

SI 670 Notes

Suggested books

- Introduction to Machine Learning with Python A Guide for Data Scientists By Andreas C. Müller and Sarah Guido
- Deep Learning with Python by Francois Chollet

Top libraries

- Scikit-learn
- SciPy
- Numpy
- Pandas
- Matplotlib

Cycle

- Feature representation
- Training
- Evaluation
- Refine cycle (hyperparameterization)

Data quality checks

- Min/max summaries
- Wrong data type, units
- Equal class representation
- Outliers
- Data distribution
- Correlations among variables

KNN Notes

Category

- Supervised
 - Classification
 - Regression

High level algorithm

Given a training set X_{train} with labels y_{train} , and given a new instance x_{test}

1. Find the observations that resemble x_{test} that are in X_{train} . Call this set of observation(s) X_{nn}
2. Get the labels of Y_{nn} for the instances in X_{nn}
3. Predict label for x_{test} by combining the labels Y_{nn} (majority vote).

Parameters

- Distance metric (Euclidian)
- Choice of k (k=1 very flexibel, k=100 rigid)
- Weighting function (neighbors that are far less influence on final prediction)

Evaluation

- Accuracy (correctly predicted / total observations) (for classification)
- R^2 (for regression, measure how does the data fit the model 0-1)

Extras

- Ensure that all observations are on the same scale
 - if not standardize them (standard scalar)

Classification

```
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X_C1, y_C1, random_state = 0)
knnnc = KNeighborsClassifier(n_neighbors = 5).fit(X_train, y_train)
print(knnnc.predict(X_test))
print('Accuracy test score: {:.3f}'.format(knnnc.score(X_test, y_test)))
```

Regression

```
from sklearn.neighbors import KNeighborsRegressor
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X_R1, y_R1, random_state = 0)
knnreg = KNeighborsRegressor(n_neighbors = 5).fit(X_train, y_train)
print(knnreg.predict(X_test))
print('R-squared test score: {:.3f}'.format(knnreg.score(X_test, y_test)))
```