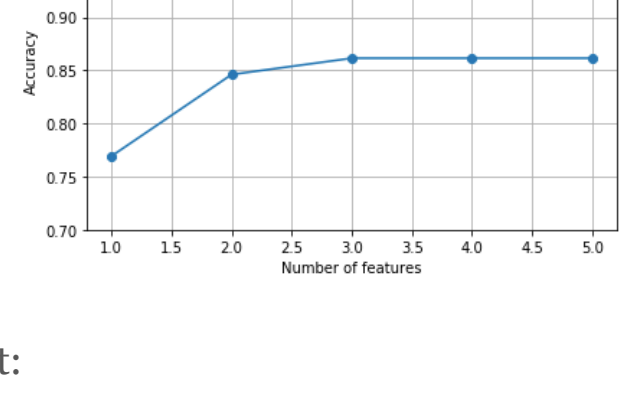


## Sequential Forward Selection – Python Example

July 30, 2020 by [Ajitesh Kumar](#) · [Leave a comment](#)



In this post, you will learn about one of **feature selection** techniques namely **sequential forward selection** with **Python** code example. Refer to my earlier post on [sequential backward selection](#) technique for feature selection. Sequential forward selection algorithm is a part of sequential feature selection algorithms. Some of the following topics will be covered in this

post:

- Introduction to sequential feature selection algorithms
- Sequential forward selection algorithm
- Python example using sequential forward selection

### Table of Contents

1. Introduction to Sequential Feature Selection
2. Sequential Forward Selection & Python Code
3. Python example using sequential forward selection

## Introduction to Sequential Feature Selection

Sequential feature selection algorithms including sequential forward selection algorithm belongs to the family of **greedy search algorithms** which are used to reduce an initial  $d$ -dimensional feature space to a  $k$ -dimensional feature subspace where  $k < d$ . The idea is to select a **subset of features** that is most relevant to the problem, which results in **optimal computation efficiency** while achieving **reduced generalization error** by filtering out irrelevant features (that acts as a noise).

Some of the common techniques used for feature selection includes **regularization techniques (L1 / L2 norm)** and **sequential forward / backward feature selection**. For those algorithms such as K-Nearest Neighbours which do not support regularisation techniques, sequential feature selection is commonly used. Here is a [good read on sequential feature selection](#) technique. Note that the sequential feature selection techniques are based on greedy search algorithms which applies combinatorial methods for feature search.

## Sequential Forward Selection & Python Code

Sequential forward selection algorithm is about execution of the following steps to search the most appropriate features out of N features to fit in K-features subset.

- First and foremost, the best single feature is selected (i.e., using some criterion function) out of all the features.
- Then, pairs of features are formed using one of the remaining features and this best feature, and the best pair is selected.
- Next, triplets of features are formed using one of the remaining features and these two best features, and the best triplet is selected.

- This procedure continues until a predefined number of features (K features) are selected.

Here is the Python code sample for the algorithm. Note some of the following in the code:

- Code block for finding the first / single feature having best score
- Code block for adding the features one by one until k\_features is reached
- Fit method to find the most appropriate indices in the feature list representing the best features
- Transform method to output the best features

- Attributes such as indices\_, subsets\_, scores\_ etc to find the best subsets and related best scores

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.linear_model import LogisticRegression
5 from itertools import combinations
6 from sklearn.base import clone
7 from sklearn.metrics import accuracy_score
8
9 class SequentialForwardSelection():
10     '''
11     Instantiate with Estimator and given number of features
12     '''
13     def __init__(self, estimator, k_features):
14         self.estimator = clone(estimator)
15         self.k_features = k_features
16
17     '''
18     X_train - Training data Pandas dataframe
19     X_test - Test data Pandas dataframe
20     y_train - Training label Pandas dataframe
21     y_test - Test data Pandas dataframe
22     '''
23     def fit(self, X_train, X_test, y_train, y_test):
24         max_indices = tuple(range(X_train.shape[1]))
25         total_features_count = len(max_indices)
26         self.subsets_ = []
27         self.scores_ = []
28         self.indices_ = []
29
30         '''
31         Iterate through the feature space to find the first
32         feature
33         which gives the maximum model performance
34         '''
35         scores = []
36         subsets = []
37         for p in combinations(max_indices, r=1):
38             score = self._calc_score(X_train.values,
39                                     X_test.values, y_train.values,
40                                     y_test.values, p)
41             scores.append(score)
42             subsets.append(p)
43
44         # Find the single feature having best score
45         #
46         best_score_index = np.argmax(scores)
47         self.scores_.append(scores[best_score_index])
48         self.indices_ = list(subsets[best_score_index])
49         self.subsets_.append(self.indices_)
50
51         # Add a feature one by one until k_features is
52         # reached
53         #
54         dim = 1
55         while dim < self.k_features:
56             scores = []
57             subsets = []
58             current_feature = dim
59
60             # Add the remaining features one-by-one from the
61             # remaining feature set
62             # Calculate the score for every feature
63             # combinations
64             for idx in combinations(max_indices, r=current_feature):
65                 while idx < total_features_count:
66                     if idx not in self.indices_:
67                         indices = list(self.indices_)
68                         indices.append(idx)
69                         score = self._calc_score(X_train.values,
70                                                 X_test.values, y_train.values,
71                                                 y_test.values, indices)
72                         scores.append(score)
73                         subsets.append(indices)
74                         idx += 1
75
76             #
77             # Get the index of best score
78             #
79             best_score_index = np.argmax(scores)
80             #
81             # Record the best score
82             #
83             self.scores_.append(scores[best_score_index])
84             #
85             # Get the indices of features which gave best
86             # score
87             #
88             self.indices_ = list(subsets[best_score_index])
89             #
90             # Record the indices of features for best score
91             #
92             self.subsets_.append(self.indices_)
93
94             dim += 1
95
96         self.k_score_ = self.scores_[-1]
97
98     '''
99     Transform training, test data set to the data set
100     having features which gave best score
101     '''
102     def transform(self, X):
103         return X.values[:, self.indices_]
104
105     '''
106     Train models with specific set of features
107     indices - indices of features
108     '''
109     def _calc_score(self, X_train, X_test, y_train, y_test,
110                    indices):
111         self.estimator.fit(X_train[:, indices],
112                             y_train.ravel())
113         y_pred = self.estimator.predict(X_test[:, indices])
114         score = accuracy_score(y_test, y_pred)
115         return score
```

## Python example using sequential forward selection

Here is the code which represents how an instance of LogisticRegression can be passed with training and test data set and the best features are derived. Although regularization technique can be used with LogisticRegression, this is just used for illustration purpose.

```
1 #
2 # Instantiate the estimator - LogisticRegression
3 #
4 lr = LogisticRegression(C=1.0, random_state=1)
5 #
6 # Number of features
7 #
8 k_features = 5
9 #
10 # Instantiate SequentialBackwardSearch
11 #
12 sfs = SequentialForwardSelection(lr, k_features)
13 #
14 # Fit the data to determine the k_features which give the
15 # most optimal model performance
16 #
17 sfs.fit(X_train, X_test, y_train, y_test)
18 #
19 # Transform the training data set to dataset having k_features
20 # giving most optimal model performance
21 #
22 X_train_sfs = sfs.transform(X_train)
23 #
24 # Transform the test data set to dataset having k_features
25 #
26 X_test_sfs = sfs.transform(X_test)
```

Here is a glimpse of the training data used in the above example:

In [147]:	X_train.head()
Out[147]:	
	hsc_p hsc_p degree_p estest_p mba_p gender_M hsc_B_Others hsc_B_Others hsc_S_Commerce hsc_S_Science degree_1_Others degree_2
202	0.249055 -0.289516 -0.114685 -0.858762 -0.439652 1 0 0 0 0 1 0
40	1.480356 2.039075 1.429653 0.411259 2.316050 0 0 1 1 1 0 0
107	1.315537 2.053830 2.228399 0.502767 1.880407 1 1 1 1 1 0 0
20	-0.475983 -0.115935 -0.114685 -1.675998 -1.020767 1 1 1 1 1 0 0
206	-2.307078 -0.112118 -0.842010 1.737400 -1.587751 1 0 0 0 0 1 0

Fig 1. Data used for sequential forward selection algorithm

Here is the plot representing the model performance vs number of features which got derived from executing **sequential forward selection algorithm**.

```
1 k_features = [len(k) for k in sfs.subsets_]
2 plt.plot(k_features, sfs.scores_, marker='o')
3 plt.ylim([0.7, 1.02])
4 plt.ylabel('Accuracy')
5 plt.xlabel('Number of features')
6 plt.grid()
7 plt.tight_layout()
8 plt.show()
```

Here is how the plot would look like:

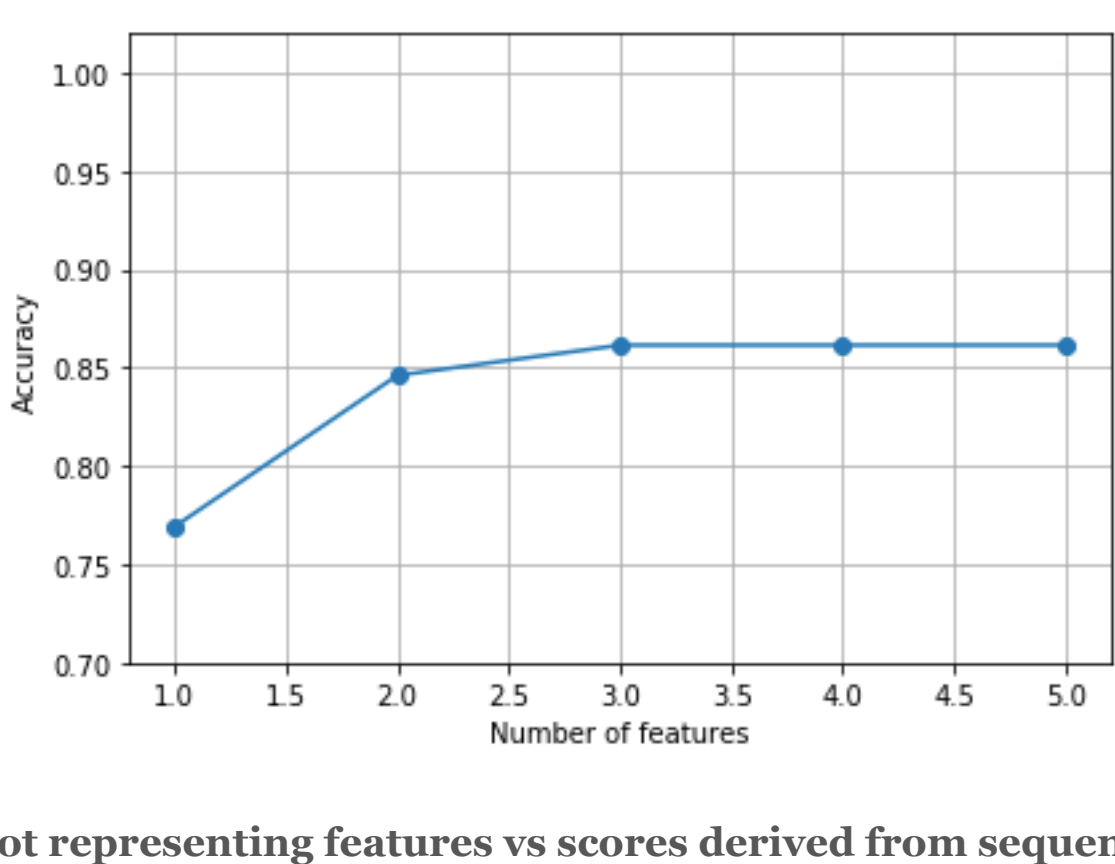



Fig 2. Plot representing features vs scores derived from sequential forward selection algorithm

Author



Recent Posts



Ajitesh Kumar

I have been recently working in the area of Data Science and Machine Learning / Deep Learning. In addition, I am also passionate about various different technologies including programming languages such as Java/JEE, Javascript, Python, R, Julia etc and technologies such as Blockchain, mobile computing, cloud-native technologies, application security, cloud computing platforms, big data etc. I would love to connect with you on [Linkedin](#) and [Twitter](#).

Follow me



Posted in [Data Science](#), [Machine Learning](#), [Python](#). Tagged with [Data Science](#), [machine learning](#), [python](#).

[← Sequential Backward Feature Selection – Python Example](#) [Feature Importance using Random Forest Classifier – Python →](#)

## Leave a Reply

Your email address will not be published. Required fields are marked \*

Comment

Name \*

Email \*

Website

× **1** = one 

About Us




Vitalflux.com is dedicated to help software engineers & data scientists get technology news, practice tests, tutorials in order to reskill / acquire newer skills from time-to-time.

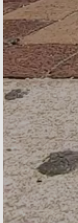
Thank you for visiting our site today. We welcome all your suggestions in order to make our website better. Please feel free to share your thoughts.

BOSTON ARCHITECTURAL COLLEGE


Transferring Made Easy

LEARN MORE






uBreakiFix - Low Price P...

uBreakiFix

Boston 10AM-6PM



### Recent Posts

- [Spend Analytics Use Cases: AI & Data Science](#)
- [Procurement: Key Advanced Analytics Use Cases](#)
- [What are Sequence Models: Types & Examples](#)
- [First Principles Thinking Explained with Examples](#)
- [Matplotlib Pyplot Plot Python Example](#)


### Tag Cloud

ai (70) Angular (50) angularjs (104) api (16) Application Security (22) artificial intelligence (20) AWS (23) bigdata (11) big data (40) blockchain (58) cloud (11) cloud computing (11) data analytics (21) Data Science (328) datascience (34) Deep Learning (37) docker (26) freshers (14) google (14) google glass (11) hyperledger (18) Interview questions (73) Java (92) javascript (103) Kubernetes (19) machine learning (324) mongodb (16) news (13) nlp (13) nosql (17) python (92) QA (12) quantum computing (13) reactjs (15) r programming (11) sklearn (30) Software Quality (11) spring framework (16) statistics (18) testing (16) tools (11) tutorials (14) UI (13) Unit Testing (18) web (16)

ASTRID LUN HOME

Seamlessly sync your *online store* inventory no matter where you sell it.

BUILD YOUR ONLINE STORE



## Everything to Sell Anything

Keep your inventory in sync with an easy-to-use interface that scales with you.

Squarespace

[Open >](#)

### Recent Comments

- [Ajitesh Kumar](#) on [Contract Management Use Cases for Machine Learning](#)
- [Prateek Madhesiya](#) on [Contract Management Use Cases for Machine Learning](#)
- [URL](#) on [What are Features in Machine Learning?](#)
- [Hold-out Method for Training Machine Learning Models - Data Analytics](#) on [What is Machine Learning? Concepts & Examples](#)
- [Hold-out Method for Training Machine Learning Models - Data Analytics](#) on [K-Fold Cross Validation – Python Example](#)