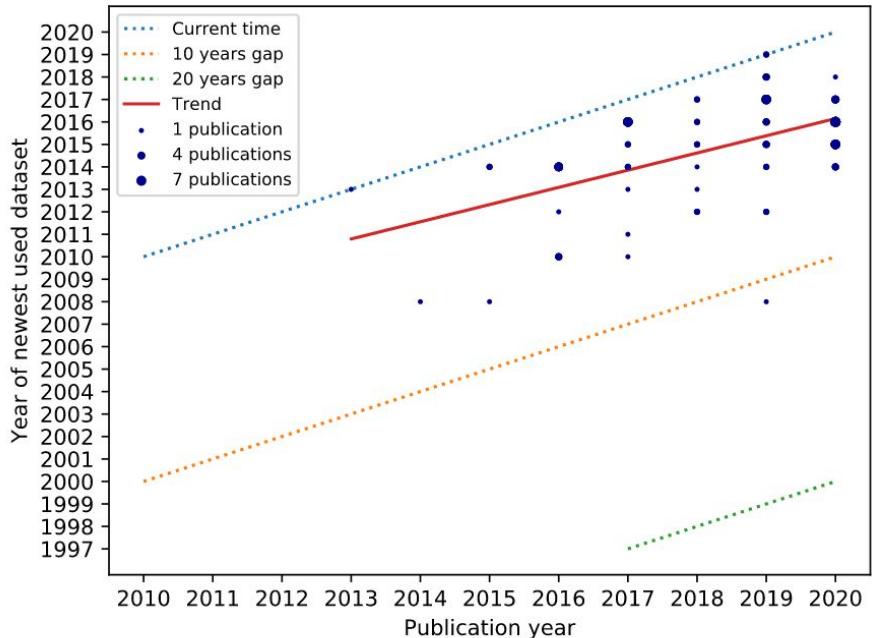




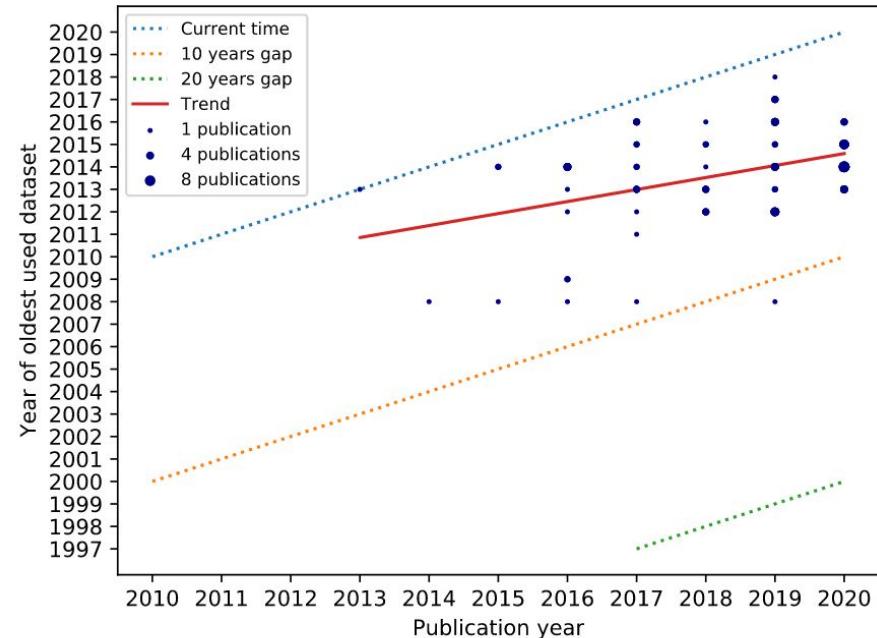
(a) Most recent microarray datasets



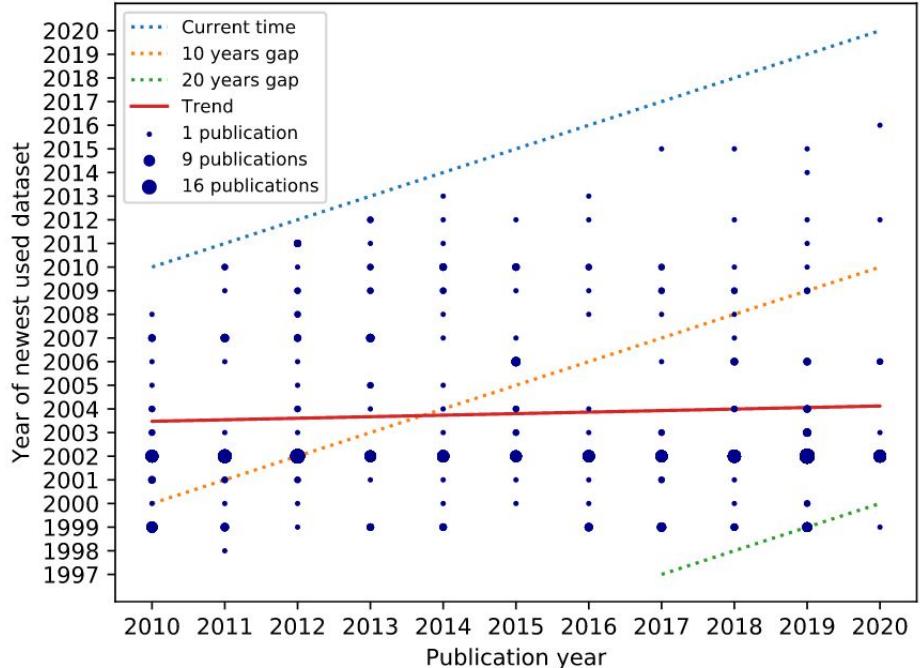
(b) Least recent microarray datasets



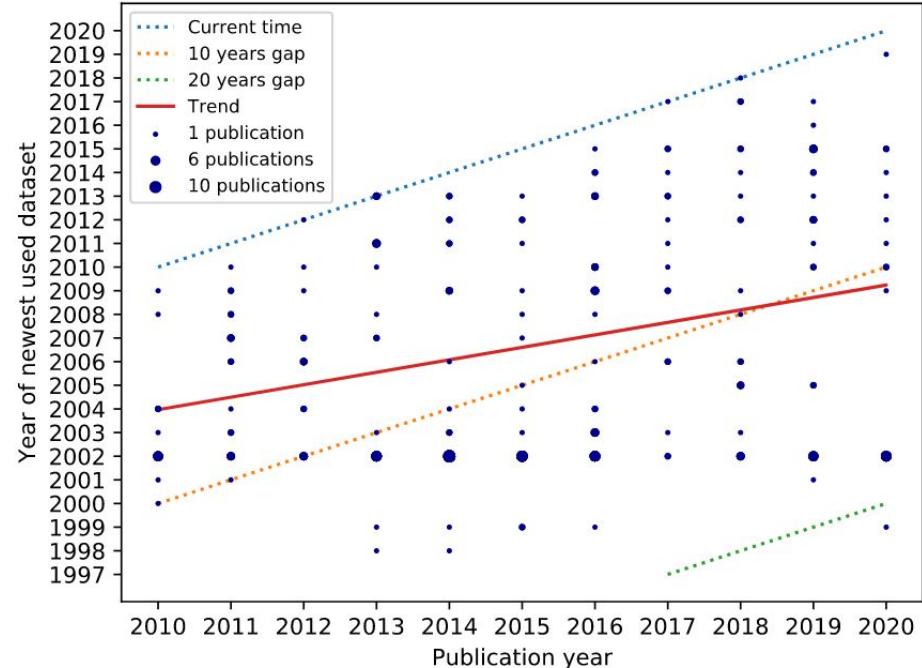
(c) Most recent RNA-seq datasets



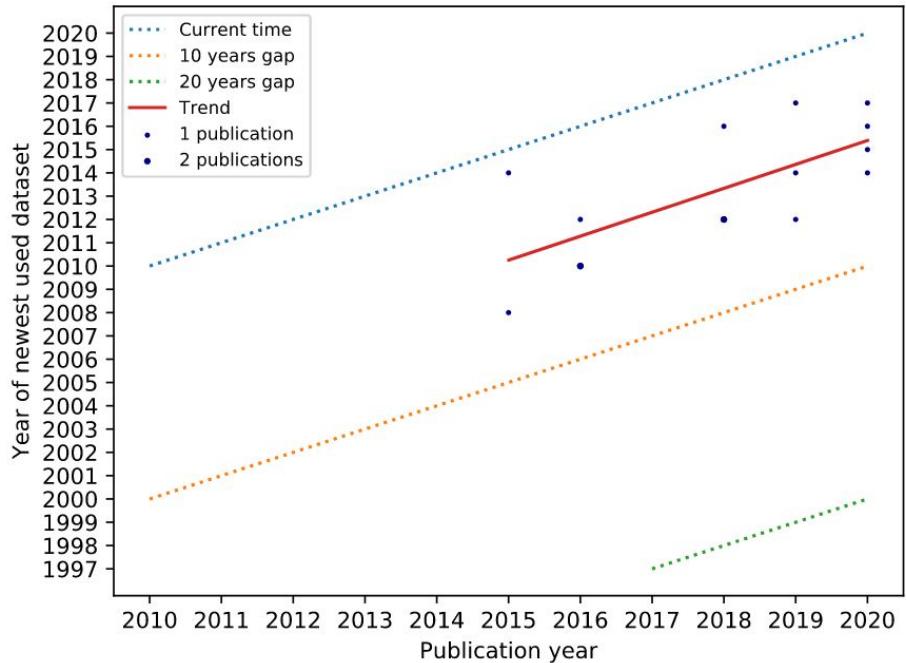
(d) Least recent RNA-seq datasets



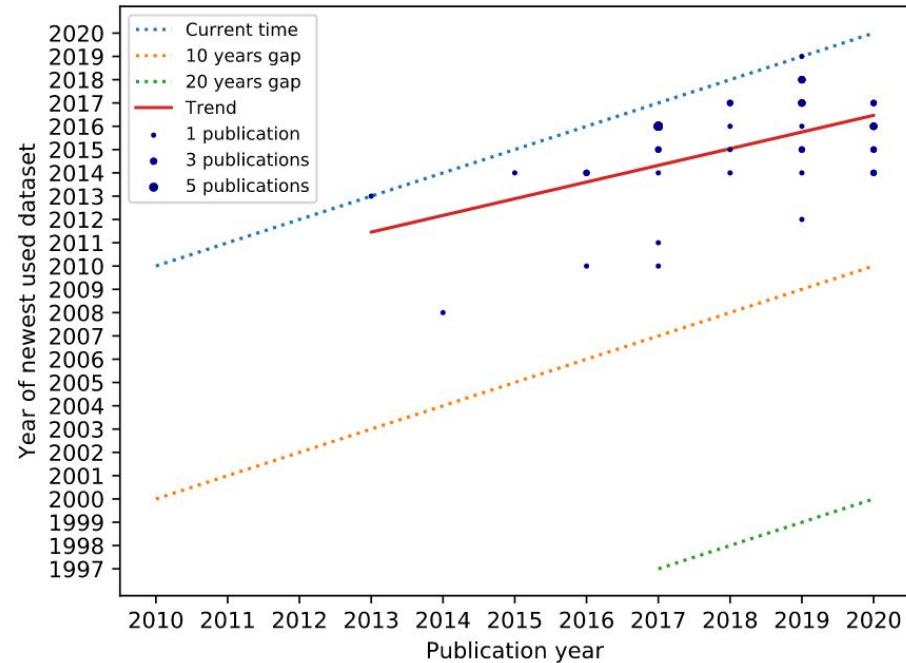
(a) IEEE Xplore — Microarray



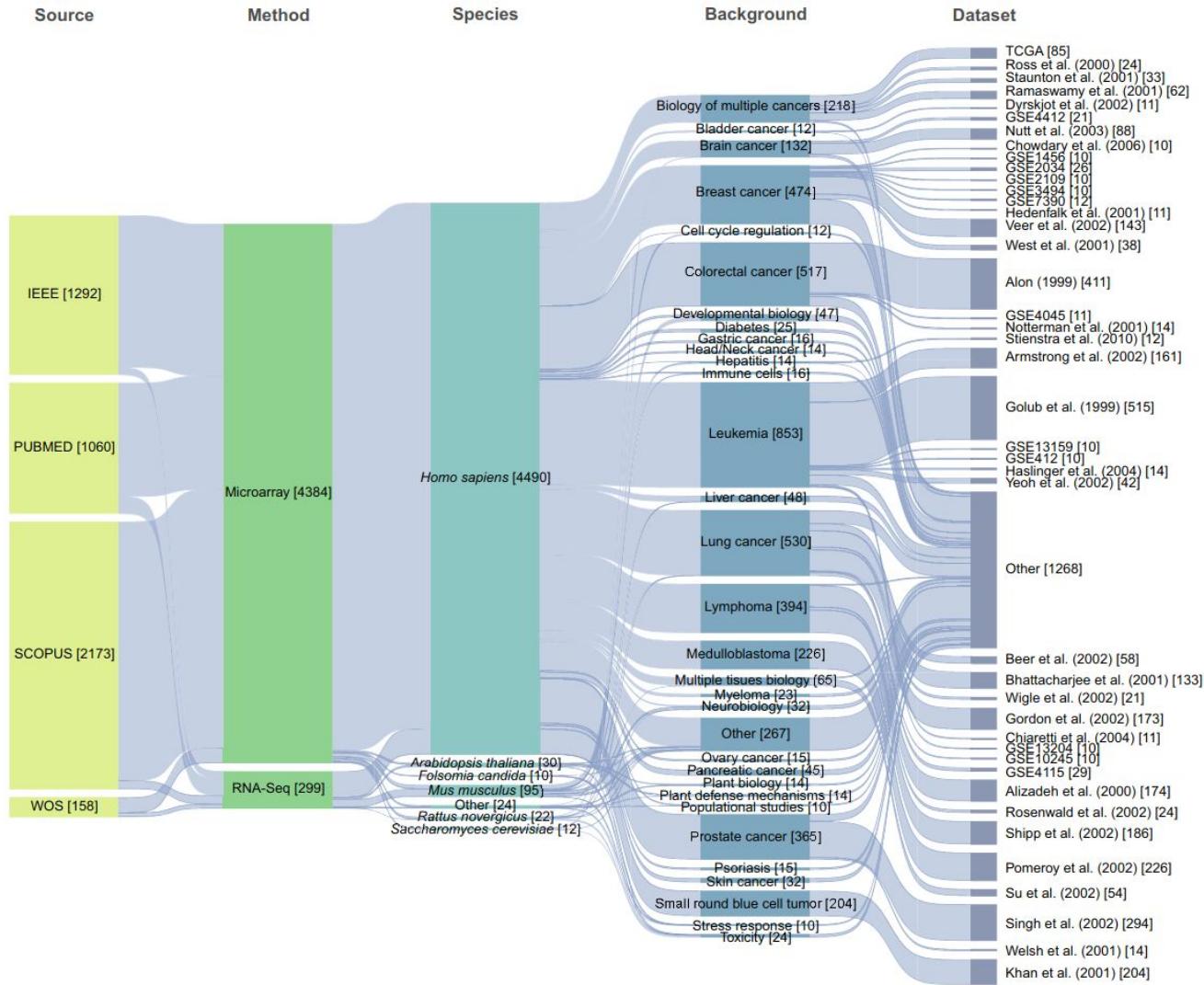
(b) PubMed — Microarray



(c) IEEE Xplore — RNA-seq



(d) PubMed — RNA-seq



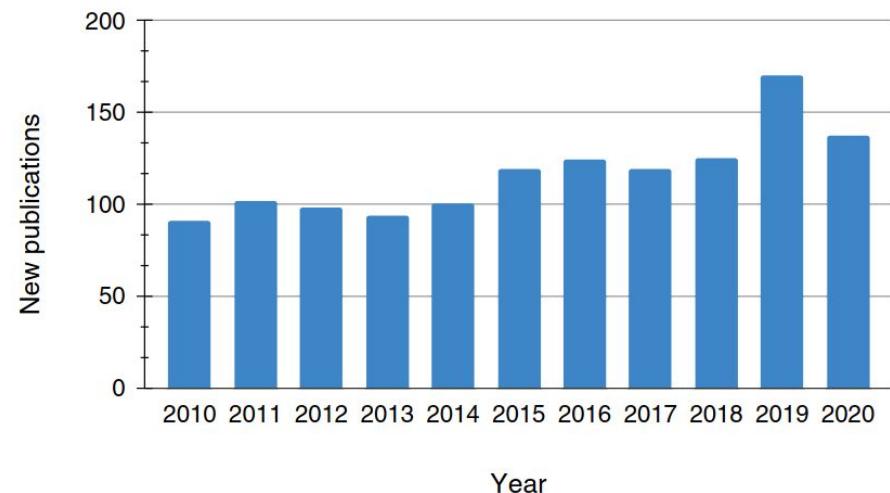


Tem motivo?

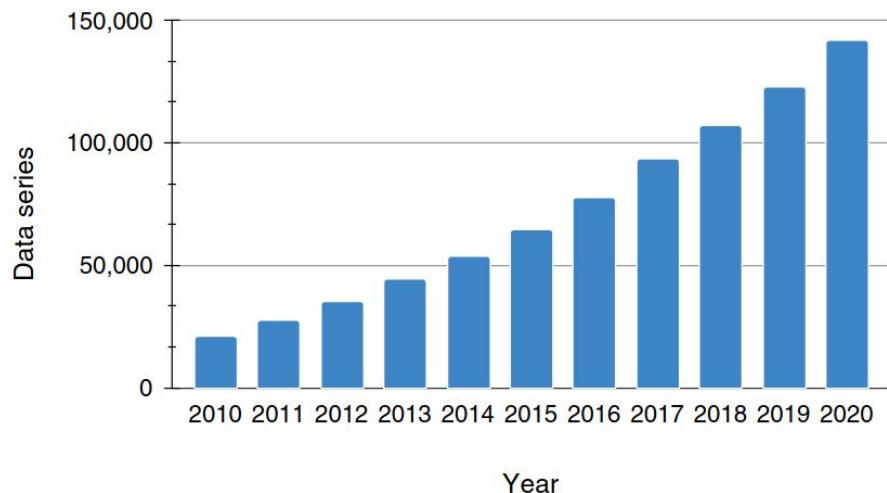


Tem motivo?

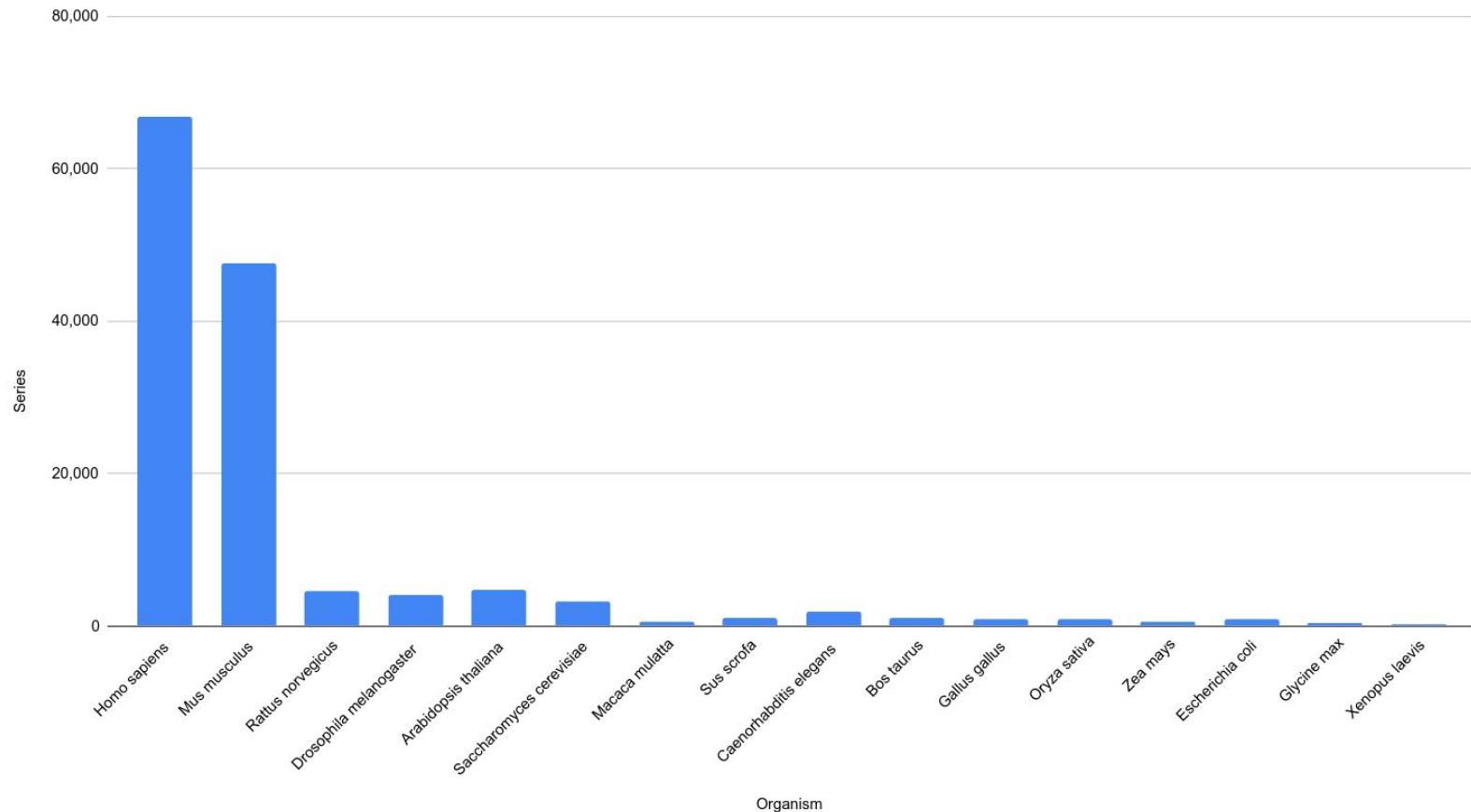
Feature selection publications using gene expression data



Growth of GEO — NCBI



GEO - NCBI: Series versus Organism





E qual o problema?

E qual o problema?

Singh et al. (2002): A distribuição de classes no conjunto de treinamento (49%/51%) é diferente do conjunto de teste (26%/74%). Além disso, este conjunto de teste tem uma diferença quase 10 vezes maior na intensidade do *microarray* em comparação com o conjunto de treinamento.

Gordon et al. (2002): divisão de 50%/50% no treinamento, mas de 90%/10% no teste. Um único gene presente nos dados de treinamento é capaz de classificar corretamente todas as amostras. No entanto, esse mesmo gene não é relevante no conjunto de teste, que, por sua vez, não é linearmente separável.

Bhattacharjee et al. (2001): Mramor et al. (2007) identificaram sete amostras rotuladas incorretamente.

E qual o problema?

Golub et al. (1999): 6817 genes e 38 amostras, sendo pequeno mesmo em comparação com dados dos anos 2000.

Altamente heterogêneo, incluindo amostras de sangue periférico, não apenas de medula óssea, além de pacientes infantis e laboratórios com diferentes protocolos de preparação de amostras.

Usando o método "análise de vizinhança" os autores relataram 100% de precisão na predição de classes. Esses resultados foram altamente insensíveis à seleção particular de genes, com preditores usando de 10 a 200 genes diferentes, todos alcançando a mesma precisão de 100%.



E qual o problema?

Alon et al. (1999) e Khan et al. (2001): experimentos de microarranjo com 6567 genes.

Os autores filtraram o número de genes para 2000 e 2308.

As versões desses conjuntos de dados empregadas por outros pesquisadores foram suas variações filtradas, não os dados brutos originais

Usar dados pré-filtrados para avaliar os selecionadores de atributos pode enviesar os resultados.

Os autores podem estar inadvertidamente misturando os resultados de um algoritmo de filtro de terceiros com os seus próprios.

E qual o problema?

Pomeroy et al. (2002) e **Shipp et al. (2002)** usaram uma plataforma de microarranjo com apenas 6817 sondas, anterior ao fim do HGP e três vezes menor do que o frequentemente usado Affymetrix U133 GeneChip.

Os links para os resultados originais de **Pomeroy et al. (2002)** e o protocolo detalhado de análise de expressão gênica de **Armstrong et al. (2002)** foram armazenados em sites privados que não estão mais funcionais, um problema comum com conjuntos de dados publicados antes do GEO.

Van't Veer et al. (2002): informação das sequências de cRNA para microarranjos, que foram extraídas antes do fim do HGP, contendo, assim, vários vieses pelos padrões atuais.



Possíveis causas?



Possíveis causas

Falta de familiaridade com os dados

Comparação com resultados anteriores

Facilidade de acesso

Bases de dados

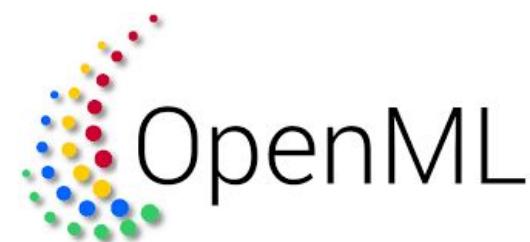
Ao longo dos anos, grandes bancos de dados de *machine learning* se tornaram disponíveis.

Esses bancos de dados contêm centenas a milhares de conjuntos de dados de vários domínios.

Objetivo principal dessas plataformas: facilitar a distribuição e acesso de conjuntos de dados.

Foco em experimentos e *benchmarks* de *machine learning*.

Podem amplificar certos vieses ao promover conjuntos de dados já populares.



Gene expression dataset (Golub et al.)

Data Code (83) Discussion (4) Metadata

▲ 246

New Notebook

Download (1 MiB)



Activity Overview

ACTIVITY STATS

VIEWS

104077

DOWNLOADS

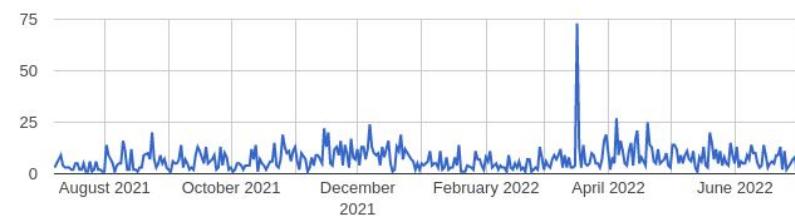
12399

DOWNLOAD PER VIEW RATIO

0.12

77

Downloads ▾



NOTEBOOKS STATS

NOTEBOOKS

83

NOTEBOOK COMMENTS

104

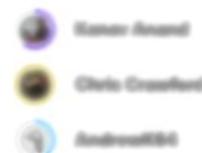
UPVOTE PER NOTEBOOK RATIO

3.96

NOTEBOOK UPVOTES

329

TOP CONTRIBUTORS



DISCUSSION STATS

TOPICS

4

TOTAL COMMENTS

9

UPVOTE PER POST RATIO

2.56

DISCUSSION UPVOTES

23

Bases de dados

Outros bancos de dados contêm dados biológicos específicos.

Voltados principalmente para biologia.

Organização, jargão e formatos de arquivo podem afastar usuários de outras áreas.

Alguns bancos de dados oferecem conjuntos de dados de expressão gênica curados para uso em pesquisas de aprendizado de máquina.



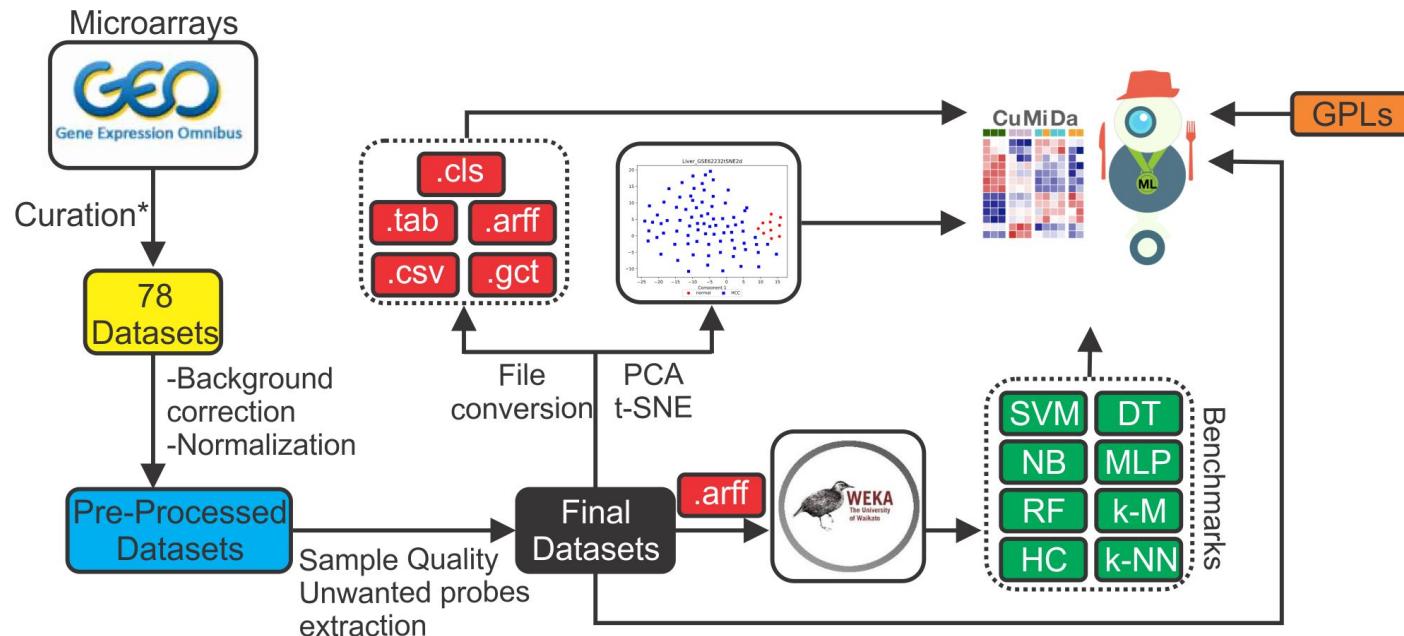
TABLE 3 The most famous public microarray data repositories according to the review by Bolón-Canedo et al. (2014).

Repositories	URL	Status
ArrayExpress European Bioinformatics Institute	http://www.ebi.ac.uk/arrayexpress/	✓
Cancer Program Data Sets Broad Institute	http://www.broadinstitute.org/cgi-bin/cancer/datasets.cgi	✗
Dataset Repository Bioinformatics Research Group of Universidad Pablo de Olavide	http://www.upo.es/eps/bigs/datasets.html	✗
Feature Selection Datasets Arizona State University	http://featureselection.asu.edu/datasets.php	✗
Gene Expression Model Selector Vanderbilt University	http://www.gems-system.org	✗
Gene Expression Omnibus National Institutes of Health	http://www.ncbi.nlm.nih.gov/geo/	✓
Gene Expression Project Princeton University	http://genomics-pubs.princeton.edu/oncology/	✓
Kent Ridge Bio-Medical Dataset Repository Agency for Science, Technology and Research	http://datam.i2r.a-star.edu.sg/datasets/krbd	✗
Stanford Microarray Database Stanford University	http://smd.stanford.edu/	✗

Note: The sources and the provided URLs are listed as informed by Bolón-Canedo et al. (2014), which accessed all repositories in January 2014. In less than a decade, most URLs are not working anymore (marked with 55). The current status of the URLs was checked on October 14, 2023.

Databases	Curated	Source	Quality control ^a	Up to date ^b	File formats ^c
ARCH4	No	Normalized; gene and transcript level counts	No	Yes	.h5
BARRA:CuRDa	Yes	Normalized	Yes	Yes	.csv; .tab; .gct; .cls
BioLab	No	Author's	No	No	.tab
CuMiDa	Yes	Normalized	Yes	Yes	.csv; .tab; .gct; .cls; .arff
Datamicroarray	No	Author's	No	No	.r; .RData
Gene Expression Project	No	Author's	No	No	.tab; .xls
InSilicodb	Yes	Varies	NS	Yes	.r
PSO-EMT datasets	No	Varies	No	No	.mat
Recount3	Yes	Gene and transcript level counts	No	Yes	.txt; .bw; .mtx; .RData
Refine.bio	No	Normalized; gene and transcript level counts	No	Yes	.sf; .json; .tsv
RNASeq-er	No	Normalized; gene and transcript level counts	No	Yes	.cram; .bw; .bedGraph

CuMiDa



CuMiDa: An Extensively Curated Microarray Database for Benchmarking and Testing of Machine Learning Approaches in Cancer Research

BRUNO CÉSAR FELTES, EDUARDO BASSANI CHANDELIER,
BRUNO IOCHINS GRISCI, and MÁRCIO DORN

ABSTRACT

The employment of machine learning (ML) approaches to extract gene expression information from microarray studies has increased in the past years, specially on cancer-related works. However, despite this continuous interest in applying ML in cancer biomedical research, there are no curated repositories focused only on providing quality data sets exclusively for benchmarking and testing of such techniques for cancer research. Thus, in this work, we present the *Curated Microarray Database* (CuMiDa), a database composed of 78 handpicked microarray data sets for *Homo sapiens* that were carefully examined from more than 30,000 microarray experiments from the *Gene Expression Omnibus* using a rigorous filtering criteria. All data sets were individually submitted to background correction, normalization, sample quality analysis and were manually edited to eliminate erroneous probes. All data sets were tested using principal component analysis (PCA) and t-distributed stochastic neighbor embedding (t-SNE) analyses to observe sample division and were additionally tested using various ML approaches to provide a base accuracy for the major techniques employed for microarray data sets. CuMiDa is a database created solely for benchmarking and testing of ML approaches applied to cancer research.

Keywords: benchmarking, cancer, classification, curation, machine learning, microarray, supervised learning, unsupervised learning.

<https://doi.org/10.1089/cmb.2018.0238>

Dataset

Cancer Type: All

GSE:

Sort by Samples Sort by Genes Sort by Classes

TYPE	GSE	GPL PLATFORM	SAMPLES	GENES	CLASSES	Download		
Pancreatic	16515	570	51	54676	2			
ZEROR	SVM	MLP	DT	NB	RF	HC	KNN	K-MEANS
0.71	0.86	0.78	0.78	0.84	0.82	0.69	0.76	0.76
TYPE	GSE	GPL PLATFORM	SAMPLES	GENES	CLASSES	Download		
Breast	33447	14550	16	36623	2			
ZEROR	SVM	MLP	DT	NB	RF	HC	KNN	K-MEANS
0.44	1	0.88	0.88	0.88	0.94	0.56	0.88	0.88
TYPE	GSE	GPL PLATFORM	SAMPLES	GENES	CLASSES	Download		
Breast	59246	13607	101	36623	2			
ZEROR	SVM	MLP	DT	NB	RF	HC	KNN	K-MEANS

<https://sbc.inf.ufrgs.br/cumida>



CuMiDa

CUMIDA

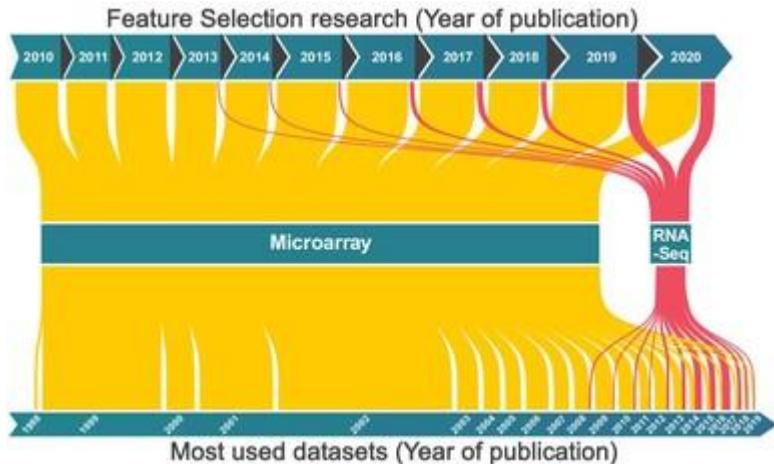
9

TABLE 3. MAJOR CRITERIA DESIGNED TO EVALUATE A PROPER BENCHMARK

<i>Criteria</i>	<i>CuMiDa</i>
Active benchmarks in research area	Results obtained from the top, most used ML tools in the field
Trustworthy evaluated tools	All tools available at WEKA software
Transparency with conducted protocols	Full description: from data mining filtering criteria to classification protocol
Availability of benchmarked in/outputs	All inputs and outputs fully available at the database
Relevance of employed metrics	All metrics are the most used in the field of cancer microarray analysis
Availability of benchmarked methods	All methods can be found in the WEKA software
No inclusion of new tools	All benchmark results are from traditional methods

ML, machine learning.

Artigo final



The use of gene expression datasets in feature selection research: 20 years of inherent bias?

Bruno I. Grisci^{1,2} | Bruno César Feltes³ | Joice de Faria Poloni³ |
 Pedro H. Narloch¹ | Márcio Dorn^{1,4,5}

¹Institute of Informatics, Federal University of Rio Grande do Sul, Porto Alegre, Rio Grande do Sul, Brazil

²Faculty of Computer Science, Dalhousie University, Halifax, Nova Scotia, Canada

³Institute of Biosciences, Federal University of Rio Grande do Sul, Porto Alegre, Rio Grande do Sul, Brazil

⁴National Institute of Science and Technology - Forensic Science, Porto Alegre, Rio Grande do Sul, Brazil

⁵Center for Biotechnology, Federal University of Rio Grande do Sul, Porto Alegre, Rio Grande do Sul, Brazil

Abstract

Feature selection algorithms are frequently employed in preprocessing machine learning pipelines applied to biological data to identify relevant features. The use of feature selection in gene expression studies began at the end of the 1990s with the analysis of human cancer microarray datasets. Since then, gene expression technology has been perfected, the Human Genome Project has been completed, new microarray platforms have been created and discontinued, and RNA-seq has gradually replaced microarrays. However, most feature selection methods in the last two decades were designed, evaluated, and validated on the same datasets from the microarray technology's infancy. In this review of over 1200 publications regarding feature selection and gene expression, published between 2010 and 2020, we found that 57% of the publications used at least one outdated dataset, 23% used only outdated data, and 32% did not cite data sources. Other issues include referencing databases that are no longer available, the slow adoption of RNA-seq datasets, and bias toward human cancer data, even for methods designed for a broader scope. In the most popular datasets, some being 23 years old, mislabeled samples, experimental biases, distribution shifts, and the absence of classification challenges are common. These problems are more predominant in publications with computer science backgrounds compared to publications from biology and can lead to inaccurate and misleading biological results.

This article is categorized under:

Algorithmic Development > Biological Data Mining
 Technologies > Machine Learning

KEY WORDS

feature selection, gene expression, machine learning, microarray, RNA-seq



Problema resolvido?

Recomendações

- Citar a fonte original dos dados.
- Mencionar as principais características dos conjuntos de dados, como o número de amostras, atributos e classes.
- Se vários conjuntos de dados forem usados, listar de forma organizada.
- Adicionar *hyperlinks* funcionais para a fonte dos dados. Se os dados estiverem em bancos de dados privados, verificar se não estão disponíveis em repositórios públicos.
- Usar vários conjuntos de dados com propriedades distintas para validação.

Recomendações

- Não usar apenas dados de *H. sapiens* e de câncer.
- Não usar conjuntos de dados pré-filtrados pelo autor.
- Evitar usar conjuntos de dados antigos.
- Usar dados de RNA-seq.
- Ao projetar um novo banco de dados de expressão gênica para aprendizado de máquina, seguir as diretrizes propostas por **Peters et al. (2018)**, **Wilkinson et al. (2016)**, **Walsh et al. (2021)** e **Hutchinson et al. (2021)**.
- Revisores e editores devem garantir a aplicação dos itens acima nas publicações.
- Usar métodos interpretáveis de aprendizado de máquina e visualização para melhorar a compreensão e a replicabilidade dos resultados de experimentos.