

# Learning from the Past: with Scikit-Learn

ANOOP THOMAS MATHEW  
Profoundis Labs Pvt. Ltd.



# Agenda

- Basics of Machine Learning
- Introduction some common techniques
- Let you know scikit-learn exists
- Some inspiration on using machine learning in daily life scenarios and live projects.

# How to draw a snake?



and ...

# How to draw a snake?



This is WEIRD!

and ...



# Introduction

A lot of Data!

What to do???

# Introduction

What is

Machine Learning  
(Data Mining)?  
*(in plain english)*

# Machine Learning

"A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ "

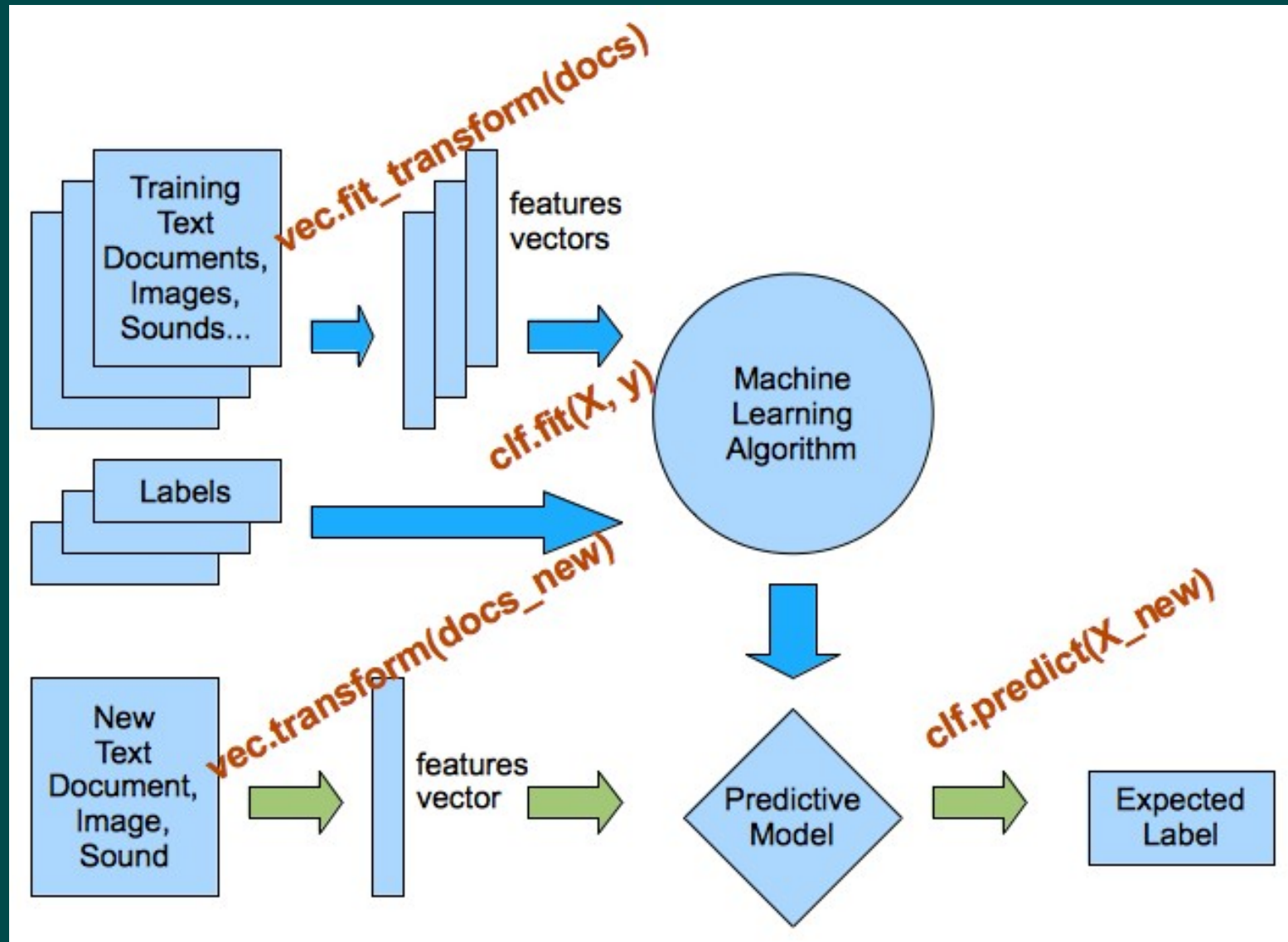
Tom M. Mitchell



# Machine Learning

- Supervised Learning – `model.fit(X, y)`
- Unsupervised Learning – `model.fit(X)`

# Supervised Learning

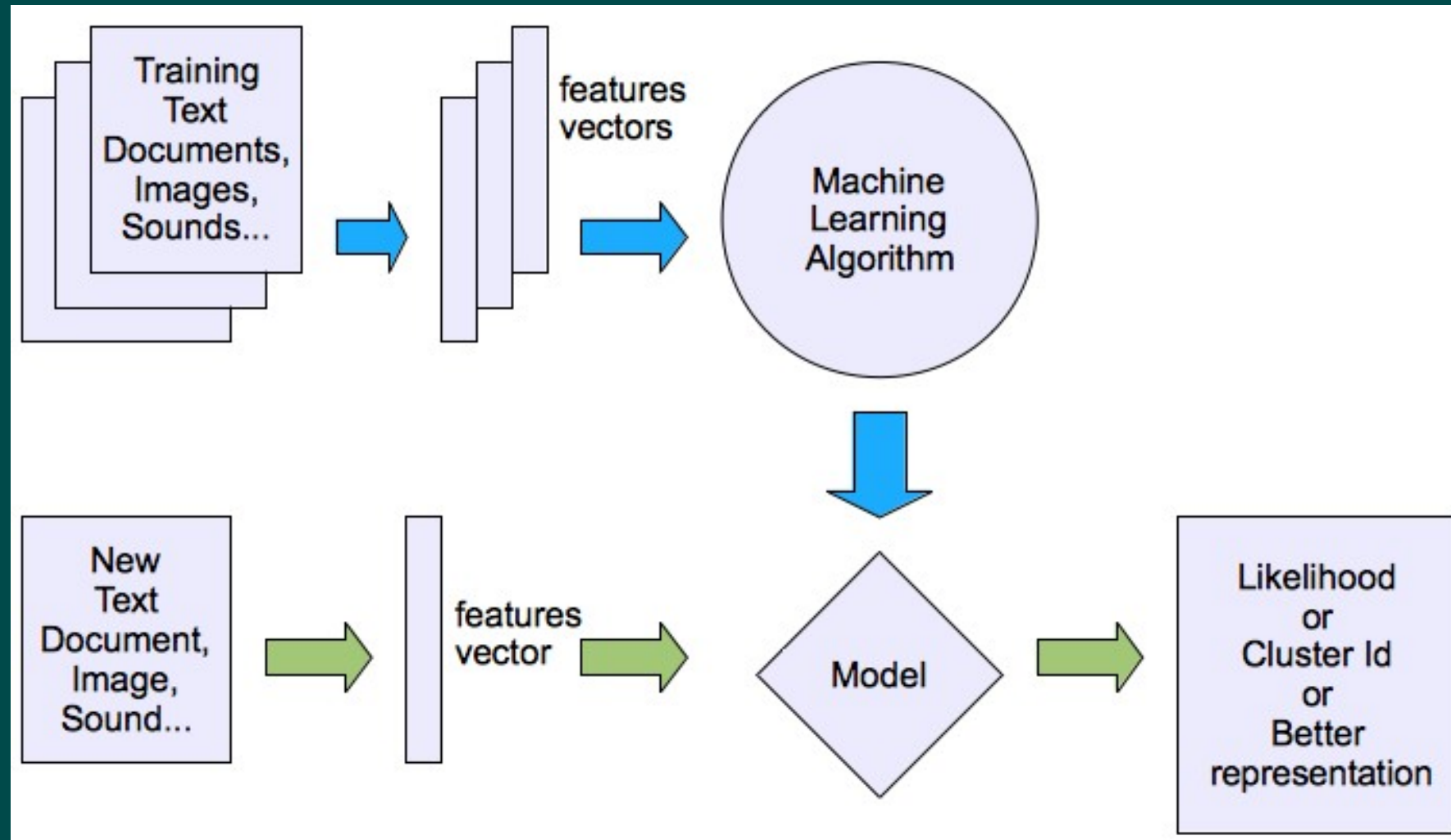


# For example ...

```
from sklearn.linear_model import Ridge as RidgeRegression
from sklearn import datasets
from matplotlib import pyplot as plt
```

```
boston = datasets.load_boston()
X = boston.data
y = boston.target
clf = RidgeRegression()
clf.fit(X, y)
clf.predict(X)
```

# Unsupervised Learning



# For example ...

```
from sklearn.cluster import KMeans  
from numpy.random import RandomState  
rng = RandomState(42)  
k_means = KMeans(3, random_state=rng)  
k_means.fit(X)
```

# What can Scikit-learn do?

Clustering  
Classification  
Regression

# Terminology

- **Model** the collection of parameters you are trying to fit
- **Data** what you are using to fit the model
- **Target** the value you are trying to predict with your model
- **Features** attributes of your data that will be used in prediction
- **Methods** algorithms that will use your data to fit a model

# Steps for Analysis

- Understand the task. See how to measure the performance.
- Choose the source of training experience.
- Decide what will be input and output.
- Choose a set of models to the output function.
- Choose a learning algorithm.

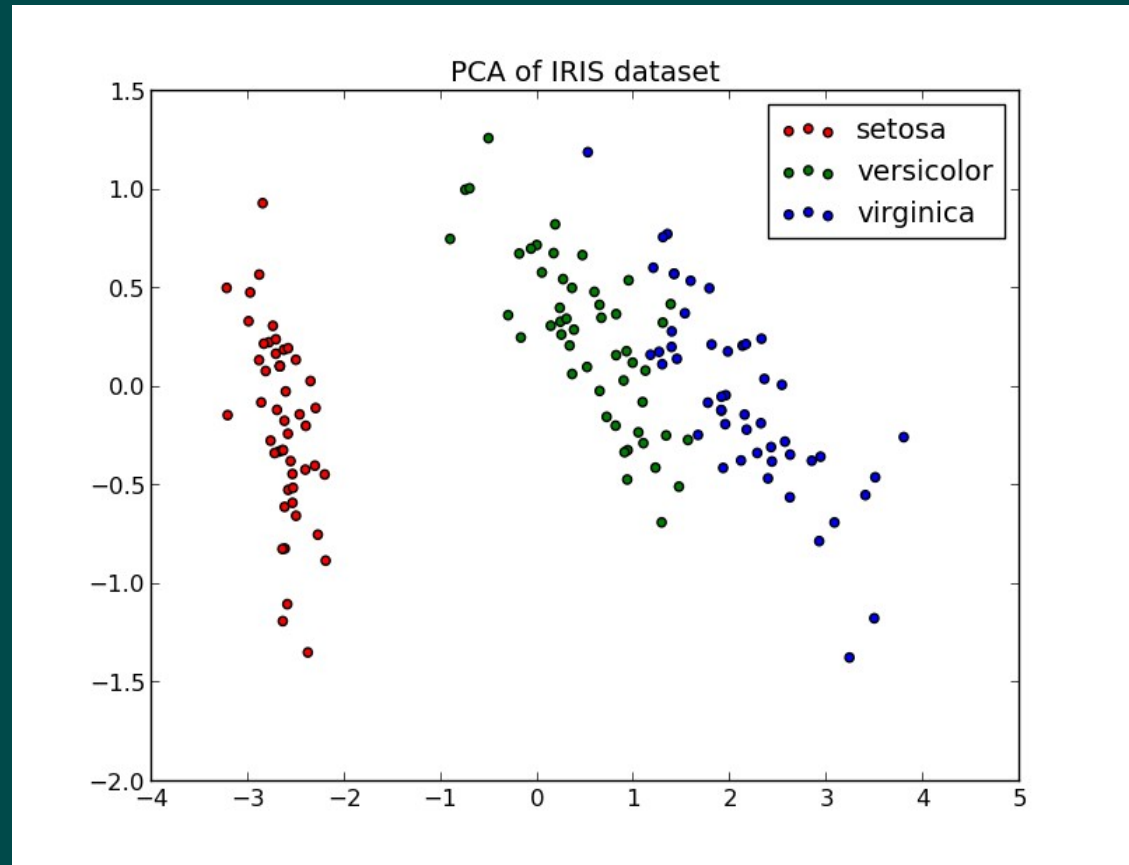
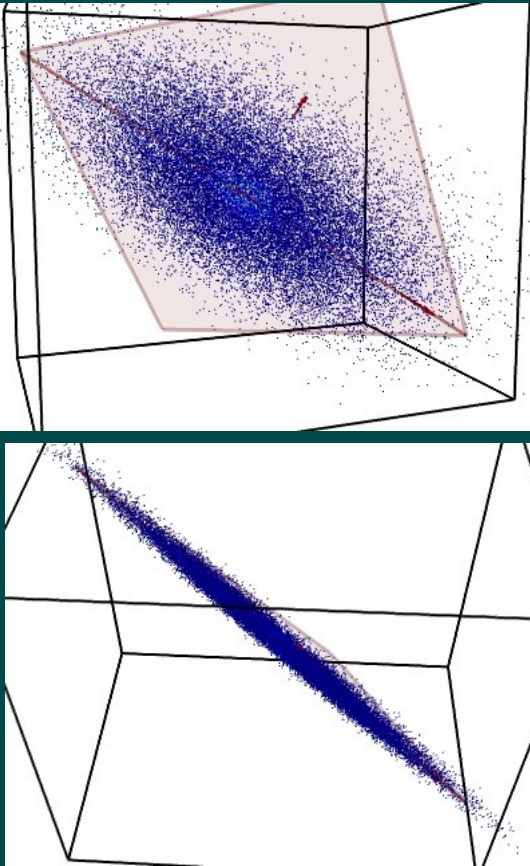


# Steps for Analysis

- Understand the task. See how to measure the performance. Find the right question to ask.
- Choose the source of training experience.
  - Keep training and testing dataset separate. Beware of *overfitting* !
- Decide what will be input and expected output.
- Choose a set of models to approximate the output function. (use dimensionality reduction)
- Choose a learning algorithm. Try different ones ;)

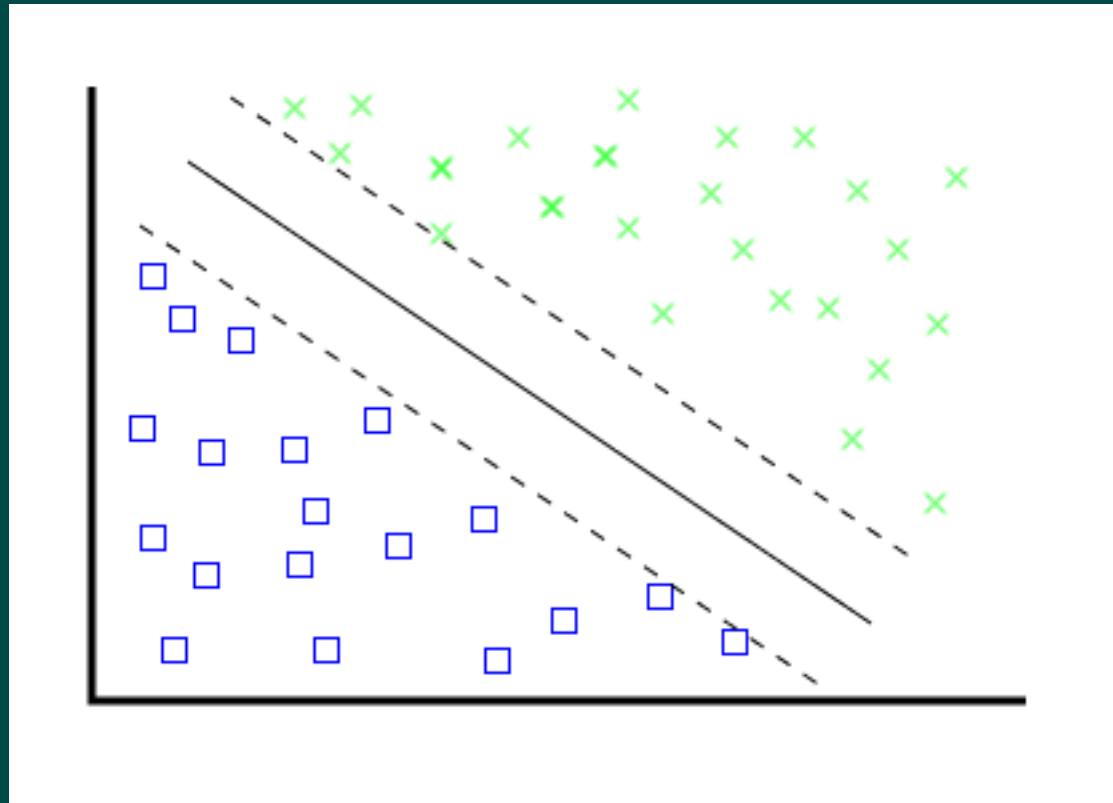
# Some Common Algorithms

## Principal Component Analysis



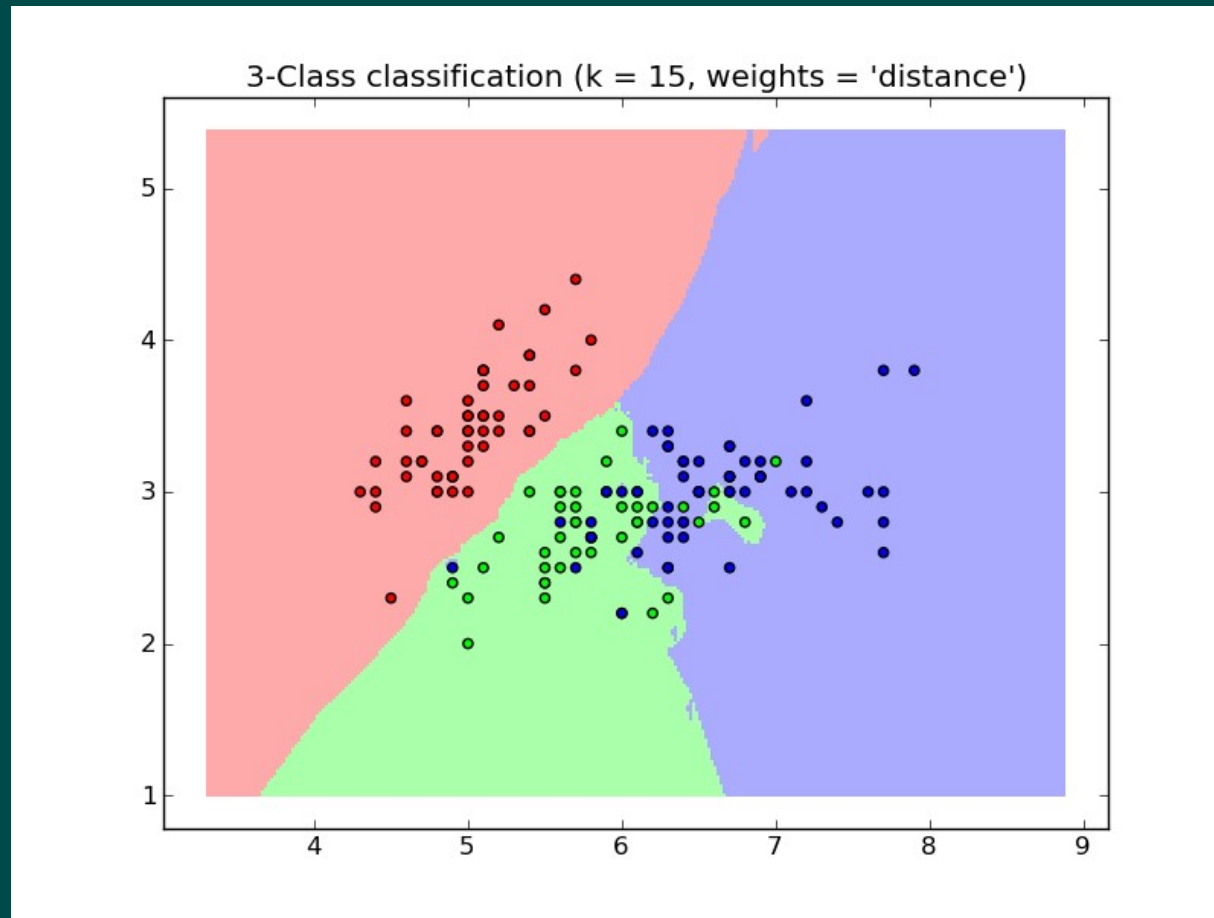
# Some Common Algorithms

- Support Vector Machine



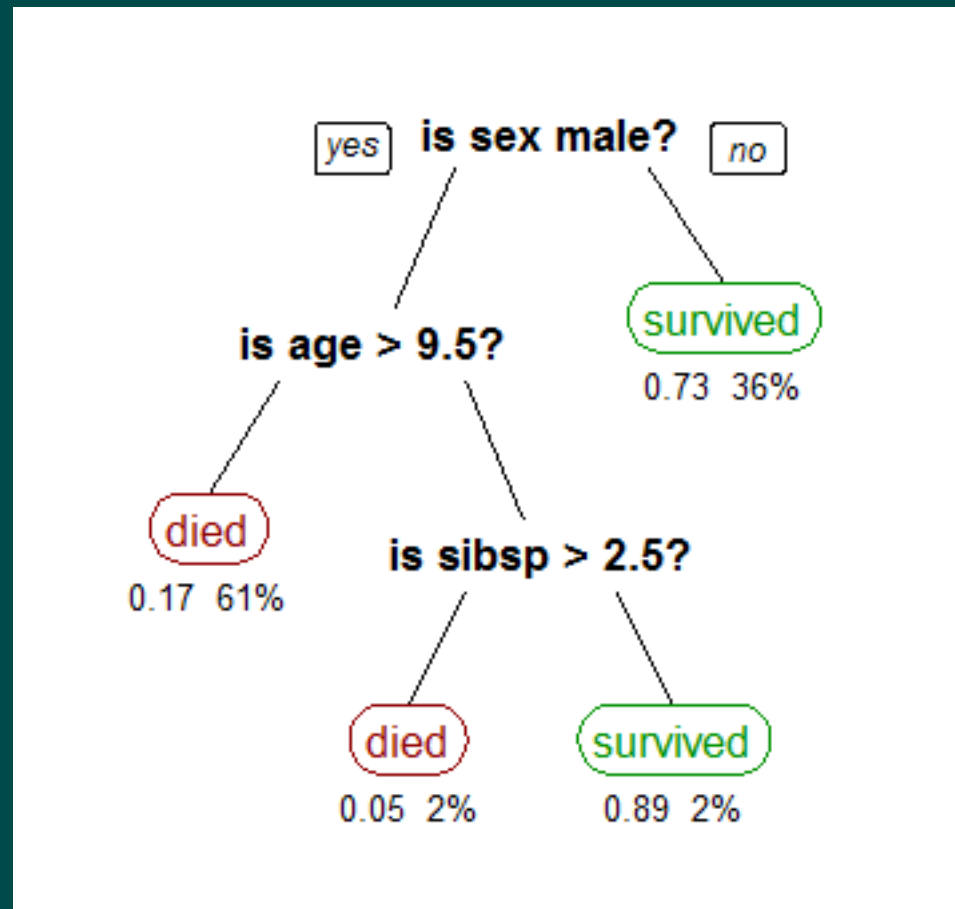
# Some Common Algorithms

- Nearest Neighbour Classifier



# Some Common Algorithms

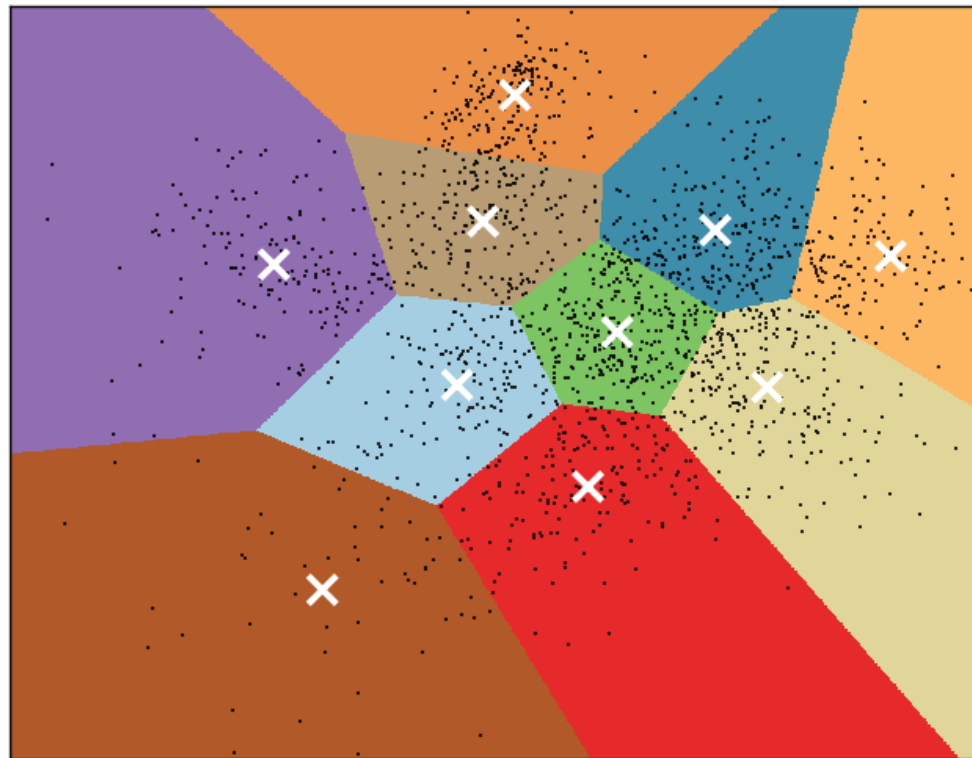
- Decision Tree Learning



# Some Common Algorithms

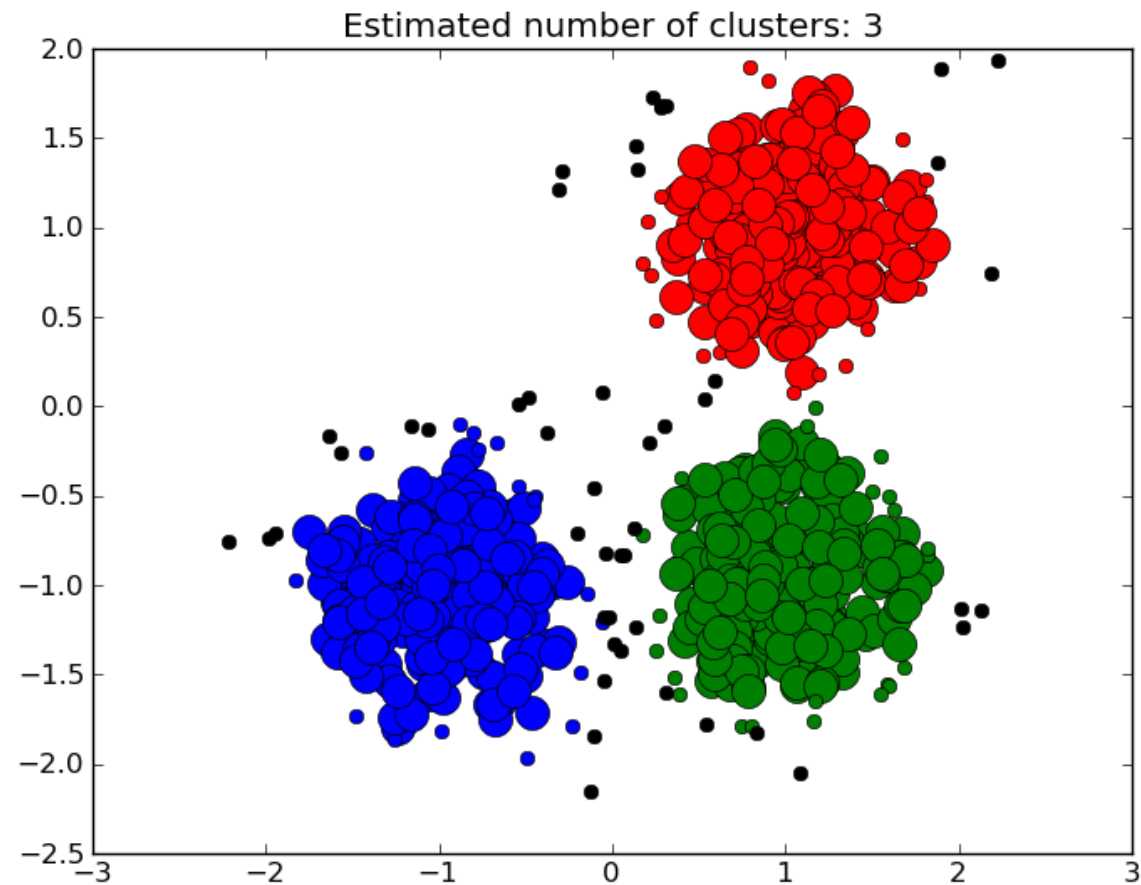
## k-means clustering

K-means clustering on the digits dataset (PCA-reduced data)  
Centroids are marked with white cross



# Some Common Algorithms

## DB SCAN Clustering



## Some Example Usecases

- Log file analysis
- Outlier detection
- Fraud Detection
- Forecasting
- User patterns



## A few comments

- **nlTK** is a good (better) for text processing
- **scikit-learn** is for medium size problems
- for humongous projects, think of **mahout**
- **matplotlib** can be used for visualization
- visualize it in browser using **d3.js**
- have a look at **pandas** for numerical analysis

## Conclusion

- This is just the tip of an iceberg.
- Scikit-learn is really cool to hack with.
- A lot of examples  
([http://scikit-learn.org/stable/auto\\_examples/index.html](http://scikit-learn.org/stable/auto_examples/index.html))

## Final words

`pip install scikit-learn`

Its all in the internet.

Happy Hacking!