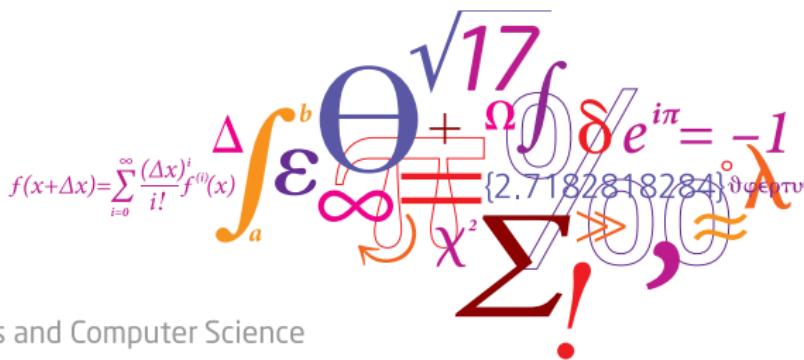


02450: Introduction to Machine Learning and Data Mining

Probability densities and data Visualization

Mikkel N. Schmidt

DTU Compute, Technical University of Denmark (DTU)



Lecture Schedule

1 Introduction

7 October: C1

Data: Feature extraction, and visualization

2 Data, feature extraction and PCA

7 October: C2, C3

3 Measures of similarity, summary statistics and probabilities

7 October: C4, C5

4 Probability densities and data Visualization

7 October: C6, C7

Supervised learning: Classification and regression

5 Decision trees and linear regression

8 October: C8, C9

6 Overfitting, cross-validation and Nearest Neighbor

8 October: C10, C12

7 Performance evaluation, Bayes, and Naive Bayes

9 October: C11, C13

Piazza online help: <https://piazza.com/dtu.dk/fall2019/october2019>

8 Artificial Neural Networks and Bias/Variance

9 October: C14, C15

9 AUC and ensemble methods

10 October: C16, C17

Unsupervised learning: Clustering and density estimation

10 K-means and hierarchical clustering

10 October: C18

11 Mixture models and density estimation

11 October: C19, C20

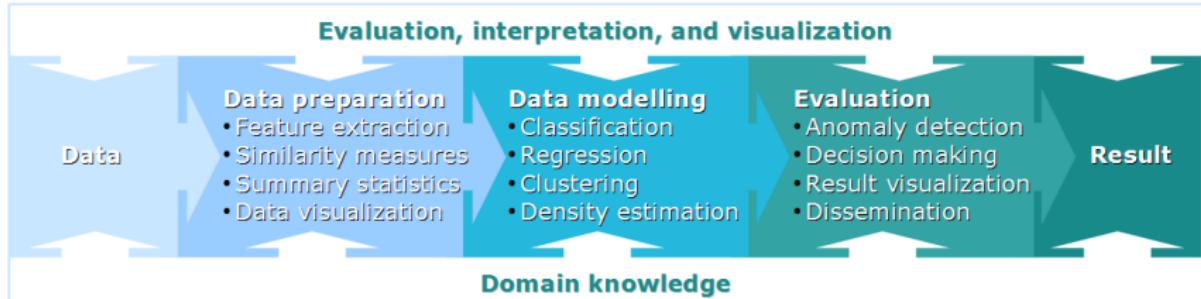
12 Association mining

11 October: C21

Recap

13 Recap

11 October: C1-C21

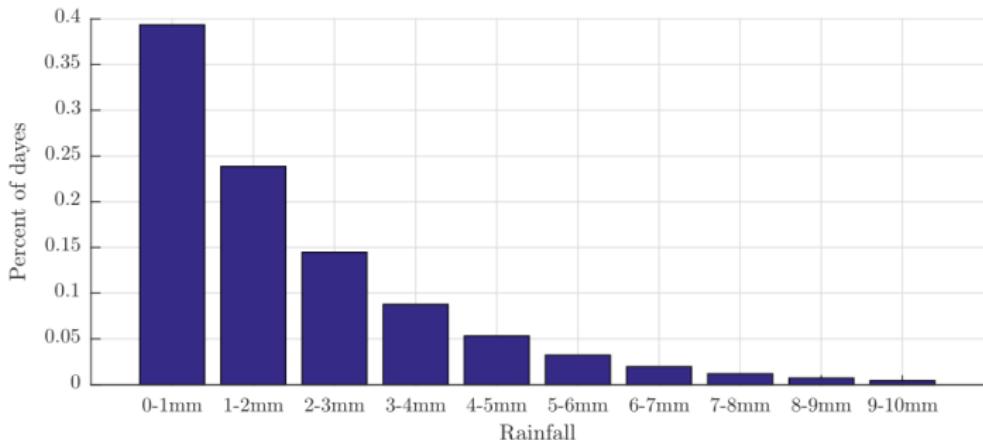


Learning Objectives

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

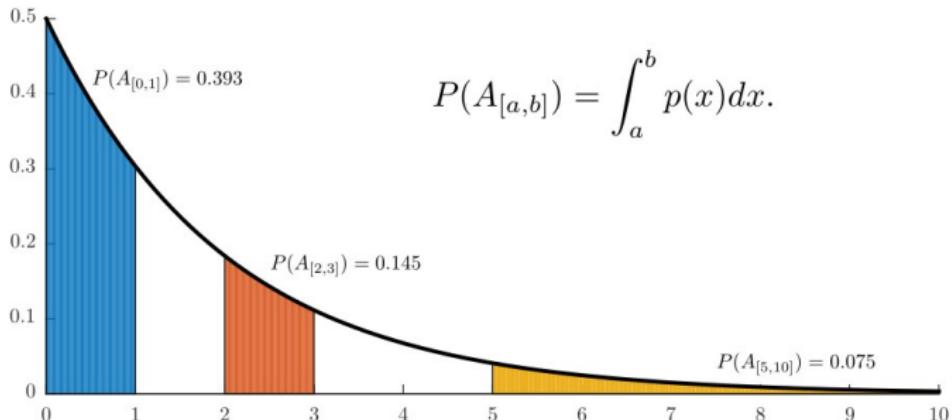
Probability vs. Density

- 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

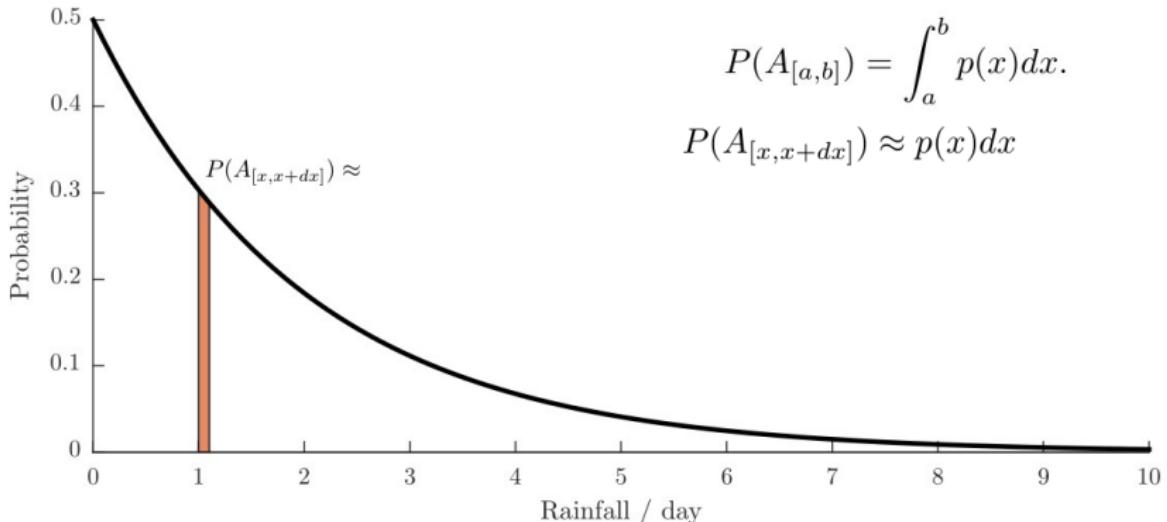
- These probabilities can be **represented** as integrals
- **Events** are **intervals**, the **probability** is the **integral**, the **curve** is the **density**



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

Probability vs. Density

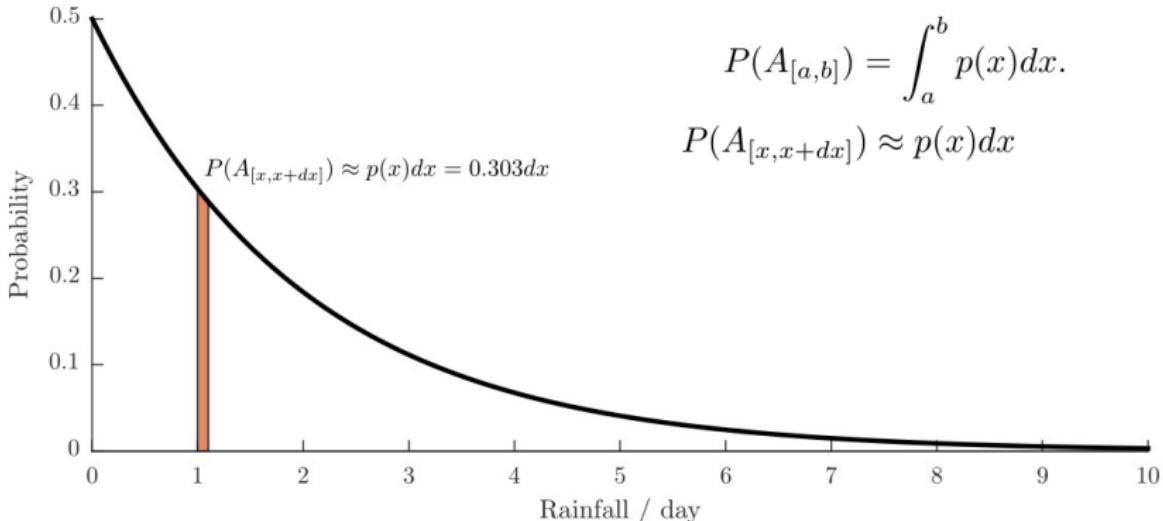
- These probabilities can be **represented** as integrals
- **Events** are **intervals**, the **probability** is the **integral**, the **curve** is the **density**
- What is the probability there will be between 1 and 1.1 mm of rain?



$A_{[a,b]}$: 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

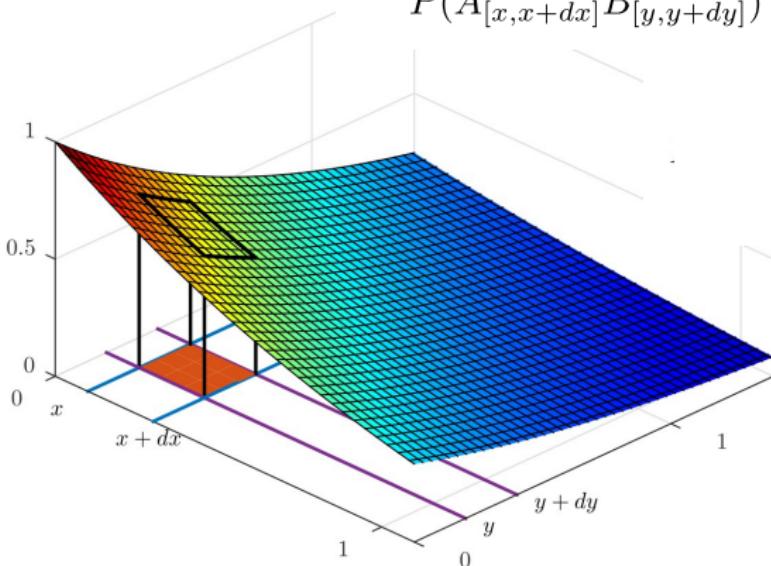
Probability vs. Density

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

$$P(A_{[x,x+dx]})$$

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

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



This implies:

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

Probability vs. Density

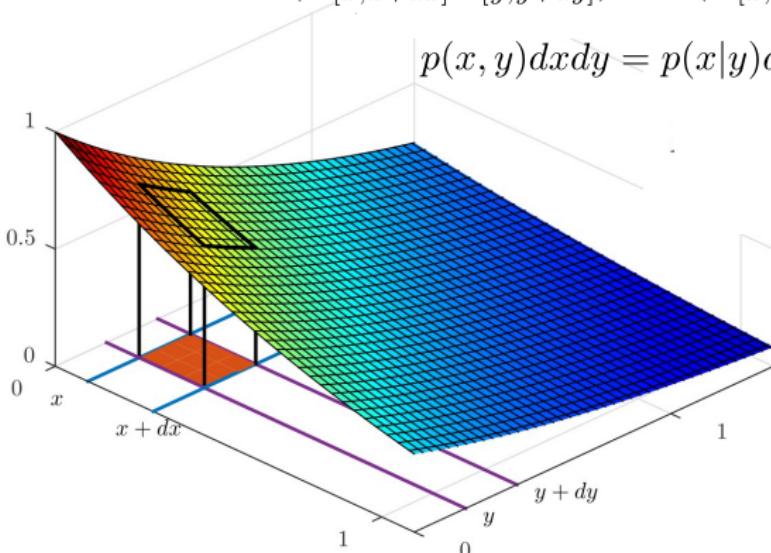
- For two variables x and y , the **probability** is an integral over an **area**

$$P(A_{[x,x+dx]})$$

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

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

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



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:

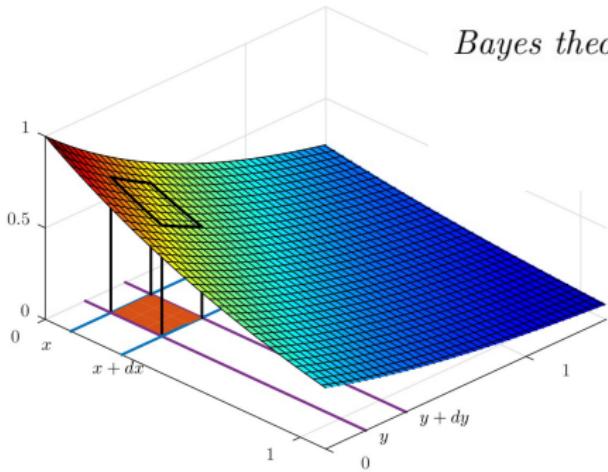
$$\int dx \ p(x|z) = 1$$

The product rule:

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

Bayes theorem:

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



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 our of simpler building blocks (densities).
In this course we will learn four:

Bernoulli density

The Catagorical density

The Beta density

The Multivariate normal density

The multivariate normal distribution

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

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

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

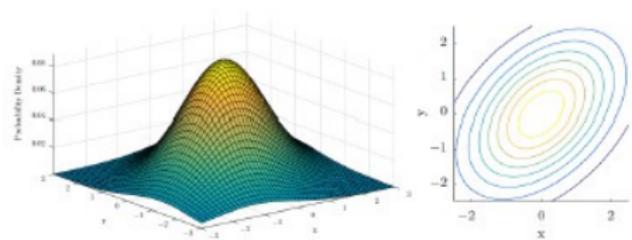
$\boldsymbol{\mu}$ is the mean vector and $\boldsymbol{\Sigma}$ is the covariance matrix:

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

- Example: 2-dimensional Normal distribution

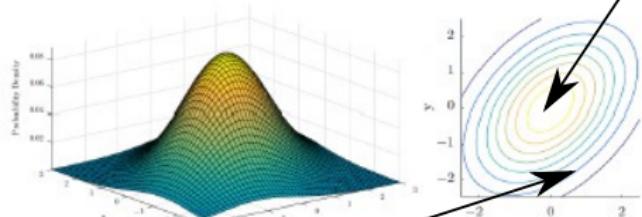
$$\boldsymbol{\mu} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \end{bmatrix}$$

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



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

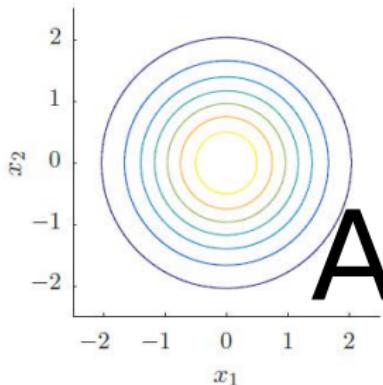
B. $\Sigma_B = \begin{bmatrix} 1 & -0.45 \\ 0.45 & 1 \end{bmatrix}$, $\Sigma_C = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$

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

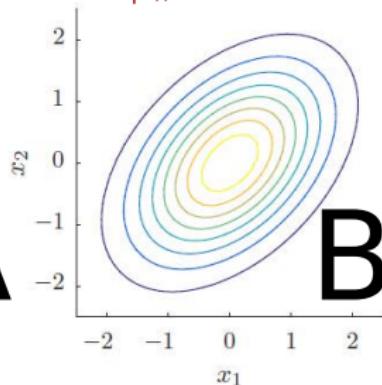
D. $\Sigma_A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$, $\Sigma_C = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$

E. Don't know.

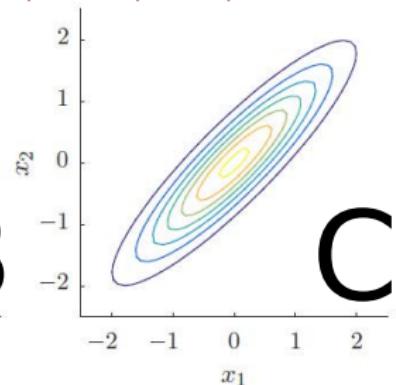
Check out the online demo <http://www2.imm.dtu.dk/courses/02450/DemoNormal.html>



A



B



C

The right answer is *D*. The covariance has to be positive (because x_1 and x_2 are positively correlated), and the variance is 1 in all cases. Furthermore, since A is axis-aligned, the covariance terms are zero. All

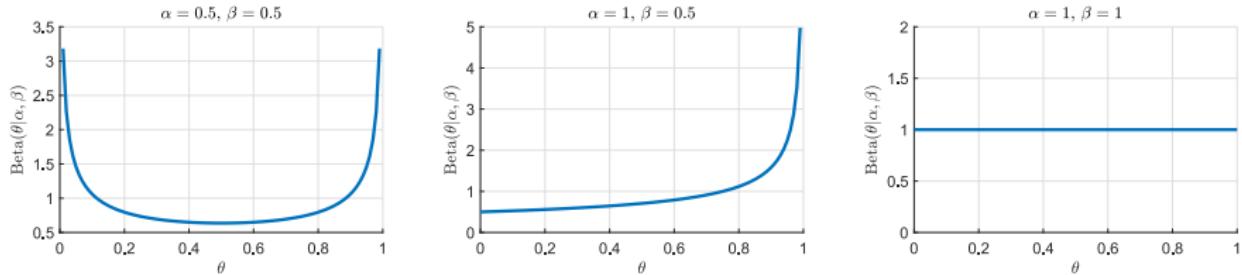
in all

$$\Sigma_A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \Sigma_B = \begin{bmatrix} 1 & 0.45 \\ 0.45 & 1 \end{bmatrix}, \quad \Sigma_C = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$$

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

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

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 last week: 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 $p(\theta) =$

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

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

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!**

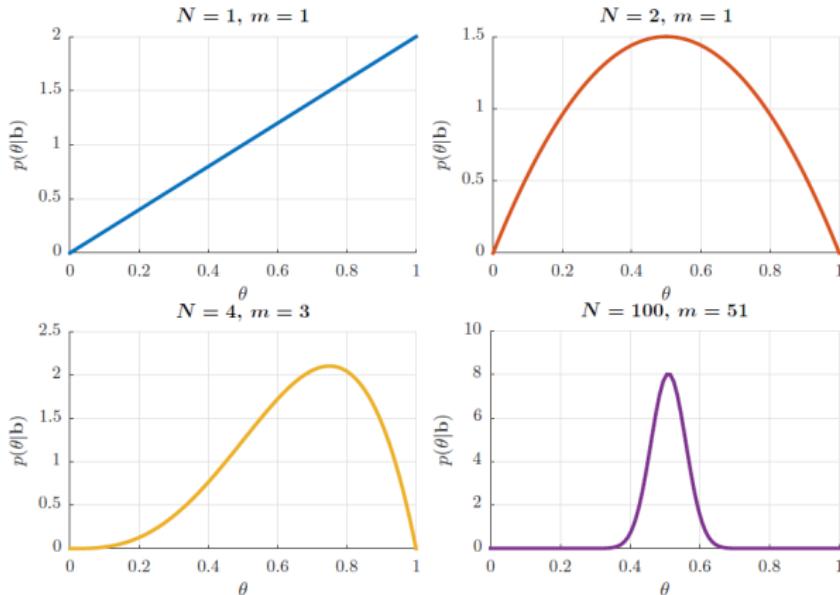
$$p(\theta|\mathbf{b}) = \frac{p(\mathbf{b}|\theta)p(\theta)}{p(\mathbf{b})} = \frac{p(\mathbf{b}|\theta)p(\theta)}{\int_0^1 p(\mathbf{b}|\theta')p(\theta')d\theta'}$$

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

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

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?

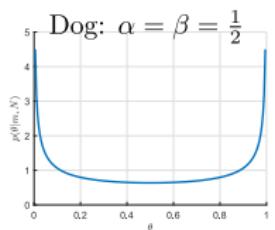
Dogs and coins



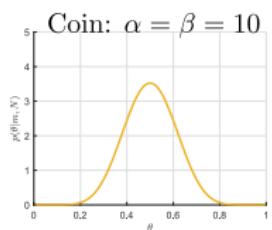
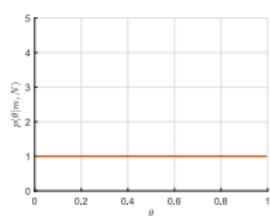
- Your friend just got a dog,
- The dog can either be in the doghouse or outside
- The first four times you come by the dog is in the doghouse
- What is the chance the dog is in the doghouse tomorrow?

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?



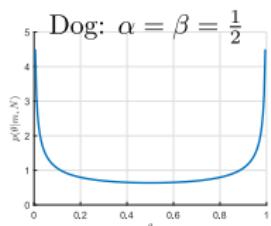
Dogs and coins



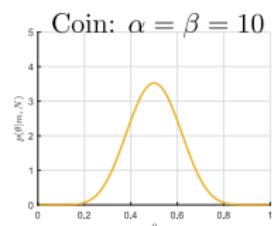
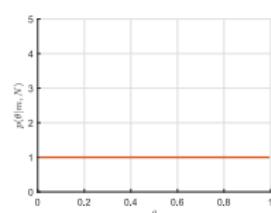
- Your friend just got a dog.
- The dog can either be in the doghouse or outside
- The first four times you come by the dog is in the doghouse
- What is the chance the dog is in the doghouse tomorrow?

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

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

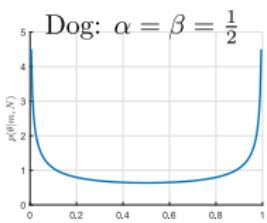
Dogs and coins



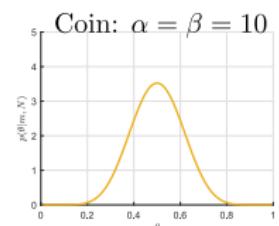
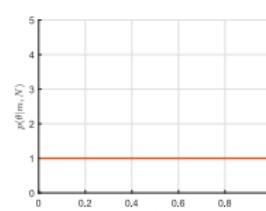
- Your friend just got a dog.
- The dog can either be in the doghouse or outside
- The first four times you come by the dog is in the doghouse
- What is the chance the dog is in the doghouse tomorrow?

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

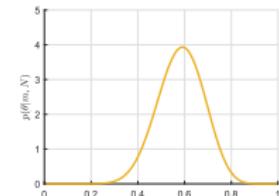
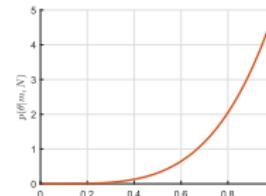
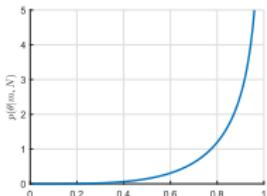
$$p(m=4, N=4|\theta) = \theta^m(1-\theta)^{N-m} = \theta^4$$

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

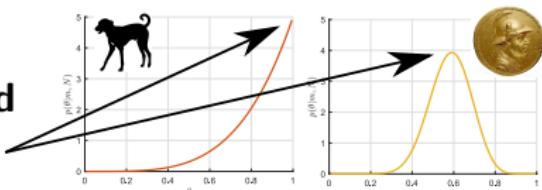
In most practical situations, we should assume as little as possible and choose $\alpha = \beta = \frac{1}{2}$

Posterior

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



Learning principle: Maximum likelihood



- 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 $f(x) > 0$:

$$\theta^* = \arg \min_{\theta} \left[-\log \frac{p(m, N | \theta) p(\theta | \alpha, \beta)}{p(m, N)} \right] \quad \begin{array}{l} \text{(likelihood)} \\ p(m, N | \theta) = \theta^m (1 - \theta)^{N-m} \end{array}$$

$$\begin{array}{l} \text{(prior)} \\ p(\theta | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1} \end{array}$$

- All in all:

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

A learning principle: Maximum likelihood

- 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 $f(x) > 0$:

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

- All in all:

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

Maximum likelihood 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 θ
- Assume**

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

Observations are not informative about each other when we know parameters

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

- Then

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

- The following are equivalent:

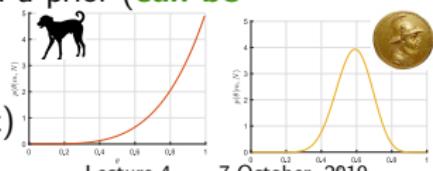
Maximum likelihood principle

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

$$\text{Minimize: } \mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w})$$

$$E(\mathbf{w}) =$$

- 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)



Maximum likelihood learning

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

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

- Then

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

- The following are equivalent:

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

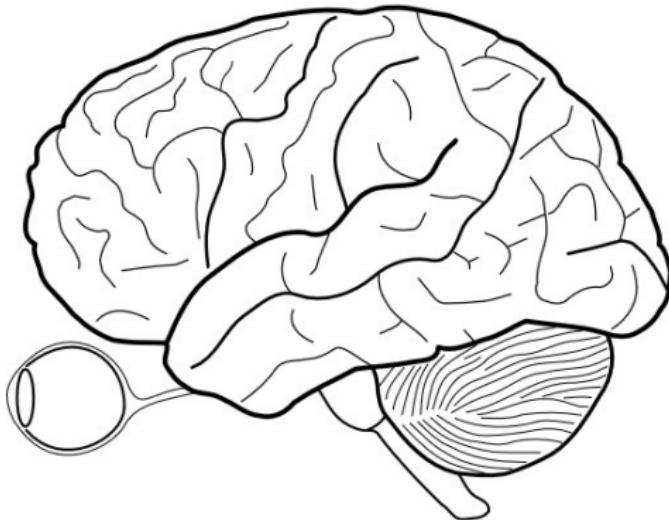
$$\text{Minimize: } \mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w})$$

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

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

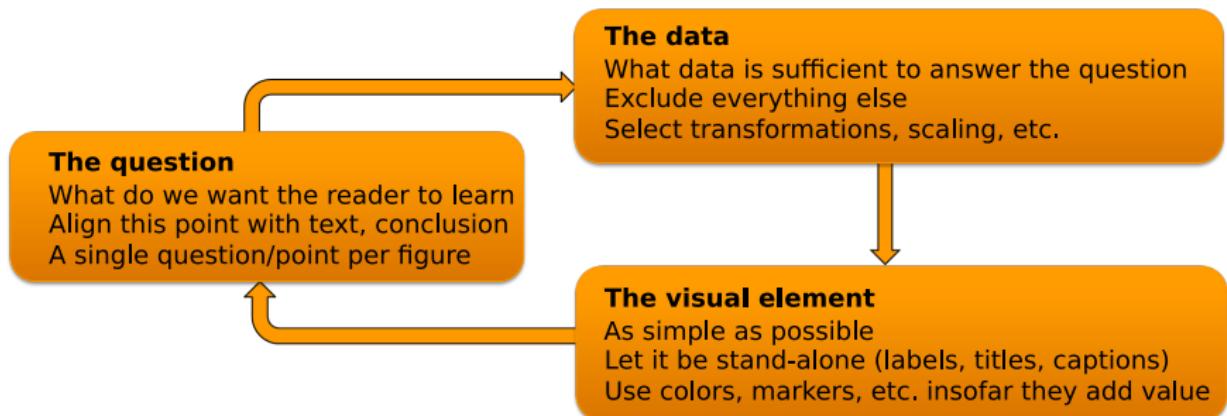
Visualization

- 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*)

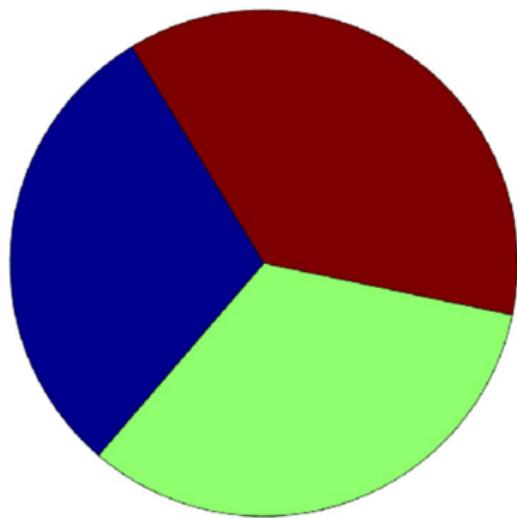


Important choices for visualizations

- **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

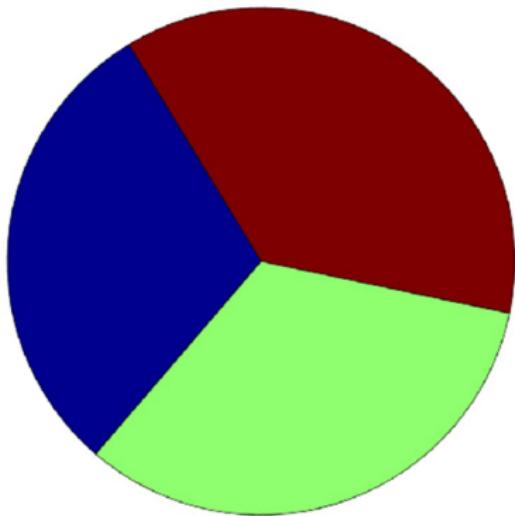
- **Area represents proportion**
 - Which is smallest, middle, and largest?
 - What are the proportions approximately?



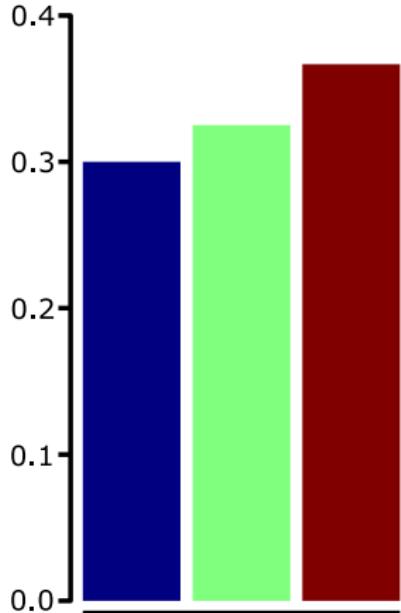
Representation

- **Area represents proportion**

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



- **Height represents proportion**



Arrangement

- **Placement of visual elements**

- Can make a great difference in how easy it is to interpret data

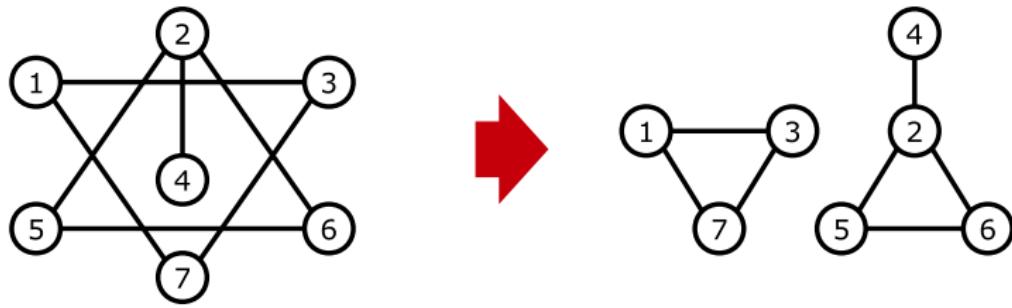
- **Example**

	1	2	3	4	5	6
1	0	1	0	1	1	0
2	1	0	1	0	0	1
3	0	1	0	1	1	0
4	1	0	1	0	0	1
5	0	1	0	1	1	0
6	1	0	1	0	0	1
7	0	1	0	1	1	0
8	1	0	1	0	0	1
9	0	1	0	1	1	0



	6	1	3	2	5	4
4	1	1	1	0	0	0
2	1	1	1	0	0	0
6	1	1	1	0	0	0
8	1	1	1	0	0	0
5	0	0	0	1	1	1
3	0	0	0	1	1	1
9	0	0	0	1	1	1
1	0	0	0	1	1	1
7	0	0	0	1	1	1

Arrangement



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**

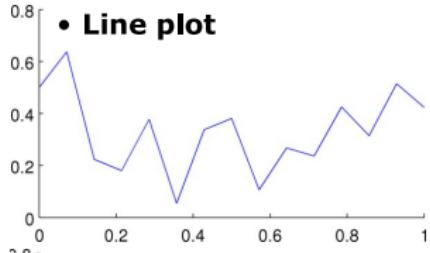
- 2D histogram
- Heat maps and contour plots
- Scatter plots

- **Visualization of high-dimensional objects**

- Matrix plots
- Parallel coordinates
- Star plots

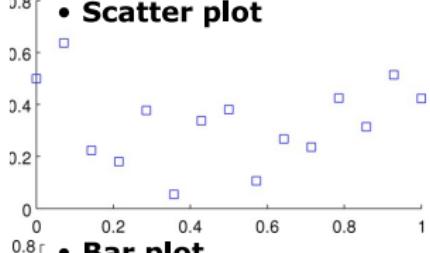
Basic plots

- **Line plot**



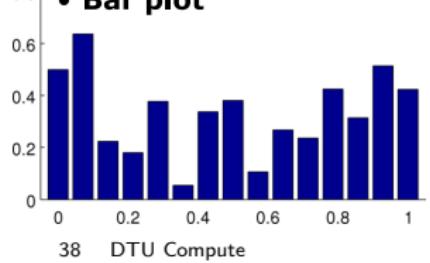
```
plot(x,y);
```

- **Scatter plot**



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

- **Bar plot**



```
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

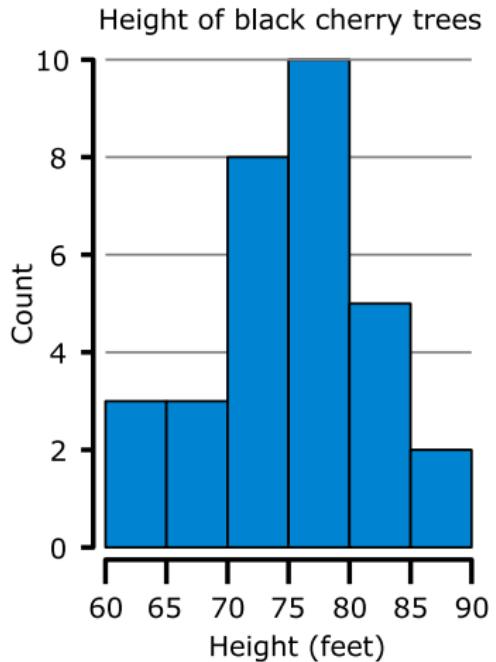
$$X^{\text{Observation} \times \text{Attribute}}$$

Distribution of a single attribute

Histograms

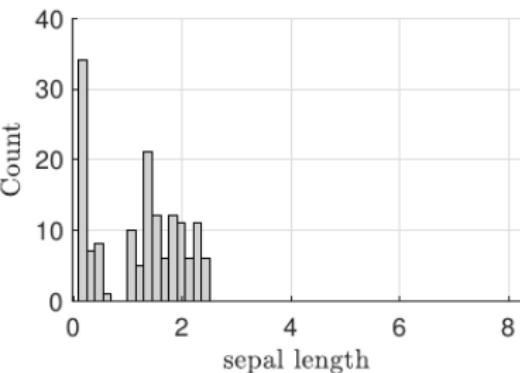
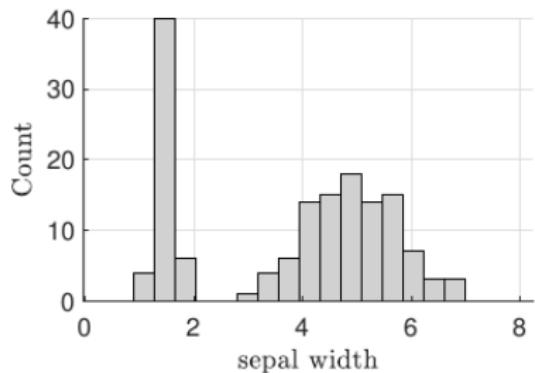
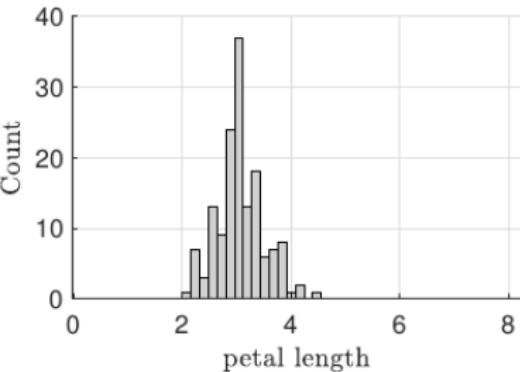
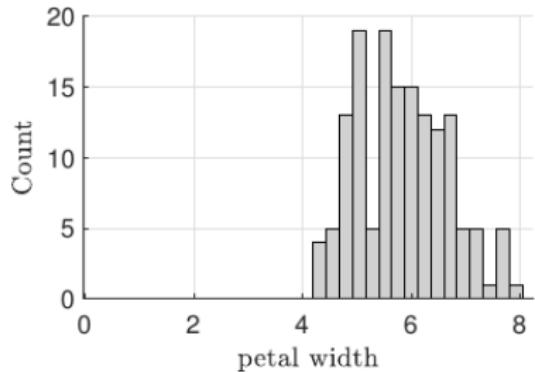
- Shows distribution of a single variable

- 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



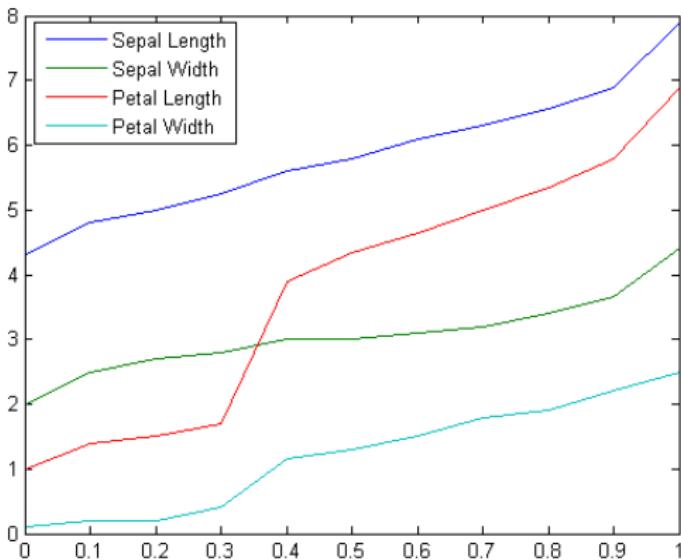
$$H = \{60, 64, 64, 66, 67, 69, 71, 72, 72, 72, 72, 73, 74, 74, 75, 75, 76, 76, 77, 77, 78, 78, 79, 80, 80, 81, 82, 84, 85, 85, 89\}$$

Histograms of the Iris data attributes



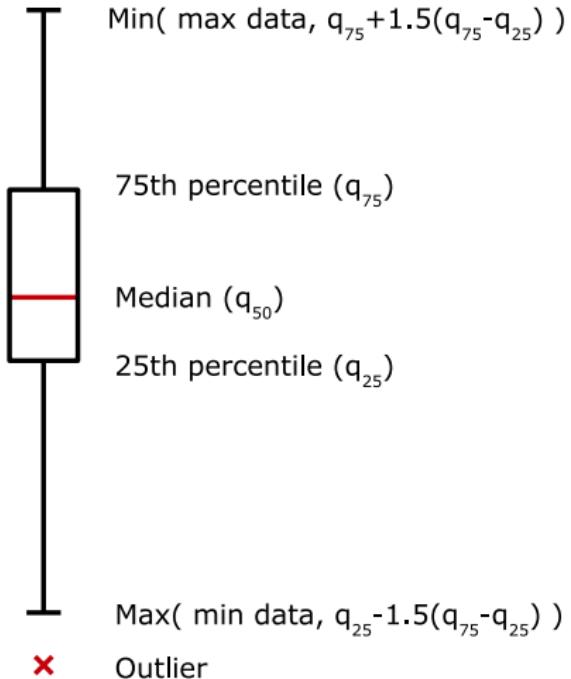
Percentile plots

Percentiles: Given an ordinal or continuous attribute x and a number p between 0 and 100, the p th percentile is a value x_p of x such that p percent of the observed values of x are less than x_p .

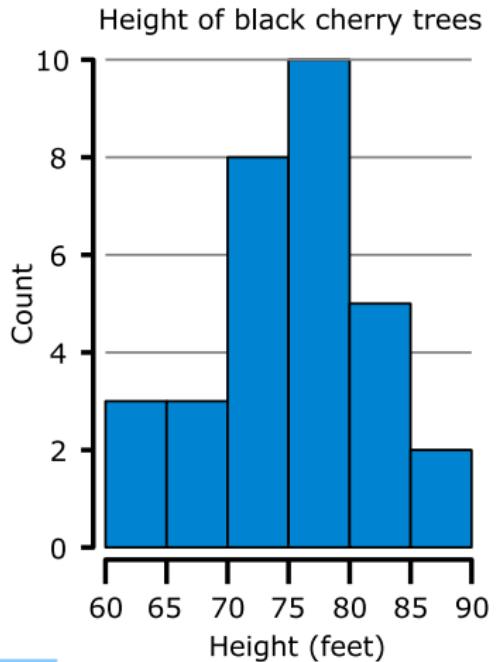


```
prctile = 0:0.1:1;
Y = quantile(X,prctile);
plot(prctile,Y);
legend(attributeNames);
```

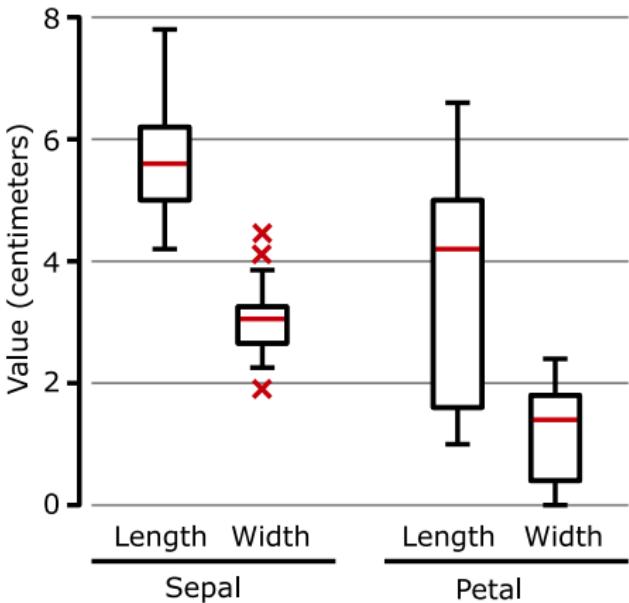
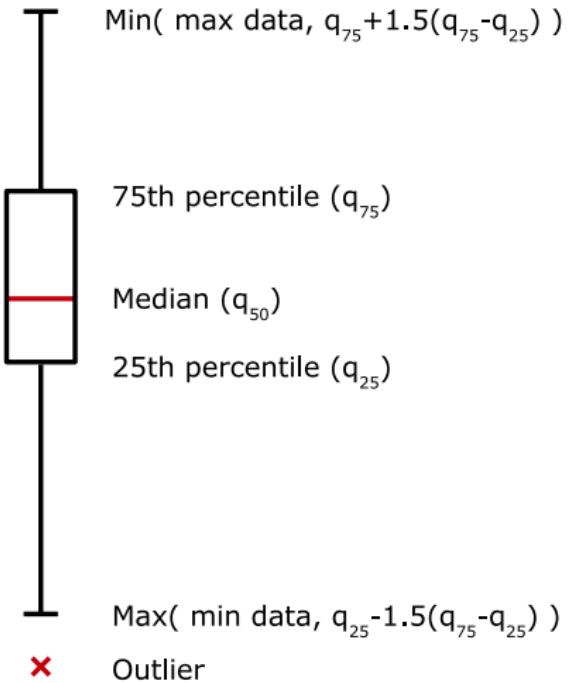
Box plots



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

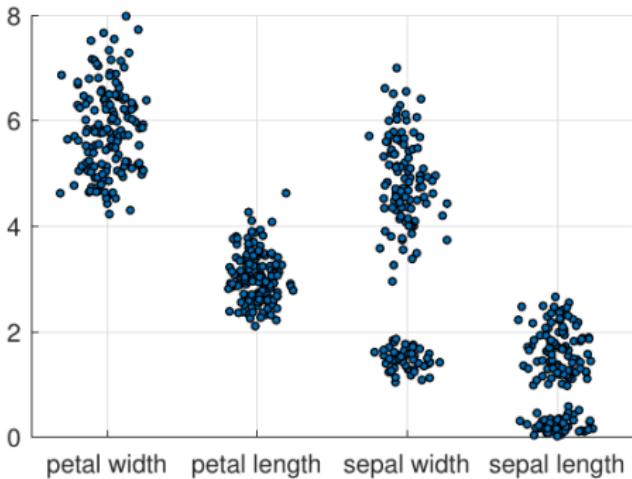
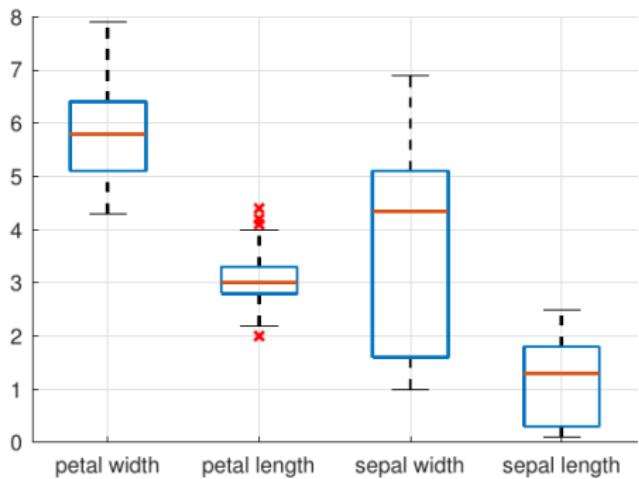


Box plots



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

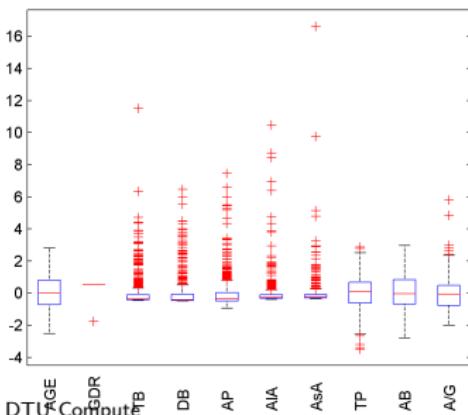
Box plots



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.

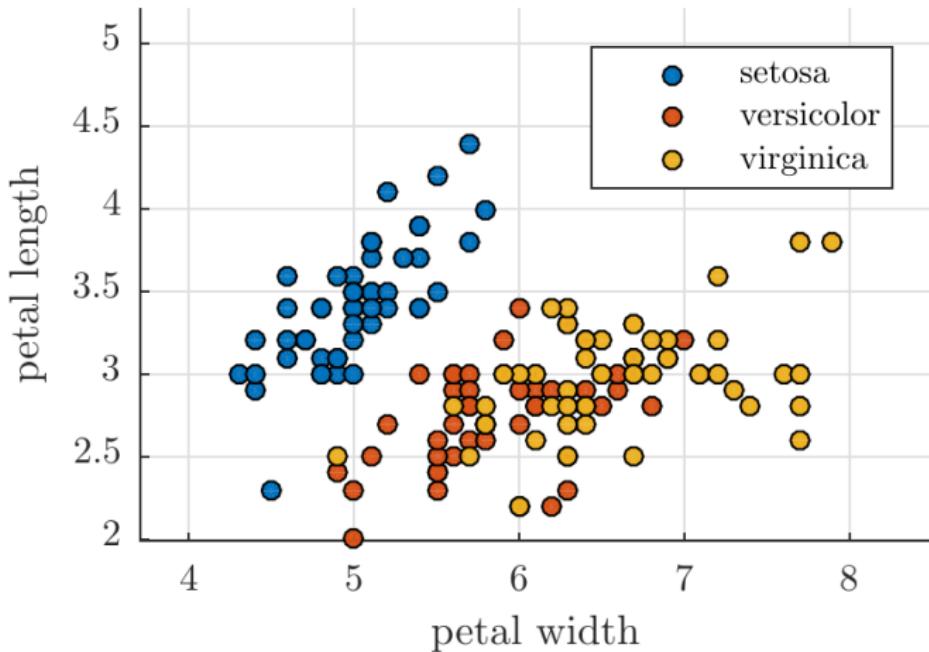
The 25th and 50th percentile but not the 50th and 75th percentiles of the attribute DB coincides. A1A and AsA will not necessarily be highly correlated even though their distributions may have a similar shape (hence, this is correct). For attributes to be correlated it is important they take on high or low values systematically, however, this can not be inspected in

a boxplot. TB is not likely to be normal distribution as this attribute does not have a symmetric but highly right skewed distribution. The attribute GDR does not have a clear outlier, in fact the outlier corresponds to the females in the dataset and all we can deduce from the plot is that more than 75 % of the observations are males.

Relation between attributes

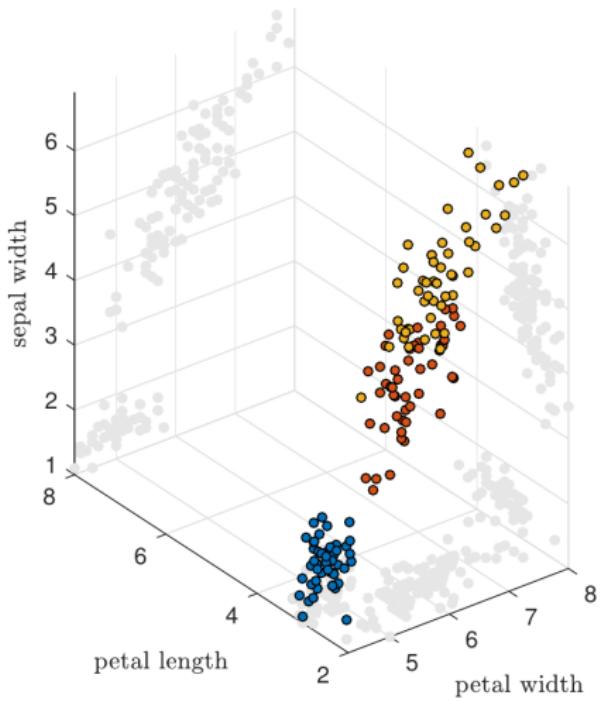
Scatter plots

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



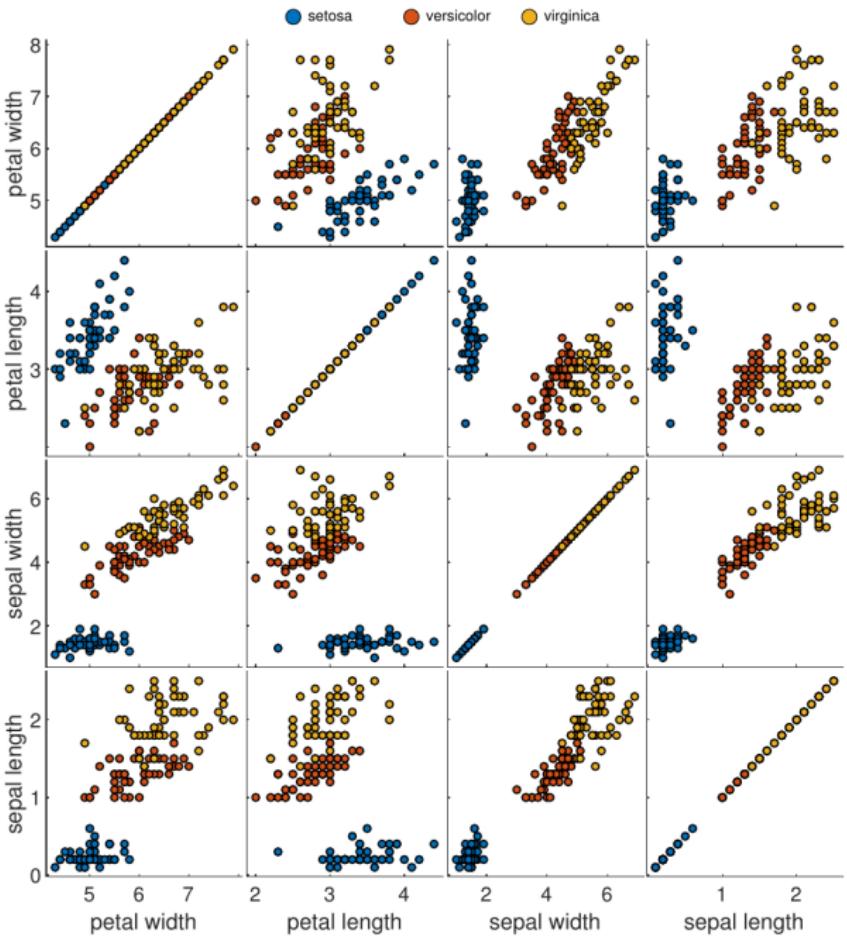
Scatter plots

- 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



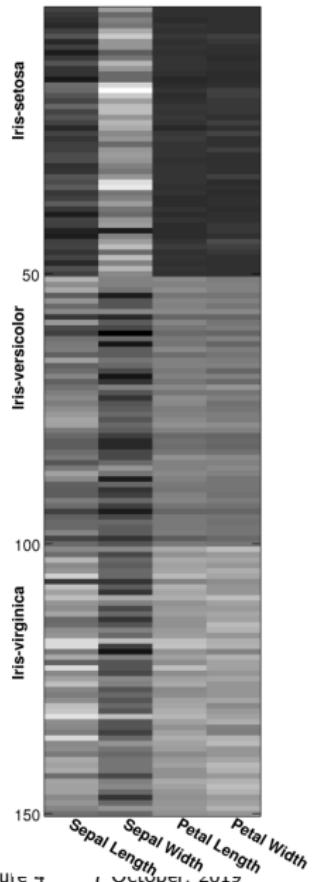
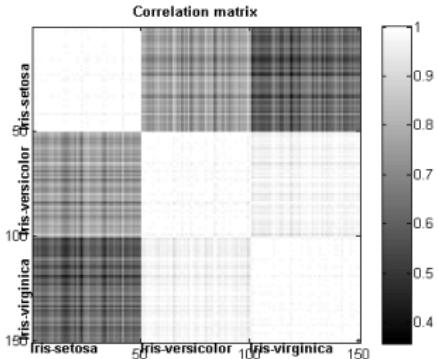
Matrix plots

- Plot of raw data matrix

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

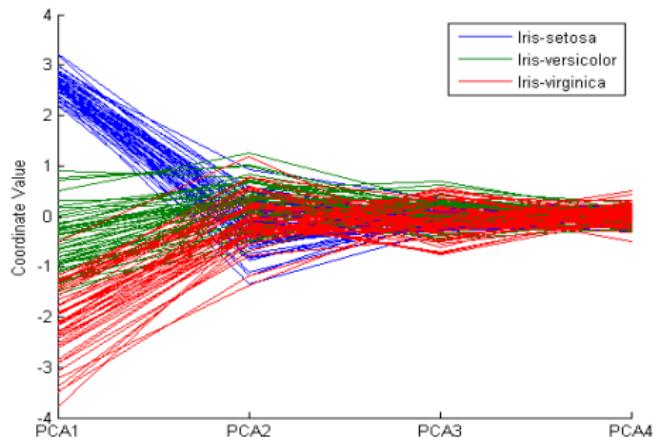
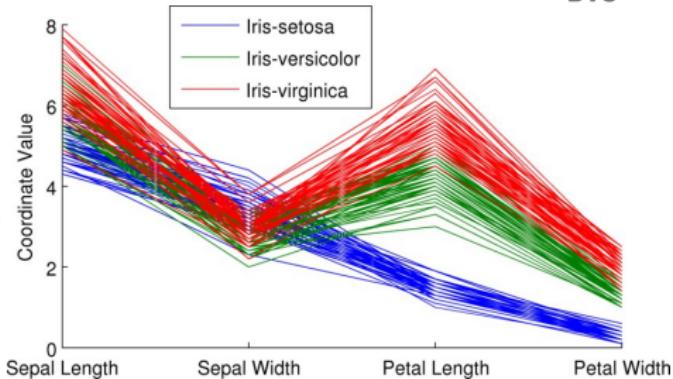
- **Plots of similarity matrices**

- 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/>



Imagine you have a dataset where some of the attributes are numeric given in the matrix X but you also have a categorical attribute given by TXT (see below). You would like to carry out a PCA on the data taking both the numeric and categoric attributes into account. How would you proceed?

Age Height Weight
(Standardized)

-0.2248	-0.4762	-0.2097
-0.5890	0.8620	0.6252
-0.2938	-1.3617	0.1832
-0.8479	0.4550	-1.0298
-1.1201	-0.8487	0.9492
2.5260	-0.3349	0.3071
X= 1.6555	0.5528	0.1352
0.3075	1.0391	0.5152
-1.2571	-1.1176	0.2614
-0.8655	1.2607	-0.9415
-0.1765	0.6601	-0.1623
0.7914	-0.0679	-0.1461
-1.3320	-0.1952	-0.5320
-2.3299	-0.2176	1.6821
-1.4491	-0.3031	-0.8757
0.3335	0.0230	-0.4838
0.3914	0.0513	-0.7120
0.4517	0.8261	-1.1742
-0.1303	1.5270	-0.1922
0.1837	0.4669	-0.2741

Nationality

'Sweden'
'Sweden'
'Sweden'
'Sweden'
'Norway'
'Norway'
'Norway'
TXT= 'Norway'
'Norway'
'Norway'
'Sweden'
'Norway'
'Denmark'
'Denmark'
'Sweden'
'Sweden'

Denmark Norway Sweden



X_tmp=

0	0	1
0	0	1
0	0	1
0	0	1
0	1	0
0	1	0
0	1	0
0	1	0
1	0	0
1	0	0
1	0	0
0	0	1
0	0	1
0	1	0
1	0	0
1	0	0
0	0	1
0	0	1
1	0	0
0	1	0
0	1	0
0	1	0
0	0	1
1	0	0

One-out-of-K coding

```
[X_tmp, attributeNames_tmp]=categoric2numeric(TXT);  
X=[X X_tmp];  
attributeNames=[attributeNames; attributeNames_tmp];
```

X=

-0.2248	-0.4762	-0.2097	0	0	1.0000
-0.5890	0.8620	0.6252	0	0	1.0000
-0.2938	-1.3617	0.1832	0	0	1.0000
-0.8479	0.4550	-1.0298	0	0	1.0000
-1.1201	-0.8487	0.9492	0	1.0000	0
2.5260	-0.3349	0.3071	0	1.0000	0
1.6555	0.5528	0.1352	0	1.0000	0
0.3075	1.0391	0.5152	0	1.0000	0
-1.2571	-1.1176	0.2614	0	1.0000	0
-0.8655	1.2607	-0.9415	0	0	1.0000
-0.1765	0.6601	-0.1623	0	1.0000	0
0.7914	-0.0679	-0.1461	1.0000	0	0
-1.3320	-0.1952	-0.5320	1.0000	0	0
-2.3299	-0.2176	1.6821	0	0	1.0000
-1.4491	-0.3031	-0.8757	0	0	1.0000
0.3335	0.0230	-0.4838	0	0	1.0000
0.3914	0.0513	-0.7120	0	0	1.0000
0.4517	0.8261	-1.1742	0	0	1.0000
-0.1303	1.5270	-0.1922	0	0	1.0000
0.1837	0.4669	-0.2741	1.0000	0	0

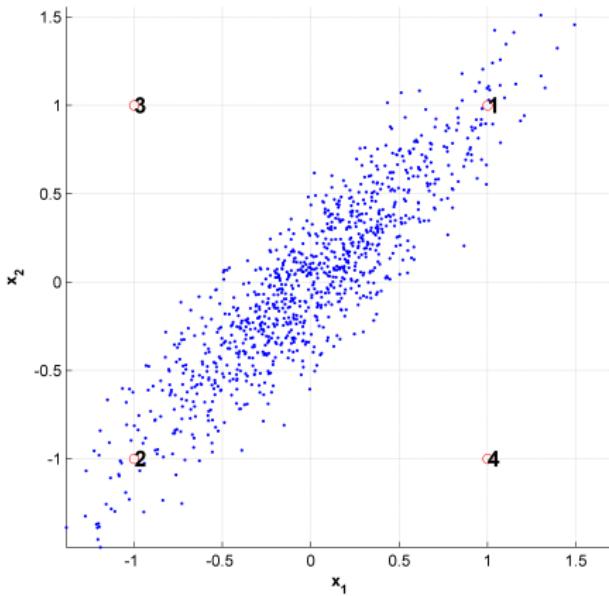
The Mahalanobis distance

How far are x_1 and x_2 apart?

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

How far are x_3 and x_4 apart?

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



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

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

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>)