# Practice 02: Intro to ML

MADMO, 2021

girafe
ai

# Outline

1. ML thesaurus and notation
2. Naïve Bayes classifier

# ML thesaurus

girafe
ai

01

# ML thesaurus

Denote the **dataset**.

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Target (passed) |
|------|-----|-------------------|---------------|-----------|-----------------|---------------|-----------------|
| John | 22 | 5 | 4 | Brown | English | 5 | TRUE |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | TRUE |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | TRUE |
| Michael | 27 | 3 | 4 | Green | French | 5 | TRUE |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | FALSE |

# ML thesaurus

**Observation** (or datum, or data point) is one piece of information.

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Target (passed) |
|---|---|---|---|---|---|---|---|
| John | 22 | 5 | 4 | Brown | English | 5 | TRUE |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | TRUE |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | TRUE |
| Michael | 27 | 3 | 4 | Green | French | 5 | TRUE |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | FALSE |

In many cases the observations are supposed to be **i.i.d.**

- **independent**
- **identically distributed**

5

# ML thesaurus

**Feature** (or predictor) represents some special property.

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Target (passed) |
|---|---|---|---|---|---|---|---|
| John | 22 | 5 | 4 | Brown | English | 5 | TRUE |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | TRUE |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | TRUE |
| Michael | 27 | 3 | 4 | Green | French | 5 | TRUE |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | FALSE |

# ML thesaurus

These all are features

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Target (passed) |
|------|-----|-------------------|---------------|-----------|-----------------|---------------|-----------------|
| John | 22 | 5 | 4 | Brown | English | 5 | TRUE |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | TRUE |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | TRUE |
| Michael | 27 | 3 | 4 | Green | French | 5 | TRUE |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | FALSE |

# ML thesaurus

These all are features

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Target (passed) |
|---|---|---|---|---|---|---|---|
| John | 22 | 5 | 4 | Brown | English | 5 | TRUE |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | TRUE |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | TRUE |
| Michael | 27 | 3 | 4 | Green | French | 5 | TRUE |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | FALSE |

# ML thesaurus

These all are features

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Target (passed) |
|------|-----|-------------------|---------------|-----------|-----------------|---------------|-----------------|
| John | 22 | 5 | 4 | Brown | English | 5 | TRUE |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | TRUE |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | TRUE |
| Michael | 27 | 3 | 4 | Green | French | 5 | TRUE |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | FALSE |

# ML thesaurus

These all are features

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Target (passed) |
|---|---|---|---|---|---|---|---|
| John | 22 | 5 | 4 | Brown | English | 5 | TRUE |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | TRUE |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | TRUE |
| Michael | 27 | 3 | 4 | Green | French | 5 | TRUE |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | FALSE |

# ML thesaurus

And even the name is a **feature**

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Target (passed) |
|------|-----|-------------------|---------------|-----------|-----------------|---------------|-----------------|
| John | 22 | 5 | 4 | Brown | English | 5 | TRUE |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | TRUE |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | TRUE |
| Michael | 27 | 3 | 4 | Green | French | 5 | TRUE |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | FALSE |

# ML thesaurus

The **design matrix** contains all the features and observations.

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Target (passed) |
|---|---|---|---|---|---|---|---|
| John | 22 | 5 | 4 | Brown | English | 5 | TRUE |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | TRUE |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | TRUE |
| Michael | 27 | 3 | 4 | Green | French | 5 | TRUE |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | FALSE |

*Features can even be multidimensional, we will discuss it later in this course.*

# ML thesaurus

**Target** represents the information we are interested in.

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Target (passed) |
|------|-----|-------------------|---------------|-----------|-----------------|---------------|-----------------|
| John | 22 | 5 | 4 | Brown | English | 5 | TRUE |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | TRUE |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | TRUE |
| Michael | 27 | 3 | 4 | Green | French | 5 | TRUE |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | FALSE |

*Target can be either a **number** (real, integer, etc.) – for **regression** problem*

# ML thesaurus

*Target* represents the information we are interested in.

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Target (passed) |
|---|---|---|---|---|---|---|---|
| John | 22 | 5 | 4 | Brown | English | 5 | TRUE |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | TRUE |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | TRUE |
| Michael | 27 | 3 | 4 | Green | French | 5 | TRUE |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | FALSE |

*Or a **label** – for **classification** problem*

# ML thesaurus

***Target*** represents the information we are interested in.

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Target (passed) |
|---|---|---|---|---|---|---|---|
| John | 22 | 5 | 4 | Brown | English | 5 | TRUE |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | TRUE |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | TRUE |
| Michael | 27 | 3 | 4 | Green | French | 5 | TRUE |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | FALSE |

*Mark can be treated as a label too (due to finite number of labels: 1 to 5). We will discuss it later.*

# ML thesaurus

Further we will work with the numerical target (mark)

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) |
|------|-----|-------------------|---------------|-----------|-----------------|---------------|
| John | 22 | 5 | 4 | Brown | English | 5 |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 |
| Michael | 27 | 3 | 4 | Green | French | 5 |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 |

# ML thesaurus

The **prediction** contains values we predicted using some **model**.

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Predicted (mark) |
|------|-----|-------------------|---------------|-----------|-----------------|---------------|------------------|
| John | 22 | 5 | 4 | Brown | English | 5 | 4.5 |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | 4.5 |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | 5 |
| Michael | 27 | 3 | 4 | Green | French | 5 | 3.5 |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | 3 |

One could notice that prediction just averages of Statistics and Python marks. So our **model** can be represented as follows:

$$\hat{\text{mark}}_{ML} = \frac{1}{2}\text{mark}_{Statistics} + \frac{1}{2}\text{mark}_{Python}$$

# ML thesaurus

The **prediction** contains values we predicted using some **model**.

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Predicted (mark) |
|------|-----|-------------------|---------------|-----------|-----------------|---------------|------------------|
| John | 22 | 5 | 4 | Brown | English | 5 | 4.5 |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | 4.5 |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | 5 |
| Michael | 27 | 3 | 4 | Green | French | 5 | 3.5 |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | 3 |

*Different models can provide different predictions:*

$$\hat{\text{mark}}_{ML} = \frac{1}{2}\text{mark}_{Statistics} + \frac{1}{2}\text{mark}_{Python}$$

# ML thesaurus

The **prediction** contains values we predicted using some **model**.

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Predicted (mark) |
|---|---|---|---|---|---|---|---|
| John | 22 | 5 | 4 | Brown | English | 5 | 1 |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | 5 |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | 2 |
| Michael | 27 | 3 | 4 | Green | French | 5 | 4 |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | 3 |

*Different models can provide different predictions:*

$$\hat{\text{mark}}_{ML} = \text{random}(\text{integer from } [1; 5])$$

# ML thesaurus

The **prediction** contains values we predicted using some **model**.

| Name | Age | Statistics (mark) | Python (mark) | Eye color | Native language | Target (mark) | Predicted (mark) |
|---|---|---|---|---|---|---|---|
| John | 22 | 5 | 4 | Brown | English | 5 | 1 |
| Aahna | 17 | 4 | 5 | Brown | Hindi | 4 | 5 |
| Emily | 25 | 5 | 5 | Blue | Chinese | 5 | 2 |
| Michael | 27 | 3 | 4 | Green | French | 5 | 4 |
| Some student | 23 | 3 | 3 | NA | Esperanto | 2 | 3 |

*Different models can provide different predictions.*

*Usually some **hypothesis** lies beneath the model choice.*

# Naïve Bayes classifier

girafe
ai

# Naïve Bayes classifier

Let's denote:

- Training set $\mathcal{L} = \{\mathbf{x}_i, y_i\}_{i=1}^{n}$ , where
  - $\mathbf{x}_i \in \mathbb{R}^p$ , $y_i \in \{C_1, \ldots, C_k\}$ for k-class classification

# Bayes' theorem

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

or, in our case

$$P(y_i = C_k|\mathbf{x}_i) = \frac{P(\mathbf{x}_i|y_i = C_k)P(y_i = C_k)}{P(\mathbf{x}_i)}$$

# Naïve Bayes classifier

Let's denote:

- Training set $\mathcal{L} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ , where
  - $\mathbf{x}_i \in \mathbb{R}^p$ , $y_i \in \{C_1, \ldots, C_K\}$  for K-class classification

$$P(y_i = C_k | \mathbf{x}_i) = \frac{P(\mathbf{x}_i | y_i = C_k) P(y_i = C_k)}{P(\mathbf{x}_i)}$$

Naïve assumption: features are **independent**

# Naïve Bayes classifier

$$P(y_i = C_k | \mathbf{x}_i) = \frac{P(\mathbf{x}_i | y_i = C_k) P(y_i = C_k)}{P(\mathbf{x}_i)}$$

Naïve assumption: features are **independent:**

$$P(\mathbf{x}_i | y_i = C_k) = \prod_{l=1}^{p} P(x_i^l | y_i = C_k)$$

# Naïve Bayes classifier

$$P(y_i = C_k | \mathbf{x}_i) = \frac{P(\mathbf{x}_i | y_i = C_k) P(y_i = C_k)}{P(\mathbf{x}_i)}$$

Optimal class label:

$$C^* = \arg\max_k P(y_i = C_k | \mathbf{x_i})$$

To find maximum we even do not need the denominator

But we need it to get probabilities

# Let's Practice

Thanks for attention!

girafe
ai

# Model validation and evaluation

# Supervised learning problem statement

Let's denote:

- Training set $\mathcal{L} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ , where

  - $(\mathbf{x} \in \mathbb{R}^p, y \in \mathbb{R})$ for regression

  - $\mathbf{x}_i \in \mathbb{R}^p$, $y_i \in \{+1, -1\}$ for binary classification

Model $f(\mathbf{x})$ predicts some value for every object

Loss function $Q(\mathbf{x}, y, f)$ that should be minimized

29

# Overfitting vs. underfitting



**Under-fitting**

(too simple to explain the variance)

**Appropriate-fitting**

**Over-fitting**

(forcefitting -- too good to be true)
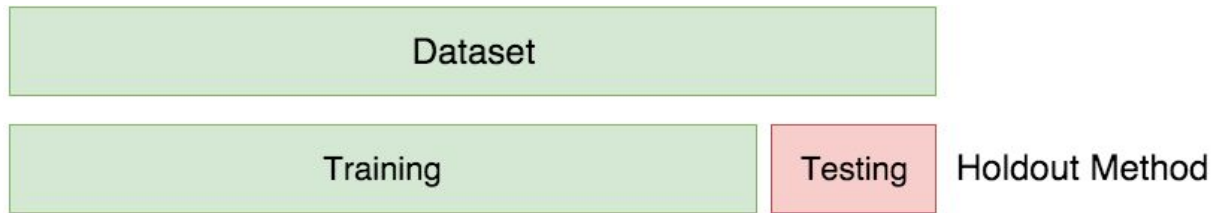
# Overfitting vs. underfitting

# Overfitting vs. underfitting

- We can control overfitting / underfitting by altering model's capacity (ability to fit a wide variety of functions):
- select appropriate hypothesis space
- learning algorithm's effective capacity may be less than the representational capacity of the model family
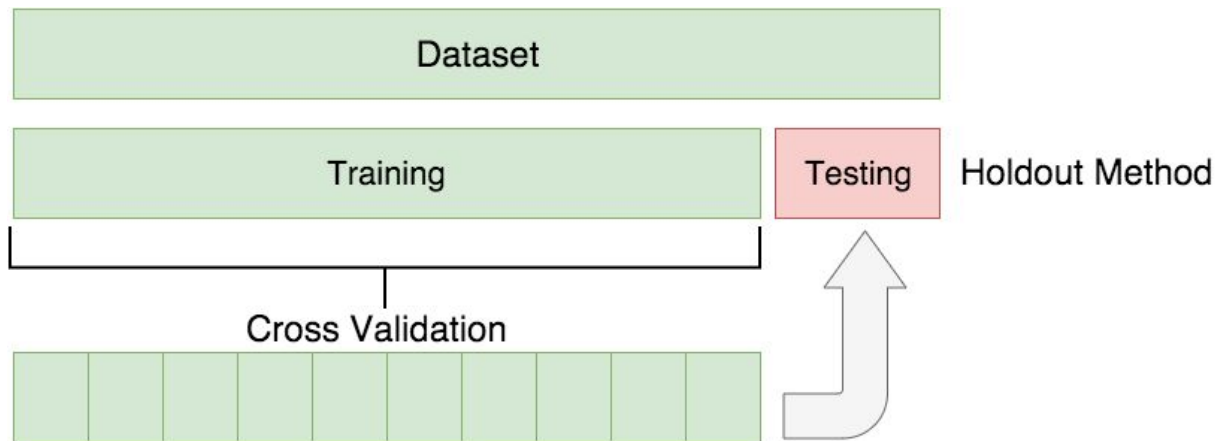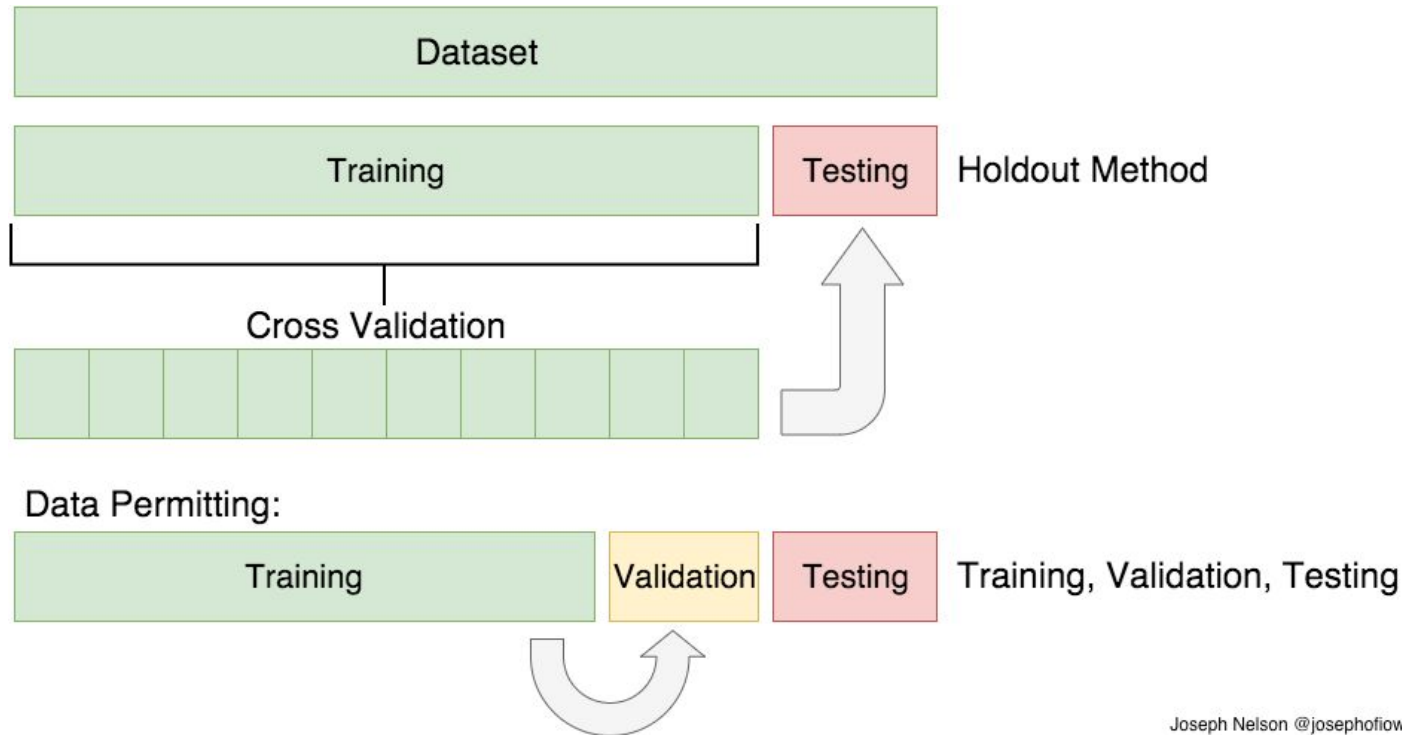
# Evaluating the quality



Dataset

Training | Testing | Holdout Method

# Evaluating the quality



Is it good enough?

# Evaluating the quality

# Evaluating the quality



Image credit: Joseph Nelson

# Cross-validation



Training set

Training folds      Test fold

1st iteration $\Rightarrow E_1$

2nd iteration $\Rightarrow E_2$

3rd iteration $\Rightarrow E_3$

...

10th iteration $\Rightarrow E_{10}$

$$E = \frac{1}{10}\sum_{i=1}^{10} E_i$$