# California Dreamin'
## An analysis on housing in California

**Gabriele Cialdea, Leonardo Filippone, Pietro Pianini**

# Introduction

## Research Question ?

How can we predict the **expected value of an house** basing on the **features of the neighbourhood**?

Analysis of the **structure** and the **variability of the sample**

Selecting a **prediction model** assessing several **supervised classification methodologies**

# Agenda

**California Dreamin'**

Gabriele Cialdea, Leonardo Filippone, Pietro Pianini

# Data Description

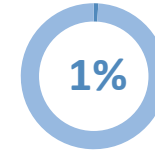## Source and dimension of the dataset and attributes of the observations

### Source

Kaggle dataset based on 1990 census on housing in California

### # observations

## 20.640[1]

### 1%   Missing values (artificially added)

For each observation the following **attributes**:

- Longitude
- Latitude
- Housing median age
- Total rooms
- Total bedrooms

- Population
- Households median income
- Median house value
- Ocean proximity (artificially added)

---

**California Dreamin'**

Gabriele Cialdea, Leonardo Filippone, Pietro Pianini

**Note**: (1) Observations relative to single census block groups

# Data Manipulation and Selection

## Selection and manipulation process

### Step 1

**Cleaning** the dataset from **missing values**

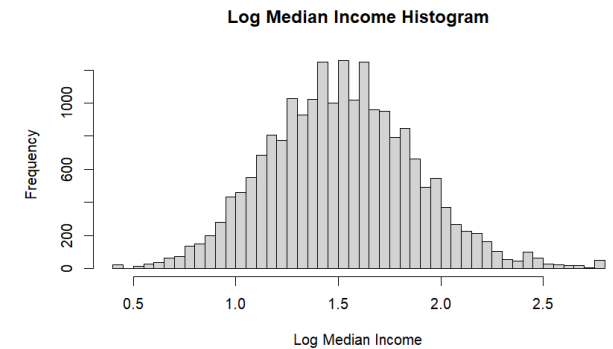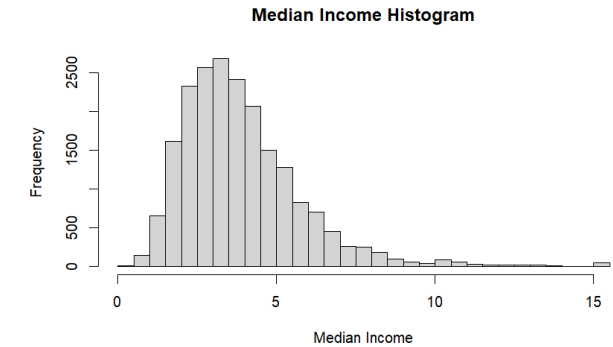**Final number of observations:** **20.433**

### Step 2

**Log-transformed** the values of the **most skewed attributes**

### Step 3

**Scaled values** for all the observations and **cathegorized the target attribute**

## 🔍 Focus on Step 2



Median Income Histogram



Log Median Income Histogram

**California Dreamin'**

Gabriele Cialdea, Leonardo Filippone, Pietro Pianini

# Agenda

**California Dreamin'**

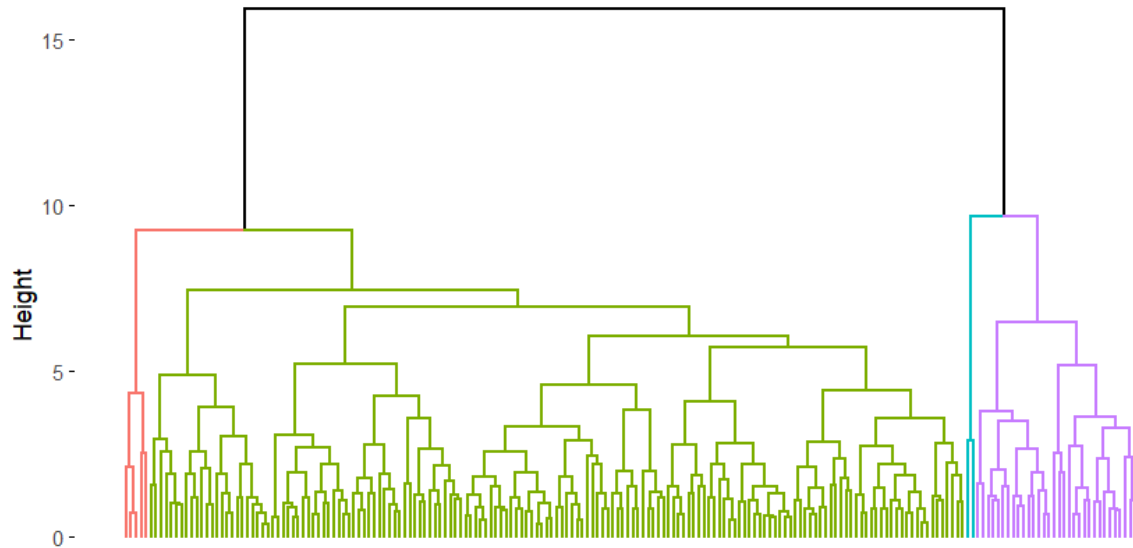Gabriele Cialdea, Leonardo Filippone, Pietro Pianini

# Hierarchical Clustering
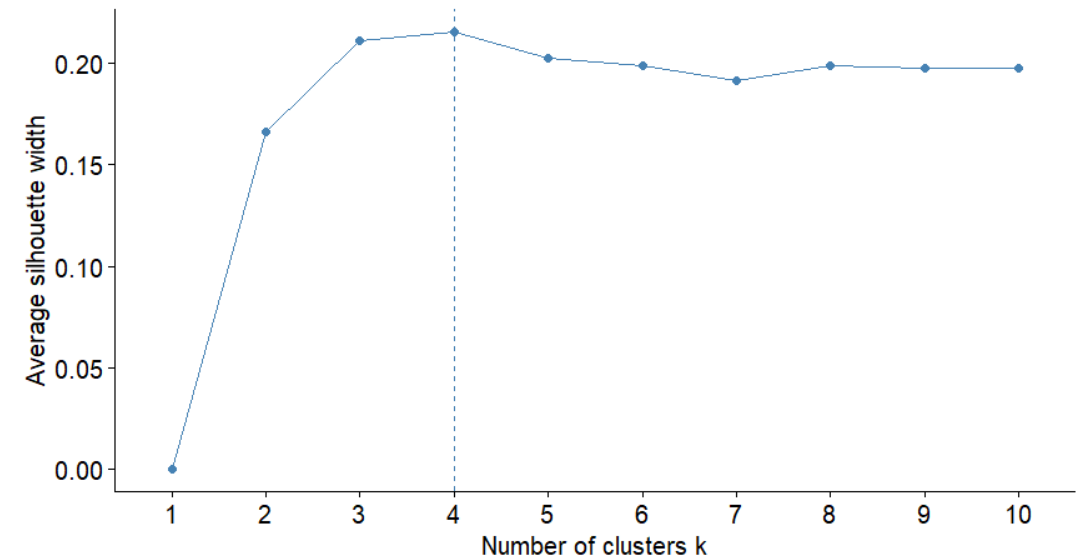
## Dendogram representation

Euclidean-Complete with eclust



| Number of clusters | 4 |
| --- | --- |

## Silhouette evaluation

Optimal number of clusters
Silhouette method AHC



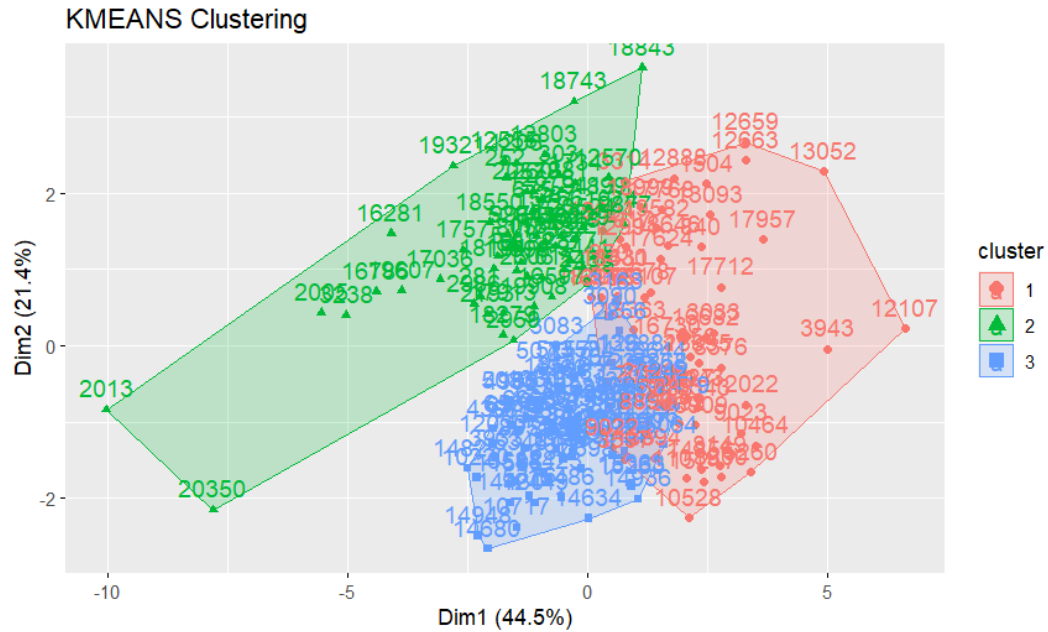| Average Silhouette Width | 0,22 |
| --- | --- |

**California Dreamin'**

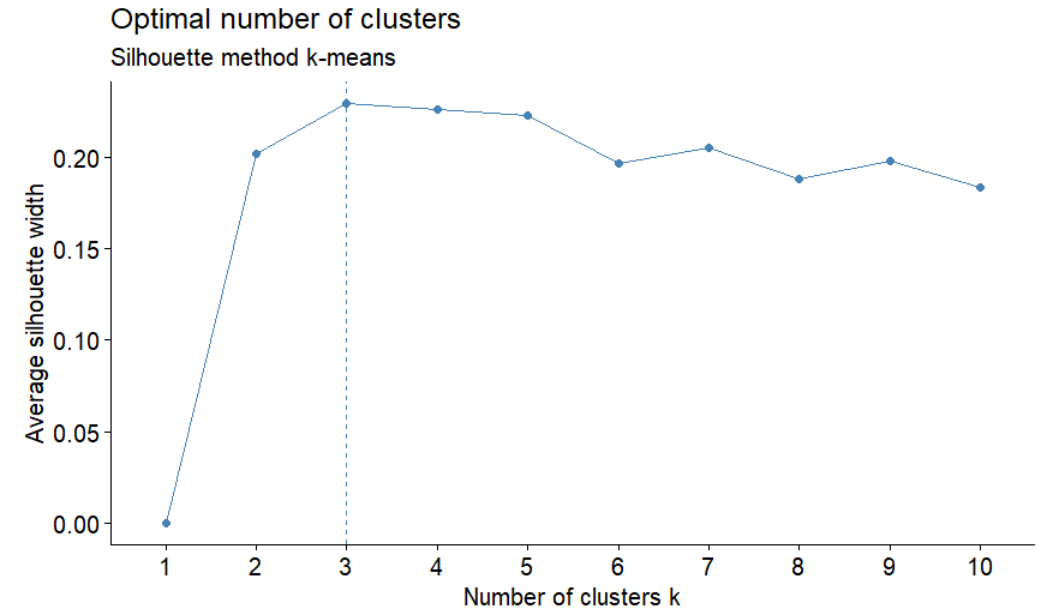Gabriele Cialdea, Leonardo Filippone, Pietro Pianini

# K-Means Clustering

## K-Means cluster representation



**Number of clusters** **3**

## Silhouette evaluation



**Average Silhouette Width** **0,23**

**California Dreamin'**

Gabriele Cialdea, Leonardo Filippone, Pietro Pianini

# Agenda

**California Dreamin'**

Gabriele Cialdea, Leonardo Filippone, Pietro Pianini

# Optimal number of components

## Screeplot

## Cumulative PVE

# Principal components representation

## Biplot
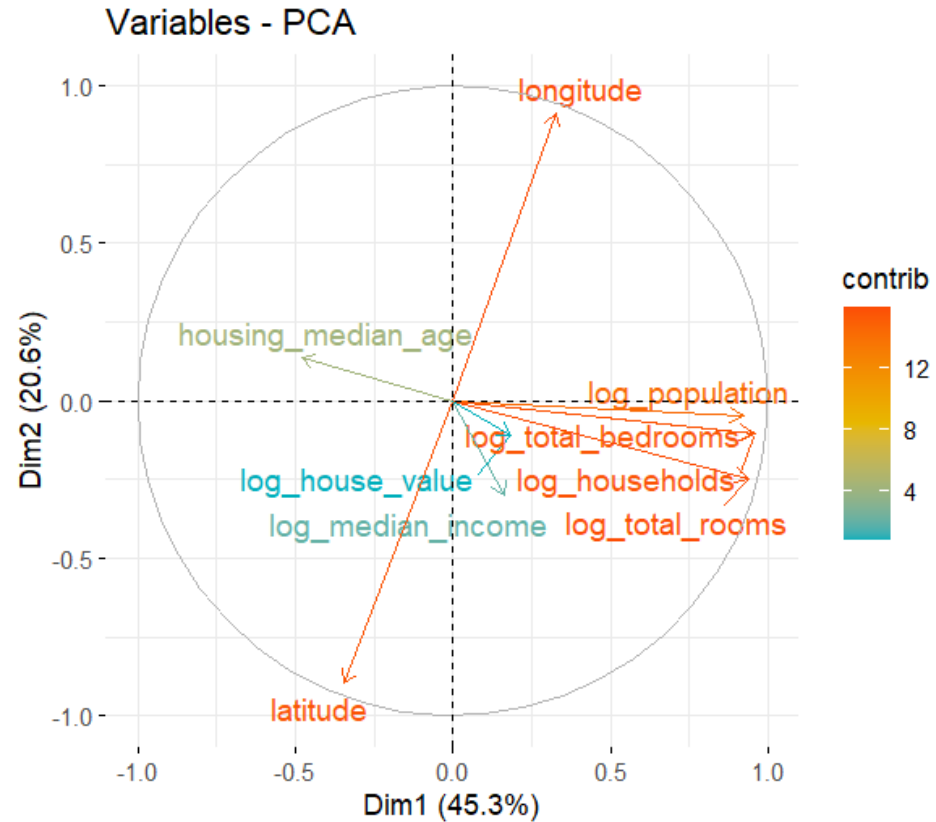


## Corrplot

# Agenda

**California Dreamin'**

Gabriele Cialdea, Leonardo Filippone, Pietro Pianini

# Main Classifiers Overview

**Decision tree**

① 

Decision Tree

**Random forest**

② 

Data set

Decision Tree-1    Decision Tree-2    Decision Tree-N

**Gradient boosting**

③ 

Error

Iterations

**Neural networks**

④ 

**Logistic regression**

⑤ 

Y

1

0.5

0

S-Curve

y=0.8

Threshold Value

y=0.3

X

**LDA**

⑥ 

**California Dreamin'**

Gabriele Cialdea, Leonardo Filippone, Pietro Pianini

# Cross Validation Methodology and Evaluation Metric

## Cross validation



Dataset

Training | Testing | Holdout Method

Cross Validation

Data Permitting:

Training | Validation | Testing | Training, Validation, Testing

| Training Set | 80% |
|---|---|
| Test Set | 20% |

## Evaluation metric

| Metric | Formula |
|---|---|
| True positive rate, recall | $\dfrac{TP}{TP+FN}$ |
| False positive rate | $\dfrac{FP}{FP+TN}$ |
| Precision | $\dfrac{TP}{TP+FP}$ |
| Accuracy | $\dfrac{TP+TN}{TP+TN+FP+FN}$ |
| F-measure | $\dfrac{2 \cdot precision \cdot recall}{precision + recall}$ |

**California Dreamin'**

Gabriele Cialdea, Leonardo Filippone, Pietro Pianini

# Results

## Performance scores and ranking of the main classifiers

**CLASSIFIERS PERFORMANCES**



- **Gradient Boosting** is recognized as the **best-fitting classifier**

- These the **optimal hyperparameters**
  - Learning rate: 0.01
  - Number of estimators: 300
  - Max depth: 5

**California Dreamin'**

Gabriele Cialdea, Leonardo Filippone, Pietro Pianini

# Agenda

**California Dreamin'**

Gabriele Cialdea, Leonardo Filippone, Pietro Pianini

# Further research

More recent data and/or extension of the analysis to **other geographical areas** (e.g. robustness in other States of the U.S.)

**Different and more complex methodologies** for the supervised classification

Employment of **more explanatory attributes** for the observations in the prediction model

# References

- Boehmke, Bradley; Greenwell, Brandon (2019). Hands-On Machine Learning with R. Chapman & Hall.

- James, G., Witten, D., Hastie, T., Tibshirani, R. (2021). An Introduction to Statistical Learning. Springer Texts in Statistics. Springer, New York, NY

- Pace, R. Kelley, and Ronald Barry. "Sparse spatial autoregressions." Statistics & Probability Letters 33.3 (1997): 291-297.