

# Demystifying Artificial Intelligence Sorcery

(Part 2: Machine Learning)<sup>a</sup>

---

Abdelbacet Mhamdi  
abdelbacet.mhamdi@bizerte.r-iset.tn

*Dr.-Ing. in Electrical Engineering*  
*Senior Lecturer at ISET Bizerte*

---

<sup>a</sup>Available @ <https://github.com/a-mhamdi/jlai/>



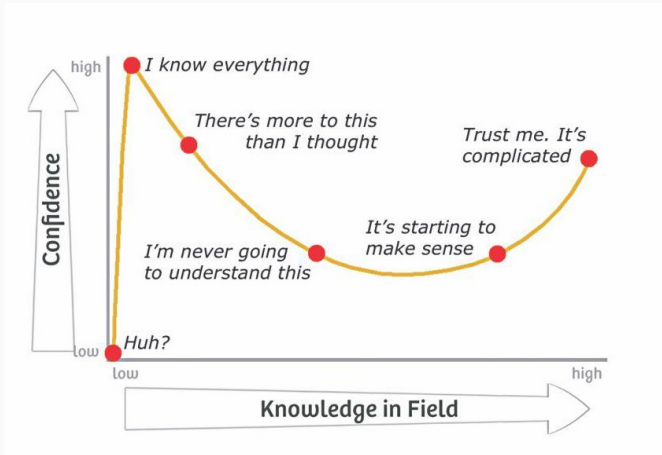
## Disclaimer

This document features some materials gathered from multiple online sources.

Please note no copyright infringement is intended, and I do not own nor claim to own any of the original materials. They are used for educational purposes only.

I have included links solely as a convenience to the reader. Some links within these slides may lead to other websites, including those operated and maintained by third parties. The presence of such a link does not imply a responsibility for the linked site or an endorsement of the linked site, its operator, or its contents.

## DUNNING-KRUGER EFFECT



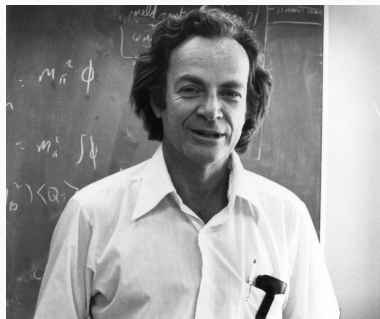
Kruger, J. and Dunning, D. (1999) *Unskilled and unaware of it: How difficulties in recognizing one's own incompetence lead to inflated self-assessments*. **J Pers Soc Psychol.** 77(6) pp. 1121–1134.

 10.1037/0022-3514.77.6.1121

“Knowledge isn’t free. You have to pay attention.”

---

**Richard P. Feynman**



## ROADMAP

1. An overview
2. Supervised Learning
3. Unsupervised Learning
4. Complementary Lab. Project
5. ML Landscape through Quizzes

## **An overview**

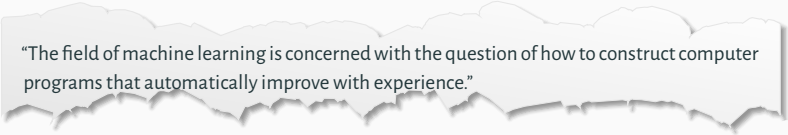
---

## GLOBAL DATA TRAFFIC



Update on the internet in real time is available [here](#).

## LITERATURE REVIEW (1/3)



“The field of machine learning is concerned with the question of how to construct computer programs that automatically improve with experience.”

Mitchell, T. (1997) *Machine Learning*. **McGraw-Hill International Editions. McGraw-Hill.**



## LITERATURE REVIEW (2/3)

“Machine learning (ML) is a scientific discipline that concerns developing learning capabilities in computer systems. Machine learning is one of central areas of Artificial Intelligence (AI). It is an interdisciplinary area that combines results from statistics, logic, robotics, computer science, computational intelligence, pattern recognition, data mining, cognitive science, and more.”

Wojtusiak, J. (2012) Machine learning. In *Encyclopedia of the Sciences of Learning*, pages 2082–2083. Springer US.

## LITERATURE REVIEW (3/3)

“Machine learning is an evolving branch of computational algorithms that are designed to emulate human intelligence by learning from the surrounding environment. They are considered the working horse in the new era of the so-called big data. Techniques based on machine learning have been applied successfully in diverse fields ranging from pattern recognition, computer vision, spacecraft engineering, finance, entertainment, and computational biology to biomedical and medical applications. [...] The ability of machine learning algorithms to learn from current context and generalize into unseen tasks would allow improvements in both the safety and efficacy of radiotherapy practice leading to better outcomes.”

El Naqa, I. and Murphy, M. J. (2015) *What Is Machine Learning?*, pages 3–11. **Springer International Publishing.**

## DEBRIEF

### Arthur Samuel (1959)

Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.

### Tom Mitchell (1998)

Well-posed Learning Problem: A computer is said to learn from experience  $\mathcal{E}$  with respect to some task  $\mathcal{T}$  and some performance measure  $\mathcal{P}$ , if its performance on  $\mathcal{T}$ , as measured by  $\mathcal{P}$ , improves with experience  $\mathcal{E}$ .

### Task #1

Suppose your email program watches which emails you do or do not mark as spam, and based on that learns how to better filter spam. What is the task  $\mathcal{T}$  in this setting?

1. Classifying emails as spam or not spam;
2. Watching you label emails as spam or not spam;
3. The number (or fraction) of emails correctly classified as spam/not spam;
4. None of the above-this not a machine learning problem.

# DEBRIEF

## Arthur Samuel (1959)

Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.

## Tom Mitchell (1998)

Well-posed Learning Problem: A computer is said to learn from experience  $\mathcal{E}$  with respect to some task  $\mathcal{T}$  and some performance measure  $\mathcal{P}$ , if its performance on  $\mathcal{T}$ , as measured by  $\mathcal{P}$ , improves with experience  $\mathcal{E}$ .

## Task #1

Suppose your email program watches which emails you do or do not mark as spam, and based on that learns how to better filter spam. What is the task  $\mathcal{T}$  in this setting?

1. Classifying emails as spam or not spam;
2. Watching you label emails as spam or not spam;
3. The number (or fraction) of emails correctly classified as spam/not spam;
4. None of the above-this not a machine learning problem.

# DEBRIEF

## Arthur Samuel (1959)

Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.

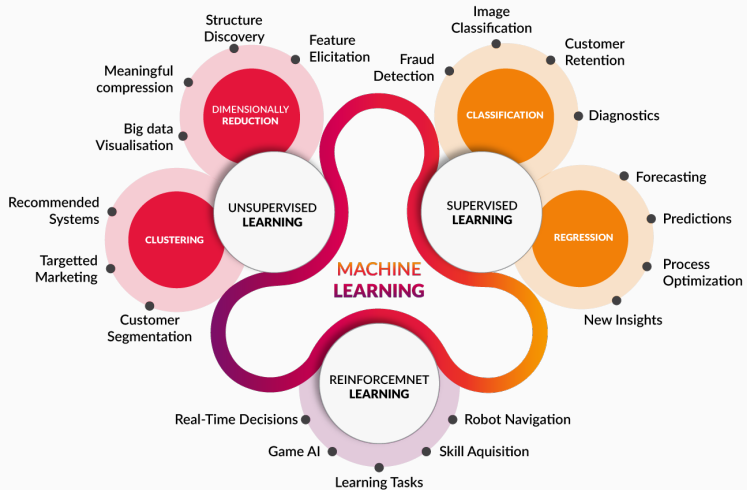
## Tom Mitchell (1998)

Well-posed Learning Problem: A computer is said to learn from experience  $\mathcal{E}$  with respect to some task  $\mathcal{T}$  and some performance measure  $\mathcal{P}$ , if its performance on  $\mathcal{T}$ , as measured by  $\mathcal{P}$ , improves with experience  $\mathcal{E}$ .

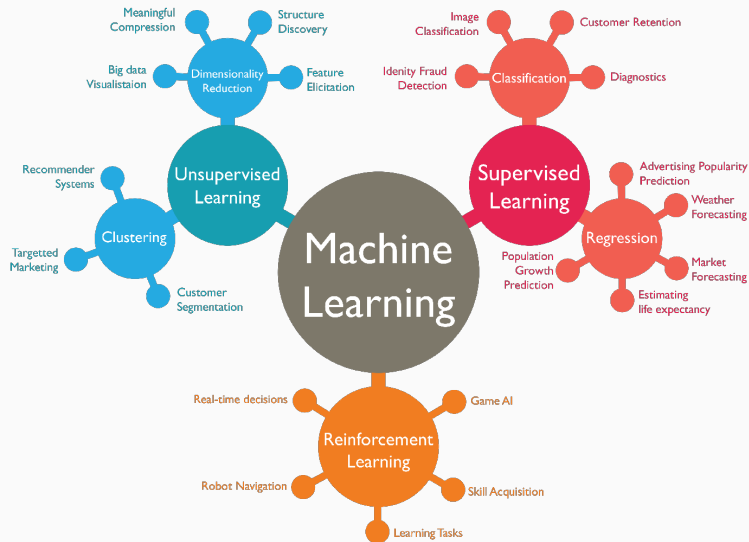
## Task #1

Suppose your email program watches which emails you do or do not mark as spam, and based on that learns how to better filter spam. What is the task  $\mathcal{T}$  in this setting?

1. Classifying emails as spam or not spam;
2. Watching you label emails as spam or not spam;
3. The number (or fraction) of emails correctly classified as spam/not spam;
4. None of the above-this not a machine learning problem.

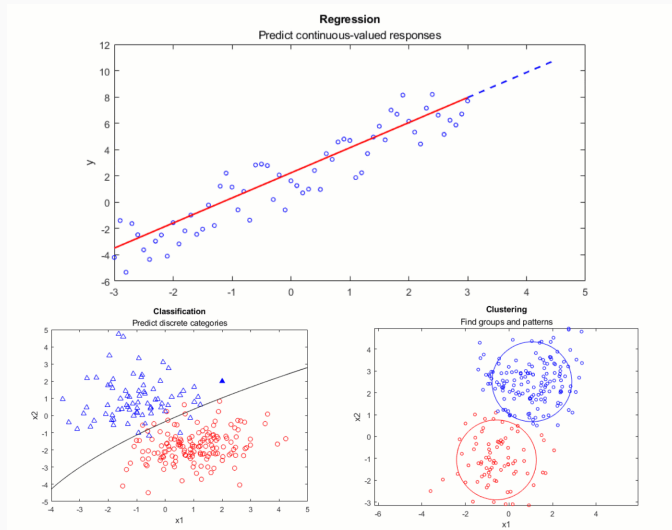


<https://www.cognub.com/index.php/cognitive-platform/>



<https://vitalflux.com/great-mind-maps-for-learning-machine-learning/>

# REGRESSION | CLASSIFICATION | CLUSTERING



<https://github.com/MathWorks-Teaching-Resources/Machine-Learning-for-Regression>



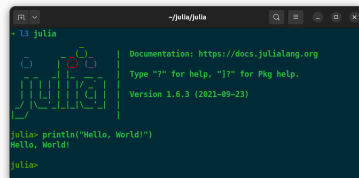


**REMINDER**

# PROGRAMMING LANGUAGE



[julialang.org/](https://julialang.org/)

A screenshot of a terminal window titled "~julia/julia". The prompt is "+ julia". The output shows a stylized "julia" logo made of green and red characters, followed by the text: "Documentation: https://docs.julialang.org", "Type '?' for help, ']' for Pkg help.", and "Version 1.6.3 (2021-09-23)". Below this, the user enters "julia> println(\"Hello, World!\")" and the output "Hello, World!" is displayed. The prompt "julia>" is shown again at the bottom.

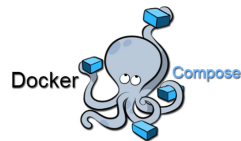
```
+ julia

julia
 |
 | Documentation: https://docs.julialang.org
 | Type '?' for help, ']' for Pkg help.
 | Version 1.6.3 (2021-09-23)
 |
julia> println("Hello, World!")
Hello, World!
julia>
```

## DEVELOPMENT ENVIRONMENTS



**Pluto.jl** 

The Pluto.jl logo consists of the text "Pluto.jl" in a bold, black, sans-serif font, followed by a small icon of three stacked circles in green, purple, and red.

▲ \$ docker compose up

▼ \$ docker compose down

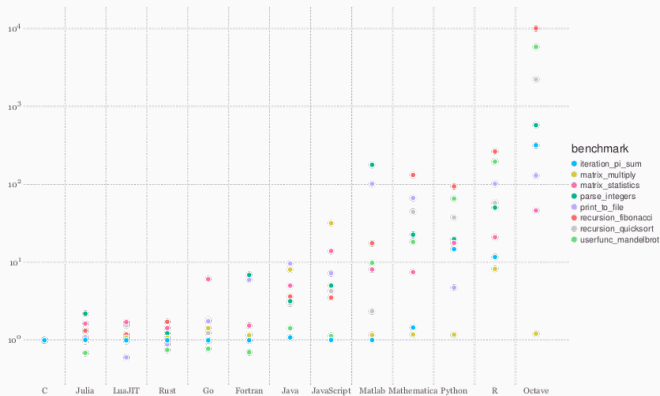


# JULIA IN A NUTSHELL

- ▲ **Fast:** native code for multiple platforms via LLVM;
- ▲ **Dynamic:** good support for interactive use (*like a scripting language*);
- ▲ **Reproducible:** environment recreation across platforms, with pre-built binaries;
- ▲ **Composable:** multiple dispatch as a paradigm (*oop & functional programming*);
- ▲ **General:** asynchronous I/O, metaprogramming, debugging, logging; profiling, pkg, ...
- ▲ **Open Source:** GitHub repository at <https://github.com/JuliaLang/julia>.



# JULIA MICRO-BENCHMARKS (1/2)



<https://julialang.org/benchmarks>



## JULIA MICRO-BENCHMARKS (2/2)

### Geometric Means<sup>1</sup> of Micro-Benchmarks by Language

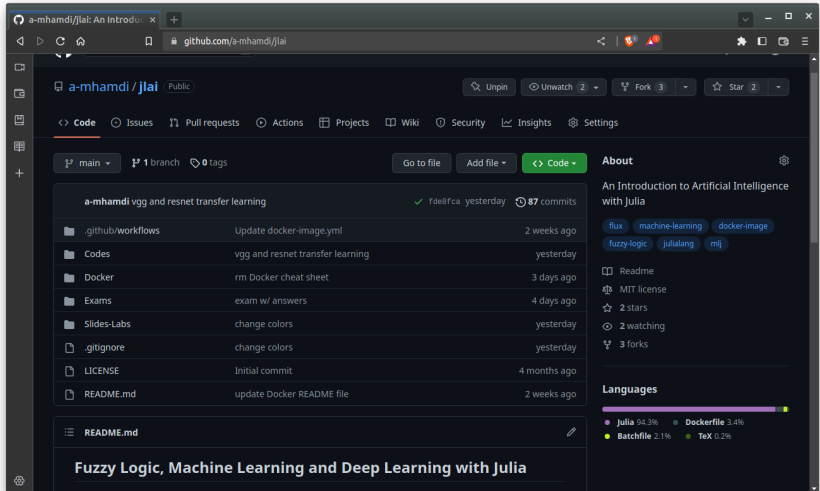
1	C	1.0
2	Julia	1.17006
3	LuaJIT	1.02931
4	Rust	1.0999
5	Go	1.49917
6	Fortran	1.67022
7	Java	3.46773
8	JavaScript	4.79602
9	Matlab	9.57235
10	Mathematica	14.6387
11	Python	16.9262
12	R	48.5796
13	Octave	338.704



<sup>1</sup>Measure of central tendency expressed as  $(x_1 \times x_2 \times \dots \times x_n)^{1/n}$



## SOURCE CONTROL MANAGEMENT (SCM)



<https://github.com/a-mhamdi/jlai>



# DOCKER IMAGE

The screenshot shows a web browser window displaying the Docker Hub page for the `abmhamdi/jlai-p2` Docker image. The page has a dark theme. At the top, the browser tab is titled "abmhamdi/jlai-p2 - Docker Image | Docker Hub - Brave". The address bar shows the URL `hub.docker.com/r/abmhamdi/jlai-p2`. The Docker Hub navigation bar is visible with links for Explore, Repositories, Organizations, and Usage. A search bar and a "Manage Repository" button are also present. The main content area shows the repository name `abmhamdi/jlai-p2` with a Docker logo icon. Below the name, it says "By `abmhamdi` · Updated 6 minutes ago" and "Artificial Intelligence Labs - Part 2 @ ISETBZ". There are tags for "DATA SCIENCE", "LANGUAGES & FRAMEWORKS", and "MACHINE LEARNING & AI". The repository has 0 stars and 64 downloads. The "Overview" tab is selected, showing a description: "Machine Learning with Julia". The description states: "This repository contains slides, labs and code examples for using Julia to implement some artificial intelligence related algorithms. Codes run on top of a Docker image, ensuring a consistent and reproducible environment." To the right, the "Docker Pull Command" is displayed as `docker pull abmhamdi/jlai-p2` with a "Copy" button.


abmhamdi/jlai-p2 - Docker Image | Docker Hub - Brave

abmhamdi/jlai-p2 - Docker Image | Docker Hub - Brave

hub.docker.com/r/abmhamdi/jlai-p2

dockerhub Explore Repositories Organizations Usage Search Dock... [ctrl+K] ? ⚙️ 🔔 🌙 ☰ A

Explore / abmhamdi/jlai-p2

 **abmhamdi/jlai-p2**

By `abmhamdi` · Updated 6 minutes ago

Artificial Intelligence Labs - Part 2 @ ISETBZ

IMAGE

DATA SCIENCE LANGUAGES & FRAMEWORKS MACHINE LEARNING & AI

☆0 ↓ 64

Manage Repository

Overview Tags

**Machine Learning with Julia**

This repository contains slides, labs and code examples for using Julia to implement some artificial intelligence related algorithms. Codes run on top of a Docker image, ensuring a consistent and reproducible environment.

Docker Pull Command

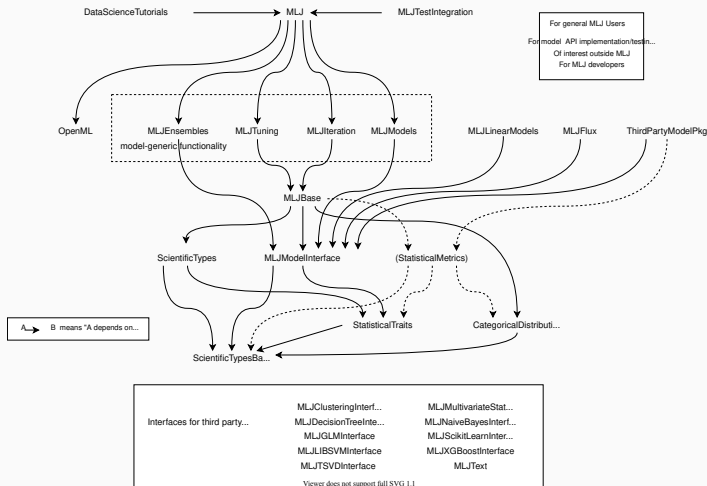
```
docker pull abmhamdi/jlai-p2
```

Copy

<https://hub.docker.com/r/abmhamdi/jlai-p2>



# A MACHINE LEARNING FRAMEWORK FOR JULIA

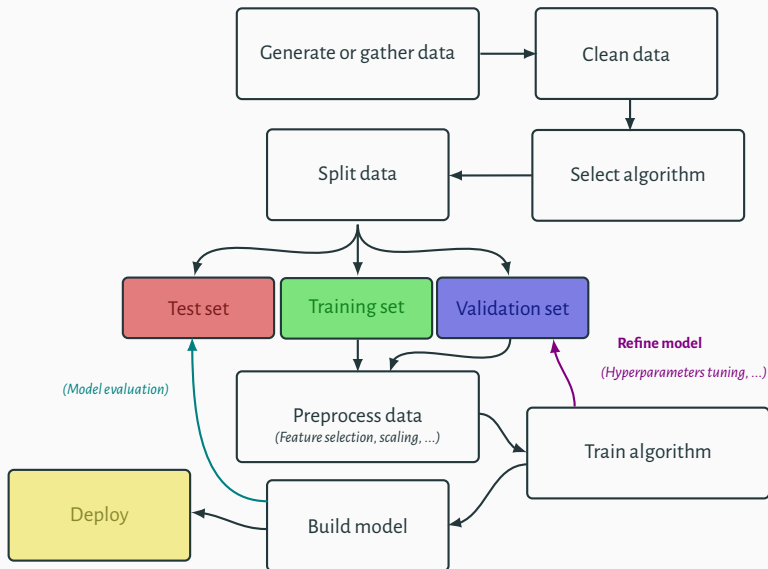


<https://docs.juliahub.com/MLJ/>

## Supervised Learning

---

# WORKFLOW IN MACHINE LEARNING



# DATA PREPROCESSING

How?

**Cleaning** Identifying and correcting or removing inaccuracies and inconsistencies in the data.

**Transformation** Converting data from one format or structure to another.

**Normalization** Scaling the data so that it fits within a specific range. This is often done to make the data more amenable to certain operations or algorithms.

# DATA PREPROCESSING

## WHY?

- ▶ Raw data is often messy and may need to be cleaned and formatted before it can be used for machine learning.  
*(This may involve removing missing or invalid data, handling outliers, and encoding categorical variables.)*
- ▶ Normalizing the data can help to scale the features so that they are on the same scale.  
*(This can be important for algorithms that use distance measures, as features on different scales can dominate the distance measure.)*
- ▶ Preprocessing techniques such as feature selection and feature extraction can help to reduce the dimensionality of the data.  
*(This may improve the performance of the model and reduce the risk of overfitting.)*
- ▶ Preprocessing techniques such as feature selection can help to identify the most important features in the data.  
*(This can make the model more interpretable and easier to understand.)*

# DATA PREPROCESSING

## FEATURE SCALING

### Normalization

---

$$X \triangleq \frac{X - \text{minimum}(X)}{\text{maximum}(X) - \text{minimum}(X)}$$

- ▲ No assumption on data distribution

### Standardization (*Standardizer*)

---

$$X \triangleq \frac{X - \mu}{\sigma}$$

- ▲ More recommended when following normal distribution
-

# DATA PREPROCESSING TEMPLATE



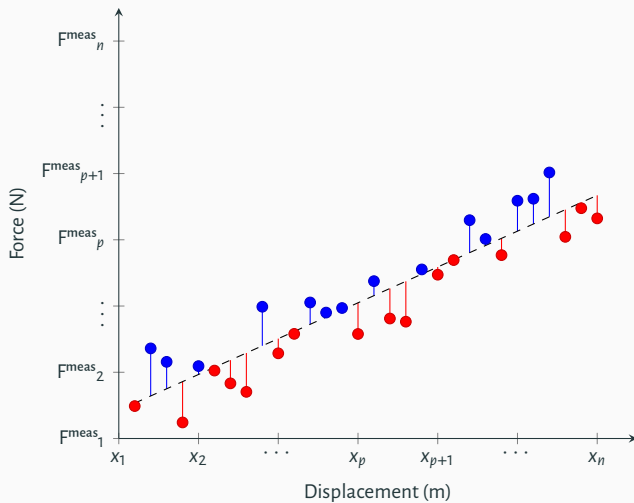
The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → *Codes* → *Julia* → *Part-2*

→ *Pluto* → *ml-workflows.jl*

**Pluto.jl** 

→ *Jupyter* → *ml-workflows.ipynb*







Consider the example of a spring. Our main goal is to determine the stiffness  $k$  of this spring, given some experimental data. The mathematical model (*Hooke's law*):

$$F = kx \quad (1)$$

Restoring force is proportional to displacement.

**Table 1:** Measurements of couple  $(x_i, F^{\text{meas}}_i)$

$x_i$	$x_1$	$\dots$	$x_p$	$\dots$	$x_n$
$F^{\text{meas}}_i$	$F^{\text{meas}}_1$	$\dots$	$F^{\text{meas}}_p$	$\dots$	$F^{\text{meas}}_n$

$$\begin{aligned} F^{\text{meas}}_i &= F_i + \varepsilon_i \\ &= kx_i + \varepsilon_i, \end{aligned} \quad (2)$$

where  $F_i$  denotes the unknown real value of the force applied to the spring. In order to estimate the stiffness value  $k$ , we can consider the quadratic criterion:

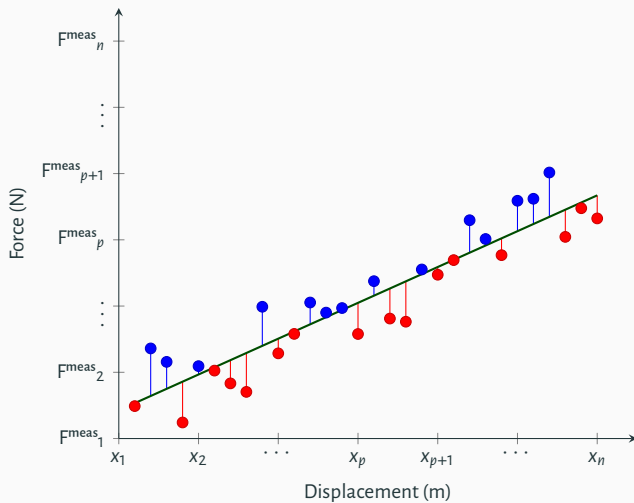
$$\begin{aligned} \mathcal{J} &= \sum_{i=1}^n \varepsilon_i^2 \\ &= \sum_{i=1}^n (F^{\text{meas}}_i - kx_i)^2 \end{aligned}$$

$$\frac{\partial \mathcal{J}}{\partial k} = 0 \quad (3)$$

$$2 \sum_{i=1}^n (F^{\text{meas}}_i - kx_i) \sum_{i=1}^n \frac{\partial (F^{\text{meas}}_i - kx_i)}{\partial k} = 0$$

$$\sum_{i=1}^n (F^{\text{meas}}_i - kx_i) \sum_{i=1}^n x_i = 0$$

$$\sum_{i=1}^n F^{\text{meas}}_i x_i = k \sum_{i=1}^n x_i^2 \iff \hat{k} = \frac{\sum_{i=1}^n F^{\text{meas}}_i x_i}{\sum_{i=1}^n x_i^2}$$



## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → Codes → Julia → Part-2

→ Pluto → *simple-regression.jl*

**Pluto.jl** 

→ Jupyter → *simple-regression.ipynb*



This example consists on determining the unknown couple  $(y_0, v_0)$  of a mobile solid. We assume that the trajectory is linear. The mathematical model that relates the position  $y$  to time  $t$  is given by this equation:

$$y = y_0 + v_0 t \quad (4)$$

**Table 2:** Measurements of position  $y$

$t_k$	$t_1$	$\dots$	$t_p$	$\dots$	$t_n$
$y_k^{\text{meas}}$	$y_1^{\text{meas}}$	$\dots$	$y_p^{\text{meas}}$	$\dots$	$y_n^{\text{meas}}$

$$\begin{aligned} y_k^{\text{meas}} &= y_k + \varepsilon_k \\ &= y_0 + v_0 t_k + \varepsilon_k, \end{aligned} \quad (5)$$

where  $y_k$  denotes the unknown real value of the position  $y$  at time point  $t_k$ .

In order to estimate the values taken by the couple  $\begin{bmatrix} y_0, & v_0 \end{bmatrix}^T$ , we consider the quadratic criterion again, as follows:

$$\begin{aligned} \mathcal{J} &= \sum_{k=1}^n \varepsilon_k^2 \\ &= \boldsymbol{\varepsilon}^T \times \boldsymbol{\varepsilon} \end{aligned}$$

The vector  $\boldsymbol{\varepsilon}$  is set by  $\varepsilon_k, \forall k \geq 1$ :

$$\boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 & \cdots & \varepsilon_n \end{bmatrix}^T$$

$$\frac{\partial \mathcal{J}}{\partial \begin{bmatrix} y_0 \\ v_0 \end{bmatrix}} = 0 \quad (6)$$

## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → Codes → Julia → Part-2

→ Pluto → *multivariable-regression.jl*

**Pluto.jl** 

→ Jupyter → *multivariable-regression.ipynb*



Consider the following multivariable equation:

$$y = \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_m x_m \quad (7)$$

For a particular single measurement, eq. (7) can be updated as

$$y_k = \theta_1 x_{(1,k)} + \theta_2 x_{(2,k)} + \cdots + \theta_m x_{(m,k)} + \varepsilon_k$$

We denote hereafter by  $\boldsymbol{\theta}$  the vector  $\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_m \end{bmatrix}$ . The function  $y_k$  becomes:

$$y_k = \underbrace{[x_{(1,k)}, x_{(2,k)}, \cdots, x_{(m,k)}]}_{\mathbf{x}_k^T} \boldsymbol{\theta} + \varepsilon_k$$

We assume that we have  $n$  measurements for  $y$ . Then we can transform the previous equation into

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\varepsilon},$$

$$\text{where } \mathbf{y}^T = [y_1, y_2, \cdots, y_n], \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}, \text{ and } \boldsymbol{\varepsilon}^T = [\varepsilon_1, \varepsilon_2, \cdots, \varepsilon_n].$$



We can consider the mean squared error or quadratic criterion in order to compute the approximated value of  $\theta$ :

$$\begin{aligned}\mathcal{J} &= \sum_{k=1}^n \varepsilon_k^2 \\ &= \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon}\end{aligned}$$

The best well estimated value of  $\hat{\theta}$  corresponds to the absolute minimum of  $\mathcal{J}$ . This leads to calculate the gradient of  $\mathcal{J}$  with respect to  $\theta$ :

$$\frac{\partial \mathcal{J}}{\partial \theta} = \frac{\partial (\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon})}{\partial \theta}$$

$$\frac{\partial (\boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon})}{\partial \theta} = 2 \left( \frac{\partial \boldsymbol{\varepsilon}}{\partial \theta} \right)^T \boldsymbol{\varepsilon}$$

Recall that  $\boldsymbol{\varepsilon} = \mathbf{y} - \mathbf{X}\boldsymbol{\theta}$ , the term  $\frac{\partial \boldsymbol{\varepsilon}}{\partial \theta}$  hence becomes:

$$\frac{\partial \boldsymbol{\varepsilon}}{\partial \theta} = -\mathbf{X}$$

$$\begin{aligned}\frac{\partial J}{\partial \theta} &= 2(-\mathbf{X})^T (\mathbf{y} - \mathbf{X}\theta) \\ &= 0\end{aligned}$$

The vector  $\hat{\theta}$  is given by

$$\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$



$\mathbf{X}^T \mathbf{X}$  is not invertible (singular/degenerate)

#### ▼ Redundant Features

Some features are linearly dependent, i.e,  $\exists$  some  $x_p \propto$  some  $x_l$ , e.g.,  $x_p$  in feet and  $x_l$  in m.

#### ▼ Too many features

Fewer observations compared to the number of features, i.e,  $m \geq n$ .

- ▲ Delete some features
- ▲ Add extra observations

▲ Use regularization:

$\lambda \sum_{i=2}^m  \theta_i $	$\frac{1}{2} \lambda \sum_{i=2}^m \theta_i^2$	$r\lambda \sum_{i=2}^m  \theta_i  + \frac{(1-r)}{2} \lambda \sum_{i=2}^m \theta_i^2$
LASSO	RIDGE	ELASTIC NET

$$\hat{\theta} = \left( \mathbf{X}^T \mathbf{X} + \lambda \begin{bmatrix} 0 & \cdots & \cdots & \cdots & 0 \\ \vdots & 1 & \ddots & & \vdots \\ \vdots & \ddots & 1 & \ddots & \vdots \\ \vdots & & \ddots & 1 & 0 \\ 0 & \cdots & \cdots & 0 & 1 \end{bmatrix}_{(m,m)} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → Codes → Julia → Part-2

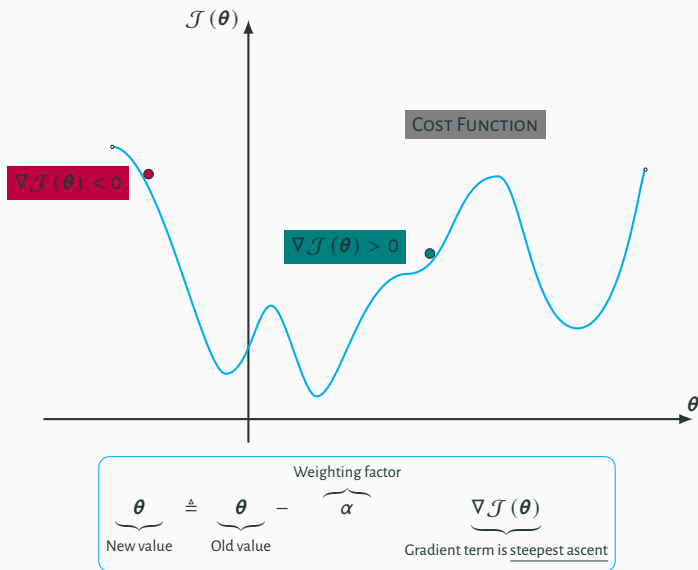
→ Pluto → *polynomial-regression.jl*

**Pluto.jl** 

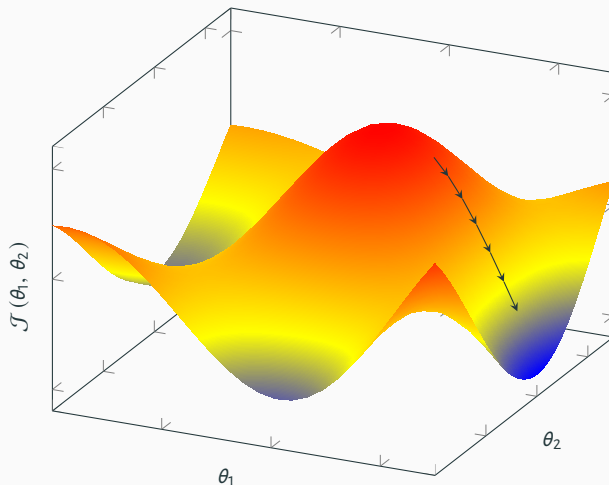
→ Jupyter → *polynomial-regression.ipynb*



## GRADIENT DESCENT



# GRADIENT DESCENT



- ① Start with some random values of  $\theta_1$  and  $\theta_2$
- ② Keep changing  $\theta_1$  and  $\theta_2$  to reduce  $\mathcal{J}(\theta_1, \theta_2)$  until we hopefully end up at minimum

# GRADIENT DESCENT

$$\theta_i \triangleq \theta_i - \underbrace{\alpha}_{\text{LEARNING RATE}} \frac{\partial \mathcal{J}}{\partial \theta_i}$$

Recall that

$$\mathcal{J} = \frac{1}{2n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k))^2 \implies \frac{\partial \mathcal{J}}{\partial \theta_i} = -\frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(i,k)}$$

$$\theta_1 \triangleq \theta_1 + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(1,k)}$$

$$\theta_2 \triangleq \theta_2 + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(2,k)}$$

⋮

$$\theta_m \triangleq \theta_m + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(m,k)}$$

## GRADIENT DESCENT

$$\theta_i \triangleq \theta_i - \underbrace{\alpha}_{\text{LEARNING RATE}} \frac{\partial \mathcal{J}}{\partial \theta_i}$$

Recall that with  $L_2$  regularization term

$$\mathcal{J} = \frac{1}{2n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k))^2 + \frac{\lambda}{2n} \sum_{i=2}^m \theta_i^2 \implies \frac{\partial \mathcal{J}}{\partial \theta_i} = -\frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(i,k)} + \frac{\lambda}{n} \theta_i \text{ iff } i \neq 1$$

$$\theta_1 \triangleq \left(1 - \cancel{\alpha \frac{\lambda}{n}}\right) \theta_1 + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(1,k)}$$

$$\theta_2 \triangleq \left(1 - \alpha \frac{\lambda}{n}\right) \theta_2 + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(2,k)}$$

$$\vdots$$

$$\theta_m \triangleq \left(1 - \alpha \frac{\lambda}{n}\right) \theta_m + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(m,k)}$$



## HANDS-ON EXAMPLE (1/3)

**Task #2**

The yield  $y$  of a chemical process is a random variable whose value is considered to be a linear function of the temperature  $x$ . The following data of corresponding values of  $x$  and  $y$  is found:

Temperature in °C ( $x$ )	0	25	50	75	100
Yield in grams ( $y$ )	14	38	54	76	95

The linear regression model  $y = \theta_0 + \theta_1 x$  is used. Determine the values of  $\theta_0$ ,  $\theta_1$ .

1. Using normal equation,
2. Using gradient descent for 5 iterations, given the following initial settings:

$$\alpha = 0.01 \quad \text{and} \quad \theta = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}$$

## HANDS-ON EXAMPLE (2/3)

## ① Normal Equation

$$\mathbf{y} = \begin{bmatrix} 14 \\ 38 \\ 54 \\ 76 \\ 95 \end{bmatrix} \quad \text{and} \quad \mathbf{X} = \begin{bmatrix} 1 & 0 \\ 1 & 25 \\ 1 & 50 \\ 1 & 75 \\ 1 & 100 \end{bmatrix} \implies \hat{\boldsymbol{\theta}} = \begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{bmatrix} = \begin{bmatrix} 15.4 \\ 0.8 \end{bmatrix}$$

## ② Stochastic Gradient Descent

$k$	1	2	3	4	5
$y$	14	38	54	76	95
$h_{\boldsymbol{\theta}}(\mathbf{x}_k)$	1	13.63	330.999	-9894.410	734688.376
$\hat{\boldsymbol{\theta}} = \begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{bmatrix}$	$\begin{bmatrix} 1.13 \\ 0.5 \end{bmatrix}$	$\begin{bmatrix} 1.374 \\ 6.592 \end{bmatrix}$	$\begin{bmatrix} -1.396 \\ -131.907 \end{bmatrix}$	$\begin{bmatrix} 98.308 \\ 7345.901 \end{bmatrix}$	$\begin{bmatrix} -7247.626 \\ -727247.475 \end{bmatrix}$

## HANDS-ON EXAMPLE (3/3)

```
1  X = [1 0; 1 25; 1 50; 1 75; 1 100] # Features
2  y = [14, 38, 54, 76, 95] # Target
3
4   $\alpha$ , n,  $\theta$  = 0.01, 5, [1; .5]
5
6  J = []
7  for k in 1:5
8      h_th = X[k, :]' *  $\theta$ 
9      println("h_th = $(h_th)")
10     cost = (y[k] - h_th)^2
11     push!(J, cost);
12      $\theta$  +=  $\alpha$  * (y[k] - h_th) * X[k, :]
13     println(" $\theta$  = $( $\theta$ )\n")
14 end
```



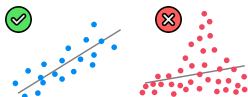
<https://github.com/a-mhamdi/journey-into-ML/blob/main/Julia/gradient-descent.jl>

# Assumptions of Linear Regression



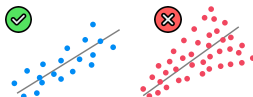
## 1. Linearity

(Linear relationship between Y and each X)



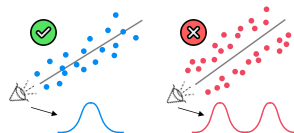
## 2. Homoscedasticity

(Equal variance)



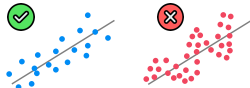
## 3. Multivariate Normality

(Normality of error distribution)



## 4. Independence

(of observations. Includes "no autocorrelation")



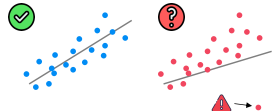
## 5. Lack of Multicollinearity

(Predictors are not correlated with each other)



## 6. The Outlier Check

(This is not an assumption, but an "extra")



## EVALUATION METRICS (1/2)

**Mean Absolute Error (MAE)** measures the average difference of absolute values between predicted and actual targets.

$$\text{MAE} = \frac{1}{n} \sum_{k=1}^n |y_k - \hat{y}_k|$$

**Root Mean Squared Error (RMSE)** measures the root of the average squared difference between predicted and actual values.

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{k=1}^n (y_k - \hat{y}_k)^2}$$

**Mean Absolute Percentage Error (MAPE)** is a measure of the prediction quality. It is equivalent to doing weighted **MAE**.

$$\text{MAPE} = \frac{1}{n} \sum_{k=1}^n \left| \frac{y_k - \hat{y}_k}{y_k} \right| \times 100\%$$

👍 *A lower error indicates a better fit of the model to the data.*

## EVALUATION METRICS (2/2)

**R-squared** is a statistical measure that quantifies the proportion of the variance in the dependent variable that is explained by the independent variables in the model.

$$\mathcal{R}^2 = 1 - \frac{SS_{\text{residuals}}}{SS_{\text{total}}} = 1 - \frac{\sum_{k=1}^n (y_k - \hat{y}_k)^2}{\sum_{k=1}^n (y_k - \bar{y})^2}$$

👍 1 indicates that the model explains **ALL** the variance in the dependent variable

👎 0 indicates that the model explains **NONE** of the variance in the dependent variable

**Adjusted R-squared** is a modified version of R-squared that accounts for the number of independent variables in the model.

$$\text{Adjusted } \mathcal{R}^2 = 1 - (1 - \mathcal{R}^2) \frac{n - 1}{n - m - 1}$$

## HANDS-ON EXAMPLE (1/3)

### Task #3

We consider a univariate regression problem with only two predictors  $m = 2$ . Compute the error metrics for the given data.

$y$	1	1	-2	5	-3.5	1
$\hat{y}$	0.9	0.85	-2.2	4.8	-3.3	1.2

## HANDS-ON EXAMPLE (2/3)

$$\text{MAE: } \frac{1}{6} \sum_{k=1}^6 |y_k - \hat{y}_k| \approx 0.175$$

$$\text{RMSE: } \sqrt{\frac{1}{6} \sum_{k=1}^6 (y_k - \hat{y}_k)^2} \approx 0.179$$

$$\text{MAPE: } \frac{1}{6} \sum_{k=1}^6 \left| \frac{y_k - \hat{y}_k}{y_k} \right| \times 100\% \approx 10.786$$

$$\mathcal{R}^2: 1 - \frac{\sum_{k=1}^6 (y_k - \hat{y}_k)^2}{\sum_{k=1}^6 (y_k - \bar{y})^2} \approx 0.996$$

$$\text{Adjusted } \mathcal{R}^2: 1 - (1 - \mathcal{R}^2) \frac{6 - 1}{6 - 2 - 1} \approx 0.992$$



## HANDS-ON EXAMPLE (3/3)

```

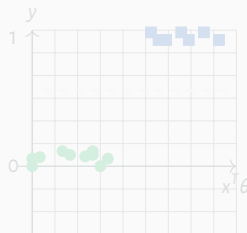
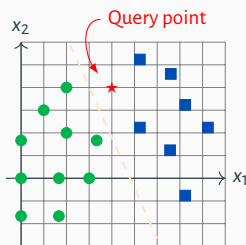
1  y = [1 1 -2 5 -3.5 1]
2  ŷ = [.9 .85 -2.2 4.8 -3.3 1.2]
3
4  mae = 1/6 * sum( abs.(y .- ŷ) ) # 0.17500000000000007
5  rmse = √(1/6 * sum( (y .- ŷ).^2 )) # 0.1791182105017057
6  mape = 1/6 * sum( abs.((y .- ŷ)./y) ) * 100 # 10.785714285714288
7
8  y_bar = 1/6 * sum(y); # 0.41666666666666663 (y_bar = ...)
9  # using Statistics; y_bar = mean(y) (y_bar = mean(y))
10 r2 = 1 - sum( (y .- ŷ).^2 )/sum( (y .- y_bar).^2 ) # 0.9955448408871745
11 adj_r2 = 1 - (1-r2) * (6-1)/(6-2-1) # 0.9925747348119576

```



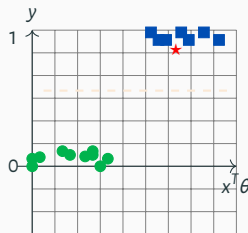
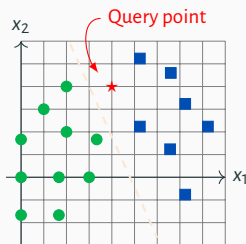
# INTRODUCTION

Classification is a type of supervised machine learning algorithm. A model is trained on a set of *labeled data*, where each data point is associated with a known class or category. The goal of the algorithm is to learn the relationship between the *input features*  $x$  and the corresponding *output classes*  $y$ , so that it can accurately predict the class of **new, unseen query points**.



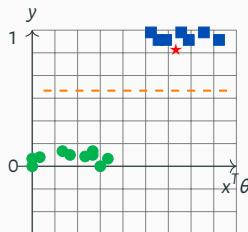
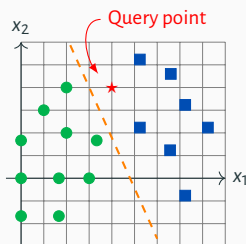
# INTRODUCTION

Classification is a type of supervised machine learning algorithm. A model is trained on a set of *labeled data*, where each data point is associated with a known class or category. The goal of the algorithm is to learn the relationship between the *input features*  $x$  and the corresponding *output classes*  $y$ , so that it can accurately predict the class of **new, unseen query points**.



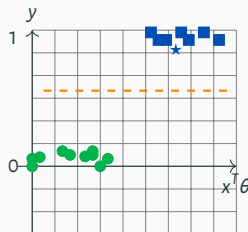
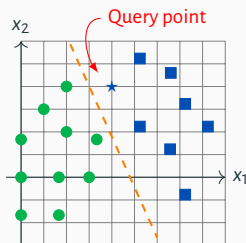
# INTRODUCTION

Classification is a type of supervised machine learning algorithm. A model is trained on a set of *labeled data*, where each data point is associated with a known class or category. The goal of the algorithm is to learn the relationship between the *input features*  $x$  and the corresponding *output classes*  $y$ , so that it can accurately predict the class of **new, unseen query points**.

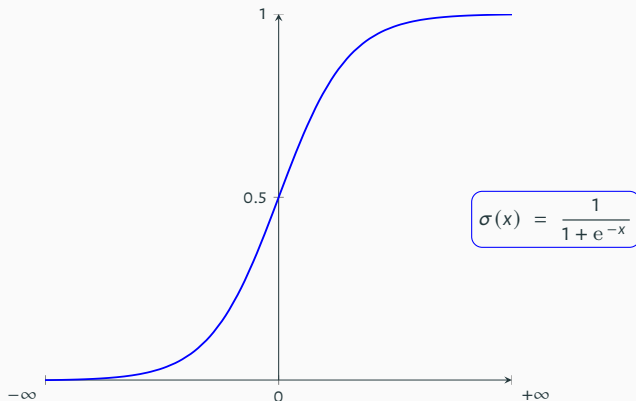


# INTRODUCTION

Classification is a type of supervised machine learning algorithm. A model is trained on a set of *labeled data*, where each data point is associated with a known class or category. The goal of the algorithm is to learn the relationship between the *input features*  $x$  and the corresponding *output classes*  $y$ , so that it can accurately predict the class of **new, unseen query points**.



## LOGISTIC OR S-SHAPED FUNCTION $\sigma$



▲  $\sigma$  squashes range of distance from  $]-\infty, +\infty[$  to  $[0, 1]$

▲  $\sigma$  is differentiable and easy to compute:  $\dot{\sigma} = \sigma \times (1 - \sigma)$

## DECISION BOUNDARY

$$y = \sigma (\theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_m x_m)$$

$$y = \frac{1}{1 + e^{-\mathbf{x}^T \boldsymbol{\theta}}}$$

### Hypothesis

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = P(y = 1 | \mathbf{x}; \boldsymbol{\theta}) = \frac{1}{1 + e^{-\mathbf{x}^T \boldsymbol{\theta}}}$$

For some given  $x_k$

$$h_{\boldsymbol{\theta}}(\mathbf{x}_k) = P(y = 1 | \mathbf{x}_k; \boldsymbol{\theta}) = \frac{1}{1 + e^{-\mathbf{x}_k^T \boldsymbol{\theta}}}$$

### Cost function

$$\mathcal{J} = \begin{cases} -\ln(h_{\boldsymbol{\theta}}(\mathbf{x})) & \text{if } y = 1 \\ -\ln(1 - h_{\boldsymbol{\theta}}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

$$\mathcal{J} = -y \ln(h_{\boldsymbol{\theta}}(\mathbf{x})) - (1 - y) \ln(1 - h_{\boldsymbol{\theta}}(\mathbf{x}))$$

# GRADIENT DESCENT

$$\theta_i \triangleq \theta_i - \underbrace{\alpha}_{\text{LEARNING RATE}} \frac{\partial \mathcal{J}}{\partial \theta_i}$$

Generalizing  $\mathcal{J}$  yields:

$$\mathcal{J} = -\frac{1}{n} \sum_{k=1}^n (y_k \ln(h_{\theta}(\mathbf{x}_k)) + (1 - y_k) \ln(1 - h_{\theta}(\mathbf{x}_k)))$$

$$\Rightarrow \frac{\partial \mathcal{J}}{\partial \theta_i} = -\frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(i,k)}$$

$$\theta_1 \triangleq \theta_1 + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(1,k)}$$

$$\theta_2 \triangleq \theta_2 + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(2,k)}$$

$\vdots$

$$\theta_m \triangleq \theta_m + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(m,k)}$$



## GRADIENT DESCENT

$$\theta_i \triangleq \theta_i - \underbrace{\alpha}_{\text{LEARNING RATE}} \frac{\partial \mathcal{J}}{\partial \theta_i}$$

Generalizing  $\mathcal{J}$  with  $\mathbf{L}_2$  regularization term yields:

$$\mathcal{J} = -\frac{1}{n} \sum_{k=1}^n (y_k \ln(h_{\theta}(\mathbf{x}_k)) + (1 - y_k) \ln(1 - h_{\theta}(\mathbf{x}_k))) + \frac{\lambda}{2n} \sum_{i=2}^m \theta_i^2$$

$$\Rightarrow \frac{\partial \mathcal{J}}{\partial \theta_i} = -\frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(i,k)} + \frac{\lambda}{n} \theta_i \text{ iff } i \neq 1$$

$$\theta_1 \triangleq \left(1 - \cancel{\alpha \frac{\lambda}{n}}\right) \theta_1 + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(1,k)}$$

$$\theta_2 \triangleq \left(1 - \alpha \frac{\lambda}{n}\right) \theta_2 + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(2,k)}$$

$$\vdots$$

$$\theta_m \triangleq \left(1 - \alpha \frac{\lambda}{n}\right) \theta_m + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(\mathbf{x}_k)) x_{(m,k)}$$

## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → *Codes* → *Julia* → *Part-2*

→ *Pluto* → *logistic-regression.jl*

**Pluto.jl** 

→ *Jupyter* → *logistic-regression.ipynb*



# CONFUSION MATRIX

		Actual	
		Positive	Negative
Predicted	Positive	TP	FP
	Negative	FN	TN

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$f1 - score = \frac{2}{\frac{1}{Precision} + \frac{1}{Recall}}$$

183	141
13	663

$$Accuracy = 0.846$$

$$Precision = 0.565$$

$$Recall = 0.934$$

$$f1 - score = 0.704$$

320	20
43	538

$$Accuracy = 0.932$$

$$Precision = 0.941$$

$$Recall = 0.882$$

$$f1 - score = 0.910$$

# CONFUSION MATRIX

		Actual	
		Positive	Negative
Predicted	Positive	<b>TP</b>	<b>FP</b>
	Negative	<b>FN</b>	<b>TN</b>

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$f1 - score = \frac{2}{\frac{1}{Precision} + \frac{1}{Recall}}$$

183	141
13	663

$$Accuracy = 0.846$$

$$Precision = 0.565$$

$$Recall = 0.934$$

$$f1 - score = 0.704$$

320	20
43	538

$$Accuracy = 0.932$$

$$Precision = 0.941$$

$$Recall = 0.882$$

$$f1 - score = 0.910$$

## EVALUATION METRICS

**ACCURACY**

**PRECISION**

**RECALL**

**f1-SCORE**

*Accuracy* denotes the ratio of how many we got right over all cases:

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

*Precision* designates how many positives do we get right over all positive predictions:

$$Precision = \frac{TP}{TP + FP}$$

*Recall* is the ratio of how many positives we got right over all actual positive cases:

$$Recall = \frac{TP}{TP + FN}$$

f1 – score denotes the Harmonic Mean of *Precision* & *Recall*:

$$f1 - score = \frac{2}{\frac{1}{Precision} + \frac{1}{Recall}}$$

# EVALUATION METRICS

## AN ANALOGY

I just explained recall/precision to a non-DS, and he got it immediately : r/datascience - Brave

reddit.com/r/datascience/comments/qdai89/i\_j...

**r/datascience** • 3 yr. ago question\_23

### I just explained recall/precision to a non-DS, and he got it immediately

Discussion

Explain it like fishing with a net. You use a wide net, and catch 80 of 100 total fish in a lake. That's 80% recall. But you also get 80 rocks in your net. That means 50% precision, half of the net's contents is junk. You could use a smaller net and target one pocket of the lake where there are lots of fish and no rocks, but you might only get 20 of the fish in order to get 0 rocks. That is 20% recall and 100% precision.

Seriously, it made me so happy since I've butted against this for years. Equations make people's eyes glaze over, but my PM understood this immediately over a voice call, without diagrams or anything.

Also I googled this and found it's a common explanation, but I'd never heard of it in my 4 years working as a DS.

Archived post. New comments cannot be posted and votes cannot be cast.

1.4K 64 Share

**r/datascience** Data Science  
A space for data science professionals to engage in discussions and debates on the subject of data science.  
Created Aug 6, 2011  
Public  
2.7M Data Scientists 114 Online Top 1% Rank by size

**SUBREDDIT NEWS**  
[Check the Wiki](#)  
[How to get user flair](#)

**RULES**

- 1 Be Respectful and Supportive
- 2 Stay On Topic
- 3 Use the Weekly Thread

## EVALUATION METRICS

### FOLLOW UP



$$f_{\beta} - \text{score} = \frac{1 + \beta^2}{\frac{1}{\text{Precision}} + \frac{\beta^2}{\text{Recall}}}$$

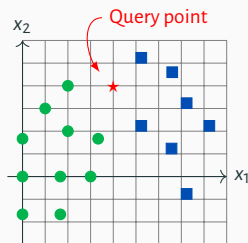
**Case #1:** Prioritize *Precision* over *Recall*, e.g.,  $\beta = 0.5$

- ▶ Mail spam detection
- ▶ Predicting appropriate day to launch a satellite

**Case #2:** Prioritize *Recall* over *Precision*, e.g.,  $\beta = 2$

- ▶ Detection of life threatening diseases like cancer
- ▶ Fraud detection

## *k*-NEAREST NEIGHBORS (1/2)



Minkowski distance

$$d(x; y) = \left( \sum_{i=1}^n |y_i - x_i|^p \right)^{1/p}$$

Manhattan distance (p=1)

$$d(x; y) = \sum_{i=1}^n |y_i - x_i|$$

Euclidean distance (p=2)

$$d(x; y) = \sqrt{\sum_{i=1}^n (y_i - x_i)^2}$$



## $k$ -NEAREST NEIGHBORS (2/2)

► Evelyn Fix and Joseph Hodges, 1951

► Thomas Cover, 1966

---

### Algorithm 1 Summary Construction

---

1: **procedure** HOW DOES  $k$ -NN WORK? (Finding Nearest Neighbors)

**Input:** A query point;

**Output:** Assign a class label to that point.

2:     Define how many neighbors will be checked to classify the specific query point;

3:     Compute the distance  $d(x; y)$  of the query point to other data points;

4:     Count the number of the data points in each category;

5:     Assign the query point to the class with most frequent neighbors.

6: **end procedure**

---

## HANDS-ON EXAMPLE (1/2)

### Task #4

Let be the following coordinate points:

$A(1, 6)$ ;  $B(2, 6)$ ;  $C(3, 1)$ ;  $D(4, 2)$ ;  $E(6, 0)$ ;  $F(7, 5)$ ;  $G(7, 3)$ ;  $H(10, 3)$ ;  $I(-4, -1)$

Using the Euclidean distance, what are the two closest neighbors of point  $P(5, 5)$ ?

$$d(A; P) = \sqrt{17} \approx 4.12 \quad d(B; P) = \sqrt{10} \approx 3.16 \quad d(C; P) = \sqrt{20} \approx 4.47$$

$$d(D; P) = \sqrt{10} \approx 3.16 \quad d(E; P) = \sqrt{26} \approx 5.1 \quad d(F; P) = \sqrt{4} = 2$$

$$d(G; P) = \sqrt{8} \approx 2.83 \quad d(H; P) = \sqrt{29} \approx 5.38 \quad d(I; P) = \sqrt{117} \approx 10.82$$

## HANDS-ON EXAMPLE (2/2)

```
1  function dds(a, b) # 'a' and 'b' are coordinates of some point
2      d_squared = (a-5)^2+(b-5)^2
3      (d_squared, sqrt(d_squared))
4  end
5
6  dds(1, 6) # Point 'A'
7  dds(2, 6) # Point 'B'
```



## HANDS-ON EXAMPLE (1/3)

## Task #5

We try to predict the color of a fruit according to its width ( $w$ ) and height ( $h$ ). The following training data is available:

Fruit	$F_1$	$F_2$	$F_3$	$F_4$	$F_5$	$F_6$	$F_7$	$F_8$
$w$	2	5	2	6	1	4	2	6
$h$	6	6	5	5	2	2	1	1
Color	Red	Yellow	Orange	Purple	Red	Blue	Violet	Green

The goal here is to study the influence of neighbors on the color property of a fruit. Let  $U$  be the new fruit of width  $w = 1$  and height  $h = 4$

1. What is its color if we consider 1 neighbor?
2. What is its color if we consider 3 neighbors?
3. Rather than majority voting, we would like to consider the vote of neighbors weighted by the distance. Each neighbor votes according to a weight inversely proportional to the square of its distance:  $\frac{1}{d^2}$ . We take 3 neighbors, what is the color of  $U$ ? Compare your results to those in question 2.

## HANDS-ON EXAMPLE (2/3)

$$d(U; F_1) = \sqrt{5} \approx 2.24 \quad d(U; F_2) = \sqrt{20} \approx 4.47 \quad d(U; F_3) = \sqrt{2} \approx 1.41$$

$$d(U; F_4) = \sqrt{26} \approx 5.1 \quad d(U; F_5) = \sqrt{4} = 2 \quad d(U; F_6) = \sqrt{13} \approx 3.6$$

$$d(U; F_7) = \sqrt{10} \approx 3.16 \quad d(U; F_8) = \sqrt{34} \approx 5.83$$

1. Color of  $U$  is Orange because  $d(U; F_3)$  is the smallest.
2. Color of  $U$  is Red:  $F_1$  and  $F_5$  (+2 to Red class),  $F_3$  (+1 to Orange class)
3. Color of  $U$  is Orange

$$S(\text{Red}) = \frac{1}{d^2(U; F_1)} + \frac{1}{d^2(U; F_5)} = 0.45$$

$$S(\text{Orange}) = \frac{1}{d^2(U; F_3)} = 0.5$$

## HANDS-ON EXAMPLE (3/3)

```
1  function dds(w, h) # `w` and `h` are width and height of some fruit
2      d_squared = (w-1)^2+(h-4)^2
3      (d_squared, sqrt(d_squared))
4  end
5
6  dds(2, 6) # Fruit `F_1`
7  dds(5, 6) # Fruit `F_2`
```



## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → *Codes* → *Julia* → *Part-2*

→ *Pluto* → *knn.jl*

**Pluto.jl** 

→ *Jupyter* → *knn.ipynb*



## RULE OF THUMB TO CHOOSE $k$

$k$  is **even** if the number of classes is odd

$k$  is **odd** if the number of classes is even

$k$  is an important hyperparameter that can affect the performance of the model.

1. Larger values of  $k$  will result in a smoother decision boundary, which can lead to a more generalized model.
2. Smaller values of  $k$  will result in a more complex decision boundary, which can lead to a model that is more prone to overfitting.
3. The optimal value of  $k$  may depend on the specific dataset and the characteristics of the data.



## SUPPORT VECTOR MACHINES (SVMs)

Support Vector Machines (SVMs) are supervised learning algorithms for classification and regression. They identify the optimal hyperplane to separate data into distinct classes.

**Hyperplane** A decision boundary separating classes in feature space.

**Support Vectors** Closest data points to the hyperplane, determining its position and orientation.

**Margin** Distance between the hyperplane and the nearest data points. SVM maximizes this margin.

For a dataset  $\{(\mathbf{x}_k, y_k)\}_{k=1}^n$ , where  $\mathbf{x}_k \in \mathbb{R}^d$  and  $y_k \in \{-1, +1\}$ :

- Find a hyperplane  $\mathbf{x}^T \mathbf{w} + b = 0$  that maximizes the margin.
- Constrained optimization:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad y_k (\mathbf{x}_k^T \mathbf{w} + b) \geq 1 \quad \forall k$$

### Types of SVMs

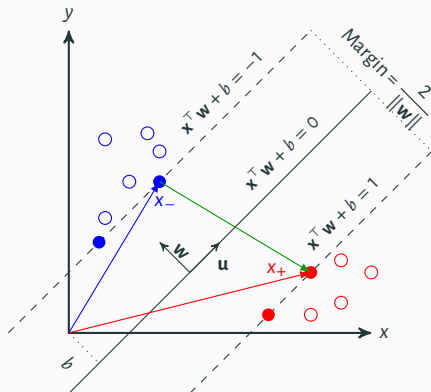
- **Linear SVM:** For linearly separable data.
- **Kernel SVM:** Maps data to higher dimensions for non-linear separation.

# VISUALIZATION OF SVM

$\mathbf{u}$  is a direction vector of  $\Delta$ .

$$\Delta : \mathbf{x}^T \mathbf{w} + b = 0 \implies \mathbf{w} \perp \mathbf{u}$$

$$\begin{aligned} \text{Margin} &= \frac{|(\mathbf{x}_+ - \mathbf{x}_-) \cdot \mathbf{w}|}{\|\mathbf{w}\|} \\ &= \frac{|\mathbf{x}_+^T \mathbf{w} - \mathbf{x}_-^T \mathbf{w}|}{\|\mathbf{w}\|} \\ &= \frac{|(1 - b) - (-1 - b)|}{\|\mathbf{w}\|} \\ &= \frac{2}{\|\mathbf{w}\|} \end{aligned}$$



## SVM: LOSS FUNCTION

**Hinge Loss Function:** The hinge loss is used in SVMs to penalize misclassifications and maximize the margin:

$$L(y, h(\mathbf{x})) = \max(0, 1 - yh(\mathbf{x})),$$

where:

- ▶  $y \in \{-1, +1\}$  is the true label,
- ▶  $h(x) = \mathbf{x}^T \mathbf{w} + b$  is the predicted score.

The loss is zero if  $yh(\mathbf{x}) \geq 1$  (correct classification with margin), and increases linearly otherwise.

### Why SVMs are powerful

- ▶ **Robustness:** Focuses on support vectors, making it less sensitive to outliers.
- ▶ **High-Dimensional Data:** Performs well even with a large number of features.
- ▶ **Kernel Trick:** Enables non-linear classification by mapping data to higher dimensions.

## SOLVING THE LAGRANGIAN PROBLEM IN SVM (1/3)

### Primal Problem (*Hard-Margin SVM*)

The objective is to maximize the margin between two classes. The primal optimization problem is:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_k (\mathbf{x}_k^\top \mathbf{w} + b) \geq 1, \quad \forall k$$

### Construct the Lagrangian

Introduce Lagrange multipliers  $\alpha_k \geq 0$  for each constraint and form the Lagrangian:

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{k=1}^n \alpha_k (y_k (\mathbf{x}_k^\top \mathbf{w} + b) - 1)$$

### Take Derivatives and Set to Zero

Minimize  $\mathcal{L}$  with respect to  $\mathbf{w}$  and  $b$ :

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} - \sum_{k=1}^n \alpha_k y_k \mathbf{x}_k = 0 \implies \mathbf{w} = \sum_{k=1}^n \alpha_k y_k \mathbf{x}_k$$

## SOLVING THE LAGRANGIAN PROBLEM IN SVM (2/3)

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_{k=1}^n \alpha_k y_k = 0 \implies \sum_{k=1}^n \alpha_k y_k = 0$$

Substitute Back into the Lagrangian

Plugging  $\mathbf{w}$  and the constraint into  $\mathcal{L}$ , we obtain the **dual problem**:

$$\mathcal{L}_D(\boldsymbol{\alpha}) = \sum_{k=1}^n \alpha_k - \frac{1}{2} \sum_{k,l=1}^n \alpha_k \alpha_l y_k y_l \mathbf{x}_k^\top \mathbf{x}_l$$

Solve the Dual Problem

Maximize  $\mathcal{L}_D(\boldsymbol{\alpha})$  subject to:

$$\alpha_k \geq 0 \quad \text{and} \quad \sum_{k=1}^n \alpha_k y_k = 0$$

## SOLVING THE LAGRANGIAN PROBLEM IN SVM (3/3)

### Obtain the Solution

- The optimal  $\mathbf{w}^*$  is a linear combination of support vectors:

$$\mathbf{w}^* = \sum_{k \in SV} \alpha_k y_k \mathbf{x}_k$$

- The bias  $b^*$  is computed using any support vector  $\mathbf{x}_k$ :

$$b^* = y_k - \mathbf{w}^{*\top} \mathbf{x}_k$$

### Final Decision Function

For a new input  $\mathbf{x}$ , the prediction is:

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^{*\top} \mathbf{x} + b^*)$$

## HANDS-ON EXAMPLE (1/9)

## Task #6

You are given a dataset with two classes: **Class A** ( $y = +1$ ) and **Class B** ( $y = -1$ ). The features ( $x_1, x_2$ ) and corresponding labels are:

Data Point	1	2	3	4	5	6
$x_1$	2	3	1	6	7	8
$x_2$	2	4	1	2	3	2
$y$	+1	+1	+1	-1	-1	-1

1. Plot the data points on a 2D plane.
2. Identify the support vectors.
3. Determine the optimal hyperplane separating the two classes using SVM.
4. Calculate the margin width if the weight vector is  $\mathbf{w} = [-1/2, 1/4]^T$ .

## Hint!

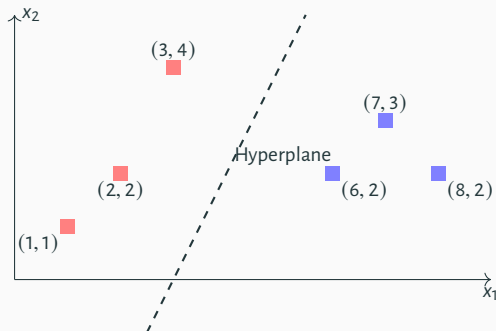
Support vectors are the data points closest to the hyperplane. The margin width is given by:

$$\text{Margin Width} = \frac{2}{\|\mathbf{w}\|}$$

## HANDS-ON EXAMPLE (2/9)

### Q. #1: Plot the Data Points

The data points are visualized on the plane  $(x_1, x_2)$ :





## HANDS-ON EXAMPLE (3/9)

### Q. #2: Support Vectors

The support vectors are the points closest to the hyperplane:

(2, 2) from Class A,   (3, 4) from Class A,   (6, 2) from Class B.

### Q. #3: Optimal Hyperplane

$$\begin{aligned}
 \mathcal{L}_D(\boldsymbol{\alpha}) &= \sum_{k=1}^6 \alpha_k - \frac{1}{2} \sum_{k,l=1}^6 \alpha_k \alpha_l y_k y_l \mathbf{x}_k^T \mathbf{x}_l \\
 &= \alpha_1 + \alpha_2 + \alpha_4 \\
 &\quad - \frac{1}{2} (\alpha_1 [8\alpha_1 + 14\alpha_2 - 16\alpha_4] + \alpha_2 [14\alpha_1 + 25\alpha_2 - 26\alpha_4] - \alpha_4 [16\alpha_1 + 26\alpha_2 - 40\alpha_4])
 \end{aligned}$$

Given

$$\sum_{k=1}^6 \alpha_k y_k = 0 \therefore \alpha_1 + \alpha_2 - \alpha_4 = 0$$

$$\mathcal{L}_D(\boldsymbol{\alpha}) = 2\alpha_1 + 2\alpha_2 - \frac{1}{2} (16\alpha_1^2 + 13\alpha_2^2 + 24\alpha_1\alpha_2)$$

## HANDS-ON EXAMPLE (4/9)

$$\frac{\partial \mathcal{L}_D(\boldsymbol{\alpha})}{\partial \alpha_1} = 0 \quad \therefore \quad 16\alpha_1 + 12\alpha_2 = 2$$

$$\frac{\partial \mathcal{L}_D(\boldsymbol{\alpha})}{\partial \alpha_2} = 0 \quad \therefore \quad 12\alpha_1 + 13\alpha_2 = 2$$

The resulting linear system is:

$$\begin{bmatrix} 16 & 12 \\ 12 & 13 \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$\boldsymbol{\alpha}$  is:

$$\boldsymbol{\alpha}^T = \begin{bmatrix} 1/32 & 1/8 & 0 & 5/32 & 0 & 0 \end{bmatrix}$$

- The optimal  $\mathbf{w}^*$  is a linear combination of support vectors:

$$\mathbf{w}^* = 1/32 \begin{bmatrix} 2 \\ 2 \end{bmatrix} + 1/8 \begin{bmatrix} 3 \\ 4 \end{bmatrix} - 5/32 \begin{bmatrix} 6 \\ 2 \end{bmatrix} = \begin{bmatrix} -1/2 \\ 1/4 \end{bmatrix}$$

## HANDS-ON EXAMPLE (5/9)

- The bias  $b^*$  is computed using any support vector  $\mathbf{x}_k$ :

$$b^* = y_k - \mathbf{x}_k^T \mathbf{w}^* = 3/2$$

The separating hyperplane is given by:

$$x_1 w_1 + x_2 w_2 + b = 0 \quad \text{where} \quad \mathbf{w} = \begin{bmatrix} -1/2 \\ 1/4 \end{bmatrix}^T, \quad b = 3/2$$

Thus, the equation becomes:

$$-2x_1 + x_2 + 6 = 0$$

### Q. #4: Margin Width

The margin width is calculated as:

$$\text{Margin Width} = \frac{2}{\|\mathbf{w}\|} = \frac{2}{\sqrt{(1/2)^2 + (1/4)^2}} = \frac{8}{\sqrt{5}}$$

## HANDS-ON EXAMPLE (6/9)

```

1  using DataFrames, MLJ
2  using Plots
3
4  df = DataFrame(
5      x1 = [2, 3, 1, 6, 7, 8],
6      x2 = [2, 4, 1, 2, 3, 2],
7      y = [1, 1, 1, -1, -1, -1]
8  )
9
10 scatter(df.x1, df.x2, group=df.y)
11
12 schema(df)
13 #=
14
15 |  names  |  scitypes  |  types  |
16 |-----|-----|-----|
17 |  x1     |  Count     |  Int64  |
18 |  x2     |  Count     |  Int64  |
19 |  y      |  Count     |  Int64  |
20 |-----|-----|-----|

```

## HANDS-ON EXAMPLE (7/9)

```

21  =#
22
23  coerce!(df, :x1 => Continuous, :x2 => Continuous, :y => OrderedFactor)
24  scitype(df.y) <: AbstractVector{<:OrderedFactor} # true
25  schema(df)
26  #=
27  |-----|
28  | names | scitypes          | types          |
29  |-----|-----|-----|
30  | x1    | Continuous          | Float64       |
31  | x2    | Continuous          | Float64       |
32  | y     | OrderedFactor{2}    | CategoricalValue{Int64, UInt32} |
33  |-----|-----|-----|
34  =#
35
36  X = select(df, [:x1, :x2])
37  y = df.y
38
39  "[LIBSVM](https://www.csie.ntu.edu.tw/~cjlin/libsvm/)"
40  LinearSVC = @load LinearSVC pkg=LIBSVM

```

## HANDS-ON EXAMPLE (8/9)

```

41  # doc("LinearSVC", pkg="LIBSVM")
42
43  svc = LinearSVC()
44  #=
45  LinearSVC(
46    solver = LIBSVM.LinearSolver.L2R_L2LOSS_SVC_DUAL,
47    tolerance = Inf,
48    cost = 1.0,
49    bias = -1.0)
50  =#
51
52  mach_svc = machine(svc, X, y) |> fit!
53  fitted_params(mach_svc)
54  #=
55  (libsvm_model = LIBLINEAR.LinearModel{UInt32}(1, 2, 2, [-0.486189670915792, 0.
    ↳8956996119492998], Int32[1, 2], UInt32[0x00000002, 0x00000001], -1.0, 0.0),
56    encoding = Dict{UInt32, CategoricalArrays.CategoricalValue{String, UInt32}}
    ↳((0x00000002 => "1", 0x00000001 => "-1"),)
57  =#
58

```

## HANDS-ON EXAMPLE (9/9)

```

59  using LIBSVM
60  Data = Matrix(X)'
61  mdl = svmtrain(D, y, kernel=Kernel.Linear)
62  mdl.SVs
63  #=
64  LIBSVM.SupportVectors{CategoricalArrays.CategoricalVector{Int64, UInt32, Int64,
CategoricalArrays.CategoricalValue{Int64, UInt32}, Union{}}}, Matrix{Float64}}(3,
Int32[2, 1], CategoricalArrays.CategoricalValue{Int64, UInt32}[1, 1, -1], [2.0 3.
↪ 0 6.0; 2.0 4.0 2.0], Int32[1, 2, 4], LIBSVM.SVMNode[LIBSVM.SVMNode(1, 2.0),
LIBSVM.SVMNode(1, 3.0), LIBSVM.SVMNode(1, 6.0)])
65  =#

```



# HARD MARGIN VS. SOFT MARGIN SVM

## ► Data Assumption

**Hard:** Strict linear separability

**Soft:** Tolerates noise/overlap

## ► Primal Form

$$\textbf{Hard: } \min \frac{1}{2} \|w\|^2 \quad \text{s.t. } y_k (\mathbf{x}_k^\top \mathbf{w} + b) \geq 1$$

$$\textbf{Soft: } \min \frac{1}{2} \|w\|^2 + C \sum \xi_k \quad \text{s.t. } y_k (\mathbf{x}_k^\top \mathbf{w} + b) \geq 1 - \xi_k$$

## ► Key Features

**Hard:** No slack variables ( $\xi$ ),  $C \rightarrow \infty$

**Soft:** Slack variables ( $\xi_i \geq 0$ ), finite  $C$  (trade-off)

## ► Dual Form

**Hard:**  $0 \leq \alpha_i$

**Soft:**  $0 \leq \alpha_i \leq C$

**Both:**  $\sum \alpha_i y_i = 0$

## ► Practicality

**Hard:** Rarely used (overfits)

**Soft:** Robust to noise, works with kernels



# KERNEL TRICK IN SVMs

Dual Form with Kernel:

$$\mathcal{L}_D(\boldsymbol{\alpha}) = \sum_{k=1}^n \alpha_k - \frac{1}{2} \sum_{k,l=1}^n \alpha_k \alpha_l y_k y_l \underbrace{\mathbf{x}_k^\top \mathbf{x}_l}_{\mathcal{K}(\mathbf{x}_k, \mathbf{x}_l)}$$

$$\mathcal{K}(\mathbf{x}_k, \mathbf{x}_l) = \Phi^\top(\mathbf{x}_k) \Phi(\mathbf{x}_l)$$

**Kernel Trick:** Replace inner products with kernel functions - no need to compute  $\Phi(\mathbf{x})$  explicitly.

Common Kernels:

- ▶ **Linear:**  $\mathcal{K}(\mathbf{x}_k, \mathbf{x}_l) = \mathbf{x}_k^\top \mathbf{x}_l$
- ▶ **Polynomial:**  $\mathcal{K}(\mathbf{x}_k, \mathbf{x}_l) = \left( \gamma \mathbf{x}_k^\top \mathbf{x}_l + r \right)^d$
- ▶ **RBF (Gaussian):**  $\mathcal{K}(\mathbf{x}_k, \mathbf{x}_l) = \exp \left( -\gamma \|\mathbf{x}_k - \mathbf{x}_l\|^2 \right)$
- ▶ **Sigmoid:**  $\mathcal{K}(\mathbf{x}_k, \mathbf{x}_l) = \tanh \left( \gamma \mathbf{x}_k^\top \mathbf{x}_l + r \right)$

Parameters

$d$ : Polynomial degree (#non-linear terms)

$\gamma$ : Kernel scale (RBF/poly)

$C$ : Regularization (controls margin vs. errors)

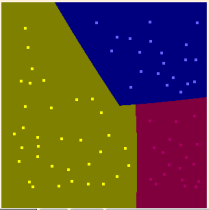
# LIBSVM LIBRARY

LIBSVM -- A Library for Support Vector Machines

Please read the [COPYRIGHT](#) notice before using LIBSVM.

## Graphic Interface

Here is a simple applet demonstrating SVM classification and regression.  
Click on the drawing area and use "Change" to change class of data. Then use "Run" to see the results.



Change Run Clear -t 2 -c 100

Examples of options: -s 0 -c 10 -t 1 -g 1 -r 1 -d 3  
Classify a binary data with polynomial kernel  $(u^v+1)^3$  and  $C = 10$

<https://www.csie.ntu.edu.tw/~cjlin/libsvm/>

## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → *Codes* → *Julia* → *Part-2*

→ *Pluto* → *svc.jl*

**Pluto.jl** 

→ *Jupyter* → *svc.ipynb*



# NAIVE BAYES

## COVID-19 TESTING

### Task #7

Imagine we have the following information:

- ▶ 1% of the population has COVID-19 (prior probability)
- ▶ A COVID test has 95% sensitivity (true positive rate) - if you have COVID, there's a 95% chance the test will be positive
- ▶ The test has 90% specificity (true negative rate) - if you don't have COVID, there's a 90% chance the test will be negative

Now, if someone receives a positive test result, what's the probability they actually have COVID-19?

# NAIVE BAYES

## COVID-19 TESTING

### Task #7

Imagine we have the following information:

- ▶ 1% of the population has COVID-19 (prior probability)
- ▶ A COVID test has 95% sensitivity (true positive rate) - if you have COVID, there's a 95% chance the test will be positive
- ▶ The test has 90% specificity (true negative rate) - if you don't have COVID, there's a 90% chance the test will be negative

Now, if someone receives a positive test result, what's the probability they actually have COVID-19?

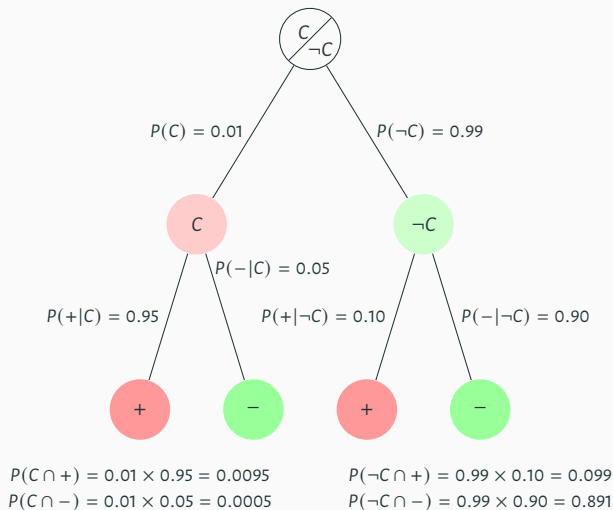
$$P(+) = P(C \cap +) + P(\neg C \cap +) = 0.0095 + 0.099 = 0.1085$$

$$P(C|+) = \frac{P(C \cap +)}{P(+)} = \frac{0.0095}{0.1085} \approx 0.088 \text{ (8.8\%)}$$

∴ Despite a positive test, there's only an 8.8% chance the person actually has COVID-19.

# NAIVE BAYES

## COVID-19 TESTING



# NAIVE BAYES (1/2)

## DEFINITION AND KEY CONCEPT

Naive Bayes is a probabilistic machine learning algorithm based on Bayes' Theorem. It is commonly used for classification tasks, such as spam detection and text classification.

### Bayes' Theorem:

$$P(y|X) = \frac{P(X|y) \cdot P(y)}{P(X)}$$

where:

- ▶  $P(y|X)$ : Posterior probability of class  $y$  given features  $X$ .
- ▶  $P(X|y)$ : Likelihood of features  $X$  given class  $y$ .
- ▶  $P(y)$ : Prior probability of class  $y$ .
- ▶  $P(X)$ : Marginal probability of features  $X$ .

$$\hat{y} = \arg \max_{y \in Y} P(y|X) = \arg \max_{y \in Y} P(X|y) \cdot P(y)$$

## NAIVE BAYES (2/2)

### DEFINITION AND KEY CONCEPT



#### NAIVE BAYES FOR SPAM DETECTION

**Features:** Words in an email (e.g., “free”, “offer”, “money”).

**Classes:** “Spam” ( $S$ ) or “Not Spam” ( $\neg S$ ).

**Steps:**

- ① Calculate  $P(S)$  and  $P(\neg S)$  from the training data.
- ② Calculate  $P(\text{Word}|S)$  and  $P(\text{Word}|\neg S)$  for each word.
- ③ Use Bayes' Theorem to predict the class of a new email.



## HANDS-ON EXAMPLE (1/7)

### Task #8

You are given a dataset of emails with the following word frequencies and labels:

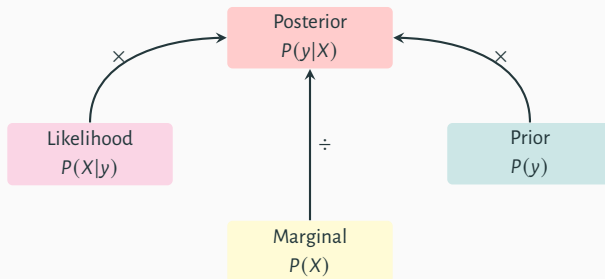
Email	“free”	“money”	Label
1	1	0	$S$
2	0	1	$S$
3	0	0	$\neg S$
4	0	1	$\neg S$

1. Calculate the prior probabilities  $P(S)$  and  $P(\neg S)$ .
2. Calculate the likelihoods  $P(\text{“free”}|S)$ ,  $P(\text{“money”}|S)$ ,  $P(\text{“free”}|\neg S)$ , and  $P(\text{“money”}|\neg S)$ .
3. Use Naive Bayes to classify a new email with the words “free” and “money”.

### Hint!

- Use Laplace smoothing (add +1 smoothing) to handle zero probabilities.

## HANDS-ON EXAMPLE (2/7)



### Step 1: Calculate Prior Probabilities

- $P(S) = \frac{\text{Number of Spam Emails}}{\text{Total Emails}} = \frac{2}{4} = 0.5$
- $P(\neg S) = \frac{\text{Number of Not Spam Emails}}{\text{Total Emails}} = \frac{2}{4} = 0.5$

## HANDS-ON EXAMPLE (3/7)

## Step 2: Calculate Likelihoods with Laplace Smoothing

$$P(w_i|c) = \frac{\text{count}(w_i, c) + k}{\sum_{w \in V} \text{count}(w, c) + k|V|}$$

$k$ : Smoothing parameter ( $k \geq 0$ )

$|V|$ : Vocabulary size

**count** ( $w_i, c$ ): Frequency of word  $w_i$  in class  $c$

$$\blacktriangleright P(\text{"free"}|S) = \frac{\text{Count of "free" in Spam} + 1}{\text{Count of Words in Spam} + |\text{Vocabulary}|} = \frac{1 + 1}{2 + 2} = \frac{2}{4} = 0.5$$

$$\blacktriangleright P(\text{"money"}|S) = \frac{\text{Count of "money" in Spam} + 1}{\text{Count of Words in Spam} + |\text{Vocabulary}|} = \frac{1 + 1}{2 + 2} = \frac{2}{4} = 0.5$$

$$\blacktriangleright P(\text{"free"}|\neg S) = \frac{\text{Count of "free" in Not Spam} + 1}{\text{Count of Words in } \neg \text{Spam} + |\text{Vocabulary}|} = \frac{0 + 1}{1 + 2} = \frac{1}{3} \approx 0.333$$

$$\blacktriangleright P(\text{"money"}|\neg S) = \frac{\text{Count of "money" in Not Spam} + 1}{\text{Count of Words in } \neg \text{Spam} + |\text{Vocabulary}|} = \frac{1 + 1}{1 + 2} = \frac{2}{3} \approx 0.666$$

## HANDS-ON EXAMPLE (4/7)

### Step 3: Classify a New Email

- For the email with words “free” and “money”:

$$\begin{aligned}P(S|\text{“free” and “money”}) &\propto P(\text{“free”}|S) \cdot P(\text{“money”}|S) \cdot P(S) \\&= 1/2 \cdot 1/2 \cdot 1/2 = 1/6\end{aligned}$$

$$\begin{aligned}P(\neg S|\text{“free” and “money”}) &\propto P(\text{“free”}|\neg S) \cdot P(\text{“money”}|\neg S) \cdot P(\neg S) \\&= 1/3 \cdot 2/3 \cdot 1/2 = 1/9\end{aligned}$$

- Since  $1/6 > 1/9$ , the email is classified as **Spam**.

## HANDS-ON EXAMPLE (5/7)

```

1  using DataFrames, MLJ
2
3  df = DataFrame(
4      free = [1, 0, 0, 0],
5      money = [0, 1, 0, 1],
6      label = ["S", "S", "NS", "NS"]
7  )
8  schema(df)
9  #=
10
11  | names | scitypes | types |
12  |-----|-----|
13  | free  | Count    | Int64  |
14  | money | Count    | Int64  |
15  | label | Textual  | String |
16  |-----|-----|
17  =#
18
19  coerce!(df, :label => OrderedFactor)
20  scitype(df.label) <: AbstractVector{<:OrderedFactor} # true

```

## HANDS-ON EXAMPLE (6/7)

```

21  schema(df)
22  #=
23  |
24  | names | scitypes | types |
25  |-----|-----|-----|
26  | free  | Count    | Int64 |
27  | money | Count    | Int64 |
28  | label | OrderedFactor{2} | CategoricalValue{String, UInt32} |
29  |-----|-----|-----|
30  =#
31
32  X = select(df, [:free, :money])
33  y = df.label
34
35  MultinomialNBClassifier = @load MultinomialNBClassifier pkg=NaiveBayes
36  mnb = MultinomialNBClassifier()
37  #=
38  MultinomialNBClassifier(
39    alpha = 1)
40  =#

```

## HANDS-ON EXAMPLE (7/7)

```
41
42 mach_nb = machine(mnb, X, y) |> fit!
43 predict(mach_nb, DataFrame(free = 1, money = 1))
44 #=
45 1-element CategoricalDistributions.UnivariateFiniteVector{Multiclass{2}, String,
46 UInt32, Float64}:
47  UnivariateFinite{Multiclass{2}}(NS=>0.471, S=>0.529)
48 =#
```



# NAIVE BAYES: KEY ASSUMPTIONS AND TYPES

## Key Assumptions

**Independence:** Features are assumed to be conditionally independent given the class.

**Equal Importance:** All features contribute equally to the prediction.

## Types

**Gaussian Naive Bayes:** Assumes features follow a normal distribution.

**Multinomial Naive Bayes:** Used for discrete data (e.g., word counts).

**Bernoulli Naive Bayes:** Used for binary features.

- ▲ Simple and fast to train.
- ▲ Works well with high-dimensional data (e.g., text).
- ▲ Requires less training data compared to other algorithms.
- ▼ Assumes feature independence, which is rarely true in real-world data.
- ▼ Struggles with zero probabilities (requires smoothing techniques).
- ▼ Less accurate for complex relationships between features.



## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → *Codes* → *Julia* → *Part-2*

→ *Pluto* → *naive-bayes.jl*

**Pluto.jl** 

→ *Jupyter* → *naive-bayes.ipynb*



## DECISION TREE: DEFINITION AND KEY CONCEPTS

A decision tree is a supervised learning algorithm used for classification and regression. It partitions the feature space into regions by applying a series of decision rules.

### Components

**Root Node:** Represents the entire dataset.

**Decision Nodes:** Split the dataset based on a feature and a threshold.

**Leaf Nodes:** Represent the final output (class label or regression value).

### Key Concepts

**Splitting Criteria:** Choose the feature and threshold that maximize information gain or minimize impurity.

**Pruning:** Remove unnecessary branches to prevent overfitting.

**Feature Importance:** Rank features based on their contribution to the model.

## GINI INDEX VS ENTROPY

### Gini Index

- ▶ Measures the impurity of a node.
- ▶ Lower values indicate purer nodes.

$$G = 1 - \sum_{i=1}^k p_i^2$$

### Entropy

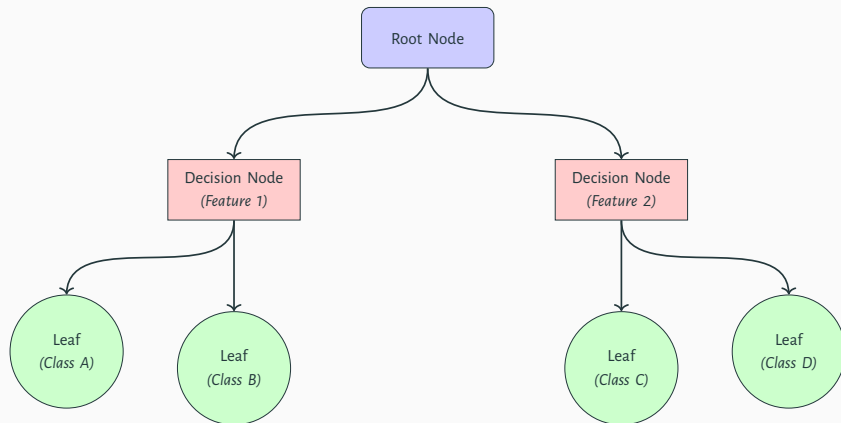
- ▶ Measures the uncertainty of a node.
- ▶ Lower values indicate more certainty.

$$H = - \sum_{i=1}^k p_i \log(p_i)$$

### Comparison

- ▶ Both are used for splitting criteria in decision trees;
- ▶ Gini Index is faster to compute;
- ▶ Entropy provides more precise splits in some cases.

## DECISION TREE VISUALIZATION



## HANDS-ON EXAMPLE (1/13)

### Task #9

Considering the following restaurant wait time dataset, we need to predict if a reservation will be honored (Yes/No). Construct a **decision tree** using Gini Index:

Party Size	Day of Week	Special Occasion?	Reservation Honored?
2	Weekday	No	Yes
4	Weekend	Yes	No
6	Weekend	Yes	No
2	Weekend	No	Yes
4	Weekday	No	Yes
8	Weekend	Yes	No
2	Weekday	Yes	Yes
4	Weekend	No	No
6	Weekday	No	Yes
2	Weekend	Yes	No

## HANDS-ON EXAMPLE (2/13)

### Step 1: Root Node Gini Impurity

$$G_{\text{Root}} = 1 - \left( \left( \frac{5}{10} \right)^2 + \left( \frac{5}{10} \right)^2 \right) = 1 - (0.25 + 0.25) = 0.5$$

### Step 2: Feature Split Analysis

Split 1: "Party Size"

$$G_{\text{ps}=2} = 1 - \left( \left( \frac{3}{4} \right)^2 + \left( \frac{1}{4} \right)^2 \right) = 0.375$$

$$G_{\text{ps}=4} = 1 - \left( \left( \frac{1}{3} \right)^2 + \left( \frac{2}{3} \right)^2 \right) = 0.444$$

$$G_{\text{ps}=6} = 1 - \left( \left( \frac{1}{2} \right)^2 + \left( \frac{1}{2} \right)^2 \right) = 0.5$$

$$G_{\text{ps}=8} = 1 - \left( \left( \frac{0}{1} \right)^2 + \left( \frac{1}{1} \right)^2 \right) = 0$$

$$G_{\text{split}} = \frac{4}{10} \times 0.375 + \frac{3}{10} \times 0.444 + \frac{2}{10} \times 0.5 + \frac{1}{10} \times 0 = 0.383$$

## HANDS-ON EXAMPLE (3/13)

*Split 2: "DOTW"*

$$G_{\text{weekday}} = 1 - \left( \left( \frac{4}{4} \right)^2 + \left( \frac{0}{4} \right)^2 \right) = 0$$

$$G_{\text{weekend}} = 1 - \left( \left( \frac{1}{6} \right)^2 + \left( \frac{5}{6} \right)^2 \right) = 0.278$$

$$G_{\text{split}} = \frac{4}{10} \times 0 + \frac{6}{10} \times 0.278 = 0.167$$

*Split 3: "Special Occasion?"*

$$G_{\text{special}} = 1 - \left( \left( \frac{1}{5} \right)^2 + \left( \frac{4}{5} \right)^2 \right) = 0.32$$

$$G_{\text{no special}} = 1 - \left( \left( \frac{4}{5} \right)^2 + \left( \frac{1}{5} \right)^2 \right) \approx 0.32$$

$$G_{\text{split}} = \frac{5}{10} \times 0.32 + \frac{5}{10} \times 0.32 \approx 0.32$$

## HANDS-ON EXAMPLE (4/13)

### Step 3: Optimal Split Selection

- ▶ **DOTW** has the lowest Gini impurity ( $G = 0.167$ )
- ▶ **Party Size?** ( $G = 0.383$ ) and **Special Occasion?** ( $G \approx 0.32$ ) are less optimal

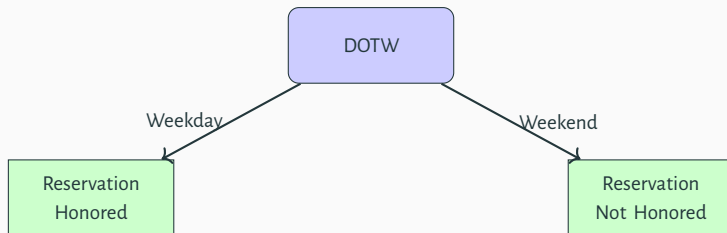
### RESULTING DECISION TREE

**1<sup>st</sup> split:** "Weekend?"

**Weekday:** Predict "Reservation Honored" (4/4 samples) (*pure leaf node*)

**Weekend:** Predict "Reservation Not Honored" (5/6 samples)

**Subsequent splits:** Refine with "Party Size" or "Special Occasion?"





## HANDS-ON EXAMPLE (5/13)

Party Size	SO?	RH?
2	No	Yes
4	No	Yes
2	Yes	Yes
6	No	Yes

## Weekday Node Gini Impurity

$$G_{\text{Weekday}} = 1 - \left( \left( \frac{4}{4} \right)^2 + \left( \frac{0}{4} \right)^2 \right) = 0$$

Party Size	SO?	RH?
4	Yes	No
6	Yes	No
2	No	Yes
8	Yes	No
4	No	No
2	Yes	No

## Weekend Node Gini Impurity

$$G_{\text{Weekend}} = 1 - \left( \left( \frac{1}{6} \right)^2 + \left( \frac{5}{6} \right)^2 \right) = 0.278$$

## HANDS-ON EXAMPLE (6/13)

*“Party Size?” Split*

$$G_{ps=2} = 1 - \left( \left( \frac{1}{2} \right)^2 + \left( \frac{1}{2} \right)^2 \right) = 0.5$$

$$G_{ps=4} = 1 - \left( \left( \frac{0}{2} \right)^2 + \left( \frac{2}{2} \right)^2 \right) = 0$$

$$G_{ps=6} = 1 - \left( \left( \frac{0}{1} \right)^2 + \left( \frac{1}{1} \right)^2 \right) = 0$$

$$G_{ps=8} = 1 - \left( \left( \frac{0}{1} \right)^2 + \left( \frac{1}{1} \right)^2 \right) = 0$$

$$G_{split} = \frac{2}{6} \times 0.5 + \frac{2}{6} \times 0 + \frac{1}{6} \times 0 + \frac{1}{6} \times 0 = 0.167$$

## HANDS-ON EXAMPLE (7/13)

*“Special Occasion?” Split*

$$G_{\text{special}} = 1 - \left( \left( \frac{0}{4} \right)^2 + \left( \frac{4}{4} \right)^2 \right) = 0$$

$$G_{\text{no special}} = 1 - \left( \left( \frac{1}{2} \right)^2 + \left( \frac{1}{2} \right)^2 \right) \approx 0.5$$

$$G_{\text{split}} = \frac{4}{6} \times 0 + \frac{2}{6} \times 0.5 \approx 0.167$$

## HANDS-ON EXAMPLE (8/13)

### Optimal Split Selection

Both “Party Size” and “Special Occasion?” have the same Gini Index (0.1667) for the “Weekend” branch.

In this case, we can choose either one. Let's use “Special Occasion?” as our second split.

### RESULTING DECISION TREE

**2<sup>nd</sup> split:** Special Occasion?

**Yes:** Predict “Reservation Not Honored” (4 samples)

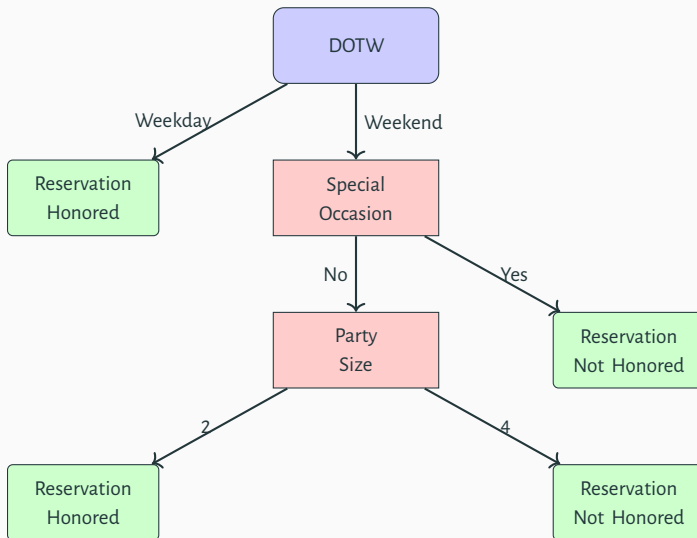
**No:** Refine with “Party Size” (2 samples)

**3<sup>rd</sup> split:** Party Size (*pure leaf*)

**2:** Predict “Reservation Honored” (1/2 sample)

**4:** Predict “Reservation Not Honored” (1/2 sample)

## HANDS-ON EXAMPLE (9/13)



## HANDS-ON EXAMPLE (10/13)

```

1  using DataFrames, MLJ
2
3  df = DataFrame(
4      PS = [2, 4, 6, 2, 4, 8, 2, 4, 6, 2],
5      DOTW = ["Weekday", "Weekend", "Weekend", "Weekend", "Weekday", "Weekend",
6      ↪ "Weekday", "Weekend", "Weekday", "Weekend"],
7      SO = ["No", "Yes", "Yes", "No", "No", "Yes", "Yes", "No", "No", "Yes"],
8      RH = ["Yes", "No", "No", "Yes", "Yes", "No", "Yes", "No", "Yes", "No"]
9  )
10 schema(df)
11
12  #=
13
14  | names | scitypes | types |
15  |-----|-----|-----|
16  | PS    | Count    | Int64  |
17  | DOTW  | Textual  | String |
18  | SO    | Textual  | String |
19  | RH    | Textual  | String |
20  |-----|-----|-----|
21
22  =#

```

## HANDS-ON EXAMPLE (11/13)

```

20
21 coerce!(df,
22     :PS => Continuous,
23     :DOTW => OrderedFactor,
24     :SO => OrderedFactor,
25     :RH => OrderedFactor
26 )
27 schema(df)
28 #=
29 | names | scitypes | types |
30 |-----|-----|-----|
31 | PS    | Continuous | Float64 |
32 | DOTW  | OrderedFactor{2} | CategoricalValue{String, UInt32} |
33 | SO    | OrderedFactor{2} | CategoricalValue{String, UInt32} |
34 | RH    | OrderedFactor{2} | CategoricalValue{String, UInt32} |
35 |-----|-----|-----|
36
37 =#
38
39 X = select(df, Not(:RH))

```

## HANDS-ON EXAMPLE (12/13)

```

40 y = df.RH
41
42 ohe = OneHotEncoder(drop_last=True)
43 mach_ohe = machine(ohe, X) |> fit!
44 W = MLJ.transform(mach_ohe, X)
45 first(W, 2)
46
47 #=
48 2×3 DataFrame
49 Row | PS          DOTW__Weekday  SO__No
50    | Float64  Float64          Float64
51 1 | 2.0          1.0          1.0
52 2 | 4.0          0.0          0.0
53 =#
54
55 DecisionTreeClassifier = @load DecisionTreeClassifier pkg=DecisionTree
56 dt = DecisionTreeClassifier() # an instance of 'DecisionTreeClassifier'
57
58 mach_tree = machine(dt, W, y)
59 fit!(mach_tree, verbosity=2)

```



## HANDS-ON EXAMPLE (13/13)

```

60 mach_tree.report[:fit].print_tree(mach_tree.fitresult[1])
61 #=
62 Feature 2: "DOTW_Weekday" < 0.5 ?
63 └─ Feature 3: "SO_No" < 0.5 ?
64     └─ 1 : 4/4
65         └─ Feature 1: "PS" < 3.0 ?
66             └─ 2 : 1/1
67                 └─ 1 : 1/1
68         └─ 2 : 4/4
69     =#

```



## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → Codes → Julia → Part-2

→ Pluto → *decision-tree\*.jl*

**Pluto.jl** 

→ Jupyter → *decision-tree\*.ipynb*



## ENSEMBLE APPROACH: RANDOM FOREST (1/2)

A random forest is an ensemble method that combines multiple decision trees to improve classification accuracy. It uses bagging (bootstrap aggregation) and random feature selection.

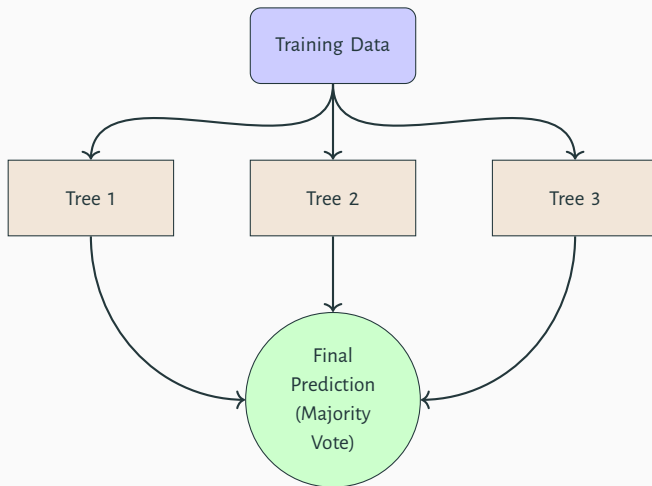
### Key Features

- ▶ **Diversity:** Trees are trained on different subsets of data and features.
- ▶ **Voting Mechanism:** For classification, the final output is the mode of the tree predictions.

### Loss Function

- ▶ Combines the loss functions of individual trees, e.g., Gini Index or Entropy, during splits.

## ENSEMBLE APPROACH: RANDOM FOREST (2/2)



## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → Codes → Julia → Part-2

→ Pluto → *random-forest-\*.jl*

**Pluto.jl** 

→ Jupyter → *random-forest-\*.ipynb*

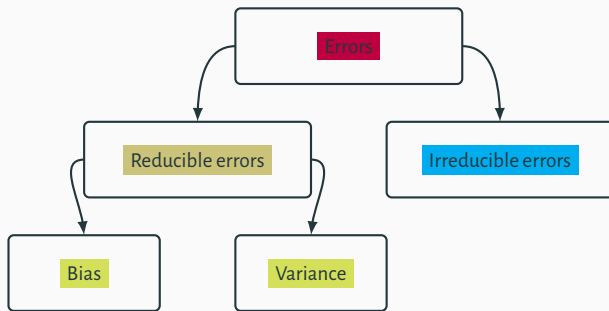


## SUMMARY

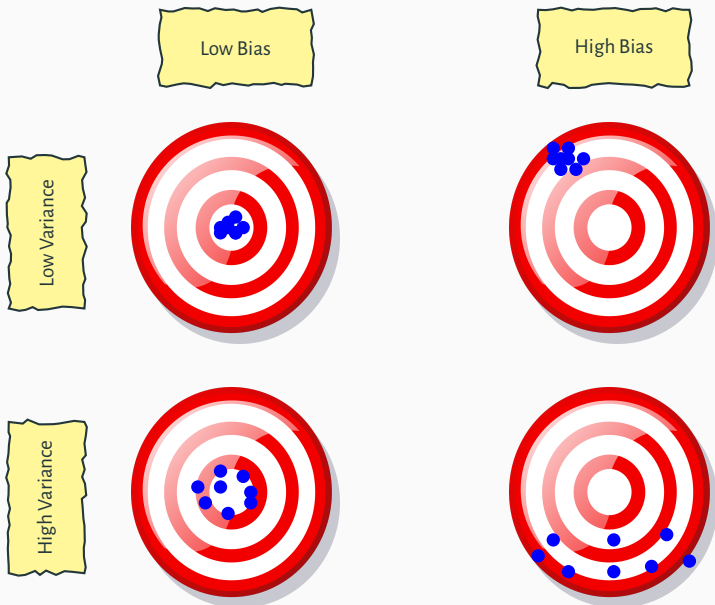
Method		Pros		Cons
<i>Logistic Regression</i>	▲	Probabilistic	▼	Almost linearly separable data
<i>k-NN</i>	▲	Fast and efficient	▼	# of neighbors $k$
			▼	Detecting outliers <sup>2</sup>
<i>SVM</i>	▲	Memory efficient	▼	Kernel's choice
	▲	Versatile	▼	Large datasets
	▲	Noise and outliers	▼	Overlapping classes
	▲	High dimension	▼	Interpretability
<i>Naive Bayes</i>	▲	Simplicity and efficiency	▼	Independence between features
	▲	High dimension	▼	∃ of irrelevant features
<i>Decision Tree</i>	▲	Interpretability	▼	Overfitting
	▲	Numerical and categorical data	▼	Unstable
	▲	Robust to outliers	▼	Continuous variables
	▲	High accuracy	▼	# of input features
<i>Random Forest</i>	▲	Less prone to overfitting	▼	Computation
	▲	High dimension	▼	Interpretability

<sup>2</sup>Points that differ significantly from the rest of the data points.

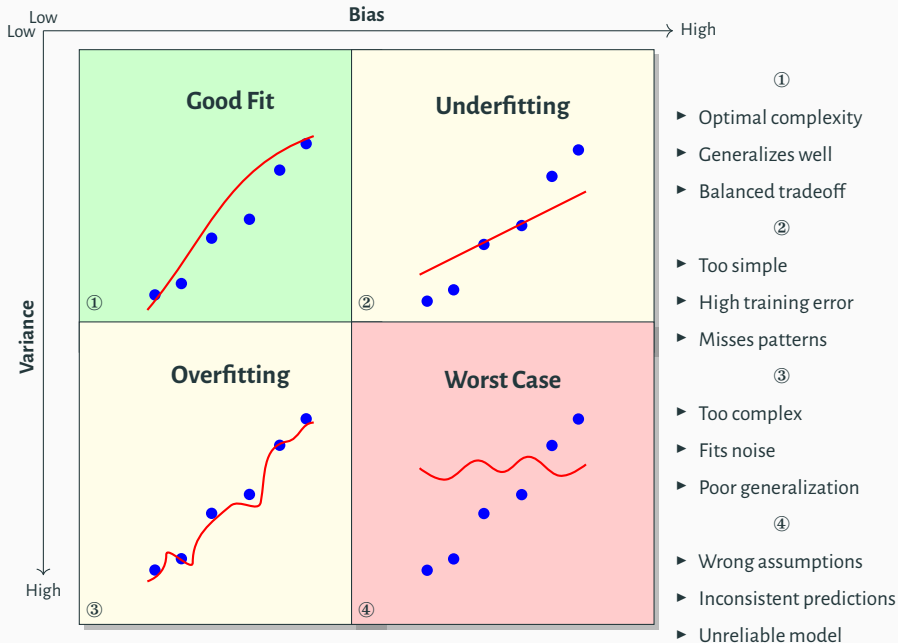
## ERRORS IN ML



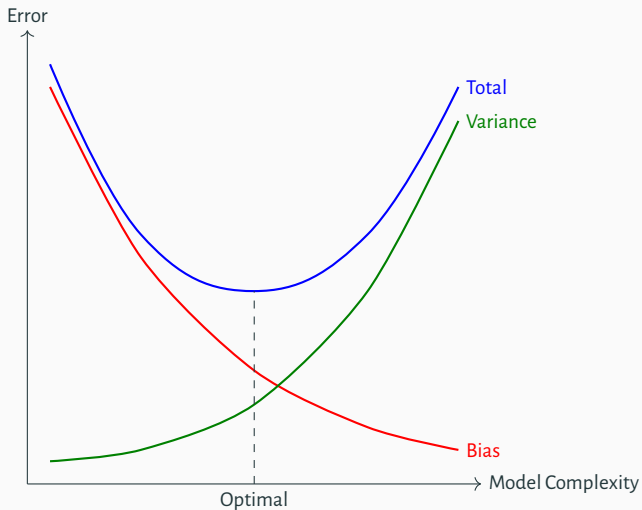
# BIAS-VARIANCE TRADEOFF







## ERROR CURVE VS. MODEL COMPLEXITY



## BIAS-VARIANCE TRADEOFF: UNDERSTANDING & RESOLVING

Underfitting & High Bias	Overfitting & High Variance
Model is too simple to capture data patterns	Model memorizes noise/outliers in training data
<ul style="list-style-type: none"> <li>▶ High training error</li> <li>▶ Poor performance on both seen/unseen data</li> </ul>	<ul style="list-style-type: none"> <li>▶ Low training error but high test error</li> <li>▶ Excellent on training data, poor on new data</li> </ul>
<ul style="list-style-type: none"> <li>▶ Oversimplified model architecture</li> <li>▶ Insufficient features</li> <li>▶ Excessive regularization</li> </ul>	<ul style="list-style-type: none"> <li>▶ Overly complex model</li> <li>▶ Too many features</li> <li>▶ Insufficient training data</li> </ul>

Fix Underfitting (High Bias)	Fix Overfitting (High Variance)
<ul style="list-style-type: none"> <li>▶ Increase model complexity</li> <li>▶ Add relevant features</li> <li>▶ Reduce regularization</li> <li>▶ Train longer with better optimization</li> </ul>	<ul style="list-style-type: none"> <li>▶ Get more training data</li> <li>▶ Use regularization (L1/L2)</li> <li>▶ Implement early stopping</li> <li>▶ Use cross-validation</li> </ul>

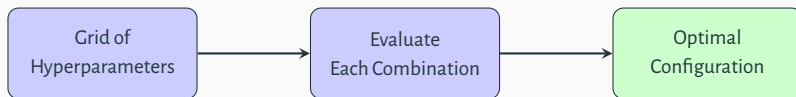
## GRID SEARCH: OPTIMIZING HYPERPARAMETERS

Grid Search is a systematic approach to finding the best hyperparameter combination for a machine learning model by evaluating all possible configurations.

### Key Concepts

- ▶ Defines a **grid** of hyperparameter values.
- ▶ Evaluates all combinations using a performance metric (e.g., accuracy, precision).
- ▶ Identifies the **optimal configuration** for the model.

### Process



### Challenges

**Computationally expensive:** Large grids require significant resources.

**Alternative:** Randomized Search samples combinations to reduce computation.

## CROSS-VALIDATION: ENSURING ROBUSTNESS

Cross-validation is a resampling technique used to evaluate model performance by dividing the dataset into multiple subsets (folds) and validating the model on each subset.

1. Split the dataset into  $k$  subsets (folds).
2. Train the model on  $k - 1$  folds and validate on the remaining fold.
3. Repeat  $k$  times, using a different fold for validation each time.
4. Compute the average performance metric across all folds.

*Validation Fold*

*Training Folds*



EXAMPLE: 5-FOLD CROSS-VALIDATION

### Benefits

- ▶ Reduces overfitting by validating on multiple subsets.
- ▶ Provides a reliable estimate of model performance.

## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → *Codes* → *Julia* → *Part-2*

→ *Pluto* → *hp.jl*

**Pluto.jl** 

→ *Jupyter* → *hp.ipynb*

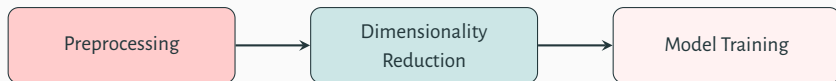


## PIPELINE MECHANISM: AUTOMATING WORKFLOW

A pipeline automates the sequential application of preprocessing steps and model training, ensuring a consistent and efficient workflow.

### Steps in a Pipeline

1. Preprocess data (e.g., scaling, normalization).
2. Apply dimensionality reduction or feature selection.
3. Train a machine learning model.



### Advantages

- ▶ Ensures consistency across training and testing data.
- ▶ Simplifies hyperparameter optimization.

# INTEGRATING PIPELINES, GRID SEARCH, AND CROSS-VALIDATION

## Overview

**Pipeline:** Automates preprocessing and training.

**Grid Search:** Finds the best hyperparameters.

**Cross-Validation:** Ensures robust performance evaluation.



## End-to-End Workflow

1. Define preprocessing steps in a pipeline.
2. Use Grid Search to optimize hyperparameters.
3. Validate results using K-Folds Cross-Validation.



## OVERALL METHODOLOGY (1/2)

### 1. Problem Definition

- Understand the business or research objective.
- Define success metrics (e.g., accuracy, F1-score, ROI).

### 2. Data Collection

- Gather raw data from databases, APIs, or files.
- Ensure data quality, relevance, and representativeness.

### 3. Data Preprocessing & Exploratory Analysis (EDA)

- Handle missing values, duplicates, and outliers.
- Perform statistical and visual analysis.

### 4. Feature Engineering & Selection

- Transform/create meaningful features (e.g., scaling, encoding).
- Select optimal features (e.g., using PCA, feature importance).

### 5. Model Selection & Training

- Split data into training, validation, and test sets.
- Train multiple algorithms (e.g., regression, neural networks).
- Apply cross-validation techniques.

## OVERALL METHODOLOGY (2/2)

### 6. Model Evaluation

- Test on unseen data (validation/test sets).
- Compare performance metrics (e.g., precision, RMSE, AUC-ROC).

### 7. Hyperparameter Tuning & Optimization

- Fine-tune models using GridSearch, RandomSearch, or Bayesian methods.

### 8. Deployment & Monitoring

- Deploy the model (e.g., as an API, cloud service, or embedded system).
- Continuously monitor performance and data drift.

### 9. Communication & Maintenance

- Document methodology, results, and limitations.
- Schedule retraining and updates as needed.

## Unsupervised Learning

---

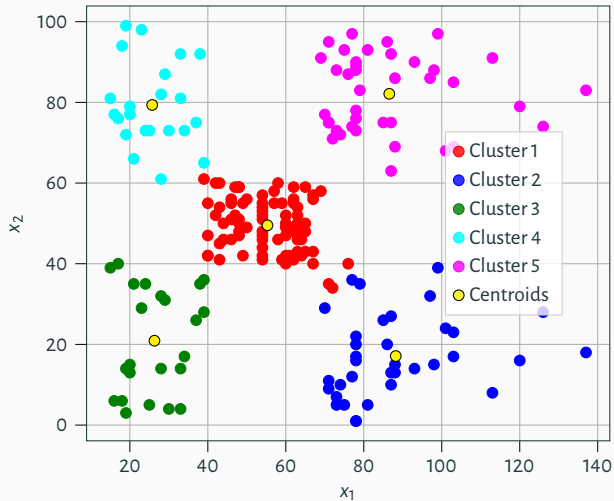
## K-MEANS CLUSTERING (1/3)

The algorithm *K-Means* allows to display regularities or patterns in unlabeled data.

- ▶ The term 'means' refers to averaging the data when computing each centroid;
- ▶ A centroid is the arithmetic mean of all the data points belonging to a particular cluster.

This technique identifies a certain number of centroids within a data set. The algorithm then allocates every data point to the nearest cluster as it attempts to keep the clusters as small as possible. At the same time, *K-Means* attempts to keep the other clusters as different as possible.

## K-MEANS CLUSTERING (2/3)



## K-MEANS CLUSTERING (3/3)

---

### Algorithm 2 Summary Construction

---

1: **procedure** HOW DOES K-MEANS WORK? (Discovering similarities)

**Input:** Unlabeled data sets;

**Output:** Grouping into clusters.

2:     Define how many clusters will be used to group the data sets;

3:     Initialize all the coordinates of the  $k$  cluster centers

4:     **repeat**

5:         Assign each point to its nearest cluster;

6:         Update the centroids coordinates;

7:     **until** No changes to the centers of the clusters

8:     Assign new cases to one of the clusters

9:     **end procedure**

---

**Task #10<sup>3</sup>**

Of the following examples, which would you address using an unsupervised learning algorithms? (*Check all that apply.*)

1. Given email labeled as spam/not spam, learn a spam filter
2. Given a set of news articles found on the web, group them into set of articles about the same story
3. Given a database of customer data, automatically discover market segments and group customers into different market segments
4. Given a dataset of patients diagnosed as either having diabetes or not, learn to classify new patients as having diabetes or not.

---

<sup>3</sup>From 'Machine Learning' course on 'Coursera'

**Task #10<sup>3</sup>**

Of the following examples, which would you address using an unsupervised learning algorithms? (*Check all that apply.*)

1. Given email labeled as spam/not spam, learn a spam filter
2. Given a set of news articles found on the web, group them into set of articles about the same story
3. Given a database of customer data, automatically discover market segments and group customers into different market segments
4. Given a dataset of patients diagnosed as either having diabetes or not, learn to classify new patients as having diabetes or not.

---

<sup>3</sup>From 'Machine Learning' course on 'Coursera'



**Task #11<sup>4</sup>**

Use K-Means algorithm to cluster the following eight points into three clusters:

$A(2, 10)$ ;  $B(2, 5)$ ;  $C(8, 4)$ ;  $D(5, 8)$ ;  $E(7, 5)$ ;  $F(6, 4)$ ;  $G(1, 2)$  and  $H(4, 9)$ .

- Initial cluster centers are:  $\alpha(2, 10)$ ;  $\beta(5, 8)$  and  $\gamma(1, 2)$
- The distance between two points:  $M(x_m, y_m)$  and  $N(x_n, y_n)$  is defined as

$$d(M; N) = |x_m - x_n| + |y_m - y_n|$$

---

<sup>4</sup>Credit: Shokoufeh Mirzaei, PhD

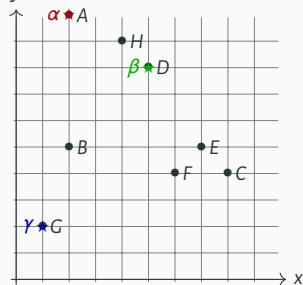
**Task #11<sup>4</sup>**

Use K-Means algorithm to cluster the following eight points into three clusters:

$A(2, 10)$ ;  $B(2, 5)$ ;  $C(8, 4)$ ;  $D(5, 8)$ ;  $E(7, 5)$ ;  $F(6, 4)$ ;  $G(1, 2)$  and  $H(4, 9)$ .

- Initial cluster centers are:  $\alpha(2, 10)$ ;  $\beta(5, 8)$  and  $\gamma(1, 2)$
- The distance between two points:  $M(x_m, y_m)$  and  $N(x_n, y_n)$  is defined as

$$d(M; N) = |x_m - x_n| + |y_m - y_n|$$



<sup>4</sup>Credit: Shokoufeh Mirzaei, PhD

**Task #11<sup>4</sup>**

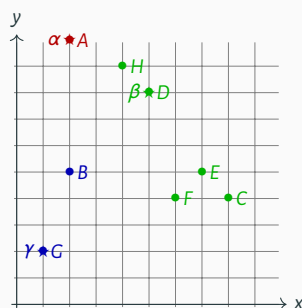
Use K-Means algorithm to cluster the following eight points into three clusters:

A(2, 10); B(2, 5); C(8, 4); D(5, 8); E(7, 5); F(6, 4); G(1, 2) and H(4, 9).

- Initial cluster centers are:  $\alpha(2, 10)$ ;  $\beta(5, 8)$  and  $\gamma(1, 2)$
- The distance between two points:  $M(x_m, y_m)$  and  $N(x_n, y_n)$  is defined as

$$d(M; N) = |x_m - x_n| + |y_m - y_n|$$

Point	$\alpha(2, 10)$	$\beta(5, 8)$	$\gamma(1, 2)$	#
A(2, 10)	0	5	9	1
B(2, 5)	5	6	4	3
C(8, 4)	12	7	9	2
D(5, 8)	5	0	10	2
E(7, 5)	10	5	9	2
F(6, 4)	10	5	7	2
G(1, 2)	9	10	0	3
H(4, 9)	3	2	10	2



<sup>4</sup>Credit: Shokoufeh Mirzaei, PhD

**Task #11<sup>4</sup>**

Use K-Means algorithm to cluster the following eight points into three clusters:

$A(2, 10)$ ;  $B(2, 5)$ ;  $C(8, 4)$ ;  $D(5, 8)$ ;  $E(7, 5)$ ;  $F(6, 4)$ ;  $G(1, 2)$  and  $H(4, 9)$ .

- Initial cluster centers are:  $\alpha(2, 10)$ ;  $\beta(5, 8)$  and  $\gamma(1, 2)$
- The distance between two points:  $M(x_m, y_m)$  and  $N(x_n, y_n)$  is defined as

$$d(M; N) = |x_m - x_n| + |y_m - y_n|$$

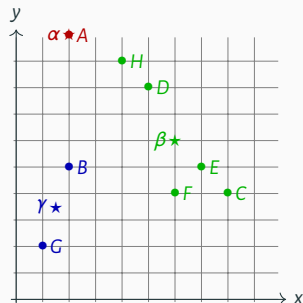
Point	$\alpha(2, 10)$	$\beta(5, 8)$	$\gamma(1, 2)$	#
$A(2, 10)$	0	5	9	1
$B(2, 5)$	5	6	4	3
$C(8, 4)$	12	7	9	2
$D(5, 8)$	5	0	10	2
$E(7, 5)$	10	5	9	2
$F(6, 4)$	10	5	7	2
$G(1, 2)$	9	10	0	3
$H(4, 9)$	3	2	10	2

$\alpha(2, 10)$

$\beta(6, 6)$

$\gamma(1.5, 3.5)$



<sup>4</sup>Credit: Shokoufeh Mirzaei, PhD

**Task #11<sup>4</sup>**

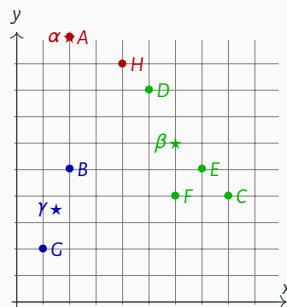
Use K-Means algorithm to cluster the following eight points into three clusters:

A(2, 10); B(2, 5); C(8, 4); D(5, 8); E(7, 5); F(6, 4); G(1, 2) and H(4, 9).

- Initial cluster centers are:  $\alpha(2, 10)$ ;  $\beta(5, 8)$  and  $\gamma(1, 2)$
- The distance between two points:  $M(x_m, y_m)$  and  $N(x_n, y_n)$  is defined as

$$d(M; N) = |x_m - x_n| + |y_m - y_n|$$

Point	$\alpha(2, 10)$	$\beta(5, 8)$	$\gamma(1.5, 3.5)$	#
A(2, 10)	0	8	7	1
B(2, 5)	5	5	2	3
C(8, 4)	12	4	7	2
D(5, 8)	5	3	8	2
E(7, 5)	10	2	7	2
F(6, 4)	10	2	5	2
G(1, 2)	9	9	2	3
H(4, 9)	3	5	8	1



<sup>4</sup>Credit: Shokoufeh Mirzaei, PhD

**Task #11<sup>4</sup>**

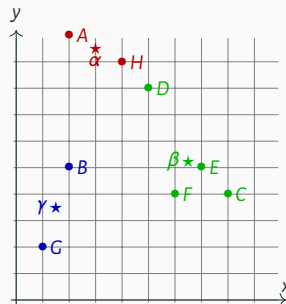
Use K-Means algorithm to cluster the following eight points into three clusters:

$A(2, 10)$ ;  $B(2, 5)$ ;  $C(8, 4)$ ;  $D(5, 8)$ ;  $E(7, 5)$ ;  $F(6, 4)$ ;  $G(1, 2)$  and  $H(4, 9)$ .

- Initial cluster centers are:  $\alpha(2, 10)$ ;  $\beta(5, 8)$  and  $\gamma(1.5, 2)$
- The distance between two points:  $M(x_m, y_m)$  and  $N(x_n, y_n)$  is defined as

$$d(M; N) = |x_m - x_n| + |y_m - y_n|$$

Point	$\alpha(2, 10)$	$\beta(6, 6)$	$\gamma(1.5, 3.5)$	#
$A(2, 10)$	0	8	7	1
$B(2, 5)$	5	5	2	3
$C(8, 4)$	12	4	7	2
$D(5, 8)$	5	3	8	2
$E(7, 5)$	10	2	7	2
$F(6, 4)$	10	2	5	2
$G(1, 2)$	9	9	2	3
$H(4, 9)$	3	5	8	1
<div> <math>\alpha(3, 9.5)</math> <math>\beta(6.5, 5.25)</math> <math>\gamma(1.5, 3.5)</math> </div>				



<sup>4</sup>Credit: Shokoufeh Mirzaei, PhD

**Task #11<sup>4</sup>**

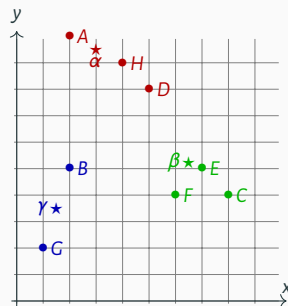
Use K-Means algorithm to cluster the following eight points into three clusters:

$A(2, 10)$ ;  $B(2, 5)$ ;  $C(8, 4)$ ;  $D(5, 8)$ ;  $E(7, 5)$ ;  $F(6, 4)$ ;  $G(1, 2)$  and  $H(4, 9)$ .

- Initial cluster centers are:  $\alpha(2, 10)$ ;  $\beta(5, 8)$  and  $\gamma(1, 2)$
- The distance between two points:  $M(x_m, y_m)$  and  $N(x_n, y_n)$  is defined as

$$d(M; N) = |x_m - x_n| + |y_m - y_n|$$

Point	$\alpha(3, 9.5)$	$\beta(6.5, 5.25)$	$\gamma(1.5, 3.5)$	#
$A(2, 10)$	1.5	9.25	7	1
$B(2, 5)$	5.5	4.75	2	3
$C(8, 4)$	10.5	2.75	7	2
$D(5, 8)$	3.5	4.25	8	1
$E(7, 5)$	8.5	0.75	7	2
$F(6, 4)$	8.5	1.75	5	2
$G(1, 2)$	9.5	8.75	2	3
$H(4, 9)$	1.5	6.25	8	1



<sup>4</sup>Credit: Shokoufeh Mirzaei, PhD

**Task #11<sup>4</sup>**

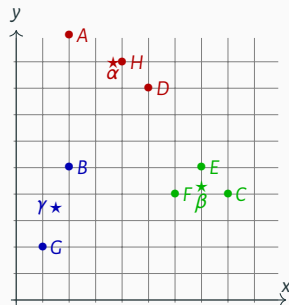
Use K-Means algorithm to cluster the following eight points into three clusters:

$A(2, 10)$ ;  $B(2, 5)$ ;  $C(8, 4)$ ;  $D(5, 8)$ ;  $E(7, 5)$ ;  $F(6, 4)$ ;  $G(1, 2)$  and  $H(4, 9)$ .

- Initial cluster centers are:  $\alpha(2, 10)$ ;  $\beta(5, 8)$  and  $\gamma(1, 2)$
- The distance between two points:  $M(x_m, y_m)$  and  $N(x_n, y_n)$  is defined as

$$d(M; N) = |x_m - x_n| + |y_m - y_n|$$

Point	$\alpha(3, 9.5)$	$\beta(6.5, 5.25)$	$\gamma(1.5, 3.5)$	#
$A(2, 10)$	1.5	9.25	7	1
$B(2, 5)$	5.5	4.75	2	3
$C(8, 4)$	10.5	2.75	7	2
$D(5, 8)$	3.5	4.25	8	1
$E(7, 5)$	8.5	0.75	7	2
$F(6, 4)$	8.5	1.75	5	2
$G(1, 2)$	9.5	8.75	2	3
$H(4, 9)$	1.5	6.25	8	1
<div> <math>\alpha(3.67, 9)</math> <math>\beta(7, 4.3)</math> <math>\gamma(1.5, 3.5)</math> </div>				



<sup>4</sup>Credit: Shokoufeh Mirzaei, PhD



**Task #11<sup>4</sup>**

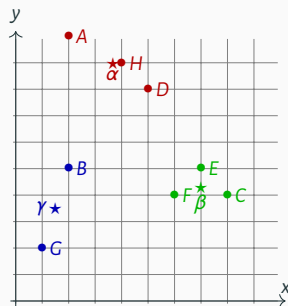
Use K-Means algorithm to cluster the following eight points into three clusters:

$A(2, 10)$ ;  $B(2, 5)$ ;  $C(8, 4)$ ;  $D(5, 8)$ ;  $E(7, 5)$ ;  $F(6, 4)$ ;  $G(1, 2)$  and  $H(4, 9)$ .

- Initial cluster centers are:  $\alpha(2, 10)$ ;  $\beta(5, 8)$  and  $\gamma(1, 2)$
- The distance between two points:  $M(x_m, y_m)$  and  $N(x_n, y_n)$  is defined as

$$d(M; N) = |x_m - x_n| + |y_m - y_n|$$

Point	$\alpha(3.67, 9)$	$\beta(7, 4.3)$	$\gamma(1.5, 3.5)$	#
$A(2, 10)$	2.67	10.7	7	1
$B(2, 5)$	5.67	5.7	2	3
$C(8, 4)$	9.33	1.3	7	2
$D(5, 8)$	2.33	5.7	8	1
$E(7, 5)$	7.33	0.7	7	2
$F(6, 4)$	7.33	1.3	5	2
$G(1, 2)$	9.67	8.3	2	3
$H(4, 9)$	0.33	7.7	8	1



<sup>4</sup>Credit: Shokoufeh Mirzaei, PhD

**Task #11<sup>4</sup>**

Use K-Means algorithm to cluster the following eight points into three clusters:

$A(2, 10)$ ;  $B(2, 5)$ ;  $C(8, 4)$ ;  $D(5, 8)$ ;  $E(7, 5)$ ;  $F(6, 4)$ ;  $G(1, 2)$  and  $H(4, 9)$ .

- Initial cluster centers are:  $\alpha(2, 10)$ ;  $\beta(5, 8)$  and  $\gamma(1, 2)$
- The distance between two points:  $M(x_m, y_m)$  and  $N(x_n, y_n)$  is defined as

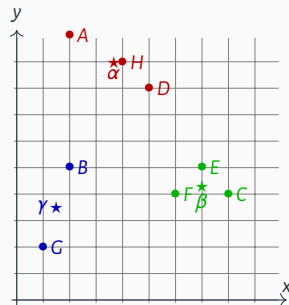
$$d(M; N) = |x_m - x_n| + |y_m - y_n|$$

Point	$\alpha(3.67, 9)$	$\beta(7, 4.3)$	$\gamma(1.5, 3.5)$	#
$A(2, 10)$	2.67	10.7	7	1
$B(2, 5)$	5.67	5.7	2	3
$C(8, 4)$	9.33	1.3	7	2
$D(5, 8)$	2.33	5.7	8	1
$E(7, 5)$	7.33	0.7	7	2
$F(6, 4)$	7.33	1.3	5	2
$G(1, 2)$	9.67	8.3	2	3
$H(4, 9)$	0.33	7.7	8	1

$\alpha(3.67, 9)$

$\beta(7, 4.3)$

$\gamma(1.5, 3.5)$



<sup>4</sup>Credit: Shokoufeh Mirzaei, PhD

## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → *Codes* → *Julia* → *Part-2*

→ *Pluto* → *kmeans.jl*

**Pluto.jl** 

→ *Jupyter* → *kmeans.ipynb*



## DBSCAN: OVERVIEW

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a density-based clustering algorithm that groups together points closely packed in space while marking points in sparse regions as outliers.

**Epsilon ( $\epsilon$ ):** Maximum distance between two points to consider them as neighbors.

**MinPts:** Minimum number of points required to form a dense region.

**Core Point:** A point with at least MinPts neighbors within  $\epsilon$ -distance.

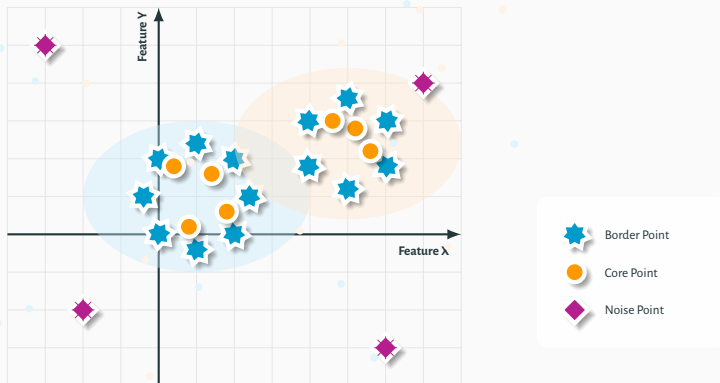
**Border Point:** A point within  $\epsilon$ -distance of a core point but with fewer than MinPts neighbors.

**Noise Point:** A point that is neither a core point nor a border point.

### Steps of DBSCAN

1. Select an unvisited point and check its  $\epsilon$ -neighborhood.
2. Mark it as a core, border, or noise point based on MinPts.
3. Expand clusters iteratively by connecting core points.
4. Continue until all points are visited.

## DBSCAN: VISUAL REPRESENTATION



- ▲ Detects clusters of arbitrary shapes.
- ▲ Robust to noise and outliers.

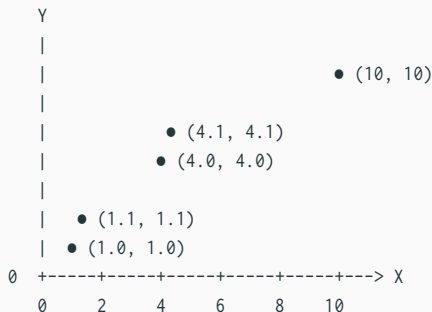
- ▼ Sensitive to the choice of  $\epsilon$  and MinPts.
- ▼ Struggles with clusters of varying densities.
- ▼ Performance degrades with high-dimensional data.

## HANDS-ON EXAMPLE (1/4)

## Task #12

Consider the following synthetic 2D dataset with 5 points:

X	Y
1.0	1.0
1.1	1.1
4.0	4.0
4.1	4.1
10.0	10.0



Apply DBSCAN with  $\epsilon = 1.5$  and  $\text{MinPts} = 2$ .

## HANDS-ON EXAMPLE (2/4)

## Compute Pairwise Distances

	(1, 1)	(1.1, 1.1)	(4, 4)	(4.1, 4.1)	(10, 10)
(1, 1)	0.00	0.141	4.243	4.384	12.728
(1.1, 1.1)	0.141	0.00	4.101	4.243	12.586
(4, 4)	4.243	4.101	0.00	0.141	8.485
(4.1, 4.1)	4.384	4.243	0.141	0.00	8.344
(10, 10)	12.728	12.586	8.485	8.344	0.00

## Identify Core Points

A point is a **core point** if  $\geq \text{MinPts}$  neighbors exist within  $\epsilon$ :

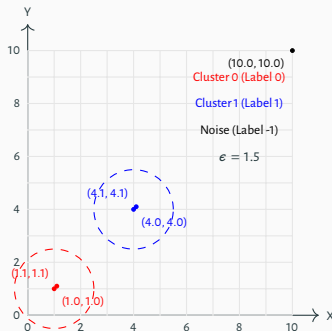
- ▶ **(1.0, 1.0)**: Neighbors = {(1.1, 1.1)}  $\rightarrow 1 (< 2) \rightarrow \text{Border}$
- ▶ **(1.1, 1.1)**: Neighbors = {(1.0, 1.0)}  $\rightarrow 1 (< 2) \rightarrow \text{Border}$
- ▶ **(4.0, 4.0)**: Neighbors = {(4.1, 4.1)}  $\rightarrow 1 (< 2) \rightarrow \text{Border}$
- ▶ **(4.1, 4.1)**: Neighbors = {(4.0, 4.0)}  $\rightarrow 1 (< 2) \rightarrow \text{Border}$
- ▶ **(10.0, 10.0)**: No neighbors within  $\epsilon \rightarrow \text{Noise}$

## HANDS-ON EXAMPLE (3/4)

## Cluster Assignment

- **Cluster 1:**  $\{(1.0, 1.0), (1.1, 1.1)\}$  (mutually reachable via  $\epsilon$ )
- **Cluster 2:**  $\{(4.0, 4.0), (4.1, 4.1)\}$  (mutually reachable via  $\epsilon$ )
- **Noise:**  $(10.0, 10.0)$  (no nearby points)

Point (X, Y)	Label	Type
(1.0, 1.0)	0	Border
(1.1, 1.1)	0	Border
(4.0, 4.0)	1	Border
(4.1, 4.1)	1	Border
(10.0, 10.0)	-1	Noise





## HANDS-ON EXAMPLE (4/4)

```

1  X = [1 1; 1.1 1.1; 4 4; 4.1 4.1; 10 10]; X'
2  #=
3  2×5 adjoint(::Matrix{Float64}) with eltype Float64:
4    1.0  1.1  4.0  4.1  10.0
5    1.0  1.1  4.0  4.1  10.0
6  =#
7
8  f(x) = sqrt.( ( x[1] .- X[:, 1] ).^2 + ( x[2] .- X[:, 2] ).^2 )
9  f(X[1, :])'
10 #=
11 1×5 adjoint(::Vector{Float64}) with eltype Float64:
12 0.0  0.141421  4.24264  4.38406  12.7279
13 =#
14
15 f(X[2, :])
16 f(X[3, :])
17 f(X[4, :])
18 f(X[5, :])

```

## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → *Codes* → *Julia* → *Part-2*

→ *Pluto* → *dbscan.jl*

**Pluto.jl** 

→ *Jupyter* → *dbscan.ipynb*



## PRINCIPAL COMPONENT ANALYSIS (PCA): OVERVIEW

PCA is a dimensionality reduction technique that transforms a high-dimensional dataset into a lower-dimensional space by identifying the directions of maximum variance.

- ▶ Reduce the number of features while retaining most of the dataset's variability.
- ▶ Identify patterns in data by capturing principal components.
- ▶ Remove redundant or irrelevant information.

### How PCA Works

1. Standardize the dataset (mean = 0, variance = 1).
2. Compute the covariance matrix.
3. Calculate eigenvalues and eigenvectors of the covariance matrix.
4. Project the data onto the eigenvectors with the largest eigenvalues (principal components).

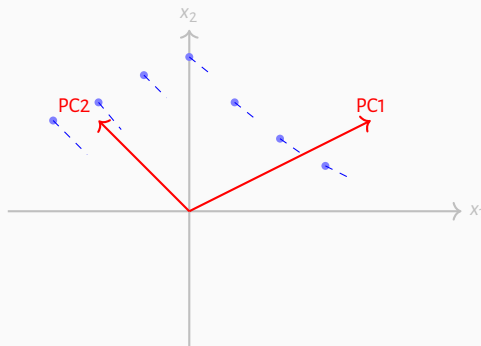
### Applications

- ▶ Image compression.
- ▶ Data visualization in 2D or 3D.
- ▶ Noise reduction.

# VISUALIZING PCA

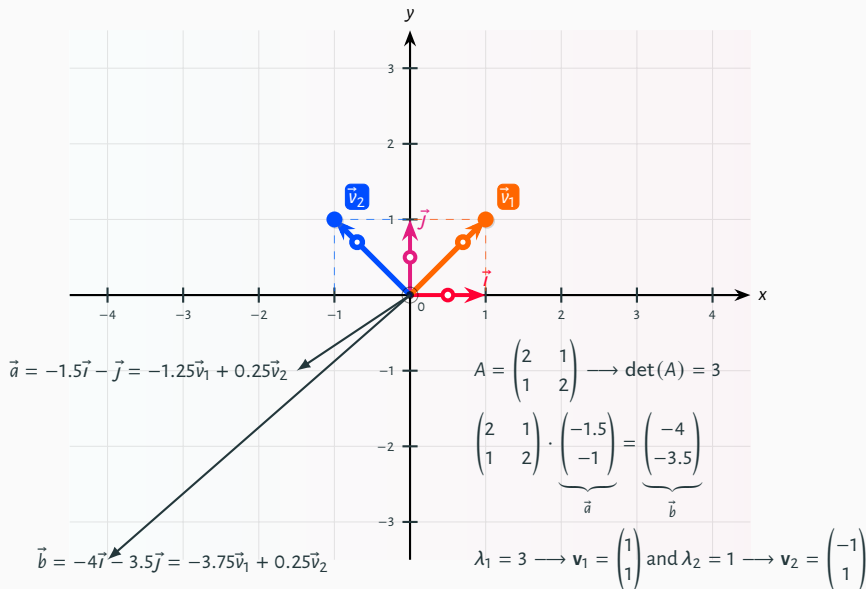
## UNDERSTANDING PCA THROUGH A 2D EXAMPLE:

- ▶ Dataset with two features ( $x_1, x_2$ ).
- ▶ PCA identifies the principal axes of variance (principal components).
- ▶ Data is projected onto these axes.



- ① PC1: Captures the most variance in the data.
- ② PC2: Captures the remaining variance, orthogonal to PC1.

## REMINDER



## HANDS-ON EXAMPLE (1/7)

### STEP-BY-STEP PCA EXAMPLE

$$\mathbf{X} = \begin{bmatrix} 2.5 & 2.4 \\ 0.5 & 0.7 \\ 2.2 & 2.9 \\ 1.9 & 2.2 \\ 3.1 & 3.0 \\ 2.3 & 2.7 \\ 2.0 & 1.6 \\ 1.0 & 1.1 \\ 1.5 & 1.6 \\ 1.1 & 0.9 \end{bmatrix}$$

## HANDS-ON EXAMPLE (2/7)

### STEP-BY-STEP PCA EXAMPLE

#### Mean Normalization

$$\bar{\mathbf{X}} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \end{bmatrix} = \begin{bmatrix} 1.81 \\ 1.91 \end{bmatrix}, \quad \mathbf{X}_{\text{normalized}} = \mathbf{X} - \bar{\mathbf{X}}$$

#### Covariance Matrix

$$\mathbf{C} = \frac{1}{n-1} \mathbf{X}_{\text{normalized}}^T \mathbf{X}_{\text{normalized}} \approx \begin{bmatrix} 0.6165 & 0.6152 \\ 0.6152 & 0.7165 \end{bmatrix}$$

#### Eigenvalues and Eigenvectors

$$\text{Eigenvalues: } \lambda_1 \approx 1.284, \quad \lambda_2 \approx 0.049$$

$$\text{Eigenvectors: } \mathbf{v}_1 \approx \begin{bmatrix} 0.678 \\ 0.735 \end{bmatrix}, \quad \mathbf{v}_2 \approx \begin{bmatrix} -0.735 \\ 0.677 \end{bmatrix}$$

# HANDS-ON EXAMPLE (3/7)

## STEP-BY-STEP PCA EXAMPLE

Project Data onto Principal Components

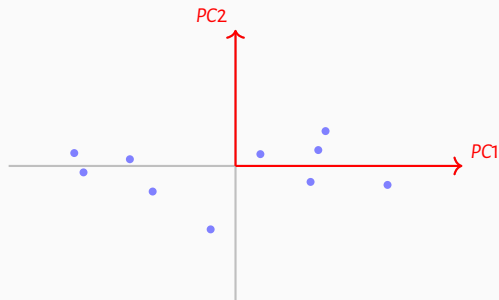
$$\mathbf{Z} = \mathbf{X}_{\text{normalized}} \mathbf{V}, \quad \therefore \mathbf{Z} \approx \begin{bmatrix} 0.827 & -0.175 \\ -1.777 & 0.142 \\ 0.992 & 0.384 \\ 0.274 & 0.130 \\ 1.675 & -0.209 \\ 0.912 & 0.175 \\ -0.099 & -0.35 \\ -1.144 & 0.046 \\ -0.438 & 0.018 \\ -1.224 & -0.163 \end{bmatrix}$$



## HANDS-ON EXAMPLE (4/7)

### STEP-BY-STEP PCA EXAMPLE

#### Visualization of Principal Components



#### Key Observations

- ▶ Most variance is along  $PC1$  (horizontal axis).
- ▶ Data projected onto  $PC2$  (vertical axis) shows minimal variance.

## HANDS-ON EXAMPLE (5/7)

## STEP-BY-STEP PCA EXAMPLE

```

1  using LinearAlgebra
2  X = [
3      2.5 2.4; 0.5 0.7; 2.2 2.9; 1.9 2.2; 3.1 3.0;
4      2.3 2.7; 2.0 1.6; 1 1.1; 1.5 1.6; 1.1 0.9
5  ];
6  X'
7  #=
8  2×10 adjoint(::Matrix{Float64}) with eltype Float64:
9      2.5  0.5  2.2  1.9  3.1  2.3  2.0  1.0  1.5  1.1
10     2.4  0.7  2.9  2.2  3.0  2.7  1.6  1.1  1.6  0.9
11  =#
12
13  X = 1/10 * sum(X, dims=1) #  $\bar{X}$  = ...
14  #=
15  1×2 Matrix{Float64}:
16      1.81  1.91
17  =#
18
19  Xnorm = X .- X ; #  $X - \bar{X}$ 

```

## HANDS-ON EXAMPLE (6/7)

## STEP-BY-STEP PCA EXAMPLE

```

20  Xnorm'
21  #=
22  2×10 adjoint(::Matrix{Float64}) with eltype Float64:
23    0.69  -1.31  0.39  0.09  1.29  0.49  0.19  -0.81  -0.31  -0.71
24    0.49  -1.21  0.99  0.29  1.09  0.79  -0.31  -0.81  -0.31  -1.01
25  =#
26
27  C = 1/(9)*(Xnorm'*Xnorm) # C = 1/(n-1) * ...
28  #=
29  2×2 Matrix{Float64}:
30    0.616556  0.615444
31    0.615444  0.716556
32  =#
33
34  λ, V = eigen(C)
35  #=
36  Eigen{Float64, Float64, Matrix{Float64}, Vector{Float64}}
37  values:
38  2-element Vector{Float64}:

```

# HANDS-ON EXAMPLE (7/7)

## STEP-BY-STEP PCA EXAMPLE

```

39     0.049083398938327194
40     1.2840277121727839
41     vectors:
42     2×2 Matrix{Float64}:
43     -0.735179  0.677873
44     0.677873  0.735179
45     =#
46
47     Z = Xnorm * V;
48     Z'
49     #=
50     2×10 adjoint(::Matrix{Float64}) with eltype Float64:
51     -0.175115  0.142857  0.384375  ...  0.0464173  0.0177646  -0.162675
52     0.82797   -1.77758  0.992197   -1.14457   -0.438046  -1.22382
53     =#

```



## PRACTICAL IMPLEMENTATION



The code is available @ [github.com/a-mhamdi/jlai](https://github.com/a-mhamdi/jlai) → *Codes* → *Julia* → *Part-2*

→ *Pluto* → *pca.jl*

**Pluto.jl** 

→ *Jupyter* → *pca.ipynb*

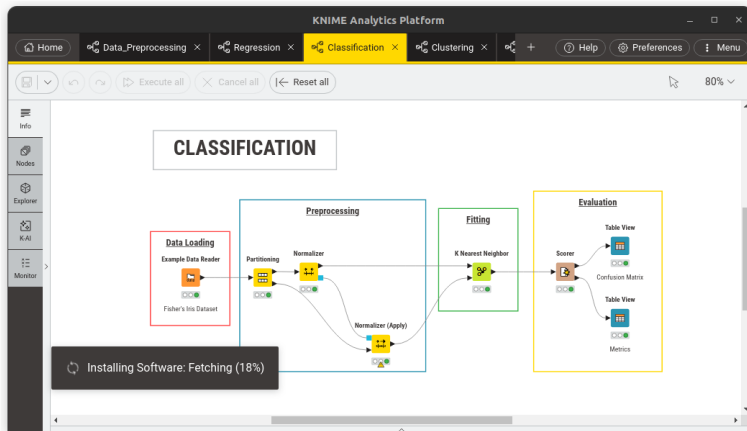


## **Complementary Lab. Project**

---

On the day of assignment, you will be informed about the **dataset to use** (*source: database*), **specific features to retain**, and the **machine learning model to implement**. You will be asked to:

- ① ETL(*Extract, Transform, Load*) data from a database;
- ② Build and evaluate the **ML** model (*pipeline, featurization, split, etc.*).



## **ML Landscape through Quizzes**

---



## KNOWLEDGE CHECK



1

Go to [wooclap.com](https://wooclap.com)

2

Enter the event code in the top banner

Event code  
**JLAI2**<https://app.wooclap.com/JLAI2>

### References

---

- [Bur19] A. Burkov. ***The Hundred-Page Machine Learning Book***. Andriy Burkov, Jan. 1, 2019. 160 pp.
- [Bur20] A. Burkov. ***Machine Learning Engineering***. True Positive Inc., Sept. 8, 2020. 310 pp.
- [DFO20] M. P. Deisenroth, A. A. Faisal, and C. S. Ong. ***Mathematics for Machine Learning***. Cambridge University Pr., Apr. 1, 2020. 398 pp.
- [ENM15] I. El Naqa and M. J. Murphy. “**What Is Machine Learning?**” In: *Machine Learning in Radiation Oncology: Theory and Applications*. Ed. by I. El Naqa, R. Li, and M. J. Murphy. Cham: Springer International Publishing, 2015, pp. 3–11. DOI: 10.1007/978-3-319-18305-3\_1.
- [Fla12] P. Flach. “**References**”. In: *Machine Learning: The Art and Science of Algorithms that Make Sense of Data*. Cambridge University Press, Sept. 2012, pp. 367–382. DOI: 10.1017/CB09780511973000.017.
- [GBC16] I. Goodfellow, J. Bengio, and A. Courville. ***Deep Learning***. MIT Press Ltd, Nov. 18, 2016. 800 pp.
- [Gé19] A. Géron. ***Hands-on Machine Learning with Scikit-Learn, Keras, and TensorFlow***. O'Reilly Media, Oct. 15, 2019. 819 pp.
- [HYU21] T. J. Hui (York University). ***Machine Learning Fundamentals***. Cambridge University Press, Nov. 25, 2021. 420 pp.

## FURTHER READING (2/2)

- [Jia22] H. Jiang. ***Machine Learning Fundamentals***. Cambridge University Pr., Jan. 31, 2022.
- [JPM21] L. M. John Paul Mueller. ***Machine Learning For Dummies***. Wiley John + Sons, Apr. 8, 2021. 464 pp.
- [Mit97] T. Mitchell. ***Machine Learning***. McGraw-Hill International Editions. McGraw-Hill, 1997.
- [Pra18] M. L. de Prado. ***Advances in Financial Machine Learning***. John Wiley & Sons Inc, May 4, 2018. 400 pp.
- [SG16] A. C. M. Sarah Guido. ***Introduction to Machine Learning with Python***. O'Reilly Media, July 31, 2016.
- [Woj12] J. Wojtusiak. **“Machine Learning”**. In: *Encyclopedia of the Sciences of Learning*. Springer US, 2012, pp. 2082–2083. DOI: 10.1007/978-1-4419-1428-6\_1927.