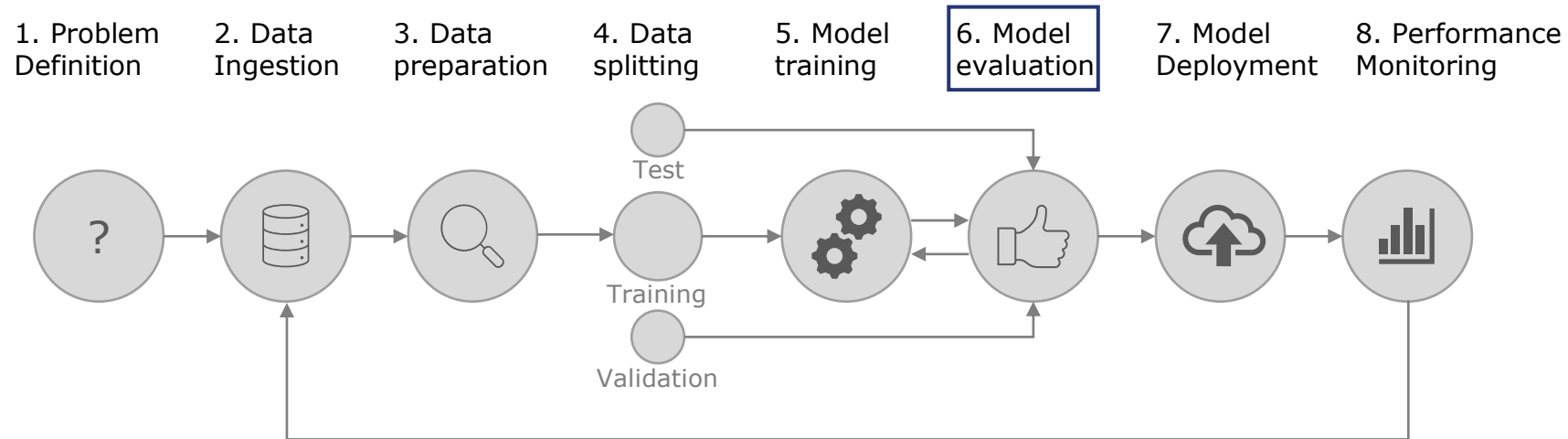


The background features a large, abstract graphic of a human head in profile, facing right. The head is composed of a blue wireframe mesh. Inside the head, there are glowing blue lines and dots, resembling a circuit board or neural network. The overall color scheme is dark blue with bright blue highlights.

# Machine Learning, Artificial Intelligence, and Big Data Analytics (IL, 4th Semester)

## Lecture 2

# Focus of this lecture



4 Data splitting= training test cross validation and separating the data

# Topics of today

- How to evaluate a ML model
- Coding:
  - Loading data
  - Manipulating with data.table
  - Data splitting and model evaluation

# Evaluating a model

What to use for evaluating a **regression** model?

MAE = Mean absolute error  
MAPE = Mean absolute percentage error  
RMSE = Root mean squared error  
R2score =

- $MAE = \frac{1}{n} \sum_{i=1}^n |e_i|$  ← Just the **average absolute error** (0 means **perfect fit**)
- $MAPE = \frac{1}{n} \sum_{i=1}^n \left| \frac{e_i}{y_i} \right|$  ← The **average error** in **relation** to the **actual values** (0% means **perfect fit**) it is not about values but percentages
- $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n |e_i|^2}$  ← The **average error** but **penalizes larger errors** more severely (0 means perfect fit)
- $R^2\text{-score} = 1 - \frac{RSS}{TSS}$  ← The **degree to which the model explains the variance in the data** (1 means **perfect** fit. 0 is **no better than the mean**. <0 is worse than the mean)
- Very easy to compute. R, Python, and Julia also provide built-in functions and usually include these metrics in the model object (from the training data).
- You should know these from the statistics lecture!
- What about **classification**?

# Evaluating a model

What to use for evaluating a **Classification** model?

- Back to the spam detection example

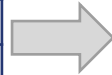
Actual	Prediction
No spam	Spam
No spam	No spam
No spam	Spam
No spam	No spam
No spam	No spam
Spam	No spam
Spam	Spam
Spam	Spam
Spam	Spam
Spam	Spam

# Evaluating a model

What to use for evaluating a **Classification** model?

- Back to the spam detection example

Actual	Prediction
No spam	Spam
No spam	No spam
No spam	Spam
No spam	No spam
No spam	No spam
Spam	No spam
Spam	Spam
Spam	Spam
Spam	Spam
Spam	Spam



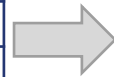
Confusion Matrix		
	Predicted	
	<b>1</b>	<b>0</b>
Actual	<b>1</b>	True positive False Negative
	<b>0</b>	False Positive True negative

# Evaluating a model

What to use for evaluating a **Classification** model?

- Back to the spam detection example

Actual	Prediction
No spam	Spam
No spam	No spam
No spam	Spam
No spam	No spam
No spam	No spam
Spam	No spam
Spam	Spam
Spam	Spam
Spam	Spam
Spam	Spam



Confusion Matrix

		Predicted	
		1	0
Actual	1	True positive	False Negative
	0	False Positive	True negative

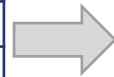
	Spam	No spam
Spam	4	1
No spam	2	3

# Evaluating a model

What to use for evaluating a **Classification** model?

- Back to the spam detection example

Actual	Prediction
No spam	Spam
No spam	No spam
No spam	Spam
No spam	No spam
No spam	No spam
Spam	No spam
Spam	Spam
Spam	Spam
Spam	Spam
Spam	Spam



Confusion Matrix

		Predicted	
		1	0
Actual	1	True positive	False Negative
	0	False Positive	True negative

	Spam	No spam
Spam	4	1
No spam	2	3

Precision = Y axis  
Recall = X axis

**Accuracy:** What fraction does it get right  
 $(\#TP + \#TN) / \#Total$

**Precision:** When it says 1 how often is it right Sensitivity  
 $\#TP / (\#TP + \#FP)$

**Recall:** What fraction of 1s does it get right Specificity  
 $\#TP / (\#TP + \#FN)$

**FP Rate:** What fraction of 0s are called 1s  
 $\#FP / (\#FP + \#TN)$

**FN Rate:** What fraction of 1s are called 0s  
 $\#FN / (\#TP + \#FN)$

**F1-Score:**  $2 * \frac{precision * recall}{precision + recall}$

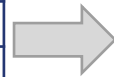


# Evaluating a model

What to use for evaluating a **Classification** model?

- Back to the spam detection example

Actual	Prediction
No spam	Spam
No spam	No spam
No spam	Spam
No spam	No spam
No spam	No spam
Spam	No spam
Spam	Spam
Spam	Spam
Spam	Spam
Spam	Spam



Confusion Matrix

		Predicted	
		1	0
Actual	1	True positive	False Negative
	0	False Positive	True negative

	Spam	No spam
Spam	4	1
No spam	2	3

**Accuracy:** What fraction does it get right  
 $(\#TP + \#TN) / \#Total = 7/10 = 70\%$

**Precision:** When it says 1 how often is it right  
 $\#TP / (\#TP + \#FP) = 4/6 = 66\%$

**Recall:** What fraction of 1s does it get right  
 $\#TP / (\#TP + \#FN) = 4/5 = 80\%$

**FP Rate:** What fraction of 0s are called 1s  
 $\#FP / (\#FP + \#TN) = 2/5 = 40\%$

**FN Rate:** What fraction of 1s are called 0s  
 $\#FN / (\#TP + \#FN) = 1/5 = 20\%$

**F1-Score:**  $2 * \frac{precision * recall}{precision + recall} = 0.72$

# Evaluating a model

The importance of looking at different metrics

- Imagine the following

Actual	Prediction
No spam	No spam
No spam	No spam
No spam	No spam
No spam	No spam
No spam	No spam
No spam	No spam
No spam	No spam
No spam	No spam
No spam	No spam
Spam	No spam
Spam	Spam

		Predicted	
		Spam	No spam
Actual	Spam	TP=1	FN=1
	No spam	FP=0	TN=8

**Accuracy:** What fraction does it get right  
 $(\#TP + \#TN) / \#Total = 9/10 = 90\%$

**Precision:** When it says 1 how often is it right  
 $\#TP / (\#TP + \#FP) = 1/1 = 100\%$

**FP Rate:** What fraction of 0s are called 1s  
 $\#FP / (\#FP + \#TN) = 0\%$

We also need to

Exam question: Given a table compute all the metrics, typical error = verify where is predicted and where is actual  
Given that a model has to detect a disease known accuracy, can you say if the model is good or bad?

# Evaluating a model

The importance of looking at different metrics

- Imagine the following

Actual	Prediction
No spam	No spam
No spam	No spam
No spam	No spam
No spam	No spam
No spam	No spam
No spam	No spam
No spam	No spam
No spam	No spam
Spam	No spam
Spam	Spam

		Predicted	
		Spam	No spam
Actual	Spam	TP=1	FN=1
	No spam	FP=0	TN=8

It says 100% but we predicted that it is SPAM only once which is small, than we need to present Recall and F1-Score

**Accuracy:** What fraction does it get right  
 $(\#TP + \#TN) / \#Total = 9/10 = 90\%$

**Precision:** When it says 1 how often is it right  
 $\#TP / (\#TP + \#FP) = 1/1 = 100\%$

**Recall:** What fraction of 1s does it get right  
 $\#TP / (\#TP + \#FN) = 1/2 = 50\%$

**FP Rate:** What fraction of 0s are called 1s  
 $\#FP / (\#FP + \#TN) = 0\%$

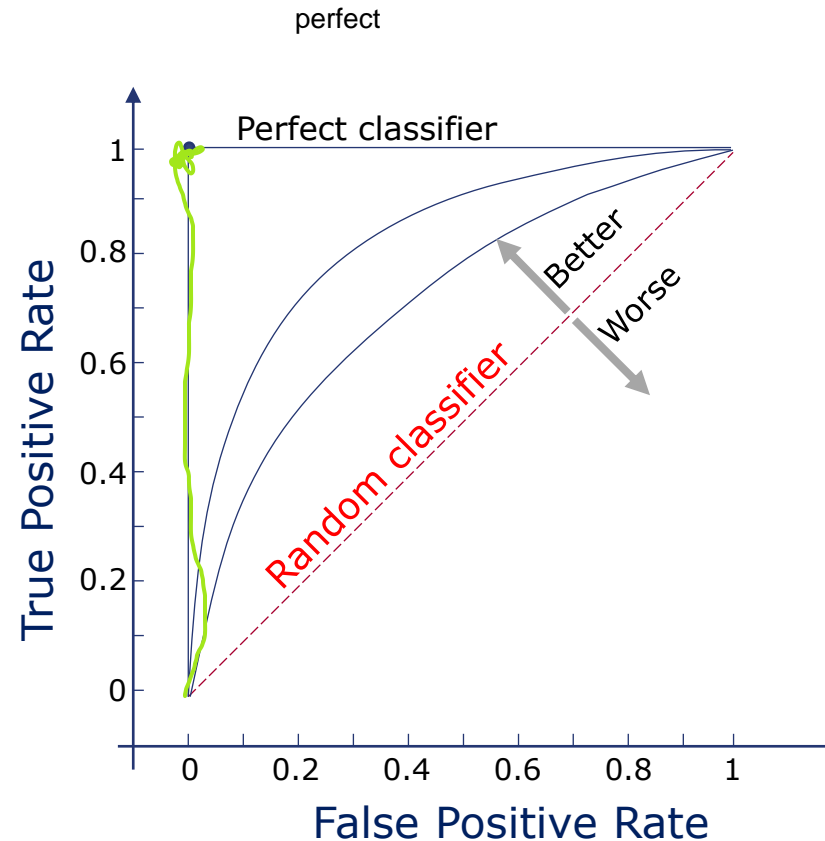
**FN Rate:** What fraction of 1s are called 0s  
 $\#FN / (\#TP + \#FN) = 1/2 = 50\%$

**F1-Score:**  $2 * \frac{precision * recall}{precision + recall} = 0.66$

Receiver operator characteristic = It plots the false positive rate (FPR), against true positive rate (TPR) X axis = TPR, Y axis = FPR

# The ROC curve and the AUC

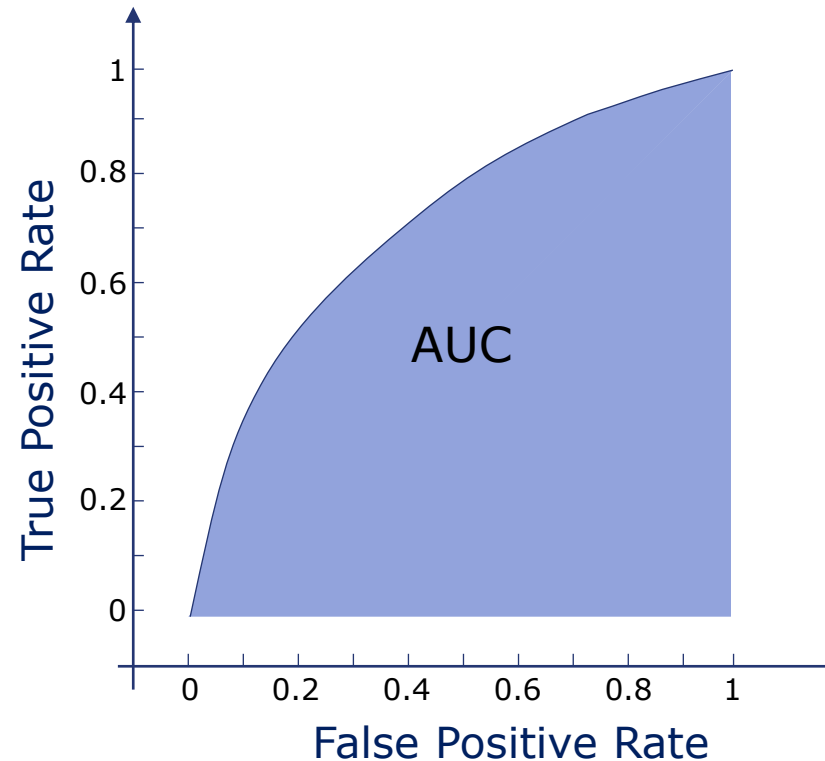
- Comparing binary classifiers
- True Positive vs. False Positive at various thresholds



AUC = area under the curve

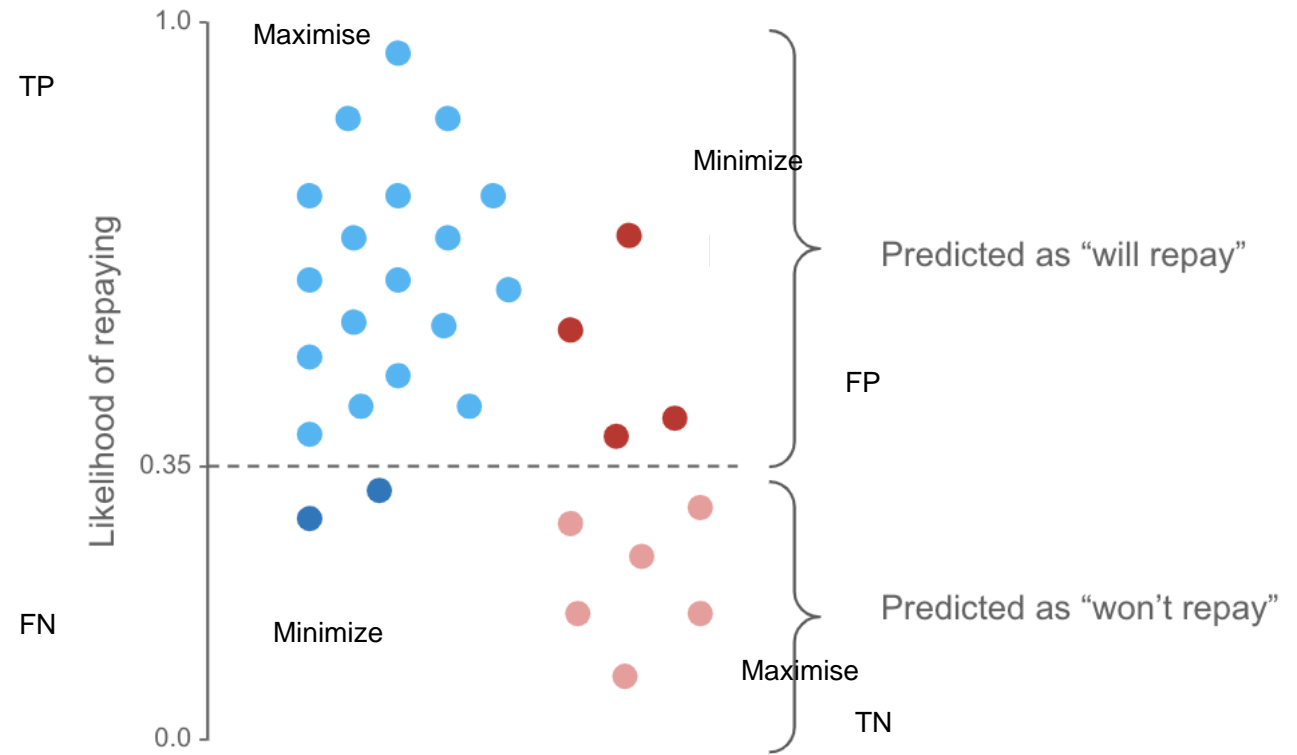
# The ROC curve and the AUC

- Comparing binary classifiers
- True Positive vs. False Positive at various thresholds
- $0 < \text{AUC} < 1$
- The larger the better



# ROC example

<https://towardsdatascience.com/understanding-the-roc-curve-in-three-visual-steps-795b1399481c>



Actual positives: users who repaid the loan

● Predicted as "will repay"

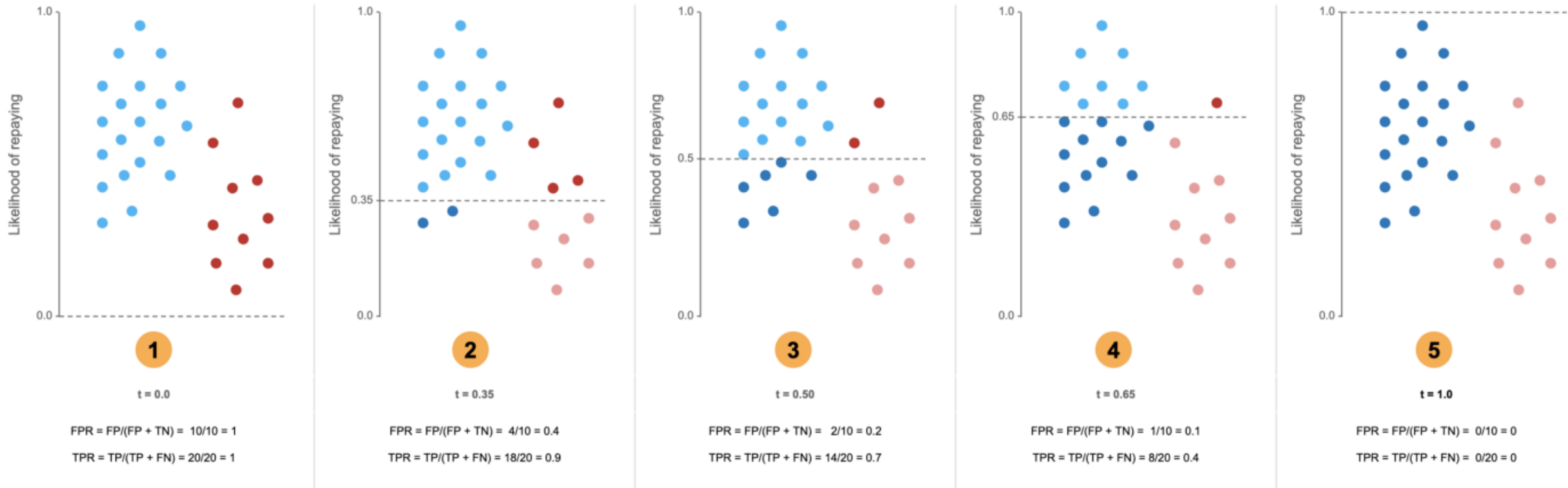
● Predicted as "won't repay"

Actual negatives: users who didn't repaid the loan

● Predicted as "won't repay"

● Predicted as "will repay"

# ROC example

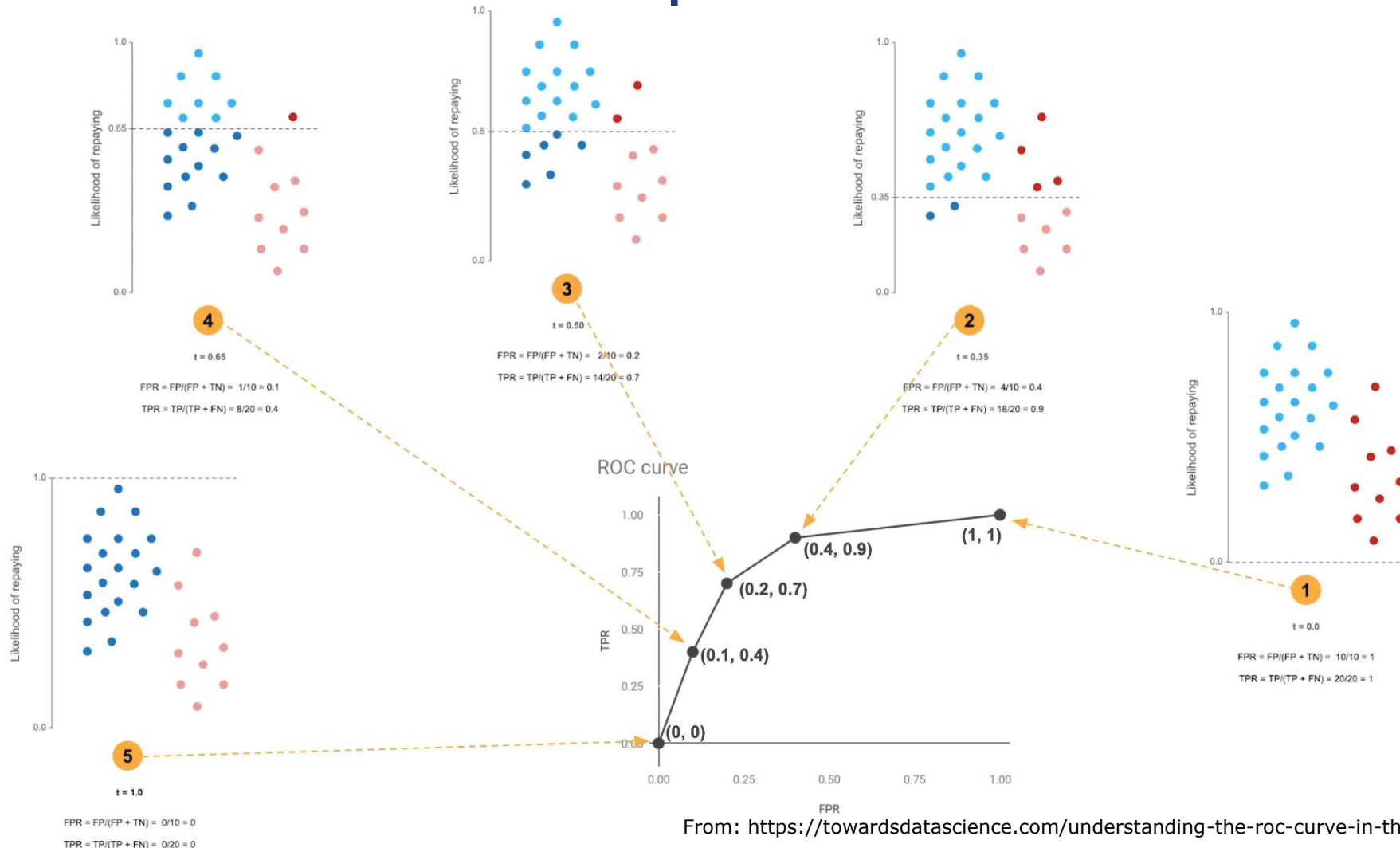


FPR and TPR decrease as the threshold gets larger

Actual positives: users who repaid the loan  
Actual negatives: users who didn't repaid the loan

● Predicted as "will repay" (TP)      ● Predicted as "won't repay" (TN)  
● Predicted as "won't repay" (FN)      ● Predicted as "will repay" (FP)

# ROC example





# Evaluating a model

## Summary

Metric	Formula	Meaning	Visual look	range								
Accuracy	(#TP+ #TN)/#Total	What fraction does it get right	<table><tr><td>TP</td><td>FN</td></tr><tr><td>FP</td><td>TN</td></tr></table> / <table><tr><td>TP</td><td>FN</td></tr><tr><td>FP</td><td>TN</td></tr></table>	TP	FN	FP	TN	TP	FN	FP	TN	0- <u>1</u>
TP	FN											
FP	TN											
TP	FN											
FP	TN											
Precision	#TP/(#TP+ #FP)	When it says 1 how often is it right	<table><tr><td>TP</td><td>FN</td></tr><tr><td>FP</td><td>TN</td></tr></table> / <table><tr><td>TP</td><td>FN</td></tr><tr><td>FP</td><td>TN</td></tr></table>	TP	FN	FP	TN	TP	FN	FP	TN	0- <u>1</u>
TP	FN											
FP	TN											
TP	FN											
FP	TN											
Recall/ Sensitivity	#TP/(#TP+ #FN)	What fraction of 1s does it get right (True Positive Rate – TPR)	<table><tr><td>TP</td><td>FN</td></tr><tr><td>FP</td><td>TN</td></tr></table> / <table><tr><td>TP</td><td>FN</td></tr><tr><td>FP</td><td>TN</td></tr></table>	TP	FN	FP	TN	TP	FN	FP	TN	0- <u>1</u>
TP	FN											
FP	TN											
TP	FN											
FP	TN											
Specificity	#TN/(#TN+ #FP)	What fraction of 0s does it get right (True Negative Rate – TNR)	<table><tr><td>TP</td><td>FN</td></tr><tr><td>FP</td><td>TN</td></tr></table> / <table><tr><td>TP</td><td>FN</td></tr><tr><td>FP</td><td>TN</td></tr></table>	TP	FN	FP	TN	TP	FN	FP	TN	0- <u>1</u>
TP	FN											
FP	TN											
TP	FN											
FP	TN											
FP Rate	#FP/(#FP+ #TN)	What fraction of 0s are called 1s	<table><tr><td>TP</td><td>FN</td></tr><tr><td>FP</td><td>TN</td></tr></table> / <table><tr><td>TP</td><td>FN</td></tr><tr><td>FP</td><td>TN</td></tr></table>	TP	FN	FP	TN	TP	FN	FP	TN	<u>0</u> -1
TP	FN											
FP	TN											
TP	FN											
FP	TN											
FN Rate	#FN/(#TP+ #FN)	What fraction of 1s are called 0s	<table><tr><td>TP</td><td>FN</td></tr><tr><td>FP</td><td>TN</td></tr></table> / <table><tr><td>TP</td><td>FN</td></tr><tr><td>FP</td><td>TN</td></tr></table>	TP	FN	FP	TN	TP	FN	FP	TN	<u>0</u> -1
TP	FN											
FP	TN											
TP	FN											
FP	TN											
F1-score	$2 * \frac{precision*recall}{precision+recall}$	How “good” are precision and recall		0- <u>1</u>								

# Things you should know

100% will be in exam

- What is underfit/overfit. What is the bias-variance tradeoff. How do they relate?
- How does cross-validation work.
- What is bootstrapping and bagging.
- How to evaluate a regression or a classification model
  - RMSE, MAE, ...
  - Accuracy, Precision, Recall, ...
  - Interpret a ROC curve

bootstrapping = sampling data with replacement from a given dataset

Bagging = training model on different parts of data

- It continues in R

# Feedback round

- Scan the barcode from your mobile phone

**OR**

- go to <http://sli.do> and insert this code: **19651**

and follow my instructions.

# Exercise 1

## Overfit

1. **Load** the dataset wines.csv (or any other regression dataset from [here](#) – i.e., Regression task, numerical variables)
  2. **Explore and visualize the dataset** (e.g., how many observations? How many features? Missing values? Are some features irrelevant?)
  3. Create a regression model (i.e., for wine: the quality by using density, chlorides, and volatile acidity).
    1. Split the data into training and test set
    2. Create a linear regression model and polynomial models with increasing degree.
    3. What's the MAE, the RMSE, and the MAPE on the training and test set for all the models?
    4. When does the model start overfitting? Which degree would you choose?
- **Due date: March 12<sup>th</sup>, 23.59 CET (Late submission +1week, 6 pts)**
  - Comment code and results (or write a notebook).
  - Use R or Python