

# Machine Learning for Official Statistics and SDGs

## Classification



# [ CLASSIFICATION ]

What is a classification problem?

# [ CLASSIFICATION ]

What is a classification problem?

- ▶ The goal is to understand why an observation belongs to a certain category

# [ CLASSIFICATION ]

What is a classification problem?

- ▶ The goal is to understand why an observation belongs to a certain category
- ▶  $y$  takes discrete values: 0/1, high school/primary school/no education; urban/rural

# [ CLASSIFICATION ]

What is a classification problem?

- ▶ The goal is to understand why an observation belongs to a certain category
- ▶  $y$  takes discrete values: 0/1, high school/primary school/no education; urban/rural
- ▶ Some variables  $x$ s may explain why  $y$  belongs to a particular category

# [ CLASSIFICATION ]

What is a classification problem?

- ▶ The goal is to understand why an observation belongs to a certain category
- ▶  $y$  takes discrete values: 0/1, high school/primary school/no education; urban/rural
- ▶ Some variables  $x$ s may explain why  $y$  belongs to a particular category

# [ CLASSIFICATION ]

What is a classification problem?

- ▶ The goal is to understand why an observation belongs to a certain category
- ▶  $y$  takes discrete values: 0/1, high school/primary school/no education; urban/rural
- ▶ Some variables  $x$ s may explain why  $y$  belongs to a particular category

A **classifier** is a tool that provides a classification for  $y$  using (*or not*) additional information from other variables

# [ SUPERVISED *vs* UNSUPERVISED CLASSIFICATION ]



# [ SUPERVISED *vs* UNSUPERVISED CLASSIFICATION ]

- ▶ In **supervised** classification, we **observe** the category for each observation

# [ SUPERVISED *vs* UNSUPERVISED CLASSIFICATION ]

- In **supervised** classification, we **observe** the category for each observation

*One may learn and estimate the impact of other variables on that classification (e.g. logit regression)*

## [ SUPERVISED *vs* UNSUPERVISED CLASSIFICATION ]

- ▶ In **supervised** classification, we **observe** the category for each observation  
*One may learn and estimate the impact of other variables on that classification (e.g. logit regression)*
- ▶ In **unsupervised** classification, we **ignore** the category (if any) of each observation

## [ SUPERVISED *vs* UNSUPERVISED CLASSIFICATION ]

- ▶ In **supervised** classification, we **observe** the category for each observation

*One may learn and estimate the impact of other variables on that classification (e.g. logit regression)*

- ▶ In **unsupervised** classification, we **ignore** the category (if any) of each observation

*The goal is to classify observations from those variables (clustering) without having any information of what a category means.*

## [ SUPERVISED *vs* UNSUPERVISED CLASSIFICATION ]

- ▶ In **supervised** classification, we **observe** the category for each observation

*One may learn and estimate the impact of other variables on that classification (e.g. logit regression)*

- ▶ In **unsupervised** classification, we **ignore** the category (if any) of each observation

*The goal is to classify observations from those variables (clustering) without having any information of what a category means.*

## [ SUPERVISED *vs* UNSUPERVISED CLASSIFICATION ]

- ▶ In **supervised** classification, we **observe** the category for each observation

*One may learn and estimate the impact of other variables on that classification (e.g. logit regression)*

- ▶ In **unsupervised** classification, we **ignore** the category (if any) of each observation

*The goal is to classify observations from those variables (clustering) without having any information of what a category means.*

- ▶ We'll focus on **supervised** classification

# [ CLASSIFICATION: AN EXAMPLE ]

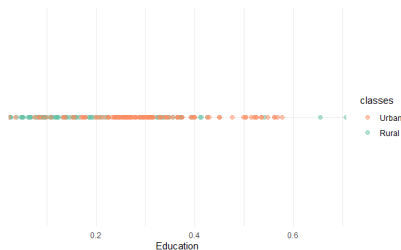
## [ CLASSIFICATION: AN EXAMPLE ]

- ▶ You observe households that are either in *Urban* or *Rural* areas (colors) and one variable (feature): *Education*.



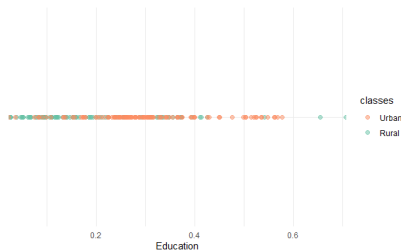
# [ CLASSIFICATION: AN EXAMPLE ]

- You observe households that are either in *Urban* or *Rural* areas (colors) and one variable (feature): *Education*.



# [ CLASSIFICATION: AN EXAMPLE ]

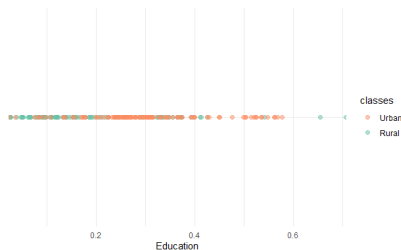
- You observe households that are either in *Urban* or *Rural* areas (colors) and one variable (feature): *Education*.



- A classifier determines the value of *Education* that separate "*Rural*" from "*Urban*"

# [ CLASSIFICATION: AN EXAMPLE ]

- You observe households that are either in *Urban* or *Rural* areas (colors) and one variable (feature): *Education*.



- A classifier determines the value of *Education* that separate "*Rural*" from "*Urban*"

Typically with a threshold rule:

"if  $x \geq t$  then category is *Urban*"

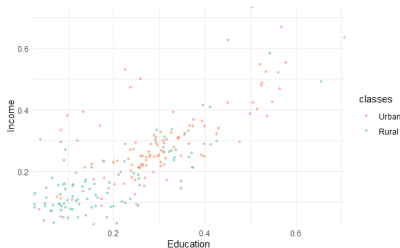
# [ CLASSIFICATION: A 2-D EXAMPLE ]

## [ CLASSIFICATION: A 2-D EXAMPLE ]

- ▶ You observe households in *Urban* or *Rural* areas and **two** variables (features): *Education* and *Income*

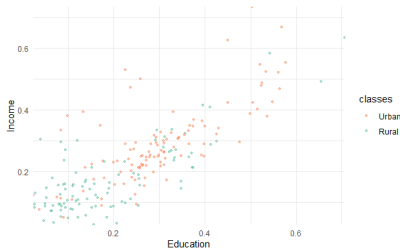
# [ CLASSIFICATION: A 2-D EXAMPLE ]

- You observe households in *Urban* or *Rural* areas and **two** variables (features): *Education* and *Income*



# [ CLASSIFICATION: A 2-D EXAMPLE ]

- You observe households in *Urban* or *Rural* areas and **two** variables (features): *Education* and *Income*



# [ CLASSIFICATION: A 2-D EXAMPLE ]

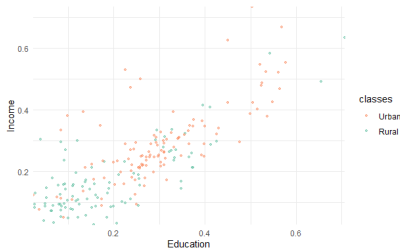


## [ CLASSIFICATION: A 2-D EXAMPLE ]

- ▶ A classifier will determine a **boundary** using both *Education* and *Income* to separate "Rural" from "Urban"

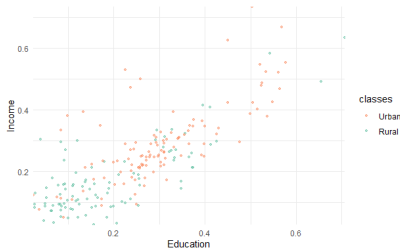
# [ CLASSIFICATION: A 2-D EXAMPLE ]

- A classifier will determine a **boundary** using both *Education* and *Income* to separate "Rural" from "Urban"



## [ CLASSIFICATION: A 2-D EXAMPLE ]

- A classifier will determine a **boundary** using both *Education* and *Income* to separate "Rural" from "Urban"



- The rule can be based on a linear relationship between *Education* and *Income* or can be non linear.

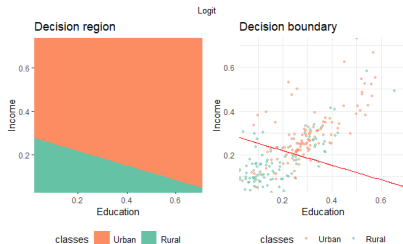
# [ CLASSIFICATION: A 2-D EXAMPLE ]

## [ CLASSIFICATION: A 2-D EXAMPLE ]

- ▶ Example of a linear classifier

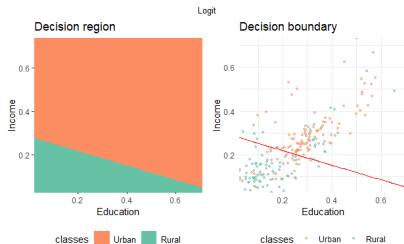
# [ CLASSIFICATION: A 2-D EXAMPLE ]

## ► Example of a linear classifier



# [ CLASSIFICATION: A 2-D EXAMPLE ]

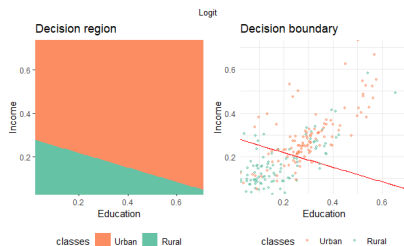
## ► Example of a linear classifier



- The separation rule is  $x'\beta \geq T_0$   
with  $T_0$  a known threshold

# [ CLASSIFICATION: A 2-D EXAMPLE ]

## ► Example of a linear classifier



- The separation rule is  $x'\beta \geq T_0$   
with  $T_0$  a known threshold  
*e.g.*  $\beta_0 + \beta_1 \text{Education} + \beta_2 \text{Income} \geq T_0 \Leftrightarrow \text{Urban}$



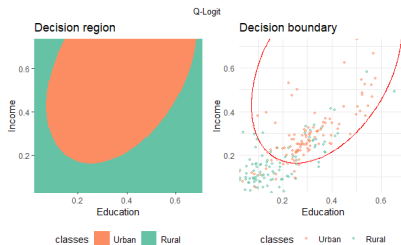
# [ CLASSIFIERS EXAMPLES]

## [ CLASSIFIERS EXAMPLES]

- ▶ Example of non-linear classifier

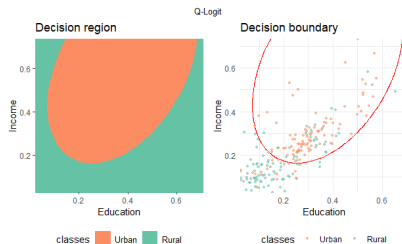
# [ CLASSIFIERS EXAMPLES ]

## ► Example of non-linear classifier



# [ CLASSIFIERS EXAMPLES ]

## ► Example of non-linear classifier



## ► The rule that separated the two classes is non linear in the variables *Education* and *Income*

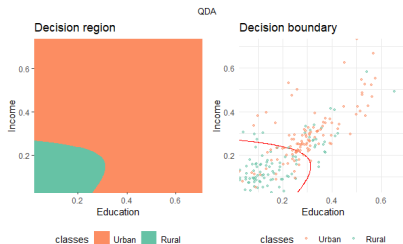
# [ CLASSIFIERS EXAMPLES ]

## [ CLASSIFIERS EXAMPLES ]

- ▶ Another non-linear example

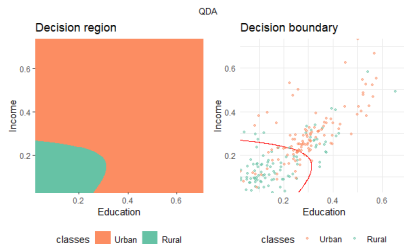
# [ CLASSIFIERS EXAMPLES ]

## ► Another non-linear example



# [ CLASSIFIERS EXAMPLES ]

## ► Another non-linear example



- The boundary is complex and uses *Education* and *Income* features.



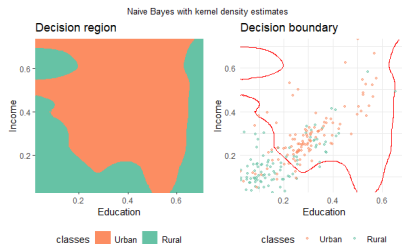
# [ CLASSIFIERS EXAMPLES ]

## [ CLASSIFIERS EXAMPLES ]

- ▶ Other examples can be very non linear

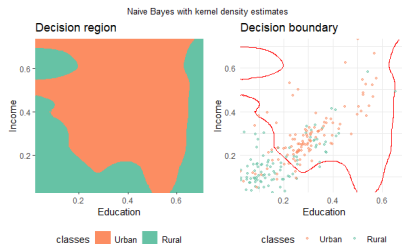
# [ CLASSIFIERS EXAMPLES ]

- Other examples can be very non linear



# [ CLASSIFIERS EXAMPLES ]

- Other examples can be very non linear



- It is hard to understand how the two classes are built using *Education* and *Income*

# [ HOW TO SELECT THE RIGHT MODEL? ]

# [ HOW TO SELECT THE RIGHT MODEL? ]

- What is the goal?

# [ HOW TO SELECT THE RIGHT MODEL? ]

- ▶ What is the goal?  
Have the “best” classification

# [ HOW TO SELECT THE RIGHT MODEL? ]

- What is the goal?

Have the “best” classification

↪ Need for a criterion to determine what is a good classifier



# [ HOW TO SELECT THE RIGHT MODEL? ]

- ▶ What is the goal?  
Have the “best” classification
- ↔ Need for a criterion to determine what is a good classifier
- ▶ Measures of fit in classification are different and specific

# [ MEASURES OF FIT IN CLASSIFICATION ]

There are several popular measures of fit, differing in their spirit and their goal

# [ MEASURES OF FIT IN CLASSIFICATION ]

There are several popular measures of fit, differing in their spirit and their goal

- ▶ Accuracy

# [ MEASURES OF FIT IN CLASSIFICATION ]

There are several popular measures of fit, differing in their spirit and their goal

- ▶ Accuracy
- ▶ Confusion matrix

# [ MEASURES OF FIT IN CLASSIFICATION ]

There are several popular measures of fit, differing in their spirit and their goal

- ▶ Accuracy
- ▶ Confusion matrix
- ▶ Sensitivity & Specificity

## [ MEASURES OF FIT IN CLASSIFICATION ]

There are several popular measures of fit, differing in their spirit and their goal

- ▶ Accuracy
- ▶ Confusion matrix
- ▶ Sensitivity & Specificity
- ▶ Kappa

# [ MEASURES OF FIT IN CLASSIFICATION ]

There are several popular measures of fit, differing in their spirit and their goal

- ▶ Accuracy
- ▶ Confusion matrix
- ▶ Sensitivity & Specificity
- ▶ Kappa
- ...

## [ MEASURES OF FIT IN CLASSIFICATION ]

There are several popular measures of fit, differing in their spirit and their goal

- ▶ Accuracy
- ▶ Confusion matrix
- ▶ Sensitivity & Specificity
- ▶ Kappa

...

Each criterion answers to a different question



# [ ACCURACY AND CONFUSION MATRIX]

Accuracy corresponds to the probability of being "accurate"

$$\Pr \left[ y_0 = \hat{f}(x_0) \right]$$

# [ ACCURACY AND CONFUSION MATRIX]

Accuracy corresponds to the probability of being "accurate"

$$\Pr \left[ y_0 = \hat{f}(x_0) \right]$$

- where  $\hat{f}(\cdot)$  is the classifier.

# [ ACCURACY AND CONFUSION MATRIX]

Accuracy corresponds to the probability of being "accurate"

$$\Pr \left[ y_0 = \hat{f}(x_0) \right]$$

- where  $\hat{f}(\cdot)$  is the classifier.
- ↪ We want the maximum possible accuracy.

# [ ACCURACY AND CONFUSION MATRIX]

Accuracy corresponds to the probability of being "accurate"

$$\Pr \left[ y_0 = \hat{f}(x_0) \right]$$

- ▶ where  $\hat{f}(\cdot)$  is the classifier.
- ↪ We want the maximum possible accuracy.
- ▶ Equivalently, we may want to minimize the *error rate* or *misclassification rate*

$$\Pr \left[ y_0 \neq \hat{f}(x_0) \right]$$

## [ CONFUSION MATRIX & ACCURACY ]

A classifier predicts in which class each observation should be:

## [ CONFUSION MATRIX & ACCURACY ]

A classifier predicts in which class each observation should be:

Predicted	Observed (True)	
	TP (True Positive)	FP (False Positive)
	FN (False Negative)	TN (True Negative)

Table: Confusion Matrix

## [ CONFUSION MATRIX & ACCURACY ]

A classifier predicts in which class each observation should be:

	Observed (True)	
	TP (True Positive)	FP (False Positive)
Predicted	FN (False Negative)	TN (True Negative)

Table: Confusion Matrix

- Accuracy is then the ratio:

$$\begin{aligned} \text{Accuracy} &= \frac{TP + TN}{TP + TN + FP + FN} \\ &= \frac{\text{TruePositives} + \text{TrueNegatives}}{N} \end{aligned}$$

## [ CONFUSION MATRIX & ACCURACY ]

A classifier predicts in which class each observation should be:

	Observed (True)	
	TP (True Positive)	FP (False Positive)
Predicted	FN (False Negative)	TN (True Negative)

Table: Confusion Matrix

- Accuracy is then the ratio:

$$\begin{aligned} \text{Accuracy} &= \frac{TP + TN}{TP + TN + FP + FN} \\ &= \frac{\text{TruePositives} + \text{TrueNegatives}}{N} \end{aligned}$$

- It is the proportion of accurate predictions



# [CONFUSION MATRIX & ACCURACY]

In practice, with a classifier we have:

		Observed (True)	
		Urban	Rural
Predicted	Urban	87 (TP)	28 (FP)
	Rural	24 (FN)	69 (TN)

Table: Confusion Matrix

# [CONFUSION MATRIX & ACCURACY]

In practice, with a classifier we have:

		Observed (True)	
		Urban	Rural
Predicted	Urban	87 (TP)	28 (FP)
	Rural	24 (FN)	69 (TN)

Table: Confusion Matrix

- Here *Urban* is the "positive" class

## [CONFUSION MATRIX & ACCURACY]

In practice, with a classifier we have:

		Observed (True)	
		Urban	Rural
Predicted	Urban	87 (TP)	28 (FP)
	Rural	24 (FN)	69 (TN)

Table: Confusion Matrix

- ▶ Here *Urban* is the "positive" class
- ▶ Accuracy is then the ratio:

$$\begin{aligned} \text{Accuracy} &= \frac{87 + 69}{87 + 69 + 28 + 24} \\ &= \frac{156}{208} = 0.75 \end{aligned}$$

## [CONFUSION MATRIX & ACCURACY]

In practice, with a classifier we have:

		Observed (True)	
		Urban	Rural
Predicted	Urban	87 (TP)	28 (FP)
	Rural	24 (FN)	69 (TN)

Table: Confusion Matrix

- ▶ Here *Urban* is the "positive" class
- ▶ Accuracy is then the ratio:

$$\begin{aligned} \text{Accuracy} &= \frac{87 + 69}{87 + 69 + 28 + 24} \\ &= \frac{156}{208} = 0.75 \end{aligned}$$

- ▶ We have an accurate prediction in 75% of the cases.

## [PROBLEM 1: ACCURACY IS ONE NUMBER ]

Accuracy is not the panacea and may be misleading

- ▶ One may be more interested in **correctly** predicting a particular outcome!

## [PROBLEM 1: ACCURACY IS ONE NUMBER ]

Accuracy is not the panacea and may be misleading

- ▶ One may be more interested in **correctly** predicting a particular outcome!
- ↪ This is often the case if the **cost** of being wrong differ

## [PROBLEM 1: ACCURACY IS ONE NUMBER ]

Accuracy is not the panacea and may be misleading

- ▶ One may be more interested in **correctly** predicting a particular outcome!
- ↪ This is often the case if the **cost** of being wrong differ
- ▶ One may need other measures, focused on one particular outcome

## [PROBLEM 1: ACCURACY IS ONE NUMBER ]

Accuracy is not the panacea and may be misleading

- ▶ One may be more interested in **correctly** predicting a particular outcome!
- ↪ This is often the case if the **cost** of being wrong differ
- ▶ One may need other measures, focused on one particular outcome
- ▶ Compute *Sensitivity & Specificity* from the confusion matrix



## [PROBLEM 1: ACCURACY IS ONE NUMBER ]

Accuracy is not the panacea and may be misleading

- ▶ One may be more interested in **correctly** predicting a particular outcome!
- ↪ This is often the case if the **cost** of being wrong differ
- ▶ One may need other measures, focused on one particular outcome
- ▶ Compute *Sensitivity & Specificity* from the confusion matrix
- ▶ They may go in different directions

# [SENSITIVITY OR *True Positive Rate*]

		Observed (True)	
		Urban	Rural
Predicted	Urban	87 (TP)	28 (FP)
	Rural	24 (FN)	69 (TN)

Table: Confusion Matrix

## [SENSITIVITY OR *True Positive Rate*]

		Observed (True)	
		Urban	Rural
Predicted	Urban	87 (TP)	28 (FP)
	Rural	24 (FN)	69 (TN)

Table: Confusion Matrix

- *Sensitivity* focuses on "positives" (here *Urban*), i.e. on predicted positives *vs* the observed positives

$$\begin{aligned} \text{Sensitivity} &= \frac{TP}{TP + FN} \\ &= \frac{87}{87 + 24} = 0.78 \end{aligned}$$

## [SENSITIVITY OR *True Positive Rate*]

		Observed (True)	
		Urban	Rural
Predicted	Urban	87 (TP)	28 (FP)
	Rural	24 (FN)	69 (TN)

Table: Confusion Matrix

- *Sensitivity* focuses on "positives" (here *Urban*), i.e. on predicted positives *vs* the observed positives

$$\begin{aligned} \text{Sensitivity} &= \frac{TP}{TP + FN} \\ &= \frac{87}{87 + 24} = 0.78 \end{aligned}$$

- On *Urban*, we correctly predict in 78% of the cases

# [SPECIFICITY OR *True Negative Rate*]

		Observed (True)	
		Urban	Rural
Predicted	Urban	87 (TP)	28 (FP)
	Rural	24 (FN)	69 (TN)

Table: Confusion Matrix

## [SPECIFICITY OR *True Negative Rate*]

		Observed (True)	
		Urban	Rural
Predicted	Urban	87 (TP)	28 (FP)
	Rural	24 (FN)	69 (TN)

Table: Confusion Matrix

- *Sensitivity* focuses on negatives (*Rural*), i.e. on predicted negatives *vs* the observed negatives

$$\text{Specificity} = \frac{TN}{TN + FP}$$

$$= \frac{69}{69 + 28} = 0.71$$

## [SPECIFICITY OR *True Negative Rate*]

		Observed (True)	
		Urban	Rural
Predicted	Urban	87 (TP)	28 (FP)
	Rural	24 (FN)	69 (TN)

Table: Confusion Matrix

- *Sensitivity* focuses on negatives (*Rural*), i.e. on predicted negatives *vs* the observed negatives

$$\begin{aligned}\text{Specificity} &= \frac{TN}{TN + FP} \\ &= \frac{69}{69 + 28} = 0.71\end{aligned}$$

- On *Rural*, we predict correctly in **only** 71% of the cases

## [PROBLEM 2: IMBALANCED OUTCOMES]

Imagine you observe much more Urban than Rural

Observed (True)	
Urban	Rural
95	5



## [PROBLEM 2: IMBALANCED OUTCOMES]

Imagine you observe much more Urban than Rural

Observed (True)	
Urban	Rural
95	5

- A "*stupid*" classifier predicting only *Urban* . . .

## [PROBLEM 2: IMBALANCED OUTCOMES]

Imagine you observe much more Urban than Rural

Observed (True)	
Urban	Rural
95	5

- A "*stupid*" classifier predicting only *Urban* . . .

		Observed (True)	
		Urban	Rural
Predicted	Urban	95 (TP)	5 (FP)
	Rural	0 (FN)	0 (TN)

## [PROBLEM 2: IMBALANCED OUTCOMES]

Imagine you observe much more Urban than Rural

Observed (True)	
Urban	Rural
95	5

- A "*stupid*" classifier predicting only *Urban* ...

		Observed (True)	
		Urban	Rural
Predicted	Urban	95 (TP)	5 (FP)
	Rural	0 (FN)	0 (TN)

... would have a very good *Accuracy* and *Sensitivity*

## [PROBLEM 2: IMBALANCED OUTCOMES]

Imagine you observe much more Urban than Rural

Observed (True)	
Urban	Rural
95	5

- A "*stupid*" classifier predicting only *Urban* ...

		Observed (True)	
		Urban	Rural
Predicted	Urban	95 (TP)	5 (FP)
	Rural	0 (FN)	0 (TN)

... would have a very good *Accuracy* and *Sensitivity*

$$\text{Accuracy} = (\text{TP} + \text{TN}) / 100 = 95 \%$$

## [PROBLEM 2: IMBALANCED OUTCOMES]

Imagine you observe much more Urban than Rural

Observed (True)	
Urban	Rural
95	5

- A "*stupid*" classifier predicting only *Urban* ...

		Observed (True)	
		Urban	Rural
Predicted	Urban	95 (TP)	5 (FP)
	Rural	0 (FN)	0 (TN)

... would have a very good *Accuracy* and *Sensitivity*

$$\text{Accuracy} = (\text{TP} + \text{TN}) / 100 = 95 \%$$

$$\text{Sensitivity} = \text{TP} / (\text{TP} + \text{FN}) = 95 / 95 = 100 \%$$

## [ THE KAPPA ( $\kappa$ ) INDEX ]

*Kappa* ( $\kappa$ ) is defined to measure the accuracy with imbalanced classes

Its formal definition is given by

$$\kappa = \frac{P_o - P_e}{1 - P_e}$$

## [ THE KAPPA ( $\kappa$ ) INDEX ]

*Kappa* ( $\kappa$ ) is defined to measure the accuracy with imbalanced classes

Its formal definition is given by

$$\kappa = \frac{P_o - P_e}{1 - P_e}$$

$P_o$  is the **current classifier accuracy** which is compared here with the accuracy of an uniformed classifier  $P_e$

## [ THE KAPPA ( $\kappa$ ) INDEX ]

*Kappa* ( $\kappa$ ) is defined to measure the accuracy with imbalanced classes

Its formal definition is given by

$$\kappa = \frac{P_o - P_e}{1 - P_e}$$

$P_o$  is the **current classifier accuracy** which is compared here with the accuracy of an uniformed classifier  $P_e$

$P_e$  is the **accuracy of an uniformed classifier** that would operate purely by chance, using no information.



## [ THE KAPPA ( $\kappa$ ) INDEX ]

*Kappa* ( $\kappa$ ) is defined to measure the accuracy with imbalanced classes

Its formal definition is given by

$$\kappa = \frac{P_o - P_e}{1 - P_e}$$

$P_o$  is the **current classifier accuracy** which is compared here with the accuracy of an uniformed classifier  $P_e$

$P_e$  is the **accuracy of an uniformed classifier** that would operate purely by chance, using no information.

NB:  $P_o$  is simple accuracy while  $P_e$  is more complex to compute.

## [ THE KAPPA ( $\kappa$ ) INDEX ]

*Kappa* ( $\kappa$ ) is defined to measure the accuracy with imbalanced classes

Its formal definition is given by

$$\kappa = \frac{P_o - P_e}{1 - P_e}$$

$P_o$  is the **current classifier accuracy** which is compared here with the accuracy of an uniformed classifier  $P_e$

$P_e$  is the **accuracy of an uniformed classifier** that would operate purely by chance, using no information.

NB:  $P_o$  is simple accuracy while  $P_e$  is more complex to compute.

- The larger  $\kappa$  is, the better the model for a given distribution of classes in a data set

# [QUIZ TIME]

# [QUIZ TIME]

## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important

## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important
- ▶ Many adjustment measures: **accuracy**, **sensitivity** and **specificity**.

## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important
- ▶ Many adjustment measures: **accuracy**, **sensitivity** and **specificity**.
  - ▶ *Sensitivity* is accuracy restricted to the positives.

## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important
- ▶ Many adjustment measures: **accuracy**, **sensitivity** and **specificity**.
  - ▶ *Sensitivity* is accuracy restricted to the positives.
  - ▶ *Specificity* is accuracy restricted to the negatives.



## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important
- ▶ Many adjustment measures: **accuracy**, **sensitivity** and **specificity**.
  - ▶ *Sensitivity* is accuracy restricted to the positives.
  - ▶ *Specificity* is accuracy restricted to the negatives.
- ▶ When outcome is *imbalanced*, one may use **kappa** has a better measure for accuracy.

## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important
- ▶ Many adjustment measures: **accuracy**, **sensitivity** and **specificity**.
  - ▶ *Sensitivity* is accuracy restricted to the positives.
  - ▶ *Specificity* is accuracy restricted to the negatives.
- ▶ When outcome is *imbalanced*, one may use **kappa** has a better measure for accuracy.

Which measure you should consider depends on the context and your goal.

# [LOGIT AS YOU KNOW IT ]

# [LOGIT AS YOU KNOW IT ]

$y$  is discrete, so no *direct* linear relationship between  $y$  and the explanatory variables  $x$  ( *i.e. Education, Income* )

# [LOGIT AS YOU KNOW IT ]

$y$  is discrete, so no *direct* linear relationship between  $y$  and the explanatory variables  $x$  ( *i.e.* *Education, Income* )

- Logit estimates the probabilities  $\pi$  ( $\in [0, 1]$ )

$$\pi = \text{Probability}[y = 1]$$

# [LOGIT AS YOU KNOW IT ]

$y$  is discrete, so no *direct* linear relationship between  $y$  and the explanatory variables  $x$  ( *i.e.* *Education, Income*)

- ▶ Logit estimates the probabilities  $\pi$  ( $\in [0, 1]$ )

$$\pi = \text{Probability}[y = 1]$$

- ▶ The definition of the logit model is:

$$\pi = \text{Pr}(y = 1) = F(x'\beta) = \frac{1}{1 + \exp(-x'\beta)}$$

## [LOGIT AS YOU KNOW IT ]

$y$  is discrete, so no *direct* linear relationship between  $y$  and the explanatory variables  $x$  ( *i.e.* *Education, Income*)

- ▶ Logit estimates the probabilities  $\pi$  ( $\in [0, 1]$ )

$$\pi = \text{Probability}[y = 1]$$

- ▶ The definition of the logit model is:

$$\pi = \text{Pr}(y = 1) = F(x'\beta) = \frac{1}{1 + \exp(-x'\beta)}$$

- ▶ This can be transformed into:

$$\pi = \frac{\exp(x'\beta)}{1 + \exp(x'\beta)}$$

# [LOGIT AS YOU **DON'T** KNOW IT ]



# [LOGIT AS YOU **DON'T** KNOW IT ]

From this equation

$$\pi = \frac{\exp(x'\beta)}{1 + \exp(x'\beta)}$$

# [LOGIT AS YOU **DON'T** KNOW IT ]

From this equation

$$\pi = \frac{\exp(x'\beta)}{1 + \exp(x'\beta)}$$

one gets the **linear** nature of the logit:

$$\log\left(\frac{\pi}{1 - \pi}\right) = x'\beta$$

# [LOGIT AS YOU **DON'T** KNOW IT ]

From this equation

$$\pi = \frac{\exp(x'\beta)}{1 + \exp(x'\beta)}$$

one gets the **linear** nature of the logit:

$$\log\left(\frac{\pi}{1 - \pi}\right) = x'\beta$$

where  $\frac{\pi}{1-\pi}$  is the odd ratio  $\in [0, \infty]$  with values indicating high or low probability that  $y = 1$

# [LOGIT AS YOU **DON'T** KNOW IT ]

From this equation

$$\pi = \frac{\exp(x'\beta)}{1 + \exp(x'\beta)}$$

one gets the **linear** nature of the logit:

$$\log\left(\frac{\pi}{1-\pi}\right) = x'\beta$$

where  $\frac{\pi}{1-\pi}$  is the odd ratio  $\in [0, \infty]$  with values indicating high or low probability that  $y = 1$

$\hookrightarrow$  "The logit models log of odd ratios as linear in  $x$ "

# [LOGIT AS A CLASSIFIER]

## [LOGIT AS A CLASSIFIER ]

- Once estimated,  $\hat{\pi}_i$  provide a simple rule for classification

## [LOGIT AS A CLASSIFIER]

- Once estimated,  $\hat{\pi}_i$  provide a simple rule for classification

$$\hat{\pi}_i > t_0 \Leftrightarrow \hat{y}_i = 1$$

## [LOGIT AS A CLASSIFIER ]

- Once estimated,  $\hat{\pi}_i$  provide a simple rule for classification

$$\hat{\pi}_i > t_0 \Leftrightarrow \hat{y}_i = 1$$

Where  $t_0$  is a threshold probability, by default 1/2.



## [LOGIT AS A CLASSIFIER]

- Once estimated,  $\hat{\pi}_i$  provide a simple rule for classification

$$\hat{\pi}_i > t_0 \Leftrightarrow \hat{y}_i = 1$$

Where  $t_0$  is a threshold probability, by default 1/2.

- If  $t_0 = 1/2$  (default), then the rule is equivalent to:

$$x'_i \hat{\beta} > 0 \Leftrightarrow \hat{y}_i = 1$$

## [LOGIT AS A CLASSIFIER]

- Once estimated,  $\hat{\pi}_i$  provide a simple rule for classification

$$\hat{\pi}_i > t_0 \Leftrightarrow \hat{y}_i = 1$$

Where  $t_0$  is a threshold probability, by default  $1/2$ .

- If  $t_0 = 1/2$  (default), then the rule is equivalent to:

$$x'_i \hat{\beta} > 0 \Leftrightarrow \hat{y}_i = 1$$

- If  $t_0 \neq 1/2$ :

$$x'_i \hat{\beta} > T_0 \Leftrightarrow \hat{y}_i = 1$$

## [LOGIT AS A CLASSIFIER]

- Once estimated,  $\hat{\pi}_i$  provide a simple rule for classification

$$\hat{\pi}_i > t_0 \Leftrightarrow \hat{y}_i = 1$$

Where  $t_0$  is a threshold probability, by default  $1/2$ .

- If  $t_0 = 1/2$  (default), then the rule is equivalent to:

$$x'_i \hat{\beta} > 0 \Leftrightarrow \hat{y}_i = 1$$

- If  $t_0 \neq 1/2$ :

$$x'_i \hat{\beta} > T_0 \Leftrightarrow \hat{y}_i = 1$$

↪ The logit classifier depends on the linear combination of the  $x$ 's

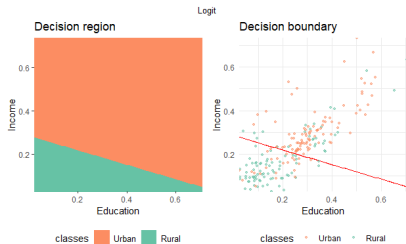
# [IMPORTANCE OF THE THRESHOLD ]

## [IMPORTANCE OF THE THRESHOLD ]

- ▶ The rule  $x'\beta \geq T_0$  defines the partition of the space

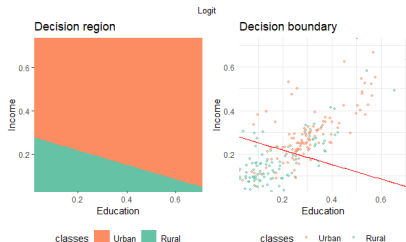
# [IMPORTANCE OF THE THRESHOLD ]

- The rule  $x'\beta \geq T_0$  defines the partition of the space



# [IMPORTANCE OF THE THRESHOLD ]

- The rule  $x'\beta \geq T_0$  defines the partition of the space



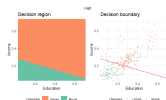
- This partition is sensitive to the choice of the threshold  $T_0$  (and the  $t_0$ )

# [IMPORTANCE OF THE THRESHOLD ]



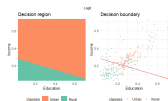


# [IMPORTANCE OF THE THRESHOLD ]



- Changing  $t_0$  will change the predictions & the classification

# [IMPORTANCE OF THE THRESHOLD ]



- Changing  $t_0$  will change the predictions & the classification  
A higher  $t_0$  will allocate less observations to the  $y = 1$  category (Urban)

# [IMPORTANCE OF THE THRESHOLD ]



- Changing  $t_0$  will change the predictions & the classification  
A higher  $t_0$  will allocate less observations to the  $y = 1$  category (Urban)  
A lower  $t_0$  will allocate more observations to the  $y = 1$  category

# [IMPORTANCE OF THE THRESHOLD ]



- Changing  $t_0$  will change the predictions & the classification  
A higher  $t_0$  will allocate less observations to the  $y = 1$  category (Urban)  
A lower  $t_0$  will allocate more observations to the  $y = 1$  category
- The choice of  $t_0$  should be done according to the data and observed classes repartition

# [IMPORTANCE OF THE THRESHOLD ]



- ▶ Changing  $t_0$  will change the predictions & the classification  
A higher  $t_0$  will allocate less observations to the  $y = 1$  category (Urban)  
A lower  $t_0$  will allocate more observations to the  $y = 1$  category
- ▶ The choice of  $t_0$  should be done according to the data and observed classes repartition
- ▶ *Specificity* and *Sensitivity* are affected by  $t_0$

# [THE ROC CURVE ]

# [THE ROC CURVE ]

- We want the *Specificity* and *Sensitivity* to be both **maximized** (ideally both would be 1)

# [THE ROC CURVE ]

- ▶ We want the *Specificity* and *Sensitivity* to be both **maximized** (ideally both would be 1)
- ▶ The ROC curve help visualize the best choice



## [THE ROC CURVE ]

- ▶ We want the *Specificity* and *Sensitivity* to be both **maximized** (ideally both would be 1)
- ▶ The ROC curve help visualize the best choice
- ▶ The ROC plots both Sensitivity and Specificity values for different thresholds

## [THE ROC CURVE ]

- ▶ We want the *Specificity* and *Sensitivity* to be both **maximized** (ideally both would be 1)
- ▶ The ROC curve help visualize the best choice
- ▶ The ROC plots both Sensitivity and Specificity values for different thresholds

*Be careful of the axes*

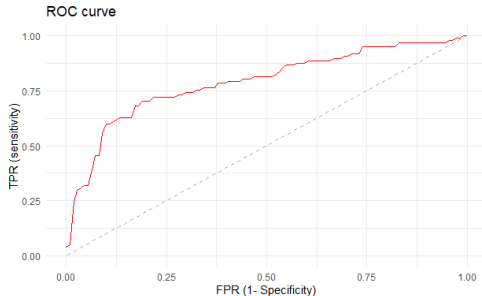
# [THE ROC CURVE]

# [THE ROC CURVE]

The ROC represents values of  $1 - \text{Specificity} = \text{FPR}$  vs  
 $\text{Sensitivity} = \text{TPR}$  for many values of the threshold  $t_0$

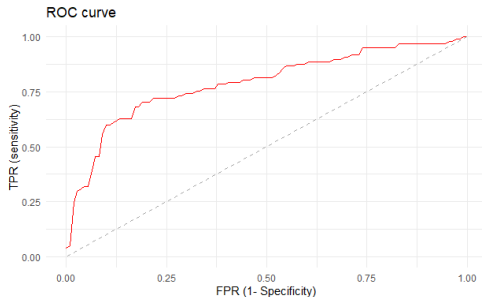
# [THE ROC CURVE]

The ROC represents values of  $1 - \text{Specificity} = \text{FPR}$  vs  
 $\text{Sensitivity} = \text{TPR}$  for many values of the threshold  $t_0$



# [THE ROC CURVE]

The ROC represents values of  $1 - \text{Specificity} = \text{FPR}$  vs  $\text{Sensitivity} = \text{TPR}$  for many values of the threshold  $t_0$



- (sometimes  $x$  is sensitivity with inverted  $x$ -axis)

# [THE ROC CURVE: HOW TO READ?]

# [THE ROC CURVE: HOW TO READ?]

Changing  $t_0$  changes the classification

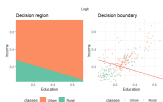




# [THE ROC CURVE: HOW TO READ?]

Changing  $t_0$  changes the classification

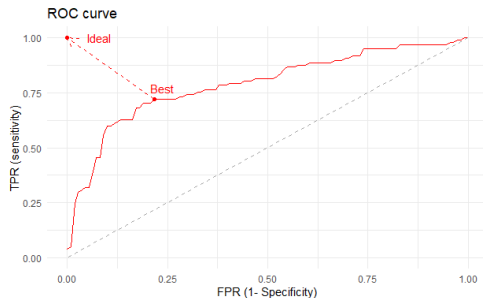
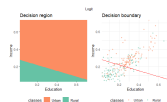
- Optimally, the curve should touch top-left corner



# [THE ROC CURVE: HOW TO READ?]

Changing  $t_0$  changes the classification

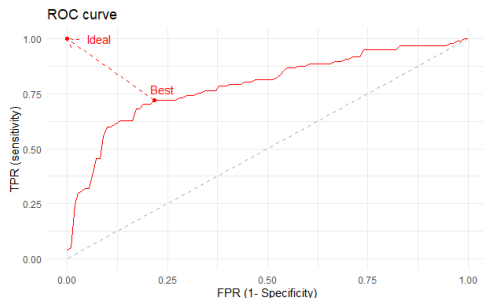
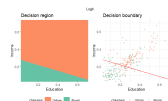
- Optimally, the curve should touch top-left corner



# [THE ROC CURVE: HOW TO READ?]

Changing  $t_0$  changes the classification

- Optimally, the curve should touch top-left corner

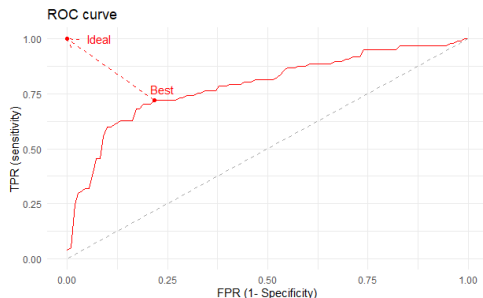
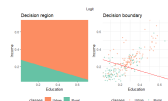


- If  $t_0 \nearrow$ , more cases classified as *Negatives*, less *Positives*

# [THE ROC CURVE: HOW TO READ?]

Changing  $t_0$  changes the classification

- Optimally, the curve should touch top-left corner



- If  $t_0 \nearrow$ , more cases classified as *Negatives*, less *Positives*
- If  $t_0 \nearrow$ , specificity  $\nearrow$  and sensitivity  $\searrow$

# [AUC AS A MEASURE OF FIT]

## [AUC AS A MEASURE OF FIT]

*A model that works well, whatever the threshold is certainly desirable*

## [AUC AS A MEASURE OF FIT]

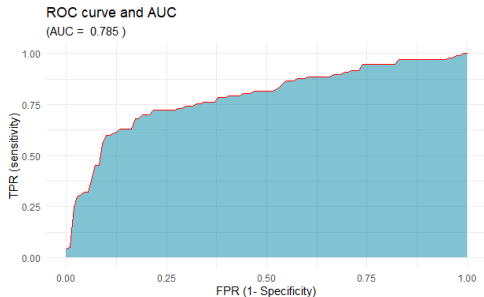
*A model that works well, whatever the threshold is certainly desirable*

- Using the AUC is also a measure of fit of a model

## [AUC AS A MEASURE OF FIT]

*A model that works well, whatever the threshold is certainly desirable*

- Using the AUC is also a measure of fit of a model

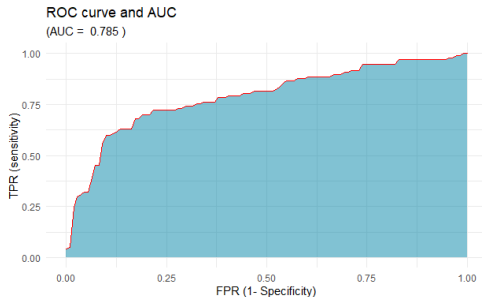




## [AUC AS A MEASURE OF FIT]

*A model that works well, whatever the threshold is certainly desirable*

- Using the AUC is also a measure of fit of a model



- The greater the area, the better the model

# [ COMPARING THE MEASURES]

## [ COMPARING THE MEASURES ]

- ▶ We have several measures at hand

## [ COMPARING THE MEASURES]

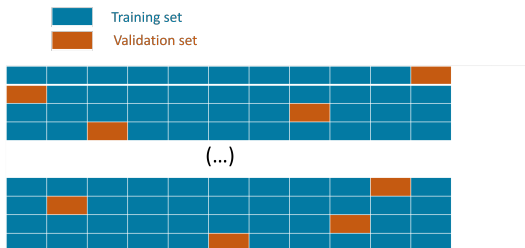
- ▶ We have several measures at hand
- ▶ We should evaluate those models on their predictive performance on a new "*unseen*" data set

## [ COMPARING THE MEASURES ]

- ▶ We have several measures at hand
  - ▶ We should evaluate those models on their predictive performance on a new "*unseen*" data set
- ↪ This is what Cross-Validation can do

# [ COMPARING THE MEASURES ]

- ▶ We have several measures at hand
  - ▶ We should evaluate those models on their predictive performance on a new "*unseen*" data set
- ↪ This is what Cross-Validation can do



# [ COMPARING THE MEASURES]

## [ COMPARING THE MEASURES]

- For any model, CV gives several classifications



## [ COMPARING THE MEASURES]

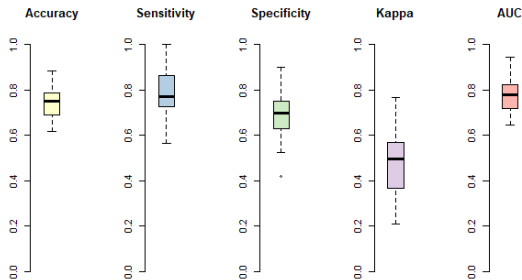
- ▶ For any model, CV gives several classifications
- ▶ All the criteria derive from the confusion matrix

## [ COMPARING THE MEASURES]

- ▶ For any model, CV gives several classifications
  - ▶ All the criteria derive from the confusion matrix
- ↪ Examine them all!

## [ COMPARING THE MEASURES ]

- ▶ For any model, CV gives several classifications
  - ▶ All the criteria derive from the confusion matrix
- ↪ Examine them all!



# [ HOW TO CHOSE THE BEST MODEL? ]

# [ HOW TO CHOSE THE BEST MODEL? ]

- ▶ We have several criteria for one model

## [ HOW TO CHOSE THE BEST MODEL? ]

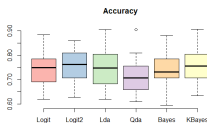
- ▶ We have several criteria for one model
- ▶ We should again evaluate the classifier based on "*unseen*" data set

## [ HOW TO CHOSE THE BEST MODEL? ]

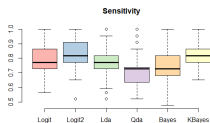
- ▶ We have several criteria for one model
  - ▶ We should again evaluate the classifier based on "*unseen*" data set
- ↪ Run Cross-Validation an all!

# [ HOW TO CHOSE THE BEST MODEL? ]

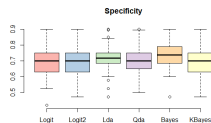
- ▶ We have several criteria for one model
  - ▶ We should again evaluate the classifier based on "*unseen*" data set
- ↪ Run Cross-Validation an all!



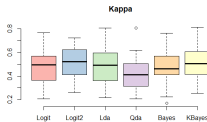
Accuracy of all CV validation sets, on all models



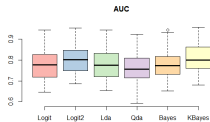
Sensitivity of all CV validation set, on all models



Specificity of all CV validation set, on all models



Kappa of all CV validation sets, on all models



AUC of all CV validation sets, on all models



# [QUIZ TIME]

# [QUIZ TIME]

## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important

## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important
- ▶ Many adjustment measures: **accuracy, sensitivity** and **specificity**.

## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important
- ▶ Many adjustment measures: **accuracy**, **sensitivity** and **specificity**.
  - ▶ *Sensitivity* is accuracy restricted to the positives.

## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important
- ▶ Many adjustment measures: **accuracy**, **sensitivity** and **specificity**.
  - ▶ *Sensitivity* is accuracy restricted to the positives.
  - ▶ *Specificity* is accuracy restricted to the negatives.

## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important
- ▶ Many adjustment measures: **accuracy**, **sensitivity** and **specificity**.
  - ▶ *Sensitivity* is accuracy restricted to the positives.
  - ▶ *Specificity* is accuracy restricted to the negatives.
- ▶ When outcome is *imbalanced*, one may use **kappa** as a better measure for accuracy.

## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important
- ▶ Many adjustment measures: **accuracy**, **sensitivity** and **specificity**.
  - ▶ *Sensitivity* is accuracy restricted to the positives.
  - ▶ *Specificity* is accuracy restricted to the negatives.
- ▶ When outcome is *imbalanced*, one may use **kappa** as a better measure for accuracy.

Which measure you should consider depends on the context and your goal.



## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important
- ▶ Many adjustment measures: **accuracy**, **sensitivity** and **specificity**.
  - ▶ *Sensitivity* is accuracy restricted to the positives.
  - ▶ *Specificity* is accuracy restricted to the negatives.
- ▶ When outcome is *imbalanced*, one may use **kappa** as a better measure for accuracy.

Which measure you should consider depends on the context and your goal.
- ▶ **Logit** is a benchmark parametric model for classification

## [TAKEAWAYS]

- ▶ In classification, the **Confusion matrix** is important
- ▶ Many adjustment measures: **accuracy**, **sensitivity** and **specificity**.
  - ▶ *Sensitivity* is accuracy restricted to the positives.
  - ▶ *Specificity* is accuracy restricted to the negatives.
- ▶ When outcome is *imbalanced*, one may use **kappa** as a better measure for accuracy.

Which measure you should consider depends on the context and your goal.

- ▶ **Logit** is a benchmark parametric model for classification  
One may use the **ROC** to change the threshold parameter

## [TAKEAWAYS]

- Use Training-Validation set to **select** parameters within a model

## [TAKEAWAYS]

- ▶ Use Training-Validation set to **select** parameters within a model
- ▶ Use Training-Validation set to **compare** models on the same criteria

## [TAKEAWAYS]

- ▶ Use Training-Validation set to **select** parameters within a model
- ▶ Use Training-Validation set to **compare** models on the same criteria
- ▶ Several criteria / measures of fit / cost functions are available

## [TAKEAWAYS]

- ▶ Use Training-Validation set to **select** parameters within a model
- ▶ Use Training-Validation set to **compare** models on the same criteria
- ▶ Several criteria / measures of fit / cost functions are available
- ▶ Time is the limit...