

Demystifying Artificial Intelligence Sorcery

(Part 2: Machine Learning)^a

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

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

^aAvailable @ <https://github.com/a-mhamdi/isetbz/>



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.

1. An overview
2. Supervised Learning
3. Unsupervised Learning
4. 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} .

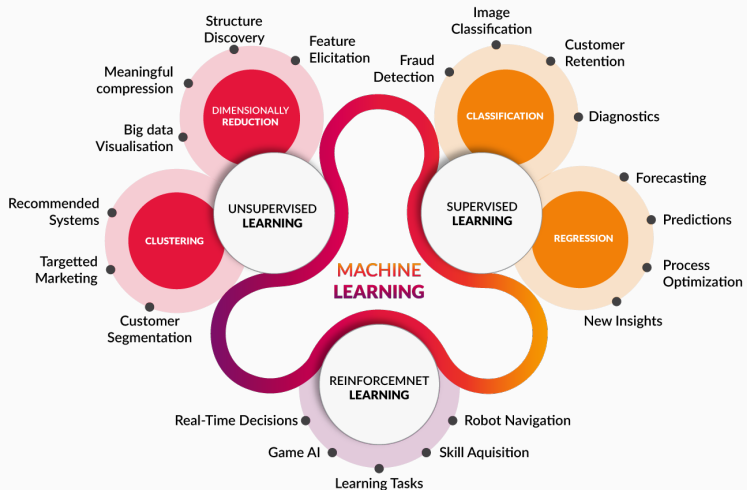
Task #1

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.

OVERALL METHODOLOGY

1. Define the problem;
2. Gather dataset;
3. Choose measure of success;
4. Decide evaluation protocol;
5. Prepare the data;
6. Develop a model;
7. Iterate models.

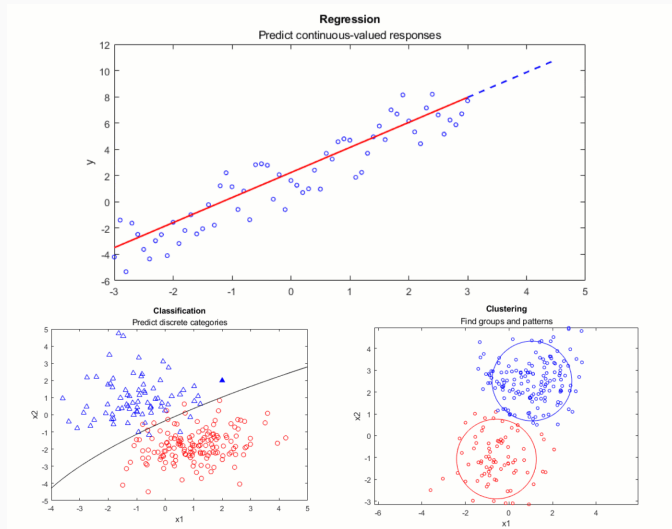


<https://www.cognub.com/index.php/cognitive-platform/>



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

REGRESSION | CLASSIFICATION | CLUSTERING

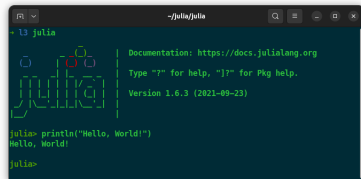


<https://github.com/MathWorks-Teaching-Resources/Machine-Learning-for-Regression>



REMINDER

PROGRAMMING LANGUAGE

julia-lang.org/

DEVELOPMENT ENVIRONMENTS



Pluto.jl



▲ \$ docker compose up

▼ \$ docker compose down



JULIA IN A NUTSHELL

- ▲ Fast
- ▲ Dynamic
- ▲ Reproducible
- ▲ Composable
- ▲ General
- ▲ Open Source



JULIA MICRO-BENCHMARKS (1/2)



<https://julialang.org/benchmarks>



JULIA MICRO-BENCHMARKS (2/2)

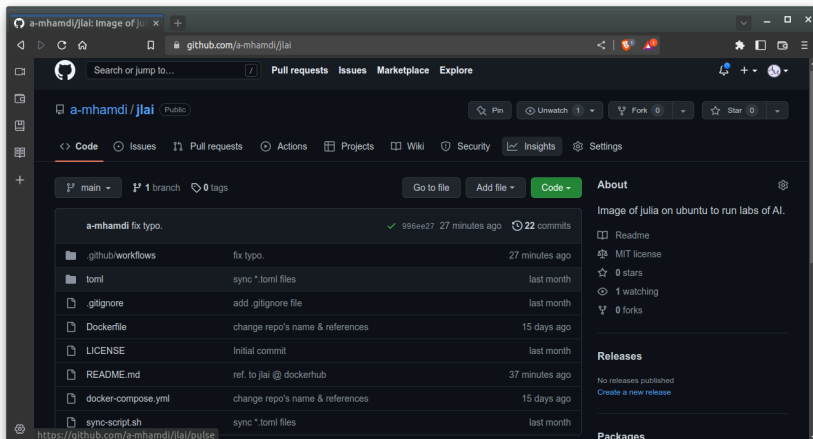
Geometric Means of Micro-Benchmarks by Language

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





SOURCE CONTROL MANAGEMENT (SCM)



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



CONTINUOUS INTEGRATION (CI)

The screenshot shows the Docker Hub interface for the repository `abmhamdi/jlai`. The page includes a search bar, navigation tabs (General, Tags, Builds, Collaborators, Webhooks, Settings), and a description of the repository as 'Artificial Intelligence Labs @ ISETBZ'. It also displays Docker commands for pushing a new tag, a table of tags and scans, and information about automated builds.

abmhamdi /jlai

Description
Artificial Intelligence Labs @ ISETBZ
Last pushed: 2 minutes ago

Docker commands
To push a new tag to this repository,
`docker push abmhamdi/jlai:tagname`

Tags and scans
This repository contains 1 tag(s).
VULNERABILITY SCANNING - DISABLED [Enable](#)

Tag	OS	Type	Pulled	Pushed
latest	linux	Image	—	2 minutes ago

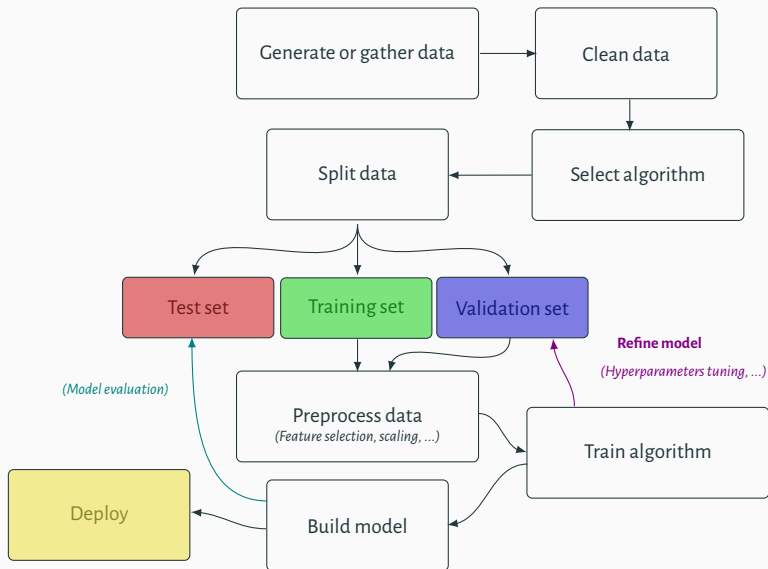
[See all](#) [Go to Advanced Image Management](#)

Automated Builds
Manually pushing images to Hub? Connect your account to GitHub or Bitbucket to automatically build and tag new images whenever your code is updated, so you can focus your time on creating.
Available with Pro, Team and Business subscriptions.
[Upgrade](#) [Learn more](#)

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

Supervised Learning

WORKFLOW IN MACHINE LEARNING



DATA PREPROCESSING

- ▶ 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.)

FEATURE SCALING

Normalisation

$$X = \frac{X - \min(X)}{X.\max() - X.\min()}$$

▲ No assumption on data distribution

Standardisation

$$X = \frac{X - X.\text{mean}()}{X.\text{std}()}$$

▲ More recommended when following normal distribution

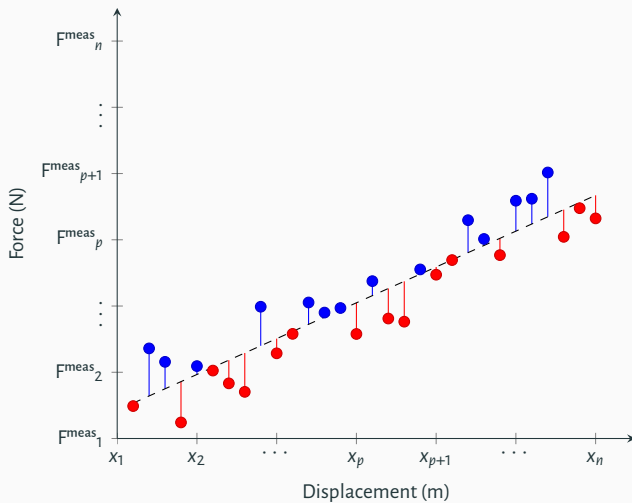
DATA PREPROCESSING TEMPLATE

CODE SNIPPET



The notebook is available at <https://github.com/a-mhamdi/isetbz/>
→ Artificial Intelligence → Codes → Julia → data-preprocessing-template.jl





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^{meas}_i)

x_i	x_1	\dots	x_p	\dots	x_n
F^{meas}_i	F^{meas}_1	\dots	F^{meas}_p	\dots	F^{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:

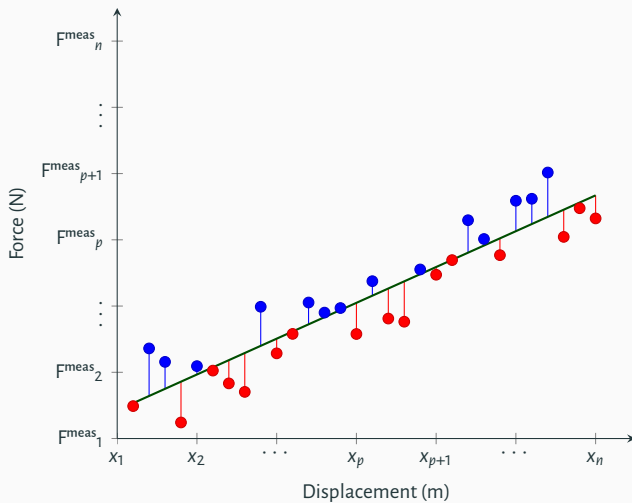
$$\begin{aligned} \mathcal{J} &= \sum_{i=1}^n \varepsilon_i^2 \\ &= \sum_{i=1}^n (F^{\text{meas}}_i - kx_i)^2 \end{aligned}$$

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

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

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

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



SIMPLE LINEAR REGRESSION

CODE SNIPPET



The notebook is available at <https://github.com/a-mhamdi/isetbz/>
→ Artificial Intelligence → Codes → Julia → simple-lr.jl



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_i	t_1	\dots	t_p	\dots	t_n
y^{meas}_i	y^{meas}_1	\dots	y^{meas}_p	\dots	y^{meas}_n

$$\begin{aligned} y^{\text{meas}}_i &= y_i + \varepsilon_i \\ &= y_0 + v_0 t_i + \varepsilon_i, \end{aligned} \quad (5)$$

where y_i denotes the unknown real value of the position y at time point t_i .

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:

$$\begin{aligned} \mathcal{J} &= \sum_{i=1}^n \varepsilon_i^2 \\ &= \varepsilon^T \times \varepsilon \end{aligned}$$

The vector ε is set by $\varepsilon_i, \forall i \geq 1$:

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

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

MULTIPLE LINEAR REGRESSION

CODE SNIPPET



The notebook is available at <https://github.com/a-mhamdi/isetbz/>
→ Artificial Intelligence → Codes → Julia → multiple-lr.jl



Consider the following multivariate 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

$$y_k = \theta_1 x_{(1,k)} + \theta_2 x_{(2,k)} + \cdots + \theta_m x_{(m,k)} + \varepsilon_k \quad (8)$$

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

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

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

$$y = H\theta + \varepsilon,$$

$$\text{where } y^T = [y_1, y_2, \cdots, y_n], X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}, \text{ and } \varepsilon^T = [\varepsilon_1, \varepsilon_2, \cdots, \varepsilon_n].$$

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

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

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 θ :

$$\frac{\partial \mathcal{J}}{\partial \theta} = \frac{\partial (\varepsilon^T \varepsilon)}{\partial \theta} \quad (9)$$

$$\frac{\partial (\varepsilon^T \varepsilon)}{\partial \theta} = 2 \left(\frac{\partial \varepsilon}{\partial \theta} \right)^T \varepsilon \quad (10)$$

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

$$\frac{\partial \varepsilon}{\partial \theta} = -X \quad (11)$$

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

The regressor is given by

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



$X^T X$ is not invertible (singular/degenerate)

▼ Redundant Features

Some features are linearly dependant, i.e, \exists some $x_p \propto$ some x_l for instance 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

GRADIENT DESCENT

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

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

$$\theta_i \triangleq \theta_i + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(x_k)) x_{(i,k)}$$

$$\theta_0 \triangleq \theta_0 + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(x_k)) x_{(0,k)}$$

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

$$\vdots$$

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

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 θ_0 , θ_1 .

1. Using normal equation,
2. Using gradient descent for 5 iterations.

$$y = \begin{bmatrix} 14 \\ 38 \\ 54 \\ 76 \\ 95 \end{bmatrix} \quad \text{and} \quad X = \begin{bmatrix} 1 & 0 \\ 1 & 25 \\ 1 & 50 \\ 1 & 75 \\ 1 & 100 \end{bmatrix} \quad \Rightarrow \quad X^T X = \begin{bmatrix} 5 & 250 \\ 250 & 18750 \end{bmatrix}$$

$$\hat{\theta} = \begin{bmatrix} \hat{\theta}_0 \\ \hat{\theta}_1 \end{bmatrix} = \begin{bmatrix} 15.4 \\ 0.8 \end{bmatrix}$$

POLYNOMIAL LINEAR REGRESSION

CODE SNIPPET



The notebook is available at <https://github.com/a-mhamdi/isetbz/>
→ Artificial Intelligence → Codes → Julia → polynomial-lr.jl



F1-Score, Accuracy, Recall and **Precision** are calculated as follow:

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

$f1 - score$ denotes the *Harmonic Mean of Recall & Precision*

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

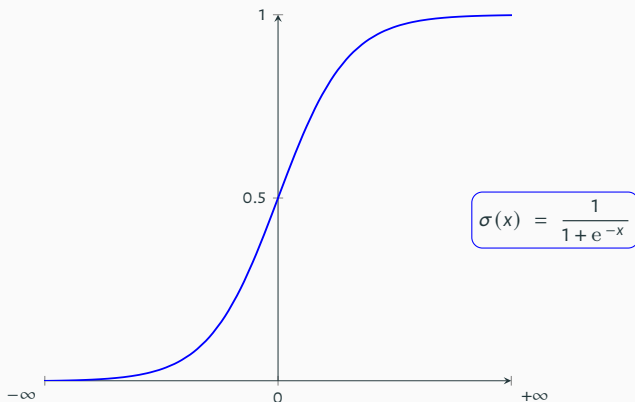
It denotes the ratio of how much we got right over all cases. Recall, on the other hand, designates the ratio of how much positives do we got right over all actual positive cases.

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

Precision, at last, is how much positives we got right over all positive predictions. It is given by:

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

LOGISTIC OR S-SHAPED FUNCTION σ



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

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

DECISION BOUNDARY

$$y = \sigma (\theta_1 x_{(1)} + \theta_2 x_{(2)} + \cdots + \theta_m x_{(m)})$$

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

Hypothesis:

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \quad h_{\theta}(x_k) = \frac{1}{1 + e^{-\theta^T x_k}}$$

Cost function:

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

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

GRADIENT DESCENT

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

Generalizing \mathcal{J} yields: $\mathcal{J} = -\frac{1}{n} \sum_{k=1}^n (y_k \ln(h_{\theta}(x_k)) + (1 - y_k) \ln(1 - h_{\theta}(x_k)))$

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

$$\theta_i \triangleq \theta_i + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(x_k)) x_{(i,k)}$$

$$\theta_0 \triangleq \theta_0 + \alpha \frac{1}{n} \sum_{k=1}^n (y_k - h_{\theta}(x_k)) x_{(0,k)}$$

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

\vdots

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

LOGISTIC REGRESSION

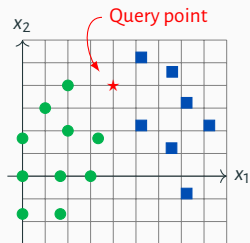
CODE SNIPPET



The notebook is available at <https://github.com/a-mhamdi/isetbz/>
→ Artificial Intelligence → Codes → Julia → logistic-regression.jl



k -NEAREST NEIGHBORS (1/6)



► Evelyn Fix and Joseph Hodges, 1951

► Thomas Cover, 1966

k-NEAREST NEIGHBORS (2/6)

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**

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 (3/6)

Task #3

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 $P(5, 5)$?

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

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

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

```
function dds(a, b) # `a` and `b` are coordinates of some point
    d_squared = (a-5)^2+(b-5)^2
    (d_squared, sqrt(d_squared))
end
```

```
dds(1, 6) # Point `A`
```

```
dds(2, 6) # Point `B`
```

k -NEAREST NEIGHBORS (4/6)

Task #4¹

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.

k-NEAREST NEIGHBORS (5/6)

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

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

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

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

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

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

k-NEAREST NEIGHBORS (6/6)

```
function dds(w, h) # `w` and `h` are width and height of some fruit
    d_squared = (w-1)^2+(h-4)^2
    (d_squared, sqrt(d_squared))
end

dds(2, 6) # Fruit `F_1`
dds(5, 6) # Fruit `F_2`
```

¹From Prof. Winston's book



The notebook is available at <https://github.com/a-mhamdi/isetbz/>
→ Artificial Intelligence → Codes → Julia → knn.jl



RULE OF THUMB TO CHOOSE k

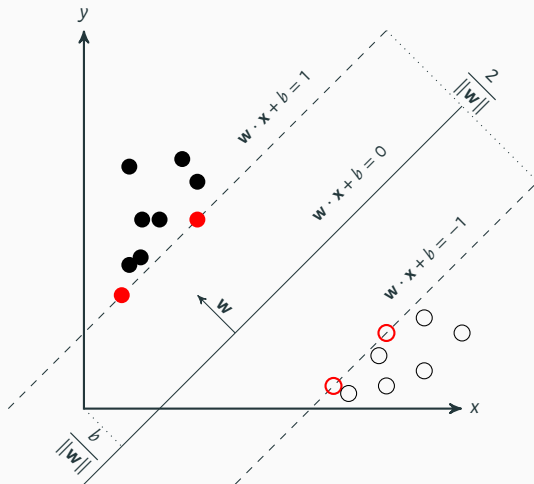
k is **even** if the number of classes is odd

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 MACHINE (SVM)



SVM FOR CLASSIFICATION

CODE SNIPPET



The notebook is available at <https://github.com/a-mhamdi/isetbz/>
→ Artificial Intelligence → Codes → Julia → svc.jl



OUTRODUCTION

Method		Pros		Cons
<i>Logistic Regression</i>	▲	Probabilistic	▼	Almost linearly separable data
<i>k-NN</i>	▲	Simple	▼	Number of neighbors k
	▲	Fast	▼	Detecting outliers ²
	▲	Efficient		

²Points that differ significantly from the rest of the data points.

Unsupervised Learning

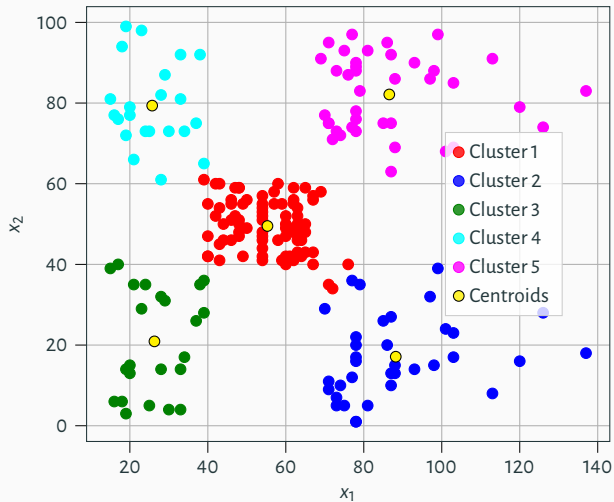
K-MEANS CLUSTERING (1/3)

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 #5³

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.

³From 'Machine Learning' course on 'Coursera'

Task #5³

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.

³From 'Machine Learning' course on 'Coursera'

Task #6⁴

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

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

⁴Credit: Shokoufeh Mirzaei, PhD

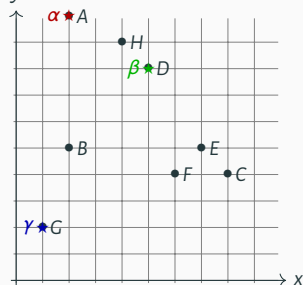
Task #6⁴

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

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



⁴Credit: Shokoufeh Mirzaei, PhD

Task #6⁴

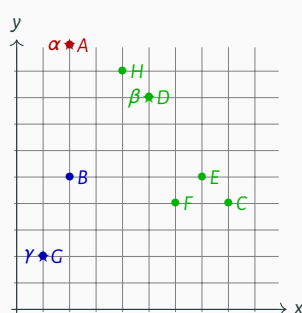
Use K-Means algorithm to cluster the following eight points into three clusters:

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



⁴Credit: Shokoufeh Mirzaei, PhD

Task #6⁴

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

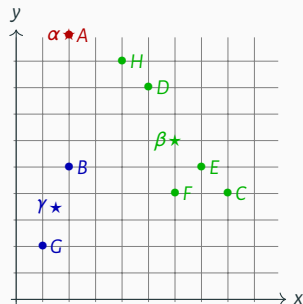
$$d(M; N) = |x_m - x_n| + |y_m - y_n|$$

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

$\alpha(2, 10)$

$\beta(6, 6)$

$\gamma(1.5, 3.5)$



⁴Credit: Shokoufeh Mirzaei, PhD

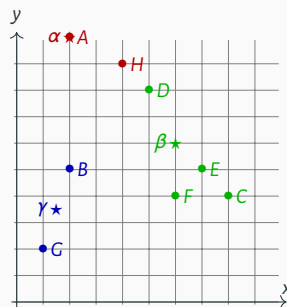
Task #6⁴

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	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

⁴Credit: Shokoufeh Mirzaei, PhD

Task #6⁴

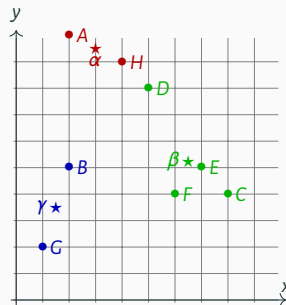
Use K-Means algorithm to cluster the following eight points into three clusters:

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



⁴Credit: Shokoufeh Mirzaei, PhD

Task #6⁴

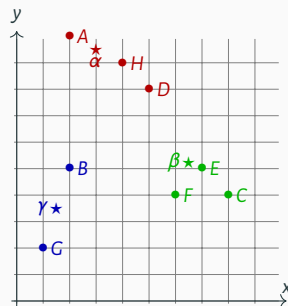
Use K-Means algorithm to cluster the following eight points into three clusters:

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



⁴Credit: Shokoufeh Mirzaei, PhD

Task #6⁴

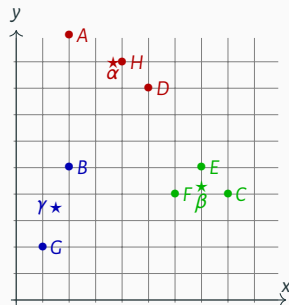
Use K-Means algorithm to cluster the following eight points into three clusters:

$A(2, 10)$; $B(2, 5)$; $C(8, 4)$; $D(5, 8)$; $E(7, 5)$; $F(6, 4)$; $G(1, 2)$ 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> $\alpha(3.67, 9)$ $\beta(7, 4.3)$ $\gamma(1.5, 3.5)$ </div>				



⁴Credit: Shokoufeh Mirzaei, PhD

Task #6⁴

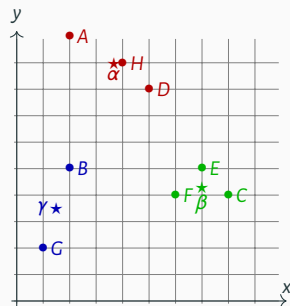
Use K-Means algorithm to cluster the following eight points into three clusters:

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



⁴Credit: Shokoufeh Mirzaei, PhD

Task #6⁴

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

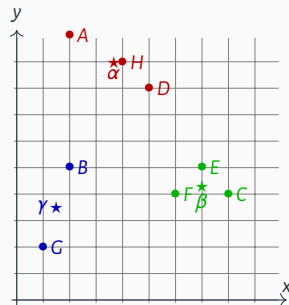
$$d(M; N) = |x_m - x_n| + |y_m - y_n|$$

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

$\alpha(3.67, 9)$

$\beta(7, 4.3)$

$\gamma(1.5, 3.5)$



⁴Credit: Shokoufeh Mirzaei, PhD

K-MEANS

CODE SNIPPET



The notebook is available at <https://github.com/a-mhamdi/isetbz/>
→ Artificial Intelligence → Codes → Julia → kmeans.jl



ML Landscape through Quizzes

MCQ (1/10)

1. ... is the machine learning algorithm that can be used with labeled data.

- ✓ Regression algorithm
- ✓ Clustering algorithm
- ✓ Association algorithm

2. What is Machine Learning (ML)?

- × The selective acquisition of knowledge through the use of computer programs
- × The selective acquisition of knowledge through the use of manual programs
- ✓ The autonomous acquisition of knowledge through the use of computer programs
- × The autonomous acquisition of knowledge through the use of manual programs

3. Successful applications of ML

- × Learning to recognize spoken words
- × Learning to drive an autonomous vehicle
- × Learning to classify new astronomical structures
- × Learning to play world-class backgammon
- ✓ All of the above

MCQ (2/10)

4. Features of Machine Learning are ...

- ✓ Automation
- ✓ Improved customer experience
- ✓ Business intelligence

5. Replace missing values with mean/median/mode helps to handle missing or corrupted data in a dataset. True/False?

- ✓ True
- × False

6. Which among the following algorithms are used in Machine learning?

- ✓ Naive Bayes
- ✓ Support Vector Machines
- ✓ k -Nearest Neighbors

MCQ (3/10)

7. Overfitting is a type of modelling error which results in the failure to predict future observations effectively or fit additional data in the existing model. Yes/No?
- ☐ Probably
 - ☒ Yes
 - ☐ No
 - ☐ Can not say
8. ... is the scenario when the model fails to decipher the underlying trend in the input data.
- ☒ Underfitting
 - ☐ Overfitting
 - ☐ All of the above
 - ☐ None of the above
9. Machine learning approaches can be traditionally categorized into ... categories.
- ☒ 3
 - ☐ 4
 - ☐ 7
 - ☐ 9

MCQ (4/10)

10. The categories in which Machine learning approaches can be traditionally categorized are ...

- × Supervised learning
- × Unsupervised learning
- × Reinforcement learning
- ✓ All of the above

11. In general, to have a well-defined learning problem, we must identify which of the following

- × The class of tasks
- × The measure of performance to be improved
- × The source of experience
- ✓ All of the above

12. The average positive difference between computed and desired outcome values

- × Root Mean Squared Error
- × Mean Squared Error
- × Mean Absolute Error
- ✓ Mean Positive Error

MCQ (5/10)

13. ... is used as an input to the machine learning model for training and prediction purposes.
- × Target variable
 - ✓ Feature vector
 - × All of the above
 - × None of the above
14. Simple regression assumes a ... relationship between the input attribute and output attribute.
- ✓ linear
 - × quadratic
 - × reciprocal
 - × inverse
15. The correlation between the number of years an employee has worked for a company and the salary of the employee is 0.75. What can be said about employee salary and years worked?
- × There is no relationship between salary and years worked.
 - ✓ Individuals that have worked for the company the longest have higher salaries.
 - × Individuals that have worked for the company the longest have lower salaries.
 - × The majority of employees have been with the company a long time.

MCQ (6/10)

16. Which machine learning models are trained to make a series of decisions based on the rewards and feedback they receive for their actions?
- × Supervised learning
 - × Unsupervised learning
 - ✓ Reinforcement learning
 - × All of the above
17. Which of the following is not a type of supervised learning?
- × Classification
 - × Regression
 - ✓ Clustering
 - × None of the above
18. As the amount of training data increases
- × Training error usually increases and generalization error usually increases
 - ✓ Training error usually increases and generalization error usually decreases
 - × Training error usually decreases and generalization error usually decreases
 - × Training error usually decreases and generalization error usually increases

MCQ (7/10)

19. Which of the following are not classification tasks?

- × Find the gender of a person by analyzing his writing style
- × Detect Pneumonia from Chest X-ray images
- ✓ Predict the price of a house based on floor area, number of rooms, etc.
- × Predict whether there will be abnormally heavy rainfall next year

20. Which of the following is a categorical feature?

- × Height of a person
- × Price of petroleum
- × Amount of rainfall in a day
- ✓ Mother tongue of a person

21. What is the use of validation dataset in Machine Learning?

- × To train the machine learning model.
- ✓ To tune the hyperparameters of the machine learning model
- × To evaluate the performance of the machine learning model
- × None of the above

MCQ (8/10)

22. When there is noise in data, which of the following options would improve the performance of the k -NN algorithm?
- ✓ Increase the value of k
 - × Decrease the value of k
 - × Changing value of k will not change the effect of the noise
 - × None of these
23. Which of the following criteria is typically used for optimizing in linear regression.
- × Maximize the number of points it touches.
 - × Minimize the number of points it touches.
 - ✓ Minimize the squared distance from the points.
 - × Minimize the maximum distance of a point from a line.
24. Logistic Regression is used for ...
- × regression purposes
 - ✓ classification purposes
 - × all of the above
 - × none of the above

MCQ (9/10)

25. The supervised learning problems can be grouped as ...

- ☐ Regression problems
- ☐ Classification problems
- ☒ All of the above
- ☐ None of the above

26. The unsupervised learning problems can be grouped as ...

- ☐ Clustering
- ☐ Association
- ☒ All of the above
- ☐ None of the above

27. Which of the following methods do we use to best fit the data in Logistic Regression?

- ☐ Least Squared Error
- ☒ Maximum Likelihood
- ☐ Jaccard distance

MCQ (10/10)

28. The term machine learning was coined by ...

- ☐ James Gosling
- ☒ Arthur Samuel
- ☐ Guido van Rossum
- ☐ None of the above

29. For two runs of K-Means clustering, is it expected to get same clustering results?

- ☐ Yes
- ☒ No

30. Which of the following can act as possible termination conditions in K-Means?

- a. For a fixed number of iterations
 - b. Assignment of observations to clusters does not change between iterations. Except for cases with a bad local minimum.
 - c. Centroids do not change between successive iterations
 - d. Terminate when RSS falls below a threshold
- ☐ a, c & d
 - ☐ a, b & c
 - ☐ a, b & d
 - ☒ All of the above

SOME USEFUL LINKS

1. <https://setosa.io/ev/>
2. <https://karpathy.ai/>
3. <http://yann.lecun.com/>
4. <https://www.hackingnote.com/>
5. <https://machinelearningmastery.com/>
6. <https://stanford.edu/~shervine/teaching/>
7. <https://www.ibm.com/downloads/cas/GB8ZMQZ3>
8. <https://colah.github.io/posts/2014-03-NN-Manifolds-Topology/>

FURTHER READING (1/2)

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.

FURTHER READING (2/2)

- [HYU21] T. J. Hui (York University). *Machine Learning Fundamentals*. Cambridge University Press, Nov. 25, 2021. 420 pp.
- [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.