

02/06/2020

Supervised Learning with SciKit-Learn

[1] Classification

Supervised Learning

• what is machine learning?

- The art and science of:

- giving computers the ability to learn to make decisions from data.

- without being explicitly programmed.

- Examples:

- learning to predict whether an email is spam or not.

- clustering wikipedia entries into different categories.

- Supervised learning: Uses labeled data.

- Unsupervised learning: Uses unlabelled data.

03/06/2020

• unsupervised learning

- uncovering hidden patterns from unlabelled data.
- Example
 - Grouping customers into distinct categories (clustering).

• Reinforcement learning

- Software agents interact with an environment
 - learn how to optimize their behavior
 - given a system of rewards & punishments.
 - Draws inspiration from behavioral psychology
- Applications
 - Economics
 - Genetics
 - Game Playing
- AlphaGo: first computer to defeat the world champion in Go.

2. Supervised learning

- Predictor variables / features & a target variable.
- Aim: Predict the target variable, given the predictor variables.
 - Classification: Target variable consists of categories.
 - Regression: Target variable is continuous.

2. Naming conventions

- Features = predictor variables = independent variables.
- Target variable = dependent variable = response variable.

2. The Iris Dataset

Features:

- Petal length
- Petal width
- Sepal length
- Sepal width

A) Target variable: Species

- Versicolor
- Virginica
- Setosa

→ The iris dataset in scikit-learn

```
from sklearn import datasets  
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
plt.style.use('ggplot')
```

```
iris = datasets.load_iris()  
type(iris)
```

```
iris.keys()
```

```
# EDA
```

```
X = iris.data  
y = iris.target  
df = pd.DataFrame(X, columns=iris.  
feature_names)  
print(df.head())
```

```
# Visual EDA
```

```
= pd.plotting.scatter_matrix(df, c=y,  
figsize=[8, 15], s=150,  
marker='D')
```

* In this iris dataset we have used scatter plot. But if we have a dataset in the form of "Yes" or "No" then we can use seaborn's countplot().

```
→ plt.figure()  
sns.countplot(x='education', hue='party'  
               data=df, palette='RdBu')  
plt.xticks([0,1], ['No', 'Yes'])  
plt.show()
```

* K-Nearest Neighbors

- Basic idea: Predict the label of a data point by its k closest labeled data points.
- looking at the k closest labeled data points.
- taking a majority vote.

* what is classification?

- A supervised learning approach.
- Categorizing some unknown items into a discrete set of categories or "classes".
- These target attribute is a categorical variable.

- classification algorithms in machine learning:
 1. Decision Trees
 2. Naïve Bayes
 3. Linear Discriminant Analysis
 4. K-Nearest Neighbor
 5. Logistic Regression
 6. Neural Networks
 7. Support Vector Machines (SVM)

• The K-Nearest Neighbor algorithm

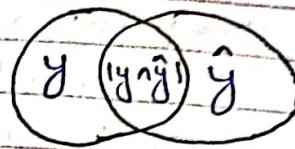
1. Pick a value for k .
2. Calculate the distance of unknown case from all cases.
3. Select the k -observations in the training data that are "nearest" to the unknown data point.
4. Predict the response of one unknown data point using the most popular response value from the k -nearest neighbors.

2. Evaluation matrix in classification

1) Jaccard index

y : Actual labels

\hat{y} : Predicted labels



$$J(y, \hat{y}) = \frac{|y \cap \hat{y}|}{|y \cup \hat{y}|} = \frac{|y \cap \hat{y}|}{|y| + |\hat{y}| - |y \cap \hat{y}|}$$

e.g.: $y: [0, 0, 0, 0, 0, 1, 1, 1, 1, 1]$

$\hat{y}: [1, 1, 0, 0, 0, 1, 1, 1, 1, 1]$

8 correctly predicted

$$J(y, \hat{y}) = \frac{8}{10 + 10 - 8} = 0.66$$

2) Confusion matrix

TP	FN
FP	TN

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

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

F1-Score = $\frac{2 \times (\text{precision} \times \text{recall})}{\text{precision} + \text{recall}}$

3) Log loss

- Performance of a classifier whether the predicted o/p is a probability value b/w 0 & 1.

- Predicting log loss for each row:

$$(y \times \log(\hat{y}) + (1-y) \times \log(1-\hat{y}))$$

- Avg. log loss

$$\text{logloss} = -\frac{1}{n} \sum (y \times \log(\hat{y}) + (1-y) \times \log(1-\hat{y}))$$

→ Scikit-learn fit and predict

- All machine learning models implemented as Python classes

- They implement the algorithm for learning and predicting.

- store the information learned from the data.

- Training a model on the data = "fitting" a model to the data
 - .fit() method.

- To predict the labels of new data:
 - `predict()` method

⇒ using scikit-learn to fit a classifier

```
→ from sklearn.neighbors import  
KNeighborsClassifier  
knn = KNeighborsClassifier(n_neighbors=6)  
knn.fit(iris['data'], iris['target'])  
→ iris['data'].shape
```

⇒ Predicting on unlabelled data

```
→ x_new = np.array([[5.6, 2.8, 3.9, 1.1],  
[5.7, 2.6, 3.8, 1.3],  
[4.7, 3.2, 1.3, 0.2]])
```

```
→ prediction = knn.predict(x_new)
```

```
→ x_new.shape
```

```
→ print('Prediction: {}' .format(prediction))
```

↳ output

Prediction: [1 1 0]

→ Measuring model performance

- could compute accuracy on data used to fit classifier.
- NOT indicative of ability to generalize.
- Split data into training & test set.
- fit/train the classifier on training set.
- Make predictions on the test set.
- compare predictions with the known labels.

Train/Test split

```
→ from sklearn.model_selection import  
train_test_split  
X_train, X_test, y_train, y_test =  
train_test_split(X, y, test_size=0.3,  
random_state=21,  
stratify=y)
```

```
knn = KNeighborsClassifier(n_neighbors=8)  
knn.fit(X_train, y_train)  
y_pred = knn.predict(X_test)  
print('Test set predictions: \n',  
format(y_pred))
```

To check accuracy we use score method.

knn.score(x-test, y-test)

O/P: 0.95555

04/06/2020

Code on K-NN

About the dataset

- Imagine a telecommunication provider has segmented its customer base by service usage patterns, categorizing the customers into 4 groups.

- The example focuses on using demographic data, such as region, age and marital status to predict usage patterns.

- The target field called "custcat", has 4 possible values that correspond to the 4 customer groups, as follows:

1 - Base Service

2 - E-service

3 - Plus-Service

4 - Total Service

```
# let's load required libraries
```

```
import itertools  
import numpy as np  
import matplotlib.pyplot as plt  
from matplotlib.ticker import NullFormatter  
import pandas as pd  
import matplotlib.ticker as ticker  
from sklearn import preprocessing  
%matplotlib inline
```

```
import itertools:
```

This module is used for implementing number of iteration tools.

```
from matplotlib.ticker import NullFormatter:
```

```
import matplotlib.ticker as ticker:
```

```
matplotlib.ticker
```

This module contains classes to support completely configurable tick locating and formatting.

```
# load data from csv file
```

```
df = pd.read_csv('telecomroot.csv')  
df.head()
```

```
# Data visualization & analysis
```

```
# let's see how many of each class is in our dataset.
```

```
df['custcat'].value_counts()
```

281 - Plus service

266 - Basic service

236 - Total service

217 - E-service

```
# feature set
```

Let's define feature sets, X :

```
x = df[['region', 'tenure', 'age', 'marital',  
        'address', 'income', 'ed', 'employ',  
        'retire', 'gender', 'reside']].values
```

```
X[0:5]
```

```
y = df['custcat'].values  
y[0:5]
```

Normalize Data

→ Data standardization give data zero mean and unit variance, it is good practice, especially for algorithms such as KNN which is based on distance of cases.

```
x = preprocessing.StandardScaler().fit(x).  
transform(x.astype(float))
```

```
x[0:5]
```

Train/Test Split

→ Out of Sample Accuracy is the percentage of correct predictions that the model makes on data that the model has not been trained on.

→ Doing a train and test on the same dataset will most likely have low out-of-sample accuracy, due to the likelihood of being over-fit.

→ It is important that our models have a high, out-of-sample accuracy; because the purpose of any model, of course, is to make correct predictions on unknown data.

```
from sklearn.model_selection import  
train_test_split  
X_train, X_test, y_train, y_test =  
train_test_split(X, y, test_size=0.2, random_state=4)  
print('Train set:', X_train.shape, y_train.shape)  
print('Test set:', X_test.shape, y_test.shape)
```

Classification

K Nearest Neighbor (KNN)

```
# import library
```

```
from sklearn.neighbors import KNeighborsClassifier
```

Training

→ let's start the algorithm with $K=4$ for now.

$K=4$

```
# Train model
```

```
neigh = KNeighborsClassifier(n_neighbors = K).
```

```
neigh.fit(X_train, y_train)
```

neigh

Predicting

→ we can use the model to predict the test set:

```
yhat = neigh.predict(x-test)  
yhat[0:5]
```

Accuracy Evaluation

→ In multilabel classification, accuracy classification score is a function that computes subset accuracy.

```
from sklearn import metrics  
print("Train set Accuracy:", metrics.accuracy_  
score(y-train, neigh.predict(x-train)))  
print("Test set Accuracy:", metrics.accuracy_  
score(y-test, yhat))
```

what about other k?

→ k in KNN is the number of nearest neighbors to examine. It is supposed to be specified by user.
So, how can we choose right value for k?
↳ The general solution is to reserve a part of your data for testing the accuracy of the model.

- ↳ Then choose $k=1$, use the training part for modelling, and calculate the accuracy of prediction using all samples in your test set.
- ↳ Repeat this process increasing the k , and see which k is the best for your model.
- We can calculate the accuracy of KNN for different k s.

$ks = 10$

mean_acc = np.zeros((ks-1))

std_acc = np.zeros((ks-1))

confusionMX = [[0] * n for n in range(1, ks+1)]

for n in range(1, ks+1):

Train model & predict

neigh = KNeighborsClassifier(n_neighbors=n).

fit(x_train, y_train)

yhat = neigh.predict(x_test)

mean_acc[n-1] = metrics.accuracy_score(y_test, yhat)

std_acc[n-1] = np.std(yhat == y_test) /
np.sqrt(yhat.shape[0])

mean_acc

Plot model accuracy for different number of Neighbors

```
plt.plot (range(1,ks), mean-acc, 'g')
plt.fill_between (range(1,ks), mean-acc-1*std-acc, mean-acc + 1 * std-acc, alpha=0.10)
plt.legend ('Accuracy', '+(-3xstd)')
plt.ylabel ('Accuracy')
plt.xlabel ('Number of Neighbors (K)')
plt.tight_layout()
plt.show()
```

```
print ("The best accuracy was with",  
      mean-acc.argmax(), "with k =",  
      mean-acc.argmax() + 1)
```

The best accuracy was with 0.34 with
k = 9