

# Handling Multi-collinearity

## 1. Remove Highly Correlated Predictors

- **Identify Correlated Variables:** Use a correlation matrix to identify pairs of variables that are highly correlated (e.g., with a correlation coefficient greater than 0.8 or 0.9).
- **Remove One of the Variables:** If two variables are highly correlated, consider removing one of them from the model, particularly if they convey similar information.

```
import pandas as pd

# Assuming df is your DataFrame with the predictors

corr_matrix = df.corr().abs()

# Select upper triangle of correlation matrix

upper = corr_matrix.where(pd.np.triu(pd.np.ones(corr_matrix.shape),
k=1).astype(bool))

# Find features with correlation greater than 0.8

to_drop = [column for column in upper.columns if any(upper[column] > 0.8)]

# Drop features

df_reduced = df.drop(to_drop, axis=1)
```

## 2. Combine Predictors

- **Create Composite Variables:** If two or more variables are highly correlated and represent similar concepts, you can combine them into a single variable. For example, you might average the scores or create an index.
- **Principal Component Analysis (PCA):** PCA is a dimensionality reduction technique that transforms correlated variables into a set of uncorrelated components. These components can then be used as predictors in your regression model.

```
from sklearn.decomposition import PCA

# Assuming X is your matrix of predictors

pca = PCA(n_components=0.95) # Keep 95% of the variance

X_reduced = pca.fit_transform(X)
```

### 3. Regularization Techniques

- **Ridge Regression:** Ridge regression adds a penalty to the regression coefficients proportional to their size (L2 regularization). This penalty can reduce the impact of multicollinearity by shrinking the coefficients of correlated predictors.
- **Lasso Regression:** Lasso regression adds a penalty equal to the absolute value of the coefficients (L1 regularization). It can shrink some coefficients to zero, effectively selecting a simpler model that may eliminate multicollinearity.

```
from sklearn.linear_model import Ridge, Lasso

# Ridge Regression

ridge = Ridge(alpha=1.0) # alpha is the regularization strength

ridge.fit(X_train, y_train)

# Lasso Regression

lasso = Lasso(alpha=0.1) # alpha is the regularization strength

lasso.fit(X_train, y_train)
```

### 4. Centering the Variables

- **Mean-Centering:** Subtract the mean of each predictor from its values to center the data. Centering can reduce multicollinearity, particularly when interactions between variables are included in the model.
- **Standardization:** Scaling the variables so that they have a mean of zero and a standard deviation of one can sometimes reduce the effects of multicollinearity.

```
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
```

```
X_scaled = scaler.fit_transform(X)
```

## 5. Variance Inflation Factor (VIF) Threshold

- **VIF Thresholding:** Calculate the VIF for each predictor. If the VIF for a variable exceeds a certain threshold (commonly 5 or 10), consider removing it from the model. This approach directly targets variables that contribute most to multicollinearity.

```
import pandas as pd

from statsmodels.stats.outliers_influence import variance_inflation_factor

import statsmodels.api as sm

# Add a constant
X = sm.add_constant(X)

# Calculate VIF for each predictor
vif_data = pd.DataFrame()
vif_data["feature"] = X.columns
vif_data["VIF"] = [variance_inflation_factor(X.values, i) for i in range(X.shape[1])]

# Drop features with VIF > 10
features_to_keep = vif_data[vif_data["VIF"] <= 10]["feature"]
X_reduced = X[features_to_keep]
```

## 6. Increase Sample Size

- **Collect More Data:** Increasing the sample size can sometimes reduce the impact of multicollinearity. With more data, the variance of the estimates can decrease, making multicollinearity less problematic.

## 7. Stepwise Regression

- **Forward or Backward Selection:** Use stepwise regression techniques to automatically include or exclude variables based on certain criteria (like AIC

or BIC). This can help in selecting a subset of predictors that have less multicollinearity.

```
from mlxtend.feature_selection import SequentialFeatureSelector as SFS

from sklearn.linear_model import LinearRegression

# Forward Selection

sfs = SFS(LinearRegression(),

         k_features='best',

         forward=True,

         floating=False,

         scoring='r2',

         cv=5)

sfs = sfs.fit(X, y)

X_selected = X[list(sfs.k_feature_names_)]
```

## 8. Partial Least Squares Regression (PLS)

- **PLS Regression:** This method reduces the predictors to a smaller set of uncorrelated components, similar to PCA, but also considers the dependent variable in the dimensionality reduction process. It is especially useful when dealing with multicollinearity.

```
from sklearn.cross_decomposition import PLSRegression

# Assuming you want to keep 2 components

pls = PLSRegression(n_components=2)

X_pls = pls.fit_transform(X, y)
```

## 9. Domain Knowledge

- **Leverage Subject Matter Expertise:** Use your understanding of the subject area to decide which variables are most important. Sometimes, even with high multicollinearity, certain variables are essential for theoretical or practical reasons.

```
from sklearn.ensemble import RandomForestRegressor

model = RandomForestRegressor()

model.fit(X, y)

importances = model.feature_importances_

# Select features with importance above a certain threshold

important_features = [X.columns[i] for i in range(len(importances)) if importances[i] > 0.05]

X_important = X[important_features]
```

## 10. Removing Interaction Terms

- **Review Interaction Terms:** Interaction terms can introduce multicollinearity, particularly if the main effects are also included in the model. Consider removing unnecessary interaction terms or using other methods to handle them.