## Baysian Classifier:

It is statistical classifier. It can predict class membership probabilities, such as the probability that a given tuple belongs to a particular class. This is based on Baye's **(18th-century British mathematician Thomas Bayes)** theorem. It shows high accuracy and speed when it is applied to a large database.

# Example

- Example: Play Tennis

*PlayTennis*: training examples

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

# Example

- Learning Phase

| Outlook | Play=*Yes* | Play=*No* |
|---------|------------|-----------|
| *Sunny* | 2/9 | 3/5 |
| *Overcast* | 4/9 | 0/5 |
| *Rain* | 3/9 | 2/5 |

| Temperature | Play=*Yes* | Play=*No* |
|-------------|------------|-----------|
| *Hot* | 2/9 | 2/5 |
| *Mild* | 4/9 | 2/5 |
| *Cool* | 3/9 | 1/5 |

| Humidity | Play=*Yes* | Play=N*o* |
|----------|------------|-----------|
| *High* | 3/9 | 4/5 |
| *Normal* | 6/9 | 1/5 |

| Wind | Play=*Yes* | Play=*No* |
|------|------------|-----------|
| *Strong* | 3/9 | 3/5 |
| *Weak* | 6/9 | 2/5 |

$$P(\text{Play}=Yes) = 9/14 \qquad P(\text{Play}=No) = 5/14$$

# Example

- **Test Phase**
  - Given a new instance,

    $\mathbf{x}'=$(Outlook=*Sunny,* Temperature=*Cool,* Humidity=*High,* Wind=*Strong*)

  - Look up tables

    P(Outlook=*Sunny*|Play=*Yes*) = 2/9

    P(Temperature=*Cool*|Play=*Yes*) = 3/9

    P(Huminity=*High*|Play=*Yes*) = 3/9

    P(Wind=*Strong*|Play=*Yes*) = 3/9

    P(Play=*Yes*) = 9/14

    P(Outlook=S*unny*|Play=*No*) = 3/5

    P(Temperature=*Cool*|Play==*No*) = 1/5

    P(Huminity=*High*|Play=*No*) = 4/5

    P(Wind=*Strong*|Play=*No*) = 3/5

    P(Play=*No*) = 5/14

  - MAP rule

    P(*Yes*|$\mathbf{x}'$): [P(*Sunny*|*Yes*)P(*Cool*|*Yes*)P(*High*|*Yes*)P(*Strong*|*Yes*)]P(Play=*Yes*) = 0.0053

    P(*No*|$\mathbf{x}'$): [P(*Sunny*|N*o*) P(*Cool*|N*o*)P(*High*|*No*)P(*Strong*|*No*)]P(Play=*No*) = 0.0206

    Given the fact P(*Yes*|$\mathbf{x}'$) < P(*No*|$\mathbf{x}'$), we label $\mathbf{x}'$ to be "*No*".

# Confusion Matrix

It is a table that is often used to **describe the performance of a classification model** (or "classifier") on a set of test data for which the true values are known. The confusion matrix itself is relatively simple to understand, but the related terminology can be confusing.

Let's start with an **example confusion matrix for a binary classifier** (though it can easily be extended to the case of more than two classes):

| n=165 | Predicted: NO | Predicted: YES |
|---|---|---|
| Actual: NO | 50 | 10 |
| Actual: YES | 5 | 100 |

What can we learn from this matrix?

- There are two possible predicted classes: "yes" and "no". If we were predicting the presence of a disease, for example, "yes" would mean they have the disease, and "no" would mean they don't have the disease.
- The classifier made a total of 165 predictions (e.g., 165 patients were being tested for the presence of that disease).
- Out of those 165 cases, the classifier predicted "yes" 110 times, and "no" 55 times.
- In reality, 105 patients in the sample have the disease, and 60 patients do not.

- Let's now define the most basic terms, which are whole numbers (not rates):

| n=165 | Predicted: NO | Predicted: YES |
|---|---|---|
| Actual: NO | 50 | 10 |
| Actual: YES | 5 | 100 |

- **true positives (TP):** These are cases in which we predicted yes (they have the disease), and they do have the disease.
- **true negatives (TN):** We predicted no, and they don't have the disease.
- **false positives (FP):** We predicted yes, but they don't actually have the disease. (Also known as a "Type I error.")
- **false negatives (FN):** We predicted no, but they actually do have the disease. (Also known as a "Type II error.")

| n=165 | Predicted: NO | Predicted: YES | |
|---|---|---|---|
| **Actual: NO** | TN = 50 | FP = 10 | 60 |
| **Actual: YES** | FN = 5 | TP = 100 | 105 |
| | 55 | 110 | |

| n=165 | Predicted: NO | Predicted: YES | |
|---|---|---|---|
| **Actual: NO** | TN = 50 | FP = 10 | 60 |
| **Actual: YES** | FN = 5 | TP = 100 | 105 |
| | 55 | 110 | |

- **Accuracy:** Overall, how often is the classifier correct?

$$Accuracy = (TP+TN)/total = (100+50)/165 = 0.91$$

- **Misclassification Rate: Overall, how often is it wrong?**
  (FP+FN)/total = (10+5)/165 = 0.09
  equivalent to 1 minus Accuracy
  also known as "Error Rate"

| n=165 | Predicted: NO | Predicted: YES | |
|---|---|---|---|
| Actual: NO | TN = 50 | FP = 10 | 60 |
| Actual: YES | FN = 5 | TP = 100 | 105 |
| | 55 | 110 | |

- **True Positive Rate:** When it's actually yes, how often does it predict yes?
  - TP/actual yes = 100/105 = 0.95
  - also known as "Sensitivity" or "Recall"

- **False Positive Rate:** When it's actually no, how often does it predict yes?
- FP/actual no = 10/60 = 0.17

- **True Negative Rate:** When it's actually no, how often does it predict no?
  - TN/actual no = 50/60 = 0.83
  - equivalent to 1 minus False Positive Rate
  - also known as "Specificity"

| n=165 | Predicted: NO | Predicted: YES | |
|---|---|---|---|
| Actual: NO | TN = 50 | FP = 10 | 60 |
| Actual: YES | FN = 5 | TP = 100 | 105 |
| | 55 | 110 | |

- **Precision:** When it predicts yes, how often is it correct?
  TP/predicted yes = 100/110 = 0.91

- **Prevalence:** How often does the yes condition actually occur in our sample?
  actual yes/total = 105/165 = 0.64

**ROC(Receiver Operating Characteristics)** curve plots **TPR** (on the **y**-axis) against **FPR** (on the **x**-axis)

(True +ve rate) $\quad TPR = \dfrac{TP}{TP+FN}$

(False +ve rate) $\quad FPR = \dfrac{FP}{FP+TN}$

# Problem 1

**Confusion Matrix for Multi clssification:**

|   | A | B | C |
|---|---|---|---|
| A | 25 | 5 | 2 |
| B | 3 | 32 | 4 |
| C | 1 | 0 | 15 |

Accuracy = (25+32+15)/(25+5+2+3+32+4+1+0+15)

$P_A$=25/(25+3+1)

$R_A$ = 25/(25+5+2)