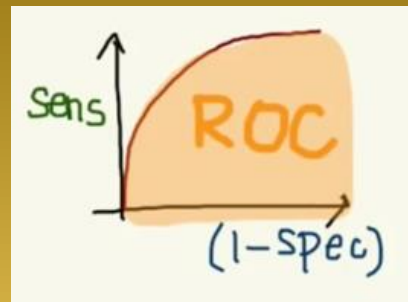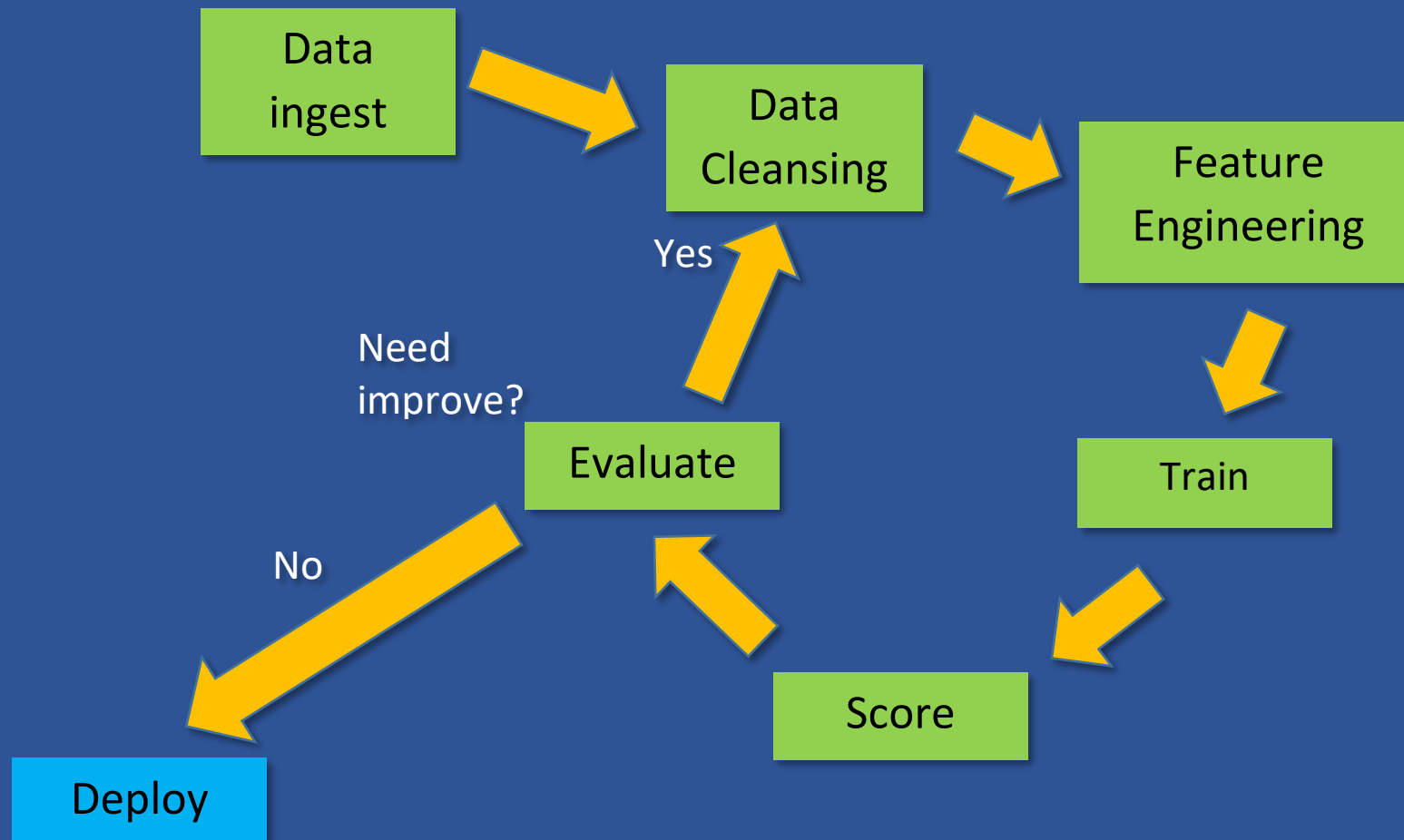# ML EVALUATION

# In this session

- ML train and evaluation circle
- How to read Histogram
- How to read Box Plot
- Adding Evaluate Model
- How to read ROC curve
- Area Under the Curve (AUC)
- How to read Evaluation metrics

## ML Evaluation
# ML evaluation circle

Data ingest → Data Cleansing → Feature Engineering

Feature Engineering → Train

Train → Score

Score → Evaluate

Evaluate → **Yes** → Data Cleansing

**Need improve?**

Evaluate → **No** → Deploy

## ML Evaluation
# How to read Scoring results



- This table = Scored dataset
- Row = 267 / Columns = 10
- Total column = 10 / Left 8 = features / Right 2 = prediction results
- Scored Label 0 = dead 1 = survived
- Scored Probabilities (SP)  SP <=0.5 == dead /  SP > 0.5 == survived

## ML Evaluation
## How to read Scoring Statistics

**Statistics**

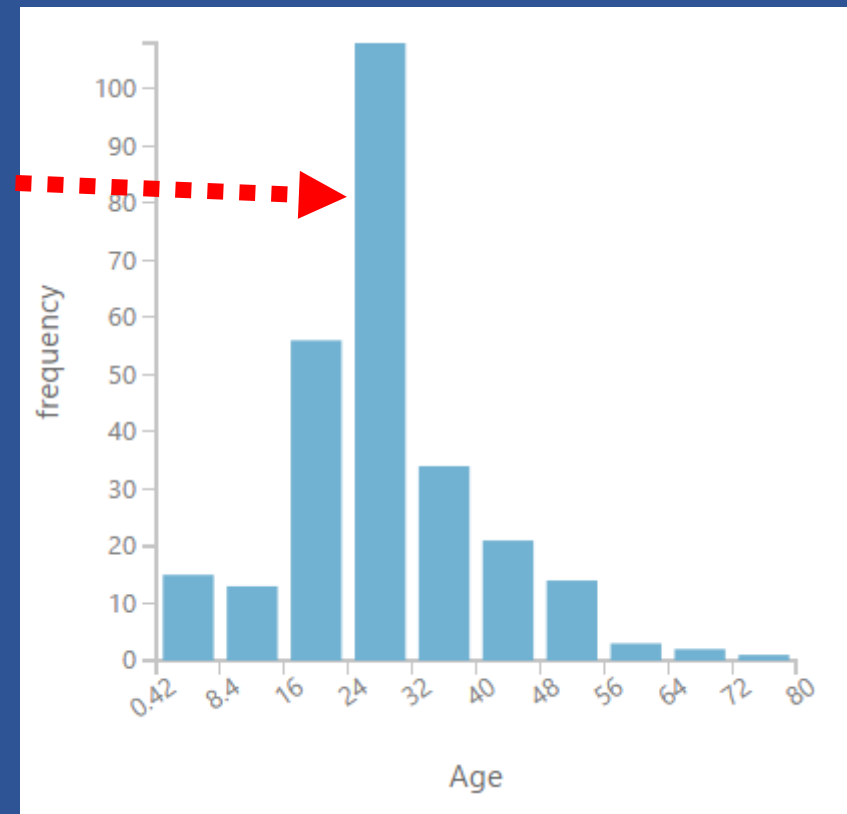| | |
|---|---|
| Mean | 28.8265 |
| Median | 28 |
| Min | 0.42 |
| Max | 80 |
| Standard Deviation | 12.3791 |
| Unique Values | 61 |
| Missing Values | 0 |
| Feature Type | Numeric Feature |

Show Statistics of the Scored dataset

- Mean = Sum of all the values divided by the number of values
- Median = The midpoint of the data after being ranked
- Standard Deviation = The square root of the variance
- Unique Values
- Missing Value

## ML Evaluation
# How to read Score Histogram

## Histogram

- Representation: distribution of numerical data
- Bin: series of intervals (bin)
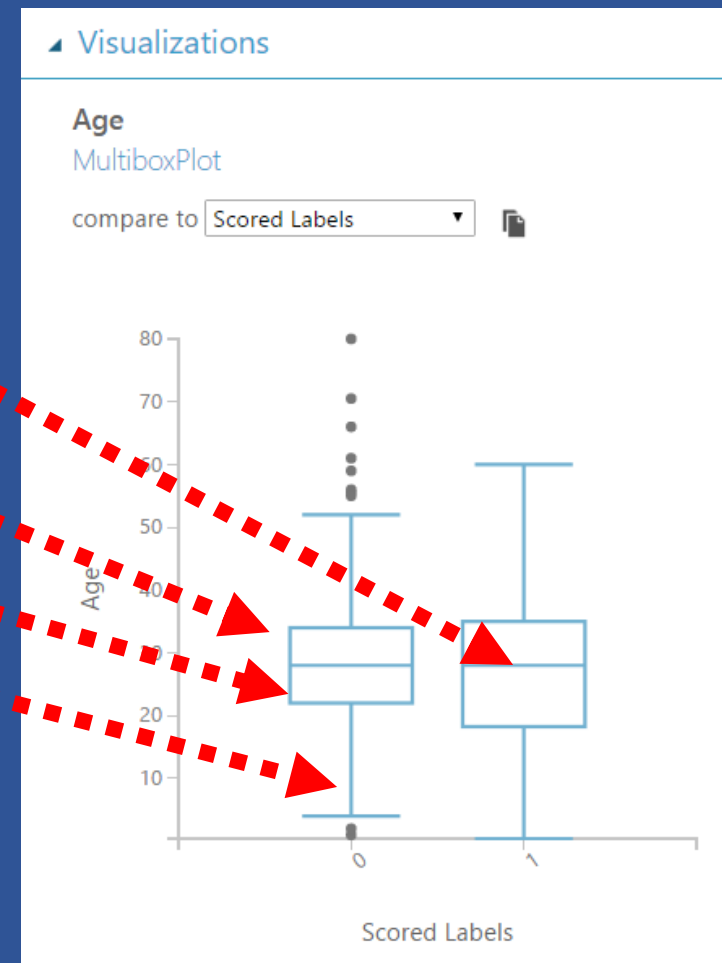- Count: values fall into each interval

ML Evaluation
# How to read Box Plot

## Box Plot

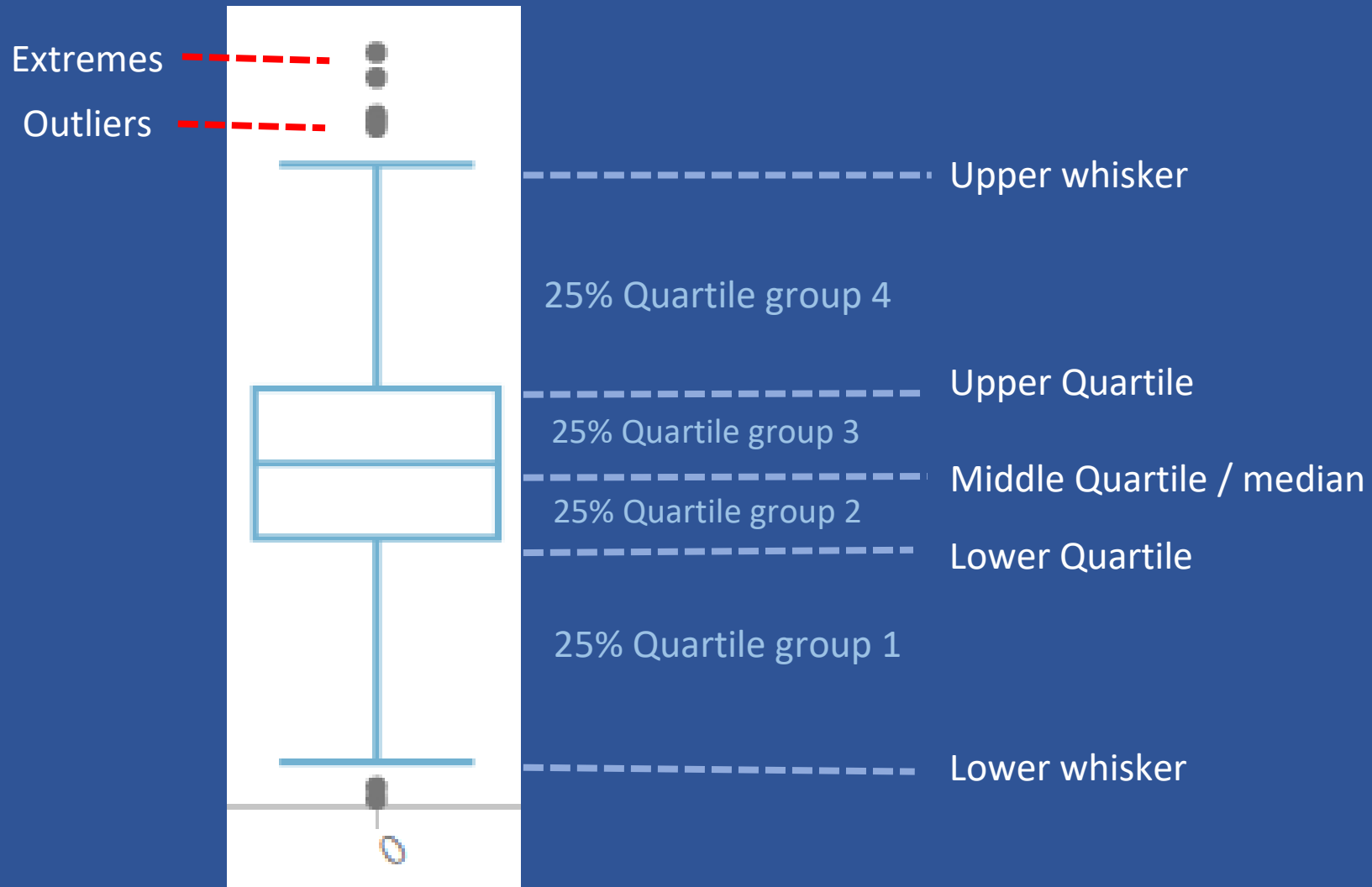Box Plot (whisker) is a standardized way of displaying the distribution of data

- **Median:** marks the mid-point of the data
- **Box:** middle 50% of scores for the group.
- **Upper quartile:** 75% of the scores fall below the upper quartile.
- **Lower quartile:** 25% of scores fall below the lower quartile.
- **Whiskers:** scores outside the middle 50%
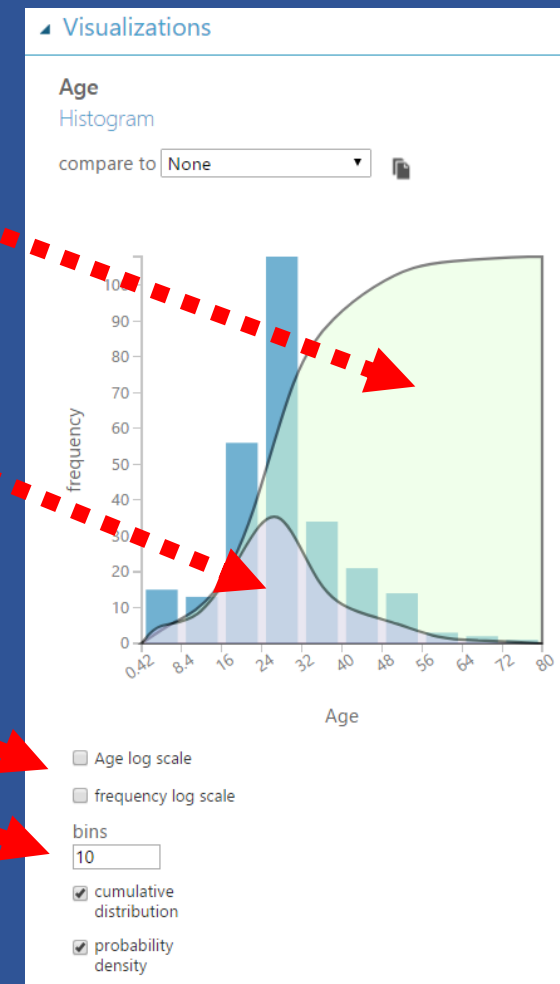
0 = dead

ML Evaluation
# Box Plot Definitions

Extremes

Outliers

Upper whisker

25% Quartile group 4

Upper Quartile

25% Quartile group 3

Middle Quartile / median

25% Quartile group 2

Lower Quartile

25% Quartile group 1

Lower whisker

**GreatFriends.Biz**

ML Evaluation
# Histogram option

## Histogram options

- **Cumulative distribution function (cdf):** shows "How common are samples that are *less than or equal* to this value?"
- **Probability density function (pdf):** shows "How common are samples at exactly this value?"
- **Scale:** scaling the distribution
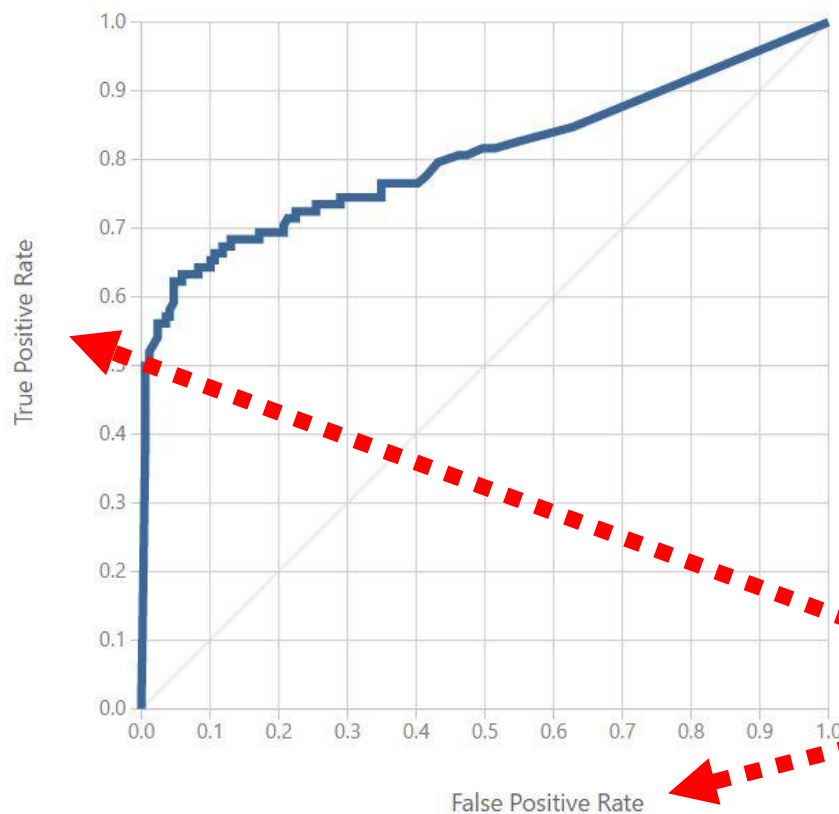- **bins:** number of bin

## ML Evaluation
# Receiver Operating Characteristic (ROC) Curve



Titanic Evaluate > Evaluate Model > Evaluation results

ROC PRECISION/RECALL LIFT

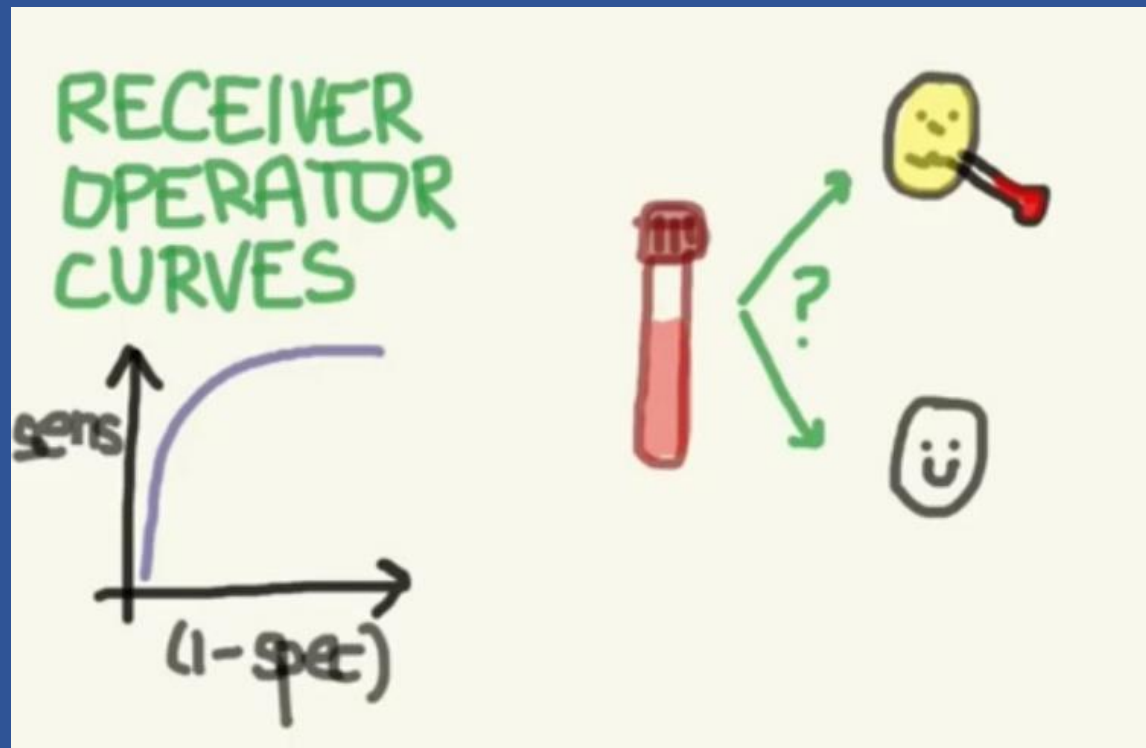| | | | |
|---|---|---|---|
| True Positive | False Negative | Accuracy | Precision |
| 64 | 34 | 0.805 | 0.780 |
| False Positive | True Negative | Recall | F1 Score |
| 18 | 151 | 0.653 | 0.711 |
| Positive Label | Negative Label | | |
| 1 | 0 | | |

Threshold
0.5

AUC
0.817

True Positive Rate (TPR)

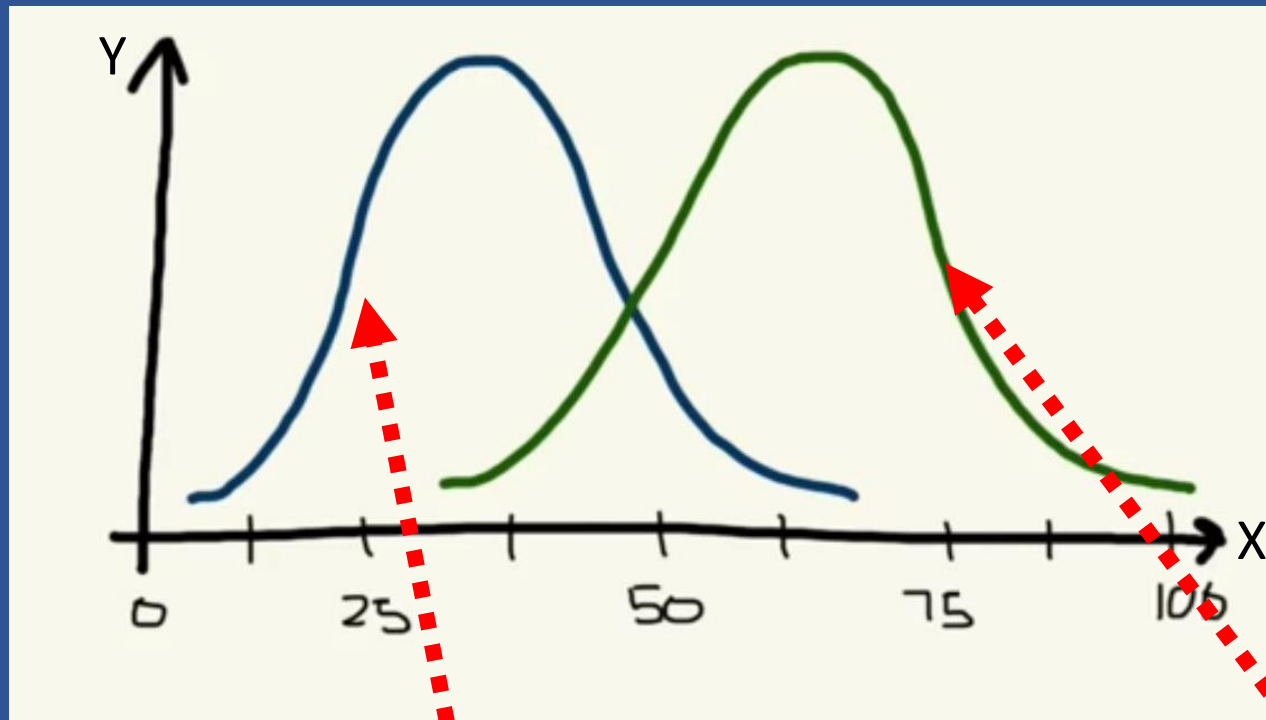False Positive Rate (FPR)

## ML Evaluation
## How to read ROC curve

ROC curve is a graphical plot that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is varied.



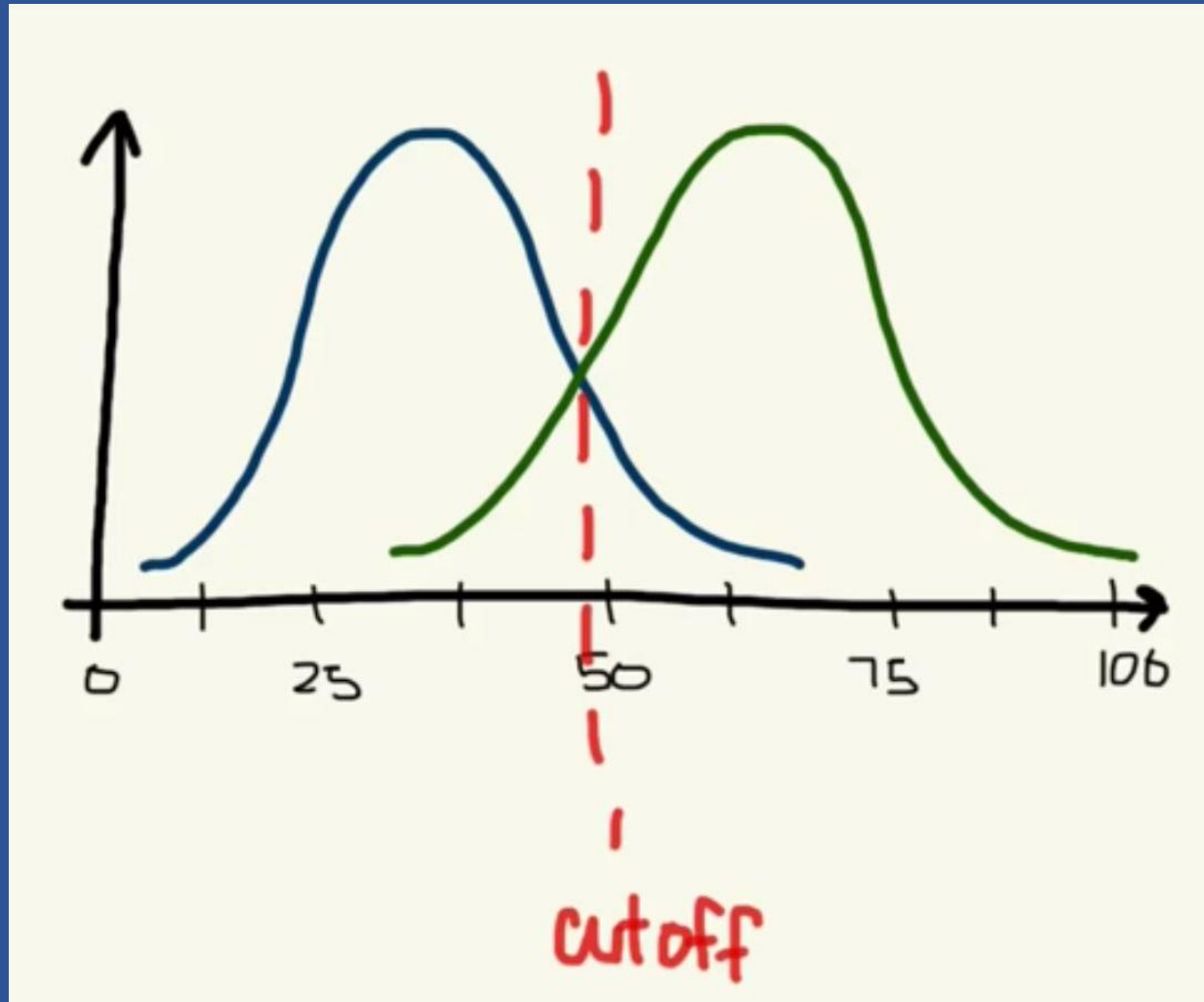ROC curve prediction result who have disease who don't

ML Evaluation
# Distribution score



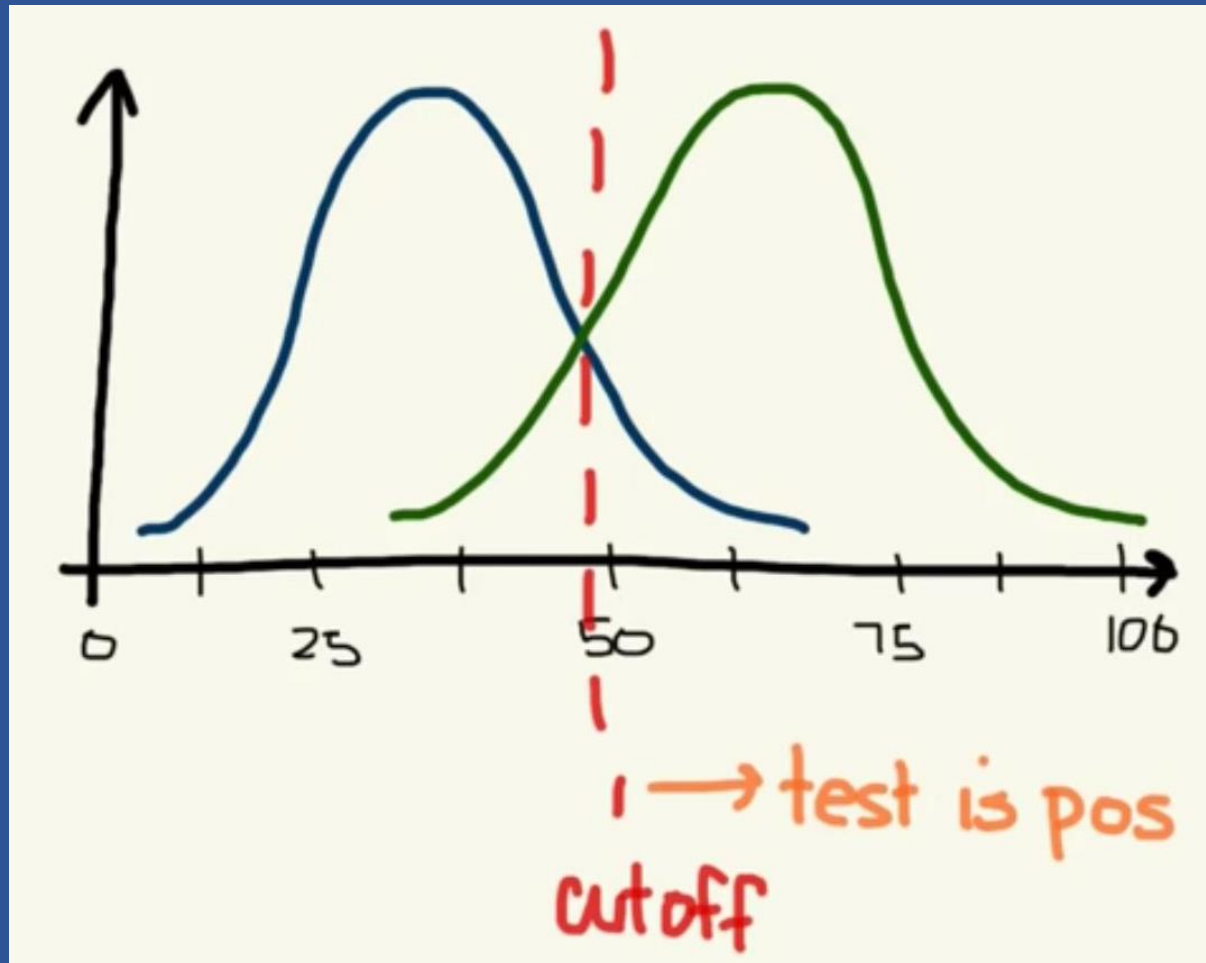Left distribution = patient who do NOT have disease (survived)  / Right = have disease (dead)

x axis = score / y axis = number of patient

ML Evaluation
# Cutoff line

ML Evaluation
# Area where the test is positive

**GreatFriends.Biz**
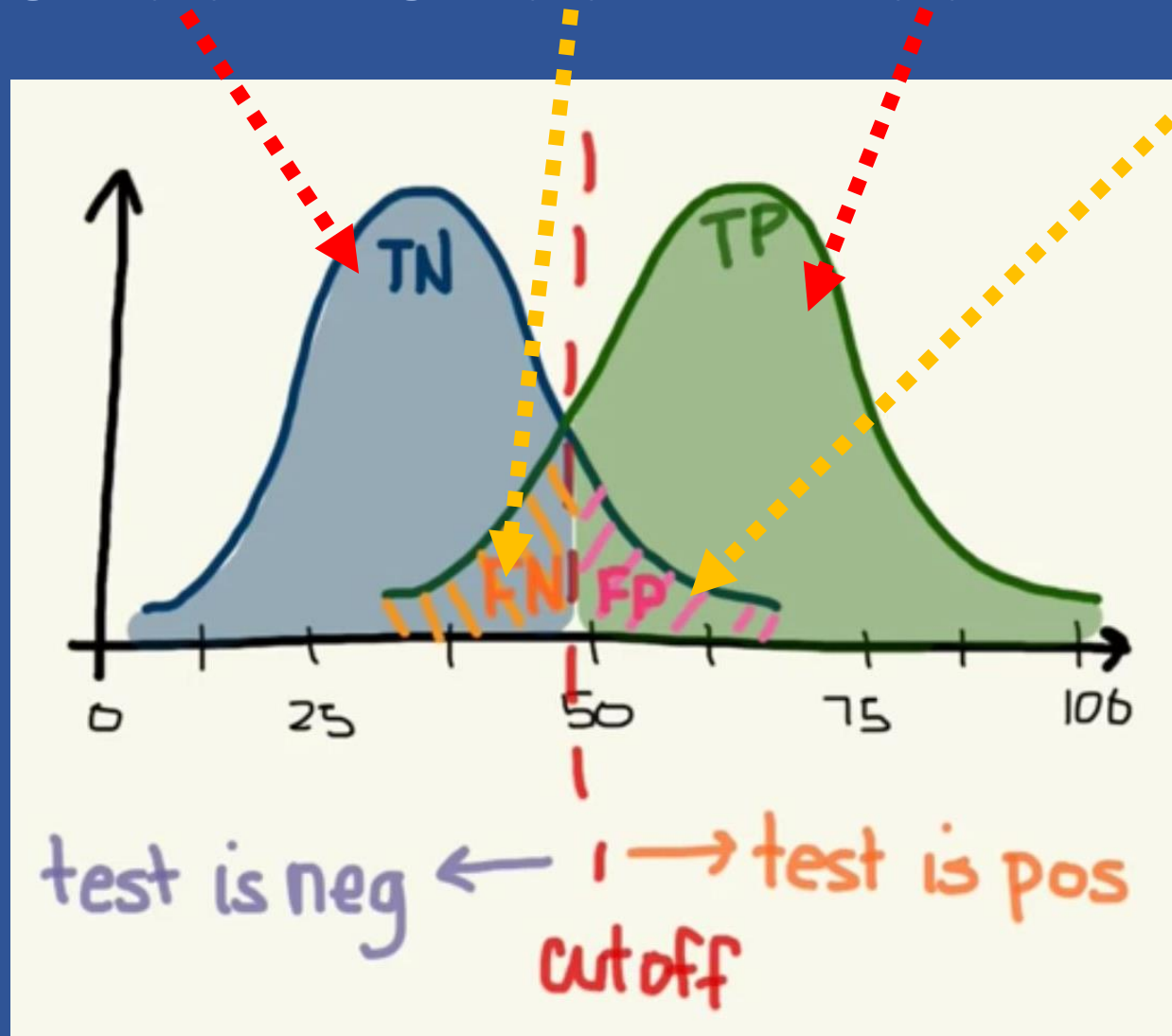
ML Evaluation
# Area where the test is negative

# ML Evaluation
## True Negative (TN), False Negative (FN) / True Positive (TP), False Positive (FP)

ML Evaluation
ROC Specificity / Sensitivity
## Specificity = True Negative Rate
## Sensitivity (Recall) = True Positive Rate

# ML Evaluation
## Move cutoff to the left Sens++ / Spec--

## ML Evaluation
## Move cutoff to right Sens-- / Spec++

ML Evaluation
# Chart proportion of Sens / Spec

ML Evaluation
# ROC curve = proportion of Sens / (1 – Spec)

ML Evaluation
# Area Under the Curve (AUC)

AUC is used to determine which of the used models predicts the classes best.

ML Evaluation
# AUC score

ML Evaluation
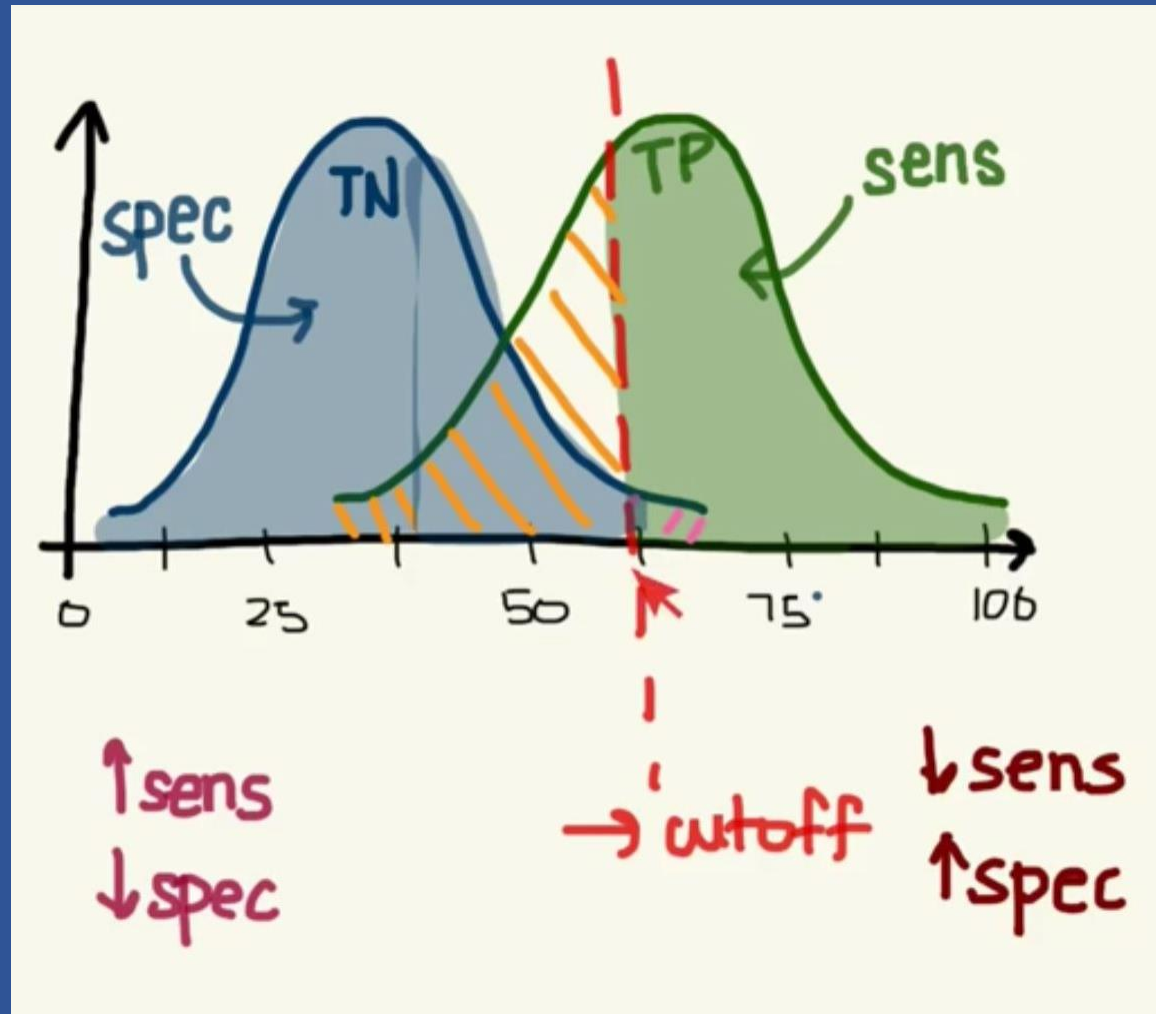# PRECISION/RECALL

**Precision:** the number of items correctly predicted as belonging to that class divided by the total number of items predicted as belonging to the class. TP / (TP + FP)

**Recall:** the number of items correctly predicted as belonging to that class divided by the total number of items that actually belong to the class. TP / (TP + FN)

ML Evaluation
# Evaluation metrics

| Score Bin | Positive Examples | Negative Examples | Fraction Above Threshold |
|---|---|---|---|
| (0.900,1.000] | 59 | 8 | 0.251 |
| (0.800,0.900] | 3 | 4 | 0.277 |
| (0.700,0.800] | 0 | 1 | 0.281 |
| (0.600,0.700] | 0 | 1 | 0.285 |

| Accuracy | F1 Score | Precision | Recall | Negative Precision | Negative Recall | Cumulative AUC |
|---|---|---|---|---|---|---|
| 0.824 | 0.715 | 0.881 | 0.602 | 0.805 | 0.953 | 0.023 |
| 0.820 | 0.721 | 0.838 | 0.633 | 0.813 | 0.929 | 0.038 |
| 0.816 | 0.717 | 0.827 | 0.633 | 0.813 | 0.923 | 0.041 |
| 0.813 | 0.713 | 0.816 | 0.633 | 0.812 | 0.917 | 0.045 |

# Evaluation metrics variable

- **True Positive (TP):** Correctly identified e.g. Sick people correctly diagnosed as sick
- **False Positive (FP):** Incorrectly identified e.g. healthy people incorrectly identified as sick
- **True Negative (TN):** Correctly rejected e.g. healthy people correctly identified as healthy
- **False Negative (FN):** Incorrectly rejected e.g. Sick people incorrectly identified as healthy
- **Accuracy :** The proportion of the total number of predictions that is correct. (TP + TN) / (TP + TN + FP + FN)

- Precision: is the proportion of positive cases that were correctly identified. TP / (TP + FP)
- Recall: Sensitivity or Recall is the proportion of actual positive cases which are correctly identified. TP / (TP + FN)
- F1 Score: is the harmonic mean of precision and Recall. 2TP / (2TP + FP + FN)
- Threshold: Threshold is the value above which it belongs to first class and all other values to the second class. E.g. if the threshold is 0.5 then any patient scored more than or equal to 0.5 is identified as sick else healthy.

ML Evaluation
# Sentiment evaluation results

- Positive Label: 1 = Good Text (GT)
- Negative Label: 0 = Bad Text  (BT)
- True Positive (TP): correctly predict GT
- True Negative (TN): correctly predict BT
- False Positive (FP):  incorrectly predict GT
- False Negative (FN): incorrectly predict BT

AUC 0.761



| True Positive | False Negative | Accuracy | Precision | Threshold |
|---|---|---|---|---|
| 106 | 46 | 0.690 | 0.693 | 0.5 |
| False Positive | True Negative | Recall | F1 Score | |
| 47 | 101 | 0.697 | 0.695 | |
| Positive Label | Negative Label | | | |
| 1 | 0 | | | |

ML Evaluation
# Metrics for Binary Classification

| METRICS | DESCRIPTION | LOOK FOR |
|---------|-------------|----------|
| Accuracy | proportion of correct predictions with a test data set | The closer to 1.00, the better |
| AUC | Area under the curve: This is measuring the area under the curve created by sweeping the true positive rate vs. the false positive rate. | The closer to 1.00, the better |
| AUCPR | Area under the curve of a Precision-Recall curve: Useful measure of success of prediction when the classes are very imbalanced (highly skewed datasets). | The closer to 1.00, the better |
| F1-score | the harmonic mean of the precision and recall. F1 Score is helpful when you want to seek a balance between Precision and Recall. | The closer to 1.00, the better |

# Next Step

## Create Sentiment model using AutoML