

Selected topics in Data Science



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

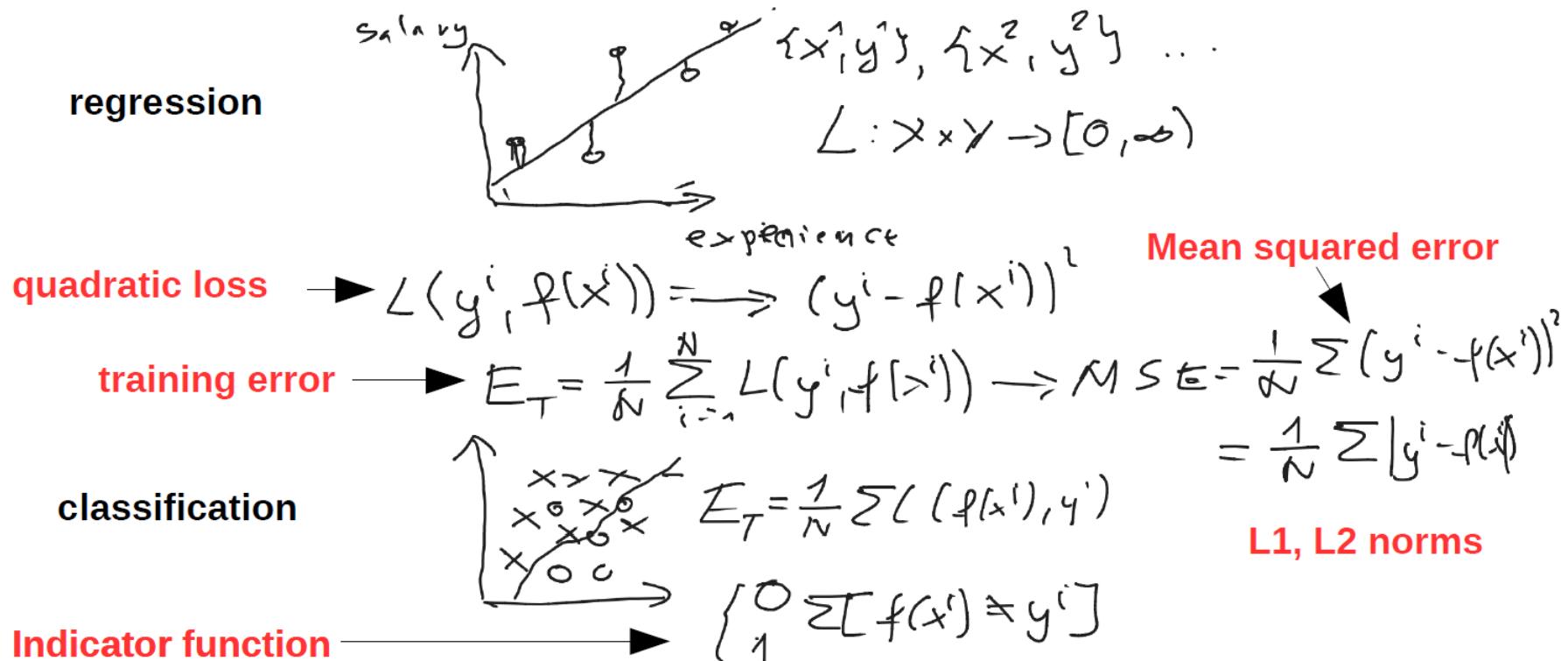
Wojciech Krzemień

30.10 2020, NCBJ

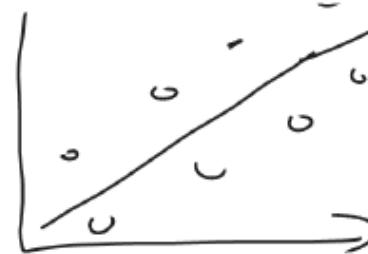
Last lecture recap

- Elements of Statistical Learning Theory:
 - loss function
 - empirical risk minimization
 - overfitting and model complexity
 - prediction vs inference
- Crash course in statistics I:
 - sample space, event space, probability measure
 - random variables, conditional probabilities and independent events
 - Bayes theorem
 - Cumulative Distribution Function (CDF), Probability Mass Function (PMF), Probability Density Function (PDF)

Loss function L to penalize the mistakes the model makes



Training expressed as **minimization** problem



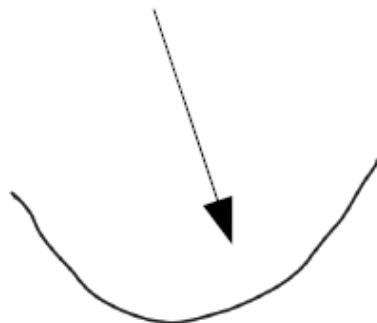
$$f_{\theta} = \theta_0 + \theta_1 x_1$$

Find parameters θ^* for which E_T has minimum

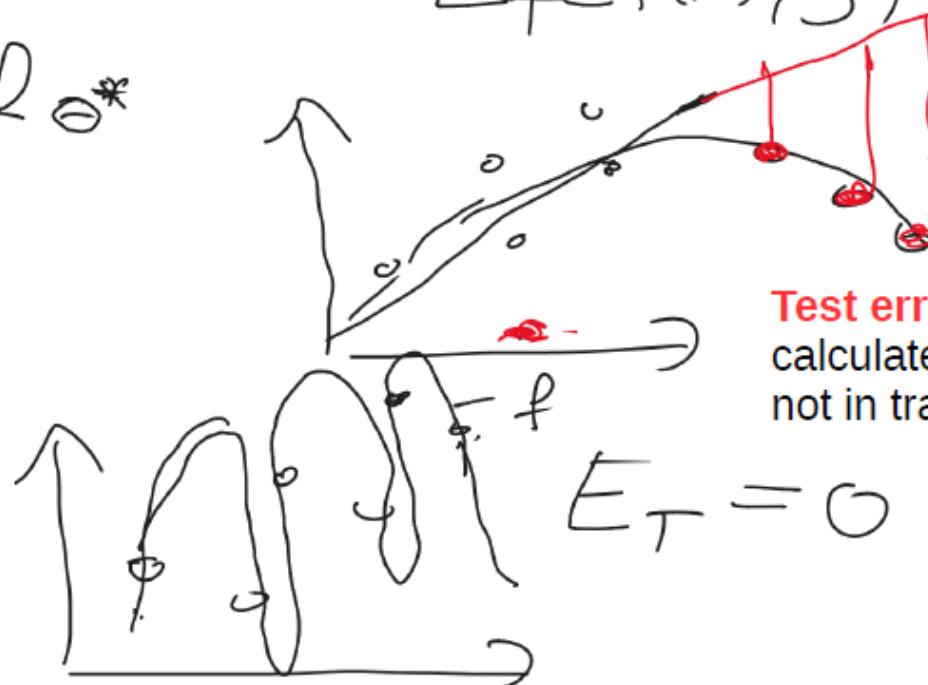
$$E[f(x), y]$$

$$f_{\theta^*}$$

Convex function



No more than one minimum



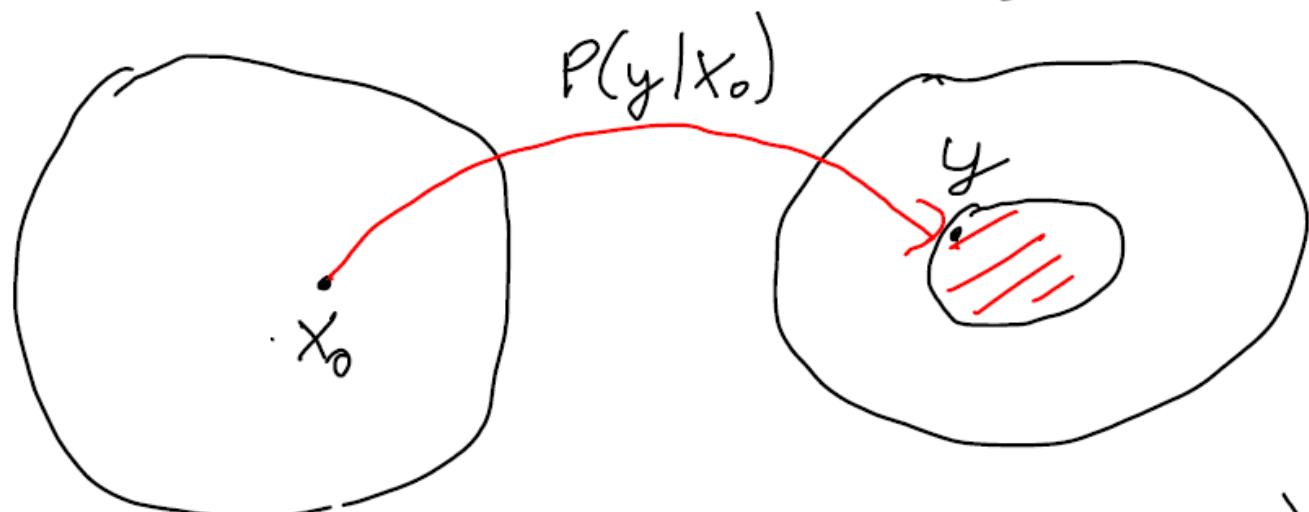
Test errors
calculated for points
not in training set

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

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

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

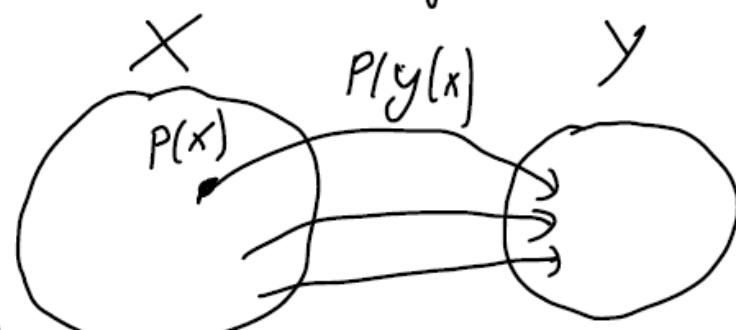
↑
Loss function

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.

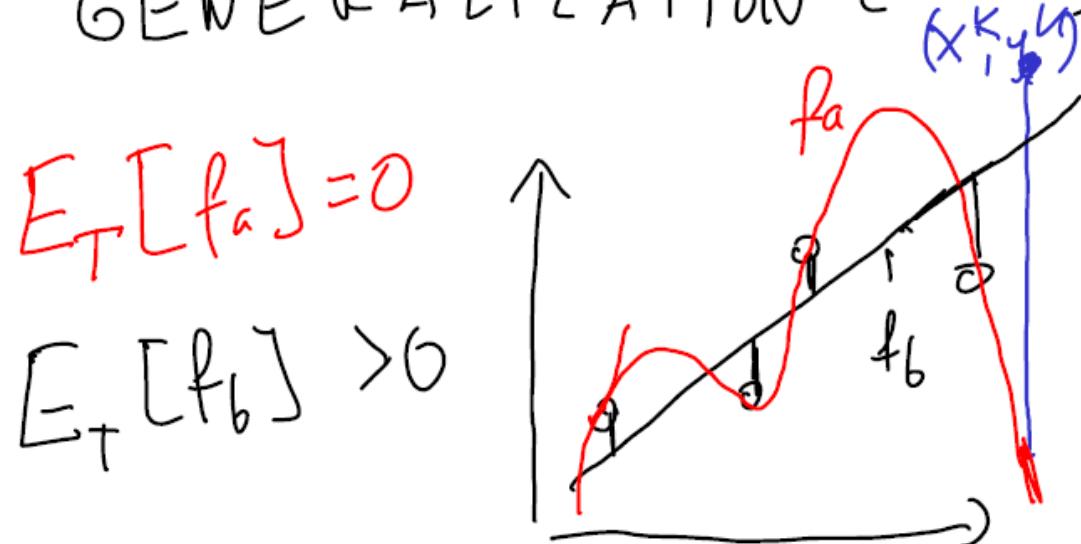


Since we don't know what is $s_{x,y}$ we cannot compute it

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

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

GENERALIZATION \leftrightarrow STABILITY

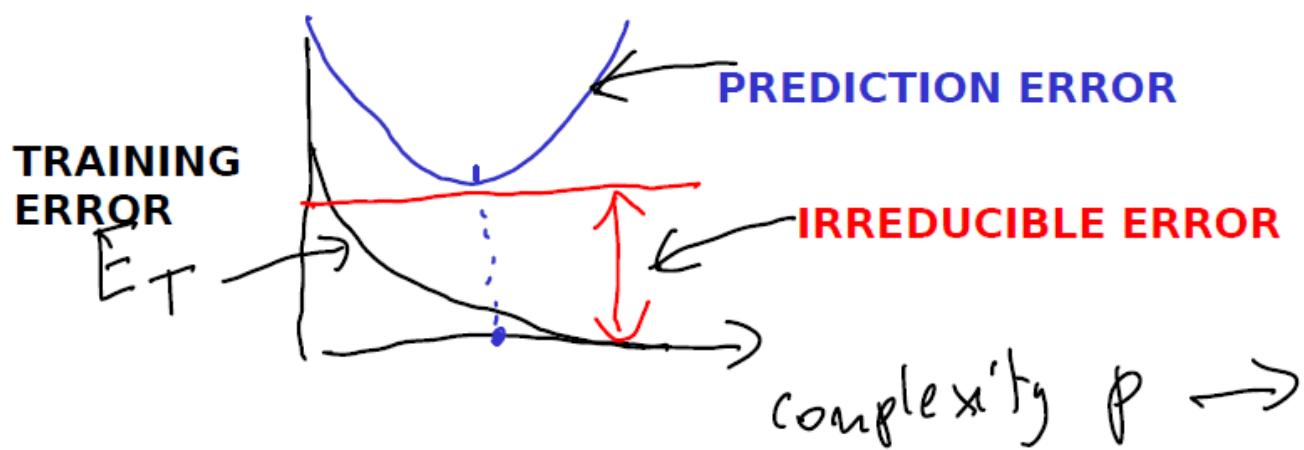


model complexity:

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

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

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



Todays plan

- Optimal algorithm and optimal classifier
- Bias-variance decomposition,
- k-NN revisited
- Crash course in Statistics II
- Work in groups:
 - Statistics problems
 - Small programs and k-NN implementations

Supervised vs unsupervised learning (Quiz)

Decide in each situation which kind of method would be better suited

- 1) A database with prices of the houses and its descriptions (e.g. area, location number of rooms). We want to predict prices for newly constructed houses.
- 2) A client database of the online shop. We want to group the customers according to their preferences.
- 3) A movie database with client profiles and a list of movies they watched, we want to develop a model to recommend movies for a given type of clients.
- 4) Samples of Internet traffic (e.g. packets), we would like to detect anomalous behaviours.

Supervised vs unsupervised learning (Quiz)

Decide in each situation which kind of method would be better suited

- 1) A database with prices of the houses and its descriptions (e.g. area, location number of rooms). We want to predict prices for newly constructed houses.

Supervised Learning

- 2) A client database of the online shop. We want to group the customers according to their preferences.

Unsupervised Learning

- 3) A movie database with client profiles and a list of movies they watched, we want to develop a model to recommend movies for a given type of clients.

Supervised Learning

- 4) Samples of Internet traffic (e.g. packets), we would like to detect anomalous behaviours.

Unsupervised Learning

Inference vs Prediction (Quiz)

Decide in each situation whether we are more interested in inference or in prediction.

- 1) The company launched an advertising campaign using different media e.g. TV, radio, or newspapers. The gathered data contains the sums of money invested in each sector and the rise in the income for different combinations. The boss would like to know what is the most optimal strategy e.g. is it worth advertising in all the sectors, is there any relation between the advertising on TV and radio, etc.

- 2) The company gathers data of their clients (e.g. age, education, preferences, and a lot more demographic information) and they want to promote their new product by a targeted mail campaign. The question is to decide which group of clients should be chosen.

Inference vs Prediction (Quiz)

Decide in each situation whether we are more interested in inference or in prediction.

- 1) The company launched an advertising campaign using different media e.g. TV, radio, or newspapers. The gathered data contains the sums of money invested in each sector and the rise in the income for different combinations. The boss would like to know what is the most optimal strategy e.g. is it worth advertising in all the sectors, is there any relation between the advertising on TV and radio, etc.

Inference / Both

- 2) The company gathers data of their clients (e.g. age, education, preferences, and a lot more demographic information) and they want to promote their new product by a targeted mail campaign. The question is to decide which group of clients should be chosen.

Prediction

Inference vs Prediction (Quiz)

Decide in each situation whether we are more interested in inference or in prediction.

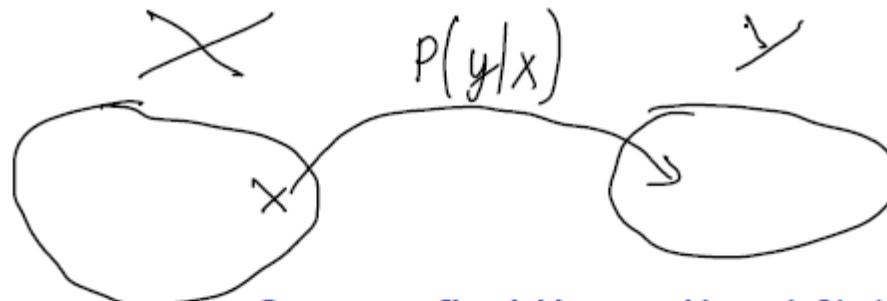
- 3) We analyze data from the university containing the results of students' exams together with the information about students and the courses that they previously took. We would like to understand what is the relation between the final scores they got and the other factors

Inference vs Prediction (Quiz)

Decide in each situation whether we are more interested in inference or in prediction.

- 3) We analyze data from the university containing the results of students' exams together with the information about students and the courses that they previously took. We would like to understand what is the relation between the final scores they got and the other factors

Inference



$$p(x,y) = p(y|x) \cdot p(x)$$

$$S_{x,y} = S_{y|x} \cdot S_x$$

Can we find the optimal $f(x)$, given quadratic loss function?

$$L[f(x), y] = (f(x) - y)^2$$

$$EPE[f] = \text{SSL}(f(x), y) \cdot S_{x,y} dx dy$$

$$\text{SSL}(f(x) - y)^2 S_{y|x} \cdot S_{x,y} dx dy =$$

$$\int S_x \int_y (f(x) - y)^2 S_{y|x} \cdot dx dy$$

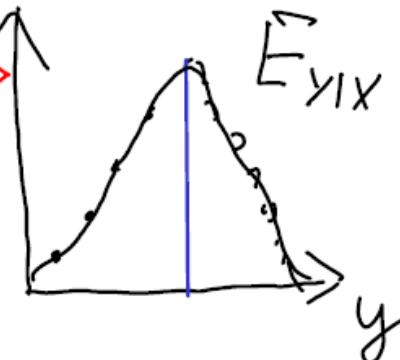
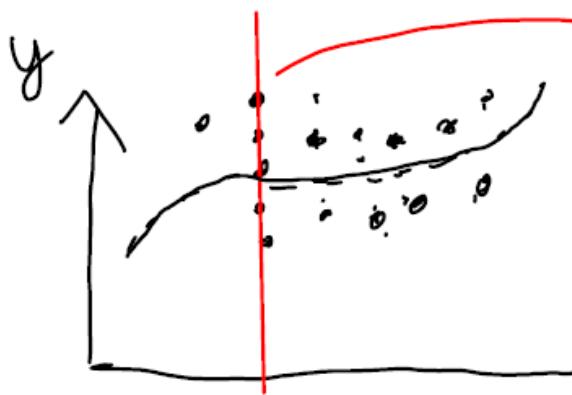
$$E[g(x)] \stackrel{\text{def}}{=} \int g(x) S_x dx$$

Expectation value
over the whole
feature space

$$\int_{\mathcal{X}} \int (f(x) - y)^2 p_{y|x}(y|x) dy dx$$

$$E_x [E_{y|x} [(f(x) - y)^2 | x]]$$

Example of the conditional expectation



It can be approximated by the average for whole y's points.

$$E_{y|x} [y|x=x_0] \approx \frac{1}{N} \sum_{i=1}^N y_i$$

x_0

We will treat the problem point-wise.
We search for the optimal solution in $x=x_0$.

$$L = E_{y|x} [(f(x_0) - y)^2 | x=x_0]$$

$$f(x_0) = c \text{ constant}$$

We want to find the minimum of L with respect to c

$$\frac{d}{dc} L = 0$$

$$\frac{d}{dc} E_{y|x} [(c - y)^2 | x=x_0] = 0$$

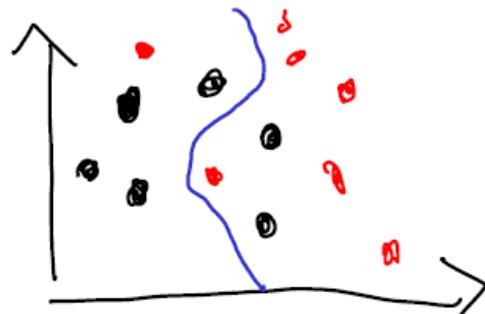
$$C \equiv f^* = E_{y|x} [y | x=x_0]$$

Regression Function

This is an optimal function for regression assuming MSE loss (!)

If we take $L(f(x), y) = |f(x) - y|$ one can find that minimization of $EPE[f]$ gives $f^* = \text{Median}(y, x=x_0)$

Classification



$$L(f(x), y) = [f(x) \neq y]$$

(1 if $f(x) \neq y$
0 if $f(x) = y$)

Again we're trying to find 'the optimal classifier' that minimizes EPE for indicator loss

$$EPE[f] = \iint L(f(x), y) g_{y|x} dy dx$$

This is the classification problem so:

$$\sum_y \text{obj} \rightarrow \sum_{k=1}^K \text{sum over all output classes}$$

Discrete output over k classes:

e.g. for binary $y \in \{0, 1\}$ classifier

$$S_{y|x} \rightarrow p(k|x) \sum_{i=1}^2$$

$$EPE[f] = E_x \left[\sum_{k=1}^K p(k|x) \cdot [f(x) \neq y] \right]$$

$$[f(x) \neq y] = 1 - [f(x) = y]$$

(...) After treating the problem point-wise $x = x_0$, we get

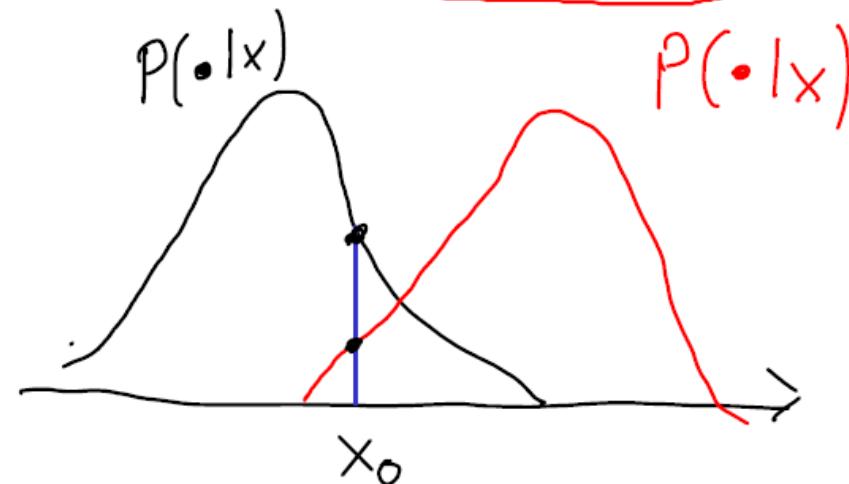
the solution

$$f^* = \arg \max_k P(k | x = x_0)$$

BAYES
CLASSIFIER

Bayes classifier is "the optimal classifier" -> with the smallest error rate.

$$f^* = \operatorname{argmax}_k P(k|x=x_0)$$



$\operatorname{argmax}[x]$ for $f(x)$ is the notation to express x value for which $f(x)$ has the maximum value:

$$\text{e.g. } \operatorname{argmax}_x \sin(x) = \frac{\pi}{2}$$

e.g. $P(\cdot|x=x_0) = 70\%$ } $P(\cdot|x=x_0) = 30\%$ } \rightarrow We choose the class with higher probability
so •, For 100 cases we count 30% errors.
(misclassifications)
 $E_{\text{Rate}} = 30\%$, but it's still the lowest possible!!!

Remark 1: Don't mix the Bayes classifier with the Naive Bayes Classifier.

Remark 2: Since normally we don't know $P(y|x)$ we cannot use it.

Optimal f^* model (algorithm)

- Let's assume that we know $\rho_{XY}(x,y)$
- $f^*(X)$ - optimal algorithm for quadratic loss function $L(f(X), Y) = (f(X) - Y)^2$
- One can show that the optimal solution is:

$$f^*(X) = E_{Y|X}[Y|X]$$

Problem:

Show that $f^*(X) = E_{Y|X}[Y|X]$ is optimal assuming $L(f(X), Y) = (f(X) - Y)^2$

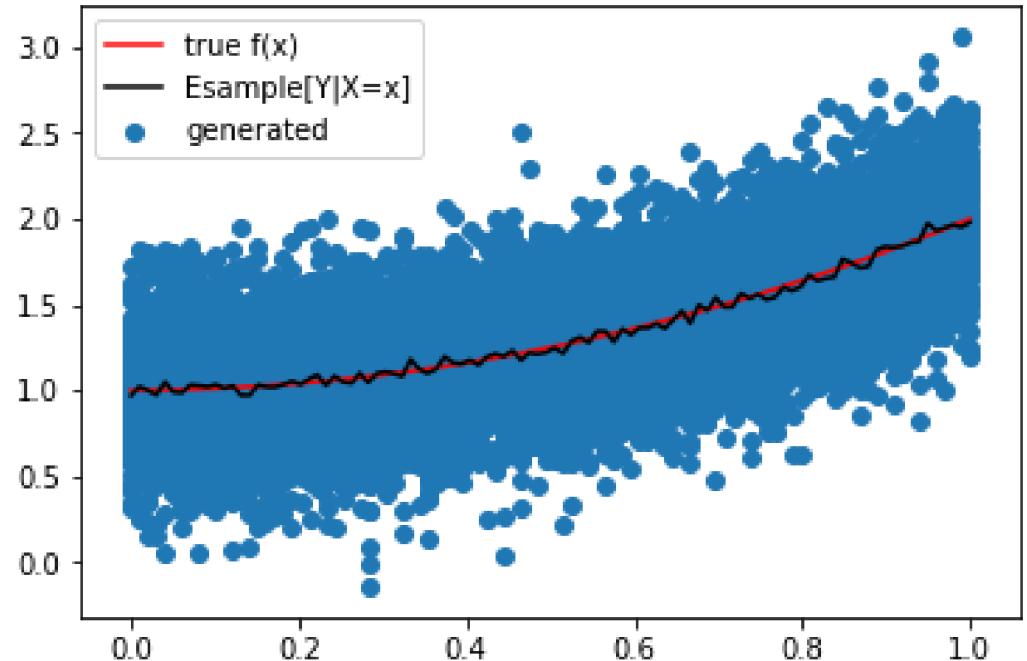
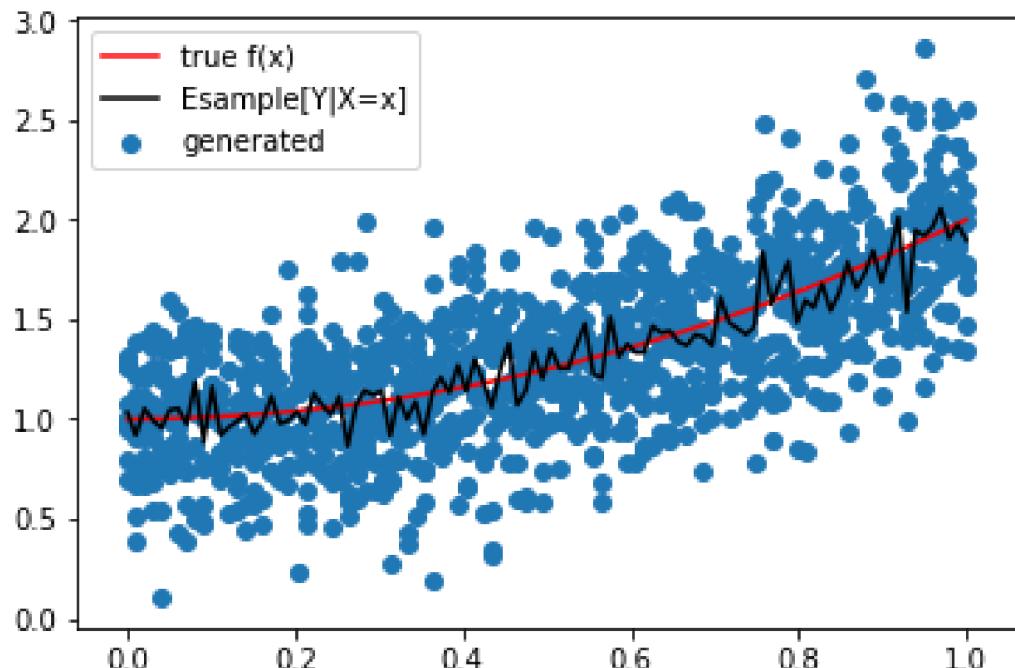
$$EPE(f) = E_{X,Y}[L(f(X), Y)] = \iint L(f(x), y) \rho_{XY}(x, y) dx dy$$

Hints:

1. Express ρ_{XY} as conditional $\rho_{Y|X}$
2. Express EPE as conditional expectations
3. Minimize the expression point-wise (for given $X = x_0$ from the feature space)

(*) regression function defined point-wise for every point x_0 of feature space X

Optimal f^* model (regression function)



$$f^*(x) = E_{Y|X}[Y|X]$$

(*) regression function defined point-wise for every point x_0 of feature space X

Problem:

- Derive the optimal algorithm $f^*(X)$ assuming $L(f(X), Y) = |(f(X) - Y)|$

Problem:

- Derive the optimal algorithm $f^*(X)$ assuming $L(f(X), Y) = (f(X) - Y)^2$

Minimizing $EPE[f]$

$$L_2 \quad L = (f(x) - y)^2$$

Indicator loss for classification

$$L_1: L = |f(x) - y|$$

$$L = [f(x) \neq y]$$

$$f^* = E_{y|x} [y|x=x_0]$$

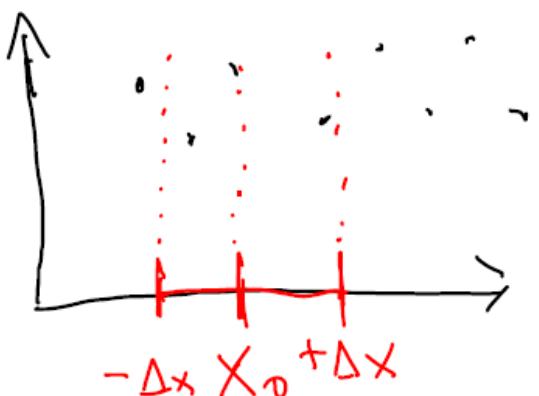
$$f^* = \text{Median}(y|x_0)$$

$$f^* = \arg \max_k P(k|x=x_0)$$

REGRESSION FUNCTION

All f^* are defined point-wise, assuming that we know $S_{Y|X}$ or $p(y|x)$.

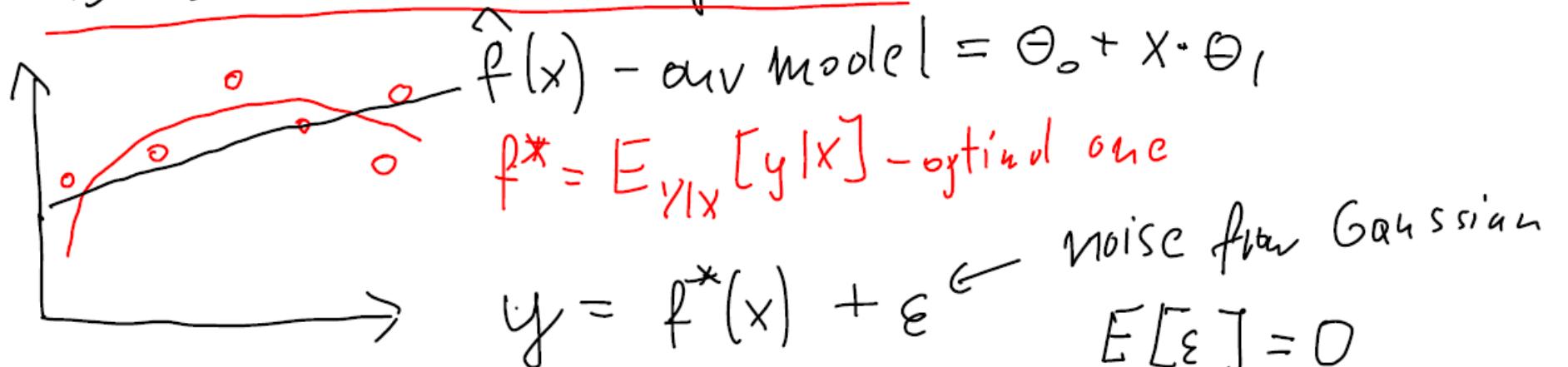
BAYES CLASSIFIER



$$f^* = E_{y|x} [y|x=x_0] \rightarrow \text{Average } y_i \\ x \in [x - \Delta x, x + \Delta x]$$

The k-NN for regression is the approximation of $E_{y|x} [y|x=x_0]$

Bias - Variance decomposition



We decompose the error for given $x = x_0$

$$E_y[(y - \hat{f}(x_0))^2] \quad [x = x_0]$$

|| (show it!)

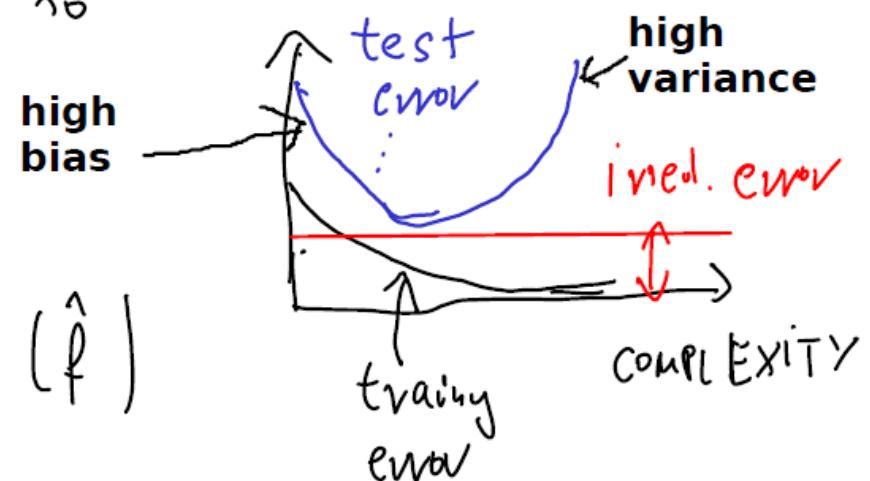
$$\text{VAR}(\epsilon) + \text{BIAS}^2(\hat{f}) + \text{VAR}(\hat{f})$$

$$(\hat{f} - f^*)^2$$

Irreducible error

How far we are from
the optimal solution

How much or model
would change for different
training sets



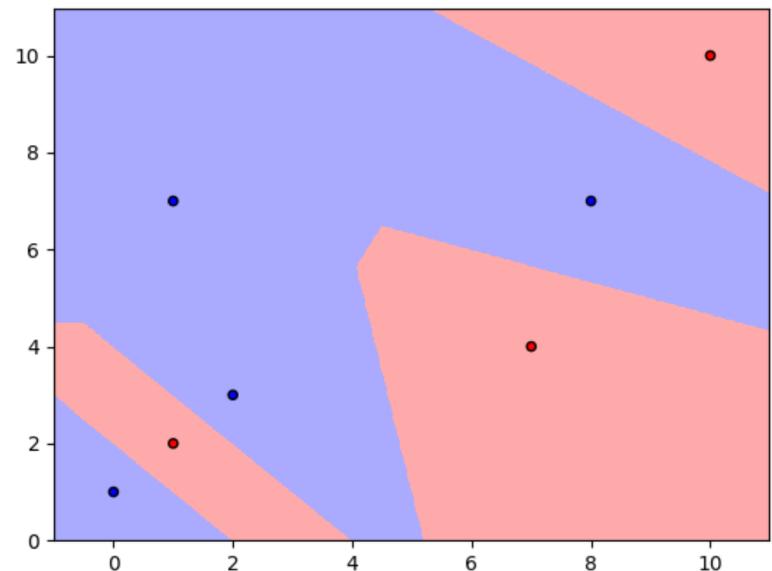
Notebook:

https://github.com/wkrzemien/dataScienceAndML2020/blob/master/notebooks/slt/EPE_train_test_errors.ipynb

k-Nearest Neighbors -model

X – space of our features, Y – space of outputs/labels of classes:

- Training set T consists of pairs of $\{X^{(i)}, Y^{(i)}\}$ $i=1,\dots,N$
- Every X is feature vector with P coordinates $X=[x_1, x_2, \dots, x_P]$
- k- number of Neighbors
- $\text{dist}(X^i, X^j)$ is a distance metric in feature space X



To classify a new X^* instance:

- we choose k closest neighbors from T e.g. $N_k(x)$
- We calculate the average output Y (e.g. majority vote)
 $y(x) = 1/k * \sum Y^{(j)}$ over j for which $x^{(j)}$ belong to $N_k(x)$
- We assign (predict) the label of X^*

- Problem Let's assume that we use kNN model to perform the classification of two type of objects (think of spam/non-spam) with $k=2$, $k=3$ and $k=10$ and $N=100$ training samples with 2 features each e.g. $\{ X^{(1)}=[0,1]^T, X^{(2)}=[-1,2]^T, X^{(3)}=[2,2]^T, X^{(4)}=[0,0.3]^T, \dots, X^{(100)}=[-100,12,5]^T \}$ and $\{ Y^{(1)}=[0], Y^{(2)}=[1], Y^{(3)}=[1], Y^{(4)}=[0] \}$
- For which k value do we expect the smallest training error?
- For which k value do we expect the highest/smallest stability?
- How would we classify $X=[0,0]$ if we use first 4 training samples for $k=1$, $k=2$, $k=3$? Calculate the majority vote result for each case.
- We increase the training sample to $N=101$. How we expect to affect the stability for different k values?
- Now instead of kNN model we use the linear regression. Discuss the stability issue.

Crash course in probability II



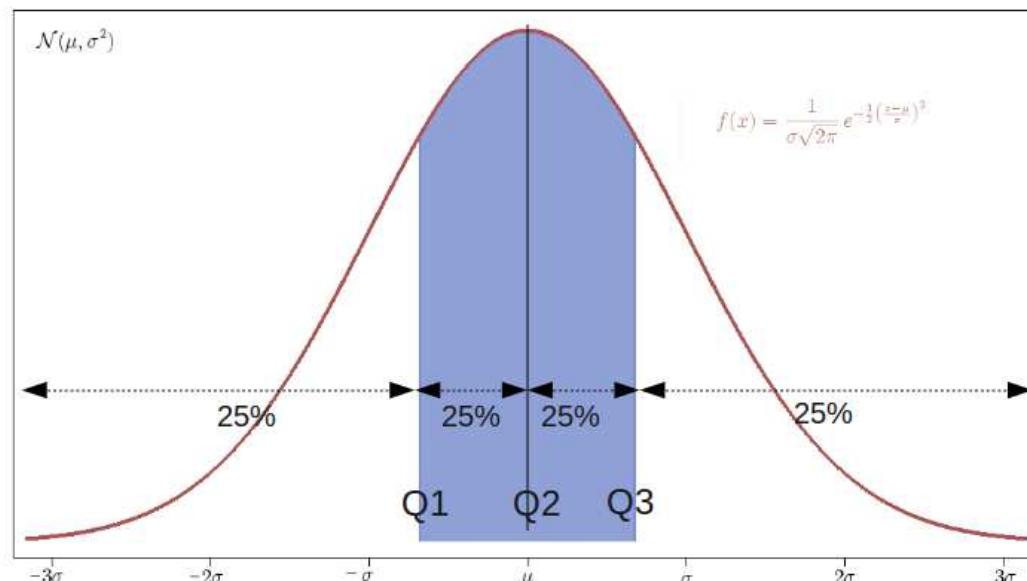
Quantiles

Quantile of the order q (q-quantile) x_q

cut points dividing the sample (distribution) into sub samples with equals probabilities:

$$F(x_q) = q \text{ for } q \in [0,1] \rightarrow F^{-1}(q) = x_q \text{ if (cdf exists)}$$

Median - $x_{1/2}$ for which there is an equal probability that x gets lower and higher value.



Expectation / Expected value

Let $g(x)$ be some arbitrary function: $R \rightarrow R$, X – random variable, $g(X)$ becomes random variable as well:

• For X being a discrete random variable with PMF $p_x(k)$:

$$E[g(X)] = \sum p_x(k) g(k) \text{ for all possible } k \text{ from } X$$

• For X being a continuous random variable with PDF $f_x(x)$:

$$E[g(X)] = \int g(x) f_x(x) dx \text{ for } x \text{ from } -\infty \text{ to } \infty$$

Properties:

• $E[f(X) + g(X)] = E[f(X)] + E[g(X)]$ - linearity

• $E[K] = K$ for K being constant

• $E[K * f(X)] = K * E[f(X)]$ for K being constant

- **Intuition:** $E[g(X)]$ is a kind of weighted average with weights defined by p_x or f_x
- **Remark:** The equations for $E[g(X)]$ are also known as the law of the unconscious statistician (LOTUS)

Variance

Variance of a random variable X expresses its concentration around the mean:

$$\text{Var}[X] = E[(E[X] - X)^2] \rightarrow \text{Var}[X] = E[X^2] - E^2[X]$$

Properties:

- $\text{Var}[K] = 0$ for K being constant
- $\text{Var}[K * f(X)] = K^2 * \text{Var}[f(X)]$ for K being constant

Problem Let $g(X) = 1$ for some set A being a subset of sample space Ω :

- What is $E[g(X)]$ if X is discrete with a given PMF or continuous with a given PDF

Problem

- What is $E[g(X)]$ for $g(X) = x$

Problem

- Show that $\text{Var}[X] = E[X^2] - E^2[X]$

Problem

- Calculate the mean and the variance for the uniform distribution

Problem:

- Implement a function that returns mean
- Implement a function that returns $\text{Var}[X]$

Problem: Similarity measures

- Implement a function that returns Euclidean distance between two vectors
- Implement a function that returns Manhattan distance between two vectors

Joint and marginal CDF

Let's consider two random variables X, Y with CDF's: $F_x(x)$, and $F_y(y)$, The joint CDF of X and Y is defined as:

$$F_{xy}(x,y) = P(X \leq x, Y \leq y)$$

F_x and F_y are called marginal CDFs and are related to F_{xy}

$$\cdot F_x(x) = \lim_{y \rightarrow \infty} F_{xy}(x,y)$$

$$\cdot F_y(y) = \lim_{x \rightarrow \infty} F_{xy}(x,y)$$

Remark: Since joint CDF is defined as probability, then obviously its values must be limited to range [0, 1].

Joint and marginal PMF and PDF

Let's consider two random variables X, Y with:

- (discrete) PMF's: $p_X(k)$, and $p_Y(r)$, The joint PMF is defined as:

$$p_{XY}(k,r) = P(X=k, Y=r)$$

- (continuous) with F_{XY} differentiable for all x,y values. The joint PDF is defined as:

$$f_{XY}(x,y) = \partial^2 F_{XY}(x,y) / \partial x \partial y$$

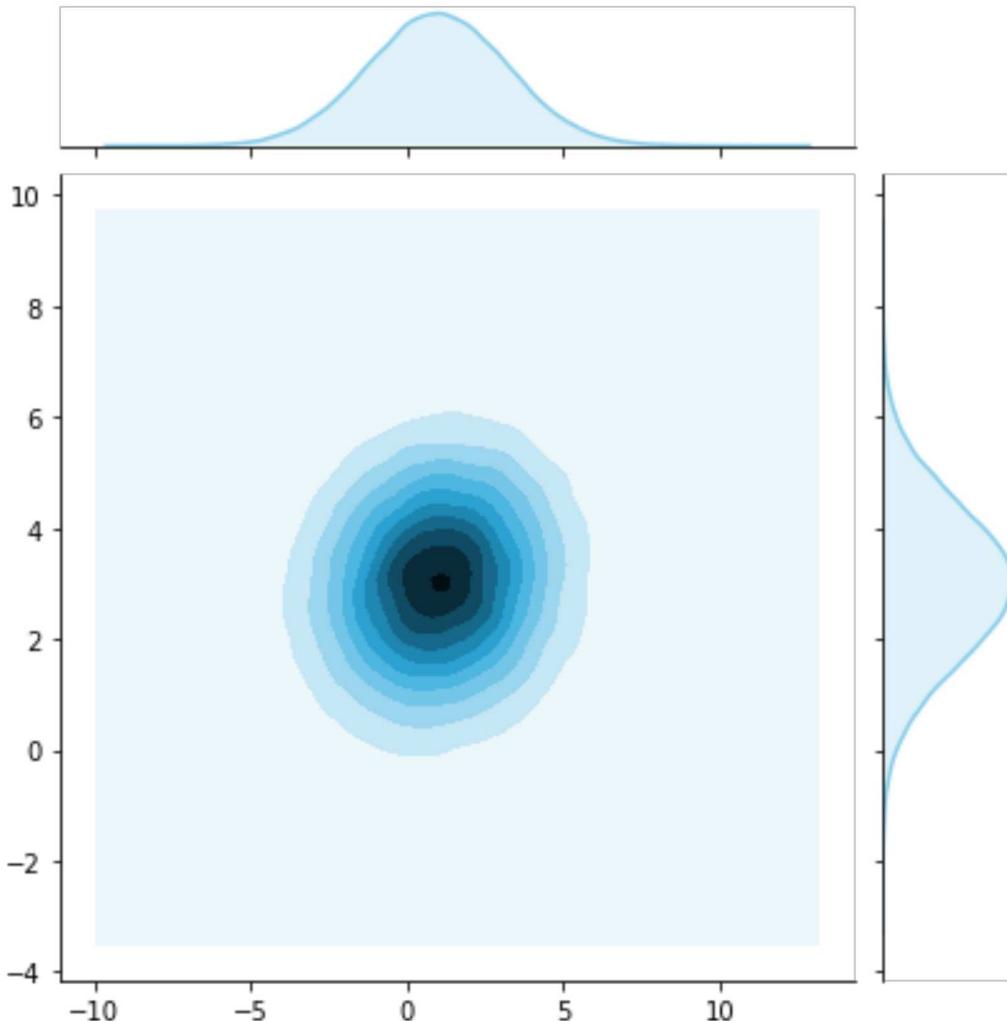
$$\iint f_{XY}(x,y) dx dy = P((X,Y) \in A) \text{ for } x,y \in A$$

p_X/f_X and p_Y/f_Y are called marginal PMFs/PDFs and are related to p_{XY}/f_{XY} :

$$\cdot p_X(k) = \sum p_{XY}(k,r) \text{ for all } r \quad f_X(x) = \int p_{XY}(x,y) dy \text{ for } y -\infty \text{ to } \infty$$

$$\cdot p_Y(r) = \sum p_{XY}(k,r) \text{ for all } k \quad f_Y(y) = \int p_{XY}(x,y) dx \text{ for } x -\infty \text{ to } \infty$$

Bivariate normal example



$$\begin{aligned}\mu_x &= 1 \\ \mu_y &= 3 \\ \text{cov} &= \begin{vmatrix} 5 & 0.3 \\ 0.3 & 2 \end{vmatrix} \\ \rho &= \text{cov}/(\sigma_x * \sigma_y)\end{aligned}$$

$$f(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)} \left[\frac{(x-\mu_X)^2}{\sigma_X^2} + \frac{(y-\mu_Y)^2}{\sigma_Y^2} - \frac{2\rho(x-\mu_X)(y-\mu_Y)}{\sigma_X\sigma_Y} \right]\right)$$

Equation taken from Wikipedia: https://en.wikipedia.org/wiki/Multivariate_normal_distribution

Problem Tossing a coin: $\Omega = \{ H, T \}$ and rolling strange die $\Omega = \{ 1, 2, 3 \}$

If H we roll the die twice, if T we roll the die once

- Calculate joint PMF
- Calculate marginal PMF based on joint ones.
- (*) Calculate joint and marginal CDF

Problem Write a program that estimates the PMFs distribution

Problem Tossing a coin: $\Omega = \{ H, T \}$ and rolling strange die $\Omega = \{ 1, 2, 3 \}$

If H we roll the die twice, if T we roll the die once

X number of heads $X = \{0, 1\}$

Y sum from die $Y = \{1, 2, 3, 4, 5, 6\}$

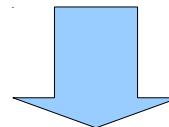
- Calculate joint PMF
- Calculate marginal PMF based on joint ones.
- (*) Calculate joint and marginal CDF

Conditional distribution

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

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

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

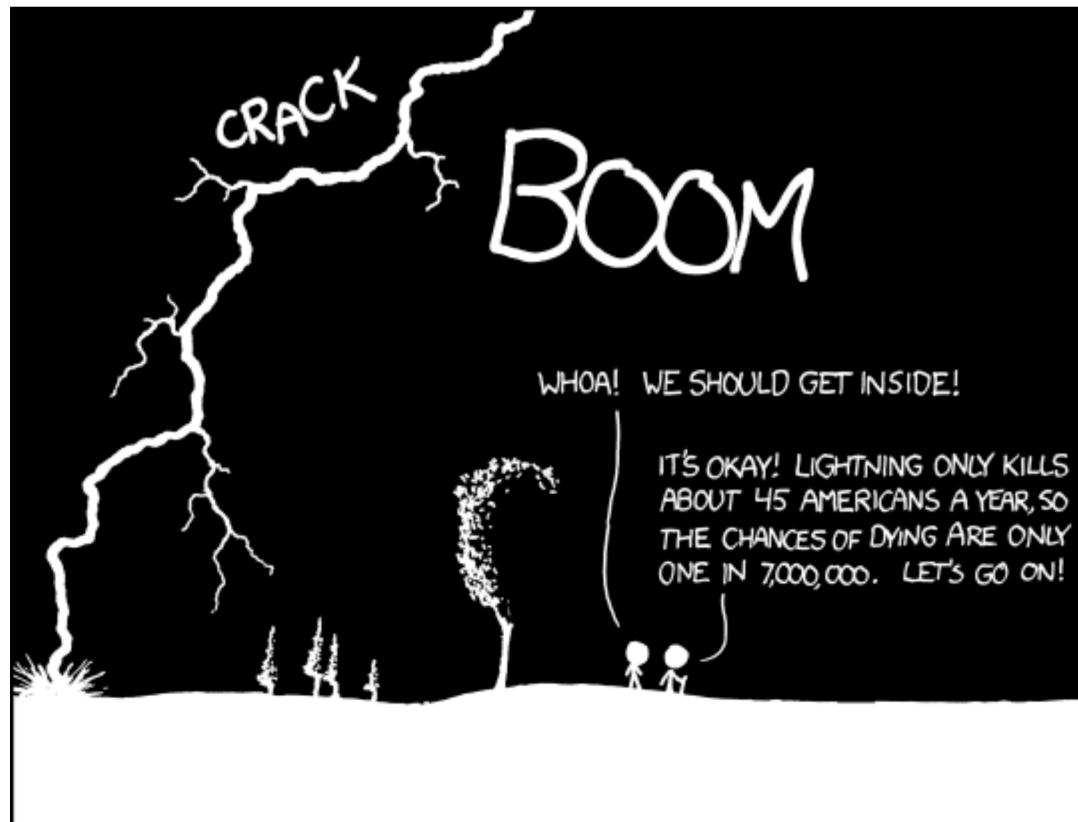


Conditional distribution

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

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

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



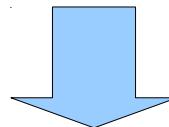
<https://xkcd.com/795/>

Conditional distribution

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

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

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



X, Y two random variables

What is the probability distribution for Y when X random variable has value set to $X=x$?

- (discrete) with the joint PMF: $p_{XY}(k,r)$, and $p_X(k) \neq 0$:

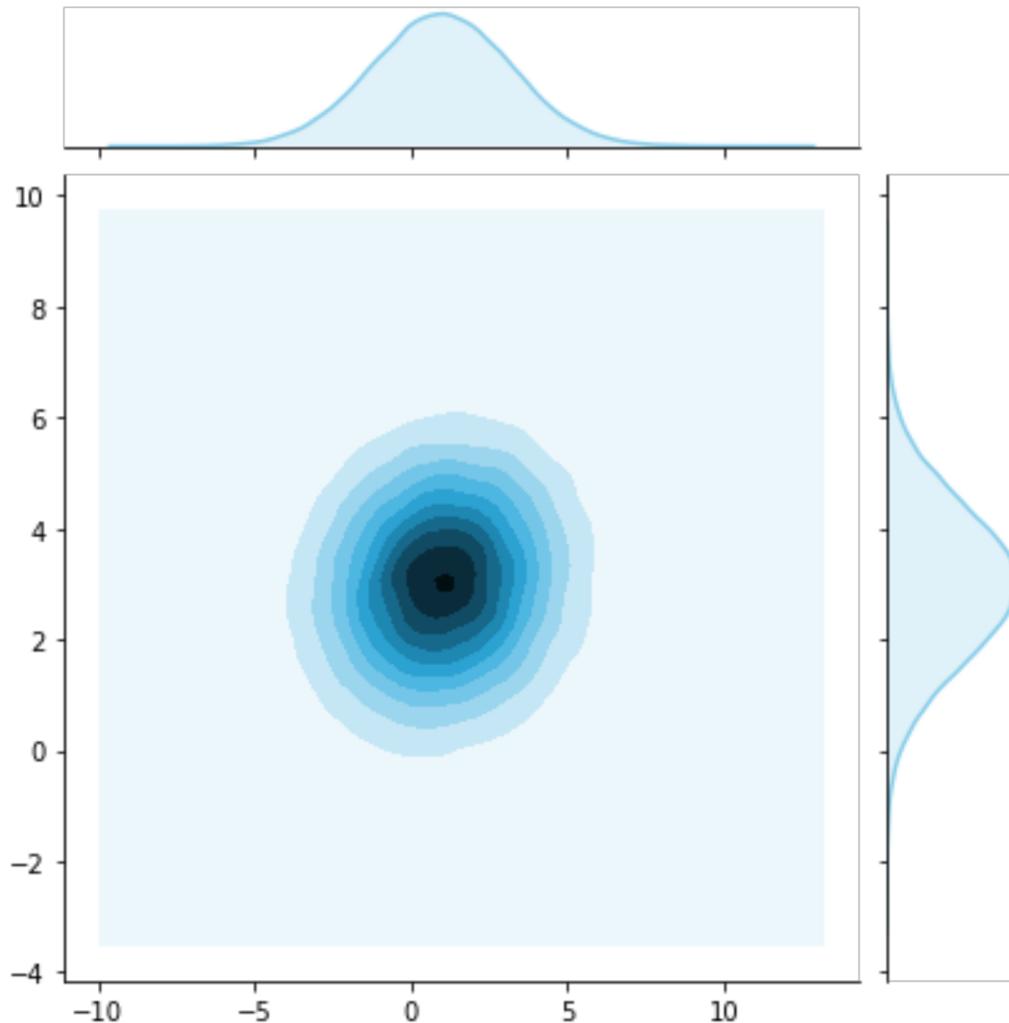
$$p_{Y|X}(r|k) = p_{XY}(k,r)/p_X(k)$$

- (continuous) with the joint PDF $f_{XY}(x,y)$ and $f_X(x) \neq 0$

$$f_{Y|X}(y|x) = f_{XY}(x,y)/f_X(x) (*)$$

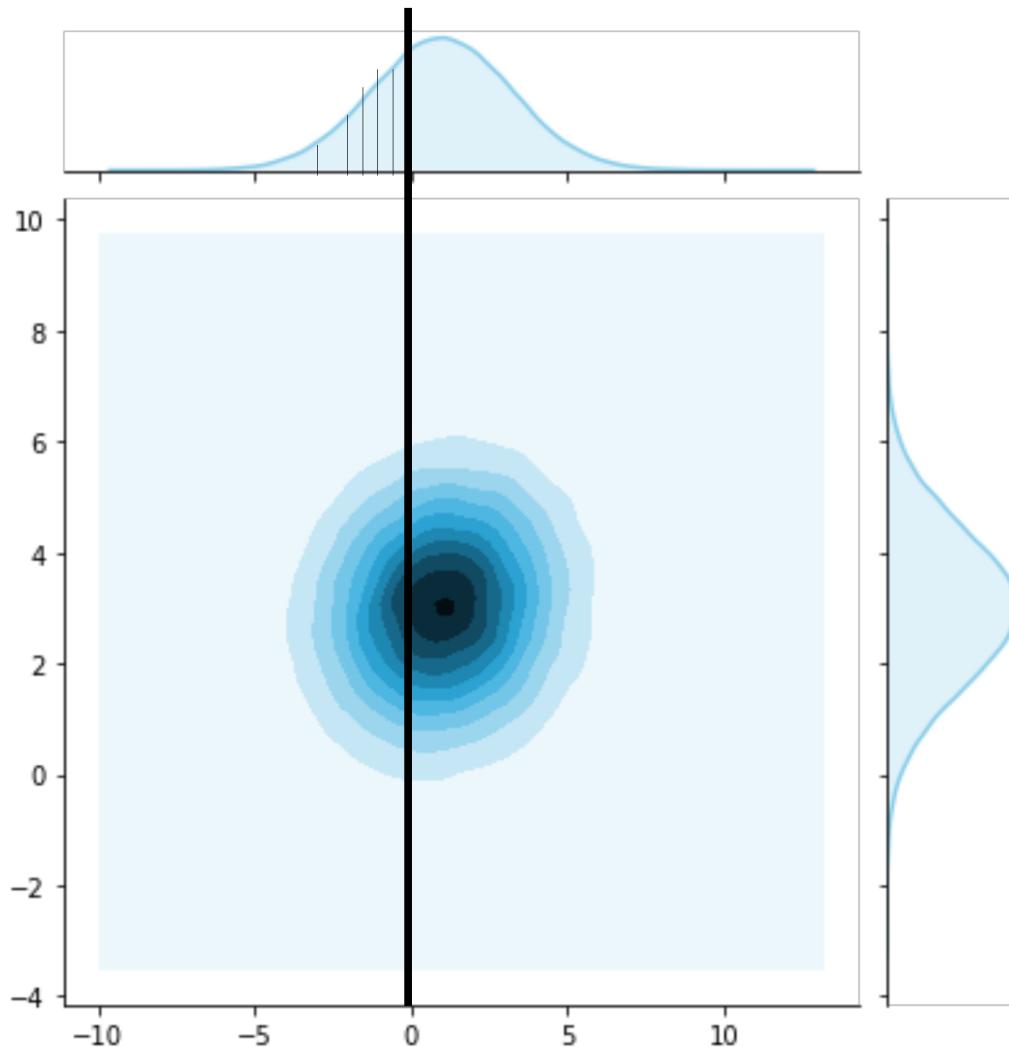
***Remark:** Formally in continuous case it should be derived using CDF and limit

Conditional probability



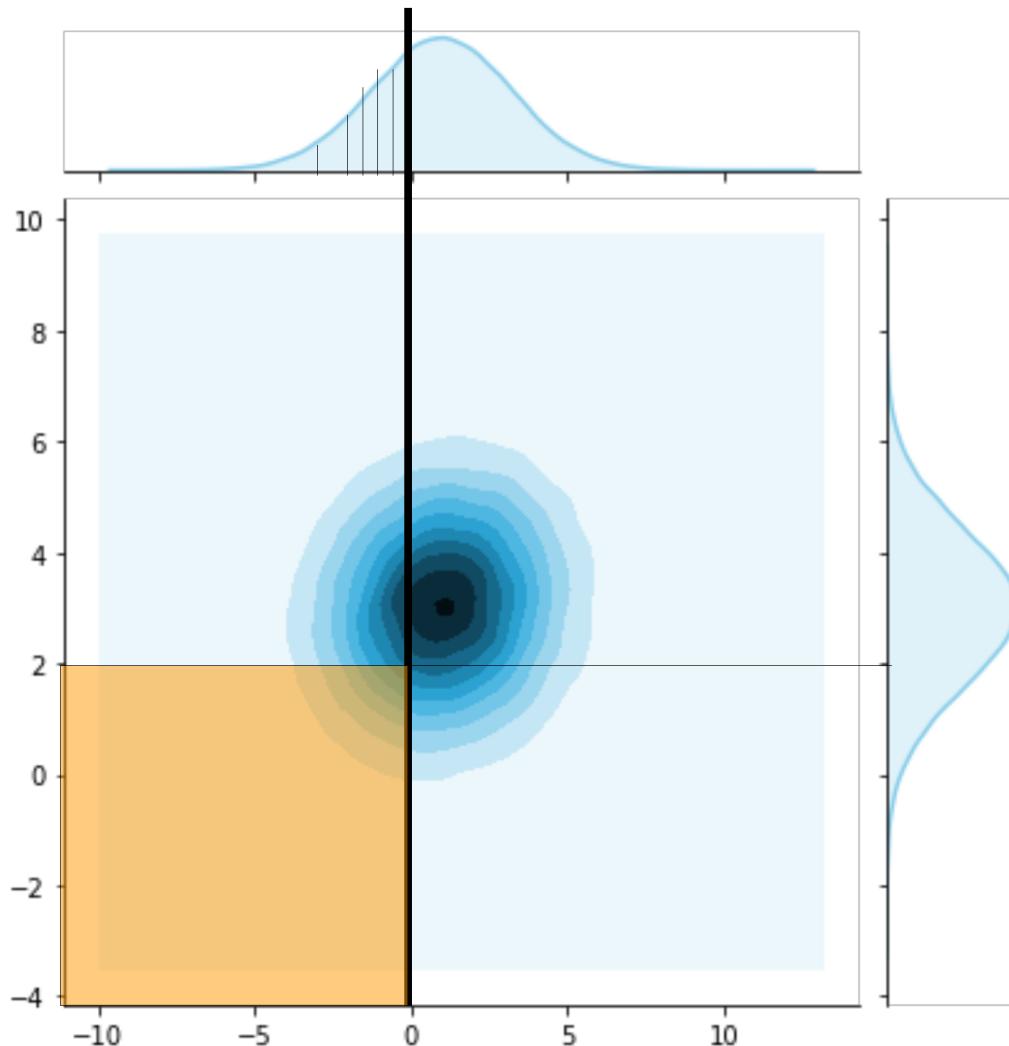
$$P(Y \leq 2 | X \leq 0) = P(X \leq 0, Y \leq 2) / P(X \leq 0)$$

Conditional probability



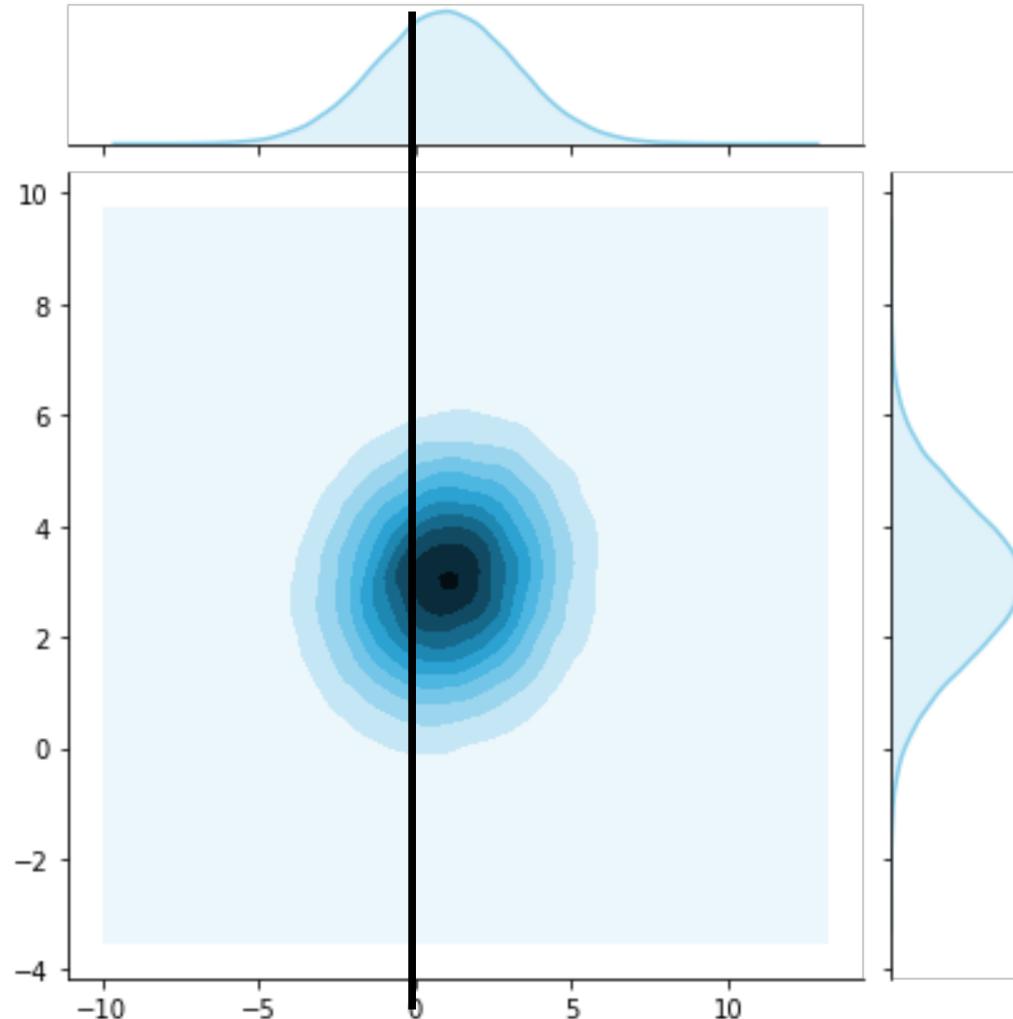
$$P(Y \leq 2 | X \leq 0) = P(X \leq 0, Y \leq 2) / P(X \leq 0)$$

Conditional probability



$$P(Y \leq 2 | X \leq 0) = P(X \leq 0, Y \leq 2) / P(X \leq 0)$$

Conditional pdf



$$f_{Y|X}(y|x=0) = f_{XY}(x=0,y)/f_X(x=0)$$

- If you think all this is trivial – look for Borel-Kolmogorov paradox

Probability function vs likelihood function

X, Y two random variables

If we treat conditional PDF as a function of y: $g(y) = f_{Y|X}(y|x)$, it is a proper **probability function**:

$$\int g(y)dy = 1$$

If we treat conditional PDF as a function of x: $h(x) = f_{Y|X}(y|x)$, it is a **likelihood function**:

$$\int h(x)dx \text{ not necessarily } = 1$$

The same remark concerns PMFs

Problem: Rolling a die $\Omega = \{1, 2, 3, 4, 5, 6\}$

X is 1 if even number 0 otherwise.

Y is 1 if prime number 0 otherwise.

- Calculate joint PMF
- Calculate marginal PMF of X and of Y
- Calculate conditional PMF $p_{Y|X}(r|X=1)$
- Check if $h(k) = p_{Y|X}(r|k)$ is a proper probability function with respect to k

Problem: Let X and Y have a joint PDF $f_{XY}(x,y) = x+y$ for $0 < x < 1, 0 < y < 1$

- Find conditional PDF $f_{Y|X}(y|x)$
- Show that the integral of $f_{Y|X}(y|x)$ over all y values is equal to 1

Expectation for X and Y

Let $g(x,y)$ be some function: $R^2 \rightarrow R$, X, Y – random variables, $g(X,Y)$ becomes random variable as well:

- For X, Y being discrete random variables with joint PMF $p_{XY}(k,r)$:

$$E[g(X,Y)] = \sum g(k,r) p_{XY}(k,r) \text{ for all possible values of } k \text{ and } r$$

- For X, Y being a continuous random variables with PDF $f_{XY}(x,y)$:

$$E[g(X,Y)] = \iint g(x,y) f_{XY}(x,y) dx dy \text{ for } x \text{ from } -\infty \text{ to } \infty, y \text{ from } -\infty \text{ to } \infty$$

Properties:

- $E[f(X,Y) + g(X,Y)] = E[f(X,Y)] + E[g(X,Y)]$ - linearity
- $E[K * f(X,Y)] = K * E[f(X,Y)]$ for K being constant

- **Intuition:** $E[g(X,Y)]$ is a kind of weighted average with weights defined by p_{XY} or f_{XY}

Covariance

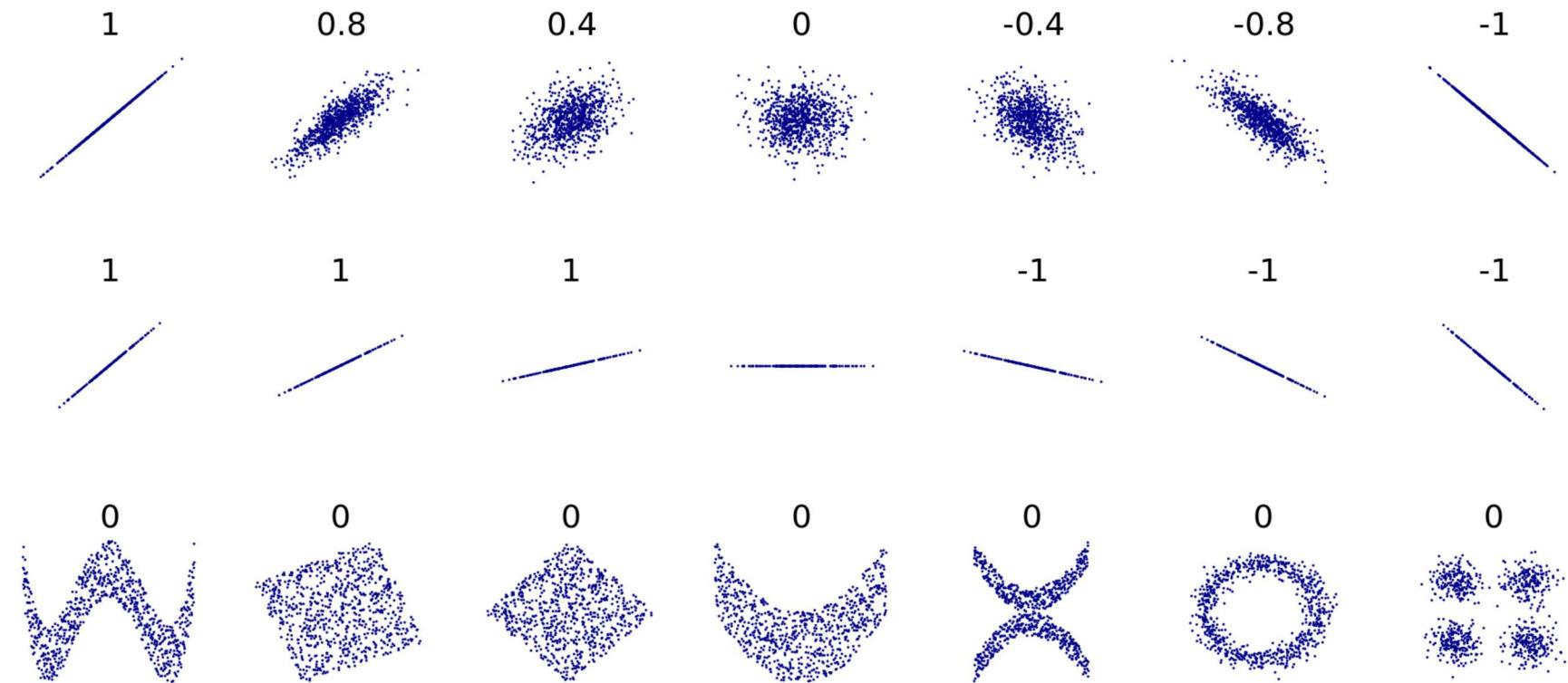
Covariance expresses the relationship between two random variables X,Y

$$\text{Cov}[X,Y] = E[(X - E[X]) * (Y - E[Y])] \rightarrow \text{Cov}[X,Y] = E[XY] - E[X]E[Y]$$

- Covariance captures only **linear** dependency between variables (!!):
 - If $\text{Cov}[X,Y] > 0$ then if X increases then Y will increase,
 - If $\text{Cov}[X,Y] < 0$ then if X increases then Y will decrease,
 - If $\text{Cov}[X,Y] = 0$ we say that X and Y are **uncorrelated**.
- $\text{Cov}[X,X] = \text{Var}[X]$
- $\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X,Y]$

Remark: If X,Y are independent, then $\text{Cov}[X,Y] = 0$ but opposite not always true

Covariance

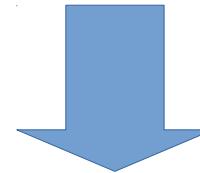


From Wikipedia by Denis Boigelot, <https://commons.wikimedia.org/w/index.php?curid=15165296>

Independence

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

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



We want to generalize the notion of independence to random distributions.

Let X, Y be two random variables.

Independent variables X, Y

Let X, Y be **independent** random variables, then it is true that:

- $E[f(X) * g(Y)] = E[f(X)] * E[g(Y)] \rightarrow$ factorization of expectation values

Independent variables X, Y

Let X, Y be **independent** random variables, then it is true that:

- $E[f(X) * g(Y)] = E[f(X)] * E[g(Y)] \rightarrow$ factorization of expectation values
- $F_{XY}(x,y) = F_x(x) * F_y(y) \rightarrow$ factorization of joint CDF

Independent variables X, Y

Let X, Y be **independent** random variables, then it is true that:

- $E[f(X) * g(Y)] = E[f(X)] * E[g(Y)] \rightarrow$ factorization of expectation values
- $F_{XY}(x,y) = F_x(x) * F_y(y) \rightarrow$ factorization of joint CDF
- For X, Y , being discrete random variables:
 - $p_{XY}(k,r) = p_X(k) * p_Y(r)$ for all k from X and r from $Y \rightarrow$ factorization of PMF
 - $p_{Y|X}(r|k) = p_Y(r)$ for all r from Y , where $p_X(k) \neq 0$

Independent variables X, Y

Let X, Y be **independent** random variables, then it is true that:

- $E[f(X) * g(Y)] = E[f(X)] * E[g(Y)] \rightarrow$ factorization of expectation values
- $F_{XY}(x,y) = F_x(x) * F_y(y) \rightarrow$ factorization of joint CDF
- For X, Y , being discrete random variables:
 - $p_{XY}(k,r) = p_X(k) * p_Y(r)$ for all k from X and r from $Y \rightarrow$ factorization of PMF
 - $p_{Y|X}(r|k) = p_Y(r)$ for all r from Y , where $p_X(k) \neq 0$
- For X, Y , being continuous random variables:
 - $f_{XY}(x,y) = f_X(x) * f_Y(y)$ for all x,y from $R \rightarrow$ factorization of PDF
 - $f_{Y|X}(y|x) = f_Y(y)$ for all y from Y , where $f_X(x) \neq 0$

Problem: Show that:

- $\text{Cov}[X,Y] = E[XY] - E[X]E[Y]$
- $\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y] + 2 \text{ Cov}[X,Y]$

Problem: Let X be uniform in $(-1,1)$ and $Y = X^2$

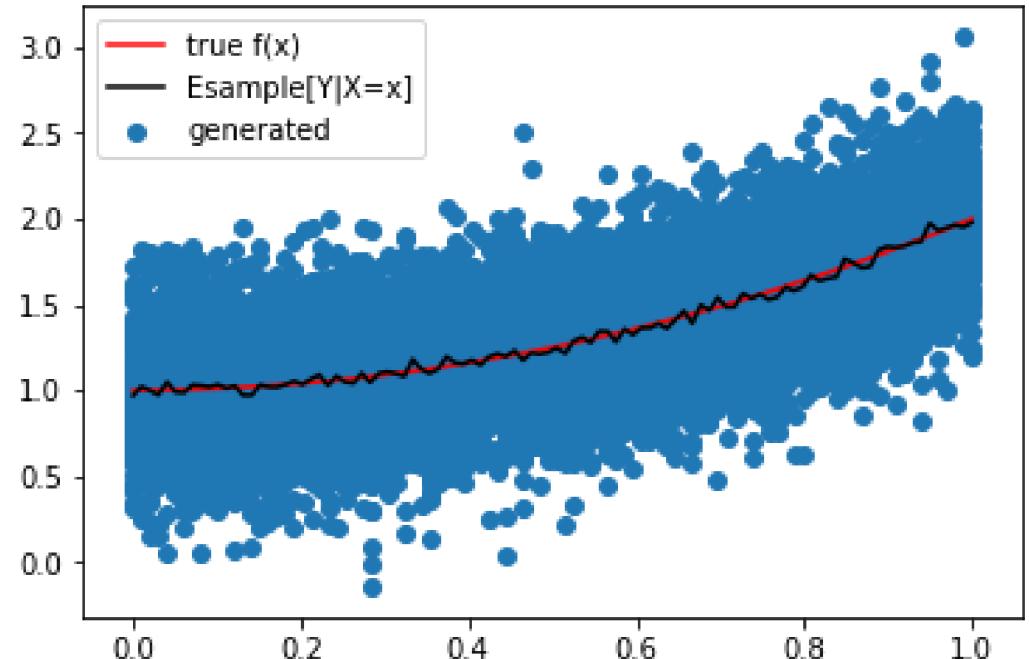
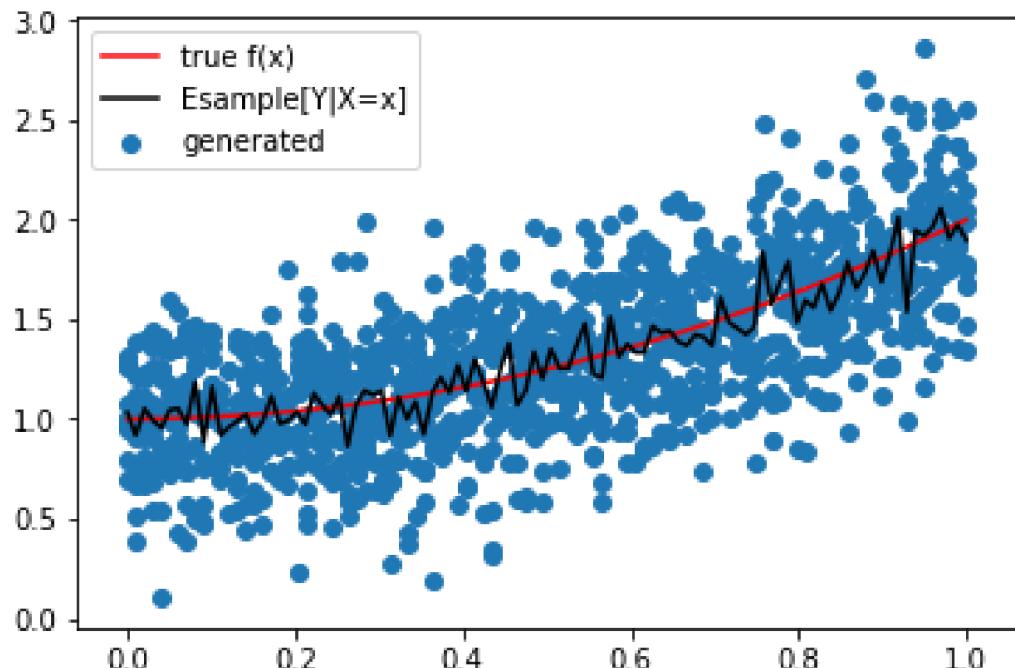
- Check if X and Y are correlated
- Check if X and Y are independent

Conditional expectation

Conditional expectation of random variable X expresses $E[X]$ under some set of conditions that “occurred”.

- For X, Y being discrete random variables with conditional PMF $p_{Y|X}(r|k)$:
$$E[Y|k] = \sum r * p_{Y|X}(r|k) \text{ for all possible values of } r$$
- For X, Y being a continuous random variables with conditional PDF $f_{Y|X}(y|x)$:
$$E[Y|x] = \int y * f_{Y|X}(y|x) dy \text{ for } y \text{ from } -\infty \text{ to } \infty$$

Conditional expectation example



Conditional variance

Conditional variance of random variable X expresses $\text{Var}[X]$ under some set of conditions that “occurred”.

- For X,Y being discrete random variables with conditional PMF $p_{Y|X}(r|k)$:

$$\text{Var}[Y|k] = \sum (r - E[Y|k])^2 * p_{Y|X}(r|k) \text{ for all possible values of } r$$

- For X,Y being a continuous random variables with conditional PDF $f_{Y|X}(y|x)$:

$$\text{Var}[Y|x] = \int (y - E[Y|x])^2 * f_{Y|X}(y|x) dy \text{ for } y \text{ from } -\infty \text{ to } \infty$$

- One can show that :

$$\text{Var}[Y|k] = E[Y^2|k] - (E[Y|k])^2$$

Problem 6 rolling a die $\Omega = \{ 1, 2, 3, 4, 5, 6 \}$

X is 1 if even number 0 otherwise.

Y is 1 if prime number 0 otherwise.

- Calculate $E[Y|1]$, $E[Y|0]$
- $\text{Var}[Y|1]$, $\text{Var}[Y|0]$

Problem:

- Implement a function that returns mean
- Implement a function that returns $\text{Var}[X]$
- Implement a function that returns $\text{Cov}[X,Y]$

Problem: Similarity measures

- Implement a function that returns Euclidean distance between two vectors
- Implement a function that returns Manhattan distance between two vectors
- Implement a function that returns Cosine similarity
- Implement a function that returns Pearson correlation coefficient

Thank you