

DTU

02450: Introduction to Machine Learning and Data Mining

Measures of similarity, summary statistics and probabilities

Georgios Arvanitidis

DTU Compute, Technical University of Denmark (DTU)

DTU Compute
Department of Applied Mathematics and Computer Science

Today

Feedback Groups of the day:
Cecilie Haslund Bertelsen, Akshat Bhardwaj, Yann Xavier Frank Birnie-Scott, Mikalle Elouise Bisgaard-Bohr, Simon Bjerrregaard, Felipe Blanco Ricco, Andrea Gonzalez Boa, Rinske Boersma, Vinzenz Bohland, Morten Rahbek Boilesen, Andrea Madelena Lang Bolvig, Simon Bondø, Peter Bonvang, Anja Lykke Borre, Bertram Ladefoged Bottos, Oliver Brandt-Olsen, Malthe August Bordin Bresler, Marie Brinck-Jensen, Hjalte Sebastian Brinkløv, Signe Wobeser Brødsbaard, Jonas Bodulv Broge, Noah Bro-Jørgensen, Jacob Møller Bruhn, Carl Bruunsmose, Lucas Brylle, Anna Bzinkowska, Irene Campillo Pereda, Paula Campiõa Monzùn, Bogdan Capsa, Francesco Centomo, Ingel Daniel Cester Sala, Lai Wun Chan, Therapol Charoensuk, AROOJ CHAUDHRY, Andrew Chen, Tonny Chen, Rong Cheng, Ion Chetaru, Shirish Reddy Chintaguntla, Erik Buur Christensen, Aksel Buur Christensen, Jasper Brodner Christensen, Joachim Christian Christensen, Jakob Friis Christensen, Clara Sofie Christiansen, Christina Herlin Christiansen, Martin Christoffersen, Michał Wiktor Chrobot, Alejandro Cirugeda Pablos, Juan Collado, Arthur Conrado Veiga Bosquetti, Anna Cecilia Cordes, Euan Thomas Cortes, Jiyuan Cui, Olga Czajkowska

Reading material:
Chapter 4, Chapter 5

Lecture 3 13 September, 2022

DTU

Lecture Schedule

1 Introduction 30 August: C1 Data: Feature extraction, and visualization	8 Artificial Neural Networks and Bias/Variance 25 October: C14, C15
2 Data, feature extraction and PCA 6 September: C2, C3	9 AUC and ensemble methods 1 November: C16, C17 Unsupervised learning: Clustering and density estimation
3 Measures of similarity, summary statistics and probabilities 13 September: C4, C5	10 K-means and hierarchical clustering 8 November: C18
4 Probability densities and data visualization 20 September: C6, C7 Supervised learning: Classification and regression	11 Mixtures models and density estimation 15 November: C19, C20 (Project 2 due before 13:00)
5 Decision trees and linear regression 27 September: C8, C9	12 Association mining 22 November: C21 Recap
6 Overfitting, cross-validation and Nearest Neighbor 4 October: C10, C12 (Project 1 due before 13:00)	13 Recap and discussion of the exam 29 November: C1-C21
7 Performance evaluation, Bayes, and Naïve Bayes 11 October: C11, C13	

Online help: Piazza
Videos of lectures: <https://video.dtu.dk>
Streaming of lectures: Zoom (link on DTU Learn)

3 DTU Compute

DTU

Evaluation, interpretation, and visualization

Data: Data preparation
• Feature extraction
• Similarity measures
• Summary statistics
• Data visualization

Data modelling: Data modelling
• Classification
• Clustering
• Density estimation

Evaluation: Evaluation
• Anomaly detection
• Decision making
• Result visualization
• Dissemination

Result

Domain knowledge

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

4 DTU Compute

Lecture 3 13 September, 2022

DTU

PCA recap: Principal component analysis on images

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

Tamara Berg "Faces in the wild"

5 DTU Compute

DTU

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
- $1000 \times 22'188$

6 DTU Compute

Lecture 3 13 September, 2022

Principal component analysis (PCA)



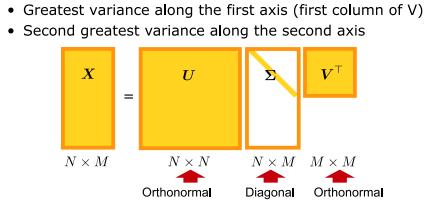
1. Subtract the mean

- 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



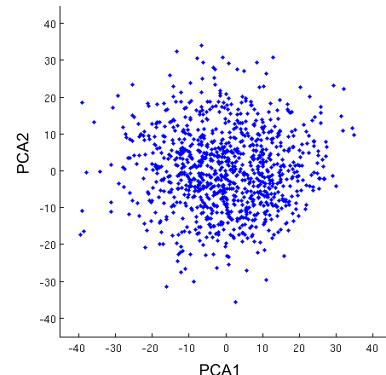
• Plot data in the transformed coordinate system

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

7 DTU Compute

Lecture 3 13 September, 2022

PCA on face images

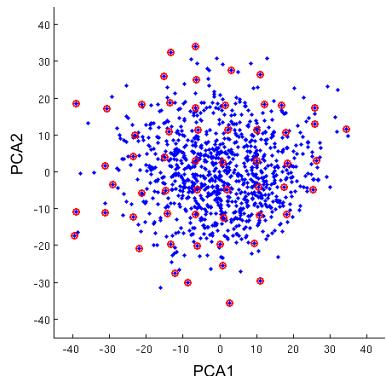


PCA1 PCA2

8 DTU Compute

Lecture 3 13 September, 2022

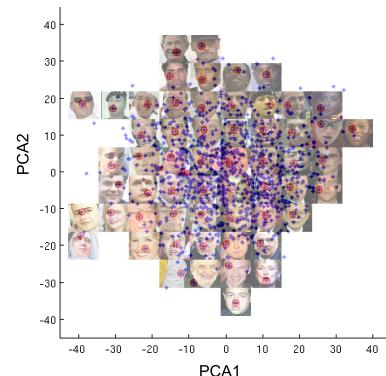
PCA on face images



9 DTU Compute

Lecture 3 13 September, 2022

PCA on face images

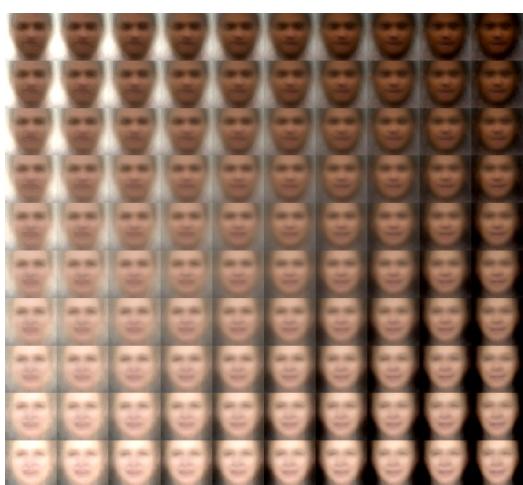


10 DTU Compute

Lecture 3 13 September, 2022



- What information do the two principal axes capture?

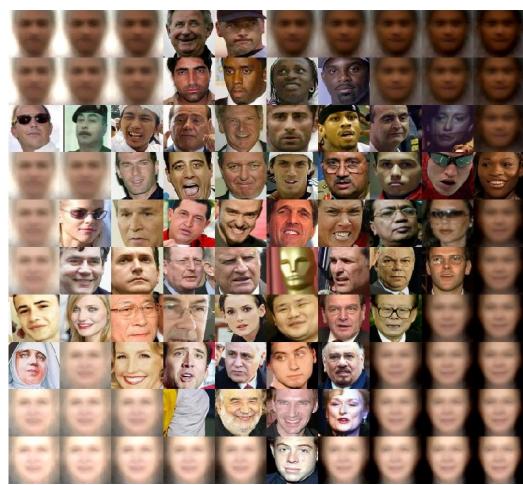


11 DTU Compute

Lecture 3 13 September, 2022



- What information do the two principal axes capture?



12 DTU Compute

Lecture 3 13 September, 2022

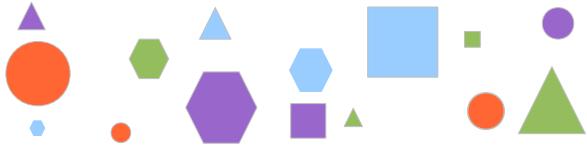
Similarity / Dissimilarity measures



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

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

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



13 DTU Compute

Lecture 3 13 September, 2022

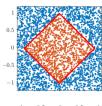
Dissimilarity measures



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

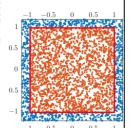
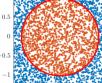
• One-norm ($p = 1$)

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



• Euclidean ($p = 2$)

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



• Max-norm distance ($p = \infty$)

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

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

14 DTU Compute

Lecture 3 13 September, 2022

Similarity measures

$$\mathbf{x} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \mathbf{y} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 1 \end{bmatrix}, K : \text{Total number of attributes}$$



f₀₀: Number of attributes where X_k=Y_k=0

f₁₁: Number of attributes where X_k=Y_k=1

Simple Matching Coefficient (SMC)

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

Jaccard Coefficient

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

Cosine similarity

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

Extended Jaccard coefficient

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

Also defined for continuous data

Lecture 3 13 September, 2022

Quiz 1, similarity measures



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

Which of the following statements are true?

- A. SMC($\mathbf{o}_1, \mathbf{o}_2$) = $\frac{3}{5}$, $J(\mathbf{o}_1, \mathbf{o}_2) = \frac{1}{2}$, $\cos(\mathbf{o}_1, \mathbf{o}_2) = \frac{2}{3}$,
- B. SMC($\mathbf{o}_1, \mathbf{o}_2$) = $\frac{3}{5}$, $J(\mathbf{o}_1, \mathbf{o}_2) = \frac{3}{4}$, $\cos(\mathbf{o}_1, \mathbf{o}_2) = \sqrt{\frac{2}{3}}$,
- C. SMC($\mathbf{o}_1, \mathbf{o}_2$) = $\frac{2}{5}$, $J(\mathbf{o}_1, \mathbf{o}_2) = \frac{1}{3}$, $\cos(\mathbf{o}_1, \mathbf{o}_2) = \frac{2}{3}$,
- D. SMC($\mathbf{o}_1, \mathbf{o}_2$) = $\frac{2}{5}$, $J(\mathbf{o}_1, \mathbf{o}_2) = \frac{1}{3}$, $\cos(\mathbf{o}_1, \mathbf{o}_2) = \sqrt{\frac{2}{3}}$,
- E. Don't know.

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

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

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

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

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

Lecture 3 13 September, 2022

Invariance

Scale invariance

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

Translation invariance

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

General invariances

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

17 DTU Compute

Lecture 3 13 September, 2022

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

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

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

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

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

Weighting Attributes have different importance

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

18 DTU Compute

Lecture 3 13 September, 2022

Empirical statistics

Given two samples $x_1, x_2, \dots, x_N \in \mathbb{R}$ and $y_1, y_2, \dots, y_N \in \mathbb{R}$:

- Empirical mean

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

- Empirical variance

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

- Empirical covariance

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

- Empirical standard deviation

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



19 DTU Compute

Lecture 3 13 September, 2022

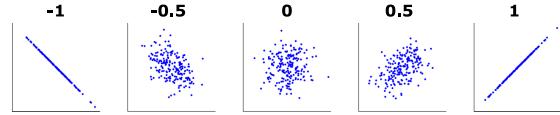
Correlation

- Measure of degree of linear relationship

$$\text{cor}[x, y] = \frac{\text{cov}[x, y]}{\hat{\sigma}_x \hat{\sigma}_y}$$

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

$$x_k = ay_k + b$$



20 DTU Compute

Lecture 3 13 September, 2022

Quantiles

Given N observations of an attribute $x_1, x_2, \dots, x_N \in \mathbb{R}$.

Quantiles describe the *points* that divide the underlying distribution into intervals that are equally probable:

- The one 2-quantile (**median**) divides the distribution in two intervals.
- The three 4-quantiles (**quartiles**) divides the distribution in four intervals.
- The 99 100-quantiles (**percentiles**) divides the distribution in 100 intervals.

The **median** is the same as the 2nd quartile or the 50th percentile.

E.g., we can (approximately) find the **median** by

- Sort the observations in ascending order $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(N)}$

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

21 DTU Compute

Lecture 3 13 September, 2022



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).

22 DTU Compute

Lecture 3 13 September, 2022

Probabilities



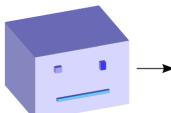
Input



$A|B$



Output



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

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

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

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

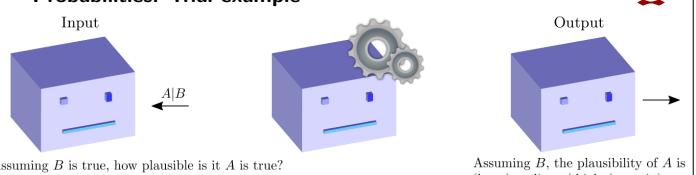
The degree-of-belief that A is true given B is accepted as true is at a level x

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

23 DTU Compute

Lecture 3 13 September, 2022

Probabilities: Trial example



G : The accused is guilty

E_1 : A car similar to his was seen at the crime scene.

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

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

$$P(G|E) > P(G|E_2)$$

24 DTU Compute

Lecture 3 13 September, 2022

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**:

$$\begin{aligned} AB &\equiv \text{True if } A \text{ and } B \text{ are both true} \\ A+B &\equiv \text{True if either } A \text{ or } B \text{ are true} \\ \bar{A} &\equiv \text{True if } A \text{ is false} \end{aligned}$$

We define two special propositions which is 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

$$A(B_1 + B_2 + \dots + B_n) = AB_1 + AB_2 + \dots + AB_n$$

25 DTU Compute

Lecture 3 13 September, 2022



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

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%?

- 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.



26 DTU Compute

Lecture 3 13 September, 2022

Rules of probability



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

$$\text{The product rule: } P(AB|C) = P(B|AC)P(A|C)$$

Interpretation:

$$\begin{aligned} P(A|B) = 0 &\quad (\text{interpretation: given } B \text{ is true, } A \text{ is certainly false}) \\ P(A|B) = 1 &\quad (\text{interpretation: given } B \text{ is true, } A \text{ is certainly true}) \end{aligned}$$

We also use the shorthand:

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

Remarkably, this is the mathematical basis for this course

27 DTU Compute

Lecture 3 13 September, 2022

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?

$$G : \text{George is guilty}, \quad D : \text{There was a positive DNA match}$$



29 DTU Compute

Lecture 3 13 September, 2022



Marginalization and Bayes' theorem



Marginalization and Bayes' theorem

Sum rule	$P(A C) + P(\bar{A} C) = 1$
Product rule	$P(AB C) = P(B AC)P(A C)$

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

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

28 DTU Compute

Lecture 3 13 September, 2022

Exclusive and exhaustive events



- A_1 : The side \square face up.
 - A_2 : The side \square face up.
 - A_3 : The side \square face up.
 - A_4 : The side \square face up.
 - A_5 : The side \square face up.
 - A_6 : The side \square 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

$$P(A + B) = P(A) + P(B) - P(AB)$$

- 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:

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

30 DTU Compute

Lecture 3 13 September, 2022

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:

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

- Stochastic variable simplify this notation by the definition:

Sum rule $P(A C) + P(\bar{A} C) = 1$	Product rule $P(AB C) = P(B AC)P(A C)$	Sum rule $\sum_i P(x_i z_k) = 1$
Marginalization $P(A C) = P(A BC)P(B C) + P(A \bar{B}C)P(\bar{B} C)$	Product rule $p(x_i, y_j z_k) = p(x_i y_j, z_k)p(y_j z_k)$	Marginalization $p(x_i z_k) = \sum_j p(x_i y_j, z_k)p(y_j z_k)$
Bayes theorem $P(A BC) = \frac{P(B AC)P(A C)}{P(B AC)P(A C) + P(B \bar{A}C)P(\bar{A} C)}$	Bayes theorem $p(y_j x_i, z_k) = \frac{p(x_i y_j, z_k)p(y_j z_k)}{\sum_j p(x_i y_j, z_k)p(y_j z_k)}$	

31 DTU Compute

Lecture 3 13 September, 2022



Quiz 3, Avila bible (Fall 2018)



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 $\hat{x}_2 = 1$ and $\hat{x}_{10} = 0$?

A. $p(y=1|\hat{x}_2=1, \hat{x}_{10}=0) = 0.25$

B. $p(y=1|\hat{x}_2=1, \hat{x}_{10}=0) = 0.313$

C. $p(y=1|\hat{x}_2=1, \hat{x}_{10}=0) = 0.262$

D. $p(y=1|\hat{x}_2=1, \hat{x}_{10}=0) = 0.298$

E. Don't know.

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

We will consider a dataset based on the Avila bible.

We wish to predict the copyist ($y = 0, 1, 2$) 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 $\hat{x}_2 = 0, 1$ and *mr/is* to $\hat{x}_{10} = 0, 1$. Suppose the probability for each of the configurations of \hat{x}_2 and \hat{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 $p(x_i, y_j z_k) = p(x_i y_j, z_k)p(y_j z_k)$
Marginalization $p(x_i z_k) = \sum_j p(x_i y_j, z_k)p(y_j z_k)$	Bayes theorem $p(y_j x_i, z_k) = \frac{p(x_i y_j, z_k)p(y_j z_k)}{\sum_j p(x_i y_j, z_k)p(y_j z_k)}$

Lecture 3 13 September, 2022

Independence



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

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

33 DTU Compute

Lecture 3 13 September, 2022

Expectations



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

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

Example: Uniform probability

$$\begin{aligned} p(x_i) &= \frac{1}{N} \\ \mathbb{E}[f] &= \frac{1}{N} \sum_{i=1}^N f(x_i) \\ \mathbb{E}[x] &= \frac{\sum_{i=1}^N x_i}{N} \\ \text{Var}[x] &= \frac{1}{N} \sum_{i=1}^N (x_i - \mathbb{E}[x])^2 \end{aligned}$$

Lecture 3 13 September, 2022

Densities and models



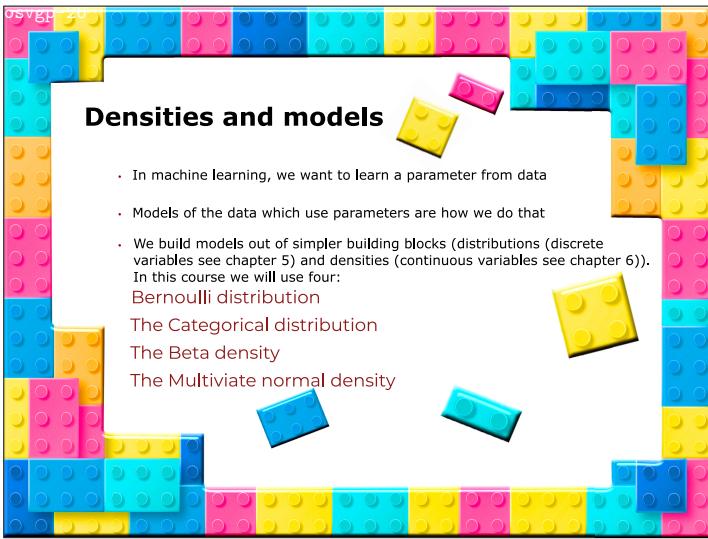
- In machine learning, we want to learn a parameter from data
- Models of the data which use parameters are how we do that
- We build models out of simpler building blocks (distributions (discrete variables see chapter 5) and densities (continuous variables see chapter 6)). In this course we will use four:

Bernoulli distribution

The Categorical distribution

The Beta density

The Multivariate normal density



The Bernoulli distribution



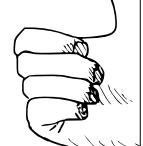
- Let $b = 0, 1$ denote a binary event.

- For instance,

- $b = 0$ corresponds to heads, and $b = 1$ to tails, or
- $b = 0$ corresponds to a person being ill, and $b = 1$ that a person is well.

- The probability of b is expressed using a parameter θ in the unit interval $[0, 1]$

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



36 DTU Compute

Lecture 3 13 September, 2022

The Bernoulli distribution, repeated events



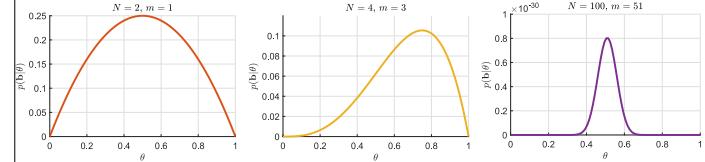
Conditional independence $p(x_i, x_j | z_k) = p(x_i | z_k)p(x_j | z_k)$

- 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 ill or well ($b_1 = 0$ 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: $p(b|\theta) = \theta^b(1-\theta)^{1-b}$.

$$\begin{aligned} p(b_1, \dots, b_N | \theta) &= \prod_{i=1}^N p(b_i | \theta) = \prod_{i=1}^N \theta^{b_i} (1-\theta)^{1-b_i} = \theta^{\sum_{i=1}^N b_i} (1-\theta)^{N - \sum_{i=1}^N b_i} \\ &= \theta^m (1-\theta)^{N-m}, \quad m = b_1 + b_2 + \dots + b_N \end{aligned}$$

The Bernoulli distribution, maximum likelihood



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

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

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

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



Resources



<https://bayes.wustl.edu> Classical textbook which treats probabilities as states-of-knowledge and discuss many practical and philosophical issues (this book converted me to ML!)

(<https://bayes.wustl.edu/etj/prob/book.pdf>)

<https://02402.compute.dtu.dk> 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=repTypePdf>)