

Machine Learning for Official Statistics and SDGs

Supervised Classification



[CLASSIFICATION]

What is a classification problem?

[CLASSIFICATION]

What is a classification problem?

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

[CLASSIFICATION]

What is a classification problem?

- ▶ The goal of classification is to understand why an observation belongs to a certain category
- ▶ The interest is on y that takes discrete values: 0/1, high school/prilmary school/no education; urban/rural

[CLASSIFICATION]

What is a classification problem?

- ▶ The goal of classification is to understand why an observation belongs to a certain category
- ▶ The interest is on y that takes discrete values: 0/1, high school/prilmary school/no education; urban/rural
- ▶ There may be variables explaining why we observe that y belongs to a particular category

[CLASSIFICATION]

What is a classification problem?

- ▶ The goal of classification is to understand why an observation belongs to a certain category
- ▶ The interest is on y that takes discrete values: 0/1, high school/prilmary school/no education; urban/rural
- ▶ There may be variables explaining why we observe that y belongs to a particular category
- ▶ A classifier is a tool that provides a classification for y

[SUPERVISED *vs* UNSUPERVISED CLASSIFICATION]

[SUPERVISED *vs* UNSUPERVISED CLASSIFICATION]

- In **supervised** classification, we observe the class for y

[SUPERVISED *vs* UNSUPERVISED CLASSIFICATION]

- In **supervised** classification, we observe the class for y
One may learn from that information and estimate the impact of other variables on that classification (*e.g.* logit regression)

[SUPERVISED *vs* UNSUPERVISED CLASSIFICATION]

- ▶ In **supervised** classification, we observe the class for y
One may learn from that information and estimate the impact of other variables on that classification (*e.g.* logit regression)
- ▶ In **unsupervised** classification, we observe a set of variables for each observation

[SUPERVISED *vs* UNSUPERVISED CLASSIFICATION]

- ▶ In **supervised** classification, we observe the class for y
One may learn from that information and estimate the impact of other variables on that classification (*e.g.* logit regression)
- ▶ In **unsupervised** classification, we observe a set of variables for 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 class for y
One may learn from that information and estimate the impact of other variables on that classification (*e.g.* logit regression)
- ▶ In **unsupervised** classification, we observe a set of variables for 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 class for y
One may learn from that information and estimate the impact of other variables on that classification (*e.g.* logit regression)
- ▶ In **unsupervised** classification, we observe a set of variables for 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

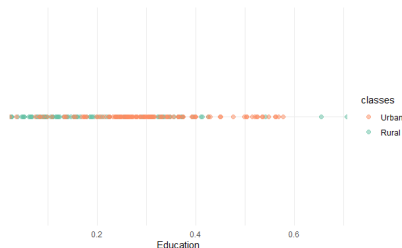
[CLASSIFICATION: AN EXAMPLE]

[CLASSIFICATION: AN EXAMPLE]

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

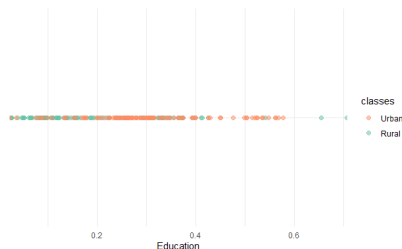
[CLASSIFICATION: AN EXAMPLE]

- 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", typically with a threshold rule : "if $x \geq x_t$ then y is categorized as Urban"

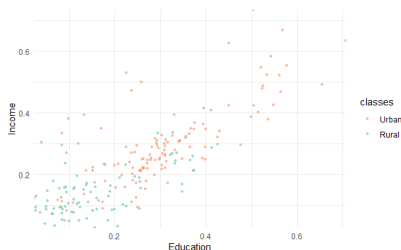
[CLASSIFICATION: A 2-D EXAMPLE]

[CLASSIFICATION: A 2-D EXAMPLE]

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

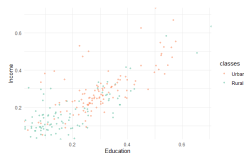
[CLASSIFICATION: A 2-D EXAMPLE]

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

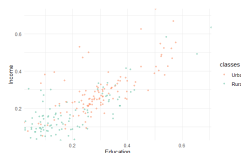


[CLASSIFICATION: A 2-D EXAMPLE]

[CLASSIFICATION: A 2-D EXAMPLE]

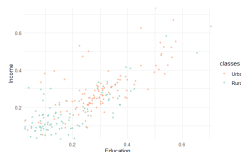


[CLASSIFICATION: A 2-D EXAMPLE]



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

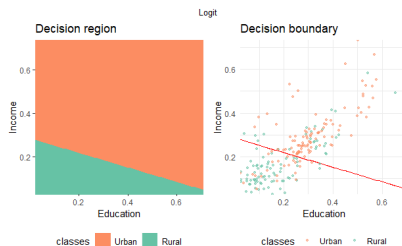
[CLASSIFICATION: A 2-D EXAMPLE]

[CLASSIFICATION: A 2-D EXAMPLE]

- ▶ The logit is an example of linear classifier

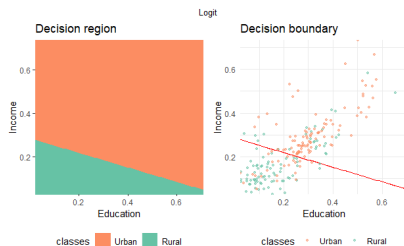
[CLASSIFICATION: A 2-D EXAMPLE]

- The logit is an example of linear classifier



[CLASSIFICATION: A 2-D EXAMPLE]

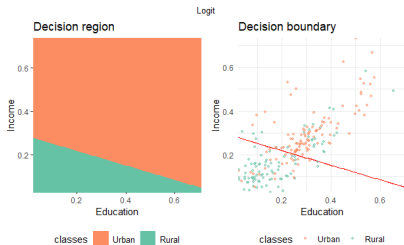
- The logit is an example of linear classifier



- The rule that separated the two classes is $x'\beta \geq T_0$ with T_0 a known threshold

[CLASSIFICATION: A 2-D EXAMPLE]

- The logit is an example of linear classifier



- The rule that separated the two classes is $x'\beta \geq T_0$ with T_0 a known threshold
e.g. $\beta_1 \text{Education} + \beta_2 \text{Income} \geq T_0$

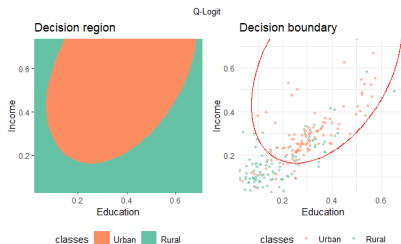
[CLASSIFICATION: A 2-D EXAMPLE]

[CLASSIFICATION: A 2-D EXAMPLE]

- The quadratic logit is an example of non-linear classifier

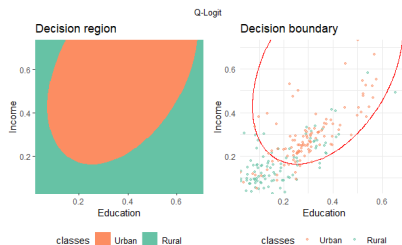
[CLASSIFICATION: A 2-D EXAMPLE]

- The quadratic logit is an example of non-linear classifier



[CLASSIFICATION: A 2-D EXAMPLE]

- The quadratic logit is an example of non-linear classifier



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

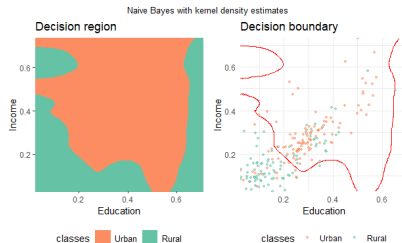
[CLASSIFICATION: A 2-D EXAMPLE]

[CLASSIFICATION: A 2-D EXAMPLE]

- ▶ Other examples can be very non linear

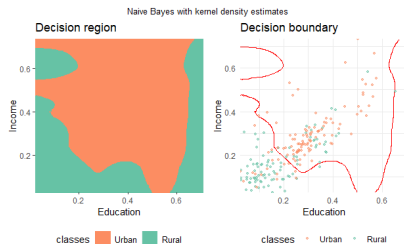
[CLASSIFICATION: A 2-D EXAMPLE]

- Other examples can be very non linear



[CLASSIFICATION: A 2-D EXAMPLE]

- Other examples can be very non linear



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

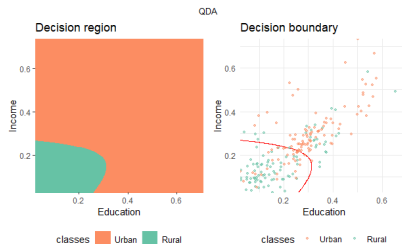
[CLASSIFICATION: A 2-D EXAMPLE]

[CLASSIFICATION: A 2-D EXAMPLE]

- ▶ Other examples can be very non linear

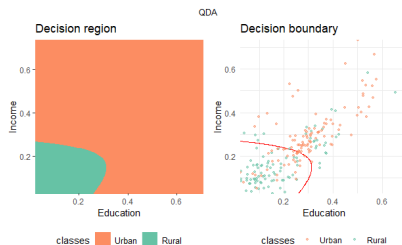
[CLASSIFICATION: A 2-D EXAMPLE]

- Other examples can be very non linear



[CLASSIFICATION: A 2-D EXAMPLE]

- Other examples can be very non linear



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

[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
- ▶ Specificity & sensitivity

[MEASURES OF FIT IN CLASSIFICATION]

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

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

[MEASURES OF FIT IN CLASSIFICATION]

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

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

[MEASURES OF FIT IN CLASSIFICATION]

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

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

...

All answer to different questions and solve different problems

[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 Positive)	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 Positive)	TN (True Negative)

Table: Confusion Matrix

- Accuracy is then the ratio:

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

[CONFUSION MATRIX & ACCURACY]

A classifier predicts in which class each observation should be:

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

- Remember that 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

- Remember that here *Urban* is the "positive" class
- Accuracy is then the ratio:

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

[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

- ▶ Remember that 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.

[CONFUSION MATRIX & ACCURACY]

Accuracy is not the panacea and may be misleading

Pb 1: Imagine a data set where 95% of the observations are in class 1 and 5% in the remaining class

[CONFUSION MATRIX & ACCURACY]

Accuracy is not the panacea and may be misleading

Pb 1: Imagine a data set where 95% of the observations are in class 1 and 5% in the remaining class

↪ A classifier that predicts always class 1 will have an accuracy of 95%!

[CONFUSION MATRIX & ACCURACY]

Accuracy is not the panacea and may be misleading

Pb 1: Imagine a data set where 95% of the observations are in class 1 and 5% in the remaining class

- ↪ A classifier that predicts always class 1 will have an accuracy of 95%!
- ▶ Accuracy is not very useful with an imbalanced data set and one may prefer another indicator: *kappa*

[CONFUSION MATRIX & ACCURACY]

Accuracy is not the panacea and may be misleading

Pb 1: Imagine a data set where 95% of the observations are in class 1 and 5% in the remaining class

↪ A classifier that predicts always class 1 will have an accuracy of 95%!

► Accuracy is not very useful with an imbalanced data set and one may prefer another indicator: *kappa*

Pb 2: One may be more interested in correctly predicting a particular outcome

[CONFUSION MATRIX & ACCURACY]

Accuracy is not the panacea and may be misleading

Pb 1: Imagine a data set where 95% of the observations are in class 1 and 5% in the remaining class

↪ A classifier that predicts always class 1 will have an accuracy of 95%!

► Accuracy is not very useful with an imbalanced data set and one may prefer another indicator: *kappa*

Pb 2: One may be more interested in correctly predicting a particular outcome

↪ This is often the case

[CONFUSION MATRIX & ACCURACY]

Accuracy is not the panacea and may be misleading

Pb 1: Imagine a data set where 95% of the observations are in class 1 and 5% in the remaining class

↪ A classifier that predicts always class 1 will have an accuracy of 95%!

► Accuracy is not very useful with an imbalanced data set and one may prefer another indicator: *kappa*

Pb 2: One may be more interested in correctly predicting a particular outcome

↪ This is often the case

► One may need other measures, such as *Specificity* or *Sensitivity*

[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 positives *vs* the observed positives

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

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

[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 positives *vs* the observed positives

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

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

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

[CONFUSION MATRIX & KAPPA (κ)]

Kappa (κ) is defined to measure the accuracy with imbalanced classes

Its formal definition is given by

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

[CONFUSION MATRIX & KAPPA (κ)]

Kappa (κ) 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

[CONFUSION MATRIX & KAPPA (κ)]

Kappa (κ) 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 a classifier that would operate purely by chance, using no information.

[CONFUSION MATRIX & KAPPA (κ)]

Kappa (κ) 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 a classifier that would operate purely by chance, using no information.

- P_o is simple (here = 0.75) while P_e is more complex to compute.

[CONFUSION MATRIX & KAPPA (κ)]

Kappa (κ) 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 a classifier that would operate purely by chance, using no information.

- ▶ P_o is simple (here = 0.75) while P_e is more complex to compute.
- ▶ The larger κ is, the better the model for a given distribution of classes in a data set

[LOGIT AS YOU KNOW IT]

[LOGIT AS YOU KNOW IT]

When dealing with a discrete outcome y we cannot use a *direct* linear relationship between y and the explanatory variables x (*i.e.* *Education, Income*)

[LOGIT AS YOU KNOW IT]

When dealing with a discrete outcome y we cannot use a *direct* linear relationship between y and the explanatory variables x (*i.e.* *Education, Income*)

- Logit aims at estimating the probabilities $\pi = \text{Probability}[y = 1]$, which are continuous ($\in [0, 1]$)

[LOGIT AS YOU KNOW IT]

When dealing with a discrete outcome y we cannot use a *direct* linear relationship between y and the explanatory variables x (*i.e. Education, Income*)

- ▶ Logit aims at estimating the probabilities $\pi = \text{Probability}[y = 1]$, which are continuous ($\in [0, 1]$)
- ▶ 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]

When dealing with a discrete outcome y we cannot use a *direct* linear relationship between y and the explanatory variables x (*i.e. Education, Income*)

- ▶ Logit aims at estimating the probabilities $\pi = \text{Probability}[y = 1]$, which are continuous ($\in [0, 1]$)
- ▶ 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 provability, 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 provability, 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 provability, 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 provability, 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

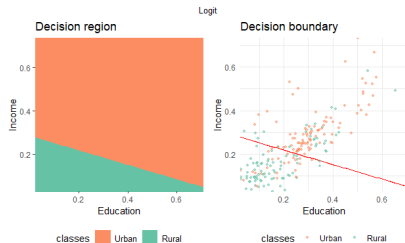
[IMPORTANCE OF THE THRESHOLD]

[IMPORTANCE OF THE THRESHOLD]

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

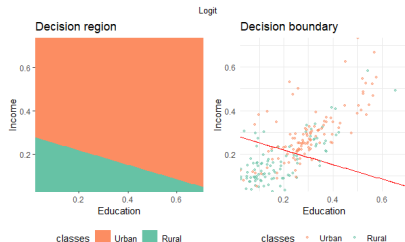
[IMPORTANCE OF THE THRESHOLD]

- 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

Be careful of the axes

[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
Be careful of the axes
- ▶

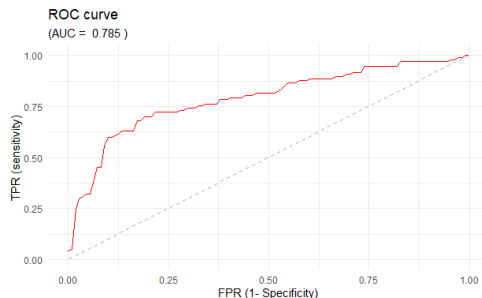
[THE ROC CURVE]

[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

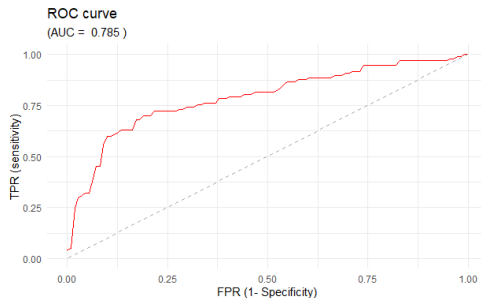
[THE ROC CURVE]

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 of observations

[THE ROC CURVE: HOW TO READ?]

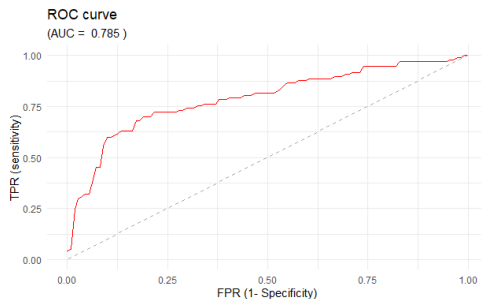
Changing t_0 changes the classification of observations

- Optimally, the curve should touch top-left corner

[THE ROC CURVE: HOW TO READ?]

Changing t_0 changes the classification of observations

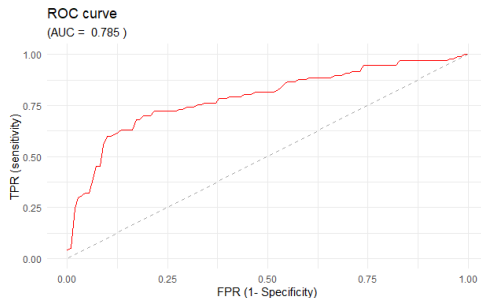
- Optimally, the curve should touch top-left corner



[THE ROC CURVE: HOW TO READ?]

Changing t_0 changes the classification of observations

- Optimally, the curve should touch top-left corner

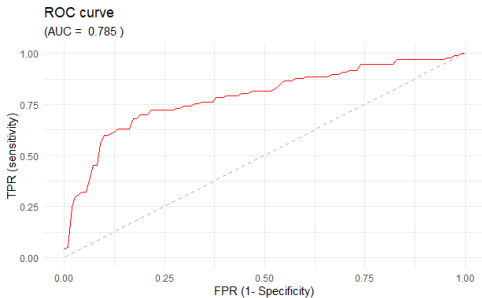


- If $t_0 \nearrow$, more cases classified as *Negatives*, less *Positives*

[THE ROC CURVE: HOW TO READ?]

Changing t_0 changes the classification of observations

- 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]

[TRAIN & VALIDATION SET]

[CV ON MEASURES OF FIT]

[MODEL COMPARISON]

[TAKEWAYS]

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

[TAKEWAYS]

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

[TAKEWAYS]

- ▶ 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.

[TAKEWAYS]

- ▶ 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.

[TAKEWAYS]

- ▶ 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

[TAKEWAYS]

- ▶ 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

[TAKEWAYS]

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

[TAKEWAYS]

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

[TAKEWAYS]

- ▶ 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

[TAKEWAYS]

- ▶ 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...