

Machine Learning Introduction

Aline Kraemer June 2020

Objective

Introduction to Machine Learning

Hands on: Examples in Python with Scikit-learn

About me





Control and Automation Engineering - UFSC



PhD student in Mechanical and Aeronautical Engineering - ITA



AceleraDev Python Women - Codenation



Backend Engineer - Loadsmart

"It is the practice of using **algorithms** to **parse data**, **learn from it**, and then **make a determination or prediction** about something in the world" – Nvidia





Applications

Speech RecognitionMedical DiagnosisTraffic PredictionImage RecognitionRecommender SystemsStock Market TradingAutonomous DrivingFraud DetectionSearch EngineFace RecognitionSpam Detection



Applications





Important concepts



Data preprocessing Feature Selection/Extraction



Model Selection Hyperparameters Tuning



Model Evaluation

Performance Metrics

Cross-validation



Overfitting vs Underfitting





Overfitting vs Underfitting





Learning strategies





Learning strategies





Learning strategies







Learning strategies





Supervised Learning







Supervised Learning

Algorithms



k-Nearest Neighbors Support Vector Machine Decision Trees Random Forests Logistic Regression Linear Regression Neural Networks



Learning strategies





Unsupervised Learning





Unsupervised Learning

Problem Types



Unsupervised Learning

Algorithms



Hands On

Scikit-learn Library

Scikit-learn Library



Open source machine learning library in Python



Supports supervised and unsupervised learning



Provides tools for: Model fitting Data preprocessing Model selection Model evaluation



























```
[1]
    ► M↓ B→B
     from sklearn import datasets
     iris = datasets.load_iris()
[2]
     ► M↓ B→B
     iris.data[:5]
     array([[5.1, 3.5, 1.4, 0.2],
            [4.9, 3., 1.4, 0.2],
            [4.7, 3.2, 1.3, 0.2],
            [4.6, 3.1, 1.5, 0.2],
            [5., 3.6, 1.4, 0.2]])
[3]
    ► M↓ B→B
     iris.feature_names
```

```
['sepal length (cm)',
'sepal width (cm)',
'petal length (cm)',
'petal width (cm)']
```



[4]	
	<pre>iris.target[:5]</pre>
	array([0, 0, 0, 0, 0])
[5]	► Mi B+B
	iris.target_names
	array(['setosa', 'versicolor', 'virginica'], dtype=' <u10')< th=""></u10')<>
[6]	► Mi Bea
	<pre>import matplotlib.pyplot as plt import seaborn as sns</pre>
	<pre>sns.set(context='talk')</pre>











k-Nearest Neighbors Classification Iris Dataset

[9] ▶ M↓ 🖁→8

from sklearn.model_selection import train_test_split
from sklearn import neighbors

X = iris.data

y = iris.target

X_train, X_test, y_train, y_test = train_test_split(

X, y, test_size=0.3)

n_neighbors = 5
clf = neighbors.KNeighborsClassifier(n_neighbors)
clf.fit(X_train, y_train)



k-Nearest Neighbors Classification Iris Dataset

[10] ⊳ M↓ B→B

from sklearn.metrics import plot_confusion_matrix





k-Nearest Neighbors Classification Iris Dataset

```
[11] ▷ ML B+8
# Using only 2 features
X = iris.data[:,:2]
y = iris.target
n_neighbors = 5
clf = neighbors.KNeighborsClassifier(n_neighbors)
clf.fit(X, y)
```

[12] ⊳ M↓ 🗄→8

```
from matplotlib.colors import ListedColormap
```

```
# Create color maps
```

```
cmap_light = ListedColormap(['#ed5e5f', '#ffa54d', '#bfbfbf'])
cmap_bold = ListedColormap(['#e41a1c', '#ff7f00', '#999999'])
```



[13]	► Mi Bea	
	<pre>import numpy as np</pre>	
	# Create a mesh	
	h = 0.02 # step size	
	<pre>x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1</pre>	
	<pre>y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1</pre>	
	<pre>xx, yy = np.meshgrid(np.arange(x_min, x_max, h),</pre>	
	np.arange(y_min, y_max, h))	
	<pre>Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])</pre>	
	# Dist the desiring boundary	
	# Plot the decision boundary	
	Z = Z.reshape(xx.shape)	
	plt.figure()	
	<pre>plt.pcolormesh(xx, yy, Z, cmap=cmap_light)</pre>	









k-Means Clustering

K = 2 (2 clusters)





K-Means Clustering

K = 2 (2 clusters)



1. Randomly pick K centroids from the sample points as initial cluster centers



K-Means Clustering

K = 2 (2 clusters)



2. Assign each sample to the nearest centroid



K-Means Clustering

K = 2 (2 clusters)



3. Move the centroids to the center of the samples that were assigned to it



K-Means Clustering

K = 2 (2 clusters)



4. Repeat steps 2 and 3 until the cluster assignments do not change or a user-defined tolerance or maximum number of iterations is reached



K-Means Clustering

[1] ▷ M↓ 🗄→8

from sklearn.datasets import make_blobs

Create dataset

[2] ▷ M↓ 🗄→8

```
import matplotlib.pyplot as plt
import seaborn as sns
```

```
sns.set(context='talk')
```



```
[3] ▷ MJ @+8
# Plot
plt.scatter(X[:, 0], X[:, 1], c='white', marker='o',
    edgecolor='black', s=50)
plt.xlabel('Feature 1')
plt.ylabel('Feature 2')
```



K-Means Clustering

[4] ▷ M↓ →8

from sklearn.cluster import KMeans

```
km = KMeans(n_clusters=3, init='random', n_init=10,
    max_iter=300, tol=1e-04, random_state=0)
```

y_km = km.fit_predict(X)



```
[5] ⊳ M↓ 8→8
     # Plot the 3 clusters
     plt.scatter(X[y_km == 0, 0], X[y_km == 0, 1],
         c='lightgreen', edgecolor='black', label='Cluster 1')
     plt.scatter(X[y_km == 1, 0], X[y_km == 1, 1],
         c='orange', edgecolor='black', label='Cluster 2')
     plt.scatter(X[y_km == 2, 0], X[y_km == 2, 1],
         c='lightblue', edgecolor='black', label='Cluster 3')
     # Plot the centroids
     plt.scatter(km.cluster_centers [:, 0], km.cluster_centers [:, 1],
         marker='s', c='red', edgecolor='black', label='Centroids')
     plt.legend(bbox_to_anchor=(1, 0.7))
     plt.xlabel('Feature 1')
     plt.ylabel('Feature 2')
```









Final Remarks



A machine learning **model** is only **as good as** the **data** it is fed.

Thank you!





We are hiring!