

A decorative graphic in the top-left corner consisting of two overlapping parallelograms. The front one is blue and the back one is a light green. They are positioned diagonally, with the blue one partially covering the green one.

Predicting Shirt Size, with Data!

Presented by:
Zack Godwin

What's the big deal?

Buying clothes online has never been easier than it is today. Businesses that provide free return shipping are faced with financial loss when the customer doesn't get the right fit.

Providing the customer with all the tools necessary to get the right fit the first time, could save money and provide a unique service to the customer.

Data Science + Apparel = \$\$\$



The Research Question:

- Can we predict chest circumference, a common measure of shirt size, using other body measurements?

The Dataset

- The dataset contained body measurement data and some demographic-related features for people age < 1 - 20 years old.
- 3,900 observations
- 122 variables

	WEIGHT	STATURE	ERECT SITTING HEIGHT	MAXIMUM HIP BREADTH (SEATED)	BUTTOCK- KNEE LENGTH	KNEE HEIGHT	HEAD CIRCUMFERENCE	HEAD BREADTH	SHOULDER BREADTH
12	499.0	1578.0	838.0	330.0	518.0	488.0	532.0	154.0	396.0
13	558.0	1618.0	826.0	325.0	568.0	511.0	549.0	152.0	402.0
15	400.0	1468.0	766.0	297.0	509.0	452.0	512.0	145.0	353.0

https://raw.githubusercontent.com/Padam-0/cluster_t-shirt_sizing/master/data.csv

Data Cleaning

```
In [2]: 1 # Cleaning the features
2 df = pd.read_csv('https://raw.githubusercontent.com/Padam-0/cluster_t-shirt_sizing/master/data.csv')
3
4 remove_cols = []
5
6 # I want to get rid of the columns with more than 2000 Nulls,
7 # because incomplete data won't be useful.
8 for i in df.columns:
9     if 3900 - df.loc[:,i].astype(bool).sum() > 2000:
10         remove_cols.append(i)
11
12 df2 = df.drop(remove_cols, axis = 'columns')
13
14 # I don't need these categorical variables
15 demographic_attributes = ['AGE IN YEARS', 'LOCATION',
16                           'BIRTH DATE', 'MEASUREMENT DATE', 'MEASUREMENT SET TP',
17                           'MEASURER NUMBER', 'COMPUTER NUMBER', 'RACE', 'GRADE LEVEL',
18                           'HANDEDNESS', 'NUMBER OF BROTHERS', 'NUMBER OF SISTERS', 'TWIN',
19                           'BIRTH ORDER', 'MOTHERS OCCUPATION', 'FATHERS OCCUPATION',
20                           'MOTHERS EDUCATION', 'FATHERS EDUCATION', 'YEARS IN COMMUNITY',
21                           'ANTHROPOMETER NO', 'CALIPER NO', 'GIRTH NO', 'PERSON #']
22
23 df2 = df2.drop(demographic_attributes, axis = 'columns')
24
25 # Recode Gender Variable
26 def binary_gender(gender):
27     if gender == 1:
28         return 0
29     else:
30         return 1
31
32 df2 = df2.replace(0, np.nan)
33 df2 = df2.dropna()
34
35 df2['gender_recode'] = df2['SEX'].apply(binary_gender)
36
37 df2 = df2.drop(['SEX'], axis = 'columns')
38
39 df2.head()
```

Creating Shirt Size Categories

```
In [4]: 1 df_men = df2[(df2.gender_recode == 0) & \
2              (df2['CHEST CIRCUMFERENCE'] >= 640)]
```

```
In [5]: 1 df_women = df2[(df2.gender_recode == 1) & \
2              (df2['CHEST CIRCUMFERENCE'] >= 640)]
```

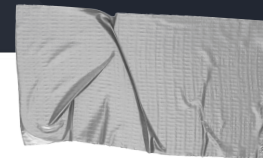
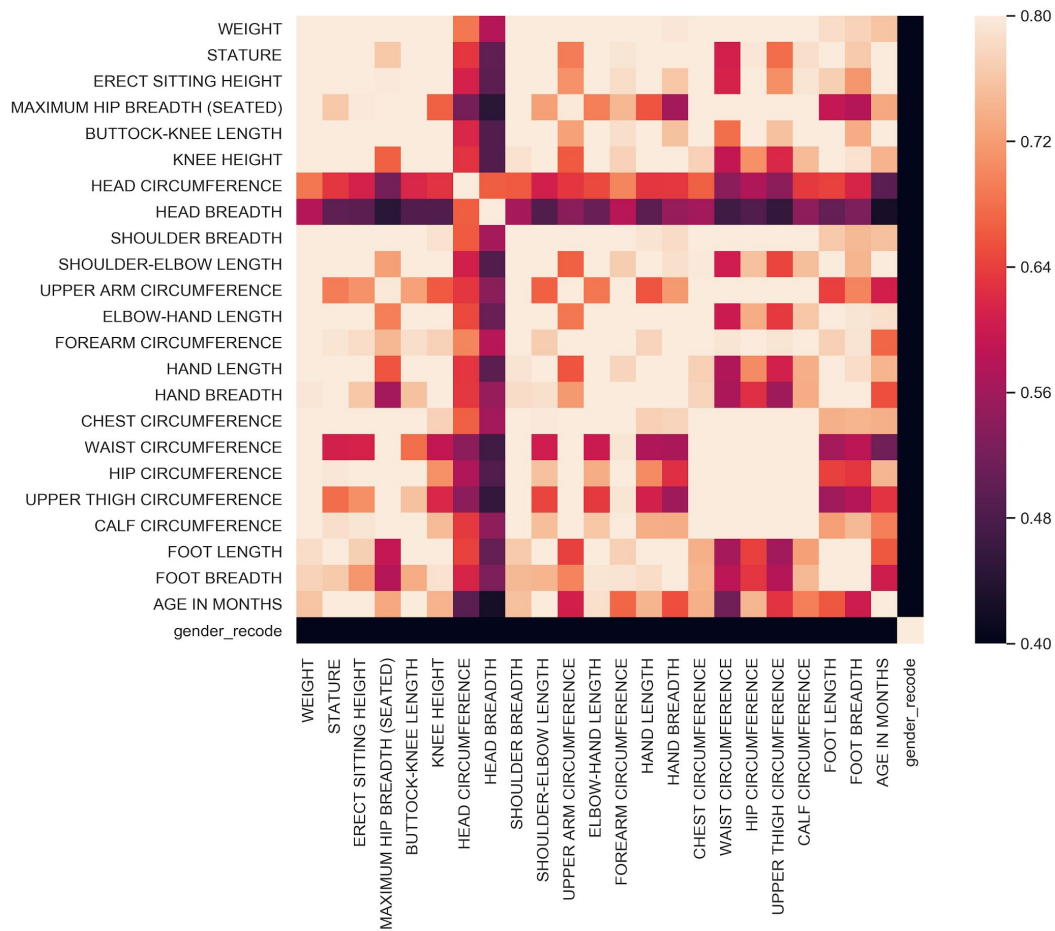
```
In [41]: 1 #Adding size category.
2 def cat_size_men(size):
3     if 660 > size and size >= 640:
4         return "YM X Small"
5     elif 690 > size and size >= 660:
6         return "YM Small"
7     elif 750 > size and size >= 690:
8         return "YM Medium"
9     elif 820 > size and size >= 750:
10        return "YM Large"
11    elif 890 > size and size >= 820:
12        return "YM Extra Large"
13    elif 960 > size and size >= 880:
14        return "AM Small"
15    elif 1040 > size and size >= 960:
16        return "AM Medium"
17    elif 1120 > size and size >= 1040:
18        return "AM Large"
19    elif 1240 > size and size >= 1120:
20        return "AM Extra Large"
21    elif 1360 > size and size >= 1240:
22        return "AM XX Large"
23    elif 1480 > size and size >= 1360:
24        return "AM XXX Large"
25    elif 1600 > size and size >= 1480:
26        return "AM XXXX Large"
27    else:
28        return "Other"
29    pd.options.mode.chained_assignment = None # default='warn'
30    df_men['cat_size'] = df_men.loc[:, 'CHEST CIRCUMFERENCE'].apply(cat_size_men)
31    df_men.head()
```

```
In [42]: 1 def cat_size_women(size):
2     if 660 > size and size >= 640:
3         return "YW X Small"
4     elif 690 > size and size >= 660:
5         return "YW Small"
6     elif 750 > size and size >= 690:
7         return "YW Medium"
8     elif 780 > size and size >= 750:
9         return "YW Large"
10    elif 830 > size and size >= 780:
11        return "AW Extra Small"
12    elif 900 > size and size >= 830:
13        return "AW Small"
14    elif 970 > size and size >= 900:
15        return "AW Medium"
16    elif 1040 > size and size >= 970:
17        return "AW Large"
18    elif 1140 > size and size >= 1040:
19        return "AW X Large"
20    elif 1240 > size and size >= 1140:
21        return "AW XX Large"
22    else:
23        return "Other"
24
25    df_women['cat_size'] = df_women.loc[:, 'CHEST CIRCUMFERENCE'].apply(cat_size_women)
26    df_women.head()
```

```
In [64]: 1 frames = [df_women, df_men]
2
3 df_full = pd.concat(frames)
4 df_full.cat_size.unique()
```

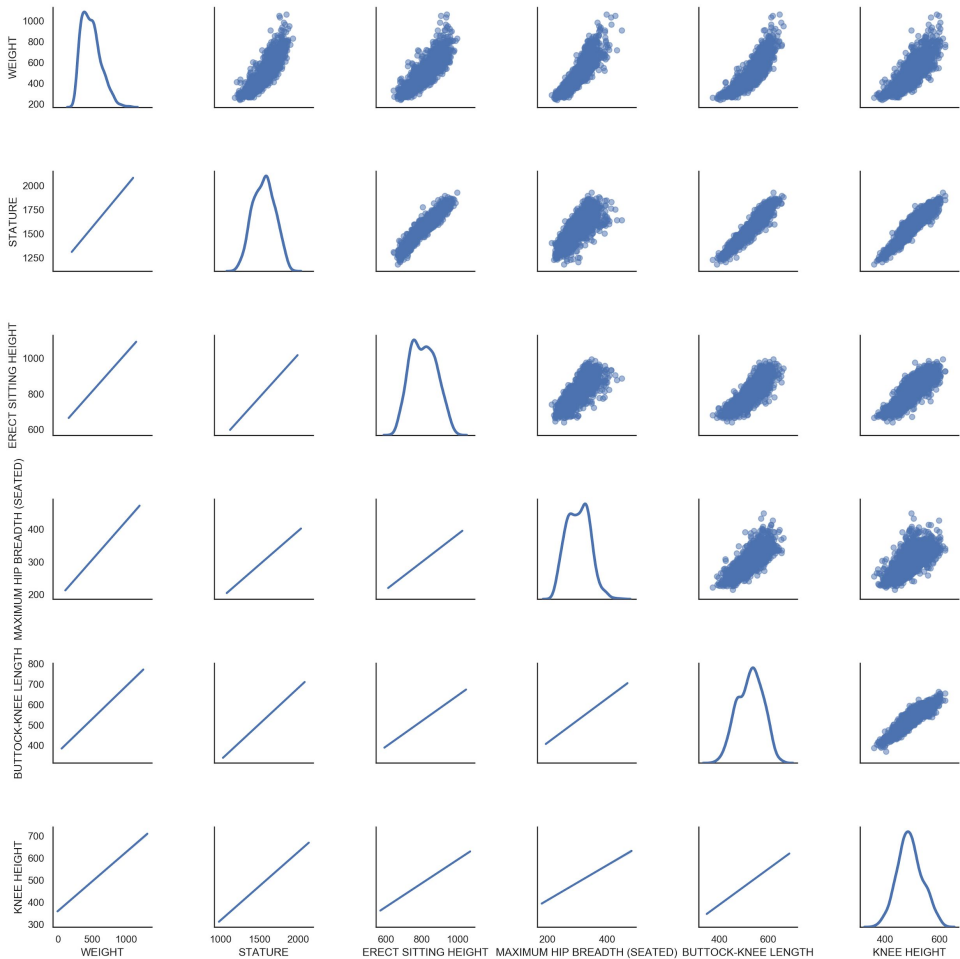
```
Out[64]: array(['YW Large', 'AW Small', 'YW Medium', 'AW Extra Small', 'YW Small',
               'AW Large', 'AW Medium', 'YW X Small', 'AW X Large', 'YM X Small',
               'YM Medium', 'YM Small', 'YM Large', 'YM Extra Large', 'AM Medium',
               'AM Small', 'AM Large', 'AM Extra Large'], dtype=object)
```

Feature correlations:



- We have very strongly correlated variables, and few weakly correlated.
- We may assume linearity in body measurement variables, but to confirm we will examine a few.

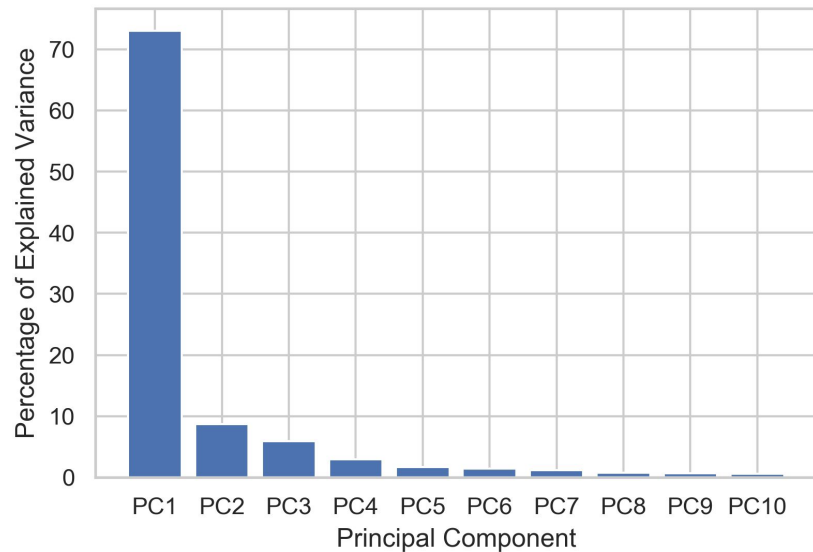
Feature exploration:



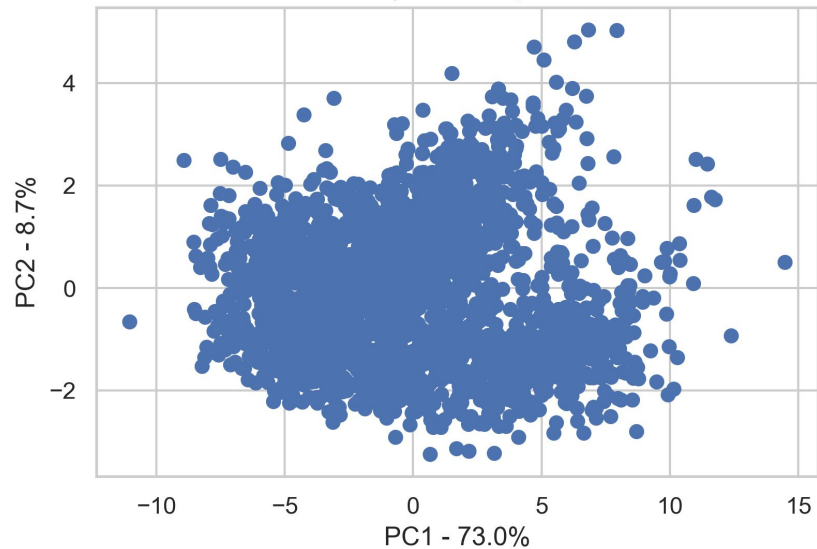
- Looking at just a few features, we see several linear relationships.
- To simplify the model, I first performed Principal Component Analysis.

PCA

Scree Plot



My PCA Graph

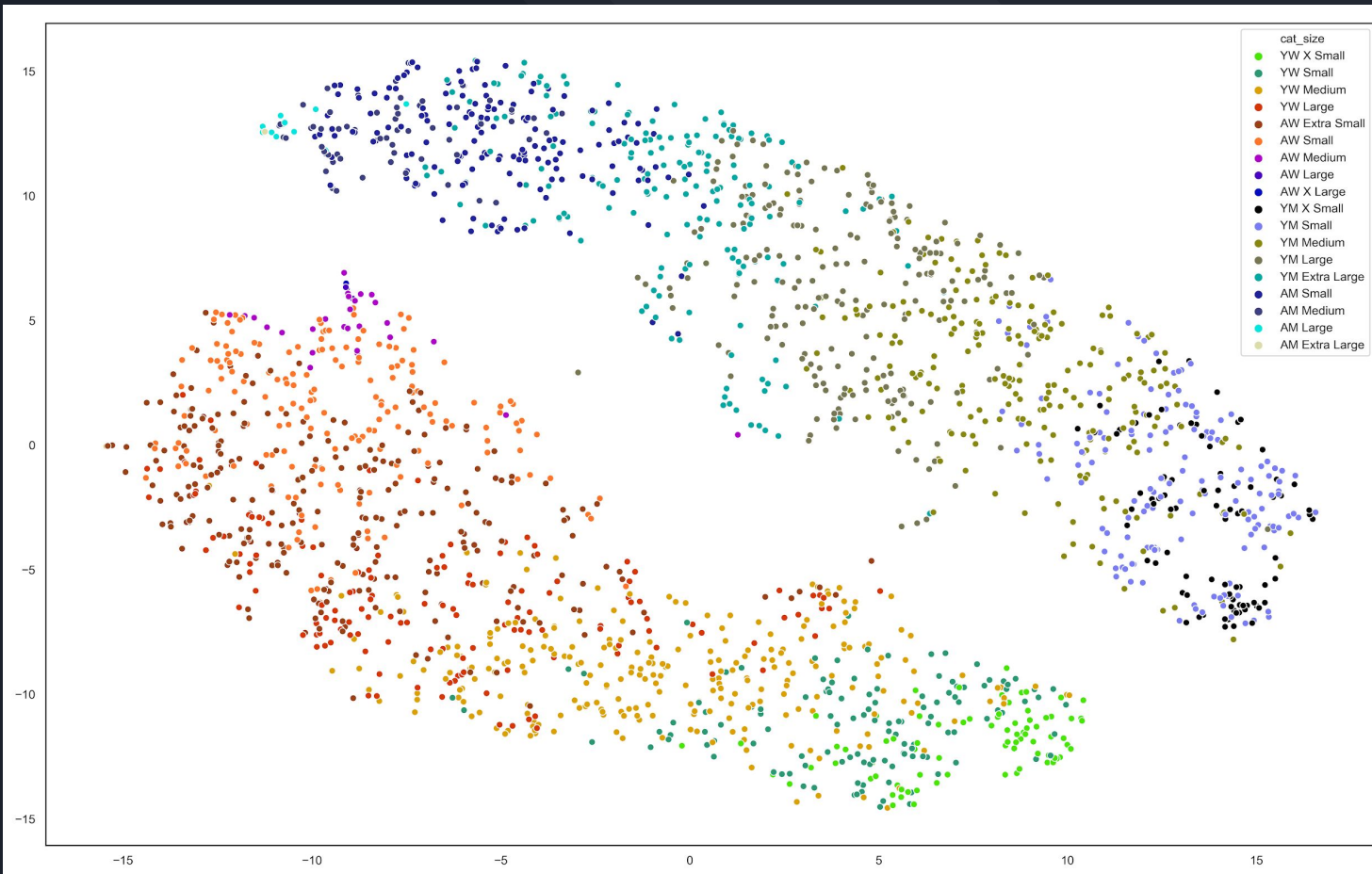


- PCA reduced 24 features down to 10 principal components.
- We'll only use PC1 and PC2.

TSNE

To visualize how the observations cluster relative to the shirt size category I created I used a t-distributed stochastic neighbor embedding algorithm.

We see clustering by gender, and clustering by size in a gradient pattern from smallest to largest!





Model Selection

OLS models trained with 80/20 split.

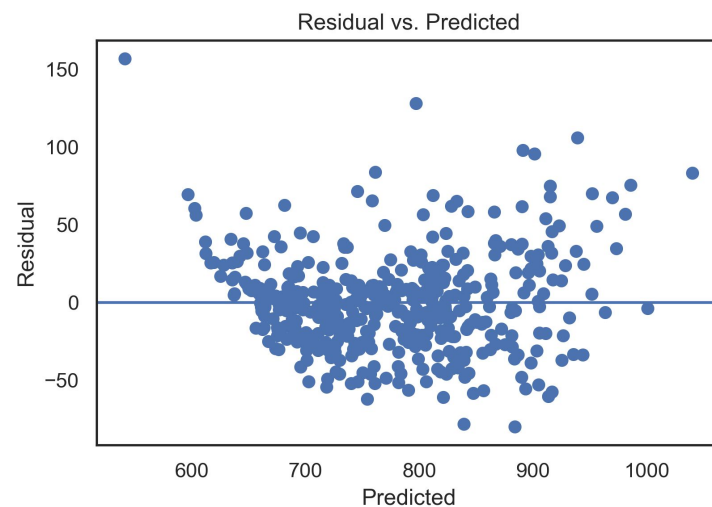
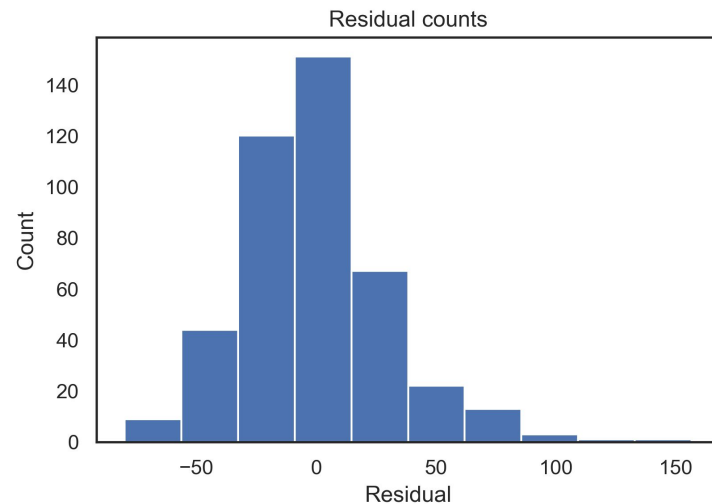
Target = Chest Circumference

Success = High R-Squared value and a high mean CV with low standard deviation.

Linear Regression on Principal Components

R-squared:
0.890

Mean (SD) Cross Validation Score:
0.94 (+/- 0.01)





Regression Model Scoreboard

Regression Model	R-Squared	Mean CV (+/- SD)
PCA Linear OLS	0.890	0.94 (+/- 0.01)



Is PCA the best method?

- I used sklearn's variance threshold to remove features with low variance.

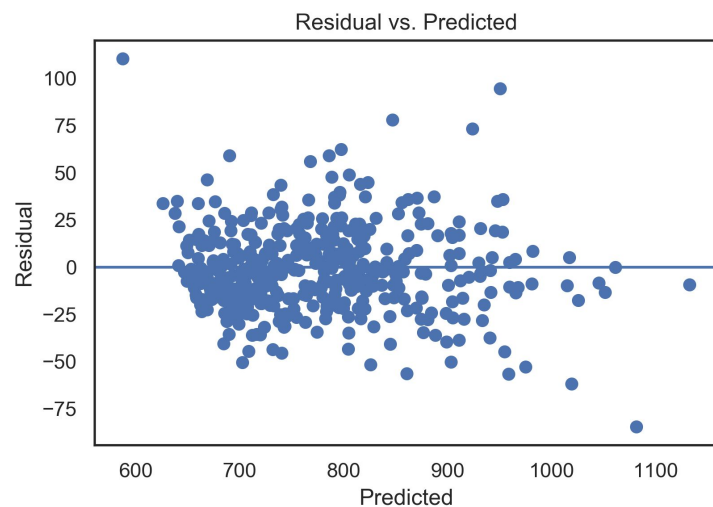
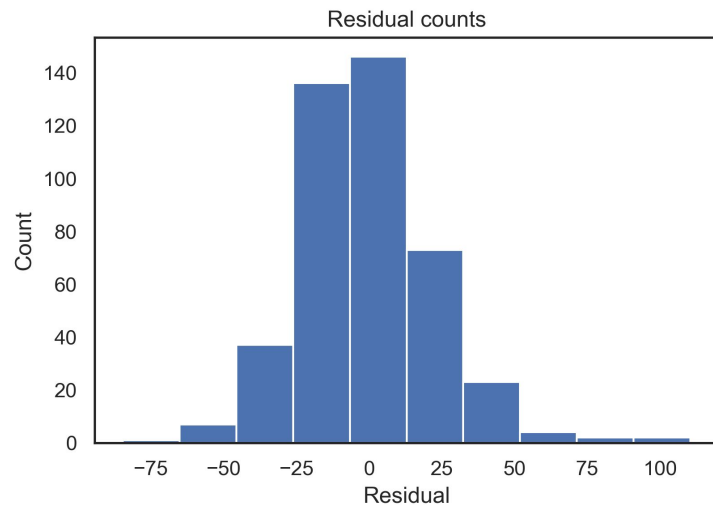
```
In [28]: 1 from sklearn.feature_selection import VarianceThreshold
          2
          3
          4 def variance_threshold_selector(data, threshold):
          5     selector = VarianceThreshold(threshold)
          6     selector.fit(data)
          7     return data[data.columns[selector.get_support(indices=True)]]
```

```
In [29]: 1 variance_threshold = variance_threshold_selector(
          2     X_scaled_df, threshold=(.8 * (1 - .8)))
```

Linear Regression on Variance Threshold

R-squared:
0.945

Mean (SD) Cross Validation Score:
0.94 (+/- 0.01)





Regression Model Scoreboard

Regression Model	R-Squared	Mean CV (+/- SD)
PCA Linear OLS	0.890	0.94 (+/- 0.01)
Variance Threshold OLS	0.945	0.94 (+/- 0.01)



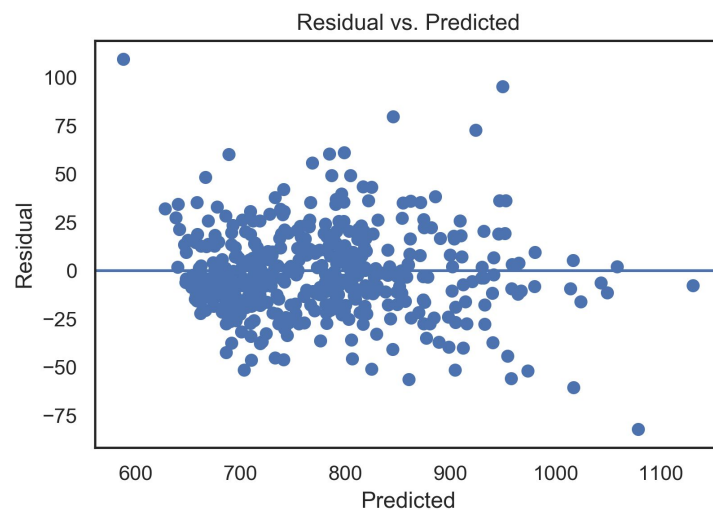
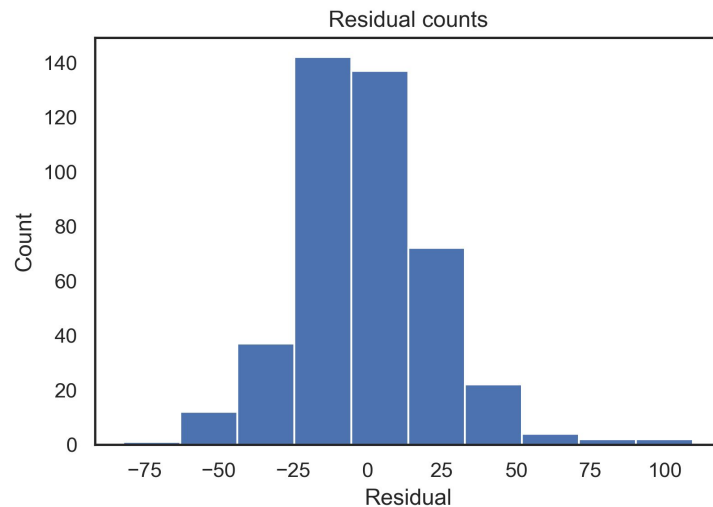
Is OLS the best model?

- I compared the results of the OLS models to Ridge and Lasso regression models.
- These models trained on the whole data, and did their own feature management.

Lasso Regression

R-squared:
0.945

Mean (SD) Cross Validation Score:
0.93 (+/- 0.05)

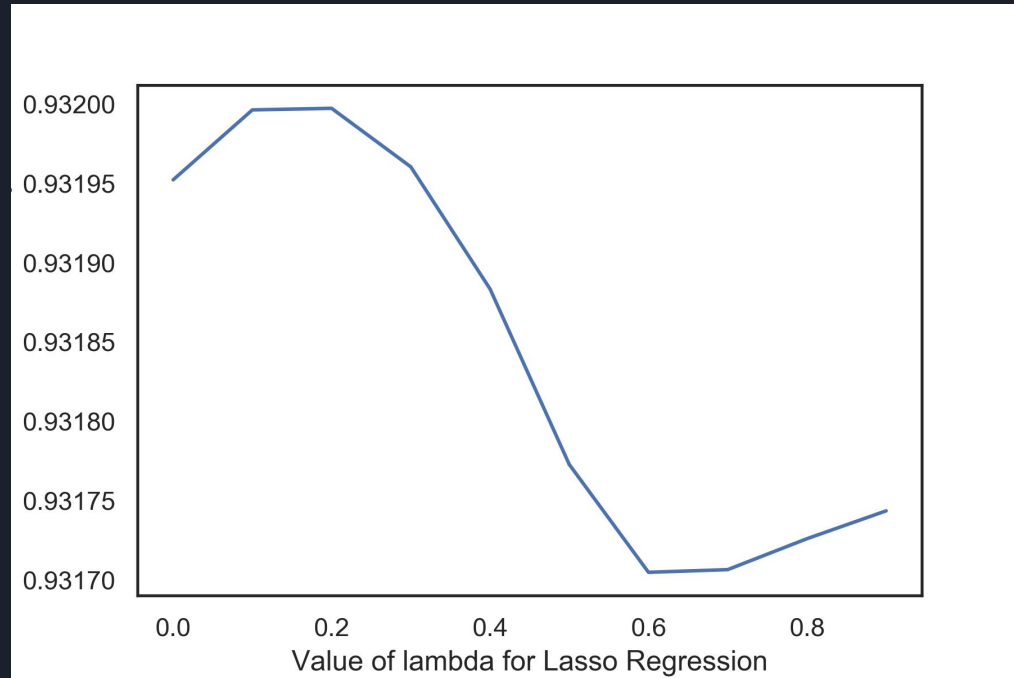




Regression Model Scoreboard

Regression Model	R-Squared	Mean CV (+/- SD)
PCA Linear OLS	0.890	0.94 (+/- 0.01)
Variance Threshold OLS	0.945	0.94 (+/- 0.01)
Lasso Regression	0.889	0.93 (+/- 0.05)

Finding Lambda...

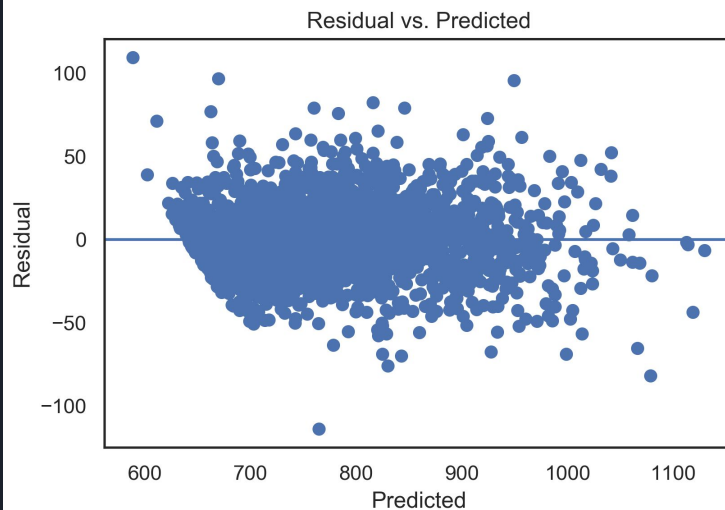
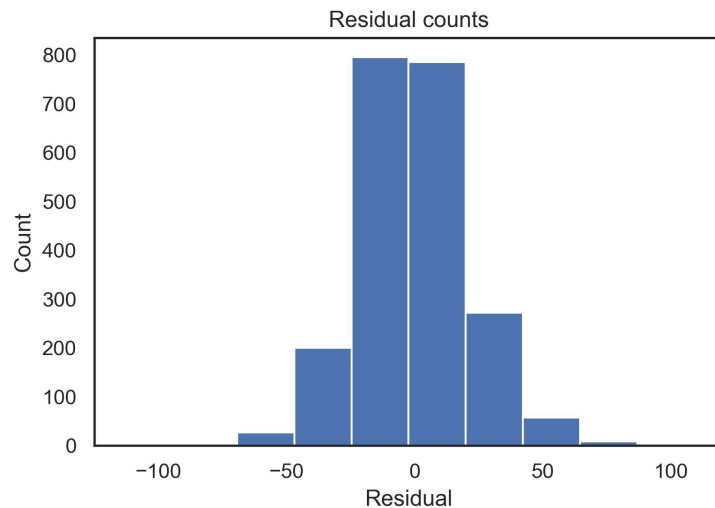


... 0.2 should work just fine.

Tuned lambda Lasso Regression

R-squared:
0.944

Mean (SD) Cross Validation Score:
0.94 (+/- 0.01)

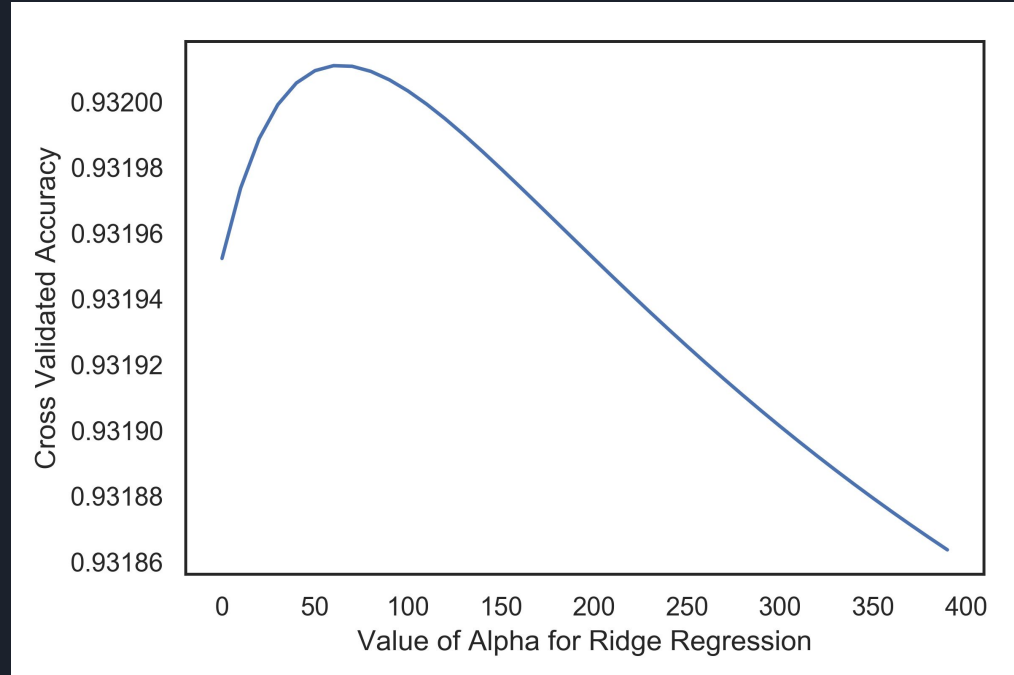




Regression Model Scoreboard

Regression Model	R-Squared	Mean CV (+/- SD)
PCA Linear OLS	0.890	0.94 (+/- 0.01)
Variance Threshold OLS	0.945	0.94 (+/- 0.01)
Lasso Regression	0.889	0.93 (+/- 0.05)
Tuned Lambda Lasso	0.944	0.94 (+/- 0.01)

Finding Alpha...

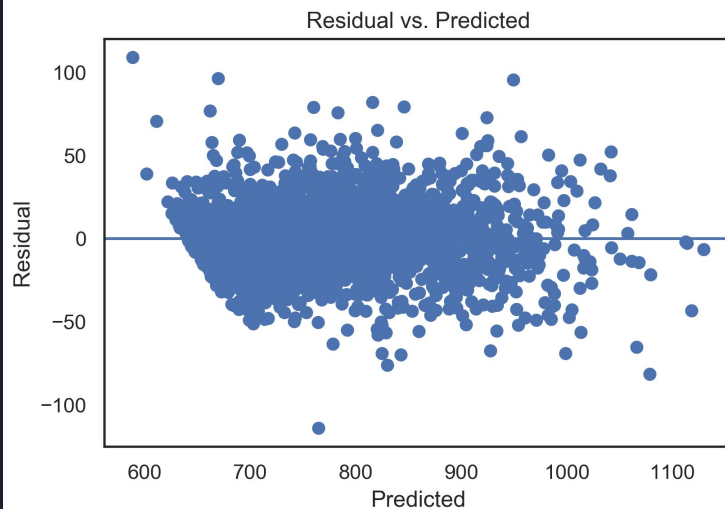
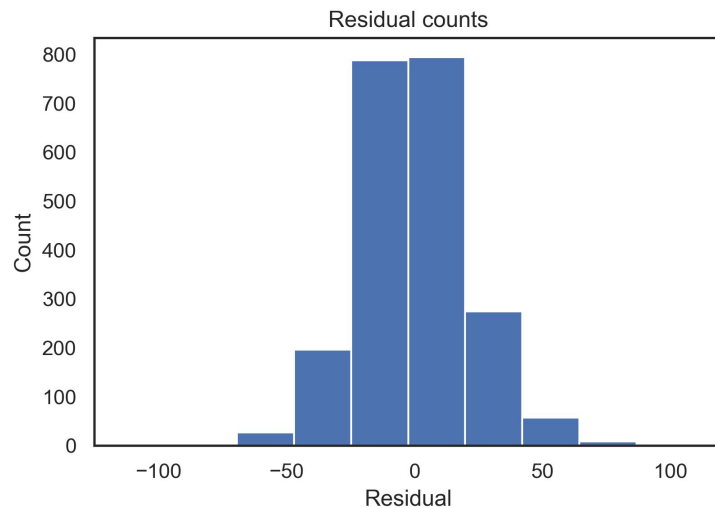


... 60 should work just fine.

Tuned Ridge Regression

R-squared:
0.944

Mean (SD) Cross Validation Score:
0.94 (+/- 0.01)





Regression Model Scoreboard

Regression Model	R-Squared	Mean CV (+/- SD)
PCA Linear OLS	0.890	0.94 (+/- 0.01)
Variance Threshold OLS	0.945	0.94 (+/- 0.01)
Lasso Regression	0.889	0.93 (+/- 0.05)
Tuned Lambda Lasso	0.944	0.94 (+/- 0.01)
Tuned Ridge Regression	0.944	0.94 (+/- 0.01)



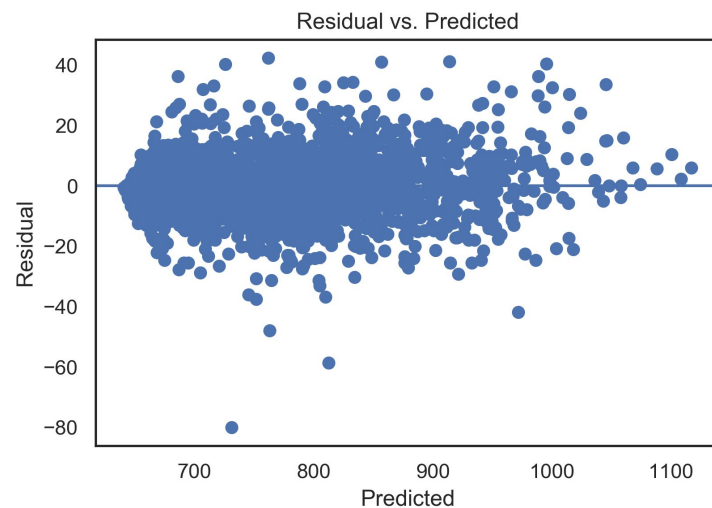
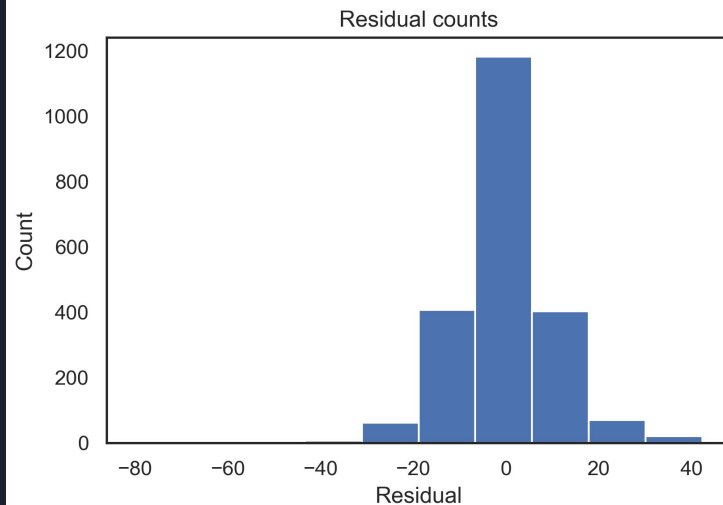
Random Forest Regression

R-squared:

0.987

Mean (SD) Cross Validation Score:

0.93 (+/- 0.01)





Regression Model Scoreboard

Regression Model	R-Squared	Mean CV (+/- SD)
PCA Linear OLS	0.890	0.94 (+/- 0.01)
Variance Threshold OLS	0.945	0.94 (+/- 0.01)
Lasso Regression	0.889	0.93 (+/- 0.05)
Tuned Lambda Lasso	0.944	0.94 (+/- 0.01)
Tuned Ridge Regression	0.944	0.94 (+/- 0.01)
Random Forest	0.979	0.93 (+/- 0.01)



Tuning Hyperparameters: Random Forest

In [42]:

```
1 from sklearn.model_selection import GridSearchCV
2
3 # Create the parameter grid based on the results of random search
4 param_grid = {
5     'bootstrap': [True, False],
6     'max_depth': [80, 90, 100, 110],
7     'min_samples_leaf': [3, 4, 5],
8     'min_samples_split': [8, 10, 12],
9     'n_estimators': [100, 200, 300, 1000]
10 }
```

In [43]:

```
1 rfr = ensemble.RandomForestRegressor()
2 rfrfit = rfr.fit(X, y)
3
4 # Instantiate the grid search model
5 grid_search = GridSearchCV(estimator=rfrfit, param_grid=param_grid,
6                             cv=3, n_jobs=-1, verbose=2)
7 grid_search.fit(X, y)
8 grid_search.best_params_
```

Fitting 3 folds for each of 288 candidates, totalling 864 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done 25 tasks      | elapsed:    16.0s
[Parallel(n_jobs=-1)]: Done 146 tasks     | elapsed:    1.3min
[Parallel(n_jobs=-1)]: Done 349 tasks     | elapsed:    3.1min
[Parallel(n_jobs=-1)]: Done 632 tasks     | elapsed:    6.6min
[Parallel(n_jobs=-1)]: Done 864 out of 864 | elapsed:    9.9min finished
```

Out[43]:

```
{'bootstrap': True,
 'max_depth': 110,
 'min_samples_leaf': 4,
 'min_samples_split': 10,
 'n_estimators': 200}
```



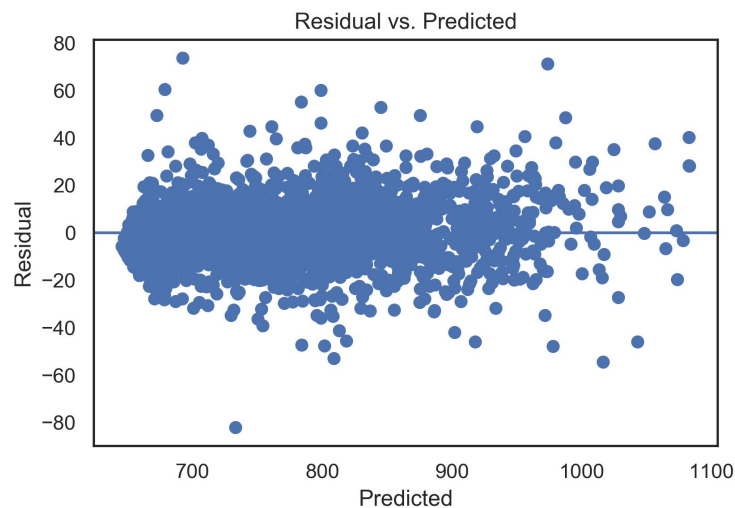
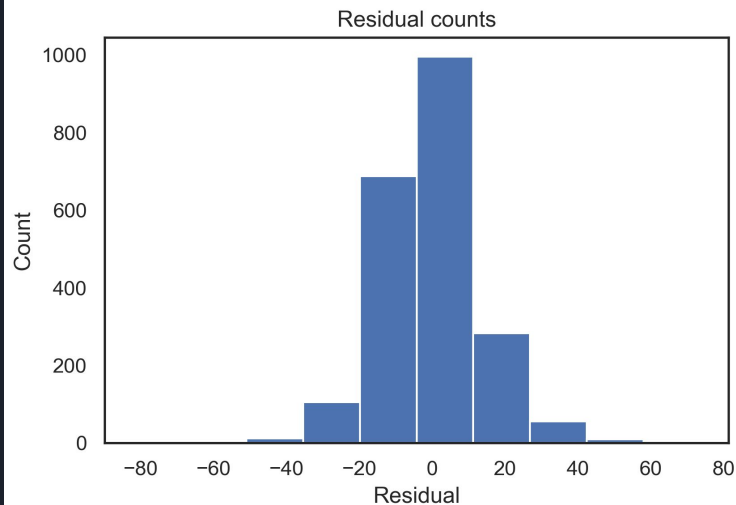
Tuned Hyperparameters Random Forest Regression

R-squared:

0.978

Mean (SD) Cross Validation Score:

0.94 (+/- 0.01)





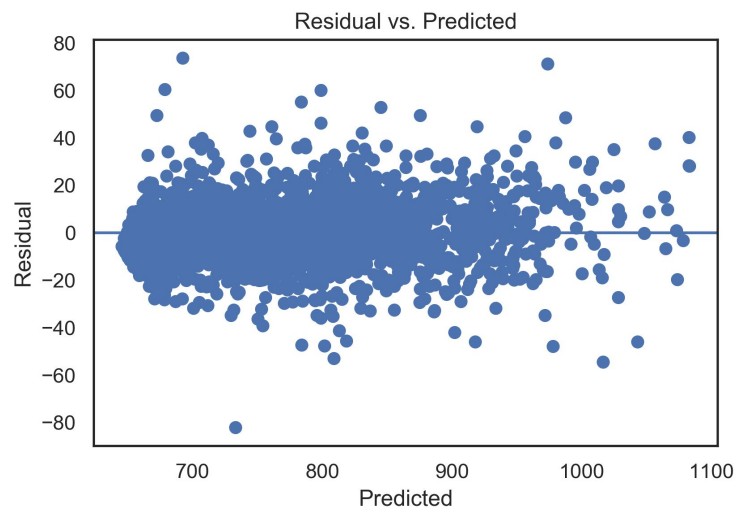
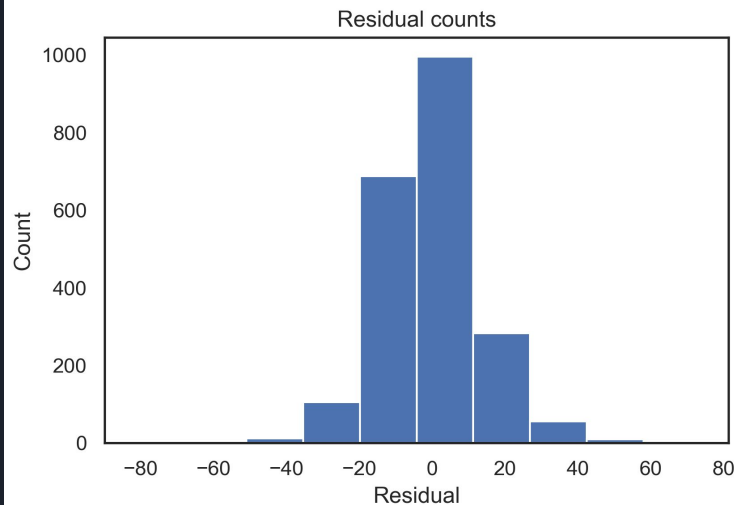
Regression Model Scoreboard

Regression Model	R-Squared	Mean CV (+/- SD)
PCA Linear OLS	0.890	0.94 (+/- 0.01)
Variance Threshold OLS	0.945	0.94 (+/- 0.01)
Lasso Regression	0.889	0.93 (+/- 0.05)
Tuned Lambda Lasso	0.944	0.94 (+/- 0.01)
Tuned Ridge Regression	0.944	0.94 (+/- 0.01)
Random Forest	0.979	0.93 (+/- 0.01)
Tuned Random Forest	0.978	0.94 (+/- 0.01)

Conclusion

→ **The Tuned-Hyperparameters Random Forest Regression** best met my success definition for this data:

- ◆ **High $R^2 = 0.978$**
- ◆ **High Mean CV of 0.94**
- ◆ **Low CV (SD) of (+/- 0.01)**





Questions?