

## Feature Selection

David Gerbing  
The School of Business  
Portland State University  
gerbing@pdx.edu

### Table of Contents

- [1\\_Preliminaries](#)
- [2\\_Data](#)
- [3\\_Feature Selection](#)
  - [3.1\\_Manual Selection](#)
  - [3.2\\_Automated Feature Selection](#)
    - [3.2.1\\_Automated Univariate Feature Selection](#)
    - [3.2.2\\_Automated Multivariate Feature Selection](#)
- [4\\_Postscript](#)

Feature selection is not always necessary for building machine learning models, but it is typically a helpful process to pursue. The goal is to reduce the number of predictor variables (features) in the model, to keep predictive accuracy at or about the same level, but with a much simpler model, with fewer predictor variables.

Two reasons to pursue feature selection:

1. Data costs money. The fewer the predictors, the less data needs to be collected.
2. Understanding the underlying relationships between predictors and target variable, which indirectly often leads to the construction of better models.

## Preliminaries

```
from datetime import datetime as dt
now = dt.now()
print ("Analysis on", now.strftime("%Y-%m-%d"), "at", now.strftime("%H:%M"))

Analysis on 2021-07-12 at 14:28

import os
os.getcwd()
'/content'

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

## Data

```
#d = pd.read_csv('data/Boston.csv')
d = pd.read_csv('http://web.pdx.edu/~gerbing/data/Boston.csv')
```

d.shape

(506, 15)

d.head()

	Unnamed: 0	crim	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	black	lstat	medv
0	1	0.00632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296	15.3	396.90	4.98	24.0
1	2	0.02731	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242	17.8	396.90	9.14	21.6
2	3	0.02729	0.0	7.07	0	0.469	7.185	61.1	4.9671	2	242	17.8	392.83	4.03	34.7
3	4	0.03237	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63	2.94	33.4
4	5	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222	18.7	396.90	5.33	36.2

Do not need the first column, so drop.

```
d = d.drop(['Unnamed: 0'], axis="columns")
d.head()

      crim  zn  indus  chas  nox  rm  age  dis  rad  tax  ptratio  black  lstat  medv
0  0.00632  18.0  2.31  0  0.538  6.575  65.2  4.0900  1  296  15.3  396.90  4.98  24.0
1  0.02731  0.0  7.07  0  0.469  6.421  78.9  4.9671  2  242  17.8  396.90  9.14  21.6
2  0.02729  0.0  7.07  0  0.469  7.185  61.1  4.9671  2  242  17.8  392.83  4.03  34.7
3  0.03237  0.0  2.18  0  0.458  6.998  45.8  6.0622  3  222  18.7  394.63  2.94  33.4
4  0.06905  0.0  2.18  0  0.458  7.147  54.2  6.0622  3  222  18.7  396.90  5.33  36.2
```

Store the features, the predictor variables, in data structure X. Store the target variable in data structure y. To run multiple regression with all possible predictor variables, one possibility defines X as the entire data frame with medv dropped, as in

X = d.drop(['medv'], axis="columns")

Alternatively, use the procedure below that manually defines a vector of the predictor variables (features) names, and then define X as the subset of d that contains just these variables.

```
y = d['medv']
pred_vars = ['crim', 'zn', 'indus', 'chas', 'nox', 'rm', 'age', 'dis', 'rad', 'tax', 'ptratio', 'black', 'lstat']
X = d[pred_vars]
```

Not necessary, but see how many features in the model, and observe the data type of the X and y data structures. The function len() provides the length of a vector, that is, the number of elements of a vector.

```
n_pred = len(pred_vars)
print("Number of predictor variables:", n_pred)
```

Number of predictor variables: 13

X and y are created as two different pandas data types: X is a data frame, y is a one-dimensional array called a series. A data frame can have a single column, but, somewhat confusingly (in my opinion), subsetting a data frame down to a single variable is no longer a data frame.

## Feature Selection

Features, the predictor variables, should be ...

- *relevant*: Predictors each correlate with the target
- *unique*: Predictors do not correlate much with each other

The problem of *collinearity* is the problem of correlated predictor variables, the features. Too much correlation and redundancy make estimating the slope coefficients difficult, though it does not harm predictive accuracy per se. Generally, improve model fit by adding new information in the form of a new predictor variable to the model to the extent that the new predictor is relevant and unique.

Do, however, be aware of the problem of *data leakage*. When testing the model on data previously unseen by the model, all aspects of that data must have been unseen, just as in a real-world forecasting scenario. Otherwise, the data is said to leak from training to testing data. Making decisions regarding the model based on *all* the data then by definition includes both training and testing data. Best to make decisions regarding model estimation only from the training data.

## Manual Selection

Base selection of the predictor variables on satisfying the two criterion: relevance and uniqueness. The goal here is to produce a single output, a table, that displays numerical indices for both criterion.

*Uniqueness*: Besides the correlation coefficient of two predictor variables, a more general indicator of collinearity is the *variance inflation factor* or *VIF*. The *VIF* assesses the linear redundancy of one predictor variable not just with one other predictor variable, but all the other predictor variables.

*Relevance*: Compute the correlation of each predictor with the target.

```
print("X is a:", type(X))
print("X.values is a:", type(X.values))
```

X is a: <class 'pandas.core.frame.DataFrame'>
X.values is a: <class 'numpy.ndarray'>

Use the statsmodels function *variance\_inflation\_factor()* to compute the variance inflation factor for each predictor. The VIF's are a property only of the X's, so the target y is not part of this analysis. The *variance\_inflation\_factor()* function does not compute all the VIF's, but only one at a time. Create a data frame named vif, then fill each row of the data frame with the corresponding name of the predictor variable and its corresponding variance inflation factor.

To systematically calculate and retrieve the VIF's, one for each feature, traverse through the variables in X one at a time with a programming structure known as a *for* loop, from the first X variable through the last X variable, where X.shape[1] is the number of rows of the data frame.

Because the loop cannot traverse through the original data frame, transfer the X data frame to a more primitive data structure, a *numpy* structure of a numeric matrix, obtained with the *values* method.

1. To begin, create an empty data frame with any valid name. Here we use *vif*. Then define a variable called *Predictor* in the data frame, filled with the names of the columns of the X data structure using the *columns* method.

2. Then create a variable called *VIF*, the variance inflation factor for each predictor variable. Loop through the data matrix (not data frame) with the *values* method for each predictor variable.

3. Calculate the correlation of each predictor (feature) with the target and store in the variable called *Relevance*. Store in the data series *cr*, then loop through *cr* for each variable to copy the value to the new *Relevance* variable.

4. Finally, display the contents of the created *vif* data frame by listing its name as the last line of code in the cell. (If we wish to display information before the last line, then need the *print()* function.)

```
from statsmodels.stats.outliers_influence import variance_inflation_factor
```

```
vif = pd.DataFrame()
```

```
vif['Predictor'] = X.columns
```

```
vif['VIF'] = [variance_inflation_factor(X.values, i)
```

```
for i in range(X.shape[1])]
```

```
vif
```

```
/usr/local/lib/python3.7/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning:
```

```
import pandas.util.testing as tm
```

```
Predictor          VIF
0      crim  2.100373
1        zn  2.844013
2    indus  14.485758
3      chas  1.152952
4      nox  73.894947
5        rm  77.948283
6      age  21.386850
7      dis  14.699652
8      rad  15.167725
9      tax  61.227274
10    ptratio  85.029547
11    black  20.104943
12    lstat  0.333
```

There is much collinearity in the data, consistent with the correlation matrix that shows many feature correlations far from 0. Many features could be deleted to yield a more parsimonious model that would be just as effective if not more so. Although *rm* has one of the highest VIF's, it is also strongly related to the target as shown by the regression coefficients analysis and has one of the highest correlations with the target. When deleted, lower the VIF on the more relevant feature because perhaps a relevant feature is correlated with other, less relevant features that, when deleted, lower the VIF on the more relevant feature.

## Automated Feature Selection

The pure machine approach seeks to automate selection. This approach makes the most sense when there are many, tens if not hundreds, of features. Otherwise, it is better to analyze, variance inflation factor, the features, and then select the best few.

The *sklearn* module *feature\_selection* has the *SelectKBest* that selects the specified number of features with the highest correlations with the target. Specify the number of retained features with the *k* parameter.

Here the logical array we name *selected* indicates which of the X feature data structure are to be retained.

```
print("X is a:", type(X))
print("X.values is a:", type(X.values))
```

X is a: <class 'pandas.core.frame.DataFrame'>
X.values is a: <class 'numpy.ndarray'>

Base selection of the predictor variables on satisfying the two criterion: relevance and uniqueness. The goal here is to produce a single output, a table, that displays numerical indices for both criterion.

*Uniqueness*: Besides the correlation coefficient of two predictor variables, a more general indicator of collinearity is the *variance inflation factor* or *VIF*. The *VIF* assesses the linear redundancy of one predictor variable not just with one other predictor variable, but all the other predictor variables.

*Relevance*: Compute the correlation of each predictor with the target.

```
print("X is a:", type(X))
print("X.values is a:", type(X.values))
```

X is a: <class 'pandas.core.frame.DataFrame'>
X.values is a: <class 'numpy.ndarray'>

Store the features, the predictor variables, in data structure X. Store the target variable in data structure y. To run multiple regression with all possible predictor variables, one possibility defines X as the entire data frame with medv dropped, as in

```
X = d.drop(['medv'], axis="columns")
```

Alternatively, use the procedure below that manually defines a vector of the predictor variables (features) names, and then define X as the subset of d that contains just these variables.

```
y = d['medv']
pred_vars = ['crim', 'zn', 'indus', 'chas', 'nox', 'rm', 'age', 'dis', 'rad', 'tax', 'ptratio', 'black', 'lstat']
X = d[pred_vars]
```

Not necessary, but see how many features in the model, and observe the data type of the X and y data structures. The function len() provides the length of a vector, that is, the number of elements of a vector.

```
n_pred = len(pred_vars)
print("Number of predictor variables:", n_pred)
```

Number of predictor variables: 13

X and y are created as two different pandas data types: X is a data frame, y is a one-dimensional array called a series. A data frame can have a single column, but, somewhat confusingly (in my opinion), subsetting a data frame down to a single variable is no longer a data frame.

## Automated Feature Selection

Features, the predictor variables, should be ...

- *relevant*: Predictors each correlate with the target
- *unique*: Predictors do not correlate much with each other

The problem of *collinearity* is the problem of correlated predictor variables, the features. Too much correlation and redundancy make estimating the slope coefficients difficult, though it does not harm predictive accuracy per se. Generally, improve model fit by adding new information in the form of a new predictor variable to the model to the extent that the new predictor is relevant and unique.

Do, however, be aware of the problem of *data leakage*. When testing the model on data previously unseen by the model, all aspects of that data must have been unseen, just as in a real-world forecasting scenario. Otherwise, the data is said to leak from training to testing data. Making decisions regarding the model based on *all* the data then by definition includes both training and testing data. Best to make decisions regarding model estimation only from the training data.

## Manual Selection

Base selection of the predictor variables on satisfying the two criterion: relevance and uniqueness. The goal here is to produce a single output, a table, that displays numerical indices for both criterion.

*Uniqueness*: Besides the correlation coefficient of two predictor variables, a more general indicator of collinearity is the *variance inflation factor* or *VIF*. The *VIF* assesses the linear redundancy of one predictor variable not just with one other predictor variable, but all the other predictor variables.

*Relevance*: Compute the correlation of each predictor with the target.

```
print("X is a:", type(X))
print("X.values is a:", type(X.values))
```

X is a: <class 'pandas.core.frame.DataFrame'>
X.values is a: <class 'numpy.ndarray'>

Store the features, the predictor variables, in data structure X. Store the target variable in data structure y. To run multiple regression with all possible predictor variables, one possibility defines X as the entire data frame with medv dropped, as in

```
X = d.drop(['medv'], axis="columns")
```

Alternatively, use the procedure below that manually defines a vector of the predictor variables (features) names, and then define X as the subset of d that contains just these variables.

```
y = d['medv']
pred_vars = ['crim', 'zn', 'indus', 'chas', 'nox', 'rm', 'age', 'dis', 'rad', 'tax', 'ptratio', 'black', 'lstat']
X = d[pred_vars]
```

Not necessary, but see how many features in the model, and observe the data type of the X and y data structures. The function len() provides the length of a vector, that is, the number of elements of a vector.

```
n_pred = len(pred_vars)
print("Number of predictor variables:", n_pred)
```

Number of predictor variables: 13

X and y are created as two different pandas data types: X is a data frame, y is a one-dimensional array called a series. A data frame can have a single column, but, somewhat confusingly (in my opinion), subsetting a data frame down to a single variable is no longer a data frame.