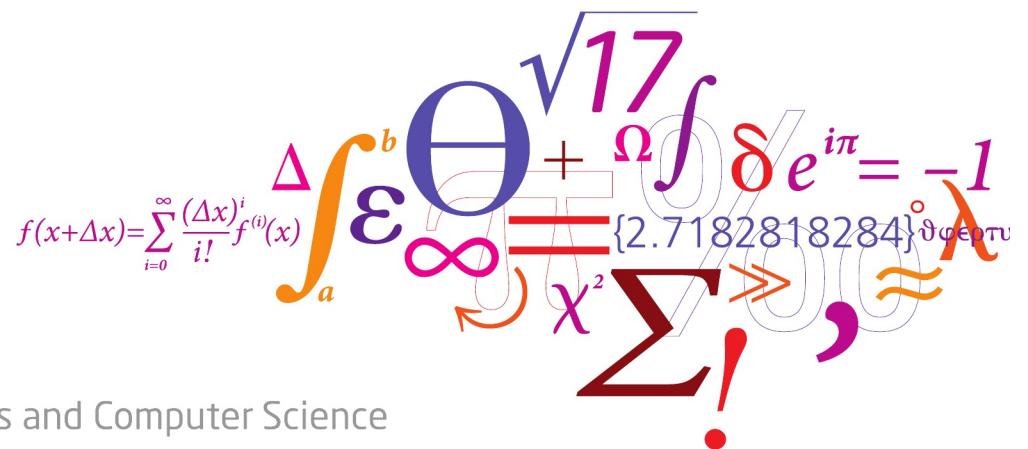


02450: Introduction to Machine Learning and Data Mining

Measures of similarity, summary statistics and probabilities

Tue Herlau

DTU Compute, Technical University of Denmark (DTU)

$$f(x+\Delta x) = \sum_{i=0}^{\infty} \frac{(\Delta x)^i}{i!} f^{(i)}(x)$$


DTU Compute

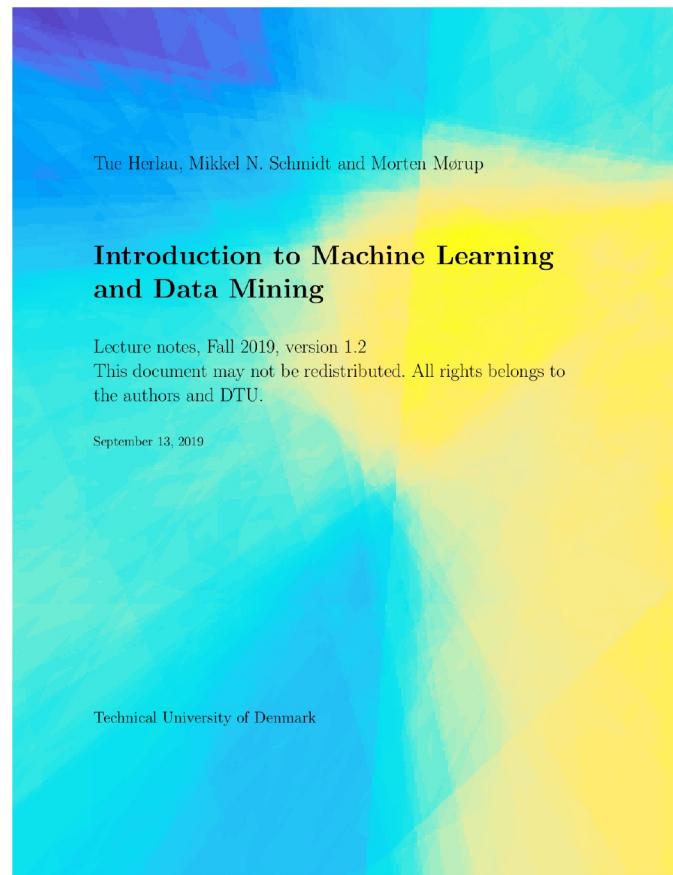
Department of Applied Mathematics and Computer Science

Today

Feedback Groups of the day:

Georgios Bekakos, Georgios Bekakos, Vasileios Boutlas, Cecilie Amalie Neijendam Thystrup, Clara Drachmann, Søren Steenstrup Zeeberg, Lisbeth Ditte Evald Sandvik, Agnete Marie Nørregaard, Robert Stig Grønlund Nielsen, Lukas Kofoed Gildhoff, Khushboo Nyman, Enrico Tolotto, Alexandros Spyropoulos, Louis Perot, Felix Gigler, David Onnen, Stefan Bîrs, Morten Bjerre, Mads Gramtorp Bjørn, Callum Blair, Kálmán Bogdán, Julian Böhm, Lisa Veibel Bonde, Julia Bonzanini, Annie Borch, Rishav Bose, Oskar Johannes Fred Bremberg, Mads Brink, Benjamin Bruun, Mikkel Bruus, Rasmus Stokholm Bryld, Razvan-Vlad Bucur, Jonathan Buhl, Philip Nielsen Butenko, Manuel Caballero Peña, Cemre Çadir, Yiyi Cao, Chiara Libera Carnevale, Dilara Eda Celepli, Emmanouil Chalvatzopoulos, Anshul Chauhan, Dominic Chen, Darius Chira, Seokhyun Choung, Emil Chrisander, Oskar Eiler Wiese Christensen, Line Wulff Christensen, Ann-Katrine Grodt Christiansen, Jonas Søbro Christophersen, Yeon Soo Chung, Sofus Rasmus Clausen, Joakim Nøddeskov Clifford

Reading material: Chapter 4, Chapter 5



Lecture 3 17 September, 2019

Lecture Schedule

1 Introduction

3 September: C1

Data: Feature extraction, and visualization

2 Data, feature extraction and PCA

10 September: C2, C3

3 Measures of similarity, summary statistics and probabilities

17 September: C4, C5

4 Probability densities and data Visualization

24 September: C6, C7

Supervised learning: Classification and regression

5 Decision trees and linear regression

1 October: C8, C9 (Project 1 due before 13:00)

6 Overfitting, cross-validation and Nearest Neighbor (Note: Tentative)

8 October: TBA

7 Bayes, Naive Bayes and performance evaluation (Note: Tentative)

22 October: TBA

Piazza online help: <https://piazza.com/dtu.dk/fall2019/02450>

Videos/streaming of lectures: <https://video.dtu.dk>

8 Artificial Neural Networks and Bias/Variance

29 October: C14, C15

9 AUC and ensemble methods

5 November: C16, C17

Unsupervised learning: Clustering and density estimation

10 K-means and hierarchical clustering

12 November: C18 (Project 2 due before 13:00)

11 Mixture models and density estimation

19 November: C19, C20

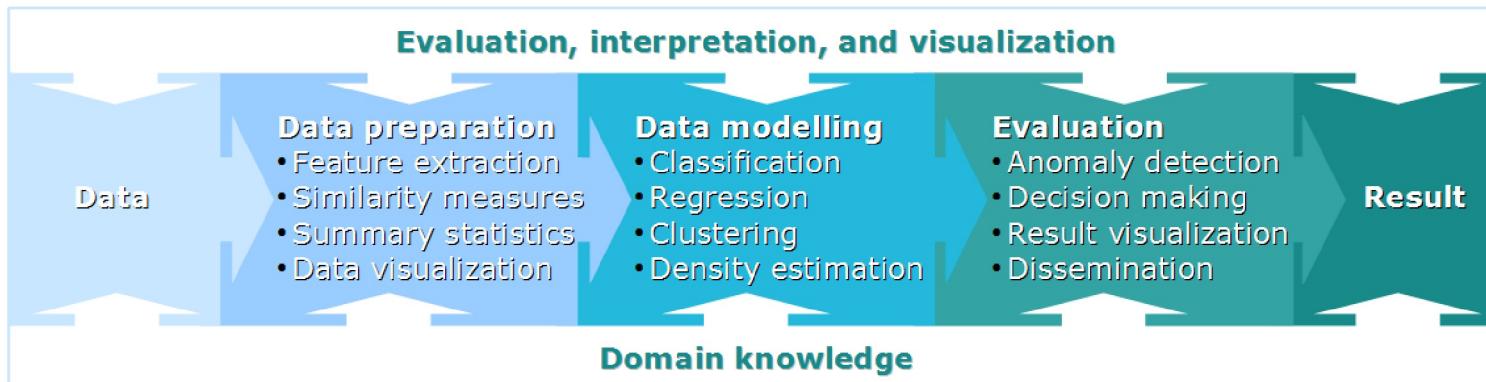
12 Association mining

26 November: C21

Recap

13 Recap and discussion of the exam

3 December: C1-C21 (Project 3 due before 13:00)



Learning Objectives

- Compute measures of similarity/dissimilarity (L_p distance, cosine distance, etc.)
- Understand basics of probability theory and stochastic variables in the discrete setting
- Understand probabilistic concepts such as expectations, independence and the Bernoulli distribution
- Understand the maximum likelihood principle for repeated binary events

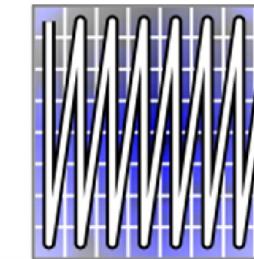
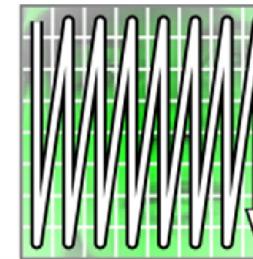
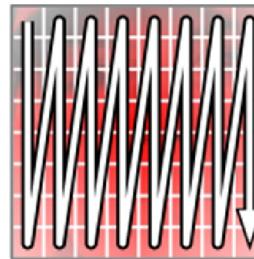
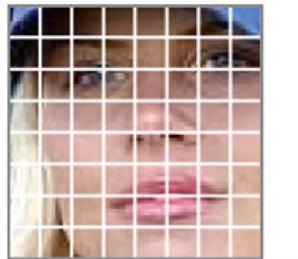
PCA recap: Principal component analysis on images



- 1000 images, 86 x 86 pixels, 3 RGB intensities

Tamara Berg "Faces in the wild"

Pre-processing



- Concatenate all pixel color values in one long vector
 - $86 \times 86 \times 3 = 22'188$
 - Image is now represented as a 22'188 dimensional vector
- Stack all 1000 images into a big matrix
 - $1000 \times 22'188 = X$

Principal component analysis (PCA)

$$X \xrightarrow[\text{(SVD)}]{\text{mean}} \tilde{X}$$

$$v_i = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}_{22188} \quad (\tilde{X}^T \tilde{X}) v_i = \lambda_i v_i$$

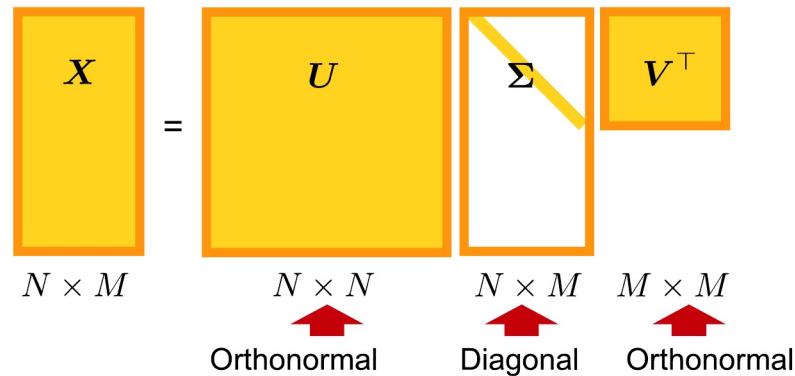
1. Subtract the mean

- Consider dividing with variance; use 1-out-of-K coding for nominal attributes

2. Compute the singular value decomposition (SVD)

- Orthogonal linear transformation
- Transforms data to a new coordinate system
 - Greatest variance along the first axis (first column of V)
 - Second greatest variance along the second axis

$$\lambda_i \cdot = \sum_{i,i} \quad V = [v_1 \dots v_M]$$



• Plot data in the transformed coordinate system

- Corresponds to looking at data from an angle where it is most spread out

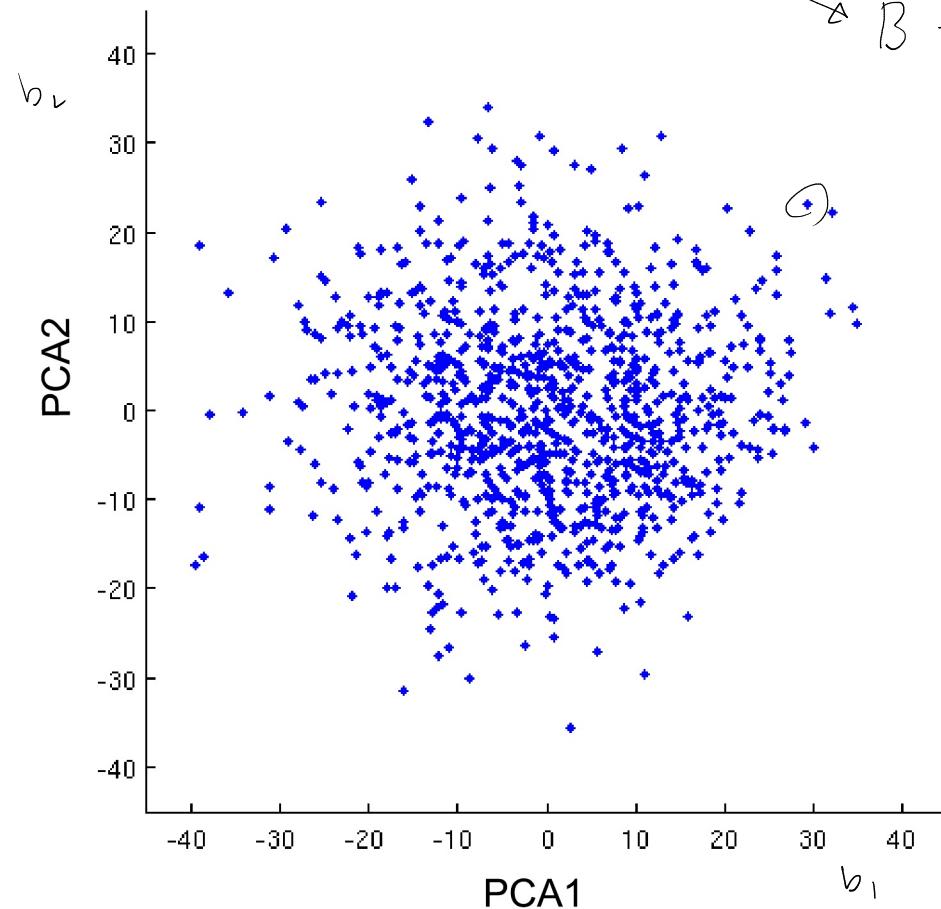
PCA on face images

$$V_n = \begin{bmatrix} V_1 & V_2 \end{bmatrix}, n=2$$

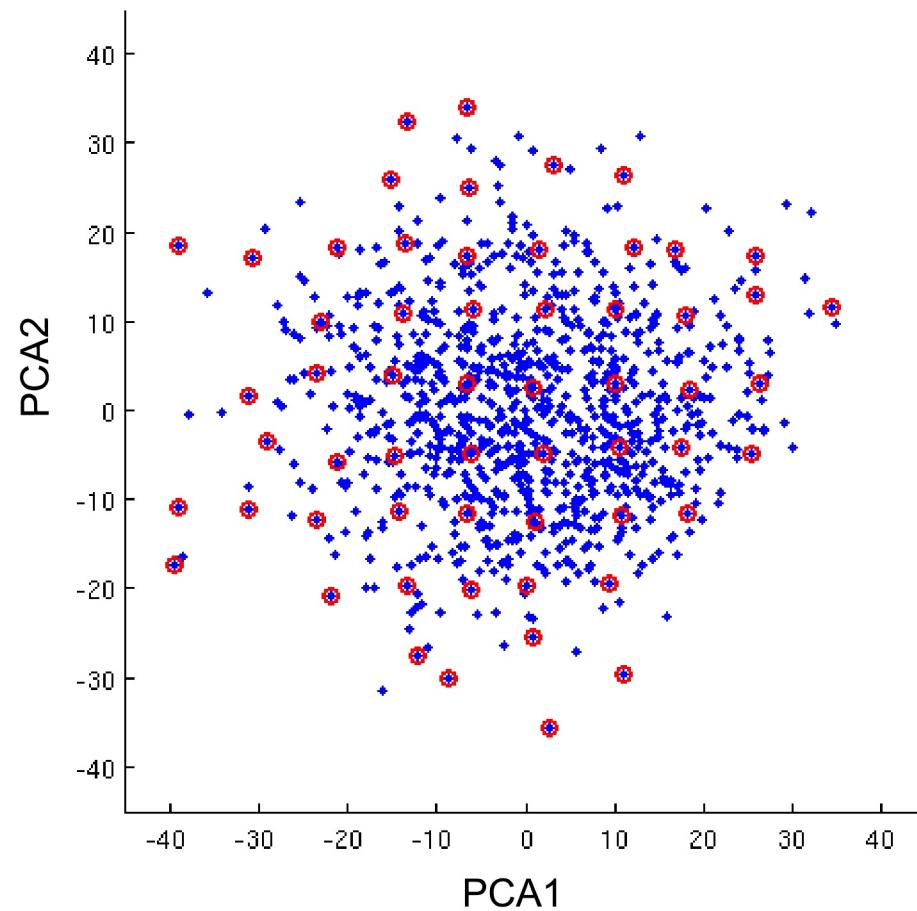
$$\begin{bmatrix} b_1 & b_2 \end{bmatrix} = \begin{bmatrix} \tilde{X}_i^T V_1 & \tilde{X}_i^T V_2 \end{bmatrix}$$

$$= \tilde{X}_i^T V_n$$

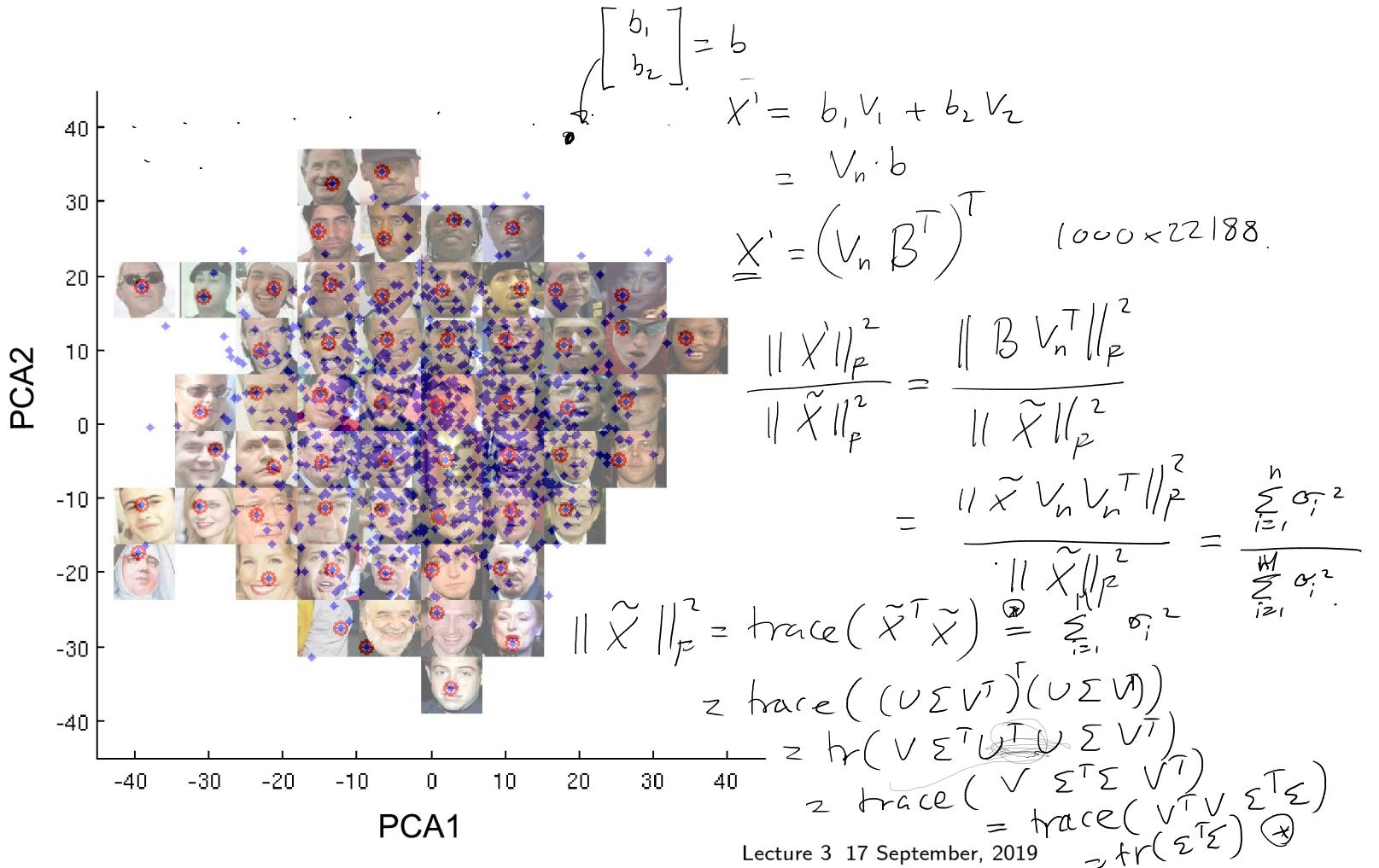
$$1000 \times 2 \quad \rightarrow \quad \mathcal{B} = \tilde{X}^T V_n \quad 22188 \times 2$$



PCA on face images



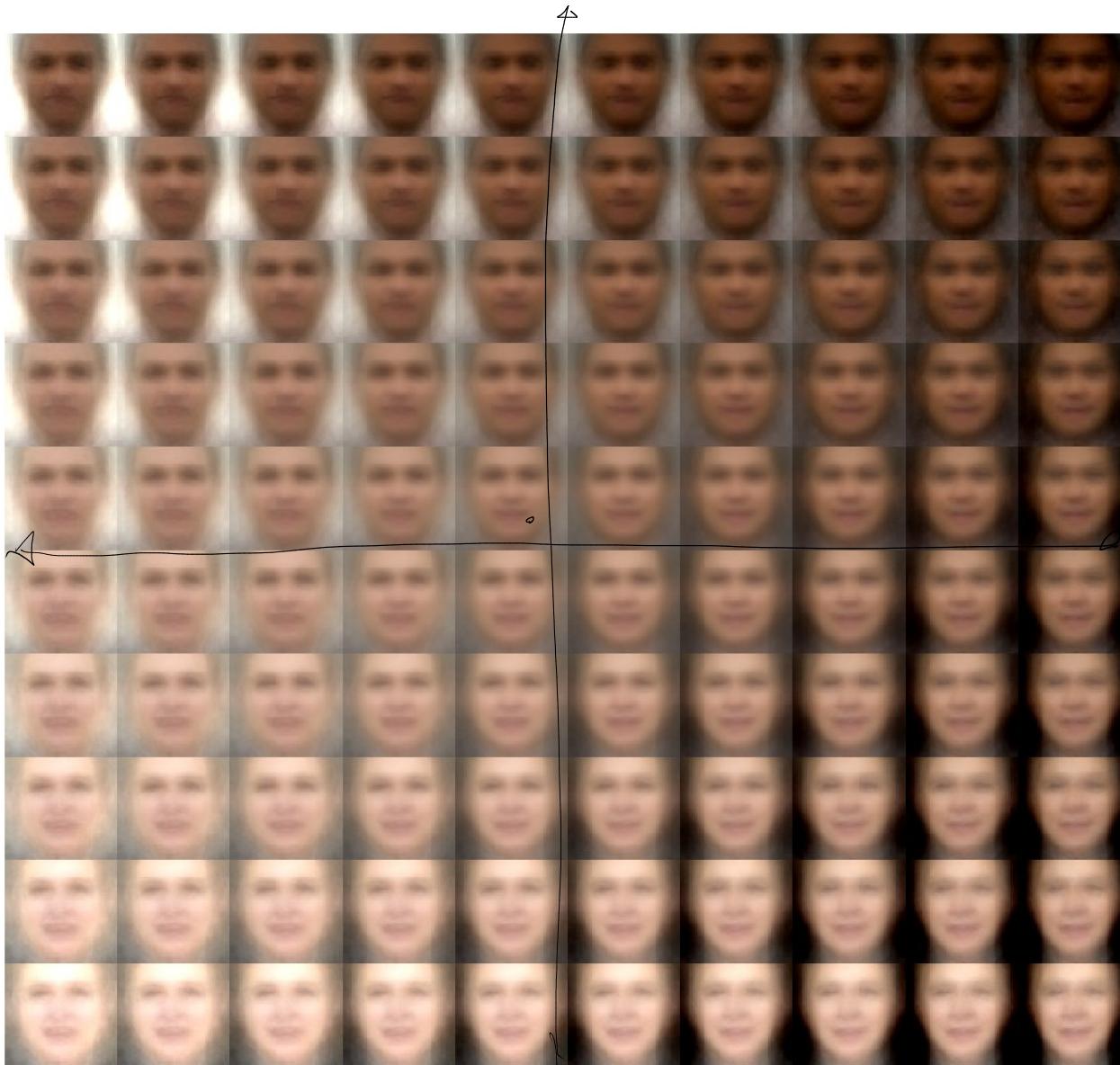
PCA on face images





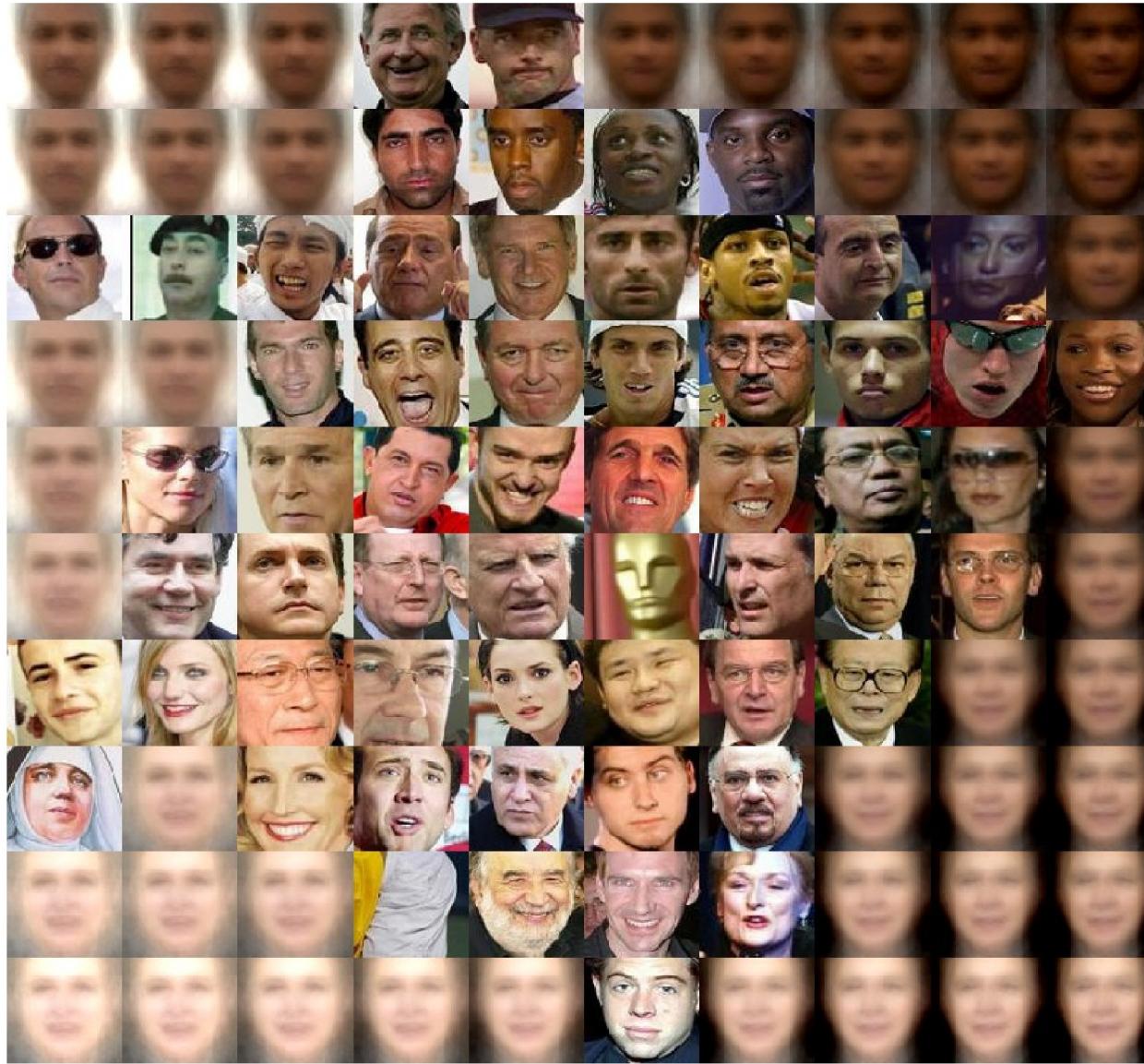
• What information do the two principal axes capture?

$$\begin{bmatrix} b_2 & b_1 \end{bmatrix}$$

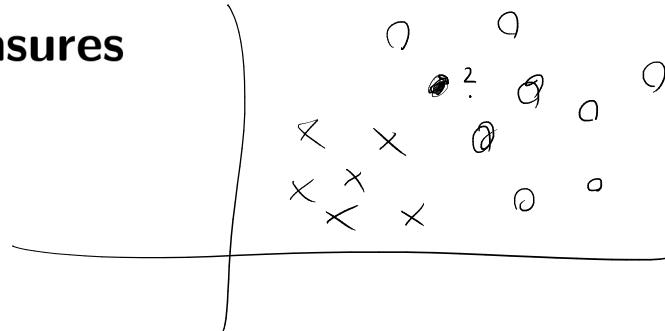




What information do the two principal axes capture?



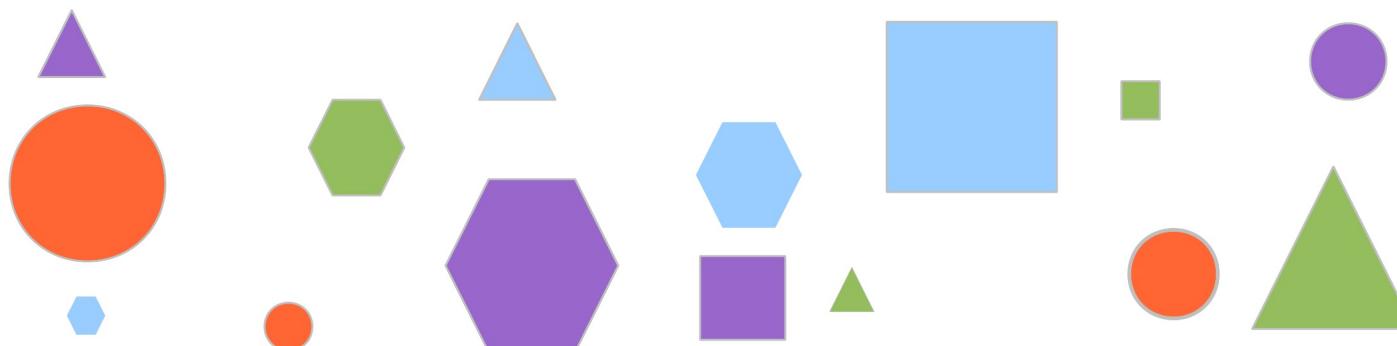
Similarity / Dissimilarity measures



- } Similarity $s(x, y)$ Often between 0 and 1. Higher means more similar
Dissimilarity $d(x, y)$ Greater than 0. Lower means more similar.

Classification Classify a document as having the same topic y as the document is **most similar/least dissimilar** to.

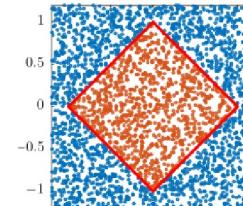
Outlier detection The observation most **dissimilar** to all other observations is an outlier



Dissimilarity measures

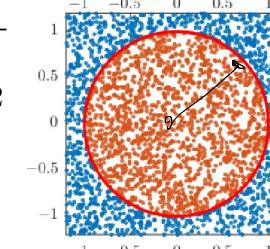
- General Minkowsky distance (p -distance) $d_p(\mathbf{x}, \mathbf{y}) = \left(\sum_{j=1}^M |x_j - y_j|^p \right)^{\frac{1}{p}}$
- One-norm ($p = 1$)

$$d_1(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^M |x_j - y_j|$$



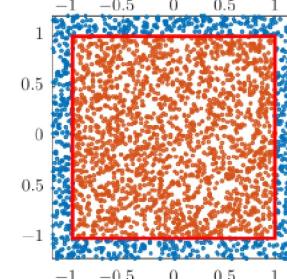
- Euclidean ($p = 2$)

$$d_2(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{j=1}^M (x_j - y_j)^2}$$



- Max-norm distance ($\underbrace{p = \infty}_{\text{}}$)

$$d_\infty(\mathbf{x}, \mathbf{y}) = \max \{|x_1 - y_1|, |x_2 - y_2|, \dots, |x_M - y_M|\}$$



Usage: Regularization and alternative optimization targets. For instance, $p = \infty$ is very affected by outliers, $p = 1$ much less so.

Similarity measures

$$x = \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

K : Total number of attributes
 f_{00} : Number of attributes where $x_k = y_k = 0$
 f_{11} : Number of attributes where $x_k = y_k = 1$

Simple Matching Coefficient (SMC)

$$\text{SMC}(x, y) = \frac{f_{00} + f_{11}}{K}$$

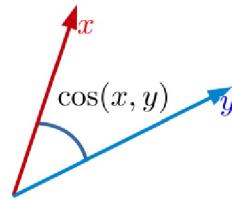
+ Symmetric

Jaccard Coefficient

$$J(x, y) = \frac{f_{11}}{K - f_{00}}$$

+ Positive matches

Cosine similarity



$$\cos(x, y) = \frac{x^\top y}{\|x\| \|y\|}$$

+ Positive matches
+ Document length

Extended Jaccard coefficient

$$\text{EJ}(x, y) = \frac{x^\top y}{\|x\|^2 + \|y\|^2 - x^\top y}$$

+ Positive matches
+ Document length

Also defined for continuous data

Quiz 1, similarity measures

Calculate the simple matching coefficient, Jaccard, cosine and extended jaccard similarity between customer 1 and customer 2 in the market basket data below.

$$f_{11} = 2 \quad f_{00} = 1$$

ID	Bread	Soda	Milk	Beer	Diaper
o_1	1	1	1	0	0
o_2	2	0	1	1	0

K : Total number of attributes

f_{00} : Number of attributes where $x_k = y_k = 0$

f_{11} : Number of attributes where $x_k = y_k = 1$

Which of the following statements are true?

- A. $SMC(o_1, o_2) = \frac{3}{5}, J(o_1, o_2) = \frac{1}{2}, \cos(o_1, o_2) = \frac{2}{3},$
- B. $SMC(o_1, o_2) = \frac{3}{5}, J(o_1, o_2) = \frac{3}{4}, \cos(o_1, o_2) = \sqrt{\frac{2}{3}},$
- C. $\cancel{SMC(o_1, o_2) = \frac{2}{5}, J(o_1, o_2) = \frac{1}{3}, \cos(o_1, o_2) = \frac{2}{3},}$
- D. $\cancel{SMC(o_1, o_2) = \frac{2}{5}, J(o_1, o_2) = \frac{1}{3}, \cos(o_1, o_2) = \sqrt{\frac{2}{3}},}$
- E. Don't know.

$$SMC(x, y) = \frac{f_{00} + f_{11}}{K}$$

$$J(x, y) = \frac{f_{11}}{K - f_{00}}$$

$$\cos(x, y) = \frac{x^\top y}{\|x\|_2 \|y\|_2}$$

$$EJ(x, y) = \frac{x^\top y}{\|x\|_2^2 + \|y\|_2^2 - x^\top y}$$

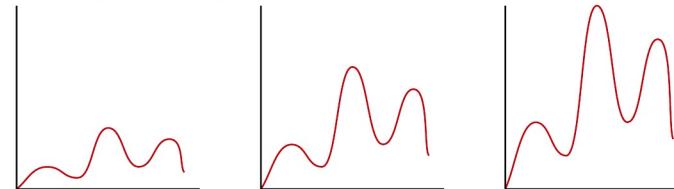
The problem is easily solved by using the inserted formula. We obtain: $\text{SMC}(o_2, o_2) = \frac{3}{5}$, $J(o_1, o_2) = \frac{1}{2}$, $\cos(o_1, o_2) = \frac{2}{3}$ and therefore the A is true. Since the

data is binary, the extended Jaccard and the jaccard coefficient agree.

Invariance

Scale invariance

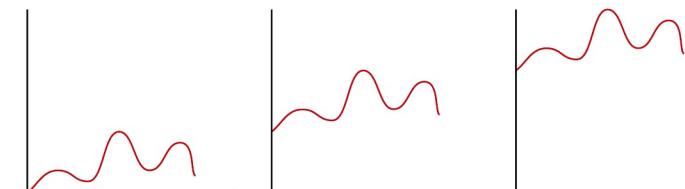
$$d(\mathbf{x}, \mathbf{y}) = d(\alpha \mathbf{x}, \mathbf{y})$$



Translation invariance

$$d(\mathbf{x}, \mathbf{y}) = d(\alpha + \mathbf{x}, \mathbf{y})$$

$$d(\mathbf{x}, \mathbf{y}) \approx d_p(\mathbf{x} - \bar{\mathbf{x}}, \mathbf{y} - \bar{\mathbf{y}})$$



General invariances

$$d(\boxed{6}, \boxed{2}) = d(\boxed{16}, \boxed{2})$$

Transformations

Standardization: Ensure a single attribute will not dominate:

$$\tilde{x}_{ik} = \frac{x_{ik} - \hat{\mu}_k}{\hat{\sigma}_k}$$

Example:

- **Number of children** ~ 0-5
- **Age** ~ 0-100 years
- **Annual income** ~ 0-50.000 €

Combining heterogeneous attributes Transform measures and combine

$$s_{\text{Edu.}} = \text{SMC}(x_{\text{Edu.}}, y_{\text{Edu.}})$$

$$s_{\text{Age.}} = a (a + d_1(x_{\text{Age.}}, y_{\text{Age.}}))^{-1}, \quad a = 1$$

$$s(x, y) = \frac{1}{2} (s_{\text{Edu.}} + s_{\text{Age.}})$$

Example:

- **Age:** Continuous
- **Education:** Binary
 - Primary (yes/no)
 - Secondary (yes/no)
 - Tertiary (yes/no)

Weighting Attributes have different importance

$$s(x, y) = 0.99s_{\text{Edu.}} + 0.01s_{\text{Age.}}$$

Empirical statistics

Given two samples x_1, x_2, \dots, x_N and y_1, y_2, \dots, y_N :

- Empirical mean

$$\hat{\mu} = \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

- Empirical variance

$$\hat{s} = \text{vár}[x] = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

- Empirical covariance

$$\text{côv}[x, y] = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu}_x)(y_i - \hat{\mu}_y)$$

- Empirical standard deviation

$$\hat{\sigma} = \hat{\text{std}}[x] = \sqrt{\hat{s}}$$

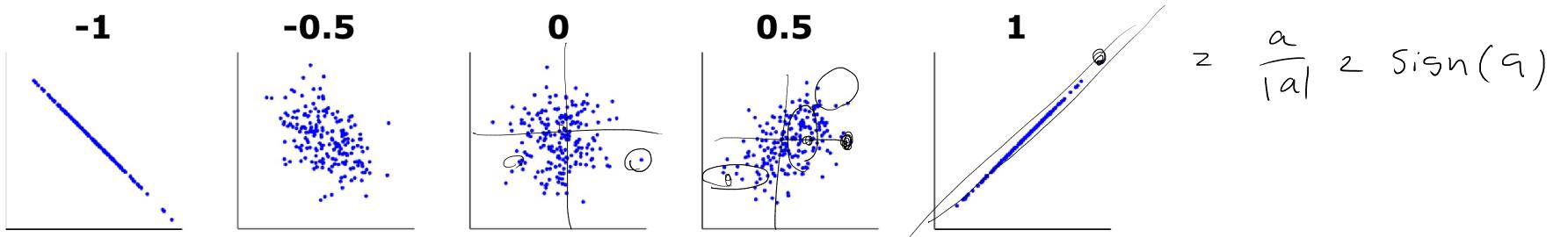
Correlation

- Measure of degree of linear relationship

$$\text{cor}[x, y] = \frac{\hat{\text{cov}}[x, y]}{\hat{\sigma}_x \hat{\sigma}_y} = \frac{\frac{1}{N-1} \sum_{i=1}^N (\hat{y}_i - \hat{\mu}_y)(\hat{a}\hat{y}_i + \hat{b} - (\hat{a}\hat{\mu}_y + \hat{b}))}{\frac{1}{N-1} a \sum_{i=1}^N (\hat{y}_i - \hat{\mu}_y)(\hat{y}_i - \hat{\mu}_y)} = \frac{\frac{a}{N-1} \sum_{i=1}^N (\hat{y}_i - \hat{\mu}_y)(\hat{y}_i - \hat{\mu}_y)}{|a| \hat{\sigma}_y \hat{\sigma}_y} = \frac{a}{|a|} = \text{sign}(a)$$

- A correlation of **1** or **-1** means there is a perfect linear relation

$$x_k = a y_k + b$$



Quantiles

Given N observations of an attribute x_1, x_2, \dots, x_N . The **q 'th quantile** is the value x_q of x such that a fraction q of the sample is **smaller** than x_q .

- Sort in ascending order $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(N)}$
- q 'th quantile is then (approximately) $x_{(\lceil Nq \rceil)}$
- Percentile is the same except q is given in percent $q = \frac{p}{100}$.
- **Median** is the $q = \frac{1}{2}$ quantile:

$$\text{median}[x] = \begin{cases} x'_{\frac{(N+1)}{2}} & \text{if } N \text{ is odd} \\ \frac{1}{2} \left(x'_{\frac{N}{2}} + x'_{\frac{N}{2}+1} \right) & \text{if } N \text{ is even.} \end{cases}$$

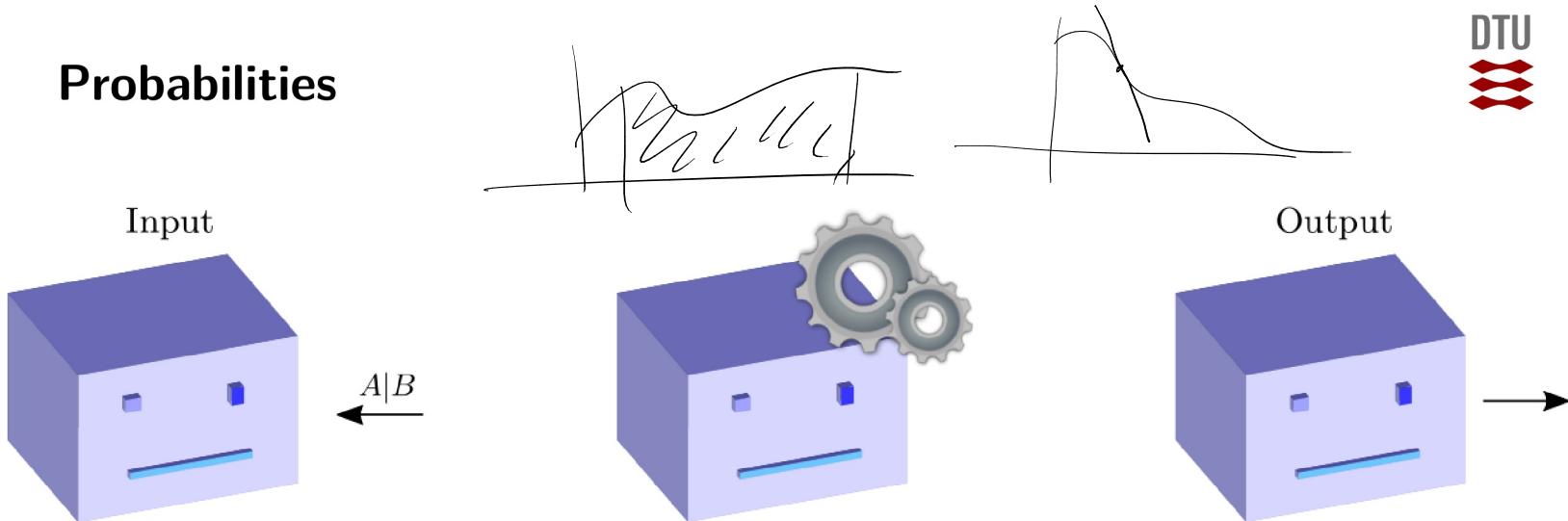
Probabilities

Pragmatically: A big part of AI is dealing with **uncertainty** and **incomplete information**. Probabilities is the formal framework for doing so

Algorithmically: If an image belongs to a particular category is a discrete event. The **probability** it belongs to a category is continuous. Algorithmically, easier to optimize continuous quantities

Convenience: There are boiler-plate ideas for transforming a **probabilistic** assumption into an algorithm (maximum likelihood)

Probabilities



Assuming B is true, how plausible is it A is true?

Assuming B , the plausibility of A is
(low / medium / high / certain)

We reason about a proposition A in light of evidence B :

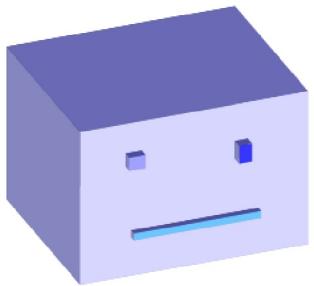
$$P(A|B) = x$$

The degree-of-belief that A is true given B is accepted as true is at a level x

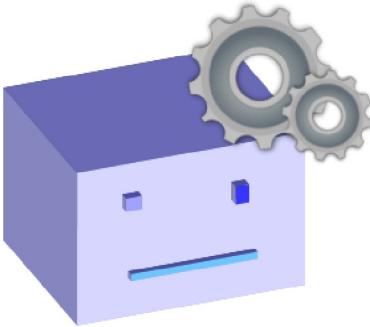
- A number between 0 and 1
- A and B are always binary (true/false) propositions
- Represents a *state of knowledge*

Probabilities: Trial example

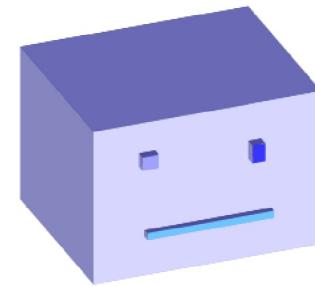
Input



$A|B$



Output



Assuming B is true, how plausible is it A is true?

Assuming B , the plausibility of A is
(low / medium / high / certain)

G : *The accused is guilty*

E_1 : *His mom says he was home on the night*

$\neg E_2$: *A large sum of money was found in his posession*

E_3 : *His fingerprints was found at the door of the bank.*

Probabilities express states-of-knowledge

$$E \equiv E_1 \text{ and } E_2 \text{ and } E_3$$

$$\underbrace{P(G|E)}_{> P(G|E_2)}$$

Binary propositions

A binary proposition is a statement which is either true or false (we might not know, but someone with complete knowledge would)

A : *In 49 BCE, Caesar crossed the Rubicon*

B : *Acceleration sensor 39 measures more than 0.85*

C : *Patient 901 has high cholesterol*

Propositions can be combined with **and**, **or** and **not**:

$AB \equiv$ True if A and B are both true

$A + B$ \equiv True if either A or B are true

$\bar{A} \equiv$ True if A is false

We define two special propositions which are always **true/false**:

1 : *A proposition which is always true*

0 : *A proposition which is always false*

...and the following identities: $A1 = A$, $A + \bar{A} = 1$, $\bar{\bar{A}} = A$ and

$$A(B_1 + B_2 + \cdots + B_n) = A\bar{B}_1 + AB_2 + \cdots + AB_n$$

Quiz 2, Probabilities

Assume we define the following 4 boolean variables.

R_1 : Handed in report 1

R_2 : Handed in report 2

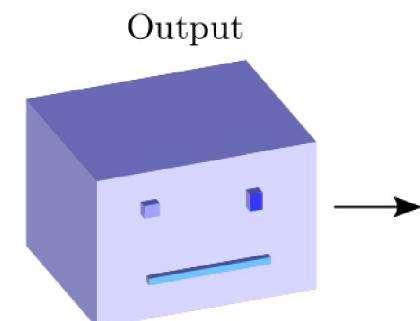
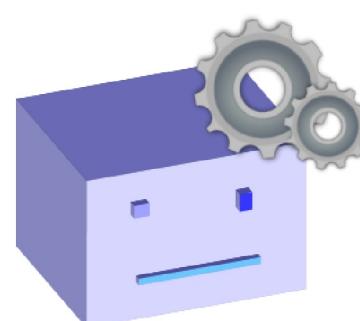
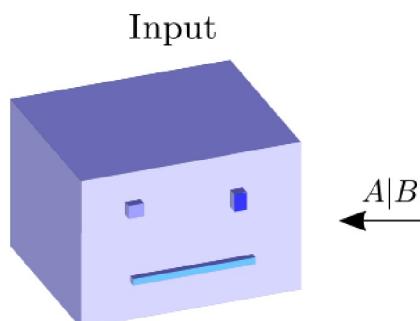
R_3 : Handed in report 3

F : Student failed 02450

- A. $P(R_1 R_2 R_3 | F) > 0.9$
- B. $P(\bar{F} | R_1 + R_2 + R_3) > 0.9$
- C. $P(\bar{F} | R_1 R_2 R_3) > 0.9$
- D. $P(R_1 + R_2 + R_3 | F) > 0.9$
- E. Don't know.

How would you express the probability of the statement:

If a student hand in report 1, 2 and 3, the chance of passing 02450 is greater than 90%?



Assuming B is true, how plausible is it A is true?

Assuming B , the plausibility of A is
(low / medium / high / certain)

Passing can be defined as not failing. Therefore, express the statement as:

$$P(\overline{F}|R_1R_2R_3) > 0.9$$

namely, if it is true report 1, 2 and 3 are all handed in, what is the chance of passing?

Rules of probability

The sum rule: $P(A|C) + P(\bar{A}|C) = 1$

The product rule: $P(AB|C) = P(B|AC)P(A|C)$

Interpretation:



$P(A|B) = 0$ (interpretation: given B is true, A is certainly false)

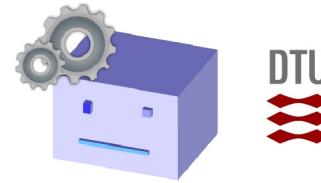
$P(A|B) = 1$ (interpretation: given B is true, A is certainly true)

We also use the shorthand:

$$P(A|1) = P(A)$$

$$\begin{aligned} P(A) + P(\bar{A}) &= 1 \\ P(AB) &= P(A|B)P(B) \end{aligned}$$

Remarkably, this is the mathematical basis for this course



Marginalization and Bayes' theorem

Sum rule

$$P(A|C) + P(\bar{A}|C) = 1$$

Product rule

$$P(AB|C) = P(B|AC)P(A|C) = p(A|BC)p(BC)$$

$$\begin{aligned} P(B|C) &= \underbrace{P(B|C)}_{\text{Product rule}} \left[\underbrace{P(A|BC)}_{\text{Sum rule}} + \underbrace{P(\bar{A}|BC)}_{\text{Sum rule}} \right] = \underbrace{P(AB|C)}_{\text{Product rule}} + \underbrace{P(\bar{A}B|C)}_{\text{Product rule}} \\ &= P(B|AC)P(A|C) + P(B|\bar{AC})P(\bar{A}|C). \end{aligned}$$

$$\text{Bayes' theorem: } P(A|BC) = \frac{P(B|AC)P(A|C)}{P(B|C)}$$

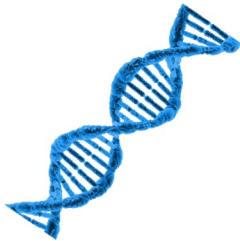
$$= \frac{P(B|AC)P(A|C)}{P(B|AC)P(A|C) + P(B|\bar{AC})P(\bar{A}|C)}.$$

Marginalization and Bayes' theorem

$$\begin{aligned} P(B|C) &= P(B|C) \left[P(A|BC) + P(\bar{A}|BC) \right] = P(AB|C) + P(\bar{A}B|C) \\ &= P(B|AC)P(A|C) + P(B|\bar{A}C)P(\bar{A}|C). \end{aligned}$$

$$\begin{aligned} P(A|BC) &= \frac{P(B|AC)P(A|C)}{P(B|C)} \\ &= \frac{P(B|AC)P(A|C)}{P(B|AC)P(A|C) + P(B|\bar{A}C)P(\bar{A}|C)}. \end{aligned}$$

DNA



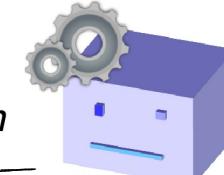
Bayes theorem

$$P(A|BC) = \frac{P(B|AC)P(A|C)}{P(B|AC)P(A|C) + P(B|\bar{A}C)P(\bar{A}|C)}$$

Crimes may be solved by matching crime-scene DNA to DNA in a database

- If the two samples are from the same person, a DNA test will always give a positive match
- If the DNA are from different persons, DNA will incorrectly give a positive match one time out of a million

A crime is committed in Racoon City by an unidentified male. Assume all 8000 possible perpetrators undergo a DNA test, and suppose the DNA test gives a positive result for George. What is the chance George is guilty?

$$G : \text{George is guilty}, \quad D : \text{There was a positive DNA match}$$
$$P(G|D) = \frac{P(D|G)P(G)}{P(D|G)P(G) + P(D|\bar{G})P(\bar{G})} = \frac{\frac{1}{8000}}{\frac{1}{8000} + 10^{-6}\left(1 - \frac{1}{8000}\right)} \approx 99\%$$


Solution:

$$\begin{aligned} P(G|D) &= \frac{P(D|G)P(G)}{P(D|G)P(G) + P(D|\bar{G})P(\bar{G})} \\ &= \frac{1 \times \frac{1}{8000}}{1 \times \frac{1}{8000} + 10^{-6} \times \left(1 - \frac{1}{8000}\right)} \\ &= 1 - \frac{1}{126} \approx 99\% \end{aligned}$$

Exclusive and exhaustive events

A_1 : The side \bullet face up.

A_3 : The side $\bullet\circ$ face up.

A_5 : The side $\circ\bullet$ face up.

A_2 : The side $\circ\circ$ face up.

A_4 : The side $\circ\bullet$ face up.

A_6 : The side $\bullet\bullet$ face up.

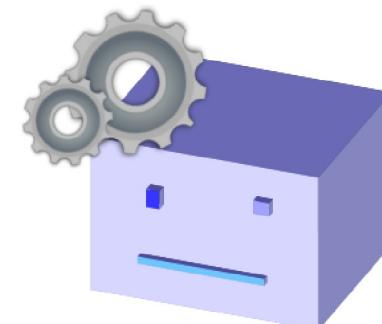
- When no two propositions can be true at the same time, they are said to be **mutually exclusive**: $\underline{A_i A_j = 0}$ for $i \neq j$

- Consider any two events A and B

$$\begin{aligned} P(A + B) &= P(A) + P(B) - P(A \cap B) \\ &\approx P(A) + P(B) - P(\text{both}) \end{aligned}$$

- In general, for n **mutually exclusive events**

$$P(\underbrace{A_1 + A_2 + \dots + A_n}_{\text{exhaustive}}) = \sum_{i=1}^n P(A_i)$$



- A set of events is **exhaustive** if one has to be true: $A_1 + \dots + A_n = 1$. Then:

$$\sum_{i=1}^n P(A_i) = P(\underbrace{A_1 + A_2 + \dots + A_n}_{\text{exhaustive}}) = 1$$

Exclusive and exhaustive events

A_1 : The side \bullet face up.

A_3 : The side $\circ\bullet$ face up.

A_5 : The side $\circ\circ$ face up.

A_2 : The side $\circ\bullet$ face up.

A_4 : The side $\bullet\bullet$ face up.

A_6 : The side $\bullet\circ$ face up.

- When no two propositions can be true at the same time, they are said to be **mutually exclusive**: $A_i A_j = 0$ for $i \neq j$
- Consider any two events A and B

$$\begin{aligned}
 P(A + B) &= 1 - P(\overline{A} \overline{B}) \\
 &= 1 - [1 - P(A|\overline{B})] P(\overline{B}) = P(B) + P(A\overline{B}) \\
 &= P(B) + P(\overline{B}|A)P(A) = P(B) + [1 - P(B|A)] P(A) \\
 &= P(A) + P(B) - P(AB).
 \end{aligned}$$

- In general, for n mutually exclusive events $P(A_1 + A_2 + \dots + A_n) = \sum_{i=1}^n P(A_i)$
- A set of events is **exhaustive** if one has to be true: $A_1 + \dots + A_n = 1$. Then:

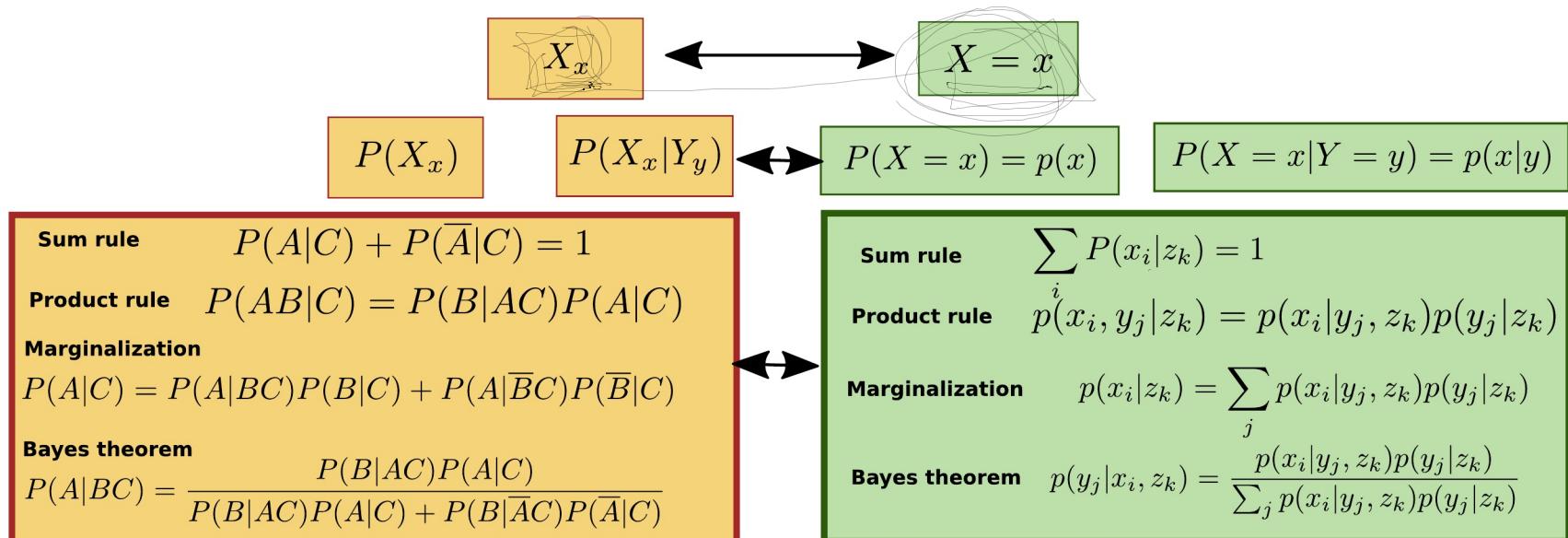
$$\sum_{i=1}^n P(A_i) = P(A_1 + A_2 + \dots + A_n) = 1$$

Stochastic variables

- Often, we will measure numerical quantities (number of children, age of a patient, etc.)
- Suppose a quantity X (number of children) takes a value $x = 3$. We can write this as the binary event X_3 and in general:

$X_x : \{ \text{The binary event that } X \text{ is equal to the number } x \}$

- Stochastic variable simplify this notation by the definition:



Quiz 3, Avila bible (Fall 2018)

$p(\tilde{x}_2, \tilde{x}_{10} y)$	$y = 1$	$y = 2$	$y = 3$
$\tilde{x}_2 = 0, \tilde{x}_{10} = 0$	0.19	0.3	0.19
$\tilde{x}_2 = 0, \tilde{x}_{10} = 1$	0.22	0.3	0.26
$\tilde{x}_2 = 1, \tilde{x}_{10} = 0$	0.25	0.2	0.35
$\tilde{x}_2 = 1, \tilde{x}_{10} = 1$	0.34	0.2	0.2

Table 1: Probability of observing particular values of \tilde{x}_2 and \tilde{x}_{10} conditional on y .

We will consider a dataset based on the Avila bible. We wish to predict the copyist ($y = 1, 2, 3$) of a bible based on the two typographic attributes *upperm* and *mr/is*. We suppose the attributes have been binarized such that *upperm* corresponds to $\tilde{x}_2 = 0, 1$ and *mr/is* to $\tilde{x}_{10} = 0, 1$. Suppose the probability for each of the configurations of \tilde{x}_2 and \tilde{x}_{10} conditional on the copyist y are as given in Table 1. and the prior probability of

Sum rule $\sum_i p(x_i|z_k) = 1$

Product rule $\prod^i p(x_i, y_j|z_k) = p(x_i|y_j, z_k)p(y_j|z_k)$

Marginalization $p(x_i|z_k) = \sum_j p(x_i|y_j, z_k)p(y_j|z_k)$

Bayes theorem $p(y_j|x_i, z_k) = \frac{p(x_i|y_j, z_k)p(y_j|z_k)}{\sum_j p(x_i|y_j, z_k)p(y_j|z_k)}$



the copyists is

$$p(y = 1) = 0.316, p(y = 2) = 0.356, p(y = 3) = 0.328.$$

Using this, what is then the probability an observation was authored by copyist 1 given that $\tilde{x}_2 = 1$ and $\tilde{x}_{10} = 0$?

A \wedge B

- A. $p(y = 1|\tilde{x}_2 = 1, \tilde{x}_{10} = 0) = 0.25$
- B. $p(y = 1|\tilde{x}_2 = 1, \tilde{x}_{10} = 0) = 0.313$
- C. $p(y = 1|\tilde{x}_2 = 1, \tilde{x}_{10} = 0) = 0.262$
- D. $p(y = 1|\tilde{x}_2 = 1, \tilde{x}_{10} = 0) = 0.298$
- E. Don't know.

$$p(y=1 | x_2, x_{10}) =$$

$$\frac{p(x_2, x_{10} | y=1) p(y=1)}{\sum_{k=1}^3 p(x_2, x_{10} | y=k) p(y=k)}$$

The problem is solved by a simple application of Bayes' theorem:

$$\begin{aligned} p(y = 1 | \tilde{x}_2 = 1, \tilde{x}_{10} = 0) \\ = \frac{p(\tilde{x}_2 = 1, \tilde{x}_{10} = 0 | y = 1)p(y = 1)}{\sum_{k=1}^3 p(\tilde{x}_2 = 1, \tilde{x}_{10} = 0 | y = k)p(y = k)} \end{aligned}$$

The values of $p(y)$ are given in the problem text and the values of $p(\tilde{x}_2 = 1, \tilde{x}_{10} = 0 | y)$ in ?? . Inserting the values we see option D is correct.

Independence

Independent: $p(x_i, y_j) = p(x_i)p(y_j)$

Conditionally independent given z_k : $p(x_i, y_j | z_k) = p(x_i | z_k)p(y_j | z_k)$

Expectations

$$\text{Expectation: } \mathbb{E}[f] = \sum_{i=1}^N f(x_i)p(x_i). \quad (2)$$

$f(x) \geq X$ $f(x) = (x - \mu)^2$

$$\text{mean: } \mathbb{E}[x] = \sum_{i=1}^N x_i p(x_i), \quad \text{Variance: } \text{Var}[x] = \sum_{i=1}^N (x_i - \mathbb{E}[x])^2 p(x_i). \quad (3)$$

Example: Uniform probability

$$p(x_i) = \frac{1}{N}$$

$$\mathbb{E}[f] = \frac{1}{N} \sum_{i=1}^N f(x_i)$$

$$\mathbb{E}[x] = \frac{\sum_{i=1}^N x_i}{N}$$

$$\text{Var}[x] = \frac{1}{N} \sum_{i=1}^N (x_i - \mathbb{E}[x])^2$$

Densities and models

- In machine learning, we want to learn a parameter from data
- Models of the data which use parameters are how we do that
- We build models our of simpler building blocks (densities).
In this course we will learn four:

Bernoulli density

The Categorical density

The Beta density

The Multivariate normal density

The Bernoulli density

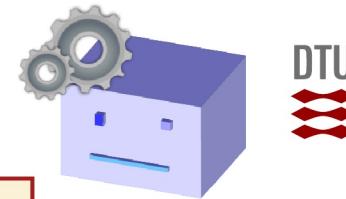
- Let $b = 0, 1$ denote a binary event.
- For instance, $b = 0$ corresponds to a person being well, and $b = 1$ that a person is ill.
- The probability of b is expressed using a parameter θ in the unit interval $[0, 1]$

Bernoulli distribution: $p(b|\theta) = \theta^b(1-\theta)^{1-b}$.

$$\boxed{p(b=1)} = p(b|\theta)$$

$$p(b=1) = \theta^1(1-\theta)^{1-1} = \theta$$

$$p(b=0) = \theta^0(1-\theta)^{1-0} = 1-\theta.$$



The Bernoulli density, repeated events

Conditional independence $p(x_i, x_j | z_k) = p(x_i | z_k)p(x_j | z_k)$

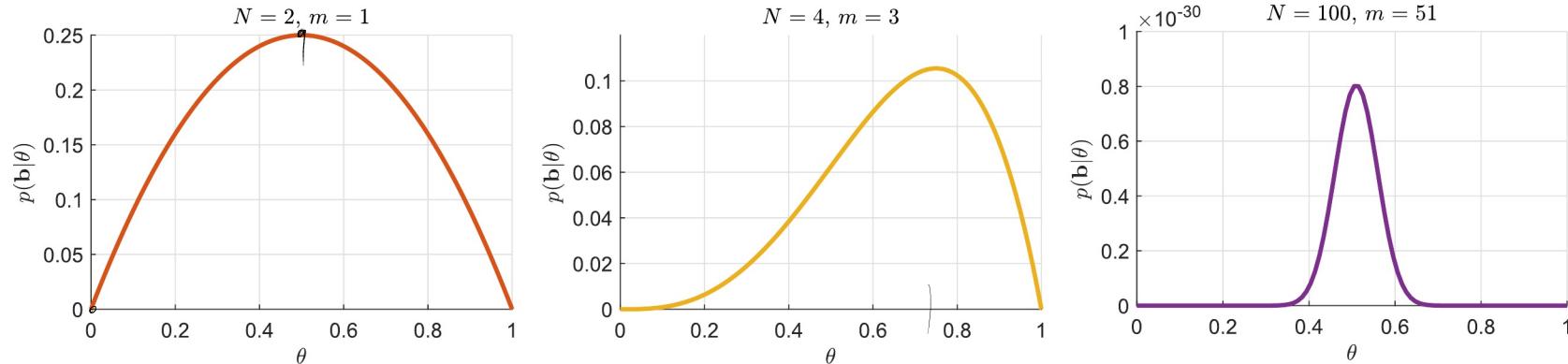
- Suppose we observe a sequence b_1, \dots, b_N of Bernoulli (binary) events.
- For instance, for N patients we record whether person 1 is well or ill ($b_1 = 0$ or $b_1 = 1$) and up to whether patient N is ill or well ($b_N = 0$ or $b_N = 1$)
- When we **know** θ (the chance a person is well or ill), the events are **independent**

Bernoulli distribution: $p(b|\theta) = \theta^b(1-\theta)^{1-b}$.
conditional independence.

$$p(b_1, \dots, b_N | \theta) = \underbrace{p(b_1 | \theta) \times \dots \times p(b_N | \theta)}_{\substack{\downarrow \\ \text{independence}}} = \prod_{i=1}^N \theta^{b_i} (1-\theta)^{1-b_i} = \theta^{\sum_{i=1}^N b_i} (1-\theta)^{N - \sum_{i=1}^N b_i}$$

$$= \theta^m (1-\theta)^{N-m}, \quad m = b_1 + b_2 + \dots + b_N$$

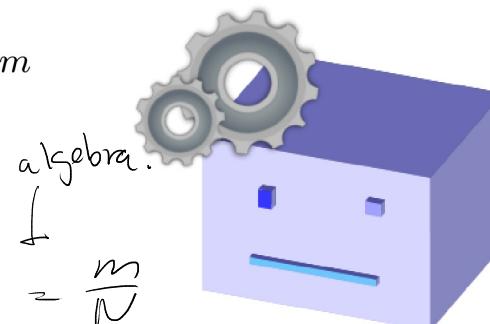
The Bernoulli density, maximum likelihood



$$p(b_1, \dots, b_N | \theta) = \theta^m (1 - \theta)^{N-m}$$

An idea for selecting θ^* is **Maximum likelihood**

$$\theta^* = \arg \max_{\theta} p(b_1, \dots, b_N | \theta) = \frac{m}{N}.$$



The value of θ according to which the data is most plausible

$$\mathcal{D} = (\mathbf{x}, \mathbf{y})$$

$$P(\mathbf{x}, \mathbf{y} | \theta)$$

$$\theta^* = \arg \max_{\theta} P(\mathbf{x}, \mathbf{y} | \theta)$$

Resources

<https://02402.compute.dtu.d> A more in-depth description of summary statistics (see chapter 1) (<https://02402.compute.dtu.dk>)

<https://www.khanacademy.org> An excellent introduction to probability theory which we recommend as a go-to resource
(<https://www.khanacademy.org/math/statistics-probability/probability-library>)

<http://citeseerx.ist.psu.edu> A more in-depth discussion of Bayes in the court room (<http://citeseerx.ist.psu.edu/viewdoc/download;jsessionid=EFE0328140036A4C668BB5B9FC76C9BE?doi=10.1.1.599.8675&rep=rep1&type=pdf>)