

DATA ANALYTICS

Statistics
AI, Machine Learning
Data Base - data warehouse OLAP

Type of Data

Record → Table, Relation, Data Matrix

Data Matrix ▶▶ มองเป็นมิติ ขนาด $m \times n$

m row → object (tuple)

n column → attribute

Document Data ▶▶ เก็บแบบ keyword

Transaction Data ▶▶ บันทึกข้อมูลทั้งชุดรายการที่สนใจ

ex. บันทึกการ shopping ของลูกค้า
คนหรืออะไรไปบ้าง

Record Data ▶▶ เก็บเป็นรายการโดยกำหนด attribute

Graph & Network → world wide web, หน้าเว็บเพจ

Ordered → sequence, time-series, ลำดับสำคัญ!

ex. ร้าน bakery เก็บข้อมูล ordered data ของช่วงเช้า, เช่น
ดูตาม Timeline (Item/Even)

Big Data 4 V

Volume High Volume data เพิ่มขึ้นแบบ expo (Big Scale)

Velocity Speed เข้าถึงข้อมูลได้รวดเร็ว ไป-กลับรวดเร็ว

+ Parallelization ทำงานแบบคู่ขนาน ตรงกับ Sequential (ตามลำดับ)

Variety complexity ex. Relation data, Text data (web), Big Public Data
many feature per item, irregular structure

Veracity trustworthy, reliability, completeness

03 Data Analytics & Introduction to Statistics

Data + IT + Statistical Analysis + Quantitative Method
+ math + computer based model

managers gain improved insight make better

Scope of Business Analytics

Descriptive "past & present" what has happen?

พรรณนา ⇒ เอาข้อมูลอดีต + ปัจจุบันมาวิเคราะห์ สร้างเช่น visualization

Predictive "future" what could happen? จะเกิดอะไรขึ้นต่อ

พยากรณ์ ⇒ วิเคราะห์ข้อมูลในอดีต เพื่อคาดการณ์ล่วงหน้า
regression, machine learning, neural networks

Prescriptive "outcome" what should we do? เราควรทำอะไร

ให้เกิดประโยชน์สูงสุด ⇒ วิเคราะห์หาวิธีกระทำให้เกิดประโยชน์สูงสุด

จำลอง Algorithms on possible

Introduction to Statistics: sample สุ่มมาแบบ NOT BIAS (SAVE TIME & MONEY)

2 Type → Descriptive : บอกแนวโน้ม-ไม่ดี

Mode สำนวนนิยม

Median ค่ากลาง สลับฐาน

Mean ค่าเฉลี่ย $\frac{\sum x}{n}$

SD ส่วนเบี่ยงเบนมาตรฐาน $\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$

SD ยิ่งมากค่าการกระจายมาก

SD ยิ่งน้อยค่าการกระจายน้อย

→ Inferential : พยากรณ์, คาดการณ์ คำนวณจาก sample

01 Introduction to Statis Machine Learning

Data → Analysis → Model → Predictions

Learn from sample (NOT BIAS)

	เรียนรู้แบบมีผู้สอน Supervised Learning	เรียนรู้แบบไม่มีผู้สอน Unsupervised Learning
Discrete data → I	Categorization	Clustering
Continuous data → R	regression สมการถดถอยเชิงเส้น	dimensionality reduction การลดขนาดมิติข้อมูล

Model ⇒ Visualization

function ⇒ $f(\text{apple}) = \{\text{หวาน}, \text{แดง}\}$; x data ที่ใส่มา, y output
 $y = \text{apple}$

Lost function ⇒ ค่าสูญเสีย / ค่าผิดพลาด; ยิ่งน้อยยิ่งข้อมูลแม่นยำ(ดี)

ex. Linear Regression plot

Machine Learning จะเรียนรู้ model ที่มาจากข้อมูลที่เราใส่มา
เพื่อ predict ที่มีค่า Loss function น้อยๆ จะเรียนรู้ได้แม่นยำ

02 Data Concept & Big Data

+ Collect Data Object & Attributes

+ Attribute "column"

rows → data objects "tuple", columns → attributes

+ Attribute Values : ทุ่ย่อยที่บอกลักษณะ

& Attribute same : ความสูง เป็นพหุตา/เมตรก็ได้
จะได้ยินหลายหน่วย

different : ID ไม่มีค่า max-min

Age มีค่า max-min

Discrete vs Continous Attribute

Discrete เช่นจำนวนเต็ม I เป็นค่าไม่ต่อเนื่อง
special case ▶ Binary Attribute

Continous เช่นจำนวนจริง R real number

Type of Data Measurement

+ Nominal "name of things" เอาไปคำนวณไม่ได้ จับได้ว่ามีเท่าไร

=, ≠ ex. eye-colour = {blue, black, brown}

Binary : 2 state "0 and 1" Yes/No

Symmetric : both เท่าเทียมสำคัญเท่ากัน ex. gender, hands

Asymmetric : ความสำคัญของ 2 สิ่งไม่เท่ากัน

ex. HIV positive - HIV negative

+ Ordinal ความเป็นอันดับ, order ได้

=, ≠, >, < ex. Size = {S, M, L}; S < M < L

+ Interval ช่วง range "Not true Zero"

=, ≠, >, <, +, - ex. Year, Temperature °C °F °K

+ Ration คำนวณเอาไปคำนวณได้จริง

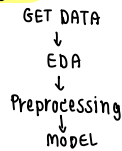
=, ≠, >, <, +, -, *, / ex. length, counts, weight, Temp K, Quantities

+ Qualitative
ปริมาณ

+ Qualitative
ค่าคงได้

05 Exploration Data Analysis

& Data Preprocessing



"EDA" การวิเคราะห์การสำรวจข้อมูล

WELL → Accuracy ความแม่นยำ

Completeness ความสมบูรณ์

Consistency ความสอดคล้อง

Timeless ไร้เวลา

Believability ความน่าเชื่อถือ

Value added คุณค่า ประโยชน์

Interpretability ข้อมูลแปลผลได้

Accessibility การเข้าถึง

Problem → Noisy & outlier: Salary = "-10", Voice flie

Dirty data → Missing value: incomplete "", n/a

→ Duplicate Data: 1 คนหลาย email

→ Inconsistent Data: data ต่าง data base

Age, Birth
Tel → String, Int

Preparation → cleaning: fill, identify, remove outlier

→ integration: เอามาเก็บในที่เดียวกัน

→ reduction & feature selection regression

→ discretization: แปลงทุกข้อมมูลเป็นจำนวนเต็ม
(continuous → discrete)

→ transformation: Z-score $z = \frac{x - \mu}{\sigma}$
จัดเรียงข้อมูล

Type of Data Format

♥ CSV: Comma Separated Value

• Excel

• JSON: JavaScript Object Notation ~ Dictionary {key: value}

• SQL: Structured Query Language

• And More! flie image, voice

PYTHON

♥ numpy → import numpy as np

♥ scipy → import scipy as sp

♥ matplotlib → import matplotlib.pyplot as plt

♥ seaborn → import seaborn as sns

♥ pandas → import pandas as pd

dataframe = pd.read_csv('file name.csv')

dataframe.columns

out: Index([...], dtype='object')

dataframe.shape

out: (2018, 18)

row column

06 Data Visualization

เข้าใจภาพ + understanding

Table

✓ specific values ค่าเฉพาะ

✓ มีความแม่นยำ precise

✓ compare related values

ต้องการเปรียบเทียบค่า

✓ different units of measure

เปรียบเทียบข้อมูลที่มีหลายหน่วย

vs.

Graph

✓ shape, bar, line

✓ relationship multiple value

ความสัมพันธ์ของหลายๆข้อมูล

✓ show trends

✓ large data sets

Bar Graph การกำหนด scale!

แกน x ข้อมูลที่ซ้ำกันบ่อยๆ

แกน y ข้อมูลที่แตกต่างกันมาก

seaborn: sns.displot(df['alcohol'])

pandas: df['alcohol'].plot.hist()

movies = ['A', 'B', 'C']

num = [2, 4, 6]

xs = range(len(movies))

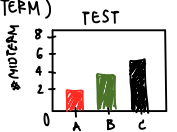
plt.bar(xs, num, color = ('r', 'g', 'b'))

plt.xticks(xs, movies)

plt.ylabel('# MIDTERM')

plt.title('TEST')

plt.show()



Pie Chart 100%

restaurant_salad = df.groupby('restaurant')['salad'].count()

plt.pie(restaurant_salad, labels = restaurant_salad.index)

df.groupby('rank')['salary'].count().plot(kind = 'pie')

bar, line, ...

Scatter Plot

sns.regplot(x = 'sodium', y = 'protein', data = df, color = 'pink')

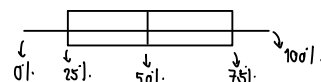
plt.scatter(df['service'], df['salary'], color = 'pink')

Bubble Plot ต้องใช้ marker

Box Plot

Quartiles → 01-min - 25% - 50% - median - 75% - 100% max

Range + Median + Percentages ต้อง outlier don



sns.boxplot(x = 'protein', y = 'restaurant', data = df)

Stacked Charts

x = ['A', 'B', 'C', 'D']

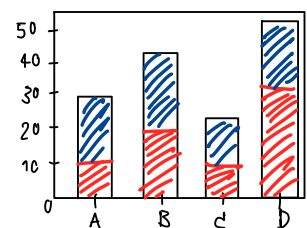
y1 = [10, 20, 10, 30]

y2 = [20, 25, 15, 25]

plt.bar(x, y1, color = 'r')

plt.bar(x, y2, bottom = y1, color = 'b')

plt.show()



Line Charts

df.groupby('rank')['salary'].count().plot(kind = 'line')

Note %matplotlib inline ช่วยจัด scale