

# Final Project Submission

Please fill out:

- Student name: Zeth Abney
- Student pace: Flex (40 week)
- Scheduled project review date/time: 04/20/2022 9:00 AM
- Instructor name: Matt Bombard
- Blog post URL: !!!TBA!!!

## Overview



The data set used in this analysis is open-source data available directly from Kings County's website (<https://kingcounty.gov/services/data.aspx>). This particular data set covers various aspects of realestate transactions including date of sale, square footage of house and lot, proximity to recreational and natural resources, etc.

The initial data set used contains 21597 total records, approximately 7% of which is eventually thrown out as a result of either data cleaning or model fitting. The dataset starts with 20 total features (i.e. columns), only 8 of which are ultimately included in the final model also with 6 additional features inferred from the original data (e.g. one-hot encoding).

For the purposes of regression modeling the data is manipulated so that every datapoint is encoded as either an integer or a decimal and there are no null or missing values. Also, the target variable 'price' is eventually log-transformed as part of the model fitting process; keep in mind that because of the log-transformation the model coefficients should be interpreted as the percentages rather than the metric's own units. This is explained further in the regression results section of this notebook.

For more details on understanding the data and statistics of the final model see the [model data dictionary \(model\\_dictionary.md\)](#).

## The Business Problem

FlipHouse, LLC. Is a profesional 'house flipping' business seeking to enter the realestate market of the pacific northwest, and specifically Seattle Washington and the surrounding area (i.e. Kings County). Before beginning any projects in the area, FlipHouse decision makers need to better understand the how to determine the opportunity cost, and potential returns for any investments made in Kings County, as well as how to maximize those returns. The oportunity cost and potential returns can indeed be determined by understanding how time, physical location, and physical condition and attributes all affect the price of a real estate property.

Therefore, this anaylisis will seek to build a statistical model that is informative as far as specifying what metrics to use and how strong each metric may be in terms of predicting the market value of a real estate property.

## Understanding the data



The data set used in this analysis is open-source data available directly from Kings County's website (<https://kingcounty.gov/services/data.aspx>). This particular data set covers various aspects of realestate transactions including date of sale, square footage of house and lot, proximity to recreational and natural resources, etc.

The initial data set used contains 21597 total records, aproximately 7% of which is eventually thrown out as a result of either data cleaning or model fitting. The dataset starts with 20 total features (i.e. columns), only 8 of which are ultimately included in the final model aslo with 6 additional features inferred from the original data (e.g. one-hot encoding).

For the purposes of regression modeling the data is manipulated so that every datapoint is encoded as either and integer or a decimal and there are no null or missing values. Also, the target variable 'price' is eventually log-transomed as part of the modelfitting process; keep in mind that because of the log-transformation the model coefficients should be interpreted as the percentages rather than the metric's own units. This is explained further in the regrssion results section of this notebook.

For more details on understanding the data and statistics of the final model see the `model_dictionary.md` file in this repo.

In [2]: ►

```
1 # raw data handling
2 import pandas as pd
3 import numpy as np
4 import datetime as dt
5
6 # data visualization
7 import matplotlib.pyplot as plt
8 import matplotlib.cm as cm
9 import seaborn as sns
10 from pywaffle import Waffle
11 import squarify
12
13 plt.style.use('fivethirtyeight')
14
15 # regression modeling
16 import statsmodels.api as sm
17 from statsmodels.formula.api import ols
18
19 # model validation
20 from sklearn.linear_model import LinearRegression
21 from sklearn.model_selection import train_test_split
22 from sklearn.model_selection import cross_val_score
23 from sklearn.metrics import mean_absolute_error, mean_squared_error
24
25 import warnings # weird sns.distplot() warnings
26 warnings.filterwarnings("ignore")
27
28 import squarify
29
30 plt.style.use('fivethirtyeight')
```

In [3]: ►

```
1 # returns OLS Linear regression model
2 def run_OLS_model(X,y):
3
4     predictors_int = sm.add_constant(X)
5     model = sm.OLS(y,predictors_int).fit()
6
7     return model
```

## Data preparation

The data is initially imported with 21597 records, 453 are immediately eliminated due to some '?' values that can not be justifiably replaced with any sort of filler such as 0 or 'NONE'. By the end of the model development process there were 19982 records remaining eliminating about 7% of the initial data overall.

In [4]:

```
1 data = pd.read_csv('data/kc_house_data.csv')
2 data.drop(data.loc[data['sqft_basement']=='?'].index,inplace=True)
```

Null values and non-numerical values are then re-encoded so that all values in the data set are numerical with no missing values. The details of how this is achieved for each data feature is explained in the comments below.

In [6]:

```
1 # convert all string types into np floats
2 data.sqft_basement = [float(sq) for sq in list(data.sqft_basement)]
3
4
5 # Replaces grade strings with numerics based on data dict.
6 grade_strings = list(data.grade.unique()) # List of unique values from column
7 grade_nums = [int(grade.split()[0]) for grade in list(data.grade.unique())]
8
9 # replaces a cell value with the int of the first character of its existing string
10 data.grade.replace(to_replace=grade_strings,value=grade_nums,inplace=True)
11
12
13 # replaces condition objects with numerics based on data dict.
14 condition_dict = {'Poor':1,'Fair':2,'Average':3,'Good':4,'Very Good':5}
15 data.condition.replace(to_replace=condition_dict,inplace=True)
16
17
18 # replace yr_built NaNs with numeric 0
19 data.yr_renovated.replace(to_replace=np.nan,value=0,inplace=True)
20
21
22 # convert waterfront into numeric boolean
23 waterfront_bool_dict = {'YES':1,'NO':0,np.nan:0}
24 data.waterfront.replace(to_replace=waterfront_bool_dict,inplace=True)
25
26
27 # convert view from string into categorical ordinal
28 view_dict = {np.nan:0,'NONE':0,'FAIR':1,'AVERAGE':2,'GOOD':3,'EXCELLENT':4}
29 data.view.replace(to_replace=view_dict,inplace=True)
30
31
32 # convert dates into ordinals, extrapolate month from date into new column
33 data.date = pd.to_datetime(data['date']) # convert date string into date
34 data['sale_month'] = data.date.apply(lambda x: x.month) # extrapolate month
35 data.date = data['date'].map(dt.datetime.toordinal) # convert original date to ordinal
36 data.rename({'date':'sale_date'},axis=1,inplace=True) # rename date to 'sale_date'
```

Next, in order to strengthen the model, outliers are eliminated. For this analysis an outlier is considered any record where the target variable (price) is greater than 3 standard deviations from the target variable mean.

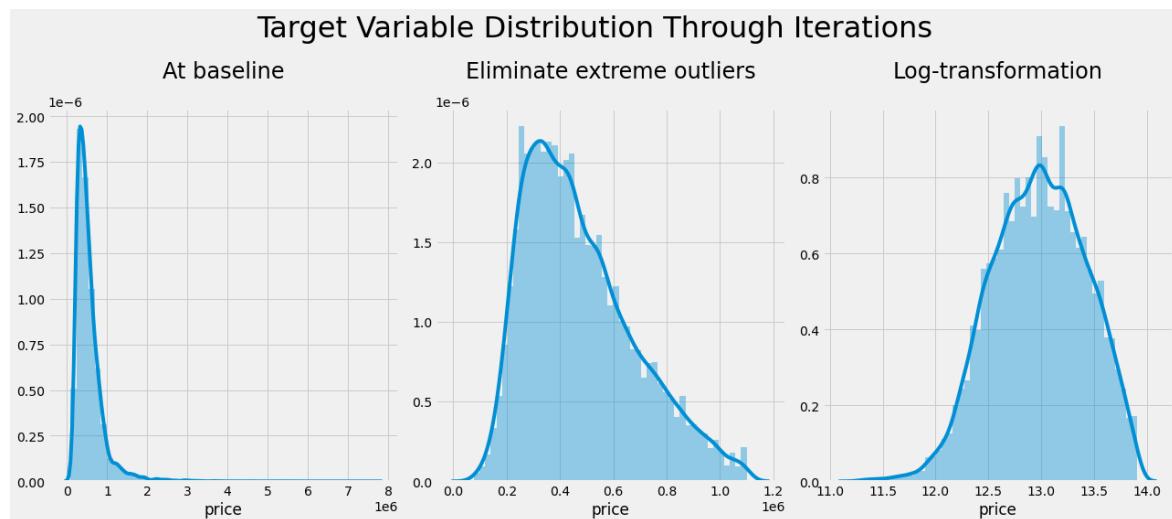
- Using statistical maximum (inter-quartile range \* 1.5) instead was also investigated but it eliminated roughly half of the data.

```
In [7]: ❶ 1 std_thresh = data.price.std()*3 # value of the third central moment of the price column
2 outlier_eliminated_df = data.loc[abs(data['price']) <= std_thresh] # slice the data
3
4 # assign X and y variables to the appropriate dataframes/series
5 y = outlier_eliminated_df.price
6 X_outliers_eliminatd = outlier_eliminated_df.drop('price',axis=1)
```

In order to increase normality in both the target variable as well as the model itself the target variable is log transformed

```
In [9]: ❶ 1 y_log = np.log(y)
```

```
In [10]: ❶ 1 fig, (ax1,ax2,ax3) = plt.subplots(1,3)
2
3 og = sns.distplot(data.price,ax=ax1).set_title('At baseline\n',fontsize=16)
4 ot = sns.distplot(y,ax=ax2).set_title('Eliminate extreme outliers\n',font-size=16)
5 lo = sns.distplot(y_log,ax=ax3).set_title('Log-transformation\n',font-size=16)
6
7 ax1.set_ylabel("")
8 ax2.set_ylabel("")
9 ax3.set_ylabel("")
10
11
12 plt.gcf().set_size_inches(18, 8)
13 plt.suptitle("Target Variable Distribution Through Iterations",font-size=16)
14 fig.tight_layout()
15 plt.show()
16
```



At this point the model was strong, the set of features was somewhat limited. The 'zipcode' feature also still remained and was essentially uninterpretable due to the fact that it is not only categorical but nominal in nature. The next successfully better iteration was to extrapolate the proximity of a particular waterfront from the zipcode feature, and then one-hot encoding each waterfront location. The details of how this was performed is outlined below.

In [11]: ►

```

1 # Dictionary with zipcodes associated with waterfronts described by the columns
2 water_loc_dict = {'Duwamish':[98168],
3 'Elliott Bay':[98119,98104,98129,98132,98127,98125,98195,98101,98134,98196],
4 'Puget Sound':[98071,98083,98013,98070,98031,98131,98063,98195,98207,98196],
5 'Lake Union':[98109],
6 'Ship Canal':[00000],
7 'Lake Washington':[98072,98077],
8 'Lake Sammamish':[98074,98075,98029],
9 'other lake':[00000],
10 'river/slough waterfronts':[00000]}
11
12 # List to contain new column data
13 waterfront_list = []
14
15 # for Loop to assign waterfront based on zipcode
16 for zipcode in X_outliers_eliminatd.zipcode:
17     for k,v in water_loc_dict.items():
18         if zipcode in v:
19             waterfront_list.append(k)
20             appended = True
21             break
22         else:
23             appended = False
24     if not appended:
25         waterfront_list.append('NONE')
26
27 # for readability and convenience create new dataframe including waterfront_loc
28 X_waters = X_outliers_eliminatd
29 X_waters['waterfront_loc'] = waterfront_list
30
31 # one-hot encoding waterfront_loc
32 waterfront_dummies = pd.get_dummies(X_waters.waterfront_loc,prefix='waterfront')
33 X_hot_waters = pd.concat([X_waters,waterfront_dummies],axis=1)
34
35 # zipcode has now been interpolated into waterfront_loc, waterfront_loc is
36 # so it is not able to be passed into the regression model
37 X_hot_waters = X_hot_waters.drop(['zipcode','waterfront_loc'],axis=1)

```

## Modeling

An extensive iterative process was conducted to achieve the final model, to see this process in detail please explore the notebook file 'EDA.ipynb' in this repository.

The iterative process included the above described data engineering techniques, but also performing feature selecting and elimination before and after every dataframe manipulation using a stepwise forward-backward selection (based on coefficient P value, community sourced function... see EDA notebook), recursive feature ranking and selection with cross validation based on coefficient strength (from scikit learn library, features selection module), elimination based on variance inflation factor.

To repeat this entire iterative process in this notebook would defeat the purpose of this notebook, so a list of column labels from the final model in the EDA journal is copy/pasted here in order to slice those particular columns from the dataframe prepared within this notebook. This should reproduce the same model as the final model in the EDA journal

```
In [12]: ┌ 1 feature_list = ['bathrooms', 'lat', 'grade', 'sqft_living15', 'view', 'condition', 'waterfront', 'grade_ix', 'sqft_living15_ix', 'lat_ix', 'view_ix', 'condition_ix', 'waterfront_ix']
  2 X = X_hot_waters[feature_list]
```

```
In [13]: ┌ 1 model = run_OLS_model(X,y_log)
  2 residuals = model.resid
```

In [14]: 1 model.summary()

Out[14]: OLS Regression Results

<b>Dep. Variable:</b>	price	<b>R-squared:</b>	0.674				
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.673				
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	2944.				
<b>Date:</b>	Sun, 17 Apr 2022	<b>Prob (F-statistic):</b>	0.00				
<b>Time:</b>	18:24:00	<b>Log-Likelihood:</b>	-944.25				
<b>No. Observations:</b>	19982	<b>AIC:</b>	1919.				
<b>Df Residuals:</b>	19967	<b>BIC:</b>	2037.				
<b>Df Model:</b>	14						
<b>Covariance Type:</b>	nonrobust						
		coef	std err	t	P> t	[0.025	0.975]
const	-173.1267	12.542	-13.804	0.000	-197.710	-148.543	
bathrooms	0.0864	0.003	26.243	0.000	0.080	0.093	
lat	1.4771	0.013	111.257	0.000	1.451	1.503	
grade	0.1527	0.003	58.332	0.000	0.148	0.158	
sqft_living15	0.0002	4.13e-06	39.495	0.000	0.000	0.000	
view	0.0734	0.003	24.402	0.000	0.068	0.079	
condition	0.0999	0.003	35.202	0.000	0.094	0.105	
waterfront_Lake Washington	-0.2440	0.013	-19.292	0.000	-0.269	-0.219	
waterfront_Duwamish	-0.2243	0.016	-14.064	0.000	-0.256	-0.193	
sale_date	0.0002	1.7e-05	9.108	0.000	0.000	0.000	
waterfront	0.3495	0.038	9.265	0.000	0.276	0.423	
waterfront_Lake Union	0.3121	0.028	11.098	0.000	0.257	0.367	
waterfront_Lake Sammamish	0.0366	0.008	4.318	0.000	0.020	0.053	
sale_month	-0.0025	0.001	-4.068	0.000	-0.004	-0.001	
waterfront_Puget Sound	-0.0569	0.013	-4.243	0.000	-0.083	-0.031	
<b>Omnibus:</b>	218.309	<b>Durbin-Watson:</b>	1.994				
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	372.712				
<b>Skew:</b>	0.041	<b>Prob(JB):</b>	1.17e-81				
<b>Kurtosis:</b>	3.664	<b>Cond. No.</b>	5.14e+09				

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 5.14e+09. This might indicate that there are strong multicollinearity or other numerical problems.

In [21]: ►

```
1 # create a train and a test split
2 X_train, X_test, y_train, y_test = train_test_split(X, y_log, test_size=0.2)
3
4 # generate predictions based on test and train samples
5 linreg = LinearRegression()
6 linreg.fit(X_train, y_train)
7
8 y_hat_train = linreg.predict(X_train)
9 y_hat_test = linreg.predict(X_test)
10
11 # calculate train and test residuals
12 train_residuals = y_hat_train - y_train
13 test_residuals = y_hat_test - y_test
```

In [22]: ►

```
1 # calculate train and test mean squared error
2 train_mse = mean_squared_error(y_train, y_hat_train)
3 test_mse = mean_squared_error(y_test, y_hat_test)
4
5 print('Train Mean Squared Error:', round(train_mse,3))
6 print('Test Mean Squared Error:', round(test_mse,3))
7 print('Difference: ', round(abs(train_mse - test_mse),3))
```

Train Mean Squared Error: 0.064

Test Mean Squared Error: 0.065

Difference: 0.0

In [23]: ►

```
1 cv_5_results = np.mean(cross_val_score(linreg, X, y_log, cv=5, scoring='neg_mean_squared_error'))
2 cv_10_results = np.mean(cross_val_score(linreg, X, y_log, cv=10, scoring='neg_mean_squared_error'))
3 cv_20_results = np.mean(cross_val_score(linreg, X, y_log, cv=20, scoring='neg_mean_squared_error'))
4
5 print('Five k-fold MSE: ', round(cv_5_results,3))
6 print('Ten k-fold MSE: ', round(cv_10_results,3))
7 print('Twenty k-fold MSE: ', round(cv_20_results,3))
```

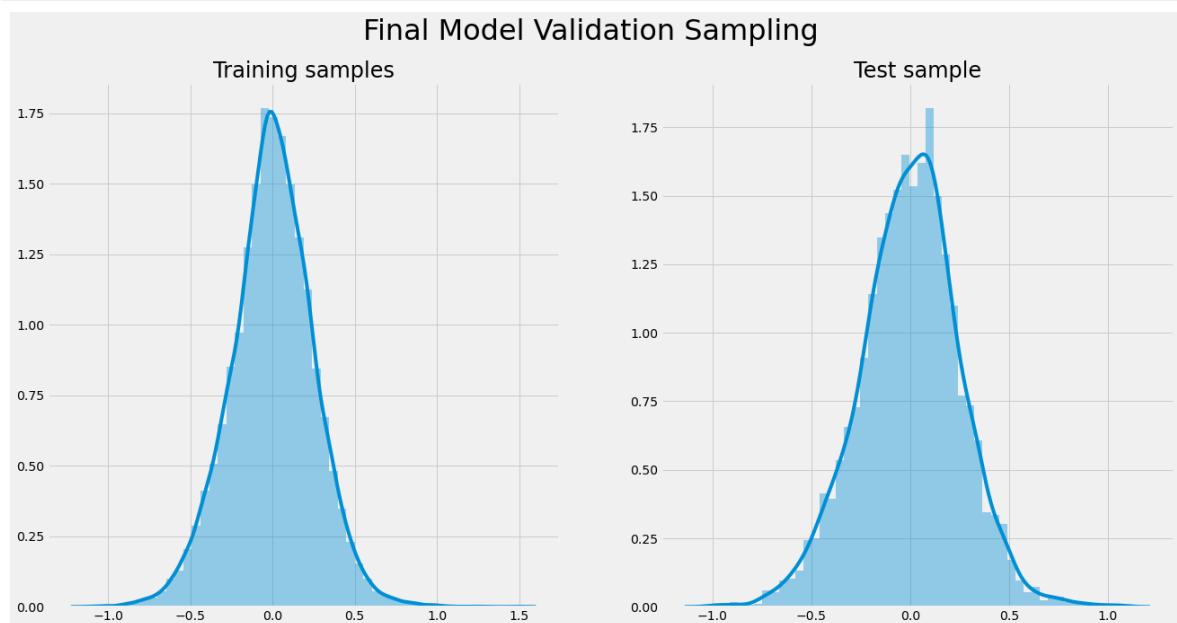
Five k-fold MSE: -0.064

Ten k-fold MSE: -0.064

Twenty k-fold MSE: -0.064

In [24]:

```
1 fig, (ax1,ax2) = plt.subplots(1,2)
2
3 sns.distplot(train_residuals,ax=ax1)
4 ax1.set_title('Training samples',fontsize=24)
5 ax1.set_ylabel('')
6 ax1.set_xlabel('')
7 sns.distplot(test_residuals,ax=ax2)
8 ax2.set_title('Test sample',fontsize=24)
9 ax2.set_ylabel('')
10 ax2.set_xlabel('')
11
12 fig.suptitle('Final Model Validation Sampling\n',fontsize=32)
13 plt.gcf().set_size_inches(20, 10)
14 plt.show()
```



The model above, as well as its mean-squared error, is identical to the final model decided on in the EDA notebook. This means that this model is in fact reproducible and not an artifact of some functional error within the notebook.

## Regression Results

```
In [25]: # data series of features and their coefficients
1 coef_series = model.summary2().tables[1]['Coef.'][1:]
2
3
4 # above dataserries as a dictionary
5 feature_effect_dict = {}
6 for item in coef_series.iteritems():
7     feature_effect_dict[item[0]] = round(item[1],3)
```

Use this function below to see in plain language how each of the model feature affects price

```
In [26]: def feature_impacts(dictionary=feature_effect_dict):
1
2
3     feature_list = sorted(list(feature_effect_dict.keys()))
4
5     print('Please select a feature from the following list...\n')
6
7     for item in feature_list:
8         print(item)
9
10    user_input = input("")
11
12    if user_input not in feature_list:
13        print('\nPlease check spelling, and ensure your input is in the list')
14
15    effect = feature_effect_dict[user_input]
16
17    if effect > 0 :
18        print(f"\nFor every unit increase of {user_input}, the price of a house increases by {effect} units")
19    else:
20        print(f"\nFor every unit increase of {user_input}, the price of a house decreases by {abs(effect)} units")
```

```
In [28]: # first uncomment the line below, press shift+enter, then type in the feature name
1 # feature_impacts()
```

## My interpretation and summary of the model

I have confidence that this is a strong model because as best I can tell it satisfies the four assumptions of linear regression as best as it can with the data available.

- When plotted the model residuals show strong resemblance to the t-distribution satisfying the normality assumption, save only for slightly long tails.
- The Drubin-Watson score is nearly 2.0 exactly, clearly demonstrating that the model satisfies the homoscedasticity assumption.
- The linearit assumption is difficult to validate with a single model alone, but I have been incresingly convinced that this model satisfies the linearity assumption thorughout the development of this model.

- The condition number is heavy which can be interpreted as an indication of multi-collinearity. However there may be a scaling issue here, when the condition numbers of each iteration of the model from the baseline up to the one in this notebook are plotted, the condition number of this model is 0 relative to its prototypical counterparts. Likewise for the Jarque Bera, and indication of normality.
- The kurtosis of this model is slightly outside the normally accepted range, through the iterative process any attempts to address this issue drastically worsened all other indicators of model strength. Additionally the kurtosis of 3.664 is the lowest of any iteration of the model so I have to run with it.

The intention of this analysis is to help Fliphouse, LLC. understand where to purchase properties as well as what physical aspects of the properties to consider before buying.

In regard to where to purchase properties, the latitude that a property lies on seems to be the strongest predictor of price. For every additional degree north, a property's price will increase by 1.47%. However the properties in this data set cover a range less than 1 degree. One degree of latitude is roughly 69 miles, and the range covered by the model data is only about 43 miles. So we can consider this to mean that approximately for every additional mile north, the price of a property increases by 0.034%

Whether or not a property lies on a waterfront is a strong predictor as well, if it in fact does the price increases by 0.35%. Additionally if a property is even in the same zipcode as Lake Union the price of the property will increase by 0.31%. However property prices decrease if the property is in the same zipcodes as Lake Washington, Duwamish, or Puget Sound.

In regard to what physical aspects of a property to consider, grade is the strongest predictor of price. Grade is on a scale of 1-11 and indicates how well the property satisfies the building codes as well as the quality and level of luxury of the property; for every additional point on the 1-11 scale the price of a property increases by 0.15% meaning that refurbishing a home that is not even up to code into a luxury home could increase the price by more than an entire percentage point.

Condition is another strong predictor and is similar to grade, it is a more qualitative version and does not refer to its satisfaction of building codes, it is on a scale of 1-5 and according to the model for every additional point here the price of the property will increase by 0.1%.

View is worth considering as well, the quality of the view from a property is also on a scale of 1-5 and for every additional point the price of a property increases by 0.07%. It is also worth noting that the price of a property increases by 0.09% for every additional bathroom in the house.

When is the best time to buy is somewhat unclear. The model suggests that prices decrease slightly month-to-month if you begin in January (i.e. sale\_month, -0.0025). However prices tend to increase the more recent the sale was (i.e. sale\_date, 0.0002). Considering that the dataset covers a range of dates from May 02, 2004 to May 24, 2015, my intuition is that spring/early summer is the when prices are highest, but I'm not confident that I can empirically support that claim with this model.

## Conclusion

For the purposes of Fliphouse, LLC penetrating the King County realestate market. I recommend looking at properties on Lake Union, or north of downtown Seattle that are also on a waterfront, preferably with a view of said waterfront. Furthermore I recommend finding properties that satisfy the aforementioned geographic stipulations that are in need of repair and potentially do not satisfy city building codes. The best course of action according to this analysis would be to bring the building and property up to code and beyond and additionally perhaps add one or two bathrooms to the home (perhaps an outdoor shower near the waterfront access).

In [29]: 1 model.summary()

Out[29]: OLS Regression Results

<b>Dep. Variable:</b>	price	<b>R-squared:</b>	0.674				
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.673				
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	2944.				
<b>Date:</b>	Sun, 17 Apr 2022	<b>Prob (F-statistic):</b>	0.00				
<b>Time:</b>	16:19:37	<b>Log-Likelihood:</b>	-944.25				
<b>No. Observations:</b>	19982	<b>AIC:</b>	1919.				
<b>Df Residuals:</b>	19967	<b>BIC:</b>	2037.				
<b>Df Model:</b>	14						
<b>Covariance Type:</b>	nonrobust						
		coef	std err	t	P> t	[0.025	0.975]
const	-173.1267	12.542	-13.804	0.000	-197.710	-148.543	
bathrooms	0.0864	0.003	26.243	0.000	0.080	0.093	
lat	1.4771	0.013	111.257	0.000	1.451	1.503	
grade	0.1527	0.003	58.332	0.000	0.148	0.158	
sqft_living15	0.0002	4.13e-06	39.495	0.000	0.000	0.000	
view	0.0734	0.003	24.402	0.000	0.068	0.079	
condition	0.0999	0.003	35.202	0.000	0.094	0.105	
waterfront_Lake Washington	-0.2440	0.013	-19.292	0.000	-0.269	-0.219	
waterfront_Duwamish	-0.2243	0.016	-14.064	0.000	-0.256	-0.193	
sale_date	0.0002	1.7e-05	9.108	0.000	0.000	0.000	
waterfront	0.3495	0.038	9.265	0.000	0.276	0.423	
waterfront_Lake Union	0.3121	0.028	11.098	0.000	0.257	0.367	
waterfront_Lake Sammamish	0.0366	0.008	4.318	0.000	0.020	0.053	
sale_month	-0.0025	0.001	-4.068	0.000	-0.004	-0.001	
waterfront_Puget Sound	-0.0569	0.013	-4.243	0.000	-0.083	-0.031	
<b>Omnibus:</b>	218.309	<b>Durbin-Watson:</b>	1.994				
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	372.712				
<b>Skew:</b>	0.041	<b>Prob(JB):</b>	1.17e-81				
<b>Kurtosis:</b>	3.664	<b>Cond. No.</b>	5.14e+09				

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 5.14e+09. This might indicate that there are strong multicollinearity or other numerical problems.

## Next steps

It is worth one-hot encoding some of the other categorical features such as view and grade. Also, investigating eliminating outliers based on category (e.g. there are only 27 data points with a grade of 4, and only 1 with a grade of 3 or less.). Some of the leptokurtosis may be address by continuing to eliminate outliers using this method.

This dataset covers a timespan of roughly one year, it would certainly strengthen the model to include analogous data from other years both before and after.

## Visualizations

Some visualizations for presentation purposes meant to hint at how well the model can answer various questions they would likely like to ask.

```
In [27]: ┏ 1 coef_series = coef_series.sort_values(ascending=False)
2
3 quest_dict = {
4 'How far north?':coef_series['lat'],
5 'On a waterfront?':coef_series['waterfront'],
6 'Which waterfront?':sum(abs(coef_series[['waterfront_Lake Union','waterf
7 'How much repair needed?':sum(abs(coef_series[['condition','grade']]))),
8 'How many bathrooms?':coef_series['bathrooms'],
9 'Good view?':coef_series['view'],
10 'When?':sum(abs(coef_series[['sale_date','sale_month']])))
11 }
12
13 quest_dict
```

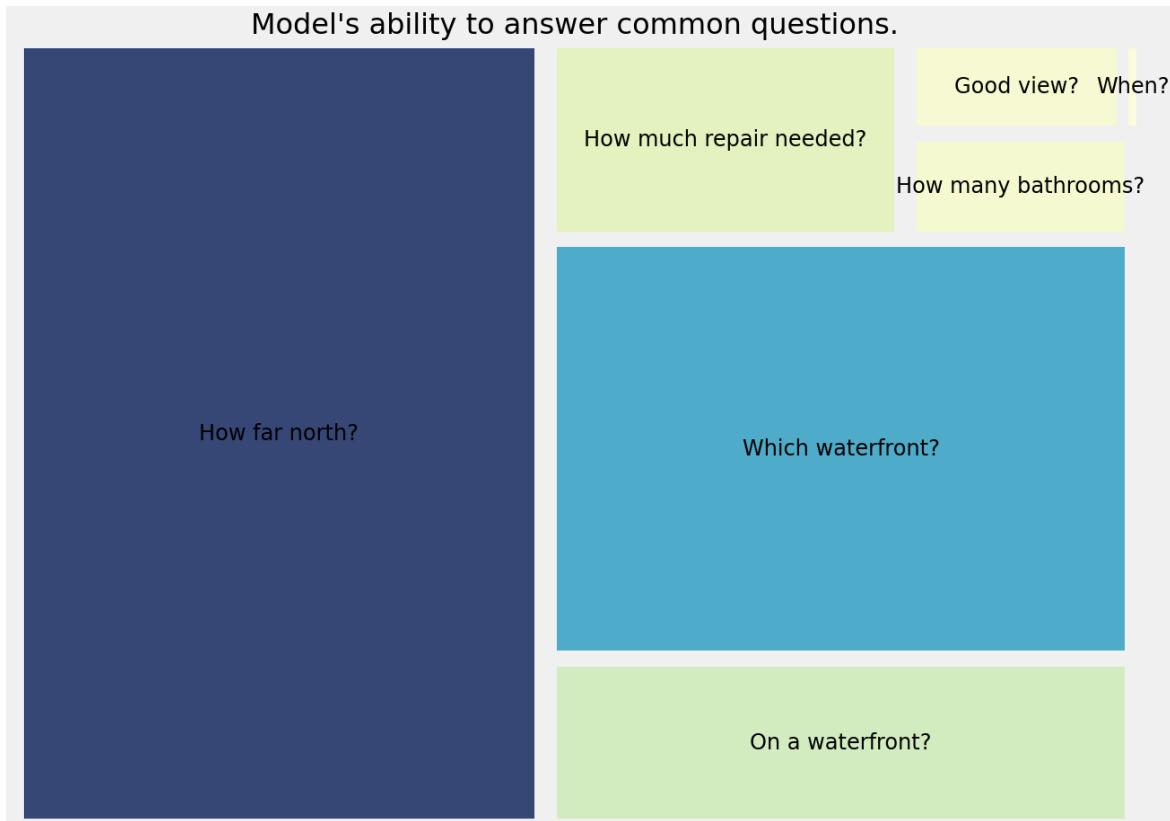
```
Out[27]: {'How far north?': 1.4770811449849992,
'On a waterfront?': 0.34949918084450954,
'Which waterfront?': 0.8739535489846464,
'How much repair needed?': 0.25260992113104447,
'How many bathrooms?': 0.08639102479362598,
'Good view?': 0.07343130156894614,
'When?': 0.0026533397723983573}
```

In [28]:

```
1 minima = min(quest_dict.values())
2 maxima = max(quest_dict.values())
3
4 norm = cm.colors.Normalize(vmin=minima, vmax=maxima, clip=True)
5 mapper = cm.ScalarMappable(norm=norm, cmap=cm.YlGnBu)
6 color_map = []
7
8 for k,v in quest_dict.items():
9     color_map.append(mapper.to_rgba(v))
```

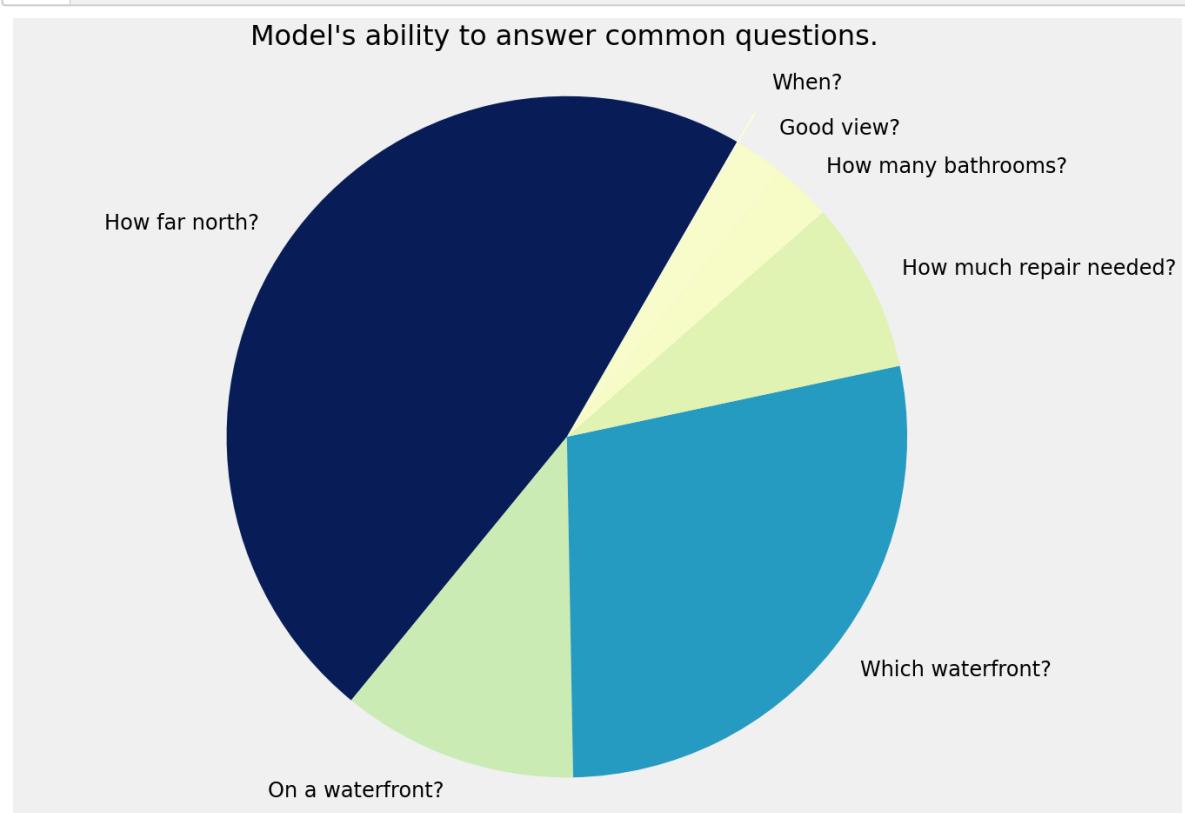
In [29]:

```
1 squarify.plot(sizes=quest_dict.values(), label=quest_dict.keys(),
2                 alpha=.8, text_kwargs={'fontsize':24}, color=color_map, pad=.1)
3
4
5 plt.gcf().set_size_inches(20, 15)
6 plt.axis('off')
7 plt.show()
```



In [48]:

```
1 fig, ax = plt.subplots()
2
3 explode = (0, 0, 0, 0, 0, 0, 0, 0.1)
4
5 ax.pie(quest_dict.values(), labels=quest_dict.keys(), explode=explode, star
6 ax.set_title("Model's ability to answer common questions.", fontsize=32)
7 ax.axis('equal')
8
9 plt.gcf().set_size_inches(20, 15)
10 plt.show()
```



For additional info, contact Zeth Abney at [\(mailto:zethusabney@gmail.com\)](mailto:zethusabney@gmail.com)

