

# Feature Selection

Daniel Cerdán, Fernando Freire

April 7, 2019

## Contents

<b>1</b>	<b>Feature Selection</b>	<b>2</b>
1.1	Question 1 . . . . .	3
1.1.1	Response 1 . . . . .	3
1.2	Question 2 . . . . .	6
1.2.1	Response 2 . . . . .	6
1.3	Question 3 . . . . .	8
1.3.1	Response 3 . . . . .	8
1.4	Question 4 . . . . .	10
1.4.1	Response 4 . . . . .	10

# 1 Feature Selection

In this practical, you will become familiarized with some basic feature selection methods implemented in scikit-learn. Consider the prostate dataset that is attached to this practical. You are asked to:

1. Estimate the performance of the nearest neighbor classifier on this dataset using 10-fold cross validation when all the features are used for prediction. The number of neighbors should be chosen using an inner cross-validation procedure. You can use 5-fold cross validation for this.
2. Estimate the performance of the nearest neighbor classifier on the same dataset when using a feature selection technique based on the F-score (ANOVA) that picks up the 10 most relevant features. Use the same cross-validation methods as in the previous step.
3. Repeat the previous experiment but when a random forest is used to pick up the 10 most relevant features. Use an initial filtering method based on the F-score to keep only the 20% most promising features.
4. What feature selection method performs best? Can you explain why?

Now we will address the problem of analyzing the trade-off between interpretability and prediction accuracy. For this, you are asked to:

1. Estimate the performance of the nearest neighbor classifier with  $K=3$  as a function of the features used for prediction. Use a 10-times 10-fold cross-validation method and plot the results obtained. That is prediction error vs. the number of features used for prediction. Use the F-score for feature selection. Report results from 1 feature to 200 features. Not all features need to be explored. Use a higher resolution when you are closer to 1 feature.
2. Repeat that process when the feature selection is done externally to the cross-validation loop using all the available data. Include these results in the previous plot.
3. Are the two estimates obtained similar? What are their differences? If they are different try to explain why this is the case.
4. By taking a look at these results, what is the optimal number of features to use in this dataset?
5. Given the results obtained in this part of the practical, you are asked to indicate which particular features should be used for prediction on this dataset. Include a list with them. Take a look at the documentation of SelectKBest from scikit-learn to understand how to do this. Use all available data to provide such a list of features.

Script 1.0.1 (python)

```
1 import warnings
2 warnings.filterwarnings("ignore")
3 %matplotlib inline
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 import matplotlib.lines as mlines
8 import matplotlib as mpl
9 from matplotlib import colors
10 import seaborn as sns; sns.set()
```

```

11 import scipy.stats as stats
12 import scipy as sp
13 from scipy import linalg
14 from sklearn.neighbors import KNeighborsClassifier
15 from sklearn.ensemble import RandomForestClassifier
16 from sklearn.pipeline import Pipeline
17 from sklearn.model_selection import train_test_split, RepeatedStratifiedKFold, KFold,
    → GridSearchCV, cross_val_score
18 from sklearn import preprocessing
19 from sklearn.feature_selection import SelectKBest, f_classif, SelectFromModel
20 from sklearn.metrics import accuracy_score, make_scorer, confusion_matrix,
    → classification_report, precision_score

```

## 1.1 Question 1

Estimate the performance of the nearest neighbor classifier on this dataset using 10-fold cross validation when all the features are used for prediction. The number of neighbors should be chosen using an inner cross-validation procedure. You can use 5-fold cross validation for this.

### 1.1.1 Response 1

We obtain a mean accuracy of 0.74 with std 0.15

#### Script 1.1.1 (python)

```

1 # Load data
2 data = pd.read_csv('prostate.csv', header = 0)
3 X = data.values[:, 0:-1].astype(np.float)
4 y = (data.values[:, -1] == 1).astype(np.int)

```

#### Script 1.1.2 (python)

```

1 X_train, X_test, y_train, y_test = train_test_split( \
2     X, y, test_size=0.3, random_state=1)
3 scaler = preprocessing.StandardScaler().fit(X_train)
4 X_train_scaled = scaler.transform(X_train)
5 X_test_scaled = scaler.transform(X_test)
6
7 #Check
8 print(np.var(X_train[:,0]))
9 print(np.var(X_train[:,1]))
10 print(np.var(X_train_scaled[:,0]))
11 print(np.var(X_train_scaled[:,1]))

```

#### Output

```

0.019316280607701256
0.021687892317377073
0.9999999999999994

```

## Script 1.1.3 (python)

```

1 pipeline = Pipeline([ ('knn', KNeighborsClassifier()) ])
2 N_neighbors = [1, 3, 5, 7, 11, 15, 20, 25, 30]
3 param_grid = { 'knn__n_neighbors': N_neighbors }
4 skfold = RepeatedStratifiedKFold(n_splits=10, n_repeats=1, random_state=0)
5 gridcv = GridSearchCV(pipeline, cv=skfold, n_jobs=1, param_grid=param_grid, \
6     scoring=make_scorer(accuracy_score))
7 result = gridcv.fit(X_train_scaled, y_train)
8 accuracies = gridcv.cv_results_['mean_test_score']
9 std_accuracies = gridcv.cv_results_['std_test_score']
10 test_accuracies = np.ones(len(N_neighbors))
11
12 for i in range(len(N_neighbors)):
13     knn = KNeighborsClassifier(n_neighbors = N_neighbors[ i ])
14     knn.fit(X_train_scaled, y_train)
15     test_accuracies[ i ] = accuracy_score(knn.predict(X_test_scaled), y_test)

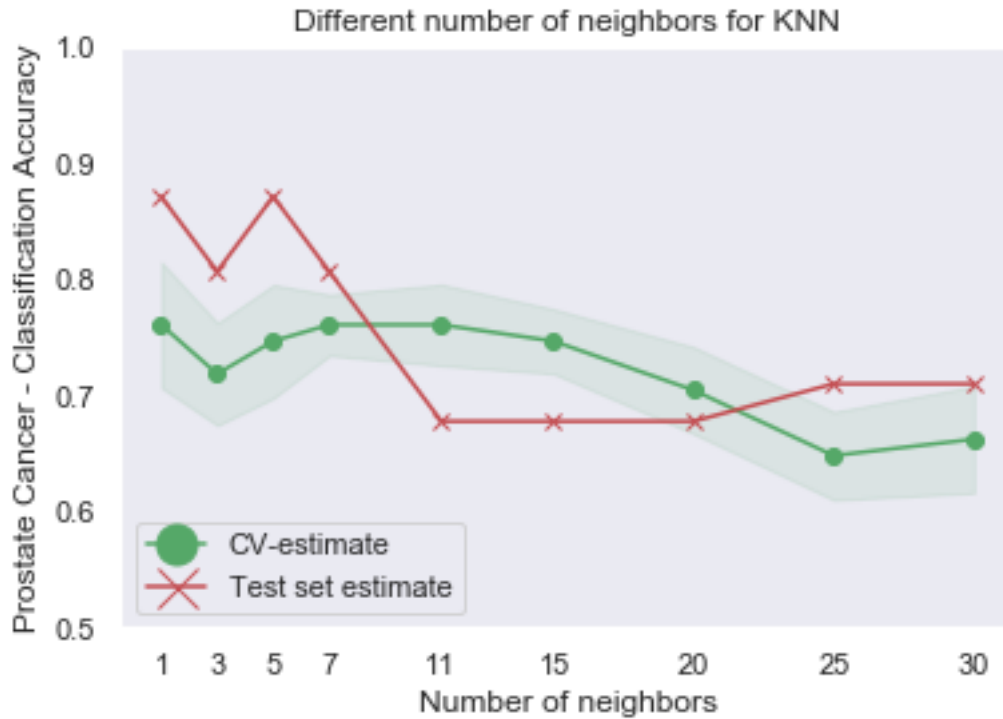
```

## Script 1.1.4 (python)

```

1 plt.figure()
2 line1, = plt.plot(N_neighbors, accuracies, 'o-', color="g")
3 line2, = plt.plot(N_neighbors, test_accuracies, 'x-', color="r")
4 plt.fill_between(N_neighbors, accuracies - std_accuracies / np.sqrt(10), \
5     accuracies + std_accuracies / np.sqrt(10), alpha=0.1, color="g")
6 plt.grid()
7 plt.title("Different number of neighbors for KNN")
8 plt.xlabel('Number of neighbors')
9 plt.xticks(np.array(N_neighbors))
10 plt.ylabel('Prostate Cancer - Classification Accuracy')
11 plt.ylim((0.5, 1.0))
12 legend_handles = [ mlines.Line2D([], [], color='g', marker='o', \
13     markersize=15, label='CV-estimate'), \
14     mlines.Line2D([], [], color='r', marker='x', \
15     markersize=15, label='Test set estimate')]
16 plt.legend(handles=legend_handles, loc = 3)
17 plt.show()
18
19 N_neighbors_selected = N_neighbors[np.argmax(accuracies)]
20 print("Selected optimal number of neighbors", N_neighbors_selected)

```



## Output

Selected optimal number of neighbors 1

## Script 1.1.5 (python)

```
1 # Training
2 knn = KNeighborsClassifier(n_neighbors = N_neighbors_selected)
```

## Script 1.1.6 (python)

```
1 # Evaluating performance for all features
2 (_, n_features) = X.shape
3 print("Number of features", n_features)
4 n_splits = 10
5 random_state = 1
6 rkf = KFold(n_splits=n_splits, random_state=random_state)
7 split = 0
8 knn_scores = np.zeros(n_splits)
9 for train_index, test_index in rkf.split(X, y):
10     print("Partition", split)
11     X_train, X_test = X[train_index], X[test_index]
12     y_train, y_test = y[train_index], y[test_index]
13
14 # Standardization
```

```

15     scaler.fit(X_train, y_train)
16     X_train = scaler.transform(X_train)
17     X_test = scaler.transform(X_test)
18
19     # Training & scoring
20     knn.fit(X_train, y_train)
21     knn_scores[split] = knn.score(X_test, y_test)
22     split += 1
23 print()
24 print("Accuracy scores", knn_scores)
25 print()
26 print("Mean", knn_scores.mean(), "Std", knn_scores.std())

```

## Output

```

Number of features 12625
Partition 0
Partition 1
Partition 2
Partition 3
Partition 4
Partition 5
Partition 6
Partition 7
Partition 8
Partition 9

Accuracy scores [0.54545455 0.81818182 0.6          1.          0.8          0.5
0.8          0.7          0.7          0.9          ]

Mean 0.7363636363636364 Std 0.1498759817894628

```

## 1.2 Question 2

Estimate the performance of the nearest neighbor classifier on the same dataset when using a feature selection technique based on the F-score (ANOVA) that picks up the 10 most relevant features. Use the same cross-validation methods as in the previous step.

### 1.2.1 Response 2

We obtain a mean accuracy score of 0.94 with std mean error 0.02

#### Script 1.2.1 (python)

```

1 n_features = 10
2 print("Number of features", n_features)
3 n_splits = 10
4 random_state = 2
5 np.random.seed(0)
6 rkf = KFold(n_splits=n_splits, random_state=random_state)

```

```

7 split = 0
8 knn_scores = np.zeros(n_splits)
9 knn_errors = np.zeros(n_splits)
10 knn_scores_train = np.zeros(n_splits)
11 for train_index, test_index in rkf.split(X, y):
12     print("Partition", split)
13     X_train, X_test = X[train_index], X[test_index]
14     y_train, y_test = y[train_index], y[test_index]
15
16     # Standardization
17     scaler.fit(X_train, y_train)
18     X_train = scaler.transform(X_train)
19     X_test = scaler.transform(X_test)
20
21     # Filtering
22     filtering = SelectKBest(f_classif, k = n_features)
23     filtering.fit(X_train, y_train)
24     X_train_filtered = filtering.transform(X_train)
25     X_test_filtered = filtering.transform(X_test)
26
27     # Training & scoring
28     knn.fit(X_train_filtered, y_train)
29     knn_scores[split] = knn.score(X_test_filtered, y_test)
30     knn_scores_train[split] = knn.score(X_train_filtered, y_train)
31     knn_errors[split] = 1.0 - np.mean(knn.predict(X_test_filtered) == y_test)
32     split += 1
33 print()
34 print("Accuracy scores", knn_scores)
35 print()
36 print("Mean", knn_scores.mean(), "std mean error", (np.std(knn_scores) /
37     → np.sqrt(len(knn_scores))))
38 print()
39 print("Errors", knn_errors)
40 print()
41 print("Mean", knn_errors.mean(), "std mean error", (np.std(knn_errors) /
42     → np.sqrt(len(knn_errors))))
43 print()
44 print("Accuracy scores train", knn_scores_train)
45 print()
46 print("Mean", knn_scores_train.mean(), "std mean error", (np.std(knn_scores_train) /
47     → np.sqrt(len(knn_scores_train))))
48 print()

```

## Output

```

Number of features 10
Partition 0
Partition 1
Partition 2
Partition 3
Partition 4

```

```

Partition 5
Partition 6
Partition 7
Partition 8
Partition 9

Accuracy scores [1.          0.90909091 1.          0.8          1.          0.9
0.9          1.          0.9          1.          ]

Mean 0.9409090909090908 std mean error 0.02081996816472506

Errors [0.          0.09090909 0.          0.2          0.          0.1
0.1          0.          0.1          0.          ]

Mean 0.05909090909090908 std mean error 0.02081996816472506

Accuracy scores train [1. 1. 1. 1. 1. 1. 1. 1. 1. 1.]

Mean 1.0 std mean error 0.0

```

### 1.3 Question 3

Repeat the previous experiment but when a random forest is used to pick up the 10 most relevant features. Use an initial filtering method based on the F-score to keep only the 20% most promising features.

#### 1.3.1 Response 3

We obtain a mean accuracy score of 0.92 with std mean error 0.02

##### Script 1.3.1 (python)

```

1 n_features = 10
2 print("Number of features", n_features)
3 n_splits = 10
4 random_state = 2
5 np.random.seed(0)
6 rkf = KFold(n_splits=n_splits, random_state=random_state)
7 split = 0
8 knn_scores = np.zeros(n_splits)
9 knn_errors = np.zeros(n_splits)
10 knn_scores_train = np.zeros(n_splits)
11 rf_selection = SelectFromModel(RandomForestClassifier(n_estimators = 2000,\
12                                                         random_state = random_state),
13                               → threshold = 0.0)
14 for train_index, test_index in rkf.split(X, y):
15     print("Partition", split)
16     X_train, X_test = X[train_index], X[test_index]
17     y_train, y_test = y[train_index], y[test_index]
18     # Standardization
19     scaler.fit(X_train, y_train)

```



```

20 X_train = scaler.transform(X_train)
21 X_test = scaler.transform(X_test)
22
23 # Filtering pass 1 via F-score, only 20%
24 _, inner_features = X_train.shape
25 features_to_retain = int(20*inner_features/100)
26 filtering = SelectKBest(f_classif, k = features_to_retain)
27 filtering.fit(X_train, y_train)
28 X_train_filtered_F = filtering.transform(X_train)
29 X_test_filtered_F = filtering.transform(X_test)
30 _, n_features_reduction_F = X_train_filtered_F.shape
31 print("Features reduced to (first pass)", n_features_reduction_F)
32
33 # Filtering random forest
34 rf_selection.fit(X_train_filtered_F, y_train)
35 rf_selection.threshold = -1.0 * np.sort(-1.0 *
    ↪ rf_selection.estimator_.feature_importances_)[n_features-1]
36 X_train_filtered_F_R = rf_selection.transform(X_train_filtered_F)
37 X_test_filtered_F_R = rf_selection.transform(X_test_filtered_F)
38 _, n_features_reduction = X_train_filtered_F_R.shape
39 print("Features reduced to (second pass)", n_features_reduction)
40
41 # Training & scoring
42 knn.fit(X_train_filtered_F_R, y_train)
43 knn_scores[split] = knn.score(X_test_filtered_F_R, y_test)
44 knn_errors[split] = 1.0 - np.mean(knn.predict(X_test_filtered_F_R) == y_test)
45 knn_scores_train[split] = knn.score(X_train_filtered_F_R, y_train)
46 split += 1
47 print()
48 print("Accuracy scores", knn_scores)
49 print()
50 print("Mean", knn_scores.mean(), "std mean error", (np.std(knn_scores) /
    ↪ np.sqrt(len(knn_scores))))
51 print()
52 print("Errors", knn_errors)
53 print()
54 print("Mean", knn_errors.mean(), "std mean error", (np.std(knn_errors) /
    ↪ np.sqrt(len(knn_errors))))
55 print()
56 print("Accuracy scores train", knn_scores_train)
57 print()
58 print("Mean", knn_scores_train.mean(), "std mean error", (np.std(knn_scores_train) /
    ↪ np.sqrt(len(knn_scores_train))))
59 print()

```

## Output

```

Number of features 10
Partition 0
Features reduced to (first pass) 2525
Features reduced to (second pass) 10

```

```

Partition 1
Features reduced to (first pass) 2525
Features reduced to (second pass) 10
Partition 2
Features reduced to (first pass) 2525
Features reduced to (second pass) 10
Partition 3
Features reduced to (first pass) 2525
Features reduced to (second pass) 10
Partition 4
Features reduced to (first pass) 2525
Features reduced to (second pass) 10
Partition 5
Features reduced to (first pass) 2525
Features reduced to (second pass) 10
Partition 6
Features reduced to (first pass) 2525
Features reduced to (second pass) 10
Partition 7
Features reduced to (first pass) 2525
Features reduced to (second pass) 10
Partition 8
Features reduced to (first pass) 2525
Features reduced to (second pass) 10
Partition 9
Features reduced to (first pass) 2525
Features reduced to (second pass) 10

Accuracy scores [0.90909091 0.90909091 1.          0.8          1.          0.9
0.9          1.          0.9          0.9          ]

Mean 0.9218181818181819 std mean error 0.018816205429811276

Errors [0.09090909 0.09090909 0.          0.2          0.          0.1
0.1          0.          0.1          0.1          ]

Mean 0.07818181818181817 std mean error 0.018816205429811276

Accuracy scores train [1. 1. 1. 1. 1. 1. 1. 1. 1. 1.]

Mean 1.0 std mean error 0.0

```

## 1.4 Question 4

What feature selection method performs best? Can you explain why?

### 1.4.1 Response 4

ANOVA performs marginally better than Random Forest, with overlapping confidence intervals:

- Anova  $0.94 \pm 0.02$ .

- Random Forest  $0.92 \pm 0.02$ .

Decision trees are prone to be more sticky to the data, because of his very hard dependences of initial selection nodes, because it chooses the best result at a given step, but does not ensure that this is the optimal decision of the whole route to the leaf node. Random forests mitigate this problem but the algorithm needs enough features in order to make more effective the randomization.

In fact, if we rise the number of selected features to 20, the numbers are just the opposite:

- Anova  $0.92 \pm 0.02$ .
- Random Forest  $0.94 \pm 0.02$ .