



# 機器學習基礎與演算法

Chapter 1 機器學習概論

講師投影片 課程投影片 資料與程式碼 今日 Playlist

#### 「版權聲明頁」

本投影片已經獲得作者授權台灣人工智慧學校得以使用於教學用途,如需取得重製權以及公開傳輸權需要透過台灣人工智慧學校取得著作人同意;如果需要修改本投影片著作,則需要取得改作權;另外,如果有需要以光碟或紙本等實體的方式傳播,則需要取得人工智慧學校散佈權。

# 課程內容

## 1. 機器學習概論

- 1-1 [理論講授] 人工智慧與 迴歸(regression)概論
- 1-2 [實作課程] 細說基礎流程
  - Step 1 定義問題
  - Step 2 蒐集, 清理資料
  - Step 3 選擇及建立模型
  - Step 4 分析結果及修正模型
- 1-3 [理論講授] 正則化(Regularization)

#### Code 放在Hub中的course內

- 為維護課程資料, courses中的檔案皆為read-only, 如需修 改請cp至自身環境中
- 打開terminal, 輸入

cp -r courses-tpe/ML/Chapter1 <存放至本機的名稱>



# Chapter 1 機器學習概論

範例程式(example)的檔名會以藍色字體顯示且旁邊附上 如此

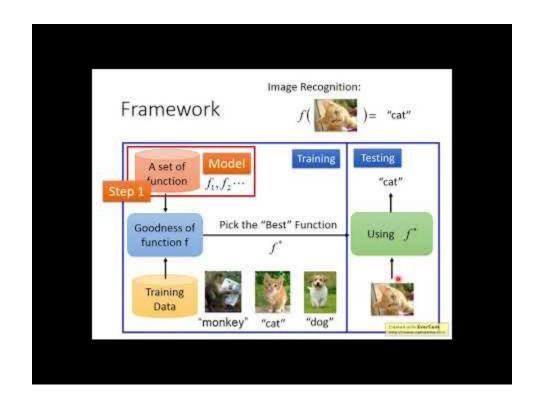


練習(exercise)的檔案以紅色字體顯示且旁邊附上



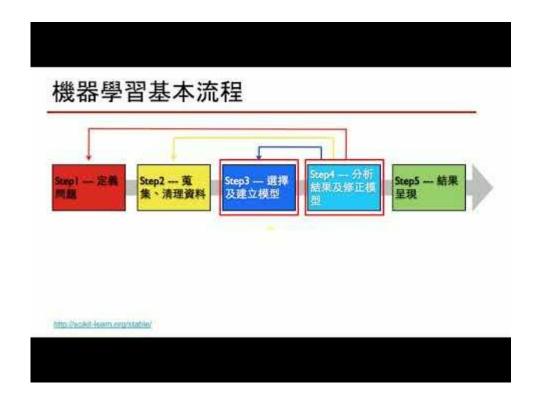
## Section 1-1 [理論講授] 人工智慧與迴歸概論







## Section 1-2 [實作課程] 細說基礎流程

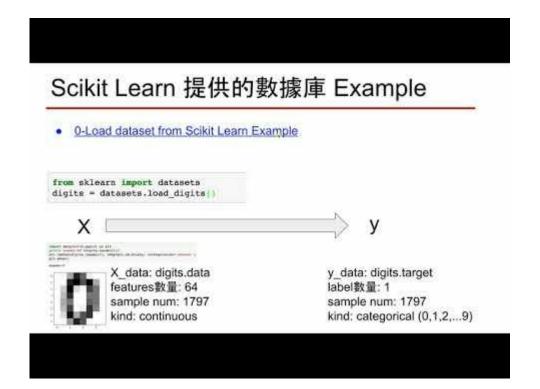




### 基礎流程

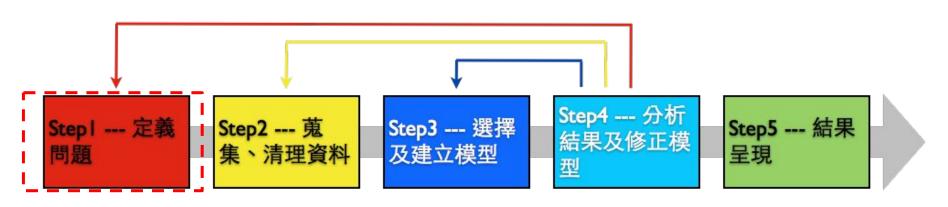


### Step 1 - 定義問題: Scikit Learn 提供的數據庫





### Scikit Learn 提供的數據庫 (1/3)



#### sklearn.datasets: Datasets

The sklearn.datasets module includes utilities to load datasets, including methods to load and fetch popular reference datasets. It also features some artificial data generators.

User guide: See the Dataset loading utilities section for further details.

#### Loaders

datasets.clear\_data\_home ([data\_home])
datasets.dump\_swalight\_file (X, y, [[, ...])
datasets.fetch\_2Onewsgroups ([data\_home, ...])
datasets.fetch\_2Onewsgroups\_vectorized ([...])

datasets.fetch\_california\_bousing([...])
datasets.fetch\_covtype[(data\_home,...))
datasets.fetch\_kdatupy9 (libsted, data\_home,...))
datasets.fetch\_ldrupy=0 (libsted,...)
datasets.fetch\_lfw\_popie (lidta\_home,...))
datasets.fetch\_ldrupe\_cove\_lddta\_home,...))
datasets.fetch\_nldate((data\_home,...))
datasets.fetch\_plote\_distare([libsted],...))
datasets.fetch\_plote\_distartibutions([...])
datasets.fetch\_popied\_distributions([...])
datasets.fetch\_popied\_distributions([...])
datasets.get\_data\_home([data\_home])
datasets.get\_data\_home([data\_home])

datasets.load breast cancer ([return\_X\_y])

Delete all the content of the data home cache.

Dump the dataset in symlight / libsym file format.

Load the filenames and data from the 20 newsgroups dataset.

Load the 20 newsgroups dataset and transform it into tf-lidf

vectors.

Loader for the California housing dataset from Statub.

Loader for the California housing dataset from Statub.

Load the covertype dataset, downloading it if necessary.

Load and return the kddcup 99 dataset (classification).

Loader for the Labeled Faces in the Wild (LFW) pairs dataset

Loader for the Labeled Faces in the Wild (LFW) people dataset

Fetch an midata.org data set Loader for the Olivetti faces data-set from AT&T. Load the RCV1 multilabel dataset, downloading it if necessary. Loader for species distribution dataset from Phillips et.

Loader for species distribution dataset from Phillips et. Return the path of the scikit-learn data dir. Load and return the boston house-prices dataset (regression). Load and return the breast cancer wisconsin dataset (classification).

#### Example:

from sklearn import datasets
digits = datasets.load\_digits()

http://scikit-learn.org/stable/modules/classes.html#module-sklearn.datasets



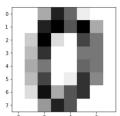
### Scikit Learn 提供的數據庫 (2/3) - Example

**0-Load dataset from Scikit Learn Example** 



```
from sklearn import datasets
digits = datasets.load digits()
```

```
import matplotlib.pyplot as plt
print('answer:%d'%digits.target[0])
plt.imshow(digits.images[0], cmap=plt.cm.binary, interpolation='nearest')
```



answer:0

X\_data: digits.data

features數量: 64

sample num: 1797

kind: continuous

y data: digits.target

label數量: 1

sample num: 1797

kind: categorical (0,1,2,...9)



### Scikit Learn 提供的數據庫 (3/3) - Exercise

#### 0-Load dataset from Scikit Learn Exercise

- 透過資料的描述對資料進行了解
- 了解4個features分別對應什麼名稱
- 了解三種labels分別對應什麼花名

	Sepal length	Sepal width	Petal length	Petal width	Class
1	5.1	3.5	1.4	0.2	Setosa
2	4.9	3.0	1.4	0.2	Setosa
50	6.4	3.5	4.5	1.2	Versicolor
150	5.9	3.0	5.0	1.8	Virginica
			1		
/	/				
					Class labels

'sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)'

> X data: iris.data features數量: 4

sample num: 150

kind: continuous

'setosa', 'versicolor', 'virginica'

y data: iris.target

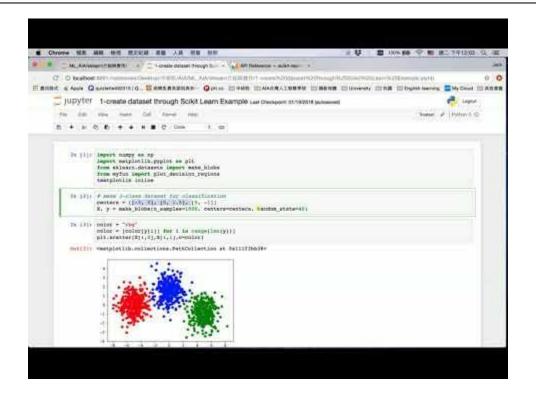
label數量: 1

sample num: 150

kind: categorical (0,1,2)

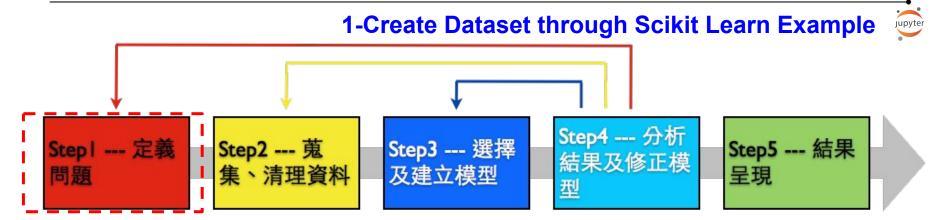


### Step 1 - 定義問題: 虛擬樣本的生成





### Scikit Learn 提供虛擬資料創建方法





biclustering. Generate isotropic Gaussian blobs for clustering. Generate an array with block checkerboard structure for Make a large circle containing a smaller circle in 2d. Generate a random n-class classification problem. Generate the "Friedman #1" regression problem Generate the "Friedman #2" regression problem Generate the "Friedman #3" regression problem Generate isotropic Gaussian and label samples by quantile Generates data for binary classification used in Hastie et al. Generate a mostly low rank matrix with bell-shaped singular Make two interleaving half circles

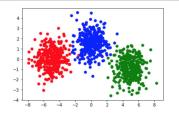
Generate a random multilabel classification problem. Generate a random regression problem Generate an S curve dataset.

Generate a sparse symmetric definite positive matrix. Generate a random regression problem with sparse uncorrelated design Generate a random symmetric, positive-definite matrix.

# make 3-class dataset for classification centers = [[-5, 0], [0, 1.5], [5, -1]]

X, y = make blobs(n samples=1000, centers=centers, random state=40)

http://scikit-learn.org/stable/modules/generated/s klearn.datasets.make blobs.html#sklearn.datase ts.make blobs

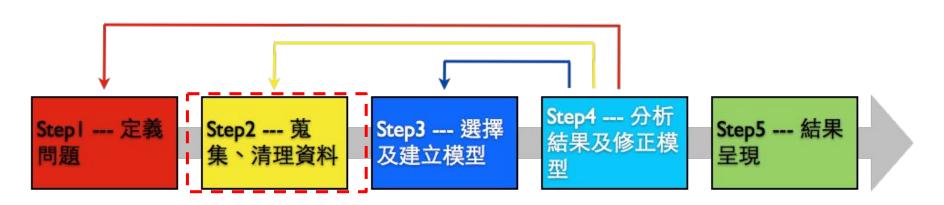


## Step 2 - 蒐集, 清理資料





#### Scikit Learn 提供資料前處理常見方法



- Training & Testing Data Split
- Data Normalization
- One Hot Encoding



### Training & Testing Data Split

# 2-Training and Testing Data Split Example

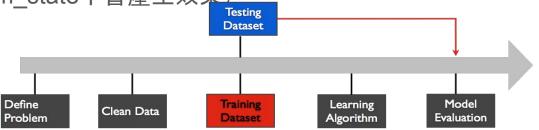


#### Import Library

調整shuffle和random state看看結果如何

from sklearn.model selection import train test split X train, X test, y train, y test = train test split(X, y, test size=0.33, random state= 42, shuffle=False) **Total Dataset** test size: 一般設定0.2~0.25 Test shuffle: default = True (將資料順序打亂) random state: 將資料打亂的方式(~random seed, 如果shuffle設定為 False, random state不會產生效果)





#### **Normalization**

#### 

 $x_2$  = number of bedrooms (1-5)  $\leftarrow$ 

 $J(\theta)$ 

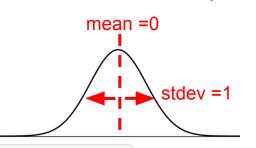


#### **Data Normalization**

#### **3- Data Normalization Example**

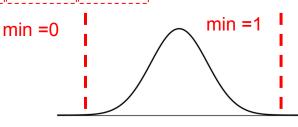


- sklearn.preprocessing.scale (Z-score)
  - (X- mean)/std
  - 將資料轉換成mean為0, stdev為1的資料分佈



	sklearn.pr	eprocessing. scale (X, axis=0, with_	_mean=True, with_std=True, copy=True)	[source
axis =0	[-100.	2.7 3.6 5. 3.6 5. 3.6 5. 3.6 3.0 3.0 3.0 3.0 3.0 3.0 3.0 3.0 3.0 3.0	[ 00.85170713 -0.55138018 -0.6 [-1.22474487 -0.55187146 -0.852133 -0.6 [ 1.22474487 1.40357859 1.40351318 1.4	57932412 <mark>]</mark>

- sklearn.preprocessing.minmax\_scale
  - $\circ \quad (X \min(X))/(\max(X) \min(X))$
  - 將資料轉換成max和min轉換到feature\_range





sklearn.preprocessing. minmax\_scale (X, feature\_range=(0, 1), axis=0, copy=True)

[source]

### Why Data Normalization?

- 提升預測準確度
- 提升模型效率

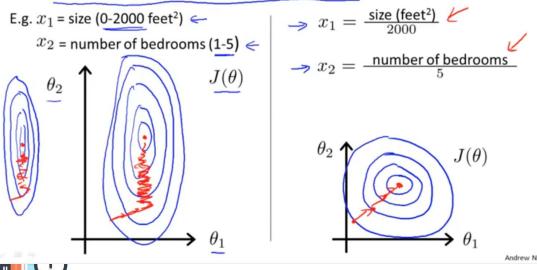
#### 4- Why Data Normalization Example in the interest of the inter

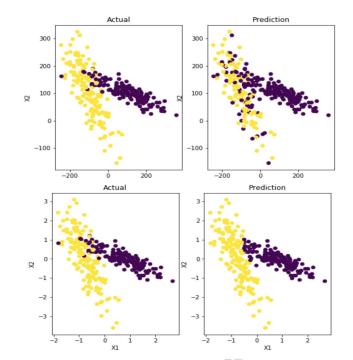
比較正規化前後的預測準確度



#### **Feature Scaling**

Idea: Make sure features are on a similar scale.



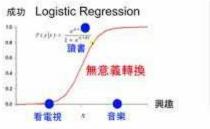


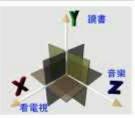
### **Label Encoding**

#### One Hot Encoding(無序類別)

 不同的興趣並沒有強度上的相關性,但regression的模型會將之視為有強度上的 差異,會導致模型無法做精準的預測。

興趣-feature(X)	看電視(1)	讀書(2)	音樂(3)
成功與否- label(y)	否(0)	是(1)	香(Ó)







#### **One Hot Encoding**

#### **5-One Hot Encoding Example**



將categorical feature轉換成數字以當作模型的輸入。

ex: 最大興趣預測未來成功與否

興趣- feature(X)	看電視	讀書	音樂
成功與否- label(y)	是	否	否

from sklearn.preprocessing import LabelEncoder

		Χ	у			Χ	У
		興趣	成功與否			興趣	成功與否
/]	、明	看電視	是		小明	1	1
/]	林	讀書	否	Label Encoder	小林	2	0
/]	英	音樂	否		小英	3	0
/	<b>、</b> 陳	看電視	是		小陳	1	1



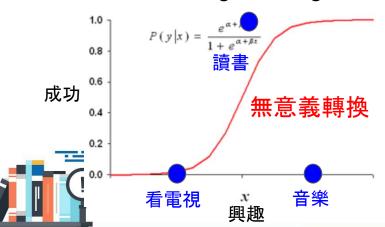
http://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.LabelEncoder

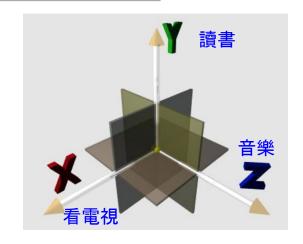
### One Hot Encoding (無序類別)(1/2)

不同的興趣並沒有強度上的相關性,但regression的模型會將之視為有強度上的 差異,會導致模型無法做精準的預測。

興趣- feature(X)	看電視(1)	讀書(2)	音樂(3)
成功與否- label(y)	否(0)	是(1)	否(0)

#### Logistic Regression





### One Hot Encoding(無序類別)(2/2)

#### 6-One Hot Encoding Example in the control of the co

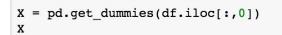


將無序的categorical feature轉換成數字以當作模型的輸入。

ex: 最大興趣預測未來成功與否

興趣- feature(X)	看電視	讀書	音樂
成功與否- label(y)	是	否	否







#### X (feature 擴充)

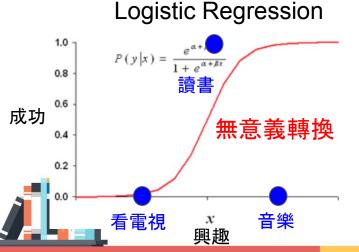
	看電視	讀書	音樂
小明	1	0	0
小林	0	1	0
小英	0	0	1
小陳	1	0	0



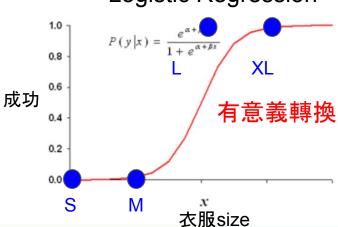
### One Hot Encoding(有序類別)(1/2)

● 衣服size具有強度上的相關性, regression的模型會將之視為有強度上的差異, 在意義上是正確的。

feature(X)	S(0)	M(1)	L(2)	XL(3)
成功與否- label(y)	否	否	是	是







台灣人工智慧學校

### One Hot Encoding(有序類別)(2/2)

#### 7-One Hot Encoding Example Jupyter



將有序的categorical feature轉換成數字以當作模型的輸入。

#### ex: 穿衣服的size預測未來成功與否

feature(X)	S	M	L	XL
成功與否- label(y)	否	否	是	是

	衣服size	成功與否
小明	XL	是
小林	S	否
小英	М	否
小陳	L	是



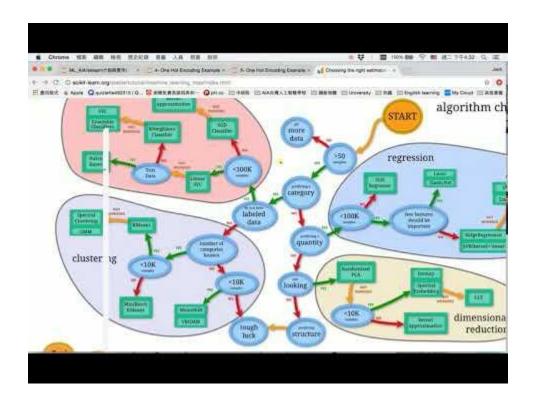


٧



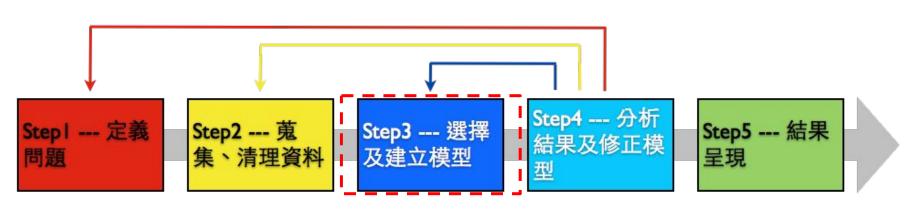
http://scikit-learn.org/stable/modules/generated/sklearn.preprocessin g.LabelEncoder.html#sklearn.preprocessing.LabelEncoder

## Step 3 - 選擇及建立模型





#### Scikit Learn 機器學習模型建立的方法



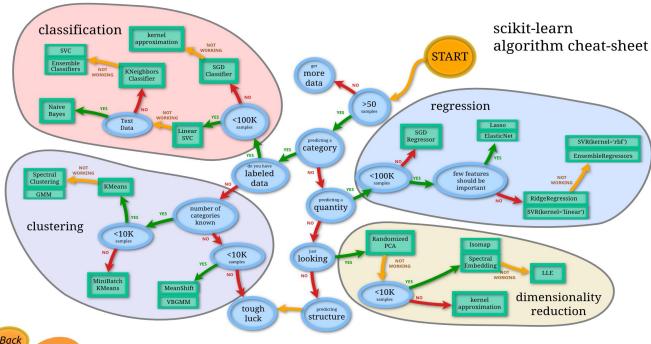
- 選擇模型 --- Machine Learning Map
- 建立模型
- 訓練模型
- 利用模型預測結果

```
# import svm classifier
from sklearn.svm import SVC
model = SVC()
model.fit(X_train, y_train)
prediction = model.predict(X)
```



### **Machine Learning Map**

http://scikit-learn.org/stable/tutorial/machine\_learning\_map/index.html







### Step 4 - 分析結果及修正模型: 迴歸型模型評估

#### Continuous

- Mean Absolute Error:
- Mean Squared Error:
  - Root Mean Squared Error:  $MSE = \frac{\sum_{i=1}^{n} e_i}{N}$
- R2 Score: coefficient of determination

 $RZ = 1 - \frac{\sum_{i=1}^{N} (e_i)^2}{\sum_{i=1}^{N} (v_i - v_i)^2} = 1 - \frac{total\ error}{data\ variance}$ 



#### 8- Model Evaluation Example

- 觀察在資料上加不同noise對預測結果造成的影響
- 了解不同metrics的意義

from sklearn import metrics

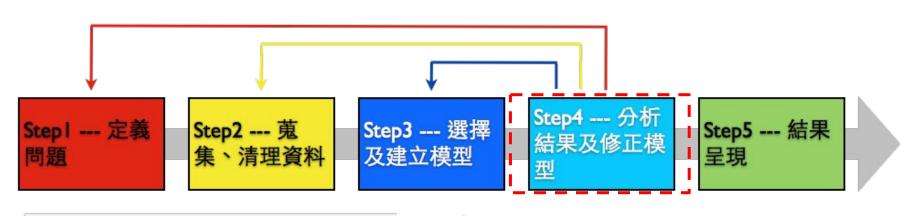
mae = metrics.mean\_absolute\_error(prediction, y)
mae = metrics.mean\_aquared\_error(prediction, y)

r2 = metrics.r2\_score(prediction, y)

partere grope annivitation in escapia organi Carlinest of assertion data State 11/10/00/1000/Lag et QCMcSqArmCatte(MATSHAR)



#### Scikit Learn 模型評估的方法



from sklearn import metrics

- Continuous
  - Mean Absolute Error
  - Root Mean Squared Error
  - R2 Score
- Binary (Classification)
  - Accuracy
  - F1 score



### 迴歸型 (Regression) 模型常用評估指標

Mean Absolute Error:

 $MAE = \frac{\sum_{i=1}^{N} |e_i|}{N}$ 

- Mean Squared Error:

Mean Squared Error: 
$$MSE = \frac{\sum_{i=1}^{N} e_i^2}{N}$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^{N} e_i^2}{N}}$$

#### 



- R2 Score:
  - coefficient of determination

$$RMSE = \frac{N}{N}$$

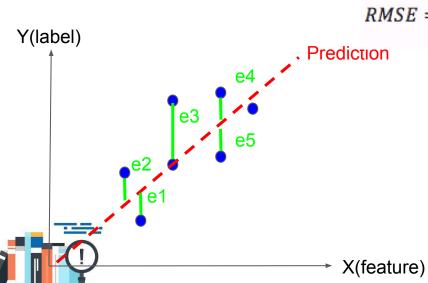
$$RMSE = \frac{\sum_{i=1}^{N} e_i^2}{N}$$

$$R2 = 1 - \frac{\sum_{i=1}^{N} (e_i)^2}{\sum_{i=1}^{N} (y - y_i)^2} = 1 - \frac{\text{total error}}{\text{data variance}}$$

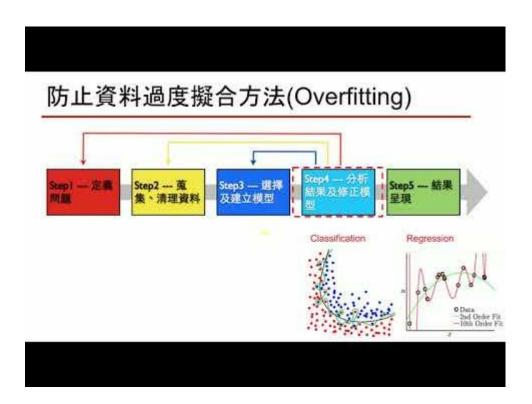
- 觀察在資料上加不同noise對預測結 果造成的影響
- 了解不同metrics的意義

from sklearn import metrics

mae = metrics.mean absolute error(y,prediction) mse = metrics.mean squared error(y,prediction) = metrics.r2 score(y,prediction)



### 防止資料過擬合





### Step 4 - 分析結果及修正模型:分類型模型評估

#### 9 Model Evaluation Example

多元分類的precision& recall & F1\_score

#### # Quantitative Measurement on the Performance

在分類問題上,所有的模型評估方法基本上都可以由confusion matrix得出來。在y方向也就是row方向代表的就是actual 0-9。x方向代表也就是column方向代表的是predicted 0-9。



### 分類型 (Classification) 常用評估指標 - Accuracy

Accuracy:

Accuracy = 
$$\frac{TN+TP}{ALL=(TN+TP+FN+FP)}$$

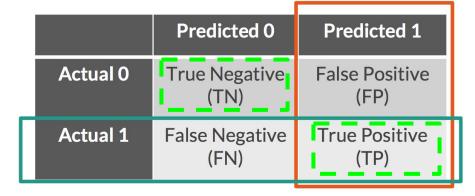


1: 壞人

0: 好人

全猜好人可以得 到高達90%的準 確率呢!

#### Confusion Matrix





樣本比例差異大















### 分類型 (Classification) 常用評估指標 - F1 Score

● F1 Score: 適用於inbalanced data ●

$$F_1 = 2 \cdot rac{1}{rac{1}{\operatorname{recall}} + rac{1}{\operatorname{precision}}} = 2 \cdot rac{\operatorname{precision} \cdot \operatorname{recall}}{\operatorname{precision} + \operatorname{recall}} \, .$$



1: 壞人

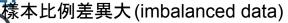
0: 好人

#### **Confusion Matrix**

	Predicted 0	Predicted 1
Actual 0	True Negative (TN)	False Positive (FP)
Actual 1	False Negative (FN)	True Positive (TP)

Recall = 
$$\frac{TP}{FN+TP}$$
 Precision =  $\frac{TP}{TP+FP}$ 

	Recall	Precision
全說壞人(1)	100%	10%
全說好人(0)	0%	
高代表的意義	大膽、信心度低	謹慎、信心度高

















#### 分類型評估練習 / 分類型評估方式 - Confusion Matrix

#### **9 Model Evaluation Example**



● 多元分類的precision& recall & F1\_score

#### # Quantitative Measurement on the Performance

在分類問題上,所有的模型評估方法基本上都可以由confusion matrix得出來。在y方向也就是row方向代表的就是actual 0~9,x方向代表也就是column方向代表的是predicted 0~9。

actual/predicted		0	1	2	3	4	5	6	7	8	9
0	[[:	37	0	0	0	0	0	0	0	0	0]
1	[	0	40	0	0	0	0	1	0	1	1]
2	[	0	0	42	2	0	0	0	0	0	0]
3	[	0	0	0	44	0	0	0	0	1	0]
4	[	0	0	0	0	37	0	0	1	0	0]
5	[	0	0	0	0	0	46	0	0	0	2]
6	[	0	1	0	0	0	0	51	0	0	0]
7	[	0	0	0	1	1	0	0	46	0	0]
8	[	0	3	1	0	0	0	0	0	44	0]
9	[	0	0	0	0	0	1	0	0	2	44]]

8是否分對的confusion matrix					
實際/預測	預測不為8	預測為8			
實際不為8					
實際為8					



### Section 1-3 [理論講授] 正則化(Regularization)

