# 下世代物聯網應用技術工作坊 - 機器學習

國立陽明交通大學MIPLab – 吳仁傑

# Course Summary

1. Validation
   - K-fold cross validation (*optional)
   - Leave-one-out cross validation
   - Holdout validation

2. Evaluation
   - Confusion matrix

# Course Summary

3. Supervised learning algorithm
   - Naive Bayes
   - Decision tree
   - Linear regression
   - Logistic regression
   - SVM
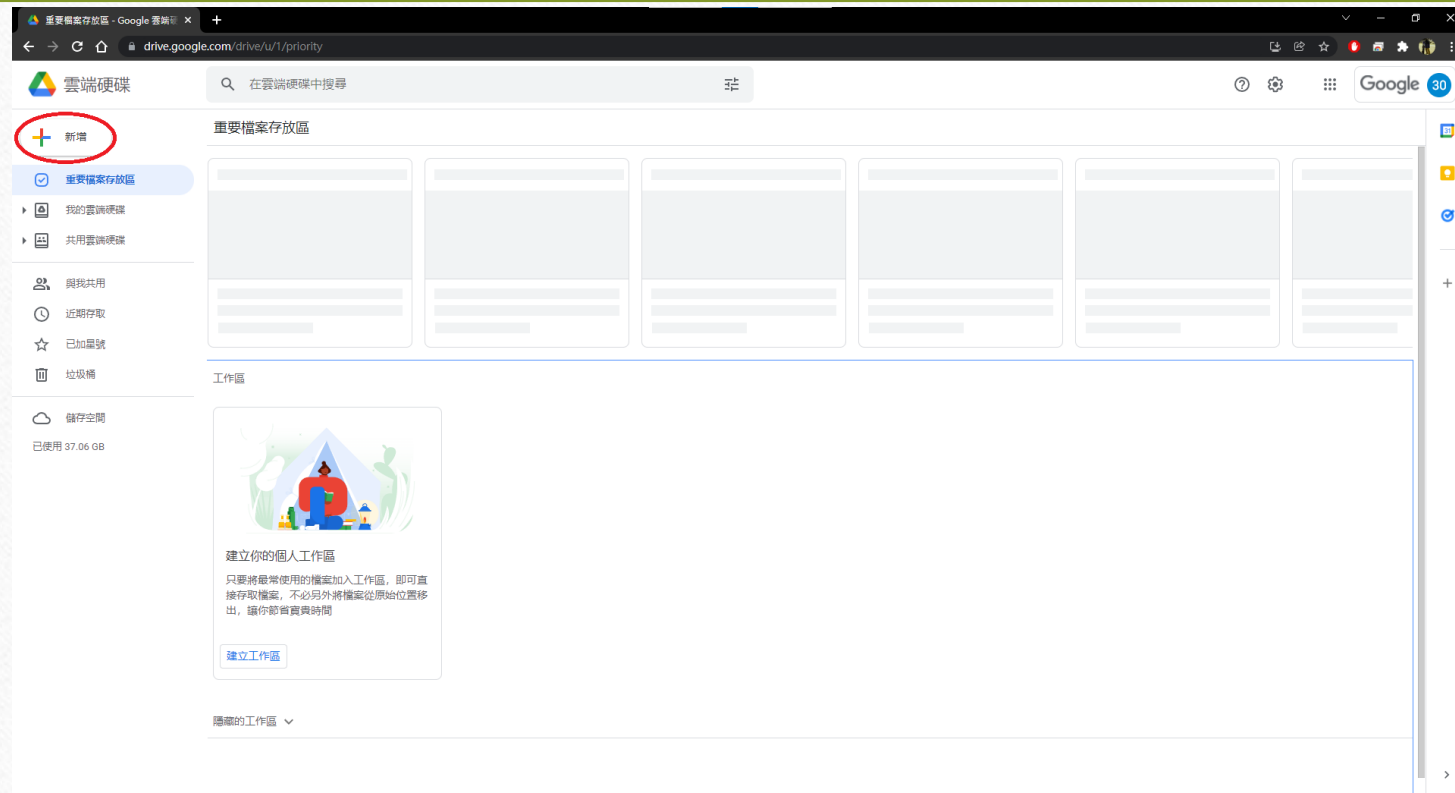
- Unsupervised learning algorithm
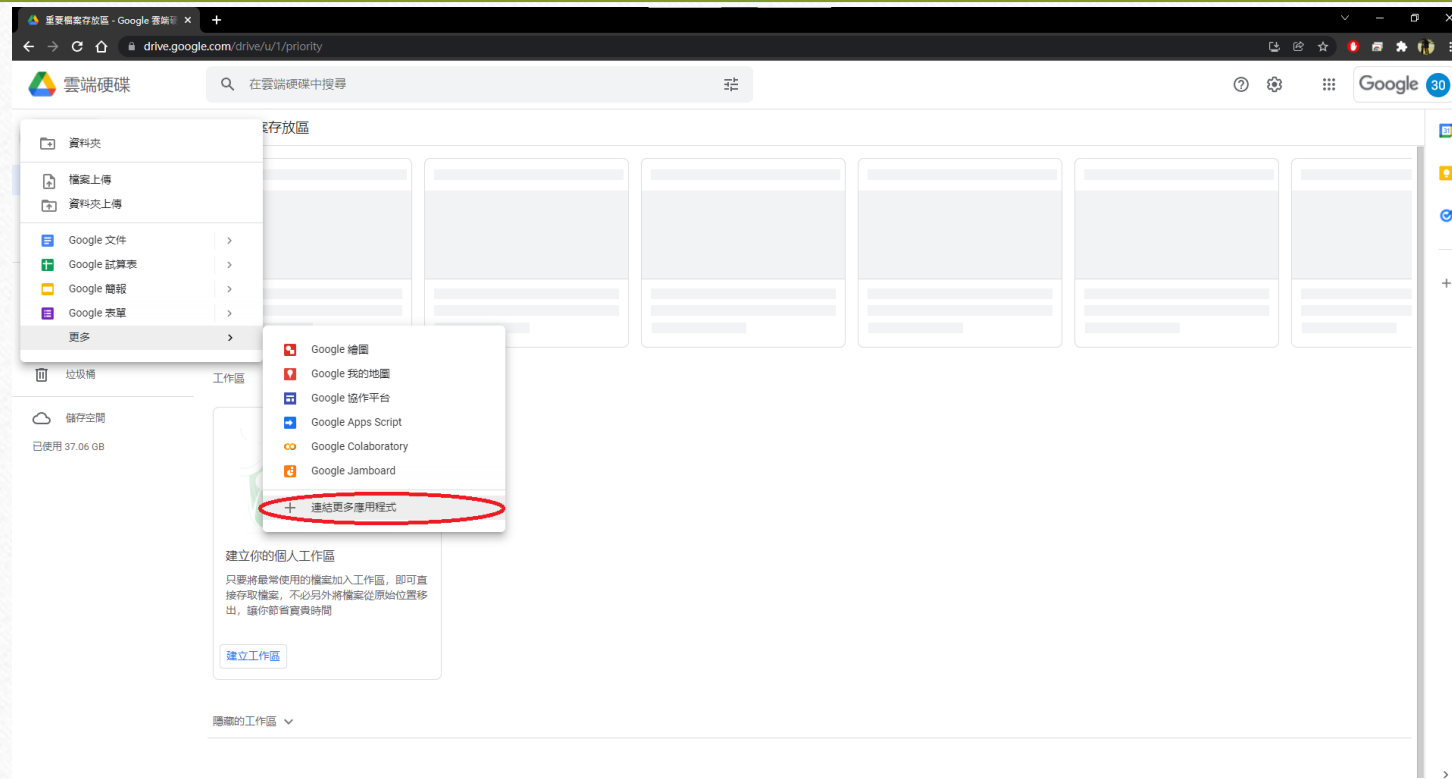  - KNN

# Setting

1. 下載教材檔案 ML_workshop.zip
   - https://drive.google.com/file/d/1GAkqaSCaGXiR1kGKrQZ4BnYFq5gFE8VL/view?usp=sharing
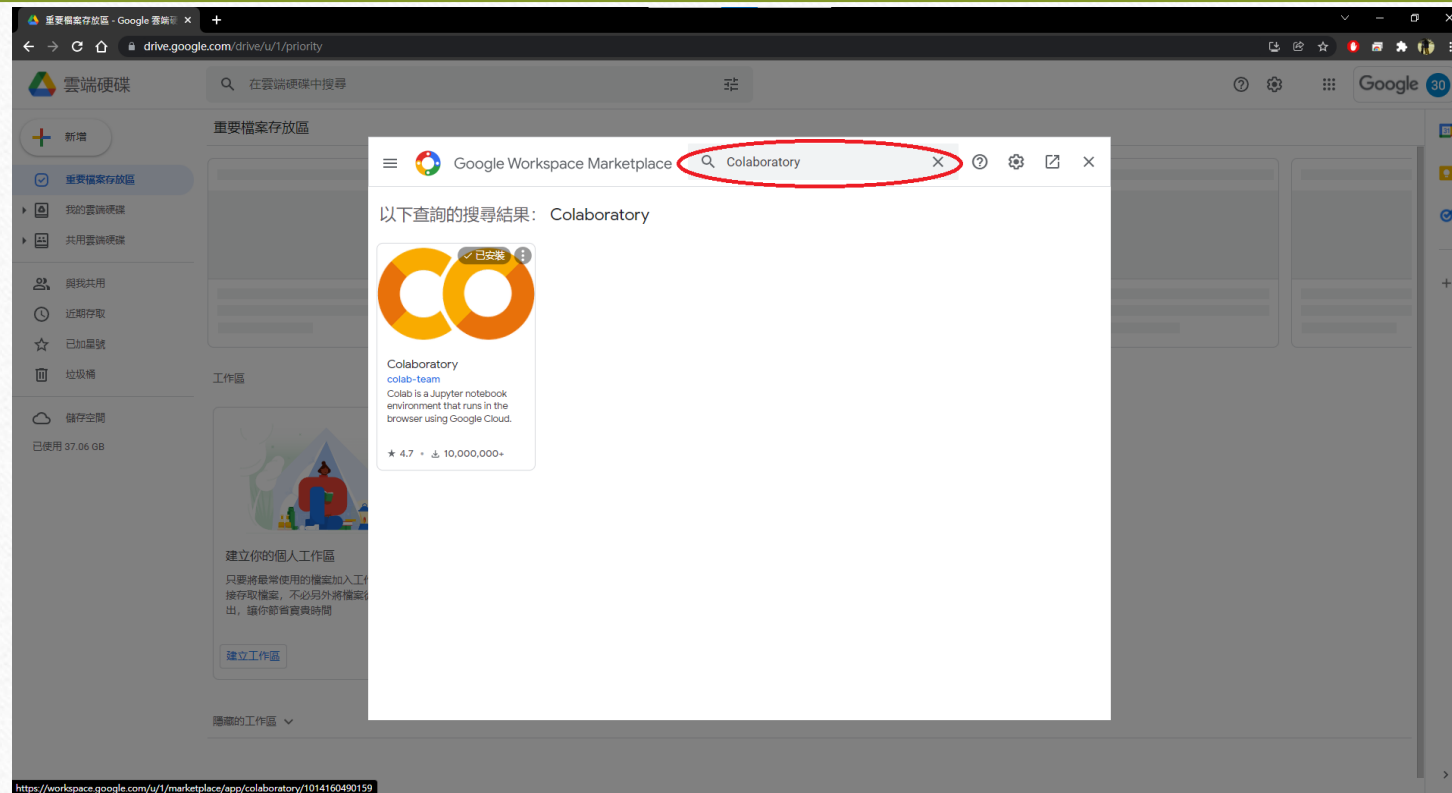2. 將檔案解壓縮後把整個ML_workshop資料夾上傳到自己的google雲端硬碟

# Setting (安裝Google Colaboratory)

# Setting (安裝Google Colaboratory)

# Setting (安裝Google Colaboratory)

# Linear regression

1. Input : $X = [x_1, \cdots, x_n]$  $(\in \mathbb{R}^n)$

   - Every $x_k$ is a quantitative variable

2. Output : $\widehat{Y} \in \mathbb{R}$

3. Parameters : $W = [w_0, w_1, \cdots, w_n]$ $(\in \mathbb{R}^{n+1})$

# Linear regression

4. Formula : $\hat{Y} = w_0 + w_1 x_1 + \cdots + w_n x_n$

5. Loss(Error) function: $L = (Y - \hat{Y})^2$

6. Optimization : $\text{argmin}_w \frac{1}{m} \sum (Y - \hat{Y})^2$

   - $m$ pieces of data

   - Find $W$ such that we have minimal average loss.

# Naive Bayes classifier

- Bayes' theorem : $P(Y|x_1, \cdots, x_n) = \frac{P(Y)P(x_1, \cdots, x_n|Y)}{P(x_1, \cdots, x_n)}$

- <span style="color:red">Naïve</span> assumption (conditional independence) :
  $P(x_i|y, x_1, \cdots, x_{i-1}, x_{i+1}, \cdots, x_n) = P(x_i|y)$

- New formula : $P(Y|x_1, \cdots, x_n) = \frac{P(Y)\Pi_{i=1}^{n}(x_i|Y)}{P(x_1, \cdots, x_n)} \propto P(Y)\Pi_{i=1}^{n}P(x_i|Y)$

# Naive Bayes classifier

- Making prediction according to the new formula :

$$\text{argmax}_Y \, P(Y|x_1, \cdots, x_n) \propto P(Y)\Pi_{i=1}^{n}P(x_i|Y)$$

1. $Y \in \{y_1, \cdots, y_k\}$ : a categorical variable

2. $x_i$ : an input feature (either categorical or quantitative)

3. $P(Y)$ : estimated by training dataset

4. $P(x_i|Y)$ : make some assumption of distribution (ex: Gaussian distribution, multinomial distribution, …)

# K-fold cross validation

- 3 roles of dataset
  - Training set : for train model
  - Validation set : for hyperparameters tuning
  - Testing set : for model performance evaluation

# K-fold cross validation

- Hyperparameters is a parameter whose value is used to control the learning process

- Find hyperparameter settings that have best average performance of each split.