



# COS10022 –DATA SCIENCE PRINCIPLES

Dr Pei-Wei Tsai (Lecturer, Unit Convenor)  
ptsai@swin.edu.au, EN508d

SWIN  
BUR  
NE

SWINBURNE  
UNIVERSITY OF  
TECHNOLOGY





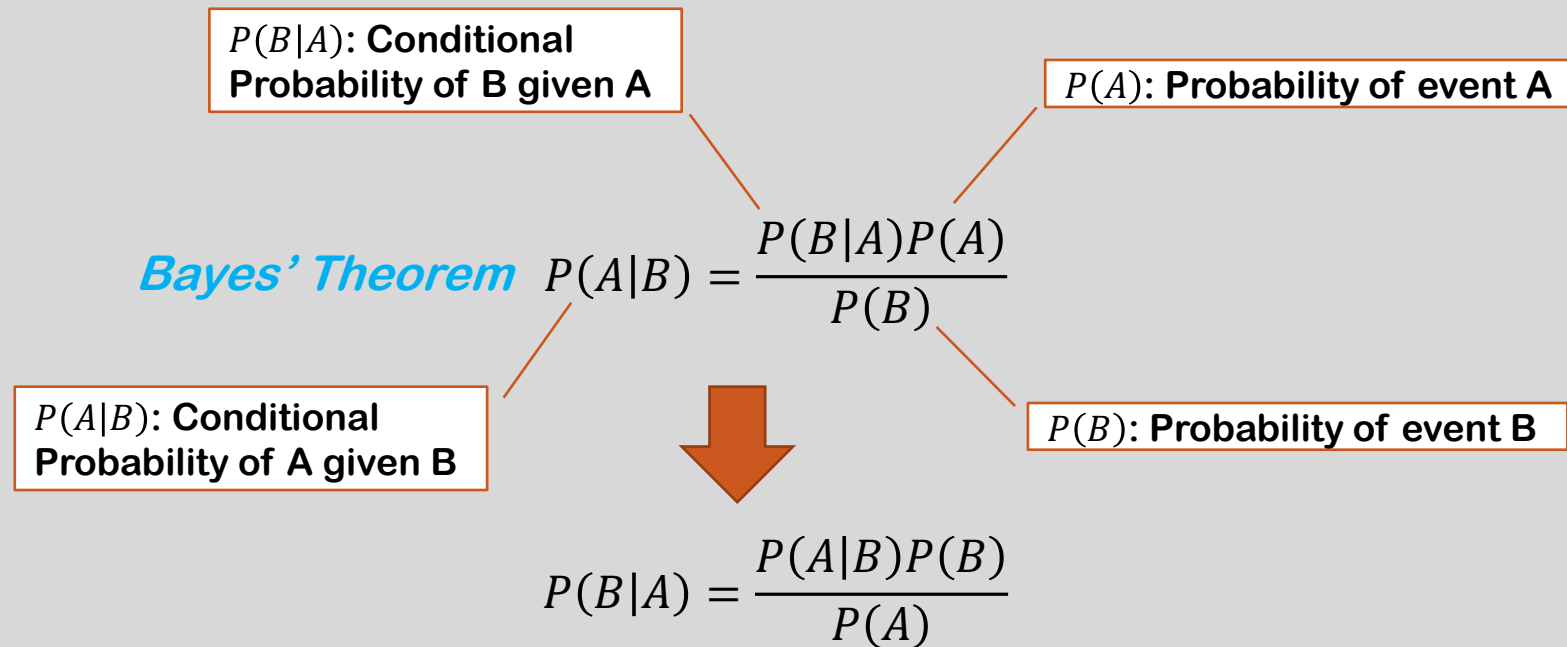
# BAYES' THEOREM

SWIN  
BUR  
NE

SWINBURNE  
UNIVERSITY OF  
TECHNOLOGY

# Bayes' Theorem

- If you know the probability of A happens conditional on B happens, the probabilities of A happens, and B happens, independently, you can derive the probability of B happens conditional on A happens.



# Why Using Bayes' Theorem?

**Cost of Data Collection.**

**Timeless.**

# Exercise 1

- If we want to find out a patient's probability of having liver disease if he or she is an alcoholic.
- **A** could mean the event “patient has liver disease.”
  - Past data tells you that 10% of patients entering your clinic have liver disease.  $\rightarrow P(A) = 0.1$
- **B** could mean the litmus test that “patient is an alcoholic.”
  - Past data tells you that 5% of patients are alcoholics.  $\rightarrow P(B) = 0.05$
- You might also know that among those patients diagnosed with liver disease, 7% are alcoholics.
  - The probability that a patient is alcoholic given that he/she has liver disease is 7%.  
 $\rightarrow P(B|A) = 0.07$
- Bayes' Theorem tells you that  $P(A|B) = \frac{(0.07 \times 0.1)}{0.05} = 0.14$ 
  - Thus, if the patient is an alcoholic, the chance of the patient having liver disease is 14%.





# MODEL EVALUATION

SWINBURNE

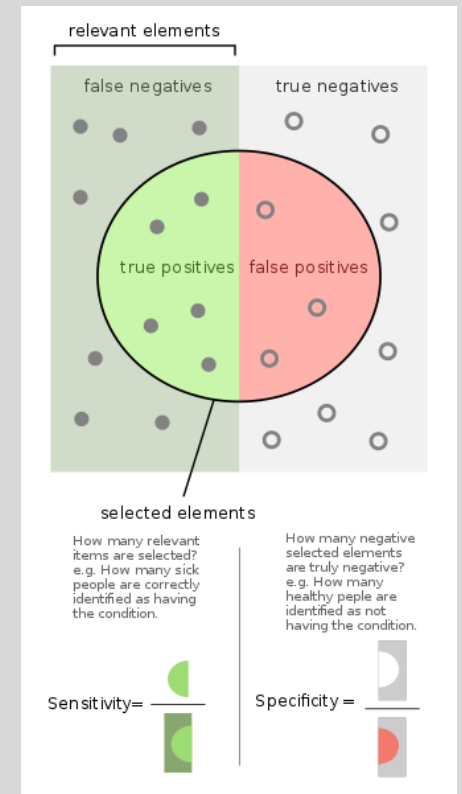
SWINBURNE  
UNIVERSITY OF  
TECHNOLOGY

# Supervised Models: Metrics and Methods

Popular metrics for evaluating the performance of supervised models:

1. **Accuracy**
2. **True Positive Rate (TPR)** – also called **Sensitivity** or **Recall**.
3. **True Negative Rate (TNR)** – also called **Specificity** or **Selectivity**.
4. **False Positive Rate (FPR)** – also called **Fall-out**.
5. **False Negative Rate (FNR)** – also called **Missing rate**.
6. **Precision:** 
$$\frac{\text{True positive}}{\text{Predictive Condition Positive}}$$
7. **Area Under the Curve (AUC)** – also called the **plasma concentration-time profile**.

These metrics can be calculated by utilising a **confusion matrix**.



# Supervised Models: Metrics and Methods

True Positives (TP):	the number of <b>positive</b> instances that a classifier correctly classifies as <b>positive</b> .
False Positives (FP):	the number of instances that a classifier identified as <b>positive</b> but in reality, are <b>negative</b> .
True Negatives (TN):	the number of <b>negative</b> instances that a classifier correctly identifies as <b>negative</b> .
False Negatives (FN):	the number of instances classified as <b>negative</b> but in reality, are <b>positive</b> .

**TP and TN are correct predictions. A good classifier should have large TP and TN, and small (ideally, zero) numbers of FP and FN.**



# Supervised Models: Metrics and Methods

**Example 1.** A confusion matrix of Naïve Bayes classifier for 100 customers in predicting whether they would subscribe to the term deposit.

		Predicted Class		Total
		Subscribed	Not Subscribed	
Actual Class	Subscribed	3 ( <b>correct prediction</b> )	8 ( <b>wrong prediction</b> )	11
	Not Subscribed	2 ( <b>wrong prediction</b> )	87 ( <b>correct prediction</b> )	89
Total		5	95	100



# QUESTION

**Which type of error is more important?**

# COVID-19 CRP Test

- Even in different stage, the important index is also different.



**Contain and Control Phase**

Missing rate (Type II error) is more important.  
FNR should be as low as possible.

**Herd Immunity  
(Community Immunity) Phase**

FPR (Type I error) is more important.  
Low FPR reduces the waste of resources.