

Unsupervised learning: clustering

Supervised learning: decision-trees

25.11.2024

EPFL

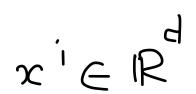
Outline

- Hour 1
 - Brief review of PCA
 - Unsupervised learning: clustering

- Hour 2: back to supervised learning
 - Decision-trees
 - For regression
 - For classification



PCA - example



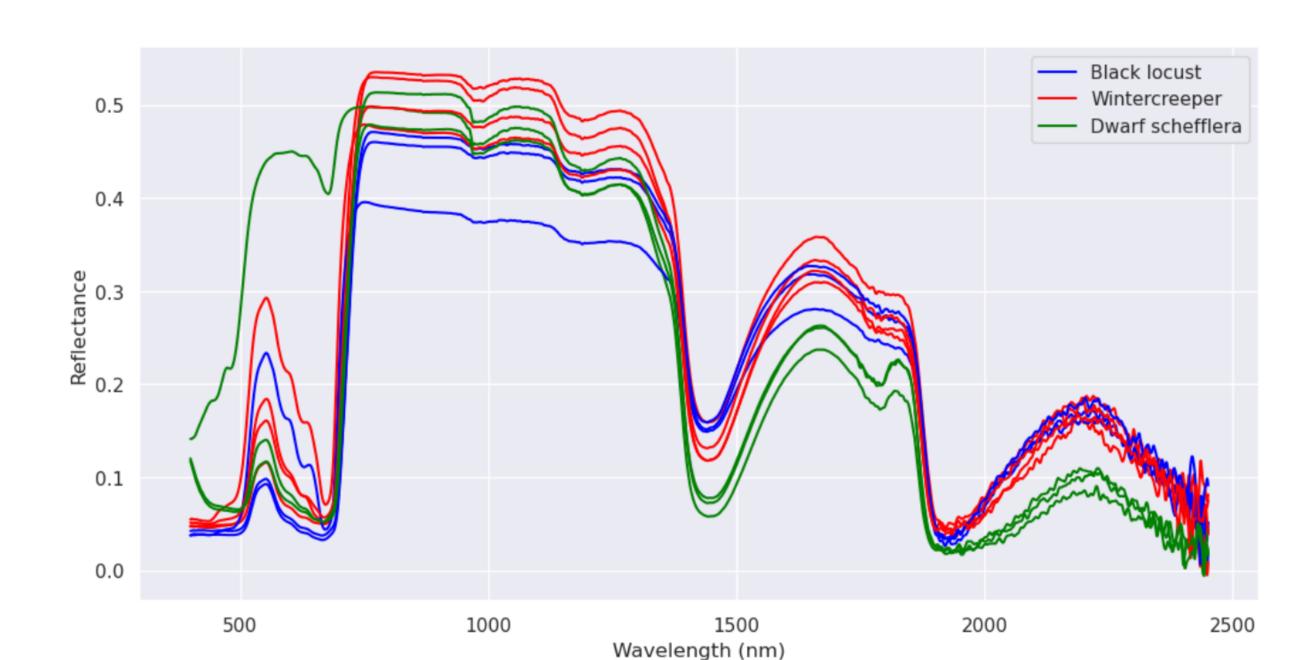
- **Dataset:** consider an unlabelled normalised data set containing N sample points and d features. $\left\{ \times \right\}_{i=1}^{N}$
- How many eigenvalues would the data covariance matrix have? $\chi \in \mathbb{R}^n$

• Suppose we perform PCA and find approximate the data based on its 3 principal components. Let $\Theta \in \mathbb{R}^{d \times r}$ be the associated eigenvectors. What is the projection of a given data point $x^i \in \mathbb{R}^d$ on the space spanned by these eigenvectors?

$$6 = [0, 10, 10] \in \mathbb{R}^{c' \times r}$$

$$r = 3$$

$$\hat{x} = 00^{T}x i \in \mathbb{R}$$



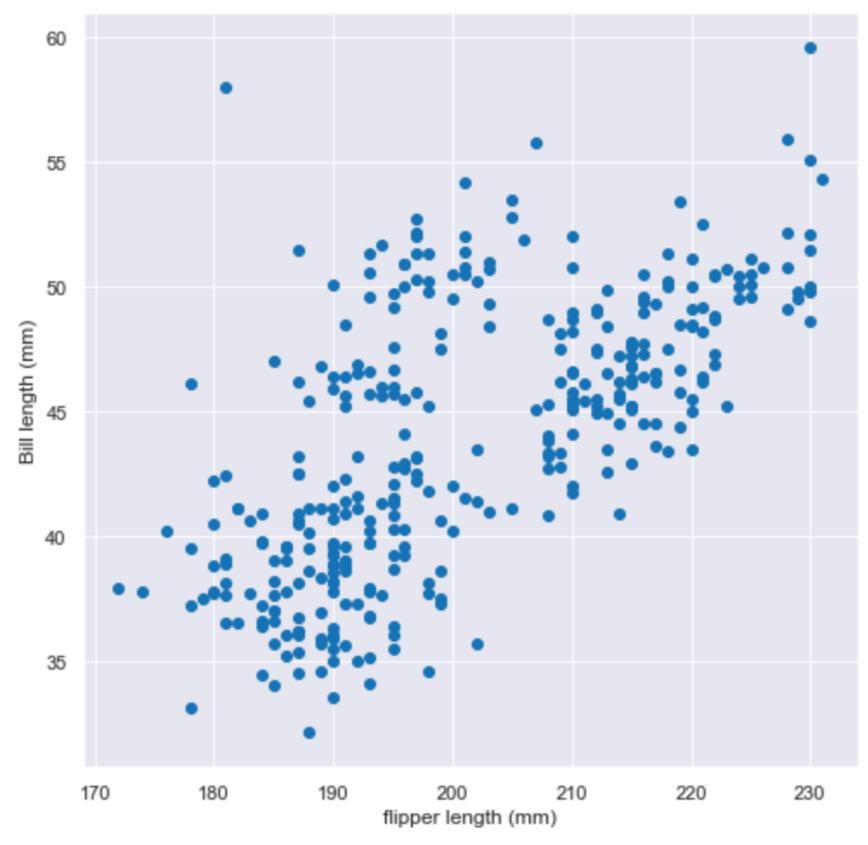


Clustering



k-means Example

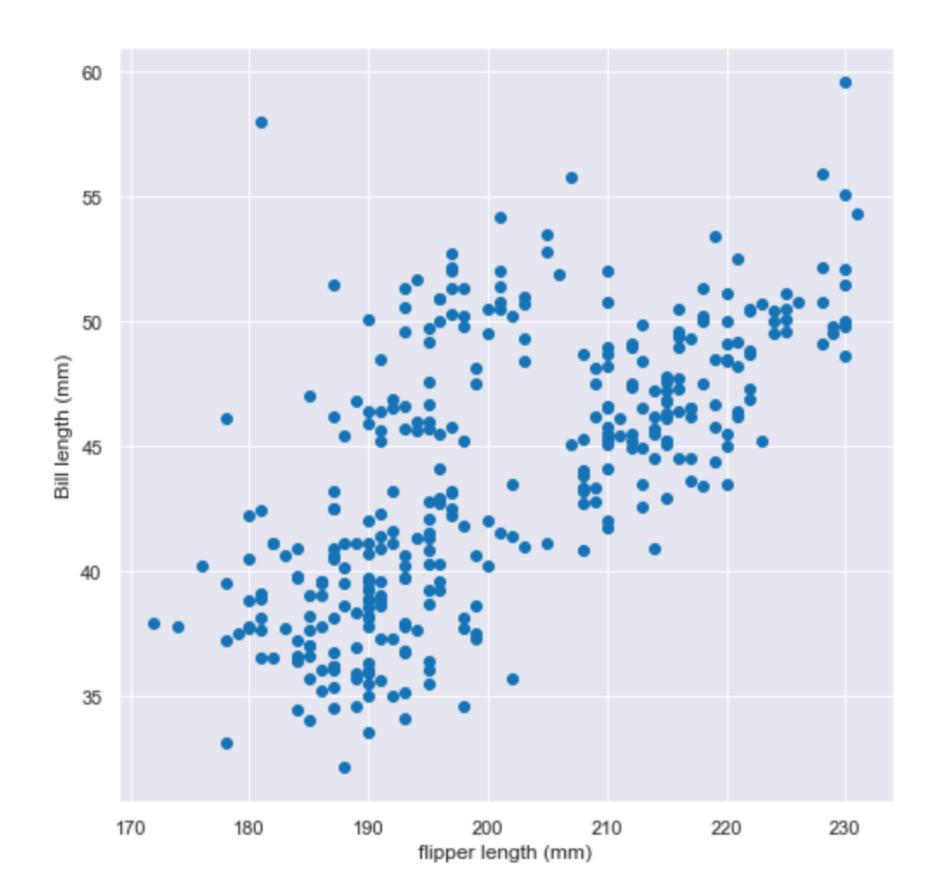
- Palmer Penguins dataset
 - Features: Flipper length and Bill length
 - Can we identify k species based on this data?

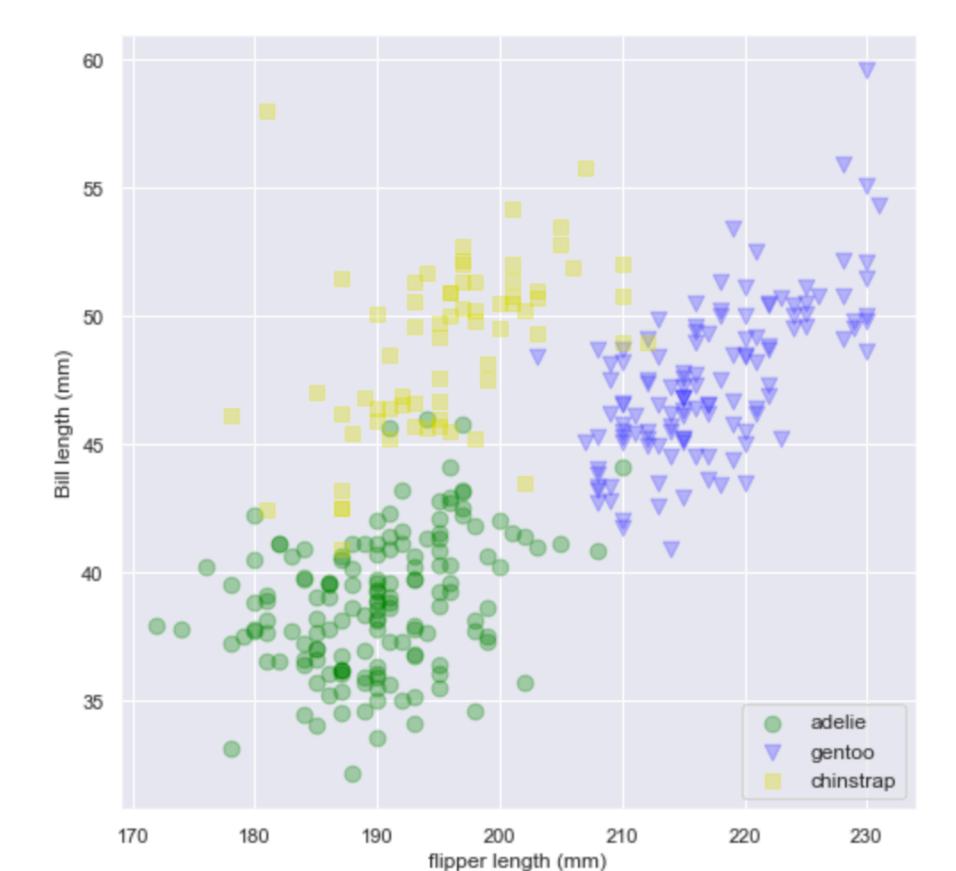




k-means Example

 Note that when we apply clustering, we don't have the labels. Our goal is to be able to best way to cluster the data (a bit vague but we will formalise one way to do this)







k-means approach to clustering

How to cluster data without labels?

Try to find similarity between groups of points

k-means: Group points based on their proximity (in terms of distance in the feature space)

- Given a set of unlabelled input samples, group the samples into k clusters ($k \in \mathbb{N}$)
- k-Means idea:
- Identify k cluster of data points given N samples.
- Find prototype points $\mu^1, \mu^2, \dots, \mu^k \in \mathbb{R}^d$ representing the center of each cluster and add the data points to the nearest cluster



k-means

- A single representative point for data: w and b h n d e e R d that

• Example: suppose we want to have one representative point
• we use a mean-squared
$$\int_{\mathcal{N}} \mathcal{N} = \int_{\mathcal{N}} \mathcal{N} = \int_{\mathcal{N}$$



k-Means

Approach
$$\langle x, y \rangle$$

$$C(i) : cluster \times i \text{ belonge to } . \quad c(i) = \underset{j=1,2,-,k}{\text{argmin}} \left(\left| x^{i} - u^{j} \right| \right|_{2}^{2}$$

Determining the cluster centres to minimize the distance of each point to its

assigned cluster

k-Means heuristic algorithm Loss is non-convex - heurish

Algorithm -

- 1. Initialize $\{\mu^1, \mu^2, \dots, \mu^k\}$ (e.g., randomly)
- 2. While not converged
 - 1. Assign each point x.. to the nearest center (see previous slicke)
 - 2. Update each center μ^{j} based on the points assigned to it $\mu^{j} = \frac{1}{N_{j}} \times \frac{1}{N_{j}}$ and $\mu^{j} = \frac{1}{N_{j}} \times \frac{1}{N_{j}}$
- Step 2.1: For each point x', compute the Euclidean distance to every center

$$\{\mu^1,\mu^2,\ldots,\mu^k\}$$
• Find the smallest distance
$$\lim_{k\to\infty} \|x^k-\mu^k\|_{2}^2 \qquad \text{for each closed points}$$

- The point is said to be assigned to the corresponding cluster (note that each point is assigned to a single cluster)
- Step 2.2:
 - Recompute each center μ^{J} as the mean of the points that were assigned to it



k-means

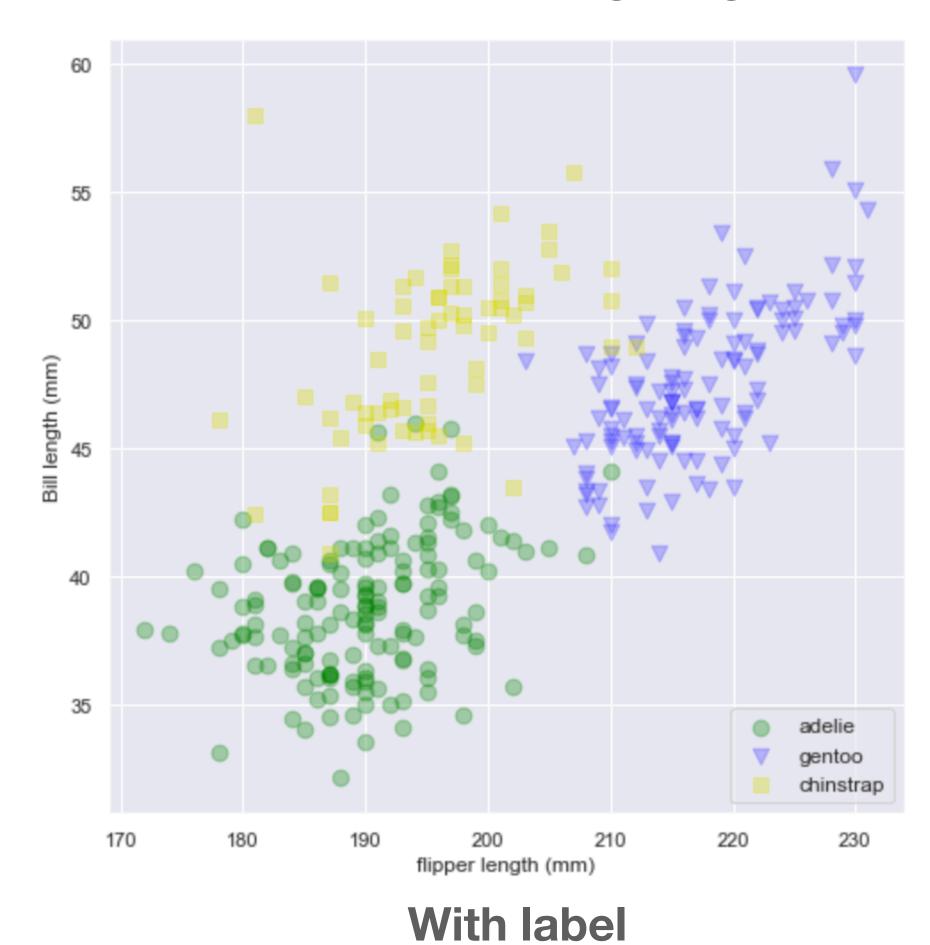
Algorithm - Convergence

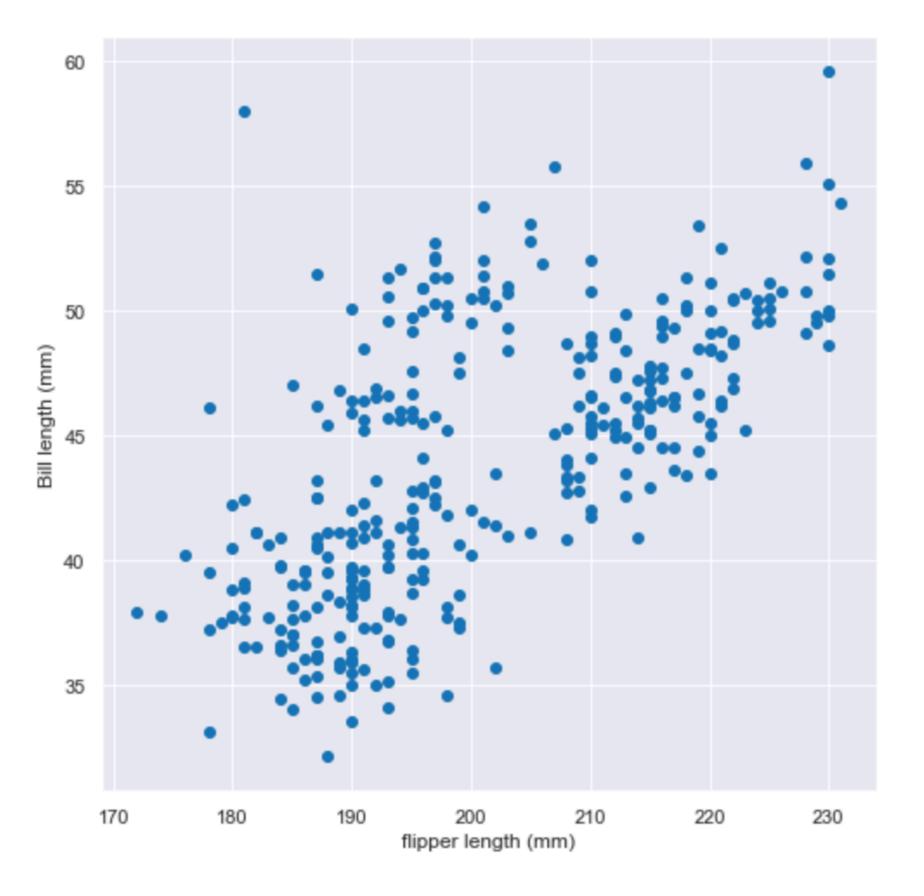
- Step 2 is repeated while k-means has not converged
 - What criteria to stop iterating?
 - Fixed number of iterations? It's arbitrary and a too small number can lead to bad results
 - The difference in assignments or center locations between two iterations can be used as criteria to stop the algorithm
- k-means does not always converge to the best solution
 - Non-convex optimization problem



k-means Example

- Use the Palmer Penguins dataset
 - With Flipper length against Bill length





Without label

k-means Example

(k = 3)3-meani

Jul] 3

Center initialisation

doman

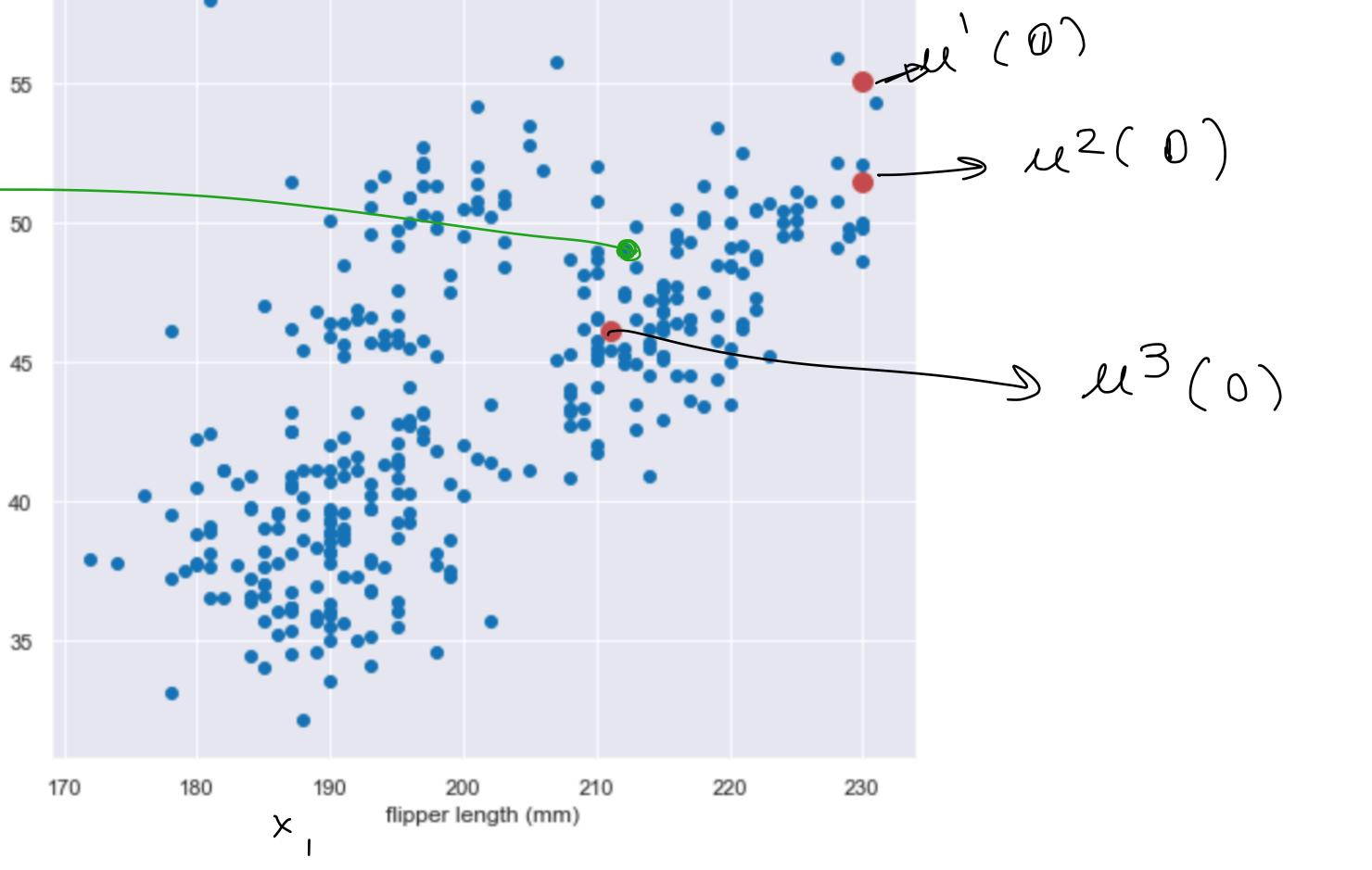
 $Min \int \left(\left(x - M \right) \right)^{3}$ (1 x, - re s// 5 / $\left(\left| x_{1}-\pi_{3}\right| \right)$

× 2 40

60

55

= ((xi - 123 [/ 2

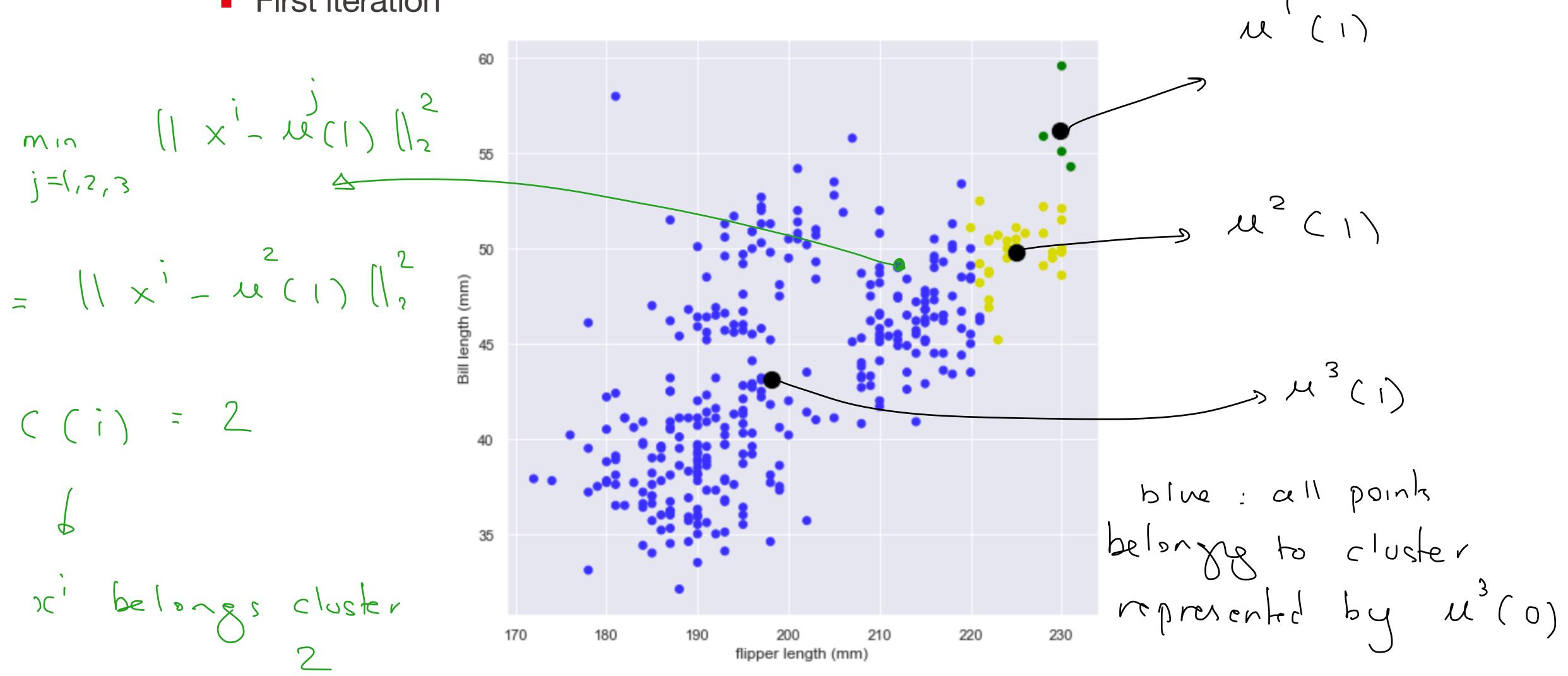


EPFL

k-Means Example

t = 1

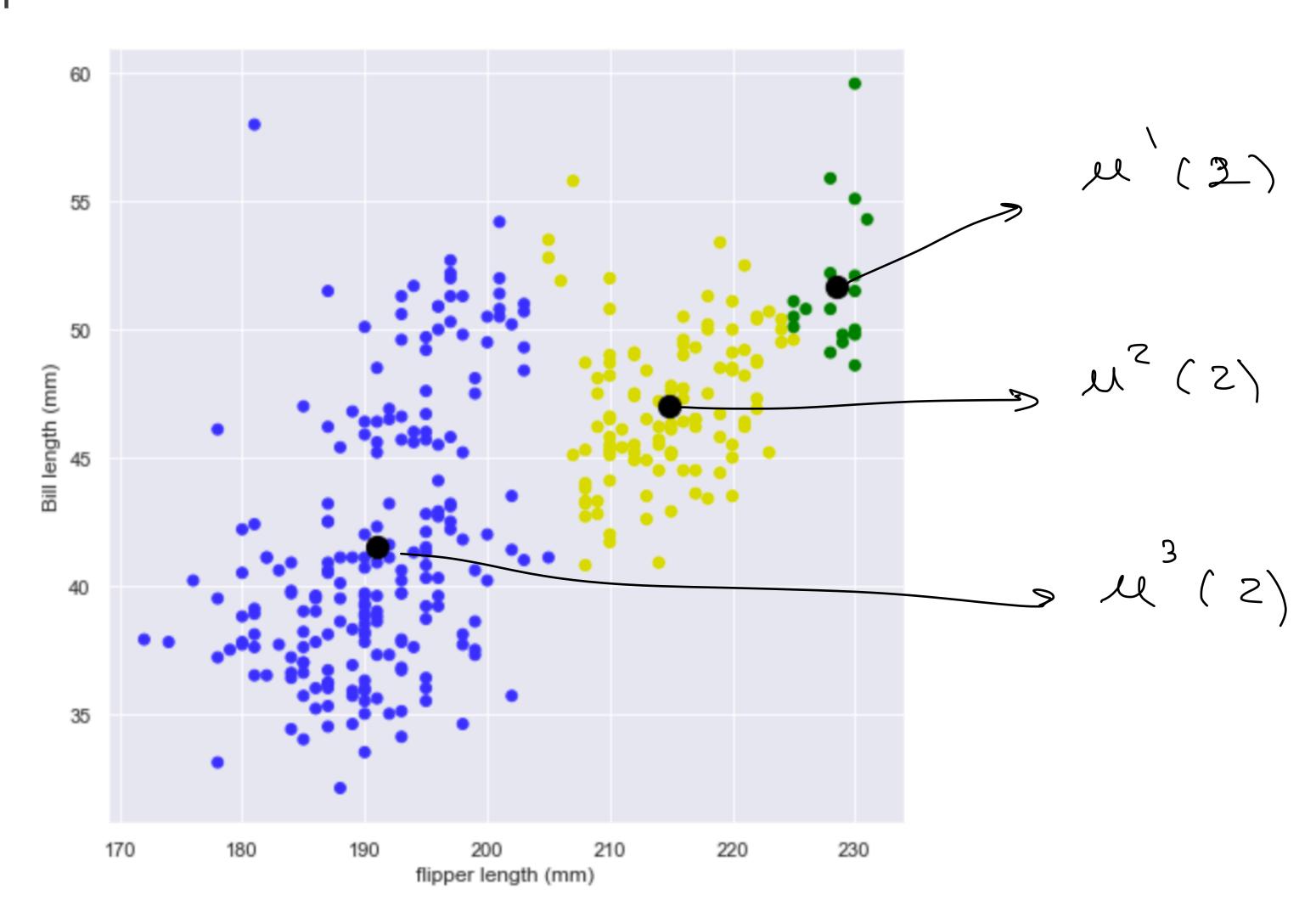
First iteration





k-Means Example

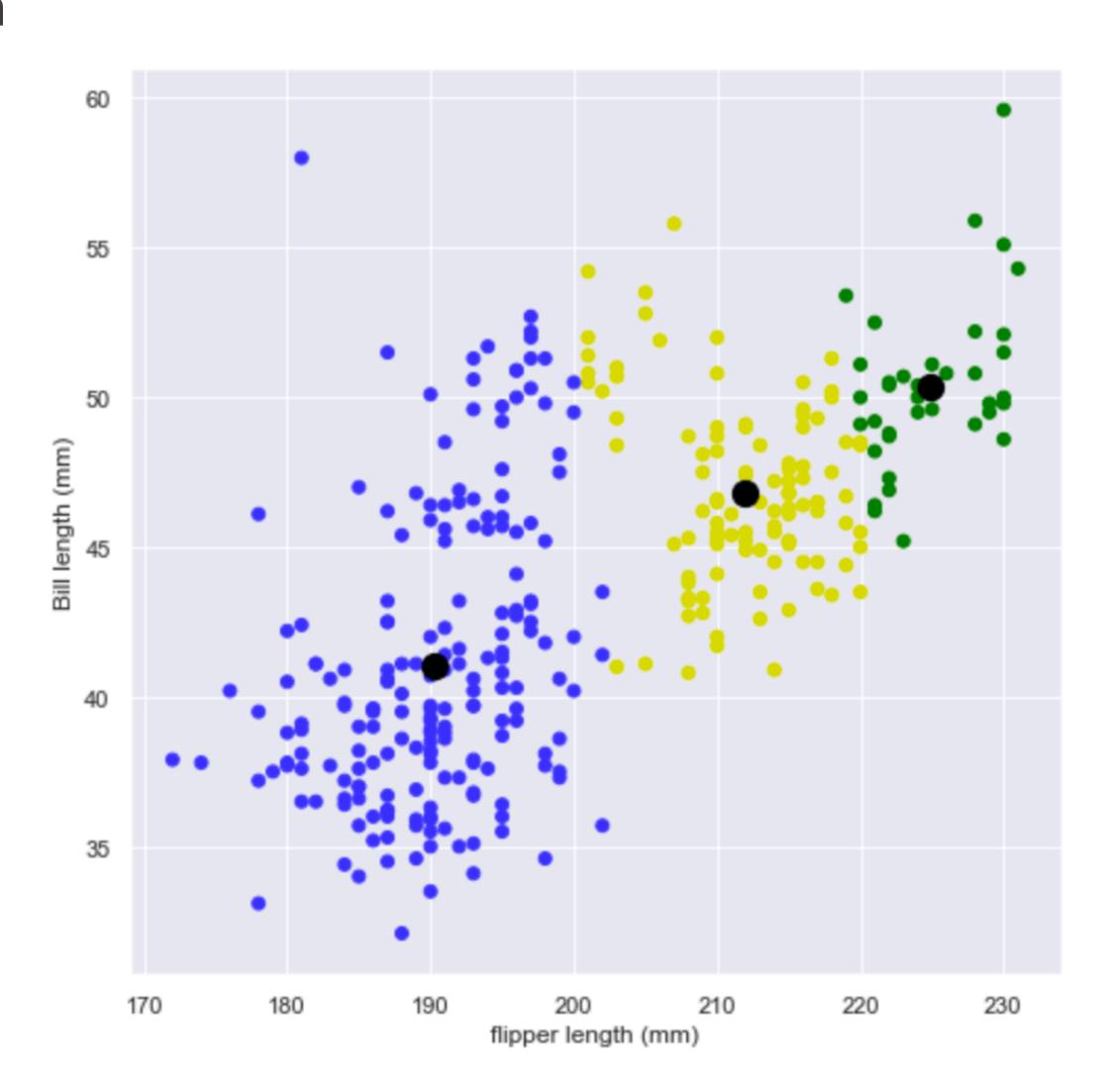
2nd iteration





k-Means Example

3rd iteration



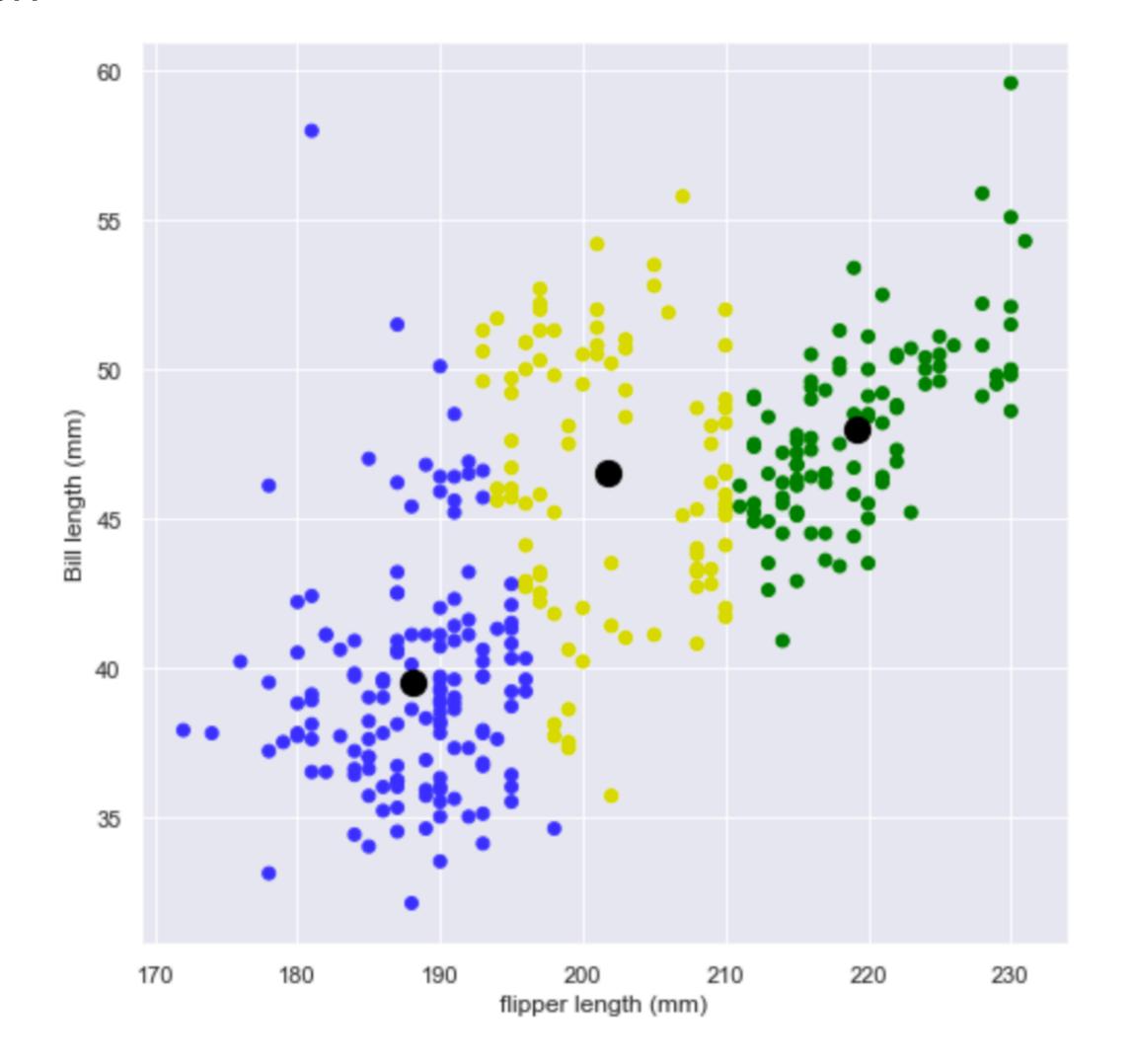


k-Means Example

Last iteration



for j = 1,2,... C





Summary - Clustering

Used for understanding data

Examples:

Topic discovery in a large set of documents

Recommendation engines

Guessing missing entries

k-means: an approach to clustering

Easy to implement and to interpret k-means algorithm usually converges, but possibly to local minima



Brief recap - data statistics Empirical distribution

$$S = \{ blue(b), yellow(y), green(g), red(r) \}$$

empirical dishibution $p: S \rightarrow IR$
 $\sum_{i=1}^{4} p(s_i) = 1, p(s_i) \ge 0$

$$P(b) = \frac{3}{12} = \frac{1}{7}, P(y) = \frac{2}{12}, P(y) = \frac{3}{12}, P(r) = \frac{4}{12}$$

Brief recap - data statistics \(\times \), \(\times \) \\ \(\times \) \\ \(\times \)

$$M_j = \frac{1}{N} \sum_{i=1}^{N} x_i^i$$
 : empirical mean of feature j

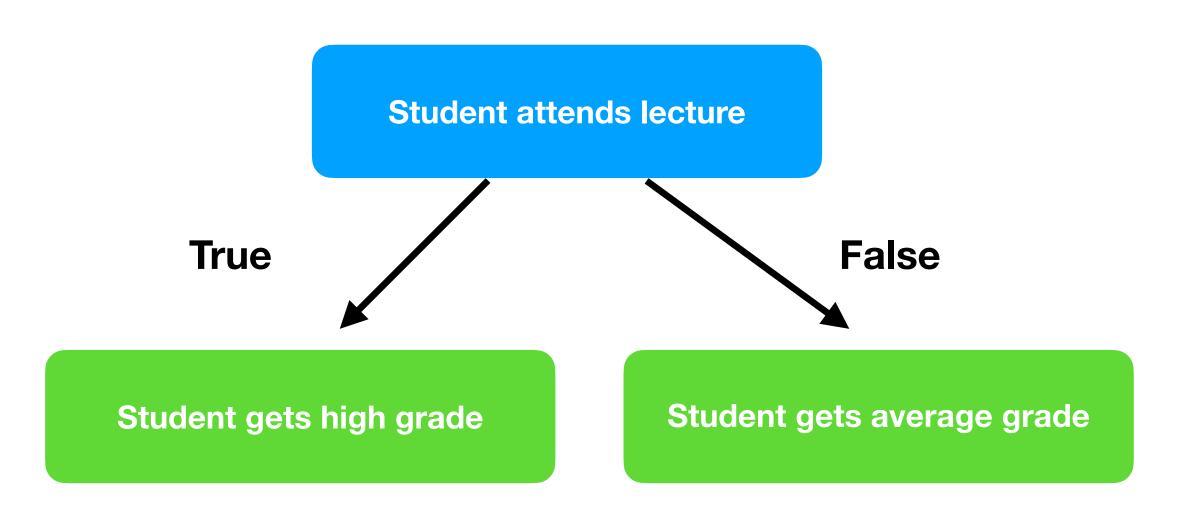
$$Cov(x_j, x_j) = \frac{1}{N-1} \sum_{i=1}^{N} (x_j^i - \mu_j^i) (x_j^i - \mu_j^i) : empirical covariance$$

Corr
$$(x_j, x_j) = \frac{(ov(x_j, x_j))}{(ov(x_j, x_j))}$$
 : empirical concludes between feature $j \not\in j$



Decision Trees: an approach to supervised learning Introduction

What's a decision tree?



Split data based on the answer to a series of Yes/No questions

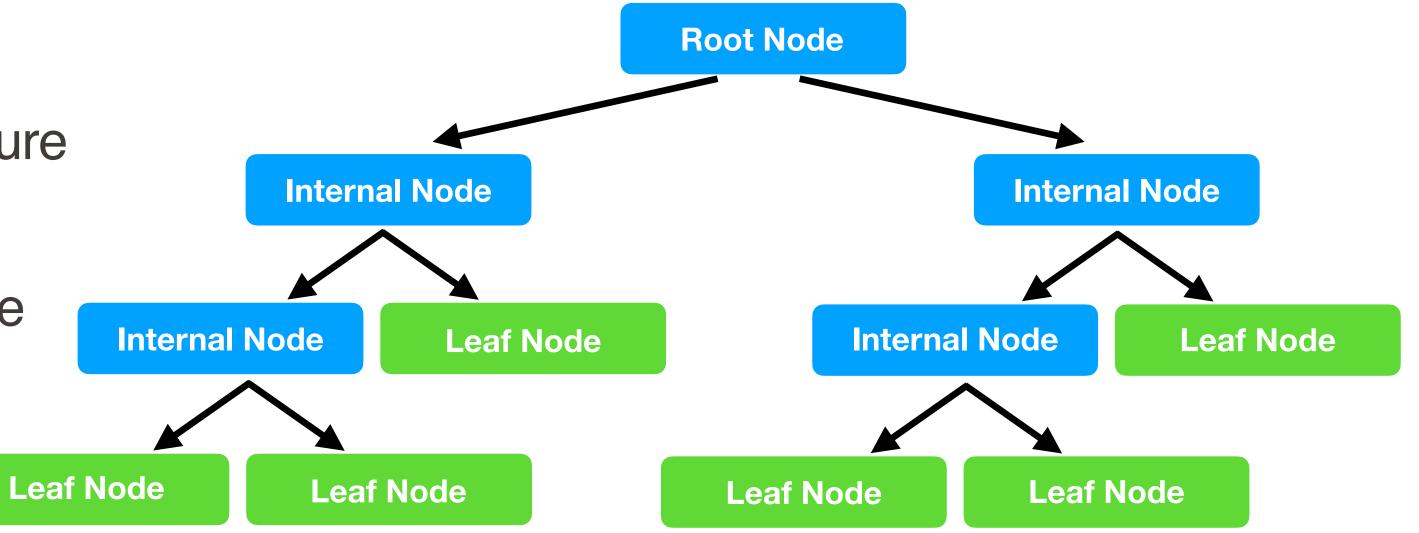


Decision Trees Terminology

Nodes are checked on a single feature

Branches are feature values

Leaves indicate prediction of the tree



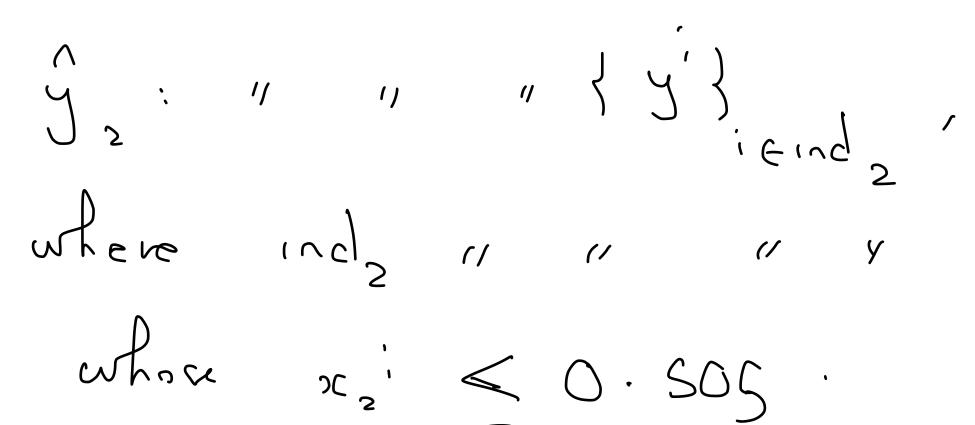


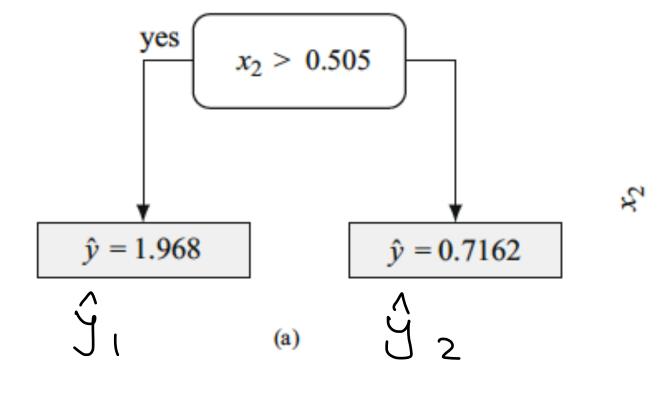
Regression tree

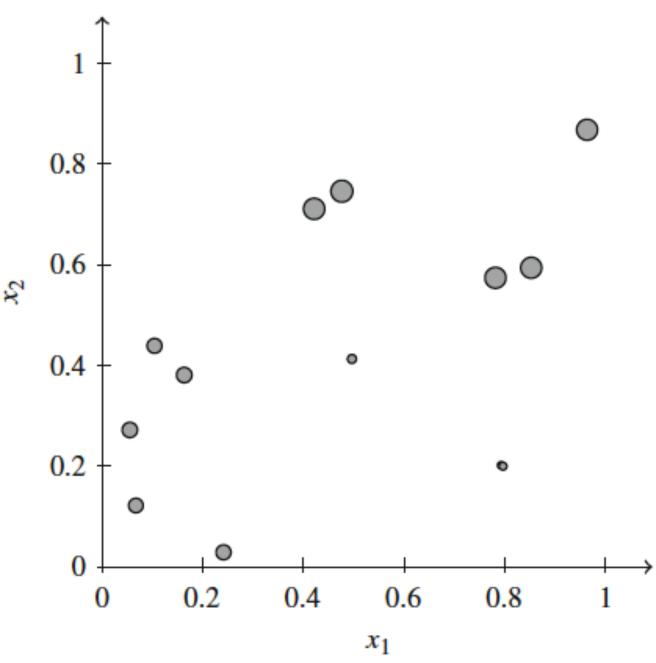
Example: how to predict the label of an unseen point based on the data?

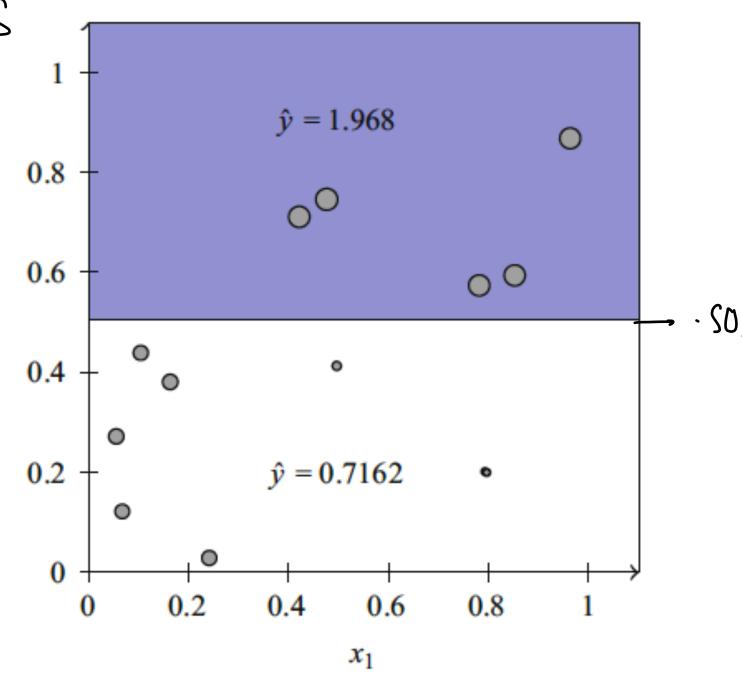
ý. : average values of 1y'}, eind,

where ind, is inclined of points whose x' > 0.50S







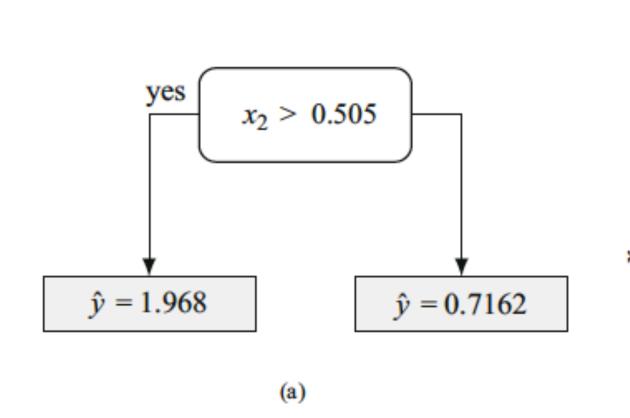


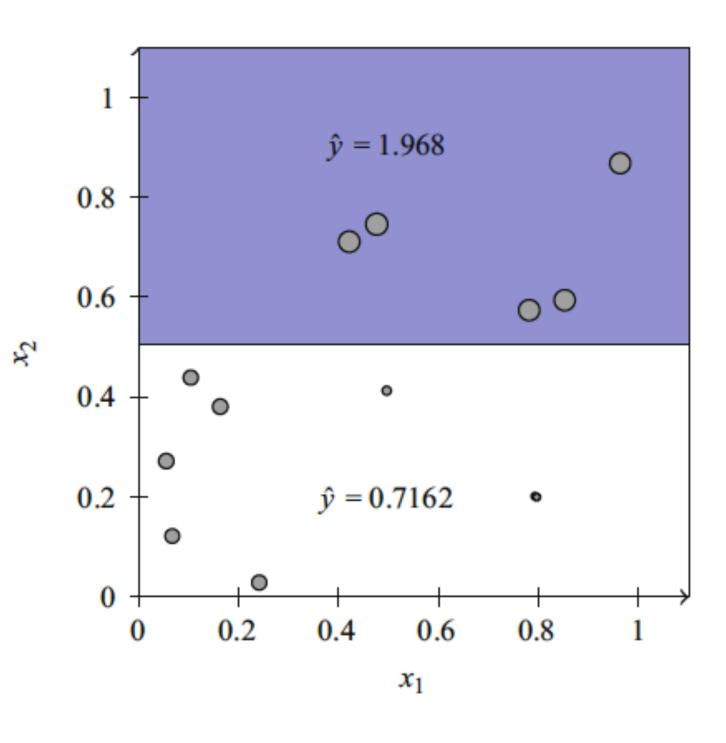
Loss function

Suppose this is our prediction model. What would be the mean-square loss?

Let ind1 and ind2 denote the set of indices of data points in the top and the bottom regions, respectively

$$\frac{1}{N} \lesssim (y'-1.968)^2 + 1.968$$

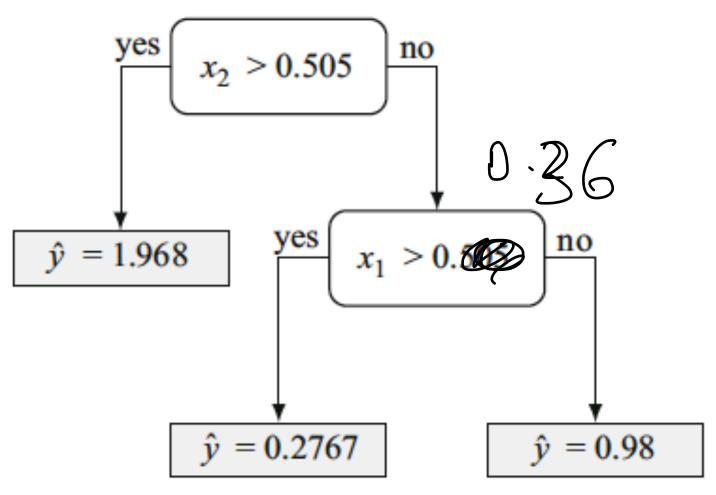


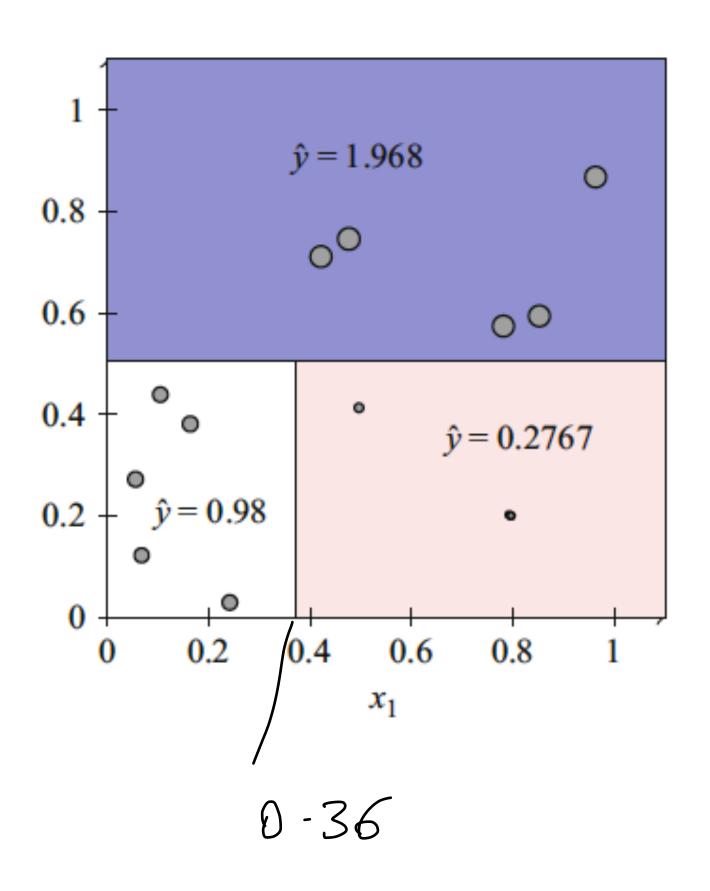


Regression tree example

Let ind2 ind3 denote the set of indices of data points in the left and right region, and ind1, as before, denote the set of indices of data points in the top region

Loss:
$$\sqrt{\frac{1}{N}} \left[\frac{1}{2} + \frac{1}{16} \cdot \frac{1}{100} + \frac{1}{100} \cdot \frac{1}{100} \right]^2 + \frac{1}{16} \cdot \frac{1}{100} \cdot \frac{1}{10$$





Regression tree Loss function

To construct the tree, we have to decide the depth of the tree. And at each node, we have to decide which feature to use for a split

what value of the feature to use for the split

min
$$\geq (y'-y_{NO})^2 + \geq (y'-y_{Ves})^2$$

i.e., $(y'-y_{Ves})^2$

where indq: \landices of x; \le S\rangle, \frac{1}{y} \text{ average of correspondy}

indz: \landices \text{ inclines } \text{ x; } > S\rangle

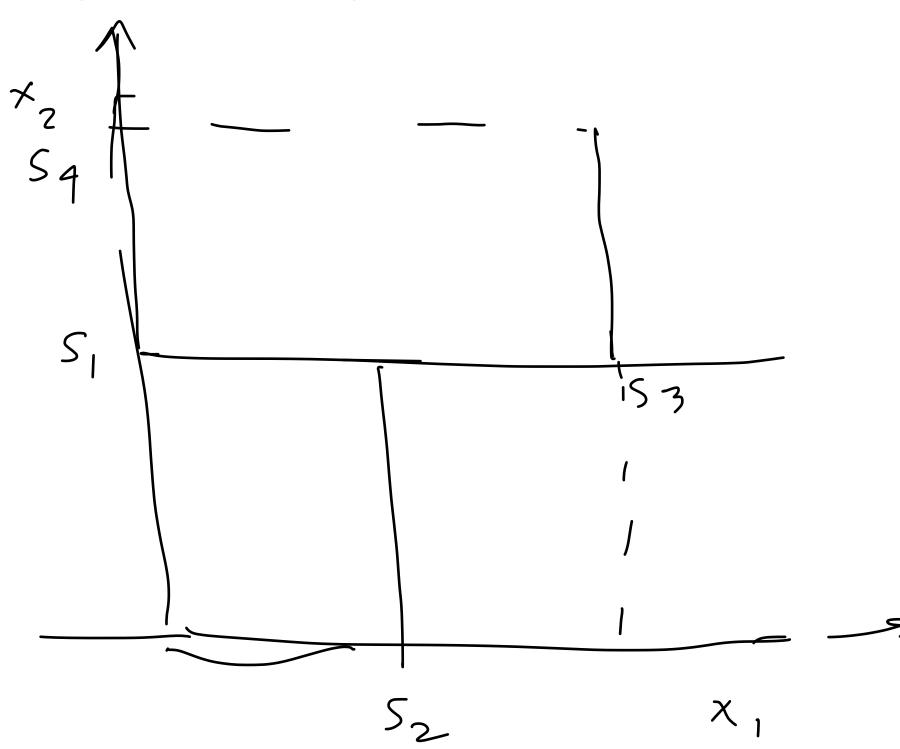
to ophnize over j. s. its non-convex.



Regression Tree Optimising the loss function - greedy algorithm

The global loss function would tell us which features and what threshold values to use for the entire tree of a given depth. This function is non-convex and hard to optimize.

We apply a greedy approach, where we start adding nodes based on optimizing the loss function at each node. Note that even the loss function at a given node is non-convex and we can only approximately find a solution





Regression tree Stopping criteria

At which depth should we stop growing the tree?

If the number of data points at leaf nodes become too small, then we often stop

If the loss at a depth d1 is close to the loss at depth d2 = d1+1, then we often stop

As the objective is non-convex and greedy approach gives us an approximate solution, it's possible that

the loss would decrease after a few iterations, even though in the next iteration it did not decrease

EPFL

Decision Trees for classification

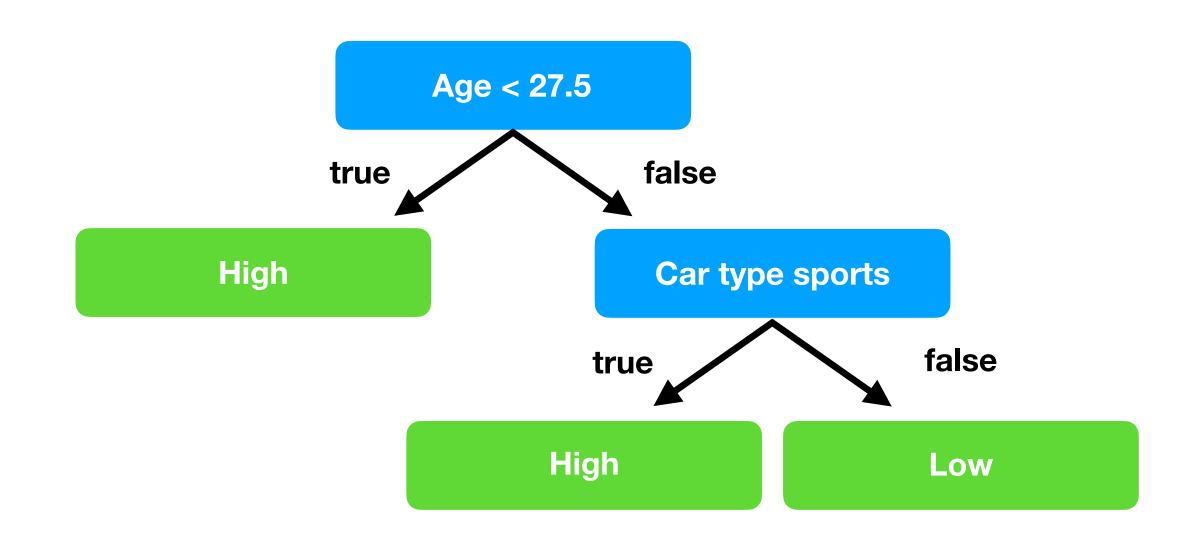
Called classification trees

$$\left\{\begin{array}{c} X \\ X \end{array}\right.$$

$$y' \in \{h, gh, low\}$$

| Continuous | Categorical | Class |
|----------------|-------------|-------|
| Feature | Feature | label |

| Age | Car type | Risk |
|-----|----------|------|
| 23 | family | high |
| 17 | sports | high |
| 43 | sports | high |
| 68 | family | low |
| 32 | family | low |
| 20 | family | high |



EPFL Classification trees

Performance metric and loss function

Given a classification tree, we can evaluate its performance by forming the confusion matrix We can measure the accuracy, error rate, etc.

However, finding the tree that minimises a given metric is a hard optimisation problem

To address this, we use an alternative measure of performance and use a greedy approach based on this measure



Classification Trees

Example - constructing the tree

Let's consider a tree of depth 2. We have to address

Which feature to use at each depth to do a split?

For the continuous feature, at what value to do a split?

For the categorical feature, which category to use for the split?

| Continuous | Categorical | Class |
|----------------|----------------|-------|
| Feature | Feature | label |

| Age | Car type | Risk |
|-----|----------|------|
| 23 | family | high |
| 17 | sports | high |
| 43 | sports | high |
| 68 | family | low |
| 32 | family | low |
| 20 | family | high |

EPFL Classification trees

Greedy approach: choose a feature and the split sequentially based on minimising a performance metric, for example, the **Gini impurity** of a node

Gini impurity of a leaf node: based on empirical probability of class

$$9 = \sum_{l=1}^{K} \sum_{l'\neq l} P_{l} P_{l'} = \sum_{l=1}^{K} P_{l} (1-P_{l})$$

$$g = (P_1P_2 + P_1P_3) + (P_2P_1 + P_2P_3) + (P_3P_1 + P_3P_2)$$

$$g = P_1P_2 + P_2P_1 = 2P_1P_2$$

Gini impurity of a node

Weighted sum of the gini impurity of the two leaf nodes associated with the node

Note: Criteria other than gini index (such as entropy) are also used for node split



Classification Trees

Example - constructing the tree

We will apply Gini impurity to construct the classification tree

| Continuous Feature | Categorical Feature | Class label |
|-----------------------|------------------------|----------------|
| Age | Car type | Risk |
| 23 | family | high |
| 17 | sports | high |
| 43 | sports | high |
| 68 | family | low |
| 32 | family | low |
| 20 | family | high |

EPFL Classification trees

Criteria for choosing a feature and a split

| Continuous | Categorical | Class |
|----------------|----------------|-------|
| Feature | Feature | label |

| Age | Car type | Risk |
|-----|----------|------|
| 23 | family | high |
| 17 | sports | high |
| 43 | sports | high |
| 68 | family | low |
| 32 | family | low |
| 20 | family | high |

