

Selected topics in Data Science



Grandjean, Martin (2014).
"La connaissance est un réseau".
Les Cahiers du Numérique

Wojciech Krzemień

23.10 2020, NCBJ

Last lecture recap

- Machine learning (ML) → learning from data without saying explicitly how.
- Main types of ML:
 - **Supervised Learning**- training set contains pairs {X,Y}, where Y is the expected output value/expected label
 - **Unsupervised Learning** - only feature vectors {X} given, the computer must classify/cluster by itself.
 - **Reinforcement Learning** - every action with environment provides „reward” looped-back to your model
- Dependent on outcome type:
 - Continuous output Y – **regression** model
 - Discrete output (labels) – **classification** model
- Introduction to **K-Nearest Neighbors (KNN)** algorithm

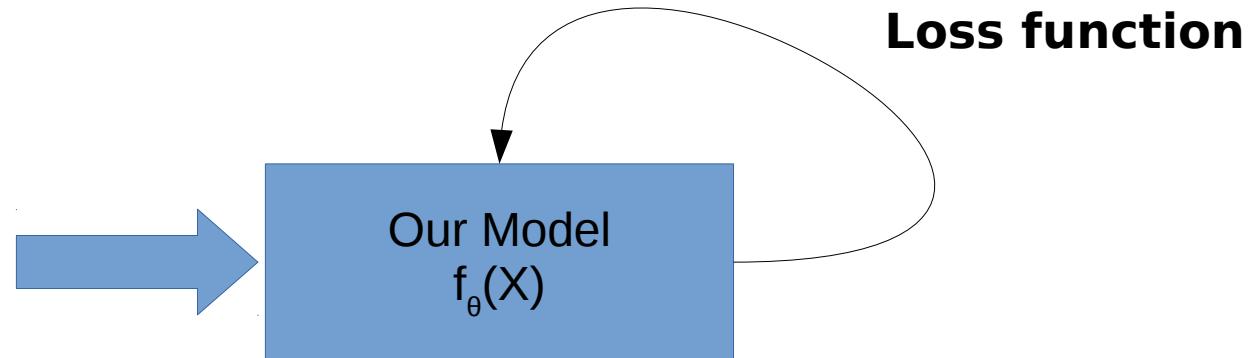
Todays plan

- Elements of **Statistical Learning Theory**:
 - Loss function
 - Notions of generalization and stability
 - Empirical risk minimization
 - overfitting and model complexity
 - Prediction vs Inference
- Crash course in **Statistics**

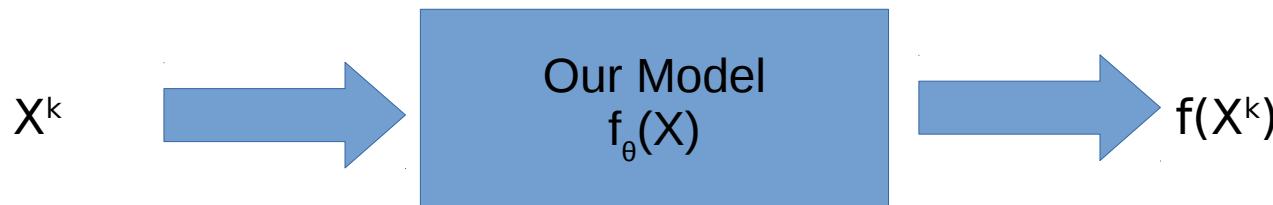
Supervised learning

1. Training

$\{X^1, Y^1\}, \{X^2, Y^2\}, \dots$
training data



2. Prediction:



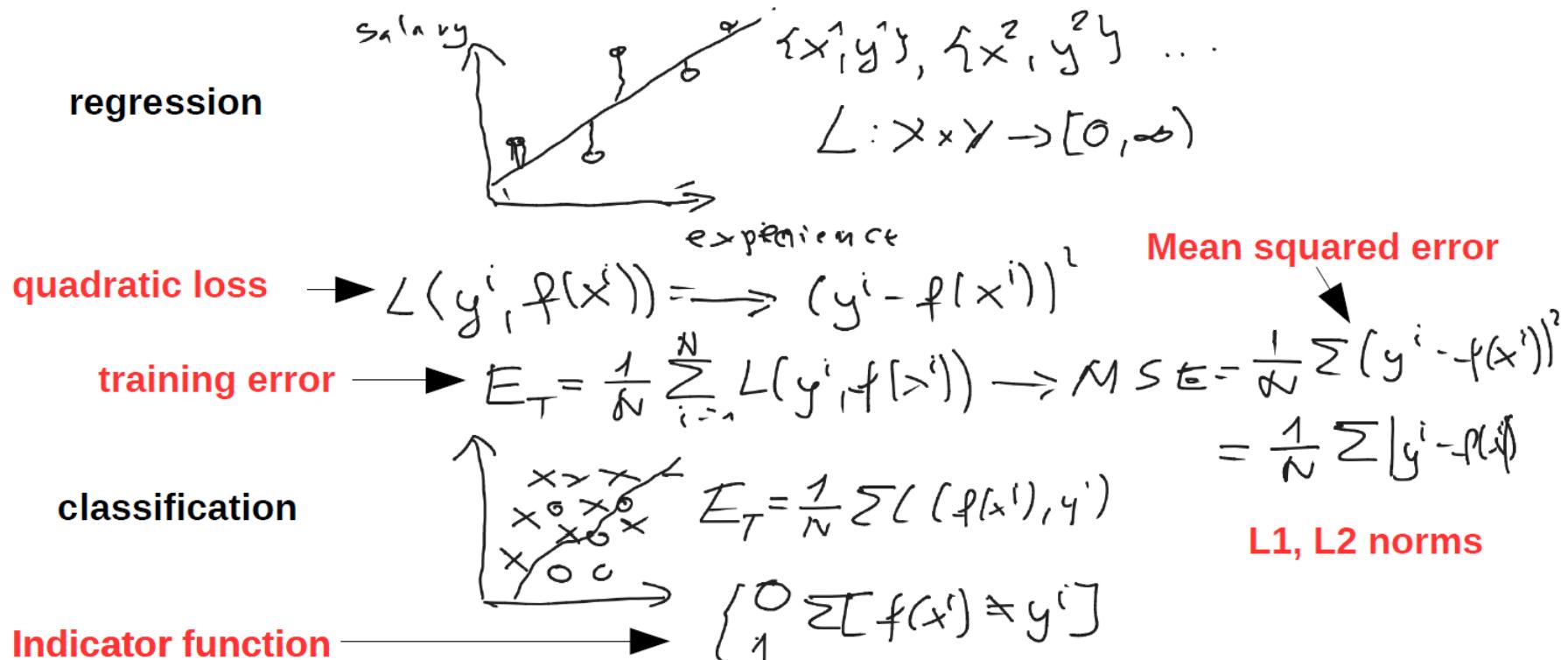
X - input, set of features (attributes)

Y - output

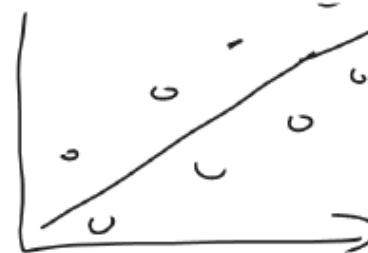
Training samples - pairs of $\{X^i, Y^i\}$ $i=1,\dots,N$

θ - some parameters of the model

Loss function L to penalize the mistakes the model makes



Training expressed as **minimization** problem



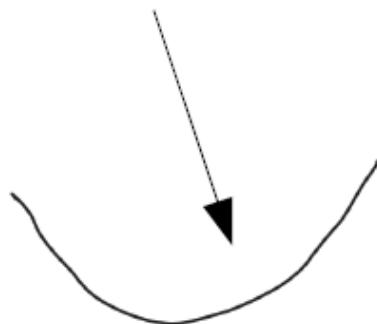
$$f_{\theta} = \theta_0 + \theta_1 x_1$$

Find parameters θ^* for which E_T has minimum

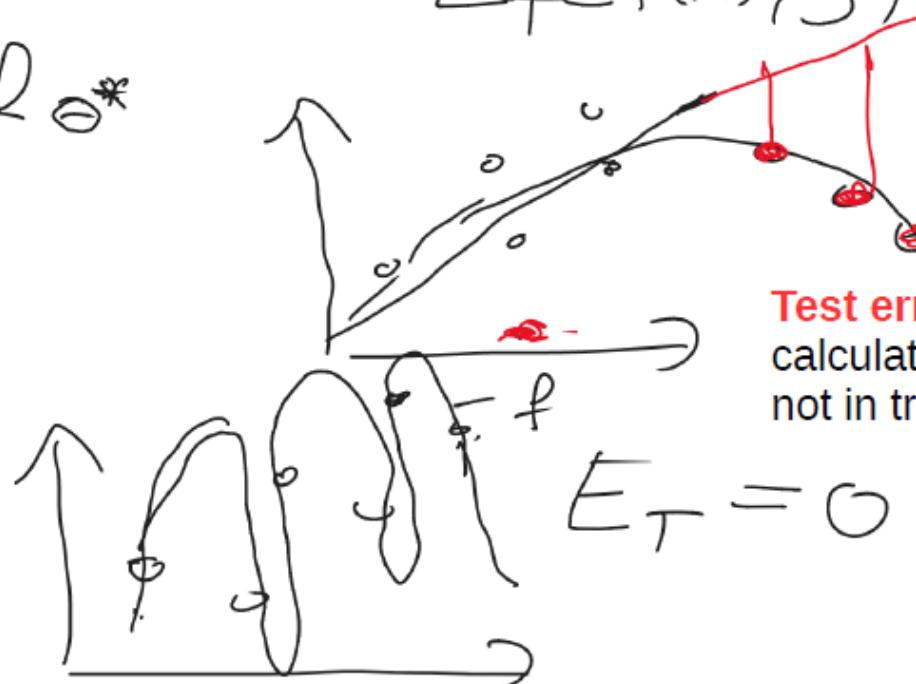
$$E[f(x), y]$$

$$f_{\theta^*}$$

Convex function



No more than one minimum



Test errors
calculated for points
not in training set

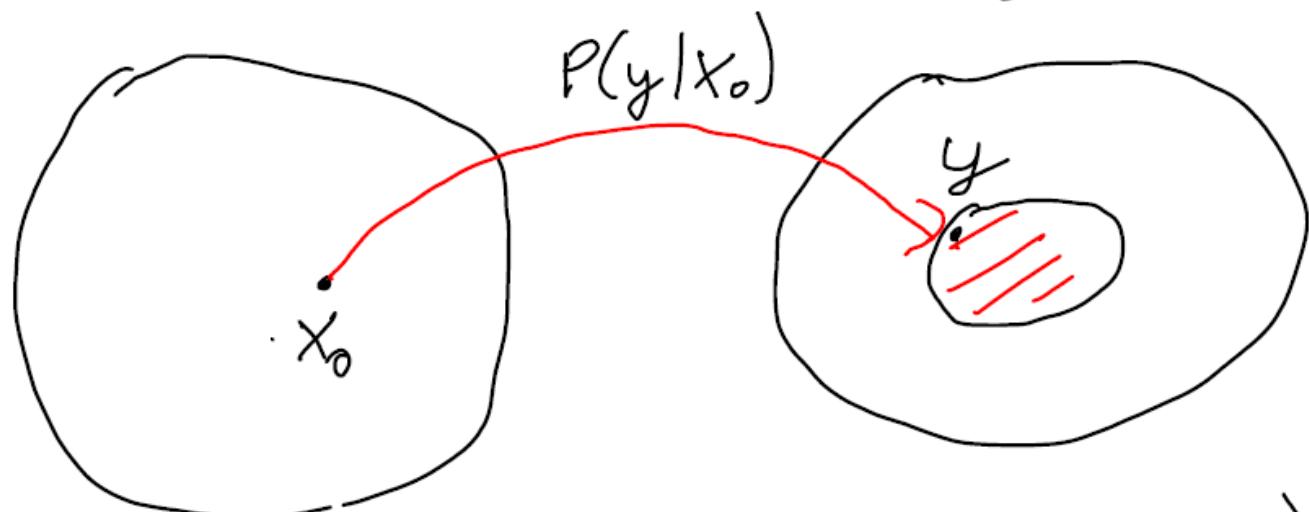
$$E_T = 0$$

OVERFITTING

Generalization

Model should describe all the data

Statistical Learning Theory



feature space X

e.g. \mathbb{R}^p

target space Y

e.g. $y \in \mathbb{R}$ - regression

$y \in \{0,1\}$ - binary classification

$Z = X \times Y$, Z is a random variable with a joint probability density distribution

$$S_Z = S_{X \times Y}$$

$$\mathcal{L}: Y \times Y \rightarrow [0, \infty)$$

↑
Loss function

Distribution $p(z)$ (\mathcal{S}_{xy}) is unknown.

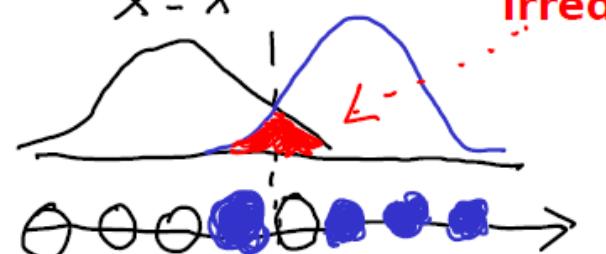
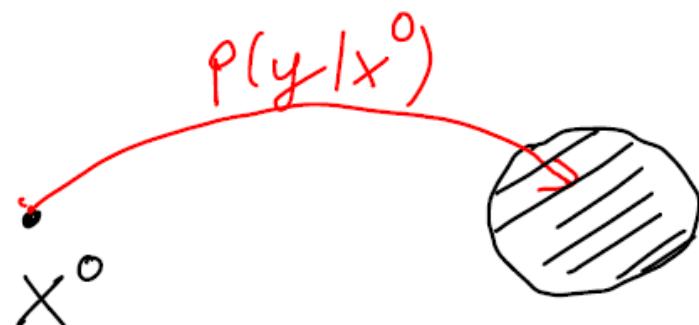
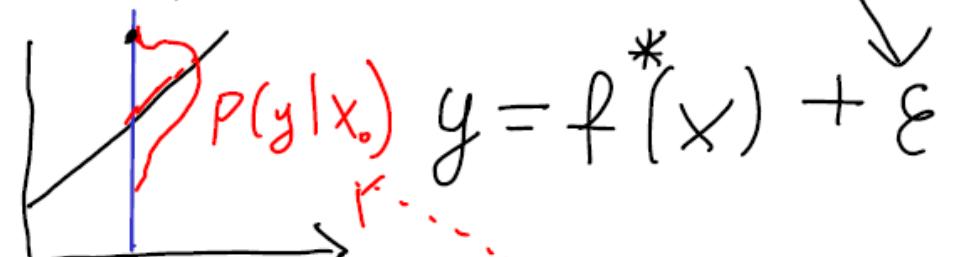
We have only a training sample:

$$T = \{(x^1, y^1), \dots, (x^N, y^N)\} \quad N - \text{number of samples}$$

Assumption: Training pairs are independent and identically distributed (iid) from Z

$$P(z) = P_X(x) \cdot P_{Y|X}(y|x)$$

gaussian noise



A set of features x_0 can correspond to different answers

$P_{\hat{x}}(x)$ this probability incorporates e.g. uncertainty due to the sampling procedure.

Key question: Starting from given training set T , and without the knowledge about the underlying distribution $p(z) = p(x) \cdot p_{y|x}(y|x)$ how to find a model f that makes 'good' predictions.

GENERALIZATION

Function (hypothesis) space H : $f : X \rightarrow Y$, contains all possible models we allow to be checked.

What conditions / constraints one should apply on function space H to ensure the generalization?

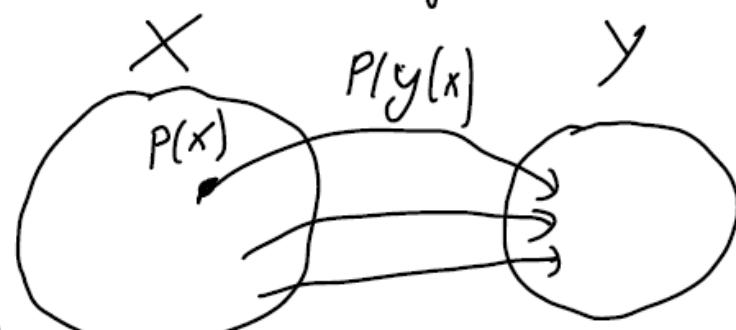
Vapnik's theory gave the answer.

Expected prediction error (true error, risk function).

Joint probability density function

$$E_{PE}[f] = \iint L(f(x), y) \underbrace{s_{x,y}}_{\text{Joint PDF}} dx dy$$

The error that given model f would commit calculated over the whole $Z=X \times Y$ space.



Since we don't know what is $s_{x,y}$ we cannot compute it

$$E_f[f] = \frac{1}{N} \sum_i^N L(f(x^i, y^i)) - \text{training error, empirical risk}$$

For the 'good' model we expect: $E_f \approx E_{PE}$

GENERALIZATION:

$$\forall_{T, S_{XY}} \lim_{N \rightarrow \infty} |E_{PE}(f_T) - E_T[f_T]| = 0$$

convergence in probability

Convergence in probability - simple example:

$x_1 \quad x_2$
•
 \vdots
 x_n

mean_T = $\frac{1}{N} \sum x_i$ \longrightarrow mean_T \rightarrow population expectation
going with $N \rightarrow \infty$

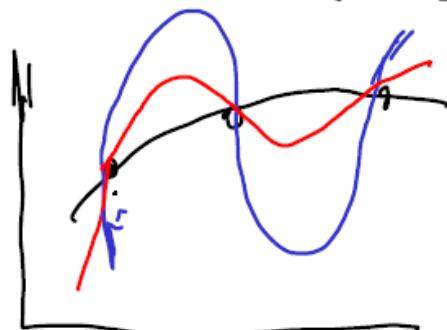
sample $T \in \{x_1, \dots, x_n\}$
 \uparrow
N numbers

Formally: $\forall_{\varepsilon > 0} \lim_{N \rightarrow \infty} \Pr(|\text{mean}_N - E| > \varepsilon) = 0$

Empirical risk minimisation (ERM) by V. Vapnik et al.

**He is the creator of
Support Vector Machines algos**

If we use training set to calculate E_T and find the optimal model f^* how we must restrict H space to fulfil generalization condition?



H space must be uniform Glivenko - Cantelli class.

For details see:

Poggio et al. proved that the ERM approach is equivalent to the requirements for a learning algorithm to be stable

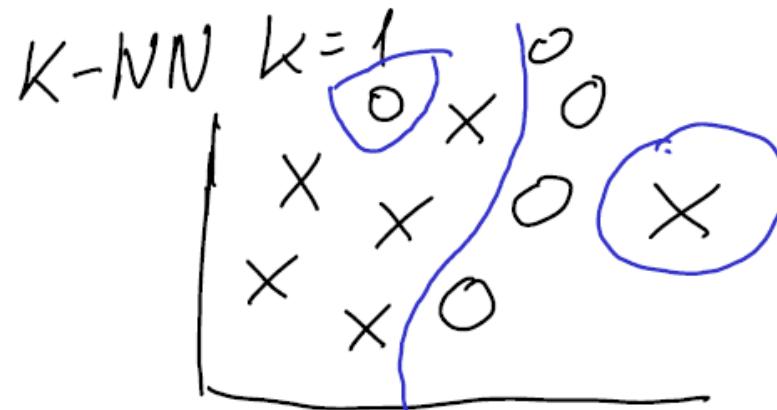
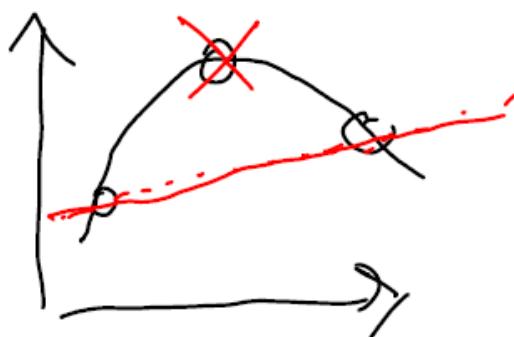
Approach can be applied to analyze broader class of ML algos.

CV_{loo} - Cross-Validation leave-one-out stability

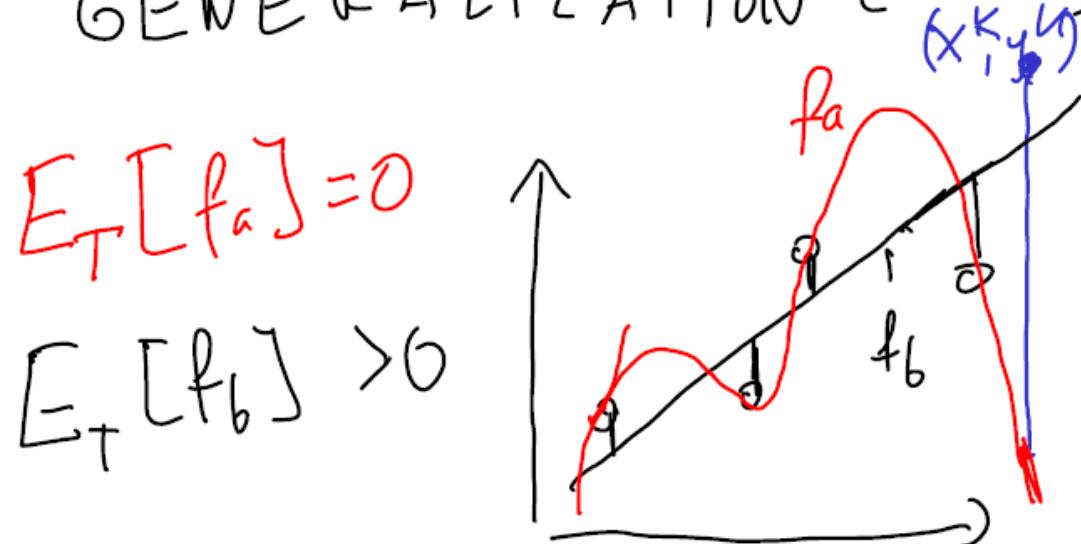
$$\forall \beta_{xy} \lim_{N \rightarrow \infty} \sup_{i \in \{1, \dots, N\}} |L[f_T(x_i), y_i] - L[\hat{f}_{T^i}(x_i), y_i]| = 0 \text{ in probability}$$



If we remove one point from our training sample how much does it perturb the loss



GENERALIZATION \leftrightarrow STABILITY

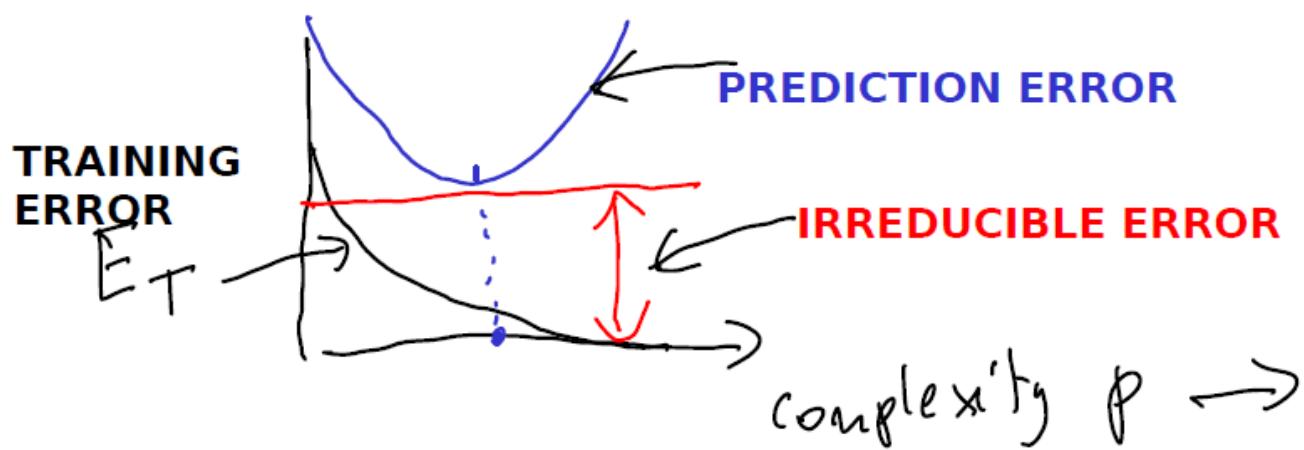


model complexity:

$$f_1(x) = \theta_0 + \theta_1 x$$

$$f_2(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

$E[f_a(x^k), y^k] \gg [f_b(x^k), y^k]$ given of polynomial complexity



Statistical Learning Theory

Generalization \leftrightarrow Stability

- V. Vapnik „Principles of Risk Minimization for Learning Theory ”([link](#))
- P. Poggio et al. „General conditions for predictivity in learning theory”,
Nature vol. 428 (2004) – ([link](#))

Prediction vs inference

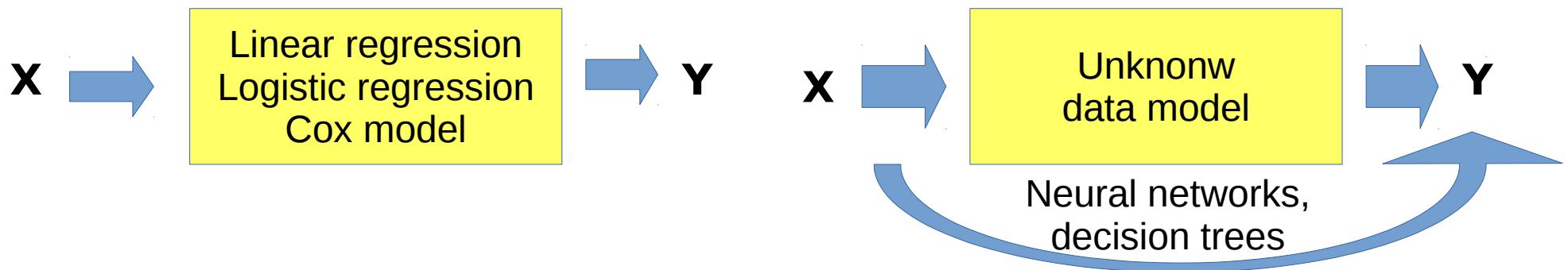


- **Prediction:** for a model $f_\theta(\mathbf{X})$ we want to **predict**, with the best possible accuracy output \mathbf{y}^i for the unknown vector \mathbf{X}^i
- **Inference** – for a model $f_\theta(\mathbf{X})$, where $\mathbf{X} = [x_1, x_2, x_3, \dots, x_p]$ we would like to **understand** how the data is generated, what are the relations between the input and output variables e.g. how the change of given input variable x_i (element of the vector X) affects the change of y

Prediction vs inference

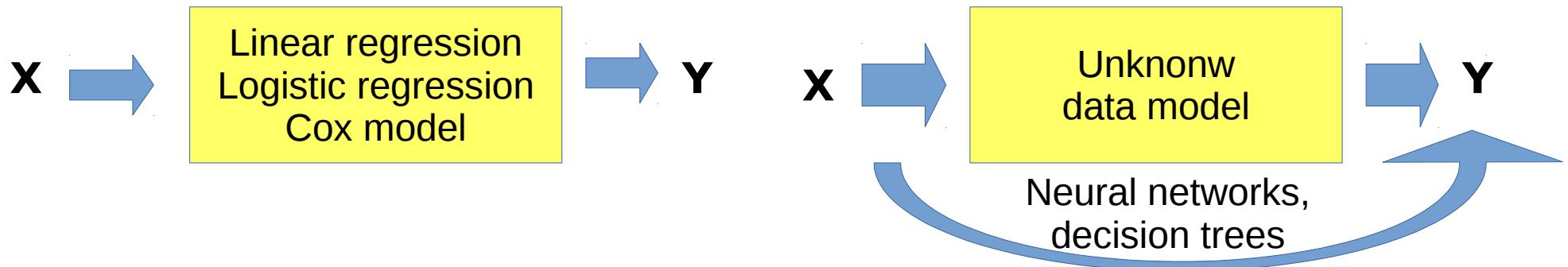
- **Statisticians' approach:** First we construct the statistical model that describes the **data generation**. Then we fit the parameters of the model and we can perform predictions and inference.
- **ML ('algorithmic') approach:** Who cares about how the data is generated ? The final goal is to have a an **algorithm** that accurately predicts the output.

Prediction vs inference



Leo Breiman „Statistical Modelling: The two cultures” (link)

Prediction vs inference



Accuracy vs Interpretability

Leo Breiman „Statistical Modelling: The two cultures” (link)

Quiz

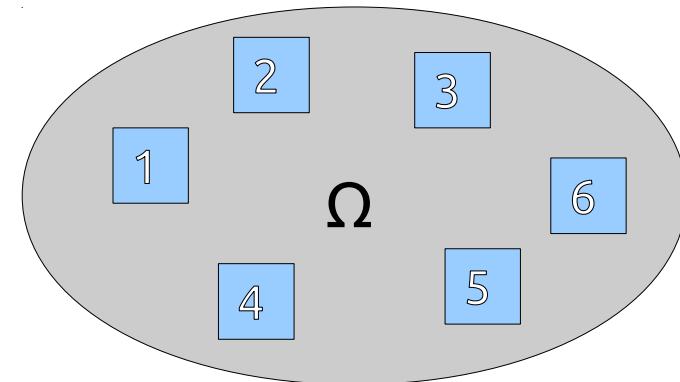
Crash course in probability (reminder)



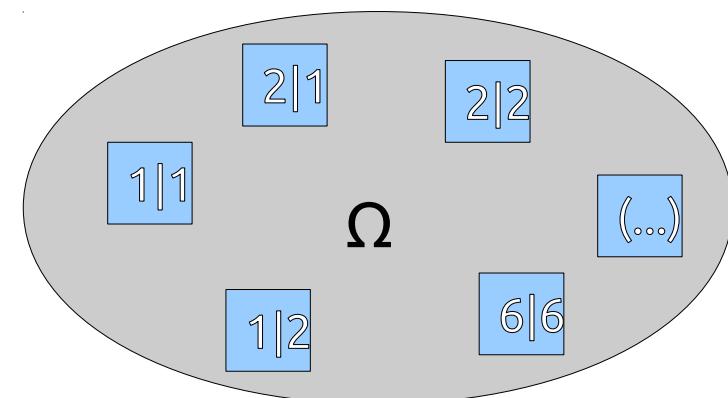
Probability reminder

Sample space Ω - set of all possible outcomes from a random experiment

Ex1. Rolling one die: $\Omega = \{1, 2, 3, 4, 5, 6\}$



Ex2. Rolling two dice: $\Omega = \{1|1, 1|2, 2|1, 2|2, \dots, 6|6\}$



Probability reminder

Event space F - set of events (outcome of experiments) E_i , being a subsets of Ω .

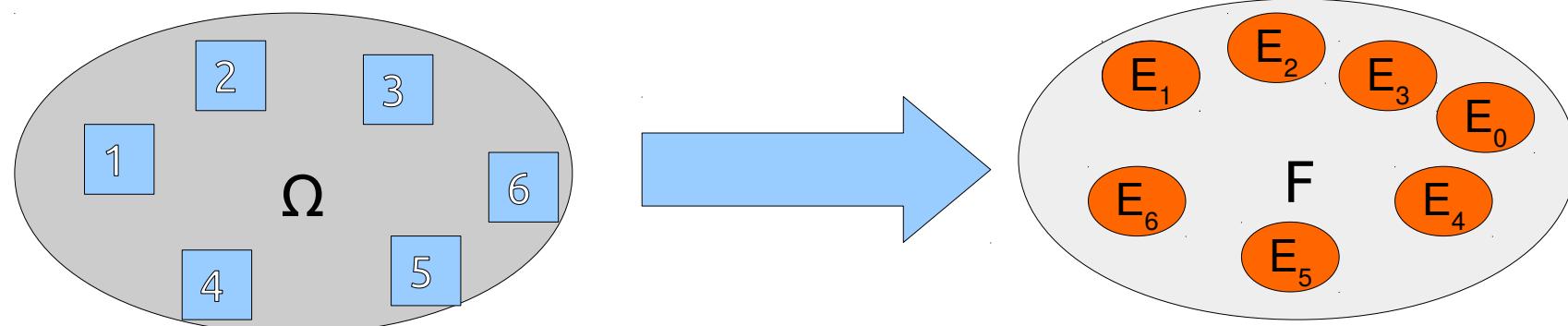
- We add an „impossible” event $E_0 = \{\emptyset\}$
- The way you define events is not unique

Ex3. Rolling one die: $\Omega = \{1, 2, 3, 4, 5, 6\}$

We are interested in probability of selecting given number:

$$E_1 = \{1\}, E_2 = \{2\}, E_3 = \{3\}, E_4 = \{4\}, E_5 = \{5\}, E_6 = \{6\}$$

$$F = \{E_1, E_2, E_3, E_4, E_5, E_6, E_0\}$$



Probability reminder

Event space F - set of events (outcome of experiments) E_i , being a subsets of Ω .

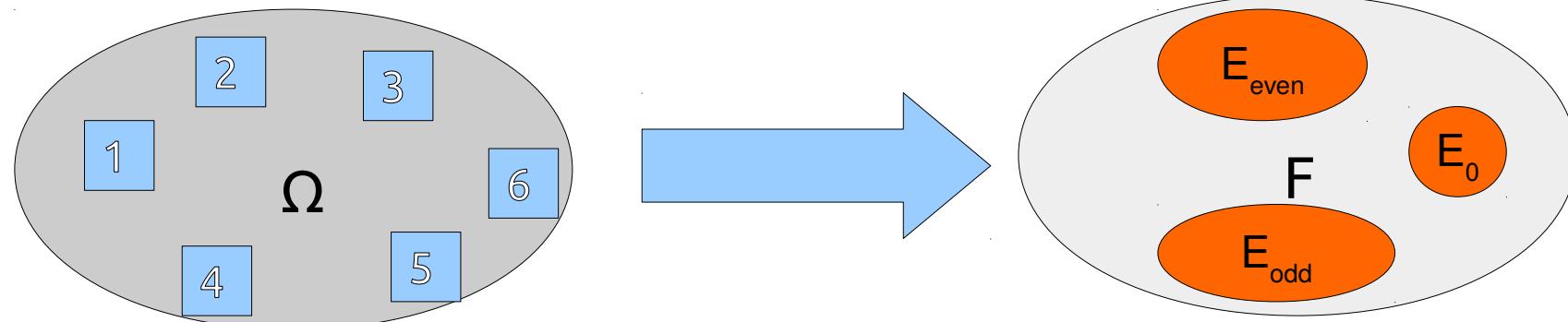
- We add an „impossible” event $E_0 = \{\emptyset\}$
- The way you define events is not unique

Ex4. Rolling one die: $\Omega = \{1, 2, 3, 4, 5, 6\}$

We are interested in probability of selecting even/odd number:

$$E_{\text{odd}} = \{1, 3, 5\}, E_{\text{even}} = \{2, 4, 6\}$$

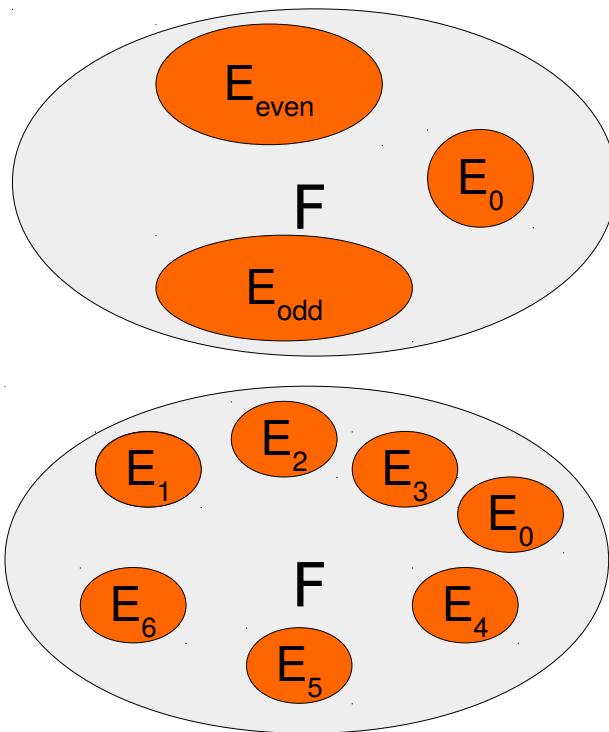
$$F = \{E_{\text{odd}}, E_{\text{even}}, E_0\}$$



Probability reminder

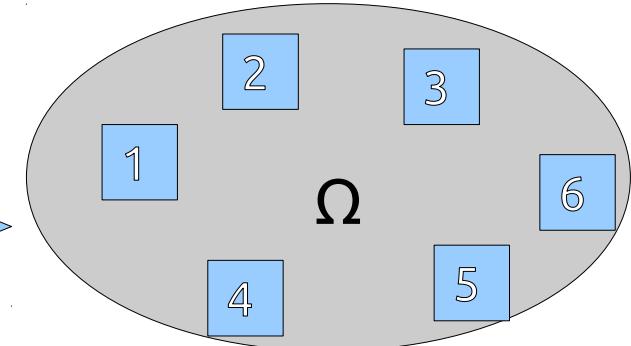
Event space \mathbf{F} - set of events (outcome of experiments) \mathbf{E}_i being a subsets of Ω .

- We add an „impossible” event $E_0 = \{\emptyset\}$
- The way you define events is not unique



$$\Omega = E_{\text{even}} \cup E_{\text{odd}}$$

$$\Omega = E_1 \cup E_2 \cup E_3 \cup E_4 \cup E_5 \cup E_6$$



***Remark:** Formally we should define σ -algebra $\sim F$ must be closed under complement and under countable unions and intersections

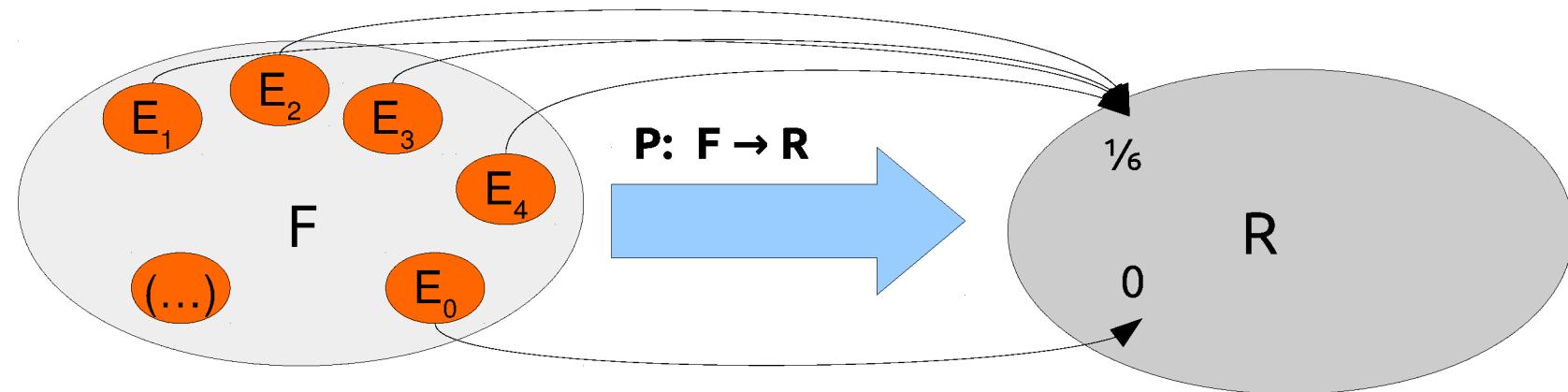
Problem Tossing a coin three times:

H - head,

T - tail

- What is the sample space Ω ?
- What is the event space F , for selection of at least two heads?

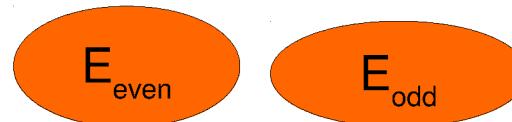
Probability measure P - Function that maps events from space \mathbf{F} to real numbers



Kolmogorov's Axioms:

- 1) Probability is non-negative: $P(E) \geq 0$
- 2) Probability of the entire sample set : $P(\Omega) = 1$
- 3) If E_1, \dots, E_n are disjoint events: $E_i \cap E_j = \emptyset$ for $i \neq j$

$$P(U E_i) = \sum P(E_i)$$



Some properties of probability measure

Properties: $E_i, E_j \in F$

- $P(\emptyset) = 0$
- $E_i \subseteq E_j \rightarrow P(E_i) \leq P(E_j)$
- $P(E_i \cap E_j) = \min(P(E_i), P(E_j))$
- $P(E_i \cup E_j) \leq P(E_i) + P(E_j)$
- $P(\Omega \setminus E) = 1 - P(E)$
- If E_1, \dots, E_n are disjoint events: $E_i \cap E_j = \emptyset$ for $i \neq j$ and $\bigcup E_i = \Omega \rightarrow \sum P(E_i) = 1$

Conditional probability

$A, B \in F$ such that $P(B) \neq 0$

Probability that event A happens “after” event B occurred:

$$P(A|B) = P(A \cap B)/P(B)$$

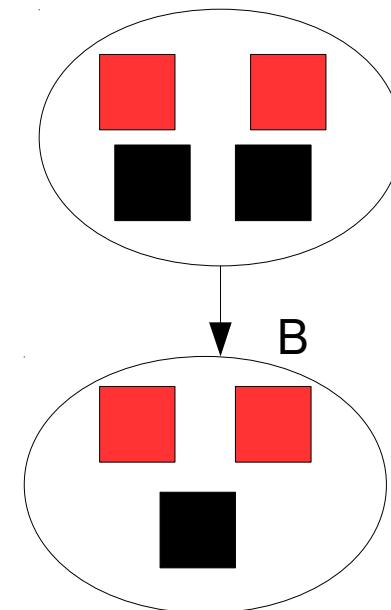
Ex5. Selecting a square from a bag with red (R) and black (B) squares

- A – selection of a red square
- B – selection of a black square

$$P(A) = \frac{1}{2}, P(B) = \frac{1}{2}$$

$A|B$ selection of a red square „after” selecting the black one:

$$P(A|B) = 2/3$$



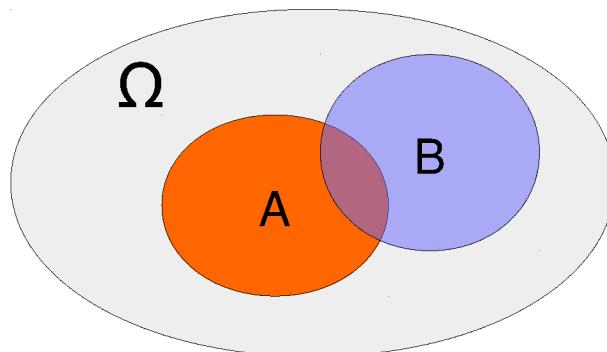
Conditional probability

$A, B \in F$ such that $P(B) \neq 0$

Probability that event A happens “after” event B occurred:

$$P(A|B) = P(A \cap B)/P(B)$$

- $P(A|B)$ fulfills Kolmogorov's axiom \rightarrow its a proper probability :-)



$$P(A|B) = P(A \cap B)/P(B)$$

- $P(A)$ can be treated as conditional $P(A) = P(A|\Omega)$

Independent events

$A, B \in F$ are independent if:

$$P(A \cap B) = P(A) * P(B) \Leftrightarrow P(A|B) = P(A), P(B|A) = P(B)$$

Ex6. Selecting (with returning) a square from a bag with red (R) and black (B) squares

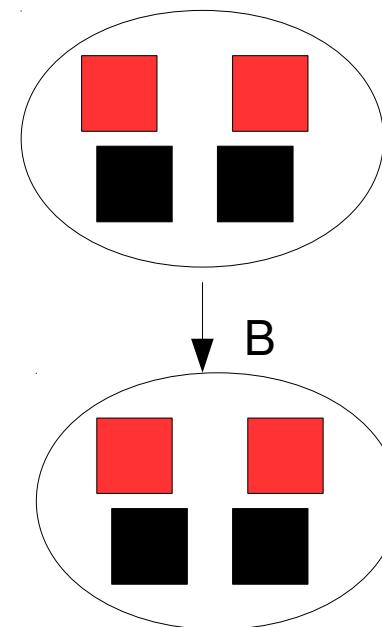
- A – selection of a red square
- B – selection of a black square

$$P(A) = \frac{1}{2}, P(B) = \frac{1}{2}$$

$A|B$ selection of a red square „after” selecting the black one:

$$P(A|B) = 2/4 = \frac{1}{2} = P(A)$$

A and B are independent



Bayes's theorem

- A,B - events from Event space

$$P(A|B) = P(A \cap B)/P(B) \rightarrow P(A \cap B) = P(A|B) * P(B)$$

$$P(B|A) = P(B \cap A)/P(A) \rightarrow P(B \cap A) = P(B|A) * P(A)$$



but, $P(A \cap B) = P(B \cap A) \rightarrow P(A|B) * P(B) = P(B|A) * P(A)$



$$P(A|B) = [P(B|A) * P(A)] / P(B) \text{ for } P(B) \neq 0$$

Bayes's theorem



Law of total probabilities:

$P(B) = \sum_i P(B|A_i) * P(A_i)$ for A_i being disjoint subsets
such that $\Omega = \cup A_i$

$$P(A|B) = [P(B|A) * P(A)] / P(B) \text{ for } P(B) \neq 0$$

$$P(A|B) = [P(B|A) * P(A)] / [\sum_i P(B|A_i) * P(A_i)] \text{ for } P(B) \neq 0$$

Problem :

$$\text{Show that } P(B) = \sum_i P(B|A_i) * P(A_i)$$

for A_i being disjoint sets, partitioning the whole sample space ($\Omega = \cup A_i$) and

$P(A_i) > 0$ for all i

Problem :

- A pharmaceutical company developed a test for detecting the rare disease, which is carried by 0.5 % cases of the whole population . Let's assume that the test gives the positive results for 96% of the cases if the patient is ill, but it also gives the positive results in 5% of the healthy patient. What is the probability that a patient is ill if his test gave a positive result?

Bayes's theorem



$$P(A|B) = [P(B|A) * P(A)] / P(B) \text{ for } P(B) \neq 0$$

posterior probability = [observation * prior probability] / evidence

How to interpret the posterior probability $P(A|B)$?

Experiment measured a mass of some particle B what can one say about the “true” mass A?

Frequentist vs Bayesian interpretations of probability

- **Frequentist approach:** probability is a **relative frequency** of occurrences

$P(A|B) = P(A \cap B)/P(B) = 95\% \rightarrow$ If we repeat the same experiment N times ($N \rightarrow \text{infinity}$), the ratio between occurrences of A together with B divided by the occurrences of B independently of A will be equal to 95%.

Frequentist vs Bayesian interpretations of probability

- **Frequentist approach:** probability is a **relative frequency** of occurrences

$P(A|B) = P(A \cap B)/P(B) = 95\% \rightarrow$ If we repeat the same experiment N times ($N \rightarrow \text{infinity}$), the ratio between occurrences of A together with B divided by the occurrences of B independently of A will be equal to 95%.

- **Bayesian approach:** probability is a measure of **degree of belief** of given hypothesis:

$P(A|B) = 95\% \rightarrow$ Probability that hypothesis is true if we measured(obtained) B is 95%.

$$P(\text{theory} | \text{data}) = P(\text{data} | \text{model}) * P(\text{model})$$

In most of the cases equivalent results independent of the approach.

Interpretation problems leads to war ...

DID THE SUN JUST EXPLODE?
(IT'S NIGHT, SO WE'RE NOT SURE.)

THIS NEUTRINO DETECTOR MEASURES
WHETHER THE SUN HAS GONE NOVA.

THEN, IT ROLLS TWO DICE. IF THEY
BOTH COME UP SIX, IT LIES TO US.
OTHERWISE, IT TELLS THE TRUTH.

LET'S TRY.

DETECTOR! HAS THE
SUN GONE NOVA?

(ROLL)

YES.



FREQUENTIST STATISTICIAN:

THE PROBABILITY OF THIS RESULT
HAPPENING BY CHANCE IS $\frac{1}{36} = 0.027$.
SINCE $p < 0.05$, I CONCLUDE
THAT THE SUN HAS EXPLODED.



BAYESIAN STATISTICIAN:

BET YOU \$50
IT HASN'T.



Physicists(at least particle) should be aware of this difference

- How to interpret unique phenomena in the frequentist framework e.g. Big Bang ?
- Different results can be obtained for upper limits calculations (depending on approach) → see Particle Data Group website

Further reading:

- B. Efron *Bayesians, Frequentists, and Scientists*

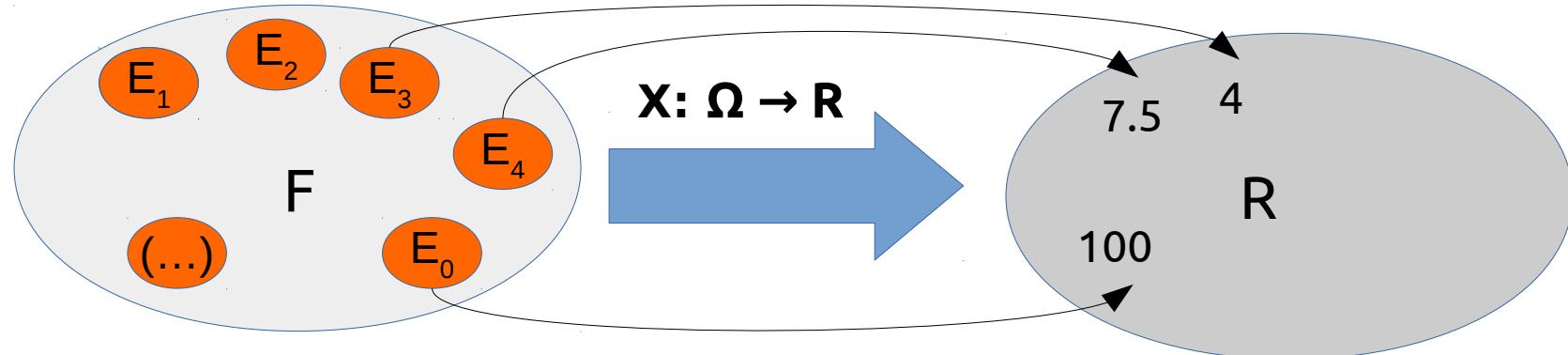
<http://statweb.stanford.edu/~ckirby/brad/papers/2005BayesFreqSci.pdf>

- L. Lyons, *Bayes and Frequentism: A particle physicist perspective:*

<https://arxiv.org/pdf/1301.1273.pdf>

Random variables

Random variable - function that maps events to some real numbers(*)

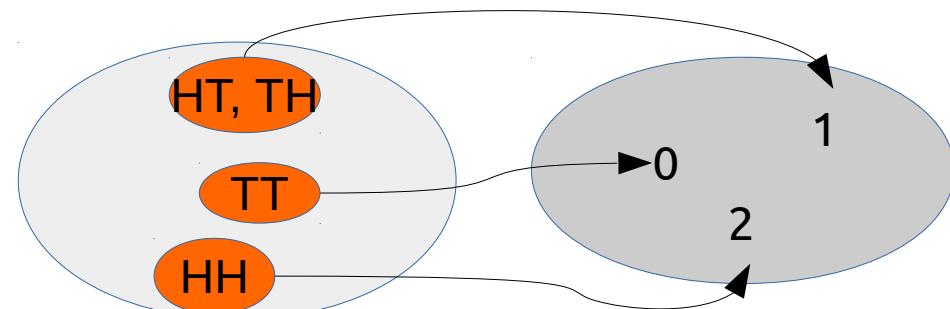


- Values are not limited to $[0,1]$ range as probability
- $X(E_i) = x$, where: X - random variable, x – value of the variable

Ex7. Tossing two coins: $\Omega = \{ HT, TH, HH, TT \}$

X is a random variable - number of heads

$$\begin{aligned} X(E) = 1 &\text{ if } E = \{ HT, TH \} \\ &0 \text{ if } E = \{ TT \} \\ &2 \text{ if } E = \{ HH \} \end{aligned}$$



*Remark: Formally we should allow only Borel functions

Random variables II

Random variables:

- **Discrete:** $P(X = k) = P(\{E: X(E) = k\})$
- **Continuous:** $P(a \leq X \leq b) = P(\{E: a \leq X \leq b\})$

Ex8. X numeric output of a die :

$$\Omega = \{1, 2, 3, 4, 5, 6\} \quad X = \{1, 2, 3, 4, 5, 6\}$$

Ex9. Y - height measurement from a group of people:

$$\Omega = \{P1, P2, P3, \dots\}, \quad Y = \{105 \text{ cm}, 183 \text{ cm}, 160,5 \text{ cm}, \dots\}$$

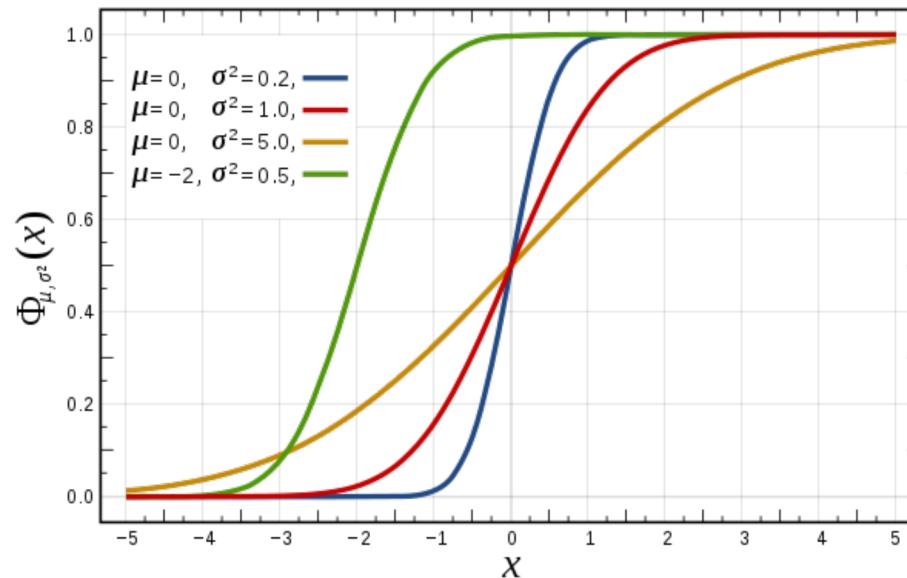
CDF

Cumulative Distribution Function (CDF) – describes probability measure for given random variable

- $F_x: \mathbb{R} \rightarrow [0,1]$
- $F_x(x) = P(X \leq x)$

Properties:

- $0 \leq F_x(x) \leq 1$
- $\lim_{x \rightarrow -\infty} F_x(x) = 0$
- $\lim_{x \rightarrow +\infty} F_x(x) = 1$
- $p \leq q \rightarrow F_x(p) \leq F_x(q)$

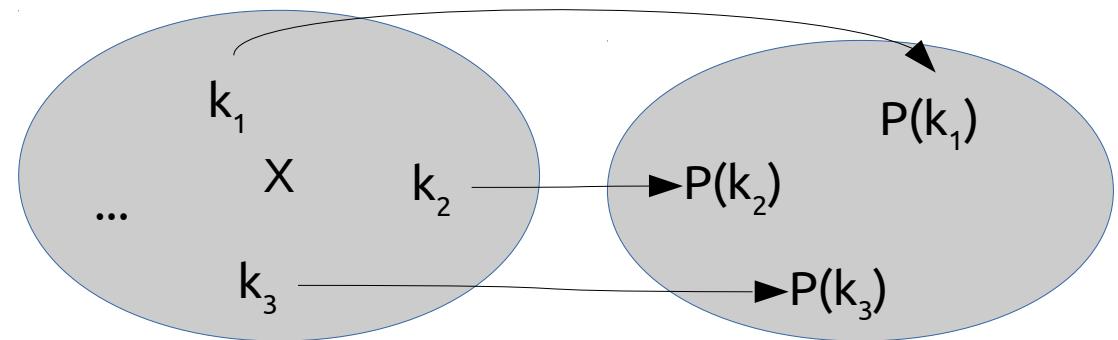


https://en.wikipedia.org/wiki/Cumulative_distribution_function

PMF

Probability Mass Function (PMF) – describes probability measure for a given discrete random variable:

- $p_x: C \rightarrow \mathbb{R}$
- $p_x(k) = P(X=k)$



Properties:

- $0 \leq p_x(x) \leq 1$
- $\sum p_x(k) = 1$ for all possible k
- $\sum p_x(k) = P(X \in A)$ for $k \in A$

PDF

Probability Density Function (PDF) – describes probability measure for a given continuous random variable (F_x must be a differentiable function):

- $f_x : \mathbb{R} \rightarrow \mathbb{R}$
- $f_x(x) = dF_x(x)/dx$

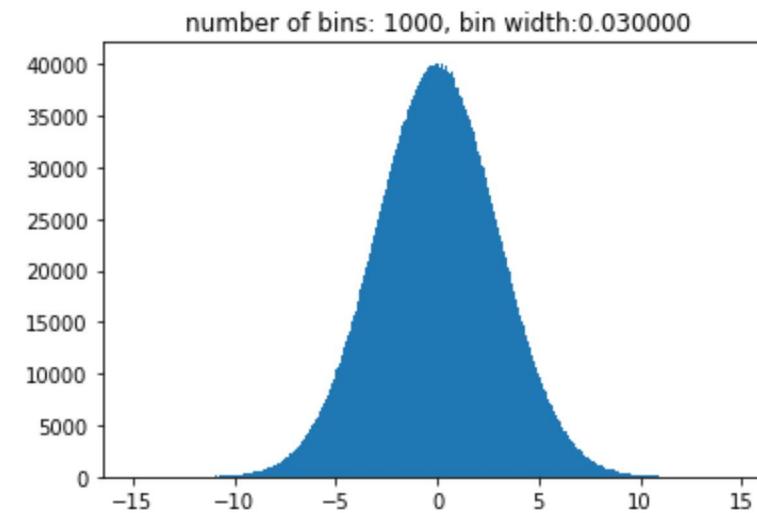
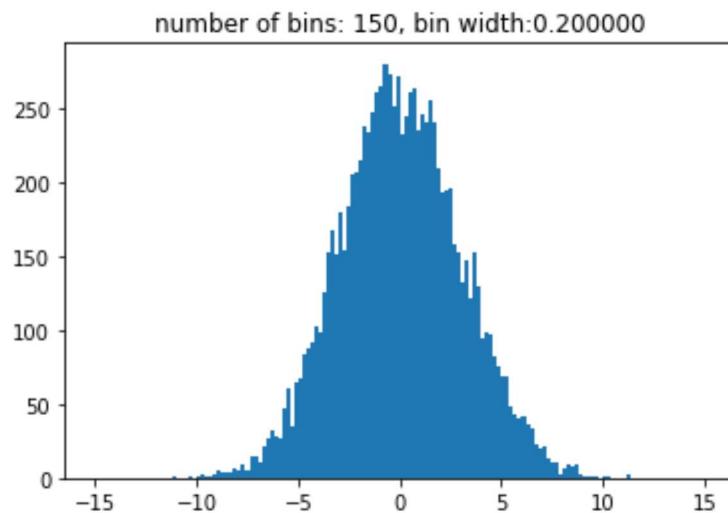
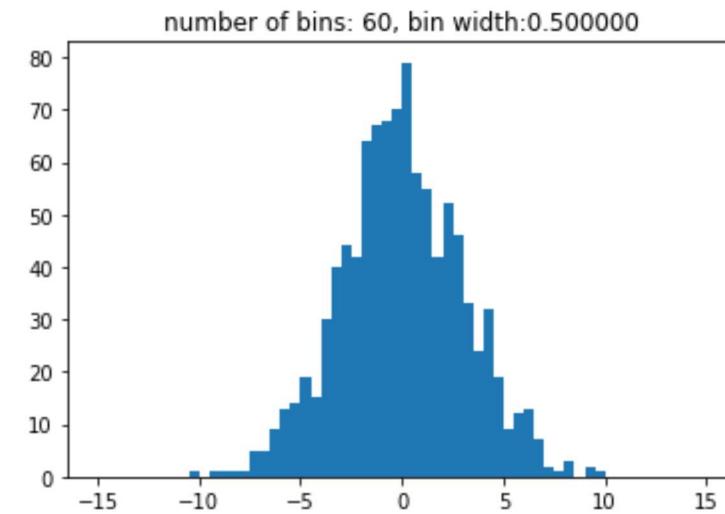
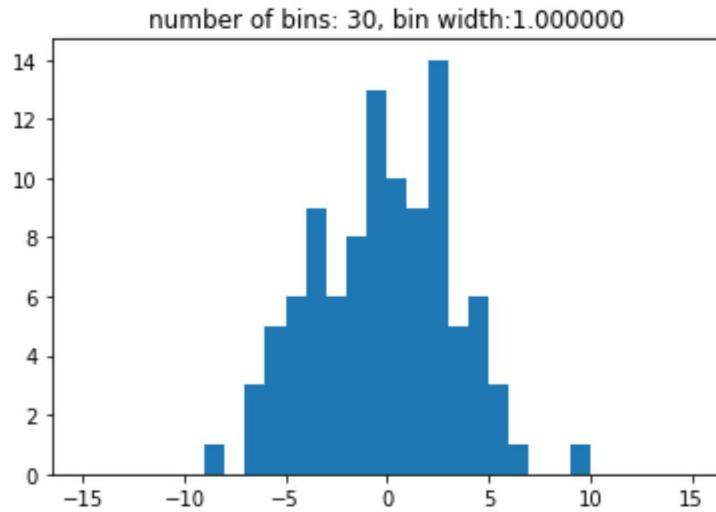
$$P(x < X < x + \Delta x) = f_x(x) \Delta x$$

Properties:

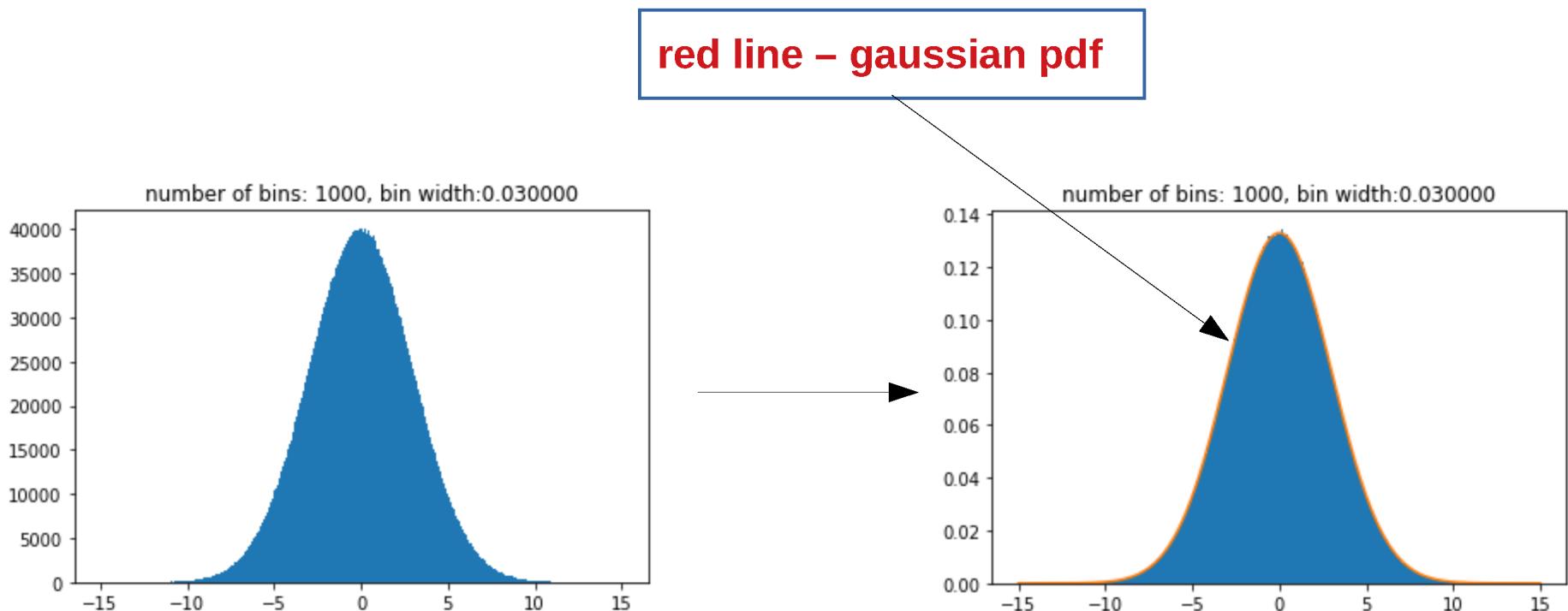
- $f_x(x) \geq 0$
- $\int f_x(x) dx = 1$ for $x -\infty$ to ∞
- $\int f_x(x) dx = P(X \in A)$ for $x \in A$

- **Remark:** f_x cannot express probability since can have values > 1
- **Remark2:** Especially, we know that for any continuous random variable $P(X=x) = 0$

Histograms vs PDFs



Histograms vs PDFs



Dividing every bin by the total number of events and by the bin width

Notebook:

https://github.com/wkrzemien/dataScienceAndML2020/blob/master/notebooks/statistics/pdfs_histograms.ipynb

Problem Tossing two coins: $\Omega = \{ HT, TH, HH, TT \}$

X is a random variable - number of heads

$$\begin{aligned} X(E) = & 1 \text{ if } E = \{ HT, TH \} \\ & 0 \text{ if } E = \{ TT \} \\ & 2 \text{ if } E = \{ HH \} \end{aligned}$$

- What are CDF and PMF functions?
- Draw it

Problem What are the PDF and CDF functions for the uniform distribution defined for $X = [0, a]$?

Problem *

Write a program that simulates the tossing of two coins, and estimate the CDF and PMF functions for the first problem above

Thank you