

CT4031
Maths for Data Science

Week 5 – Probability Part 1

### Probability: An overview

Probability provides means to measure uncertainty of the phenomena around us.

The probability of any given phenomena is a value between 0 (cannot occur) and 1 (can occur).

Probability is used in data science for predicting how likely a phenomenon will occur, based on current observed data.



### **Probability: An overview**

The probability of an event, P(A), occurring is calculated as:

$$P(A) = \frac{\text{No. of favourable outcomes}}{\text{No. of all possible outcomes}}$$
 (1)

If P(A) represents the probability of an event occurring, its complement  $P(\tilde{A})$  is calculated as

$$P(\tilde{A}) = 1 - P(A) \tag{2}$$

# Probability: An overview

- ☐ The best example for understanding probability is flipping a coin:
- ☐ There are two possible outcomes—heads or tails.
- ☐ What's the probability of the coin landing on Heads? We can find out using the equation.
- ☐ In this case:
- ☐ Probability of an event = (# of ways it can happen) / (total number of outcomes)=1/2=0.5



## **Practice Examples of Probability**



- ☐ Example: Flipping a dice
  - ☐ There are six different outcomes. What's the probability of rolling a one?
  - □ Probability of an event = (# of ways it can happen) / (total number of outcomes)=1/6
  - □ P(1)=0.1666
- ☐ What's the probability of rolling a one or a six?
  - □ Probability of an event = (# of ways it can happen) / (total number of outcomes)=1/6
  - □ P(1 or 6)=2/6=1/3=0.3333333333333
- ☐ What's the probability of rolling an even number (i.e., rolling a two, four or a six)?
  - Probability of an event = (# of ways it can happen) / (total number of outcomes)
  - □ P(2,4,6)=3/6=1/2=0.5



### **Practice Exercise**

- 1. There are 6 beads in a bag, 3 are red, 2 are yellow and 1 is blue. What is the probability of picking a yellow?
- 2. A bag contains three bananas and nothing else. What is the probability of reaching into the bag and pulling out an apple?
- 3. Example: there are 5 marbles in a bag: 4 are blue, and 1 is red. What is the probability that a blue marble gets picked?



Assume that we have a box containing 5 balls (red, green, blue, yellow, orange). If a ball is to be drawn out of the box at random, what is the probability of:

- Getting a red ball?
- Getting a ball which is not red?



# **Important Tips**

- ☐ The probability of an event can only be between 0 and 1 and can also be written as a percentage.
- ☐ The probability of event A is often written as P(A)
- $\square$ If P(A) > P(B), then event A has a higher chance of occurring than event B.
- $\square$ If P(A) = P(B), then events A and B are equally likely to occur.



## **Mutual Exclusivity**

Given an observed phenomenon, events are mutually exclusive *if* at most one of them can occur at any one time.

This means that two events A and B are mutually exclusive if the occurrence of one of them eliminates the possibility of the other one occurring. (Think of flipping a coin.)

The probability of either A or B occurring then becomes:

$$P(A \text{ or } B) = P(A) + P(B) \tag{3}$$



### **Probabilistic Independence**

Events are said to be independent *if* the probability of each occurring is independent of each other.

The probability of two independent events A and B then becomes:

$$P(A \text{ and } B) = P(A) \times P(B)$$
 (4)



Imagine that we conduct an experiment which involves flipping a fair coin and recording the number of times we get a head (H) or a tail (T).

- What is the probability of getting a tail?
- Now imagine that we conduct the experiment three times. Now what is the probability of
  - Getting a T on the first turn, a H on the second turn, and a T on the third turn?
  - Getting at least one H during the experiments?



Now imagine that we conduct the experiment three times. Now what is the probability of

- Getting a T on the first turn, a H on the second turn, and a T on the third turn?
- Getting at least one H during the experiments?

### Independent:

2 \* 2 \* 2 = 8 combinations



Now imagine that we conduct the experiment three times. Now what is the probability of

- Getting a T on the first turn, a H on the second turn, and a T on the third turn?
  - 1 in 8 = 12.5%
- Getting at least one H during the experiments?
  - 7 in 8 = 87.5%
  - T, T, T
  - T,T,H
  - T,H,T
  - T,H,H
  - H,H,H
  - H,H,T
  - H,T,H
  - H,T,T



# **Conditional Probability**

The probability calculations which we have done so far involve looking at distinct events.

However, we tend to use conditional probability in data analytics to determine the probability of an event occurring.

- We determine the probability of an event occurring, given another event occurring.
- For example, we want to calculate the probability of raining tomorrow, given the weather has been hot for the past three days.

One of the most widely-used tools is Bayes' rule.



## **Conditional Probability**

What is the probability of it raining today, given that it has been sunny for the past three days?

What is the probability of Man Utd winning the next match, given it lost the previous match?

What are the chances of the refectory selling a new lunch menu, based on the fact that it has been selling fish & chips for the past 2 weeks?



Discovered by Rev. Thomas Bayes

One of the most frequently-used approach to determining conditional probability.

Involves the use of known (prior) probabilities in determining the likelihood of an event occurring.

Here likelihood indicates the amount of confidence we have on an event occurring.



- Traditional approaches to probability usually involve determining a single event (e.g., coin toss, getting a particular-coloured ball, etc)
- Bayes' rule allows us to incorporate the probability of other events in determining the probability of a particular event occurring
- Frequently used in machine learning (i.e., classification) algorithms for these reasons



### Bayes' Rule Formula

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

- ☐ Which tells us:
  - □ how often A happens given that B happens (Posterior probability), written P(A|B),
  - □ When we know: how often B happens given that A happens (Likelihood probability), written P(B|A)
  - □ and how likely A is on its own (Prior probability), written P(A)
  - □ and how likely B is on its own (Marginal probability), written P(B)
- ☐ The reverse is also true; for example:
- $\square$  P(B|A) = "P(B) P(A|B)" /"P(A) "
- ☐ Let us say P(Fire) means how often there is fire, and P(Smoke) means how often we see smoke, then:
- □ P(Fire|Smoke) means how often there is fire when we can see smoke
- ☐ P(Smoke|Fire) means how often we can see smoke when there is fire
- ☐ So the formula kind of tells us "forwards" P(Fire|Smoke) when we know "backwards" P(Smoke|Fire)



Traditional approaches to probability usually involve determining a single event (e.g., coin toss, getting a particular-coloured ball, etc)

Does not allow to take into account of other event(s) occurring

Bayes' rule allows us to incorporate the probability of other events in determining the probability of a particular event occurring

It also allows us to incorporate our own biases in determining the probability of an event occurring

Frequently used in machine learning (i.e., classification) algorithms for these reasons



In the annual Cheltenham Open, there are 40 horses which take part in the racecourses. Based on the statistics from the previous years, the following probabilities have been discovered:

- Probability of a participating horse winning a race: 20%.
- Probability of a horse getting an injury: 15%.
- Probability of an injured horse winning a race: 10%.

Question: What is the probability of a horse getting injured, given it has won a race?



Question: What is the probability of a horse getting **injured**, **given** it has **won** a race?

$$P(injured|winning) = \frac{0.1 \times 0.15}{0.2} = 0.075 = 7.5\%$$
 (6)

Probability of a horse winning a race, P(winning) = 0.2

Probability of getting injured, P(injury) = 0.15

P(winning|injured) = 0.1

The probability we want to determine becomes

$$P(injured|winning) = \frac{0.1 \times 0.15}{0.2} = 0.075 = 7.5\%$$
 (6)



#### **Exercise**

Probability of having a zero-trust architecture?

- 10,000 cases
- 1,000 zero-trust architectures

Probability of trojan infecting a network?

- 10,000 cases
- 3,000 trojan infections

Probability of a trojan infecting a network with zero-trust architecture?

- 10,000 cases
- 20 trojan infections

What is the probability that a zero-trust architecture is infected by a trojan?



# Bayes' Rule (just a bit more...)

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

Probability of A given B



## **Bays Rule Example 1**

- □ Dangerous fires are rare (1%)
- ☐But smoke is fairly common (10%) due to barbecues,
- □and 90% of dangerous fires make smoke
- ■We can then discover the probability of dangerous Fire when there is Smoke:
- □P(Fire|Smoke) = P(Dangerous Fire) P(Smoke|Dangerous Fire)/P(Smoke)
- $\Box = 1\% \times 90\%/10\%$
- □ = 9%
- ☐So it is still worth checking out any smoke to be sure.



### **Bays Rule Example 2**

- ☐ You are planning a picnic today, but the morning is cloudy. Oh no! 50% of all rainy days start off cloudy!, But cloudy mornings are common (about 40% of days start cloudy), and this is usually a dry month (only 3 of 30 days tend to be rainy, or 10%). What is the chance of rain during the day?
  - □ We will use Rain to mean rain during the day, and Cloud to mean cloudy morning.
  - ☐ The chance of Rain given Cloud is written P(Rain|Cloud)
  - So let's put that in the formula:
  - □ P(Rain|Cloud) = P(Rain) P(Cloud|Rain)/P(Cloud)
  - □ P(Rain) is Probability of Rain = 10%
  - □ P(Cloud|Rain) is Probability of Cloud, given that Rain happens = 50%
  - □ P(Cloud) is Probability of Cloud = 40%
  - $\square$  P(Rain|Cloud) = 0.1 x 0.5/0.4= = .125
  - Or a 12.5% chance of rain.
  - Not too bad, let's have a picnic!



## **Bays Rule Practice Exercise 1**

In the annual Cheltenham Open, there are 40 horses which take part in the racecourses. Based on the statistics from the previous years, the following probabilities have been discovered:

- Probability of a participating horse winning a race: 20%.
- Probability of a horse getting an injury: 15%.
- Probability of an injured horse winning a race: 10%.

Question: What is the probability of a horse getting injured, given it has won a race?



### **Bays Rule Practice Exercise 2**

- Probability of England winning the World Cup = 1 / 21
- Probability of Italy not being on the World Cup = 3/21
- Probability of Italy not being on the World Cup given that England has won = 0/21
- What is the Probability of England winning the World Cup given that has Italy is no longer on the competition?



### **Bays Rule Practice Exercise 3**

• In a particular pain clinic, 10% of patients are prescribed narcotic pain killers. Overall, five percent of the clinic's patients are addicted to narcotics (Including pain killers and illegal substances). Out of the all the people prescribed pain pills, 8% are addicts. If a patient is an addict, what is the probability that they will be prescribed pain pills?



Probability of England winning the World Cup given that has Italy is no longer on the competition?

There is a brand new virus known as XoviD. You do a test for XoviD and get a positive result, what are the chances that you have XoviD?



Probability of England winning the World Cup given that has Italy is no longer on the competition?

- Probability of England winning the World Cup = 1 / 21 = P(A)
- Probability of Italy not being on the World Cup = 3/21 = P(B)
- Probability of Italy not being on the World Cup given that England has won = 0/21 = P (B|A)

$$P(A|B) = P(B|A) * P(A) / P(B)$$
  
 $P = (0 * 0.047) / 0.142 = 0$ 

It might not come home this time....



You get a positive test, what are the chances that you have XoviD? 1% of people in the world has XoviD (thus, 99% have not) 80% of PCR tests detect a positive case correctly (thus, 20% are false-negative)

90.4% of Lateral Flow Tests detect a negative case correctly (thus, 9.6% are false-positive)



You get a positive test, what are the chances that you have XoviD? 1% of people in the world has XoviD (thus, 99% have not) 80% of PCR tests detect a positive case correctly (thus, 20% are false-positive)

90.4% of Lateral Flow Tests detect a negative case correctly (thus, 9.6% are false-negative)

Test Type	Has XoviD (1%)	Does NOT have XoviD (99%)
PCR	Correct Positive – 80%	False Positive – 20%
LFT	False Negative – 9.6%	Correct Negative – 90.4%



You get a positive test, what are the chances that you have XoviD? 80%? 90%? 1%?

Test Type	Has XoviD (1%)	Does NOT have XoviD (99%)
PCR	Correct Positive – 80%	False Positive – 20%
LFT	False Negative – 9.6%	Correct Negative – 90.4%



You get a positive test, what are the chances that you have XoviD? 80%? 90%? 1%?

Test Type	Positive	Negative
PCR	Correct Positive 1% x 80% = 0.008	False Positive – 20% 1% x 20% = 0.002
LFT	False Negative – 9.6% 99% x 9.6% = 0.095	Correct Negative – 90.4% 99% x 90.4% = 0.894



You get a positive test, what are the chances that you have XoviD? 80%? 90%? 1%?

P(XoviD|P) = P(P|XoviD) \* P(XoviD) / P(P)

Test Type	Positive	Negative
PCR	Correct Positive 1% x 80% = 0.008	False Positive – 20% 1% x 20% = 0.002
LFT	False Negative – 9.6% 99% x 9.6% = 0.095	Correct Negative – 90.4% 99% x 90.4% = 0.894



$$P(XoviD|P) = P(P|XoviD) * P(XoviD) / P(P)$$

$$P(XoviD|P) = ?$$

$$P(P|XoviD) = 80\% = 0.8$$

$$P(XoviD) = 1\% = 0.01$$

$$P(P) = 0.008 + 0.095 = 0.103$$

Test Type	Positive	Negative
PCR	Correct Positive 1% x 80% = 0.008	False Positive – 20% 1% x 20% = 0.002
LFT	False Negative – 9.6% 99% x 9.6% = 0.095	Correct Negative – 90.4% 99% x 90.4% = 0.894



$$P(XoviD|P) = 0.8 * 0.01 / 0.103 = aprox. 7.8\%$$

$$P(XoviD|P) = 7.8\%$$

$$P(P|XoviD) = 80\% = 0.8$$

$$P(XoviD) = 1\% = 0.01$$

$$P(P) = 0.008 + 0.095 = 0.103$$

Test Type	Positive	Negative
PCR	Correct Positive 1% x 80% = 0.008	False Positive – 20% 1% x 20% = 0.002
LFT	False Negative – 9.6% 99% x 9.6% = 0.095	Correct Negative – 90.4% 99% x 90.4% = 0.894



### Post-sessional work

Write a brief report discussing the differences between mutual exclusive and probabilistic independence. Make sure to include references using Harvard referencing style.



### References

Albright, S and Wayne Winston (2014). Business Analytics: Data Analysis & Decision Making. Nelson Education.

University celebrates \betting man" Thomas Bayes (2016). http://www.bbc.co.uk/news/uk-scotland-edinburgh-east-fife-14708583.



### **Next Session!**

Probability Part 2





