# Variance Inflation Factor (VIF) Explained

08 March 2017  /  Python  /  Data Wrangling

Colinearity is the state where two variables are highly correlated and contain similiar information about the variance within a given dataset. To detect colinearity among variables, simply create a correlation matrix and find variables with large absolute values. In R use the `corr` (http://www.sthda.com/english/wiki/correlation-matrix-a-quick-start-guide-to-analyze-format-and-visualize-a-correlation-matrix-using-r-software) function and in python this can by accomplished by using numpy's `corrcoef` (https://docs.scipy.org/doc/numpy/reference/generated/numpy.corrcoef.html) function.

Multicolinearity (https://en.wikipedia.org/wiki/Multicollinearity) on the other hand is more troublesome to detect because it emerges when three or more variables, which are highly correlated, are included within a model. To make matters worst multicolinearity can emerge even when isolated pairs of variables are not colinear.

A common R function used for testing regression assumptions and specifically multicolinearity is "VIF()" and unlike many statistical concepts, its formula is straightforward:

$$ V.I.F. = 1 / (1 - R^2). $$

The Variance Inflation Factor (VIF) is a measure of colinearity among predictor variables within a multiple regression. It is calculated by taking the the ratio of the variance of all a given model's betas divide by the variane of a single beta if it were fit alone.

## Steps for Implementing VIF

1. Run a multiple regression.
2. Calculate the VIF factors.
3. Inspect the factors for each predictor variable, if the VIF is between 5-10, multicolinearity is likely present and you should consider dropping the variable.

```python
#Imports
import pandas as pd
import numpy as np
from patsy import dmatrices
import statsmodels.api as sm
from statsmodels.stats.outliers_influence import variance_inflation_factor

df = pd.read_csv('loan.csv')
df.dropna()
df = df._get_numeric_data() #drop non-numeric cols

df.head()
```

| | id | me mbe r_id | loa n_a mnt | fund ed_a mnt | funded _amnt _inv | int _r at e | inst allm ent | ann ual_ inc | d ti | deli nq_ 2yrs | ... | tota l_ba l_il | il _u ti l | open _rv_ 12m | open _rv_ 24m | max _bal _bc | al l_ut il | total_r ev_hi _lim | i n q _ fi | tota l_c u_tl | inq_l ast_1 2m |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 10 77 50 1 | 129 659 9 | 500 0.0 | 500 0.0 | 4975.0 | 1 0. 65 | 162. 87 | 240 00.0 | 2 7. 6 5 | 0.0 | ... | NaN | N a N | NaN | NaN | NaN | N a N | NaN | N a N | Na N | NaN |
| 1 | 10 77 43 0 | 131 416 7 | 250 0.0 | 250 0.0 | 2500.0 | 1 5. 27 | 59.8 3 | 300 00.0 | 1. 0 0 | 0.0 | ... | NaN | N a N | NaN | NaN | NaN | N a N | NaN | N a N | Na N | NaN |
| 2 | 10 77 17 5 | 131 352 4 | 240 0.0 | 240 0.0 | 2400.0 | 1 5. 96 | 84.3 3 | 122 52.0 | 8. 7 2 | 0.0 | ... | NaN | N a N | NaN | NaN | NaN | N a N | NaN | N a N | Na N | NaN |
| 3 | 10 76 86 3 | 127 717 8 | 100 00. 0 | 1000 0.0 | 10000. 0 | 1 3. 49 | 339. 31 | 492 00.0 | 2 0. 0 0 | 0.0 | ... | NaN | N a N | NaN | NaN | NaN | N a N | NaN | N a N | Na N | NaN |
| 4 | 10 75 35 8 | 131 174 8 | 300 0.0 | 300 0.0 | 3000.0 | 1 2. 69 | 67.7 9 | 800 00.0 | 1 7. 9 4 | 0.0 | ... | NaN | N a N | NaN | NaN | NaN | N a N | NaN | N a N | Na N | NaN |

5 rows × 51 columns

```python
df = df[['annual_inc','loan_amnt', 'funded_amnt','annual_inc','dti']].dropna() #subset the dat
```

## Step 1: Run a multiple regression

```python
%%capture
#gather features
features = "+".join(df.columns - ["annual_inc"])

# get y and X dataframes based on this regression:
y, X = dmatrices('annual_inc ~' + features, df, return_type='dataframe')
```

## Step 2: Calculate VIF Factors

```
# For each X, calculate VIF and save in dataframe
vif = pd.DataFrame()
vif["VIF Factor"] = [variance_inflation_factor(X.values, i) for i in range(X.shape[1])]
vif["features"] = X.columns
```

## Step 3: Inspect VIF Factors

```
vif.round(1)
```

|   | VIF Factor | features |
|---|---|---|
| **0** | 5.1 | Intercept |
| **1** | 1.0 | dti |
| **2** | 678.4 | funded_amnt |
| **3** | 678.4 | loan_amnt |

As expected, the total funded amount for the loan and the amount of the loan have a high variance inflation factor because they "explain" the same variance within this dataset. We would need to discard one of these variables before moving on to model building or risk building a model with high multicolinearity.

### Find an error or bug? Have a suggestion?

Everything on this site is avaliable on GitHub. Head on over and submit an issue. (https://github.com/etav/) You can also message me directly on Twitter (https://twitter.com/etav).

This project contains 16 pages and is available on GitHub (https://github.com/etav/).
Copyright © Ernest Tavares III, 2017.