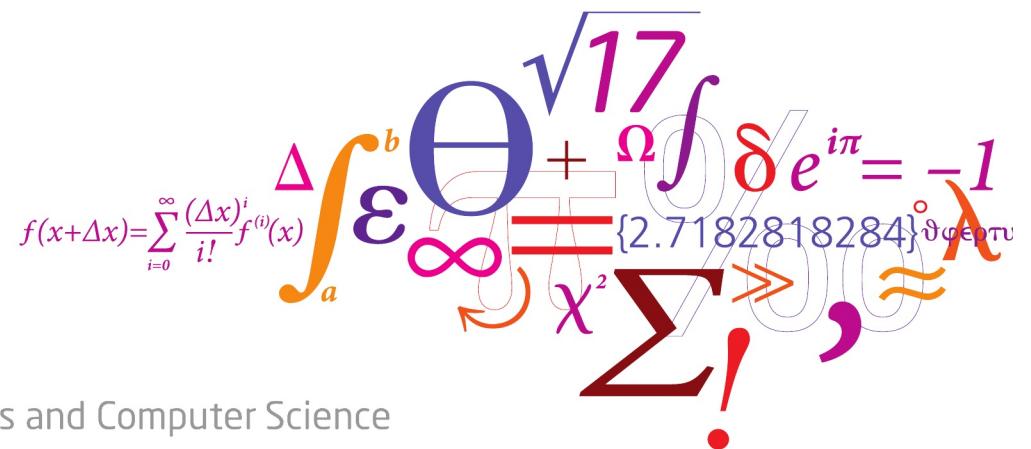


02450: Introduction to Machine Learning and Data Mining

Probability densities and data Visualization

Tue Herlau

DTU Compute, Technical University of Denmark (DTU)

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


DTU Compute

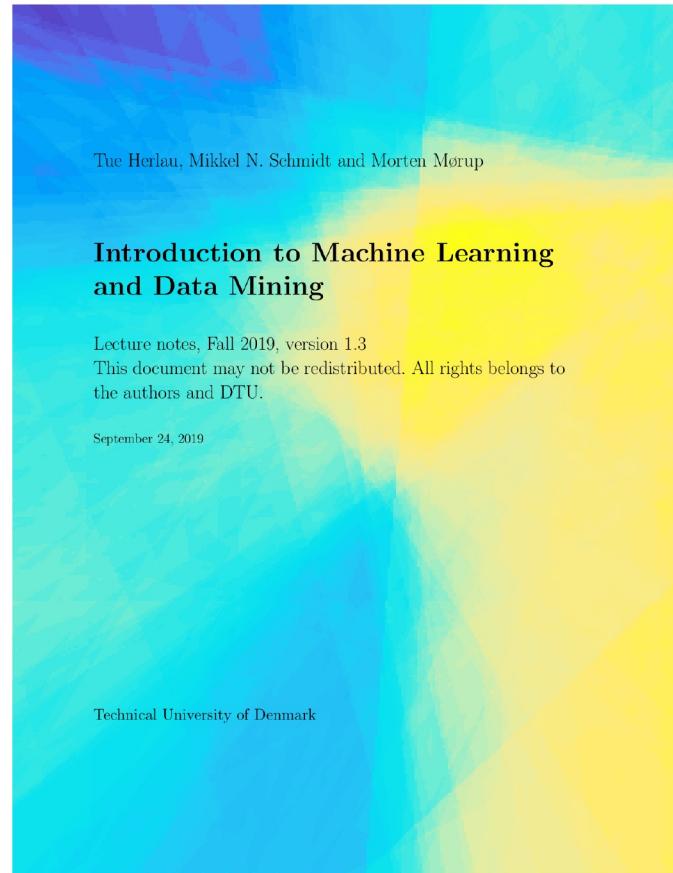
Department of Applied Mathematics and Computer Science

Today

Feedback Groups of the day:

Fabian Felix Bertoldo, Jørgen Taule, Raül Pérez i Gonzalo, Michelle Anker Pihl, Lene Christiansen, Mie Guldhammer Andersen, Nikolaj Geertinger, William Peter Settnes, Wictor Lang Jensen, Hadeel Moustafa, Abdulstar Shihada Kousa, Nojus Mickus, Andreas Goltermann, Jakob Malte Skou Lindstad, Christian Rømer Thulstrup, Niklaes Dino Robbin Jacobsen, Sebastian Stokkebro Sørensen, Kilian Waldemar Conde Frieboes, Thomas Rudolph Sparre Conrad, Asger Conradsen, Lukas Fredrik Cronje, Oskar Steentoft Dahlberg, Matteo D'Andrea, Benjamin Danielsen, Nichlas Davidsen, Seiran Davoudi, Sander de Boer, Matilde Maria de Place, Agathe De Vulpian, Sebastian Nyman Deleuran, Blanca Robledo Diaz, Rebecca Dürmüller, Emma Cathrine Rud Egelund, Pil Skov Eghoff, Lawrence Kofi Asane Egyir, Yosef Khodr El-Fil, Niels Kjær Ersbøll, Mahsa Eskandarzadeh, Kimmo Lilholt Bang Feldbæk, Þórður Páll Fjalarsson, Fredrik Hole Flaa, Maria Mandrup Fogh, Anton Palm Folkmann, Martin Trangbæk Forsingdal, Edvard Foss, Aslak Rønnow Franzen, Andreas Mørk Frendorf, Ditlev Reimer Frickmann, Tina Funck, Annika Lund Gade, Ruairí James Gallagher, Emilie Dybdahl Gandrup, Imogen Garner

Reading material: Chapter 6, Chapter 7



Lecture Schedule

① Introduction

3 September: C1

Data: Feature extraction, and visualization

② Data, feature extraction and PCA

10 September: C2, C3

③ Measures of similarity, summary statistics and probabilities

17 September: C4, C5

④ Probability densities and data Visualization

24 September: C6, C7

Supervised learning: Classification and regression

⑤ Decision trees and linear regression

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

⑥ Overfitting, cross-validation and Nearest Neighbor

8 October: C10, C12

⑦ Performance evaluation, Bayes, and Naive Bayes

22 October: C11, C13

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

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

⑧ Artificial Neural Networks and Bias/Variance

29 October: C14, C15

⑨ AUC and ensemble methods

5 November: C16, C17

Unsupervised learning: Clustering and density estimation

⑩ K-means and hierarchical clustering

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

⑪ Mixture models and density estimation

19 November: C19, C20

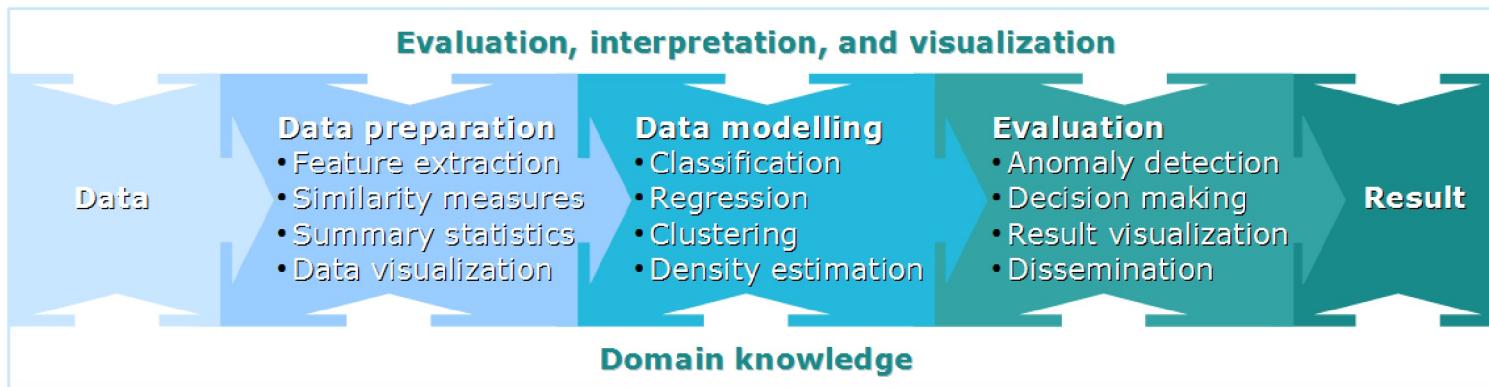
⑫ Association mining

26 November: C21

Recap

⑬ Recap and discussion of the exam

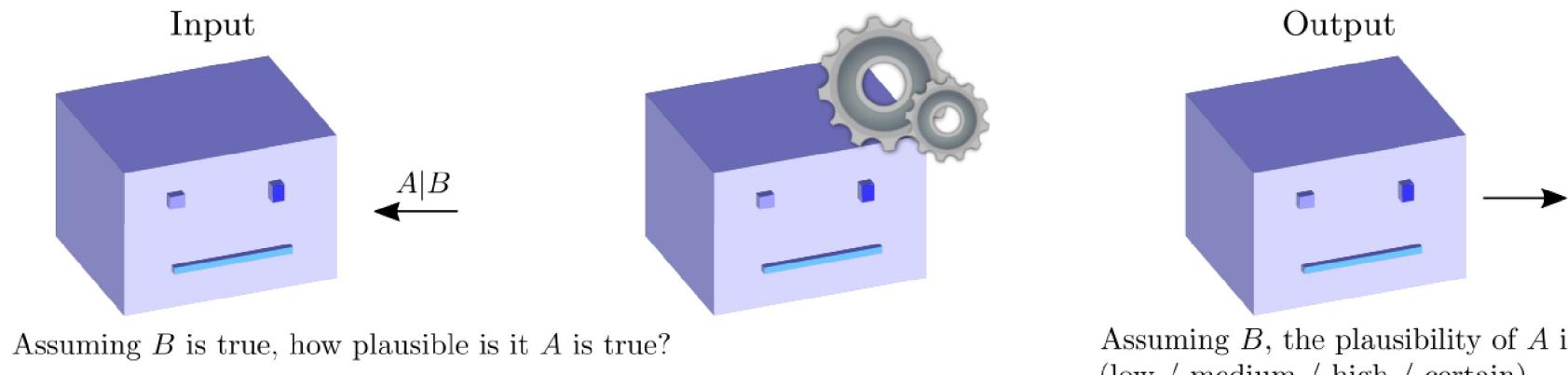
3 December: C1-C21 (Project 3 due before 13:00)



Learning Objectives

- 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

Probabilities recap



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

$$P(A|B) = x \in [0, 1]$$

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

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

Probabilities recap

- Probabilities $P(A|B)$ obey

Sum rule

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

Product rule

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

$$A_i A_j = 0$$

$$A_1 + A_2 + \dots + A_n = 1$$

$$P(A) \leq P(A|B)$$

- If we consider **mutually exclusive** and **exhaustive** events A_1, \dots, A_n :

Bayes' theorem

$$P(A_i|BC) = \frac{P(B|A_i C)P(A_i|C)}{\sum_j P(B|A_j C)P(A_j|C)}$$

$$P(A_i|B) \approx \frac{P(B|A_i)P(A_i)}{\sum_j P(B|A_j)P(A_j)}$$

Stochastic variables

Suppose X, Y , etc. are quantities we measure. $X = x$ is a binary event.
For instance

- $N = \{1, 2, \dots\}$: Number of children
- $T = 1, \dots, 24$: Hour of day
- $C = 1, \dots, K$ favorite cereal

Bayes theorem is then:

$$P(Y = y|X = x, Z = z) = \frac{P(X = x|Y = y, Z = z)P(Y = y|Z = z)}{\sum_{y'} P(X = x|Y = y', Z = z)P(Y = y'|Z = z)}$$

Using $\underbrace{p(x|y, z)}_{\text{ }} = P(X = x|Y = y, Z = z)$ and ignoring z this becomes:

$$p(y|x, \cancel{z}) = \frac{p(x|y)p(y)}{\sum_{y'} p(x|y')p(y')}$$

Quiz 1, Probabilities (Spring 2014)

Consider the Wholesale 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).

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?

D

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

$$\begin{aligned}
 R &= L, O_p, T \\
 P(R=L | D) &= \frac{P(D | R=L) P(R=L)}{\sum_{r=1}^3 P(D | R=r) P(R=r)} = \frac{(44.1\%) \times (17.5\%)}{---} \\
 &= \frac{0.441 \times .175}{0.441 \times .175 + 0.489 \times .107 + 0.516 \times 0.718} = \underline{0.1544}
 \end{aligned}$$

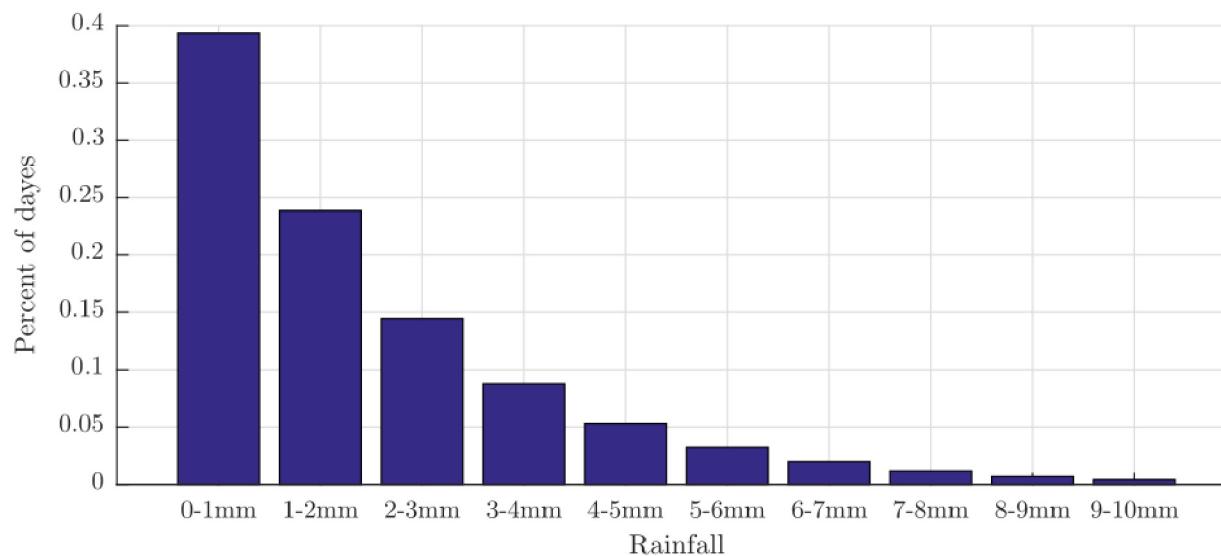
We will let $LISBON$, $OPORTO$ and $OTHER$ denote coming from Lisbon, Oporto and the other region respectively. $DELI_H$ will denote above the median value of delicatessen consumption. We now

find using Bayes theorem:

$$\begin{aligned}
 P(LISBON|DELI_H) &= \frac{P(DELI_H|LISBON)P(LISBON)}{P(DELI_H)} \\
 &= \frac{P(DELI_H|LISBON)P(LISBON)}{P(DELI_H|LISBON)P(LISBON) \\
 &\quad + P(DELI_H|OPORTO)P(OPORTO) \\
 &\quad + P(DELI_H|OTHER)P(OTHER)} \\
 &= \frac{44.1 \cdot 0.175}{0.441 \cdot 0.175 + 0.489 \cdot 0.107 + 0.516 \cdot 0.718} \\
 &= 0.1544 \approx 15.4\%
 \end{aligned}$$

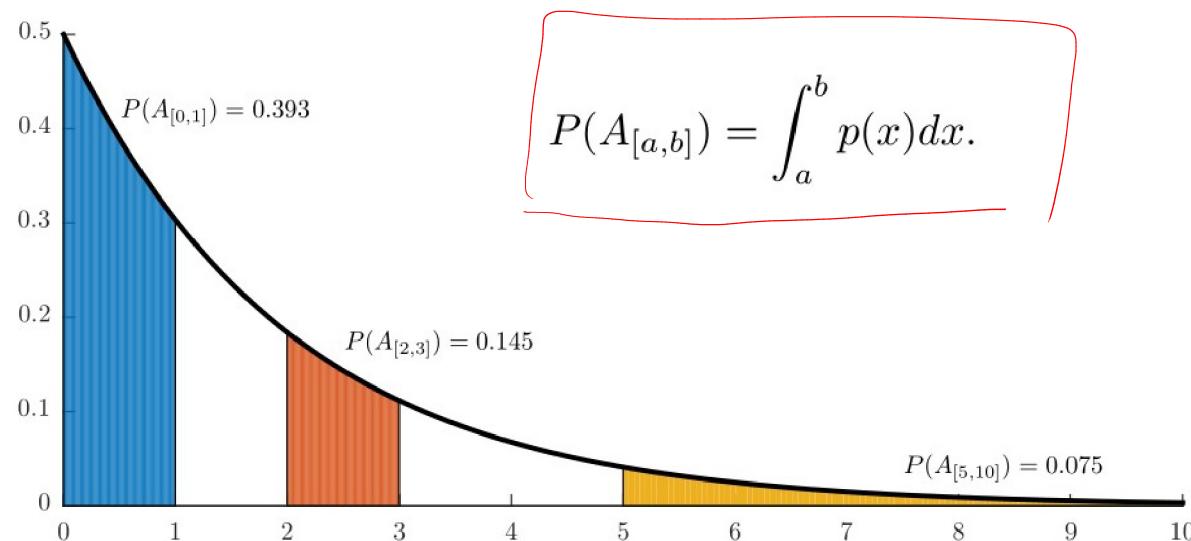
Probability vs. Density

- Suppose we consider the rainfall on an average day $r \in [0, \infty]$
- **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

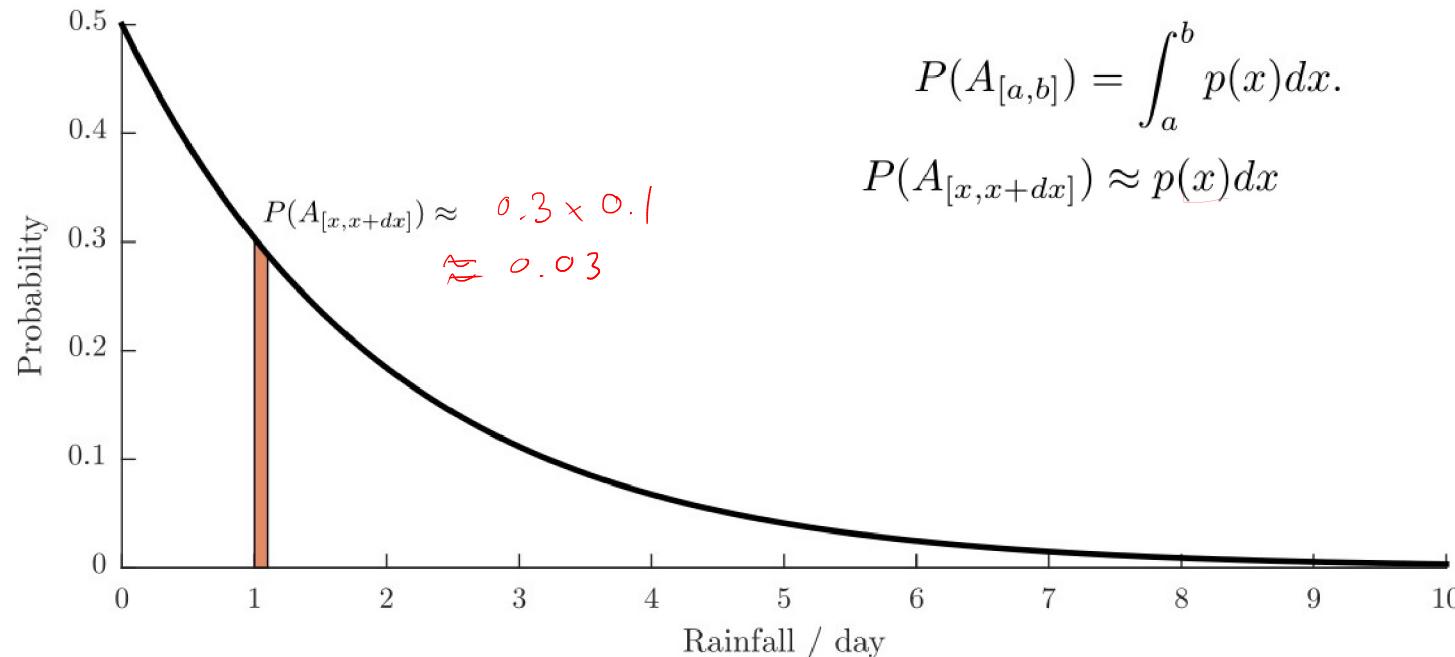
- These probabilities can be **represented** as integrals
- **Events** are **intervals**, the **probability** is the **integral**, the **curve** is the **density**



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

Probability vs. Density

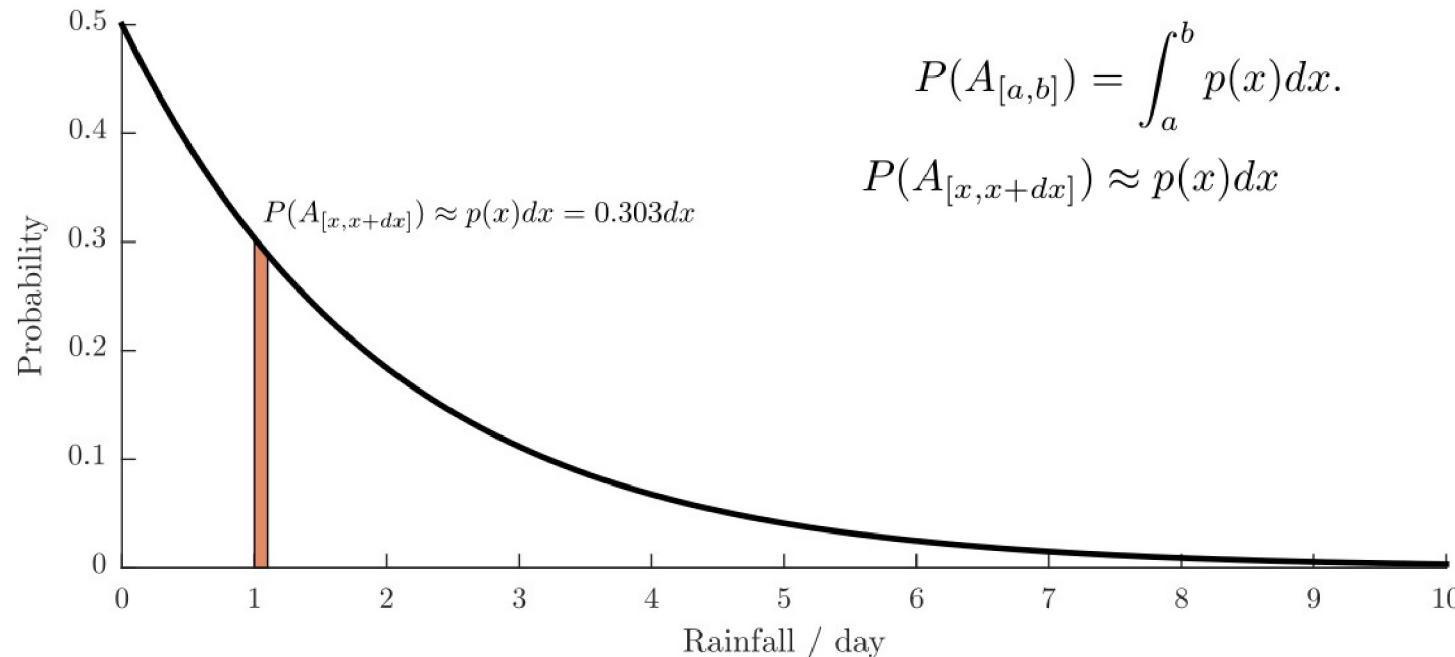
- These probabilities can be **represented** as integrals
- **Events** are **intervals**, the **probability** is the **integral**, the **curve** is the **density**
- What is the probability there will be between 1 and 1.1 mm of rain?



$A_{[a,b]}$: 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?



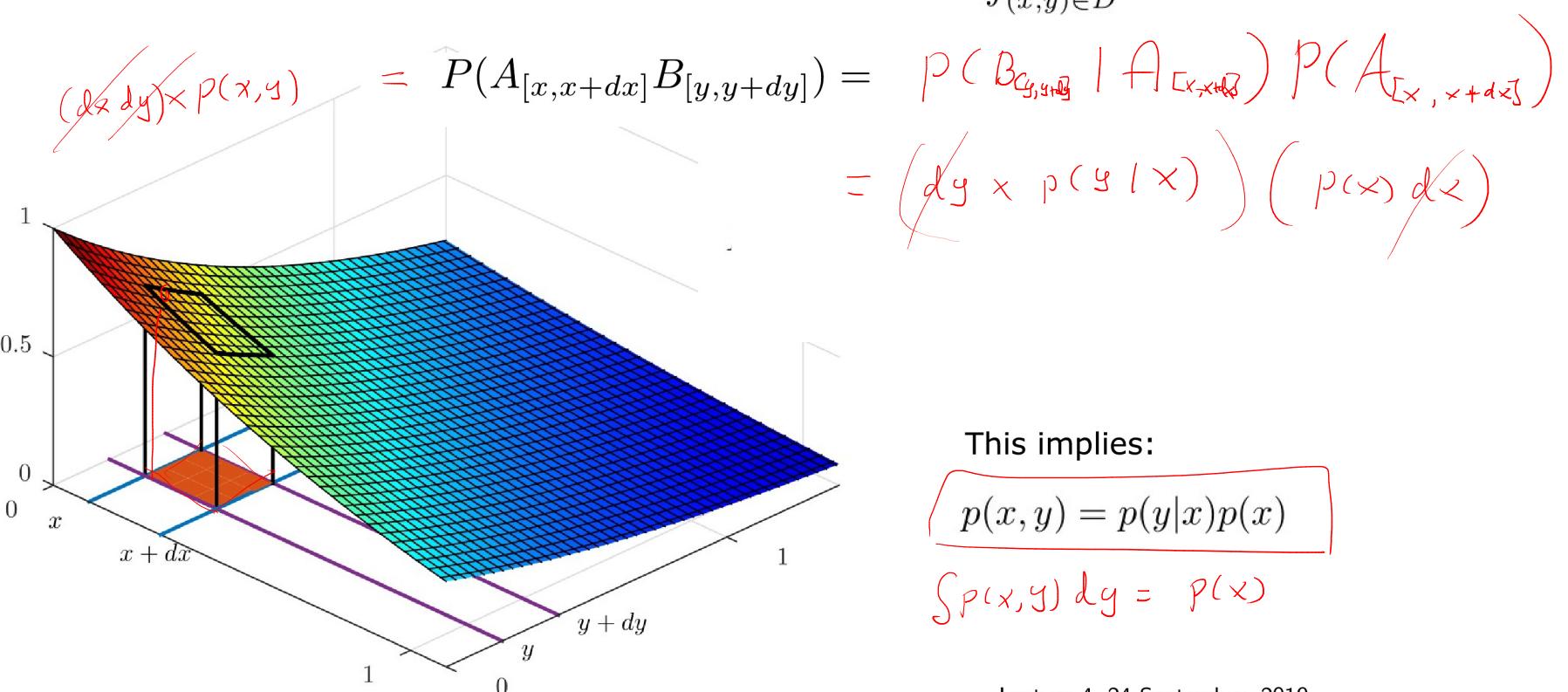
$A_{[a,b]} : \text{There will be between } a \text{ and } b \text{ mm of rain}$

Probability vs. Density

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

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

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



Probability vs. Density

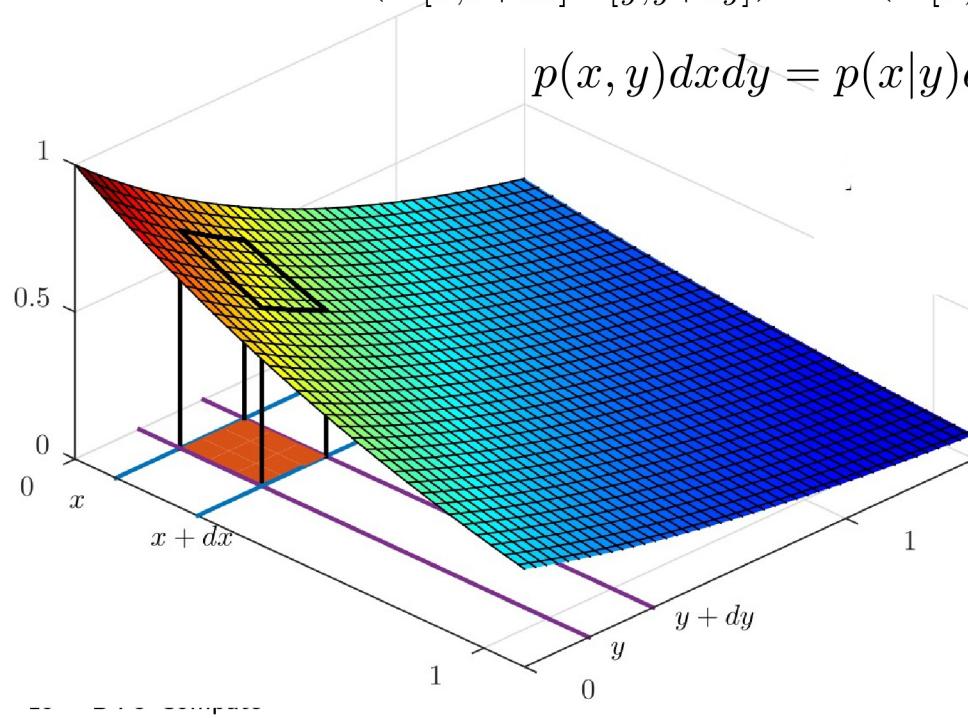
- For two variables x and y , the **probability** is an integral over an **area**

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

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

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

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



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:

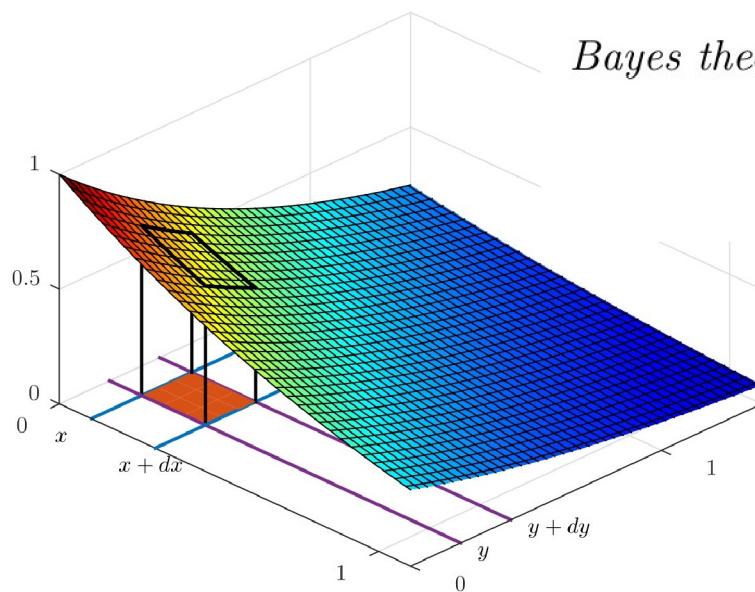
$$\int dx \ p(x|z) = 1$$

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:

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

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

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**

$$\bar{x} = \mathbb{E}[x] = \int x p(x) dx.$$

- **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:

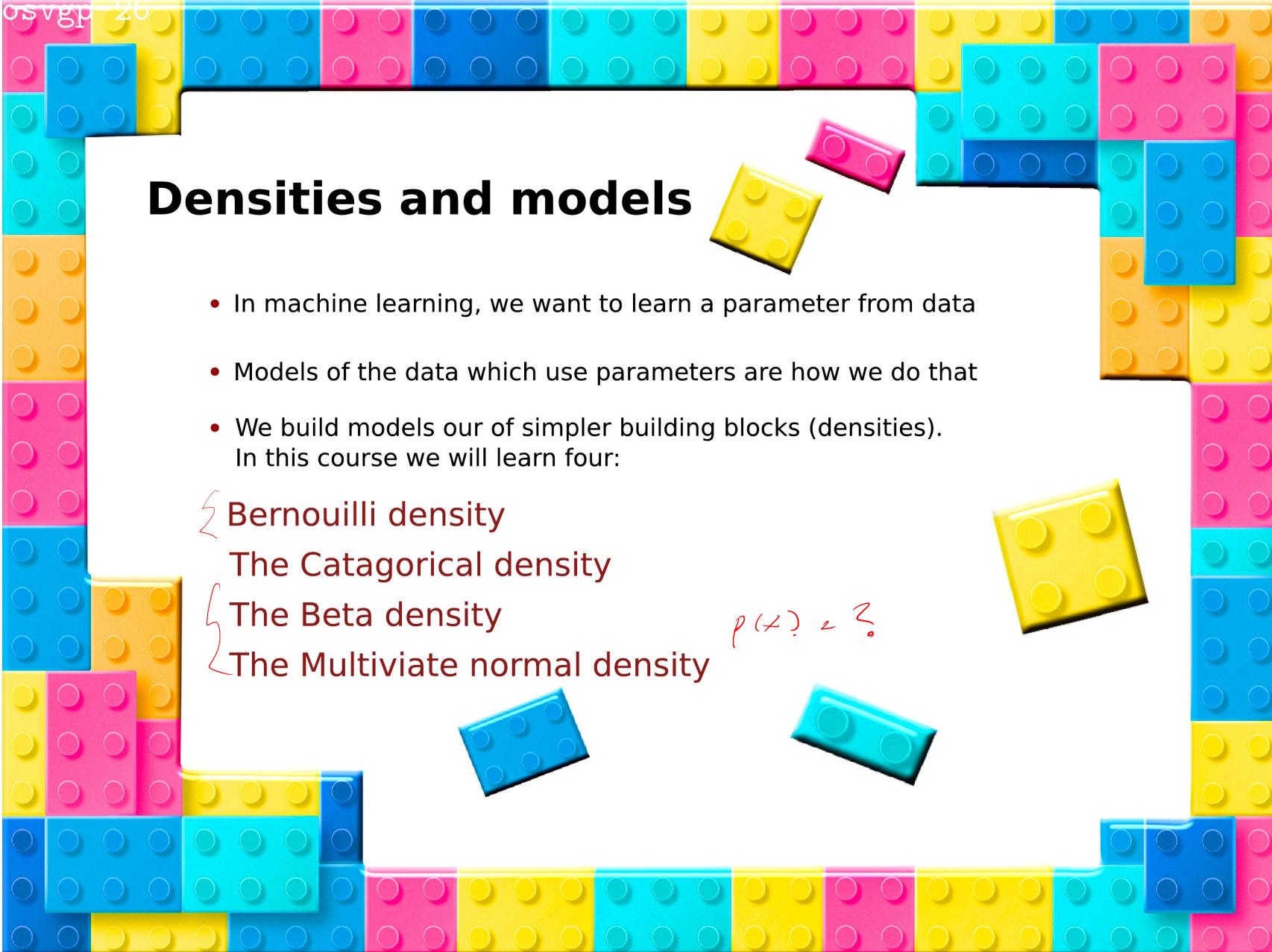
The Bernoulli density

The Categorical density

The Beta density

The Multivariate normal density

$$p(x) \propto ?$$



The multivariate normal distribution

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

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

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

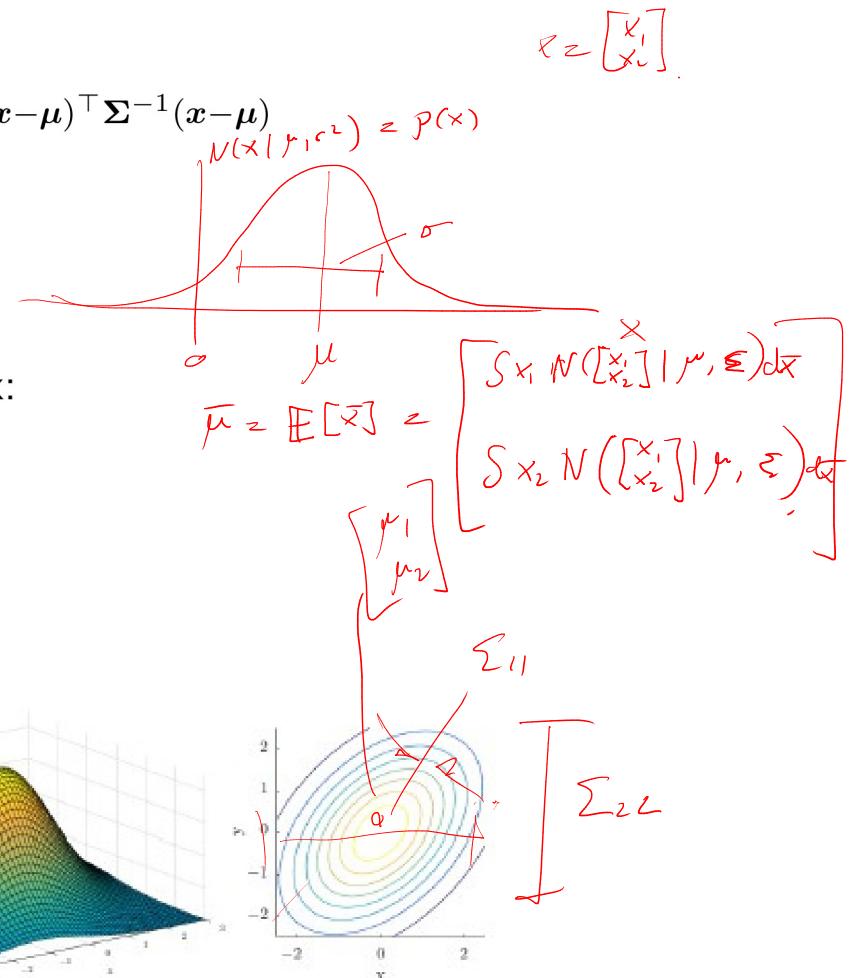
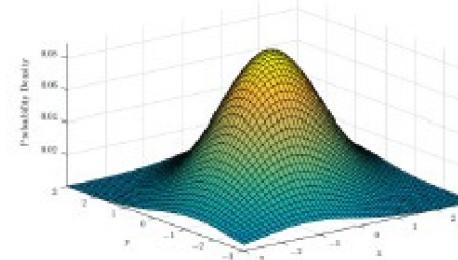
$\boldsymbol{\mu}$ is the mean vector and $\boldsymbol{\Sigma}$ is the covariance matrix:

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

- Example: 2-dimensional Normal distribution

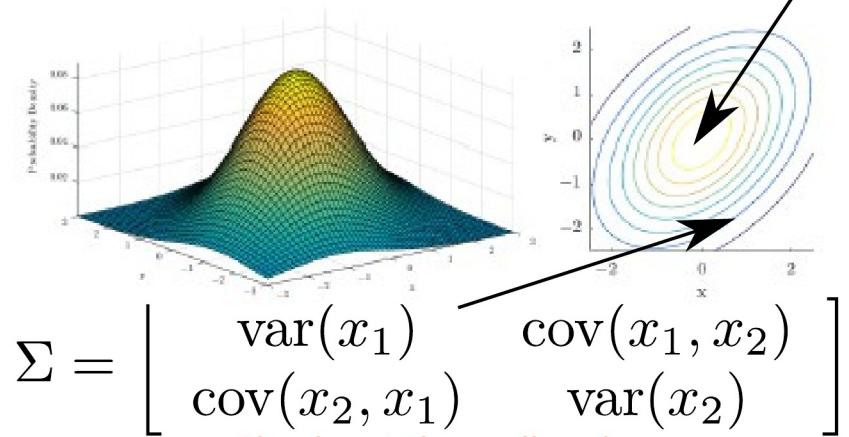
$$\boldsymbol{\mu} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \end{bmatrix}$$

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



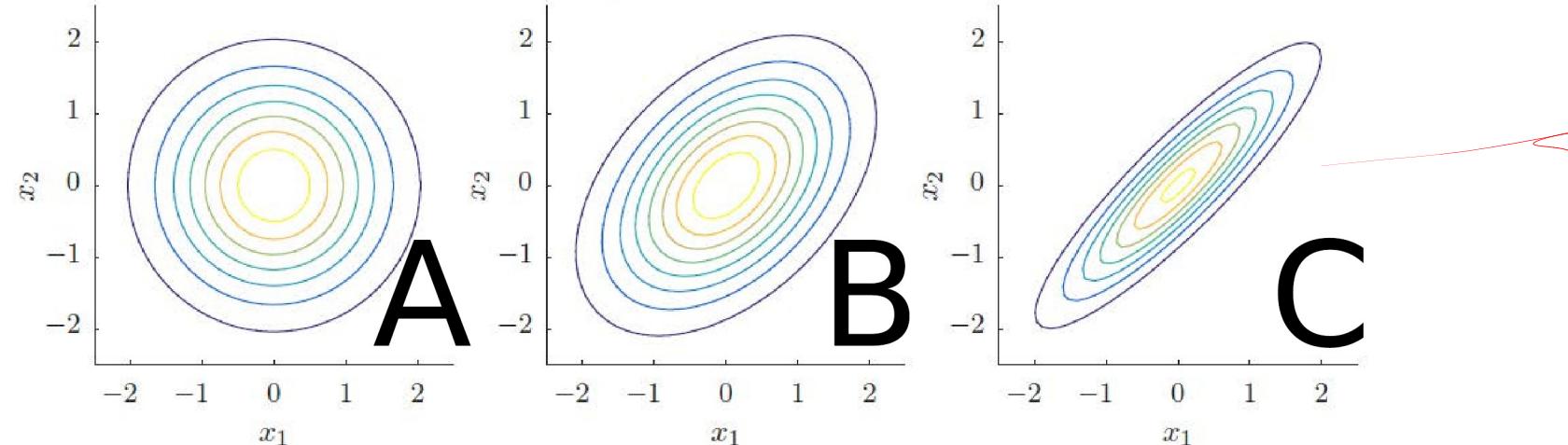
Quiz 2, Covariance

- Match the covariances to the contour plots



- 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>



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

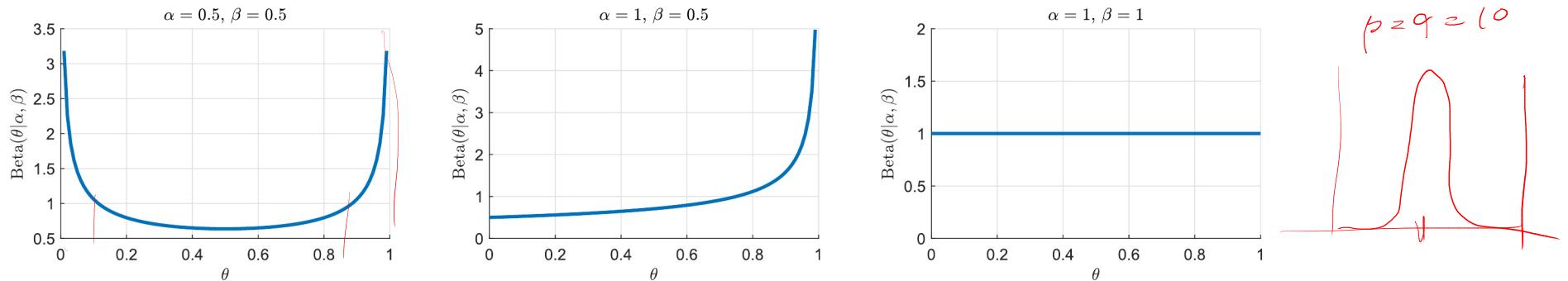
in all

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

Beta distribution

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

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



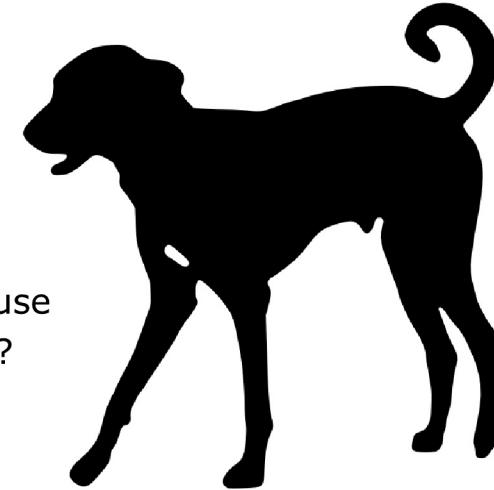
$\alpha, \beta > 0$ are related to the variance and mean

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

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?

b



- 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?

b



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**

$$p(b|\theta) = \underbrace{\theta^b(1-\theta)^{1-b}}_{\text{y}}$$

- 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

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

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

$$\underline{p(\theta|b)} = \frac{p(b|\theta) p(\theta)}{\int_{\theta} p(b|\theta) p(\theta) d\theta}$$

- Assume $p(\theta) = \text{Beta}(\theta | \alpha, \beta)$

$$p(\theta|b, \alpha, \beta) = \frac{\theta^m (1-\theta)^{N-m} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}}{\int \theta^m (1-\theta)^{N-m} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta}$$

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

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!**

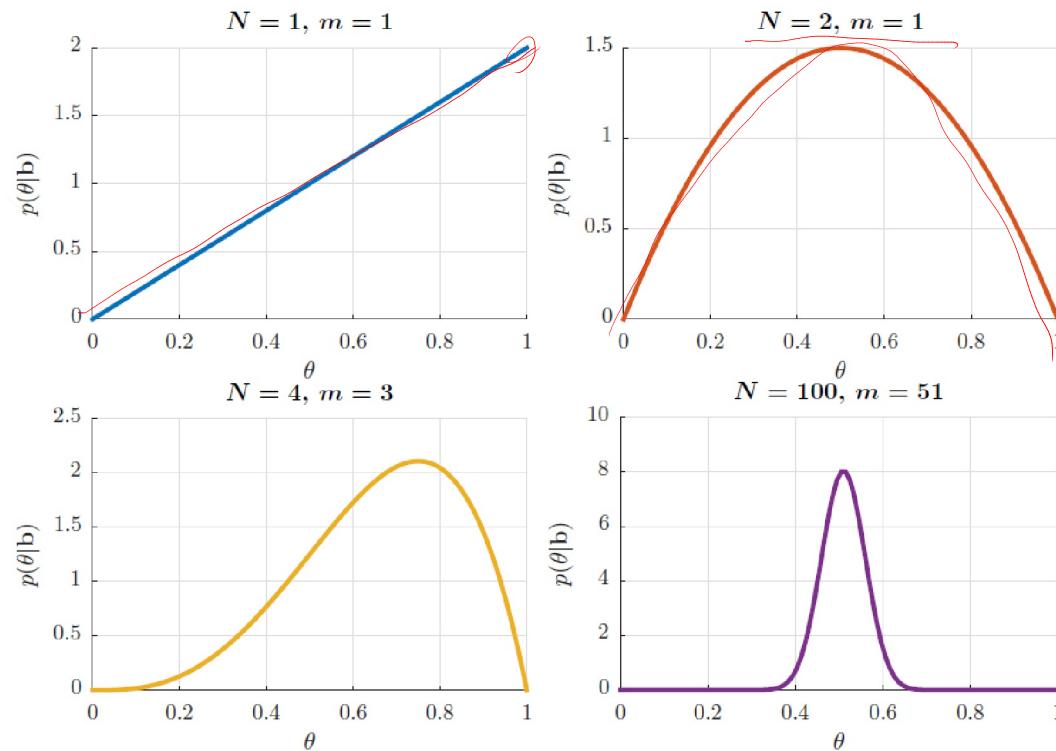
$$p(\theta|\mathbf{b}) = \frac{p(\mathbf{b}|\theta)p(\theta)}{p(\mathbf{b})} = \frac{p(\mathbf{b}|\theta)p(\theta)}{\int_0^1 p(\mathbf{b}|\theta')p(\theta')d\theta'}$$

- 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



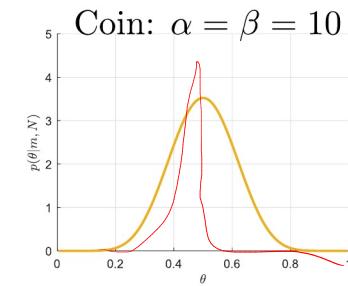
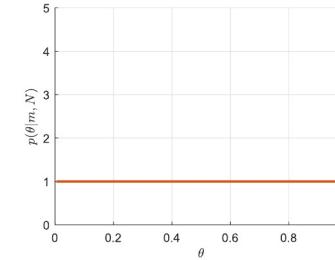
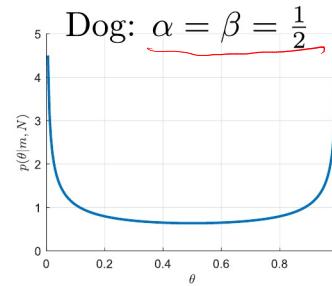
- Your friend just got a dog.
- The dog can either be in the doghouse or outside
- The first four times you come by the dog is in the doghouse
- What is the chance the dog is in the doghouse tomorrow?



- 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?

Prior

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



Dogs and coins



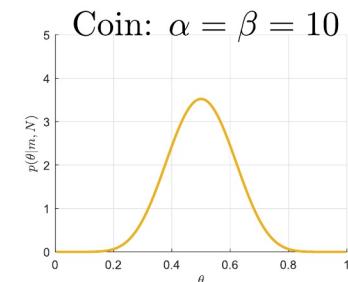
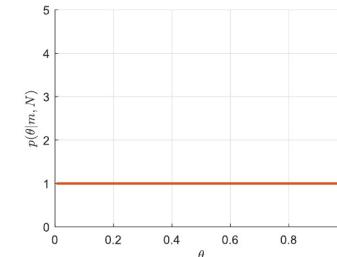
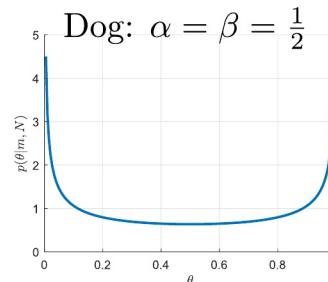
- Your friend just got a dog.
- The dog can either be in the doghouse or outside
- The first four times you come by the dog is in the doghouse
- What is the chance the dog is in the doghouse tomorrow?



- 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?

Prior

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



Likelihood

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

Dogs and coins



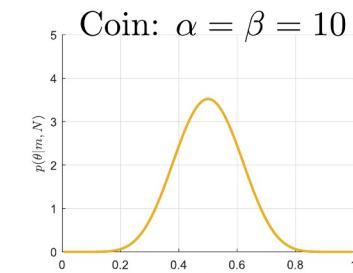
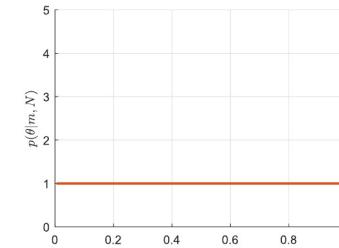
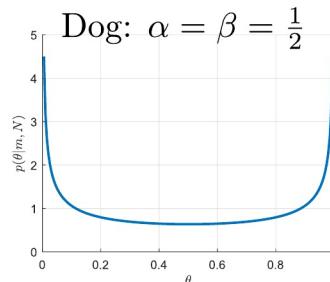
- Your friend just got a dog.
- The dog can either be in the doghouse or outside
- The first four times you come by the dog is in the doghouse
- What is the chance the dog is in the doghouse tomorrow?



- 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?

Prior

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



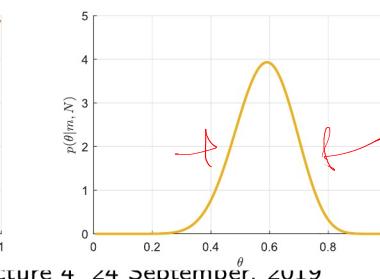
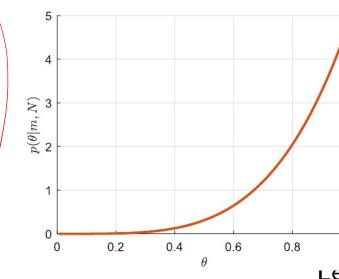
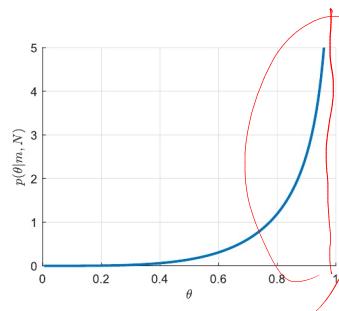
Likelihood

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

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

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

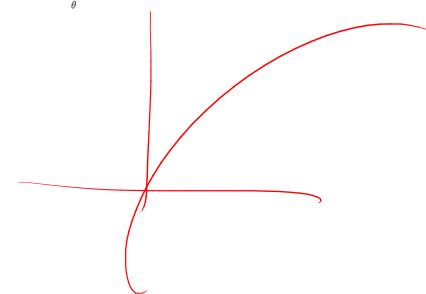
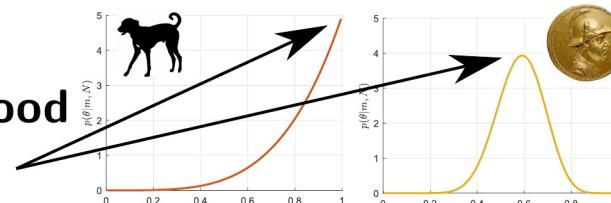


14.15

Learning principle ii: 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 $f(x) > 0$:

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

(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:

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

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^w
- 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

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

- The following are equivalent:

Maximum likelihood principle

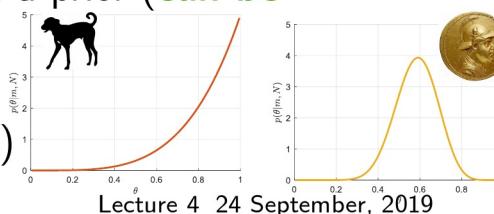
$$\text{Maximize: } \mathbf{w}^* = \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \arg \max_{\mathbf{w}} \left(\log \left(\prod_{i=1}^N p(y_i|\mathbf{x}_i, \mathbf{w}) \right) + \log p(\mathbf{w}) \right)$$

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

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

- 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)



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**

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

- 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:

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

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

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

The drawing shows me at one glance what might be spread over ten pages in a book."

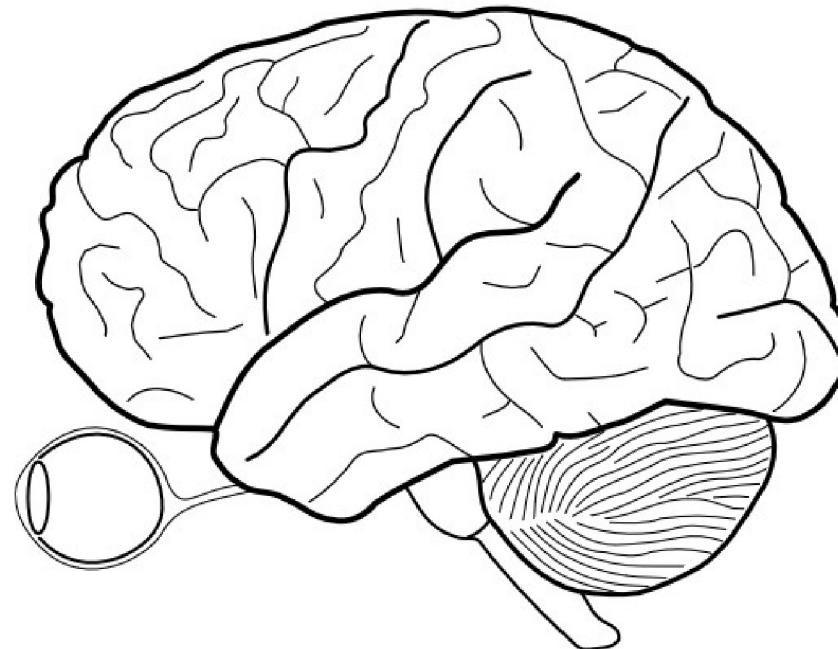
- Ivan S. Turgenev's novel Fathers and Sons, 1862.

Use a picture. It's worth a thousand words."

- Arthur Brisbane to the Syracuse Advertising Men's Club, in March 1911

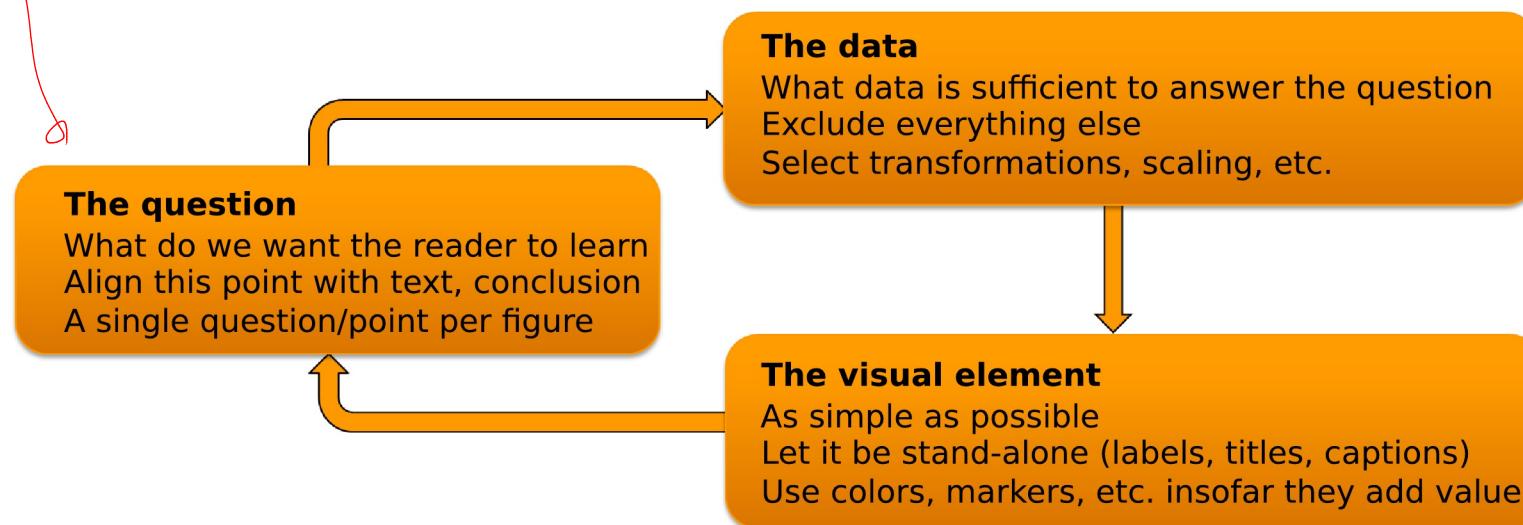
Visualization

- 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*)

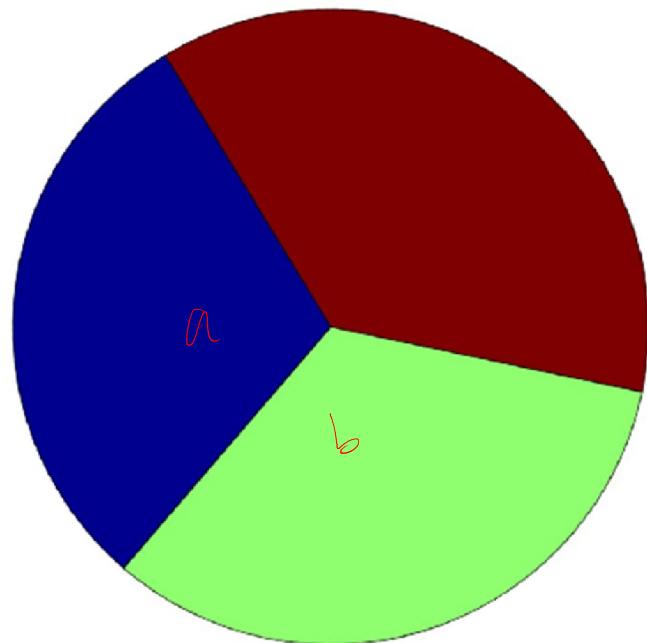


Important choices for visualizations

- **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

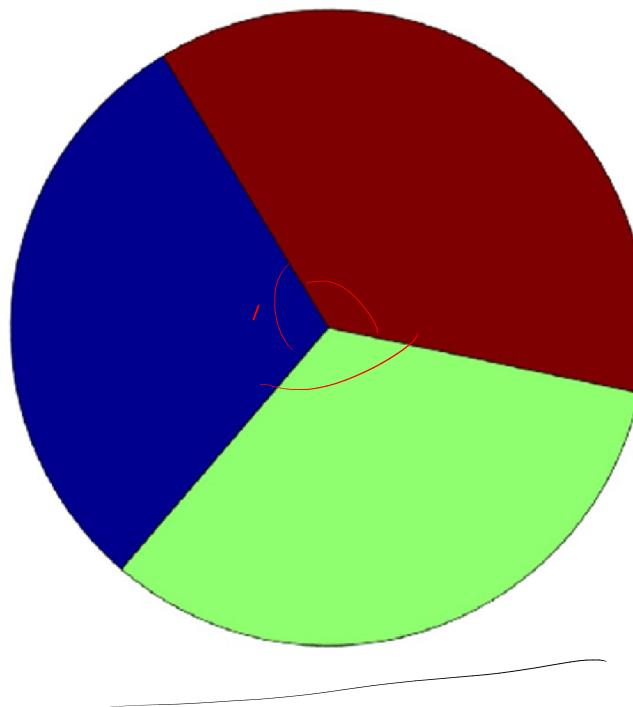
- **Area represents proportion**
 - Which is smallest, middle, and largest?
 - What are the proportions approximately?



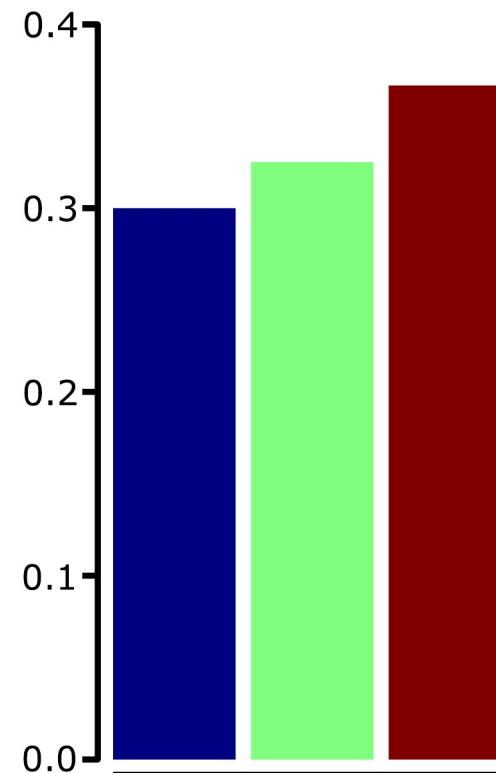
Representation

- **Area represents proportion**

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



- **Height represents proportion**



Selection

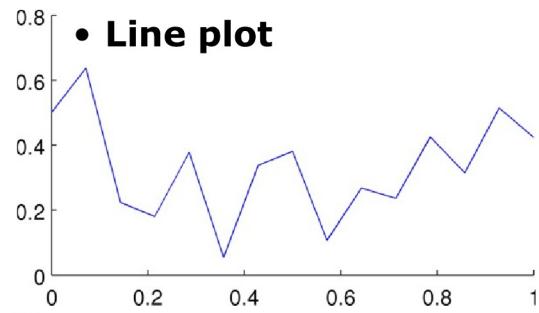
- Elimination or de-emphasis of certain objects or attributes
- A subset of **attributes**
 - **Why?** A graph can only show so many attributes – focus on the relevant
 - **How?**
 - Dimensionality reduction
 - Plot pairs of attributes
- A subset of **objects**
 - **Why?** A graph can only show so many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

Types of plots

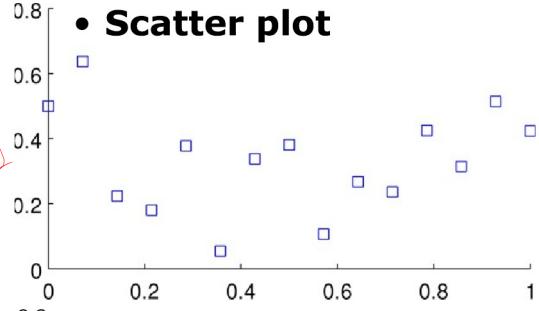
- **Distribution of a single attribute**
 - Histogram
 - Empirical cumulative distribution
 - Percentile plots
 - Box plot
- **Relation between attributes**
 - 2D histogram
 - Heat maps and contour plots
 - Scatter plots
- **Visualization of high-dimensional objects**
 - Matrix plots
 - Parallel coordinates
 - Star plots

(x, y)

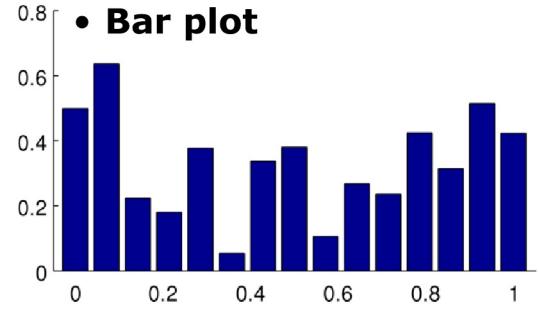
Basic plots



~~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

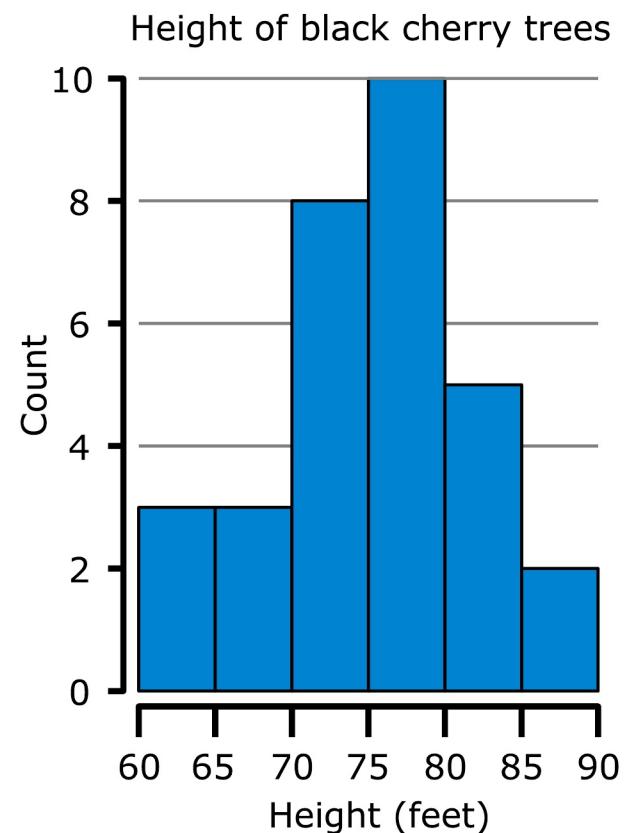
Flower ID	Attribute			
	Sepal Length	Sepal Width	Petal Length	Petal Width
1	5.1	3.5	1.4	0.2
2	4.9	3.0	1.4	0.2
3	4.7	3.2	1.3	0.2
4	4.6	3.1	1.5	0.2
.
.
150	5.9	3.0	5.1	1.8

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

Distribution of a single attribute

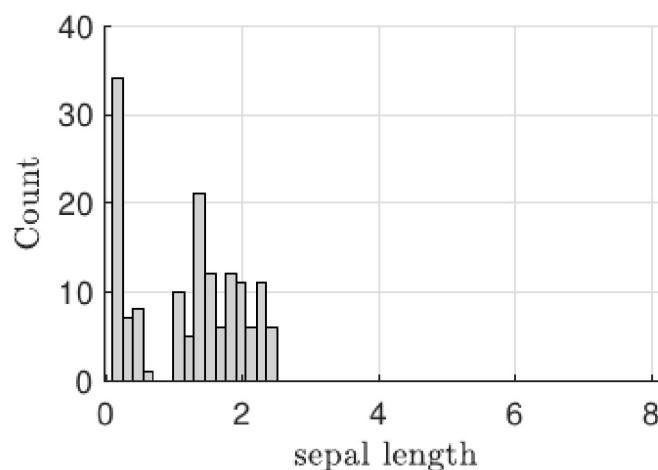
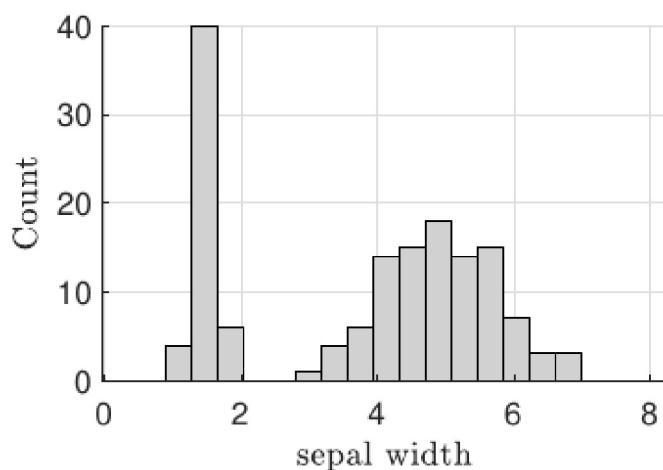
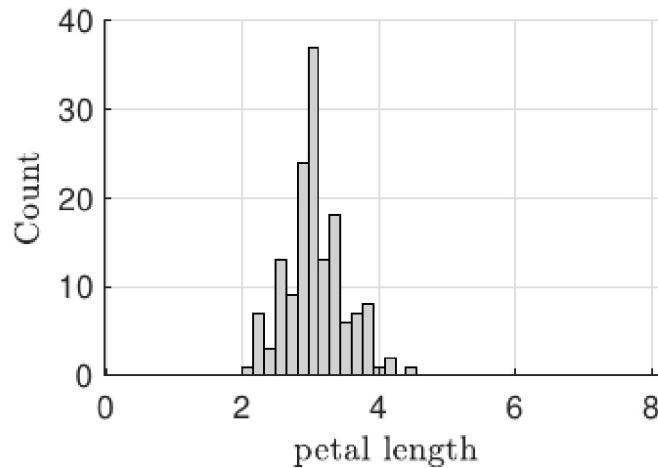
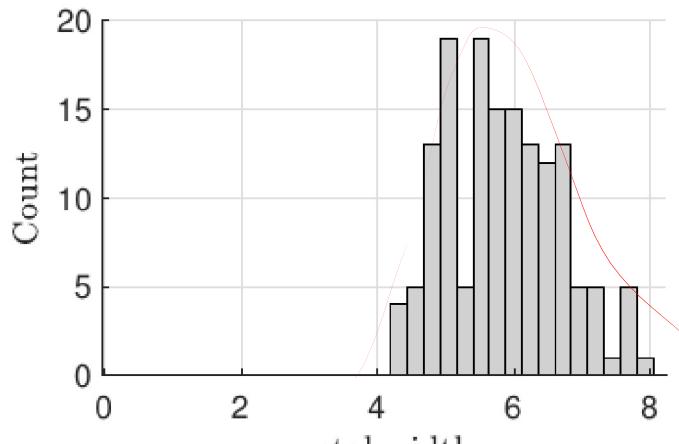
Histograms

- **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



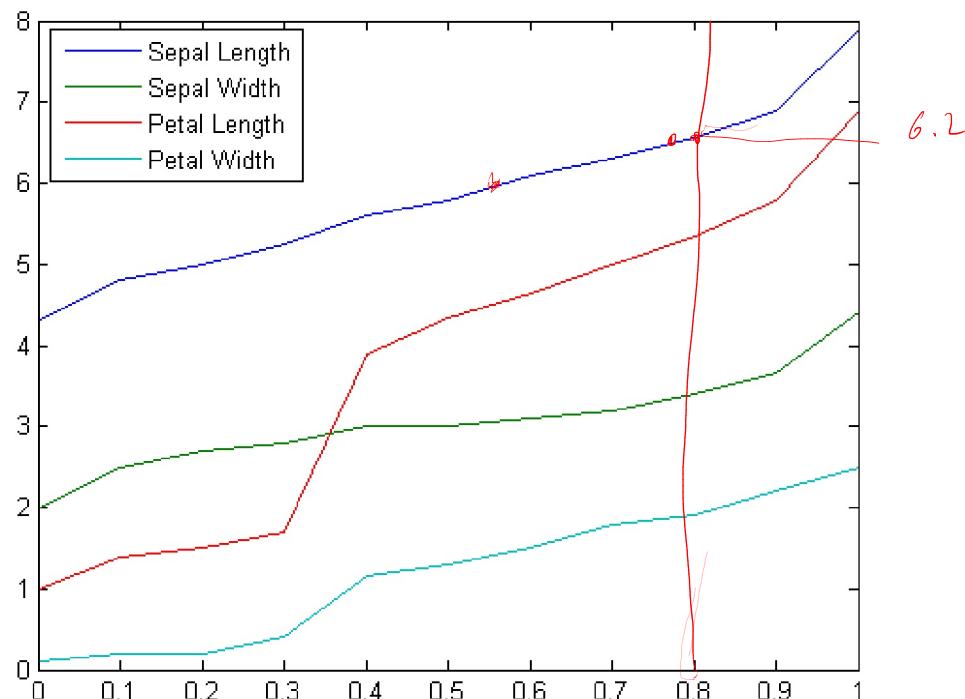
$$H = \{60, 64, 64, 66, 67, 69, 71, 72, 72, 72, 72, 73, 74, 74, 75, 75, 76, 76, 76, 77, 77, 78, 78, 79, 80, 80, 81, 82, 84, 85, 85, 89\}$$

Histograms of the Iris data attributes



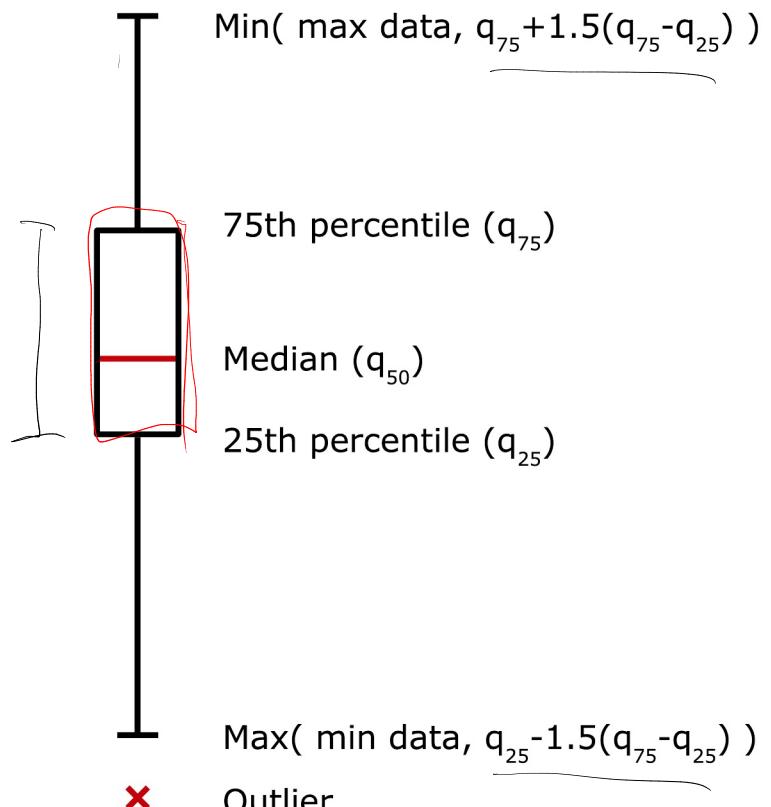
Percentile plots

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

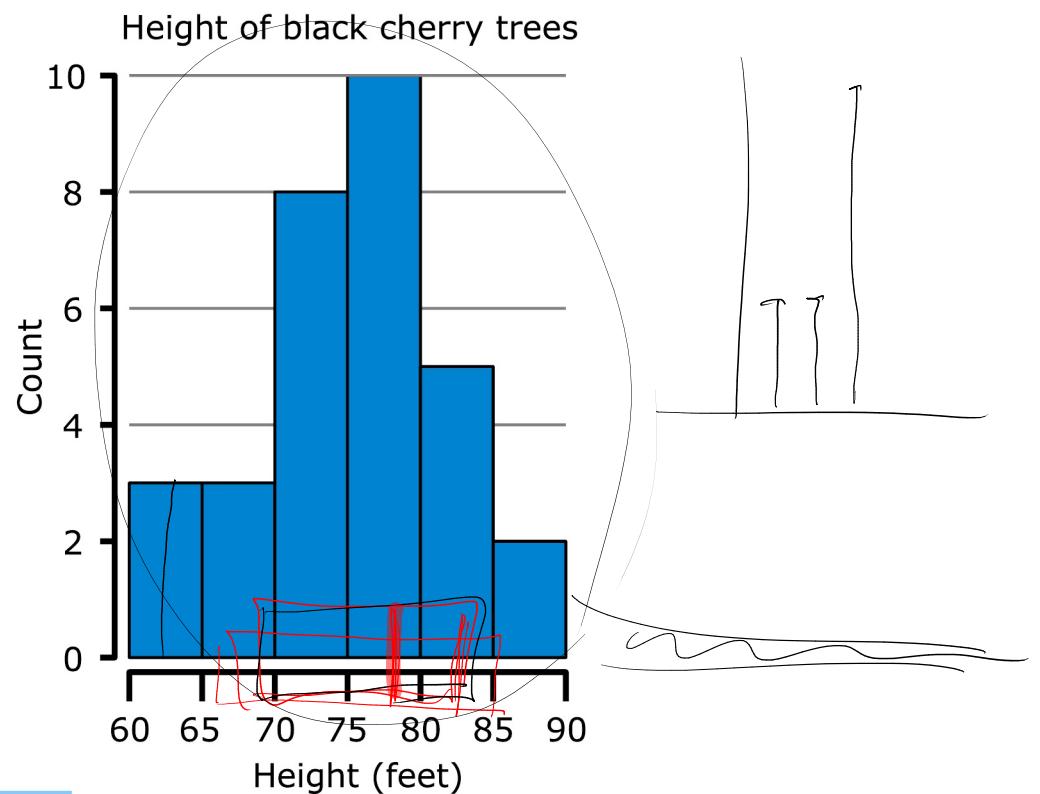


```
prctile = 0:0.1:1;
Y = quantile(X,prctile);
plot(prctile,Y);
legend(attributeNames);
```

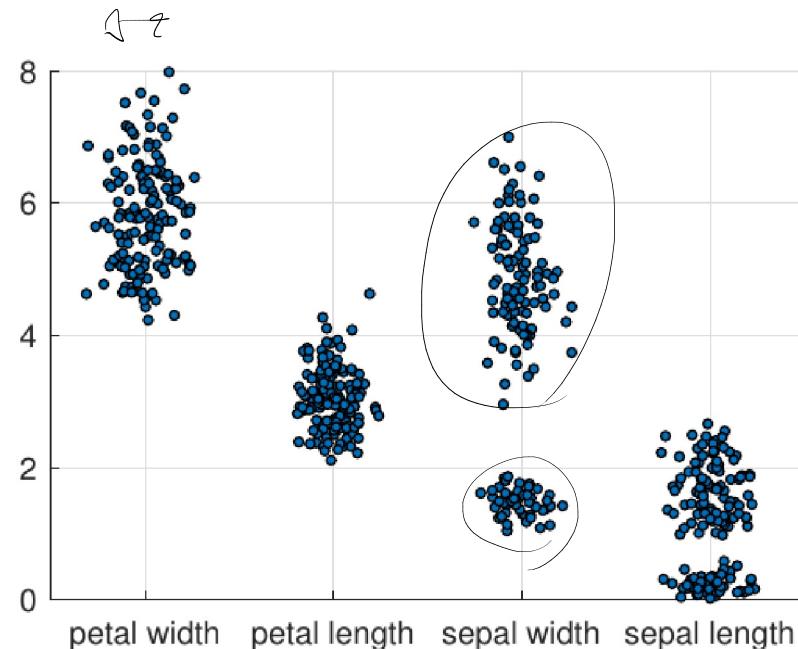
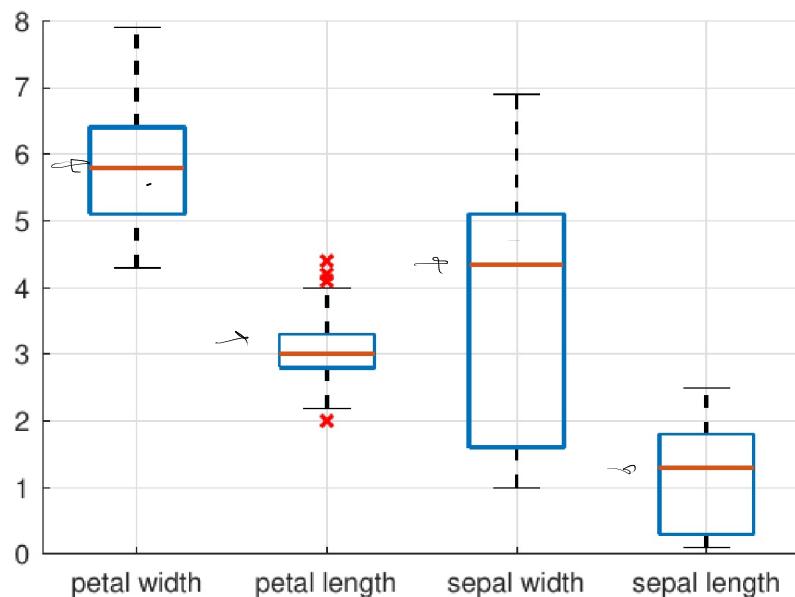
Box plots



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



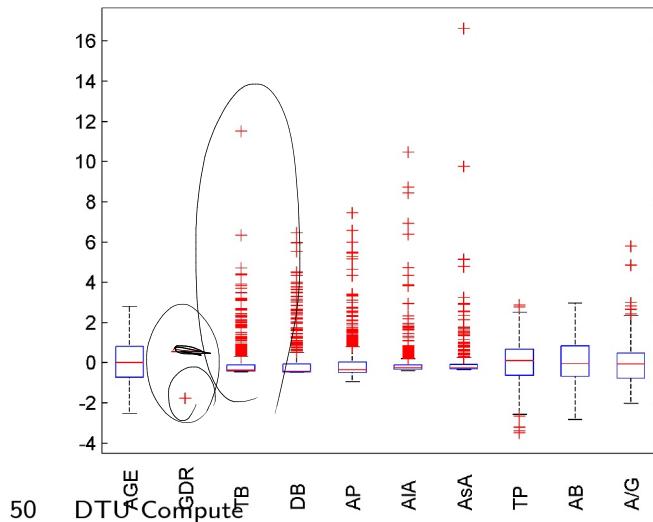
Box plots



Quiz 3, Boxplots (Fall 2012)

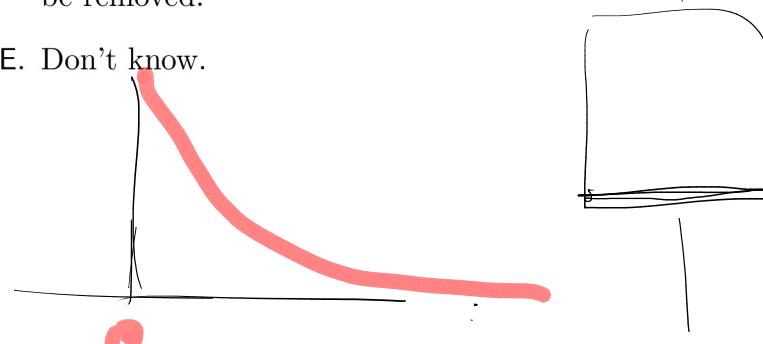
No.	Attribute description	Abbrev.
x_1	Age (in years)	AGE
x_2	Gender (Female=0, Male=1)	GDR
x_3	Total Bilirubin	TB
x_4	Direct Bilirubin	DB
x_5	Alkaline Phosphotase	AP
x_6	Alamine Aminotransferase	AlA
x_7	Aspartate Aminotransferase	AsA
x_8	Total Proteins	TP
x_9	Albumin	AB
x_{10}	Albumin to Globulin ratio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

Table 1: Liver disease dataset.



The attributes x_1-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 AlA 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.



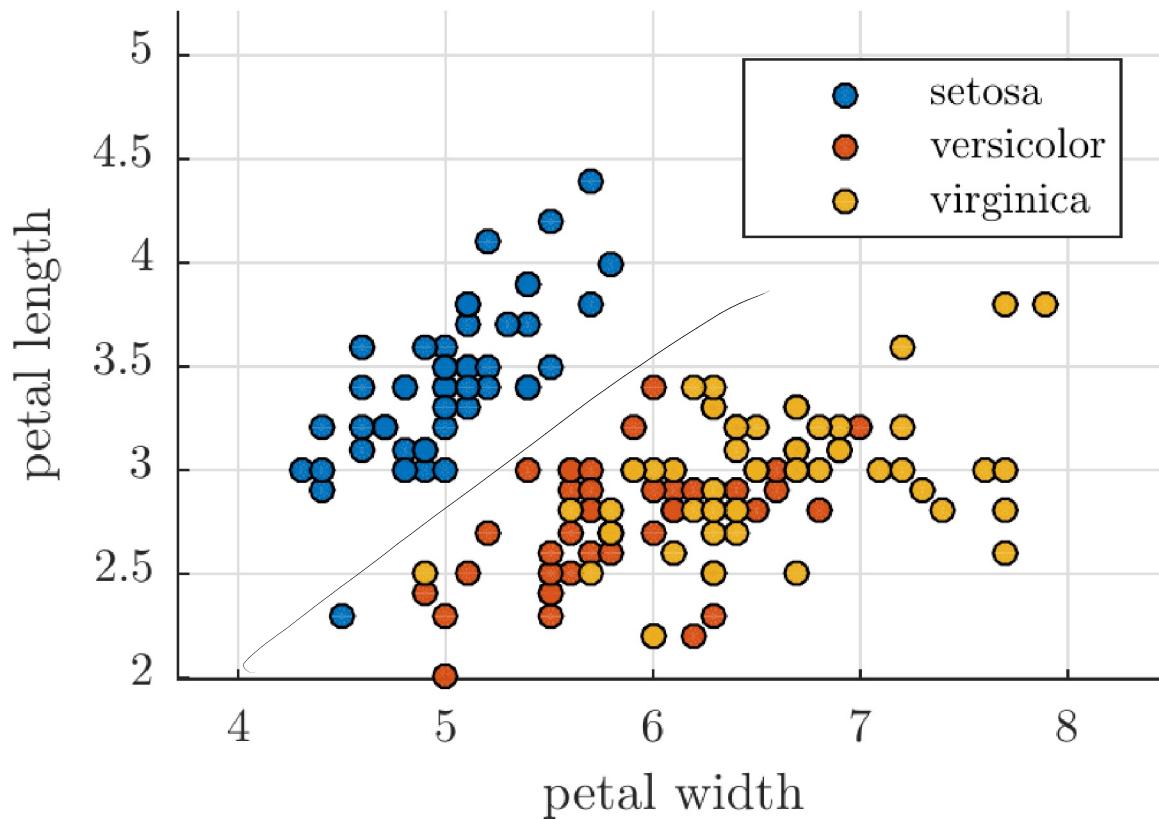
The 25th and 50th percentile but not the 50th and 75th percentiles of the attribute DB coincides. AlA and AsA will not necessarily be highly correlated even though their distributions may have a similar shape (hence, this is correct). For attributes to be correlated it is important they take on high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

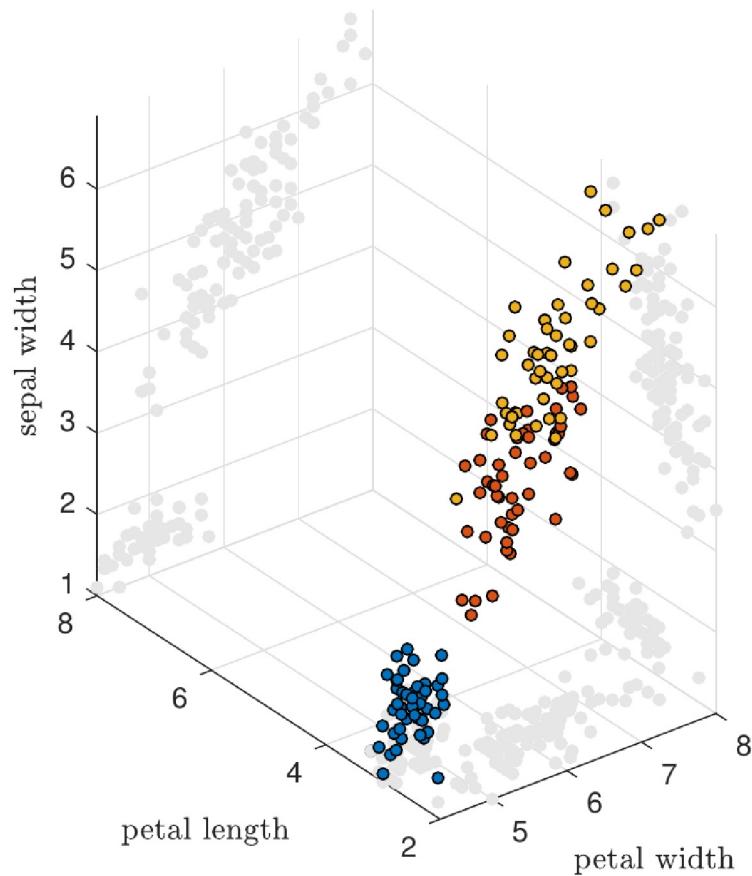
Scatter plots

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



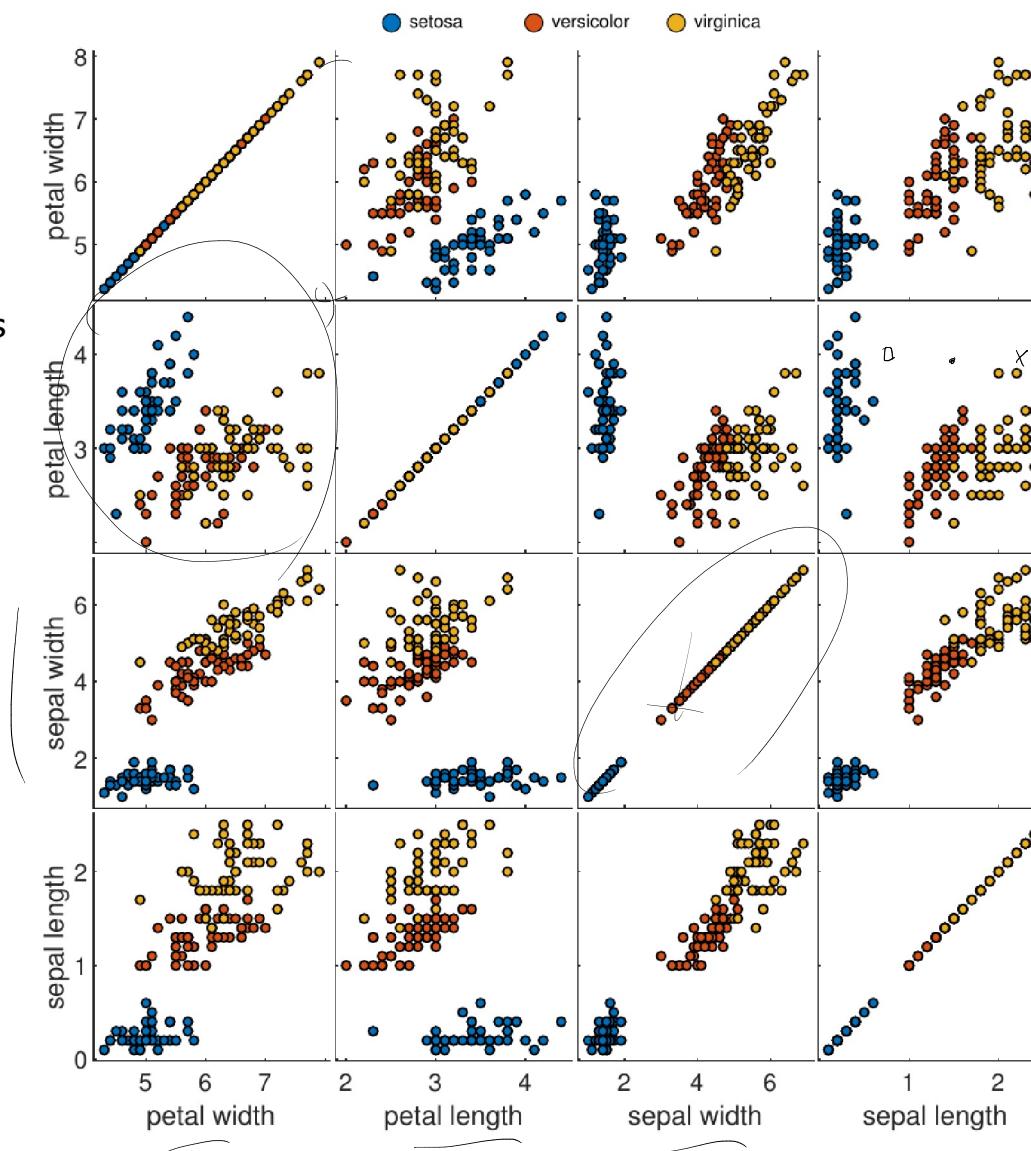
Scatter plots

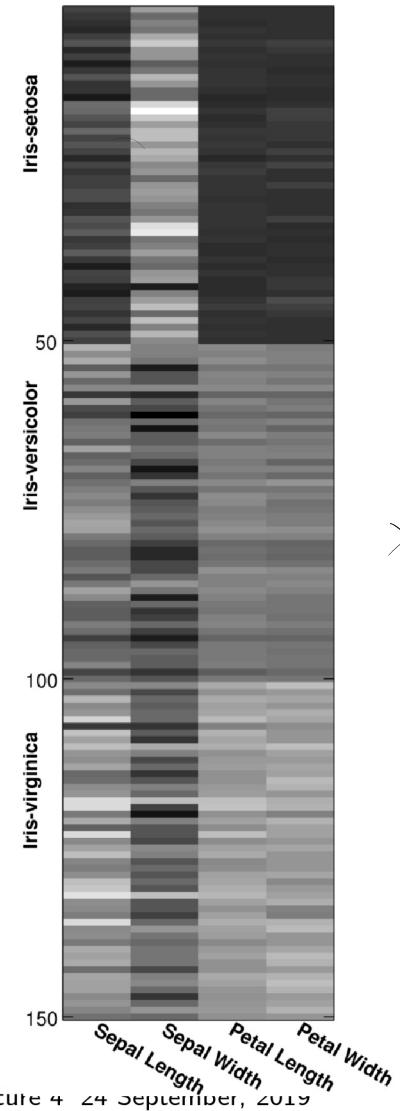
- Shows **relation** between attributes
 - Assess dependence between attributes
 - Used with classes to assess separability
 - 3D plots are often confusing;
 avoid if possible



Scatter plots

- Scatter plot matrix
 - All pairs of attributes





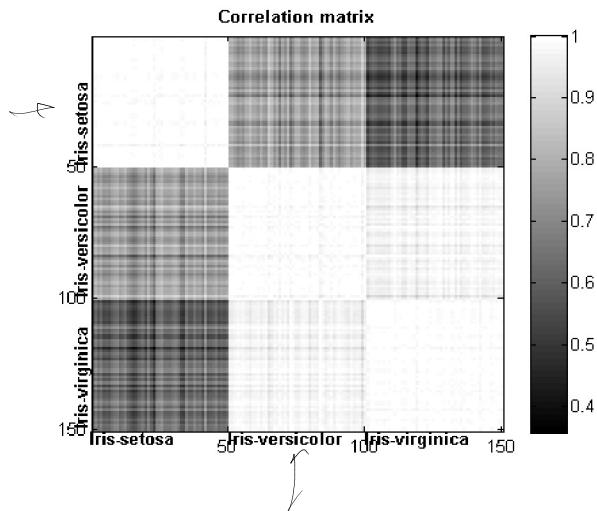
Matrix plots

- **Plot of raw data matrix**

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

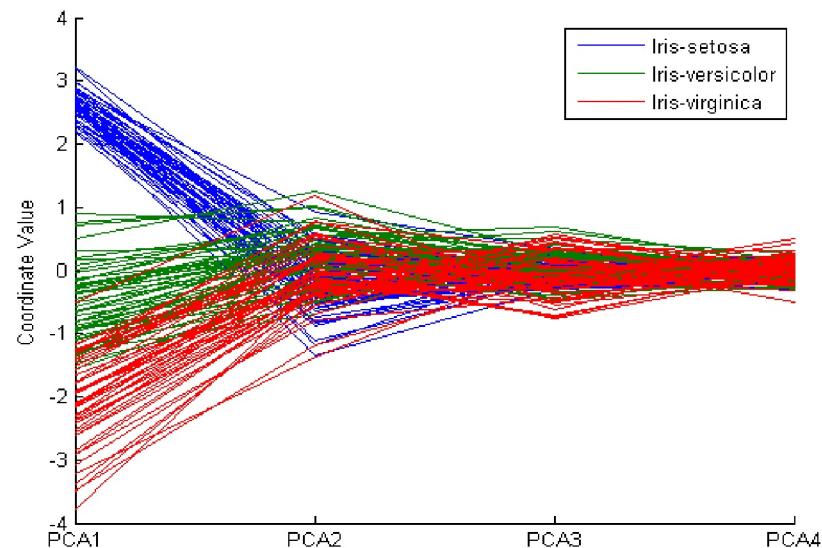
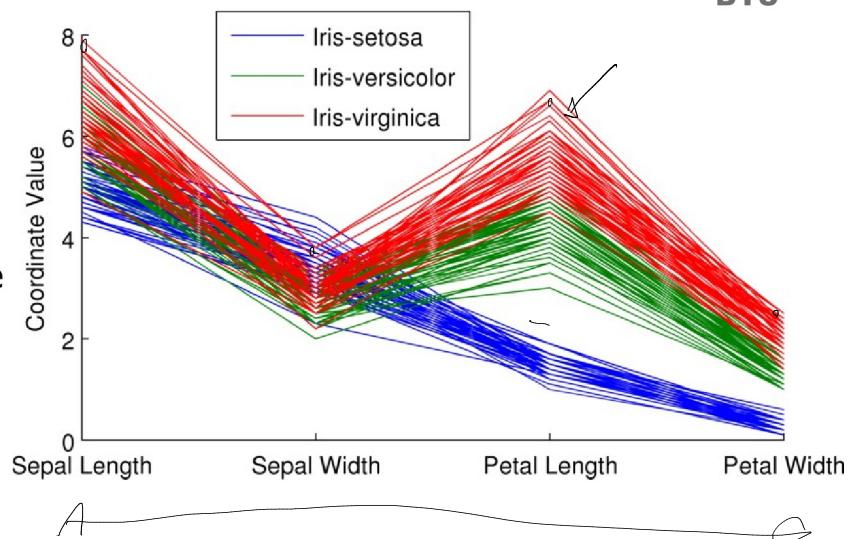
- **Plots of similarity matrices**

- 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

Lecture 4 24 September, 2019

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/>



The Mahalanobis distance

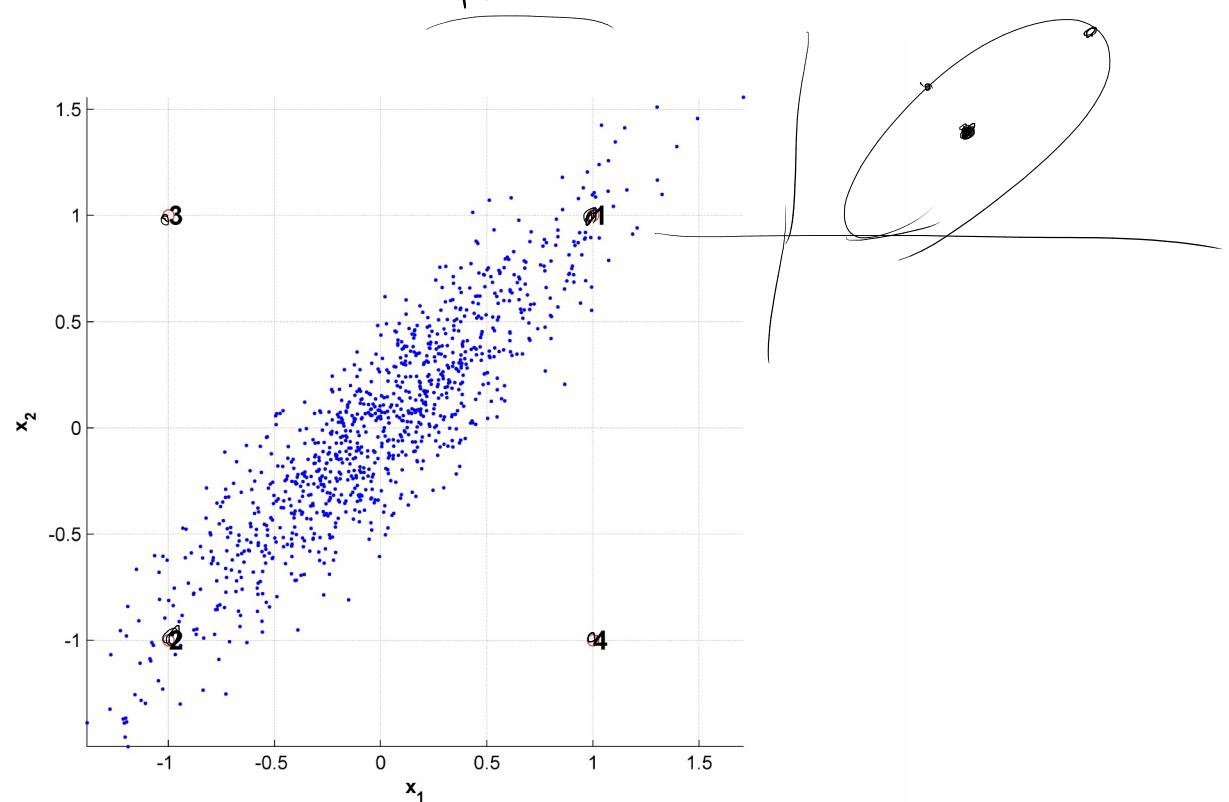
$$N(x|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$$

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?

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



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

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

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>)