



COMP20008 Elements of Data Processing

Semester 2 2018

Lecture 10: Mutual Information



THE UNIVERSITY OF
MELBOURNE

Announcements

- Project Phase 1 due next Friday 31st August at 11:59am
- Consultation sessions about Project Phase 1
 - Fri 17/08/2018 Room 09.02 Doug McDonell 9.30am-10:30am
 - Thu 23/08/2018 Room 07.02 Doug McDonell 11am-12pm
 - Wed 29/08/2018 Room 07.02 Doug McDonell 10:30am-11:30am



THE UNIVERSITY OF
MELBOURNE

Plan today

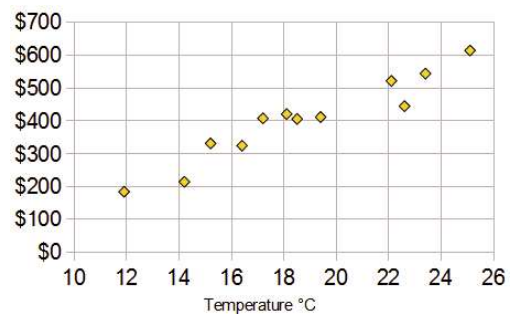
- Recap – correlations
 - Pearson correlation
- Another measure for correlation
 - Mutual information
 - Entropy
 - Conditional entropy



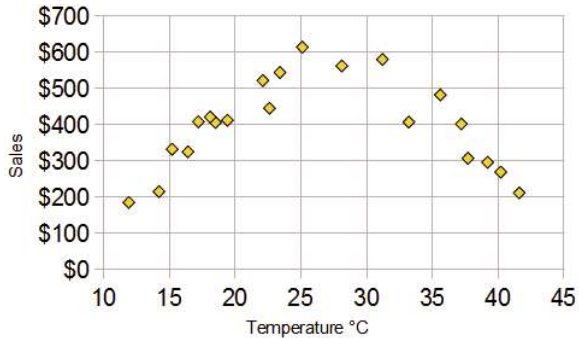
THE UNIVERSITY OF
MELBOURNE

Pearson correlation – assess degree of linear correlation between two features

$$r = r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$



<https://www.mathsisfun.com/data/correlation.html>



Pearson correlation is not suitable for this scenario (value less than 0.1)

<https://www.mathsisfun.com/data/correlation.html>

- A correlation measure that can detect **non-linear relationships**
 - It operates with **discrete features**
 - Preprocessing**: continuous features are first discretised into bins (categories). E.g. small [0,1.4], medium (1.4,1.8), big [1.8,3.0]

Object	Height	Discretised Height
1	2.03	big
2	1.85	big
3	1.23	small
4	1.31	small
5	1.72	medium
6	1.38	small
7	0.94	small

- Domain knowledge**: assign thresholds manually
 - Car speed:
 - 0-40: slow
 - 40-60: mid
 - >60: high
- Equal-width bin**
 - Divide the range of the continuous feature into equal length intervals (bins).
 - If speed ranges from 0-100, then the 10 bins are [0,10), [10,20), [20,30), ..., [90,100]
- Equal frequency bin**
 - Divide range of continuous feature into equal frequency intervals (bins). Sort the values and divide so that each bin has same number of objects.

- Given the values 2, 2, 3, 10, 13, 15, 16, 17, 19, 19, 20, 20, 21
 - Show a 3 bin equal width discretisation
 - Show a 3 bin equal frequency discretisation

- Given the values 2, 2, 3, 10, 13, 15, 16, 17, 19, 19, 20, 20, 21
 - Show a 3 bin equal width discretisation
 - $width = \frac{max-min}{N} \rightarrow min = 2, max = 21, N = 3$
 - $width = \frac{21-2}{3} = 6.33$
 - bin1 – range: [2, 8.33) \rightarrow values: 2, 2, 3
 - bin2 – range: [8.33, 14.667) \rightarrow values: 10, 13
 - bin3 – range: [14.667, 21) \rightarrow values: 15, 16, 17, 19, 19, 20, 20, 21
 - Show a 3 bin equal frequency discretization
 - $freq = \frac{number-of-object}{N} \rightarrow \frac{13}{3} = 4$
 - bin1 – values: 2, 2, 3, 10
 - bin2 – values: 13, 15, 16, 17
 - bin3 – values: 19, 19, 20, 20, 21

- $2*2*2*2=16$
 - $\log_2 16 = 4$ (16 is 2 to the power 4)
- $y = \log_2 x$ (y is the solution to the question “To what power do I need to raise 2, in order to get x?”)
- $\log_2 32 = 5$
- $\log_2 30 = 4.9$
- $\log_2 1.2 = 0.26$
- $\log_2 0.5 = -1$
- In what follows, we’ll write log instead of \log_2

- Entropy is a measure used to assess the amount of **uncertainty** in an outcome
- E.g. Randomly select an element from {1,1,1,1,1,1,1,2}, versus randomly select an element from {1,2,2,3,3,4,5}
 - In which case is the value selected more “predictable”? Why?

- E.g. Randomly select an element from {1,1,1,1,1,1,1,2}, versus randomly select an element from {1,2,2,3,3,4,5}
 - In which case is the value selected more “predictable”? Why?
 - The former case **more certain** \Rightarrow low entropy
 - The latter case is **less certain** \Rightarrow higher entropy
 - Entropy is used to quantify this degree of uncertainty (surprisingness)

- Consider the sample of all people in this lecture theatre. Each person is labelled young (<30 years) or old (>30 years)
- Randomly select a person and inspect whether they are young or old.
 - How surprised am I likely to be by the outcome?
- Suppose I repeat the experiment using a random sample of people catching the train to the city in peak hour?
 - How surprised am I likely to be by the outcome?

- Given a feature **X**. Then $H(\mathbf{X})$ is its entropy. Assuming X uses a number of categories (bins)

$$H(\mathbf{X}) = - \sum_{i=1}^{\#bins} p_i \log p_i$$

- p_i : proportion of points in the i -th bin
- May sometimes write $p(i)$ instead of p_i
- E.g. Suppose there are 3 bins, each bins contains exactly one third of the objects (points)
 - $H(\mathbf{X})=?$

$$H(\mathbf{X}) = - \sum_{i=1}^{\#bins} p_i \log p_i$$

- p_i : proportion of points in the i -th bin
- E.g. Suppose there are 3 bins, each bins contains exactly one third of the objects (points)
 - bin1: $p_1 \log p_1 = \frac{1}{3} \log \frac{1}{3}$
 - bin2: $p_2 \log p_2 = \frac{1}{3} \log \frac{1}{3}$
 - bin3: $p_3 \log p_3 = \frac{1}{3} \log \frac{1}{3}$
 - $H(\mathbf{X}) = - [0.33 \log(0.33) + 0.33 \log(0.33) + 0.33 \log(0.33)]$
- The log can be any base, we will assume base 2 (\log_2)

A	B	B	A	C	C	C	C	A
---	---	---	---	---	---	---	---	---

We have 3 categories/bins (A,B,C) for a feature **X**
 9 objects, each in exactly one of the 3 bins
 What is the entropy of this sample of 9 objects?

- $H(\mathbf{X}) = ???$

Entropy example

A B B A C C C C A

We have 3 categories/bins (A,B,C) for a feature **X**

9 objects, each in exactly one of the 3 bins

What is the entropy of this sample of 9 objects?

- bin A: $p_1 \log p_1 = \frac{3}{9} \log \frac{3}{9}$
- bin B: $p_2 \log p_2 = \frac{2}{9} \log \frac{2}{9}$
- bin C: $p_3 \log p_3 = \frac{4}{9} \log \frac{4}{9}$
- $H(X) = -[\frac{3}{9} \log(\frac{3}{9}) + \frac{2}{9} \log(\frac{2}{9}) + \frac{4}{9} \log(\frac{4}{9})]$

Answer: $H(\mathbf{X})=1.53$

Entropy example

A B B A C C C C A

We have 3 categories/bins (A,B,C) for a feature **X**

9 objects, each in exactly one of the 3 bins

What is the entropy of this sample of 9 objects?

- bin A: $p_1 \log p_1 = \frac{3}{9} \log \frac{3}{9}$
- bin B: $p_2 \log p_2 = \frac{2}{9} \log \frac{2}{9}$
- bin C: $p_3 \log p_3 = \frac{4}{9} \log \frac{4}{9}$
- $H(X) = -[\frac{3}{9} \log(\frac{3}{9}) + \frac{2}{9} \log(\frac{2}{9}) + \frac{4}{9} \log(\frac{4}{9})]$

Answer: $H(\mathbf{X})=1.53$

- $H(X)$ is the amount of randomness in the above feature vector
- High value for $H(x)$
→ high randomness
→ high uncertainty
→ less predictability

Properties of the entropy

- $H(X) \geq 0$
- Entropy – when using log base 2 – measures uncertainty of the outcome in bits. This can be viewed as the information associated with learning the outcome
- Entropy is **maximized** for uniform distribution (**highly uncertain** what value a randomly selected object will have)

Conditional entropy intuition

- Suppose I randomly sample a person in the class. Look to see whether they wear glasses – how surprised am I by their age?

Object	WearGlasses(X)	Age (Y)
1	No	young
2	No	young
3	No	young
4	No	young
5	Yes	old
6	Yes	old
7	Yes	old

Conditional entropy $H(Y|X)$

- Measures how much information needed to **describe outcome Y**, given that **outcome X is known**. Suppose X is Height and Y is Weight.

Object	Height (X)	Weight (Y)
1	big	light
2	big	heavy
3	small	light
4	small	light
5	small	light
6	small	light
7	small	heavy

$$H(Y|X) = \sum_{x \in \mathcal{X}} p(x) H(Y|X = x)$$

Conditional entropy

Object	Height (X)	Weight (Y)
1	big	light
2	big	heavy
3	small	light
4	small	light
5	small	light
6	small	light
7	small	heavy

$$H(Y|X) = \sum_{x \in \mathcal{X}} p(x) H(Y|X = x)$$

- $H(Y|X) = 2/7 * H(Y|X = big) + 5/7 * H(Y|X = small)$
 - $H(Y|X = big) = -0.5 \log 0.5 - 0.5 \log 0.5$
 - $H(Y|X = small) = -0.8 \log 0.8 - 0.2 \log 0.2$
- $H(Y|X) = \frac{2}{7} (-0.5 \log 0.5 - 0.5 \log 0.5) + \frac{5}{7} (-0.8 \log 0.8 - 0.2 \log 0.2) = 0.801$

Conditional entropy (another example)

Object	Height (X)	Weight (Y)
1	small	light
2	big	heavy
3	small	light
4	small	light
5	small	light
6	small	light
7	small	heavy

$$H(Y|X) = \sum_{x \in \mathcal{X}} p(x) H(Y|X = x)$$

- $H(Y|X) = 0.5572$
 - $H(Y) = 0.8631205$
 - $H(Y) - H(Y|X) = 0.306$
 - (how much information about Y is gained by knowing X)
- What 0.306 means?
 - How it is different than 0.7?

Conditional entropy another example

Object	Height (X)	Weight (Y)
1	big	light
2	big	heavy
3	small	light
4	small	jumbo
5	medium	light
6	medium	light
7	small	heavy

$$H(Y|X) = \sum_{x \in \mathcal{X}} p(x) H(Y|X = x)$$

- $H(Y|X) = 0.965$
- $H(Y) = 1.379$
- $H(Y) - H(Y|X) = 0.414$

$$MI(X, Y) = H(Y) - H(Y|X)$$

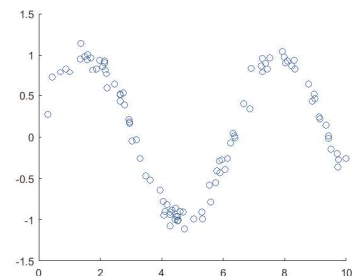
$$= H(X) - H(X|Y)$$

- Where X and Y are features (columns) in a dataset.
- MI (mutual information) is a measure of correlation
 - The amount of information about X we gain by knowing Y, or
 - The amount of information about Y we gain by knowing X

- The amount of information shared between two variables X and Y
- $MI(X, Y)$
 - large: X and Y are highly correlated (more dependent)
 - small: X and Y have low correlation (more independent)
- $0 \leq MI(X, Y)$

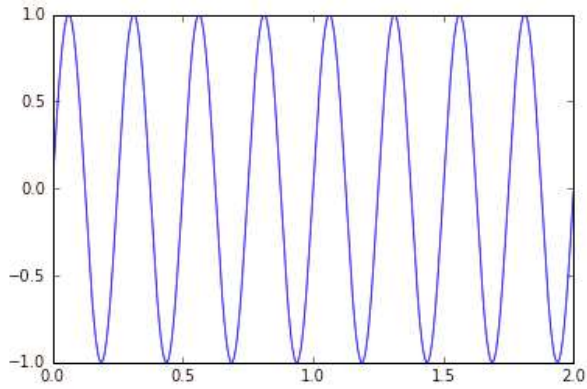
- What large means?
 - Is 0.7 large value?
 - Is 3.7 large value?

- $MI(X, Y)$ is always at least zero, may be larger than 1
- In fact, one can show it is true that
 - $0 \leq MI(X, Y) \leq \min(H(X), H(Y))$
 - (where $\min(a, b)$ indicates the minimum of a and b)
- Thus if want a measure in the interval $[0, 1]$, we can define normalized mutual information (NMI)
 - $NMI(X, Y) = \frac{MI(X, Y)}{\min(H(X), H(Y))}$
- $NMI(X, Y)$
 - large: X and Y are highly correlated (more dependent)
 - small: X and Y have low correlation (more independent)



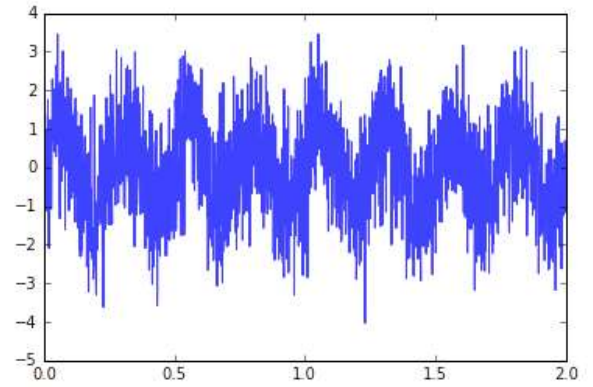
- Pearson: -0.0864
- NMI: 0.43 (3-bin equal frequency discretization)

Examples (2)



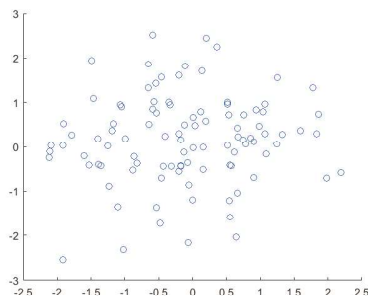
- Pearson: -0.1
- NMI: 0.84

Examples (3)



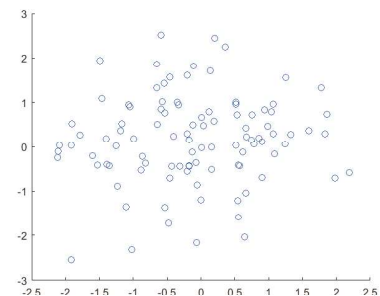
- Pearson: -0.05
- NMI: 0.35

Examples (4)



- Pearson?
- NMI?

Examples (4)



- Pearson: 0.08
- NMI: 0.009

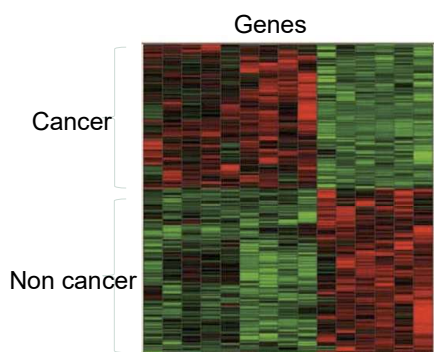
- Identifying features that are highly correlated with a class feature

HoursSleep	HoursExercise	HairColour	HoursStudy	Happy (class feature)
12	20	Brown	low	Yes
11	18	Black	low	Yes
10	10	Red	medium	Yes
10	9	Black	medium	Yes
10	10	Red	high	No
7	11	Red	high	No
6	15	Brown	high	No
2	13	Brown	high	No

- Compute $MI(\text{HoursSleep}, \text{Happy})$, $MI(\text{HoursExercise}, \text{Happy})$, and $MI(\text{HoursStudy}, \text{Happy})$, $MI(\text{HairColour}, \text{Happy})$. Retain most predictive feature(s)

HoursSleep	HoursExercise	HairColour	HoursStudy	Happy (class feature)
12	20	Brown	low	Yes
11	18	Black	low	Yes
10	10	Red	medium	No
10	9	Black	medium	Yes
10	10	Red	high	No
7	11	Black	high	No
6	15	Brown	high	No
2	13	Brown	high	No

- $MI(\text{HairColour}, \text{Happy})=0.27$ (NMI=0.28)
- $MI(\text{HoursStudy}, \text{Happy})=0.70$ (NMI=0.74)
-
- Can rank features according to their predictiveness –then focus further on just these



- $\text{Cancer} = f(\text{gene1}, \text{gene2}, \dots, \text{gene } n)$ #f is some unknown function
- Use correlation to reduce the number of variables

	Gene 1	Gene 2	Gene 3	...	Gene n	Cancer
Person 1	2.3	1.1	0.3	...	2.1	1
Person 2	3.2	0.2	1.2	...	1.1	1
Person 3	1.9	3.8	2.7	...	0.2	0
...
Person m	2.8	3.1	2.5	...	3.4	0

- Use relevant genes only: improving accuracy & performance

- Advantage
 - Can detect both **linear and non linear** dependencies (unlike Pearson)
 - Applicable and very effective for use with **discrete features** (unlike Pearson correlation)
- Disadvantage
 - If feature is continuous, it first must be discretised to compute mutual information. This involves making choices about **what bins to use**.
 - This may not be obvious. Different bin choices will lead to **different estimations of mutual information**

Instance ID	Predicted class	Actual class
1	X	X
2	X	Y
3	Y	Y
4	X	X
5	X	Y
6	Y	X
7	X	X
8	Y	Y
9	Y	X
10	Y	Y

- a) (1 mark) Would Pearson correlation be suitable to compute the correlation between the *Predicted class* and *Actual class*? Why or why not?

- a) Richard is a data wrangler. He does a survey and constructs a dataset recording average time/day spent studying and average grade for a population of 1000 students:

Student Name	Average time per day spent studying	Average Grade
...

- i) Richard computes the Pearson correlation coefficient between *Average time per day studying* and *Average grade* and obtains a value of **0.85**. He concludes that **more time spent studying causes a student's grade to increase**. Explain the **limitations** with this reasoning and suggest **two alternative explanations** for the 0.85 result.

- a) Richard is a data wrangler. He does a survey and constructs a dataset recording average time/day spent studying and average grade for a population of 1000 students:

Student Name	Average time per day spent studying	Average Grade
...

- ii) Richard separately discretises the two features *Average time per day spent studying* and *Average grade*, each into 2 bins. He then computes the normalised mutual information between these two features and obtains a value of 0.1, which seems surprisingly low to him. Suggest two reasons that might explain the mismatch between the normalised mutual information value of 0.1 and the Pearson Correlation coefficient of 0.85. Explain any assumptions made.

- understand the advantages and disadvantages of using mutual information for computing correlation between a pair of features. Understand the main differences between this and Pearson correlation.
- understand the meaning of the variables in the mutual information and how they can be calculated. Be able to compute this measure on a simple pair of features. The formula for mutual information will be provided on the exam.
- understand the role of data discretization in computing mutual information
- understand the meaning of the entropy of a random variable and how to interpret an entropy value. Understand its extension to conditional entropy
- be able to interpret the meaning of the mutual information between two features

- understand the use of mutual information for computing correlation of some feature with a class feature and why this is useful. Understand how this provides a ranking of features, according to their predictiveness of the class
- understand that normalised mutual information can be used to provide a more interpretable measure of correlation than mutual information. The formula for normalised mutual information will be provided on the exam