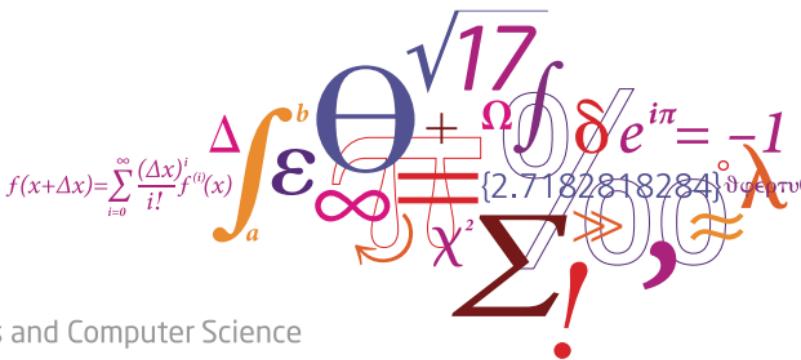


02450: Introduction to Machine Learning and Data Mining

Probability densities and data visualization

Bjørn Sand Jensen

DTU Compute, Technical University of Denmark (DTU)



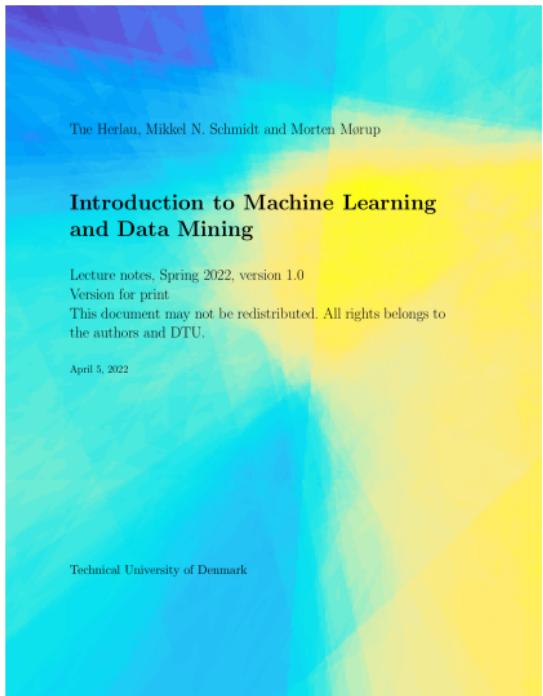
Today

Feedback Groups of the day:

Christoffer Kofoed Christensen, Jorge Mauricio Cisneros Caballero, Celina Lykke Clausen, Jens Kristian Dahl, Pietro Dal Cin, Jacob Quist Dalgaard, Snædís Lilja Daníelsdóttir, Lasse Schnell Danielsen, Barak Davidi, Yahir Alejandro de Leon Dominguez, Carlos De Santiago Leon, Larissa de Souza Dos Santos, Gaia Dellarole, Gemma Di Federico, Ann-Kathrin Dilfer, Denitsa Diyanova Dimitrova, Minh Tam Doan, Tobias Doll, Filipe Dos Santos Branco, Stavroula Douva, Tobias Drue, Kamma Dyhr, Marcus Alexander Dyhr, Mathias Hadi Dyhr, Adikrishna Murali mohan Efternavn, Kazi Ejajul, Sixten Matias Skovsted Eld, Marina Epitropaki, Ezel Daglar Erguden, Mads Dan Eriksen, Magnus Odd Olsen Erler, Sofie Ertel, Anton Espholm, Rolando Esquivel-Sancho, Ejnar Billeskov Exsteen, Latif Faghiri

Reading material:

Chapter 6, Chapter 7



Lecture Schedule

1 Introduction

31 January: C1

Data: Feature extraction, and visualization

2 Data, feature extraction and PCA

7 February: C2, C3

3 Measures of similarity, summary statistics and probabilities

14 February: C4, C5

4 Probability densities and data visualization

21 February: C6, C7

Supervised learning: Classification and regression

5 Decision trees and linear regression

28 February: C8, C9

6 Overfitting, cross-validation and Nearest Neighbor

7 March: C10, C12 (Project 1 due before 13:00)

7 Performance evaluation, Bayes, and Naive Bayes

14 March: C11, C13

8 Artificial Neural Networks and Bias/Variance

21 March: C14, C15

9 AUC and ensemble methods

28 March: C16, C17

Unsupervised learning: Clustering and density estimation

10 K-means and hierarchical clustering

11 April: C18

11 Mixture models and density estimation

18 April: C19, C20 (Project 2 due before 13:00)

12 Association mining

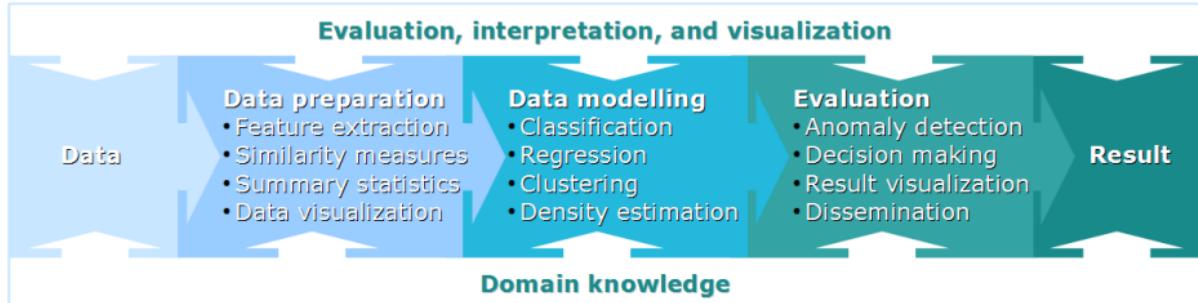
25 April: C21

Recap

13 Recap and discussion of the exam

2 May: C1-C21

Online 24/7 help: Discussion Forum/Piazza
Streaming & Videos: <https://panopto.dtu.dk/>
Online exercises: MS Teams



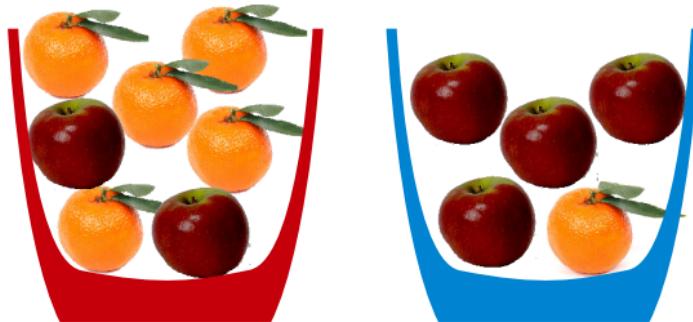
Learning Objectives

- Understand probability densities and related concepts
- Derive cost-functions from likelihood functions using Bayes' theorem
- Understand and apply a wide range of data visualization approaches
- Understand good practice in plotting including Tufte's guidelines

Probabilities recap

Example: Computing with probabilities

- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?



Apple taken from: https://upload.wikimedia.org/wikipedia/commons/3/32/Dark_apple.png
Orange (clementine) taken from: https://commons.wikimedia.org/wiki/File:Clementine_orange.jpg

Probabilities

- In more common notation we have

- Sum rule

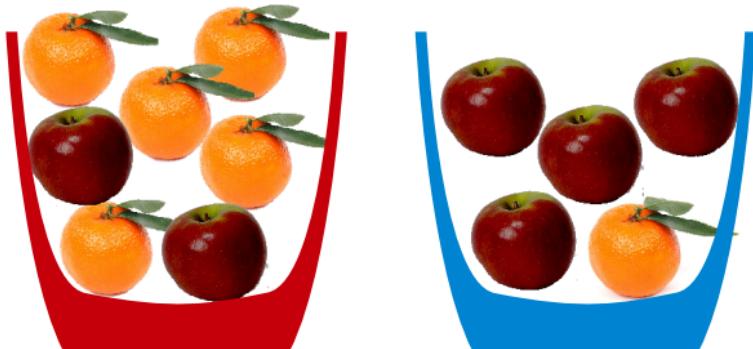
$$p(x) = p(x, y=0) + p(x, y=1)$$

- Product rule

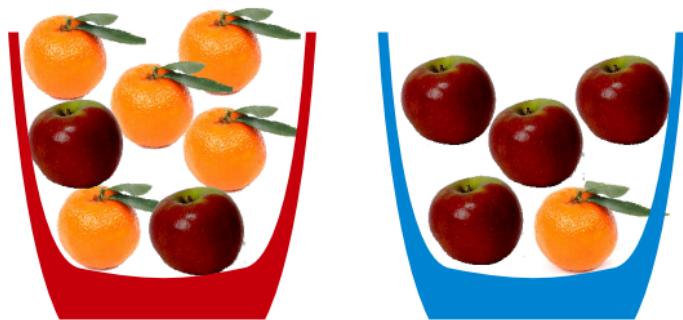
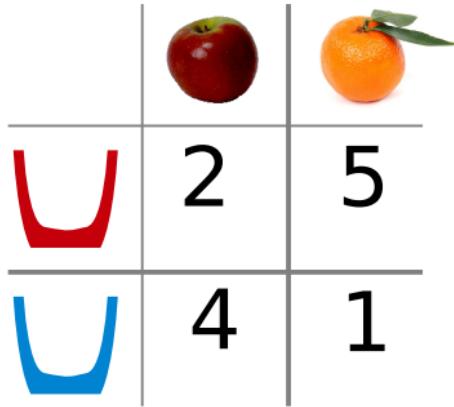
$$p(x, y) = p(x|y)p(y)$$

- Bayes' rule

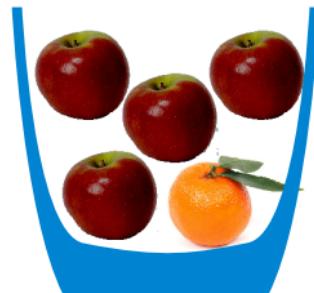
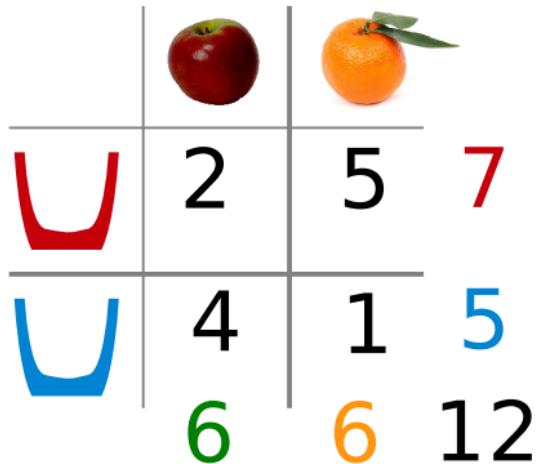
$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$



- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?



- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?



- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?

| | | | |
|---|---|---|---|
| | |  |  |
| U | 2 | 5 | 7 |
| U | 4 | 1 | 5 |
| 6 | 6 | 12 | |

$$p(o|r) = \frac{p(r,o)}{p(r)} = \frac{5/12}{7/12} = 5/7$$

- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?

| | | | |
|---|---|---|---|
| | |  |  |
| U | 2 | 5 | 7 |
| U | 4 | 1 | 5 |
| 6 | 6 | 12 | |

$$p(o|r) = \frac{p(r,o)}{p(r)} = \frac{5/12}{7/12} = 5/7$$

$$p(r|o) = \frac{p(r,o)}{p(o)} = \frac{5/12}{6/12} = 5/6$$

- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?

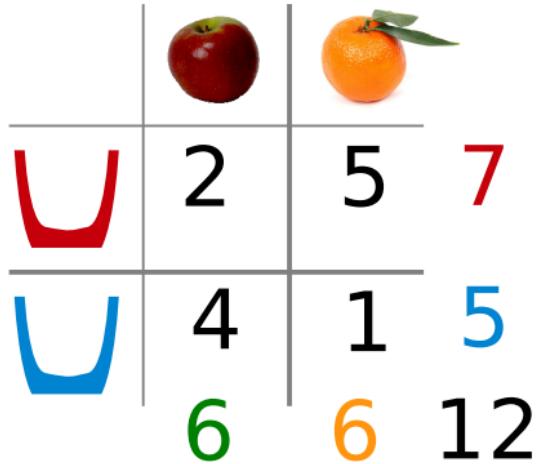
| | | | |
|---|---|---|---|
| | |  |  |
| U | 2 | 5 | 7 |
| U | 4 | 1 | 5 |
| 6 | 6 | 12 | |

$$p(o|r) = \frac{p(r,o)}{p(r)} = \frac{5/12}{7/12} = 5/7$$

$$p(r|o) = \frac{p(r,o)}{p(o)} = \frac{5/12}{6/12} = 5/6$$

$$= \frac{p(o|r)p(r)}{p(o)}$$

- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?



$$p(o|r) = \frac{p(r,o)}{p(r)} = \frac{5/12}{7/12} = 5/7$$

$$p(r|o) = \frac{p(r,o)}{p(o)} = \frac{5/12}{6/12} = 5/6$$

$$= \frac{p(o|r)p(r)}{p(o)}$$

$$= \frac{5/7 \cdot 7/12}{6/12} = 5/6$$



Medical test

$$\begin{aligned} P(D|+) &= P(D,+)/P(+) \\ &= P(+|D)*P(D) / (P(+|D)*P(D) + P(+|N)*P(N)) \\ &= 0.99*0.01 / (0.99*0.01 + 0.02*0.99) = 1/3 \end{aligned}$$

A medical test for a given disease

- Correctly identifies the disease 99% of the time (true positives), and
- Incorrectly turns out positive 2% of the time (false positives).

You know that

- 1% of the population suffers from the disease.

$$P(+|D) = 0.99$$

$$P(+|N) = 0.02$$

$$P(-|D) = 0.01$$

$$P(-|N) = 0.98$$

$$P(D) = 0.01$$

$$P(N) = 0.99$$

You go to the doctor to get tested, and the test turns out to be positive.

What is the probability you have the disease?

Hints:

- Identify from the text: ($x=Positive$,
 $y=0: no\ disease$, $y=1: Disease$)

$$p(Positive|Disease)$$

$$p(Positive|No\ Disease)$$

$$p(Disease)$$

$$p(No\ Disease)$$

- Use the basic rules of probability given to the right to find:

$$p(Disease|Positive)$$

$$\begin{aligned} p(y) &= \sum_x p(y, x) \\ &= p(y|x)p(x) + p(y|\bar{x})p(\bar{x}) \end{aligned}$$

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

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

Quiz 1, Probabilities (Spring 2014)

Consider a dataset which describe the consumption of delicatessen products in different cities. Each observation in the dataset is a customer, and we record the city the customer is from as well as their consumption of delicatessen. Suppose you are told:

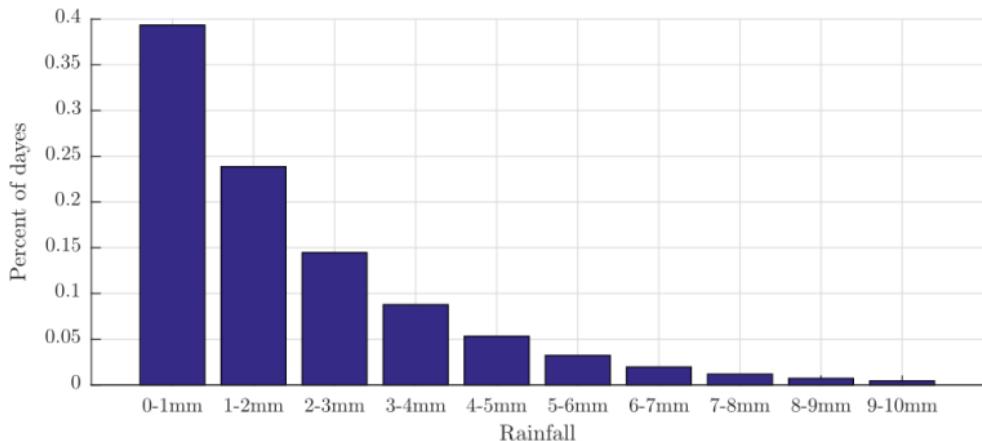
- 17.5 % were from Lisbon, 10.7 % were from Oporto and 71.8 % from the Other region.
- 44.1 % of the costumers from Lisbon spent above the median consumption on delicatessen (DELI).
- 48.9 % of the costumers from Oporto spent above the median consumption on delicatessen (DELI).
- 51.6 % of the costumers from the Other region spent above the median consumption on delicatessen (DELI).

What is the probability based on the wholesale data that a costumer that spent above the median consumption on delicatessen (DELI) come from Lisbon?

- A. 7.7 %
- B. 15.4 %**
- C. 44.1 %
- D. 59.6 %
- E. Don't know.

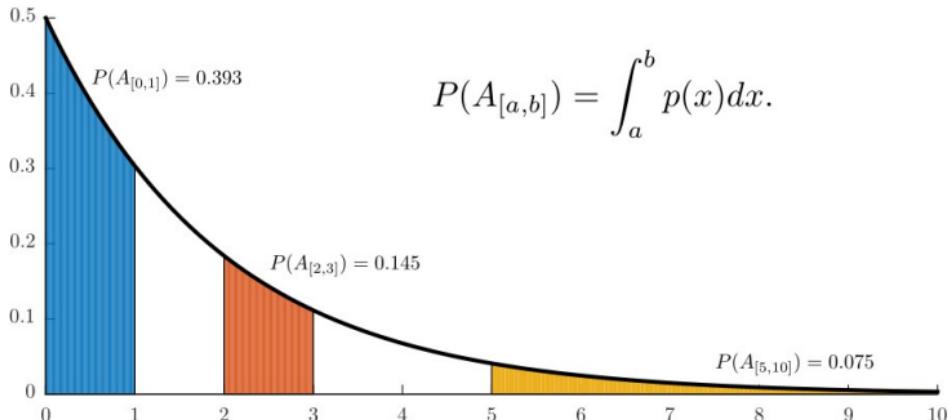
Probability vs. Density

- Suppose we consider the rainfall on an average day r
- **Can't** talk about the probability there will be **exactly** $r=2.3$ mm of rain, $P(r=2.3\text{mm})$
- **Can** talk about the probability there will be **between** 1 and 2 mm of rain



Probability vs. Density

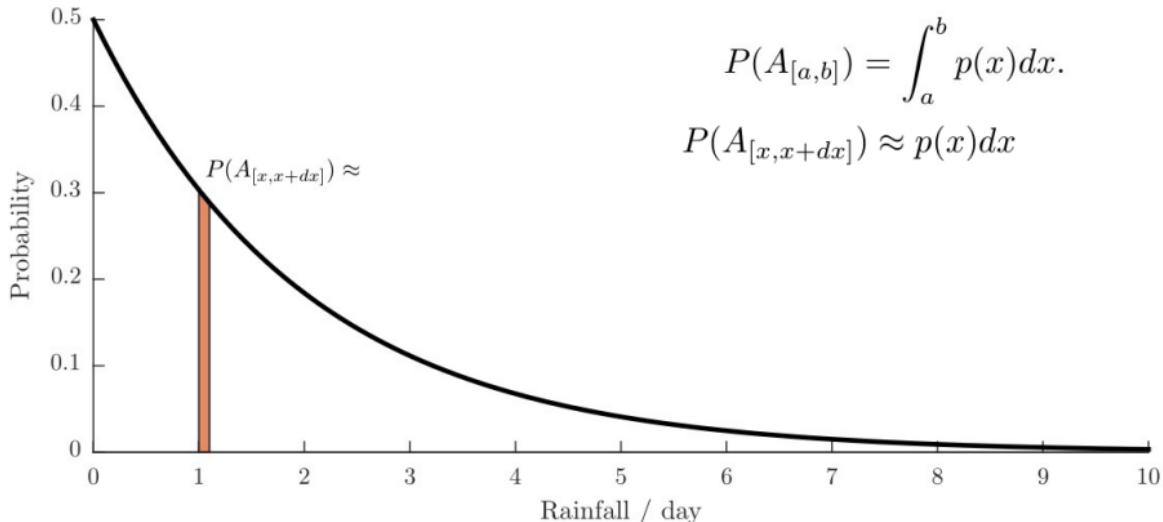
- These probabilities can be **represented** as integrals
- **Events** are **intervals**, the **probability** is the **integral**, the **curve** is the **density**



$A_{[a,b]}$: There will be between a and b mm of rain

Probability vs. Density

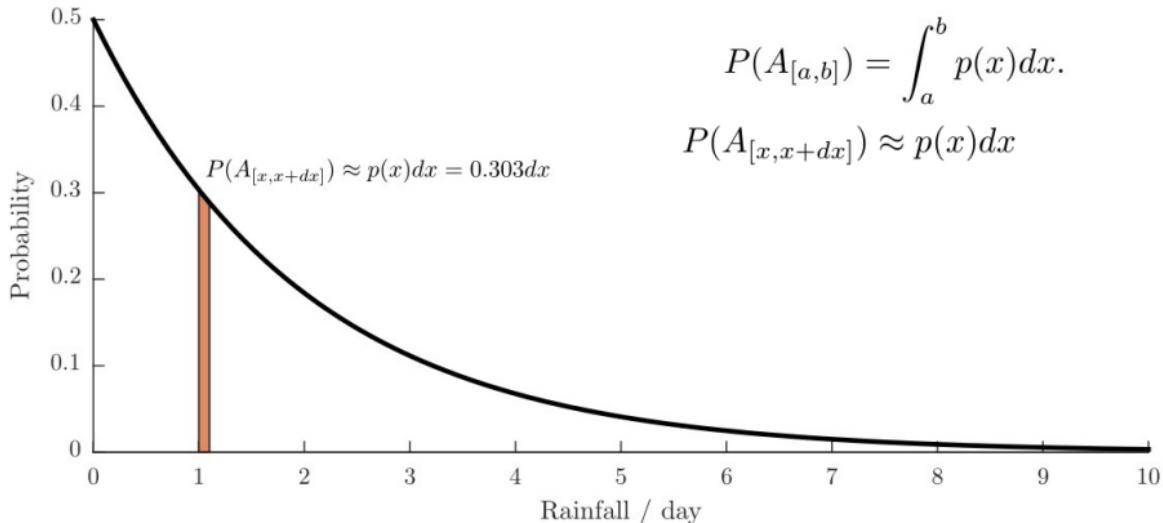
- These probabilities can be **represented** as integrals
- **Events** are **intervals**, the **probability** is the **integral**, the **curve** is the **density**
- What is the probability there will be between 1 and 1.1 mm of rain?



$A_{[a,b]}$: There will be between a and b mm of rain

Probability vs. Density

- These probabilities can be **represented** as integrals
- **Events** are **intervals**, the **probability** is the **integral**, the **curve** is the **density**
- What is the probability there will be between 1 and 1.1 mm of rain?



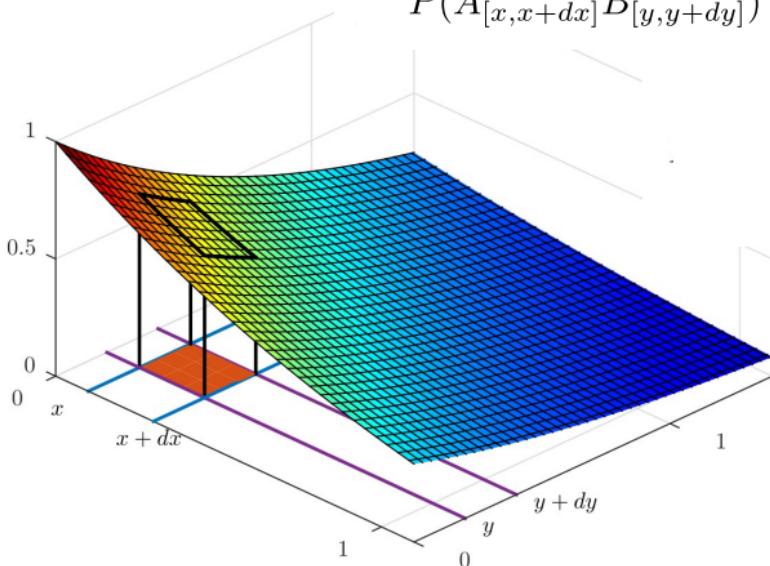
$A_{[a,b]}$: There will be between a and b mm of rain

Probability vs. Density

- For two variables x and y , the **probability** is an integral over an **area**

$$P((x, y) \in D) = \int_{(x,y) \in D} p(x, y) dx dy$$

$$\hat{P}(A_{[x,x+dx]} B_{[y,y+dy]}) =$$



This implies:

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

Probability vs. Density

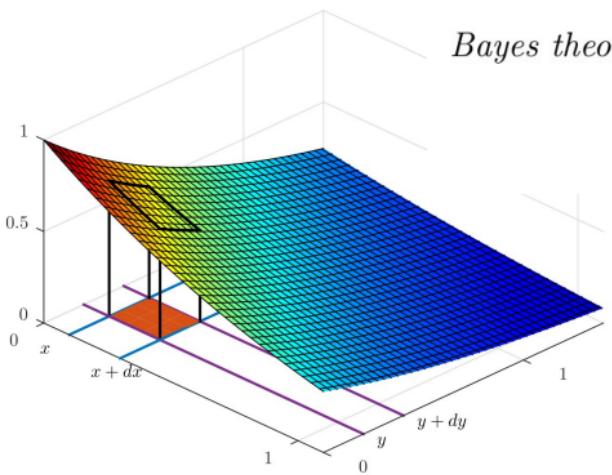
- Thus, we have shown the rules of probability theory also holds for densities

The sum rule:

$$\int dx \ p(x|z) = 1$$

The product rule:

$$p(x, y|z) = p(y|x, z)p(x|z)$$



$$\begin{aligned} p(x|y, z) &= \frac{p(y|x, z)p(x|z)}{p(y|z)} \\ &= \frac{p(y|z)p(x|y, z)}{\int p(y|x', z)p(x'|z)dx'}. \end{aligned}$$

Collecting all of this we obtain:

- Rules of probability for densities

Marginalization:

$$\int p(x, y|z)dx = p(y|z)$$

The product rule:

$$p(x, y|z) = p(y|x, z)p(x|z)$$

Bayes theorem:

$$p(x|y, z) = \frac{p(y|x, z)p(x|z)}{\int p(y|x', z)p(x'|z)dx'}.$$

- Rules of probability for discrete variables

Marginalization:

$$\sum_c p(x = c, y|z) = p(y|z)$$

The product rule:

$$p(x, y|z) = p(y|x, z)p(x|z)$$

Bayes theorem:

$$p(x|y, z) = \frac{p(y|x, z)p(x|z)}{\sum_c p(y|x = c, z)p(x = c|z)}.$$

Expected values

- Discrete random variable

$$\mathbb{E}[g] = \sum_i g(x_i)P(x_i)$$

- Continuous random variable

$$\mathbb{E}[g] = \int_{-\infty}^{\infty} g(x)p(x)dx$$

Statistics

- **Mean**

$$\bar{x} = \mathbb{E}[x]$$

- **Covariance**

$$\text{cov}(x, y) = \mathbb{E}[(x - \bar{x})(y - \bar{y})]$$

- **Variance**

$$\text{var}(x) = \text{cov}(x, x) = \mathbb{E}[(x - \bar{x})^2]$$

- **Standard deviation**

$$\text{std}(x) = \sqrt{\text{var}(x)}$$

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 multivariate normal distribution

A distribution for M -dimensional vectors \boldsymbol{x} :

$$\mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{M}{2}} \sqrt{|\boldsymbol{\Sigma}|}} e^{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})}$$

$$M=1: \quad \mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

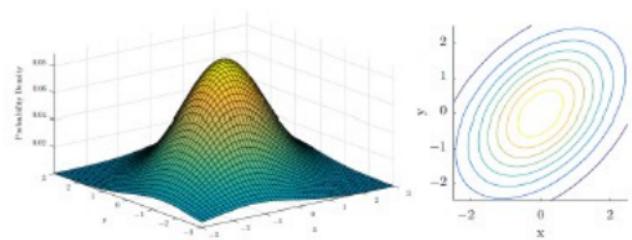
$\boldsymbol{\mu}$ is the mean vector and $\boldsymbol{\Sigma}$ is the covariance matrix:

$$\boldsymbol{\mu} = \mathbb{E}[\boldsymbol{x}], \quad \Sigma_{ij} = \text{cov}[x_i, x_j]$$

- Example: 2-dimensional Normal distribution

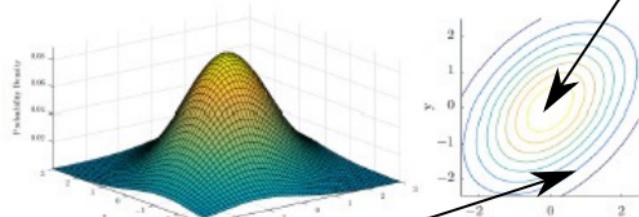
$$\boldsymbol{\mu} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \end{bmatrix}$$

$$\boldsymbol{\Sigma} = \begin{bmatrix} \text{var}(x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{var}(x_2) \end{bmatrix}$$



Quiz 2, Covariance

- Match the covariances to the contour plots



$$\Sigma = \begin{bmatrix} \text{var}(x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{var}(x_2) \end{bmatrix}$$

$$\mu = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \end{bmatrix}$$

A. Covariance of A is $\Sigma_A = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$, $\Sigma_C = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

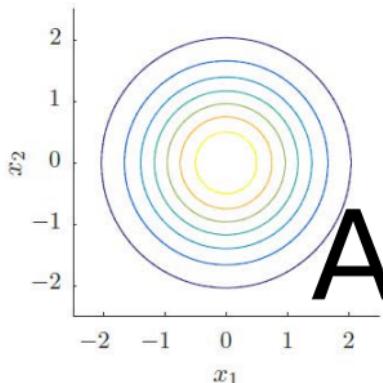
B. $\Sigma_B = \begin{bmatrix} 1 & -0.45 \\ 0.45 & 1 \end{bmatrix}$, $\Sigma_C = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$

C. $\Sigma_B = \begin{bmatrix} 10 & 4.5 \\ 4.5 & 10 \end{bmatrix}$, $\Sigma_C = \begin{bmatrix} 10 & 9 \\ 9 & 10 \end{bmatrix}$

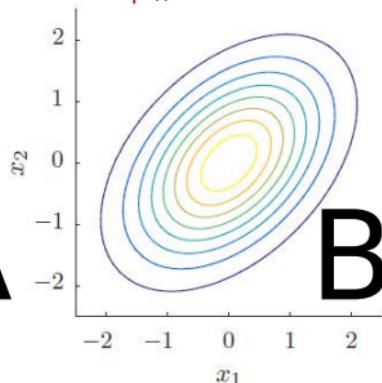
D. $\Sigma_A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$, $\Sigma_C = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$

E. Don't know.

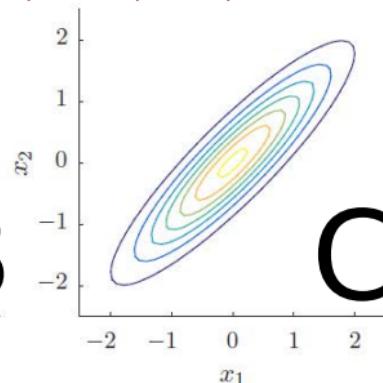
Check out the online demo <http://www2.imm.dtu.dk/courses/02450/DemoNormal.html>



A



B



C

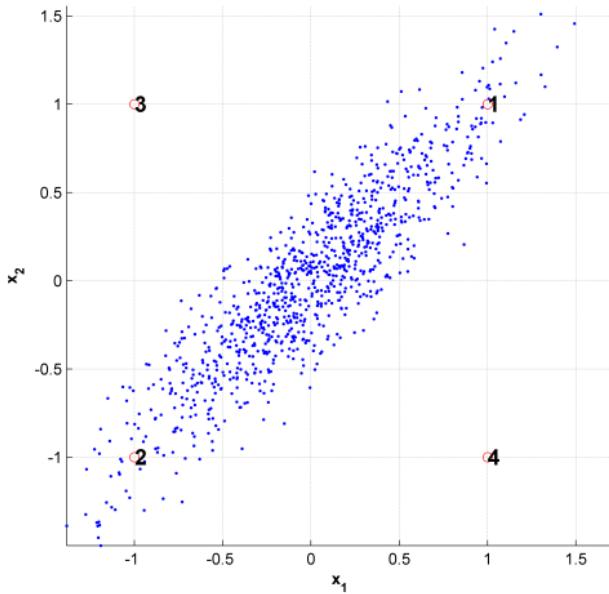
The Mahalanobis distance

How far are x_1 and x_2 apart?

- mahalanobis(x_1, x_2) = 4.2
- $d_{\text{Euclidean}}(x_1, x_2)^2$ = 8.0

How far are x_3 and x_4 apart?

- mahalanobis(x_3, x_4) = 80
- $d_{\text{Euclidean}}(x_3, x_4)^2$ = 8.0



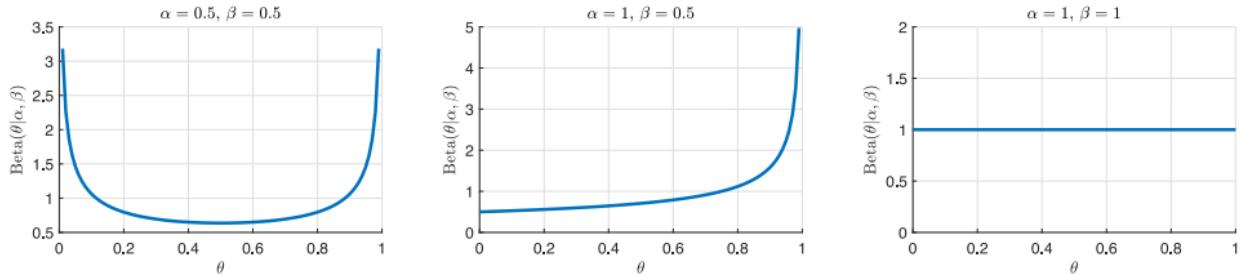
$$\text{mahalanobis}(x, y)^2 = (x - y)^\top \Sigma^{-1} (x - y)$$

$$d_{\text{euclidian}}(x, y)^2 = (x - y)^\top I^{-1} (x - y)$$

Beta distribution

Suppose θ is defined on the unit interval $[0, 1]$

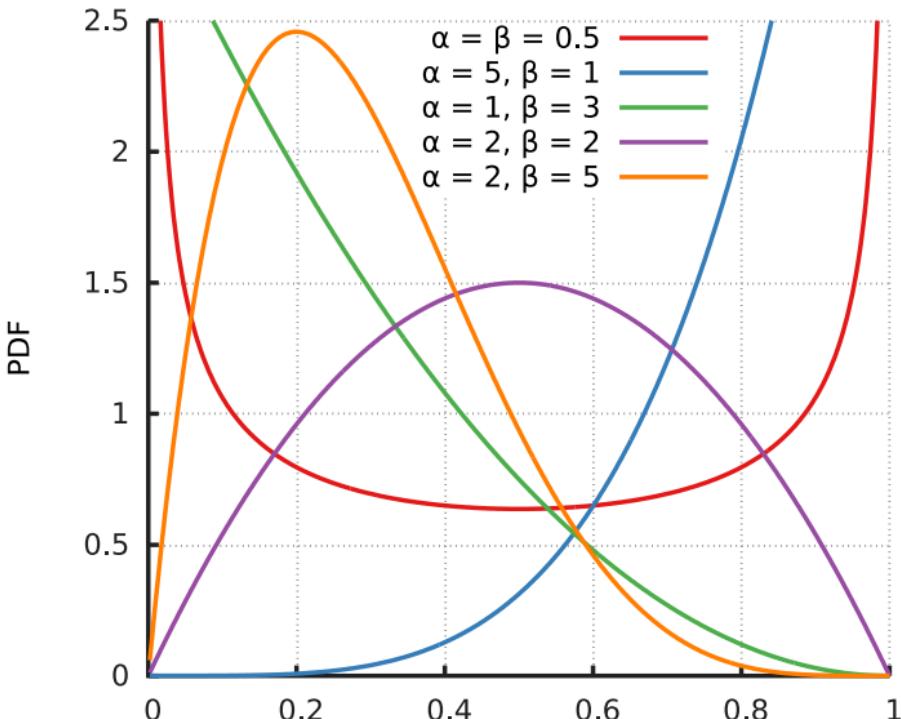
Beta density: $p(\theta|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)}\theta^{\alpha-1}(1-\theta)^{\beta-1}.$



$\alpha, \beta > 0$ are related to the variance and mean

$$\mathbb{E}_{p(\theta|\alpha, \beta)}[\theta] = \frac{\alpha}{\alpha + \beta}, \quad \text{Var}_{p(\theta|\alpha, \beta)}[\theta] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}.$$

Beta distribution



Probabilities and learning

- Your friend just got a dog.
- The dog can either be in the doghouse or outside
- The first four times you come by the dog is in the doghouse
- What is the chance the dog is in the doghouse tomorrow?



- Your friend buys a coin.
- The coin can either come up heads or tails
- The first four times you flip the coin it comes up tails
- What is the chance the coin comes up tails in the next flip?



Intuition tells us the answers are different, but the situation seems similar...

Recall from lecture 3: The Bernoulli distribution

- Suppose a coin come up heads with probability θ
 - Suppose $b = 1$ is the event the coin land heads
 - $b = 0$ is the event the coin land tails
- The density is given by the **Bernoulli distribution**

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

- For a sequence of N flips b_1, b_2, \dots, b_N **Independence**

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

- **What is θ ?**

The Bernoulli distribution

- A magic coin is a coin that comes up heads with probability θ
 - Suppose $b = 1$ is the event the coin land heads
 - $b = 0$ is the event the coin land tails
- For a sequence of N flips b_1, b_2, \dots, b_N

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

- **What is θ ?** Answer: **Use Bayes' Theorem!**

$$p(\theta|\mathbf{b}) =$$

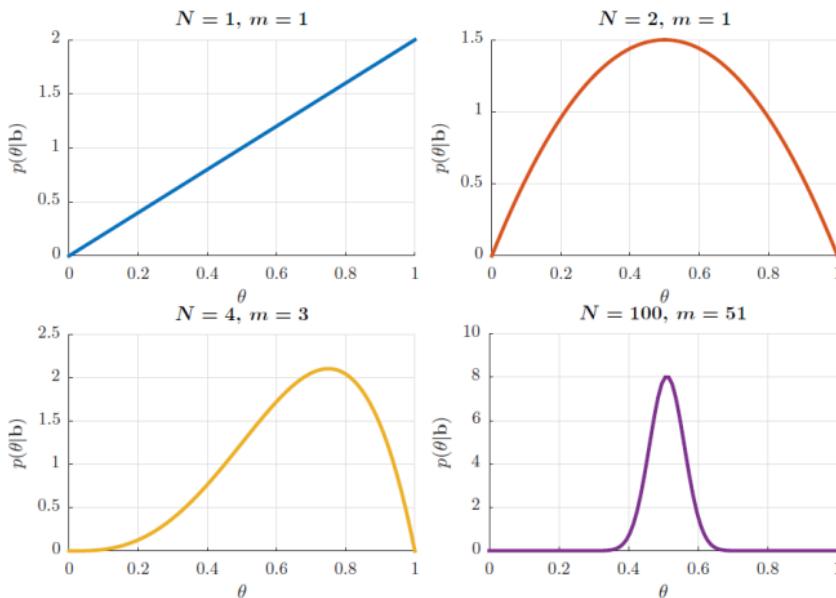
- Assume $p(\theta) =$

$$p(\theta|\mathbf{b}, \alpha, \beta) =$$

$$= \frac{\Gamma(\alpha + \beta + N)}{\Gamma(\alpha + m)\Gamma(\beta + N - m)} \theta^{\alpha+m-1} (1 - \theta)^{\beta+N-m-1}$$

Example: $\alpha = \beta = 1$

$$\begin{aligned}
 p(\theta | \mathbf{b}, \alpha, \beta) &= \frac{\Gamma(\alpha + \beta + N)}{\Gamma(\alpha + m)\Gamma(\beta + N - m)} \theta^{\alpha+m-1} (1-\theta)^{\beta+N-m-1} \\
 &= \frac{(N+1)!}{m!(N-m)!} \theta^m (1-\theta)^{N-m}
 \end{aligned}$$



Dogs and coins



- Your friend just got a dog.
- The dog can either be in the doghouse or outside
- The first four times you come by the dog is in the doghouse
- What is the chance the dog is in the doghouse tomorrow?

- Your friend buys a coin.
- The coin can either come up heads or tails
- The first four times you flip the coin it comes up tails
- What is the chance the coin comes up tails in the next flip?



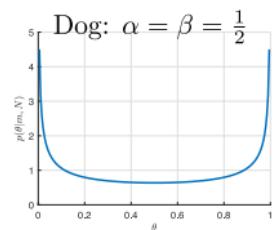
Dogs and coins



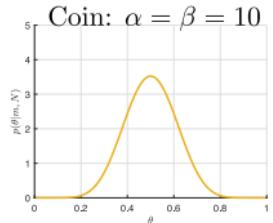
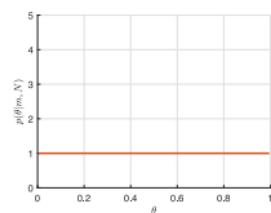
- Your friend just got a dog.
- The dog can either be in the doghouse or outside
- The first four times you come by the dog is in the doghouse
- What is the chance the dog is in the doghouse tomorrow?

Prior

$$p(\theta|\alpha, \beta) =$$



- Your friend buys a coin.
- The coin can either come up heads or tails
- The first four times you flip the coin it comes up tails
- What is the chance the coin comes up tails in the next flip?



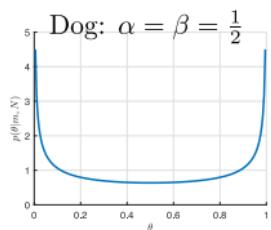
Dogs and coins



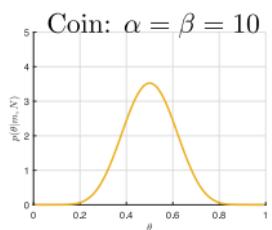
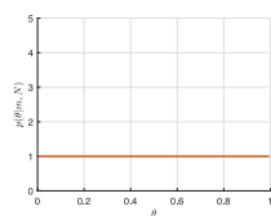
- Your friend just got a dog.
- The dog can either be in the doghouse or outside
- The first four times you come by the dog is in the doghouse
- What is the chance the dog is in the doghouse tomorrow?

Prior

$$p(\theta|\alpha, \beta) =$$



- Your friend buys a coin.
- The coin can either come up heads or tails
- The first four times you flip the coin it comes up tails
- What is the chance the coin comes up tails in the next flip?



Likelihood

$$p(m=4, N=4|\theta) = \theta^m(1-\theta)^{N-m} = \theta^4$$

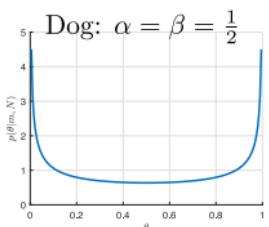
Dogs and coins



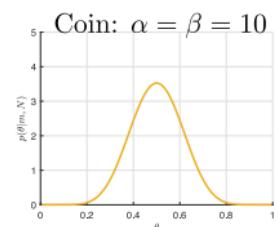
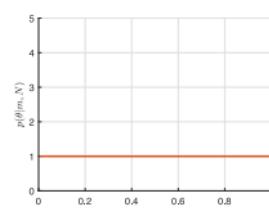
- Your friend just got a dog.
- The dog can either be in the doghouse or outside
- The first four times you come by the dog is in the doghouse
- What is the chance the dog is in the doghouse tomorrow?

Prior

$$p(\theta|\alpha, \beta) =$$



- Your friend buys a coin.
- The coin can either come up heads or tails
- The first four times you flip the coin it comes up tails
- What is the chance the coin comes up tails in the next flip?



Likelihood

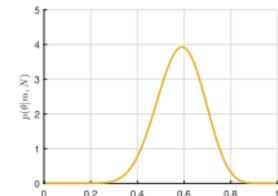
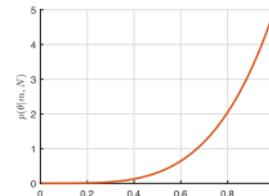
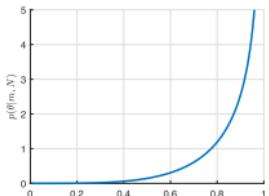
$$p(m=4, N=4|\theta) = \theta^m(1-\theta)^{N-m} = \theta^4$$

The difference between the two cases is that we have prior knowledge which tell us most coins are fair, and this affects our conclusions.

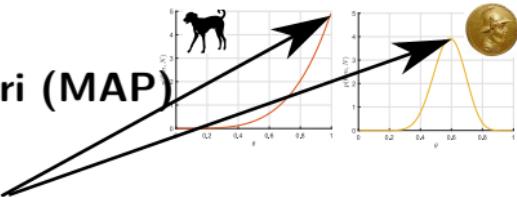
In most practical situations, we should assume as little as possible and choose $\alpha = \beta = \frac{1}{2}$

Posterior

$$p(\theta|m, N) = \frac{p(m, N|\theta)p(\theta|\alpha, \beta)}{p(m, N)} =$$



Learning principle: Maximum a posteriori (MAP)



- Another idea: Select θ which is "most probably"

$$\theta^* = \arg \max_{\theta} p(\theta | M, N) = \arg \max_{\theta} \left[\frac{p(m, N | \theta) p(\theta | \alpha, \beta)}{p(M, N)} \right]$$

- Use that $\arg \max_x f(x) = \arg \min_x [-\log f(x)]$ if $f(x) > 0$:

$$\theta^* = \arg \min_{\theta} \left[-\log \frac{p(m, N | \theta) p(\theta | \alpha, \beta)}{p(m, N)} \right]$$

(likelihood)

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

(prior)

$$p(\theta | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

- All in all:

$$\theta^* = \arg \min E(\theta), \quad E(\theta) = -\log p(m, N | \theta) - \log p(\theta | \alpha, \beta)$$

Maximum a posteriori (MAP) learning

- Consider some data $\mathbf{X} = \mathbf{x}_1, \dots, \mathbf{x}_N$ and $\mathbf{y} = y_1, \dots, y_N$
- Suppose we think x_i relates to y_i by some parameters $\boldsymbol{\theta}$
- Assume**

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \prod_{i=1}^N p(y_i|\mathbf{x}_i, \mathbf{w}), \quad p(\mathbf{w}|\mathbf{X}) = p(\mathbf{w})$$

Observations are not informative about each other when we know parameters

Without \mathbf{y} , we cannot learn the parameters

- Then

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})}$$

- The following are equivalent:

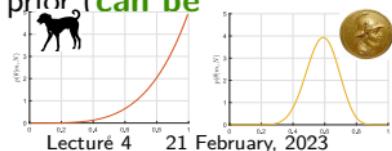
$$\text{Maximize: } \mathbf{w}^* = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathbf{X}, \mathbf{y})$$

$$\text{Minimize: } \mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w}), \quad E(\mathbf{w}) = \left[\frac{1}{N} \sum_{i=1}^N -\log p(y_i|\mathbf{x}_i, \mathbf{w}) \right] - \frac{1}{N} \log p(\mathbf{w})$$

- All we need is a likelihood (**usually pretty simple**) and a prior (**can be omitted**) and we have a machine-learning method.

- Pro:** Easy, conceptually simple, efficient

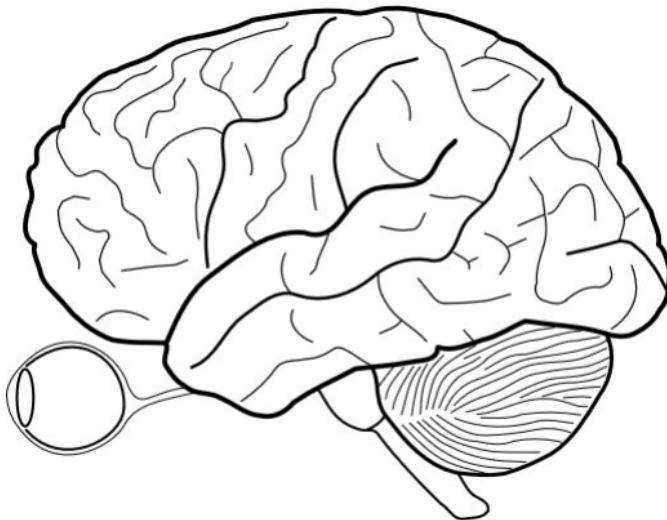
- Con:** Can sometimes give spurious results (overfit)



The drawing shows me at one glance what might be spread over ten pages in a book."
- Ivan S. Turgenev's novel Fathers and Sons, 1862.
Use a picture. It's worth a thousand words."
- Arthur Brisbane to the Syracuse Advertising Men's Club, in March 1911

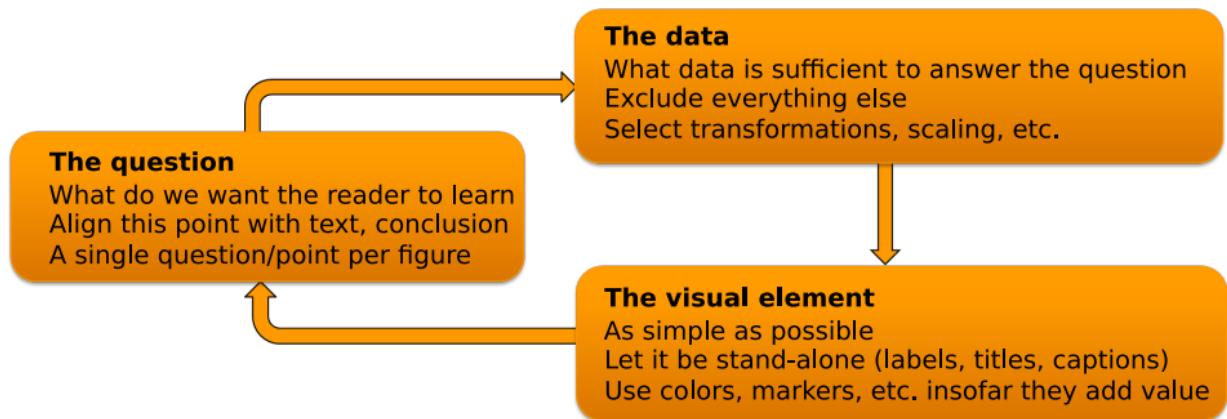
Visualization

- A main function of the brain is to process visual information
- We can exploit this capacity using visualization of the data in order to:
 - Detect new patterns, i.e. exploratory data analysis)
 - **Dissiminate results, i.e. visualizations/plots in written work (today)**
- We should take into account how the brains visual system works



Illustrations as technical writing

- The purpose of the text is to communicate an idea (*vs. plots has a purpose*)
- Be grammatically correct (*vs. elementary "rules" of good plotting*)
- Ensure the text is readable (*vs. labels, legends or lines nobody can read*)
- Avoid long/complicated paragraph (*vs. plots that are overly complicated*)
- Dont lie or exaggerate. (*vs. distort data in a plot*)

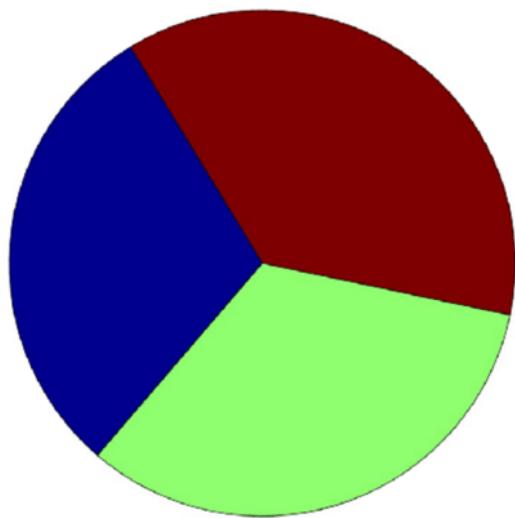


Important choices for visualizations

- **Representation:** How will you map objects, attributes, and relations to visual elements?
 - Positions, lengths, colors, areas, orientation
- **Arrangement:** How will you display the visual elements?
 - Viewpoint, transparency, separation, grouping
- **Selection:** How will you handle a large number of attributes and data objects?
 - Display a subset, focus on a region of interest, show summaries

Representation

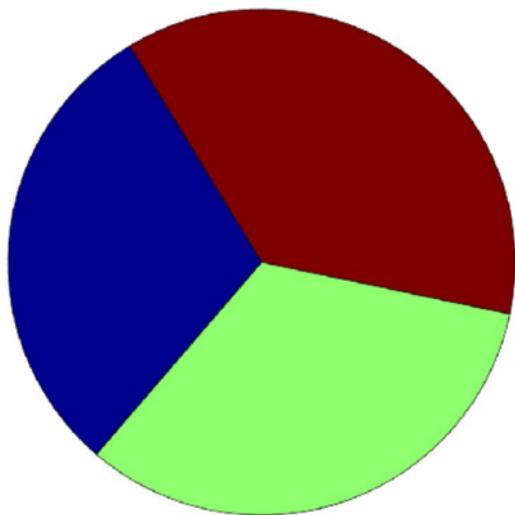
- **Area represents proportion**
 - Which is smallest, middle, and largest?
 - What are the proportions approximately?



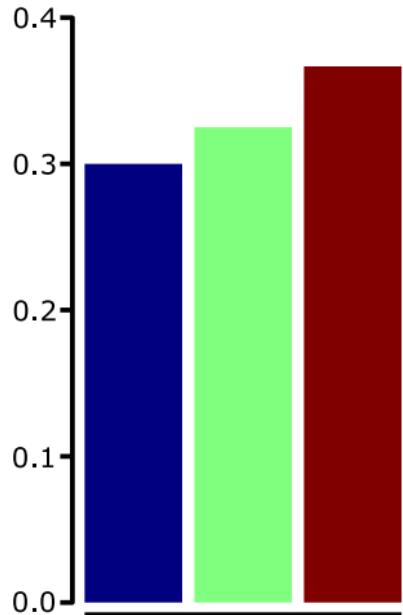
Representation

- **Area represents proportion**

- Which is smallest, middle, and largest?
- What are the proportions approximately?



- **Height represents proportion**



Selection

- Elimination or de-emphasis of certain objects or attributes
- A subset of **attributes**
 - **Why?** A graph can only show so many attributes – focus on the relevant
 - **How?**
 - Dimensionality reduction
 - Plot pairs of attributes
- A subset of **objects**
 - **Why?** A graph can only show so many objects – focus on the relevant
 - **How?**
 - Random sampling
 - Display of region of interest
 - Use density estimation

Types of plots

- **Distribution of a single attribute**

- Histogram
- Empirical cumulative distribution
- Percentile plots
- Box plot

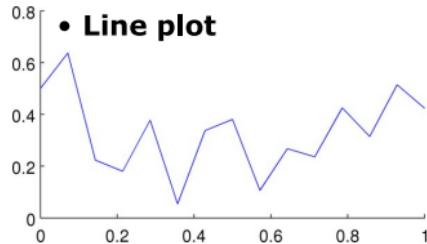
- **Relation between attributes**

- 2D histogram
- Heat maps and contour plots
- Scatter plots

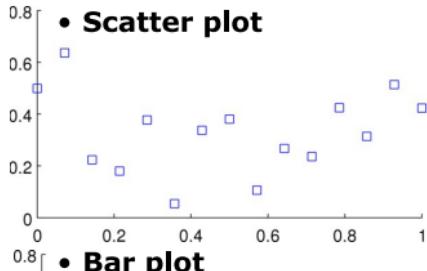
- **Visualization of high-dimensional objects**

- Matrix plots
- Parallel coordinates
- Star plots

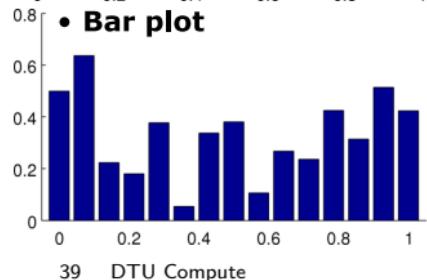
Basic plots



```
plot(x,y);
```



```
plot(x,y, 's');  
scatter(x,y, 's')
```



```
bar(x,y);
```

The iris data set

- **Three flowers**

- 50 instances of each class, 150 in total

- **Attributes**

- Sepal (outermost leaves)

- length in cm
- width in cm

- Petal (innermost leaves)

- length in cm
- width in cm

- Class of flower

- Iris Setosa
- Iris Versicolour
- Iris Virginica

| Flower ID | Attribute | | | |
|-----------|--------------|-------------|--------------|-------------|
| | Sepal Length | Sepal Width | Petal Length | Petal Width |
| 1 | 5.1 | 3.5 | 1.4 | 0.2 |
| 2 | 4.9 | 3.0 | 1.4 | 0.2 |
| 3 | 4.7 | 3.2 | 1.3 | 0.2 |
| 4 | 4.6 | 3.1 | 1.5 | 0.2 |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| 150 | 5.9 | 3.0 | 5.1 | 1.8 |

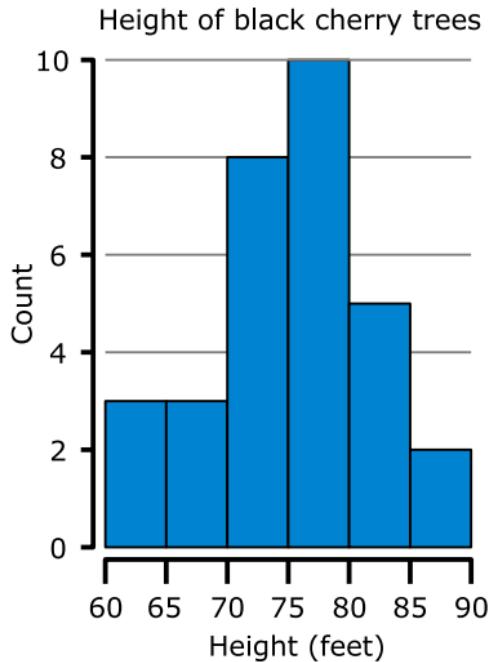
$X^{\text{Observation} \times \text{Attribute}}$

Distribution of a single attribute

Histograms

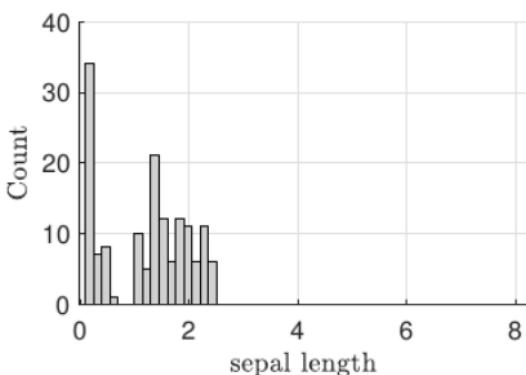
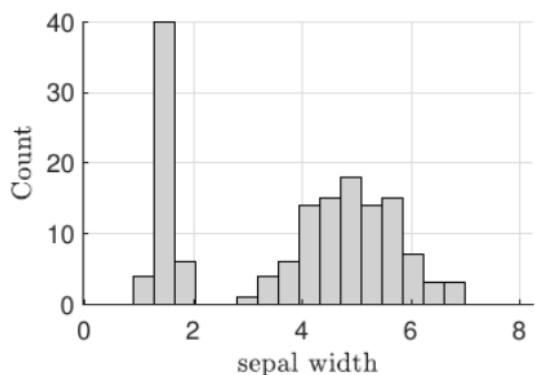
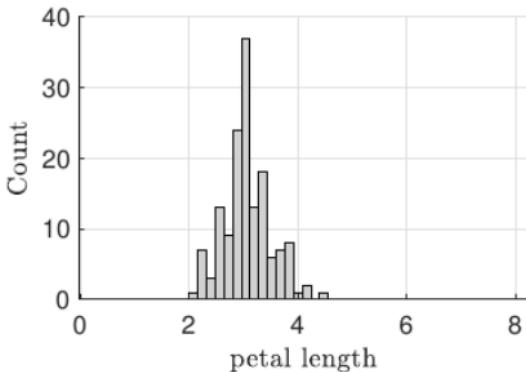
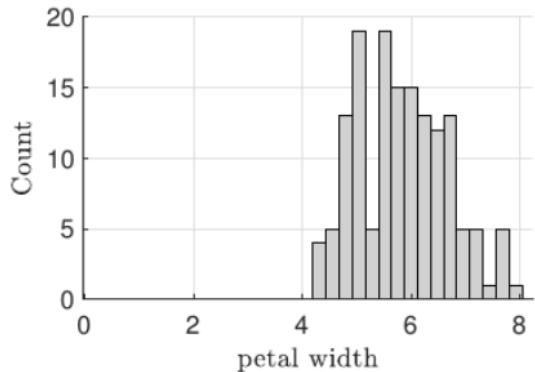
- Shows distribution of a single variable

- Divide the values into bins
- Bar plot of the number of values in bin
- Height indicates count of values
- Shape determined by
 - Distribution of data
 - Number of bins / bin width

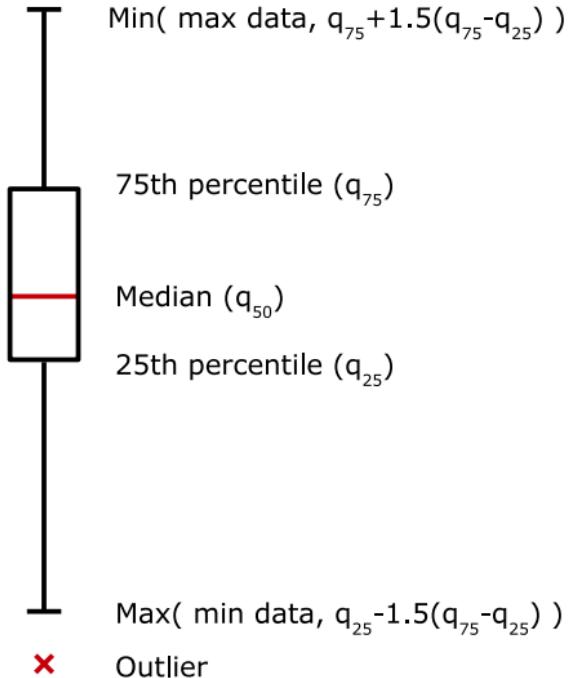


$$H = \{60, 64, 64, 66, 67, 69, 71, 72, 72, 72, 72, 73, 74, 74, 75, 75, 76, 76, 77, 77, 78, 78, 79, 80, 80, 81, 82, 84, 85, 85, 89\}$$

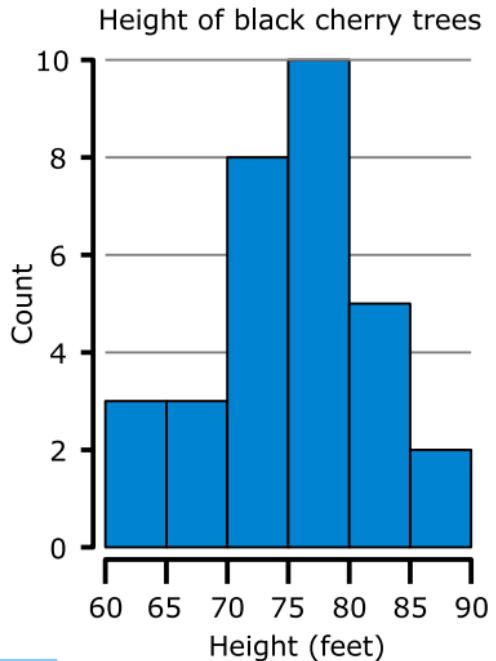
Histograms of the Iris data attributes



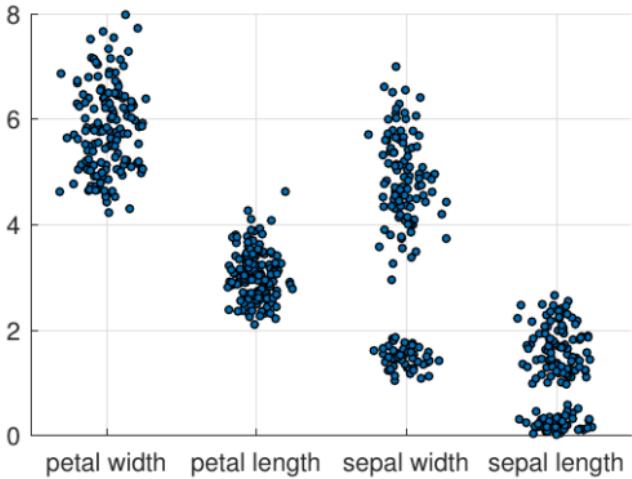
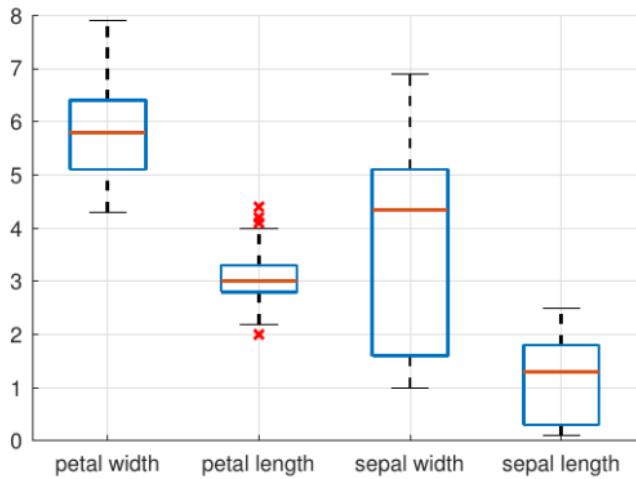
Box plots



The plotted whisker extends to the adjacent value, which is the most extreme data value that is not an outlier.



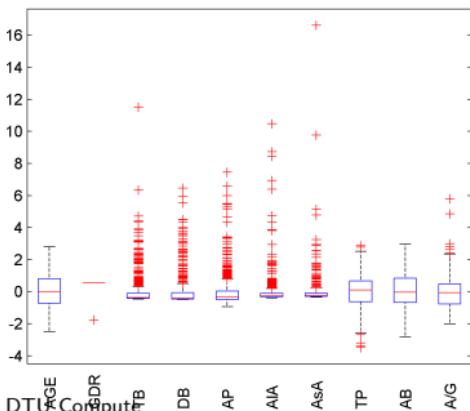
Box plots



Quiz 3, Boxplots (Fall 2012)

| No. | Attribute description | Abbrev. |
|----------|-------------------------------------|---------|
| x_1 | Age (in years) | AGE |
| x_2 | Gender (Female=0, Male=1) | GDR |
| x_3 | Total Bilirubin | TB |
| x_4 | Direct Bilirubin | DB |
| x_5 | Alkaline Phosphotase | AP |
| x_6 | Alamine Aminotransferase | AIA |
| x_7 | Aspartate Aminotransferase | AsA |
| x_8 | Total Proteins | TP |
| x_9 | Albumin | AB |
| x_{10} | Albumin to Globulin ratio | A/G |
| y | 0=No liver disease, 1=Liver disease | LD |

Table 1: Liver disease dataset.



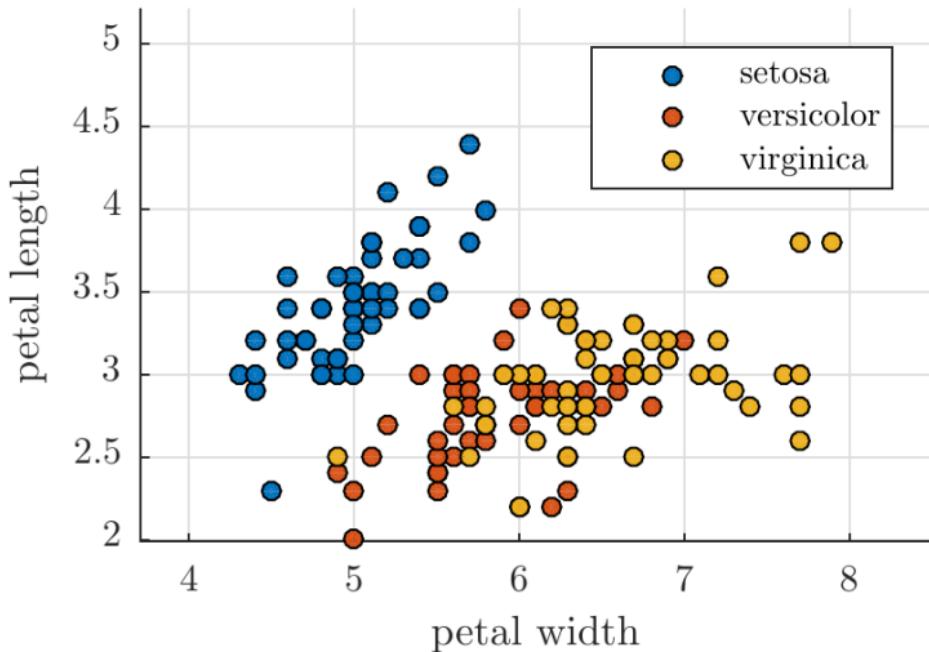
The attributes $x_1 \dots x_{10}$ are standardized (i.e., the mean has been subtracted each attribute and the attributes divided by their standard deviations). The figure shows a boxplot for the standardized data. Which of the following statements is *correct*?

- A. The value of the 50th and 75th percentiles of the attribute DB coincides.
- B. Even though the distribution of AIA and AsA may have a similar shape this does not imply that the two attributes are correlated.
- C. The attribute TB is likely to be normal distributed.
- D. The attribute GDR has a clear outlier that should be removed.
- E. Don't know.

Relation between attributes

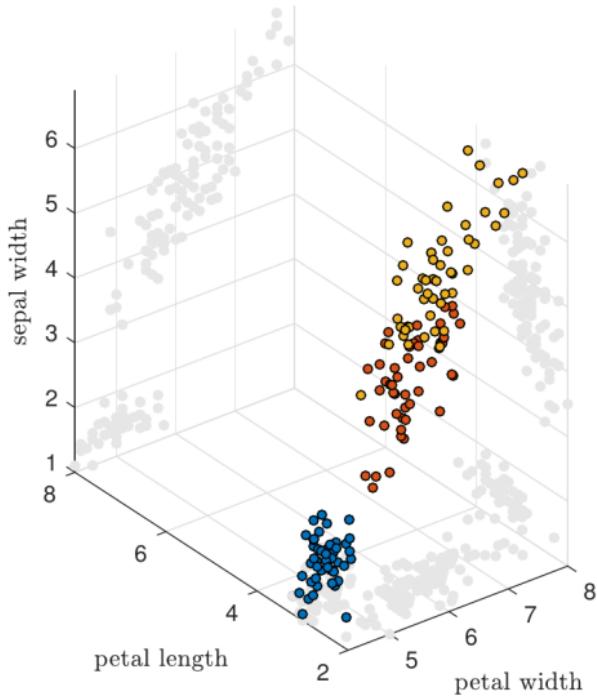
Scatter plots

- Shows **relation** between attributes
 - Assess dependence between attributes
 - Used with classes to assess separability



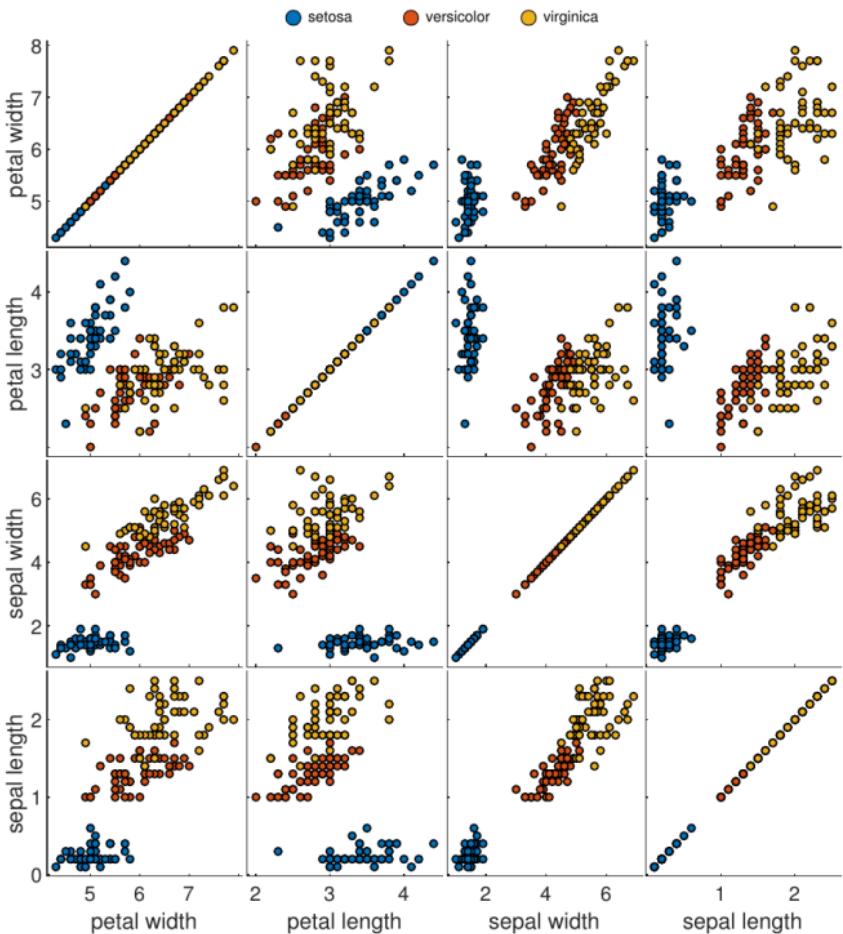
Scatter plots

- Shows **relation** between attributes
 - Assess dependence between attributes
 - Used with classes to assess separability
 - 3D plots are often confusing;
avoid if possible



Scatter plots

- Scatter plot matrix
 - All pairs of attributes



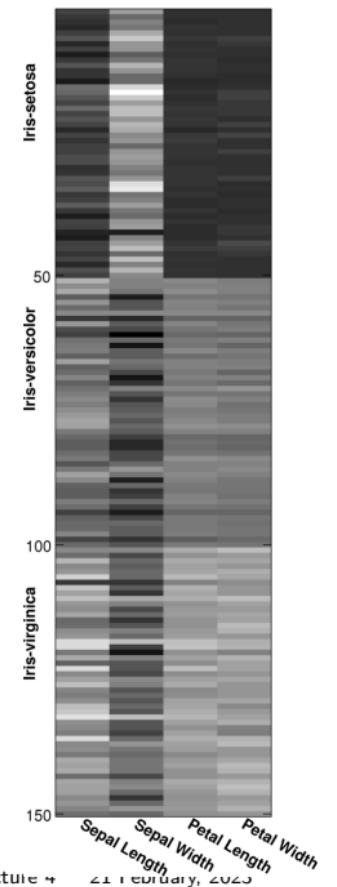
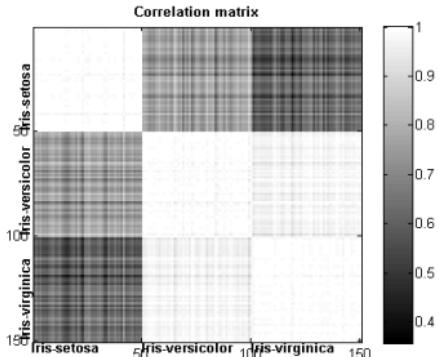
Matrix plots

- **Plot of raw data matrix**

- Useful when objects are sorted according to class
- Typically, attributes are normalized

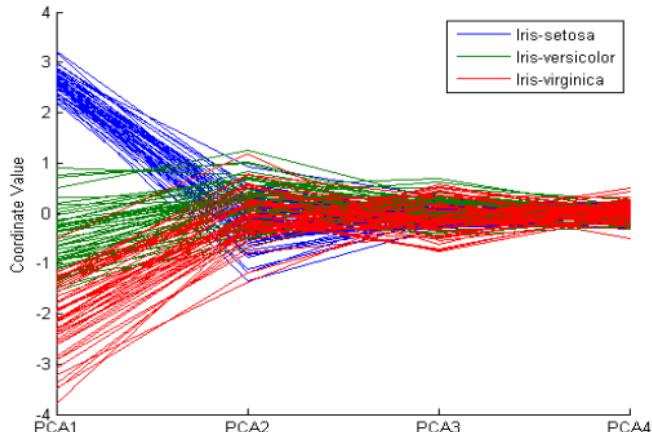
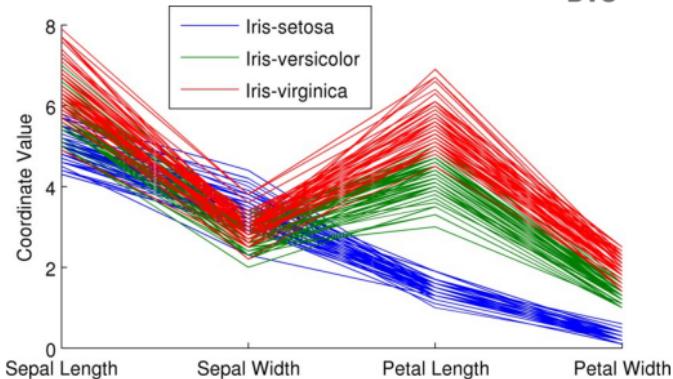
- **Plots of similarity matrices**

- Useful for visualizing the relation between objects



Parallel coordinates

- Plot high-dimensional data
- Instead of perpendicular axes
 - Use parallel axes
- Attribute values are plotted as a point
 - and the points are connected by a line
- Each object is represented as a line
- Lines representing a group of objects
 - Are similar in some sense
 - Ordering of attributes is important in seeing such groupings



ACCENT

- **Apprehension**

- Is it easy to see what is important in the graph?

- **Clarity**

- Are the most important elements visually most prominent?

- **Consistency**

- Have you used the same colors, shapes, etc. as in other graphs?

- **Efficiency**

- Does it convey its information in the most simple and efficient way?

- **Necessity**

- Are all elements of the graph necessary to represent data?

- **Truthfulness**

- Does the graph represent the data correctly?

Tufte's guidelines

- **Graphical excellence**

- Well-designed presentation of interesting data – a matter of
 - substance, statistics, and design
- Complex ideas communicated with
 - clarity, precision, and efficiency
- Gives the viewer
 - the greatest number of ideas
 - in the shortest time
 - with the least ink
 - in the smallest place.
- Nearly always multivariate
- Requires telling the truth about the data
- Maximise Data-ink ratio:

$$\text{Data-ink ratio} = \frac{\text{Data-ink}}{\text{Total ink used}}$$



Edward Tufte

https://commons.wikimedia.org/wiki/File:Edward_Tufte_-_cropped.jpg
Made available by Keegan Peterzell

Making good data visualizations is an art

For some interesting data visualizations see also

http://www.ted.com/talks/david_mccandless_the_beauty_of_data_visualization.html

<http://www.informationisbeautiful.net/>

<http://www.junkcharts.typepad.com/>

Resources

<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>)

<https://junkcharts.typepad.com> Excellent resource on creating good visualizations (https://junkcharts.typepad.com/junk_charts/)

<http://www2.imm.dtu.dk> Our demo of the multivariate normal distribution which illustrates the effect of the covariance matrix

(<http://www2.imm.dtu.dk/courses/02450/DemoNormal.html>)