

02450 Introduction to machine learning and data mining

DTU Informatics

Department of Informatics and Mathematical Modeling

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

$$\Theta^{\sqrt{17}} + \Omega \int_a^b \delta e^{inx} dx = [2.7182818284 \dots] \sum_{n=0}^{\infty} x^n$$

Reading material

Tan, Steinbach and Kumar

"Introduction to Data Mining"

Section 10.1-10.4

Feedback Groups of the day:

Wen Hsin Li, Robert Lyck, Viivi Halla-aho

Oliver Naaby, Rasmus Olsen

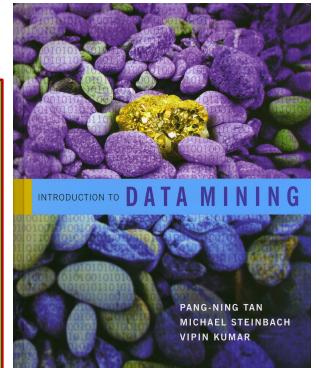
Simon Hemmingsen, Thomas Hammerbak

Henrik Thirstrup, Mette Marie Nielsen
Thomas Kristensen

Nicolai Pedersen, Jeppe Thagaard, Mathias Brade,
Christian Kragh

Oscar Petersen, Søren Vørne Nielsen

Andreas Jacobsen, Asger Anker Sørensen, Phong
Le Trung



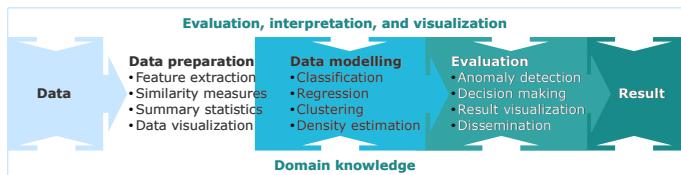
If possible, please (1) stay to give me feedback after the second lecture today (ca. 15:00) and (2) send an email or contact me at the exercises next week with feedback/suggestions on the exercises for today.

Lecture schedule

1. Introduction
(Tan 1.1-1.4)
2. Data and feature extraction
(Tan 2.1-2.3 + B1 (+ A))
3. Measures of similarity and summary statistics
(Tan 2.4 + 3.1-3.2 + C1-C2)
4. Data visualization
(Tan 3.3)
5. Decision trees and linear regression
(Tan 4.1-4.3 + D)
6. Overfitting and performance evaluation
(Tan 4.4-4.6)
7. Nearest neighbor, naive Bayes, and artificial neural networks
(Tan 5.2-5.4)
8. Ensemble methods and multi class classifiers
(Tan 5.6-5.8)
9. K-means and hierarchical clustering
(Tan 8.1-8.3+8.5.7)
10. Mixture models and association mining
(Tan 9.2.2 + 6.1-6.3)
11. Density estimation and anomaly detection
(Tan 10.1-10.4)
12. Putting it all together: Summary and overview
13. Mini project

3 DTU Informatics, Technical University of Denmark

Data modeling framework

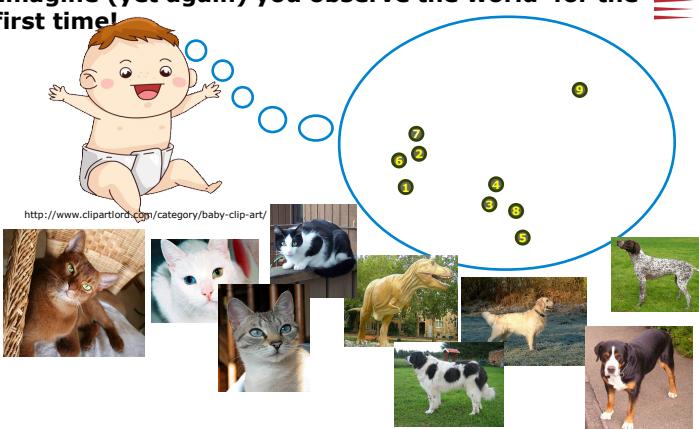


After today you should be able to:

Understand and apply a variety of outlier/anomaly detection approaches, including density estimation approaches and proximity-based techniques.

4 DTU Informatics, Technical University of Denmark

Imagine (yet again) you observe the world for the first time!



How do we detect anomalous objects
(i.e., the dinosaur in the world of cats and dogs?)

5 DTU Informatics, Technical University of Denmark

Anomaly detection: Definition

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

6 DTU Informatics, Technical University of Denmark

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

7 DTU Informatics, Technical University of Denmark



Group exercise

- Come up with **your own definition** of an outlier / anomaly
- How can we detect outliers using some of the methods you have already learned in the course?

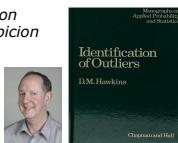
8 DTU Informatics, Technical University of Denmark



Group exercise

- Come up with **your own definition** of an outlier / anomaly
- How can we detect outliers using some of the methods you have already learned in the course?

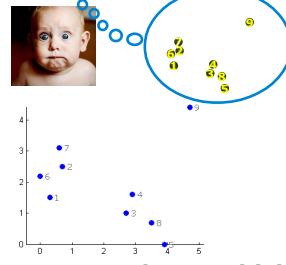
Hawkins' definition of an outlier: An outlier is an observation that differs so much from other observations as to arouse suspicion that it was generated by a different mechanism



Probabilistic definition of an outlier: An outlier is an object that has a low probability with respect to a probability distribution model of the data.

9 DTU Informatics, Technical University of Denmark

Data example I: Cats, Dogs and Dinosaurs



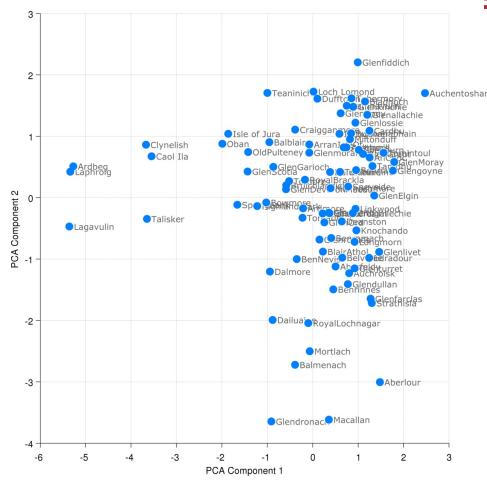
Data example II: Whisky

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



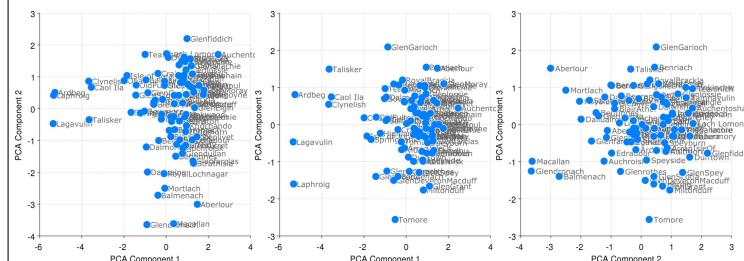
10 DTU Informatics, Technical University of Denmark

PCA plot



11 DTU Informatics, Technical University of Denmark

PCA plot

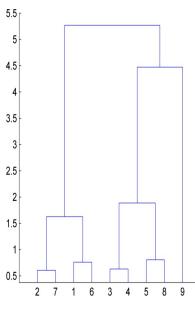


12 DTU Informatics, Technical University of Denmark

Dendrogram

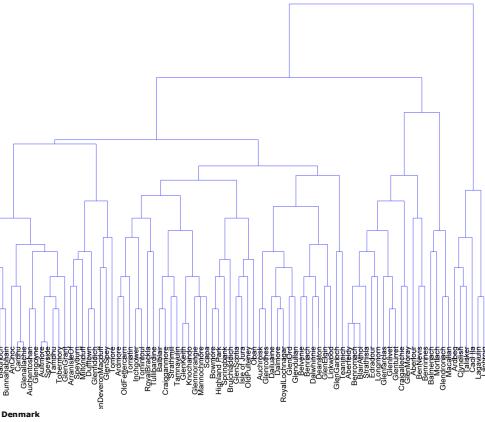
- Dendograms can be used to visualize relative distances between the observations

Data I: Cats, Dogs and Dinosaurs



13 DTU Informatics, Technical University of Denmark

Data II: Whisky



Approaches to anomaly detection

Density-based techniques

- Estimate the density of data objects
- Outliers are:
 - Data objects in low density area

Approaches we will consider:

- Univariate normal distribution
- Kernel density estimation

Proximity-based techniques

- Measure the distance between data objects
- Outliers are:
 - Data objects far from the other data objects

Approaches we will consider:

- Mahalanobis distance to center of data
- Distance to Kth nearest neighbour
- Inverse average distance to K nearest neighbours (KNN density)
- Average relative KNN density

14 DTU Informatics, Technical University of Denmark

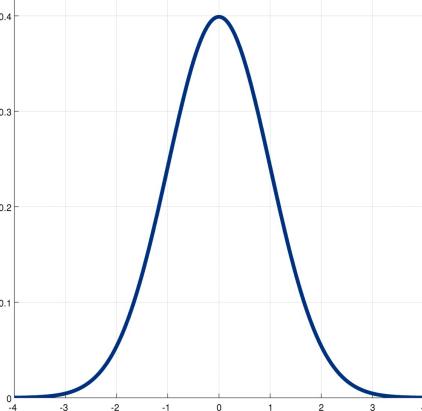
Density based techniques: Univariate normal distribution

- 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



15 DTU Informatics, Technical University of Denmark

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

16 DTU Informatics, Technical University of Denmark

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

17 DTU Informatics, Technical University of Denmark

Medicinal: z-score



Normal distribution

- Map attribute to standard Normal variable

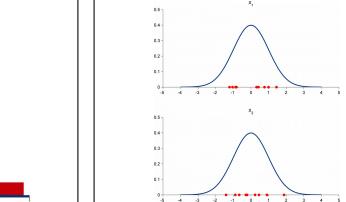
$$z = \frac{x - \mu}{\sigma}$$

- Choose a threshold

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

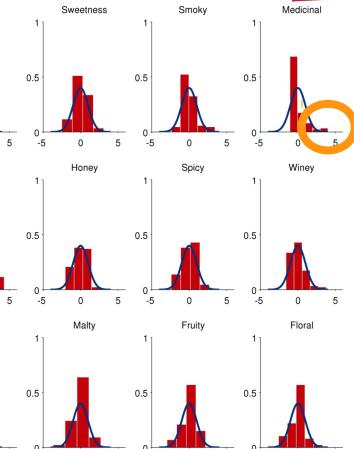
$$c = 3.2905$$

Data I: Cats, Dogs and Dinosaurs



18 DTU Informatics, Technical University of Denmark

Data II: Whisky



Normal distribution

- Map attribute to standard Normal variable

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

- Choose a threshold

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

Medicinal: z-score

Ardbeg
Lagavulin
Laphroig

19 DTU Informatics, Technical University of Denmark

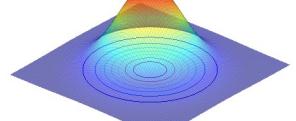
Density based techniques: Kernel Density Estimator

Remember from Last week:
Gaussian Mixture Model (GMM)

$$\text{Data density} \quad \sum_{k=1}^K w_k \mathcal{N}(x|\mu_{(k)}, \Sigma_{(k)}) \\ (\text{s.t. } \sum_{k=1}^K w_k = 1, \quad w_k \geq 0)$$

$\mu_{(k)}$: Cluster center (prototypical example in cluster)
 $\Sigma_{(k)}$: 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 x_n and co-variance $\sigma^2 I$ around each Observation.



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

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

Only free parameter σ !

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

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

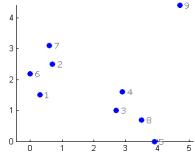
This may be useful if x is discrete or non-negative.

20 DTU Informatics, Technical University of Denmark

How do we determine σ ?



Data I: Cats, Dogs and Dinosaurs

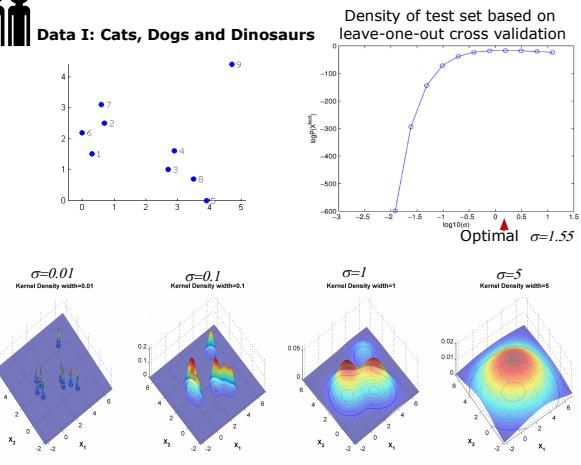


21 DTU Informatics, Technical University of Denmark

How do we determine σ ?

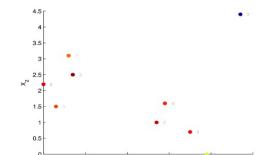


Data I: Cats, Dogs and Dinosaurs



Estimated leave-one-out density evaluated at each observation

Estimated leave-one-out density evaluated at each observation

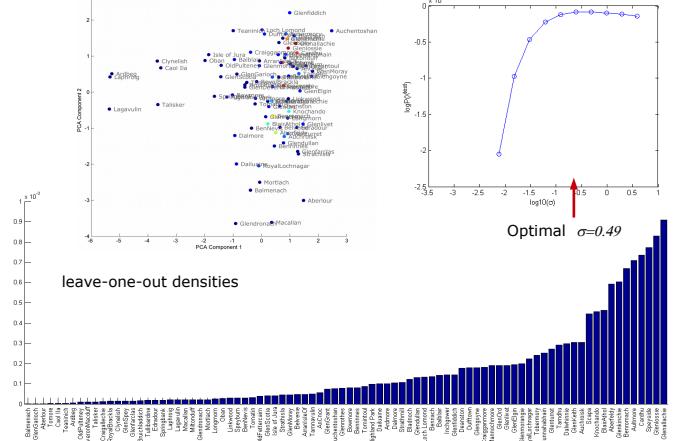


23 DTU Informatics, Technical University of Denmark

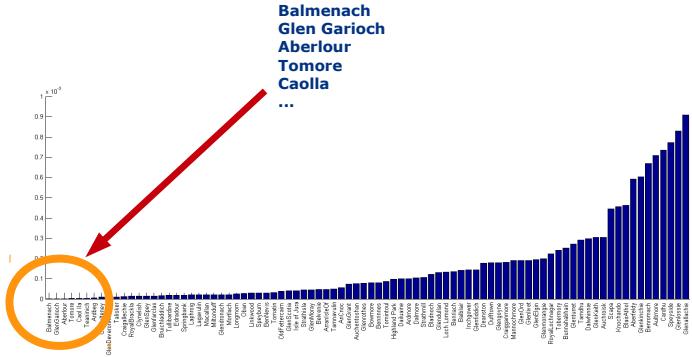
Data II: Whisky



Data II: Whisky



Data II: Whisky



Proximity-based techniques

- Mahalanobis distance to center of data
- Distance to Kth nearest neighbour
- Inverse average distance to KNN

26 DTU Informatics, Technical University of Denmark

Mahalanobis distance

- Distance to the mean of data, taking covariance into account

$$d_{\text{mahalanobis}}(\mathbf{x}, \bar{\mathbf{x}}) = (\mathbf{x} - \bar{\mathbf{x}})^T S^{-1} (\mathbf{x} - \bar{\mathbf{x}})$$

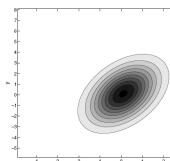


Figure 10.2. Probability density of points for the Gaussian distribution used to generate the points of Figure 10.3.

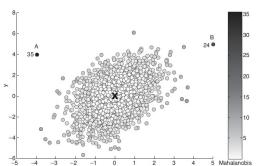


Figure 10.3. Mahalanobis distance of points from the center of a two-dimensional set of 2002 points.

$\bar{\mathbf{x}}$ Mean vector

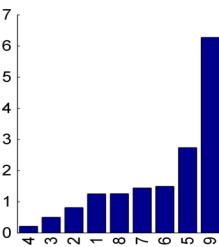
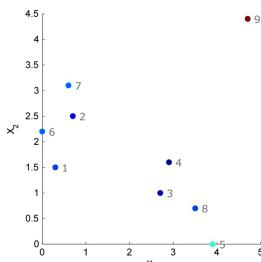
S Covariance matrix

27 DTU Informatics, Technical University of Denmark

Mahalanobis distance

- Distance to the mean of data, taking covariance into account

Data I: Cats , dogs and dinosaurs

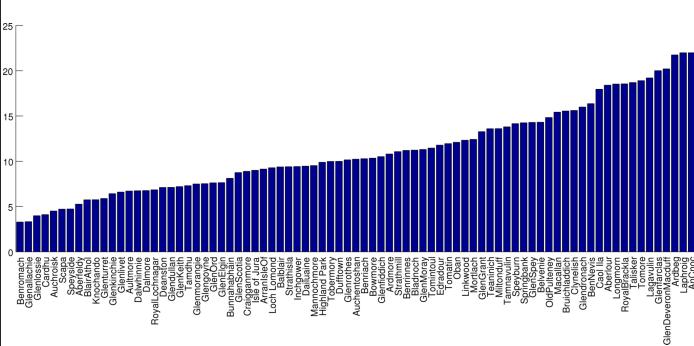


28 DTU Informatics, Technical University of Denmark

Mahalanobis distance

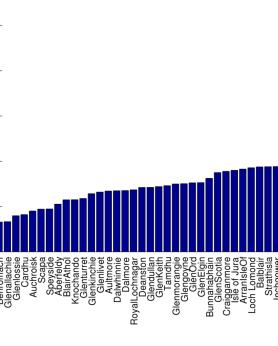
- Distance to the mean of data, taking covariance into account

Data II: Whisky



Mahalanobis distance

- Distance to the mean of data, taking covariance into account



30 DTU Informatics, Technical University of Denmark

Distance to k-nearest neighbor

- Measure the distance to the k'th nearest neighbor

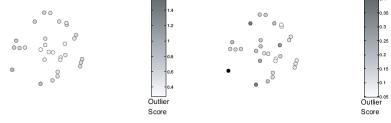


Figure 10.4. Outlier score based on the distance to fifth nearest neighbor.

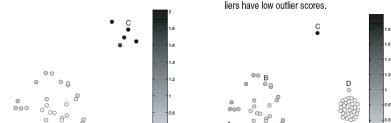


Figure 10.5. Outlier score based on the distance to the first nearest neighbor. Nearby outliers have low outlier scores.

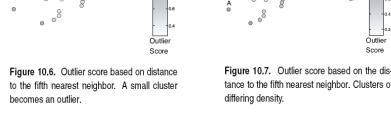


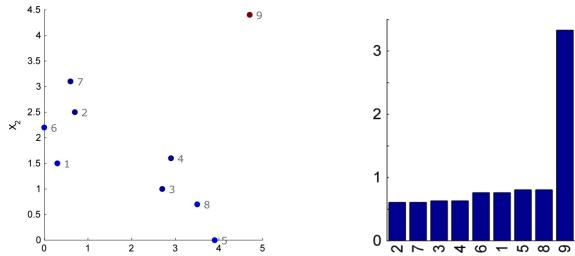
Figure 10.6. Outlier score based on distance to the fifth nearest neighbor. A small cluster becomes an outlier.

Figure 10.7. Outlier score based on the distance to the fifth nearest neighbor. Clusters of differing density.

Distance to k-nearest neighbor

- Measure the distance to the 1st nearest neighbor

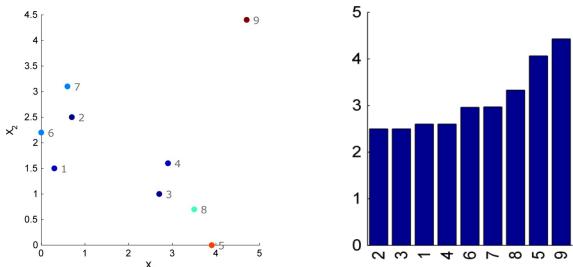
Data I: Cats , dogs and dinosaurs



Distance to kth Nearest neighbour

- Measure the distance to the 5th nearest neighbor

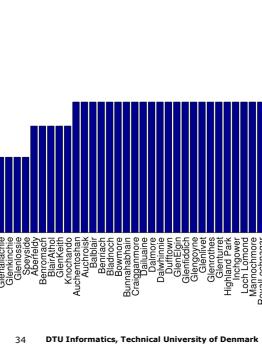
Data I: Cats , dogs and dinosaurs



Distance to k-nearest neighbor

- Measure the distance to the 1st nearest neighbor

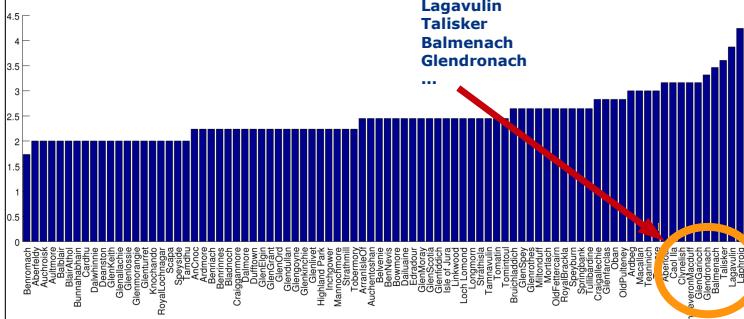
Glen Garroch
Balmenach
Aberlour
Tomore
Teaninich



Distance to k-nearest neighbor

- Measure the distance to the 5th nearest neighbor

Laphroig
Lagavulin
Talisker
Balmenach
Glendronach
...



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

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

Relative density

- Density compared to density at nearest neighbors

$$\text{average relative density}(\mathbf{x}, k) = \frac{\text{density}(\mathbf{x}, k)}{\frac{1}{k} \sum_{\mathbf{y} \in N(\mathbf{x}, k)} \text{density}(\mathbf{y}, k)}$$

$N(\mathbf{x}, k)$ The set of k nearest neighbors



Consider the pairwise distance matrix given to the left. What is the density and average relative density of the first observation for $k=3$?

$d(\mathbf{x}_i, \mathbf{x}_j)$	\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_3	\mathbf{x}_4	\mathbf{x}_5	\mathbf{x}_6	\mathbf{x}_7
\mathbf{x}_1	0	2.5	2.4	4.0	0.8	0.6	3.3
\mathbf{x}_2	2.5	0	0.6	1.6	2.9	3.0	1.1
\mathbf{x}_3	2.4	0.6	0	1.9	3.0	2.7	1.0
\mathbf{x}_4	4.0	1.6	1.9	0	4.5	4.6	3.8
\mathbf{x}_5	0.8	2.9	3.0	4.5	0	1.1	3.9
\mathbf{x}_6	0.6	3.0	2.7	4.6	1.1	0	3.8
\mathbf{x}_7	3.3	1.1	1.0	3.8	3.9	3.8	0

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

$$\text{average relative density}(\mathbf{x}, k) = \frac{\text{density}(\mathbf{x}, k)}{\frac{1}{k} \sum_{y \in N(\mathbf{x}, k)} \text{density}(\mathbf{y}, k)}$$



Consider the pairwise distance matrix given to the left. What is the density and average relative density of the first observation for $k=3$?

$d(\mathbf{x}_i, \mathbf{x}_j)$	\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_3	\mathbf{x}_4	\mathbf{x}_5	\mathbf{x}_6	\mathbf{x}_7
\mathbf{x}_1	0	2.5	2.4	4.0	0.8	0.6	3.3
\mathbf{x}_2	2.5	0	0.6	1.6	2.9	3.0	1.1
\mathbf{x}_3	2.4	0.6	0	1.9	3.0	2.7	1.0
\mathbf{x}_4	4.0	1.6	1.9	0	4.5	4.6	3.8
\mathbf{x}_5	0.8	2.9	3.0	4.5	0	1.1	3.9
\mathbf{x}_6	0.6	3.0	2.7	4.6	1.1	0	3.8
\mathbf{x}_7	3.3	1.1	1.0	3.8	3.9	3.8	0

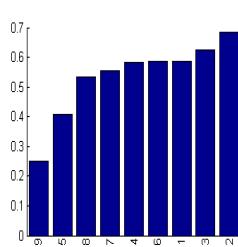
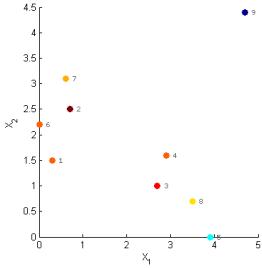
$$\text{density}(\mathbf{x}, k) = \left(\frac{1}{k} \sum_{y \in N(\mathbf{x}, k)} \text{distance}(\mathbf{x}, \mathbf{y}) \right)^{-1}$$

$$\text{average relative density}(\mathbf{x}, k) = \frac{\text{density}(\mathbf{x}, k)}{\frac{1}{k} \sum_{y \in N(\mathbf{x}, k)} \text{density}(\mathbf{y}, k)}$$

Inverse distance density estimation

- KNN density (5 nearest neighbors)

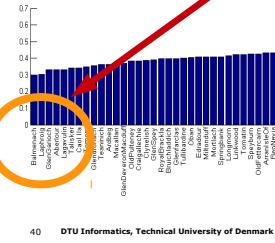
Data I: Cats , dogs and dinosaurs



Inverse distance density estimation

- KNN density (5 nearest neighbors)

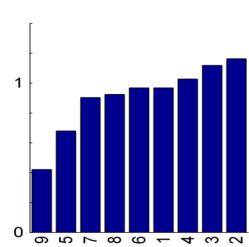
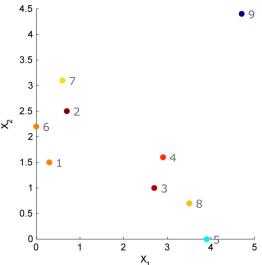
Balmenach
Laphroig
GlenGarioch
Aberlour
Lageveulin
...



Average Relative density

- Average Relative KNN density (5 nearest neighbors)

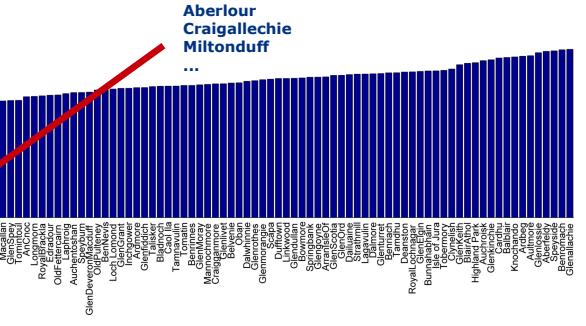
Data I: Cats , dogs and dinosaurs



Average relative density

- Average relative KNN density (5 nearest neighbors)

Balmenach
Glen Garioch
Aberlour
Craigallechie
Miltonduff
...



Results using different methods

- **Univariate Normal distribution**

- Ardbeg
- Lagavulin
- Laphroig

- **Kernel Density Estimation**

- Balmenach
- Glen Garioch
- Aberlour
- Tomore
- Caolla

- **Mahalanobis distance**

- Tullibardine
- Craigallechie
- Glen Garioch
- Old Fettercairn
- Balmenach

- **Distance to nearest neighbor**

- Glen Garioch
- Balmenach
- Aberlour
- Tomore
- Teaninich

43 DTU Informatics, Technical University of Denmark

- **Distance to 5th nearest neighbor**

- Laphroig
- Lagavulin
- Talisker
- Balmenach
- Glendronach

- **KNN density**

- Balmenach
- Laphroig
- Glen Garioch
- Aberlour
- Lagavulin

- **KNN average relative density**

- Balmenach
- Glen Garioch
- Aberlour
- Craigallechie
- Miltonduff

Common: Balmenach, Glen Garioch, Laphroig, Aberlour, Tomore, Lagavulin, Craigallechie

Example of exam questions

QI: What is the average relative density for observation 2 (i.e. \mathbf{x}_2) for k=2 nearest neighbours?

A:1/5
B:3/10
C:7/10
D:1

d($\mathbf{x}_i, \mathbf{x}_j$)	\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_3	\mathbf{x}_4	\mathbf{x}_5
)	0	2.0	0.2	0.9	0.2
\mathbf{x}_1	2.0	0	1.5	0.5	2.0
\mathbf{x}_2	0.2	1.5	0	1.2	1.4
\mathbf{x}_3	0.9	0.5	1.2	0	1.0
\mathbf{x}_4	0.2	2.0	1.4	1.0	0
\mathbf{x}_5	0.2	2.0	1.4	1.0	0

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

$$\text{avg. rel. den.}(\mathbf{x}, k) = \frac{\text{density}(\mathbf{x}, k)}{\sum_{y \in N(\mathbf{x}, k)} \text{density}(\mathbf{y}, k)}$$

44 DTU Informatics, Technical University of Denmark

QI: C