# 從零開始的關聯式學習

## Pandas 與 Mlxtend

**張家瑋** 博士

副教授

國立臺中科技大學資訊工程系

A.I.
Big Data
Images
Videos
IoT
Audios
Texts

**Work Experience**

➢ 2022/2~ Now
   **Associate Professor**
   National Taichung University of
   Science and Technology
➢ 2018/2 ~ Now
   **Adjunct Assistant Professor**
   National Cheng Kung University
➢ 2015/8 ~ 2017/11
   **Project Manager & Data Scientist**
   NEXCOM International Co., Ltd.

# About Me

- [Since Jan. 2019] Young Professionals Chair, IET Taipei Local Network.
- [Since Dec. 2017] Consultant, NEXCOM Industry 4.0 Center.
- [Jan. 2017] Ph.D. degree, National Cheng Kung University.

# Research Topics

a) Natural Language Processing
   ✓ Natural Language Understanding
   ✓ Chatbot
   ✓ Text Summarization / Classification
b) Deep Learning
c) Data Mining
d) Internet of Things
   ✓ Smart Speaker
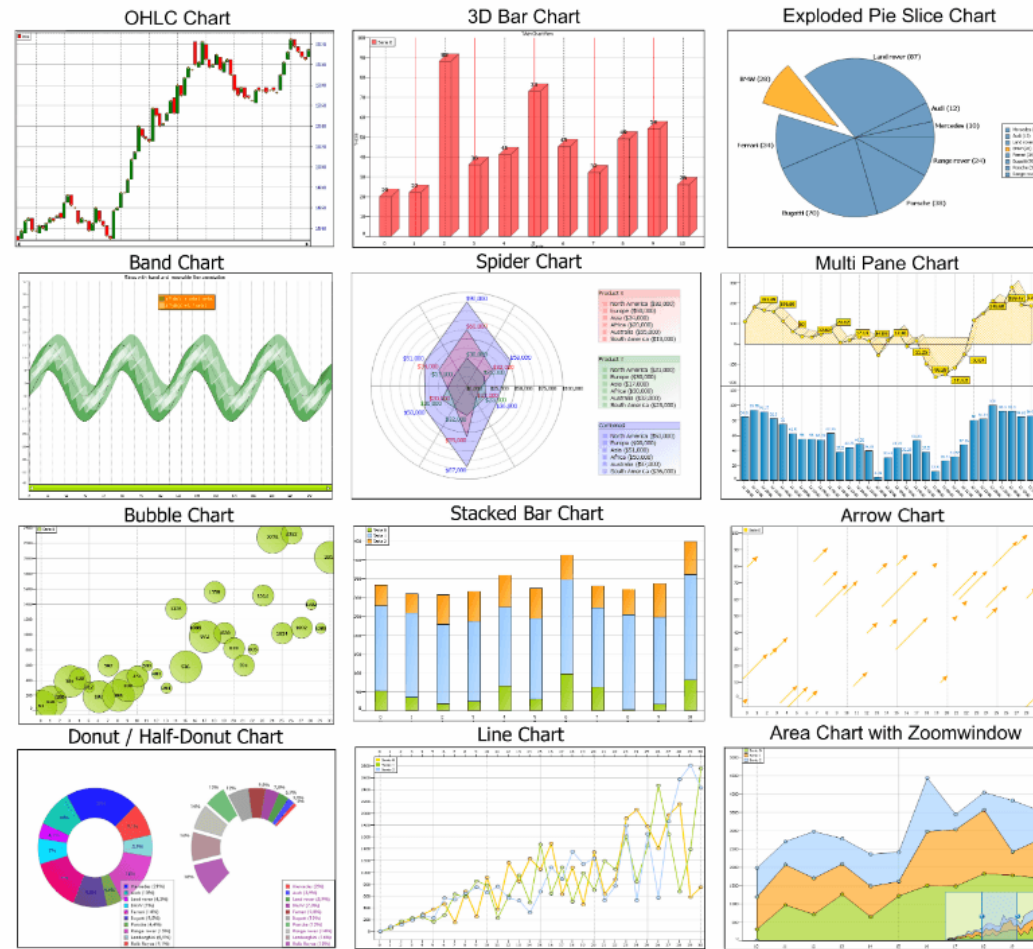
# AI

**Basic Background
of Data Science**

# Data?



- 未經過處理的原始記錄。
- 資料**缺乏組織及分類**，無法明確的表達事物代表的意義。
- 資料是關於事件的一組**離散**且**客觀**的事實描述，是構成資訊和知識的原始材料。

# Information?

- 資訊是經過處理後的資料。

- 資訊是有用的或有意義的資料。

- 對接受者**有意義的資料能使接受者產生資訊**。

# Information?

# Knowledge?

- 知識是資訊、文化脈絡以及經驗的**整合**。

- 知識是對**某個主題**確信的認識，並且這些認識擁有潛在的能力為**特定目的**而使用。

- 藉由**專業技能**或**豐富經驗**用以分析資訊的結果。

# Intelligence?

- 以知識為基礎，**運用**個人能力，

  **實踐**能力來開創價值。

- 分析、判斷、創造、思考的能力。

- 智慧具有反應能力與價值判斷。

# 人工智慧？

- 指由人製造出來的機器所表現出來的智慧。(Wiki)

- 弱人工智慧 → 專家系統

  – 處理特定的問題

- 強人工智慧

  – 通用人工智慧

# 過去

- Small Data
  - 針對某一個問題，只能獲得**小量數據**。

  - 數百筆到數萬筆。

  - 花費**大量人工編碼**。

# 過去

- Small Data 統計分析

  – 樣本推論母體**(抽樣)**

  – 在小樣本中，需要發展一系列理

    論來解釋事物的原理**(學說)**

  – [啟示] 1936 羅斯福與藍頓 的民調

# 過去

- 小數據

- Rule-based AI

- 類神經網路 (1980s)

# 現在與未來

- Big Data

  - 由**"母體"**來分析數據

  - 數萬筆到幾近無限

  - **雜亂**的原始資料

# 上一世代

- 大數據

- 分類：SVM → 機器學習

- 分群：Kmeans

- 關聯式法則：Apriori

# What's difference?

- Small Data vs Big Data

  – 都有目的或待解的問題

But

  – **減少假設**

  – **力求呈現真實世界**

# What's difference?

- 資料可重組與檢視關聯。

- 接受「**數據的雜亂性**」，不再追求「精確」的數據。

- **重「相關」而輕「因果」**。

# 現在

- 大數據

- 運算力的提升

- 深度學習 (強AI的可能性)
  - 類神經網路的文藝復興

- 演進趨勢
  - 腦神經科學
  - 認知科學
  - 認知心理學

# Big Data 的沿革 (1/3)

- Data Mining

    - 資料探勘是利用分析技術來發掘資料間**未知的關聯性與規則**。

    - [少女未婚懷孕　購物商場比老爸還早知道？！](#)

        - https://www.nownews.com/news/20120223/42676

# Data Mining

✓ 分群

 – 用於沒有標籤的資料，又通常為非監督式演算法。

✓ 分類

 – 用於有標籤的資料，又通常為監督式演算法。

✓ 關聯式法則

 – 有序性規則的資料

# Data Mining

✓ 分群

– 用於沒有標籤的資料，又通常為非監督式演算法。

# Data Mining

✓ 分群

 − 用於沒有標籤的資料，又通常為非監督式演算法。

# Data Mining

✓ 分類

  – 用於有標籤的資料，又通常為監督式演算法。

# Data Mining

✓ 關聯式法則

– 有序性 (尿布與啤酒)

# Big Data 的沿革 (2/3)

- Machine Learning

  - 人工智慧的分支，可用於資料探勘。

  - 讓機器可以**自動學習**、從巨量資料中找到規則，
  進而有能力做出**分類或預測**。

    - 判斷出類別

    - 估計出數值

# Big Data 的沿革 (3/3)

- Deep Learning

  – 是機器學習的分支

  – **類神經網路**的文藝復興

  – 從**大規模未標記資料**中建立更好的預測模型

  – 建立強 AI 的可能性

# 資料分析的基本步驟

1. 資料清除：去除極端、遺失值資料、不重要的屬性

2. 資料整合：因應用目的或特性，整合不同來源的資料

3. 資料選擇：揀選重要的屬性來逼近目的之最佳成效

4. 資料轉換：基於領域知識進行特徵縮放、數值類別轉換等

5. 資料探勘：選用合適的分析演算法得到目的之結果

6. 樣式評估：評估結果的樣式，是否如預期

7. 知識表示：因應目的將樣式轉換成合適的表達方法

# 資料分析的演算法重點

- 預處理（Preprocessing）

- 降維（Dimensionality Reduction）

- 模型選擇（Model Selection）
  - 監督式學習（Supervised learning）
    - 分類（Classification）：機器給出一個類別
    - 迴歸（Regression）：機器給出一個數值
  - 非監督式學習（Unsupervised learning）
    - 分群（Clustering）
  - 關聯式學習（Association Rule Learning）

# 資料分析的常見角色

- 資料產品經理人：將真實世界的問題轉換成資料可以解決的問題，通常是該問題領域的專業人士

- 資料工程師：蒐集、整理、清理資料，通常是具備程式技術能力的工程師

- 資料分析師：負責資料建模和分析，通常由擅長找出資料關聯的統計人擔當

- 資料視覺化設計師：將報表變得簡明易懂

# 資料驅動創新應用

- 文字、聲音、影像

  – 自然語言處理

  – 語音辨識

  – 影像辨識

- 數值與非數值

  – 連續性

  – 離散性、類別

# AI

**Open Datasets**

# UC Irvine Machine Learning Repository

http://archive.ics.uci.edu/ml/datasets.html

# Kaggle DataSets

https://www.kaggle.com/datasets

# 臺南市開放資料

http://data.tainan.gov.tw/dataset

# AI

關聯規則學習
Association Rule Learning

# 概念

- 在大型資料庫中發現項目間關聯的方法。

  - {牛奶, 麵包}→{可樂}：代表某人同時買了牛奶和麵包，就可能會買可樂。

- 該方法常使用於電子商務上，通常可為**促銷**、**產品推薦**等行銷活動的決策依據。

# 定義

- 商品的項目集合(itemset)，$I = \{ I_1, I_2, ..., I_m \}$。#Item

- 交易資料庫(Database)，$D = \{ t_1, t_2, ..., t_n \}$。 #Transaction

- 關聯規則(Association Rule)，$X \rightarrow Y$

Apriori

# 概念

- 逐層搜索的迭代方法。

- $k$-itemset 用於探索（$k+1$）- itemset。

  1. 找出 frequent 1-itemset，$F_1$。$F_1$ 用來找 frequent 2-itemset，$F_2$。而 $F_2$ 用來找到 $F_3$。直到不能找到 $k$-itemset。

  2. 每找一個 $F_k$ 需要掃描一次資料庫。為提高頻繁項集逐層產生的效率，Apriori 性質則可減少搜索。

- Apriori 性質：frequent itemset 的所有非空子集都必須是頻繁的。

  - 若某個 $k$-itemset 的 candidate 的 subsets 不在 $(k$-1$)$-itemset 時，這個 candidate 就可以直接刪除。

# 當最小支持度為 2 時的情況

**Database**

| TID | Items |
|-----|-------|
| 10 | A, C, D |
| 20 | B, C, E |
| 30 | A, B, C, E |
| 40 | B, E |

**1st scan**

$C_1$

| Itemset | Support |
|---------|---------|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {D} | 1 |
| {E} | 3 |

$F_1$

| Itemset | Support |
|---------|---------|
| {A} | 2 |
| {B} | 3 |
| {C} | 3 |
| {E} | 3 |

$C_2$

| Itemset |
|---------|
| {A,B} |
| {A,C} |
| {A,E} |
| {B,C} |
| {B,E} |
| {C,E} |

$C_2$

| Itemset | Support |
|---------|---------|
| {A,B} | 1 |
| {A,C} | 2 |
| {A,E} | 1 |
| {B,C} | 2 |
| {B,E} | 3 |
| {C,E} | 2 |

**2nd scan**

$F_2$

| Itemset | Support |
|---------|---------|
| {A,C} | 2 |
| {B,C} | 2 |
| {B,E} | 3 |
| {C,E} | 2 |

$C_3$

| Itemset |
|---------|
| {B,C,E} |

**3rd scan**

$F_3$

| Itemset | Support |
|---------|---------|
| {B,C,E} | 2 |

# 方法

*1. $C_3 = F_2$ 的組合*

- $F_2$ = {{A, C}, {B, C}, {B, E}, {C, E}}

  $C_3$ = {{A, B, C}, {A, C, E}, {B, C, E}}

2. 使用 Apriori 性質剪枝：某個 frequent itemset 的所有 subsets 必須是頻繁的，對 candidate itemset $C_3$，我們可以刪除其非頻繁的 subsets：

- {A, B, C} 的 2-itemset 是 {A, B}, {A, C}, {B, C}，其中 {A, B} 不是 $F_2$ 的元素，所以刪除;

- {A, C, E} 的 2-itemset 是 {A, C}, {A, E}, {C, E}，其中 {A, E} 不是 $F_2$ 的元素，所以刪除;

- {B, C, E} 的 2-itemset 是 {B, C}, {B, E}, {C, E}，所有 2-itemset 都是 $F_2$ 的元素，因此保留。

3. 剪枝後得到 $C_3$ ={{B, C, E}}

# 剪枝

if {A,B} is infrequent

# 案例

| TID | 網球拍 | 網 球 | 運動鞋 | 羽毛球 |
|-----|--------|-------|--------|--------|
| 1 | 1 | 1 | 1 | 0 |
| 2 | 1 | 1 | 0 | 0 |
| 3 | 1 | 0 | 0 | 0 |
| 4 | 1 | 0 | 1 | 0 |
| 5 | 0 | 1 | 1 | 1 |
| 6 | 1 | 1 | 0 | 0 |

- 顧客購買記錄的資料庫 *D*，包含 6 個 Transactions
- 項目集 *I* = {網球拍, 網球, 運動鞋, 羽毛球}

觀察關聯規則，網球拍 → 網球。
1. Transaction 1, 2, 3, 4, 6 包含網球拍。
2. Transaction 1, 2, 6 同時包含網球拍和網球。
3. 支持度 = 3/6 = 0.5，信心度 = 3/5 = 0.6。

- 若最小支持度為 0.5，最小信心度為 0.6。
- 關聯規則 "網球拍→網球" 是存在強關聯的。

- 1-itemset (4): {網球拍}, {網球}, {運動鞋}, {羽毛球}
- 2-itemset (7): {網球拍, 網球}, {網球拍, 運動鞋}, {網球拍, 羽毛球}, {網球, 運動鞋}, {網球, 羽毛球}, {運動鞋, 羽毛球}
- 3-itemset (4): {網球拍, 網球, 運動鞋}, {網球拍, 網球, 羽毛球}, {網球拍, 運動鞋, 羽毛球}, {網球, 運動鞋, 羽毛球}

實作開始

# Google Colab



https://colab.research.google.com/

# MIxtend

min_support : float (default: 0.5)
✓ A float between 0 and 1 for minumum support of the itemsets returned.

transactions_where_item(s)_occur / total_transactions

$$\text{support}(A \rightarrow C) = \text{support}(A \cup C), \quad \text{range: } [0, 1]$$

$$\text{confidence}(A \rightarrow C) = \frac{\text{support}(A \rightarrow C)}{\text{support}(A)}, \quad \text{range: } [0, 1]$$

# MIxtend

$$\text{lift}(A \rightarrow C) = \frac{\text{confidence}(A \rightarrow C)}{\text{support}(C)}, \quad \text{range: } [0, \infty] \quad = \quad \frac{Support}{Supp(X) \times Supp(Y)}$$

$$\text{levarage}(A \rightarrow C) = \text{support}(A \rightarrow C) - \text{support}(A) \times \text{support}(C), \quad \text{range: } [-1, 1]$$

$$\text{conviction}(A \rightarrow C) = \frac{1 - \text{support}(C)}{1 - \text{confidence}(A \rightarrow C)}, \quad \text{range: } [0, \infty] \quad = \quad P(A)P(B')/P(A \cap B')$$

# UCI - Online Retail Data Set



https://archive.ics.uci.edu/ml/datasets/online+retail

**Machine Learning Repository**
Center for Machine Learning and Intelligent Systems

Check out the beta version of the new UCI Machine Learning Repository we are currently testing! Contact us if you have any issues, questions, or concerns. Click here to try out the new site.

## Online Retail Data Set

Download: Data Folder, Data Set Description

Abstract: This is a transnational data set which contains all the transactions occurring between 01/12/2010 and 09/12/2011 for a UK-based and registered non-store online retail.

| Data Set Characteristics: | Multivariate, Sequential, Time-Series | Number of Instances: | 541909 | Area: | Business |
|---|---|---|---|---|---|
| Attribute Characteristics: | Integer, Real | Number of Attributes: | 8 | Date Donated | 2015-11-06 |
| Associated Tasks: | Classification, Clustering | Missing Values? | N/A | Number of Web Hits: | 764798 |

### Source:

Dr Daqing Chen, Director: Public Analytics group. chend '@' lsbu.ac.uk, School of Engineering, London South Bank University, London SE1 0AA, UK.

### Data Set Information:

This is a transnational data set which contains all the transactions occurring between 01/12/2010 and 09/12/2011 for a UK-based and registered non-store online retail.The company mainly sells unique all-occasion gifts. Many customers of the company are wholesalers.

### Attribute Information:

InvoiceNo: Invoice number. Nominal, a 6-digit integral number uniquely assigned to each transaction. If this code starts with letter 'c', it indicates a cancellation.
StockCode: Product (item) code. Nominal, a 5-digit integral number uniquely assigned to each distinct product.
Description: Product (item) name. Nominal.
Quantity: The quantities of each product (item) per transaction. Numeric.
InvoiceDate: Invice Date and time. Numeric, the day and time when each transaction was generated.
UnitPrice: Unit price. Numeric, Product price per unit in sterling.
CustomerID: Customer number. Nominal, a 5-digit integral number uniquely assigned to each customer.
Country: Country name. Nominal, the name of the country where each customer resides.

# AI

懒人包
[ link ]

Thank you