Lecture 08
# Machine Learning 2:
*classification (part 1)*

2024-10-09

Sébastien Valade
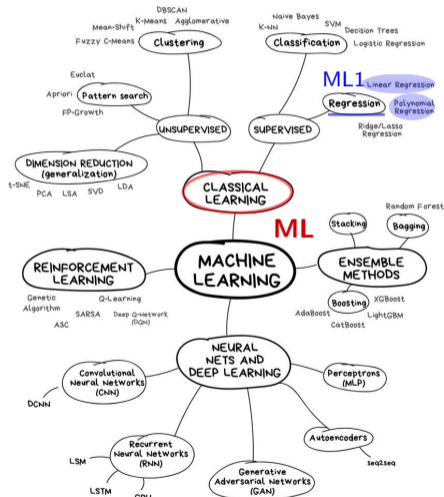
VNIVER₷DAD NAÇIONAL
AVFN°MA DE
MEXIC₽

1. Introduction

2. Probabilistic classification

3. Non-probabilistic classification
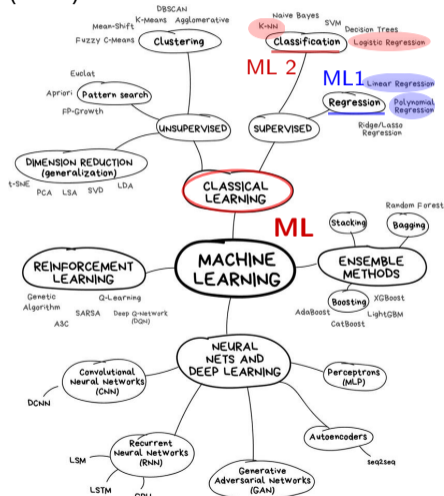
Last week: regression (ML1)

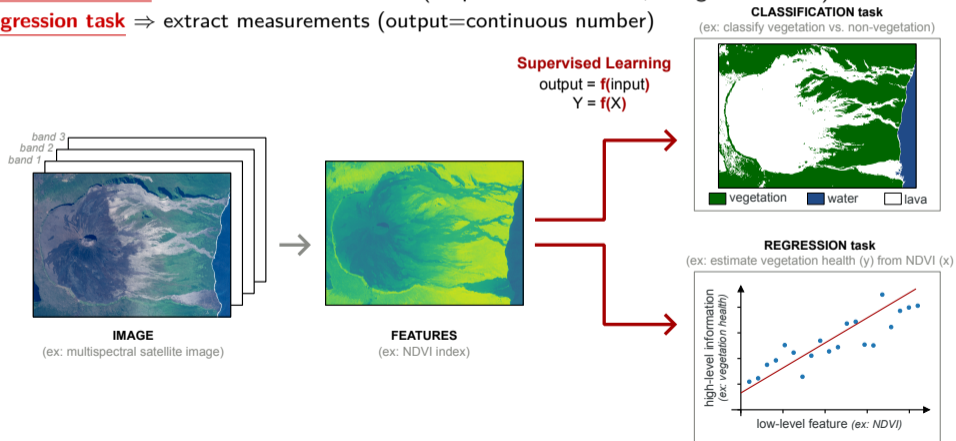Last week: regression (ML1)
This week: classification *part-1* (ML2)



source

## Reminder:

$\Rightarrow$ the goal of **supervised learning** is to learn a function $f$ which maps low-level image features ($X$) to high-level image information ($Y$), using training data (i.e. known pairs of ($X_i$, $Y_i$)):

→ **classification task** $\Rightarrow$ extract semantic classes (output=discrete labels, categorical values)

→ **regression task** $\Rightarrow$ extract measurements (output=continuous number)



**CLASSIFICATION task**
(ex: classify vegetation vs. non-vegetation)

vegetation   water   lava

**REGRESSION task**
(ex: estimate vegetation health (y) from NDVI (x))

high-level image information
*(ex: vegetation health)*

low-level feature *(ex: NDVI)*

**Supervised Learning**
output = **f**(input)
Y = **f**(X)

band 3
band 2
band 1

$\rightarrow$

**IMAGE**
(ex: multispectral satellite image)

**FEATURES**
(ex: NDVI index)

*the term "feature" is here used in a broad sense, referring to any information extracted from the image*
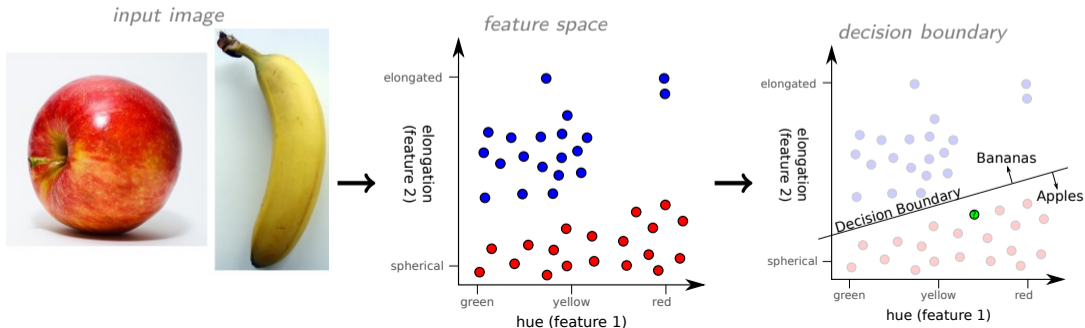
## What is classification?

⇒ the goal of **classification** is to assign **class labels** $Y$ (i.e., discrete categorical values) to data points (e.g., pixels, images)

⇒ extracting **features** from the data is useful to find a space where samples from different classes are well separable (feature crafting can be either manual, or learned from the data using *unsupervised learning*, e.g. PCA)

⇒ the classification algorithm will have to learn the **decision boundary** in an $N$-dimensional **feature space**

**What is classification?**

⇒ the goal of **classification** is to assign **class labels** $Y$ (i.e., discrete categorical values) to data points (e.g., pixels, images)

⇒ extracting **features** from the data is useful to find a space where samples from different classes are well separable (feature crafting can be either manual, or learned from the data using *unsupervised learning*, e.g. PCA)

⇒ the classification algorithm will have to learn the **decision boundary** in an *N*-dimensional **feature space**

**Toy example:** *classify fruit images into either bananas or apples*

### How is the decision boundary learned?

⇒ the **decision boundary** is the surface that separates the feature space into different regions corresponding to different classes

⇒ many algorithms exist to learn this boundary:

- *probabilistic approaches*:
  - **Logistic Regression** ⇒ estimates the probability of a class using a logistic function, fitting a linear decision boundary (binary classification)
  - **Softmax Regression** ⇒ a multi-class extension of logistic regression that assigns probabilities to each class and fits linear boundaries between them
  - **Naive Bayes** ⇒ based on Bayes' theorem, uses probabilistic reasoning to calculate the likelihood of class membership

- *deterministic approaches*:
  - **Perceptron** ⇒ a linear classifier that finds a hyperplane to separate classes, adjusting weights based on misclassified points. Similar to logistic regression but non-probabilistic
  - **k-Nearest Neighbors (kNN)** ⇒ non-parametric method that classifies based on the majority class of the nearest neighbors, leading to non-linear boundaries
  - **Support Vector Machines (SVM)** ⇒ next lectures
  - **Random Forest** ⇒ next lectures
  - **Convolutional Neural Networks (CNNs)** ⇒ next lectures

**How is the decision boundary learned?**

⇒ the **decision boundary** is the surface that separates the feature space into different regions corresponding to different classes

⇒ many algorithms exist to learn this boundary:

- *probabilistic approaches*:
  - **Logistic Regression** ⇒ *estimates the probability of a class using a logistic function, fitting a linear decision boundary (binary classification)*
  - **Softmax Regression** ⇒ *a multi-class extension of logistic regression that assigns probabilities to each class and fits linear boundaries between them*
  - **Naive Bayes** ⇒ *based on Bayes' theorem, uses probabilistic reasoning to calculate the likelihood of class membership*

- *deterministic approaches*:
  - **Perceptron** ⇒ *a linear classifier that finds a hyperplane to separate classes, adjusting weights based on misclassified points. Similar to logistic regression but non-probabilistic*
  - **k-Nearest Neighbors (kNN)** ⇒ *non-parametric method that classifies based on the majority class of the nearest neighbors, leading to non-linear boundaries*
  - **Support Vector Machines (SVM)** ⇒ *next lectures*
  - **Random Forest** ⇒ *next lectures*
  - **Convolutional Neural Networks (CNNs)** ⇒ *next lectures*

**How is the decision boundary learned?**

⇒ the **decision boundary** is the surface that separates the feature space into different regions corresponding to different classes

⇒ many algorithms exist to learn this boundary:

- *probabilistic approaches*:
  - **Logistic Regression** ⇒ *estimates the probability of a class using a logistic function, fitting a linear decision boundary (binary classification)*
  - **Softmax Regression** ⇒ *a multi-class extension of logistic regression that assigns probabilities to each class and fits linear boundaries between them*
  - **Naive Bayes** ⇒ *based on Bayes' theorem, uses probabilistic reasoning to calculate the likelihood of class membership*

- *deterministic approaches*:
  - **Perceptron** ⇒ *a linear classifier that finds a hyperplane to separate classes, adjusting weights based on misclassified points. Similar to logistic regression but non-probabilistic*
  - **k-Nearest Neighbors (kNN)** ⇒ *non-parametric method that classifies based on the majority class of the nearest neighbors, leading to non-linear boundaries*
  - **Support Vector Machines (SVM)** ⇒ *next lectures*
  - **Random Forest** ⇒ *next lectures*
  - **Convolutional Neural Networks (CNNs)** ⇒ *next lectures*

## Logistic Regression

- Linear Regression (recap)
  - ⇒ used to predict continuous values of Y given X
  - ⇒ models the relationship between X and Y as a *linear equation*:
    $Y = \beta_0 + \beta_1 X$
  - ⇒ best model parameters $(\beta_0, \beta_1)$ are found by
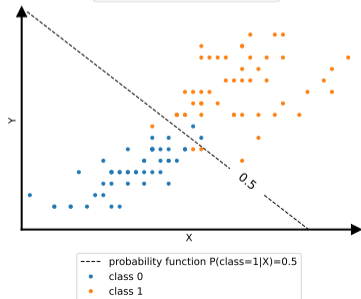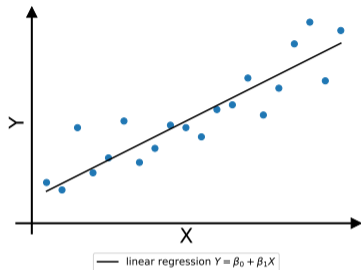    *minimizing Mean Squared Error (MSE)*



— linear regression $Y = \beta_0 + \beta_1 X$

- Logistic Regression
  - ⇒ used to predict binary class values (discrete categorical values
    $y \in 0, 1$) of a data point, given features (X, Y)
  - ⇒ models a probability function using the *logistic function*:
    $p(y = 1|X) = \frac{1}{1+e^{-(\beta_0 + \beta_1 X)}}$
  - ⇒ best model parameters $(\beta_0, \beta_1)$ are found by
    *maximizing Likelihood*



----- probability function P(class=1|X)=0.5
• class 0
• class 1

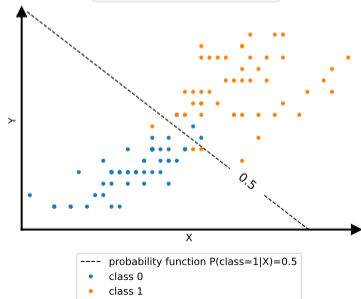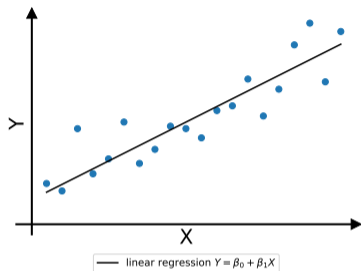## Logistic Regression

- Linear Regression (recap)
  - ⇒ used to predict continuous values of Y given X
  - ⇒ models the relationship between X and Y as a *linear equation*:
    $Y = \beta_0 + \beta_1 X$
  - ⇒ best model parameters $(\beta_0, \beta_1)$ are found by
    *minimizing Mean Squared Error (MSE)*

- Logistic Regression
  - ⇒ used to predict binary class values (discrete categorical values
    $y \in 0, 1$) of a data point, given features (X, Y)
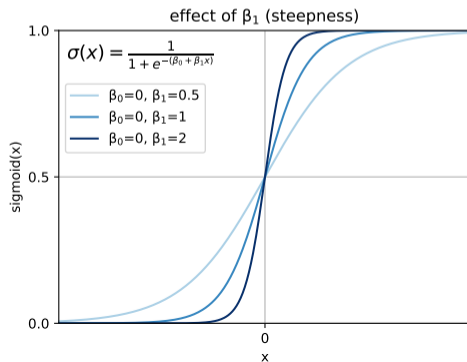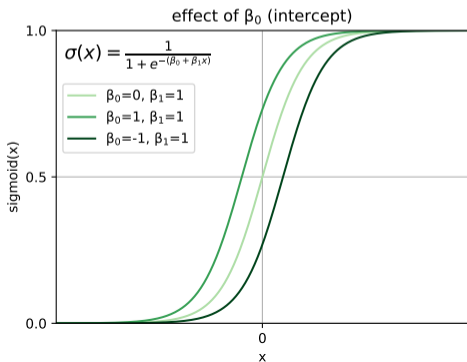  - ⇒ models a probability function using the *logistic function*:
    $p(y = 1|X) = \frac{1}{1+e^{-(\beta_0+\beta_1 X)}}$
  - ⇒ best model parameters $(\beta_0, \beta_1)$ are found by
    *maximizing Likelihood*



—— linear regression $Y = \beta_0 + \beta_1 X$



----- probability function P(class=1|X)=0.5
· class 0
· class 1

## Logistic Regression

$\Rightarrow$ in order to model the binary class probabilities, we use the **logistic function** (a.k.a. sigmoid function, S-shaped curve), which maps any real value of feature $X$ to the range [0,1]:

$$\sigma(X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X)}} = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$



effect of $\beta_0$ (intercept)

$\sigma(x) = \frac{1}{1+e^{-(\beta_0 + \beta_1 x)}}$

$\beta_0=0, \beta_1=1$
$\beta_0=1, \beta_1=1$
$\beta_0=-1, \beta_1=1$

effect of $\beta_1$ (steepness)

$\sigma(x) = \frac{1}{1+e^{-(\beta_0 + \beta_1 x)}}$

$\beta_0=0, \beta_1=0.5$
$\beta_0=0, \beta_1=1$
$\beta_0=0, \beta_1=2$

## Logistic Regression

$\Rightarrow$ in order to model the binary class probabilities, we use the **logistic function** (a.k.a. sigmoid function, S-shaped curve), which maps any real value of feature $X$ to the range [0,1]:

$$\sigma(X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X)}} = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}} = \frac{1}{1 + e^{-\theta^\top X}}$$

with: $\quad \theta = [\beta_0, \beta_1] \quad$ and $\quad X = [1, X]$

$\Rightarrow$ the probability $p$ is then calculated as:
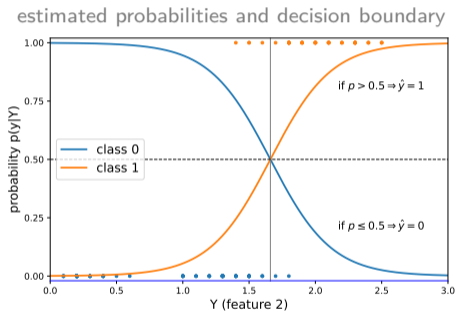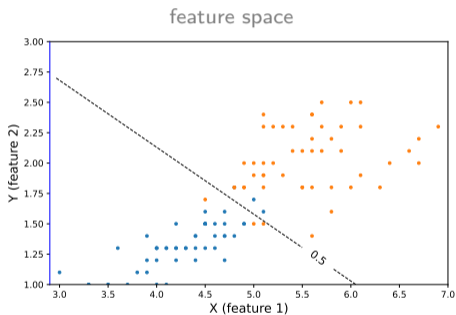
$$p(y = 1|X) = \sigma(X)$$
$$p(y = 0|X) = 1 - \sigma(X)$$

$\Rightarrow$ the prediction of the class $\hat{y} \in \{0, 1\}$ is made by comparing the probability $p(X)$ to a threshold (e.g. 0.5):

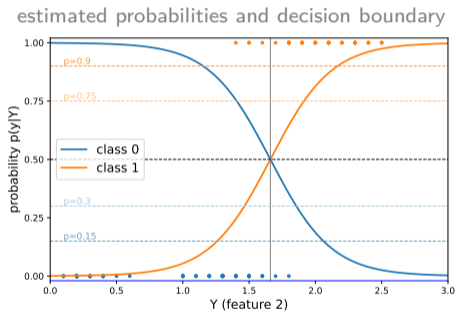$$\hat{y} = \begin{cases} 1 & \text{if } p(X) \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$
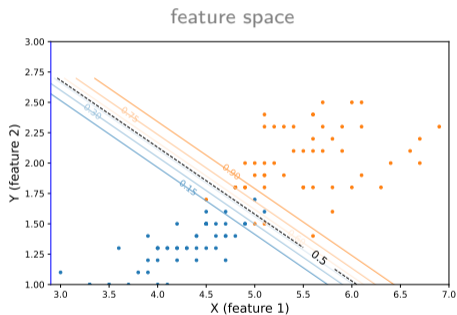
## Logistic Regression

<u>EX</u>: estimate data point class $\hat{y} \in \{0|1\}$ from estimated probabilities $p(y|Y)$



feature space

estimated probabilities and decision boundary

## Logistic Regression

<u>EX</u>: estimate data point class $\hat{y} = \{0|1\}$ from estimated probabilities $p(y|Y)$



$\longleftrightarrow$

## Logistic Regression

$\Rightarrow$ the coefficients $\beta_0$ and $\beta_1$ are unknown, and must be estimated based on the available training data

$\Rightarrow$ the coefficients are estimated by using the **maximum likelihood function**

  $\rightarrow$ find best estimates of $\beta_0$ and $\beta_1$, such that the predicted probability $\hat{p}(x_i)$ for each data point returns as closely as possible to the expected class

  $\rightarrow$ this can be formalized using a mathematical equation called a likelihood function:

$$L(\beta_0, \beta_1) = \prod_{i:y_i=1} p(x_i) \prod_{i:y_i=0} (1 - p(x_i))$$

  $\rightarrow$ the maximization of the likelihood function allows for the estimation of the best parameters $\beta_0$ and $\beta_1$:

$$\hat{\beta}_0, \hat{\beta}_1 = \arg \max_{\beta_0, \beta_1} L(\beta_0, \beta_1)$$

  NB: we could also use least squares to fit the model, but the maximum likelihood is preferred because it has better statistical properties

## Logistic Regression

$\Rightarrow$ the coefficients $\beta_0$ and $\beta_1$ are unknown, and must be estimated based on the available training data

$\Rightarrow$ the coefficients are estimated by using the **maximum likelihood function**

$\rightarrow$ *find best estimates of $\beta_0$ and $\beta_1$, such that the predicted probability $\hat{p}(x_i)$ for each data point returns as closely as possible to the expected class*

$\rightarrow$ *this can be formalized using a mathematical equation called a likelihood function:*

$$L(\beta_0, \beta_1) = \prod_{i:y_i=1} p(x_i) \prod_{i:y_i=0} (1 - p(x_i))$$

$\rightarrow$ *the maximization of the likelihood function allows for the estimation of the best parameters $\beta_0$ and $\beta_1$:*

$$\hat{\beta}_0, \hat{\beta}_1 = \arg\max_{\beta_0, \beta_1} L(\beta_0, \beta_1)$$

*NB: we could also use least squares to fit the model, but the maximum likelihood is preferred because it has better statistical properties*

## Logistic Regression

$\Rightarrow$ the coefficients $\beta_0$ and $\beta_1$ are unknown, and must be estimated based on the available training data

$\Rightarrow$ the coefficients are estimated by using the **maximum likelihood function**

  $\rightarrow$ *find best estimates of $\beta_0$ and $\beta_1$, such that the predicted probability $\hat{p}(x_i)$ for each data point returns as closely as possible to the expected class*

  $\rightarrow$ *this can be formalized using a mathematical equation called a likelihood function:*

$$L(\beta_0, \beta_1) = \prod_{i:y_i=1} p(x_i) \prod_{i:y_i=0} (1 - p(x_i))$$

  $\rightarrow$ *the maximization of the likelihood function allows for the estimation of the best parameters $\beta_0$ and $\beta_1$:*

$$\hat{\beta}_0, \hat{\beta}_1 = \arg\max_{\beta_0, \beta_1} L(\beta_0, \beta_1)$$

  <u>NB</u>: we could also use least squares to fit the model, but the maximum likelihood is preferred because it has better statistical properties

## Logistic Regression

$\Rightarrow$ the coefficients $\beta_0$ and $\beta_1$ are unknown, and must be estimated based on the available training data

$\Rightarrow$ the coefficients are estimated by using the **maximum likelihood function**

$\rightarrow$ *find best estimates of $\beta_0$ and $\beta_1$, such that the predicted probability $\hat{p}(x_i)$ for each data point returns as closely as possible to the expected class*

$\rightarrow$ *this can be formalized using a mathematical equation called a likelihood function:*

$$L(\beta_0, \beta_1) = \prod_{i:y_i=1} p(x_i) \prod_{i:y_i=0} (1 - p(x_i))$$

$\rightarrow$ *the maximization of the likelihood function allows for the estimation of the best parameters $\beta_0$ and $\beta_1$:*

$$\hat{\beta}_0, \hat{\beta}_1 = \arg \max_{\beta_0, \beta_1} L(\beta_0, \beta_1)$$

*<u>NB</u>: we could also use least squares to fit the model, but the maximum likelihood is preferred because it has better statistical properties*

**Logistic Regression**

⇒ Logistic Regression can also be used to to predict a <u>binary response but using multiple predictors</u>
(in the previous slides, the binary class prediction was done using just the 1 feature $Y$)

⇒ if we consider $k$ predictors, the model is then defined as:

$$p(X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k)}}$$
$$= \frac{e^{\beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k}}{1 + e^{\beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k}}$$

<u>NB</u>: the classification results obtained using one predictor or several may be quite different, especially when there is correlation among the predictors!

⇒ nevertheless, **Logistic Regression** is limited as it can only model **binary classification**!
⇒ for **multi-class classification**, we need to use **Softmax Regression**

**Logistic Regression**

⇒ Logistic Regression can also be used to to predict a binary response but using multiple predictors
(in the previous slides, the binary class prediction was done using just the 1 feature $Y$)

⇒ if we consider $k$ predictors, the model is then defined as:

$$p(X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k)}}$$

$$= \frac{e^{\beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k}}{1 + e^{\beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k}}$$

NB: the classification results obtained using one predictor or several may be quite different, especially when there is correlation among the predictors!

⇒ nevertheless, **Logistic Regression** is limited as it can only model **binary classification**!
⇒ for **multi-class classification**, we need to use **Softmax Regression**

## Softmax Regression

$\Rightarrow$ the **Softmax Regression** (a.k.a. *Multinomial Logistic Regression*) generalizes the logistic regression for multiple classes, by calculating probabilities $p(y)$ for each class $y \in \{1, ..., K\}$

$$p(y = k|X; \theta) = \frac{\exp(\theta^{(k)\top} X)}{\sum_{j=1}^{K} \exp(\theta^{(j)\top} X)}$$

where $\begin{cases} \text{the numerator gives the exponentiated score for each class k} \\ \text{the denominator normalizes the scores into a valid probability distribution} \end{cases}$

In other words:

$$\begin{bmatrix} p(y=1|X;\theta) \\ p(y=2|X;\theta) \\ \vdots \\ p(y=K|X;\theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^{K} \exp(\theta^{(j)\top} X)} \begin{bmatrix} exp(\theta^{(1)\top} X) \\ exp(\theta^{(2)\top} X) \\ \vdots \\ exp(\theta^{(K)\top} X) \end{bmatrix}$$

where $\theta^{(1)}, \theta^{(2)}, \ldots, \theta^{(K)}$ are the model parameters

$\Rightarrow$ like with the *Logistic Regression classifier*, the *Softmax Regression classifier* predicts the class with the highest estimated probability:

$$\hat{y} = \arg \max_{k} p(y = k|X; \theta)$$

## Softmax Regression

⇒ the **Softmax Regression** (a.k.a. *Multinomial Logistic Regression*) generalizes the logistic regression for multiple classes, by calculating probabilities $p(y)$ for each class $y \in \{1, ..., K\}$

$$p(y = k|X; \theta) = \frac{\exp(\theta^{(k)\top} X)}{\sum_{j=1}^{K} \exp(\theta^{(j)\top} X)}$$

where $\begin{cases} \text{the numerator gives the exponentiated score for each class k} \\ \text{the denominator normalizes the scores into a valid probability distribution} \end{cases}$

In other words:

$$\begin{bmatrix} p(y = 1|X; \theta) \\ p(y = 2|X; \theta) \\ \vdots \\ p(y = K|X; \theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^{K} \exp(\theta^{(j)\top} X)} \begin{bmatrix} exp(\theta^{(1)\top} X) \\ exp(\theta^{(2)\top} X) \\ \vdots \\ exp(\theta^{(K)\top} X) \end{bmatrix}$$
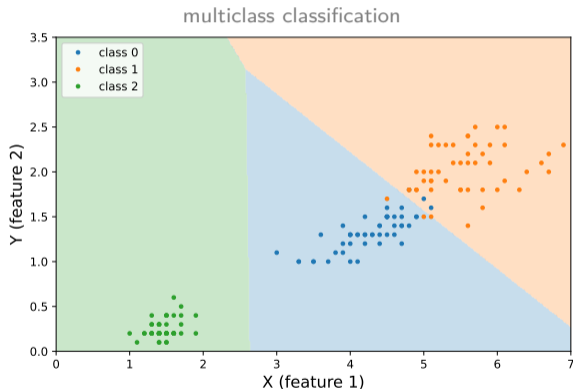
where $\theta^{(1)}, \theta^{(2)}, \ldots, \theta^{(K)}$ are the model parameters

⇒ like with the *Logistic Regression classifier*, the *Softmax Regression classifier* predicts the class with the highest estimated probability:

$$\hat{y} = \arg \max_{k} p(y = k|X; \theta)$$

## Softmax Regression

$\Rightarrow$ the **Softmax Regression** (a.k.a. *Multinomial Logistic Regression*) generalizes the logistic regression for multiple classes, by calculating probabilities $p(y)$ for each class $y \in \{1, ..., K\}$

$$p(y = k|X; \theta) = \frac{\exp(\theta^{(k)\top} X)}{\sum_{j=1}^{K} \exp(\theta^{(j)\top} X)}$$

where $\begin{cases} \text{the numerator gives the exponentiated score for each class k} \\ \text{the denominator normalizes the scores into a valid probability distribution} \end{cases}$

In other words:

$$\begin{bmatrix} p(y = 1|X; \theta) \\ p(y = 2|X; \theta) \\ \vdots \\ p(y = K|X; \theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^{K} \exp(\theta^{(j)\top} X)} \begin{bmatrix} exp(\theta^{(1)\top} X) \\ exp(\theta^{(2)\top} X) \\ \vdots \\ exp(\theta^{(K)\top} X) \end{bmatrix}$$

where $\theta^{(1)}, \theta^{(2)}, \ldots, \theta^{(K)}$ are the model parameters

$\Rightarrow$ like with the *Logistic Regression classifier*, the *Softmax Regression classifier* predicts the class with the highest estimated probability:

$$\hat{y} = \arg \max_k p(y = k|X; \theta)$$

## Softmax Regression

EX: estimate data point class $\hat{y} \in \{0, 1, 2\}$



multiclass classification

⇒ **Softmax Regression** finds linear boundaries between multiple classes
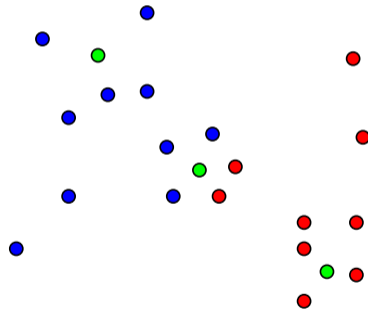→ it is commonly used as the final layer in **neural networks** for classification tasks

## kNN classification

⇒ Idea: **classify** a data point based on the majority class of its *k* **Nearest Neighbors**

⇒ Method:

1. take random data points in the training dataset

2. for a sample find the *k* (e.g. 5) closest data points in the training dataset (*k* is a *hyperparameter*)

3. look at the neighbor labels, return/assign the mode
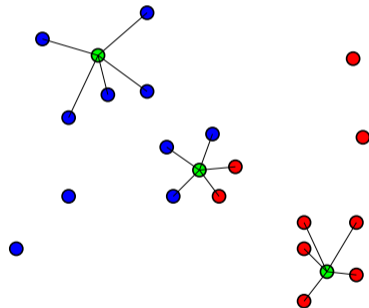
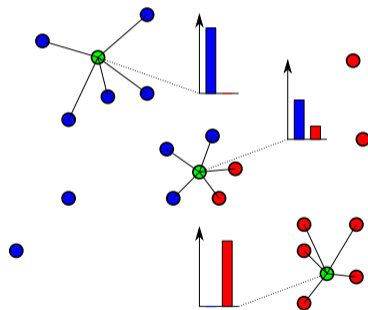4. non-linear decision boundary can be recovered
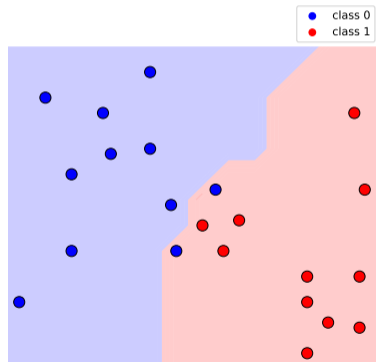
## kNN classification

⇒ Idea: **classify** a data point based on the majority class of its *k* **Nearest Neighbors**

⇒ Method:

1. take random data points in the training dataset

2. for a sample find the *k* (e.g. 5) closest data points in the training dataset *(k is a hyperparameter)*

3. look at the neighbor labels, return/assign the mode

4. non-linear decision boundary can be recovered

**kNN classification**

$\Rightarrow$ Idea: **classify** a data point based on the majority class of its $k$ **Nearest Neighbors**

$\Rightarrow$ Method:

1. take random data points in the training dataset
2. for a sample find the $k$ (e.g. 5) closest data points in the training dataset *(k is a hyperparameter)*
3. look at the neighbor labels, return/assign the mode
4. non-linear decision boundary can be recovered

**kNN classification**

$\Rightarrow$ Idea: **classify** a data point based on the majority class of its $k$ **Nearest Neighbors**
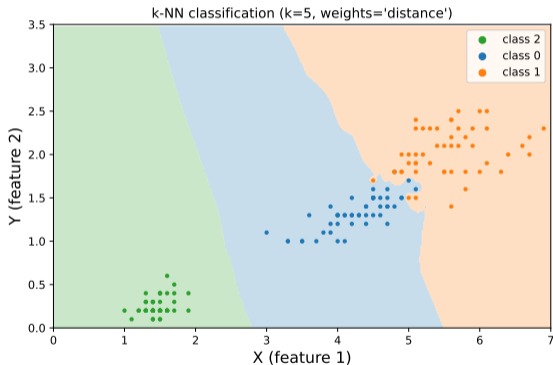
$\Rightarrow$ Method:

1. take random data points in the training dataset
2. for a sample find the $k$ (e.g. 5) closest data points in the training dataset *(k is a hyperparameter)*
3. look at the neighbor labels, return/assign the mode
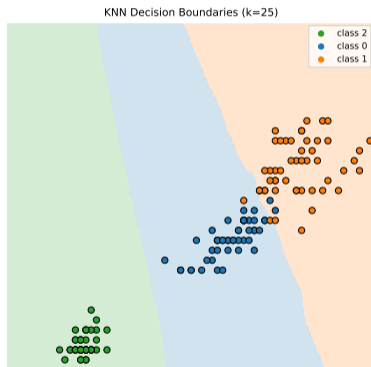4. non-linear decision boundary can be recovered

## kNN classification

⇒ Idea: **classify** a data point based on the majority class of its *k* **Nearest Neighbors**

⇒ Method:



1. take random data points in the training dataset
2. for a sample find the *k* (e.g. 5) closest data points in the training dataset *(k is a hyperparameter)*
3. look at the neighbor labels, return/assign the mode
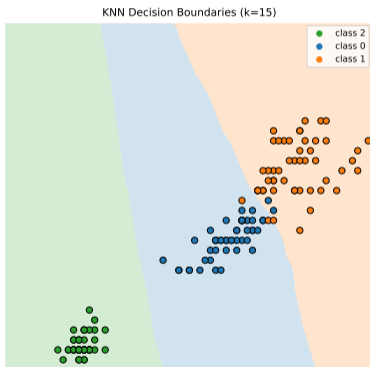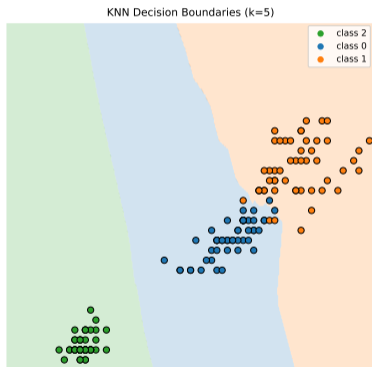4. non-linear decision boundary can be recovered

### **kNN classification**

<u>EX</u>: estimate classification boundary using the kNN algorithm



k-NN classification (k=5, weights='distance')

## kNN classification

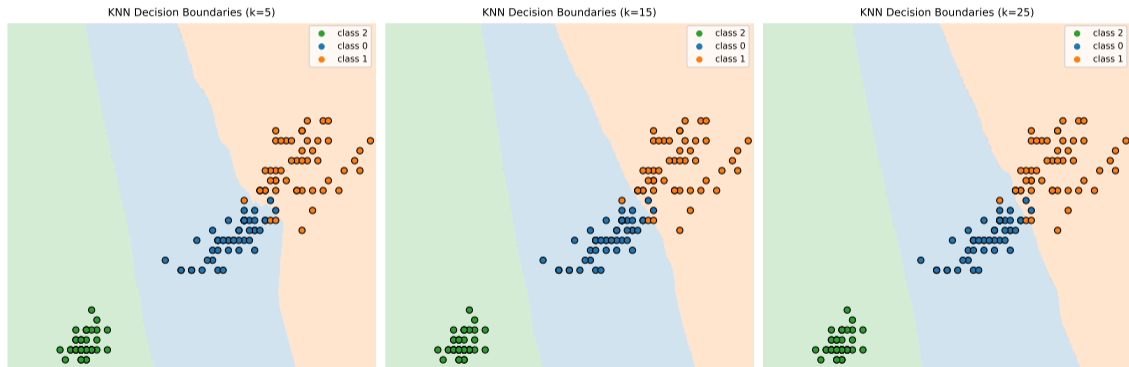EX: estimate classification boundary using the kNN algorithm

## kNN classification

EX: estimate classification boundary using the kNN algorithm



---

→ Pros: simple, easy to understand, works well with small datasets
→ Cons: slow for large datasets, sensitive to choice of distance metric and the value of k

**more on classification next week !**

(perceptron, SVM, PCA)