

# **Business Analytics**

## **Lecture 2**

### **Probability**

Dr Yufei Huang

## Q&A

**Q: How many decimal points to keep in the answers?**

A: We usually keep 3 digits, if some calculation table gives 4 digits, then keep 4.

**Q: I'm using Windows, can't find Data Analysis tool in Excel.**

A: Please follow this link:

<https://support.office.com/en-gb/article/Load-the-Analysis-ToolPak-6a63e598-cd6d-42e3-9317-6b40ba1a66b4>

**Q: I'm using mac book, can't find Data Analysis tool in Excel.**

A: Please follow this link:

<https://support.microsoft.com/en-gb/help/2431349/how-to-find-and-install-data-analysis-toolpak-or-solver-for-excel-for>

**Q: Can I hand write and scan the report?**

A: Hand writing is ok. But it is recommended to type in word, then save as a pdf file. Otherwise, Turnitin function on Moodle will not work.

# Contents

- Set Theory
  - Definitions
  - Union, Intersection, Compliment
- Probability Theory
  - Definitions
  - Probability rules
  - Conditional probability
  - Bayes theorem

# Set Theory

- **Definition.** A **set** is a collection of objects, called elements or members.
- **Examples:**
  - $\mathbb{R}$  is the set of all real numbers between  $-\infty$  and  $\infty$ .
  - $\{1,2,4,7\}$  is the set including the numbers 1,2,4,7.
  - Intervals are also sets, e.g. all numbers between 3 and 4
  - Set of daily prices of the stock of HP during the last year.

# Set Theory

- **Definition.** A **subset** is any collection of elements in a set. An **empty set** is a set which contains no elements ( $\Phi$ ).
- **Examples**
  - $(3,4]$  is a subset of  $\mathbb{R}$ .
  - The price of HP's stock on the first day of each month is a subset of the set of its daily prices during a year.
  - The empty set is a subset of any set.

# Set Theory: Union

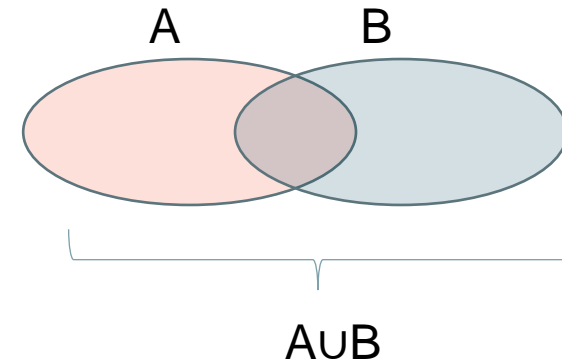
- **Union.** Let  $A, B$ , be two sets. Their union is the set which includes all elements which are in  $A$  or in  $B$ . It is denoted by  $A \cup B$ .

- **Example:**

$A$ : the set of daily prices of HP's stock between January and April 2014.

$B$ : the set of HP's daily prices between March and May 2014.

$A \cup B$ : the set of HP's daily prices between January and May 2014.



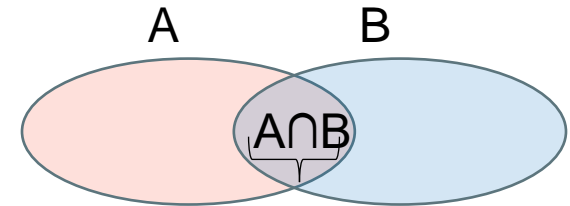
# Set Theory: Intersection

- **Intersection.** Let  $A, B$ , be two sets. Their intersection is the set which includes all elements which are in  $A$  **and** in  $B$ . It is denoted by  $A \cap B$ .
- **Example:**

$A$ : the set of daily prices of HP's stock between January and April 2014.

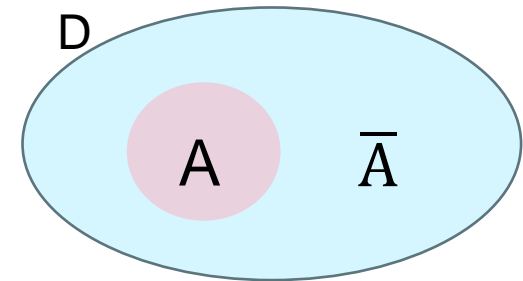
$B$ : the set of HP's daily prices between March and May 2014.

$A \cap B$ : the set of HP's prices between March and April.



# Set Theory: Compliment

- **Compliment.** Let  $A$  be a subset of a set  $D$ .  
Its complement,  $\bar{A}$ , consists of all of the elements in  $D$  apart from those in  $A$ .



- **Example.**  
 $D$  is the set of a survey results.  
 $A$  is the set of men's answers.  
 $\bar{A}$  is the set of the women's answers.
- An important set is **the set of all possible subsets** of a given set,  $A$ .

- **Example.**  
If  $A = \{1, 2, 3\}$ , the set of all subsets consists of:  
 $\{ \Phi, \{1\}, \{2\}, \{3\}, \{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\} \}$ .



# Probability Theory: Definitions

- **An experiment** is any process or procedure for which more than one outcome is possible.
- **Examples:**
  - Tossing a coin: Possible outcome: head
  - Rolling a die: Possible outcome: 6
  - The price of IBM's stock tomorrow: Possible outcome: \$187.81
- **Sample Space** is the set of all possible outcomes of an experiment.
- **Examples:**
  - Tossing a coin: {head, tail}
  - Rolling a die: {1,2,3,4,5,6}
  - The price of IBM's stock tomorrow:  $[0, \infty)$
- **An event** is a subset of a sample space.
- **Examples:**
  - The price of IBM's stock tomorrow is between \$180 and \$190.

# Probability Theory: Definitions

- **Probability** is a numerical measure of the chance that an event will occur.
- **Probability Measure** is a **function**  $P$  from the set of all of the events of a sample space  $\Omega$  to  $[0,1]$ , which satisfies for all disjoint events  $A_1, A_2, \dots, A_n$ :
  - $P(A) \geq 0$
  - $P(A_1 \cup A_2 \cup \dots \cup A_n) = P(A_1) + P(A_2) + \dots + P(A_n)$ .
  - $P(\Omega) = 1$
- A probability measure is a **function** which assigns a probability to each outcome in the sample space.
- Note. for any event  $A$ ,  $0 \leq P(A) \leq 1$ :
  - $P(A)=0$ , if  $A$  never happens.
  - $P(A)=1$ , if  $A$  always happens.

# Types of Probability

- Subjective Probability
- A-priori Classical Probability
- Empirical Classical Probability

# Subjective Probability

- **Subjective Probability:** Probability values are assigned subjectively by people.
- **Example:**
  - I have a feeling that the probability that IBM's stock price tomorrow will be higher than \$190 is 80%.
  - Someone less optimistic might estimate this probability by 50%.
- **Advantage:** Does not involve calculations
- **Disadvantage:** Subjective
  - Bilgin (2012) showed that people judge losses to have higher probabilities than gains
  - "It is very likely that this will happen" → Estimate the probability

# A-priori Classical Probability

- **A-priori classical probability:** the probability is determined according to a certain theory.
- **Example:** Random walk model predicts that the probability that IBM's stock price will increase tomorrow is 50%, because according to the random walk model, the probability of an increase in the price equals to the probability of decrease.
- **Advantages:**
  - May be used to develop theories
  - Helps calculations
- **Disadvantages:**
  - One needs a trustworthy theory

# Empirical Classical Probability

- Empirical classical probability:

- Probability is based on observed data.

$$P(\text{event}) = \frac{\text{Number of cases in which the event occurred}}{\text{The total number of cases}}$$

- **Example:** One monitors a very large number of trading days, e.g. 1000. and observes, 481 times that the price increased, then  $\text{Prob.} = 481/1000 = 0.481$ .

- Advantages:

- Does not require a theory
- Based on real data

- Disadvantages:

- Requires collection of data (time, money,...)
- Based on the assumption that probability distributions do not change.  
(Is this assumption true?)

# Exercise

In a marketing survey, 1000 consumers were asked whether they intend to buy a new product or not, and in a follow-up survey, the same group of consumers were asked whether they have bought this product. The data is shown below:

| Planned to Purchase | Actually Purchased |     |       |
|---------------------|--------------------|-----|-------|
|                     | Yes                | No  | Total |
| Yes                 | 200                | 100 | 300   |
| No                  | 250                | 450 | 700   |
| Total               | 450                | 550 | 1000  |

1. What is the probability that a person plans to purchase this product?
2. What approach did we use here?

# Solution

| Planned to Purchase | Actually Purchased |     |       |
|---------------------|--------------------|-----|-------|
|                     | Yes                | No  | Total |
| Yes                 | 200                | 100 | 300   |
| No                  | 250                | 450 | 700   |
| Total               | 450                | 550 | 1000  |

1. The probability that a person plans to purchase this product within a year is:

$$\begin{aligned}
 P(\text{planned to purchase}) &= \frac{\text{number of people who planned to purchase}}{\text{number of people in the sample}} \\
 &= \frac{200 + 100}{1000} = 0.3
 \end{aligned}$$

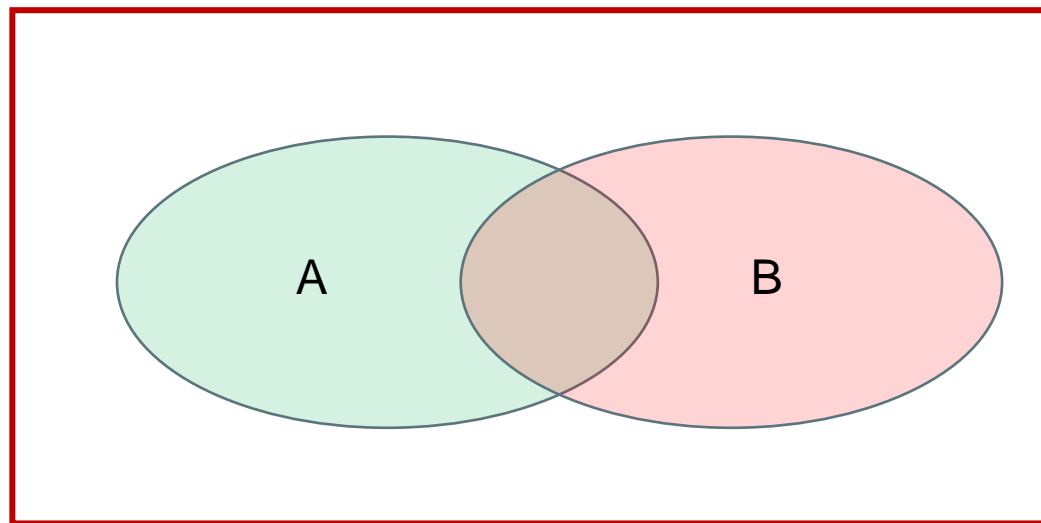
2. We used the empirical classical probability approach.



# Probability Rules

1.  $P(A) + P(\bar{A}) = 1$
2.  $P(A \cap B) + P(A \cap \bar{B}) = P(A)$
3.  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

- Can you prove the above rules graphically?



# Conditional Probability

- **Definition:** The probability of certain events **depends on** the probability of other events.
- **Example:**
  - Fruit and veges depend on the sun.
  - ([http://www.nzherald.co.nz/lifestyle/news/article.cfm?c\\_id=6&objectid=11723008](http://www.nzherald.co.nz/lifestyle/news/article.cfm?c_id=6&objectid=11723008) )
- We denote the probability of A given B by  **$P(A|B)$** .
- If  $P(B) \neq 0$ , then  $P(A | B) = \frac{P(A \cap B)}{P(B)}$

# Exercise

In a marketing survey, 1000 consumers were asked whether they intend to buy a new product or not, and in a follow-up survey, the same group of people were asked whether they have bought this product. The data is shown below:

|                     | Actually Purchased |     |       |
|---------------------|--------------------|-----|-------|
| Planned to Purchase | Yes                | No  | Total |
| Yes                 | 200                | 100 | 300   |
| No                  | 250                | 450 | 700   |
| Total               | 450                | 550 | 1000  |

- What is the probability that a consumer purchased this product **given that** he or she planned to purchase?
- What is the probability that a consumer did not purchase this product given that he or she planned to purchase?

# Solution

|                     | Actually Purchased |     |       |
|---------------------|--------------------|-----|-------|
| Planned to Purchase | Yes                | No  | Total |
| Yes                 | 200                | 100 | 300   |
| No                  | 250                | 450 | 700   |
| Total               | 450                | 550 | 1000  |

$$P(\text{Purchased} | \text{Planned}) = \frac{P(\text{Purchased and planned})}{P(\text{Planned})} = \frac{200/1000}{300/1000} = \frac{200}{300} = 0.6667$$

$$P(\text{Did not purchase} | \text{Planned}) = \frac{P(\text{Did not purchase and planned})}{P(\text{Planned})} = \frac{100/1000}{300/1000} = \frac{100}{300} = 0.3333$$

# Bayes Theorem

- **Bayes' Theorem:** 
$$P(B | A) = \frac{P(A | B)P(B)}{P(A | B)P(B) + P(A | \bar{B})P(\bar{B})}$$
- **Purpose:** If we know  $P(A|B), P(B), P(A|\bar{B})$  and  $P(\bar{B})$ , then we can calculate  $P(B|A)$ .
- **Example:** A factory manager uses quality test to identify defective products. The probability that a product is defective is 0.03. When the product is defective, the probability that this quality test will give a positive result is 0.90. When the product is not defective, the probability of a positive test result is 0.02. Suppose that the quality test has given a positive result. What is the probability that the product is actually defective?

- **Solution:** D: the product is defective,  $\bar{D}$ : the product is not defective  
T: test is positive,  $\bar{T}$ : test is negative

$$P(D)=0.03, P(\bar{D})=1-0.03=0.97, P(T|D)=0.90, P(T|\bar{D})=0.02$$

$$P(D | T) = \frac{P(T | D)P(D)}{P(T | D)P(D) + P(T | \bar{D})P(\bar{D})} = \frac{0.9 * 0.03}{0.9 * 0.03 + 0.02 * 0.97} = 0.582$$

- **Exercise:** Can you prove Bayes Theorem using conditional probability?

# Bayes Theorem in Court

<http://www.theguardian.com/law/2011/oct/02/formula-justice-bayes-theorem-miscarriage>

Edition: **UK** | US | AU | Profile | **Beta** | About us ▼

**theguardian**

News | Sport | Comment | Culture | Business | Money | Life & style |

News > Law

## A formula for justice

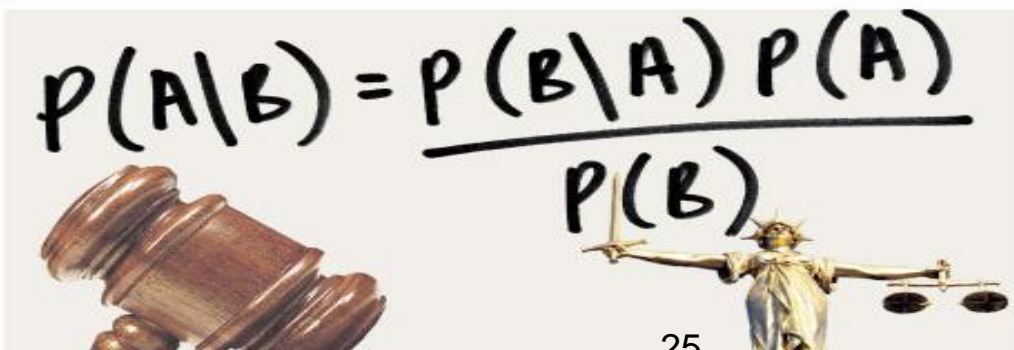
Bayes' theorem is a mathematical equation used in court cases to analyse statistical evidence. But a judge has ruled it can no longer be used. Will it result in more miscarriages of justice?



**Angela Saini**

The Guardian, Sunday 2 October 2011 21.30 BST

[Jump to comments \(...\)](#)



# Independent Events

- **Definition:** Two events, A and B are called **independent** if  $P(A|B)=P(A)$ . Namely, the outcome of B does not affect the probability of occurrence of A.
- **Example:** Two coin tosses are independent if the coins are fair.
- **Note:** The following necessary and sufficient conditions for independent events are equivalent:

$$P(B | A) = P(B) \quad P(A | B) = P(A) \quad P(A \cap B) = P(A) * P(B)$$

- **Exercise:** Can you prove the above note?

# Exercise

In a marketing survey, 1000 consumers were asked whether they intend to buy a new product or not, and in a follow-up survey, the same group of consumers were asked whether they have bought this product. The data is shown below:

| Planned to Purchase | Actually Purchased |     |       |
|---------------------|--------------------|-----|-------|
|                     | Yes                | No  | Total |
| Yes                 | 200                | 100 | 300   |
| No                  | 250                | 450 | 700   |
| Total               | 450                | 550 | 1000  |

- Given that a consumer planned to purchase, what is the probability that he/she finally purchased the product?
- Given that a consumer planned not to purchase, what is the probability that he/she finally purchased the product?
- Are consumers' willingness to purchase and their actual purchase decisions independent or not? Why?



# Solution

|                     | Actually Purchased |     |       |
|---------------------|--------------------|-----|-------|
| Planned to Purchase | Yes                | No  | Total |
| Yes                 | 200                | 100 | 300   |
| No                  | 250                | 450 | 700   |
| Total               | 450                | 550 | 1000  |

- Solution.

$$P(\text{Purchased} | \text{Planned}) = \frac{P(\text{Purchased and planned})}{P(\text{Planned})} = \frac{200/1000}{300/1000} = \frac{200}{300} = 0.6667$$

$$P(\text{Purchased} | \text{Planned Not}) = \frac{P(\text{Purchase and Planned Not})}{P(\text{Planned Not})} = \frac{250/1000}{700/1000} = \frac{250}{700} = 0.3571$$

$$P(\text{Purchased} | \text{Planned}) = 0.6667 \quad P(\text{Purchased}) = 0.45$$

So the two events are not independent.

# The Birthday Problem

To start,  $P(\text{two people share birthday}) = 1 - P(\text{no people share birthday})$ .

The probability of two people sharing birthday is difficult to get, but we can get **the probability of no people sharing birthday**:

- One person

This person can have any birthday.  $P(\text{no}) = (365/365) = 1$

- Two persons

$P(\text{no}) = (365/365) * (364/365) = 99.73\%$

- Three persons

$P(\text{no}) = ((365/365) * (364/365) * (363/365)) = 99.18\%$

.....

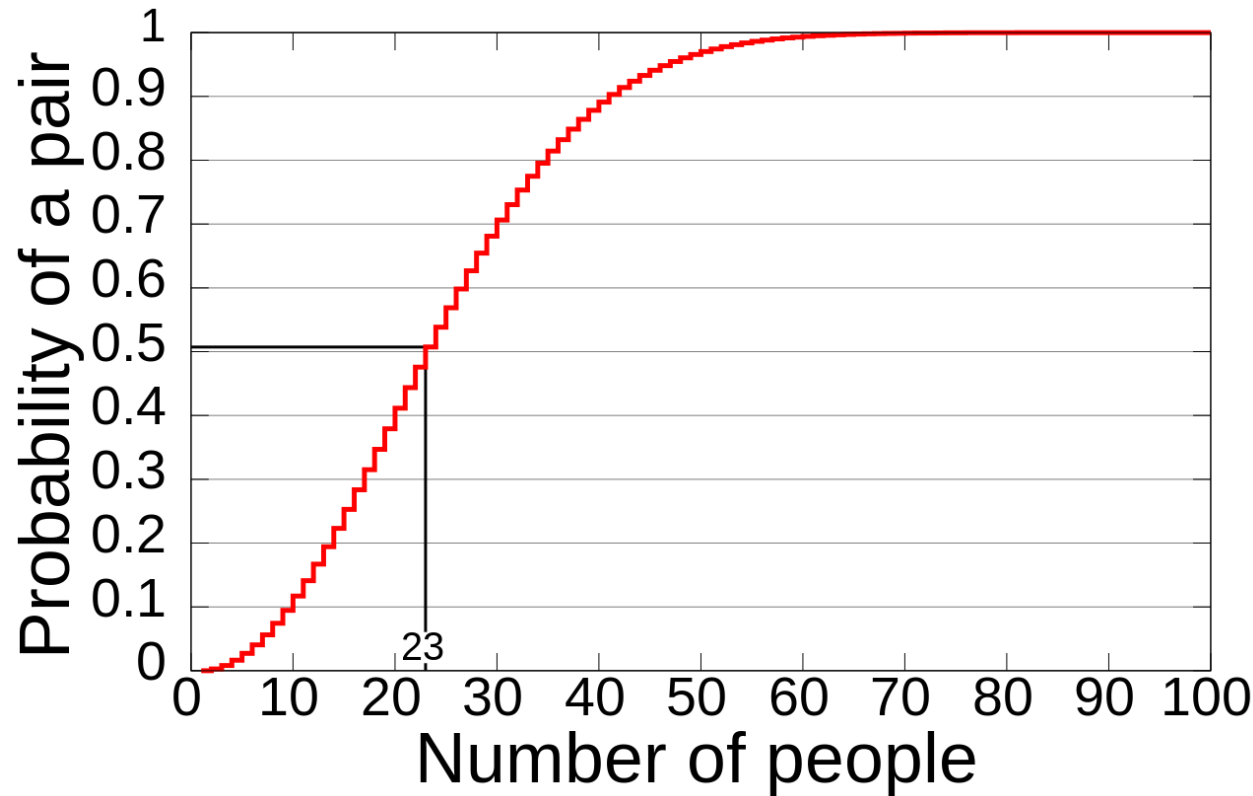
- n persons

$$P(\text{no}) = ((365/365) * ((365-1)/365) * ((365-2)/365) * \dots * ((365-n+1)/365))$$

$$= 365! / ((365-n)! * 365^n) = [365 * 364 * \dots * (365-n+1)] / 365^n$$

$P(\text{two people share birthday}) = 1 - P(\text{no people share birthday})$   
 $= 1 - [365 * 364 * \dots * (365-n+1)] / 365^n$

# The Birthday Problem



Next, we will solve the Birthday Problem using Excel.

# References

Chapter 2 of:

Aczel, A., & J. Sounderpandian. 2008. Complete Business Statistics. McGraw-Hill/Irwin, Seventh Edition.

Additional reading:

- Levine, Stephan, Krehbiel, & Berenson. Statistics for managers using Microsoft Excel. Prentice Hall, Upper Saddle River, New Jersey. Chapter 4: Basic probability.