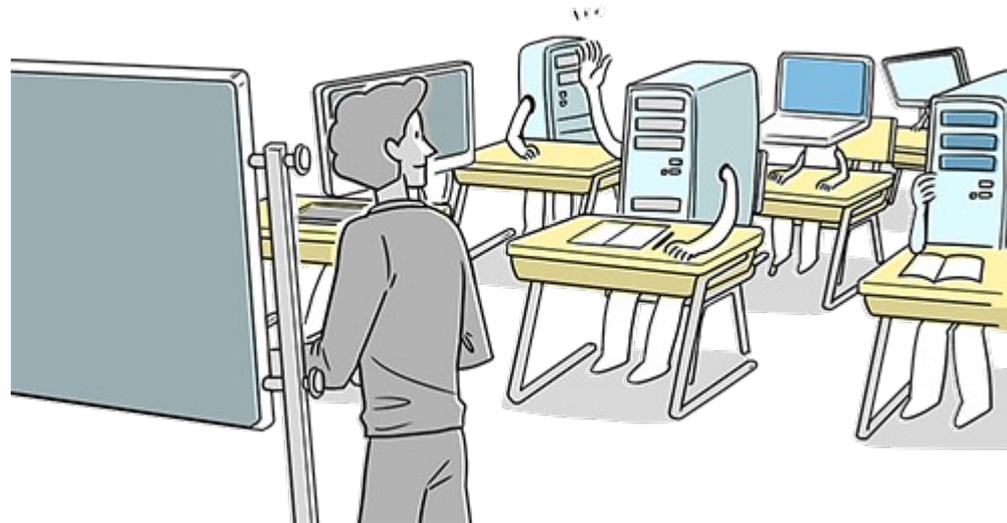


# Machine learning

## Lecture 1



Marcin Wolter  
*IFJ PAN*

*29 November 2019*

- Machine learning: what does it mean?
- Software to work with and literature.
- A little bit of mathematics and examples of simple linear classifiers.
- Some examples

All slides will be here: <https://indico.ifj.edu.pl/event/289/>



# Outline of the course

- Introduction: introduction to statistics, what does "Machine Learning" mean? A little bit of mathematics, but also examples of simple linear classifiers.
- Simple non-linear methods like Naive Bayes, k-Nearest Neighbors, Probability Density Estimators and Boosted Decision Trees (BDT).
- Neural Networks and Bayesian Neural Networks.
- Cross-validation and optimization of machine learning algorithms.
- Introduction to Deep Learning.
- Convolution Deep Neural network – classification of images.
- Some regression problems with Deep Neural Networks.
- Generative deep networks – Generative Adversary Networks (GANs)



# 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.
- Internet information about Deep Neural Networks – many things have changed in the recent years.



# Programs

- Python programming language
- <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
- <https://keras.io/>  
**Keras: The Python Deep Learning library**  
Emulates Deep Neural Network, uses google TensorFlow software
- **Google Colaboratory** as a working platform – you need a google ID  
<https://colab.research.google.com/>  
Platform to run python notebooks, gives free GPU



# About the course

- We will have 20 hours of lectures in total.
- I propose to do some exercises, run some python programs on *google colaboratory*. So please bring your laptops!
- We should make some type of exam. I propose, that instead of exam you write small projects with machine learning. This gives a chance to learn how to use machine learning tools.





# Caesar cipher

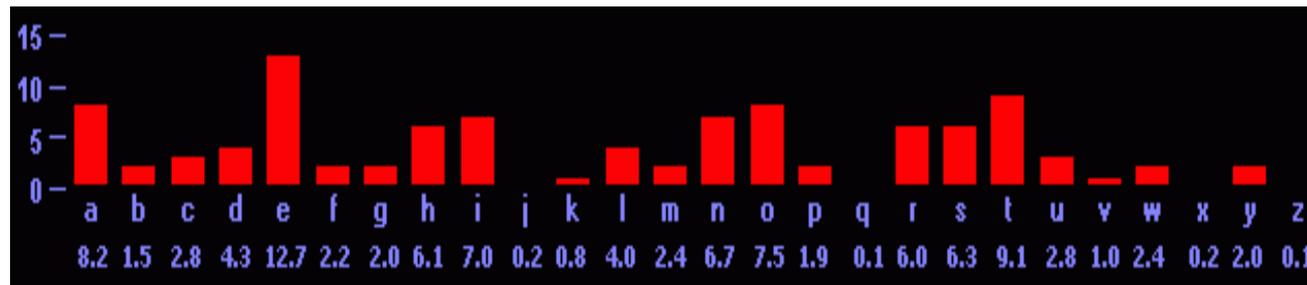
- Each letter of a text is replaced by another, shifted by n letters:

A	B	C	D	E
C	D	E	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)

A	B	C	D	E
Z	D	P	G	T

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

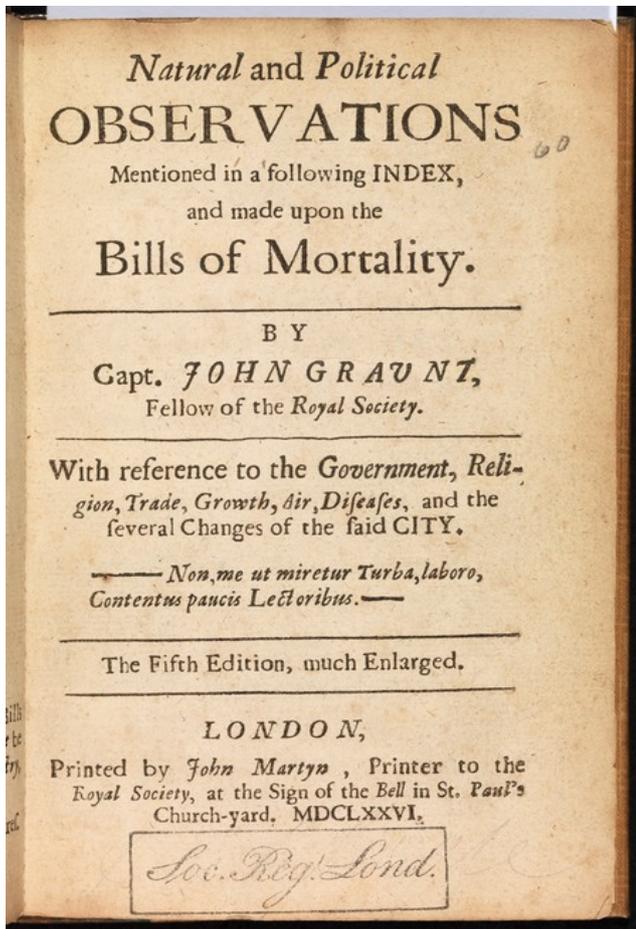


- Frequency of different letters in English.



# 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*“).



*The Table of CASUALTIES.*

The Years of our Lord	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	In 20 Years																																														
Age and Still-born	335	329	327	354	380	381	384	433	483	419	463	457	411	344	409	439	410	445	500	475	50	323	1703	2003	1841	1587	1812	1847	1855	1829	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	855		
Small Pox	918	835	859	898	780	834	864	974	743	892	889	1170	909	1095	779	712	661	671	704	623	728	714	2471	2814	3310	4432	3680	4377	4575	4631	1631	1632	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	1575
Measles and Suddenly	1260	884	751	970	1038	1112	882	1371	689	875	999	1800	303	2148	950	1091	1115	1108	953	1279	1622	2360	4418	6233	3865	4903	4363	4010	2378	1631	1632	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	2378	
Whooping Cough	68	74	64	74	100	111	118	86	92	102	113	128	91	87	22	36	17	24	35	28	4	4	54	24	4	9	1	1	1	1631	1632	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	11	
Whooping Cough	4	1	1	2	7	7	6	4	7	5	5	3	8	13	8	10	13	6	4	4	4	4	54	24	4	9	1	1	1	1631	1632	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	61	
Whooping Cough	3	2	5	1	3	4	3	2	1	7	10	3	5	7	3	10	7	5	1	3	12	3	25	19	24	31	20	19	103	1631	1632	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	103	
Whooping Cough	155	176	802	389	833	762	200	380	168	368	382	433	346	251	449	438	352	348	278	512	344	330	1587	1400	1422	1181	116	1597	1631	1632	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	7814		
Whooping Cough	3	6	10	5	11	8	5	7	10	5	7	4	6	6	3	10	7	5	1	3	12	3	25	19	24	31	20	19	103	1631	1632	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	103	
Whooping Cough	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1631	1632	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	11		
Whooping Cough	26	29	31	28	31	53	36	37	73	31	24	35	63	51	20	14	23	28	27	30	24	30	85	112	105	137	150	114	609	1631	1632	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	8	
Whooping Cough	66	28	54	42	68	51	33	72	44	81	18	27	73	68	6	4	4	1	1	1	1	74	15	29	190	244	161	133	659	1631	1632	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	3164	
Whooping Cough	161	106	114	117	200	213	158	192	177	201	236	225	225	194	150	157	112	171	132	143	163	210	590	608	498	769	859	490	2160	1631	1632	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	22106	
Whooping Cough	1369	1254	1069	890	1227	1280	1050	1343	1088	1393	1182	1144	858	1123	1598	1374	1035	2268	2130	2315	2113	1855	9279	8451	4678	4910	4788	4519	2160	1631	1632	1633	1634	1635	1636	1637	1638	1639	1640	1641	1642	1643	1644	1645	1646	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699</			

# What is an “event”?

- Elementary event is each result of an experiment (like throwing a dice), which result is random. It contains only a single outcome in the sample space.

**The numerical value of a probability may be obtained from its "classical" definition: The probability is equal to the quotient of the number of cases "favouring" a certain event to the total number of "equally possible" cases. (Laplace 1812).**

- Let's denote a set of all possible events as  $\Omega$ . Elements of a set  $\Omega$  are elementary events  $\omega$ , so  $\Omega$  is a set of elementary events. The set of events favoring  $A$  is a subset of  $\Omega$  and:

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

where  $|A|$  is a number of elements of a set  $A$ , and  $|\Omega|$  a number of elements of a set  $\Omega$ .

Example: probability of getting 6 while throwing a dice. Set of elementary events  $\Omega = \{1, 2, 3, 4, 5, 6\}$ , so the number of elementary events is  $|\Omega| = 6$ . Set of events favoring  $A = \{6\}$ , their number is  $|A| = 1$ . So  $P(6) = 1/6$



*Pierre-Simon Laplace  
(1749–1827)*

# Frequentist definition of probability

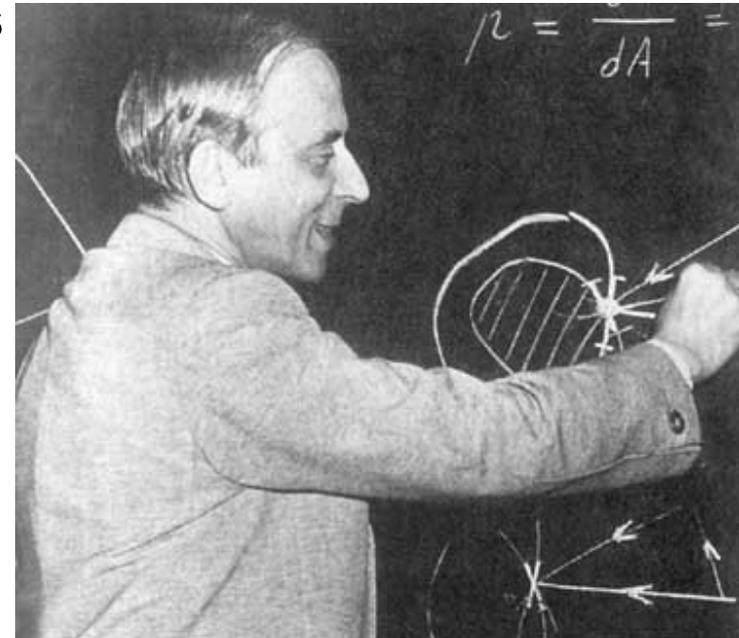
- Probability (frequentist definition) of an event  $A$  is a limit ( $N$  approaching infinity) of the ratio of  $n$  events when the event  $A$  occurred to the total number of trials  $N$ :

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

What is a probability of getting “6” while throwing a dice? It is the relative ratio of obtaining “6” in an infinite series of throws.

- The definition comes from Richard von Mises (born 19 of April 1883 in Lwów, died 14 July 1953 in Boston) – mathematician, brother of the economist Ludwig von Mises.

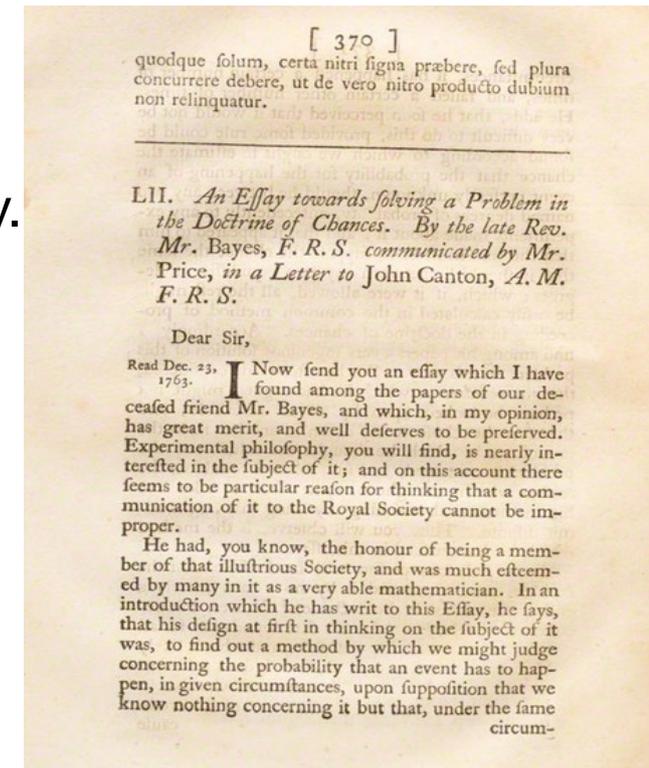
This idea of frequentist probability was used earlier, for example by de’Moivre, Bernoulli, Gauss ...



# Bayesian definition



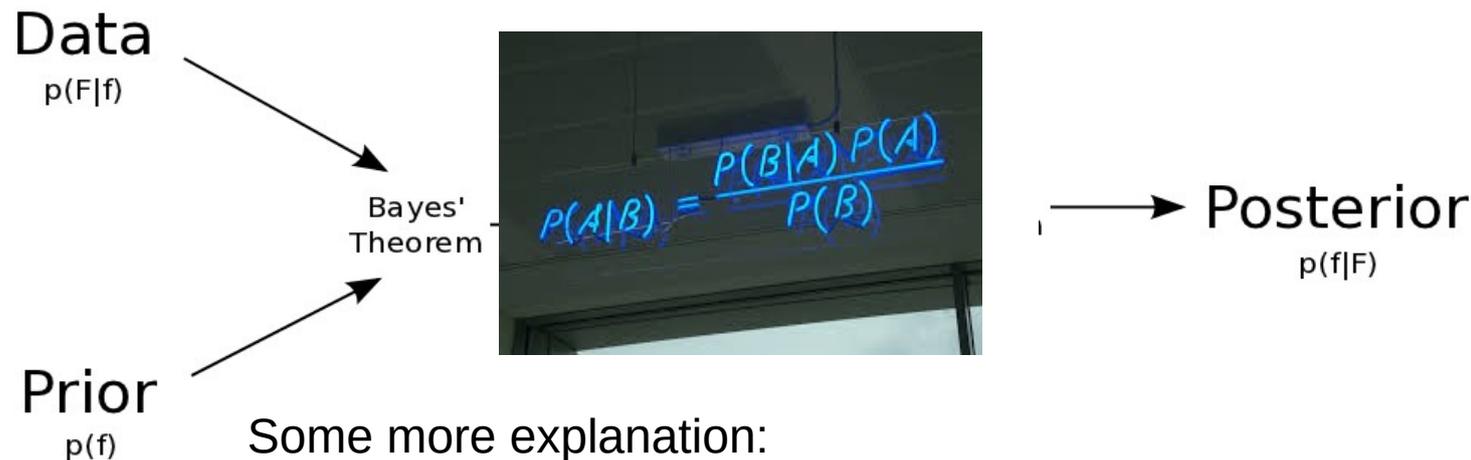
- Probability “a priori”, called unconditional, is a measure of belief, based on rational premises, that a given event will occur.
- In the next step we perform an experiment, called observation, and their results allow to modify the probability. We get the probability “a posteriori”, which is again a measure of belief, but modified by the observation.
- Supporters of the Bayesian approach were P. S. Laplace, H. Poincare and the economist John Keynes, arguing, that this is a method we use analyzing the world around us.



*Thomas Bayes (1702 - 1761) was an English statistician, philosopher and Presbyterian minister. The most important work: „Essay Towards Solving a Problem in the Doctrine of Chances”.*

# Bayesian probability

- Experiment we can't repeat many times: what is a probability to pass an exam?
- Based on our knowledge we estimate:  $\frac{1}{2}$  (probability *a priori*).
- But if all people before us didn't pass an exam (set of experiments) and we know our knowledge is not significantly higher we should verify this estimate (probability *a posteriori*).



<https://towardsdatascience.com/probability-concepts-explained-bayesian-inference-f-or-parameter-estimation-90e8930e5348>



# 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: a 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, let say 80% in our case.
- $\Pr(A)$  = Chance of having cancer (1%).
- $\Pr(\text{not } A)$  = Chance of not having cancer (99%).
- $\Pr(X|\text{not } A)$  = Chance of a positive test (X) given that you didn't have cancer (not 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.

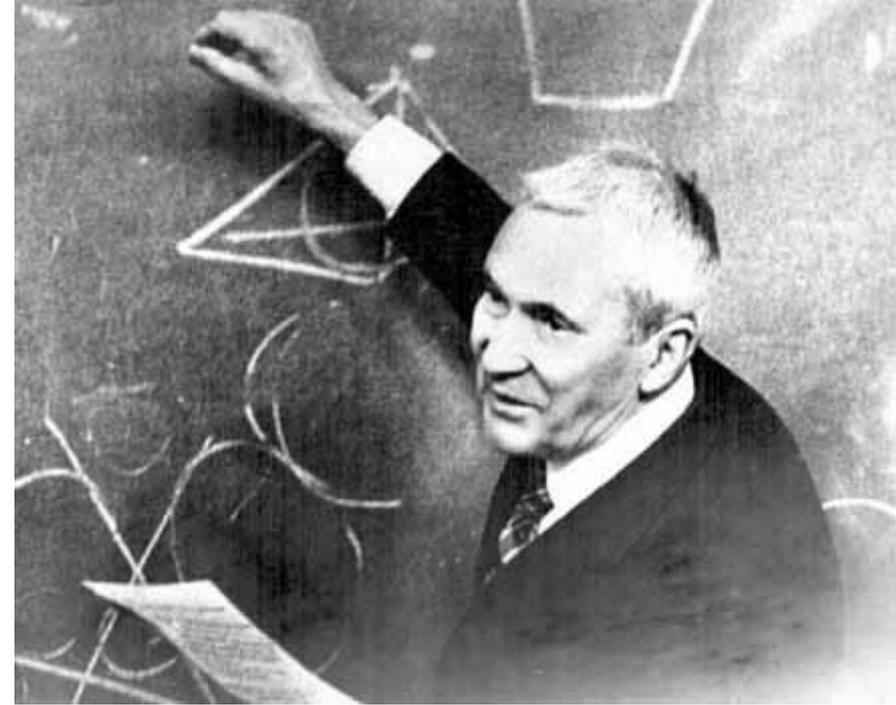
*“We don’t know all about the world to start with; our knowledge by experience consists simply of a rather scattered lot of sensations, and we cannot get any further without some a priori postulates. My problem is to get these stated as clearly as possible.”*

*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*

# Axiomatic definition

Probability can be defined in many ways...



Андрей Николаевич Колмогоров (1903-1987)

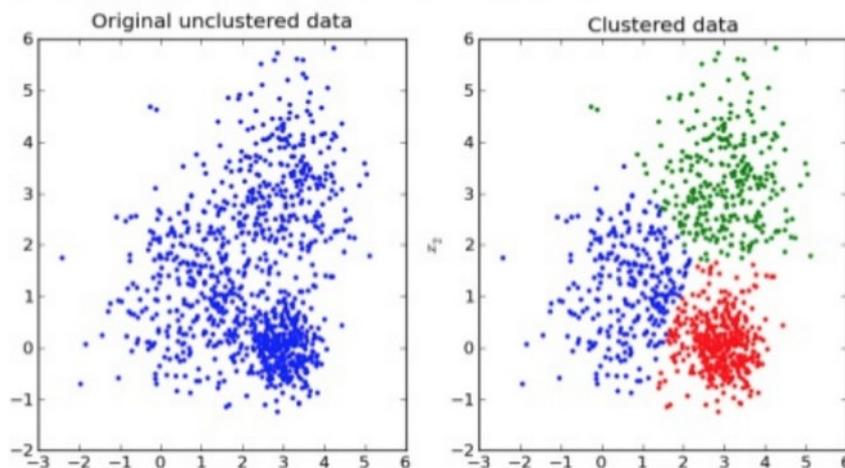
## Kolmogorov axiomatic definition:

Let  $Q_i$  denote anything subject to weighting by a normalized linear scheme of weights that sum to unity in a set  $W$ . The Kolmogorov axioms state that:

1. For every  $Q_i$  in  $W$ , there is a real number  $Q(Q_i)$  (the Kolmogorov weight of  $Q_i$ ) such that  $0 < Q(Q_i) < 1$ .
2.  $Q(Q_i) + Q(Q'_i) = 1$ , where  $Q'_i$  denotes the complement of  $Q_i$  in  $W$ .
3. For the mutually exclusive subsets  $Q_1, Q_2, \dots$  in  $W$ ,  
 $Q(Q_1 \cup Q_2 \cup Q_3 \cup \dots) = Q(Q_1) + Q(Q_2) + Q(Q_3) + \dots$

# What does “machine learning” mean?

- **Machine learning** is a field of computer science that gives computer systems the ability to "learn" (i.e. progressively improve performance on a specific task) with data, without being explicitly programmed.
- Problems:
  - Supervised learning (classification & regression)
  - Clustering (unsupervised learning)
  - Dimensionality reduction
  - Reinforcement learning
  - Many others.....



## ➤ Unsupervised Learning

- ❑ Technique of trying to find hidden structure in unlabeled data

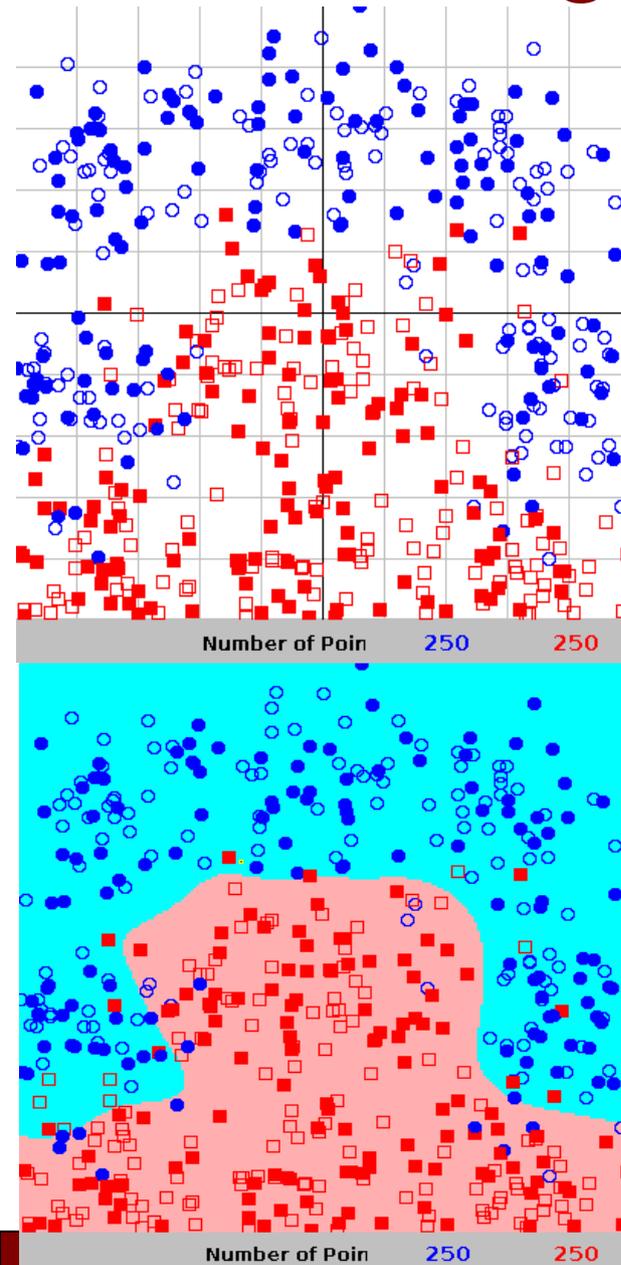
## ➤ Supervise Learning

- ❑ Technique for creating a function from training data. The training data consist of pairs of input objects (typically vectors), and desired outputs.

# How do the (supervised) machine learning algorithms work?

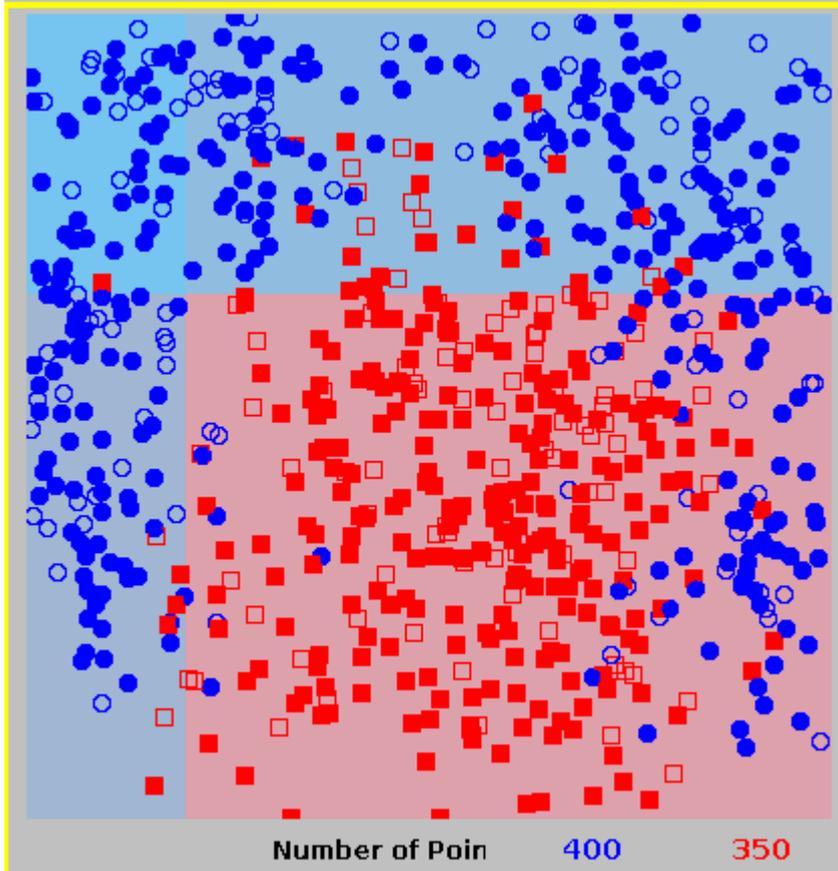
- We need **training data**, for which we know the correct answer, whether it's a signal or background. We divide the data into two samples: training and test.
- We find the best function  $f(\mathbf{x})$  which describes the probability, that a given event belongs to the class "signal". This is done by minimizing the loss function (for example  $\chi^2$ ).
- Different algorithms differ by: the class of function used as  $f(\mathbf{x})$  (linear, non-linear etc), loss function and the way it's minimized.
- All these algorithms try to approximate the unknown *Bayesian Decisive Function* (BDF) relying on the finite training sample.

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

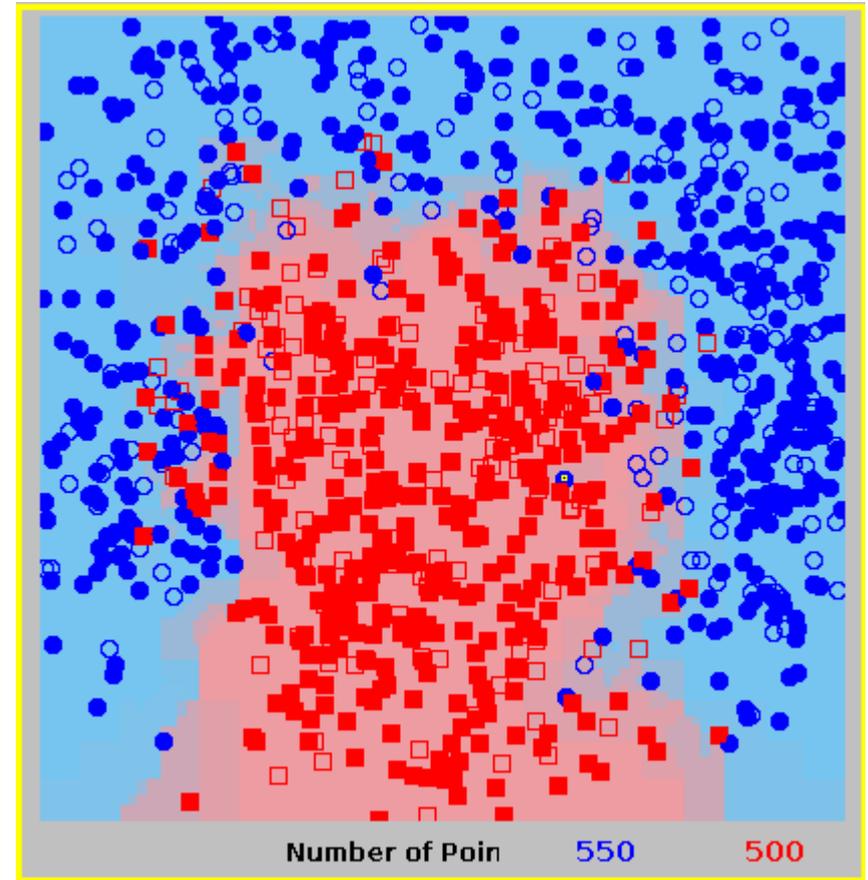


# Cuts vs non-linear separation

Cuts



Non-linear separation



Neural Networks, boosted decision trees, and so on....

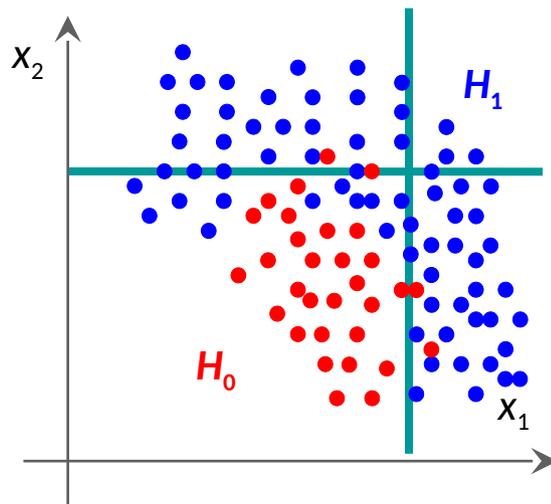
# Types of algorithms



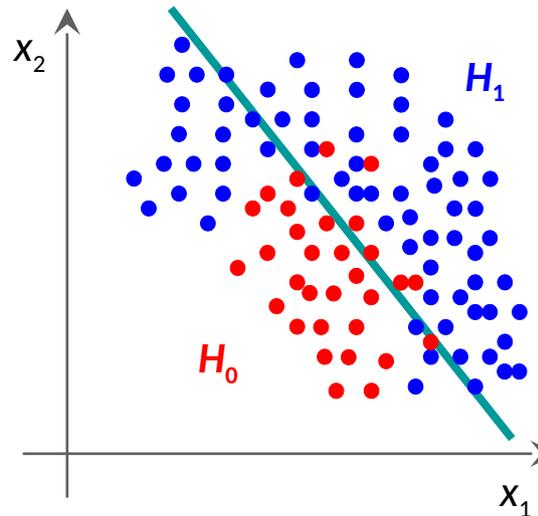
How to use the information available

**Classification:** find a function  $f(x_1, x_2)$  giving the probability, that a given data point belongs to a given class (signal vs background).

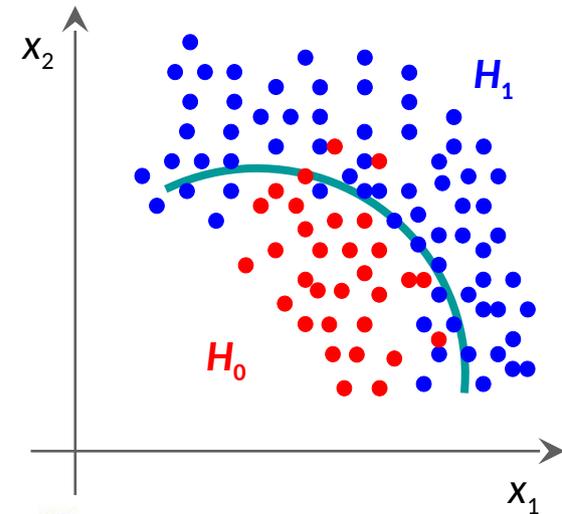
Simple cuts  
(easy and intuitive)



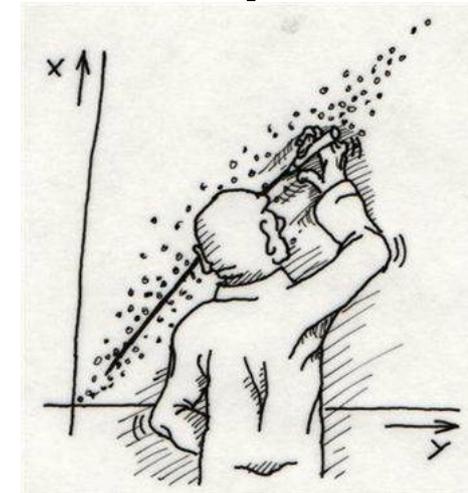
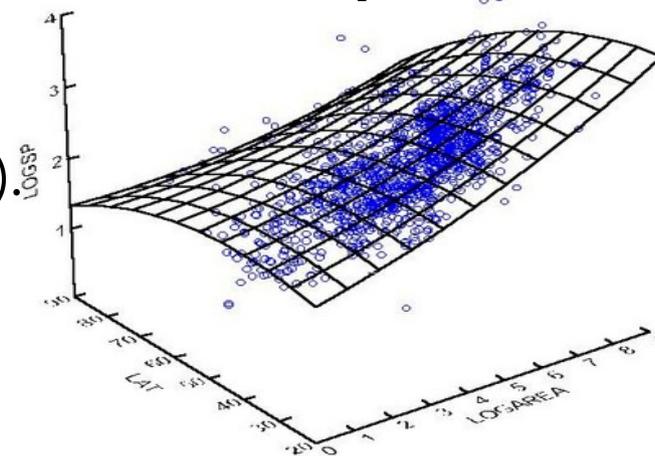
Linear  
(fast and stable)



Non-linear  
(most effective)



**Regression:** fit a continuous function  
(find particle energy from calo readouts).



# Classification

A Bayes classifier (optimal classifier):

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

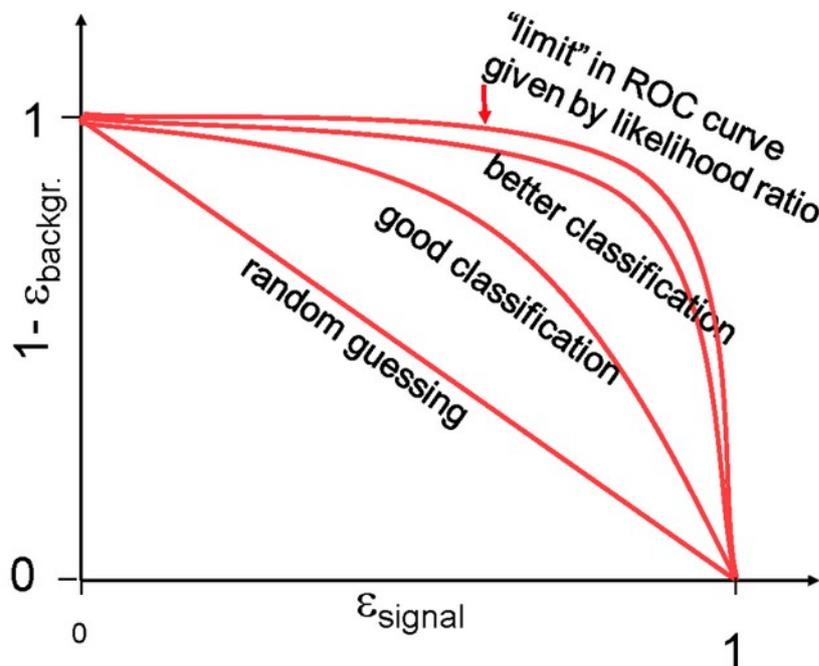
where **S** is associated with  $y = 1$  and **B** with  $y = 0$ . **Bayes classifier** accepts events  $x$  if  $p(\mathbf{S}|x) > \mathbf{cut}$  as belonging to **S**.

We need to approximate probability distributions  $P(x|\mathbf{S})$  and  $P(x|\mathbf{B})$ .

- 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 all we will be talking about, 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 ( $1-\varepsilon_B$ ) vs signal efficiency  $\varepsilon_B$ . Shows how good the classifier is.
- The integral of ROC could be a measure of the classifier quality:



Integral(ROC) =  $\frac{1}{2}$  – random

Integral(ROC) = 1 - ideal



# Practical applications

## A Short List of Multivariate Methods

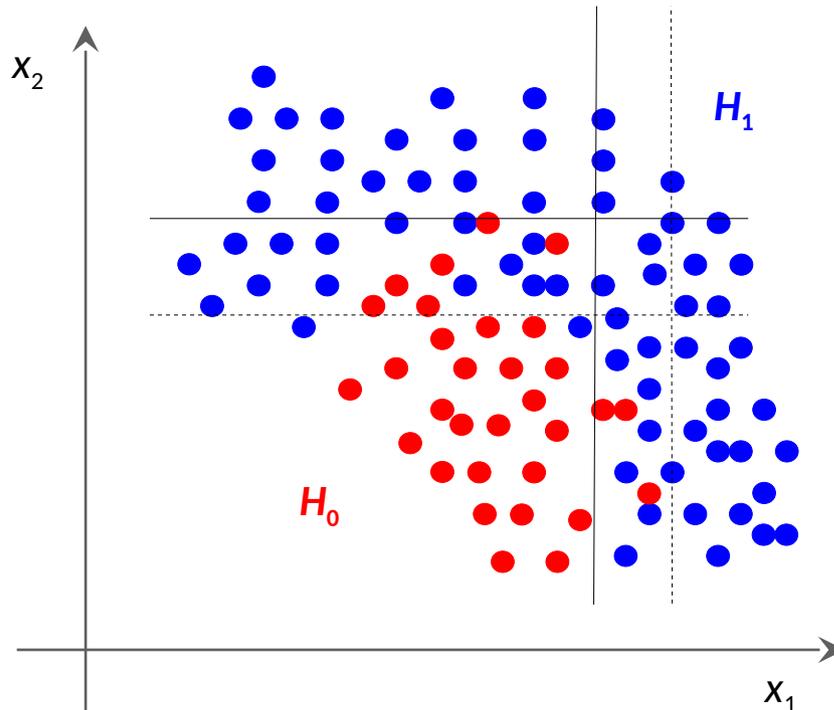
- Cuts
- Linear Discriminants (like Fisher)
- Naive Bayes (Likelihood Discriminant)
- Kernel Density Estimation
- Decision Trees
- Neural Networks
- Bayesian Neural Networks
- Genetic Algorithms
- Support Vector Machines
  
- And many, many others..... I want to present briefly just few of them.



# We will talk in the first lectures about:

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

# Cuts



## Optimization of cuts:

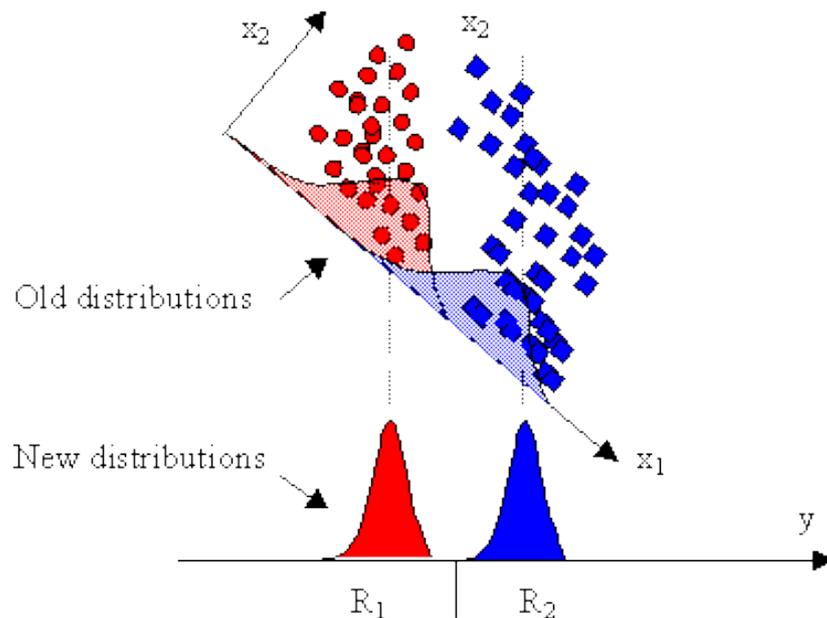
- Move cuts as long as we get the optimal signal vs. background selection. For a given signal efficiency we find the best background rejection → we get the entire ROC curve.
- Optimization methods:
  - Brute force
  - Genetic algorithms
  - Many others...

# Fisher discriminants

## LDA, Linear Discriminant Analysis

Projection to one dimension, than discrimination

Equivalent to linear separation



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

### Assumptions for new basis:

- Maximize distance between projected class means
- Minimize projected class variance

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

# Fisher Linear Discriminant Analysis

## Objective

$$\operatorname{argmax}_w J(w) = \frac{w^T S_B w}{w^T S_W w}$$

$$m_i = \frac{1}{n_i} \sum_{x \in C_i} x$$

$$S_B = (m_2 - m_1)(m_2 - m_1)^T \quad \text{Variance Between classes}$$

$$S_W = \sum_j^2 \sum_{x \in C_j} (x - m_j)(x - m_j)^T \quad \text{Variance Within class}$$

## Algorithm

1. Compute class means
2. Compute  $w = S_W^{-1}(m_2 - m_1)$
3. Project data  $y = w^T x$



# Fisher's linear discriminant

The terms *Fisher's linear discriminant* and *LDA* are often used interchangeably, although [Fisher's](#) 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  $\vec{w}^T \Sigma_{y=i} \vec{w}$  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 [signal-to-noise ratio](#) for the class labelling. It can be shown that the maximum separation occurs when

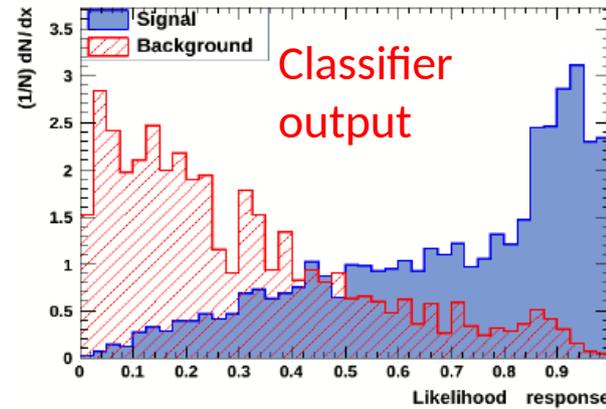
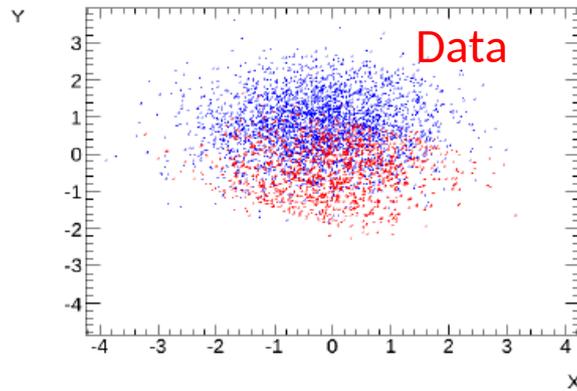
$$\vec{w} = (\Sigma_{y=0} + \Sigma_{y=1})^{-1} (\vec{\mu}_{y=1} - \vec{\mu}_{y=0})$$

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

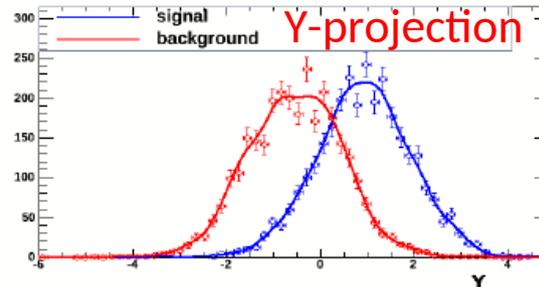
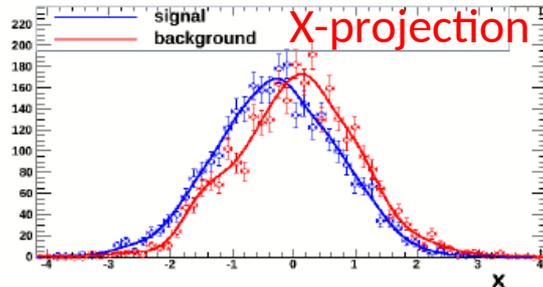
Be sure to note that the vector  $\vec{w}$  is the normal to the discriminant hyperplane. As an example, in a two dimensional problem, the line that best divides the two groups is perpendicular to  $\vec{w}$ .

Generally, the data points are projected onto  $\vec{w}$ . However, to find the actual plane that best separates the data, one must solve for the bias term  $b$  in  $w^T \mu_1 + b = -(w^T \mu_2 + b)$ .

# Naive Bayes classifier



Frequently called  
**“projected likelihood”**  
 (TMVA).



- Based on the assumption, that variables are independent (so „naive“):

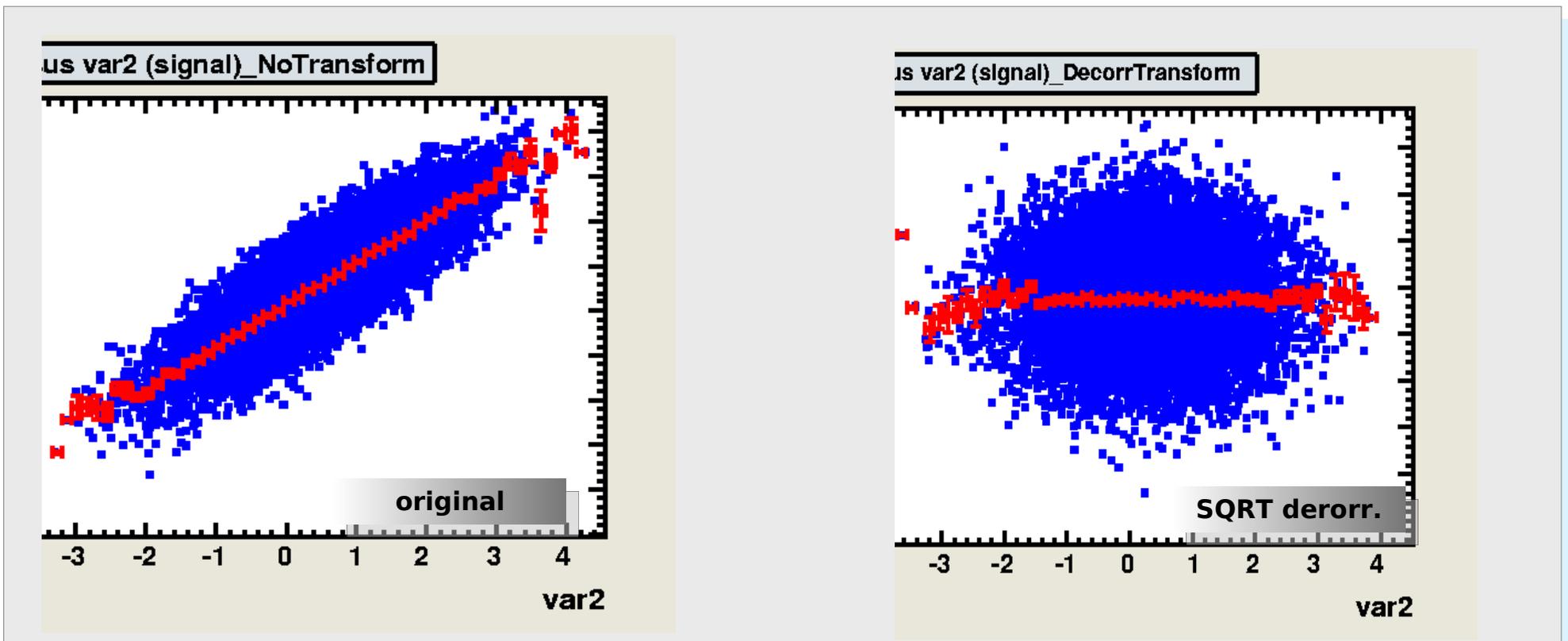
$$P(y | x_1, \dots, x_n) = \frac{P(y)P(x_1, \dots, x_n | y)}{P(x_1, \dots, x_n)} \quad \text{“Naive” assumption: } P(x_i | y, x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) = P(x_i | y),$$

$$P(y | x_1, \dots, x_n) = \frac{P(y) \prod_{i=1}^n P(x_i | y)}{P(x_1, \dots, x_n)}$$

- Output probability is a **product of probabilities for all variables**.
- Fast and stable, not optimal, but in many cases sufficient.

# Decorrelation

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

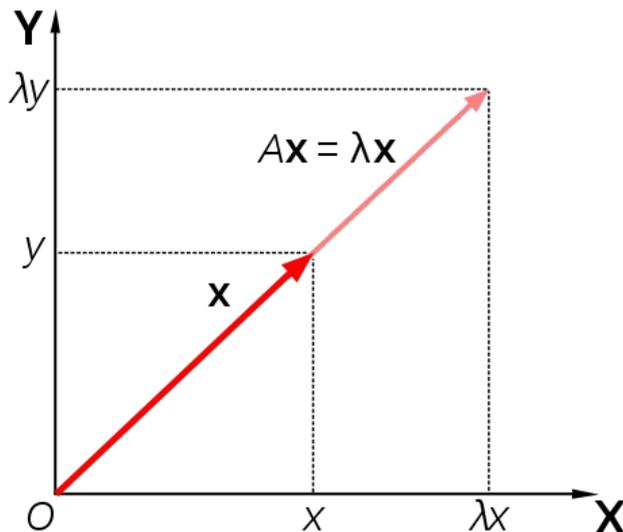


# Eigenvalues and eigenvectors

In essence, an eigenvector  $\mathbf{v}$  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  $\lambda$ , called an eigenvalue. This condition can be written as the equation

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

referred to as the eigenvalue equation or eigenequation. In general,  $\lambda$  may be any scalar. For example,  $\lambda$  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  $\mathbf{x}$ , not changing its direction, so  $\mathbf{x}$  is an eigenvector of  $A$ .

$$\begin{matrix} \mathbf{A} & & \mathbf{Q} & & \mathbf{\Lambda} & & \mathbf{Q}^{-1} \\ \left[ \begin{array}{|c|} \hline \phantom{0} \\ \hline \end{array} \right] & = & \left[ \begin{array}{|c|} \hline \mathbf{v}_1 \\ \hline \mathbf{v}_2 \\ \hline \mathbf{v}_3 \\ \hline \end{array} \right] \left[ \begin{array}{|c|} \hline \lambda_1 & 0 & 0 \\ \hline 0 & \lambda_2 & 0 \\ \hline 0 & 0 & \lambda_3 \\ \hline \end{array} \right] \left[ \begin{array}{|c|} \hline \mathbf{v}_1 \\ \hline \mathbf{v}_2 \\ \hline \mathbf{v}_3 \\ \hline \end{array} \right]^{-1} \\ \text{Eigen vectors} & & \text{Eigen values} & & \text{Eigen vectors} \\ \text{of} & & \text{of} & & \text{of} \\ \mathbf{A} & & \mathbf{A} & & \mathbf{A} \end{matrix}$$

Eigendecomposition of a matrix

# Covariance Matrix

- Let  $X$  be a  $p$ -variate random vector. The covariance matrix of  $X$  is defined as:

$$\begin{aligned}\Sigma_{XX} &= \text{Var}(X) = E\{(X - \mu)^T (X - \mu)\} \\ &= \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_p) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & \dots & \text{Cov}(X_2, X_p) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_p, X_1) & \text{Cov}(X_p, X_2) & \dots & \text{Var}(X_p) \end{pmatrix}.\end{aligned}$$

Where:

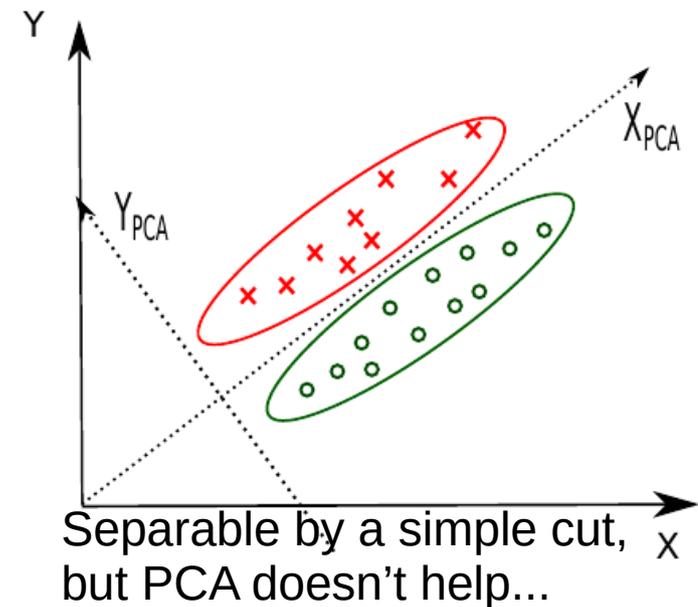
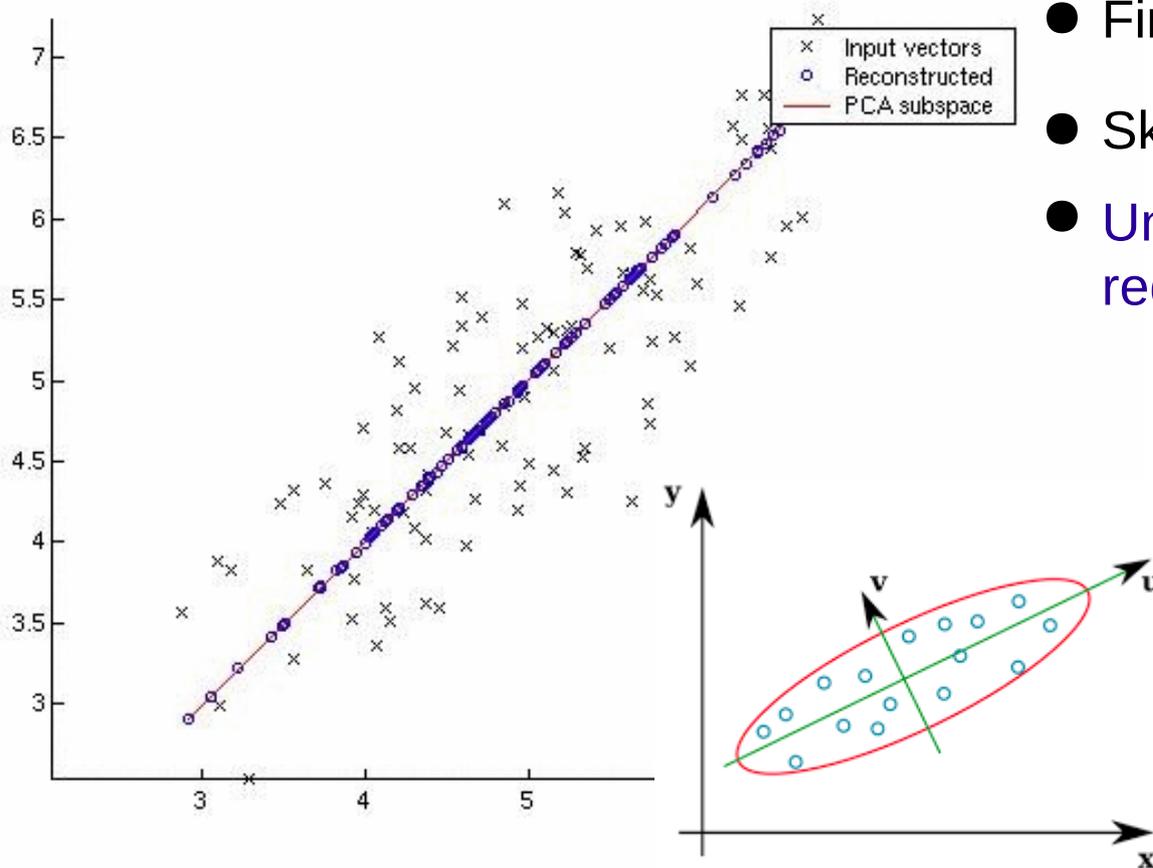
$$\begin{aligned}\text{Var}(X_i) &= E\{(X_i - \mu(X_i)) \cdot (X_i - \mu(X_i))\} \\ \text{Cov}(X_i, Y_j) &= E\{(X_i - \mu(X_i)) \cdot (Y_j - \mu(Y_j))\}\end{aligned}$$

# 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 might be skipped

## Procedure:

- Find the covariance matrix  $\text{Cov}(X)$
- Find eigenvalues  $\lambda_i$  and eigenvectors  $v_i$
- Skip smallest  $\lambda_i$
- Unsupervised learning & dimensionality reduction



# Applications

## ● Uses:

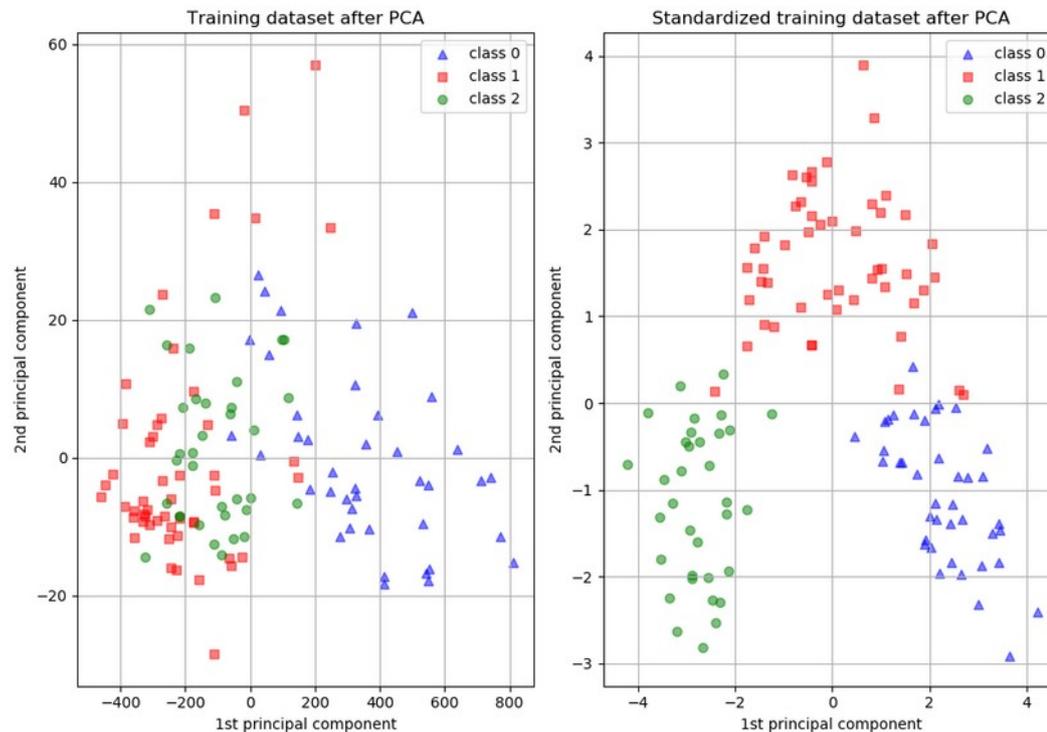
- Data Visualization
- Data Reduction
- Data Classification
- Noise Reduction

## ● Examples:

- How many unique “sub-sets” are in the sample?
- How are they similar / different?
- Which measurements are needed to differentiate?
- How to best present what is “interesting”?
- Which “sub-set” does this new sample rightfully belong?

# What would happen If I did PCA without normalization? Why do we normalize data?

In PCA we are interested in the components that maximize the variance. If one component (e.g. human height) varies less than another (e.g. weight) because of their respective scales (meters vs. kilos), PCA might determine that the direction of maximal variance more closely corresponds with the 'weight' axis, if those features are not scaled. As a change in height of one meter can be considered much more important than the change in weight of one kilogram, this is clearly incorrect.



*The dataset used is the Wine Dataset available at UCI. This dataset has continuous features that are heterogeneous in scale due to differing properties that they measure (i.e alcohol content, and malic acid).*

# Correlation matrix

The **correlation matrix** refers to the symmetric array of numbers

$$\mathbf{R} = \begin{pmatrix} 1 & r_{12} & r_{13} & \cdots & r_{1p} \\ r_{21} & 1 & r_{23} & \cdots & r_{2p} \\ r_{31} & r_{32} & 1 & \cdots & r_{3p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ r_{p1} & r_{p2} & r_{p3} & \cdots & 1 \end{pmatrix}$$

where

$$r_{jk} = \frac{s_{jk}}{s_j s_k} = \frac{\sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)}{\sqrt{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2} \sqrt{\sum_{i=1}^n (x_{ik} - \bar{x}_k)^2}}$$

is the Pearson correlation coefficient between variables  $\mathbf{x}_j$  and  $\mathbf{x}_k$ .



# Correlation or covariance matrix?

- Mean-centering is unnecessary if performing a principal components analysis on a correlation matrix, as the data are already centered after calculating correlations.
- We tend to use the covariance matrix when the variable scales are similar and the correlation matrix when variables are on different scales.

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$$

Correlation between X and Y

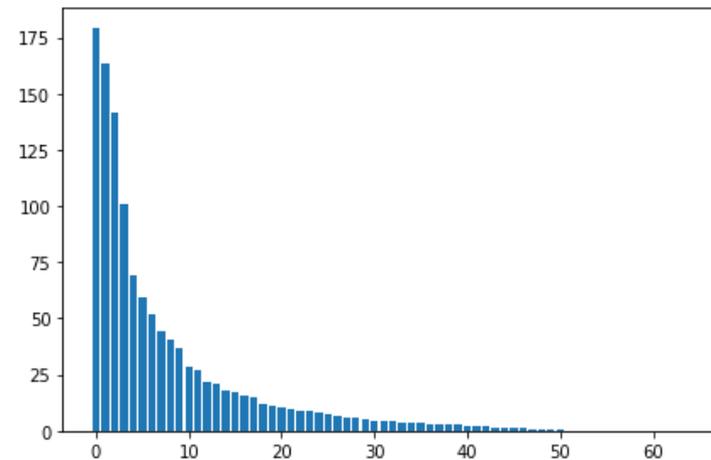
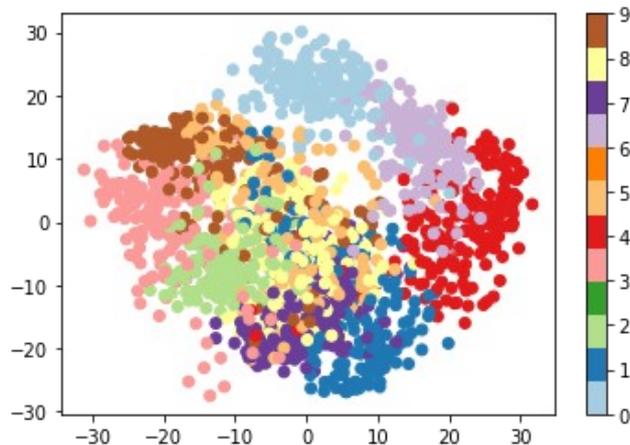
Standard deviation of X

Standard deviation of Y

Covarianced normalized by Standard Deviation

# Example

- All examples will be available here:  
<https://github.com/marcinwolter/MachineLearnin2019>
- `plot_digits_classif.ipynb` - iPython notebook prepared to run on Google Colaboratory <https://colab.research.google.com/> (for example)
  - Reads handwritten digits
  - Performs PCA
  - Displays two first principal components:



- Classification using Naive Bayes and LDA



# Summary

- What is machine learning
- Simple unsupervised methods – PCA and decorrelation
- Classifiers:
  - Cuts
  - Naive bayes
  - Fisher Linear Discriminants