

Classification: Part 1

Statistical Analysis and Document Mining
Spring 2019

Vasilii Feofanov
Université Grenoble Alpes
vasilii.feofanov@univ-grenoble-alpes.fr

1 Classification: First Sight

- 1.1 Introduction
- 1.2 k Nearest Neighbours Algorithm

2 Probabilistic Classification

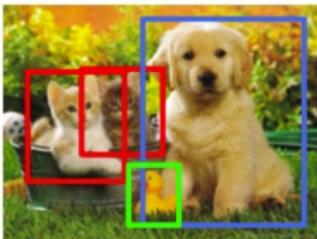
- 2.1 Problem Statement
- 2.2 Bayes Classifier
- 2.3 Linear Discriminant Analysis

3 Classification and Artificial Intelligence

Introduction

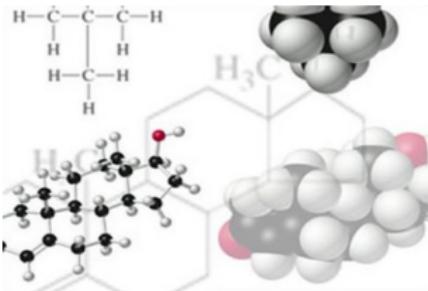


CAT



CAT, DOG, DUCK

(a) Computer Vision



(b) Bioinformatics



(c) Recommender Systems

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

(d) Pattern Recognition



(a) Iris Setosa



(b) Iris Versicolor



(c) Iris Virginica

Figure: The data set consists of 50 samples from each of 3 species of Iris. From each sample the length and the width of the sepals and petals were measured. Based on this, the goal is to build a model that distinguishes the species from each other.¹

¹https://en.wikipedia.org/wiki/Iris_flower_data_set

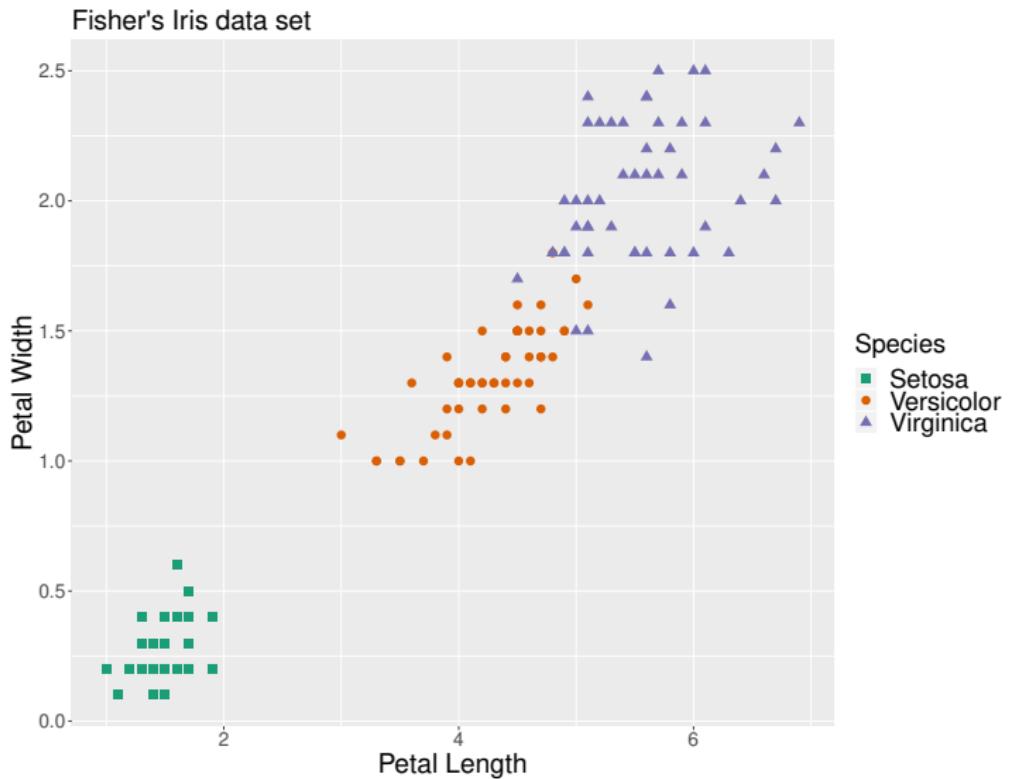
- Training set $\{\mathbf{x}_i, y_i\}_{i=1}^n$: collected samples with class labels;
- A new data point \mathbf{x} without any label;
- Build a classifier $h(\mathbf{x})$ that predicts y as accurately as possible:

| | Sepal.L | Sepal.W | Petal.L | Petal.W |
|----------------|---------|---------|---------|---------|
| $\mathbf{x} =$ | (6.2, | 2.8, | 5.6, | 2.4) |

$\xrightarrow{?}$ $y = \text{"Virginica"}$

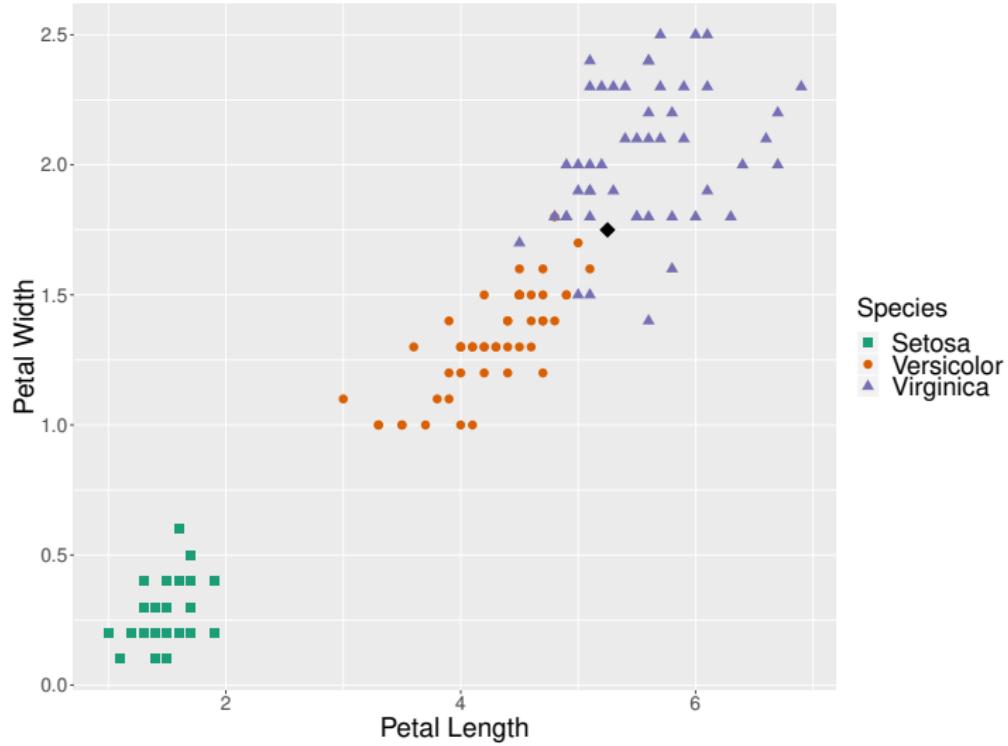
How to do that?

Iris Classification Based on the Petal Information



Iris: New Data Point

Fisher's Iris data set

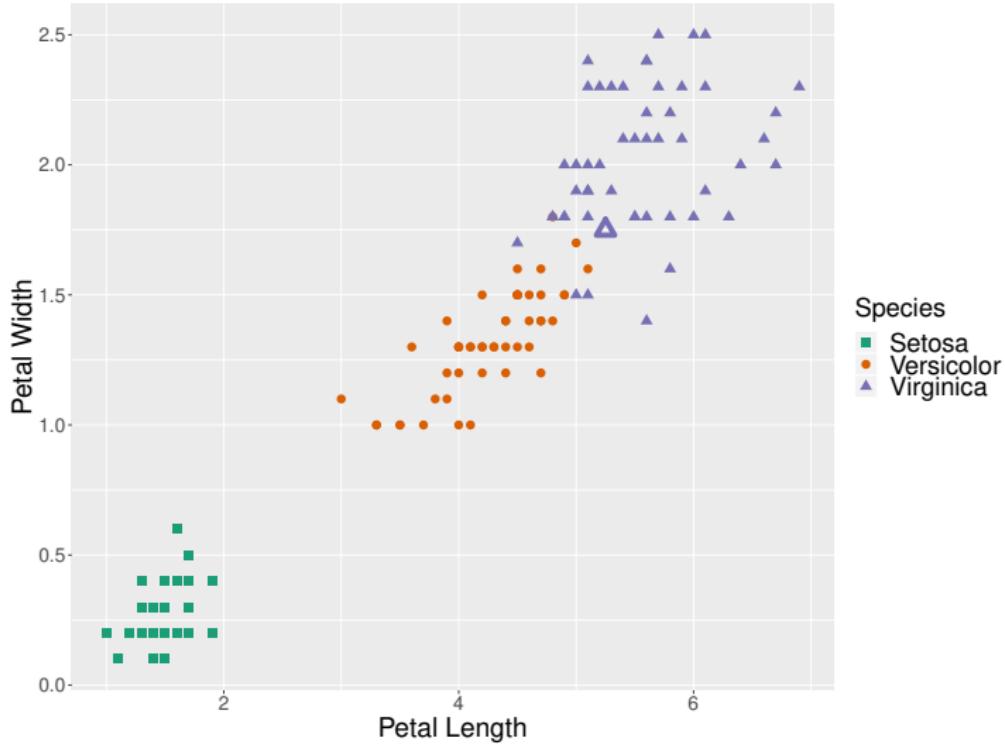


4 Nearest Neighbours



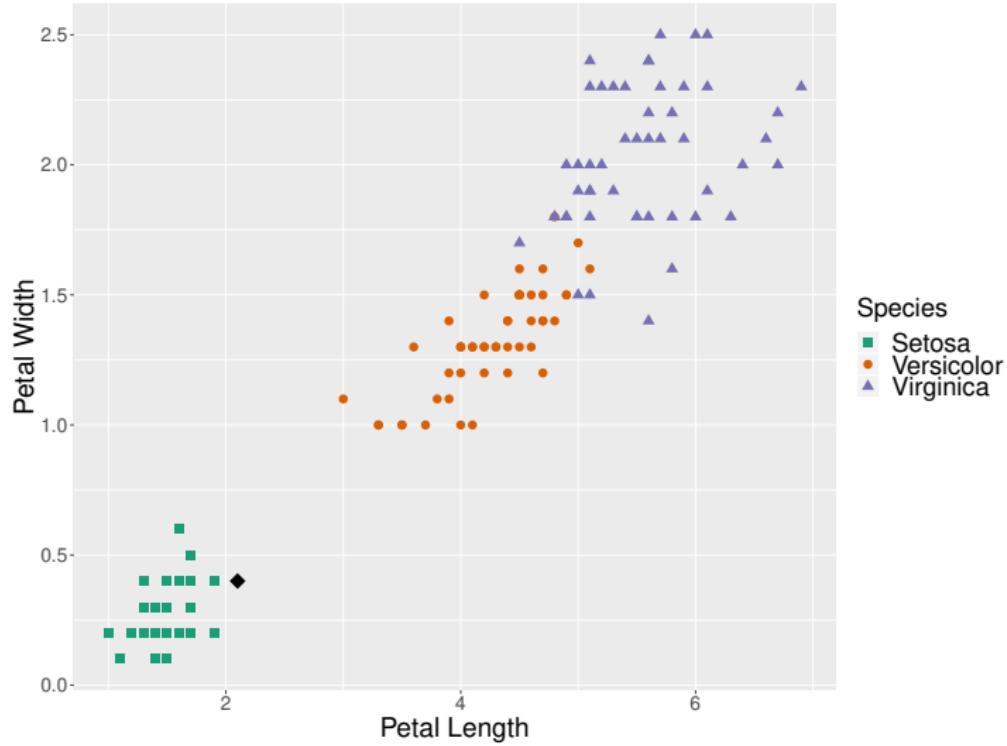
4NN Clasification

Fisher's Iris data set



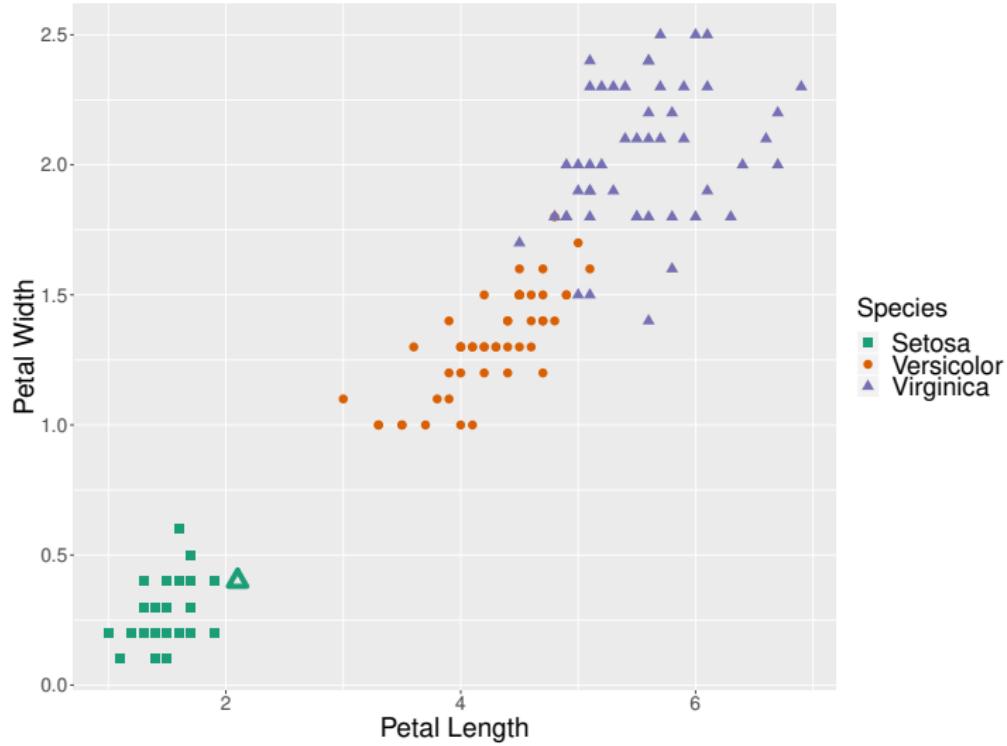
Another Data Point

Fisher's Iris data set

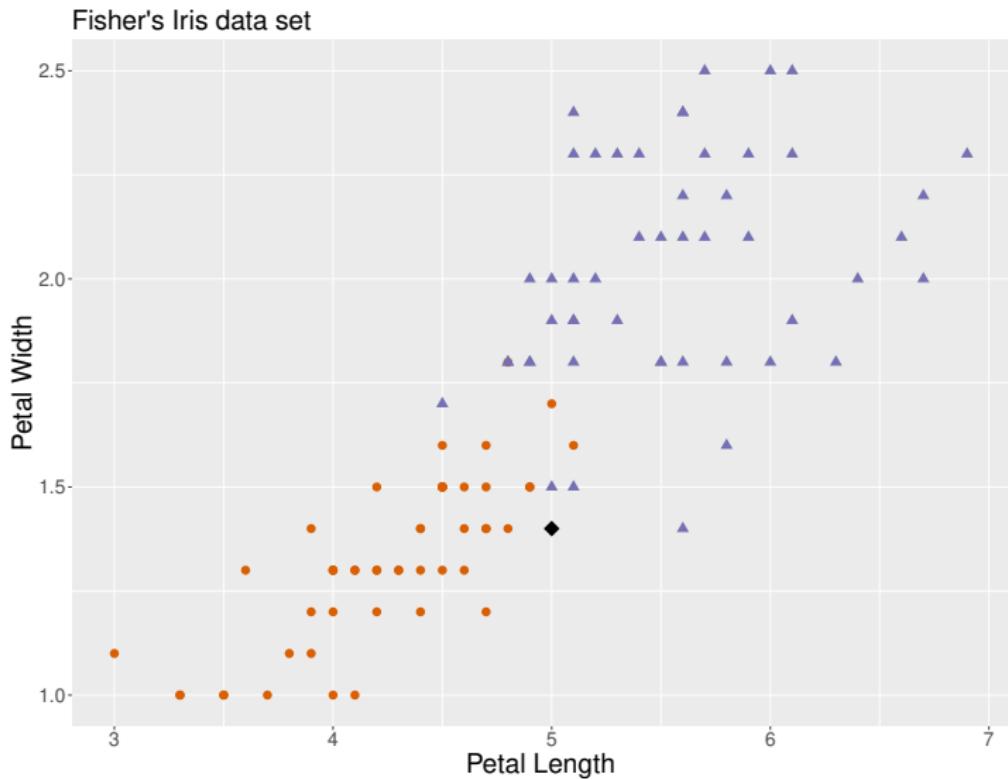


And Its Prediction

Fisher's Iris data set



One More Example



What to do in this case?



Algorithm k Nearest Neighbours (kNN)

Input: Training set $\{\mathbf{x}_i, y_i\}_{i=1}^n$;

Number of classes K ;

New data point \mathbf{x} .

1. Compute distance $d(\mathbf{x}, \mathbf{x}_i)$ for $i = \{1, \dots, n\}$.
2. Find k closest training examples to \mathbf{x} :

$$J \subset \{1, \dots, n\} \quad \text{s.t.} \quad |J| = k;$$

$$\forall j \in J, \forall t \in \{1, \dots, n\} \setminus J : d(\mathbf{x}, \mathbf{x}_j) \leq d(\mathbf{x}, \mathbf{x}_t).$$

Output: Majority class $h(\mathbf{x}) = \operatorname{argmax}_{c=1, \dots, K} \sum_{j \in J} \mathbb{I}(y_j = c)$.

Algorithm k Nearest Neighbours (kNN)

Input: Training set $\{\mathbf{x}_i, y_i\}_{i=1}^n$;

Number of classes K ;

New data point \mathbf{x} .

1. Compute distance $d(\mathbf{x}, \mathbf{x}_i)$ for $i = \{1, \dots, n\}$.
2. Find k closest training examples to \mathbf{x} :

$$J \subset \{1, \dots, n\} \quad \text{s.t.} \quad |J| = k;$$

$$\forall j \in J, \forall t \in \{1, \dots, n\} \setminus J : d(\mathbf{x}, \mathbf{x}_j) \leq d(\mathbf{x}, \mathbf{x}_t).$$

Output: Majority class $h(\mathbf{x}) = \operatorname{argmax}_{c=1, \dots, K} \sum_{j \in J} \mathbb{I}(y_j = c)$.

- What is the time complexity of the algorithm?

Algorithm k Nearest Neighbours (kNN)

Input: Training set $\{\mathbf{x}_i, y_i\}_{i=1}^n$;

Number of classes K ;

New data point \mathbf{x} .

1. Compute distance $d(\mathbf{x}, \mathbf{x}_i)$ for $i = \{1, \dots, n\}$.
2. Find k closest training examples to \mathbf{x} :

$$J \subset \{1, \dots, n\} \quad \text{s.t.} \quad |J| = k;$$

$$\forall j \in J, \forall t \in \{1, \dots, n\} \setminus J : d(\mathbf{x}, \mathbf{x}_j) \leq d(\mathbf{x}, \mathbf{x}_t).$$

Output: Majority class $h(\mathbf{x}) = \operatorname{argmax}_{c=1, \dots, K} \sum_{j \in J} \mathbb{I}(y_j = c)$.

- How does the choice of k affect kNN?

Algorithm k Nearest Neighbours (kNN)

Input: Training set $\{\mathbf{x}_i, y_i\}_{i=1}^n$;

Number of classes K ;

New data point \mathbf{x} .

1. Compute distance $d(\mathbf{x}, \mathbf{x}_i)$ for $i = \{1, \dots, n\}$.
2. Find k closest training examples to \mathbf{x} :

$$J \subset \{1, \dots, n\} \quad \text{s.t.} \quad |J| = k;$$

$$\forall j \in J, \forall t \in \{1, \dots, n\} \setminus J : d(\mathbf{x}, \mathbf{x}_j) \leq d(\mathbf{x}, \mathbf{x}_t).$$

Output: Majority class $h(\mathbf{x}) = \operatorname{argmax}_{c=1, \dots, K} \sum_{j \in J} \mathbb{I}(y_j = c)$.

- How does the choice of distance metric affect the algorithm?

Algorithm k Nearest Neighbours (kNN)

Input: Training set $\{\mathbf{x}_i, y_i\}_{i=1}^n$;

Number of classes K ;

New data point \mathbf{x} .

1. Compute distance $d(\mathbf{x}, \mathbf{x}_i)$ for $i = \{1, \dots, n\}$.
2. Find k closest training examples to \mathbf{x} :

$$J \subset \{1, \dots, n\} \quad \text{s.t.} \quad |J| = k;$$

$$\forall j \in J, \forall t \in \{1, \dots, n\} \setminus J : d(\mathbf{x}, \mathbf{x}_j) \leq d(\mathbf{x}, \mathbf{x}_t).$$

Output: Majority class $h(\mathbf{x}) = \operatorname{argmax}_{c=1, \dots, K} \sum_{j \in J} \mathbb{I}(y_j = c)$.

- Is the nearest neighbours approach applicable for the regression task?

Algorithm k Nearest Neighbours (kNN)

Input: Training set $\{\mathbf{x}_i, y_i\}_{i=1}^n$;

Number of classes K ;

New data point \mathbf{x} .

1. Compute distance $d(\mathbf{x}, \mathbf{x}_i)$ for $i = \{1, \dots, n\}$.
2. Find k closest training examples to \mathbf{x} :

$$J \subset \{1, \dots, n\} \quad \text{s.t.} \quad |J| = k;$$

$$\forall j \in J, \forall t \in \{1, \dots, n\} \setminus J : d(\mathbf{x}, \mathbf{x}_j) \leq d(\mathbf{x}, \mathbf{x}_t).$$

Output: Majority class $h(\mathbf{x}) = \operatorname{argmax}_{c=1, \dots, K} \sum_{j \in J} \mathbb{I}(y_j = c)$.

- kNN is an *instance-based* learning algorithm. What is the main drawback of such methods?

1 Classification: First Sight

- 1.1 Introduction
- 1.2 k Nearest Neighbours Algorithm

2 Probabilistic Classification

- 2.1 Problem Statement
- 2.2 Bayes Classifier
- 2.3 Linear Discriminant Analysis

3 Classification and Artificial Intelligence

Discriminant Analysis (Fisher, 1936)

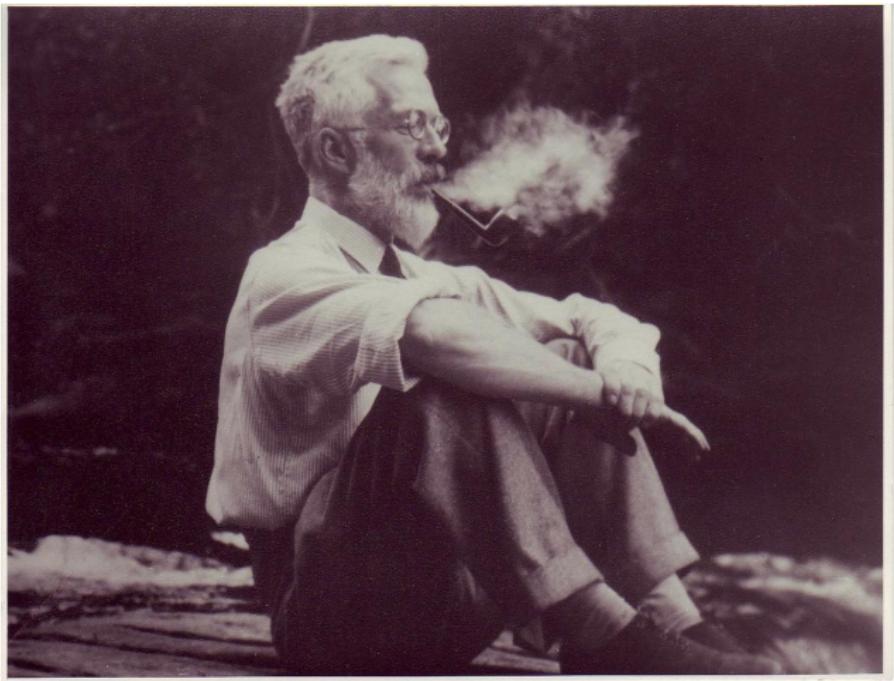


Figure: Ronald A. Fisher in 1946.

- *Input space:* $\mathcal{X} \subseteq \mathbb{R}^d$;
- *Output space:* $\mathcal{Y} = \{-1, +1\}$ (binary classification),
 $\mathcal{Y} = \{1, \dots, K\}$ (multi-class classification);
- *Assumption:* all $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$ are **i.i.d.** from \mathcal{D} with respect to a fixed unknown probability distribution $P(\mathbf{X}, Y)$;
- *Sample Data:* we observe $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$;

- *Input space:* $\mathcal{X} \subseteq \mathbb{R}^d$;
- *Output space:* $\mathcal{Y} = \{-1, +1\}$ (binary classification),
 $\mathcal{Y} = \{1, \dots, K\}$ (multi-class classification);
- *Assumption:* all $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$ are **i.i.d.** from \mathcal{D} with respect to a fixed unknown probability distribution $P(\mathbf{X}, Y)$;
- *Sample Data:* we observe $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$;
- *Loss Function:* $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$;
- *Target:* minimise the *risk* $R^\ell(h) := \mathbb{E}_{(\mathbf{X}, Y) \sim \mathcal{D}} \ell(h(\mathbf{X}), Y)$.

0/1 loss function has the following view:

$$\ell^{0/1}(h(\mathbf{x}), y) = \mathbb{I}(h(\mathbf{x}) \neq y) = \begin{cases} 1, & \text{if } h(\mathbf{x}) \neq y; \\ 0, & \text{if } h(\mathbf{x}) = y \end{cases}$$

0/1 loss function has the following view:

$$\ell^{0/1}(h(\mathbf{x}), y) = \mathbb{I}(h(\mathbf{x}) \neq y) = \begin{cases} 1, & \text{if } h(\mathbf{x}) \neq y; \\ 0, & \text{if } h(\mathbf{x}) = y \end{cases}$$

In this case, the risk is written as:

$$\begin{aligned} R(h) &= P(h(\mathbf{X}) \neq Y) \\ &= \sum_{c \in \{1, \dots, K\}} P(Y = c) P(h(\mathbf{X}) \neq c | Y = c) \end{aligned}$$

0/1 loss function has the following view:

$$\ell^{0/1}(h(\mathbf{x}), y) = \mathbb{I}(h(\mathbf{x}) \neq y) = \begin{cases} 1, & \text{if } h(\mathbf{x}) \neq y; \\ 0, & \text{if } h(\mathbf{x}) = y \end{cases}$$

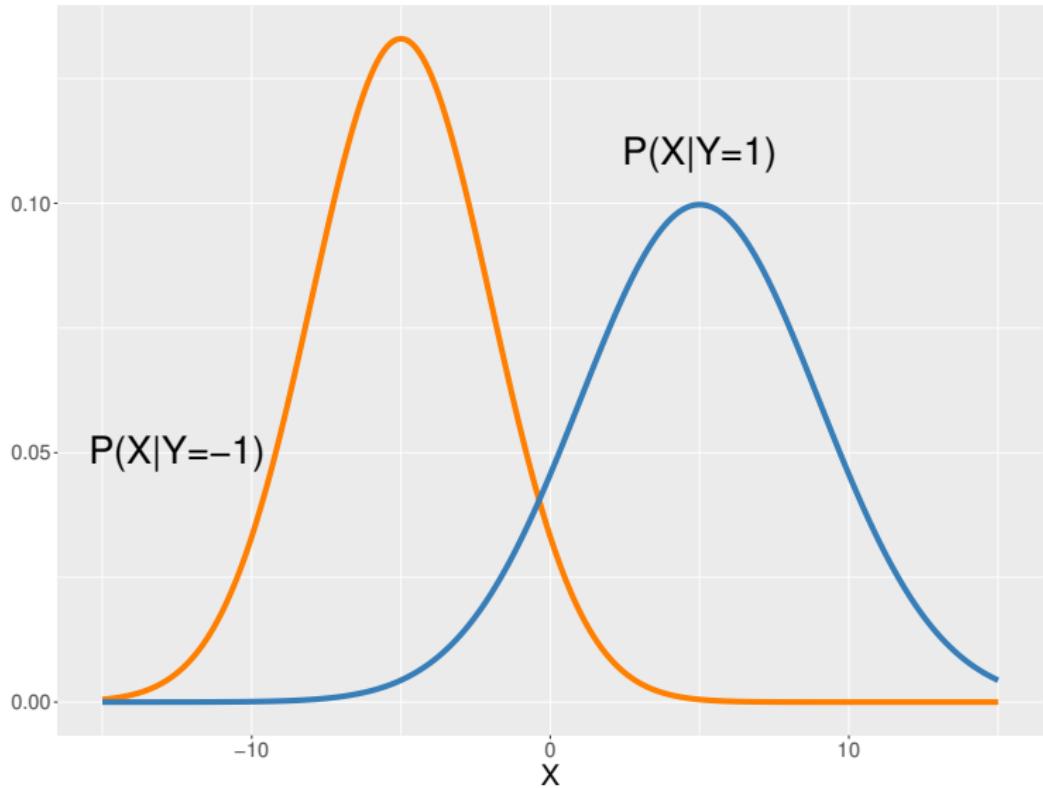
In this case, the risk is written as:

$$\begin{aligned} R(h) &= P(h(\mathbf{X}) \neq Y) \\ &= \sum_{c \in \{1, \dots, K\}} P(Y = c) P(h(\mathbf{X}) \neq c | Y = c) \end{aligned}$$

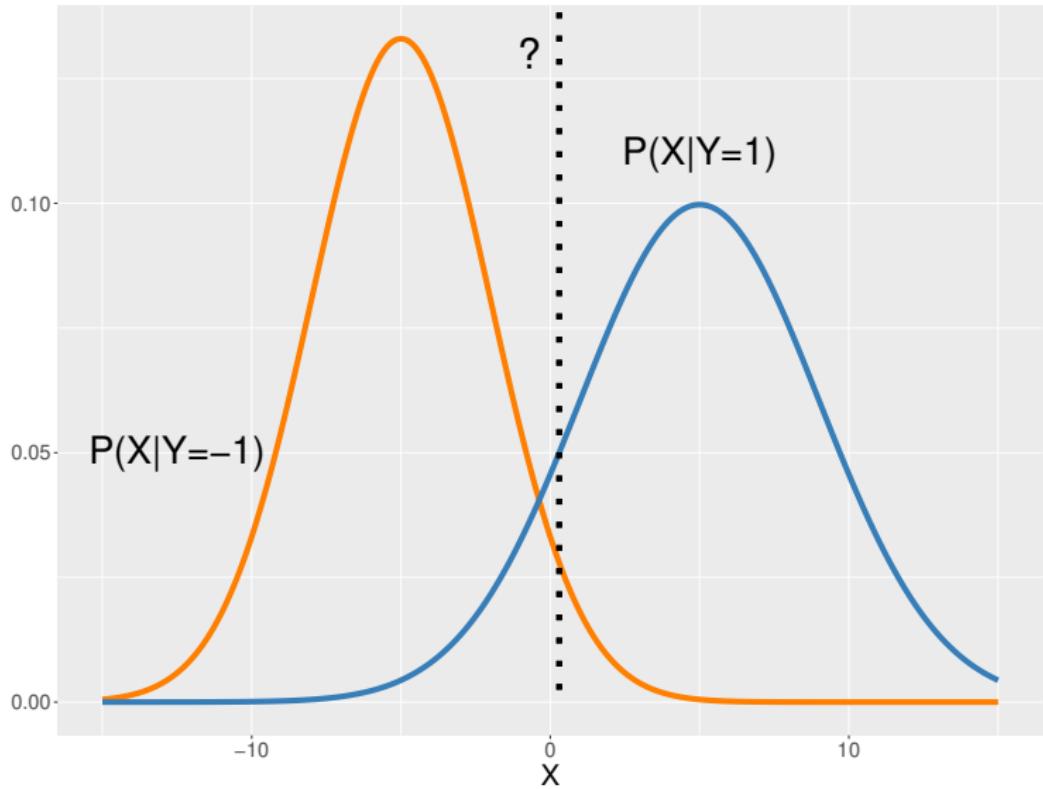
In the binary case, it is represented as:

$$R(h) = P(Y = -1)P(h(\mathbf{X}) = 1 | Y = -1) + P(Y = 1)P(h(\mathbf{X}) = -1 | Y = 1)$$

Classification in Terms of Distributions



Classification in Terms of Distributions



$$P(Y|\mathbf{X}) = \frac{\text{Posterior}}{\text{Evidence}} \cdot \frac{\text{Likelihood}}{\text{Class Prior}}$$
$$P(Y|\mathbf{X}) = \frac{P(\mathbf{X}|Y)P(Y)}{P(\mathbf{X})}$$

Idea: Classify \mathbf{x} to a class with the highest posterior probability:

$$h_B(\mathbf{x}) := \operatorname{argmax}_{y \in \mathcal{Y}} P(Y = y | \mathbf{X} = \mathbf{x}).$$

This is equivalent to:

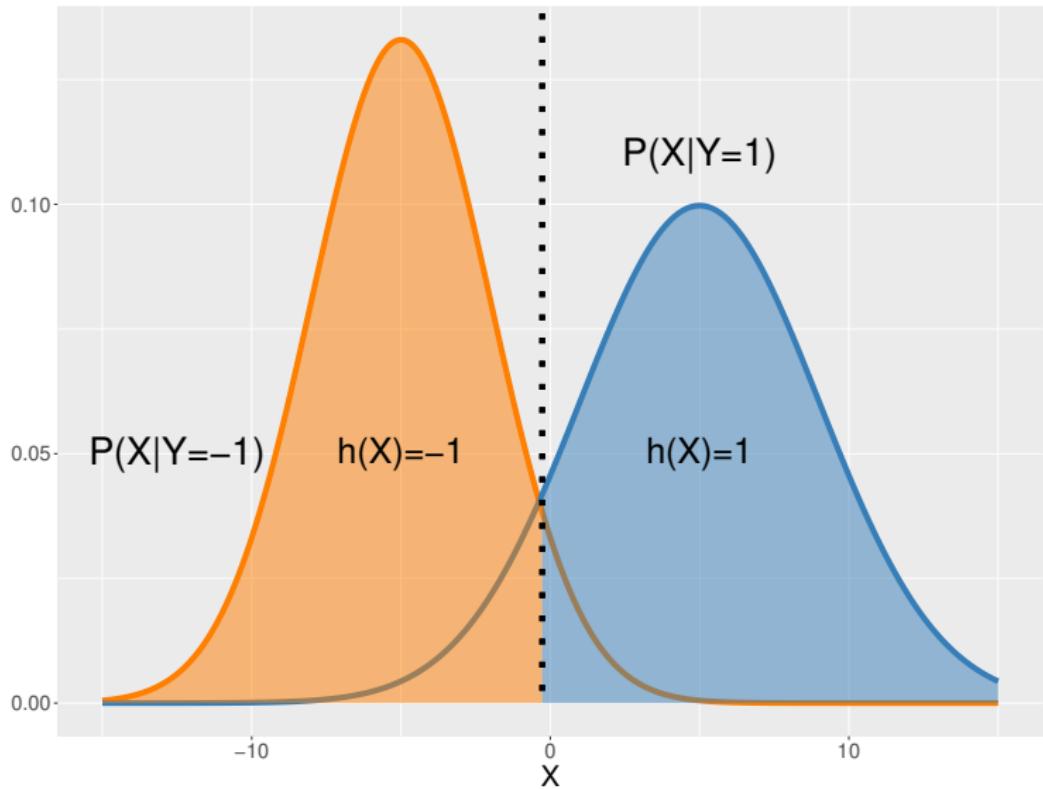
$$h_B(\mathbf{x}) \propto \operatorname{argmax}_{y \in \mathcal{Y}} P(\mathbf{X} = \mathbf{x} | Y = y) P(Y = y).$$

The method is often called the *Bayes* classifier.

In the binary case ($\mathcal{Y} = \{-1, +1\}$), the Bayes classifier is usually defined in the following way:

$$h_B(\mathbf{x}) = \begin{cases} +1, & \text{if } P(Y = +1 | \mathbf{X} = \mathbf{x}) \geq P(Y = -1 | \mathbf{X} = \mathbf{x}); \\ -1, & \text{if } P(Y = +1 | \mathbf{X} = \mathbf{x}) < P(Y = -1 | \mathbf{X} = \mathbf{x}). \end{cases}$$

Particular Case: Maximum Likelihood Estimator



Theorem

Suppose $P(Y)$ and $P(\mathbf{X}|Y)$ are given. Then the Bayes classifier yields the minimum of the misclassification error.

Exercise: Prove this theorem. For sake of simplicity, consider the binary case ($\mathcal{Y} = \{-1, +1\}$).

Theorem

Suppose $P(Y)$ and $P(\mathbf{X}|Y)$ are given. Then the Bayes classifier yields the minimum of the misclassification error.

Exercise: Prove this theorem. For sake of simplicity, consider the binary case ($\mathcal{Y} = \{-1, +1\}$).

- 1 $P(h(\mathbf{X}) \neq Y) = \int P(h(\mathbf{X}) \neq Y | \mathbf{X} = \mathbf{x}) d\mathbf{x}$. Then, what is the value of $P(h_B(\mathbf{X}) \neq Y | \mathbf{X} = \mathbf{x})$?

Theorem

Suppose $P(Y)$ and $P(\mathbf{X}|Y)$ are given. Then the Bayes classifier yields the minimum of the misclassification error.

Exercise: Prove this theorem. For sake of simplicity, consider the binary case ($\mathcal{Y} = \{-1, +1\}$).

- 1 $P(h(\mathbf{X}) \neq Y) = \int P(h(\mathbf{X}) \neq Y | \mathbf{X} = \mathbf{x}) d\mathbf{x}$. Then, what is the value of $P(h_B(\mathbf{X}) \neq Y | \mathbf{X} = \mathbf{x})$?
- 2 $P(h_B(\mathbf{x}) \neq Y | \mathbf{X} = \mathbf{x}) = 1 - \max_{c \in \mathcal{Y}} P(Y = c | \mathbf{X} = \mathbf{x})$. Does this prove the theorem?

In practice, we don't have any information about data distribution.

- Can we estimate $P(Y)$?

In practice, we don't have any information about data distribution.

- Can we estimate $P(Y)$?
- Can we estimate $P(\mathbf{X}|Y)$?

- *Assumption 1:* Observations from a class $c \in \mathcal{Y}$ are normally distributed $[\mathbf{X}|Y = c] \sim \mathcal{N}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$.
 - Remember the formula of the multivariate normal distribution.

- *Assumption 1:* Observations from a class $c \in \mathcal{Y}$ are normally distributed $[\mathbf{X}|Y = c] \sim \mathcal{N}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$.
 - Remember the formula of the multivariate normal distribution.

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_c)}.$$

- *Assumption 1:* Observations from a class $c \in \mathcal{Y}$ are normally distributed $[\mathbf{X}|Y = c] \sim \mathcal{N}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$.
 - Remember the formula of the multivariate normal distribution.

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_c)}.$$

- *Assumption 2:* The covariance matrices of the classes are equal:

$$\boldsymbol{\Sigma}_1 = \cdots = \boldsymbol{\Sigma}_K = \boldsymbol{\Sigma}.$$

Due to the logarithm properties, the Bayes classifier can be written as follows:

$$h_B(\mathbf{x}) = \operatorname{argmax}_{c \in \mathcal{Y}} [\ln P(\mathbf{X} = \mathbf{x} | Y = c) + \ln P(Y = c)].$$

Due to the logarithm properties, the Bayes classifier can be written as follows:

$$h_B(\mathbf{x}) = \operatorname{argmax}_{c \in \mathcal{Y}} [\ln P(\mathbf{X} = \mathbf{x} | Y = c) + \ln P(Y = c)].$$

Taking into account Assumption 1 and Assumption 2, we derive the algorithm called *Linear Discriminant Analysis (LDA)*:

$$h_{LDA}(\mathbf{x}) = \operatorname{argmax}_{c \in \mathcal{Y}} \delta_c(\mathbf{x}),$$

$$\delta_c(\mathbf{x}) = \boldsymbol{\mu}_c^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \ln P(Y = c);$$

δ_c is usually called the *discriminant function*.

In the binary case, the decision rule can be also written as:

$$h_{LDA}(\mathbf{x}) = \begin{cases} +1, & \text{if } \textit{Condition} \text{ is true;} \\ -1, & \text{otherwise.} \end{cases}$$

Condition:

$$(\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} (\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_{+1} + \boldsymbol{\mu}_{-1}) - \ln \frac{P(Y = -1)}{P(Y = +1)} \geq 0.$$

In the binary case, the decision rule can be also written as:

$$h_{LDA}(\mathbf{x}) = \begin{cases} +1, & \text{if } \textit{Condition} \text{ is true;} \\ -1, & \text{otherwise.} \end{cases}$$

Condition:

$$(\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} (\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_{+1} + \boldsymbol{\mu}_{-1}) - \ln \frac{P(Y = -1)}{P(Y = +1)} \geq 0.$$

- Why it is called "linear" ?

In the binary case, the decision rule can be also written as:

$$h_{LDA}(\mathbf{x}) = \begin{cases} +1, & \text{if } \textit{Condition} \text{ is true;} \\ -1, & \text{otherwise.} \end{cases}$$

Condition:

$$\underbrace{(\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1}}_{\mathbf{a}^\top} \mathbf{x} - \underbrace{\frac{1}{2} (\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_{+1} + \boldsymbol{\mu}_{-1}) - \ln \frac{P(Y = -1)}{P(Y = +1)}}_{\mathbf{b}} \geq 0.$$

- Why it is called "linear"?
- How many parameters we need to estimate?

In the binary case, the decision rule can be also written as:

$$h_{LDA}(\mathbf{x}) = \begin{cases} +1, & \text{if } \textit{Condition} \text{ is true;} \\ -1, & \text{otherwise.} \end{cases}$$

Condition:

$$\overbrace{(\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1} \mathbf{x}}^{\mathbf{a}^\top} - \overbrace{\frac{1}{2} (\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_{+1} + \boldsymbol{\mu}_{-1}) - \ln \frac{P(Y = -1)}{P(Y = +1)}}^{\mathbf{b}} \geq 0.$$

- Why it is called "linear"?
- How many parameters we need to estimate?
- What happens when Assumption 1, 2 are violated?

In the binary case, the decision rule can be also written as:

$$h_{LDA}(\mathbf{x}) = \begin{cases} +1, & \text{if } \textit{Condition} \text{ is true;} \\ -1, & \text{otherwise.} \end{cases}$$

Condition:

$$\underbrace{(\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1}}_{\mathbf{a}^\top} \mathbf{x} - \underbrace{\frac{1}{2} (\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_{+1} + \boldsymbol{\mu}_{-1}) - \ln \frac{P(Y = -1)}{P(Y = +1)}}_b \geq 0.$$

- Why it is called "linear"?
- How many parameters we need to estimate?
- What happens when Assumption 1, 2 are violated?
- What is the time complexity during training phase?

In the binary case, the decision rule can be also written as:

$$h_{LDA}(\mathbf{x}) = \begin{cases} +1, & \text{if } \textit{Condition} \text{ is true;} \\ -1, & \text{otherwise.} \end{cases}$$

Condition:

$$\underbrace{(\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1}}_{\mathbf{a}^\top} \mathbf{x} - \underbrace{\frac{1}{2} (\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_{+1} + \boldsymbol{\mu}_{-1}) - \ln \frac{P(Y = -1)}{P(Y = +1)}}_b \geq 0.$$

- Why it is called "linear"?
- How many parameters we need to estimate?
- What happens when Assumption 1, 2 are violated?
- What is the time complexity during training phase?
- What is the time complexity to predict a label for new \mathbf{x} ?

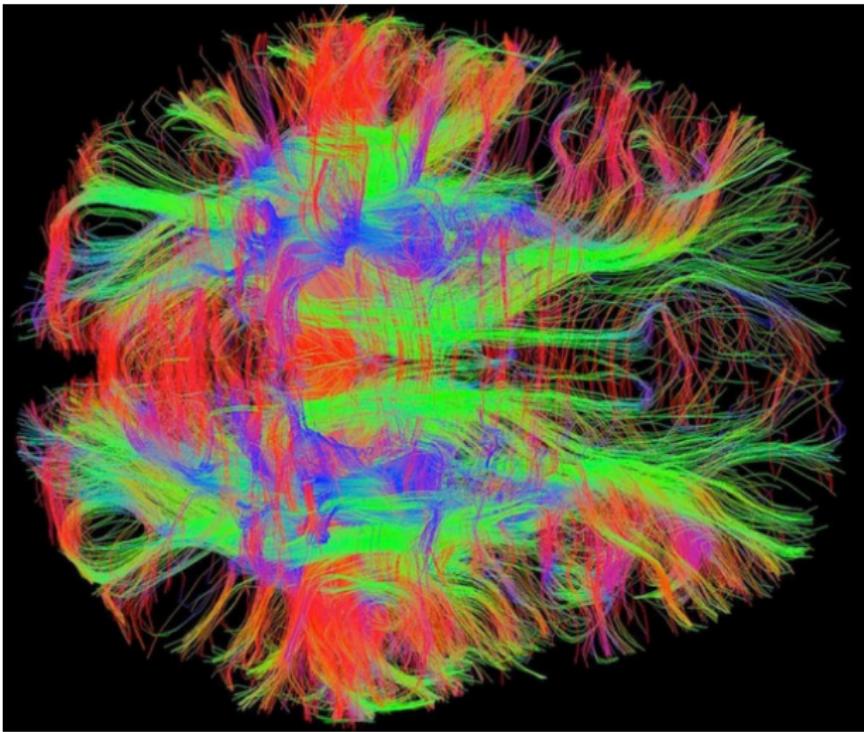
1 Classification: First Sight

- 1.1 Introduction
- 1.2 k Nearest Neighbours Algorithm

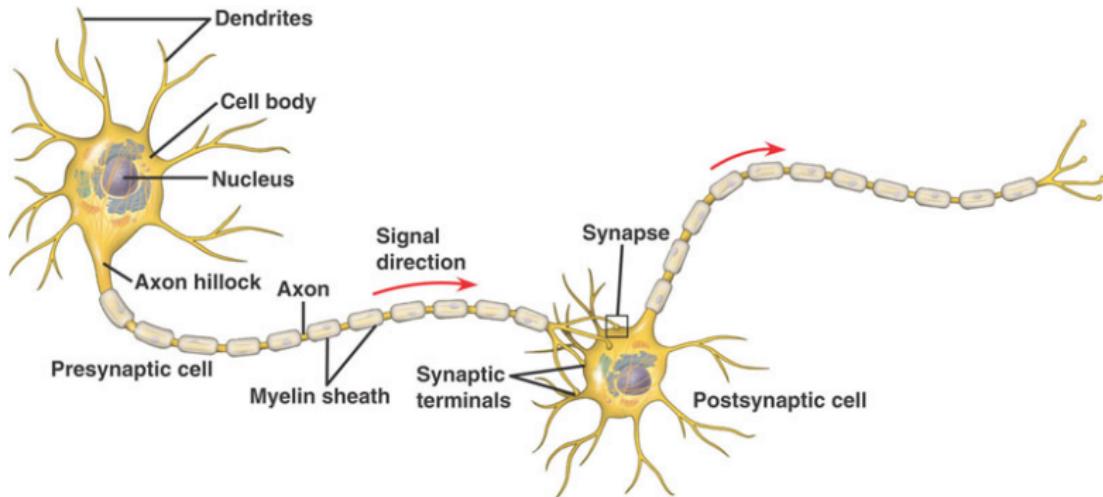
2 Probabilistic Classification

- 2.1 Problem Statement
- 2.2 Bayes Classifier
- 2.3 Linear Discriminant Analysis

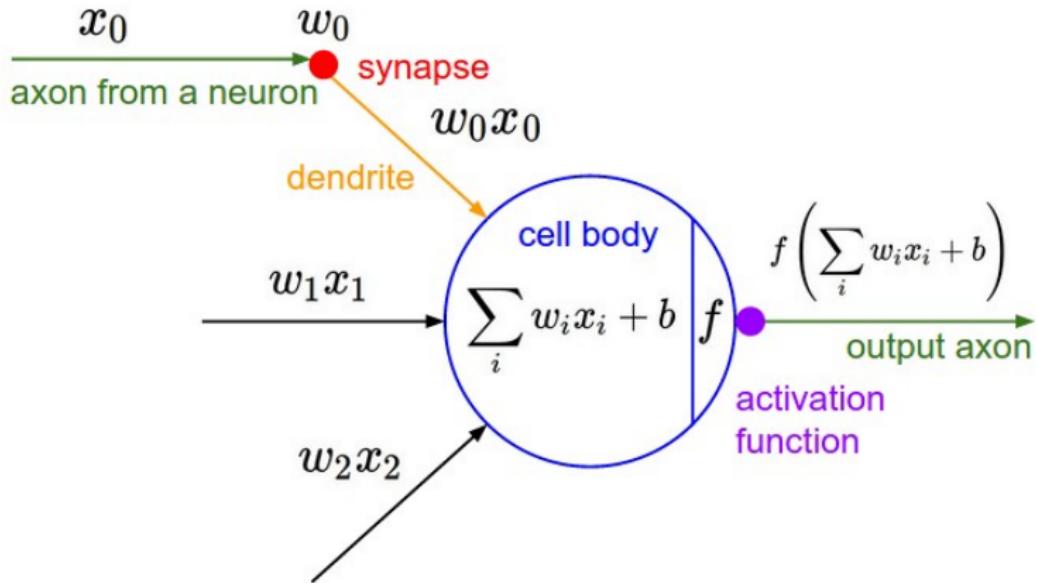
3 Classification and Artificial Intelligence



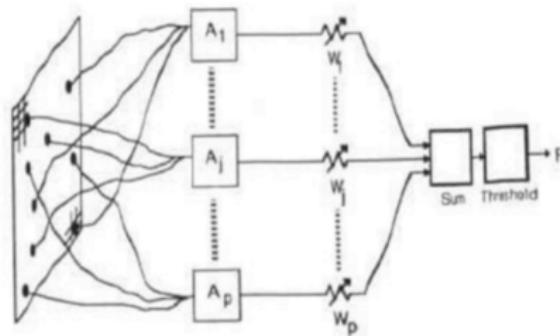
Biological Neuron



Neuron Model (McCulloch & Pitts, 1943)



Perceptron Machine (Rosenblatt, 1958)



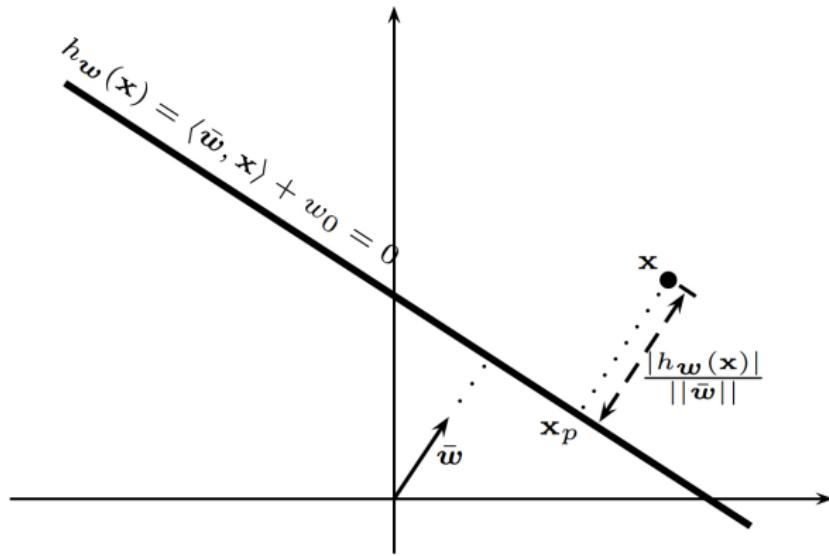
- **Summation:**

$$h_w(\mathbf{x}) := \langle \mathbf{w}, \mathbf{x} \rangle + w_0 = \sum_{j=1}^d w_j x_j + w_0.$$

- **Activation:**

$$a(\mathbf{x}) = \begin{cases} +1, & \text{if } y h_w(\mathbf{x}) \geq 0; \\ -1, & \text{if } y h_w(\mathbf{x}) < 0. \end{cases}$$

We want to find parameters (w_0, \mathbf{w}) such that the distance between the misclassified examples and the decision boundary is minimised.



- We would like to minimise the *perceptron error*:

$$\hat{L}(\mathbf{w}, w_0) = - \sum_{i' \in \mathcal{I}} y_{i'} (\langle \mathbf{w}, \mathbf{x}_{i'} \rangle + w_0).$$

- We would like to minimise the *perceptron error*:

$$\hat{L}(\mathbf{w}, w_0) = - \sum_{i' \in \mathcal{I}} y_{i'} (\langle \mathbf{w}, \mathbf{x}_{i'} \rangle + w_0).$$

- If we take derivatives with respect to the weights:

$$\frac{\partial \hat{L}(\mathbf{w})}{\partial w_0} = - \sum_{i' \in \mathcal{I}} y_{i'},$$

$$\nabla \hat{L}(\mathbf{w}) = - \sum_{i' \in \mathcal{I}} y_{i'} \mathbf{x}_{i'}.$$

- We would like to minimise the *perceptron error*:

$$\hat{L}(\mathbf{w}, w_0) = - \sum_{i' \in \mathcal{I}} y_{i'} (\langle \mathbf{w}, \mathbf{x}_{i'} \rangle + w_0).$$

- If we take derivatives with respect to the weights:

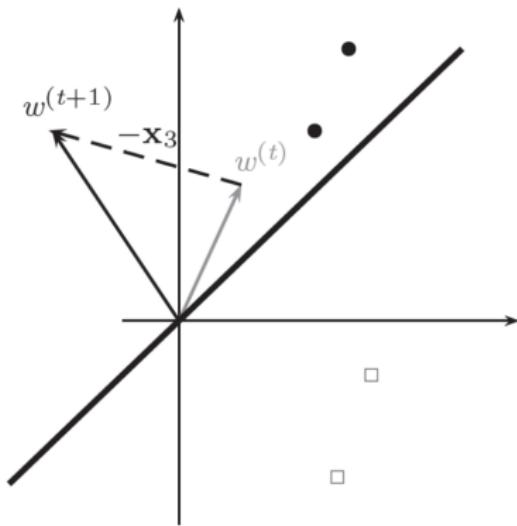
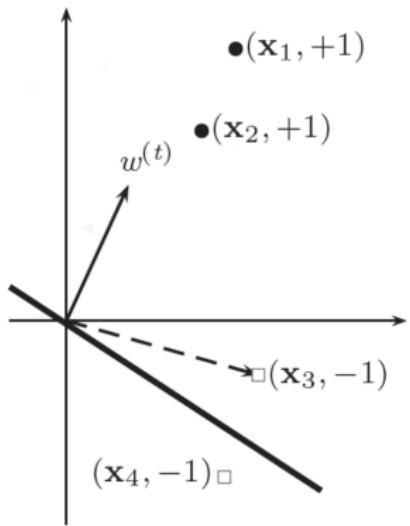
$$\frac{\partial \hat{L}(\mathbf{w})}{\partial w_0} = - \sum_{i' \in \mathcal{I}} y_{i'},$$

$$\nabla \hat{L}(\mathbf{w}) = - \sum_{i' \in \mathcal{I}} y_{i'} \mathbf{x}_{i'}.$$

- Minimisation of the perceptron error is done by stochastic gradient descent algorithm:

$$\text{if } y(\langle \mathbf{w}, \mathbf{x} \rangle + w_0) < 0, \text{ then } \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix} \leftarrow \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix} + \eta \begin{pmatrix} y \\ y\mathbf{x} \end{pmatrix}.$$

Update Rule: Illustration



Algorithm Perceptron

Input: Training set $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$.

Max number of iterations T .

Initialisation: Weights $\mathbf{w}^{(0)} \leftarrow \mathbf{0}, w_0 \leftarrow 0$;

Counter $t \leftarrow 0$, learning rate $\eta > 0$.

repeat

Choose randomly an example $(\mathbf{x}^{(t)}, y^{(t)}) \in S$

if $y \cdot (\langle \mathbf{w}^{(t)}, \mathbf{x}^{(t)} \rangle + w_0) < 0$ **then**

$$w_0^{(t+1)} \leftarrow w_0^{(t)} + \eta y^{(t)}$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta y^{(t)} \mathbf{x}^{(t)}$$

end if

$$t \leftarrow t + 1$$

until $t > T$

Output: $(w_0^{(T)}, \mathbf{w}^{(T)})$.

- What is the time complexity during training phase?

- What is the time complexity during training phase?
- What is the time complexity to predict a label for new x ?

- What is the time complexity during training phase?
- What is the time complexity to predict a label for new x ?
- How to tune learning rate η ?

- What is the time complexity during training phase?
- What is the time complexity to predict a label for new x ?
- How to tune learning rate η ?
- Does the algorithm converge to an exact solution?

- What is the time complexity during training phase?
- What is the time complexity to predict a label for new x ?
- How to tune learning rate η ?
- Does the algorithm converge to an exact solution?
- Under what condition the algorithm does not converge?

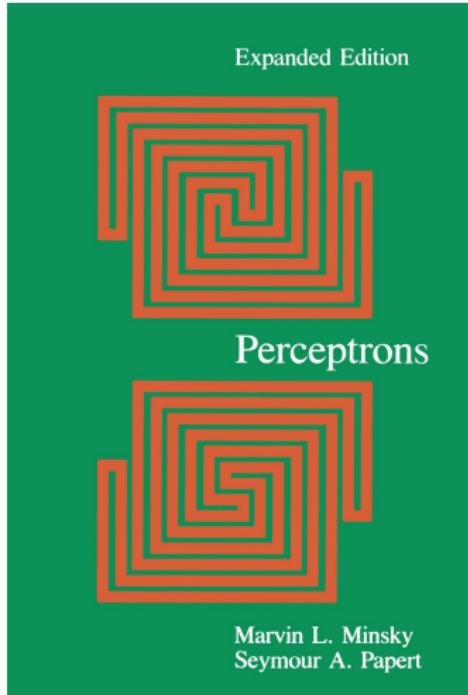


Figure: (Minsky & Papert, 1969)



Figure: Minsky and Papert in 1971.

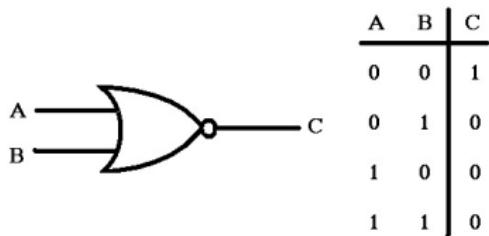
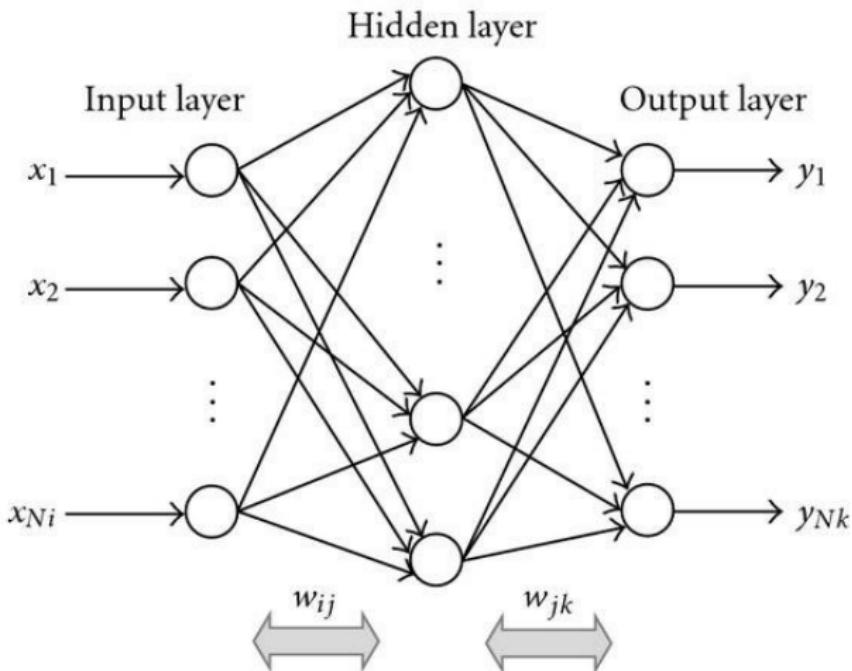
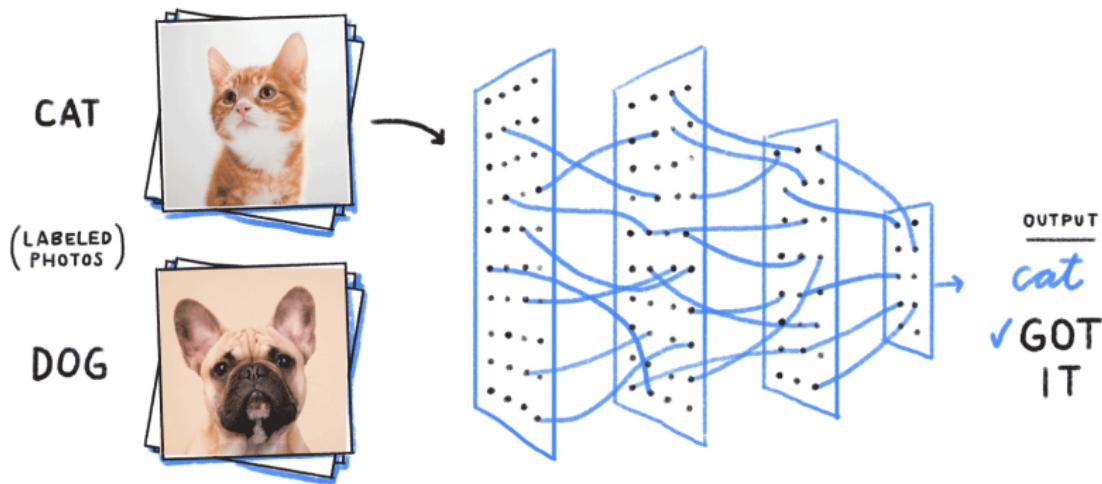


Figure: The XOR operation.





Classification: Part 2

Statistical Analysis and Document Mining
Spring 2019

Vasilii Feofanov
Université Grenoble Alpes
vasilii.feofanov@univ-grenoble-alpes.fr

1 Naive Bayes Classification

1.1 Concept

1.2 Particular Cases

2 Classification in Practice

2.1 Performance Estimation

2.2 Script

- *Input space:* $\mathcal{X} \subseteq \mathbb{R}^d$;
- *Output space:* $\mathcal{Y} = \{-1, +1\}$ (binary classification),
 $\mathcal{Y} = \{1, \dots, K\}$ (multi-class classification);
- *Assumption:* all $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$ are **i.i.d.** from \mathcal{D} with respect to a fixed unknown probability distribution $P(\mathbf{X}, Y)$;
- *Sample Data:* we observe $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$;

- *Input space:* $\mathcal{X} \subseteq \mathbb{R}^d$;
- *Output space:* $\mathcal{Y} = \{-1, +1\}$ (binary classification),
 $\mathcal{Y} = \{1, \dots, K\}$ (multi-class classification);
- *Assumption:* all $(\mathbf{X}, Y) \in \mathcal{X} \times \mathcal{Y}$ are **i.i.d.** from \mathcal{D} with respect to a fixed unknown probability distribution $P(\mathbf{X}, Y)$;
- *Sample Data:* we observe $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$;
- *Loss Function:*

$$\ell^{0/1}(h(\mathbf{x}), y) = \mathbb{I}(h(\mathbf{x}) \neq y) = \begin{cases} 1, & \text{if } h(\mathbf{x}) \neq y; \\ 0, & \text{if } h(\mathbf{x}) = y. \end{cases}$$

- *Target:* minimise the misclassification error:

$$P(h(\mathbf{X}) \neq Y) = \sum_{c \in \{1, \dots, K\}} P(Y = c) P(h(\mathbf{X}) \neq c | Y = c).$$

$$P(Y|\mathbf{X}) = \frac{\text{Posterior}}{\text{Evidence}} \cdot \frac{\text{Likelihood}}{\text{Class Prior}}$$
$$P(Y|\mathbf{X}) = \frac{P(\mathbf{X}|Y)P(Y)}{P(\mathbf{X})}$$

The Bayes classifier predicts a class with the highest posterior probability:

$$h_B(\mathbf{x}) := \operatorname{argmax}_{y \in \mathcal{Y}} P(Y = y | \mathbf{X} = \mathbf{x}).$$

This is equivalent to:

$$\begin{aligned} h_B(\mathbf{x}) &\propto \operatorname{argmax}_{y \in \mathcal{Y}} P(\mathbf{X} = \mathbf{x} | Y = y) P(Y = y) \\ &\propto \operatorname{argmax}_{y \in \mathcal{Y}} \log P(\mathbf{X} = \mathbf{x} | Y = y) + \log P(Y = y). \end{aligned}$$

- *Problem:* When d is large, parametric estimation of $P(\mathbf{X}|Y)$ requires a large number of samples.

- *Problem:* When d is large, parametric estimation of $P(\mathbf{X}|Y)$ requires a large number of samples.
- *Idea:* We assume *naively* that features are conditionally *independent* given the class.

- *Problem:* When d is large, parametric estimation of $P(\mathbf{X}|Y)$ requires a large number of samples.
- *Idea:* We assume *naively* that features are conditionally *independent* given the class.
- Then, denoting $\mathbf{X} = (X_1, \dots, X_d)$, $\mathbf{x} = (x_1, \dots, x_d)$, we obtain that:

$$P(\mathbf{X} = \mathbf{x}|Y = c) = P(X_1 = x_1|Y = c) \cdots P(X_d = x_d|Y = c).$$

- *Problem:* When d is large, parametric estimation of $P(\mathbf{X}|Y)$ requires a large number of samples.
- *Idea:* We assume *naively* that features are conditionally *independent* given the class.
- Then, denoting $\mathbf{X} = (X_1, \dots, X_d)$, $\mathbf{x} = (x_1, \dots, x_d)$, we obtain that:

$$P(\mathbf{X} = \mathbf{x}|Y = c) = P(X_1 = x_1|Y = c) \cdots P(X_d = x_d|Y = c).$$

- Thus, the *naive Bayes classifier* is defined in the following way:

$$\begin{aligned} h_B(\mathbf{x}) &:= \operatorname{argmax}_{c \in \mathcal{Y}} P(Y = c) \prod_{j=1}^d P(X_j = x_j|Y = c) \\ &\propto \operatorname{argmax}_{c \in \mathcal{Y}} \log P(Y = c) + \sum_{j=1}^d \log P(X_j = x_j|Y = c). \end{aligned}$$

- We assume that the j -th feature of observations from the class $c \in \mathcal{Y}$ is normally distributed:

$$[X_j | Y = c] \sim \mathcal{N}(\mu_{j,c}, \sigma_{j,c}).$$

- We assume that the j -th feature of observations from the class $c \in \mathcal{Y}$ is normally distributed:

$$[X_j | Y = c] \sim \mathcal{N}(\mu_{j,c}, \sigma_{j,c}).$$

- Then, the distribution of the class $c \in \mathcal{Y}$ is defined as:

$$P(\mathbf{X} = \mathbf{x} | Y = c) = \prod_{j=1}^d \frac{1}{\sqrt{2\pi}\sigma_{j,c}} e^{-\frac{(x_j - \mu_{j,c})^2}{2\sigma_{j,c}^2}}.$$

- We assume that the j -th feature of observations from the class $c \in \mathcal{Y}$ is normally distributed:

$$[X_j | Y = c] \sim \mathcal{N}(\mu_{j,c}, \sigma_{j,c}).$$

- Then, the distribution of the class $c \in \mathcal{Y}$ is defined as:

$$P(\mathbf{X} = \mathbf{x} | Y = c) = \prod_{j=1}^d \frac{1}{\sqrt{2\pi}\sigma_{j,c}} e^{-\frac{(x_j - \mu_{j,c})^2}{2\sigma_{j,c}^2}}.$$

- Finally, the Gaussian naive Bayes classifier is defined as:

$$h_B(\mathbf{x}) := \operatorname{argmax}_{c \in \mathcal{Y}} \left[\ln P(Y = c) - \sum_{j=1}^d \ln \sigma_{j,c} - \sum_{j=1}^d \frac{(x_j - \mu_{j,c})^2}{2\sigma_{j,c}^2} \right].$$

Bernoulli Naive Bayes Classifier

- Each feature X_j is a binary variable ($X_j \in \{0, 1\}$).

- Each feature X_j is a binary variable ($X_j \in \{0, 1\}$).
- Suppose that $X_j|Y = c$ is distributed according to *Bernoulli* distribution so that the probability of success and failure are defined respectively as:

$$P(X_j = 1|Y = c) = p_{j,c};$$

$$P(X_j = 0|Y = c) = 1 - p_{j,c}.$$

- Each feature X_j is a binary variable ($X_j \in \{0, 1\}$).
- Suppose that $X_j|Y = c$ is distributed according to *Bernoulli* distribution so that the probability of success and failure are defined respectively as:

$$P(X_j = 1|Y = c) = p_{j,c};$$

$$P(X_j = 0|Y = c) = 1 - p_{j,c}.$$

- Then, the Bernoulli naive Bayes classifier is defined as:

$$\begin{aligned} h_B(\mathbf{x}) &:= \operatorname{argmax}_{c \in \mathcal{Y}} P(Y = c) \prod_{j=1}^d p_{j,c}^{x_j} (1 - p_{j,c})^{1-x_j} \\ &\propto \operatorname{argmax}_{c \in \mathcal{Y}} \log P(Y = c) + \sum_{j=1}^d x_j \log p_{j,c} + \sum_{j=1}^d (1 - x_j) \log(1 - p_{j,c}). \end{aligned}$$

Multinomial Distribution

- Let T_1, \dots, T_m be i.i.d. random variables that represent trials with the output $\in \{1, \dots, d\}$ and are distributed as follows:

$$P(T_i = j) = p_j, \quad i \in \{1, \dots, m\}, \quad j \in \{1, \dots, d\}, \quad \sum_{j=1}^d p_j = 1.$$

- Let T_1, \dots, T_m be i.i.d. random variables that represent trials with the output $\in \{1, \dots, d\}$ and are distributed as follows:

$$P(T_i = j) = p_j, \quad i \in \{1, \dots, m\}, \quad j \in \{1, \dots, d\}, \quad \sum_{j=1}^d p_j = 1.$$

- Random variable X_j is a number of trials with the outcome j :

$$X_j = \sum_{i=1}^m \mathbb{I}(T_i = j), \quad j \in \{1, \dots, d\}, \quad \sum_{j=1}^d X_j = m.$$

- Let T_1, \dots, T_m be i.i.d. random variables that represent trials with the output $\in \{1, \dots, d\}$ and are distributed as follows:

$$P(T_i = j) = p_j, \quad i \in \{1, \dots, m\}, \quad j \in \{1, \dots, d\}, \quad \sum_{j=1}^d p_j = 1.$$

- Random variable X_j is a number of trials with the outcome j :

$$X_j = \sum_{i=1}^m \mathbb{I}(T_i = j), \quad j \in \{1, \dots, d\}, \quad \sum_{j=1}^d X_j = m.$$

- Then, the *multinomial* distribution $P(X_1, \dots, X_d)$ is defined as:

$$P(X_1 = x_1, \dots, X_d = x_d) = \frac{m!}{x_1! \cdots x_d!} p_1^{x_1} \cdots p_d^{x_d}.$$

- Suppose $\mathbf{X}|Y = c$ is distributed according to the multinomial distribution with parameters $\mathbf{p}_c = (p_{1,c}, \dots, p_{d,c})$. Denoting $\mathbf{X} = (X_1, \dots, X_d)$, $\mathbf{x} = (x_1, \dots, x_d)$, we obtain:

$$P(\mathbf{X} = \mathbf{x}|Y = c) \propto \prod_{j=1}^d p_{j,c}^{x_j}.$$

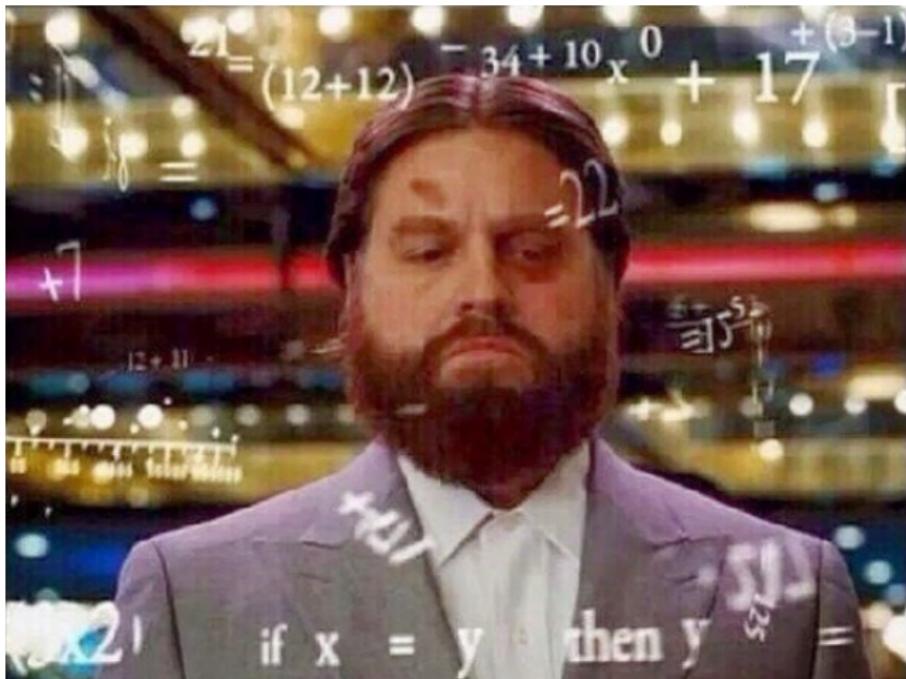
- Suppose $\mathbf{X}|Y = c$ is distributed according to the multinomial distribution with parameters $\mathbf{p}_c = (p_{1,c}, \dots, p_{d,c})$. Denoting $\mathbf{X} = (X_1, \dots, X_d)$, $\mathbf{x} = (x_1, \dots, x_d)$, we obtain:

$$P(\mathbf{X} = \mathbf{x}|Y = c) \propto \prod_{j=1}^d p_{j,c}^{x_j}.$$

- Then, we classify by the following rule:

$$\begin{aligned} h_B(\mathbf{x}) &:= \operatorname{argmax}_{c \in \mathcal{Y}} P(Y = c) \prod_{j=1}^d p_{j,c}^{x_j} \\ &\propto \operatorname{argmax}_{c \in \mathcal{Y}} P(Y = c) \sum_{j=1}^d x_j \log p_{j,c}. \end{aligned}$$

Why These Distributions?





1 Naive Bayes Classification

1.1 Concept

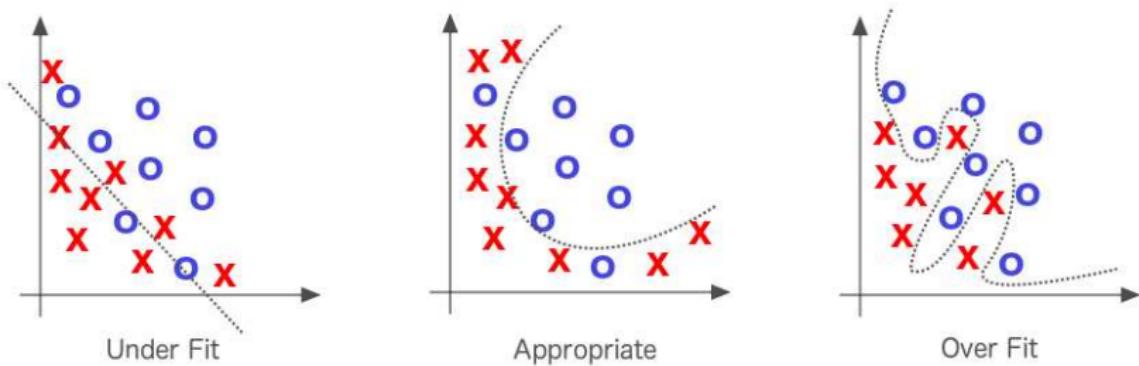
1.2 Particular Cases

2 Classification in Practice

2.1 Performance Estimation

2.2 Script

Performance of an Algorithm



Question: How to estimate the error value (e.g. $P(h(\mathbf{X}) \neq Y)$) in practice when the sample data $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$ is available only?

- Error on the training set?

Question: How to estimate the error value (e.g. $P(h(\mathbf{X}) \neq Y)$) in practice when the sample data $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$ is available only?

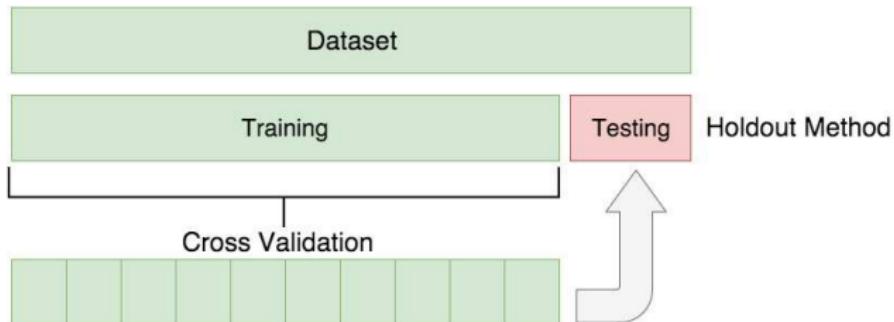
- Error on the training set? \Rightarrow Overfitting.
- Train/test split?

Question: How to estimate the error value (e.g. $P(h(\mathbf{X}) \neq Y)$) in practice when the sample data $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$ is available only?

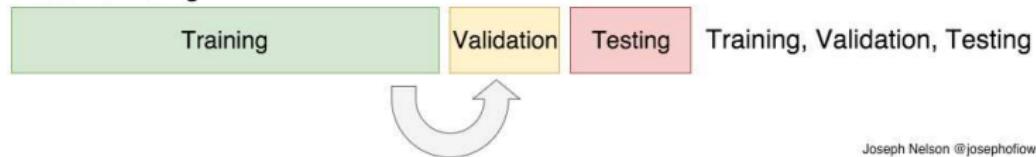
- Error on the training set? \Rightarrow Overfitting.
- Train/test split? \Rightarrow We should have enough data.
- Cross-validation?

Question: How to estimate the error value (e.g. $P(h(\mathbf{X}) \neq Y)$) in practice when the sample data $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$ is available only?

- Error on the training set? \Rightarrow Overfitting.
- Train/test split? \Rightarrow We should have enough data.
- Cross-validation? What is an impact of the number of folds?



Data Permitting:



Joseph Nelson @josephoflawa

Leave-One-Out



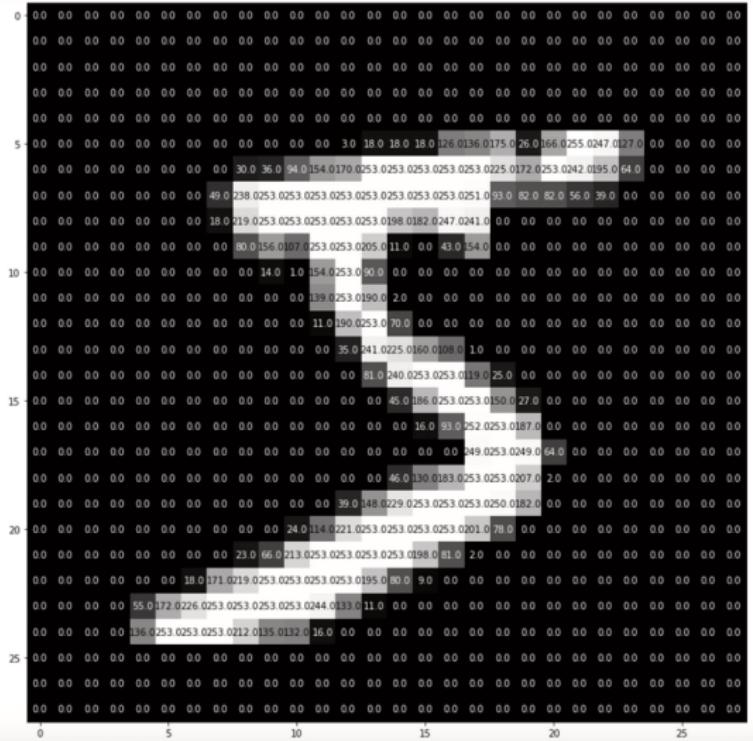
Confusion Matrix

| n=165 | Predicted: | |
|----------------|------------|-----|
| | NO | YES |
| Actual: NO | 50 | 10 |
| Actual: YES | 5 | 100 |

Please follow the link:

`https://chamilo.grenoble-inp.fr/courses/ENSIMAG4MMSADM/document/DemoR/script_classif_part_I.html`

MNIST Features



Document Classification

Statistical Analysis and Document Mining
Spring 2019

Vasilii Feofanov, Massih-Reza Amini
Université Grenoble Alpes
vasilii.feofanov@univ-grenoble-alpes.fr

1 Introduction

- 1.1 Motivation
- 1.2 Outline

2 First Look at the Problem

- 2.1 Preprocessing
- 2.2 Bag-of-Words
- 2.3 Naive Bayes Document Classification

3 Advanced Look at the Problem

- 3.1 Sparse Files
- 3.2 Two Laws of IR
- 3.3 tf-idf and Other Techniques

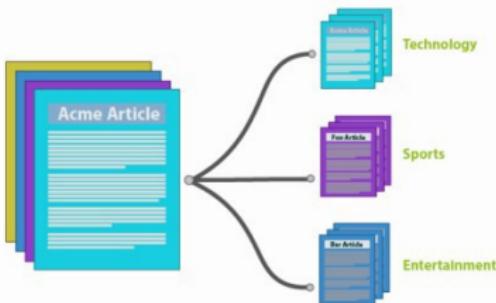
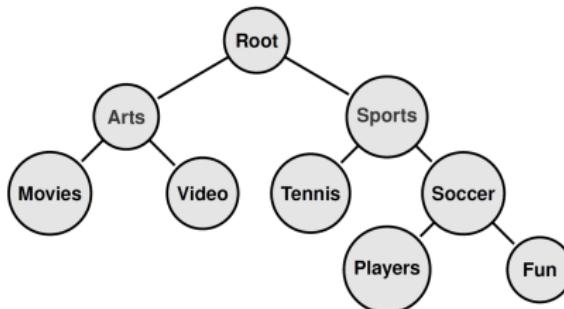
Automatic Document Classification



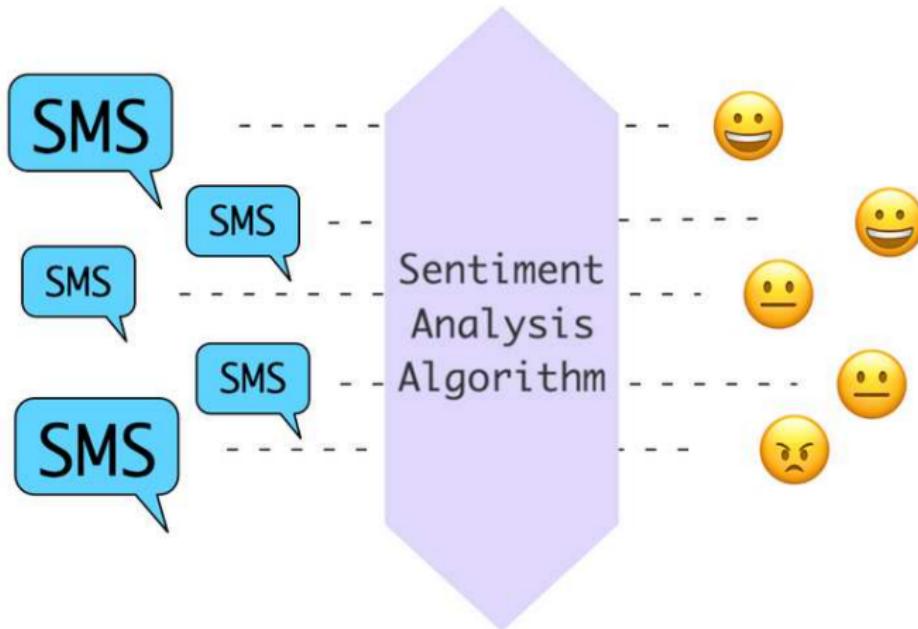
Application 1: Spam Detection



Application 2: Genre Classification



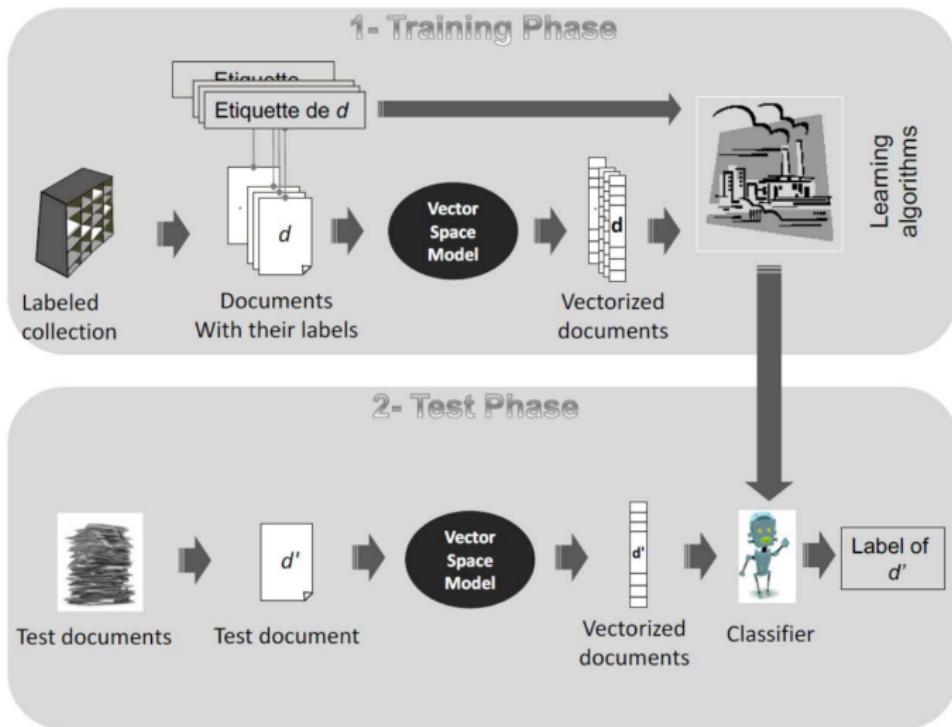
Application 3: Sentiment Analysis



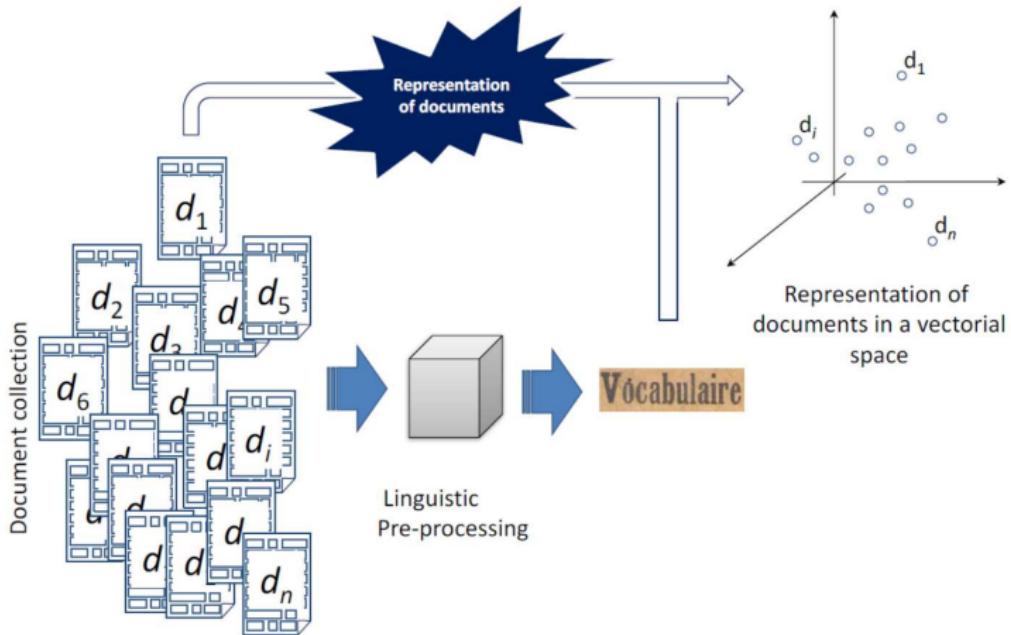
Application 4: Language Identification



Classification Scheme



From Document to Vector



1 Introduction

- 1.1 Motivation
- 1.2 Outline

2 First Look at the Problem

- 2.1 Preprocessing
- 2.2 Bag-of-Words
- 2.3 Naive Bayes Document Classification

3 Advanced Look at the Problem

- 3.1 Sparse Files
- 3.2 Two Laws of IR
- 3.3 tf-idf and Other Techniques

- *Segmentation* (tokenization): separate a sequence of characters into semantic elements, or *words*.
- *Term* (type of words): class of all words having the same sequence of characters.
- *Example*:
"The cat sat on the mat."
Words: The, cat, sat, on, the, mat
Terms: the, cat, sat, on, mat
- *Dificulty*: Tokenization is language specific.

In French, the following issues may arise during the segmentation process:

- Lexical components with hyphens:
chassé-croisé, peut-être, rendez-vous
- Lexical components with an apostrophe:
jusqu' où, aujourd'hui, prud'homme
- Idiomatic expressions:
au fait, poser un lapin, tomber dans les pommes
- Contracted forms:
j', M'sieur, Gad'zarts (les gars des Arts et Métiers)
- Acronyms:
K7, A.R., CV, càd, P.-V.

1 *Textual normalization:* consists in reducing the words of a same family to their canonical forms.

- Punctuation: suppression of points and hyphens;
- Lower-upper case: transform all upper cases to lower cases;
- Accents: suppression of accents.

2 *Linguistic normalization* consists in

- Rooting: replace each word by its root;
- Stemming: replace each word by its canonical form.

Non-spam message before preprocessing

Subject: Re: 5.1344 Native speaker intuitions The discussion on native speaker intuitions has been extremely interesting, but I worry that my brief intervention may have muddied the waters. I take it that there are a number of separable issues. The first is the extent to which a native speaker is likely to judge a lexical string as grammatical or ungrammatical per se. The second is concerned with the relationships between syntax and interpretation (although even here the distinction may not be entirely clear cut).

Non-spam message before preprocessing

Subject: Re: 5.1344 Native speaker intuitions The discussion on native speaker intuitions has been extremely interesting, but I worry that my brief intervention may have muddied the waters. I take it that there are a number of separable issues. The first is the extent to which a native speaker is likely to judge a lexical string as grammatical or ungrammatical per se. The second is concerned with the relationships between syntax and interpretation (although even here the distinction may not be entirely clear cut).

Non-spam message after preprocessing

re native speaker intuition discussion native speaker intuition extremely interest worry brief intervention muddy waters number separable issue first extent native speaker likely judge lexical string grammatical ungrammatical per se second concern relationship between syntax interpretation although even here distinction entirely clear cut

Non-spam message after preprocessing

re native speaker intuition discussion native speaker intuition extremely interest
worry brief intervention muddy waters number separable issue first extent
native speaker likely judge lexical string grammatical ungrammatical per se
second concern relationship between syntax interpretation although even here
distinction entirely clear cut

Spam message after preprocessing

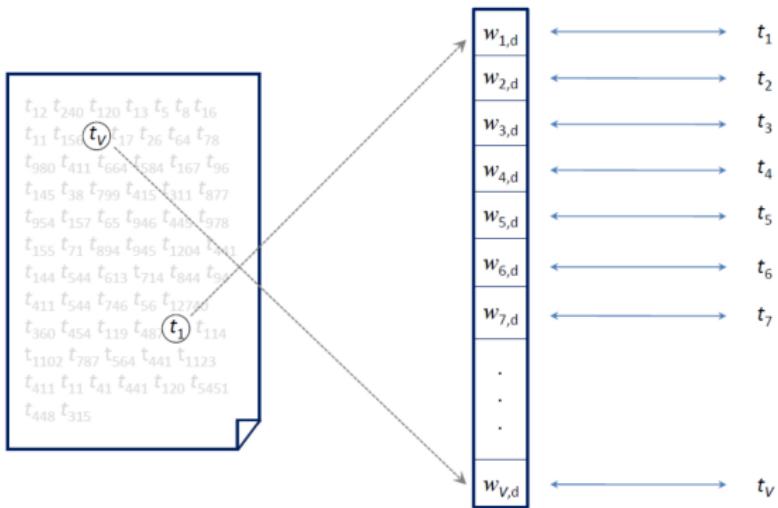
financial freedom follow financial freedom work ethic extraordinary desire earn
least per month work home special skills experience required train personal
support need ensure success legitimate homebased income opportunity put
back control finance life ve try opportunity past fail live promise

The Bag of Words Representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!

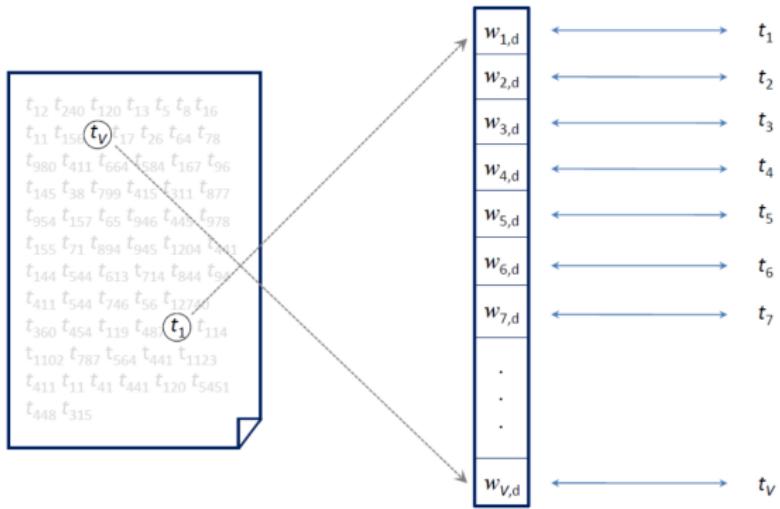


In the bag-of-words the word order in a sentence is ignored. Given a document, to each term t_j a specific value $w_{j,d}$ is assigned.



In the bag-of-words the word order in a sentence is ignored. Given a document, to each term t_j a specific value $w_{j,d}$ is assigned.

⇒ A document can be represented as a vector: $\mathbf{d} = (w_{j,d})_{j=1}^V$.



- Observation is a document $\mathbf{d} = (w_{j,d})_{i=1}^V$, where
 - V is a vocabulary size;
 - $w_{j,d}$ is a value (importance) corresponding to the term t_j .

- Observation is a document $\mathbf{d} = (w_{j,d})_{i=1}^V$, where
 - V is a vocabulary size;
 - $w_{j,d}$ is a value (importance) corresponding to the term t_j .
- Training set: $S = \{\mathbf{d}_i, y_i\}_{i=1}^n$, where
 - $\mathbf{d}_i \in \mathbb{R}^V$;
 - $y_i \in \{1, \dots, K\}$.

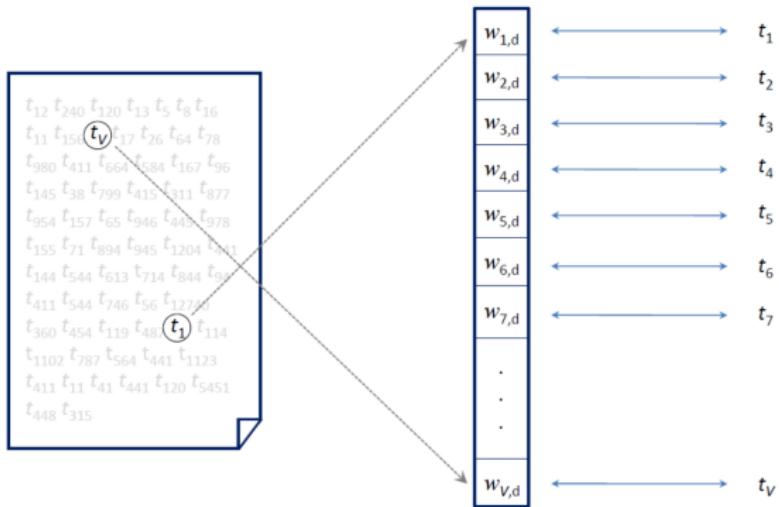
- Observation is a document $\mathbf{d} = (w_{j,d})_{i=1}^V$, where
 - V is a vocabulary size;
 - $w_{j,d}$ is a value (importance) corresponding to the term t_j .
- Training set: $S = \{\mathbf{d}_i, y_i\}_{i=1}^n$, where
 - $\mathbf{d}_i \in \mathbb{R}^V$;
 - $y_i \in \{1, \dots, K\}$.
- *Target:* minimise the misclassification error:

$$P(h(\mathbf{D}) \neq Y) = \sum_{c \in \{1, \dots, K\}} P(Y = c) P(h(\mathbf{D}) \neq c | Y = c).$$

How to Construct a Vector Space?

A document is represented as a vector: $\mathbf{d} = (w_{j,d})_{j=1}^V$.

How we define importance $w_{j,d}$, $j \in \{1, \dots, V\}$?



- If term t_j is present in the document, then $w_{j,d} = 1$.
Otherwise, $w_{j,d} = 0$.

- If term t_j is present in the document, then $w_{j,d} = 1$.
Otherwise, $w_{j,d} = 0$.
- Suppose conditional independence of terms given a class.

- If term t_j is present in the document, then $w_{j,d} = 1$.
Otherwise, $w_{j,d} = 0$.
- Suppose conditional independence of terms given a class.
- *Question:* How then the term t_j contributes into classification?

- If term t_j is present in the document, then $w_{j,d} = 1$.
Otherwise, $w_{j,d} = 0$.
- Suppose conditional independence of terms given a class.
- *Question:* How then the term t_j contributes into classification?
 - For each class, we can compute the probability that the term is present in a document from this class.
 \Rightarrow Bernoulli distribution.

- Let $W_{j,d}|Y = c$ be distributed acc. to *Bernoulli*. Then, probability that t_j is present in a document \mathbf{d} is:

$$P(W_{j,d}|Y = c) = p_{t_j|c}.$$

- Let $W_{j,d}|Y = c$ be distributed acc. to *Bernoulli*. Then, probability that t_j is present in a document \mathbf{d} is:

$$P(W_{j,d}|Y = c) = p_{t_j|c}.$$

- Naive assumption for documents:

$$P(\mathbf{D} = \mathbf{d}|Y = c) = \prod_{j=1}^V P(W_{j,d} = w_{j,d}|Y = c).$$

- Let $W_{j,d}|Y = c$ be distributed acc. to *Bernoulli*. Then, probability that t_j is present in a document \mathbf{d} is:

$$P(W_{j,d}|Y = c) = p_{t_j|c}.$$

- Naive assumption for documents:

$$P(\mathbf{D} = \mathbf{d}|Y = c) = \prod_{j=1}^V P(W_{j,d} = w_{j,d}|Y = c).$$

- Bernoulli naive Bayes decision rule for the document classification task:

$$\begin{aligned} h_{NB}(\mathbf{d}) &:= \operatorname{argmax}_{c \in \mathcal{Y}} P(Y = c) \prod_{j=1}^V P(W_{j,d} = w_{j,d}|Y = c) \\ &= \operatorname{argmax}_{c \in \mathcal{Y}} P(Y = c) \prod_{j=1}^V p_{t_j|c}^{w_{j,d}} (1 - p_{t_j|c})^{1-w_{j,d}}. \end{aligned}$$

- Estimation of priors:

$$P(Y = c) \leftarrow \frac{1}{n} \sum_{i=1}^n \mathbb{I}(y_i = c) =: \frac{n_c}{n}.$$

- Estimation of priors:

$$P(Y = c) \leftarrow \frac{1}{n} \sum_{i=1}^n \mathbb{I}(y_i = c) =: \frac{n_c}{n}.$$

- Estimation of the Bernoulli parameters:

$$\hat{p}_{t_j|c} = \frac{\sum_{i=1}^n \mathbb{I}(w_{j,d_i} = 1 \wedge y_i = c)}{\sum_{i=1}^n \mathbb{I}(y_i = c)} =: \frac{\text{df}_{t_j}(c)}{n_c}.$$

- What happens if $\hat{p}_{t_j|c} = 0$?

- Estimation of priors:

$$P(Y = c) \leftarrow \frac{1}{n} \sum_{i=1}^n \mathbb{I}(y_i = c) =: \frac{n_c}{n}.$$

- Estimation of the Bernoulli parameters:

$$\hat{p}_{t_j|c} = \frac{\sum_{i=1}^n \mathbb{I}(w_{j,d_i} = 1 \wedge y_i = c)}{\sum_{i=1}^n \mathbb{I}(y_i = c)} =: \frac{\text{df}_{t_j}(c)}{n_c}.$$

- What happens if $\hat{p}_{t_j|c} = 0$?
- It would be better to estimate $\hat{p}_{t_j|c}$ by:

$$\hat{p}_{t_j|c} = \frac{\text{df}_{t_j}(c) + 1}{n_c + 2}.$$

- Estimation of priors:

$$P(Y = c) \leftarrow \frac{1}{n} \sum_{i=1}^n \mathbb{I}(y_i = c) =: \frac{n_c}{n}.$$

- Estimation of the Bernoulli parameters:

$$\hat{p}_{t_j|c} = \frac{\sum_{i=1}^n \mathbb{I}(w_{j,d_i} = 1 \wedge y_i = c)}{\sum_{i=1}^n \mathbb{I}(y_i = c)} =: \frac{\text{df}_{t_j}(c)}{n_c}.$$

- What happens if $\hat{p}_{t_j|c} = 0$?

- It would be better to estimate $\hat{p}_{t_j|c}$ by:

$$\hat{p}_{t_j|c} = \frac{\text{df}_{t_j}(c) + 1}{n_c + 2}.$$

- Considering the log scale, we classify new \mathbf{d} as follows:

$$h_{NB}(\mathbf{d}) := \operatorname{argmax}_{c \in \mathcal{Y}} \ln \frac{n_c}{n} + \sum_{\substack{j=1, \dots, V \\ t_j \in d}} \ln \hat{p}_{t_j|c} + \sum_{\substack{j=1, \dots, V \\ t_j \notin d}} \ln (1 - \hat{p}_{t_j|c}).$$

- We assume that more often term t_j appears in a document, more important the term is.

- We assume that more often term t_j appears in a document, more important the term is.
- Then, for each t_j , we compute the *term frequency*, i.e. the number of occurrences of t_j in the document:

$$w_{j,d} = \text{tf}_{t_j,d}.$$

- We assume that more often term t_j appears in a document, more important the term is.
- Then, for each t_j , we compute the *term frequency*, i.e. the number of occurrences of t_j in the document:

$$w_{j,d} = \text{tf}_{t_j,d}.$$

- Let N_d be the number of words in the document.
Then, $\sum_{j=1}^V \text{tf}_{t_j,d} = N_d$.

- We assume that more often term t_j appears in a document, more important the term is.
- Then, for each t_j , we compute the *term frequency*, i.e. the number of occurrences of t_j in the document:

$$w_{j,d} = \text{tf}_{t_j,d}.$$

- Let N_d be the number of words in the document.
Then, $\sum_{j=1}^V \text{tf}_{t_j,d} = N_d$.
- *Naive assumption*: all N_d words are *i.i.d.* by the following law:

$$P(\text{word} = t_j) = p_{t_j|c}, \quad j = \{1, \dots, V\}, \quad \sum_{j=1}^V p_{t_j|c} = 1.$$

⇒ Then, document **d** is distributed acc. to multinomial distribution with parameters $(p_{t_j|c})_{j=1}^V$.

- Let $\mathbf{D}|Y = c$ be distributed acc. to the *multinomial law*.
Then, the likelihood of \mathbf{D} to be from the class c is:

$$\begin{aligned} P(\mathbf{D}|Y = c) &= \frac{N_d!}{w_{1,d}! \cdots w_{V,d}!} p_{t_1|c}^{w_{1,d}} \cdots p_{t_V|c}^{w_{V,d}} \\ &\propto \prod_{j=1}^V p_{t_j|c}^{w_{j,d}}. \end{aligned}$$

- Let $\mathbf{D}|Y = c$ be distributed acc. to the *multinomial law*. Then, the likelihood of \mathbf{D} to be from the class c is:

$$P(\mathbf{D}|Y = c) = \frac{N_d!}{w_{1,d}! \cdots w_{V,d}!} p_{t_1|c}^{w_{1,d}} \cdots p_{t_V|c}^{w_{V,d}}$$

$$\propto \prod_{j=1}^V p_{t_j|c}^{w_{j,d}}.$$

- Multinomial naive Bayes decision rule for the document classification task:

$$h_{NB}(\mathbf{d}) := \operatorname{argmax}_{c \in \mathcal{Y}} P(Y = c) \prod_{j=1}^V p_{t_j|c}^{w_{j,d}}$$

$$\propto \operatorname{argmax}_{c \in \mathcal{Y}} \ln P(Y = c) + \sum_{j=1}^V w_{j,d} \ln p_{t_j|c}.$$

- Estimation of priors:

$$P(Y = c) \leftarrow \frac{n_c}{n}.$$

- Estimation of priors:

$$P(Y = c) \leftarrow \frac{n_c}{n}.$$

- Estimation of the multinomial distribution parameters:

$$\hat{p}_{t_j|c} = \frac{\sum_{\substack{i=\{1,\dots,n\} \\ \text{s.t. } y_i=c}} \text{tf}_{t_j,d}}{\sum_{j'=1}^V \sum_{\substack{i=\{1,\dots,n\} \\ \text{s.t. } y_i=c}} \text{tf}_{t_{j'},d}}.$$

- Estimation of priors:

$$P(Y = c) \leftarrow \frac{n_c}{n}.$$

- Estimation of the multinomial distribution parameters:

$$\hat{p}_{t_j|c} = \frac{\sum_{\substack{i=\{1,\dots,n\} \\ \text{s.t. } y_i=c}} \text{tf}_{t_j,d}}{\sum_{j'=1}^V \sum_{\substack{i=\{1,\dots,n\} \\ \text{s.t. } y_i=c}} \text{tf}_{t_{j'},d}}.$$

- It would be better to estimate $\hat{p}_{t_j|c}$ by:

$$\hat{p}_{t_j|c} = \frac{\sum_{\substack{i=\{1,\dots,n\} \\ \text{s.t. } y_i=c}} \text{tf}_{t_j,d} + 1}{\sum_{j'=1}^V \sum_{\substack{i=\{1,\dots,n\} \\ \text{s.t. } y_i=c}} \text{tf}_{t_{j'},d} + V}.$$

- Estimation of priors:

$$P(Y = c) \leftarrow \frac{n_c}{n}.$$

- Estimation of the multinomial distribution parameters:

$$\hat{p}_{t_j|c} = \frac{\sum_{\substack{i=1, \dots, n \\ \text{s.t. } y_i=c}} \text{tf}_{t_j,d}}{\sum_{j'=1}^V \sum_{\substack{i=1, \dots, n \\ \text{s.t. } y_i=c}} \text{tf}_{t_{j'},d}}.$$

- It would be better to estimate $\hat{p}_{t_j|c}$ by:

$$\hat{p}_{t_j|c} = \frac{\sum_{\substack{i=1, \dots, n \\ \text{s.t. } y_i=c}} \text{tf}_{t_j,d} + 1}{\sum_{j'=1}^V \sum_{\substack{i=1, \dots, n \\ \text{s.t. } y_i=c}} \text{tf}_{t_{j'},d} + V}.$$

- We classify new \mathbf{d} by the following rule:

$$h_{NB}(\mathbf{d}) := \operatorname{argmax}_{c \in \mathcal{Y}} \ln \frac{n_c}{n} + \sum_{j=1}^V w_{j,d} \ln \hat{p}_{t_j|c}.$$

- What happens when the training documents of different size?
 - Bernoulli?

- What happens when the training documents of different size?
 - Bernoulli?
 - Multinomial?

- What happens when the training documents of different size?
 - Bernoulli?
 - Multinomial?
- What is the main drawback of
 - the binary weighting?

- What happens when the training documents of different size?
 - Bernoulli?
 - Multinomial?

- What is the main drawback of
 - the binary weighting?
 - the term frequency weighting?

- What happens when the training documents of different size?
 - Bernoulli?
 - Multinomial?
- What is the main drawback of
 - the binary weighting?
 - the term frequency weighting?
 - the bag of words representation?

- What happens when the training documents of different size?
 - Bernoulli?
 - Multinomial?

- What is the main drawback of
 - the binary weighting?
 - the term frequency weighting?
 - the bag of words representation?
 - the naive Bayes classifier?

1 Introduction

- 1.1 Motivation
- 1.2 Outline

2 First Look at the Problem

- 2.1 Preprocessing
- 2.2 Bag-of-Words
- 2.3 Naive Bayes Document Classification

3 Advanced Look at the Problem

- 3.1 Sparse Files
- 3.2 Two Laws of IR
- 3.3 tf-idf and Other Techniques

- Generally, there is access to a large collection of documents.

- Generally, there is access to a large collection of documents.
- At the same time, the vocabulary size is also large.

- Generally, there is access to a large collection of documents.
- At the same time, the vocabulary size is also large.
- *Problem:* Need to store a matrix of size $n \cdot V$.

- Generally, there is access to a large collection of documents.
- At the same time, the vocabulary size is also large.
- *Problem:* Need to store a matrix of size $n \cdot V$.
 - What is approximately the size of the dataset in GB, if $n = 200,000$, $V = 10,000$ and one entry needs 8 bytes?

- Generally, there is access to a large collection of documents.
- At the same time, the vocabulary size is also large.
- *Problem:* Need to store a matrix of size $n \cdot V$.
 - What is approximately the size of the dataset in GB, if $n = 200,000$, $V = 10,000$ and one entry needs 8 bytes?
 - Around 16 GB!

| Variables | Values |
|---|------------------|
| # of documents in the collection | 1,349,539 |
| Total # of occurrences of words | 696,668,157 |
| Average # of words per document | 416 |
| Size of the pre-processed collection on the disk | 4.6 GB |
| Total # of types of words | 757,476 |
| Total # of types of words after rooting | 604,244 |
| Size of the vocabulary | 604,244 |
| Average # of terms per document | 225 |
| Size of the collection after removing a stop-list | 2.8 GB |

- *Sparse matrix:* most of entries are zeros.

- *Sparse matrix*: most of entries are zeros.
- No need to store entries with zero values.

- *Sparse matrix*: most of entries are zeros.
- No need to store entries with zero values.
- *LibSVM format*: the observation is written in the following format:

| y | index-value | index-value |
|---|-------------|-------------|
| 2 | 5:0.356 | ... |
| 3 | 2:10.2 | ... |

| Variables | Values |
|---|-------------|
| # of documents in the collection | 1,349,539 |
| Total # of occurrences of words | 696,668,157 |
| Average # of distinct words per document | 416 |
| Size of the pre-processed collection on the disk | 4.6 GB |
| Total # of types of words | 757,476 |
| Total # of types of words after rooting | 604,244 |
| Size of the vocabulary | 604,244 |
| Average # of terms per document | 225 |
| Size of the collection after removing a stop-list | 2.8 GB |

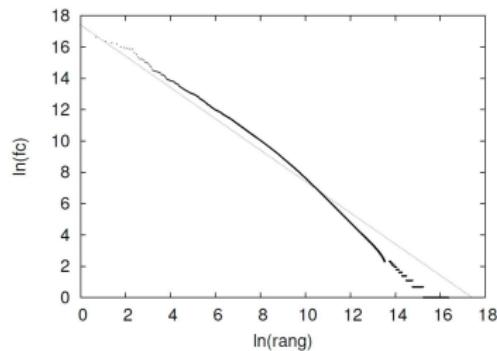
- The encyclopedia "Grand Robert" contains around 75,000 words. The most extensive record shows that the French language would contain about 700,000 words. Why in French Wikipedia we found even more words (757,476)?

- The encyclopedia "Grand Robert" contains around 75,000 words. The most extensive record shows that the French language would contain about 700,000 words. Why in French Wikipedia we found even more words (757,476)?
- When the collection was filtered by removing a stop-list ("a", "the", "of", etc.) of size 200 words, the average number of terms was reduced in documents from 416 to 225 (around 45% reduction). Why?
In addition, their filtering reduces the space on the disk of about 39% (from 4.6 GB to 2.8 GB).

The number of occurrences $fc(m)$ of a word m in a document collection is inversely proportional to its rank:

$$\forall m : fc(m) \approx \frac{\lambda}{\text{rang}(m)}.$$

⇒ The k -th most frequent word is approximately k times less present than the most frequent one.



| Rank | Word | Freq. | % |
|------|-----------|------------|-------|
| 1 | the | 22,038,615 | 4.9% |
| 2 | be | 12,545,825 | 2.79% |
| 3 | and | 10,741,073 | 2.39% |
| 4 | of | 10,343,885 | 2.3% |
| 5 | a | 10,144,200 | 2.25% |
| 6 | in | 6,996,437 | 1.56% |
| 7 | to (i.m.) | 6,332,195 | 1.41% |
| 8 | have | 4,303,955 | 0.96% |
| 9 | to (p.) | 3,856,916 | 0.86% |
| 10 | it | 3,872,477 | 0.86% |

Top 10 frequent word from the 450 million word corpus
<https://www.wordfrequency.info>.

- We suppress very frequent words that are present in all of the documents and that do not bring any information.

- We suppress very frequent words that are present in all of the documents and that do not bring any information.
- *Example:*
"The cat sat on the mat."
Before filtering: the, cat, sit, on, mat
After filtering: cat, sit, mat

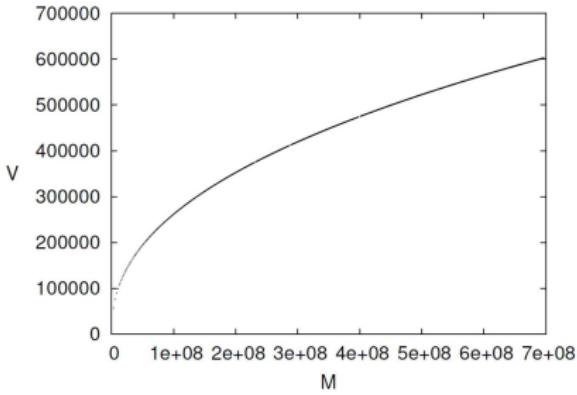
- We suppress very frequent words that are present in all of the documents and that do not bring any information.
- *Example:*
"The cat sat on the mat."
Before filtering: the, cat, sit, on, mat
After filtering: cat, sit, mat
- How many frequent words we should suppress?

The size of the vocabulary V increases sub-linearly with respect to the number of words present in a collection M :

$$V = k \cdot M^\beta,$$

where k and β are parameters that are dependent on the collection. Typically, in English text corpora $k \in [10, 100]$, and $\beta \in [0.4, 0.6]$.

⇒ Larger the collection size, larger the vocabulary size.



You are willing to analyse a document containing 1,000,000 words.

- Let $k = 10$, $\beta = 0.5$. Following the Heaps' law, What is the number of distinct words in the document?

You are willing to analyse a document containing 1,000,000 words.

- Let $k = 10$, $\beta = 0.5$. Following the Heaps' law, What is the number of distinct words in the document?
- It was found that the 7% of all words are "the" article. Following the Zipf's law, estimate the value of λ . What is the frequency of the second most frequent word? The third most frequent one?

- With filtering of stopwords the term frequency weighting may work better.

- With filtering of stopwords the term frequency weighting may work better.
- However, in many applications less frequent words also play a crucial role.

Example: Medical Prescription vs Recipe.

| take | water | glass | eat | wait | ... | paracetamol | sugar | stomach |
|------|-------|-------|-----|------|-----|-------------|-------|---------|
| 7 | 6 | 4 | 4 | 4 | ... | 0 | 2 | 0 |
| 6 | 7 | 4 | 3 | 5 | ... | 1 | 0 | 1 |

- With filtering of stopwords the term frequency weighting may work better.
- However, in many applications less frequent words also play a crucial role.

Example: Medical Prescription vs Recipe.

| take | water | glass | eat | wait | ... | paracetamol | sugar | stomach |
|------|-------|-------|-----|------|-----|-------------|-------|---------|
| 7 | 6 | 4 | 4 | 4 | ... | 0 | 2 | 0 |
| 6 | 7 | 4 | 3 | 5 | ... | 1 | 0 | 1 |

- We want to diminish the weight of terms that occur frequently in general and increase the weight of terms that occur rarely in average.

tf-idf weighting is a trade-off between term frequency and document frequency:

- Normalised term frequency (tf part):

$$\frac{\text{tf}_{t_j,d}}{\sum_{j=1}^V \text{tf}_{t_j,d}} = \frac{\text{tf}_{t_j,d}}{N_d}.$$

tf-idf weighting is a trade-off between term frequency and document frequency:

- Normalised term frequency (tf part):

$$\frac{\text{tf}_{t_j,d}}{\sum_{j=1}^V \text{tf}_{t_j,d}} = \frac{\text{tf}_{t_j,d}}{N_d}.$$

- Inverse document frequency (idf part):

$$\ln \frac{n}{\text{df}_{t_j}} := \ln \frac{n}{\sum_{i=1}^n \mathbb{I}(\text{tf}_{t_j,d_i} \neq 0)}.$$

tf-idf weighting is a trade-off between term frequency and document frequency:

- Normalised term frequency (tf part):

$$\frac{\text{tf}_{t_j,d}}{\sum_{j=1}^V \text{tf}_{t_j,d}} = \frac{\text{tf}_{t_j,d}}{N_d}.$$

- Inverse document frequency (idf part):

$$\ln \frac{n}{\text{df}_{t_j}} := \ln \frac{n}{\sum_{i=1}^n \mathbb{I}(\text{tf}_{t_j,d_i} \neq 0)}.$$

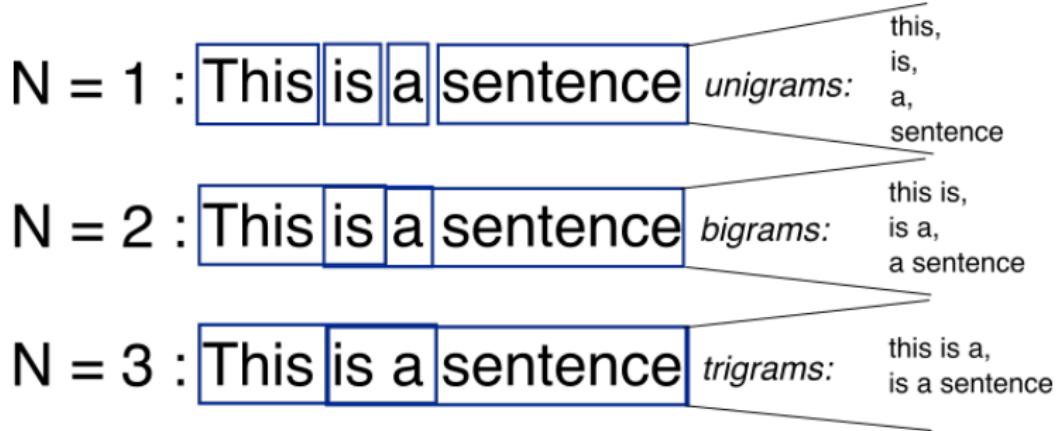
- Then, the tf-idf weight is defined as:

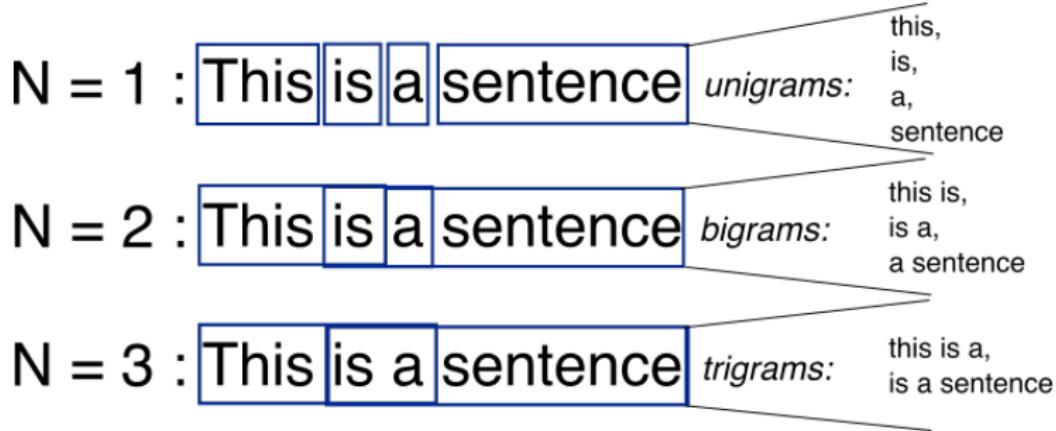
$$w_{t_j,d} = \frac{\text{tf}_{t_j,d}}{N_d} \ln \frac{n}{\text{df}_{t_j}}.$$

- We have the following training documents:

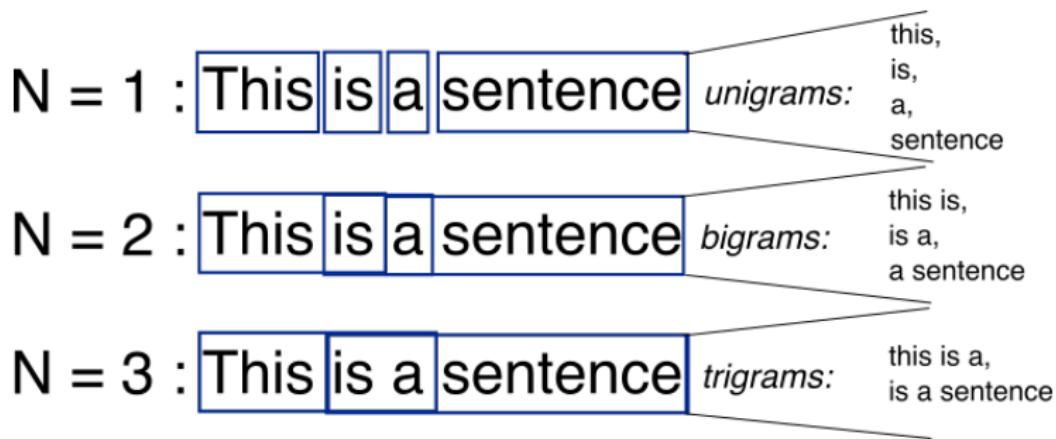
$$d_1 : \{ "cat", "sit", "mat", "cat", "jump", "bed", "cat", "good", "sit" \}$$
$$d_2 : \{ "cat", "dog", "jump", "sit", "dog", "cat", "animal" \}$$
$$d_3 : \{ "table", "sit", "write", "think", "good", "book", "dog", "sit" \}$$

- What is the tf-idf representation for this dataset?





- Several words could be more important when they appear together.



- Several words could be more important when they appear together.
- The approach is very costly when V is large.

... ATTACACGGT**GACC**AACCTATT...

| Gram | Frequency |
|-------------|-----------|
| ATTA | 4 |
| GACC | 3 |
| CGGT | 5 |
| ... | ... |

