

# Lecture 08: How well do classifiers work – ROC and AUC

Sergei V. Kalinin

This and that

Feel free to put the text of written  
homework into Google Colabs

# We have got the data... Now what?

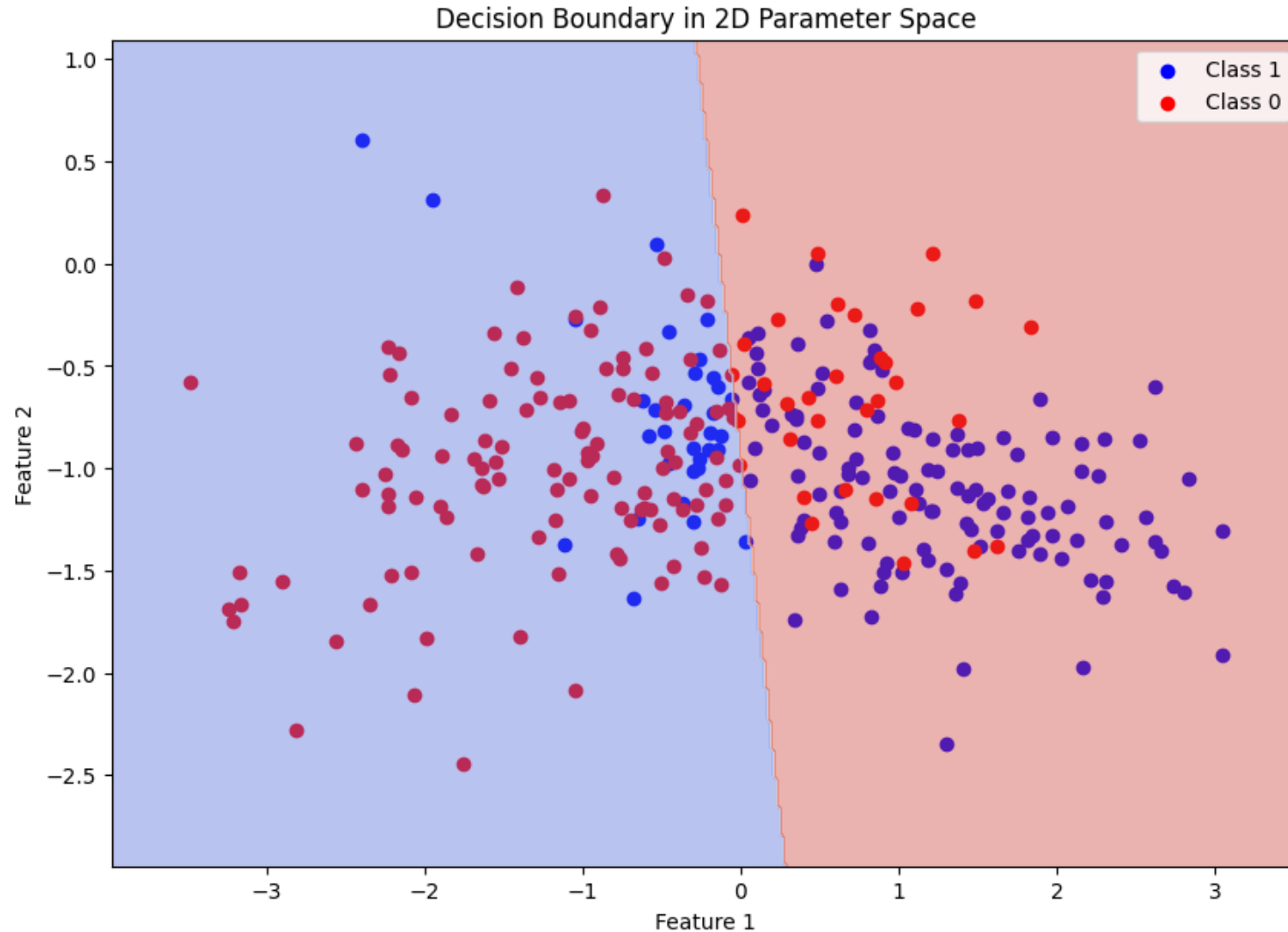
- Classification is the process of identifying and grouping objects or ideas into predetermined categories.
- In data management, classification enables the separation and sorting of data according to set requirements for various business or personal objectives.
- In machine learning (ML), classification is used in predictive modeling to assign input data with a class label.

# We have (already) several tools in the toolbox

- Classifier and Regression Trees (CART), Random Forest
- Perceptron
- Adaline
- Logistic regression

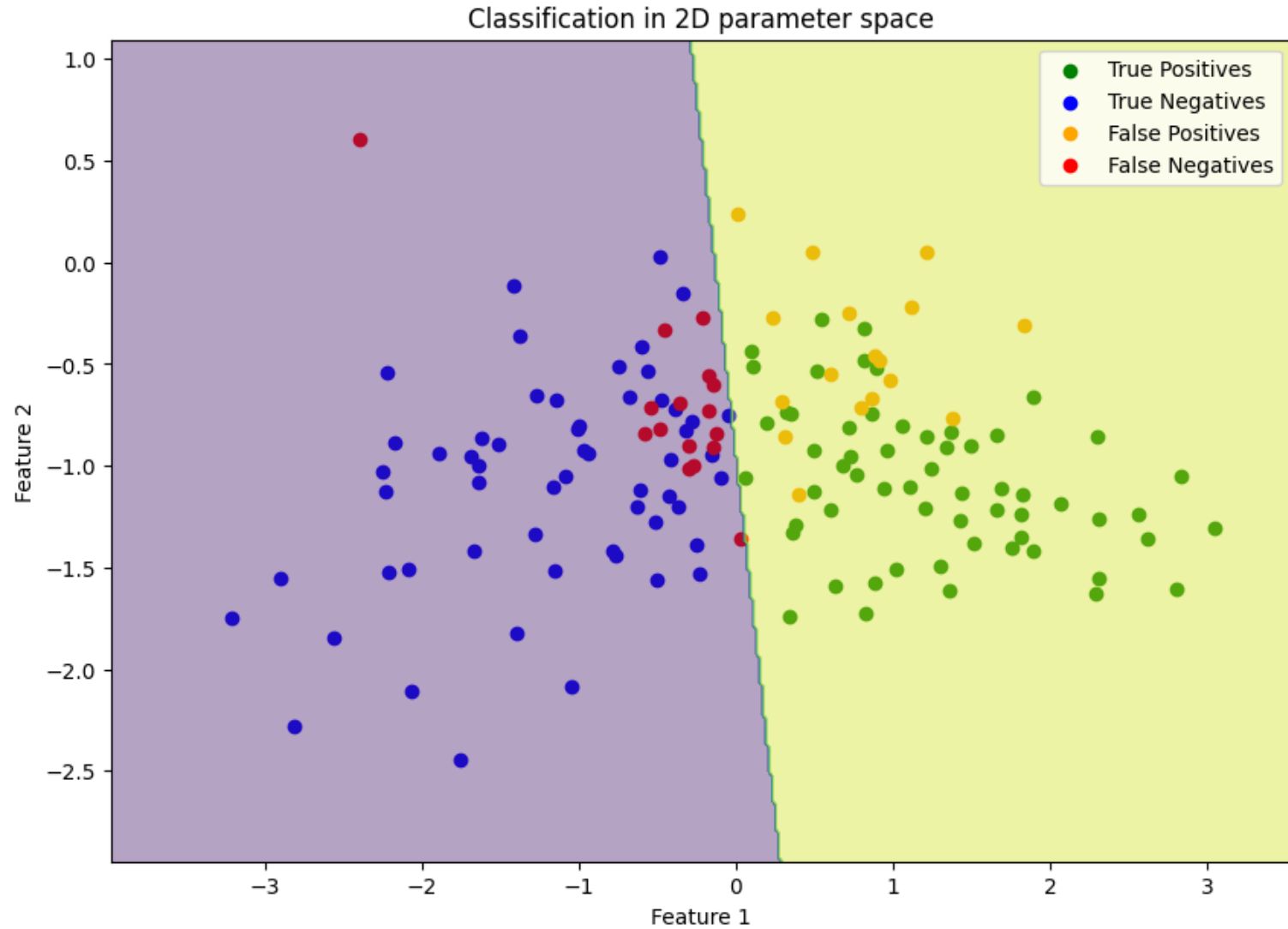
There are more classifiers in scikit-learn – and all of them use the same Pythonic methods `.fit()`, `.predict()`, and (many) yield probabilities

# Let's consider the binary classifier



- Classification is not perfect
- We have examples of Class 1 identified as 0, and Class 0 as 1

# Let's consider the binary classifier



- Does the mis-identification matter?

# It depends on the problem!

- Molecular prediction: toxicity
- Medical tests: tumor, COVID, flu, etc...
- Finance: fraud detection
- Metal detector: sensitivity threshold
- Object detection: face identification
- Quality control: ripe/rotten

# Standard assumptions

- Standard Cost Model
  - correct classification costs 0
  - cost of misclassification depends only on the class, not on the individual example
  - over a set of examples costs are additive
- Costs or Class Distributions:
  - are not known precisely at evaluation time
  - may vary with time
  - may depend on where the classifier is deployed
- True FP and TP do not vary with time or location, and are accurately estimated.



# Basic definitions

		Predicted class	
		P	N
Actual class	P	True positives (TP)	False negatives (FN)
	N	False positives (FP)	True negatives (TN)

Error:

$$ERR = \frac{FP + FN}{FP + FN + TP + TN}$$

Accuracy:

$$ACC = \frac{TP + TN}{FP + FN + TP + TN} = 1 - ERR$$

# Same definition – different names

<b>Event happens?</b>	<b>Forecast says event will happen?</b>	
	<b>Yes</b>	<b>No</b>
<b>Yes</b>	<b>Hit</b>	<b>Miss</b>
<b>No</b>	<b>False Alarm</b>	<b>Correct Rejection</b>

# Basic definitions

**False positive rate:**  $FPR = \frac{FP}{N} = \frac{FP}{FP + TN}$

**True positive rate:**  $TPR = \frac{TP}{P} = \frac{TP}{FN + TP}$

Useful for imbalanced class problems

		Predicted class	
		P	N
Actual class	P	True positives (TP)	False negatives (FN)
	N	False positives (FP)	True negatives (TN)

# Basic definitions

**Recall:**  $REC = TPR = \frac{TP}{P} = \frac{TP}{FN + TP}$

**Precision:**  $PRE = \frac{TP}{TP + FP}$

		Predicted class	
		P	N
Actual class	P	True positives (TP)	False negatives (FN)
	N	False positives (FP)	True negatives (TN)

## For medicine (cancer detection):

Optimizing for recall helps with minimizing the chance of not detecting a malignant tumor. However, this comes at the cost of predicting malignant tumors in patients although the patients are healthy (a high number of FPs).

Optimize for precision, on the other hand, we emphasize correctness if we predict that a patient has a malignant tumor. However, this comes at the cost of missing malignant tumors more frequently (a high number of FNs).

# There is always more....

**F1 Score:**

$$F1 = 2 \frac{PRE \times REC}{PRE + REC}$$

David M. W. Power *Evaluation: From Precision, Recall and F-Factor to ROC, Informedness, Markedness & Correlation*, <https://arxiv.org/abs/2010.16061>.

**Matthews correlation coefficient:**

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

D. Chicco and G. Jurman, *The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation* by *BMC Genomics*. pp. 281-305, 2012, <https://bmcbgenomics.biomedcentral.com/articles/10.1186/s12864-019-6413-7>.

From S. Raschka, Machine Learning with PyTorch and Scikit-Learn

# ... and more

Sources: [17][18][19][20][21][22][23][24][25] view · talk · edit

		Predicted condition			
		Positive (PP)	Negative (PN)	Informedness, bookmaker informedness (BM) $= \text{TPR} + \text{TNR} - 1$	Prevalence threshold (PT) $= \frac{\sqrt{\text{TPR} \times \text{FPR}} - \text{FPR}}{\text{TPR} - \text{FPR}}$
Actual condition	Total population $= P + N$				
	Positive (P)	True positive (TP), hit	False negative (FN), type II error, miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power $= \frac{\text{TP}}{P} = 1 - \text{FNR}$	False negative rate (FNR), miss rate $= \frac{\text{FN}}{P} = 1 - \text{TPR}$
	Negative (N)	False positive (FP), type I error, false alarm, overestimation	True negative (TN), correct rejection	False positive rate (FPR), probability of false alarm, fall-out $= \frac{\text{FP}}{N} = 1 - \text{TNR}$	True negative rate (TNR), specificity (SPC), selectivity $= \frac{\text{TN}}{N} = 1 - \text{FPR}$
	Prevalence $= \frac{P}{P + N}$	Positive predictive value (PPV), precision $= \frac{\text{TP}}{\text{PP}} = 1 - \text{FDR}$	False omission rate (FOR) $= \frac{\text{FN}}{\text{PN}} = 1 - \text{NPV}$	Positive likelihood ratio (LR+) $= \frac{\text{TPR}}{\text{FPR}}$	Negative likelihood ratio (LR-) $= \frac{\text{FNR}}{\text{TNR}}$
	Accuracy (ACC) $= \frac{\text{TP} + \text{TN}}{P + N}$	False discovery rate (FDR) $= \frac{\text{FP}}{\text{PP}} = 1 - \text{PPV}$	Negative predictive value (NPV) $= \frac{\text{TN}}{\text{PN}}$ $= 1 - \text{FOR}$	Markedness (MK), deltaP ( $\Delta p$ ) $= \text{PPV} + \text{NPV} - 1$	Diagnostic odds ratio (DOR) $= \frac{\text{LR}^+}{\text{LR}^-}$
	Balanced accuracy (BA) $= \frac{\text{TPR} + \text{TNR}}{2}$	$F_1$ score $= \frac{2\text{PPV} \times \text{TPR}}{\text{PPV} + \text{TPR}} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$	Fowlkes–Mallows index (FM) $= \sqrt{\text{PPV} \times \text{TPR}}$	Matthews correlation coefficient (MCC) $= \frac{\sqrt{\text{TPR} \times \text{TNR} \times \text{PPV} \times \text{NPV}}}{\sqrt{\text{FNR} \times \text{FPR} \times \text{FOR} \times \text{FDR}}}$	Threat score (TS), critical success index (CSI), Jaccard index $= \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}}$

[https://en.wikipedia.org/wiki/Receiver\\_operating\\_characteristic](https://en.wikipedia.org/wiki/Receiver_operating_characteristic)

# Limitation of the scalar measures

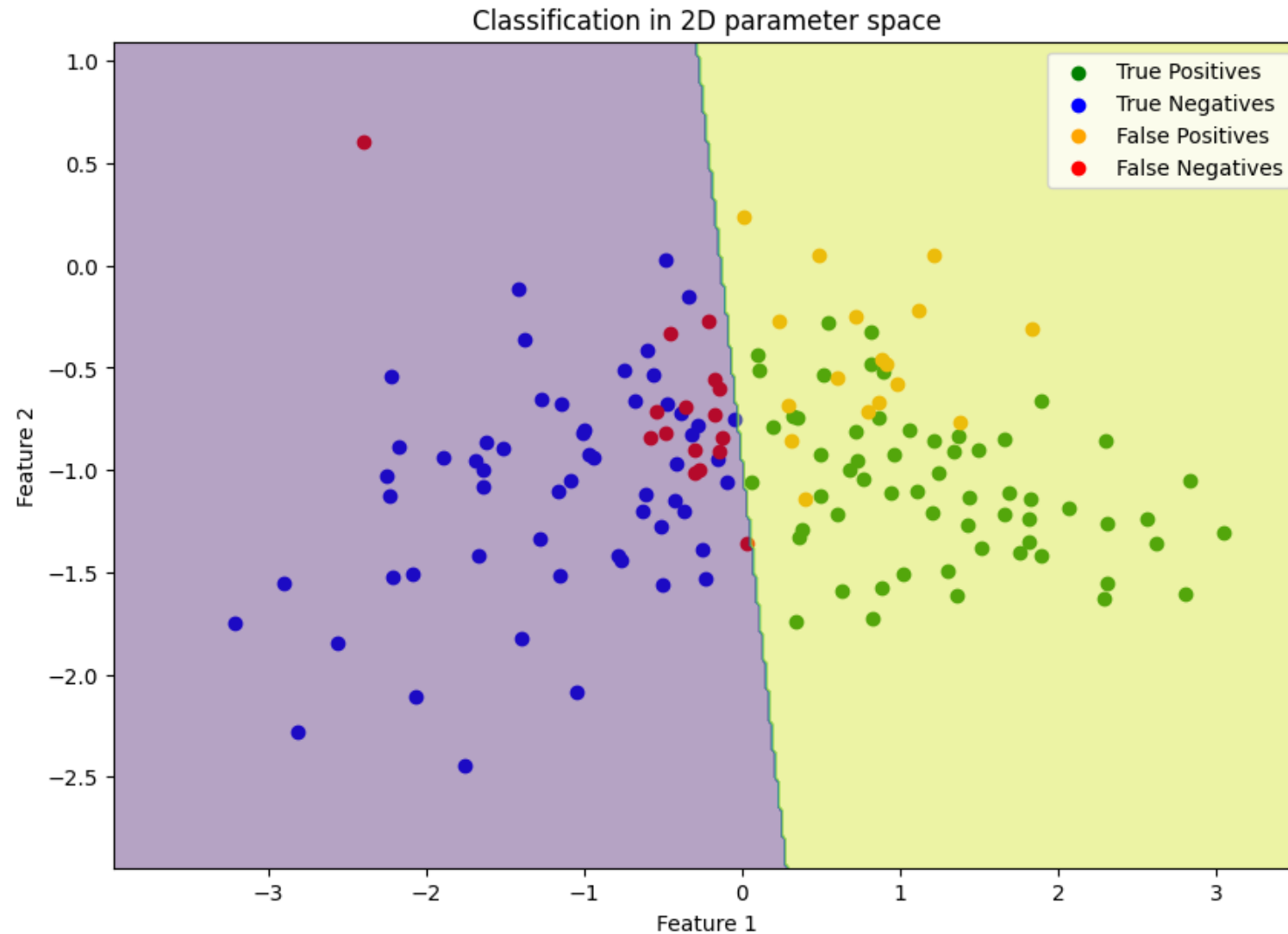
- **A scalar does not tell the whole story.**
  - There are fundamentally two numbers of interest (FP and TP), a single number invariably loses some information.
  - How are errors distributed across the classes ?
  - How will each classifier perform in different **testing** conditions (costs or class ratios other than those measured in the experiment) ?
- **A scalar imposes a linear ordering on classifiers.**
  - what we want is to identify the conditions under which each is better.

# Limitations of Scalar Measures

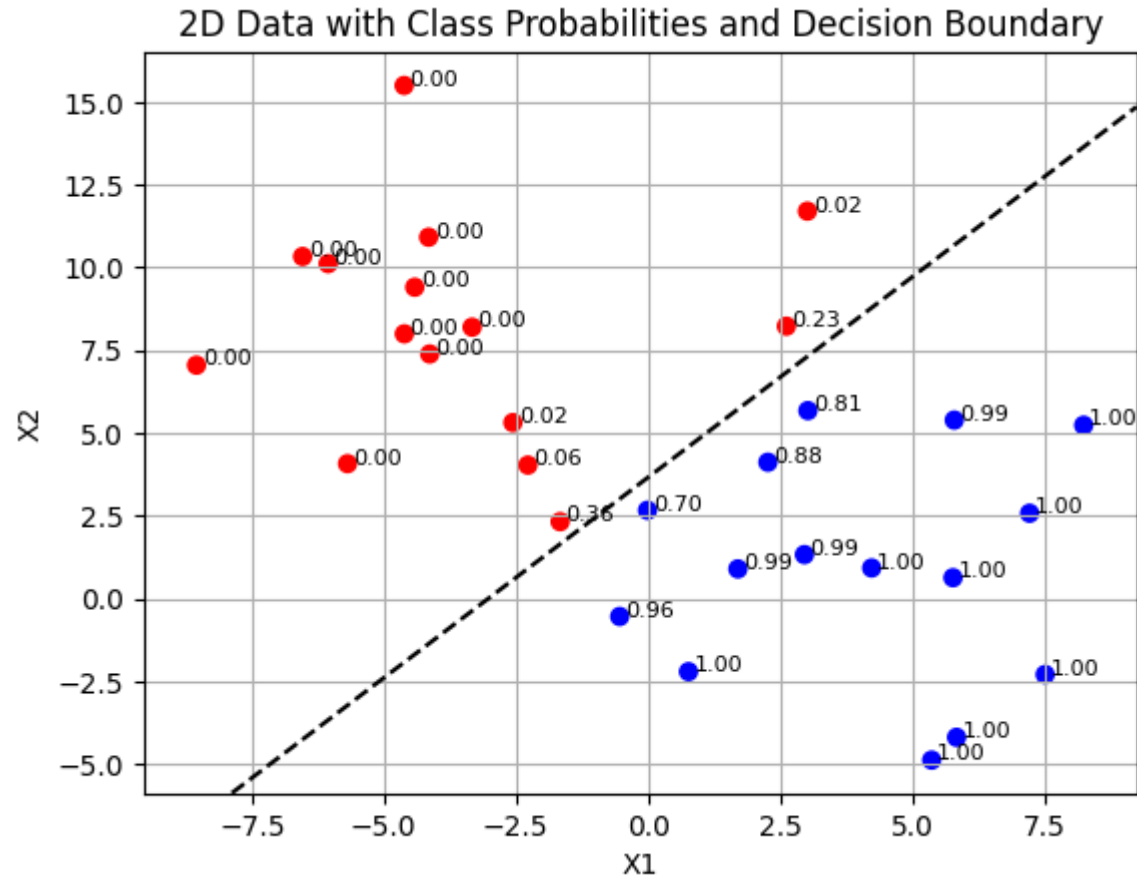
- **A table of scalars is just a mass of numbers.**
  - No immediate impact
  - Poor way to present results in a paper
  - Equally poor way for an experimenter to analyze results
- **Some scalars (accuracy, expected cost) require precise knowledge of costs and class distributions.**
  - Often these are not known precisely and might vary with time or location of deployment.



# Can we tune the classifier?



- Depending on context, cost of FP and FN errors can be different
- Can we tune the classifiers to match the business case?



**1. Train the Classifier** on train data

**2. Get Class Probabilities** for test dataset. This will give the probability of each instance belonging to the positive class.

**3. Order Probabilities:** Sort by the predicted probabilities of belonging to the positive class.

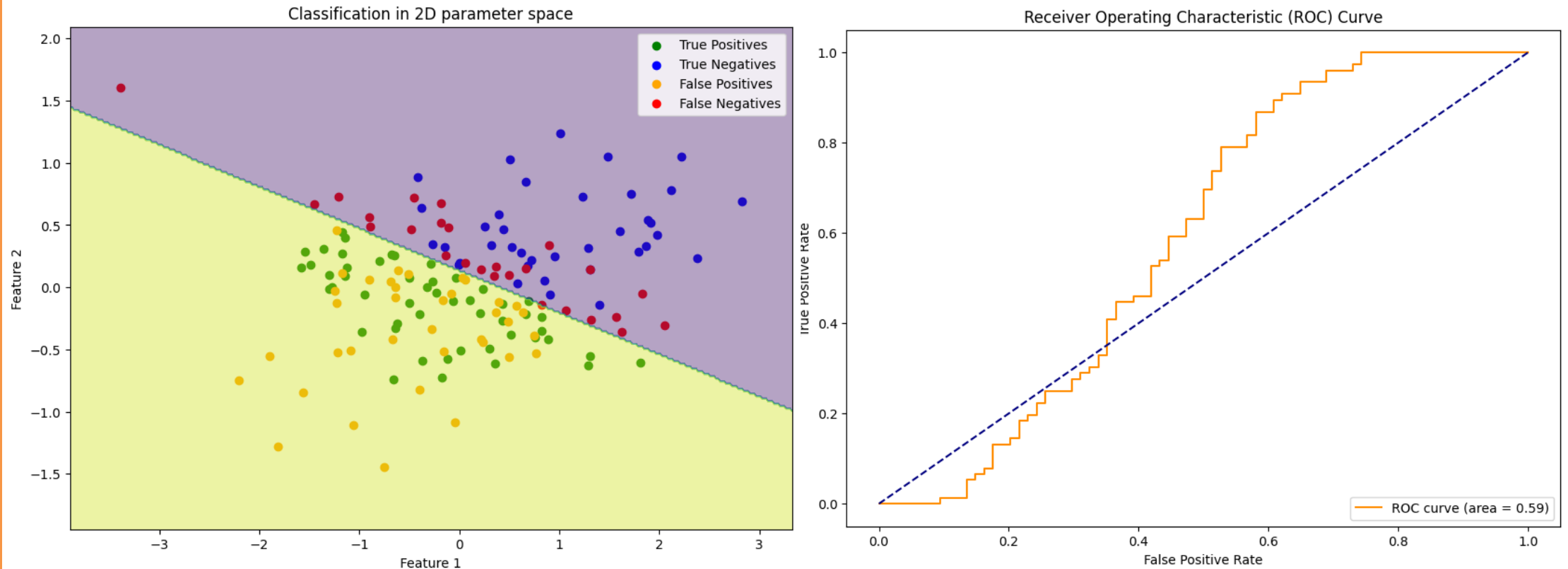
**4. Calculate TPR and FPR for Each Threshold:** For every unique probability value, treat it as a threshold. Instances with probabilities above (or equal to) the threshold are classified as positive, and those below are classified as negative. For each threshold, calculate TPR and FPR.

**5. Plot ROC Curve:** Plot TPR vs. FPR for each threshold, resulting in the ROC curve.

**6. Calculate AUC:** Compute the Area Under the Curve (AUC)

# Receiver-Operator Characteristics (ROC)

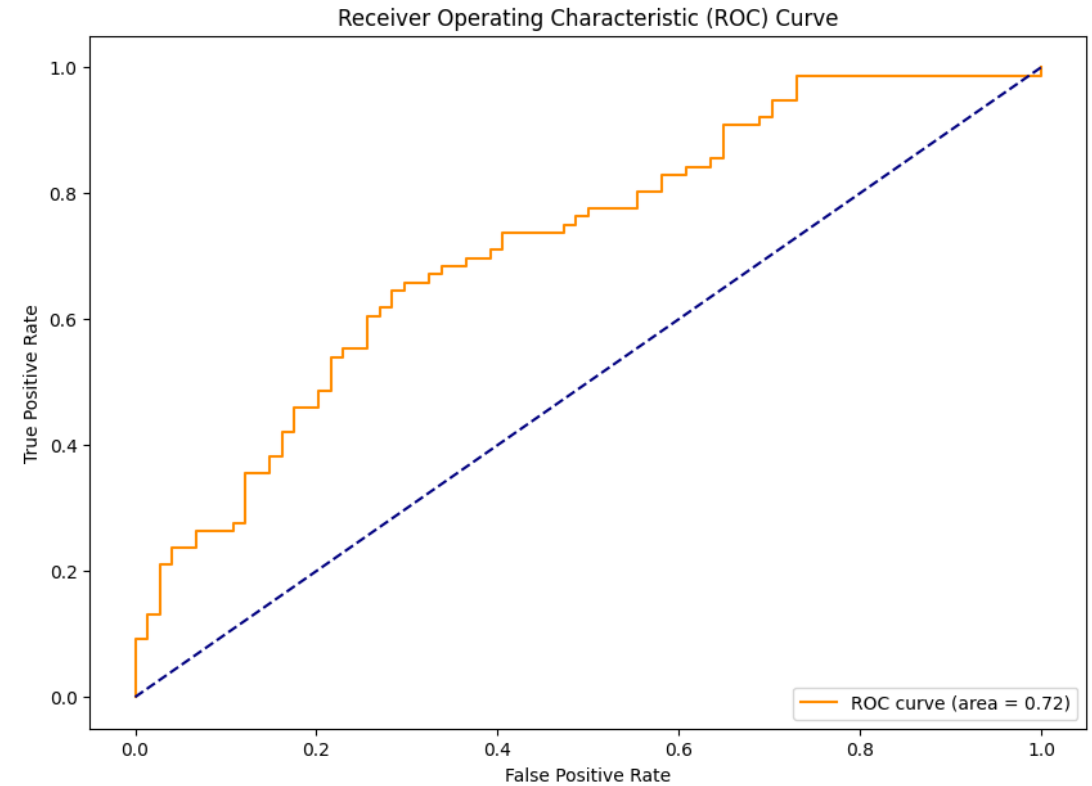
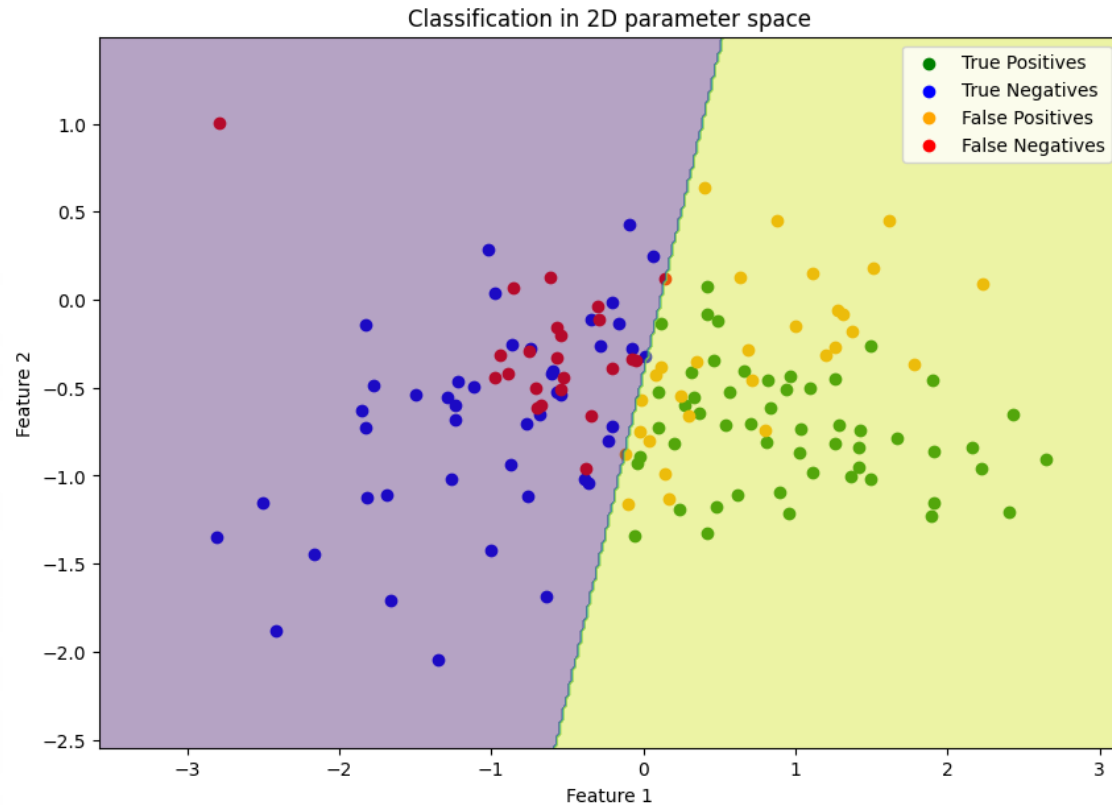
Class separation o



The Receiver Operating Characteristic (ROC) curve is a graphical representation used to assess the performance of a binary classifier system. It is a plot of the true positive rate (sensitivity) against the false positive rate (1-specificity) for different decision thresholds.

# Receiver-Operator Characteristics (ROC)

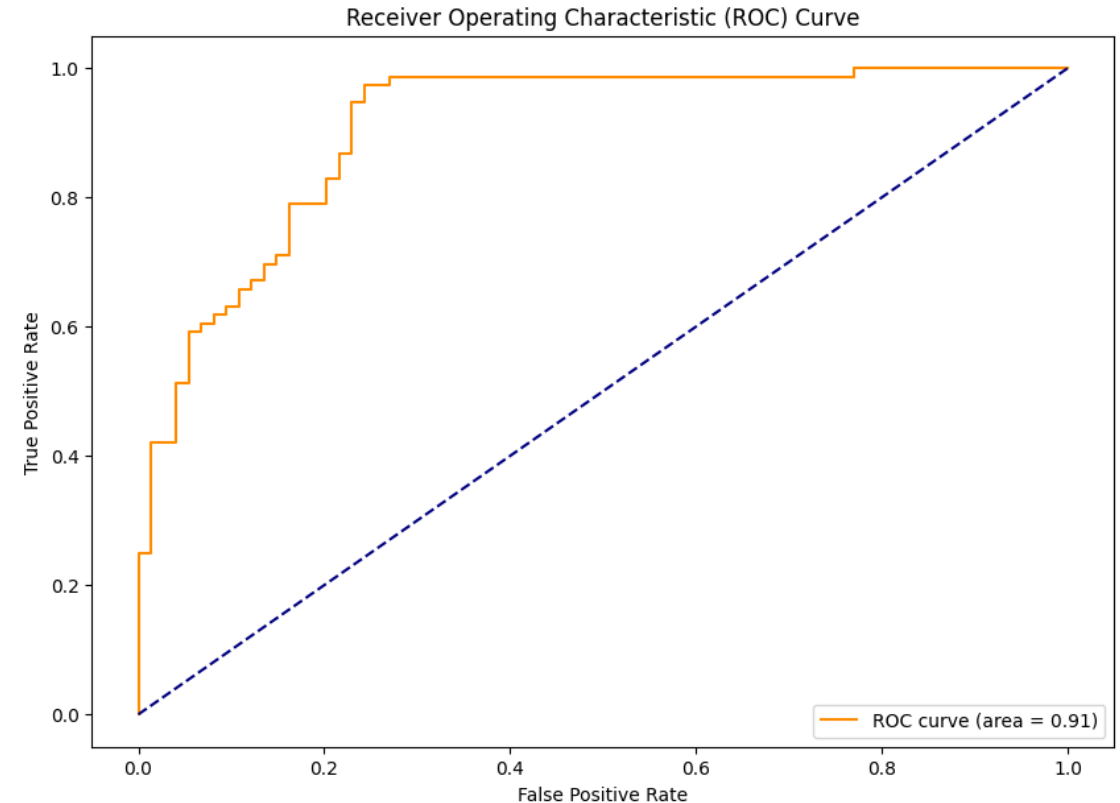
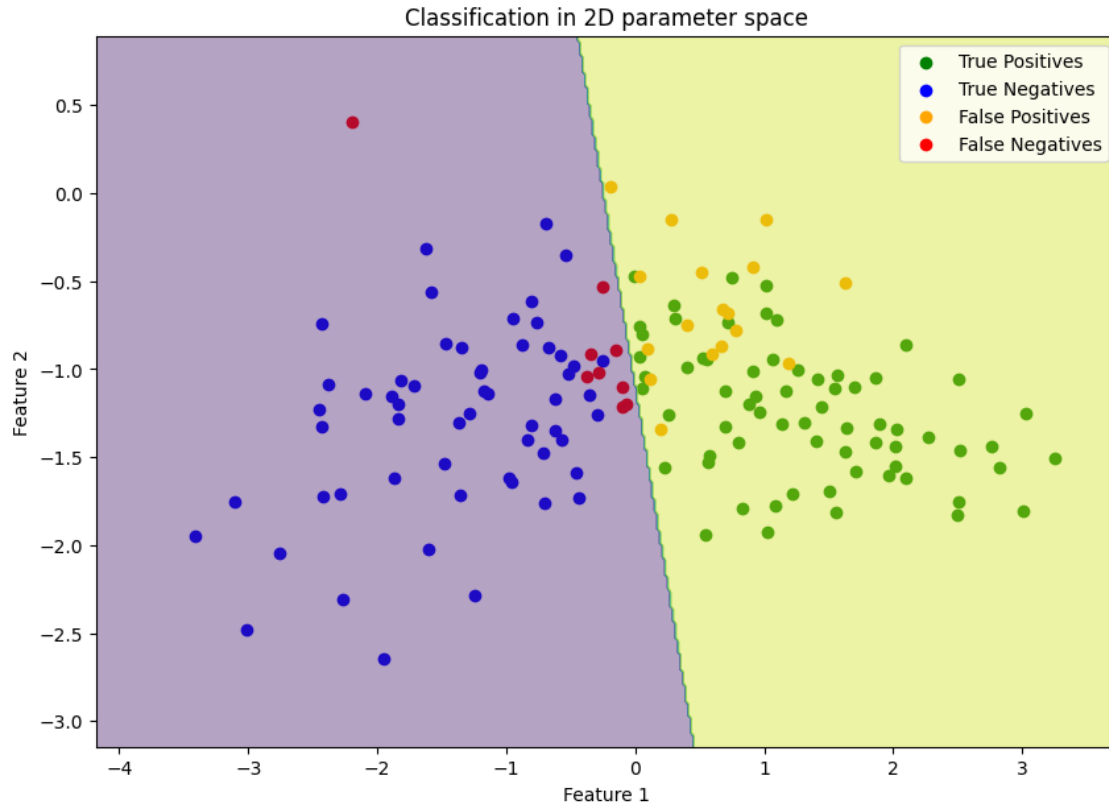
Class separation 0.6



The Receiver Operating Characteristic (ROC) curve is a graphical representation used to assess the performance of a binary classifier system. It is a plot of the true positive rate (sensitivity) against the false positive rate (1-specificity) for different decision thresholds.

# Receiver-Operator Characteristics (ROC)

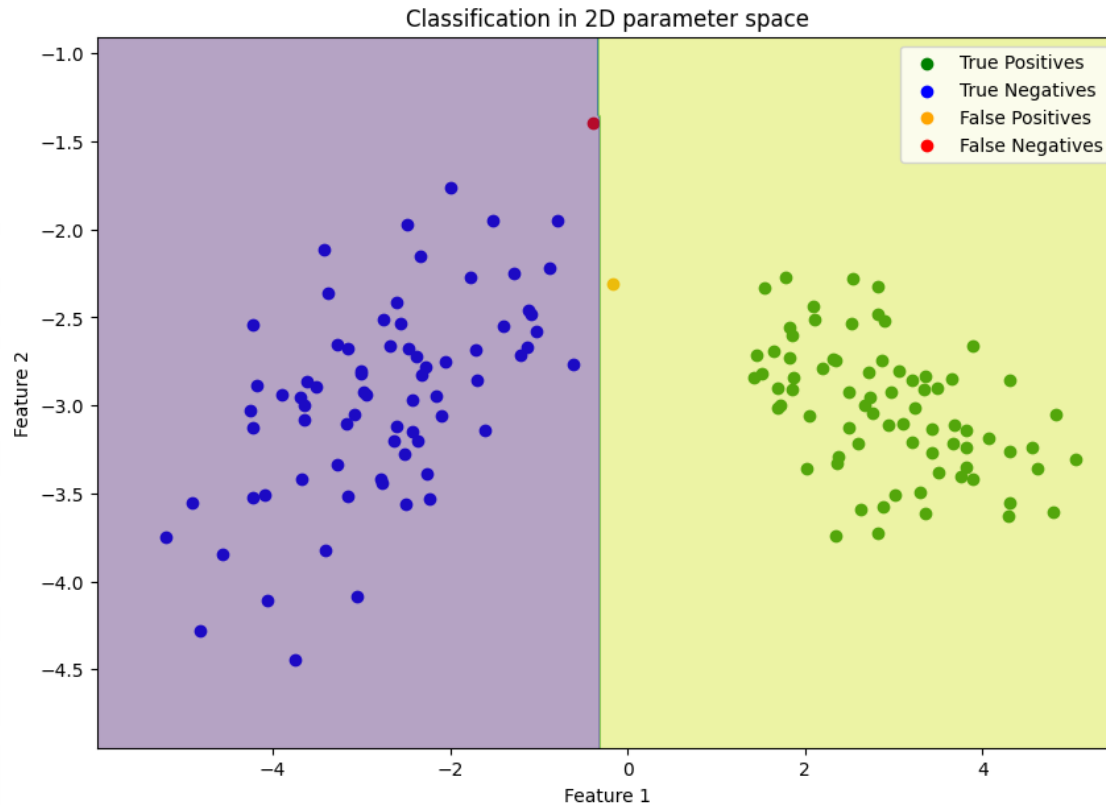
Class separation 1.2



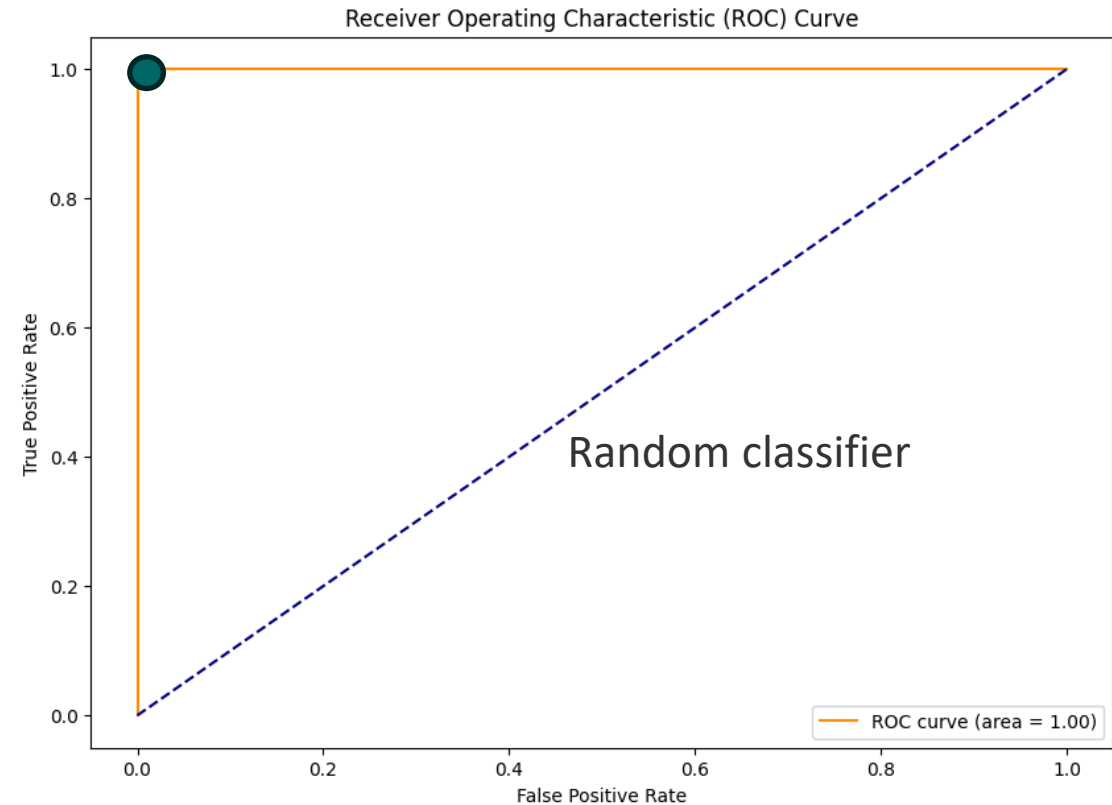
The Receiver Operating Characteristic (ROC) curve is a graphical representation used to assess the performance of a binary classifier system. It is a plot of the true positive rate (sensitivity) against the false positive rate (1-specificity) for different decision thresholds.

# Receiver-Operator Characteristics (ROC)

## Class separation 3



## Ideal classifier



The Receiver Operating Characteristic (ROC) curve is a graphical representation used to assess the performance of a binary classifier system. It is a plot of the true positive rate (sensitivity) against the false positive rate (1-specificity) for different decision thresholds.

# Advantage of ROC approach:

- A classifier produces a single ROC point.
- If the classifier has a “sensitivity” parameter, varying it produces a series of ROC points (confusion matrices).
- Alternatively, if the classifier is produced by a learning algorithm, a series of ROC points can be generated by varying the class ratio in the training set.
- ROC curves allow:
  - Visualization of all tradeoffs a model can make
  - Comparison of models across types of tradeoffs
- AUC – an aggregate measure of quality across tradeoffs
- Operating points are the tradeoff selected for specific application

# From ROC curves to Business models

$$Y = FN \cdot X + FP \cdot (1 - X)$$

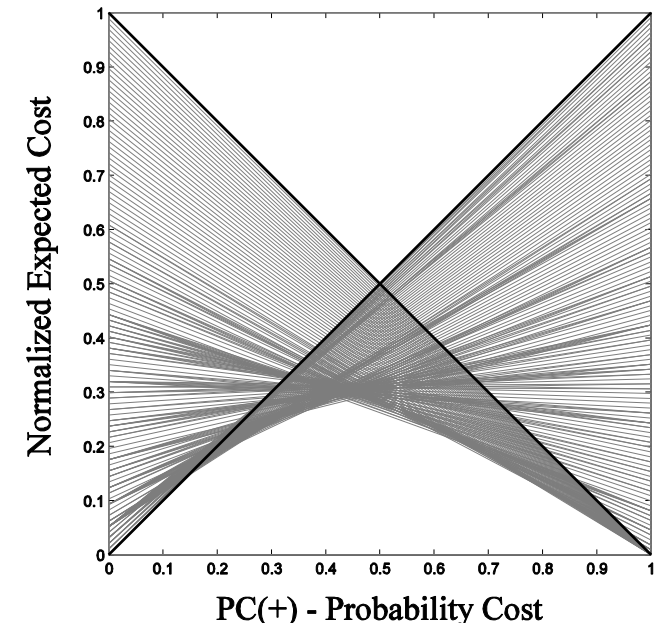
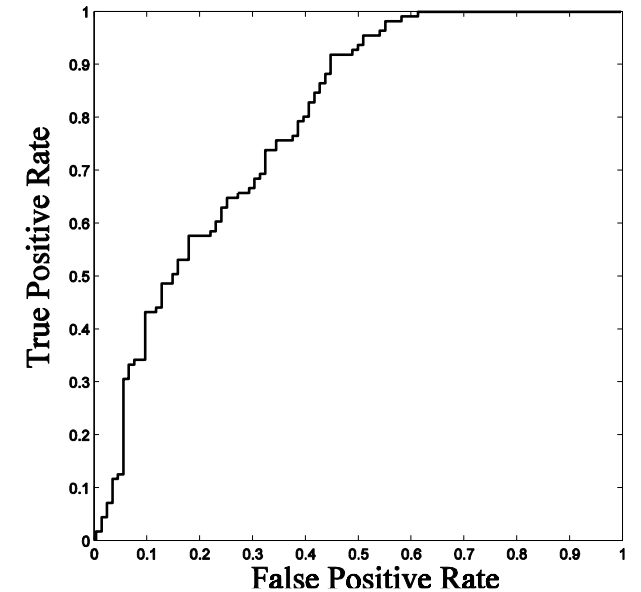
So far,  $X = p(+)$ , making  $Y = \text{error rate}$

$$X = \frac{p(+) \cdot C(-|+)}{p(+) \cdot C(-|+) + (1 - p(+)) \cdot C(+|-)}$$

$Y = \text{expected cost normalized to } [0,1]$

Cost curves enable easy visualization of

- Average performance (expected cost)
- Operating range
- Confidence intervals on performance
- Difference in performance and its significance





# What if we have multiple classes?

**1. One-vs-One (OvO) or One-vs-Rest (OvR) approach:** This involves decomposing the multi-class classification problem into multiple binary classification problems and plotting a ROC curve for each one.

**1. One-vs-Rest (OvR):** For  $K$  classes, you train  $K$  separate binary classifiers. For each classifier, one class is treated as positive while all others are treated as negative. We will have  $K$  separate ROC curves.

**2. One-vs-One (OvO):** For  $K$  classes, we train a binary classifier for every pair of classes. This results in  $K(K-1)/2$  classifiers and ROC curves.

For multi-class classification, we can compute the AUC for each class individually and then average them, providing a single scalar value that summarizes the performance.

# What if we have multiple classes?

## 1. Macro-averaging and Micro-averaging:

- 1. Macro-averaging:** Calculate the ROC curve for each class and then average them to get a single ROC curve.
- 2. Micro-averaging:** Aggregate the outcomes of individual classifications (over all classes) to compute the ROC curve. This means you'll combine all the true positives, false positives, true negatives, and false negatives across all classes and then compute the TPR and FPR.

# What if we have imbalanced classes?

- Assign a larger penalty to wrong predictions on the minority class. Via scikit-learn, adjusting such a penalty is as convenient as setting the *class\_weight* parameter to *class\_weight='balanced'*, which is implemented for most classifiers
- Upsampling the minority class or downsampling the majority class
- Generation of synthetic training examples