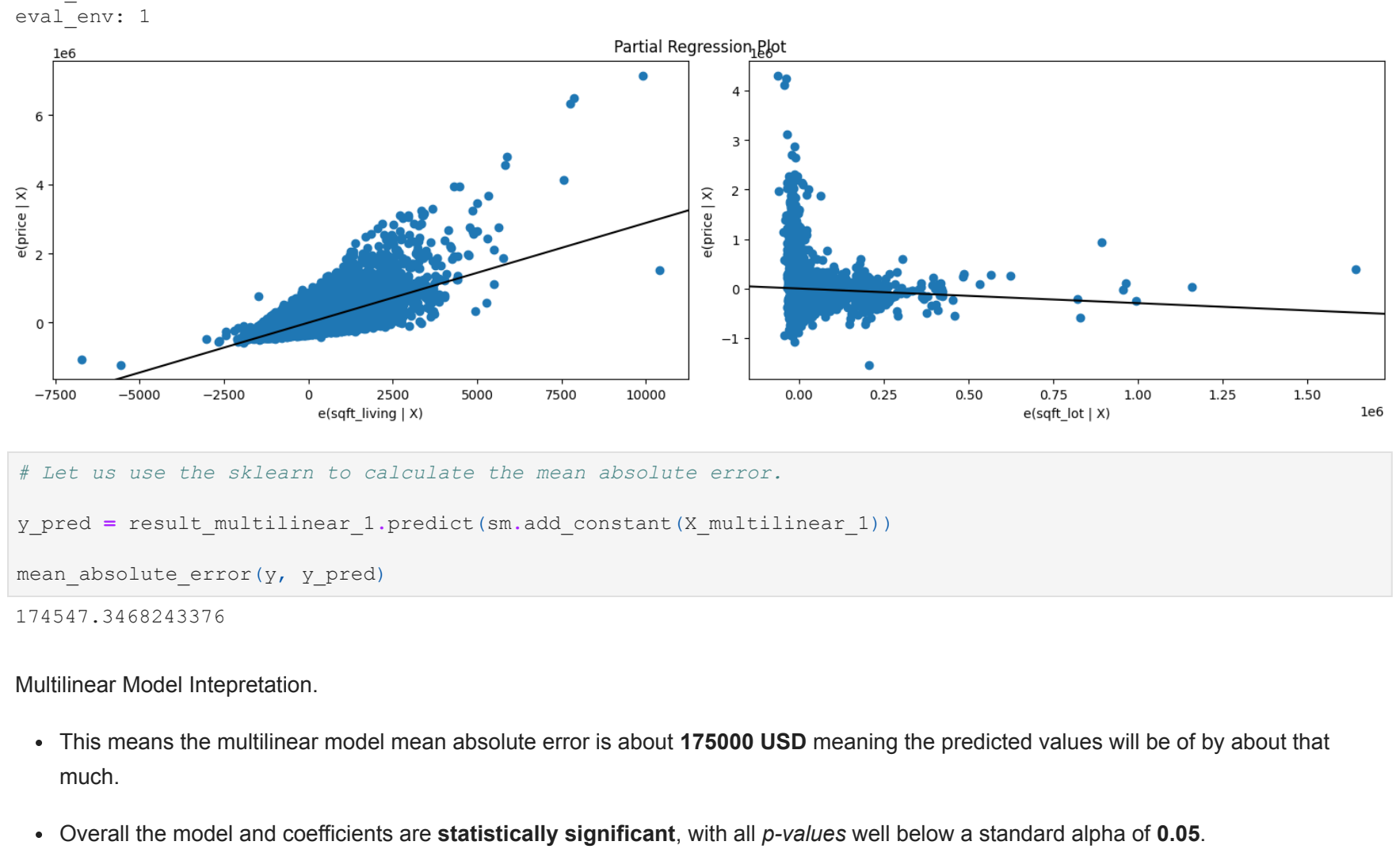






```
eval_env: 1
eval_env: 1
eval_env: 1
```



```
In [32]: # Let us use the sklearn to calculate the mean absolute error.

y_pred = result_multilinear_1.predict(sm.add_constant(X_multilinear_1))

mean_absolute_error(y, y_pred)

Out[32]: 174547.3468243376
```

- Multilinear Model Interpretation.
- This means the multilinear model mean absolute error is about **175000 USD** meaning the predicted values will be off by about that much.
  - Overall the model and coefficients are **statistically significant**, with all *p*-values well below a standard alpha of **0.05**.
  - The Adj. R-squared value is about **0.500** i.e. **50%** of the variance in the target variable can be explained by **sqft\_living** and **sqft\_lot**. It is almost similar to the baseline model.
  - The intercept is about **-6.529e+04** i.e. **-65,290 USD**, meaning that if one wanted a house with a square foot of living space of **0** and **sqft\_lot** of 0 they would expect a price of about **-55160 USD** but this can be ignored as a house with 0 sqft of living room cant exist outside the mathematical realm.
  - The **sqft\_living** coefficient is about **290**, meaning that for each additional 1 USD on price of the house we have an additional square feet of living of about **290 sqft**.
  - The **sqft\_above** coefficient is about **-0.3**, meaning that for each additional 1 USD on price of the house we have an additional square footage of the house lot of about **-0.3 sqft** which is impractical in the realm of mathematics so they can be ignored, as no house can exist with a lot of 0.

## Multilinear Model 1

$$y = ([282 - 291] \text{Sqft\_living} + [(-0.3 - 0.05)] \text{sqft\_lot} + [(-65,400 - 45,000)] \text{USD} + 175,000 \text{ USD}$$

## Multilinear Regression with a categorical column.

let us randomly pick a category to add onto the multilinear regression.  
Houses with desired aesthetic views have a price premium over similar houses without such views so from this

```
In [33]: # checking out the value counts found in the category column 'view'

X['view'].value_counts()

Out[33]:
NONE      14241
AVERAGE    688
GOOD       348
FAIR       245
EXCELLENT  240
Name: view, dtype: int64

In [34]: # Definisg a new variable to hold the independent variables to use for the new model.

X_multilinear_2 = X[['sqft_living', 'sqft_lot', 'view']]

In [35]: # before i can use the new column views, i have to change it into a binary system.

X_multilinear_2 = pd.get_dummies(X_multilinear_2, columns=['view'])

In [36]: X_multilinear_2.head()

Out[36]:
```

	sqft_living	sqft_lot	view_AVERAGE	view_EXCELLENT	view_FAIR	view_GOOD	view_NONE
6414100192	2570	7242	0	0	0	0	1
2487200875	1960	5000	0	0	0	0	1
1954400510	1680	8080	0	0	0	0	1
7237550310	5420	101930	0	0	0	0	0
1321400060	1715	6819	0	0	0	0	1

```
In [37]: # To avoid multicollinearity we have to drop one column.

# I will drop none as its the least priced category.

X_multilinear_2.drop(axis=1, inplace=True, column=['view_NONE'])

In [38]: # Preview of the dataframe.

X_multilinear_2.head()

Out[38]:
```

	sqft_living	sqft_lot	view_AVERAGE	view_EXCELLENT	view_FAIR	view_GOOD
6414100192	2570	7242	0	0	0	0
2487200875	1960	5000	0	0	0	0
1954400510	1680	8080	0	0	0	0
7237550310	5420	101930	0	0	0	0
1321400060	1715	6819	0	0	0	0

```
In [39]: # using the Ordinary Least Method from stats model we can fit our model.

Multilinear_model_2 = sm.OLS(endog=y, exog=sm.add_constant(X_multilinear_2))
result_multilinear_2 = Multilinear_model_2.fit()

# Let us print out the summary for the model

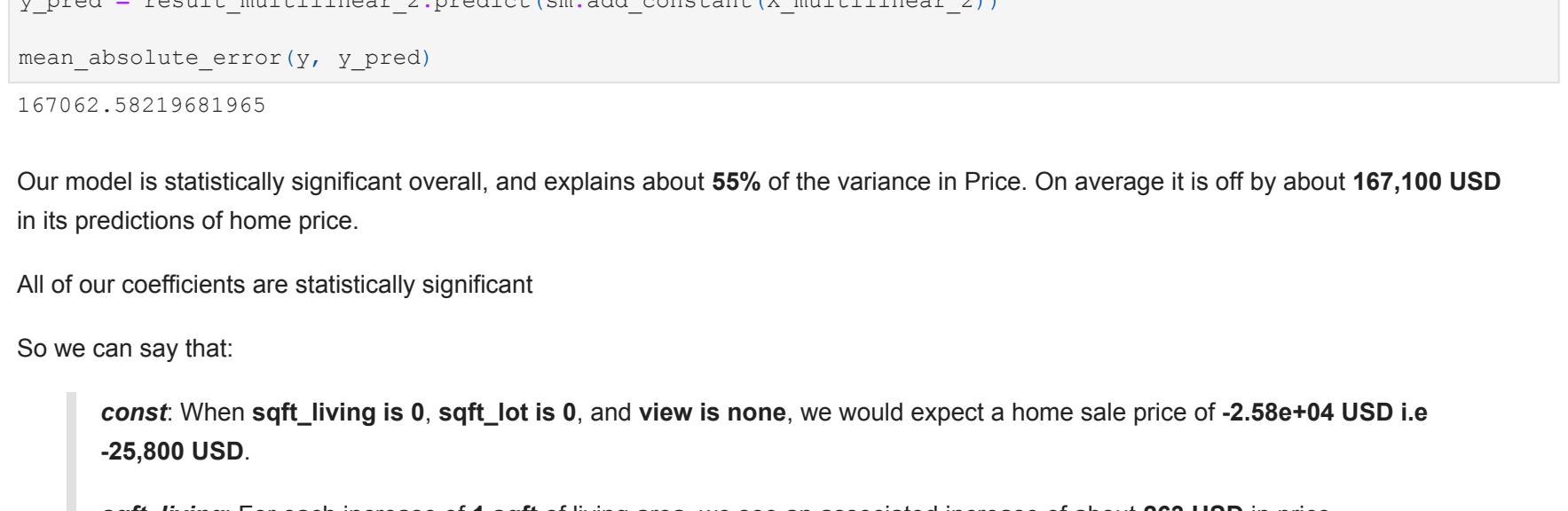
print(result_multilinear_2.summary())
```

```
OLS Regression Results
=====
Dep. Variable:      price      R-squared:      0.549
Model:              OLS      Adj. R-squared:    0.549
Method:             Least Squares      F-statistic:    3196.
Date:               Fri, 30 Sep 2022      Prob (F-statistic):    0.00
Time:               12:50:09      Log-Likelihood:    -2.1827e+05
No. Observations:   15762      AIC:      4.366e+05
Df Residuals:       15755      BIC:      4.366e+05
Df Model:           6
Covariance Type:    nonrobust

=====
coef      std err      t      P>|t|      [0.025      0.975]
-----
const      -2.58e+04      4994.742      -5.186      0.000      -3.56e+04      -1.6e+04
sqft_living    263.4652      2.290      115.063      0.000      258.977      267.953
sqft_lot      -0.3197      0.048      -6.609      0.000      -0.414      -0.225
view_AVERAGE    1.432e+05      9876.703      14.497      0.000      1.24e+05      1.63e+05
view_EXCELLENT    6.009e+05      1.66e+04      36.276      0.000      5.68e+05      6.33e+05
view_FAIR      1.749e+05      1.62e+04      10.818      0.000      1.43e+05      2.07e+05
view_GOOD      2.124e+05      1.39e+04      15.406      0.000      1.85e+05      2.39e+05
-----
Omnibus:      10453.751      Durbin-Watson:    1.972
Prob(Omnibus):    0.000      Jarque-Bera (JB):    452515.177
Skew:            2.621      Prob(JB):      0.00
Kurtosis:        28.721      Cond. No.      3.75e+05
=====
```

Notes:  
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.  
[2] The condition number is large, 3.75e+05. This might indicate that there are strong multicollinearity or other numerical problems.

```
In [40]: fig = plt.figure(figsize=(15,8))
sm.graphics.plot_partregress_grid(
    result_multilinear_2,
    exog_list=list(X_multilinear_2.columns.values),
    grid=(2,4),
    fig=fig)
plt.show()
```



```
In [41]: # Let us use the sklearn to calculate the mean absolute error.

y_pred = result_multilinear_2.predict(sm.add_constant(X_multilinear_2))

mean_absolute_error(y, y_pred)

Out[41]: 167062.58219681965
```

Our model is statistically significant overall, and explains about **55%** of the variance in Price. On average it is off by about **167,100 USD** in its predictions of home price.

All of our coefficients are statistically significant

So we can say that:

**const:** When **sqft\_living** is 0, **sqft\_lot** is 0, and **view** is none, we would expect a home sale price of **-2.58e+04 USD** i.e **-25,800 USD**.

**sqft\_living:** For each increase of 1 **sqft** of living area, we see an associated increase of about **263 USD** in price.

**sqft\_lot:** For each increase of 1 in **sqft** of the lot, we see an associated decrease of about **-0.32 USD** in price.

**view\_AVERAGE:** Compared to None, meaning no view, we see an associated increase of about **1.432e+05** i.e **143,200 USD** for an average view

**view\_FAIR:** Compared to a None, no view , we see an associated increase of about **1.749e+05** i.e **174,900 USD** for a fair view.

**view\_GOOD:** Compared to no view (None\_View), we see an associated increase of about **2.124e+05** i.e **212,400 USD** for a good view.

**view\_EXCELLENT:** Compared to no view (None\_View), we see an associated increase of about **6.009e+05** i.e **600,900 USD** for an excellent view.

Looking at the partial regression plots, the dummy variables look fairly different from the other variables. They tend to have two clusters rather than a continuous "cloud".

## Multilinear Regression with another more categorical column.

Apart from view i would like to try a multilinear model with the condition category.

The condition the house is in will have a major bearing on the price it could potentially sell for.

and since having more than two categories may lead to more errors i will not mix it with views.

```
In [42]: X['condition'].value_counts()

Out[42]:
Average    10221
Good       4137
Very Good  1254
Fair       131
Poor       19
Name: condition, dtype: int64

In [43]: # Definisg a new variable to hold the independent variables to use for the new model.

X_multilinear_3 = X[['sqft_living', 'sqft_lot', 'condition']]

In [44]: # Before i can use the new column views, i have to change it into a binary system.

X_multilinear_3 = pd.get_dummies(X_multilinear_3, columns=['condition'])

In [45]: # Preview of the dataframe.

X_multilinear_3.head()

Out[45]:
```

	sqft_living	sqft_lot	condition_Average	condition_Fair	condition_Good	condition_Poor	condition_Very Good
6414100192	2570	7242	1	0	0	0	0
2487200875	1960	5000	0	0	0	0	1
1954400510	1680	8080	1	0	0	0	0
7237550310	5420	101930	1	0	0	0	0
1321400060	1715	6819	1	0	0	0	0

```
In [46]: # To avoid multicollinearity we have to drop one column.

# I will drop condition_Average as its the worst priced category.

X_multilinear_3.drop(axis=1, inplace=True, column=['condition_Average'])

In [47]: # using the Ordinary Least Method from stats model we can fit our model.

Multilinear_model_3 = sm.OLS(endog=y, exog=sm.add_constant(X_multilinear_3))
result_multilinear_3 = Multilinear_model_3.fit()

# Let us print out the summary for the model

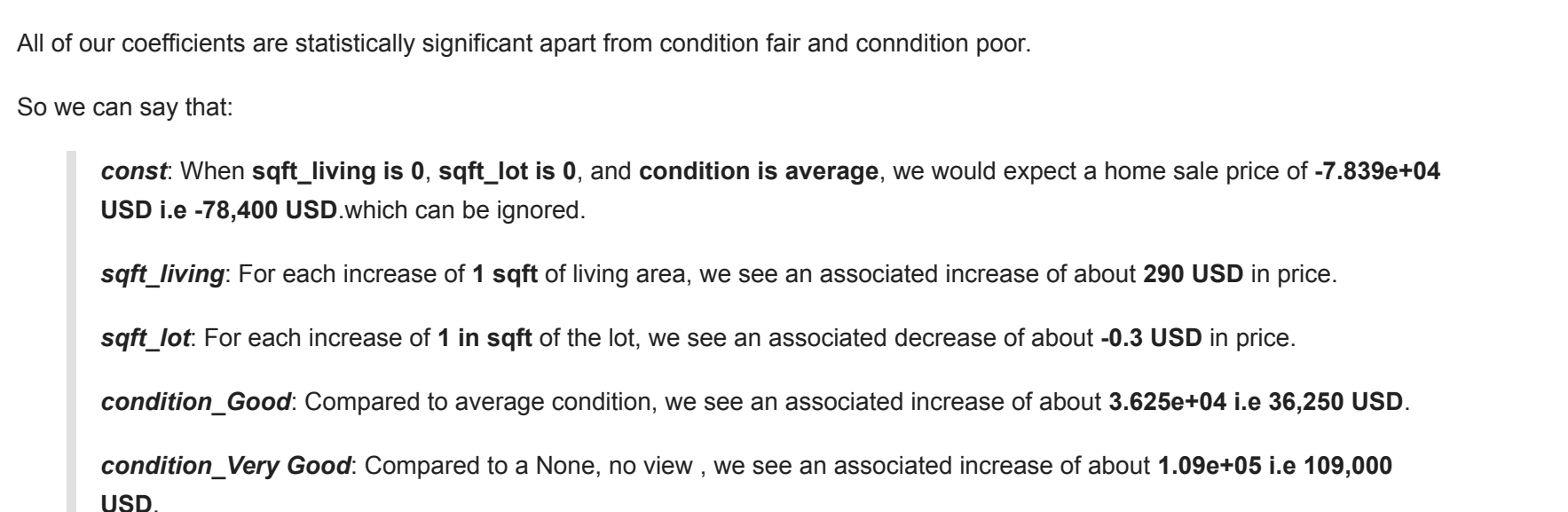
print(result_multilinear_3.summary())
```

```
OLS Regression Results
=====
Dep. Variable:      price      R-squared:      0.507
Model:              OLS      Adj. R-squared:    0.506
Method:             Least Squares      F-statistic:    2697.
Date:               Fri, 30 Sep 2022      Prob (F-statistic):    0.00
Time:               12:50:24      Log-Likelihood:    -2.1898e+05
No. Observations:   15762      AIC:      4.380e+05
Df Residuals:       15755      BIC:      4.380e+05
Df Model:           6
Covariance Type:    nonrobust

=====
coef      std err      t      P>|t|      [0.025      0.975]
-----
const      -7.839e+04      5558.834      -14.101      0.000      -8.93e+04      -6.75e+04
sqft_living    260.0568      2.318      125.404      0.000      256.114      263.200
sqft_lot      -0.2957      0.051      -5.840      0.000      -0.395      -0.196
condition_Fair    6.687e+05      2.31e+04      2.90      0.772      -3.86e+04      5.2e+04
condition_Good    3.625e+04      4841.651      7.487      0.000      2.68e+04      4.57e+04
condition_Poor    8.502e+04      6.01e+04      1.415      0.157      -3.28e+04      2.03e+05
condition_Very Good    1.09e+05      7830.791      13.916      0.000      9.36e+04      1.24e+05
-----
Omnibus:      11119.556      Durbin-Watson:    1.973
Prob(Omnibus):    0.000      Jarque-Bera (JB):    462701.290
Skew:            2.899      Prob(JB):      0.00
Kurtosis:        28.902      Cond. No.      1.28e+06
=====
```

Notes:  
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.  
[2] The condition number is large, 1.28e+06. This might indicate that there are strong multicollinearity or other numerical problems.

```
In [48]: fig = plt.figure(figsize=(15,8))
sm.graphics.plot_partregress_grid(
    result_multilinear_3,
    exog_list=list(X_multilinear_3.columns.values),
    grid=(3,4),
    fig=fig)
plt.show()
```



```
In [49]: # Let us use the sklearn to calculate the mean absolute error.

y_pred = result_multilinear_3.predict(sm.add_constant(X_multilinear_3))

mean_absolute_error(y, y_pred)

Out[49]: 173655.8755334509
```

Our model is statistically significant overall, and explains about **51%** of the variance in Price. On average it is off by about **174,000 USD** in its predictions of home price.

All of our coefficients are statistically significant apart from condition fair and condition poor.

So we can say that:

**const:** When **sqft\_living** is 0, **sqft\_lot** is 0, and **condition** is average, we would expect a home sale price of **-7.839e+04 USD** i.e **-78,400 USD** which can be ignored.

**sqft\_living:** For each increase of 1 **sqft** of living area, we see an associated increase of about **260 USD** in price.

**sqft\_lot:** For each increase of 1 in **sqft** of the lot, we see an associated decrease of about **-0.3 USD** in price.

**condition\_Good:** Compared to average condition, we see an associated increase of about **3.625e+04** i.e **36,250 USD**.

**condition\_Very Good:** Compared to a None, no view , we see an associated increase of about **1.09e+05** i.e **109,000 USD**.

Looking at the partial regression plots, the dummy variables look fairly different from the other variables. They tend to have two clusters rather than a continuous "cloud".

## RESULTS

### Baseline Model

$$y = (290) \text{Sqft\_living} - 55,160 \text{ USD} + 175,000 \text{ USD}$$

This is the formula to predict the price, but it only explains about 50% of the variance in price. Taking into consideration the 95% confidence for the coefficients;

$$\text{constant} = [-65,400 - 45,000] \\ \text{slope} = [282 - 291]$$

$$y = ([282 - 291] \text{Sqft\_living} + [(-65,400 - 45,000)] \text{USD} + 175,000 \text{ USD}$$

### Multilinear Models

From the two multilinear models made above;

**sqft\_lot** views conditions

I will favour the views as all its coefficients are statistically significant in comparison to conditions. And it explains about 60% of the variance in price in comparison to sqft\_lot and sqft\_living alone.

From the views multimedia regression model. The predictive formula is given by;

$$y = (263) \text{Sqft\_living} - (0.32) \text{Sqft\_lot} + (143,200) \text{view\_Average} + (174,900) \text{view\_Fair} + (212,400) \text{view\_Good} + (600,900) \text{view\_Excellent} - 25,800 \text{ USD} + 167,100 \text{ USD}$$

Taking into consideration the 95% confidence for the coefficients;

$$\text{constant} = [-35,600 - 16,000] \text{ slope\_sqft\_living} = [269 - 268] \text{ slope\_sqft\_lot} = [-0.41 - 0.23] \text{ slope\_view\_AVERAGE} = [143,200 - 9,880] \text{ slope\_view\_Fair} = [174,900 - 16,200] \text{ slope\_view\_Good} = [212,400 - 13,800] \text{ slope\_view\_Excellent} = [600,900 16,500]$$

$$y = ([259 - 268] \text{Sqft\_living} - [(-0.41 - 0.23)] \text{Sqft\_lot} + ([143,200 - 9,880] \text{view\_Average} + ([174,900 - 16,200] \text{view\_Fair} + ([212,400 - 13,800] \text{view\_Good} + [600,900 16,600] \text{view\_Excellent} - 25,800 \text{ USD} + 167,100 \text{ USD}$$

## CONCLUSION

I think the accuracy of the predictive formula would be more accurate if the data was grouped by the zipcode as the variance comes from mixing zipcodes of low costing houses with zipcodes with high priced areas. With zipcode specificity the accuracy i believe would be greatly increased instead of just relying on county. I got this observation by using tableau.