Machine learning Lecture 1



Marcin Wolter IFJ PAN

27 February 2017

- Machine learning: what does it mean?
- What is a difference between "machine learning" and "bayesian learning"?
- A little bit of mathematics and examples of simple linear classifiers.
- Software to work with and literature.

27.02.2017

Recommended books

- M. Krzyśko, Systemy uczące się: rozpoznawanie wzorców, analiza skupień i redukcja wymiarowości. WNT, 2008.
- C. Bishop, *Pattern recognition and machine learning*. Springer, 2009.

and maybe my thesis:

 M. Wolter, Metody analizy wielu zmiennych w fizyce wysokich energii https://www.epnp.pl/ebook/metody_analizy_wielu_zmiennych_w_fizyce_wysokich_energii

Programs

• TMVA – integrated with the ROOT package

http://tmva.sf.net

Installs together with root

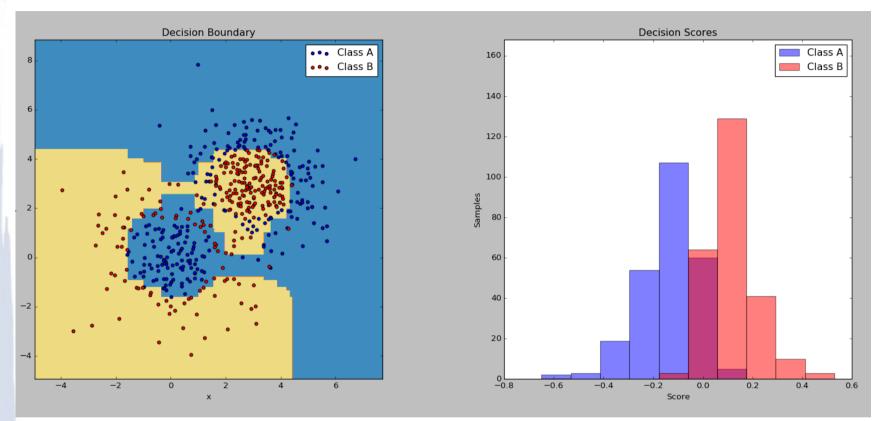
Very popular at CERN

http://scikit-learn.org

scikit-learn - Machine Learning in Python

Simple and efficient tools for data mining and data analysis Accessible to everybody, and reusable in various contexts Built on NumPy, SciPy, and matplotlib I have never used scikit, but we can learn together!

Scikit-learn example



- BDT AdaBoost separation example
- http://scikit-learn.org/stable/auto_examples/ensemble/plot_adaboost_ twoclass.html#sphx-glr-auto-examples-ensemble-plot-adaboost-twoclass -py

Computing

- https://www.cloud.ifj.edu.pl/
- Register, you can create your virtual linux box and play with it.

CC1 • Machines • Farr Home		•
CC1	Welcome to CC1!	
News	System documentation can be found at project's site. In case you encounter any troubles while using the system, please visit the Help section. Resource usage: Cpu: 4/14 Memory: 6.5 GB/28.8 GB () Storage: 126.7 GB/175.8 GB () Points: 5657/8000 ()	
27.02.2017	Machine Learning, M. Wolter	5

Statistics

- Statistics describes random events.
- First works أبو يوسف يعقوب بن إسحاق لسبّاح لكندي Al-Kindi (801-873) used statistical methods to break the Ceasar cipher by investigating the frequency in which particular letters appear.



دا معالده ما دالدم و محف ماطلومالغن احدة مرد الاالدي مردم ما العرف من ماله الد معاد معد مقط من منه محمد من ما مها مرا الموالق و معن طل High adapt aller is being hours walk ... make المدي معيد سرد على السالم لم يصعول بلد السابعد والدوسل للحرو - مراكحها كمرمادد والداد ادلية وكرد وابلسه والركد واسع السرع للساء العسا وجالوا مدسم المهادي ويلوا علوو والرحم والتعتع وحن وللم مراجعا و والاللي و

فرااداد والحداله ردالعالم وصلوا يساعلم مدمحة والمسه م

لسمالد الح الاسفالعدد والمحدالين فاسلواه يدفع لمامانه زيسه وكالمحارور لمركد وحد مزاهدان فالحوله الدرص المحالية ومندا يعجر الندجيم الدومية ومساير الفتدا اجمنه ودادان ومعالمان وعرواكما العكر والصل لما ال

Caesar cipher

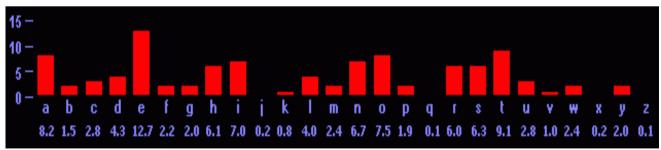
• Each letter of a text is replaced by another, shifted by n letters:

Α	В	С	D	Е
С	D	Е	F	G

 In general Al-Kindi's method can be used to break any replacement cipher (each letter is replaced by another letter, always the same)

А	В	С	D	Ε
Ζ	D	Ρ	G	Т

 Frequency analysis – the frequency of appearance of different letters is investigated.



• Frequency of different letters in English.

27.02.2017

Statistics

 Important step in the development of statistics were the first studies of demography and the games of chance (1663 John Graunt "Natural and Political Observations upon the Bills of Mortality").

Natural and Political OBSERVATIONS

Mentioned in a following INDEX, and made upon the Bills of Mortality.

B Y Capt. **70HN GRAUNT**, Fellow of the Royal Society.

With reference to the Government, Religion, Trade, Growth, dir, Difeafes, and the feveral Changes of the faid CITY.

Contentus paucis Le Et oribus.

The Fifth Edition, much Enlarged.

LONDON,

Printed by John Martyn, Printer to the Royal Society, at the Sign of the Bell in St. Paul's Church-yard, MDCLXXVI.

Joe. Real Sona

							Th	e T	able	e of	C	4	se	11	LT	11	S.						1620	1624	1648	1651	1050	1619	10:
The Years of our Loid	11547	1648	1649	1650	1651	1652	1653	1654	11055	1656	1657	1635	1659	1660	1629	1 630	1631	1632	1633	1634	163	11636	1631	1635 1635	1649	1051	1657	1649	Year
rite and Still-born																											1812		
and Fever	1160	884	7.51	650	10:5	1111	202	1371	0.89	875	1 999	1800	2202	3148	450.	1091	1115	1108	953	1279	1022	+100	4410	9.4330	10010	192035	4363	40105	237
al ma content	1 14	1	1	100	7	76		Î	4	1	5			5	43	8	10	13	6			4	54	14	5	9	14	10	. A
ing to Hex, Scouring, and Flax	3	170	801	180	841	762	300	385	165	362	362	4	7	251	449	438	352	345	173	512	348	330	14	400 1	11	181	19	17	781
t mi Scalded	3	6	10	5	11	5	5	7	10	5	7	•	6		3	10	7		1	3	4	E E	-2	1	2	4	3	19	11
n,Gingrene and Fiftula	2.6	29	31	19	5.0	53	36	37	73	31	34	:35	63	52	3.0	-14	23	38	\$7	30	24	30	85	8		4.2.8	150		60
r, Sore-mouth and Thruth	66	1.20	100	41	1000	1.2.1	100.00	10.000	1000	3.0.1	3.2.6	100.00	1.00	68 194	Sec.	+ 157	112	171	132	143	163	74	590	608	408	769	161 859 1758 4	490	65 330
insand Infinits .	1309	1254	1005	990	1497	13.80	1020	*3.45	1000	393	1000	1144	030	11231	2392	-17-	1035		1130			1.5	101	85	100	24.	497	247	1250
and Cough	120	1.5	1	1			. 41	36		. 28	30	31	23	24	10	58	51	55	45	54	50	57 5	1578	200 8	5999	91412	140 245 7 1177 1	43	59 4448
allen .	684	491	530	4.93	569	653	606	\$28	701	1037	807	841	742	1031	52	87	18	341	221	380	418	709	01	00	01	0303	9	1	907
ardie Stone	185	434	421	505	444	356	517	704	660	705	. j 6j1	3	645	4 573	235	2.52	279	280	100	250	315	359	C48 1	714	5382	32121	982 13 215 1	20	962 82
ratd Nive drinking	47	40	35	27	49	50	20	34	12	4.9	2	90	. "	40	- 72	22	1		2/	1	32	45	02	100	14 4 97		79	1	18
ruted ted in a Bath	1	17	35	45	24	12	4	21	19	32	30	15	7	18	19	耳	11	15	73	1	-	3	37	24	100	1	. 5		
rg-Schuel) Lord fmill Fox	135	400	1150	184	\$25	3	139	812	1194	813	835	409	1323	354	71	40		531	74	1354	IJ	14	81	40 1	39	14	161 27	2.9	1057
nd dead in the Streets who Fox	15	4 25	15	18	21	20	20	20	29	33	= j	53	51	31					7	17	13	22	22	백	80	5	112	23	3.9
(hted at	4	1	1	P	17	7	5	6 17	8	7	8	13	14	1 4	81	5	32	4	4	5	17.0	0 20 7	14 74 57	50	35	25 59	57.	おりまい	いわ
ri ged, and roade-away themfoly	0 11	11	1	14	17	74 74	15	9	14	16			11	30		8	6	15		-	-	1. 2	9	0	48	14	12	46	22 05 00
d-Ach	51	1	1 21	40	41	41	\$2	71	61	41	40	77	101	75	47	59	35	1 43	35	45	34	1	47	197	100	1	0	10	- 99

How to define the probability?

- In probability theory, the sample space of an experiment or random trial is the set of all possible outcomes or results of that experiment.
- Probability of an event A (frequentist definition) is a limit by N going to infinity of n/N, where n is a number of successes and N is a number of trials:

$$P(A) = \lim_{N \to \infty} \frac{n}{N}$$

What is a probability of getting "6" by throwing a dice? The same as the fraction of "6" results in the infinite number of trials.

The definition comes from a text book of Abraham de Moivre (1667-1754) – a text book on statistics "The Doctrine of Chances" (1718).



DOCTRINE

A Method of Calculating the Probability of Events in Play.

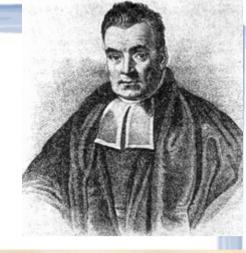


By A. Dr. Mainre. F. R. S.

L O N D O N: Printed by W. Peerfen, for the Author. M DCCXVIII.

Bayes definition

- Probability "a priori", i.e. unconditional, is understood as a measure of belief, based on rational evidence, that such an event will happen.
- In the next step we make an experiment, called "observations", and their results modify the probability. We get the probability "a posteriori", which is a measure of belief modified by the experiment.
- This idea of Thomas Bayes was supported by P. S. Laplace, H. Poincare or the well known economist John Keynes. The stated, that this is a way we recognize and study the nature.



[370] quodque folum, certa nitri figna præbere, fed plura concurrere debere, ut de vero nitro producto dubium non relinquatur.

LII. An Effay towards folving a Problem in the Dostrine of Chances. By the late Rev. Mr. Bayes, F. R. S. communicated by Mr. Price, in a Letter to John Canton, A. M. F. R. S.

Dear Sir,

Read Dec. 23, I Now fend you an effay which I have 1763. I found among the papers of our deceafed friend Mr. Bayes, and which, in my opinion, has great merit, and well deferves to be preferved. Experimental philosophy, you will find, is nearly interefted in the fubject of it; and on this account there feems to be particular reafon for thinking that a communication of it to the Royal Society cannot be improper.

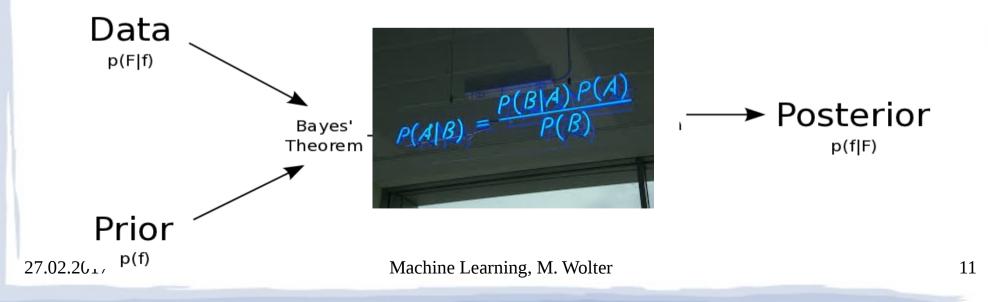
He had, you know, the honour of being a member of that illuftrious Society, and was much efteemed by many in it as a very able mathematician. In an introduction which he has writ to this Eflay, he fays, that his defign at firft in thinking on the fubject of it was, to find out a method by which we might judge concerning the probability that an event has to happen, in given circumftances, upon fuppofition that we know nothing concerning it but that, under the fame circum-

Thomas Bayes (1702 - 1761) was a British mathematician and the Presbyterian pastor. His most important work is "Essay Towards Solving a Problem in the Doctrine of Chances".

27.02.2017

Bayes definition

- The experiment we can't repeat many times: what is a probability to pass an exam?
- Based on our knowledge (we studied text books for few days), we estimate the probability to be: ¹/₂ (probability a priori).
- If all four people trying to pass the exam before us failed, and we know that their knowledge wasn't much different from our, wouldn't we modify our estimation? In this way we get the probability a posteriori.



Bayes Theorem

 Bayes' theorem relates the conditional (posterior) and marginal (prior) probabilities of events A and B:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- **P(A)** is the prior probability or marginal probability of A. It is a "prior" in the sense that it does not take into account any information about B.
- P(A|B) is the conditional probability of A, given B. It is also called the posterior probability because it is derived from or depends upon the specified value of B.
- Intuitively, Bayes' theorem in this form describes the way in which one's beliefs about observing 'A' are updated by having observed 'B'.

Bayes Theorem – an example: cancer test

$$\Pr(A|X) = \frac{\Pr(X|A) \Pr(A)}{\Pr(X)} = \frac{\Pr(X|A) \Pr(A)}{\Pr(X|A) \Pr(A) + \Pr(X|\text{not } A) \Pr(\text{not } A)}$$

- Pr(A|X) = Chance of having cancer (A) given a positive test (X). This is what we want to know: How likely is it to have cancer with a positive result? .
- Pr(X|A) = Chance of a positive test (X) given that you had cancer (A). This
 is the chance of a true positive, 80% in our case.
- Pr(A) = Chance of having cancer (1%).
- Pr(not A) = Chance of not having cancer (99%).
- Pr(X|not A) = Chance of a positive test (X) given that you didn't have cancer (~A). This is a false positive, 9.6% in our case.
- In our case Pr(A|X) is 7.8%

Bayesian vs. Frequentist approach

- PROBABILITY: degree of belief (Bayes, Laplace, Gauss, Jeffreys, de Finetti)
- **PROBABILITY: relative frequency** (Venn, Fisher, Neyman, von Mises).
- Bayesian approach: probability is degree of belief. Thus the probability p is our assessment of the probability of success at each trial, based on our current state of knowledge.
 - If our assessment, initially, is incorrect? As our state of knowledge changes, our assessment of the probability of success changes accordingly.
- Bayesian inference is statistical inference in which evidence or observations are used to update or to newly infer the probability that a hypothesis may be true.
- This allows for a *cleaner* foundation than the frequentist interpretation.

"We don't know all about the world to start with; our knowledge by experience consists simply of a rather scattered lot of sensations, and we cannot get any further without some a priori postulates. My problem is to get these stated as clearly as possible."

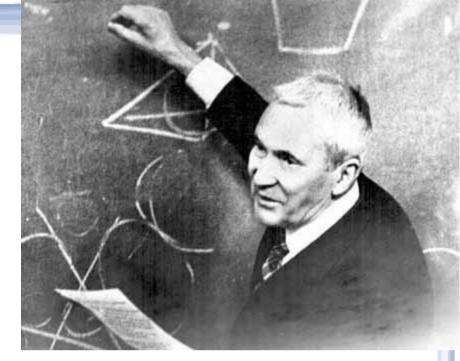
Sir Harold Jeffreys, in a letter to Sir Ronald Fisher dated 1 March, 1934

H.B. Prosper,"Bayesian Analysis", arXiv:hep-ph/0006356v1 30 Jun 2000

27.02.2017

Axiomatic definition

Probability could be defined in many ways, not necessarily like de Moivre...

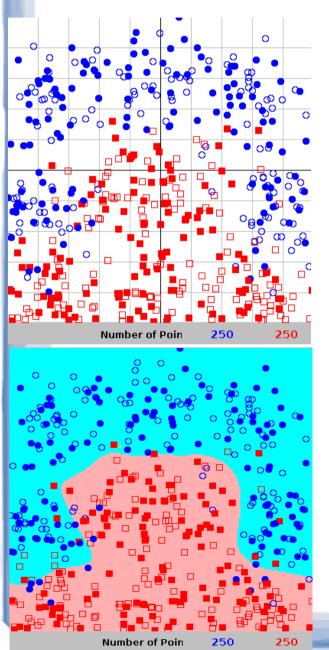


Андре́й Никола́евич Колмого́ров (1903-1987)

Axiomatic definition by Kolmogorov: probability is a function P defined on the space of elementary events, which assigns to each event A a number P(A) such, that:

- $P(A) \ge 0$ for each event A
- P(A)=1 for the sure event A
- $P(A \cup B) = P(A) + P(B)$ when the events A and B are mutually exclusive.

How do the machine learning algorithms work?



- We need training data, for which we know the correct answer, whether it's a signa or background. We divide the data into two samples: training and test.
- We find the best function *f(x)* which describes the probability, that a given event belongs to the class "signal". This is done by minimizing the loss function (for example χ²).
- Different algorithms differ by: the class of function used as *f(x)* (linear, non-linear etc), loss function and the way it's minimized.
- All these algorithms try to approximate the unknown *Bayessian Decisive Function* (BDF) relying on the finit training sample.

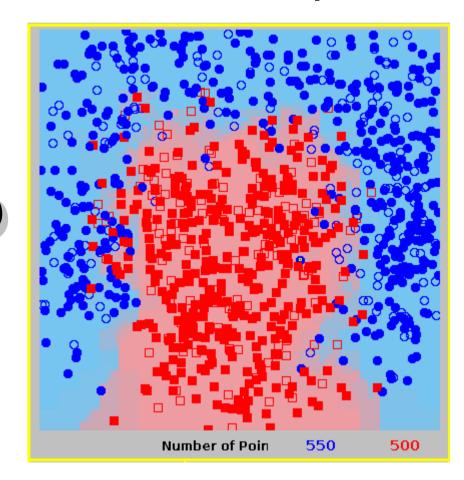
BDF -an ideal classification function given by the unknown probability densities of signal and background.

Cuts vs non-linear separation

Cuts

Number of Poin 400 350

Non-linear separation



Neural Networks, boosted decision trees, itd

27.02.2017

What are ML methods used for?

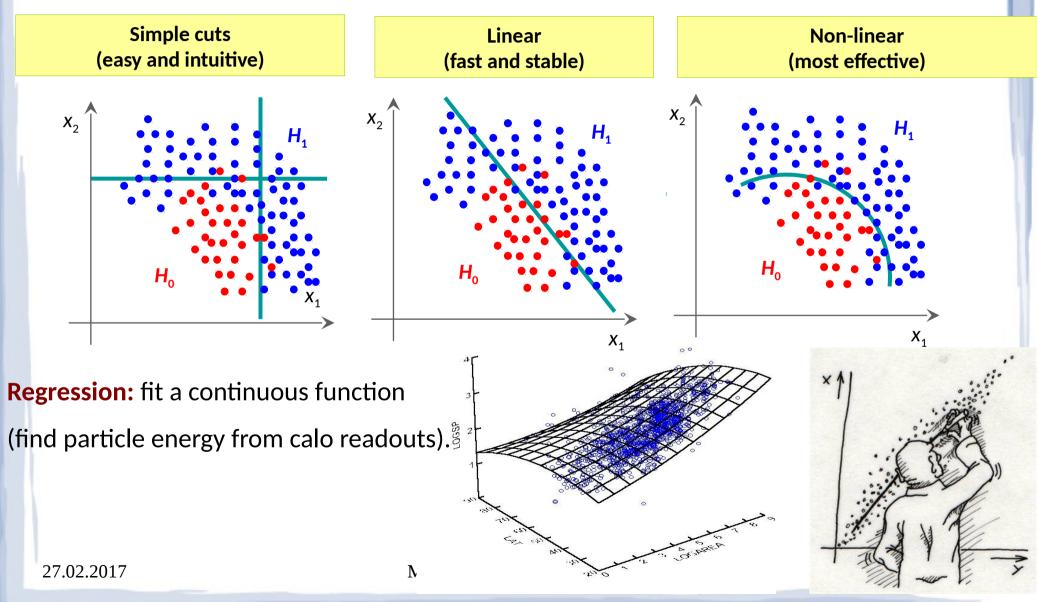
Multivariate methods (machine learning) can be used for:

- Classification (select signal out of background)
- Function approximation (regression) for example find a particle energy using signals from the calorimeter cells.
- Probability density estimation (estimate the probability distribution)
- Variable selection (find most important variables)
- Optimization (tuning, verification)
- Many others...

Types of algorithms

How to use the information available

•**Classification:** find a function *f(x1,x2)* giving the probability, that a given data point belongs to a given class (signal vs background).



Multivariate Methods

Machine Learning

Teach a machine to learn y = f(x) by feeding it training data $T = (x, y) = (x, y)_1, (x, y)_2, \dots, (x, y)_N$ and a constraint on the class of functions.

Bayesian Learning

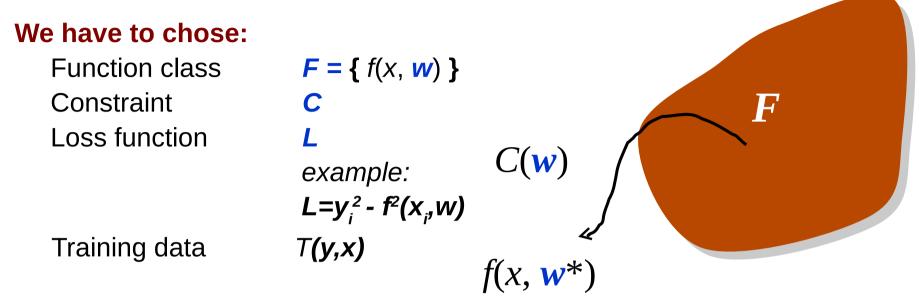
For each function f(x) in the function space F calculate the posterior probability p(f | T) using a given training sample T= (x, y).

Don't use a single function, use a bunch of functions weighted by probabilities

- The posterior probability is the conditional probability that is assigned after the relevant evidence is taken into account.
- Training sample: T = (x, y) set of input vectors x and desired outputs y.

27.02.2017

Machine Learning



Method

Find f(x, w) by minimizing the **empirical risk** *R*

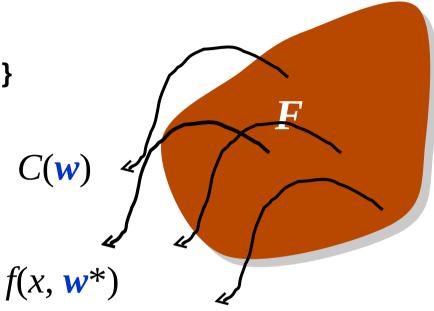
$$R(w) = \frac{1}{N} \sum_{i=1}^{N} L(y_i, f(x_i, w))$$
 subject to the constraint $C(w)$

We choose at the end a single "best" function f(x, w) (best single Neural Network, best likelihood etc.)

Bayesian Learning

Choose

Function class Prior Likelihood Training data $F = \{ f(x, w) \}$ p(w) p(y|x, w) T(y,x)



We do not pick up a single function f(x), instead p(w|T) assigns a probability density to every function in the function class.

Bayesian Learning: Why?

- <u>Probabilistic learning</u>: Calculate explicit probabilities for hypothesis, among the most practical approaches to certain types of learning problems.
- Probabilistic prediction: Predict multiple hypotheses, weighted by their probabilities.
- Incremental: Each training example can incrementally increase or decrease the probability that a hypothesis is correct. Prior knowledge can be combined with observed data.
- <u>Standard</u>: Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured.

Classification

A Bayes classifier:

$$p(S|x) = \frac{p(x|S) p(S)}{p(x|S) p(S) + p(x|B) p(B)}$$

where **S** is associated with y = 1 and **B** with y = 0. Bayes classifier accepts events x if p(S|x) > cut as belonging to **S**.

We need to approximate probability distributions P(x|S) and P(x|B).

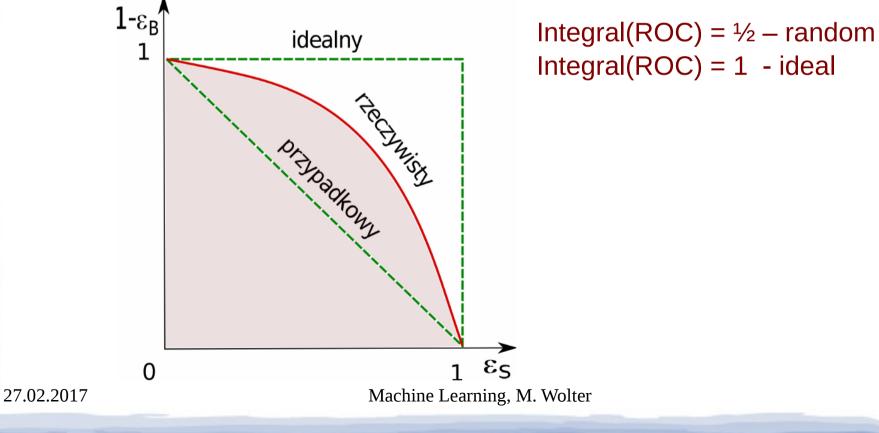
- If your goal is to classify objects with the fewest errors, then the Bayes classifier is the optimal solution.
- Consequently, if you have a classifier known to be close to the Bayes limit, then any other classifier, however sophisticated, can at best be only marginally better than the one you have.

=>If your problem is **linear** you don't gain anything by using sophisticated **Neural Network**

• All classification methods, such as the ones in TMVA, are different numerical approximations of the Bayes classifier.

ROC curve

- ROC (Receiver Operation Characteristic) curve was first used to calibrate radars.
- Shows the background rejection $(1-\epsilon_{\rm B})$ vs signal efficiency $\epsilon_{\rm B}$. Shows how good the classifier is.
- The integral of ROC could be a measure of the classifier quality:



Practical applications

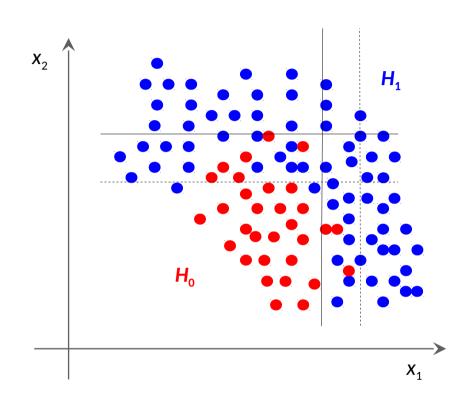
A Short List of Multivariate Methods

- Cuts
- Linear Discriminants (like Fisher)
- Support Vector Machines
- Naive Bayes (Likelihood Discriminant)
- Kernel Density Estimation
- Decision Trees
- Neural Networks
- Bayesian Neural Networks
- Genetic Algorithms
- And many, many others..... I want to present briefly just few of them.

We describe now

- Simple ML linear methods:
 - Cuts
 - Fisher linear discriminant
 - Principal Component Analysis, PCA
 - Independent Component Analysis, ICA

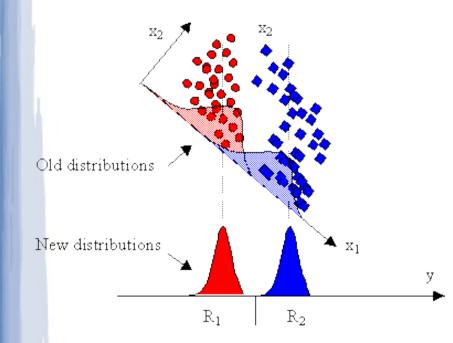
Cuts



- Optimization of cuts:
- Move cuts as long as we get the optimal signal vs. background selection. For a given signal efficiency we find the best background rejection → we get the entire ROC curve.
- Optimization methods:
 - Brute force
 - Genetic algorithms
 - Many others...

Fisher discriminants

Projection to one dimension, than discrimination



Equivalent to linear separation

We choose a projection vector in such a way, that the separation is maximized.

Method introduced by Fisher in 1936. Optimal separation for Gaussian distributions.

27.02.2017

Fisher's linear discriminant

The terms *Fisher's linear discriminant* and *LDA* are often used interchangeably, although <u>Fisher's</u> original article *The Use of Multiple Measures in Taxonomic Problems* (1936) actually describes a slightly different discriminant, which does not make some of the assumptions of LDA such as normally distributed classes or equal class covariances.

Suppose two classes of observations have means $\vec{\mu}_{y=0}, \vec{\mu}_{y=1}$ and covariances $\Sigma_{y=0}, \Sigma_{y=1}$. Then the linear combination of features $\vec{w} \cdot \vec{x}$ will have means $\vec{w} \cdot \vec{\mu}_{y=i}$ and variances $\vec{w}^T \Sigma_{y=i} \vec{w}_{i}$ for i = 0, 1. Fisher defined the separation between these two distributions to be the ratio of the variance between the classes to the variance within the classes:

$$S = \frac{\sigma_{between}^2}{\sigma_{within}^2} = \frac{(\vec{w} \cdot \vec{\mu}_{y=1} - \vec{w} \cdot \vec{\mu}_{y=0})^2}{\vec{w}^T \Sigma_{y=1} \vec{w} + \vec{w}^T \Sigma_{y=0} \vec{w}} = \frac{(\vec{w} \cdot (\vec{\mu}_{y=1} - \vec{\mu}_{y=0}))^2}{\vec{w}^T (\Sigma_{y=0} + \Sigma_{y=1}) \vec{w}}$$

This measure is, in some sense, a measure of the <u>signal-to-noise ratio</u> for the class labelling. It can be shown that the maximum separation occurs when

$$\vec{w} = (\Sigma_{y=0} + \Sigma_{y=1})^{-1} (\vec{\mu}_{y=1} - \vec{\mu}_{y=0})$$

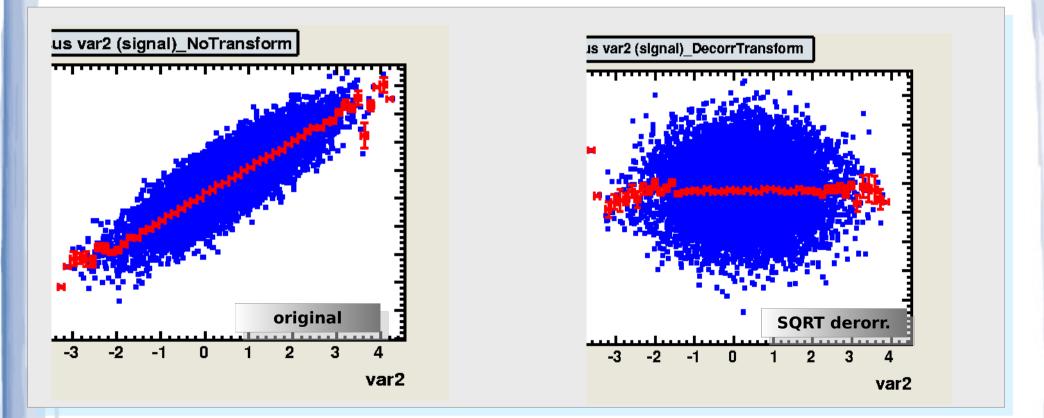
When the assumptions of LDA are satisfied, the above equation is equivalent to LDA.

Be sure to note that the vector \vec{w} is the normal to the discriminant hyperplane. As an example, in a two dimensional problem, the line that best divides the two groups is perpendicular to \vec{w} .

2' Generally, the data points are projected onto \vec{w} . However, to find the actual plane that best separates the 30 data, one must solve for the bias term b in $w^T \mu_1 + b = -(w^T \mu_2 + b)$.

Decorrelation

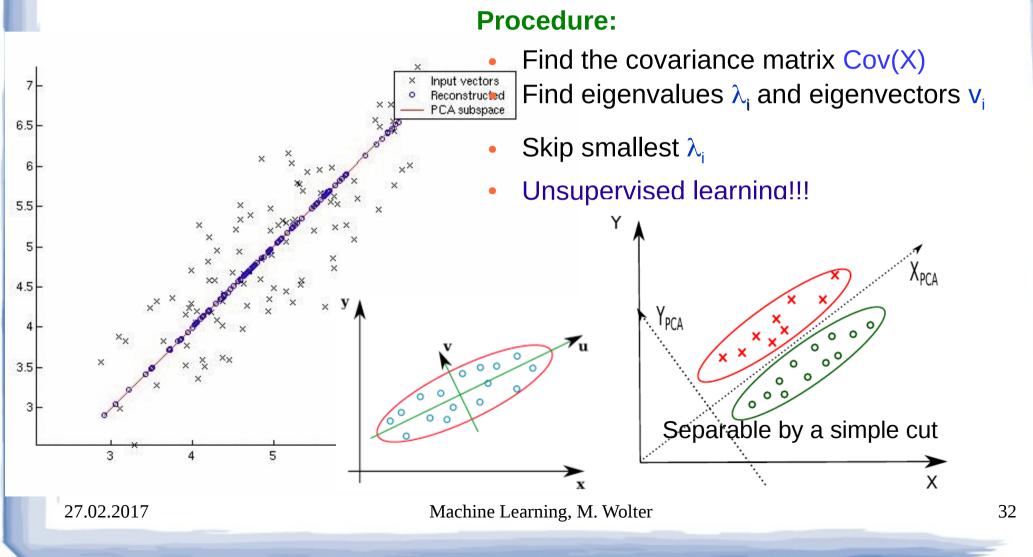
Removes correlation between variables by a rotation in the space of variables



Principal Component Analysis - PCA

Task: reduce the number of dimensions minimizing the loss of information

Finds the orthogonal base of the covariance matrix , the eigenvectors with the smallest eigenvalues migh be skipped

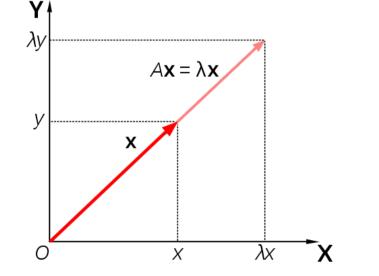


Eigenvalues and eigenvectors

In essence, an eigenvector v of a linear transformation T is a non-zero vector that, when T is applied to it, does not change direction. Applying T to the eigenvector only scales the eigenvector by the scalar value λ , called an eigenvalue. This condition can be written as the equation

$$\mathsf{T}(\mathbf{v}) = \lambda \mathbf{v}$$

referred to as the eigenvalue equation or eigenequation. In general, λ may be any scalar. For example, λ may be negative, in which case the eigenvector reverses direction as part of the scaling, or it may be zero or complex.



Matrix A acts by stretching the vector x, not changing its direction, so x is an eigenvector of A.

27.02.2017

Machine Learning, M. Wolter

Independent Component Analysis ICA

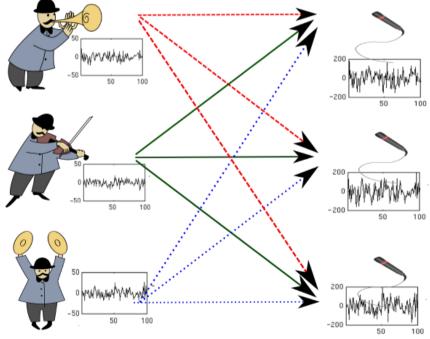
Developed at Helsinki University of Technology http://www.cis.hut.fi/projects/ica/

Problem:

- Assume, that signal X is a linear combination X = AS of independent sources
 S. The mixing matrix A and vector of sources S are unknown.
- Task: find a matrix T (inverted A), such that elements of vector U = TX are statistically independent. T is the matrix returning the original signals.

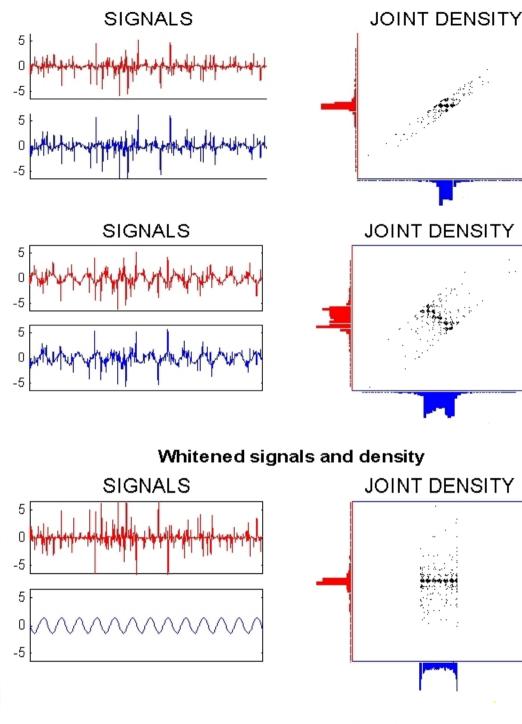
Applications:

- Filtering of one source out of many others,
- Separation of signals in telecommunication,
- Separation of signals from different regions of brain,
- Signal separation in astrophysics,
- Decomposition of signals in accelerator beam analysis in FERMILAB.



27.02.2017

Machine Learning, M.



Separated signals after 5 steps of FastICA

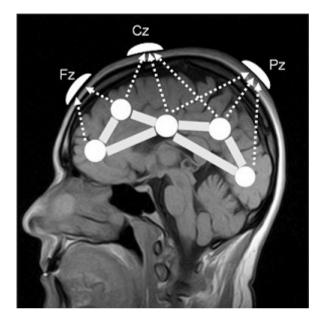
How does ICA work?

- We have two measured signals and we want to separate them into two independent sources.
- Preparing data decorrelation (correlation coefficients equal zero, σ=1).

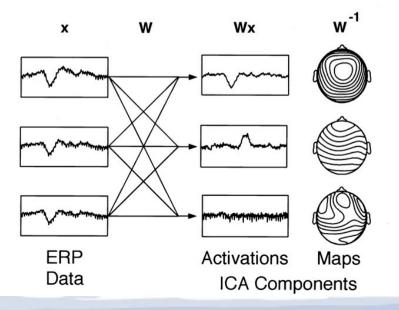
Superposition of many independent distributions gives Gaussian in the limit.

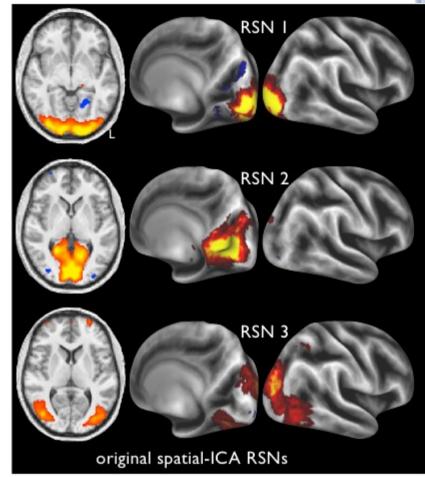
- ICA rotation, signals should be maximally non-Gaussian (measure of non-Gaussianity might be curtosis).
- Curtosis: $\operatorname{Kurt}_{n} = \frac{\frac{1}{n}\sum_{i=1}^{n}(x_{i} - \mu)^{4}}{\sigma^{4}} - 3$ where μ is i σ^{4} - 3 distribution and σ is a standard deviation.

ICA – brain research, signal separation



ICA Decomposition





3 components from 21-dimensional decomposition using the "spatial-ICA" algorithm.

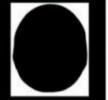
ning, M. Wolter

PNAS February 21, 2012 vol. 109 no. 8 3131-3136 36

ICA and magnetic resonance

White Matter

Skull





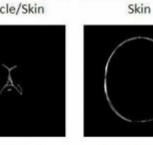


Fat





Muscle/Skin

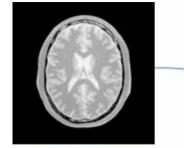


Glial Matter

Conn. Tissue

Sources of signals

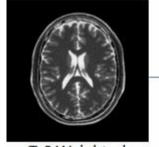
Blind Source Separation in Magnetic Resonance Images January 30, 2010 by Shubhendu Trived



Proton Density



T-1 Weighted



T-2 Weighted



Separated components

27.02.2017

Machine Learning, M. Wolter

37

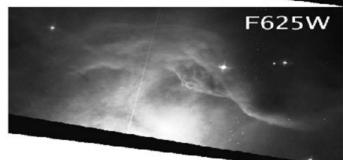
IC3

IC1

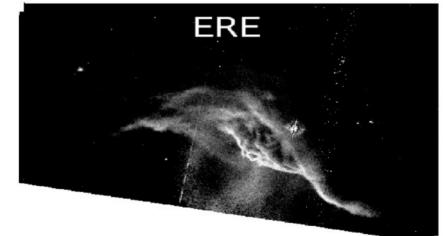
IC2

ICA – astronomy









On the left: HST images of the NGC 7023 North-West PDR in three SDSS wide-band filters. On the right: scattered light and ERE (Extended Red Emission) images extracted with FastICA from the observations.

F850LP

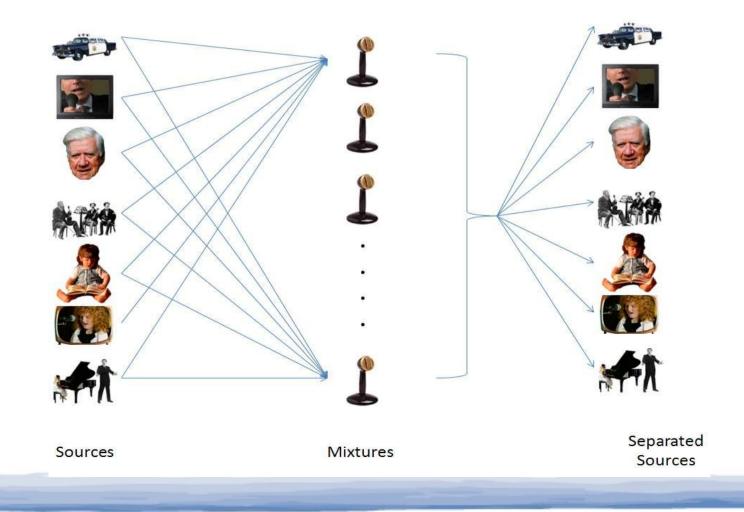
27.02.2017

Machine Learning, M. Wolter

A&A 479, L41-L44 (2008) DOI: 10.1051/0004-6361:200791588

Deser

 Cocktail Party Demo - applet pokazujący działanie algorytmu analizy składowych niezależnych na przykładzie problemu ,,przyjęcia cocktailowego".



Event Weighting

The probability p(S|x) is optimal in another sense: If one *weights* an admixture of signal and background events by the weight function

W(x) = p(S|x)

then the *signal* strength will be extracted with *zero bias* and the *smallest possible variance*, provided that our models describe the signal and background densities accurately and the signal to background ratio p(S)/p(B) is equal to the true value.

We can recover signal by event weighting with the discriminant output without cutting on it.

Roger Barlow, J. Comp. Phys. 72, 202 (1987)