

SCHOOL OF ELECTRONIC ENGINEERING AND COMPUTER SCIENCE  
QUEEN MARY UNIVERSITY OF LONDON

# CBU5201 Principles of Machine Learning

## Supervised learning: Classification II

Dr Chao Liu

Credit to Dr Jesús Requena Carrión

Oct 2023



# Agenda

The Bayes classifier

Using data to build posterior probabilities

Classification performance: Beyond accuracy

Summary

# Is a high accuracy what we want?

Our notion of **quality** is defined by a **metric**, which allows us to rank different solutions. We have used two equivalent metrics for classifiers:

- **Accuracy**: Proportion of correctly classified samples.
- **Error rate**: Proportion of misclassified samples.

Note that both metrics are **blind to the class** that misclassified samples belong to.

However, is it the same misclassifying:

- A healthy patient and an ill patient, when deciding whether to administer some treatment?
- A good business and a bad business, when receiving a loan application?

If it is not, accuracy or error rate are not the quality metrics that we need.

# A Bayesian extension

Consider a binary problem with two classes  $\circ$  and  $\bullet$ , where:

- The cost of misclassifying a  $\bullet$  sample is  $C_{\bullet}$ .
- The cost of misclassifying a  $\circ$  sample is  $C_{\circ}$ .

The Bayes classifier achieves the highest accuracy by comparing:

$$\frac{P(y = \bullet | \mathbf{x})}{P(y = \circ | \mathbf{x})} \lessgtr 1$$

and assigning the sample to  $\bullet$  if  $> 1$  or  $\circ$  if  $< 1$ .

The expected cost will be

- $C_{\bullet} \times P(y = \bullet | \mathbf{x})$ , if our classifier labels the sample as  $\bullet$ .
- $C_{\circ} \times P(y = \circ | \mathbf{x})$ , if our classifier labels the sample as  $\circ$ .

## A Bayesian extension

Consider the following posterior probabilities:

- $P(y = \text{red} | \mathbf{x}) = 0.9$ ,
- $P(y = \text{blue} | \mathbf{x}) = 0.1$ ,

and misclassification costs:

- $C_{\text{blue}} = 5000 \text{ £}$ ,
- $C_{\text{red}} = 20 \text{ £}$ .

The Bayesian classifier would label sample  $\mathbf{x}$  as red. Accordingly:

- 10% of the time you would be misclassifying blue samples.
- The expected cost will be  $C_{\text{blue}} \times P(y = \text{blue} | \mathbf{x}) = 5000 \times 0.1 = 500 \text{ £}$ .

What if we were to label sample  $\mathbf{x}$  as blue against the advice of the Bayes classifier? The accuracy would be lower (10%), but the expected cost would be  $C_{\text{red}} \times P(y = \text{red} | \mathbf{x}) = 20 \times 0.9 = 18 \text{ £}$ .

## A Bayesian extension

To account for misclassification costs, we can use the following comparison instead:

$$\frac{C_{\circ} \times P(y = \circ | \mathbf{x})}{C_{\bullet} \times P(y = \bullet | \mathbf{x})} \leq 1 \quad \text{or} \quad \frac{P(y = \circ | \mathbf{x})}{P(y = \bullet | \mathbf{x})} \leq \frac{C_{\bullet}}{C_{\circ}}$$

A classifier that follows this strategy will **minimise the expected cost**, rather than the accuracy.

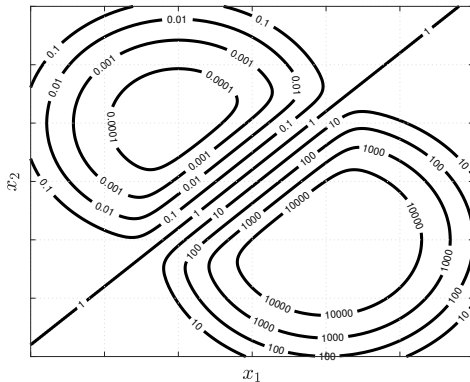
In general, our Bayesian extension will be expressed as

$$\frac{P(y = \circ | \mathbf{x})}{P(y = \bullet | \mathbf{x})} \leq T$$

where  $T$  is a threshold value corresponding to the ration of misclassification costs.

# A Bayesian extension

Changing the threshold value  $T$  changes the boundary of our classifier.



# Confusion matrix

Accuracy and error rate are not the most suitable quality metrics in **class-sensitive problems**, where the cost of misclassifying samples from different classes is different.

In addition to using the **misclassification cost** for each class, we can assess how well the classifier deals with each class individually. This is precisely the information that a **confusion** or **contingency matrix** shows.



# Confusion matrix: Counting

A confusion matrix shows for each class, the number of samples:

- **Correctly classified:** diagonal
- **Misclassified:** non-diagonal

In the following confusion matrix, 3 ○ samples are misclassified as ○, and 4 ○ samples are correctly classified. We can also learn that the dataset has 10 ○ samples, 20 ○ samples and 5 ○ samples and the **accuracy** is 24/35.

		Actual class		
		<span style="color: red;">○</span>	<span style="color: blue;">○</span>	<span style="color: green;">○</span>
Predicted class	<span style="color: red;">○</span>	5	2	0
	<span style="color: blue;">○</span>	3	15	1
	<span style="color: green;">○</span>	2	3	4

## Confusion matrix: Rates

The confusion matrix can also show rates, defined as the proportion of samples from one class that are assigned to any other class.

This is useful when working with imbalanced datasets, where the counts might be misleading. The example confusion matrix below uses rates, instead of counts.

		Actual class		
		○	○	○
Predicted class	○	0.5	0.1	0
	○	0.3	0.75	0.2
	○	0.2	0.15	0.8

# Confusion matrix: Detection problems

Many binary problems consider classes that represent the **presence** or **absence** of some property. For these problems, it is common to use the terms **positive** (presence) and **negative** (absence) for each class.

		Actual class	
		Positive	Negative
Predicted class	Positive	True positive	False positive
	Negative	False negative	True negative

In addition, we use the terms:

- True positive (TP) and true positive rate (TPR).
- False negative (FN) and false negative rate (FNR).
- False positive (FP) and false positive rate (FPR).
- True negative (TN) and true negative rate (TNR).

# Confusion matrix: Detection example

Number of samples

Predicted	Actual	
	10	11
	2	9

Rates

Predicted	Actual	
	0.83	0.55
	0.17	0.45

- $TP = 10 \rightarrow TPR = 10/12 = 0.83$
- $FP = 11 \rightarrow FPR = 11/20 = 0.55$
- $FN = 2 \rightarrow FNR = 2/12 = 0.17$
- $TN = 9 \rightarrow TNR = 9/20 = 0.45$
- $A = (10+9)/32 = 19/32 = 0.59$
- $E = (11+2)/32 = 13/32 = 0.41$

# Class-sensitive rates: Terminology

Error rate and accuracy are performance rates that do not allow us to investigate how a classifier treats each class. To do so, we can define other rates, such as the ones included in a confusion matrix.

In detection problems, the most commonly rates used are:

- **Sensitivity** (recall or true positive rate):  $TP/(TP+FN)$
- **Specificity** (true negative rate):  $TN/(TN+FP)$
- **Precision** (positive predictive value):  $TP/(TP+FP)$

These rates can be used as **quality metrics**.

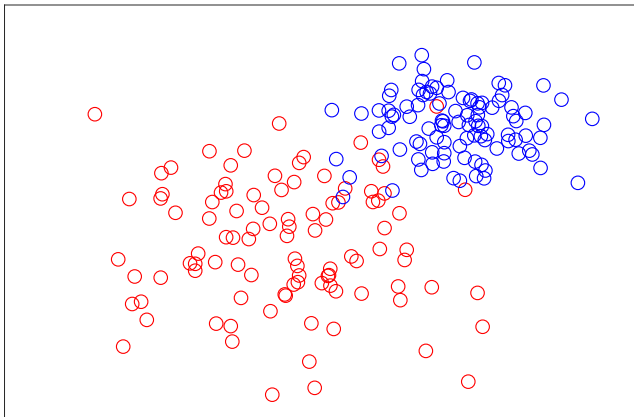
## Class-sensitive rates rates: Example

	Actual	
	10	11
Predicted	2	9

- Sensitivity =  $10/12 = 0.83$
- Specificity =  $9/20 = 0.45$
- Precision =  $10/21 = 0.48$

# Class-sensitive rates rates: Optimisation

If ● is the positive class and ● the negative class, obtain a linear boundary with the highest sensitivity and another with the highest specificity.



# Confusion matrix: Optimisation

Improving one quality metric individually is easy. For instance, if we label every sample as positive, we would achieve a perfect sensitivity. The problem is that improving **one quality metric deteriorates others**.



# Confusion matrix: Optimisation

Since improving the performance on class deteriorates the performance on the other, we usually consider simultaneously **pairs of quality metrics**:

- Sensitivity and specificity.
- Precision and recall.

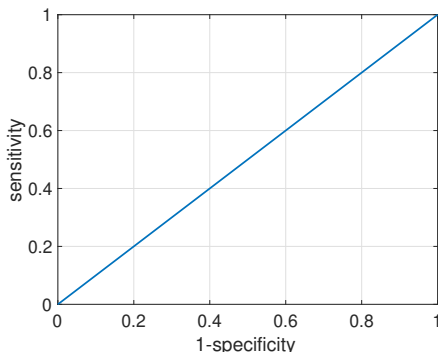
*(For instance, to travel to the UK you need to take a COVID test of 80 % sensitivity and 97 % specificity.)*

The F1-score is another widely used performance metric that provides an average between precision and recall:

$$F1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

# The ROC plane

The ROC (*receiver operating characteristic*) plane is used to represent the performance of a classifier in terms of its **sensitivity** and **1-specificity**.

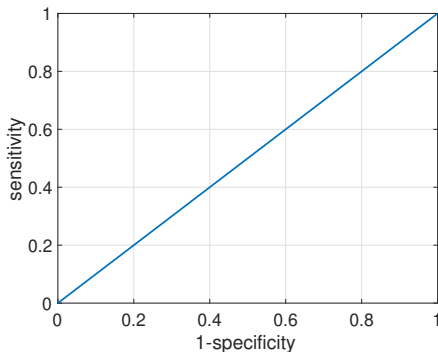


We would like the sensitivity to be close to 1 and the 1-specificity to be close to 0 (top left corner).

# The ROC plane

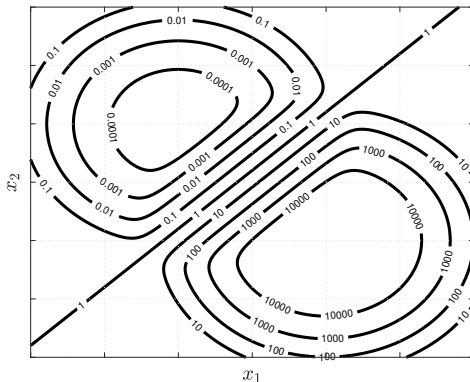
Note that we **cannot rank classifiers** using two metrics simultaneously.

The usual practice is to fix a minimum value for one of the metrics and optimise the other, for instance: obtain the highest sensitivity with a minimum specificity of 70 %.



# Back to the decision boundary

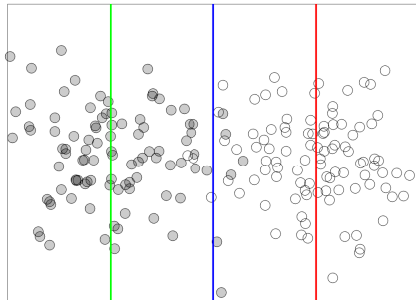
Our classifiers implement the comparison  $\frac{P(y=\text{red}|\mathbf{x})}{P(y=\text{blue}|\mathbf{x})} \lesseqgtr T$ . For each value of  $T$  we have a different boundary that defines a different classifier.



Hence,  $T$  can be **calibrated** to achieve your target performance.



# Calibration, boundaries and the confusion matrix

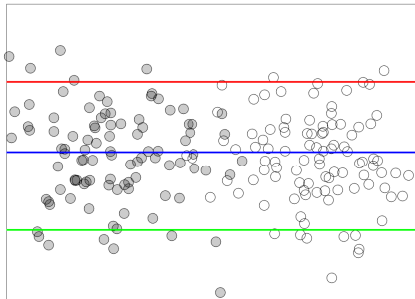


Predicted	Actual	
	1.00	0.50
	0	0.50

Predicted	Actual	
	0.95	0.05
	0.05	0.95

Predicted	Actual	
	0.50	0
	0.50	1.00

# Calibration, boundaries and the confusion matrix



Predicted	Actual	
	0.05	0.05
Predicted	0.95	0.95

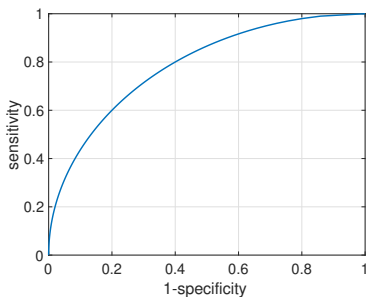
Predicted	Actual	
	0.50	0.50
Predicted	0.50	0.50

Predicted	Actual	
	0.92	0.92
Predicted	0.08	0.08

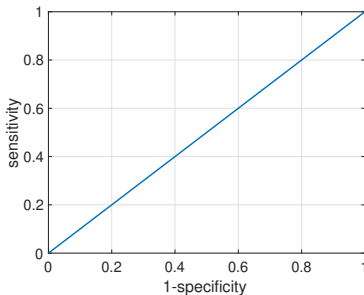
# Calibration in the ROC plane: The AUC

The **area under the curve (AUC)** is a measure of goodness for a classifier that can be calibrated.

Good classifier ( $\text{AUC} \approx 0.8$ )



Bad classifier ( $\text{AUC} = 0.5$ )



Good classifiers will have AUC close to 1, bad classifiers close to 0.5. Can you think of a classifier whose  $\text{AUC} < 0.5$ ?



# Agenda

The Bayes classifier

Using data to build posterior probabilities

Classification performance: Beyond accuracy

Summary

# The Bayes classifier

- The highest accuracy can be achieved comparing the posterior probabilities for each class and assigning a sample to the most probable class.
- The **Bayes classifier** is an ideal classifier that uses the **true posterior probabilities**.
- In general, we don't know the true posterior probabilities. In machine learning, classifiers can be seen as machines that use data to build posterior probabilities.

# Beyond accuracy

- The accuracy does not tell us how a classifier treats each class, nor accounts for different costs in misclassifying samples from each class.
- If we know the misclassification costs, we can build classifiers that **minimise the global cost**, rather than maximising the accuracy.
- We can also define class-sensitive quality metrics, using the **confusion matrix**.
- Class-sensitive quality metrics are usually **conflicting**.
- The **ROC plane** allows to explore class-sensitive quality metrics.
- We always need a **test task** to evaluate the quality of a final classifier.

# The best metric?

Don't take any quality metric for granted. Ask yourself: does this metric capture my notion of quality?

Think about the following class-sensitive classification scenarios. Which performance metrics would you use?

- Decision system in a bank offering loans.
- A security system to detect break-ins.
- A medical screening technique.
- Smoke alarm.