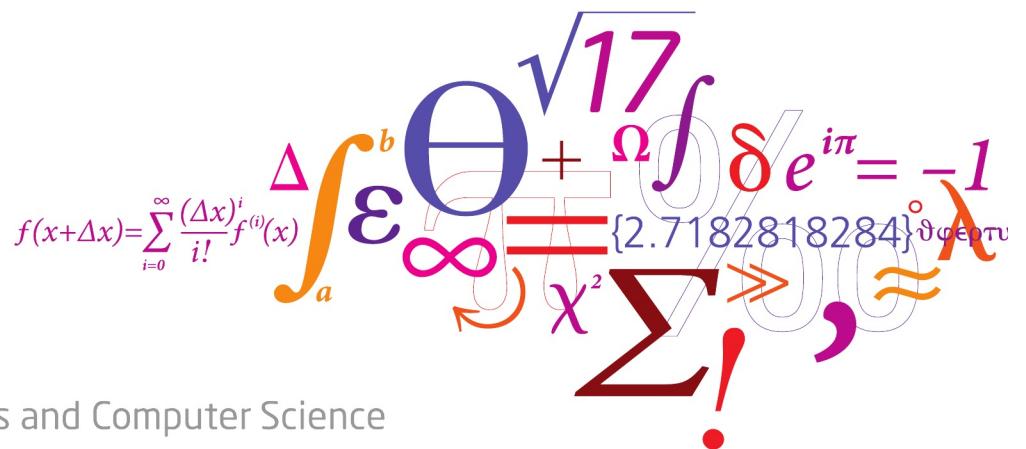


02450: Introduction to Machine Learning and Data Mining

Mixture models and density estimation

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

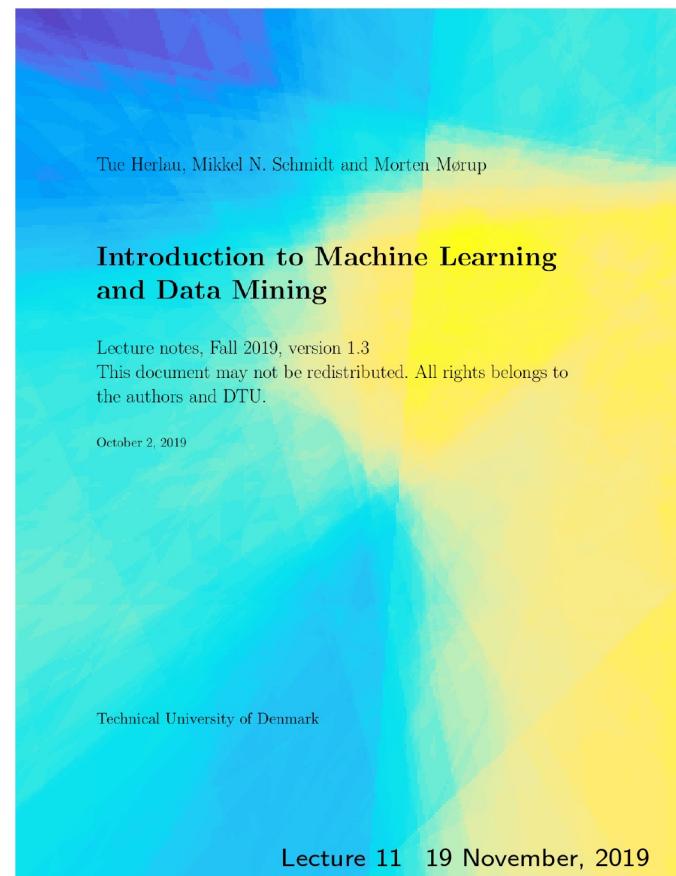
Remember you can evaluate the course on DTU inside
Today + check exam on Inside.



Feedback Groups of the day:

Lau Halkier Wandall, Mabel-Brenda Ifeoma Okikaa, Jacob Pedersen, Carla Cristina Salazar Navalón, Silas El-Azm Stryhn, Yifei Xue, Anqi Yan, Mengge Hu, Naveen Karun Somasundaram, Carlos Parra, Ashwinth Mathivanan, Søren Bojer Jørgensen, Simon Majgaard, Joachim Secher, Michael Baand Severinsen, Sofie Betzer Rossen, Joakim Nilaus Hven, Frederik David Damsgaard Popp, Johannes Schou, Emma Dam Mortensen, Stefanie Wøhler Nielsen, Anna Bjerring Jensen, Sophia Troldborg Ohmann, Jens Damholt Richardt, Richard Thyssen Nørby Larsen, Niels Christian Dahl, Kasper Fibæk Schlosser, Lasse Pærgård Kristiansen, Jonas Weile, Evangelos Kerasidis, Allan Ozvan, Casper Chris Adriaan Bekkers, Michiel Goderie, Yannick Christopher Heijne, Maaike Elgersma, Maarten Marijn van Elst, Sean Christian Lindholm, Troels Lund, Michael Thorbjøll Kristensen, Ali Jamal Jomeh, Jonathan Skovsholm Winther, Anthony Laye, Jimmy Xu, Andrew Nitu, Ruoyu Huang, Takuya Omori, Hui Jun Yap, Thor Larsen, Jakob Hammer Hedemann, Martin Illum, Hans Christian Bechsøfft Mikkelsen, Nicholas Rose, Elig Saraliev, Mikkel Nørgaard Schmidt, Mareva Ji Simon, Eskild Børsting Sørensen

Reading material: Chapter 19, Chapter 20



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

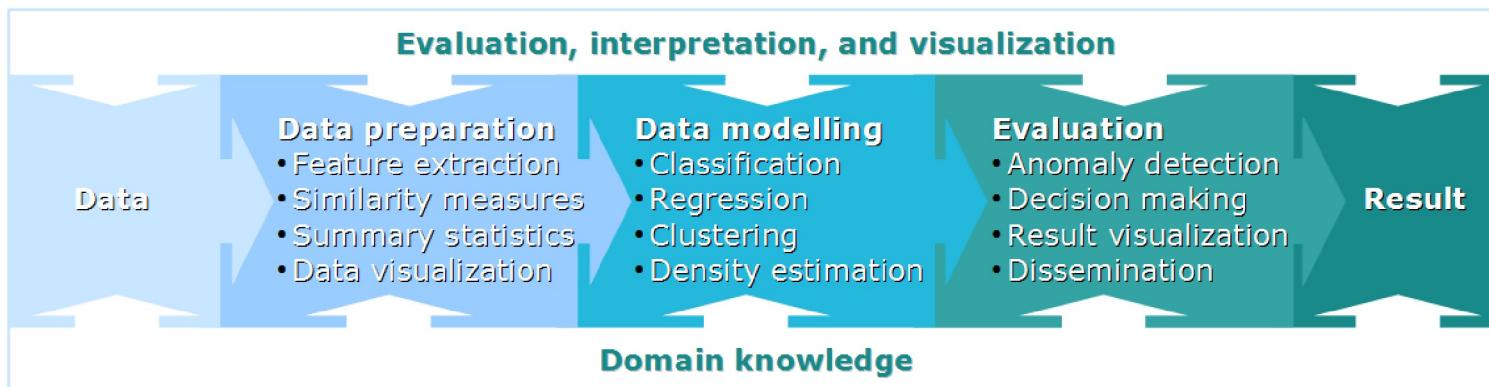
⑫ Association mining

26 November: C21

Recap

⑬ Recap and discussion of the exam

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



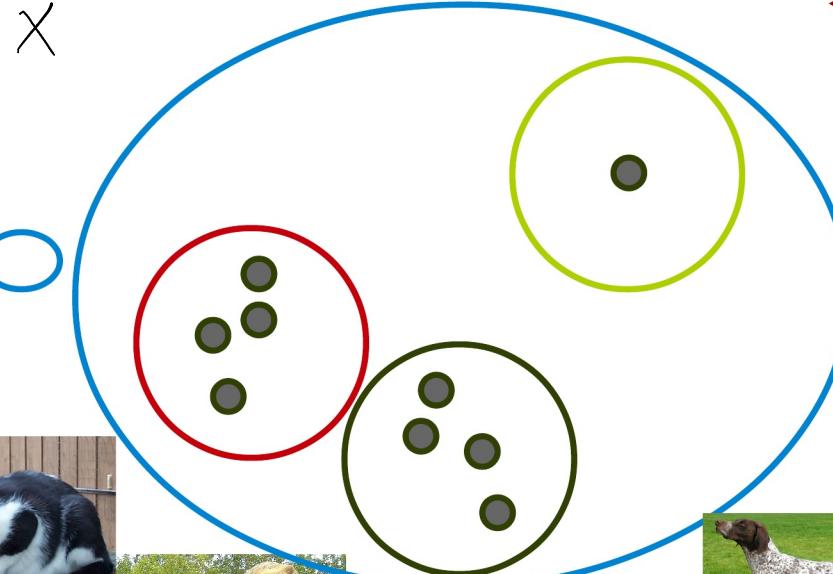
Learning Objectives

- Explain the role of the parameters in the Gaussian Mixture Model (GMM) and how the parameters are updated using the EM-algorithm
- Explain how cross-validation can be used for GMM
- Understand and apply kernel density, K-nearest neighbour density and average relative density estimation for outlier detection

Imagine you observe the world for the first time!

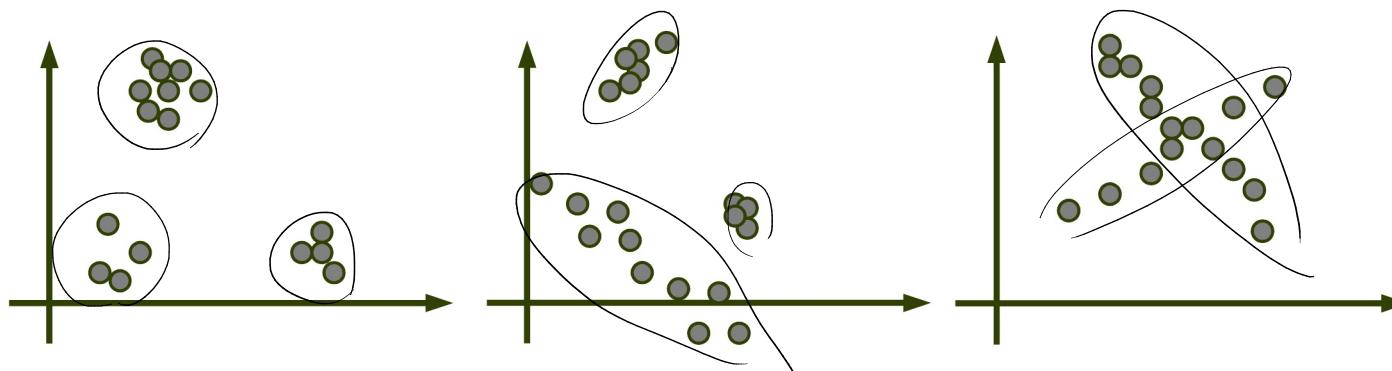


<http://www.clipartlord.com/category/baby-clip-art/>

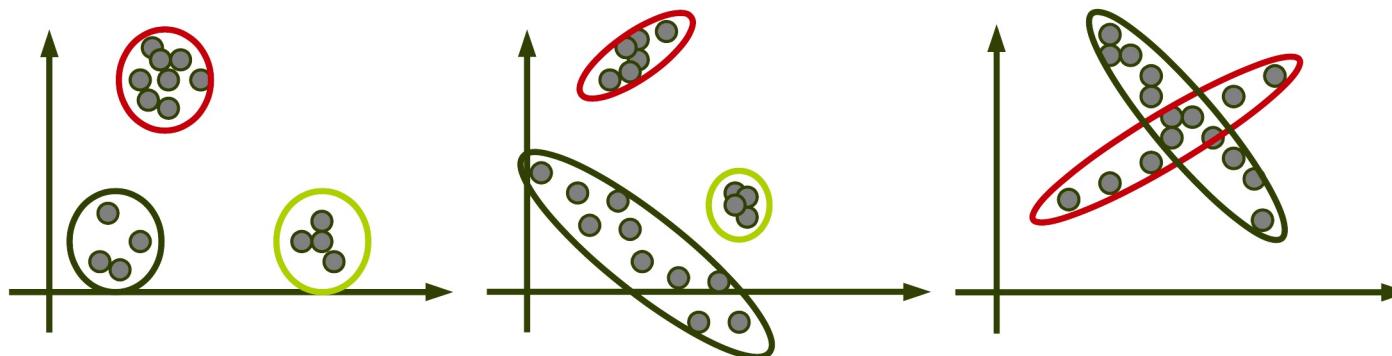


We humans are skilled at dividing objects into groups (clustering), but how do we make computers do the same?

- What are the clusters below and what characterize each cluster?
- Is **k-means** well suited for modeling the clusters below?
 - Will it always find the optimum solution?
 - Can it model the sizes of the clusters?
 - Can it model the shape of the clusters?
 - How can we determine the number of clusters?



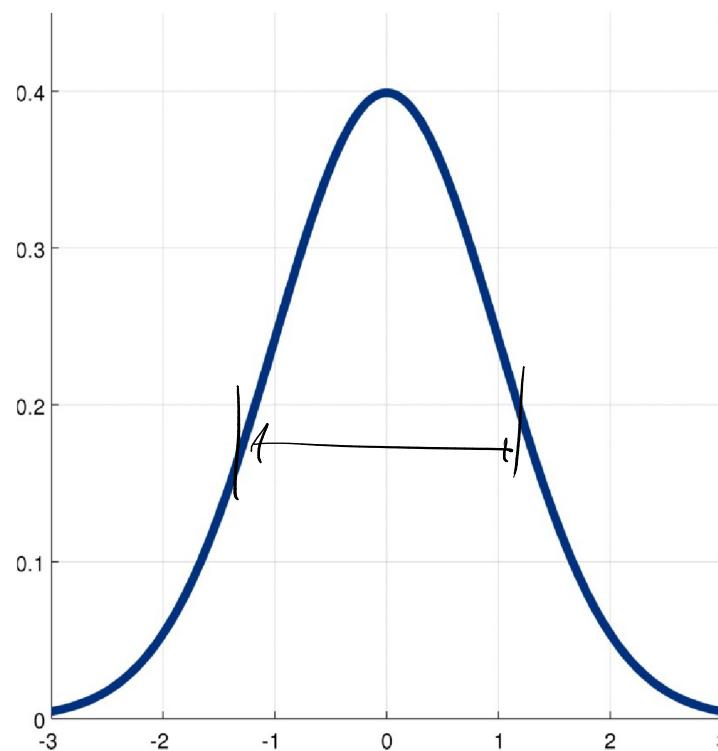
- What are the clusters below and what characterize each cluster?
- Is **k-means** well suited for modeling the clusters below?
 - Will it always find the optimum solution?
 - Can it model the sizes of the clusters?
 - Can it model the shape of the clusters?
 - How can we determine the number of clusters?



Normal distribution

- Probability density function describes the relative chance of a given value to occur
- Normal distribution characterized by
 - Mean
 - Variance

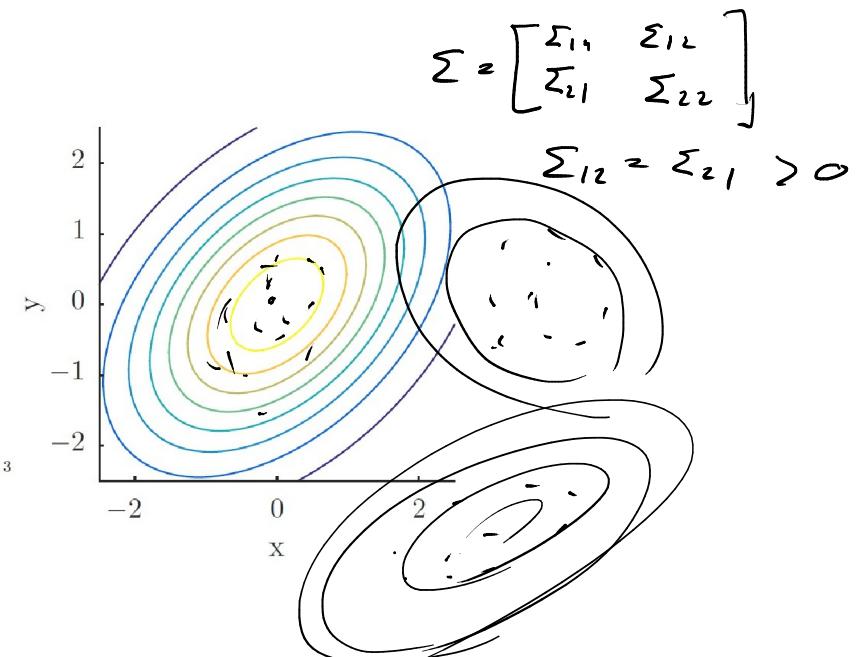
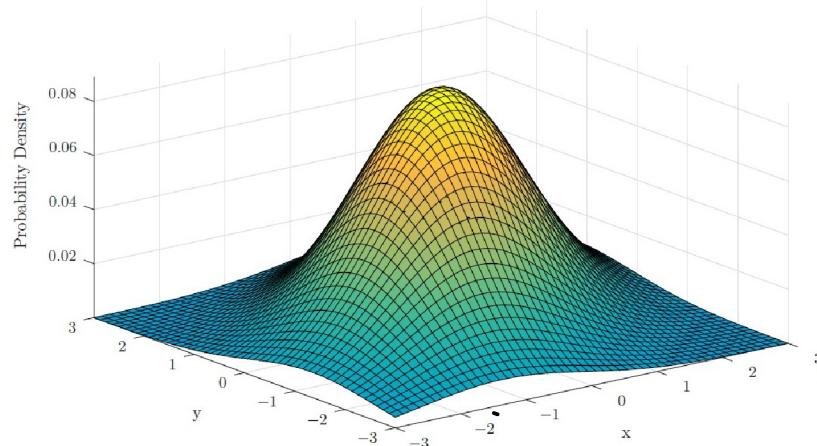
$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$



Multivariate Normal distribution

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

$p(x_1, x_2)$



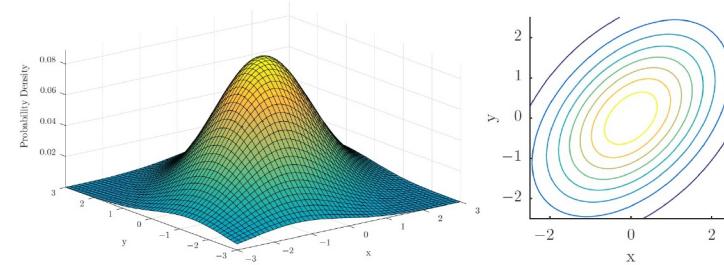
Multivariate Normal distribution

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

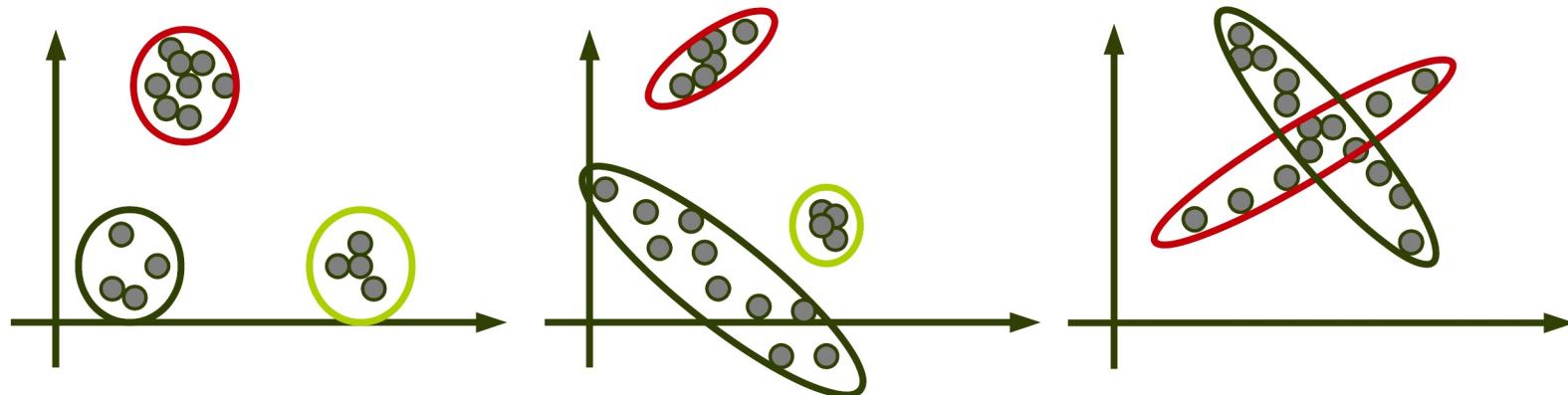
- Example: 2-dimensional Normal distribution

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

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



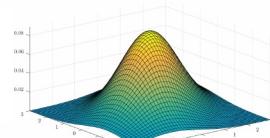
Prototypical mixture model



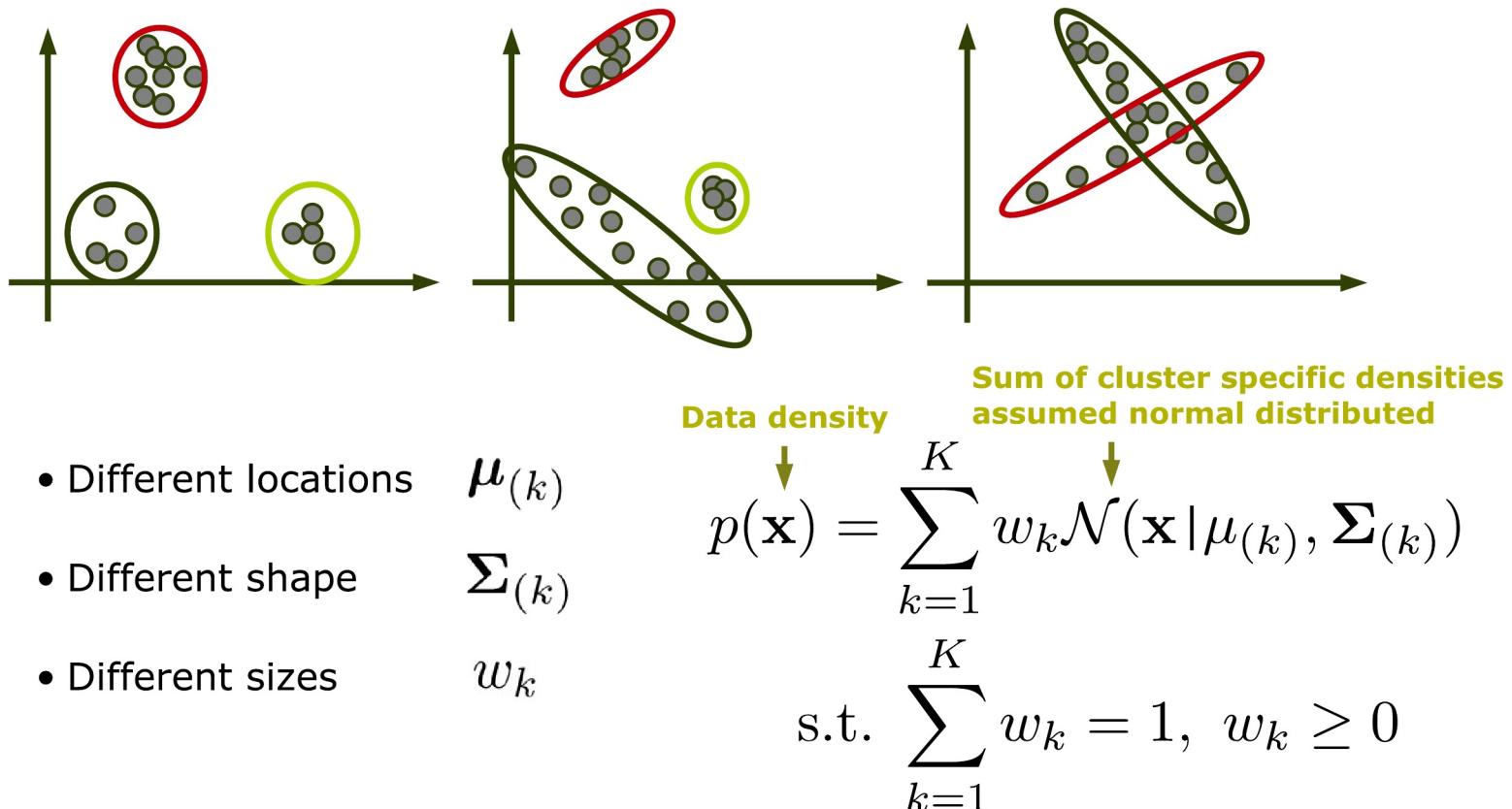
- We want a **density** $p(\mathbf{x})$ of our observations $\mathbf{x} \in \mathbb{R}^M$
- Suppose we have K clusters and let $z = k$ if \mathbf{x} belongs to cluster k
 $z = 1, \dots, K$.
- According to the basic rules of probability: $p(\mathbf{x} | z=k) \cdot p(z=k)$

$$p(\mathbf{x}) = \sum_{k=1}^K p(\mathbf{x}, z=k) = \sum_{k=1}^K p(\mathbf{x}|z=k)p(z=k)$$

- If we specify $p(\mathbf{x}|z=k)$ and $p(z=k) = w_k$ we have a model



The Gaussian Mixture Model (GMM)

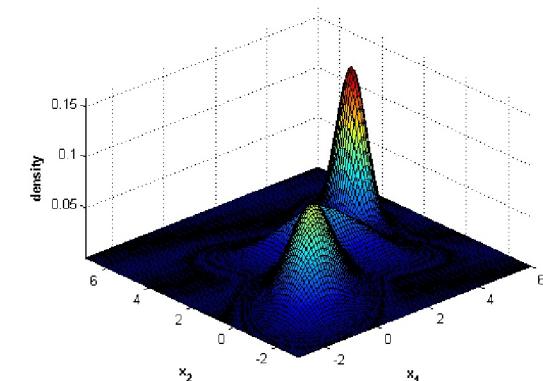
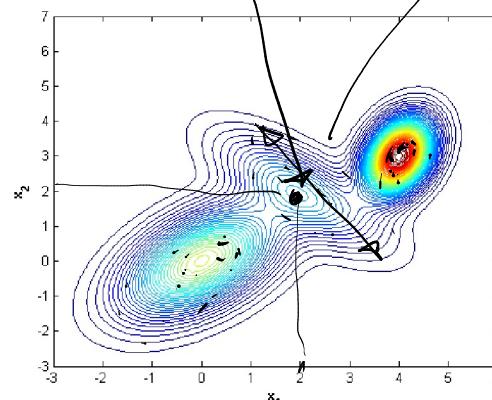
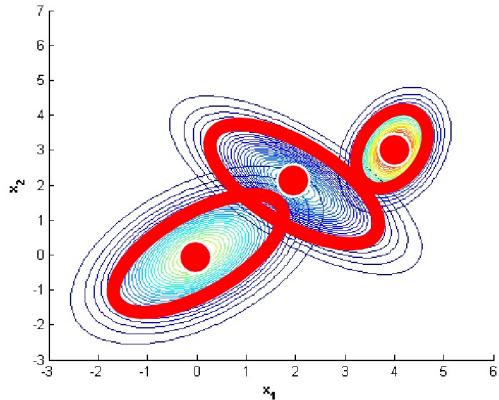


GMM example

$$k=3, \omega_1=0.5, \omega_2=0.2, \omega_3=0.3$$

$$\mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \dots$$

$$p(\mathbf{x}) = 0.5\mathcal{N}(\mathbf{x} | \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}) + 0.2\mathcal{N}(\mathbf{x} | \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 & -0.7 \\ -0.7 & 1 \end{bmatrix}) + 0.3\mathcal{N}(\mathbf{x} | \begin{bmatrix} 4 \\ 3 \end{bmatrix}, \begin{bmatrix} 0.2 & 0.1 \\ 0.1 & 0.5 \end{bmatrix})$$



$\mu_{(k)}$: Cluster center (prototypical example in cluster)

$\Sigma_{(k)}$: Shape of the cluster

w_k : Relative size/density of the cluster

Quiz 01 (please answer on Piazza): GMM

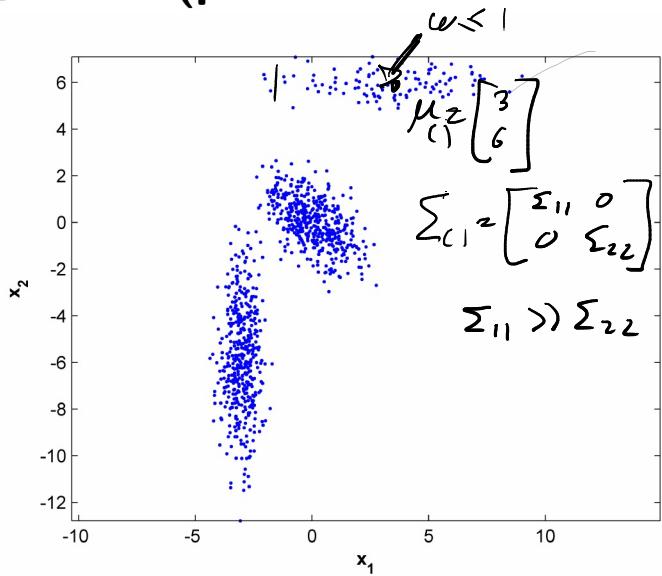


Figure 1: 1000 data observations drawn from a Gaussian Mixture Model (GMM) with three clusters.

In Figure 1 is shown 1000 observations drawn from a density defined by a Gaussian Mixture Model (GMM) with three clusters. Suppose

$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{2\pi|\boldsymbol{\Sigma}|}} \exp(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$ is the multivariate normal distribution, which one of the following GMM densities was used to generate the data?

A $p(x) = 0.1 \cdot \mathcal{N}(\mathbf{x} | \begin{bmatrix} 3 \\ 6 \end{bmatrix}, \begin{bmatrix} 5 & 0 \\ 0 & 0.25 \end{bmatrix}) + 0.45 \cdot \mathcal{N}(\mathbf{x} | \begin{bmatrix} -3 \\ -6 \end{bmatrix}, \begin{bmatrix} 0.25 & 0 \\ 0 & 5 \end{bmatrix}) + 0.45 \cdot \mathcal{N}(\mathbf{x} | \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix})$ ✓

B $p(x) = 0.45 \cdot \mathcal{N}(\mathbf{x} | \begin{bmatrix} 3 \\ 6 \end{bmatrix}, \begin{bmatrix} 5 & 0 \\ 0 & 0.25 \end{bmatrix}) + 0.1 \cdot \mathcal{N}(\mathbf{x} | \begin{bmatrix} -3 \\ -6 \end{bmatrix}, \begin{bmatrix} 0.25 & 0 \\ 0 & 5 \end{bmatrix}) + 0.45 \cdot \mathcal{N}(\mathbf{x} | \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix})$

C $p(x) = 0.1 \cdot \mathcal{N}(\mathbf{x} | \begin{bmatrix} 3 \\ 6 \end{bmatrix}, \begin{bmatrix} 0.25 & 0 \\ 0 & 5 \end{bmatrix}) + 0.45 \cdot \mathcal{N}(\mathbf{x} | \begin{bmatrix} -3 \\ -6 \end{bmatrix}, \begin{bmatrix} 5 & 0 \\ 0 & 0.25 \end{bmatrix}) + 0.45 \cdot \mathcal{N}(\mathbf{x} | \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix})$ ✓

D $p(x) = 0.1 \cdot \mathcal{N}(\mathbf{x} | \begin{bmatrix} 3 \\ 6 \end{bmatrix}, \begin{bmatrix} 0.25 & 0 \\ 0 & 5 \end{bmatrix}) + 0.45 \cdot \mathcal{N}(\mathbf{x} | \begin{bmatrix} -3 \\ -6 \end{bmatrix}, \begin{bmatrix} 5 & 0 \\ 0 & 0.25 \end{bmatrix}) + 0.45 \cdot \mathcal{N}(\mathbf{x} | \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix})$ ✓

E Don't know.

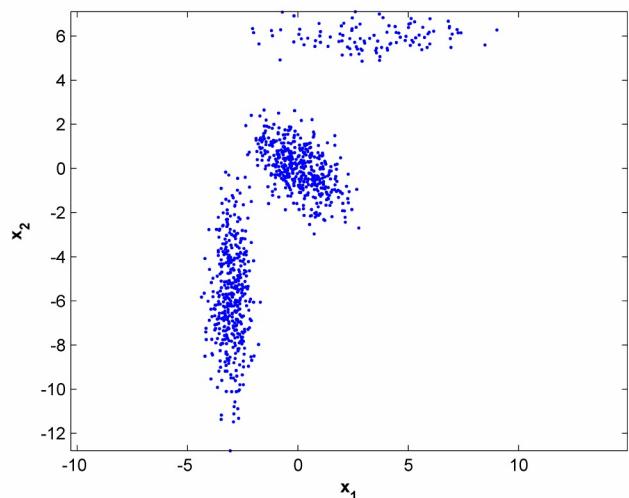
Solution:

Figure 1: 1000 data observations drawn from a Gaussian Mixture Model (GMM) with three clusters.

The centroids of the clusters are $\begin{bmatrix} 3 \\ 6 \end{bmatrix}$, $\begin{bmatrix} -3 \\ -6 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$. The cluster at $\begin{bmatrix} 3 \\ 6 \end{bmatrix}$ is not very dense and should therefore have the coefficient 0.1, it further has a large spread in the x_1 direction and small spread in the x_2 direction corresponding to a covariance matrix of $\begin{bmatrix} 5 & 0 \\ 0 & 0.25 \end{bmatrix}$, thus, answer option one is the only correct answer.

Sanity check time:

- Consider the Gaussian mixture model (GMM)

$$p(\mathbf{x}) = \sum_{k=1}^K w_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{(k)}, \boldsymbol{\Sigma}_{(k)}) \quad \text{s.t. } \sum_{k=1}^K w_k = 1, \quad w_k \geq 0$$

- What is the value of the integral?

$$\begin{aligned} \int p(\mathbf{x}) d\mathbf{x} &= \int \sum_{k=1}^K w_k N(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) d\mathbf{x} = \sum_k w_k \int N(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) d\mathbf{x} \\ &= \sum_k w_k = 1 \end{aligned}$$

Gaussian mixture models, EM algorithm

Select an initial set of model parameters
(mean and covariance for each cluster)

Repeat

- **Expectation**

- For each object, calculate the probability of belonging to each distribution

- **Maximization**

- For each probability distribution, estimate parameters by maximum likelihood

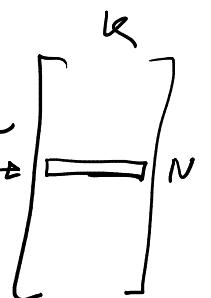
Until the parameters do not change

E-step

$$p(z_n = k|x_n) \leftarrow \frac{w_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{(k)}, \boldsymbol{\Sigma}_{(k)})}{\sum_{K=1}^K w_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{(k)}, \boldsymbol{\Sigma}_{(k)})}$$

$n = 1 \dots N$: observation nr.
 x_n : observation n .
 z_n : which cluster x_n is assigned to.

$$\frac{p(z_n = k|x_n)}{p(z_n = \dots | x_n)} \xrightarrow{n=1 \dots N}$$



M-step

$$\begin{aligned} \underline{N}_k &= \sum_{n=1}^N p(z_n = k|x_n) && \text{Diagram: A vertical bar divided into } k \text{ segments, with the } k \text{-th segment shaded.} \\ \underline{\boldsymbol{\mu}}_{(k)} &= \frac{1}{\underline{N}_k} \sum_{n=1}^N \mathbf{x}_n p(z_n = k|x_n) \\ \underline{w}_k &= \frac{\underline{N}_k}{N} \\ \underline{\boldsymbol{\Sigma}}_{(k)} &= \frac{1}{\underline{N}_k} \sum_{n=1}^N (\mathbf{x}_n - \underline{\boldsymbol{\mu}}_{(k)}) (\mathbf{x}_n - \underline{\boldsymbol{\mu}}_{(k)})^\top p(z_n = k|x_n) \end{aligned}$$

$k = 1 \dots K$.

The GMM update rules approximately implements

$$w_k \mathcal{N}(x | \mu_k, \Sigma_k) = \frac{p(z=k) p(x|z=k)}{p(z=k, x)}$$

E-step

$$p(z_n = k|x_n) \leftarrow \frac{w_k \mathcal{N}(x|\mu_{(k)}, \Sigma_{(k)})}{\sum_{k=1}^K w_k \mathcal{N}(x|\mu_{(k)}, \Sigma_{(k)})}$$

$\underbrace{\phantom{\sum_{k=1}^K w_k \mathcal{N}(x|\mu_{(k)}, \Sigma_{(k)})}}_{= p(\varphi)}$

M-step

$$\begin{aligned} N_k &= \sum_{n=1}^N p(z_n = k|x_n) \\ \mu_{(k)} &= \frac{1}{N_k} \sum_{n=1}^N x_n p(z_n = k|x_n) \\ w_k &= \frac{N_k}{N} \\ \Sigma_{(k)} &= \frac{1}{N_k} \sum_{n=1}^N (x_n - \mu_{(k)})(x_n - \mu_{(k)})^\top p(z_n = k|x_n) \end{aligned}$$

$$\approx \frac{p(z=k, x)}{p(x)} = p(z=k|x)$$

$$\left\{ \begin{array}{l} w_k = p(z = k) \\ \mu_{(k)} = \mathbb{E}_{p(x|z=k)} [x] \\ \Sigma_{(k)} = \mathbb{E}_{p(x|z=k)} [(x - \bar{x})(x - \bar{x})^\top] \end{array} \right.$$

$$\begin{aligned} w_k &= \frac{N_k}{N} = \frac{1}{N} \sum_{n=1}^N p(z_n = k|x_n) \\ &\approx \sum_{n=1}^N p(x_n) p(z_n = k|x_n) \\ &\approx \sum_{n=1}^N p(x_n, z_n = k) \\ &\approx \underline{p(z = k)} \end{aligned}$$

The GMM update rules approximately implements

$$w_k = p(z = k)$$

$$\mu_{(k)} = \mathbb{E}_{p(x|z=k)} [x]$$

$$\Sigma_{(k)} = \mathbb{E}_{p(x|z=k)} [(x - \bar{x})(x - \bar{x})^\top]$$

E-step

$$p(z_n = k|x_n) = \frac{w_k \mathcal{N}(x|\mu_{(k)}, \Sigma_{(k)})}{\sum_{k=1}^K w_k \mathcal{N}(x|\mu_{(k)}, \Sigma_{(k)})}$$

M-step

	$N_k = \sum_{n=1}^N p(z_n = k x_n)$
→	$\mu_{(k)} = \frac{1}{N_k} \sum_{n=1}^N x_n p(z_n = k x_n)$
→	$w_k = \frac{N_k}{N}$
→	$\Sigma_{(k)} = \frac{1}{N_k} \sum_{n=1}^N (x_n - \mu_{(k)})(x_n - \mu_{(k)})^\top p(z_n = k x_n)$

$$\mu_{(k)} = \frac{1}{N_k} \sum_{n=1}^N x_n p(z_n = k|x_n)$$

$$\begin{aligned}
 &= \frac{1}{N_k} \frac{1}{N} \sum_{n=1}^N x_n p(z_n = k|x_n) \\
 &= \frac{1}{\sum_{n=1}^N p(x_n)} \sum_{n=1}^N x_n p(z_n = k|x_n) \\
 &\geq \frac{\sum_{n=1}^N x_n}{\sum_{n=1}^N} \frac{p(z_n = k|x_n)p(x_n)}{p(z_n = k)} \\
 &= \sum_{n=1}^N x_n p(x_n | z_n = k) = \mathbb{E}_{p(x|z_n=k)} [x]
 \end{aligned}$$

The GMM update rules approximately implements

$$\begin{aligned} w_k &= p(z = k) \\ \mu_{(k)} &= \mathbb{E}_{p(\mathbf{x}|z=k)} [\mathbf{x}] \\ \Sigma_{(k)} &= \mathbb{E}_{p(\mathbf{x}|z=k)} [(\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^\top] \end{aligned}$$

Derivation:

$$\begin{aligned} w_k &= \frac{N_k}{N} \\ &= \sum_{n=1}^N p(z_n = k | \mathbf{x}_n) \frac{1}{N} \\ &= \sum_{n=1}^N p(z_n = k | \mathbf{x}_n) p(\mathbf{x}_n) \\ &= \sum_{n=1}^N p(z_n = k, \mathbf{x}_n) \\ &= p(z = k) \end{aligned}$$

The GMM update rules approximately implements

$$w_k = p(z = k)$$

$$\boldsymbol{\mu}_{(k)} = \mathbb{E}_{p(\mathbf{x}|z=k)} [\mathbf{x}]$$

$$\boldsymbol{\Sigma}_{(k)} = \mathbb{E}_{p(\mathbf{x}|z=k)} [(\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^\top]$$

$$\begin{aligned}\boldsymbol{\mu}_{(k)} &= \frac{1}{N_k} \sum_{n=1}^N \mathbf{x}_n p(z_n = k | \mathbf{x}_n) \\ &= \frac{1}{\frac{N_k}{N}} \sum_{n=1}^N \mathbf{x}_n p(z_n = k | \mathbf{x}_n) \frac{1}{N} \\ &= \sum_{n=1}^N \mathbf{x}_n \frac{p(z = k)p(\mathbf{x}_n)}{p(z = k)} \\ &= \sum_{n=1}^N \mathbf{x}_n \frac{p(z = k, \mathbf{x}_n)}{p(z = k)} \\ &= \sum_{n=1}^N \mathbf{x}_n p(\mathbf{x}_n | z = k)\end{aligned}$$

Quiz 02 (please answer on Piazza): GMM

DTU
14:05.

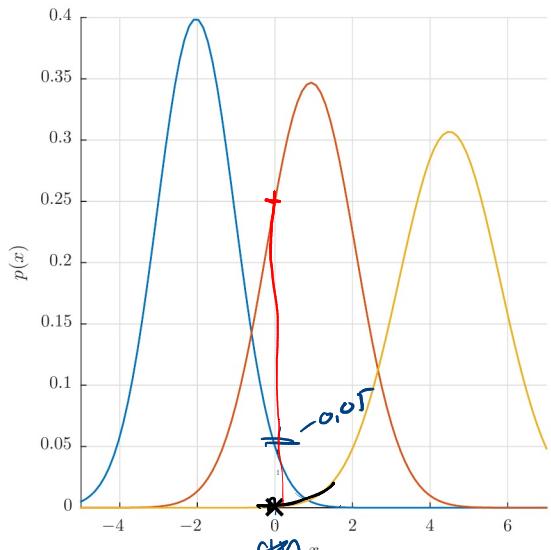


Figure 1: Mixture components in a GMM mixture model with $K = 3$

$$z_n \approx 0,$$

Consider a 1D GMM mixture model where each of the $K = 3$ (Gaussian) mixture components are illustrated in Figure 1 as the colored curves and the figure also shows a new observation indicated by the cross. Suppose we wish to apply the EM algorithm to this mixture model beginning with the E-step (i.e. assuming the mixture components has the means and variances indicated by Figure 1 and equal weights). According to the EM algorithm, what is the (approximate) probability the black cross is assigned to the blue (left-most) mixture component?

- A. 0.05
- B. 0.17
- C. 0.25
- D. 0.02
- E. Don't know.

$$\begin{aligned}
 p(z_n = \text{"blue"} | X_n) &\propto \omega_1 \cdot N(x_n | z_n = \text{"blue"}, \sigma^2) \\
 p(z_n = \text{"blue"}) &= \frac{p(x_n | z_n = \text{"blue"})}{(p(x_n | z_n = \text{"blue"}) + p(x_n | z_n = \text{"orange"}) + p(x_n | z_n = \text{"yellow"}))} \\
 &= \frac{0.05 \cdot \frac{1}{3}}{0.05 \cdot \frac{1}{3} + 0.25 \cdot \frac{1}{3} + 0} = \frac{1}{1+5} \\
 &= \frac{1}{6}.
 \end{aligned}$$

Solution:

The probability of the black cross under each of the three mixture components can be read off as approximately $p(x_0|\mu_1, \sigma_1) \approx 0.05$, $p(x_0|\mu_2, \sigma_2) \approx 0.25$, $p(x_0|\mu_3, \sigma_3) \approx 0$. Since they are weighted equally

the assignment to the left-most component is

$$p(z=1|x_0) = \frac{\frac{1}{3}p(x_0|\mu_1, \sigma_1)}{\sum_{i=1}^3 \frac{1}{3}p(x_0|\mu_i, \sigma_i)} \approx \frac{0.05}{0.05 + 0.25} = 0.17$$

$$\Sigma_k \sim \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix}$$

Gaussian mixture models, EM algorithm

Select an initial set of model parameters
(mean and covariance for each cluster)

Repeat

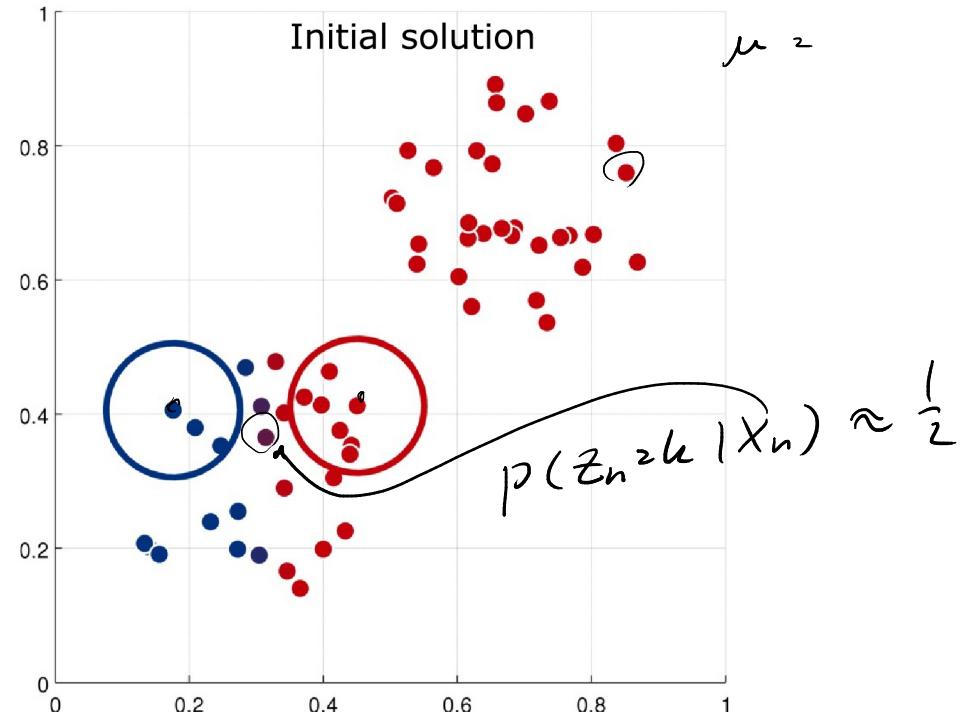
- **Expectation**

- For each object, calculate the probability of belonging to each distribution

- **Maximization**

- For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

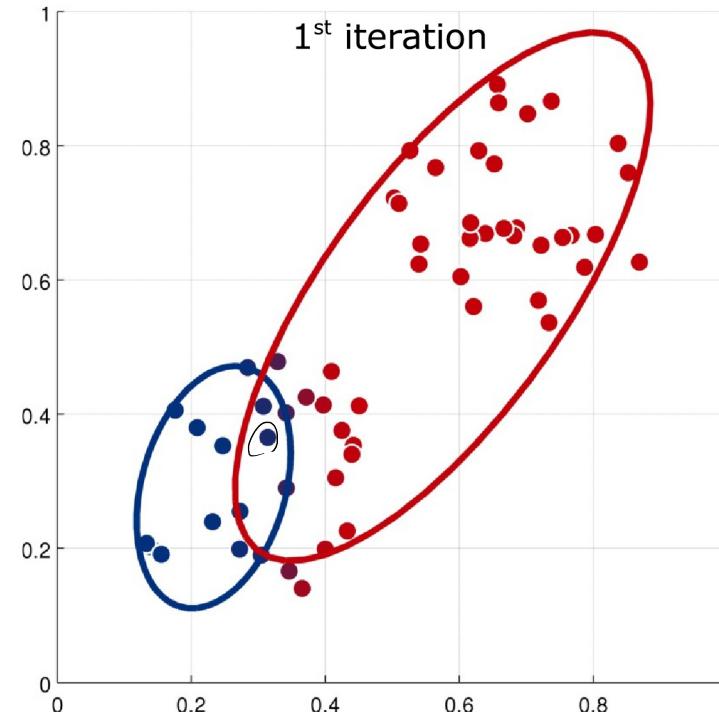
- **Expectation**

- For each object, calculate the probability of belonging to each distribution

- **Maximization**

- For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

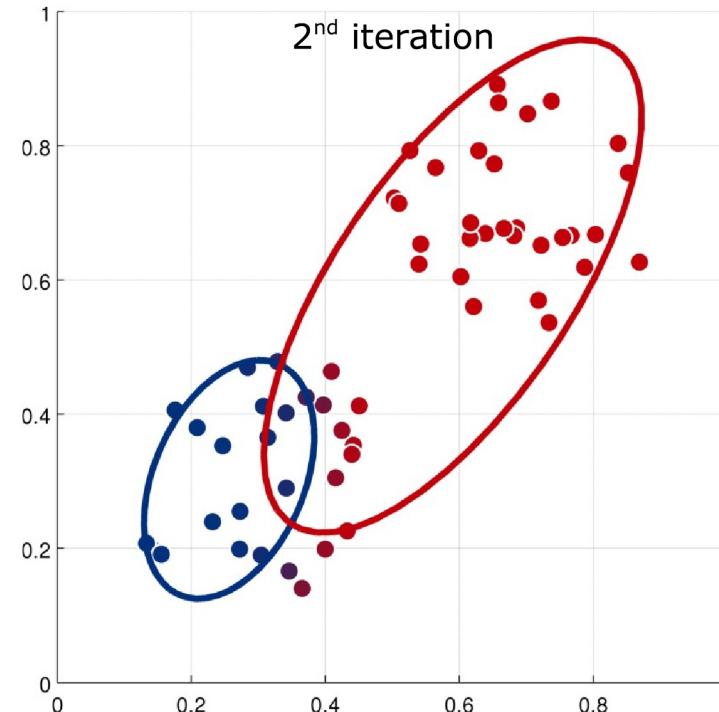
- **Expectation**

- For each object, calculate the probability of belonging to each distribution

- **Maximization**

- For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

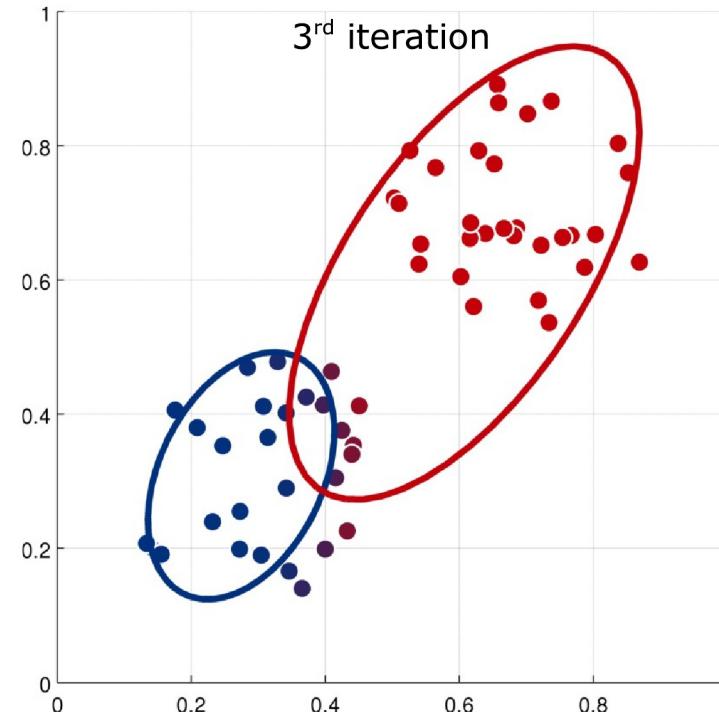
- **Expectation**

- For each object, calculate the probability of belonging to each distribution

- **Maximization**

- For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

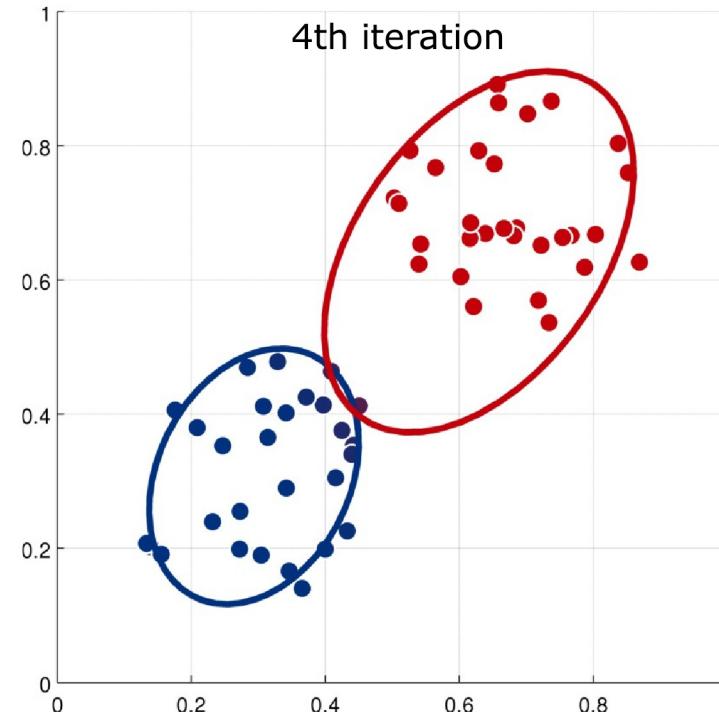
- **Expectation**

- For each object, calculate the probability of belonging to each distribution

- **Maximization**

- For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



Gaussian mixture models, EM algorithm

Select an initial set of model parameters (mean and covariance for each cluster)

Repeat

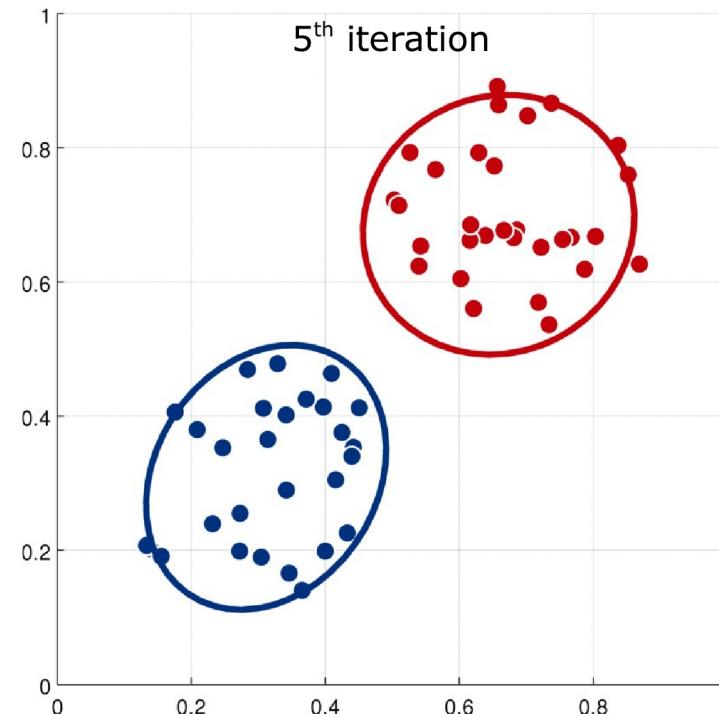
- **Expectation**

- For each object, calculate the probability of belonging to each distribution

- **Maximization**

- For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



Gaussian mixture models, EM algorithm

Select an initial set of model parameters
(mean and covariance for each cluster)

Repeat

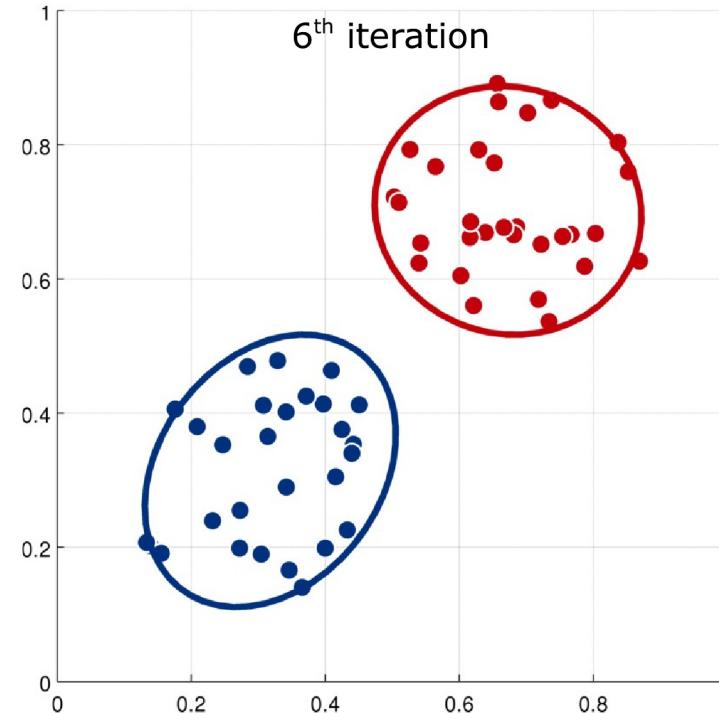
- **Expectation**

- For each object, calculate the probability of belonging to each distribution

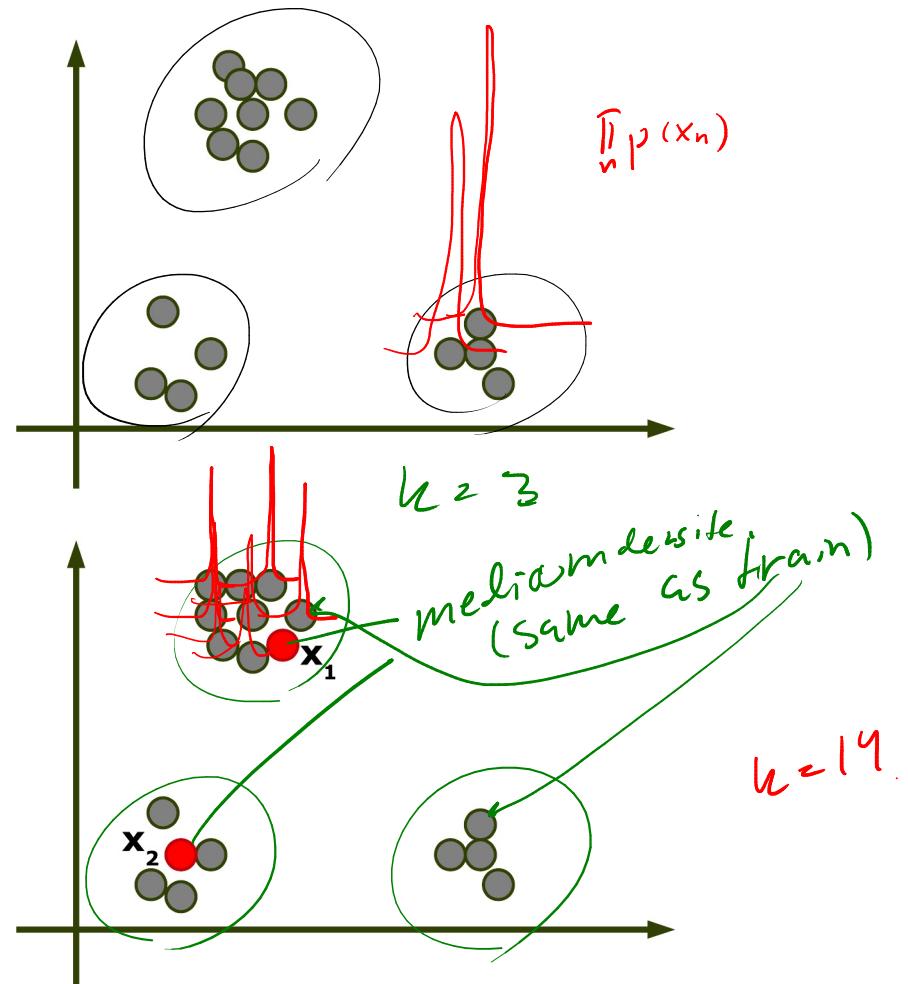
- **Maximization**

- For each probability distribution, estimate parameters by maximum likelihood

Until the parameters do not change



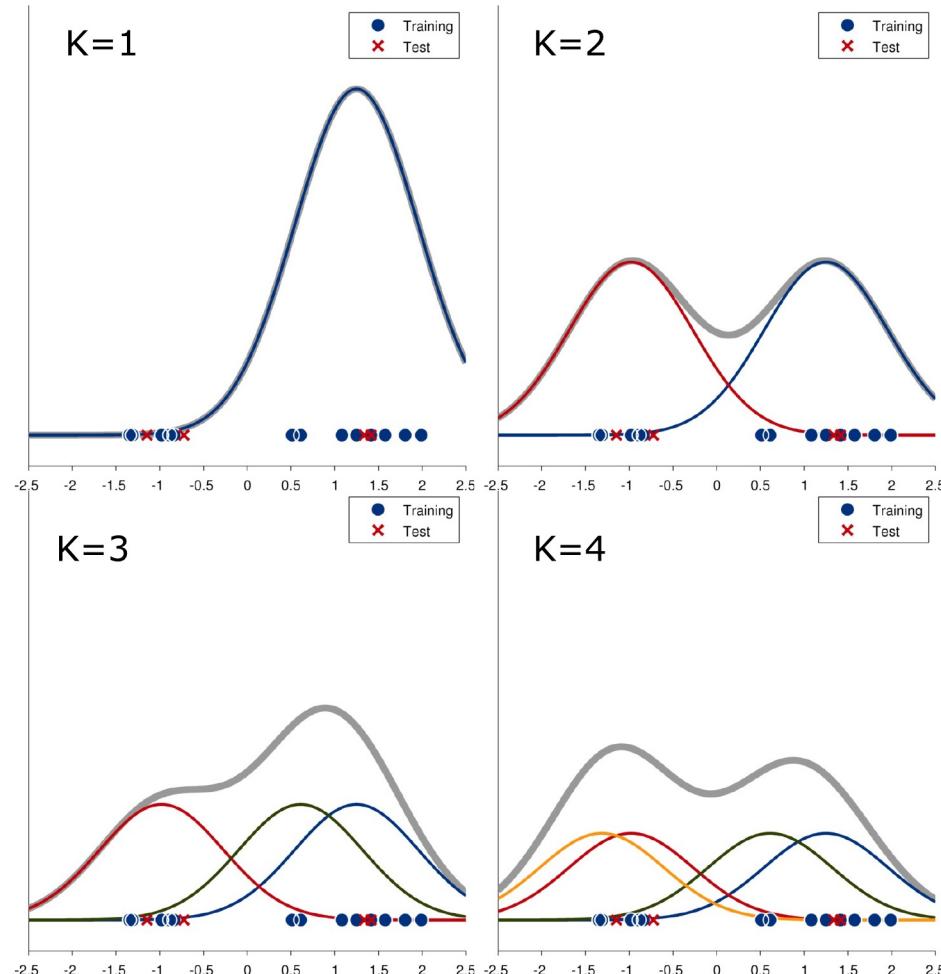
- Consider the data to the right with 16 observations.
 - What would ideally happen if we used a GMM with K=16 clusters to model the data?
- Imagine we have two **test observations** denoted \mathbf{x}_1 and \mathbf{x}_2 (red points) that are not used for training.
 - What happens to $p(\mathbf{x}_1)$ and $p(\mathbf{x}_2)$ if we use K=3 and K=16 clusters?



EM Initial solution

Mixture models

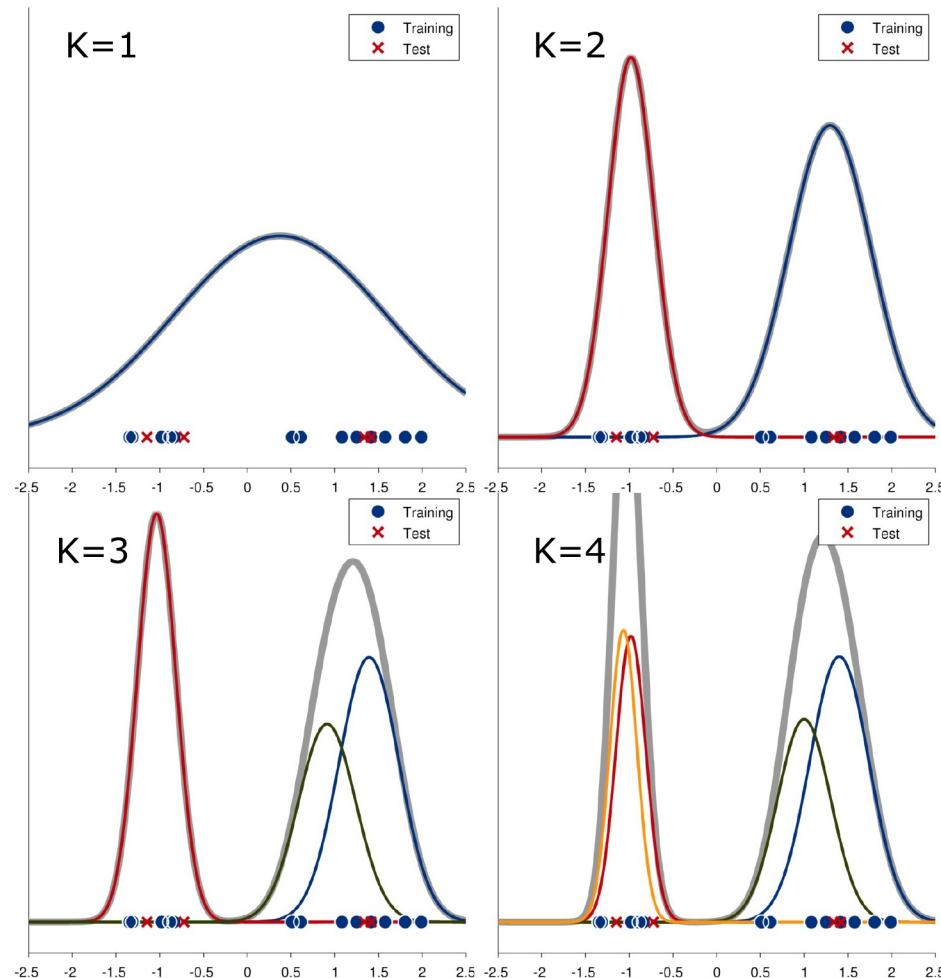
- Selecting complexity using crossvalidation



EM 1st iteration

Mixture models

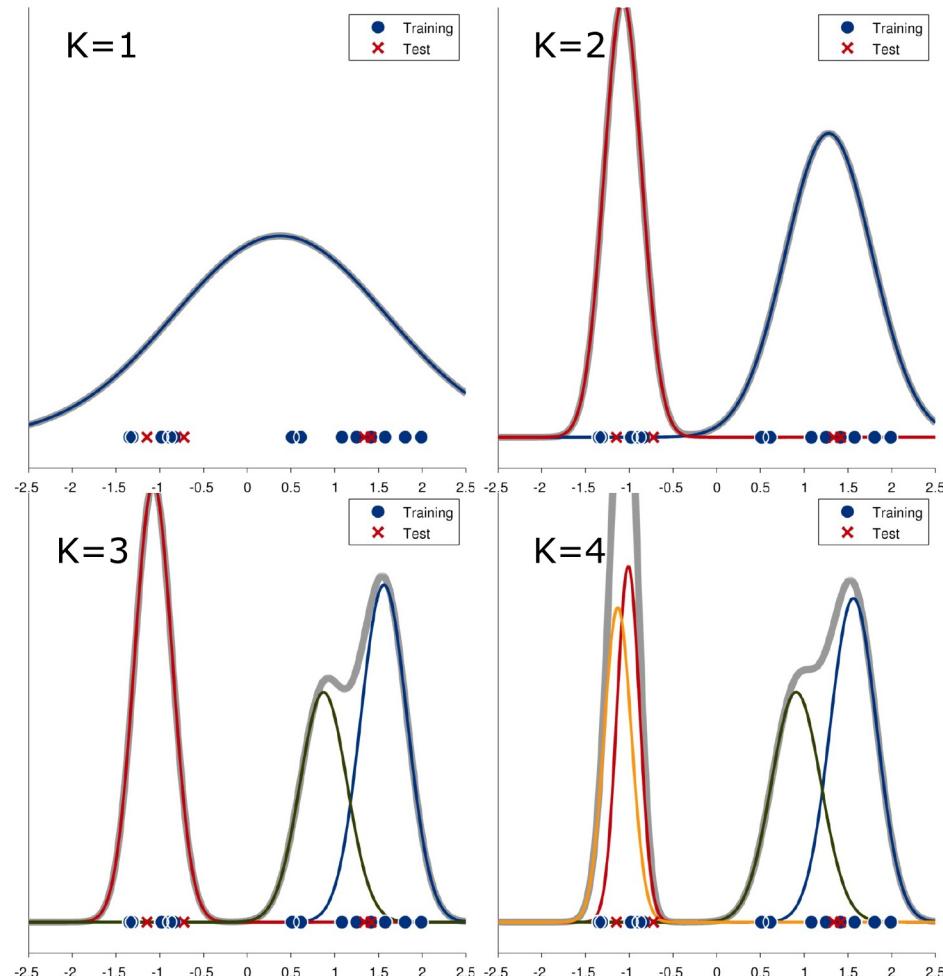
- Selecting complexity using crossvalidation



EM 2nd iteration

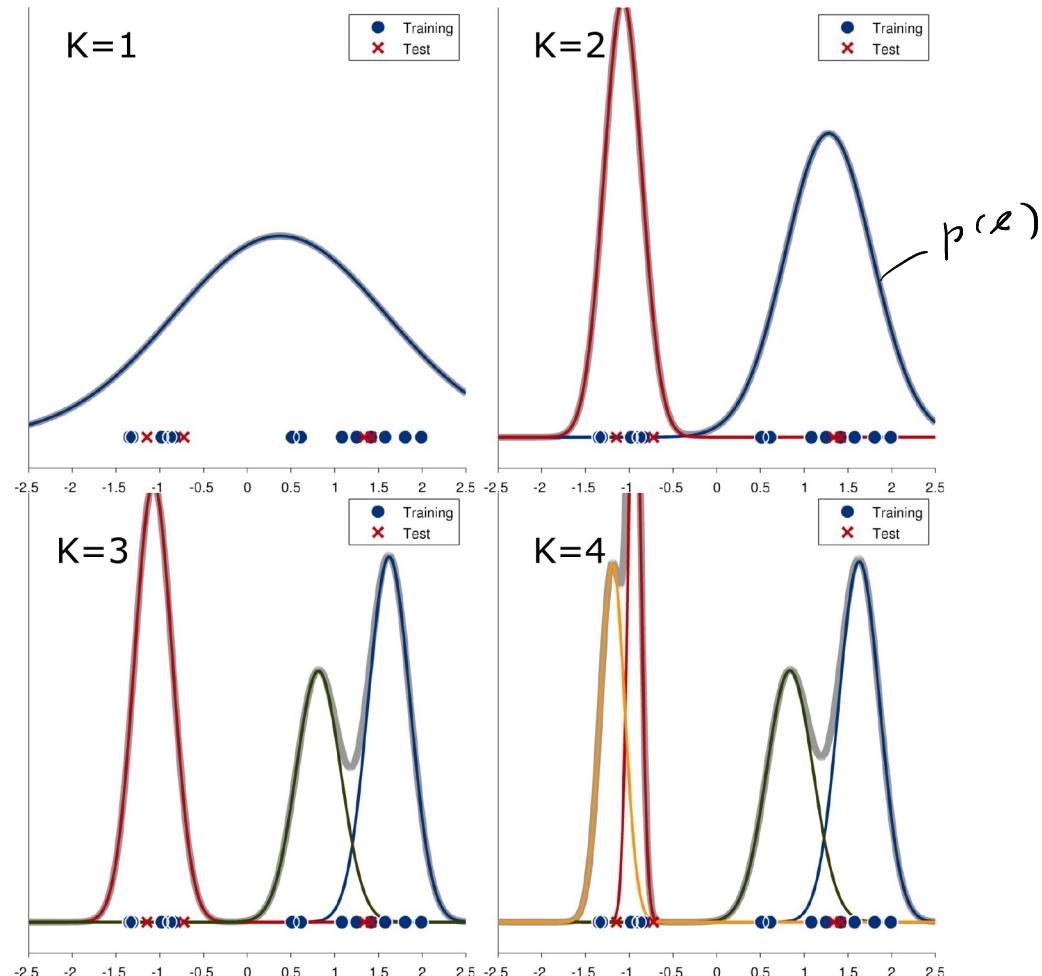
Mixture models

- Selecting complexity using crossvalidation



Mixture models

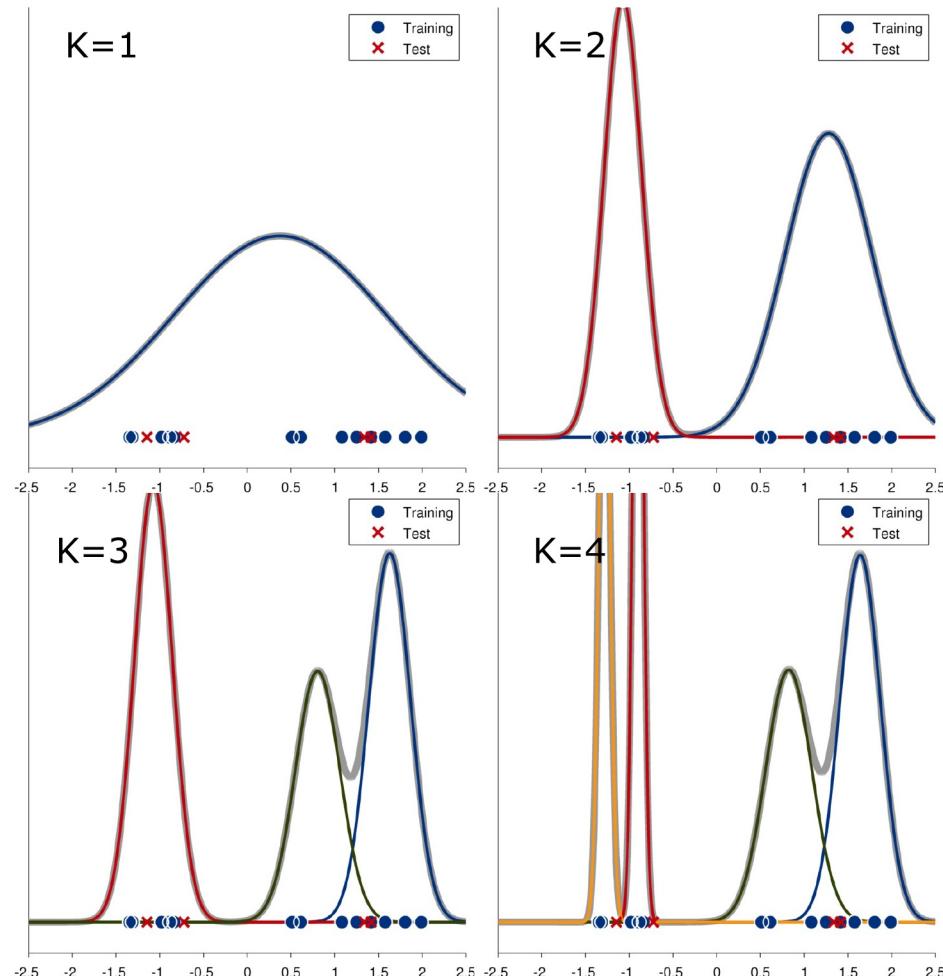
- Selecting complexity using crossvalidation



EM 4th iteration

Mixture models

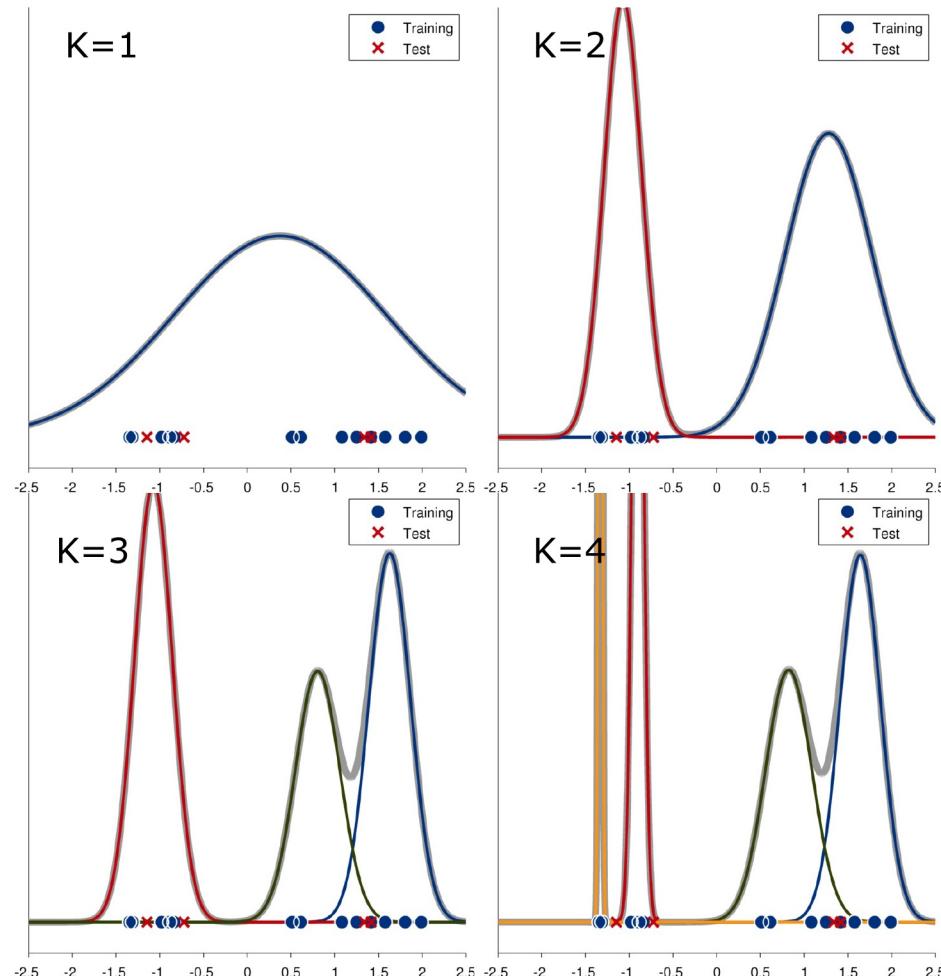
- Selecting complexity using crossvalidation



EM 5th iteration

Mixture models

- Selecting complexity using crossvalidation



Test data evaluation

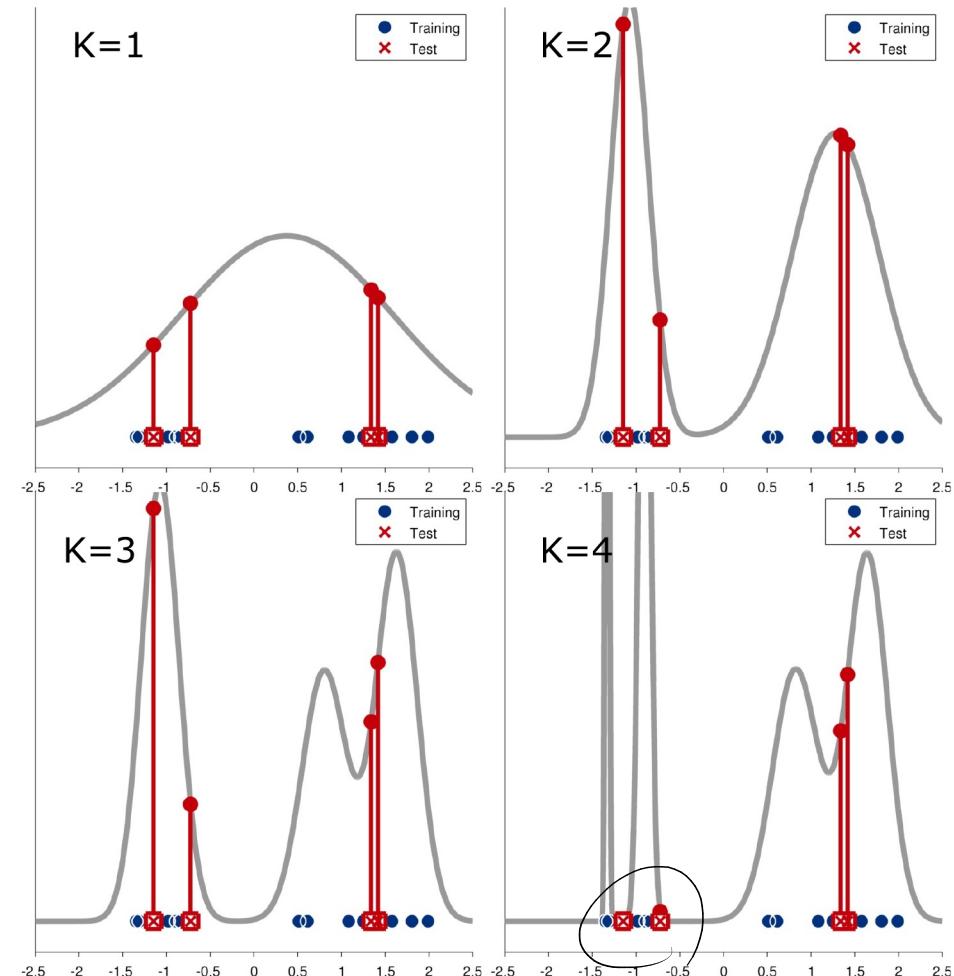
Mixture models

- Selecting complexity using crossvalidation

$$p(x^{test}) = \prod_{i=1}^{N^{test}} p(x_i^{test})$$

$$\log L \approx \log p(x^{test}) = \sum_{i=1}^{N^{test}} \log p(x_i^{test})$$

$$(E) = -(\log L)$$

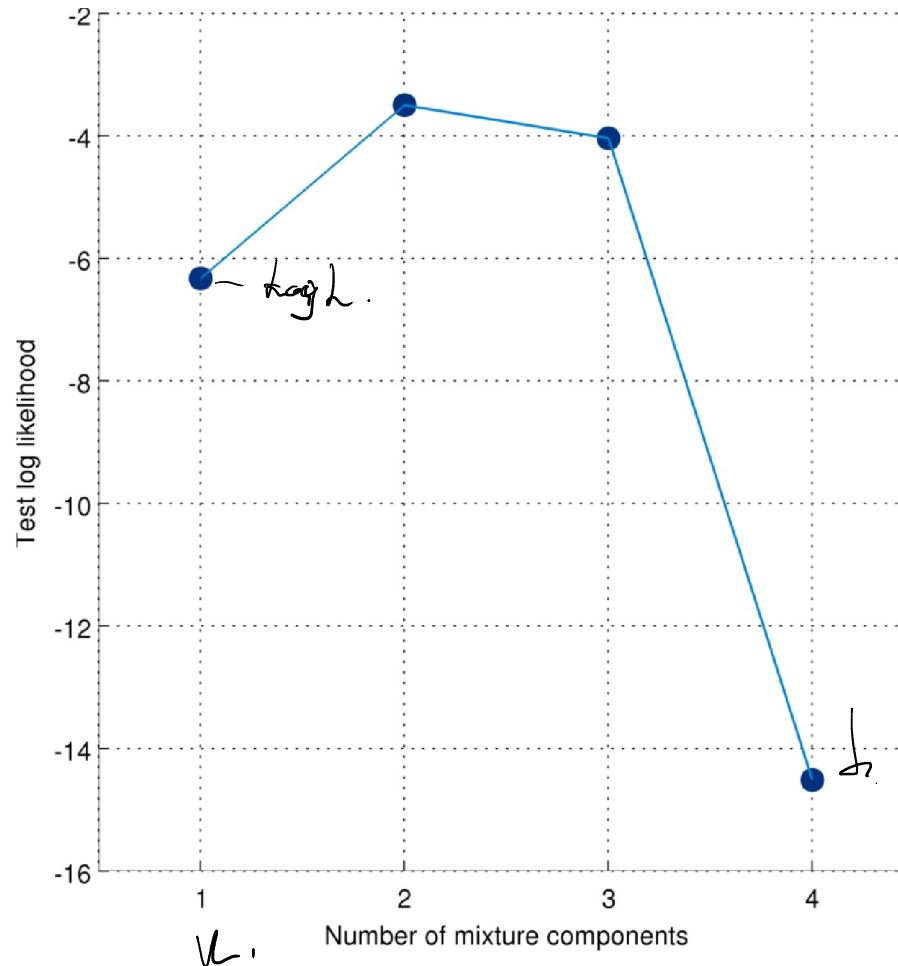


Mixture models

- Selecting complexity using crossvalidation

$$\tilde{\Sigma}_{(n)} + \Sigma_{(n)} + \lambda I$$

regularization.

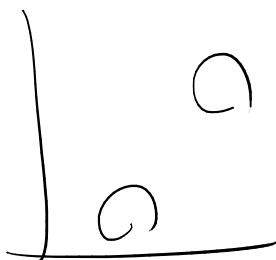


K-means versus GMM

K-means

- No guarantee of optimal solution
- Does not model shape of clusters
- Does not model the size of clusters
- Difficult to assess the number of clusters to use particularly when there is no ground truth

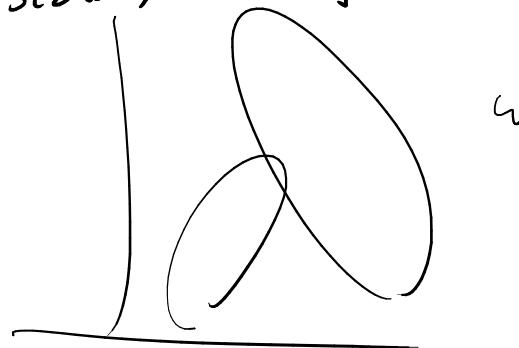
+ fast -



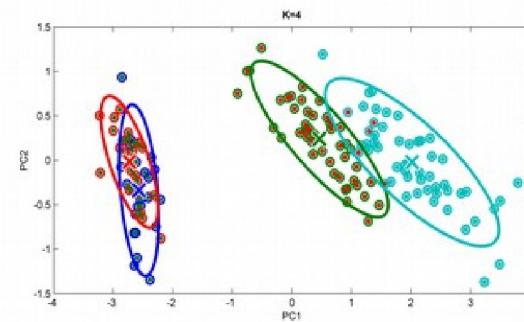
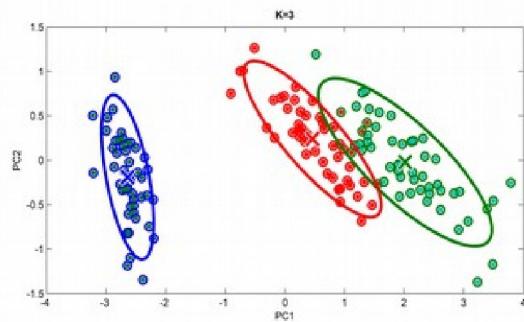
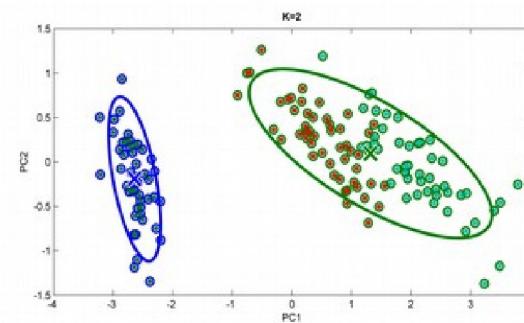
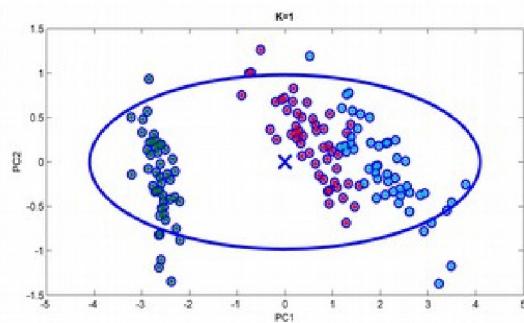
Gaussian mixture model (GMM)

- No guarantee of optimal solution (even more local minima issues due to the additional model parameters)
- Models shape of cluster as ellipsoid ✓
- Models the size of clusters ✓
- Possible to estimate the number of components by cross-validation

- slow, memory intensive.

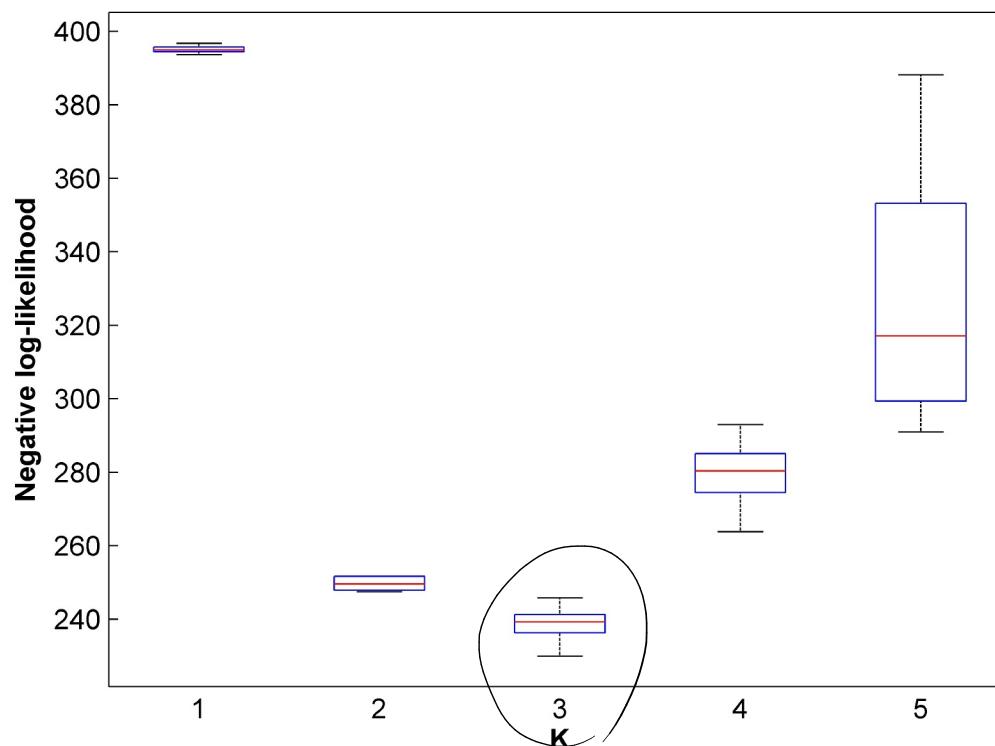


GMM on Iris data using 1,2,3 and 4 components



Recap of GMM on Iris data

GMM 10 fold cross-validation on Iris data repeated five times where the five runs are plotted using box-plots.



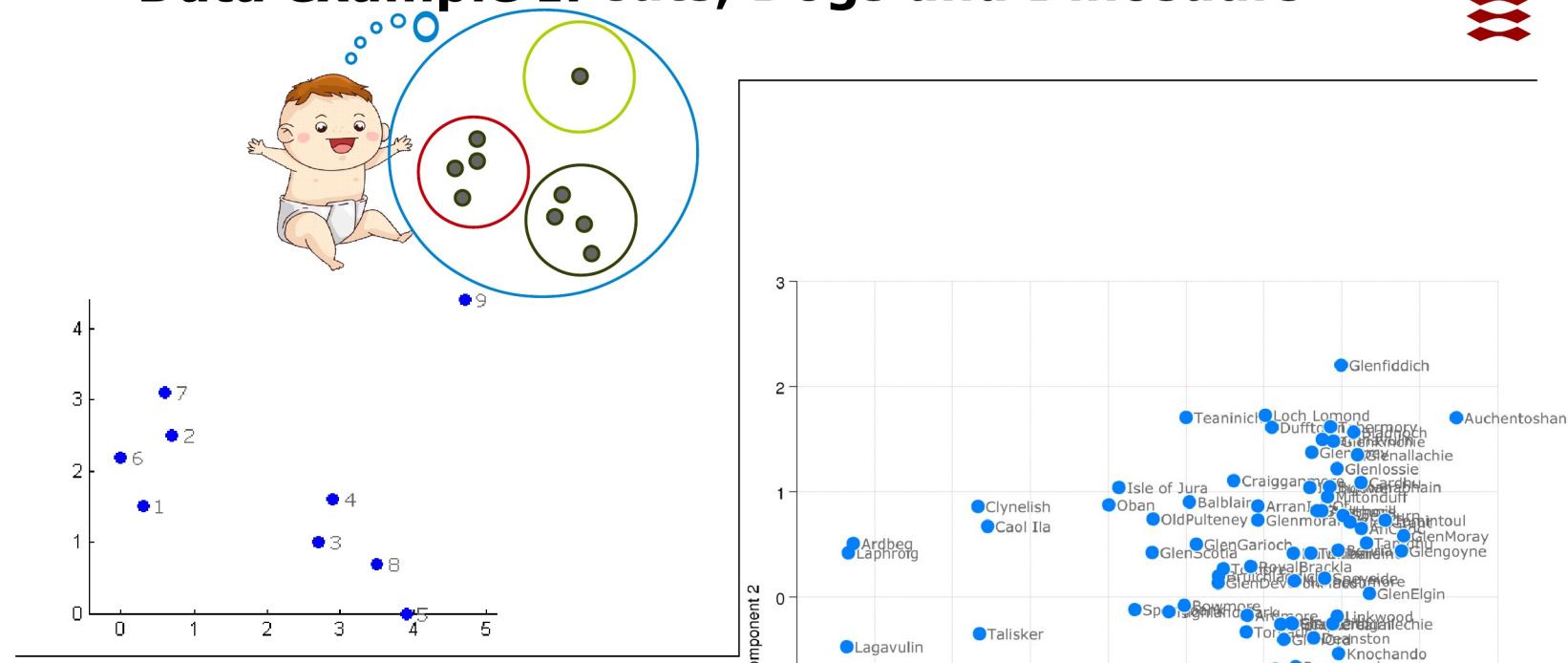
Anomaly detection: Definition

- Given a collection of data objects
 - Each object has associated a number of features
- Detect which objects **deviate from normal** behaviour

Anomaly detection: Example

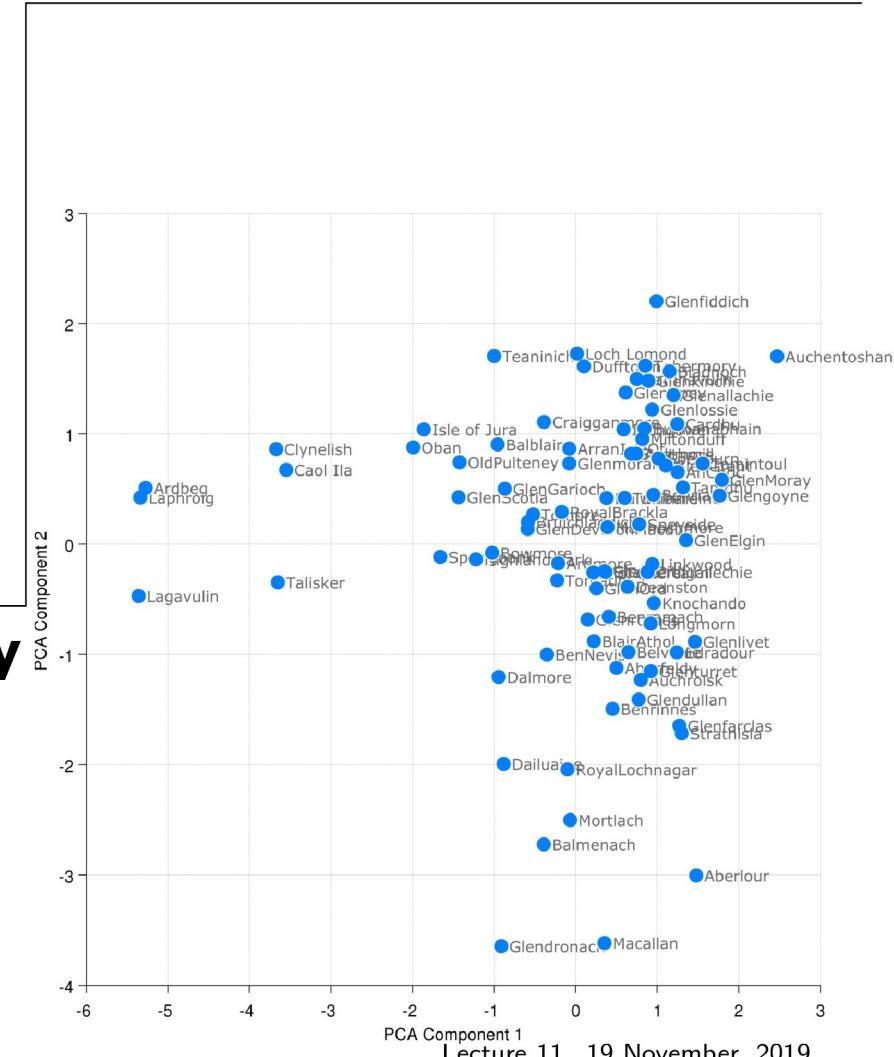
- Credit card **fraud detection**
 - Recognize dubious credit card transactions based on the transaction history of the card holder
- Network **intrusion detection**
 - Detect hacker attacks, web crawlers etc.
- **Ecosystem disturbances**
 - Detect hurricanes, floods droughts, heat waves and fires
- **Health and medicine monitoring**
 - Detect abnormal behaviour in populations and patients
- **Fault detection in industry systems**
 - Detect when a wind turbine performs poorly due to ice coating on blades
- Detection of **outliers** in data measurements
 - Remove erroneous measurements due to misreading from an instrument

Data example I: Cats, Dogs and Dinosaurs

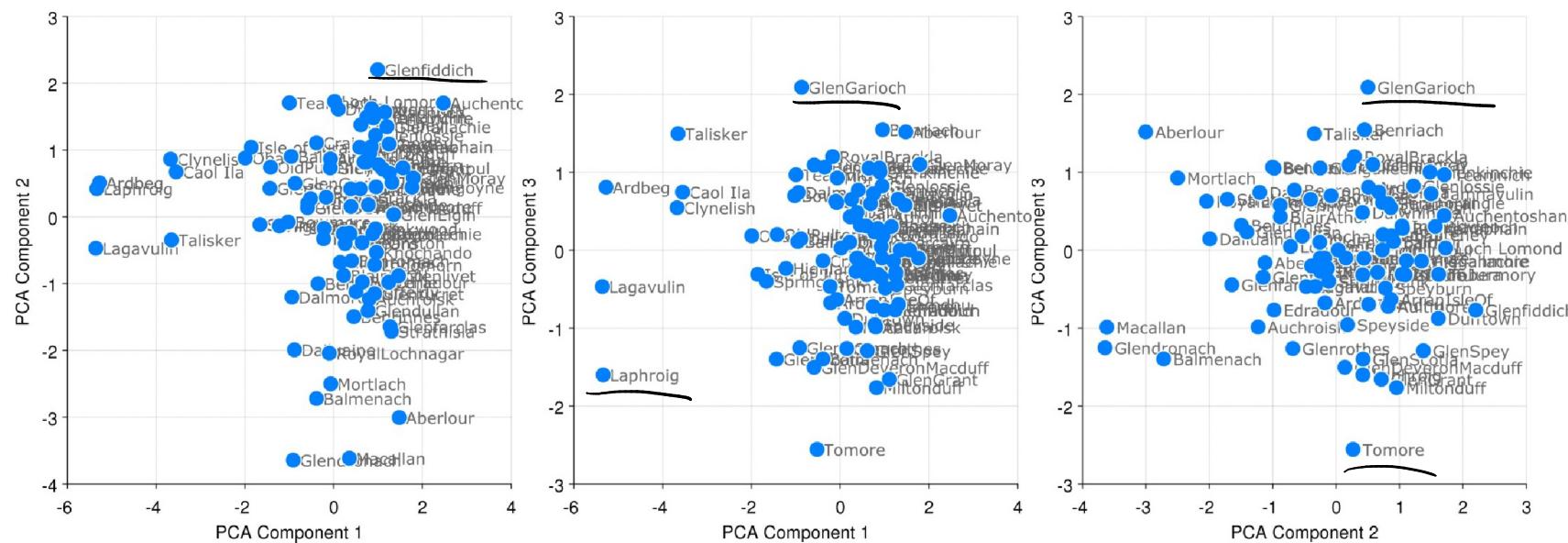


Data example II: Whisky

- 86 types of Scotch whisky
- Human ratings 1-5
- 12 taste categories
 - body, sweetness, smoky, medicinal, tobacco, honey, spicy, winey, nutty, malty, fruity, floral



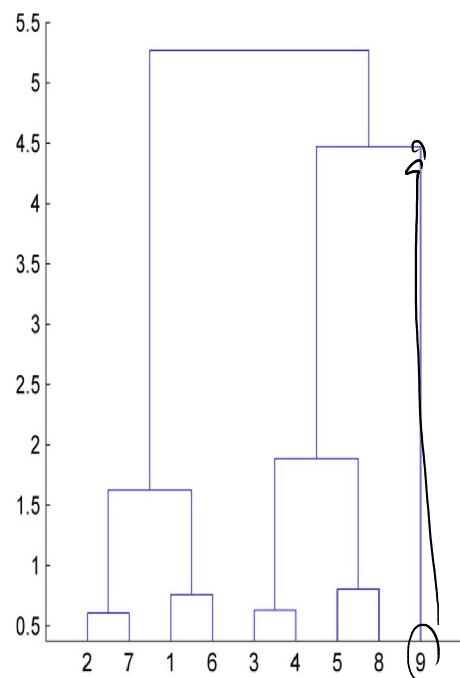
PCA plot



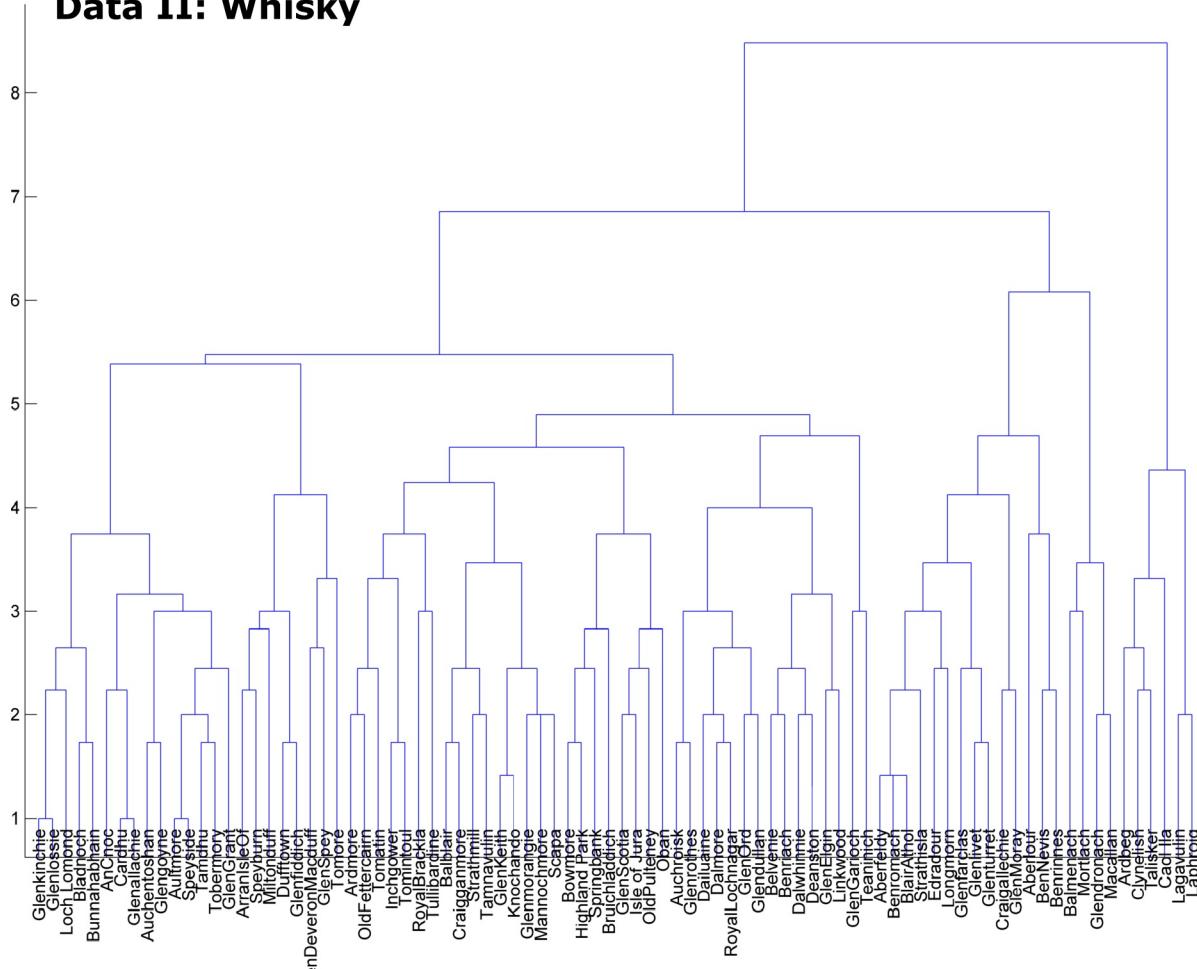
Dendrogram

- Dendograms can be used to visualize relative distances between the observations

Data I: Cats, Dogs and Dinosaurs



Data II: Whisky



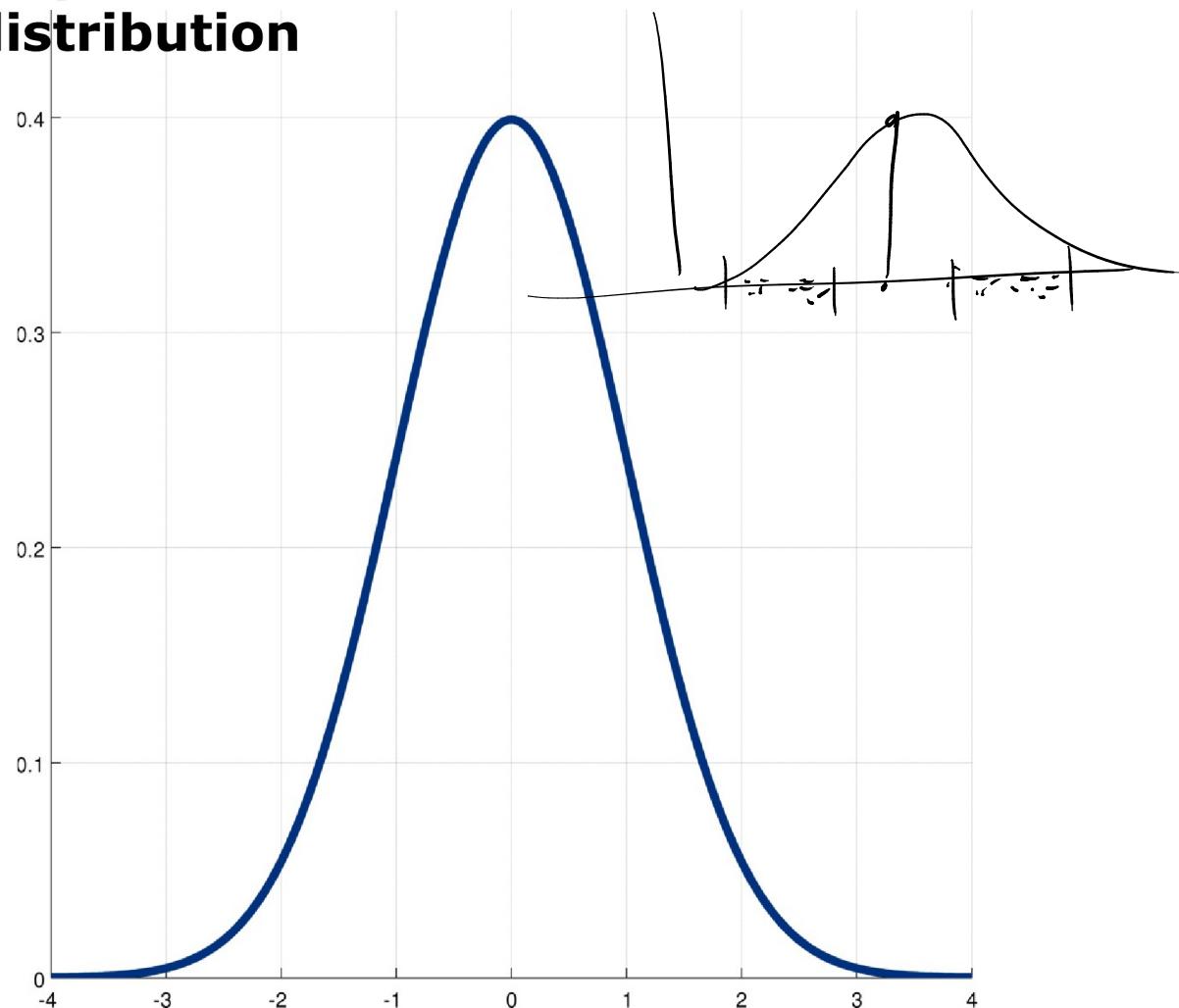
Density based techniques: Univariate normal distribution

- Map attribute to standard Normal variable

$$\hat{z} = \frac{x - \mu}{\sigma}$$

- Choose a threshold

c	p(z >c)
1.0	0.3173
1.5	0.1336
2.0	0.0455
2.5	0.0124
3.0	0.0027
3.5	0.0005
4.0	0.0001



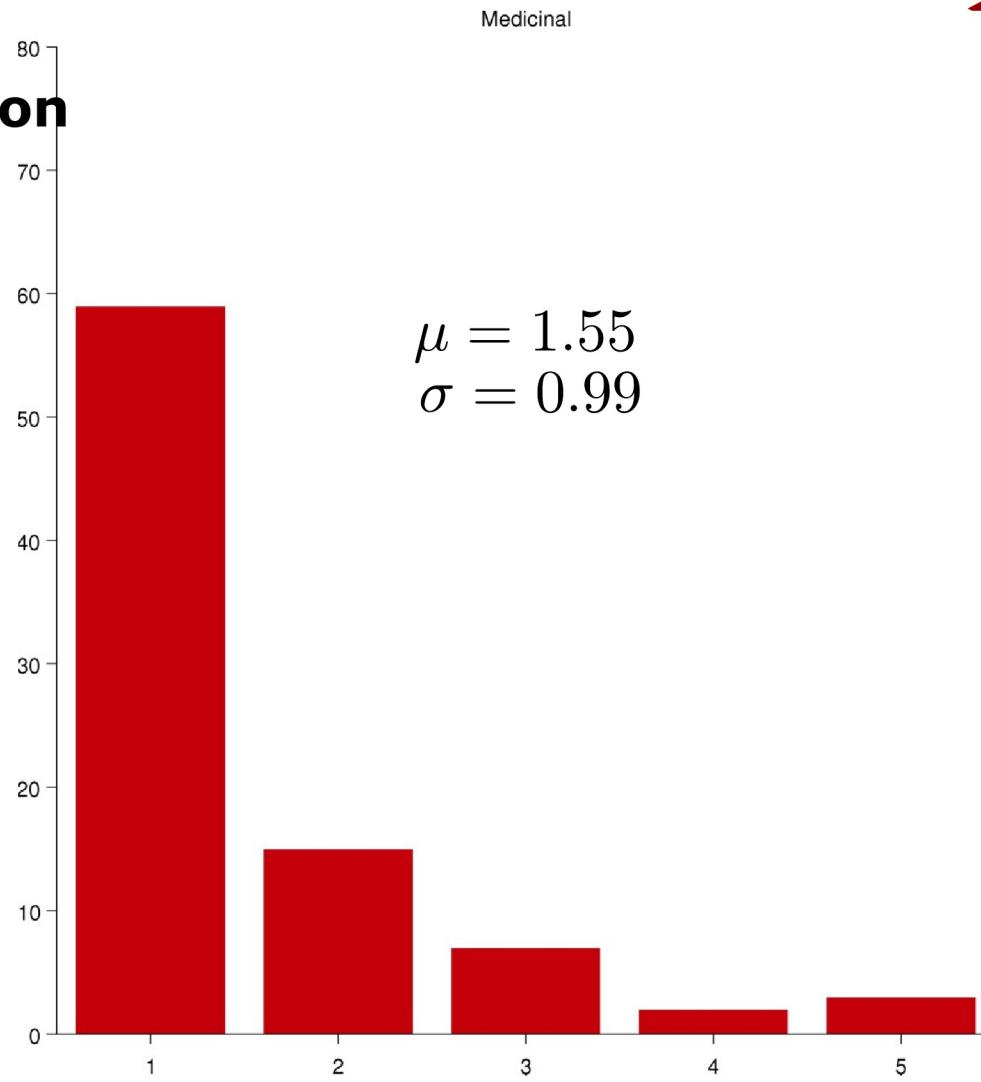
Normal distribution

- Map attribute to standard Normal variable

$$z = \frac{x - \mu}{\sigma}$$

- Choose a threshold

c	p(z >c)
1.0	0.3173
1.5	0.1336
2.0	0.0455
2.5	0.0124
3.0	0.0027
3.5	0.0005
4.0	0.0001



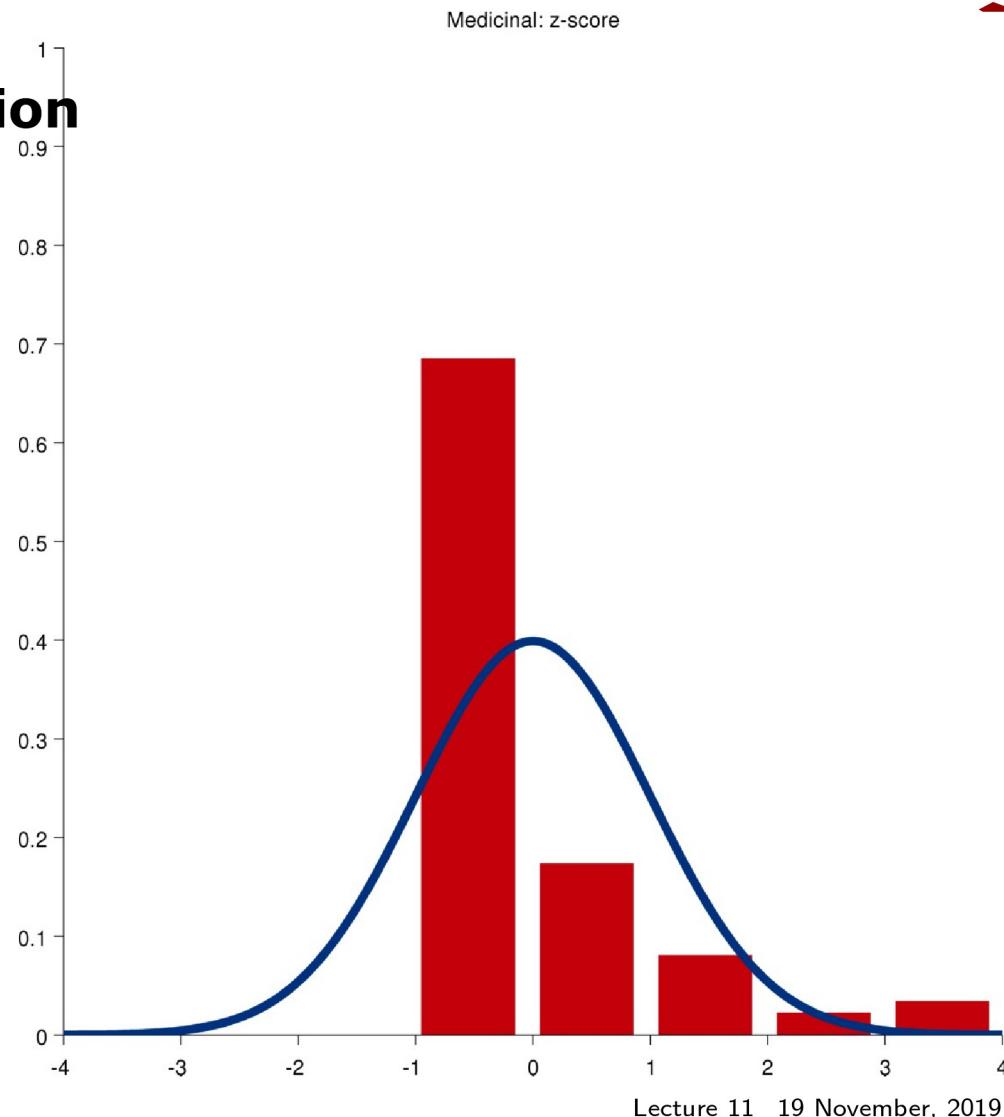
Normal distribution

- Map attribute to standard Normal variable

$$z = \frac{x - \mu}{\sigma}$$

- Choose a threshold

c	p(z >c)
1.0	0.3173
1.5	0.1336
2.0	0.0455
2.5	0.0124
3.0	0.0027
3.5	0.0005
4.0	0.0001



Normal distribution

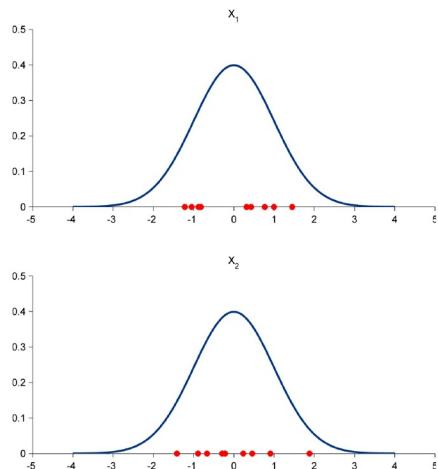
- Map attribute to standard Normal variable
- Choose a threshold

$$z = \frac{x - \mu}{\sigma}$$

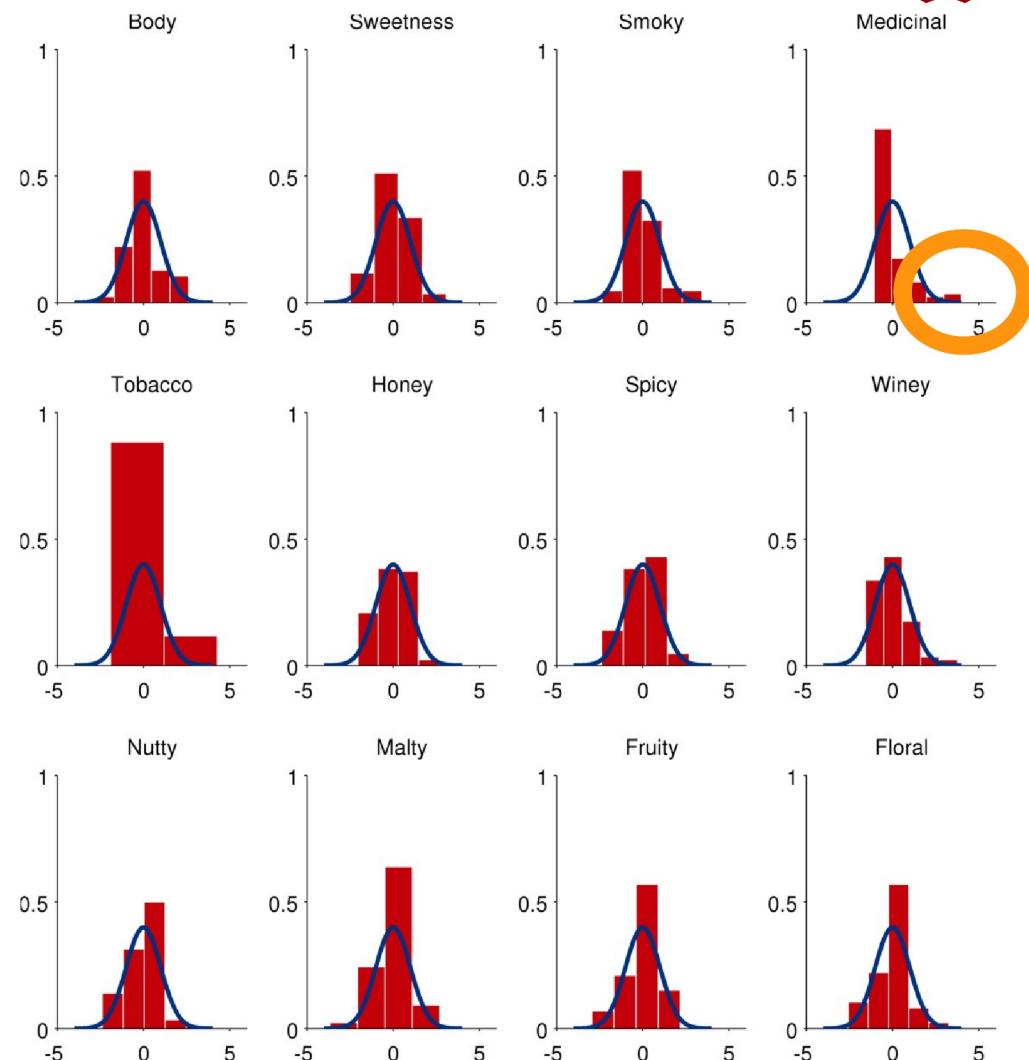
$$p(|z| > c) = 0.001$$

$$c = 3.2905$$

Data I: Cats, Dogs and Dinosaurs



Data II: Whisky



Note: Assumes attributes follow a normal distribution which may not be a valid assumption!

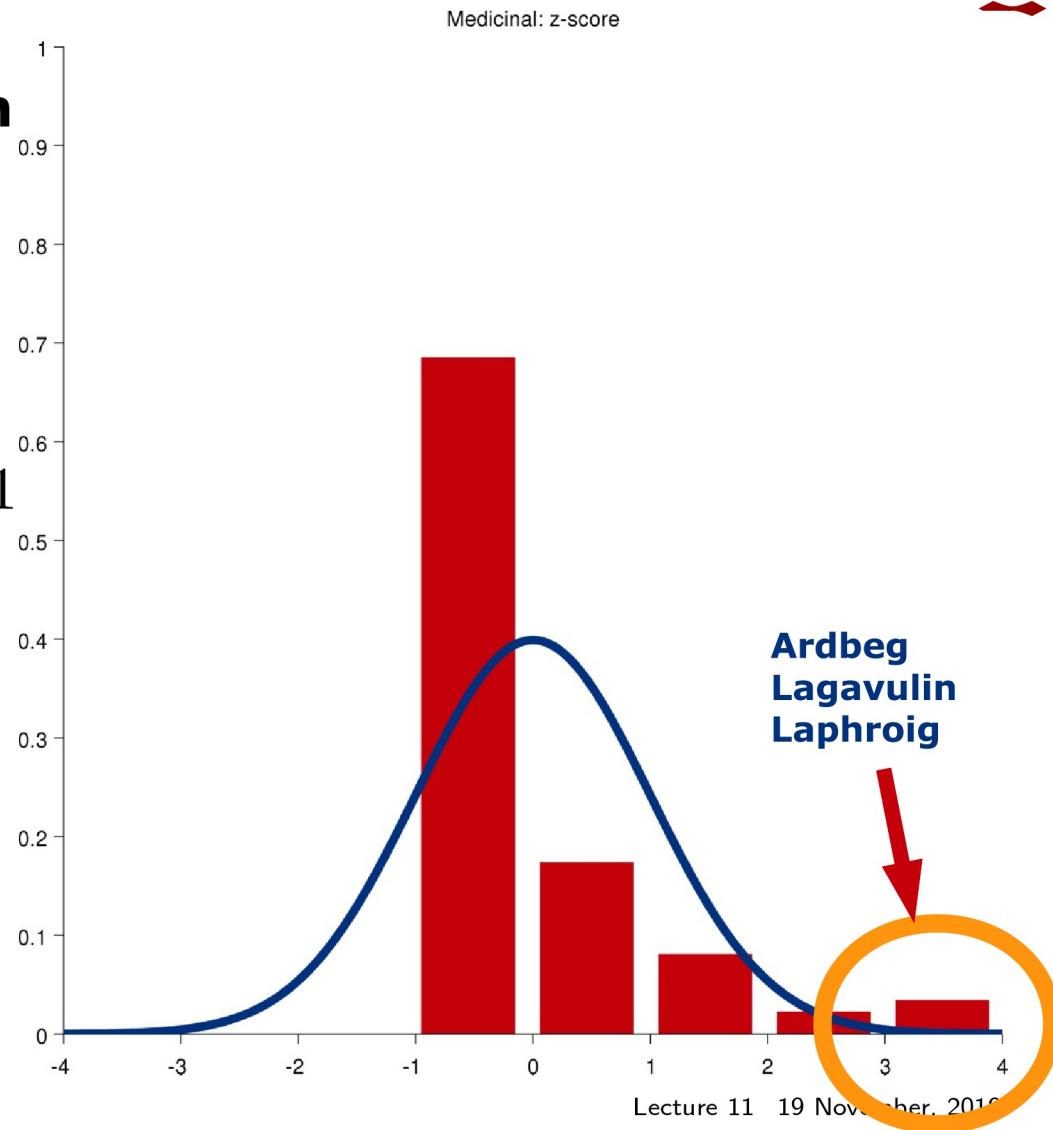
Normal distribution

- Map attribute to standard Normal variable

$$z = \frac{x - \mu}{\sigma}$$

- Choose a threshold

$$p(|z| > c) = 0.001$$
$$c = 3.2905$$



Approaches to anomaly detection

- **Density-based techniques**
 - Estimate the density of data objects
 - Outliers are:
 - Data objects in low density area
- **We can of course use the GMM to evaluate the density of test data.**
 - why not on the training data?**
- **Approaches we will presently also consider:**
 - Kernel density estimation
 - Inverse average distance to K nearest neighbours (KNN density)
 - Average relative KNN density

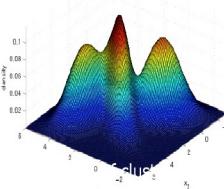
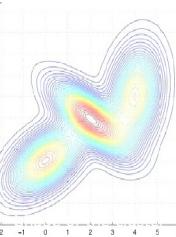
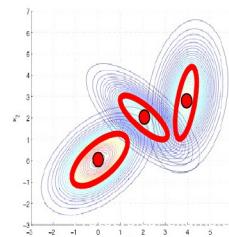
Density based techniques: Kernel Density Estimator

Recall the Gaussian Mixture Model (GMM)

Data density **Sum of cluster specific densities assumed normal distributed**

$$p(\mathbf{x}) = \sum_{k=1}^K w_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{(k)}, \boldsymbol{\Sigma}_{(k)})$$

$$(s.t. \sum_{k=1}^K w_k = 1, \quad w_k \geq 0)$$

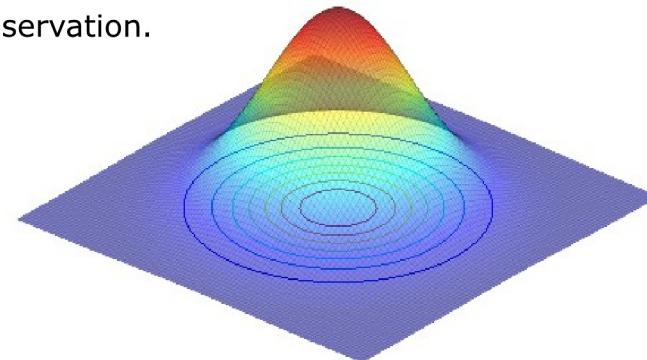


$\boldsymbol{\mu}_{(k)}$: Cluster center (prototypical example in cluster)

$\boldsymbol{\Sigma}_{(k)}$: Shape of the cluster

w_k : Relative density of the cluster

Kernel Density estimation based on Gaussian Kernel:
Consider the GMM and define a Gaussian with mean \mathbf{x}_n and co-variance $\sigma^2 I$ around each Observation.



$$\begin{aligned}\boldsymbol{\mu}_n &= \mathbf{x}_n \\ \boldsymbol{\Sigma}_n &= \sigma^2 I\end{aligned}$$

Let all observation weight the same, i.e. $w_n = 1/N$

$$p(\mathbf{x}) = \sum_{n=1}^N \frac{1}{N} \mathcal{N}(\mathbf{x} | \mathbf{x}_n, \sigma^2 \mathbf{I})$$

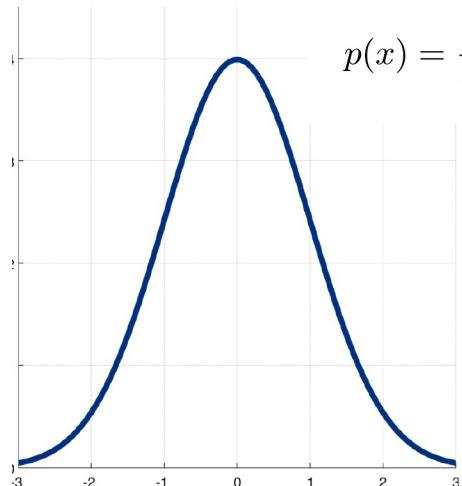
Only free parameter σ^2 !

There is nothing special about the normal distribution. For a general mixture distribution p the general form of kernel density estimator is:

$$p(\mathbf{x}) = \sum_{n=1}^N \frac{1}{N} p(\mathbf{x} | \mathbf{x}_n, \theta)$$

This may be useful if \mathbf{x} is discrete or non-negative.

Piazza quiz 03: Kernel density (Spring 2013)

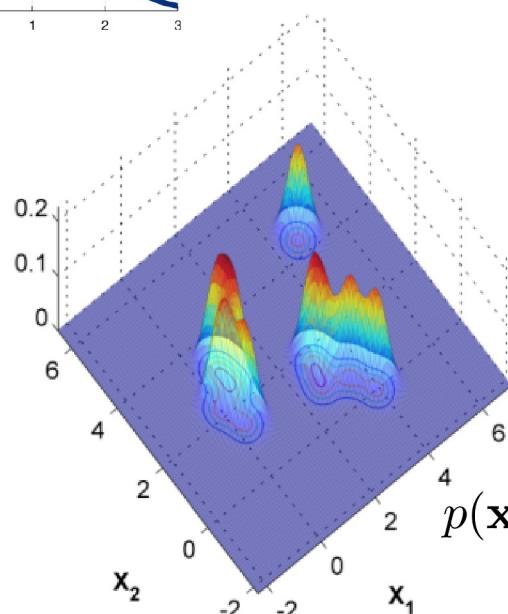


$$p(x) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Consider five observations of an attribute x given by

$$\mathbf{X} = \{2, 3, 5, 10, 12\}.$$

Based on the five observations, what is the Gaussian kernel density estimate at $x = 4$ using $\sigma^2 = 4$?



- A. $\frac{1}{\sqrt{8\pi}} \exp(-\frac{53}{4})$
- B. $\frac{1}{5\sqrt{8\pi}} \exp(-\frac{53}{4})$
- C. $\frac{1}{5\sqrt{8\pi}} (\exp(-\frac{1}{2}) + 2 \cdot \exp(-\frac{1}{8}) + \exp(-\frac{9}{2}) + \exp(-8))$
- D. $\frac{1}{5\sqrt{8\pi}} (\exp(-1) + 2 \cdot \exp(-\frac{1}{4}) + \exp(-9) + \exp(-16))$
- E. Don't know.

$$p(\mathbf{x}) = \sum_{n=1}^N \frac{1}{N} \mathcal{N}(\mathbf{x} | \mathbf{x}_n, \sigma^2 \mathbf{I})$$

Solution:

When inserting $\sigma^2 = 4$ the Gaussian kernel density is given by

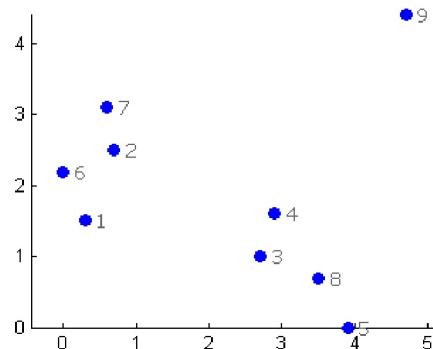
$$\begin{aligned} p(x) &= \frac{1}{N} \sum_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right) = \\ &\quad \frac{1}{5\sqrt{8\pi}} \left(\exp\left(-\frac{(x-2)^2}{2\cdot 4}\right) + \exp\left(-\frac{(x-3)^2}{2\cdot 4}\right) \right. \\ &\quad \left. + \exp\left(-\frac{(x-5)^2}{2\cdot 4}\right) + \exp\left(-\frac{(x-10)^2}{2\cdot 4}\right) + \exp\left(-\frac{(x-12)^2}{2\cdot 4}\right) \right). \end{aligned}$$

Evaluating the density at $x = 4$ we obtain

$$\begin{aligned} p(x=4) &= \frac{1}{N} \sum_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-x_n)^2}{2\sigma^2}\right) = \\ &\quad \frac{1}{5\sqrt{8\pi}} \left(\exp\left(-\frac{(2)^2}{2\cdot 4}\right) + \exp\left(-\frac{(1)^2}{2\cdot 4}\right) \right. \\ &\quad \left. + \exp\left(-\frac{(-1)^2}{2\cdot 4}\right) + \exp\left(-\frac{(-6)^2}{2\cdot 4}\right) + \exp\left(-\frac{(-8)^2}{2\cdot 4}\right) \right) = \\ &\quad \frac{1}{5\sqrt{8\pi}} \left(\exp\left(-\frac{1}{2}\right) + 2 \cdot \left(\exp\left(-\frac{1}{8}\right) + \exp\left(-\frac{9}{2}\right) + \exp(-8) \right) \right). \end{aligned}$$

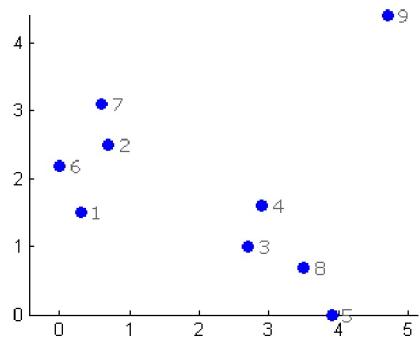
How do we determine σ^2 ?

Data I: Cats, Dogs and Dinosaurs

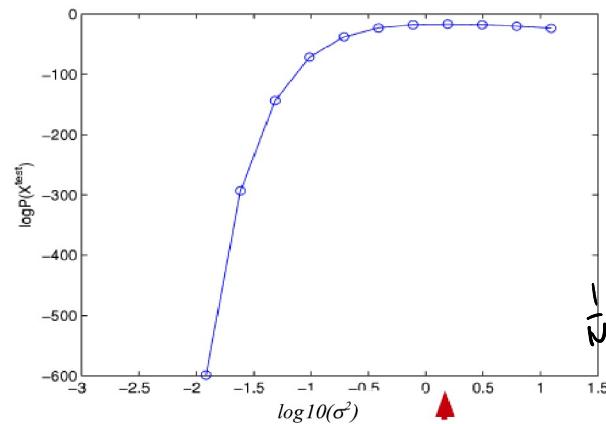


How do we determine σ^2 ? Crossvalidation!

Data I: Cats, Dogs and Dinosaurs



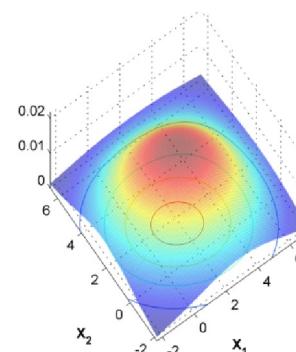
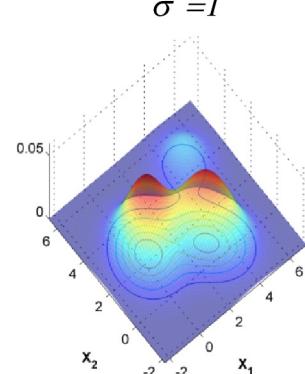
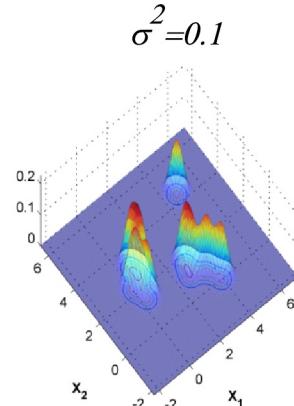
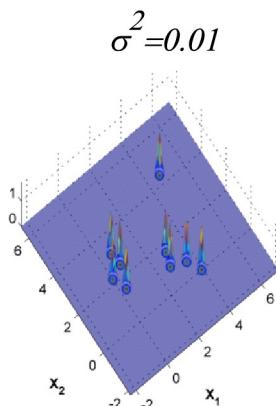
Density of test set based on leave-one-out cross validation



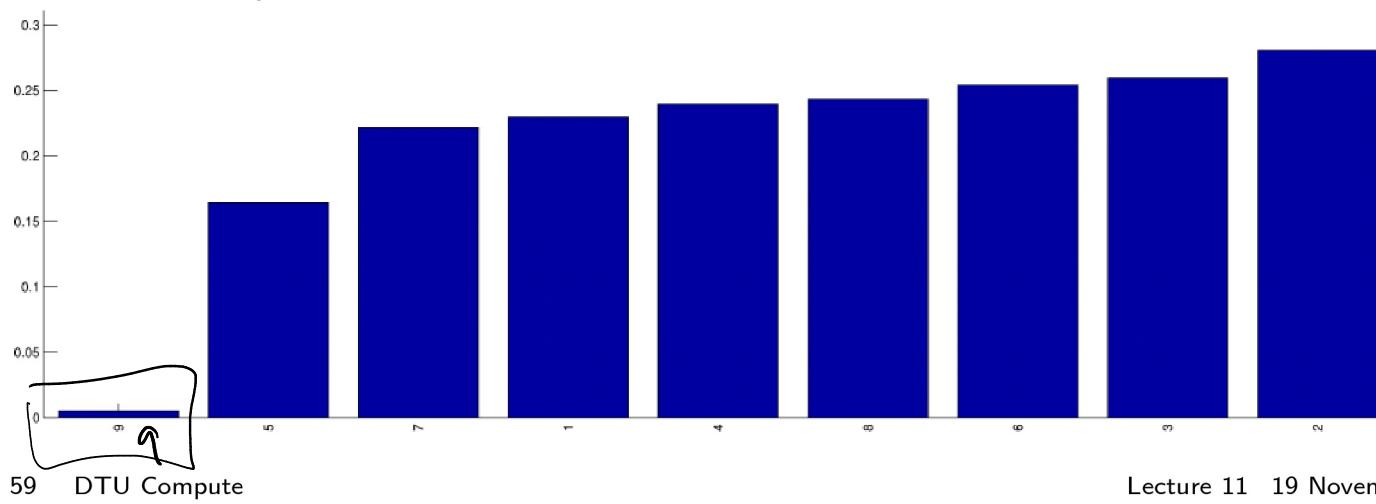
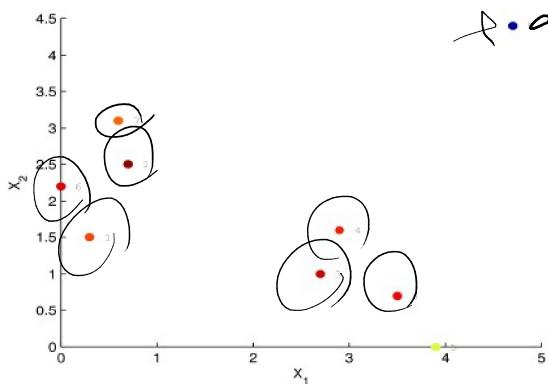
$\sigma^2 = 0.01, 0.1, 1, \dots$

$$P_{N \setminus i}^\sigma(x) \quad \sum_{i=1}^n -\log P_{Ni}^\sigma(x_i) = \text{const.}$$

Optimal $\sigma^2 = 1.55$



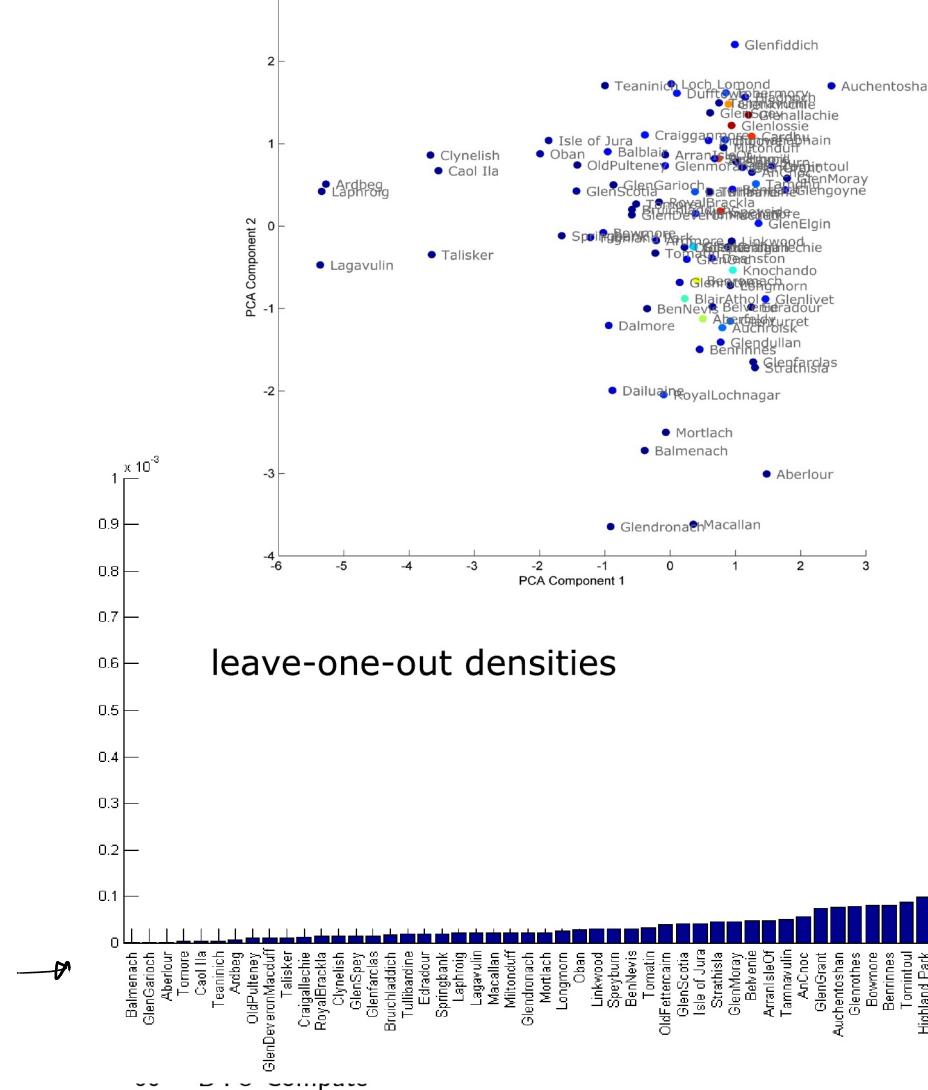
Estimated leave-one-out density evaluated at each observation



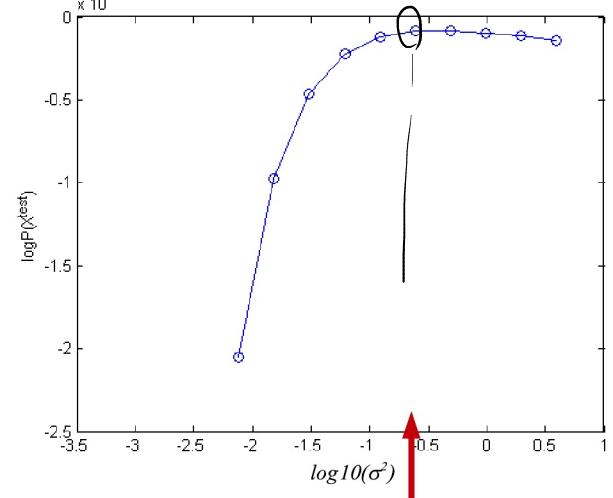
Estimated density evaluated at each observation



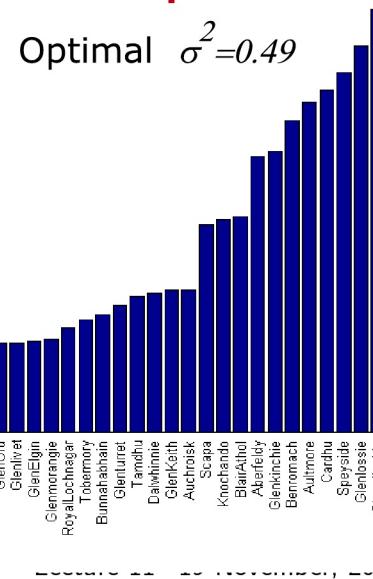
Data II: Whisky



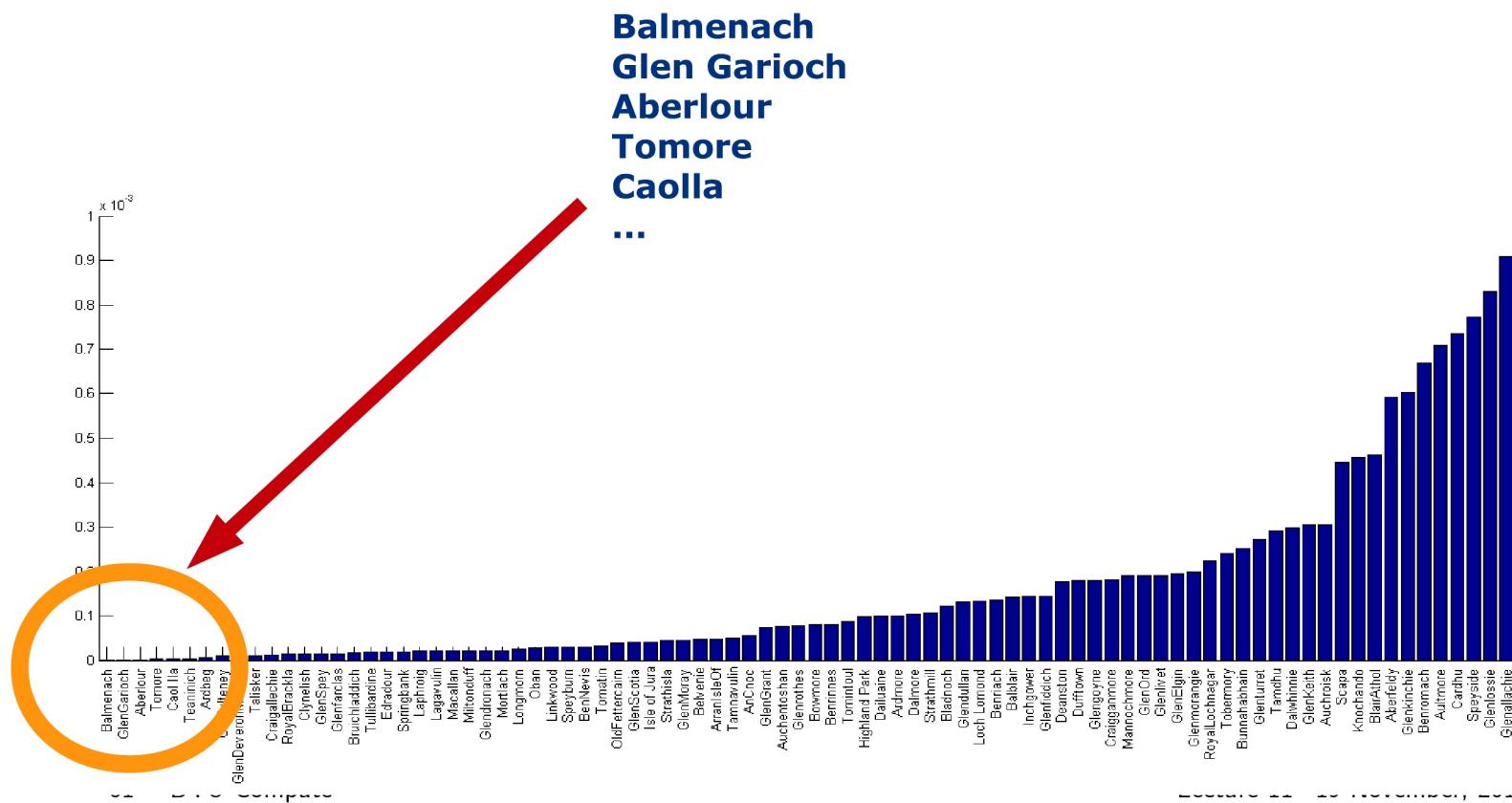
Density of test set based on leave-one-out cross validation

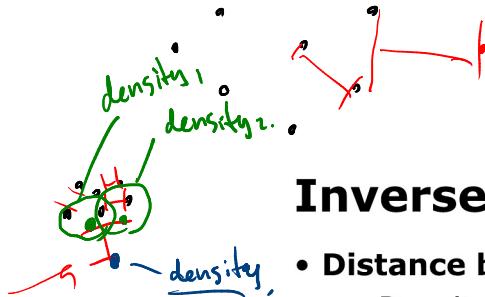


Optimal $\sigma^2 = 0.49$

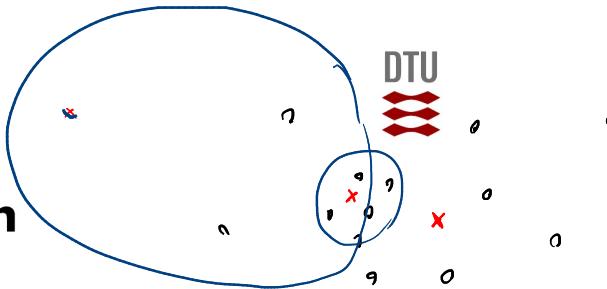


Data II: Whisky





Inverse distance density estimation



- **Distance based measure of density**

- Density is inverse proportional to average distance to k nearest neighbors
- Density is low if nearest neighbors are far away

$$\text{density}_{\mathbf{X} \setminus i}(x_i, K) = \frac{1}{\frac{1}{K} \sum_{x' \in N_{\mathbf{X} \setminus i}(x_i, K)} d(x_i, x')}$$

average distance
 to k nn. (not including x_i)

newest k observations
 to x_i who are not x_i

- **Relative density**

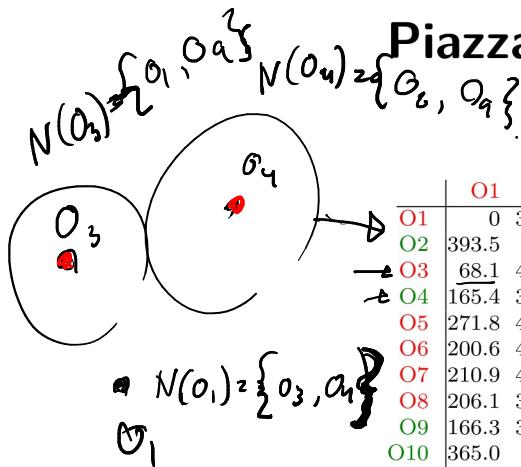
- Density compared to density at nearest neighbors

$$\text{ard}_{\mathbf{X}}(x_i, K) = \frac{\text{density}_{\mathbf{X} \setminus i}(x_i, K)}{\frac{1}{k} \sum_{x_j \in N_{\mathbf{X} \setminus i}(x_i, K)} \text{density}_{\mathbf{X} \setminus j}(x_j, K)}$$

k = 2

$N_{\mathbf{X}}(x, K) = \{\text{The } K \text{ observations in } \mathbf{X} \text{ which are nearest to } x\}$ average of density₁, density₂.

$$\mathbf{X}_{\setminus i}^T = [x_1 \ x_2 \ \cdots \ x_{i-2} \ x_{i-1} \ x_{i+1} \ x_{i+2} \ \cdots \ x_N]$$



Piazza quiz 4: ARD (Spring 2015)

	O1	O2	O3	O4	O5	O6	O7	O8	O9	O10
O1	0	393.5	68.1	165.4	271.8	200.6	210.9	206.1	166.3	365.0
O2	393.5	0	411.3	361.8	478.6	490.9	409.2	382.3	391.1	37.4
O3	68.1	411.3	0	119.8	208.4	136.6	152.8	154.3	111.1	387.1
O4	165.4	361.8	119.8	0	137.5	130.8	62.1	44.7	32.5	346.2
O5	271.8	478.6	208.4	137.5	0	99.0	76.8	101.0	116.4	468.5
O6	200.6	490.9	136.6	130.8	99.0	0	100.1	124.0	100.5	473.8
O7	210.9	409.2	152.8	62.1	76.8	100.1	0	29.5	45.2	396.8
O8	206.1	382.3	154.3	44.7	101.0	124.0	29.5	0	44.6	370.1
O9	166.3	391.1	111.1	32.5	116.4	100.5	45.2	44.6	0	375.1
O10	365.0	37.4	387.1	346.2	468.5	473.8	396.8	370.1	375.1	0

Table 1: Pairwise Euclidean distance between the 10 first observations in the PM10 (air pollution) dataset. Red/green observations corresponds to high/low pollution levels.

$$\text{density}(O_1) = \frac{1}{2}(68 + 165)$$

$$\text{density}(O_3) = \frac{1}{2}(68 + 111)$$

$$\text{density}(O_4) = \frac{1}{2}(44.7 + 32.5)$$

$$\text{Ard}(O_1) = \frac{\text{density}(O_1)}{\frac{1}{2}(\text{density}(O_3) + \text{density}(O_4))} =$$

We suspect that observation O1 in Table 1 may be an outlier. In order to assess if this is the case we will calculate the average relative density (ARD) based on the distances in the table using the definitions:

$$\text{density}(\mathbf{x}, K) = \left(\frac{1}{K} \sum_{y \in N(\mathbf{x}, K)} \text{distance}(\mathbf{x}, \mathbf{y}) \right)^{-1},$$

$$\text{a.r.d.}(\mathbf{x}, K) = \frac{\text{density}(\mathbf{x}, K)}{\frac{1}{K} \sum_{y \in N(\mathbf{x}, K)} \text{density}(\mathbf{y}, K)},$$

where $N(\mathbf{x}, K)$ is the set of K nearest neighbors of observation \mathbf{x} and $\text{a.r.d.}(\mathbf{x}, K)$ is the average relative density of \mathbf{x} using K nearest neighbors. What is ARD for observation O1 for $K = 2$ nearest neighbors?

- A. 0.01
- B. 0.02
- C. 0.23
- D. 0.46
- E. Don't know.

Solution:

	O1	O2	O3	O4	O5	O6	O7	O8	O9	O10
O1	0	393.5	68.1	165.4	271.8	200.6	210.9	206.1	166.3	365.0
O2	393.5	0	411.3	361.8	478.6	490.9	409.2	382.3	391.1	37.4
O3	68.1	411.3	0	119.8	208.4	136.6	152.8	154.3	111.1	387.1
O4	165.4	361.8	119.8	0	137.5	130.8	62.1	44.7	32.5	346.2
O5	271.8	478.6	208.4	137.5	0	99.0	76.8	101.0	116.4	468.5
O6	200.6	490.9	136.6	130.8	99.0	0	100.1	124.0	100.5	473.8
O7	210.9	409.2	152.8	62.1	76.8	100.1	0	29.5	45.2	396.8
O8	206.1	382.3	154.3	44.7	101.0	124.0	29.5	0	44.6	370.1
O9	166.3	391.1	111.1	32.5	116.4	100.5	45.2	44.6	0	375.1
O10	365.0	37.4	387.1	346.2	468.5	473.8	396.8	370.1	375.1	0

Table 1: Pairwise Euclidean distance between the 10 first observations in the PM10 (air pollution) dataset. Red/green observations corresponds to high/low pollution levels.

The intermediate computations are:

$$\text{density}(\mathbf{x}_{O1}, 2) = \left(\frac{1}{2} \cdot (68.1 + 165.4) \right)^{-1} = 0.0086$$

$$\text{density}(\mathbf{x}_{O3}, 2) = \left(\frac{1}{2} \cdot (68.1 + 111.1) \right)^{-1} = 0.0112$$

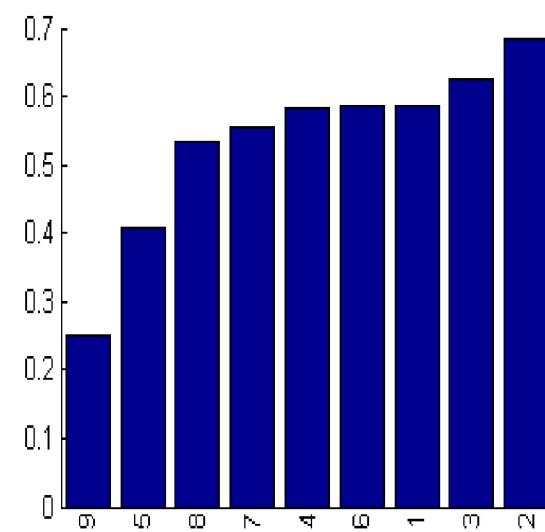
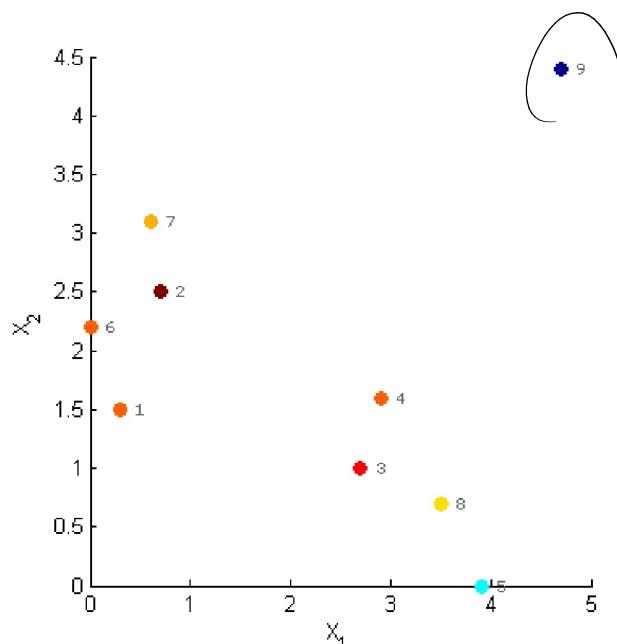
$$\text{density}(\mathbf{x}_{O4}, 2) = \left(\frac{1}{2} \cdot (32.5 + 44.7) \right)^{-1} = 0.0259$$

$$\begin{aligned} \text{a.r.d.}(\mathbf{x}, K) &= \frac{\text{density}(\mathbf{x}_{O1}, 2)}{\frac{1}{2}(\text{density}(\mathbf{x}_{O3}, 2) + \text{density}(\mathbf{x}_{O4}, 2))} \\ &= \frac{0.0086}{\frac{1}{2} \cdot (0.0112 + 0.0259)} = 0.46 \end{aligned}$$

Inverse distance density estimation

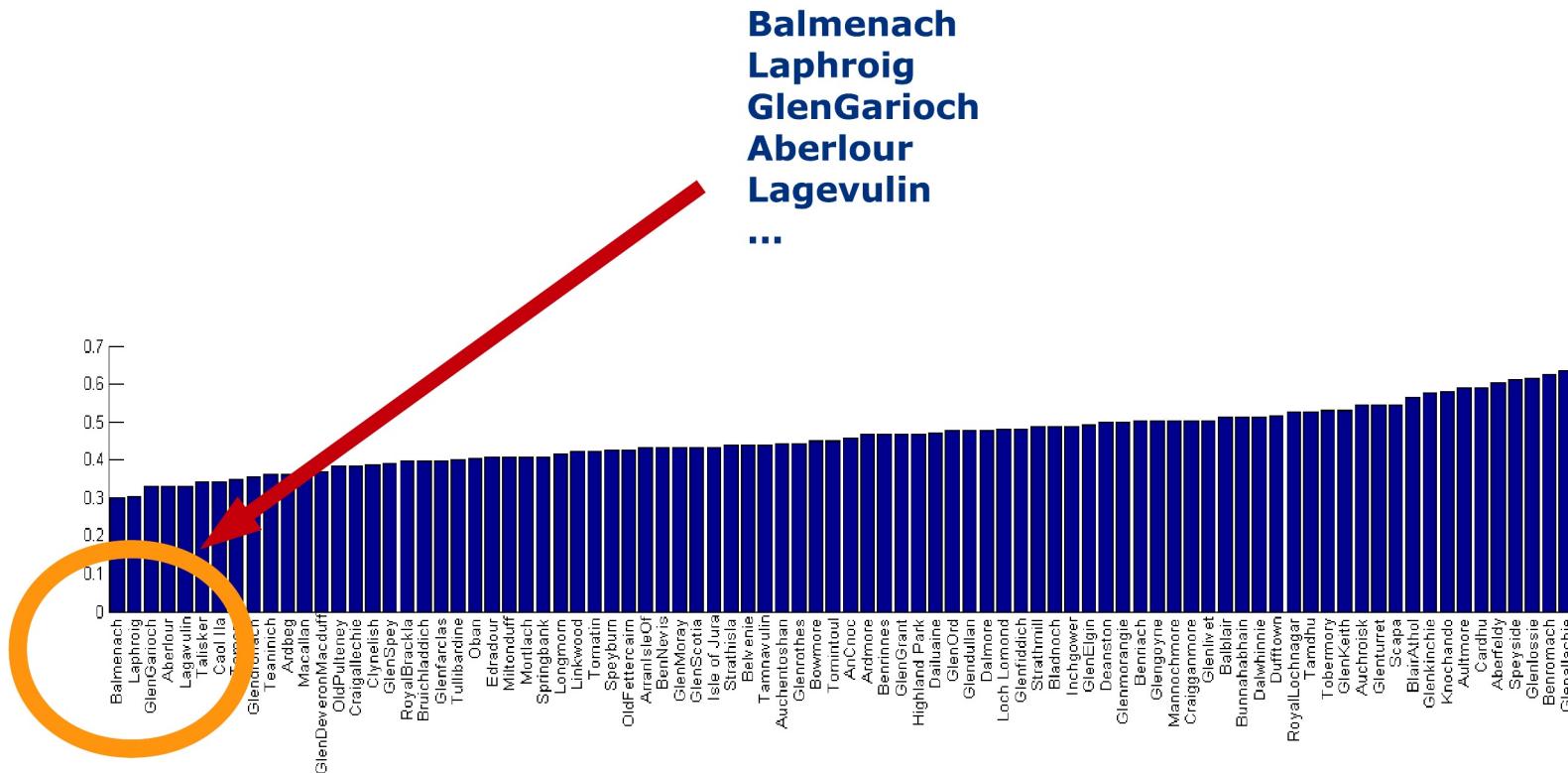
- KNN density (5 nearest neighbors)

Data I: Cats , dogs and dinosaurs



Inverse distance density estimation

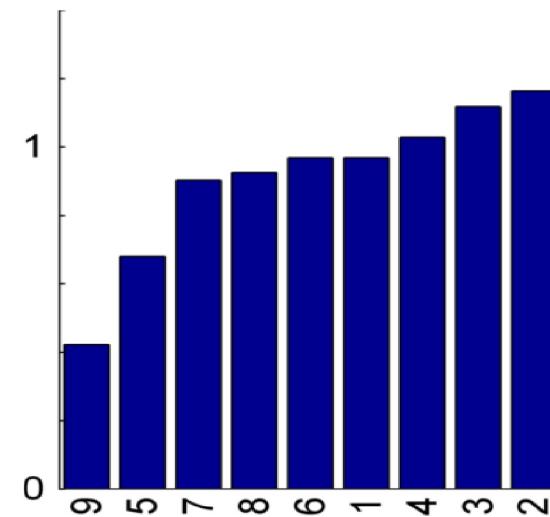
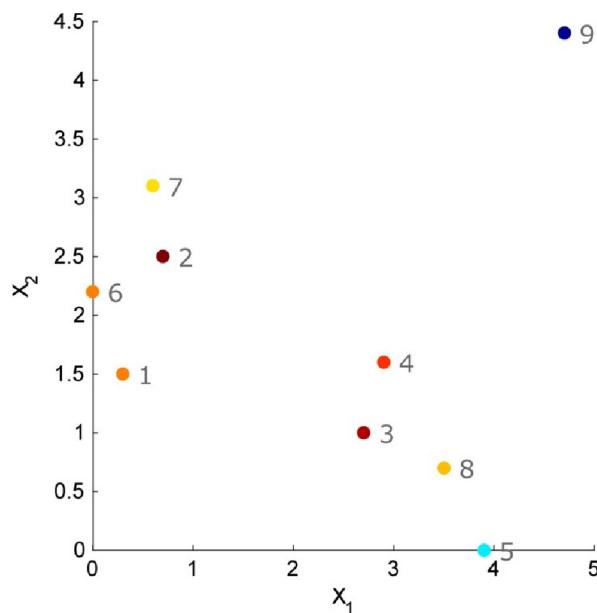
- KNN density (5 nearest neighbors)



Average Relative density

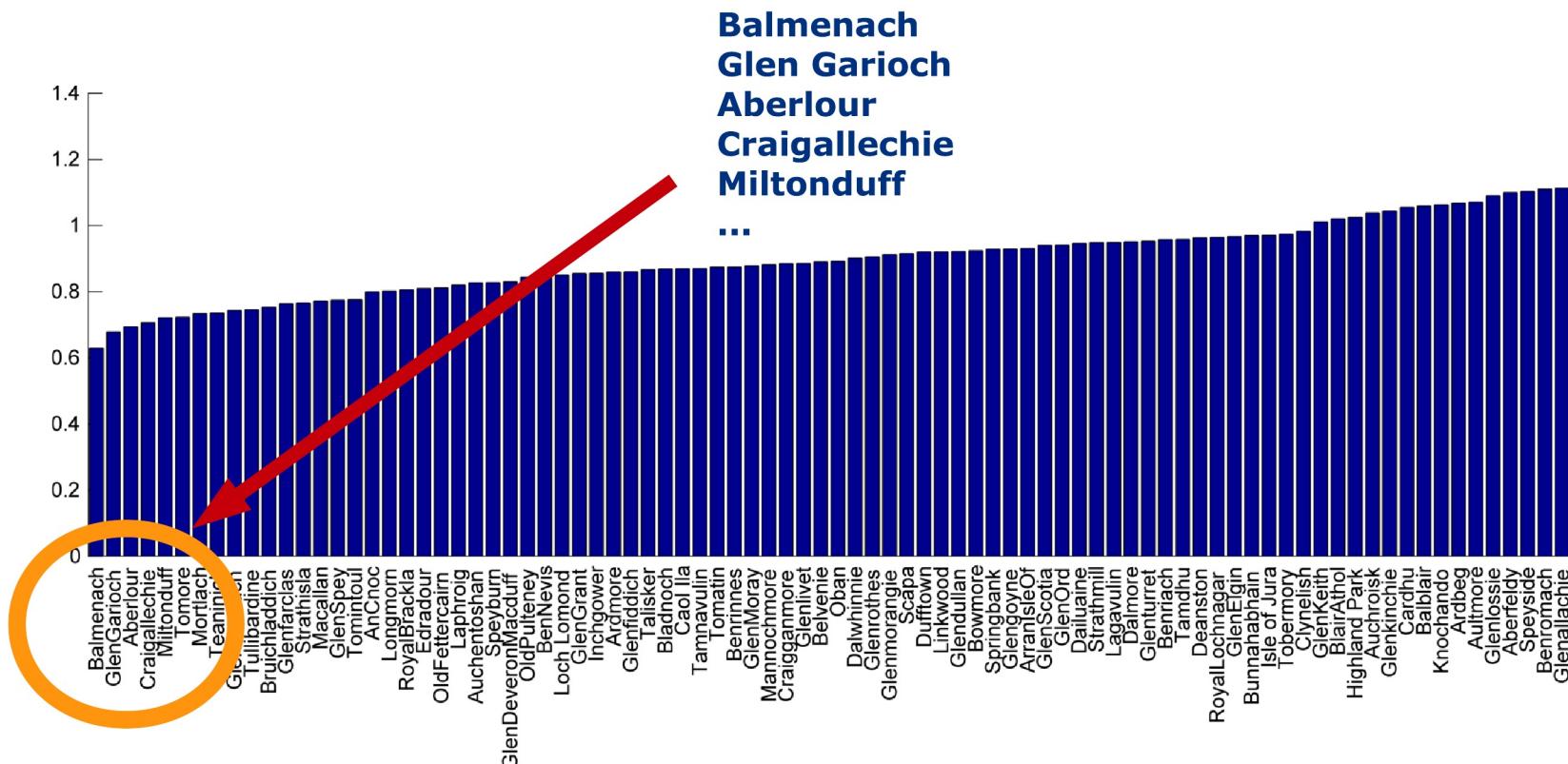
- Average Relative KNN density (5 nearest neighbors)

Data I: Cats , dogs and dinosaurs



Average relative density

- Average relative KNN density (5 nearest neighbors)



Results using different methods



- **Kernel Density Estimation**
 - Balmenach
 - Glen Garioch
 - Aberlour
 - Tomore
 - Caolla
- **KNN density**
 - Balmenach
 - Laphroig
 - Glen Garioch
 - Aberlour
 - Lagavulin
- **KNN average relative density**
 - Balmenach
 - Glen Garioch
 - Aberlour
 - Craigallechie
 - Miltonduff

Common: Balmenach, Glen Garioch,
Aberlour

About project 3

- Assuming your (true) class labels are y and your data is $X = N \times M$
- Cluster your data in the original M -dimensional space. Remember to exclude the y -variable.
- Evaluate performance by computing similarity between clustering and y (last week)
- Plot data and clusters by projecting onto e.g. 2D space using PCA



Resources

<https://www.youtube.com> Nice explanation of expectation maximization for
the Guassian Mixture Model (<https://www.youtube.com/watch?v=WaKNSBeDLTw>)