

心理與神經資訊學

(Psychoinformatics & Neuroinformatics)

課號: Psy5261

識別碼: 227U9340

教室:彷彿在雲端

時間: —789





scikit-learn也有JS版

!!

AI & ML 總論

in-tel-li-gence

noun \in- 'te-lə-jən(t)s

- (1) the ability to **learn** or understand or to deal with new or trying situations
- (2) the ability to apply **knowledge** to manipulate one's environment or to think abstractly as measured by objective criteria

人工智慧的兩種典型

演繹法vs.歸納法



機器學習

是現在主流的人工智慧方法

人工智慧

機器學習

深度學習

機器學習(Machine Learning)

電腦可以幫我們從大量的資料中
找出微細的規律性
進而對資料做自動的分類與預測

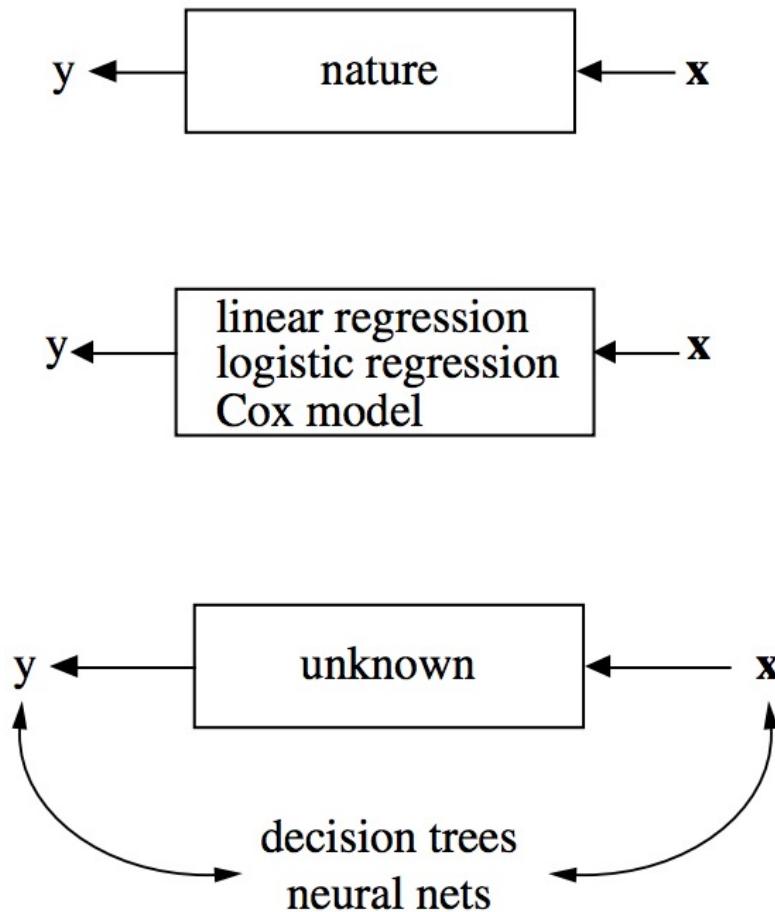


兩類統計模型建構：重解釋vs.重預測

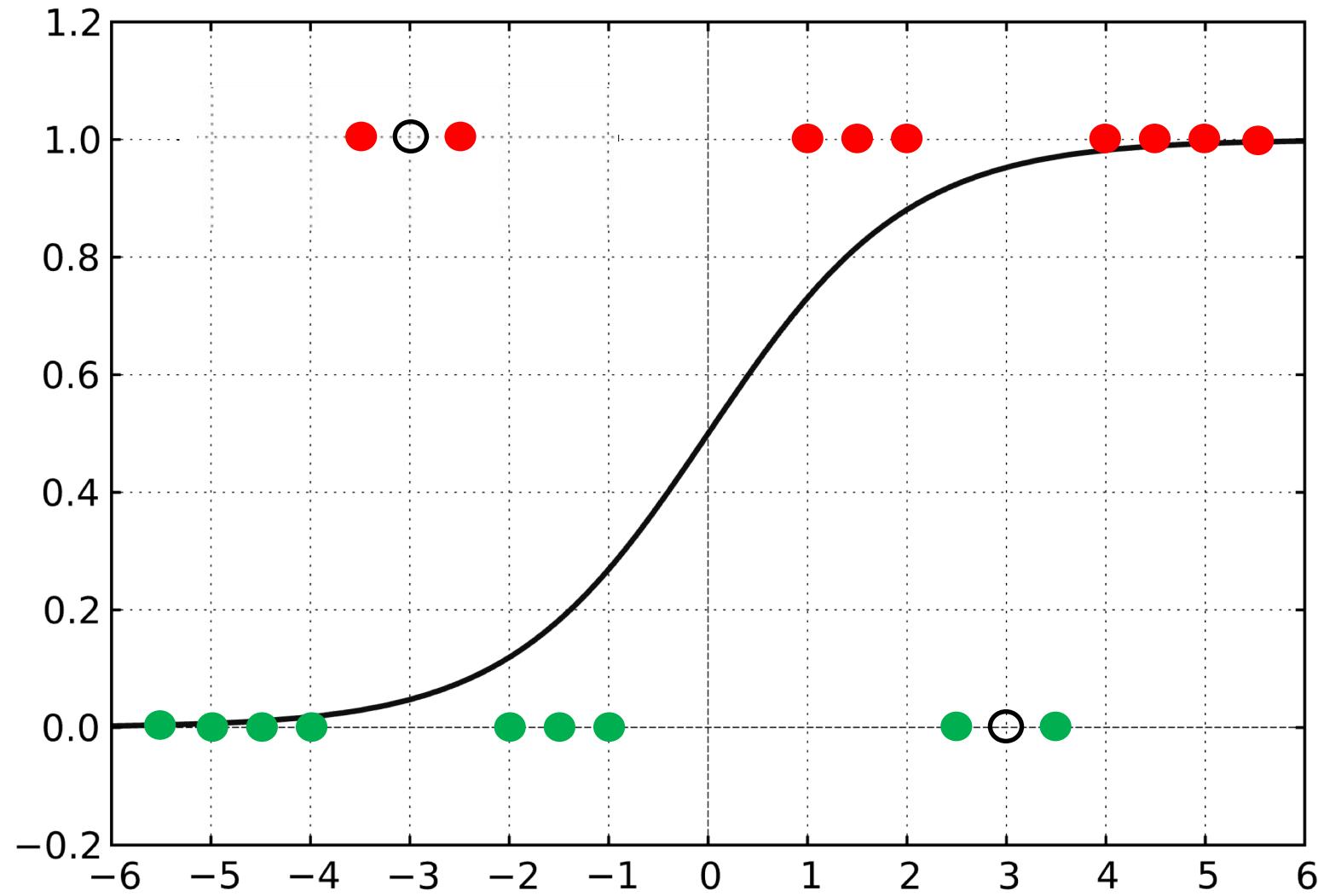
Statistical Science
2001, Vol. 16, No. 3, 199–231

Statistical Modeling: The Two Cultures

Leo Breiman



正確率：預測性模型 > 解釋性模型



邏輯迴歸與K-近鄰演算法的預測差異。邏輯迴歸對空心兩點的y值估計會根據連續的邏輯函數而預測 $y(x=-3)=0$ 與 $y(x=3)=1$ ；2-近鄰演算法則會利用 $x=-3$ 或 $x=3$ 旁左右兩個y值來投票預測 $y(x=-3)=1$ 與 $y(x=3)=0$ 。

心理學理論：解釋vs.預測

Choosing Prediction Over Explanation in Psychology: Lessons From Machine Learning

Perspectives on Psychological Science
2017, Vol. 12(6) 1100–1122

© The Author(s) 2017

Reprints and permissions:
sagepub.com/journalsPermissions.nav
DOI: 10.1177/1745691617693393
www.psychologicalscience.org/PPS



Tal Yarkoni and Jacob Westfall

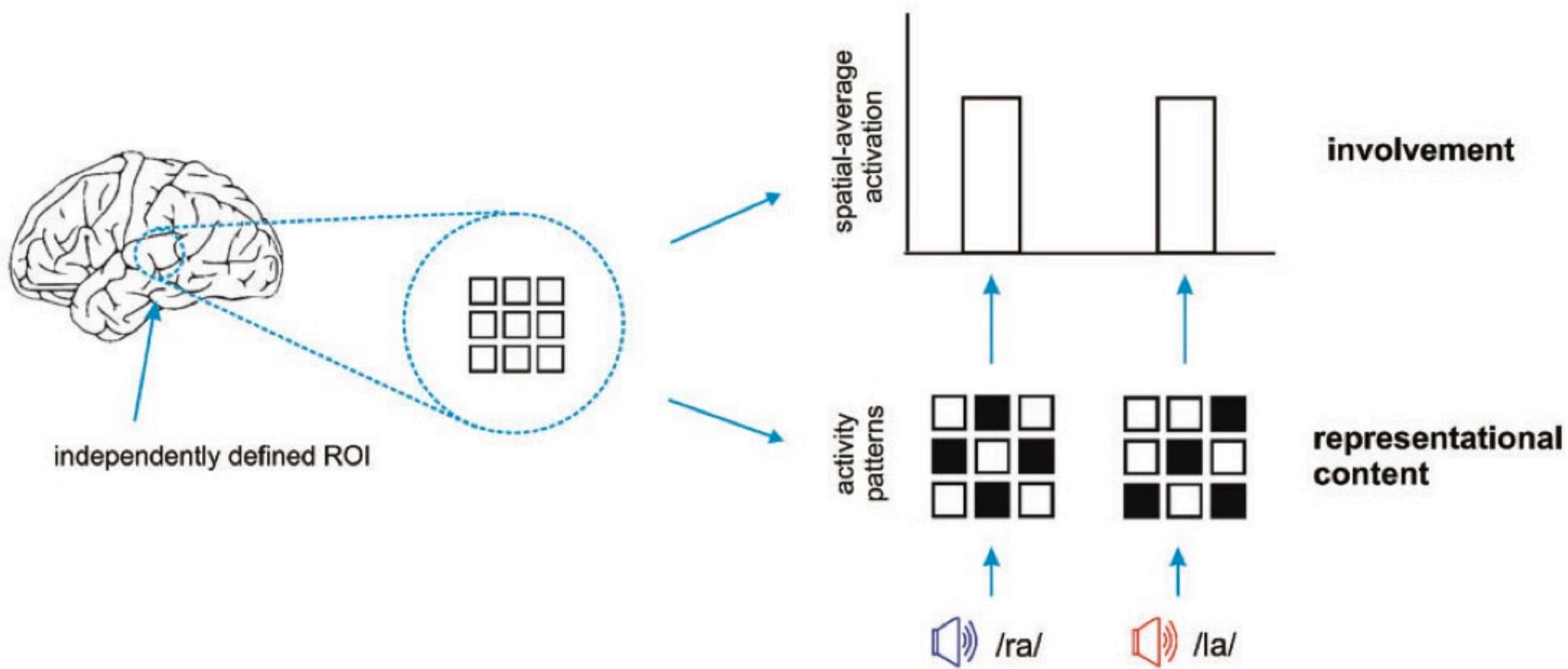
University of Texas at Austin

Abstract

Psychology has historically been concerned, first and foremost, with explaining the causal mechanisms that give rise to behavior. Randomized, tightly controlled experiments are enshrined as the gold standard of psychological research, and there are endless investigations of the various mediating and moderating variables that govern various behaviors. We argue that psychology's near-total focus on explaining the causes of behavior has led much of the field to be populated by research programs that provide intricate theories of psychological mechanism but that have little (or unknown) ability to predict future behaviors with any appreciable accuracy. We propose that principles and techniques from the field of machine learning can help psychology become a more predictive science. We review some of the fundamental concepts and tools of machine learning and point out examples where these concepts have been used to conduct interesting and important psychological research that focuses on predictive research questions. We suggest that an increased focus on prediction, rather than explanation, can ultimately lead us to greater understanding of behavior.

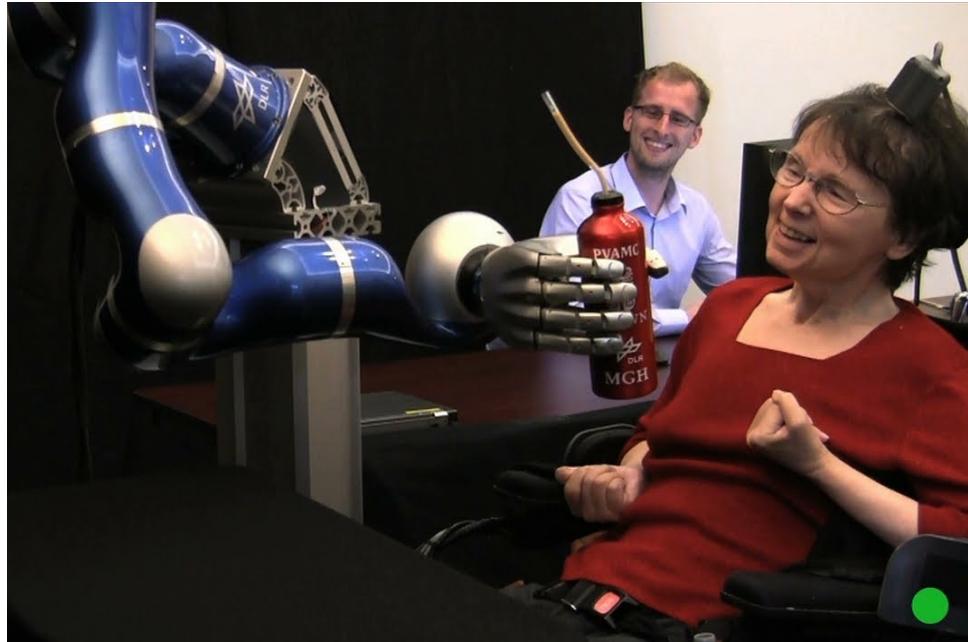
腦科學案例研究: Brain Decoding

心理學家和算命師有何不同?



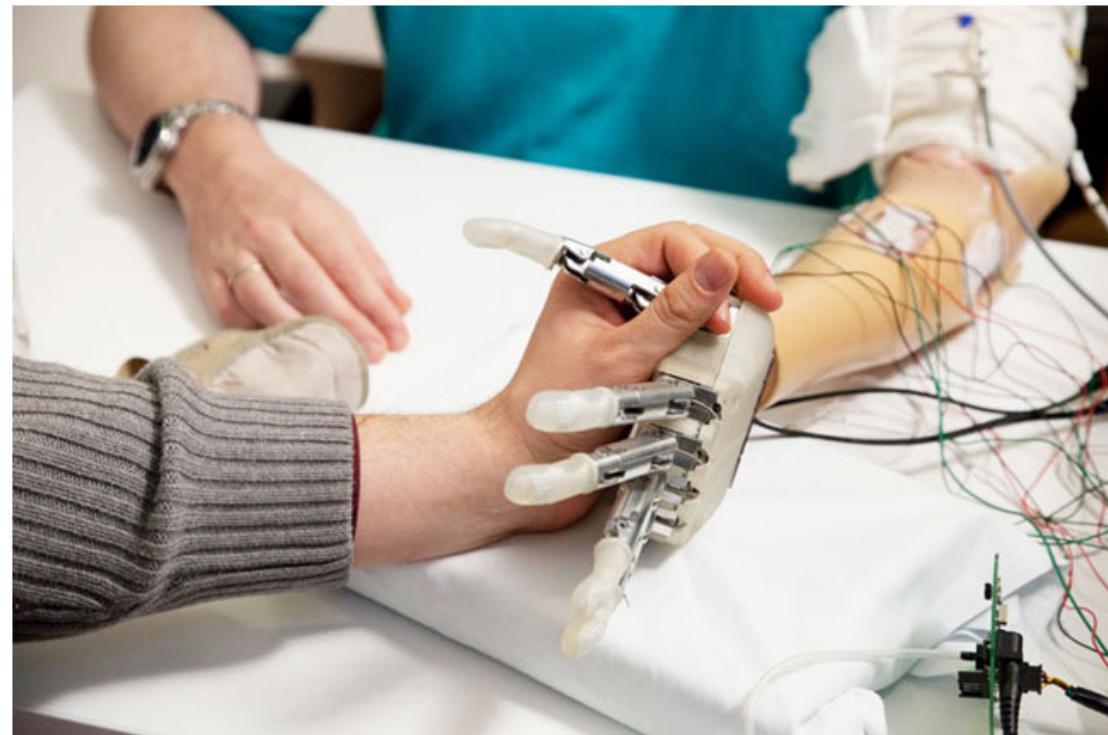
兩者都會讀心和解夢

神經義肢(Neuroprosthetics)



現在技術已可感受
物體的軟硬

可幫助身體障礙者
與世界互動



實務面 (scikit-learn)

市場需求

R&D 資料科學與工程研發	應徵職缺請 Email	hr@mobagel.com
職稱 (產業替代役可)	年薪 (含配股獎金)	需求條件
VP of Data Science	200萬~600萬	<p>年資 (參考二者年資)</p> <p>1. 5年以上資料分析專案或專業顧問經驗 2. 2年以上資料科學相關 PhD、PostDoc、教授等研究經驗</p> <p>必要條件 (至少其中三項具備) :</p> <ol style="list-style-type: none">熟悉 Python 及 Numpy, Pandas, Scikit-learn 套件熟悉 R 及統計理論熟悉資訊視覺化設計具備商業分析專案經驗具備工業設備分析專案經驗Kaggle 分析競賽排名至少前 10%英文語言能力佳 (多益860分以上) 或母語者日文語言能力佳 (日文 N1) 或母語者 <p>加分 :</p> <ol style="list-style-type: none">各項資訊專題競賽、黑客松、書卷獎資料分析相關學術論文發表
Chief Data Scientist	600萬~1200萬	<p>年資 (參考二者年資)</p> <p>1. 5年以上資料分析專案或專業顧問經驗 2. 2年以上資料科學相關 PhD、PostDoc、教授等研究經驗</p> <p>必要條件 (至少其中四項具備) :</p> <ol style="list-style-type: none">熟悉機器學習/深度學習演算法及統計方法

Scikit-Learn



Classification

Identifying to which set of categories a new observation belong to.

Applications: Spam detection, Image recognition.

Algorithms: *SVM, nearest neighbors, random forest, ...*

— Examples

Regression

Predicting a continuous value for a new example.

Applications: Drug response, Stock prices.

Algorithms: *SVR, ridge regression, Lasso, ...*

— Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: *k-Means, spectral clustering, mean-shift, ...*

— Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: *PCA, Isomap, non-negative matrix factorization.*

— Examples

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

Modules: *grid search, cross validation, metrics.*

— Examples

Preprocessing

Feature extraction and normalization.

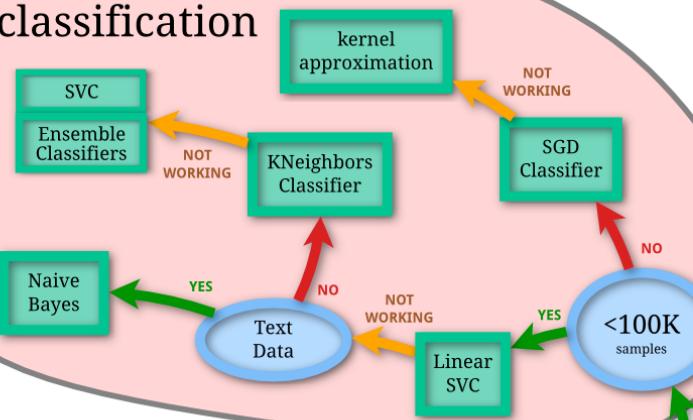
Application: Transforming input data such as text for use with machine learning algorithms.

Modules: *preprocessing, feature extraction.*

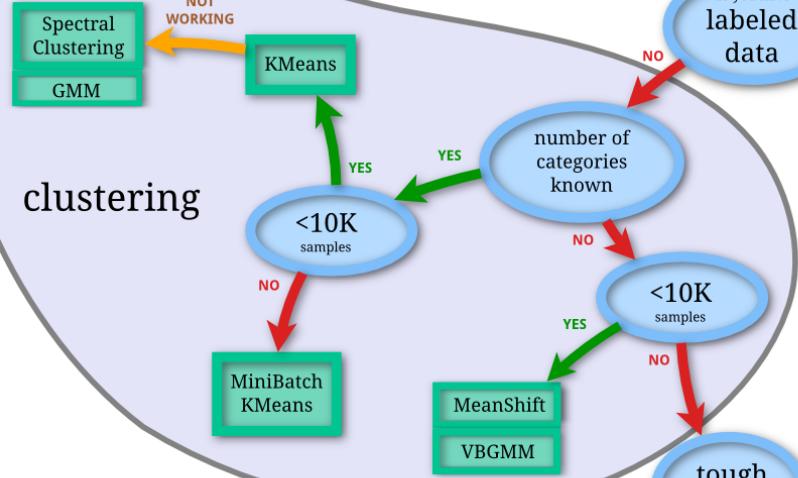
— Examples

scikit-learn推薦的workflow

classification



clustering



scikit-learn
algorithm cheat-sheet

predicting a category

predicting a quantity

just looking

predicting structure

regression

START

>50 samples

get more data

few features should be important

predicting structure

tough luck

dimensionality reduction

Randomized PCA

RidgeRegression

kernel approximation

Lasso

ElasticNet

SVR(kernel='rbf')

EnsembleRegressors

SVR(kernel='linear')

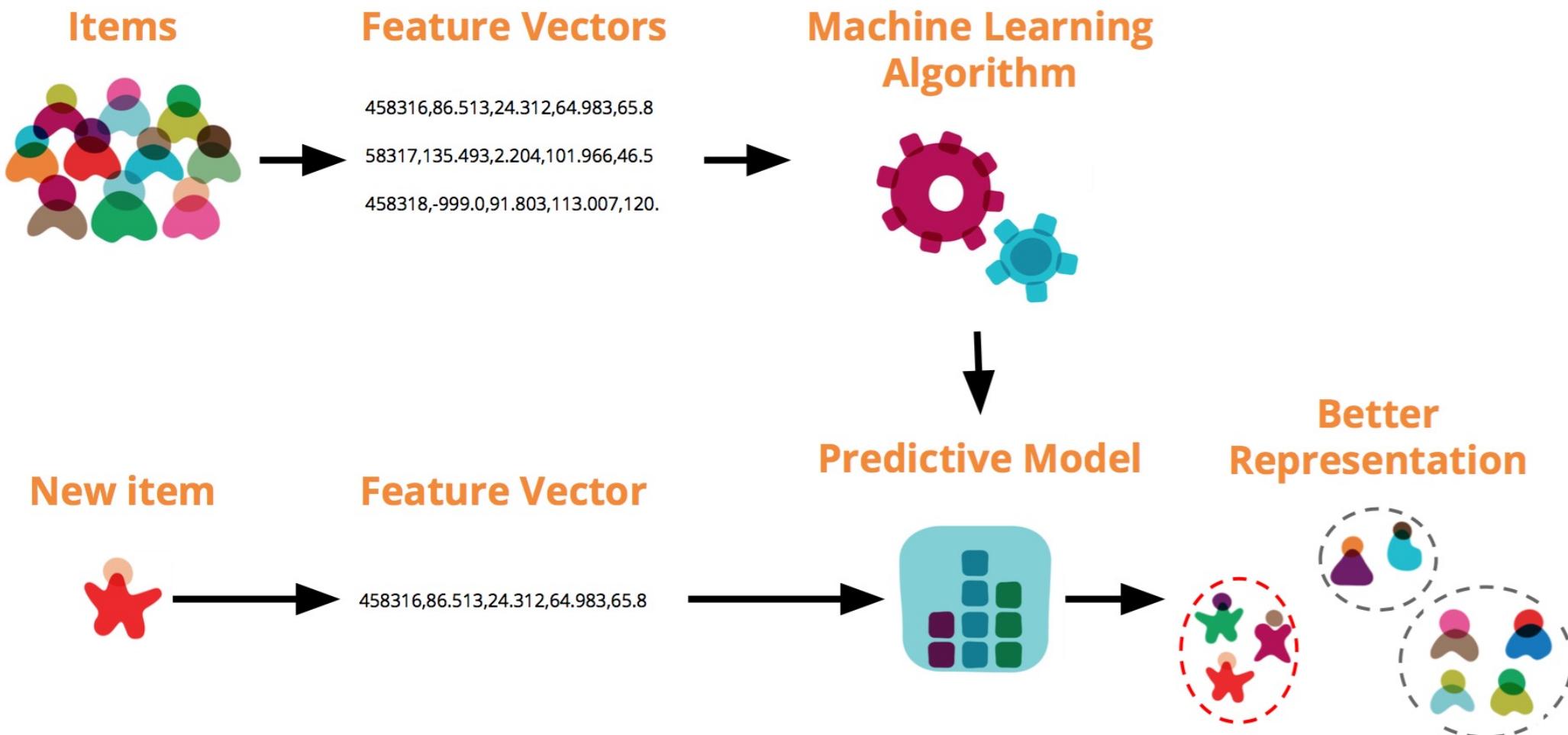
Isomap

Spectral Embedding

LLE

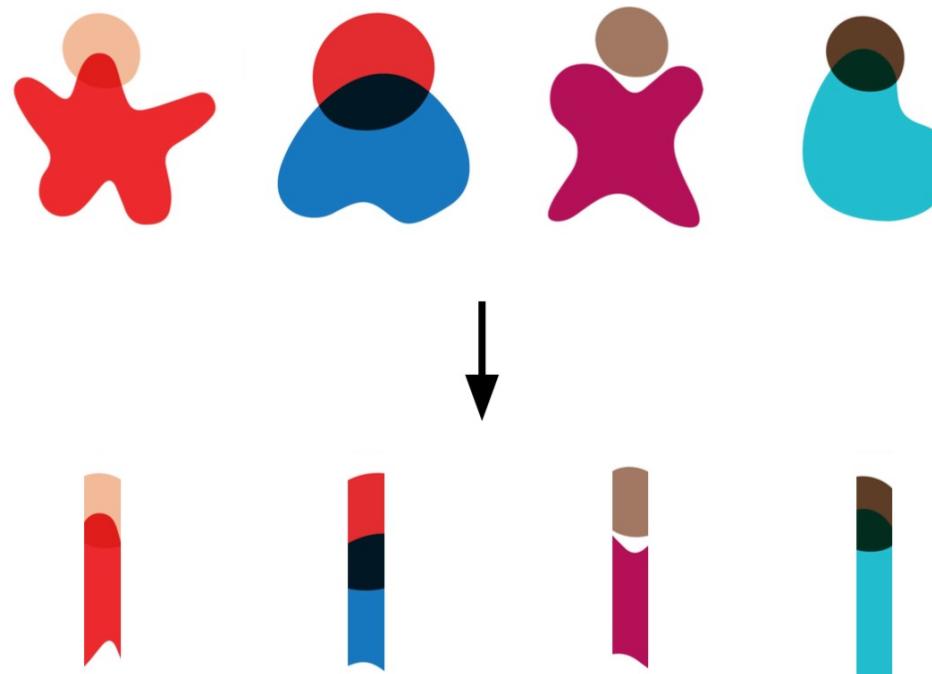
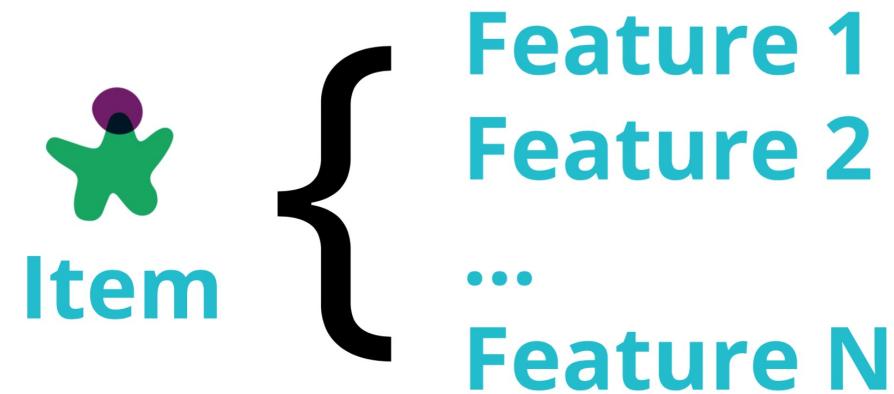
Features & Normalization

每個feature單位不同(如身高vs.體重)
因此最好能去單位化(normalization)再一起比較



Dimensionality Reduction

Feature space過大較不容易找出regularities
最好能想辦法先降低維度(如用PCA)



Unsupervised vs. Supervised Learning

非監督式的學習：沒有教師提供正確答案

監督式的學習：有教師提供正確答案

半監督式的學習：有教師提供部分正確答案



根據花瓣花萼大小的不同，鳶尾花似乎有3種

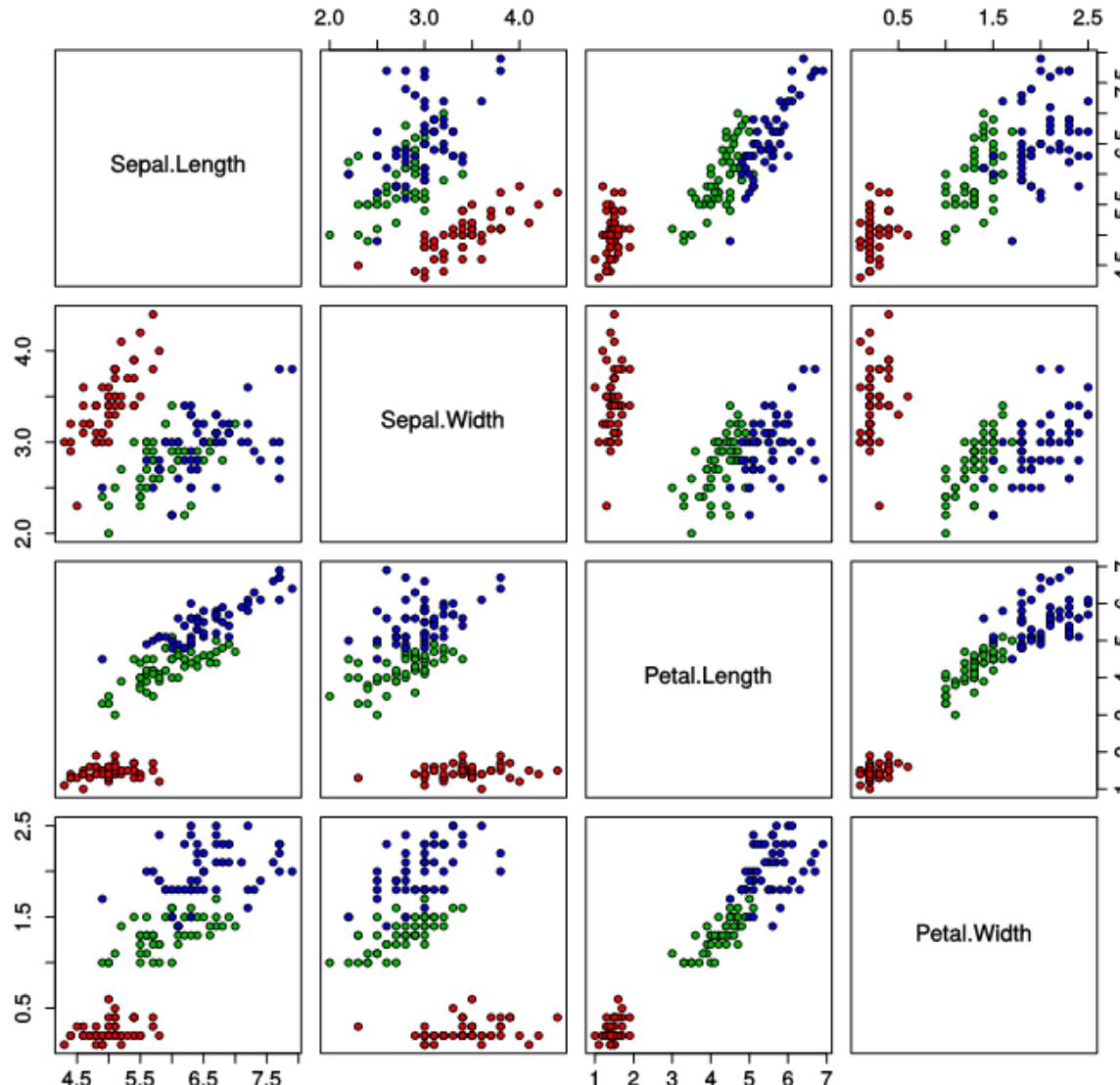
鳶尾花有3種，這個樣本是第k種

分類問題上主要是看有無提供類別標記以供學習

The Iris Dataset

常用來測試各種新的演算法

Iris Data (red=setosa,green=versicolor,blue=virginica)



非監督式學習 (Unsupervised Learning)

琳琅滿目的演算法

對演算法的原理有一些了解才知道其特性

Method name	Parameters	Scalability	Usecase	Geometry (metric used)
K-Means	number of clusters	Very large n_samples , medium n_clusters with MiniBatch code	General-purpose, even cluster size, flat geometry, not too many clusters	Distances between points
Affinity propagation	damping, sample preference	Not scalable with n_samples	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Mean-shift	bandwidth	Not scalable with n_samples	Many clusters, uneven cluster size, non-flat geometry	Distances between points
Spectral clustering	number of clusters	Medium n_samples , small n_clusters	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters	Large n_samples and n_clusters	Many clusters, possibly connectivity constraints	Distances between points
Agglomerative clustering	number of clusters, linkage type, distance	Large n_samples and n_clusters	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
DBSCAN	neighborhood size	Very large n_samples , medium n_clusters	Non-flat geometry, uneven cluster sizes	Distances between nearest points
Gaussian mixtures	many	Not scalable	Flat geometry, good for density estimation	Mahalanobis distances to centers
Birch	branching factor, threshold, optional global clusterer.	Large n_clusters and n_samples	Large dataset, outlier removal, data reduction.	Euclidean distance between points

No Free Lunch Theorem

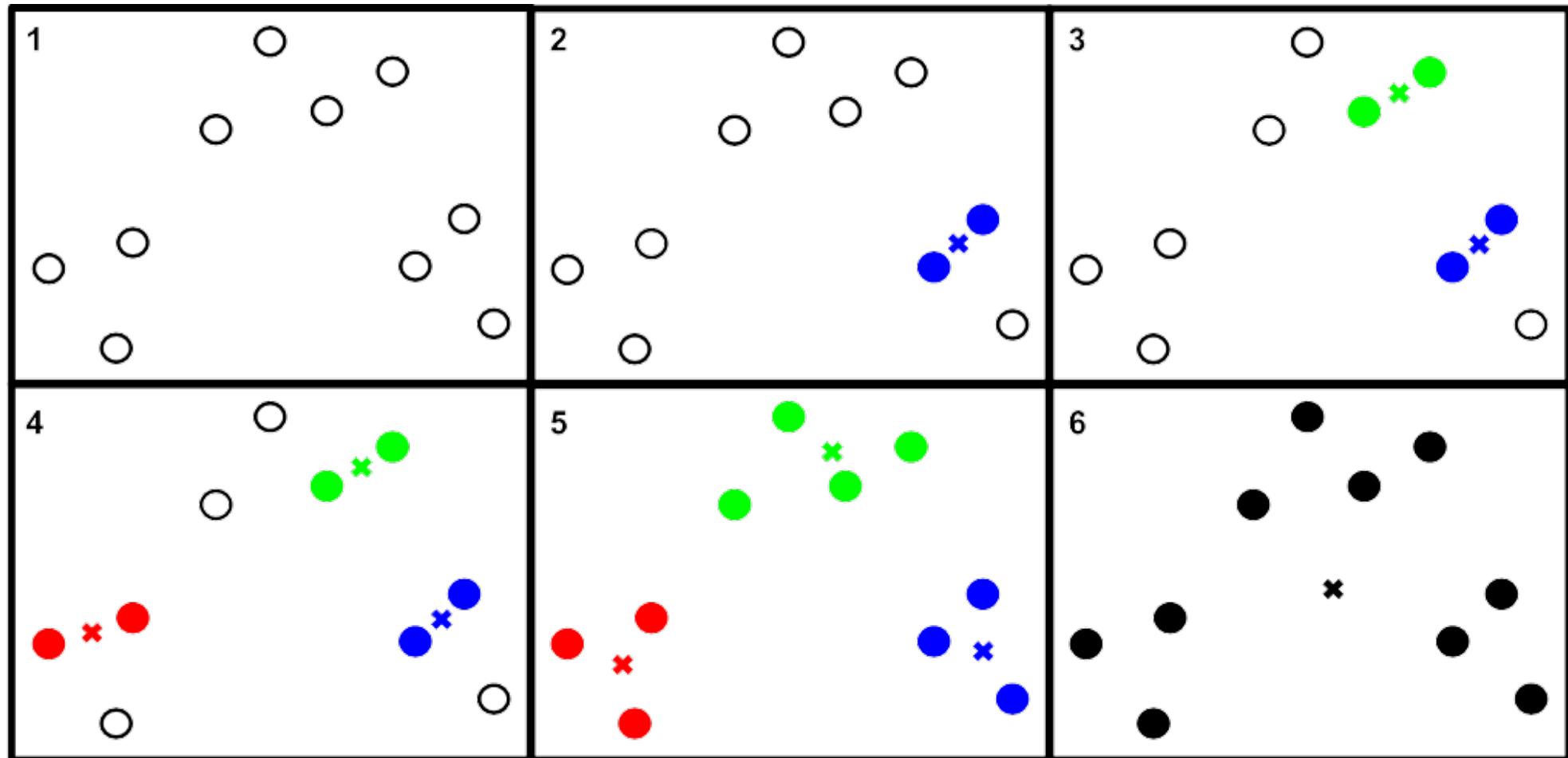
Any two optimization algorithms are equivalent when their performance is averaged across all possible problems.



每種機器學習演算法有其各自的長短處

Hierarchical Clustering

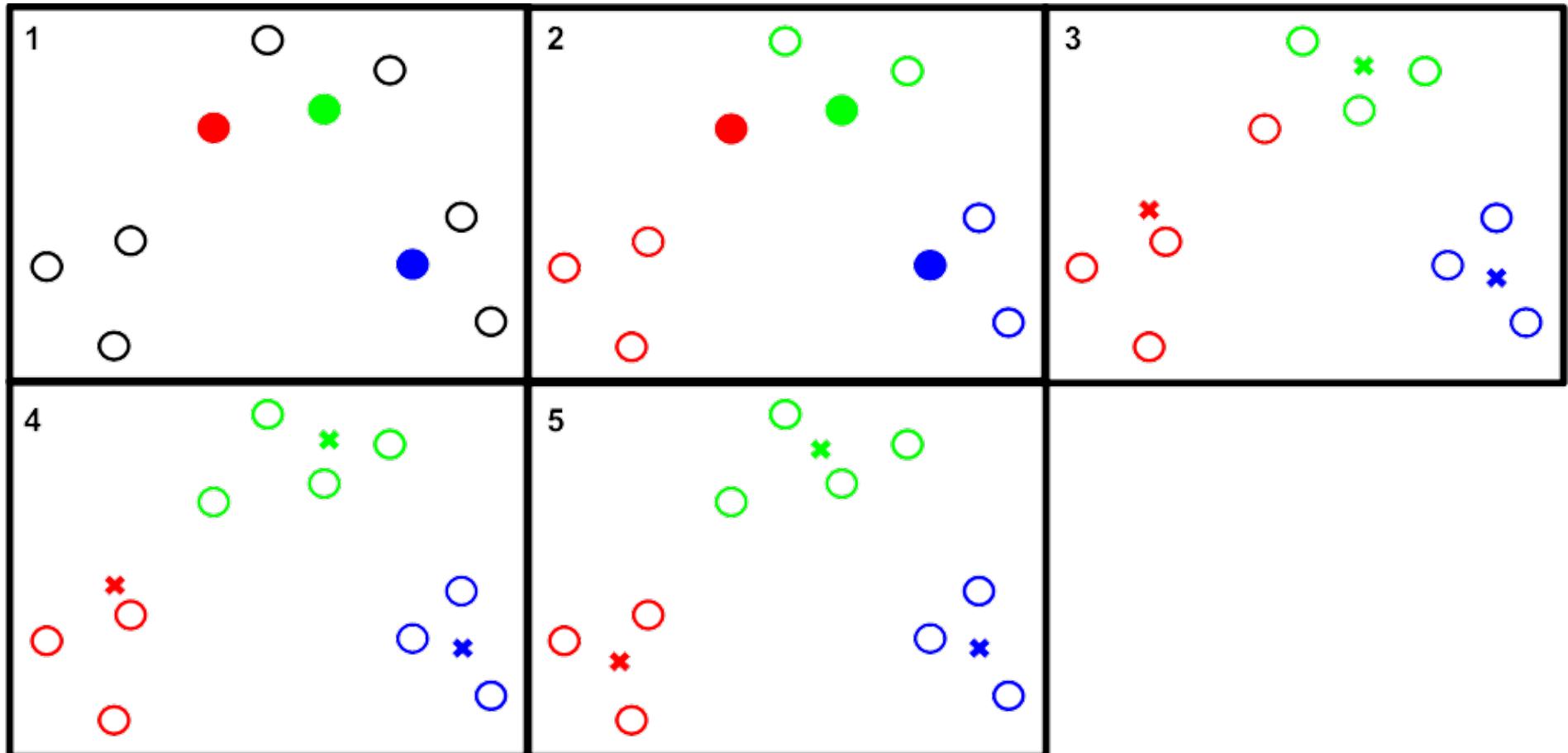
每一步驟最相近/相似的兩組合成為一群
但定義兩群組的距離的方法有很多種



```
model=AgglomerativeClustering(n_clusters=3)  
model.fit(X); print(model.labels_)
```

K-Means Clustering

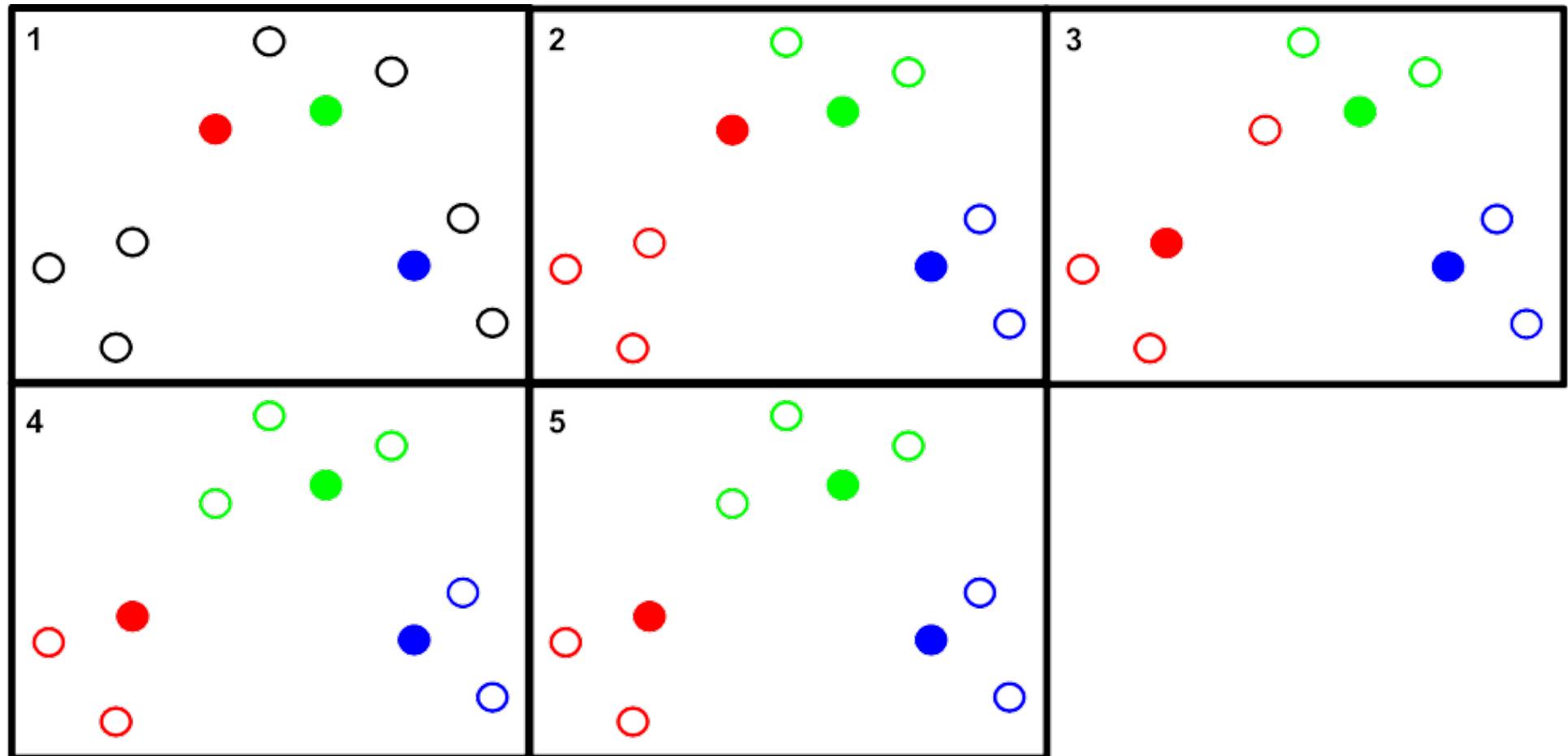
指派各點 x_j 到不同的組群 C_i ,使得
每個點到其群組中心位置 μ_i 為k種可能中最近



```
model=KMeans(n_clusters=3)  
model.fit(X); print(model.labels_)
```

K-Medoids Clustering

用某一樣本當作每群的原型以避免極值影響
此原型與其他同群內的樣本距離之和為最小

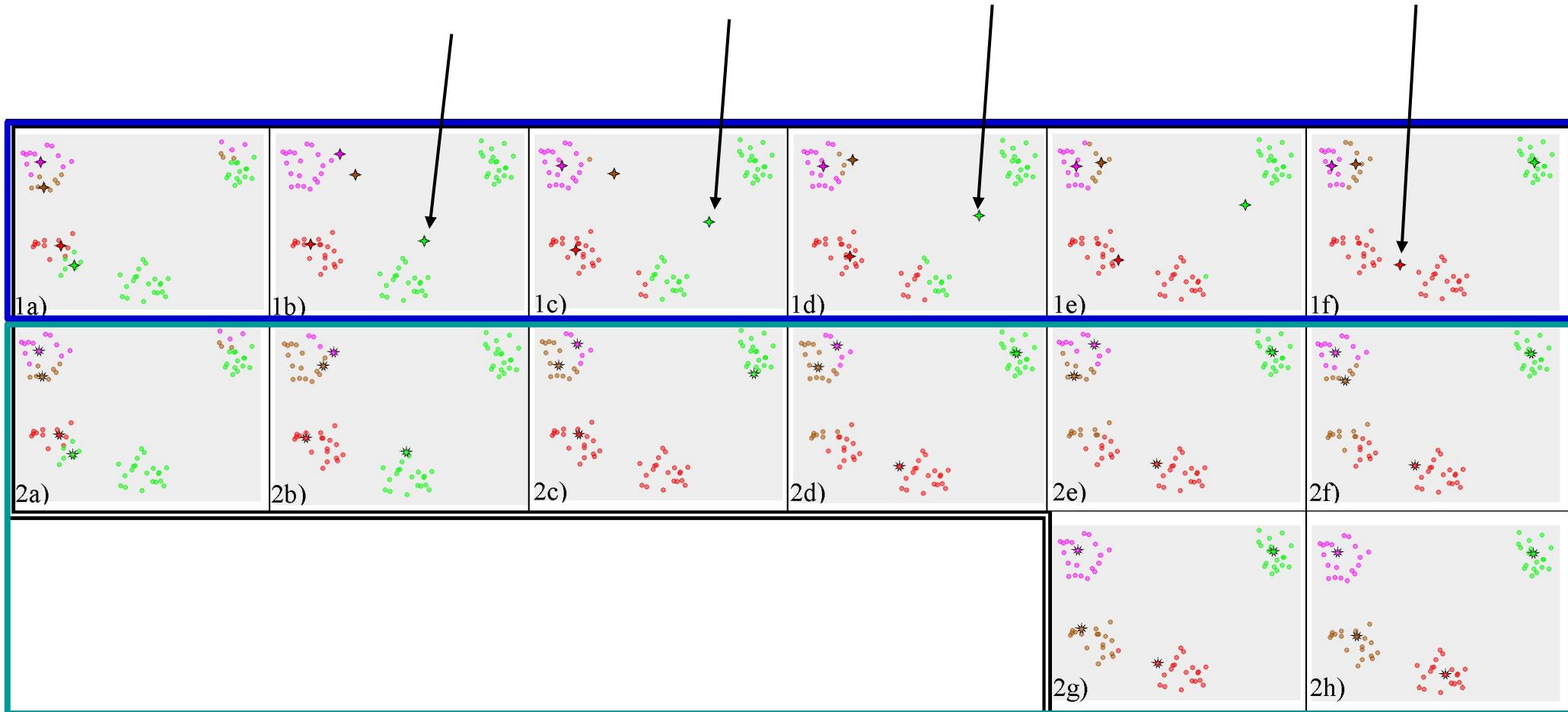


N/A in scikit-learn

Find implementations in Google Search/GitHub

K-Means vs. K-Medoids

K-Means容易產生四不像的群組原型(prototype)！

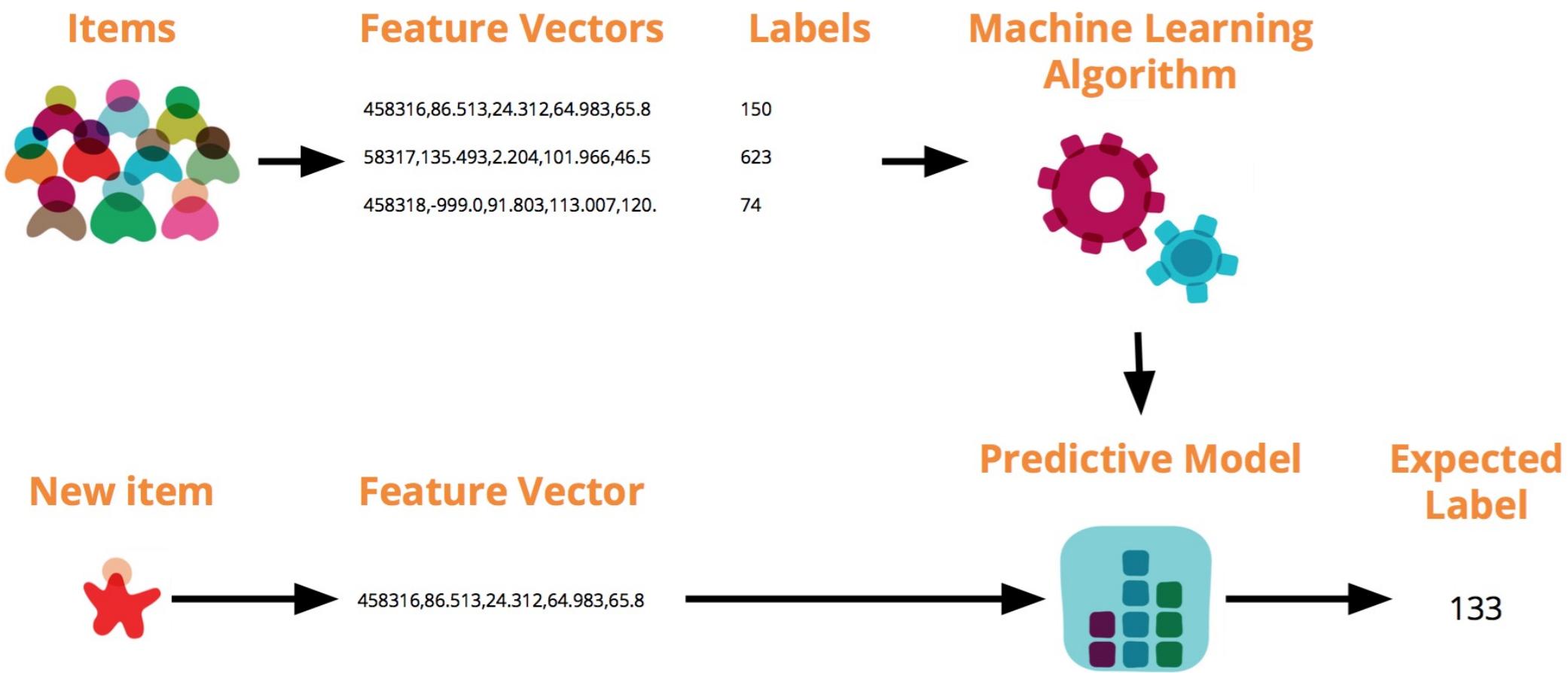


監督式學習

(Supervised Learning)

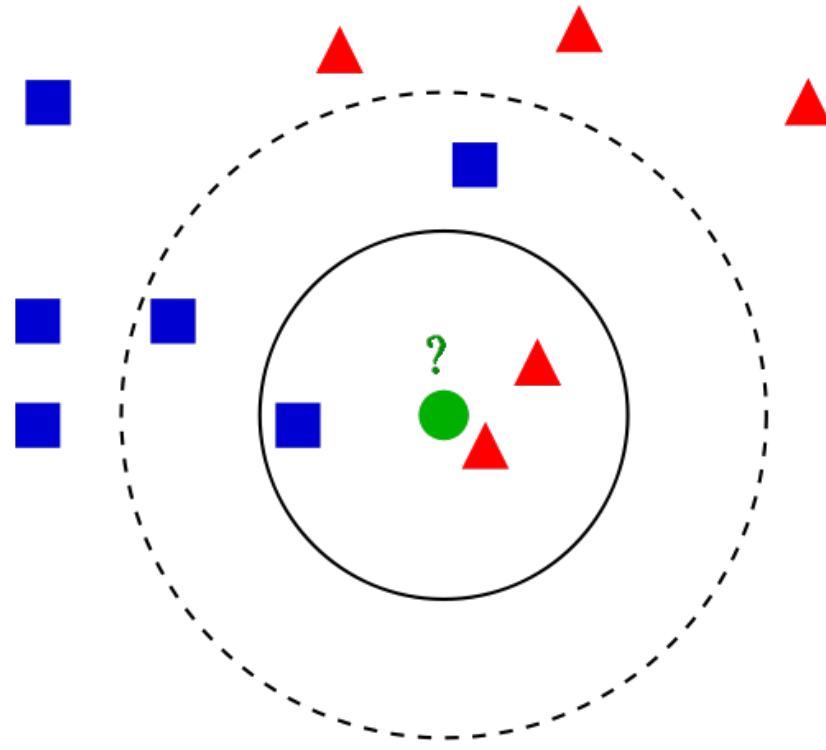
整體流程

比unsupervised learning多出了labels的部分



K-Nearest Neighbor (kNN)

用回憶中最相似的k個範例，以多數決判斷類別/數值



$k=3$ ● = ▲

$k=5$ ● = ■

$k=11$ ● = ?

```
from sklearn import *
clf=neighbors.KNeighborsClassifier(3) #try 1
clf.fit(X,Y) #training
print(np.mean(clf.predict(X)==Y)) #testing
```

Distance-Weighted KNN

$$\hat{f}(x_q) = \frac{\sum_{i=1}^k w_i f(x_i)}{\sum_{i=1}^k w_i} \quad \text{with } w_i = \frac{1}{d(x_q, x_i)^2}$$

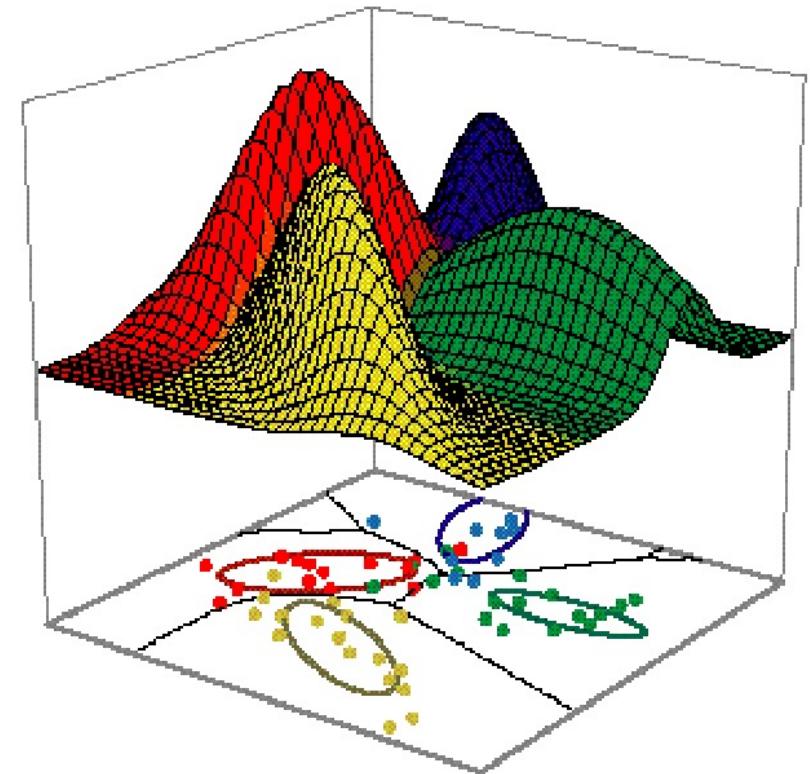

改進KNN使得投票比重和距離成反比

```
from sklearn import *
clf=neighbors.KNeighborsClassifier(3,'distance')
clf.fit(X,Y) #training
print(np.mean(clf.predict(X)==Y)) #testing
```

Naive Bayes Classifier

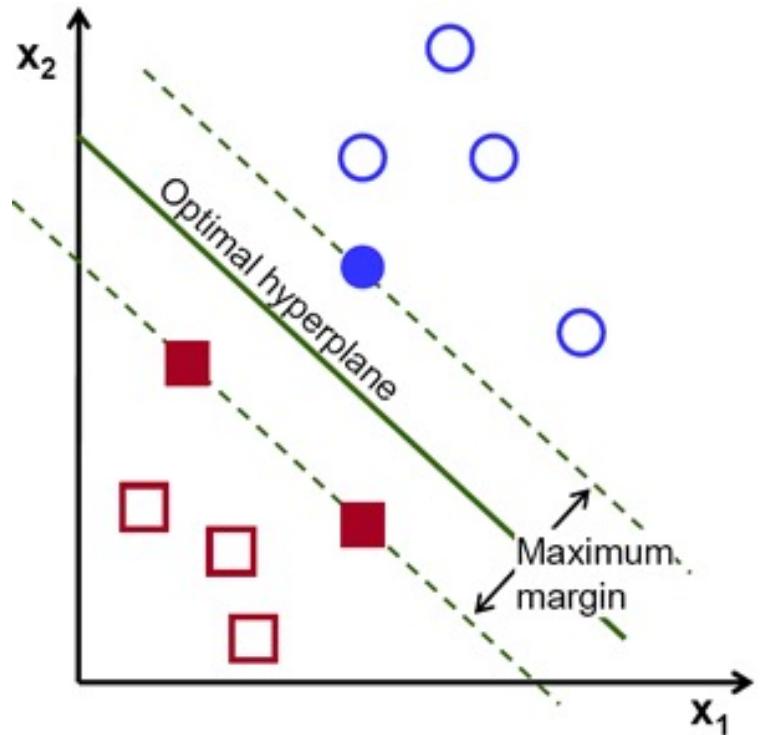
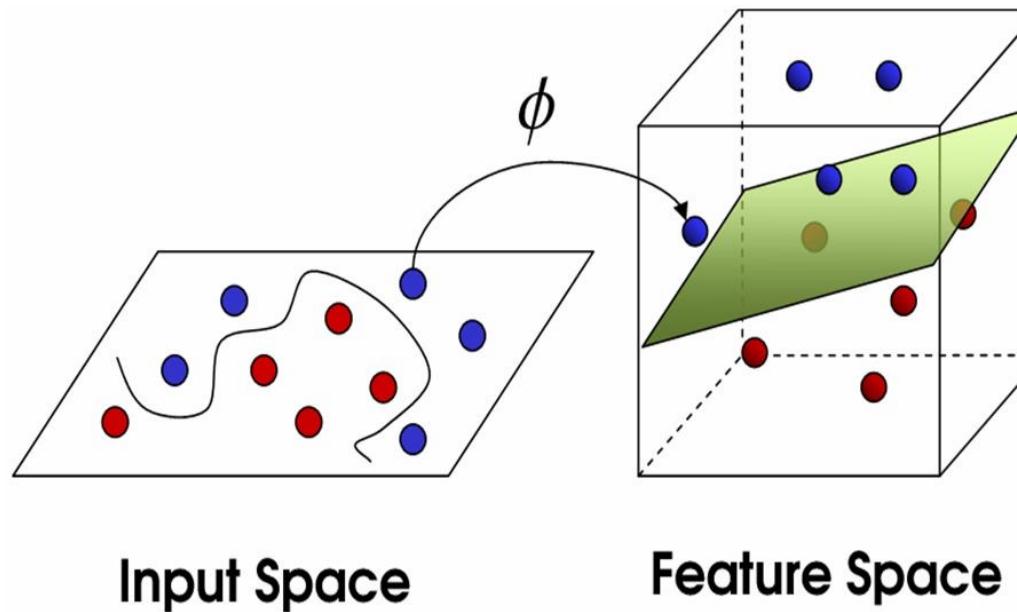
$$P(C_i|\vec{F}) = \frac{P(\vec{F}|C_i)P(C_i)}{P(\vec{F})}$$

用觀察到的範例估計母體分布，
判斷一個樣本最有可能的類別。



```
from sklearn import *
clf=naive_bayes.GaussianNB()
clf.fit(X,Y) #training
print(np.mean(clf.predict(X)==Y)) #testing
```

Support Vector Machine



```
from sklearn import *
clf=svm.SVC()
clf.fit(X,Y) #training
print(np.mean(clf.predict(X)==Y)) #testing
```

機器學習的各種亂做

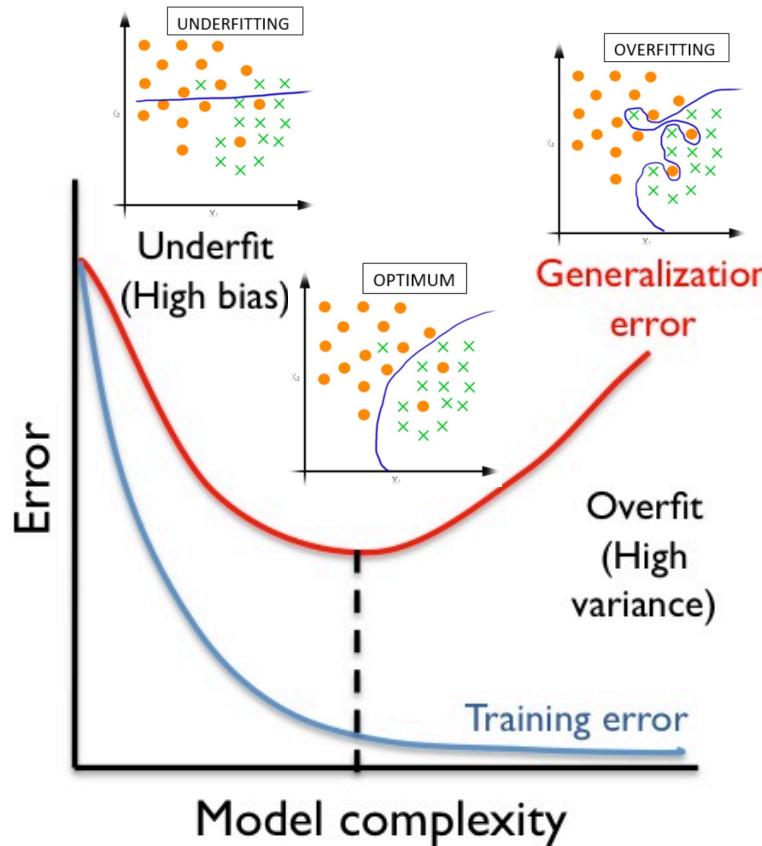
機器學習常見的錯誤(1/3)



```
X2=np.random.rand(150,4) # or randint for y  
Y=iris.target  
clf=neighbors.KNeighborsClassifier(1)  
#clf=svm.SVC()  
clf.fit(X2,Y);  
pred=clf.predict(X2)  
print(np.mean(pred==Y))  
print(metrics.confusion_matrix(Y,pred))
```

訓練(training)適可而止

training vs. testing errors
又叫In- vs. out-sample errors

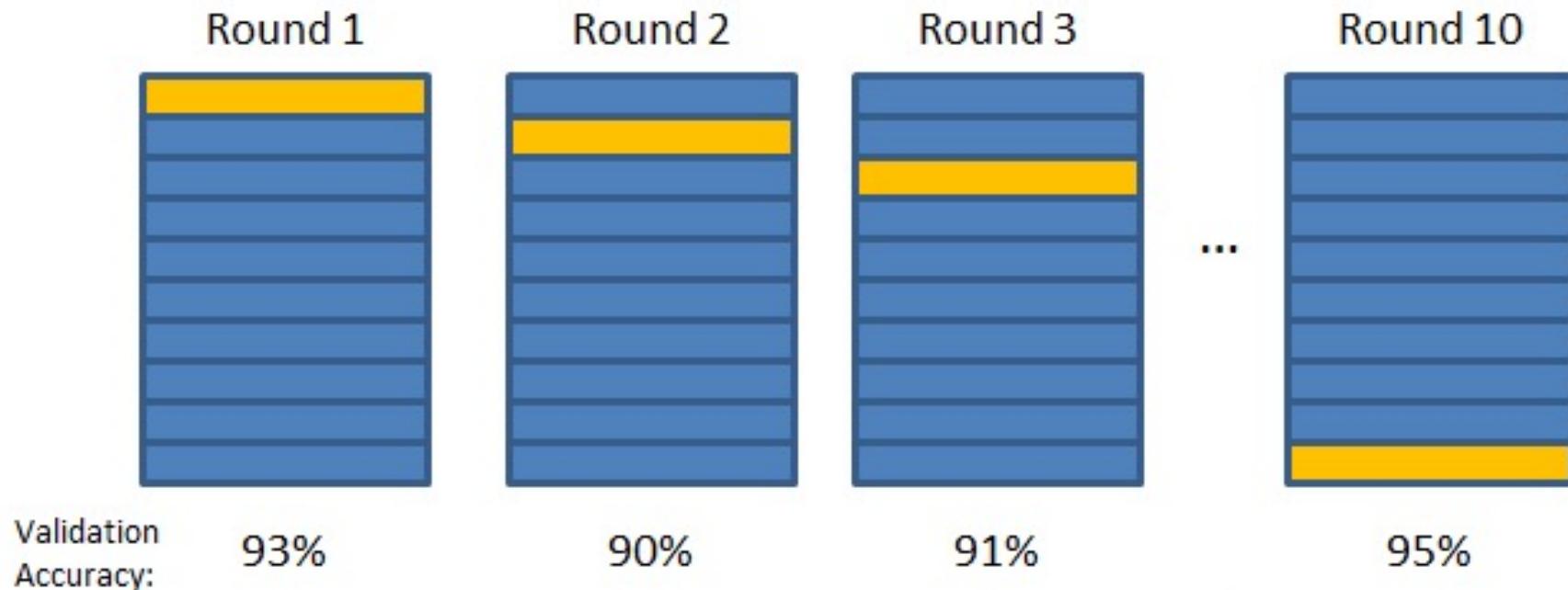


監督式學習的目標在minimize testing errors

監督式學習需要Cross-Validation

10-fold validation的例子：

- Validation Set
- Training Set



$$\text{Final Accuracy} = \text{Average}(\text{Round 1}, \text{Round 2}, \dots)$$

當 $k=N$ 時稱為leave-one-out cross-validation

機器學習常見的錯誤(2/3)



BUT

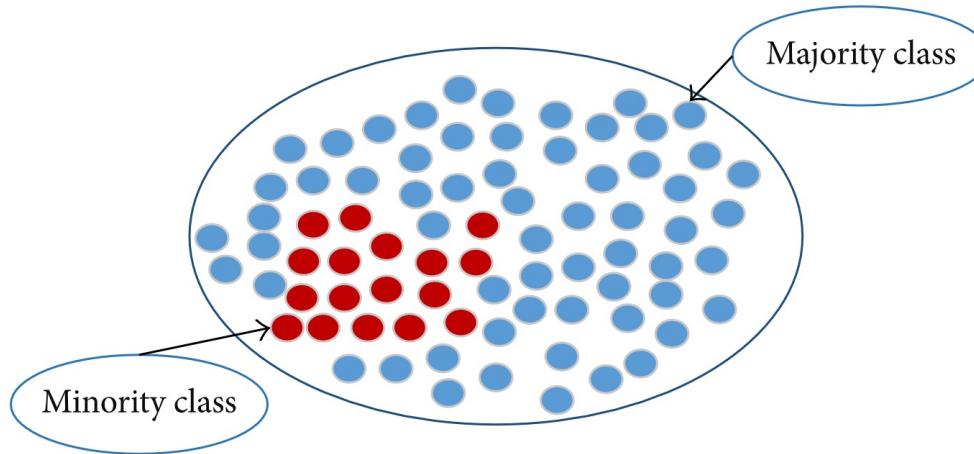
ANYTHING
THAT
CAN GO WRONG
WILL GO
WRONG



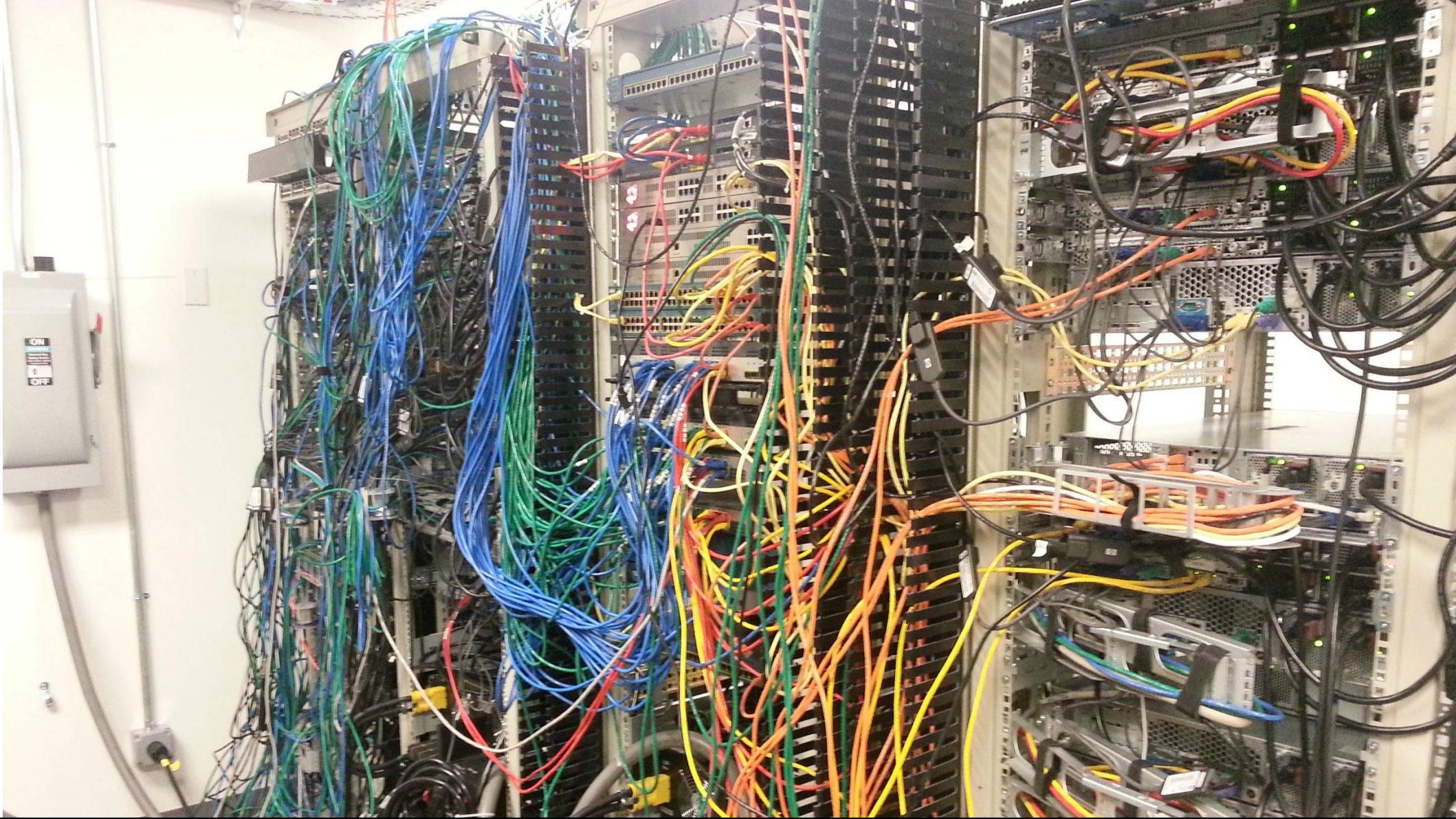
```
from sklearn import model_selection
clf=svm.SVC() # try other supervised classifiers
kf=model_selection.KFold(5)
s1=model_selection.cross_val_score(clf,X,Y,cv=kf)
s2=model_selection.cross_val_score(clf,X2,Y,cv=kf)
print(s1.mean(),s2.mean())
```

機器學習常見的錯誤(3/3)

Supervised learners常會學到prior distributions



```
from sklearn.model_selection import *
from sklearn.metrics import *
x=random.rand(100,3) # 3-d random features
y=random.permutation([0]*90+[1]*10) # 2 categories
clf=svm.SVC(); cv=KFold(100)
yp=cross_val_predict(clf,x,y,cv=cv)# leave-1-out
print('Accuracy:',mean(y==yp)) # mean accuracy
print('C. Matrix:\n',confusion_matrix(y,yp)) # c. matrix
```



聽本魯師一個勸：
每個ML pipeline都要打亂測過

Game Over

