JADS
Jheronimus
Academy
of Data Science

TU/e EINDHOVEN
UNIVERSITY OF
TECHNOLOGY

# Wine Quality Classification

Hands-on

Introduction to Data Science Elective

On behalf of the PDEng Data Science

5 November 2019

# Agenda

**Yesterday:**

Exploratory analysis of the Wine Quality dataset.

**Today:**

Machine learning: wine quality classification.

# Learning from data

# Inputs and Outputs

- **Inputs**
  - Words, sentences
  - Images, videos
  - Sensor observations, time series
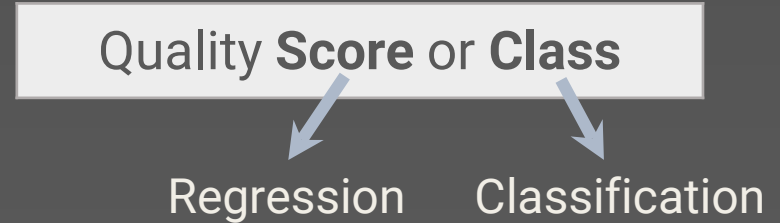  - Voice

$x$

- **Outputs**
  - Class label

$y$

# Inputs and Outputs in the Wine Data

## Inputs

| | |
|---|---|
| Fixed acidity | Total sulfur dioxide |
| Volatile acidity | Density |
| Citric acid | pH |
| Residual sugar | Sulphates |
| Chlorides | Alcohol |
| Free sulfur dioxide | |

## Outputs

Quality **Score** or **Class**

Regression     Classification

# What is learning?

- **Learning** or **training** refers to estimate the parameters of a model.

<div style="border:1px solid black; background-color:#FCD462; text-align:center; font-weight:bold;">Dataset</div>
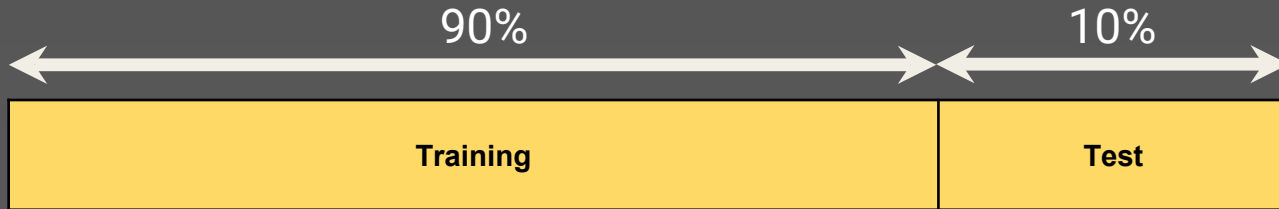
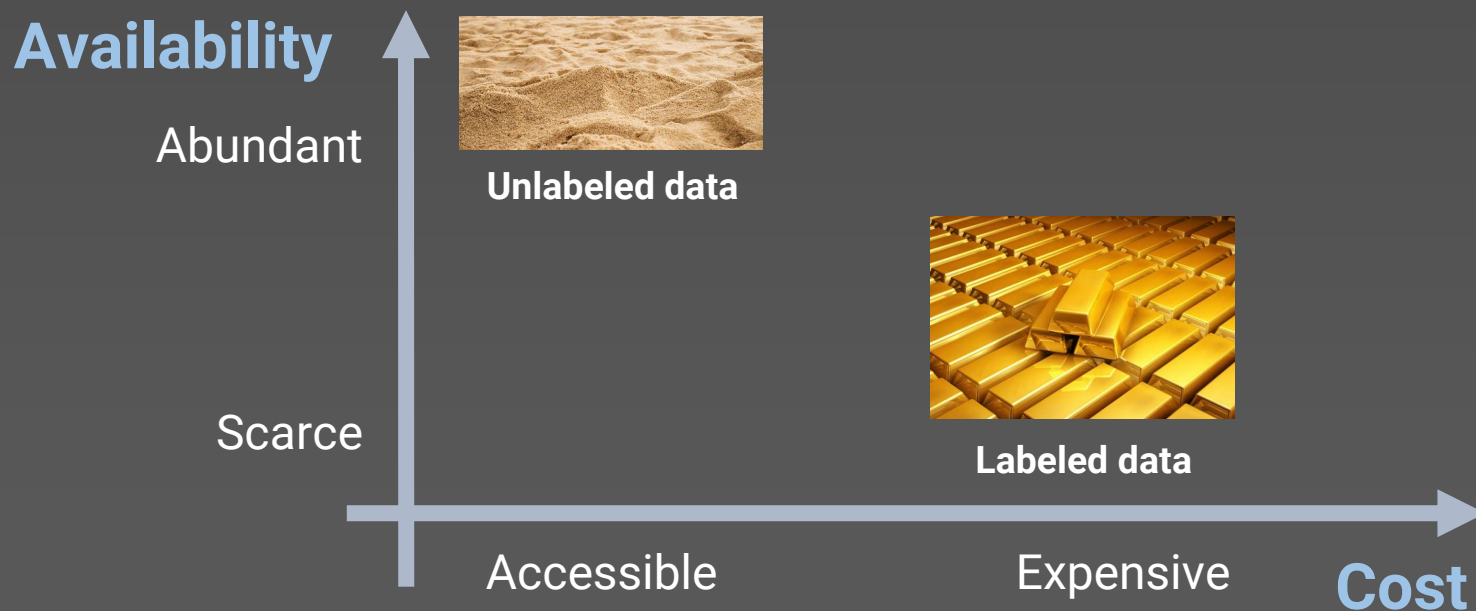# Machine learning is about generalizing to **unseen data.**

# Split Data

- Split dataset into a training and test set:

```
x_train, x_test, y_train, y_test = train_test_split(data_features,
data['grade'], test_size=0.1)
```
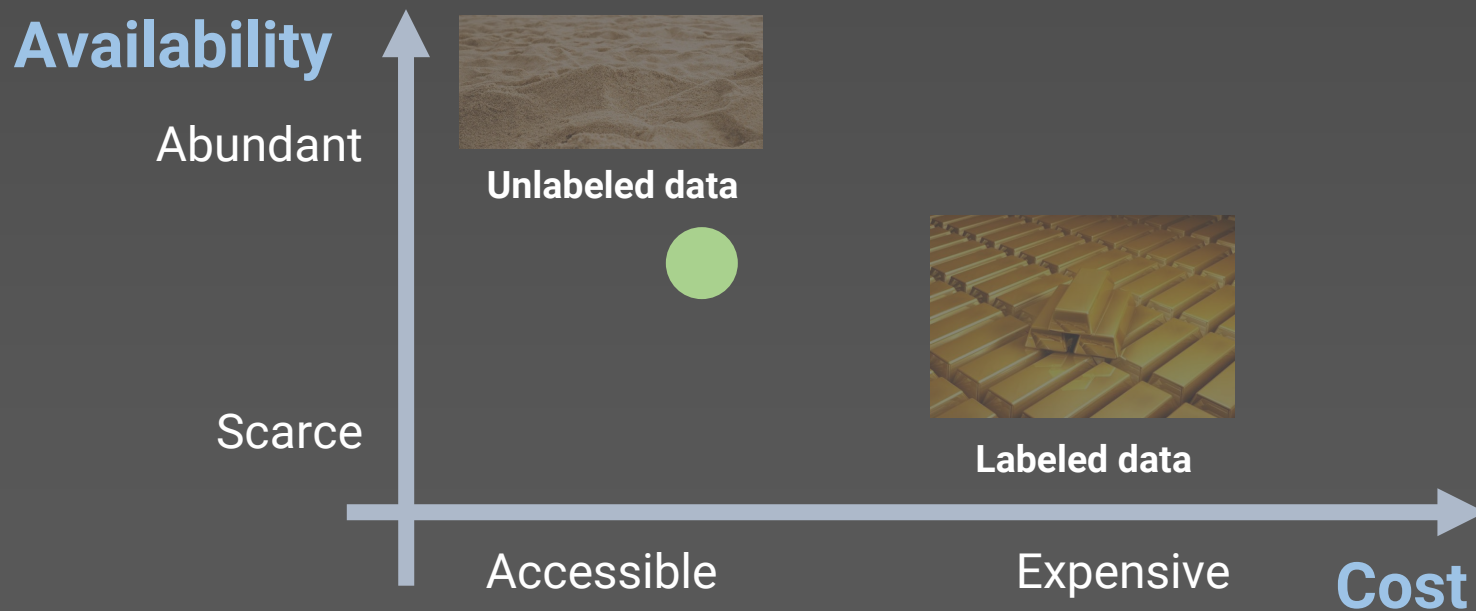
# What data is out there?

**Availability**

Abundant



**Unlabeled data**



**Labeled data**

Scarce

Accessible        Expensive        **Cost**

# Data in the PDEng Data Science

**Availability**

Abundant



**Unlabeled data**



**Labeled data**

Scarce

Accessible          Expensive          **Cost**

# Types of learning

- Supervised learning
- Unsupervised learning
- Reinforcement learning

# LeCun's Cake Analogy



How Much Information is the Machine Given during Learning?

Y. LeCun

▶ **"Pure" Reinforcement Learning (cherry)**
  ▶ The machine predicts a scalar reward given once in a while.
  ▶ **A few bits for some samples**

▶ **Supervised Learning (icing)**
  ▶ The machine predicts a category or a few numbers for each input
  ▶ Predicting human-supplied data
  ▶ **10→10,000 bits per sample**

▶ **Self-Supervised Learning (cake génoise)**
  ▶ The machine predicts any part of its input for any observed part.
  ▶ Predicts future frames in videos
  ▶ **Millions of bits per sample**

© 2019 IEEE International Solid-State Circuits Conference    1.1: Deep Learning Hardware: Past, Present, & Future    59

# Supervised Learning
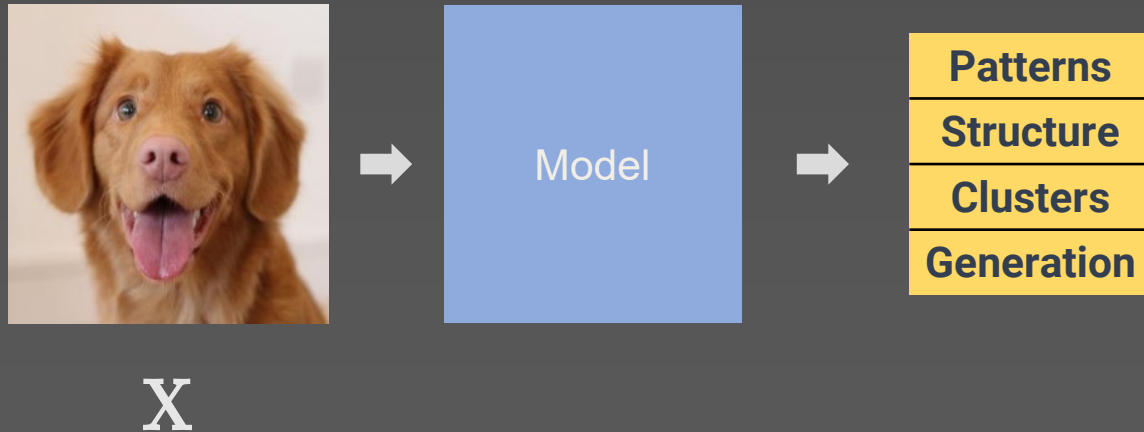
- Given a dataset D of inputs x and <u>labeled</u> targets y, *learn* **to predict** y from x.



- Most successful paradigm in machine learning.

# Unsupervised Learning
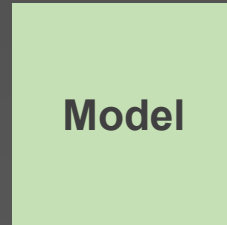
- Given only the inputs $\mathbf{x}$, models $p(\mathbf{x})$ and find



Model

**Patterns**

**Structure**

**Clusters**

**Generation**

$\mathbf{x}$

# Wine Quality Dataset

Supervised

## Inputs

| | |
|---|---|
| Fixed acidity | Total sulfur dioxide |
| Volatile acidity | Density |
| Citric acid | pH |
| Residual sugar | Sulphates |
| Chlorides | Alcohol |
| Free sulfur dioxide | ... |

Handcrafted features

## Outputs

**Model** → Quality

Good
**>6**

Poor

# Example: Decision Tree

Model

```
decisiontree = DecisionTreeClassifier(max_depth=4)

decisiontree.fit(x_train_values, y_train)

decisiontree_pred = decisiontree.predict(x_test_values)
```

*Learn* decision tree on the **training data**
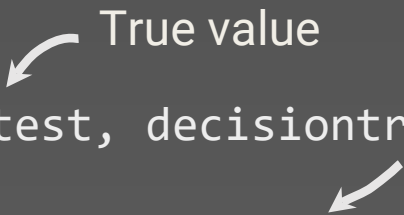
*Predict* wine quality on the **test data**

# How to evaluate our model?

**Use a performance metric (e.g., accuracy)**

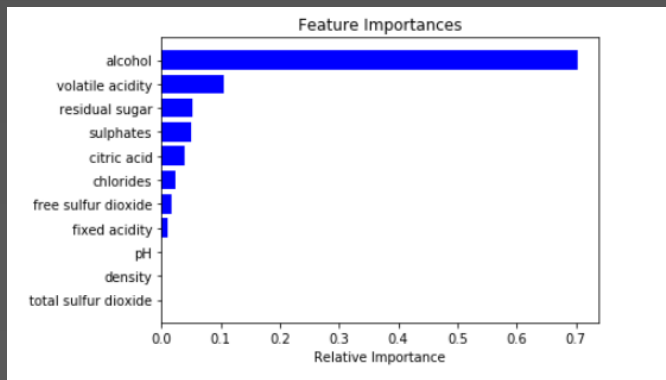Number of correct predictions made divided by the total.

True value

```
accuracy_score(y_test, decisiontree_pred)*100
```
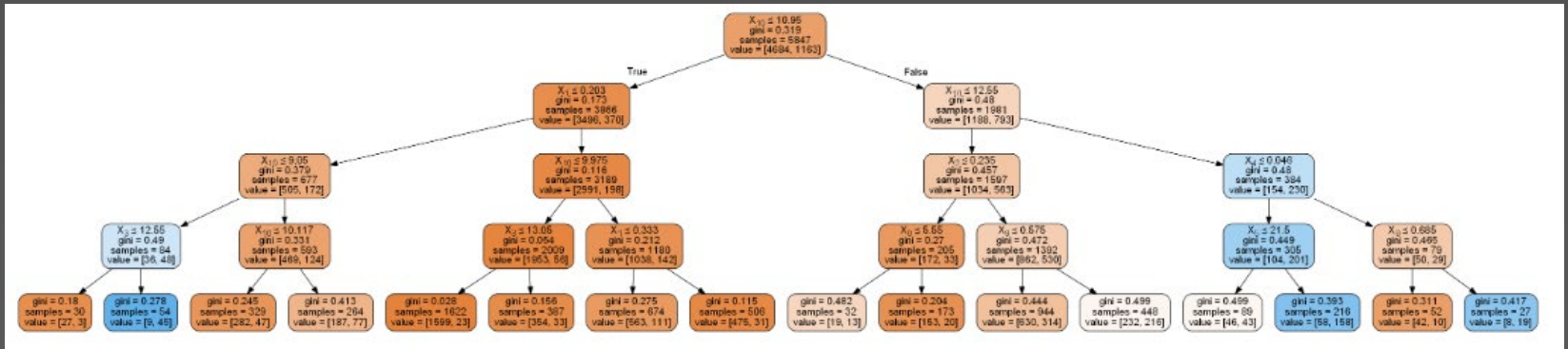
Predicted value

# How to interpret our model?

In a decision tree, compute **feature importance**.

```
importances = decisiontree.feature_importances_
indices = np.argsort(importances)
```



Feature Importances

# How to visualize our model?

Visualize the decision tree!

# Support Material

**Jupyter Notebook:**

Wine Quality classification using decision tree.

Available in the shared folder @ https://bit.ly/34r6YUs

# References

- **A Few Useful Things to Know About Machine Learning**, Domingos, 2012 (Link)

- Dataset source (Link)
- Modelling wine preferences by data mining from physicochemical properties (Link)
- Predicting wine quality using data analytics (Link)

# References

- Predicting quality of wine based on chemical attributes ([Link](#))
- Data analysis on the wine dataset ([Link](#))
- Wine Quality Classification ([Link](#))
- Vinho Verde webpage ([Link](#))

# Thank you for your attention!