



UNIVERSITY OF  
PORTSMOUTH

# Python for Data Analysis

## Modeling in Python – Bayesian and SVM (TB2 - Week 3)

Atefeh Khazaei

[atefeh.khazaei@port.ac.uk](mailto:atefeh.khazaei@port.ac.uk)



# What we will learn this week?

- ❑ Previous weeks modelling algorithms
  - ❑ K-Nearest Neighbours
  - ❑ Decision Tree
  - ❑ Random Forest
- ❑ This week modelling algorithms
  - ❑ Bayesian methods (Naïve Bayes)
  - ❑ Support Vector Machines (SVM)

# Naïve Bayes

- ❑ **A probabilistic method**

- ❑ Predicts class membership probabilities

- ❑ **Naïve Bayes classifier**

- ❑ Foundation: Based on Bayes' theorem
  - ❑ Performance: Acceptable accuracy and speed when applied to large databases

# Naïve Bayes (cont.)

## □ Bayes' Theorem:

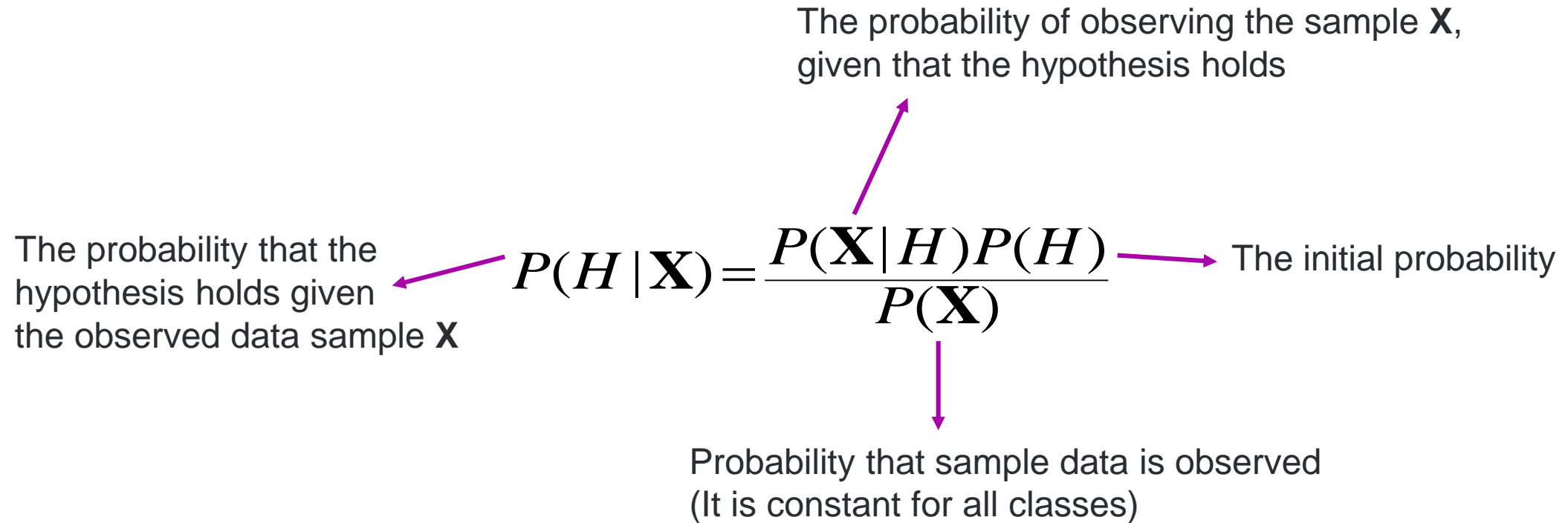
The probability of observing the sample  $\mathbf{X}$ ,  
given that the hypothesis holds

The probability that the hypothesis holds given  
the observed data sample  $\mathbf{X}$

$$P(H | \mathbf{X}) = \frac{P(\mathbf{X} | H)P(H)}{P(\mathbf{X})}$$

The initial probability

Probability that sample data is observed  
(It is constant for all classes)

A diagram illustrating Bayes' Theorem. The central equation is  $P(H | \mathbf{X}) = \frac{P(\mathbf{X} | H)P(H)}{P(\mathbf{X})}$ . Four arrows point from descriptive text to parts of the equation: an arrow from the top text points to  $P(\mathbf{X} | H)$ ; an arrow from the left text points to  $P(H | \mathbf{X})$ ; an arrow from the right text points to  $P(H)$ ; and an arrow from the bottom text points to  $P(\mathbf{X})$ .

# Naïve Bayes (cont.)

## Simple Example: One Feature, One Target

	Basketball	Cereal
1	Yes	No
2	Yes	Yes
3	Yes	No
4	No	Yes
5	No	Yes
6	No	No
...	...	...
5000	Yes	Yes

C \ B			
	YES	NO	
YES	2000	1750	3750
NO	1000	250	1250
	3000	2000	5000

# Naïve Bayes (cont.)

## Simple Example: One Feature, One Target

C \ B	YES	NO	
YES	2000	1750	3750
NO	1000	250	1250
	3000	2000	5000

Does someone who plays BASKETBALL (X, sample) eats CEREAL (H, Target)?

$$P(C_i | \mathbf{X}) = \frac{P(\mathbf{X} | C_i) P(C_i)}{P(\mathbf{X})}$$

$$P(\text{Cereal} | \text{Basketbal}) = \frac{P(\text{Basketbal} | \text{Cereal}) P(\text{Cereal})}{P(\text{Basketbal})} = \frac{\frac{2000}{3750} \frac{3750}{5000}}{\alpha} = \frac{0.4}{\alpha}$$

$$P(\overline{\text{Cereal}} | \text{Basketbal}) = 1 - P(\text{Cereal} | \text{Basketbal}) = \frac{P(\text{Basketbal} | \overline{\text{Cereal}}) P(\overline{\text{Cereal}})}{P(\text{Basketbal})} = \frac{\frac{1000}{1250} \frac{1250}{5000}}{\alpha} = \frac{0.2}{\alpha}$$

# sklearn.naive\_bayes.GaussianNB()

- ❑ You can find more details in the following link:
  - ❑ [https://scikit-learn.org/stable/modules/naive\\_bayes.html](https://scikit-learn.org/stable/modules/naive_bayes.html)
  - ❑ [https://scikit-learn.org/stable/modules/generated/sklearn.naive\\_bayes.GaussianNB.html](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.GaussianNB.html)
- ❑ About different types of NB:
  - ❑ GaussianNB,
  - ❑ MultinomialNB,
  - ❑ ComplementNB,
  - ❑ BernoulliNB,
  - ❑ CategoricalNB

# sklearn.naive\_bayes.GaussianNB()

```
from sklearn.naive_bayes import GaussianNB
from sklearn import metrics
from sklearn.model_selection import train_test_split

# Training data features, skip the first column 'Survived'
train_features = train_data[:, 1:]

# 'Survived' column values
train_target = train_data[:, 0]

# Split 80-20 train vs test data
train_x, test_x, train_y, test_y = train_test_split(train_features,
                                                    train_target,
                                                    test_size=0.20,
                                                    random_state=0)

clf = GaussianNB()
clf = clf.fit(train_x, train_y)
predict_y = clf.predict(test_x)

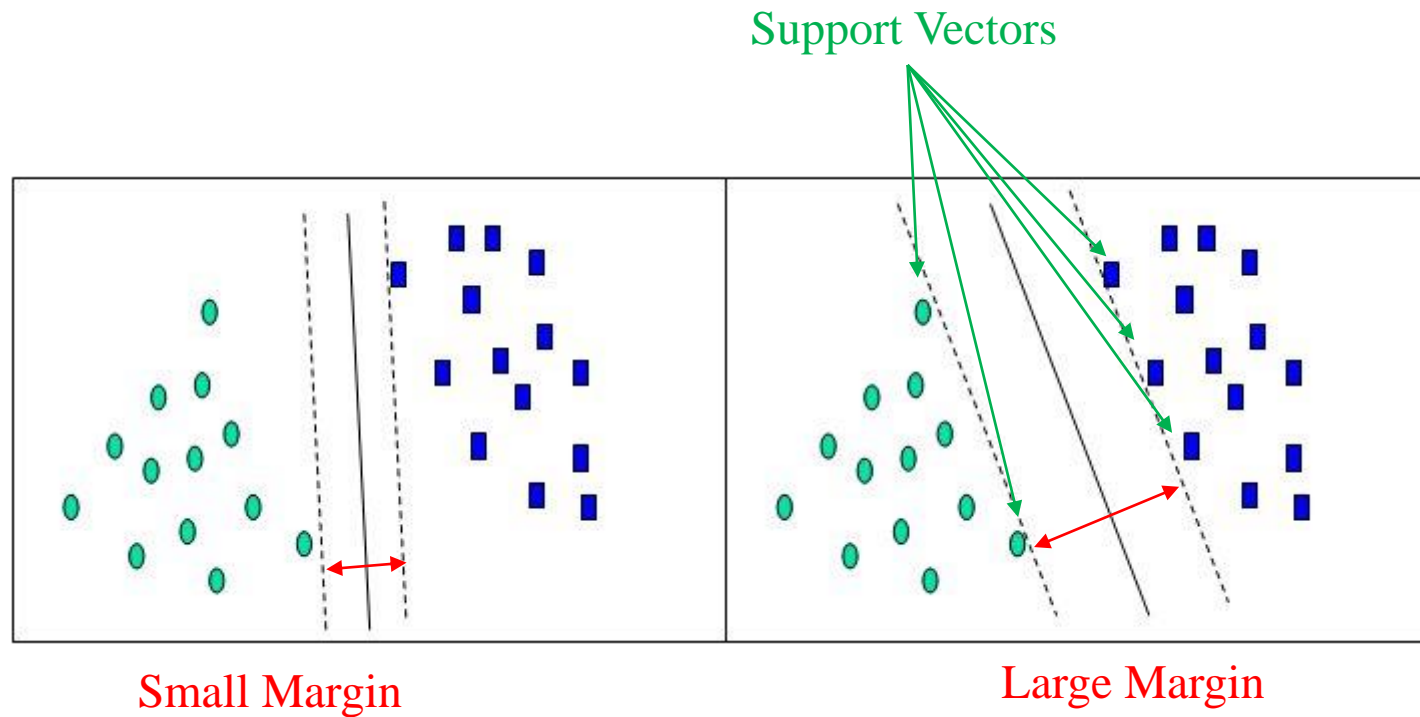
from sklearn.metrics import accuracy_score
print ("Accuracy = %.2f" % (accuracy_score(test_y, predict_y)))
```

Accuracy = 0.81



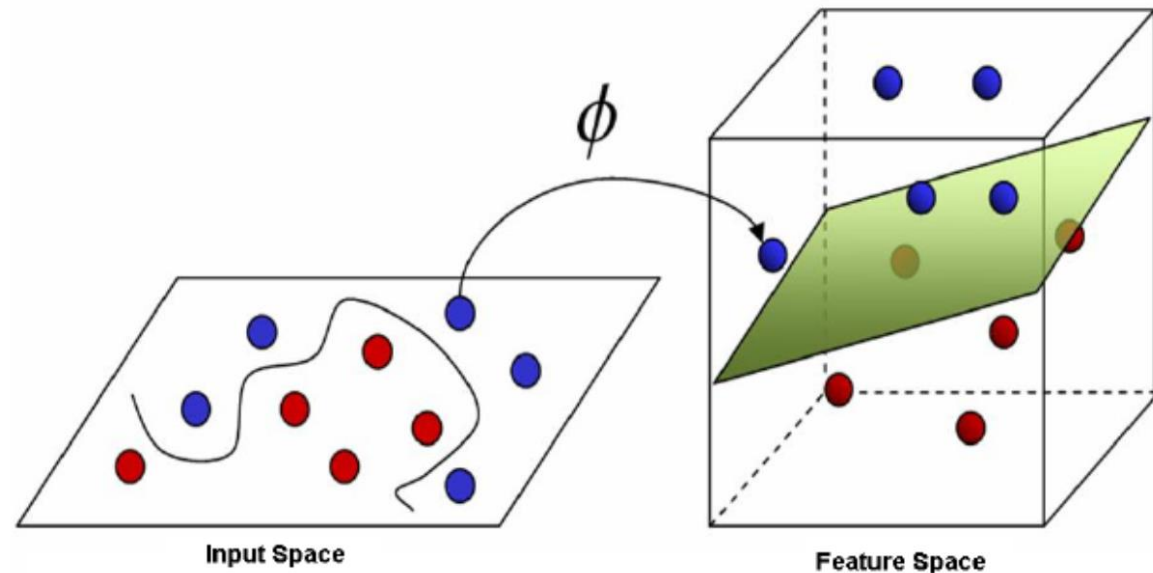
# Support Vector Machines (SVM)

- Support Vector Machine constructs a hyperplane in a high-dimensional space, which can be used for classification, regression, or other tasks like outliers detection.



# Support Vector Machines (SVM) (Cont.)

- With an appropriate nonlinear mapping to a sufficiently high dimension (**kernel** trick), data from two classes can always be separated by a hyper plane.
- With the new dimension, SVM searches for the optimal linear decision boundary.



# sklearn.svm.SVC()

- ❑ You can find more details in the following link:
  - ❑ <https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>
- ❑ About different parameters of SVM:
  - ❑ kernel (linear, poly, rbf, sigmoid, precomputed)
  - ❑ degree
  - ❑ gamma
  - ❑ ...

# sklearn.svm.SVC()

```
from sklearn.svm import SVC
from sklearn import metrics
from sklearn.model_selection import train_test_split

# Training data features, skip the first column 'Survived'
train_features = train_data[:, 1:]

# 'Survived' column values
train_target = train_data[:, 0]

# Split 80-20 train vs test data
train_x, test_x, train_y, test_y = train_test_split(train_features,
                                                    train_target,
                                                    test_size=0.20,
                                                    random_state=0)

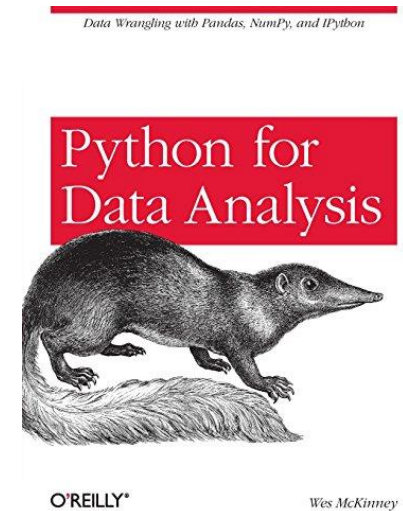
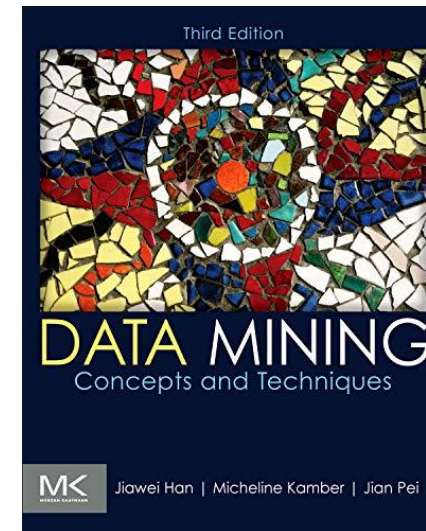
clf = SVC(kernel='linear')
clf = clf.fit(train_x, train_y)
predict_y = clf.predict(test_x)

from sklearn.metrics import accuracy_score
print ("Accuracy = %.2f" % (accuracy_score(test_y, predict_y)))
```

# References & More Resources

## References:

- McKinney, Wes. *Python for data analysis: Data wrangling with Pandas, NumPy, and Ipython*, O'Reilly Media, Inc., 2012.
- Han, Jiawei, Jian Pei, and Micheline Kamber. *Data mining: concepts and techniques*. Elsevier, 2011.



# Practical Session

- ❑ Revise the Titanic Case Study (Last session of TB1) and build some **Naïve Bayesian** and **SVM** models for Titanic. Try **different parameters** for these models and **compare them together**.