

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)
knnc = KNeighborsClassifier(n_neighbors = 5).fit(X_train, y_train)
print(knnc.predict(X_test))
print('Accuracy test score: {:.3f}'.format(knnc.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)))
```

Linear Regression

Supervised

By now it should be clear that when we have a quantitative, or qualitative response, and we want to train a model and then predict new, unseen observations with our model we are in the world of supervised learning.

Assumption

Now one of the main assumptions with this type of learning is that we expect the training data to have the same structure/relationships as the test data.

Linear Regression

Is type of linear model in which we attempt to predict the target value (quantitative response) with a weighted sum of features (predictors). The usual formula is $Y = b_0 + B_1X_1 + E$

- the goal is to minimize the residual sum of squares, which is the sum of the squared differences between y and \hat{y} , or actual minus predicted.

Evaluation

- R^2 , a.k.a *coefficient of determination* measures how well a prediction model fits a determinate dataset. The value is between 0 & 1.
- Look at the distribution of the residuals. Are they constant, or is there some systematic bias?

Cross-validation

It uses multiple test-train splits. Each split is used to train the model and get an error metric. At the end, an overall average is computed. This way we are able to get a more stable and realistic performance of our model.

- k-Fold CV, divide the dataset in k folds, run k models
- Stratified CV, so as to allow class proportions to be preserved in the splits.
- Leave-one-out CV, run n models

Linear Regression vs KNN-Regression

Linear	KNN
few parameters	Non-Parametric
Small dataset	Large dataset needed
Generalizes beyond	Limited generalization

Polynomial feature expansion

Generate new features consisting of all polynomial combinations of the original two features. The degree of the polynomial specifies how many variables participate at a time in each new feature. It is still a linear model, and can use same least-squares estimation method

Why?

- To capture interactions between the original features by adding them as features to the linear model.
- To make a classification problem easier

Pitfalls?

- Beware of polynomial feature expansion with high degree as this can lead to complex models that overfit

Generalization

Refers to an algorithm's ability to give accurate predictions for new, previously unseen data

- Models that are too complex for the amount of training data available are said to overfit and are not likely to generalize well to new examples.
- Models that are too simple, that don't even do well on the training data, are said to underfit and also not likely to generalize well.
- Refer to bias-variance trade off for an indepth analysis of why over/underfitting happens

Linear Regression

```
from sklearn.linear_model import LinearRegression
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)
linreg = LinearRegression().fit(X_train, y_train)
print("linear model intercept (b): {}".format(linreg.intercept_))
print("linear model coeff (w): {}".format(linreg.coef_))
```

Polynomial feature expansion regression

```
from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import PolynomialFeatures

poly = PolynomialFeatures(degree=degree)

#Applying standardization & reshaping(was giving an error!)
X_train_poly = poly.fit_transform(X_train.reshape(-1, 1))

X_test_poly = poly.fit_transform(X_test.reshape(-1,1))

#Fit the model
linreg = LinearRegression().fit(X_train_poly, y_train)

#Prediction
y_pred = linreg.predict(X_train_poly)

#Score
training_score = linreg.score(X_train_poly, y_train)
testing_score = linreg.score(X_test_poly, y_test)

r2_train[degree]= training_score
r2_test[degree] = testing_score
```

Logistics Regression

Bias-Variance decomposition, reducible error of model can be decomposed into $\text{Bias}^2 + \text{Variance}$. High variance our $f(\cdot)$ estimate is not stable, it produces different results for different data from same population. High bias, our model has a systematic error, which impedes it from capturing the true relationships within the data. Often, as we increase model complexity - BIAS tends to lower, & VARIANCE tends to increase.

Regularization, allows us to avoid overfitting model to data. The idea revolves around: preferring small coefficients for predictors (reduce variance) & simpler models over more complicated ones (i.e. few predictors better than all predictors).

Ridge regression, L2 regularization, alpha parameter for controlling learning, ensure to normalize inputs. Assumption is that we have many small/medium sized coefficients.

Lasso regression, L1 regularization, alpha parameter for controlling learning. Only few variables with medium/large effect.

Robust regression, fit model on all points, refit model on all points & check results for which fit works out, fit once more leaving out all points that do not fit model.

Logistics regression, it is known as a soft classifier given it “spits” out a probability for a particular observation & then it assigns that observation to the class for which has the highest probability. Model $1/1+\exp(-(BO+B1X1\dots))$. It takes as input any value & spits out a number within (0,1).

Evaluation metrics *Accuracy*, total correctly predicted/total observations. Now, this is not a good metric given it omits important information from the data (i.e. think about accuracy result/reality in a case in which we have imbalanced data).

Dummy classifier, is a sanity check against our model so as to ensure that our results are somewhat valid. It is not a real classifier just a random simulation to compare against our results.

Confusion matrix, gives a better idea about the model’s performance in which we have the predicted vs the actual & their respective ratios. Possible options true positive (class1: predicted correctly), true negatives (class 0: predicted correctly), false positives (Class 0: predicted by classifier as belonging to class1), false negatives (class 1: predicted by classifier as belonging to class 0)

Classification error, 1- accuracy, incorrect classifications

Recall, TPR, sensitivity, $TP / TP+FN$, class 1 correctly predicted/ class1 total

Precision, $TP/TP+FP$, correctly predicted class1/ total observation predicted as belonging to class1

False Positive Rate, $FP/TN+FP$, class 0 WRONGLY predicted/ class0 total

Specificity, $TN/TN+ FP$, class 0 correctly predicted/ class0 total

Precision - Recall trade off, ideally we want high recall & high precision, but as it often true we can’t have both, and we have to pick one, which one? It depends on the problem that we are trying to solve / application. High-Precision: low FP, but may have high FN High-Recall: low FN, but may have high FP

F1, $2 \text{ precision} \times \text{recall} / (\text{precision} + \text{recall})$, combines recall & precision in one number

Decision functions, given a score from each point build precision-recall thresholds..

ROC, rate of FP to Tp, ideally 0 FP & ALL TP *AUC*, summarizing in one number (0 bad, 1 good) area under curve for ROC