# EE3001 Foundations of Data Engineering

Lecturer: Dr. YUEN Shiu Yin, Kelvin

Email: kelviny.ee   Tel: x 7717   Rm: G6359
Web:   http://www.ee.cityu.edu.hk/~syyuen

# Course Aims

It aims to teach probability and statistics from a data engineering perspective.

It prepares you to take more advanced data engineering courses as well as gives you generic skills to handle other courses that requires knowledge about probability and statistics
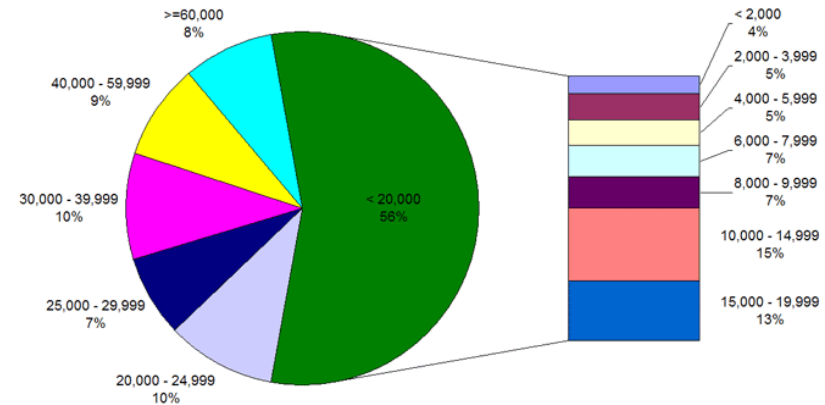
Covers M1 in local secondary school, but goes well beyond it in depth, particularly theoretical derivations. Strongly encourage you to attend ALL lessons irrespective of whether you have taken M1 or not.
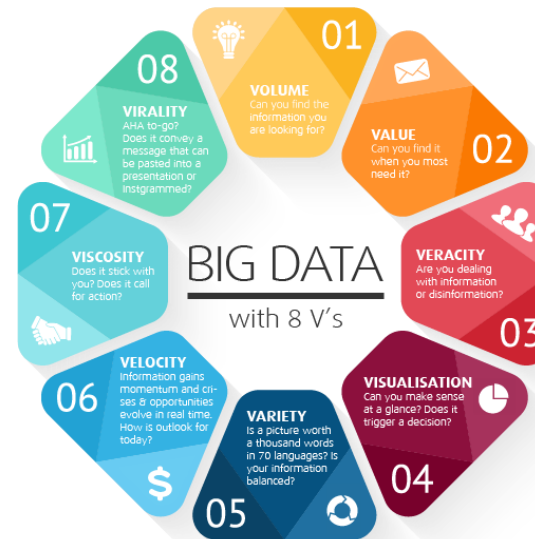
# Data helps us to identify problems





Domestic households by monthly income



港鐵沿線屋苑住戶月入中位數

# Data helps us to identify correlations and trends

# Data helps us to make predictions

Predictions: global temperature rise, warming oceans, shrinking ice sheets, sea level rise, extreme events, ocean acidification, biodiversity loss, …

Source: https://climate.nasa.gov/evidence/

# Beware of incorrect use/pitfalls/limitations in using data

Example 1  Decision making misled by meaningless correlation

Pigeon superstition (Skinner's experiment)

# Example 2  Incorrect prediction of trend



| High bias (underfit) | "Just right" | High variance (overfit) |
| --- | --- | --- |
| $\theta_0 + \theta_1 x$ | $\theta_0 + \theta_1 x + \theta_2 x^2$ | $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$ |
| Incorrect | Correct | Incorrect |

# Example 3  Identifying meaningful association while failing to find the underlying common cause

Bait Shyness in rats
(Garcia et al.'s experiment)

Pigeon superstition
(Skinner's experiment)



For more details, see Machine learning today and tomorrow: a panel discussion, 2017 youtube (8:20 – 20:05)

# Example 4  Misinterpretation of statistics



The 95% confidence interval does not mean that there is a 95% probability that the value of the variable lies within the interval

Example 5  Using a normal distribution when the distribution is not normal

Something you would understand well if you pay attention in the course:

An example of normal distribution is the grade distribution in a large class.  Do you know what assumptions are behind it?

The number of attempts required to obtain a PASS in this course do not follow the normal distribution. Why?

# Probability & Statistics

Fair Coin

Biased Coin

**Probability**
Given model, predict data

**Statistics**
Given data, predict model

also discuss the applications in
Data Engineering
as well as uses and misuses

# Motivations for studying probability and statistics

- Probability
  - It is a general mathematical tool useful in other subjects.
  - It gives you better reasoning power in making decisions that involve risk taking
- Statistics
  - There are many statistics in our daily life. It gives you an understanding of how are they arrived at and how they can be used
  - It equips you on how to correctly use the data and avoid misuse of the data

# Course Content

1. Introduction
2. Probability
3. Random Variables and Expectation
4. Statistical Descriptors
5. Special Random Variables
6. Parameter Estimation
7. Hypothesis Testing
8. Non-parametric Statistics
9. Linear Regression and Other Prediction Methods
10. Multiple Comparison Tests

# Assessment and Schedule

**Assessment Tasks/Activities (ATs)**

*(ATs are designed to assess how well the students achieve the CILOs.)*

| Assessment Tasks/Activities | CILO No. | | | | | | Weighting* | Remarks |
|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | | |
| Continuous Assessment: 50% | | | | | | | | |
| Tests (min.: 2) | ✓ | ✓ | ✓ | ✓ | | | 30% | |
| #Assignments (min.: 3) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | 20% | |
| Examination: 50% (duration:  2hrs    , if applicable) | | | | | | | | |
| Examination | ✓ | ✓ | ✓ | ✓ | ✓ | | 50% | |

\* The weightings should add up to 100%.          100%

**Remark:**

To pass the course, students are required to achieve at least 30% in course work and 30% in the examination.

# may include mini projects, in-class assignments, and homework assignments.

# Coursework Components (50%)

| Time | Item | Scope | Percentage |
|---|---|---|---|
| **Wk 6** | **Quiz 1** | everything taught in Wk 1-5 | 15% |
| **Wk 12** | **Quiz 2** | everything taught in Wk 6-11 | 15% |
| | **In-Class Assignments and/or Assignments** | | 20% |

**In-class assignments** refer to assignments conducted during lecture

# Grade distribution Report in 2018/19

# Teaching Assistant 1

Ms. CHEN Xueli

Office: FYW2386   Tel: 3442 2615

Email: xuelichen3-c@my.cityu.edu.hk

# Teaching Assistant 2

Mr. JIANG Mingjie



Office: FYW2384    Tel: 3442 4119

Email: minjiang5-c@my.cityu.edu.hk

# Related Data Engineering Courses



Structure and Flowchart for BEng in Computer and Data Engineering
2019/20 Entering Major⌗

Updated on: 2 July 2020

# Relationship with EE3211 Modelling Techniques

- You will learn the basic theory of linear regression and logistic regression in the EE3001 course

- In EE3211, you will learn how to interpret outputs from statistical programs (e.g. R, SAS). They will have hands on experience in the tutorial session to solve two real world problems below

## EE3211 Solve the following real world problem using linear regression: Pediatrics Hypertension

- Investigate how the relationship between the blood-pressure levels of newborns and infants relate to subsequent adult blood pressure.

- One problem that arises is that the blood pressure of a newborn is affected by several extraneous factors that make this relationship difficult to study.

- In particular, newborn blood pressures are affected by:
  (1) Birthweight
  (2) the day of life on which blood pressure is measured

# EE3211 Solve the following real world problem using logistic regression: Infectious Disease

- *Chlamydia trachomatis :* microorganism that has been established as an important cause of nongonococcal urethritis, pelvic inflammatory disease

- A study of risk factors for *C. trachomatis in* 431 female college students

- Because multiple risk factors may be involved, <u>several risk factors must be controlled for</u> simultaneously in analyzing variables associated with *C. trachomatis*

# Relationship with EE4146 Date Engineering and Learning Systems

- You will learn in details the Bayes' rule and also briefly learn the concept of maximum likelihood estimation

- In EE4146, you will apply Bayes' rule to classification problems and learn how to estimate the parameters of the classification model using maximum likelihood estimation

# Extract from EE4146 notes: Applying Bayes' formula for classification

## Generative vs. Discriminative Models

- In both types of models, the goal is to determine the posterior probability of different classes.
- And then assign x to a class with highest $f(C_k|\boldsymbol{x})$ (See 1.3.6)

$$f(C_1 \mid x) = \frac{f(x \mid C_1)f(C_1)}{f(x \mid C_1)f(C_1) + f(x \mid C_2)f(C_2)}$$

$$= \frac{1}{1+\exp(-a)} = \sigma(a) \quad \text{where} \quad a = \ln\frac{f(x \mid C_1)f(C_1)}{f(x \mid C_2)f(C_2)}$$

- In a generative model we estimate the class-conditionals $f(\boldsymbol{x}|C_k)$ (which are used to determine $a$).
- In the discriminative approach we directly estimate $a$ as a linear function of $\mathbf{x}$ i.e., $a = \mathbf{w}^T\mathbf{x} + b$.

6

## Example: Two class case, Gaussian: Step 1 (4)

### Maximum likelihood estimate (MLE) of $\pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2$

#### Estimates for prior probabilities

Log likelihood function that depend on $\pi$ are $\sum_{n=1}^{N} \{t_n \ln \pi + (1 - t_n) \ln(1 - \pi)\}$

MLE for $\pi$ is
Fraction of points

Setting derivative to zero and rearranging $\pi = \frac{1}{N} \sum_{n=1}^{N} t_n = \frac{N_1}{N_1 + N_2}$ where $N_1$ is no fo
data points in class $C_1$ and $N_2$ in class $C_2$.

#### Estimates for class means

Now consider maximization w.r.t. $\mu_1$. Pick log likelihood function depending only on $\mu_1$

$$\sum_{n}^{N} t_n \ln \mathcal{N}(x_n | \mu_1, \Sigma) = -\frac{1}{2} \sum_{n}^{N} t_n (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) + \text{const}$$

Setting derivative to zero and solving $\mu_1 = \frac{1}{N_1} \sum_{n=1}^{N} t_n x_n$

Mean of all input vectors
$x_n$ assigned to class $C_I$

Similarly $\mu_2 = \frac{1}{N_2} \sum_{n=1}^{N} (1 - t_n) x_n$

17

# Text book and References

- Text book

  S.M. Ross, Introduction to Probability and Statistics for Engineers and Scientists, 5th Edition, Elsevier 2014.

- References

  J.P. Marques De SÁ, Chance, the life of games & the game of life, Springer 2008.

  S.M. Ross, A First Course in Probability, 9th Edition, Pearson 2014.

**E-books** of the above are available in CityU Library. They can be found in EE3001 Course Reserve.

# Copyright Issues

The copyright of some text and images in the notes belong to the authors and should not be used elsewhere without permission. Other images are usually downloaded from Google images or from Wikipedia