

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import folium
6 import json
7 import plotly.express as px
8
9 from matplotlib import style
10 from folium import plugins
11 from plotly.subplots import make_subplots
12
13 #from mlxtend.evaluate import bias_variance_decomp
14
15 from warnings import filterwarnings
16 filterwarnings('ignore')
17
18 %matplotlib inline
```

```
In [2]: 1 # Choropleth Function
2
3 #Zipcode choropleth maps with average values per a zipcode (King County)
4 def map_choropleth_zip(df, column, title, column_name):
5     fig=px.choropleth_mapbox(data_frame=df, locations='zipcode', geojson=KC_zip_json, color=column,
6                             mapbox_style='open-street-map', zoom=8.5, height=900, featureidkey='properties.ZCTA5CE10',
7                             center={'lat': 47.403768, 'lon': -122.005863}, opacity=0.4,
8                             color_continuous_scale=px.colors.sequential.YlOrRd,
9                             title=title,
10                            template = "plotly_dark",
11                            labels={
12                                column: column_name})
13     fig.update_layout(
14         font_family="Arial",
15         font_size=16,
16         font_color="white",
17         title_font_family="Arial",
18         title_font_color="white",
19         title_font_size=20)
20
21     fig.update_layout(
22         title={
23             'y':0.98,
24             'x':0.5,
25             'xanchor': 'center',
26             'yanchor': 'top',
27         })
28
29     fig.show()
30     return None
31
```

```
In [3]: 1 # Importing raw data
2 df=pd.read_csv('data/kc_house_data.csv')
3 df_test=pd.read_csv('data/df_test.csv', dtype={'recent_renovation_new_str': str})
4 df_zipcode_viz=pd.read_csv('data/df_zipcode_vs.csv')
5 KC_zip_json=json.load(open('data/wa_washington_zip_codes_geo.min.json', 'r'))
```

### Displaying DataFrames

In [4]: 1 df\_test.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5779 entries, 0 to 5778
Data columns (total 13 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Unnamed: 0        5779 non-null   int64  
 1   sqft_living_st    5779 non-null   float64 
 2   distance_st       5779 non-null   float64 
 3   bathrooms_st      5779 non-null   float64 
 4   grade_st          5779 non-null   float64 
 5   log_price          5779 non-null   float64 
 6   recent_renovation_new  5779 non-null   float64 
 7   sqft_living         5779 non-null   int64  
 8   distance            5779 non-null   float64 
 9   grade               5779 non-null   int64  
 10  bathrooms           5779 non-null   float64 
 11  price               5779 non-null   float64 
 12  recent_renovation_new_str 5779 non-null   object  
dtypes: float64(9), int64(3), object(1)
memory usage: 587.1+ KB
```

In [5]: 1 df\_zipcode\_viz

Out[5]:

	Unnamed: 0	zipcode	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	condition	grade	yr_built
0	0	98001	2.813837e+05	3.402235	2.012570	1908.256983	14966.393855	1.432961	0.000000	0.094972	3.329609	7.296089	1981.047486
1	1	98002	2.340839e+05	3.314721	1.837563	1627.416244	7509.248731	1.332487	0.000000	0.010152	3.751269	6.695431	1967.695431
2	2	98003	2.942254e+05	3.358696	2.053442	1931.438406	10625.137681	1.309783	0.000000	0.217391	3.373188	7.550725	1976.873188
3	3	98004	1.355200e+06	3.853968	2.530159	2910.730159	13084.374603	1.431746	0.003175	0.307937	3.495238	8.688889	1971.542857
4	4	98005	8.102897e+05	3.851190	2.424107	2656.803571	19928.785714	1.279762	0.000000	0.095238	3.696429	8.488095	1969.744048
...	...	...	...	...	...	...	...	...	...	...	...	...	...
65	65	98177	6.770352e+05	3.397638	2.099409	2325.511811	11921.086614	1.277559	0.003937	0.814961	3.496063	7.976378	1960.822835
66	66	98178	3.106226e+05	3.306202	1.738372	1736.744186	8308.457364	1.186047	0.034884	0.542636	3.325581	6.829457	1955.325581
67	67	98188	2.893271e+05	3.437037	1.870370	1806.125926	10136.644444	1.225926	0.000000	0.148148	3.333333	7.044444	1965.681481
68	68	98198	3.029441e+05	3.178182	1.792727	1749.530909	10553.425455	1.225455	0.032727	0.603636	3.450909	7.109091	1966.865455
69	69	98199	7.919480e+05	3.208861	2.166930	2162.246835	5439.721519	1.466772	0.003165	0.553797	3.503165	8.009494	1956.569620

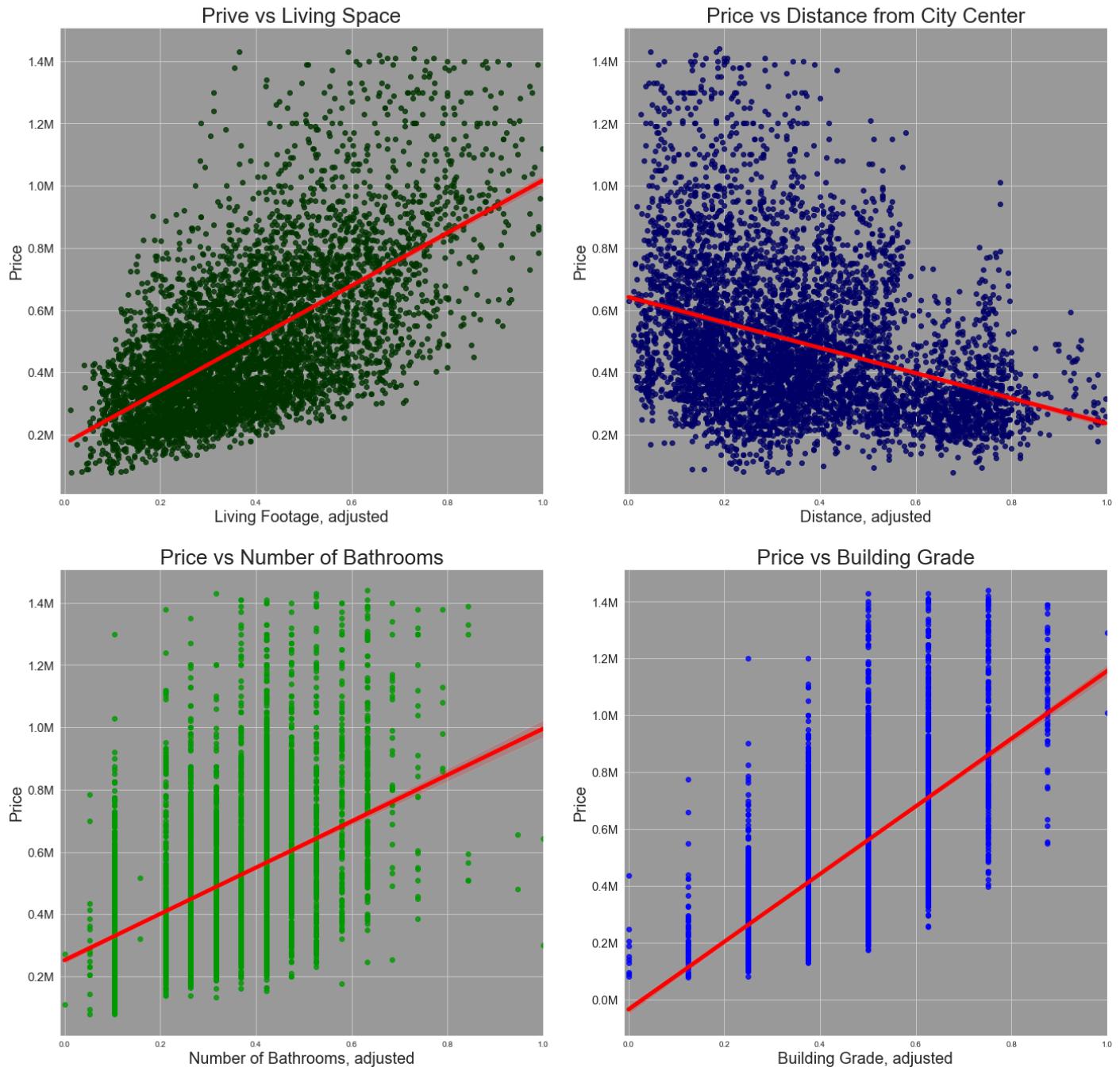
70 rows × 15 columns

Visualization

In [6]:

```
1 # Regplots for all four variables
2
3 sns.set_style("darkgrid", {"axes.facecolor": ".6"})
4 fig, axes = plt.subplots(figsize=(20,20), ncols=2, nrows=2)
5
6 g1=sns.regplot(data=df_test, x="sqft_living_st", y="price", color="#003300", fit_reg=True,
7     ax=axes[0,0], line_kws={"color": "red", "lw":5});
8 g2=sns.regplot(data=df_test, x="distance_st", y="price", color="#000066",
9     ax=axes[0,1], line_kws={"color": "red", "lw":5});
10 g3=sns.regplot(data=df_test, x="bathrooms_st", y="price", color="#009900",
11     ax=axes[1,0], line_kws={"color": "red", "lw":5});
12 g4=sns.regplot(data=df_test, x="grade_st", y="price", color="#0000ff",
13     ax=axes[1,1], line_kws={"color": "red", "lw":5});
14
15 axes[0,0].set_title("Price vs Living Space", fontsize=26);
16 axes[0,0].set_ylabel('Price', fontsize=20)
17 axes[0,0].set_xlabel('Living Footage, adjusted', fontsize=20)
18 axes[0,0].set_xlim(-0.01, 1.0)
19 ylabels = ['{:,.1f}'.format(x) + 'M' for x in g1.get_yticks()/1000000]
20 axes[0,0].set_yticklabels(ylabels, size=15)
21 axes[0,0].grid(color='lightgrey')
22
23 axes[0,1].set_title("Price vs Distance from City Center", fontsize=26);
24 axes[0,1].set_ylabel('Price', fontsize=20)
25 axes[0,1].set_xlabel('Distance, adjusted', fontsize=20)
26 axes[0,1].set_xlim(-0.01, 1.0)
27 ylabels = ['{:,.1f}'.format(x) + 'M' for x in g2.get_yticks()/1000000]
28 axes[0,1].set_yticklabels(ylabels, size=15)
29 axes[0,1].grid(color='lightgrey')
30
31 axes[1,0].set_title("Price vs Number of Bathrooms", fontsize=26);
32 axes[1,0].set_ylabel('Price', fontsize=20)
33 axes[1,0].set_xlabel('Number of Bathrooms, adjusted', fontsize=20)
34 axes[1,0].set_xlim(-0.01, 1.0)
35 ylabels = ['{:,.1f}'.format(x) + 'M' for x in g3.get_yticks()/1000000]
36 axes[1,0].set_yticklabels(ylabels, size=15)
37 axes[1,0].grid(color='lightgrey')
38
39 axes[1,1].set_title("Price vs Building Grade", fontsize=26);
40 axes[1,1].set_ylabel('Price', fontsize=20)
41 axes[1,1].set_xlabel('Building Grade, adjusted', fontsize=20)
42 axes[1,1].set_xlim(-0.01, 1.0)
43 ylabels = ['{:,.1f}'.format(x) + 'M' for x in g4.get_yticks()/1000000]
44 axes[1,1].set_yticklabels(ylabels, size=15)
45 axes[1,1].grid(color='lightgrey')
46
47
48 plt.suptitle("Regression plots of Price vs Four Independent Variables", size=30, c="Blue")
49 plt.tight_layout(pad=3)
50
```

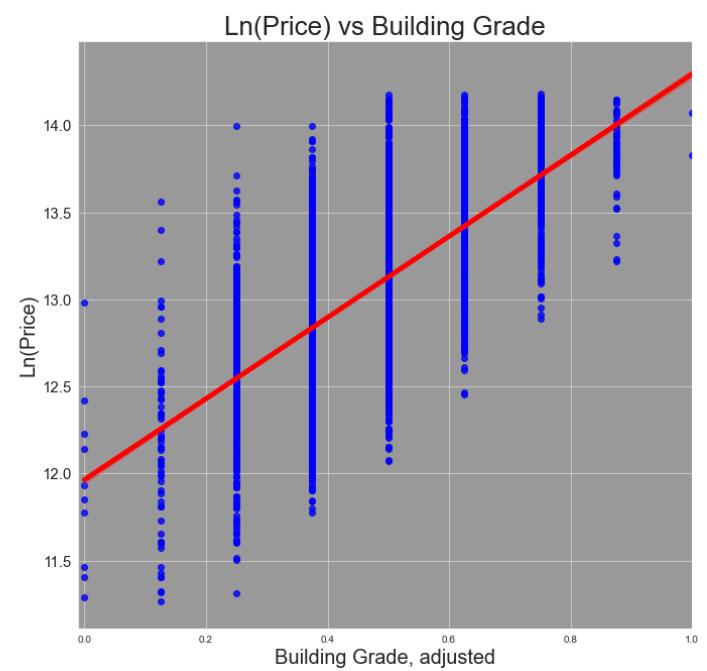
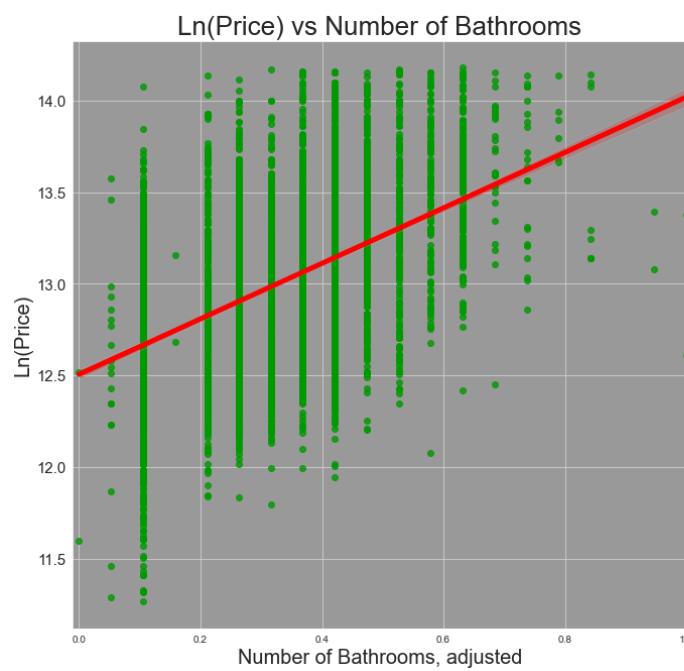
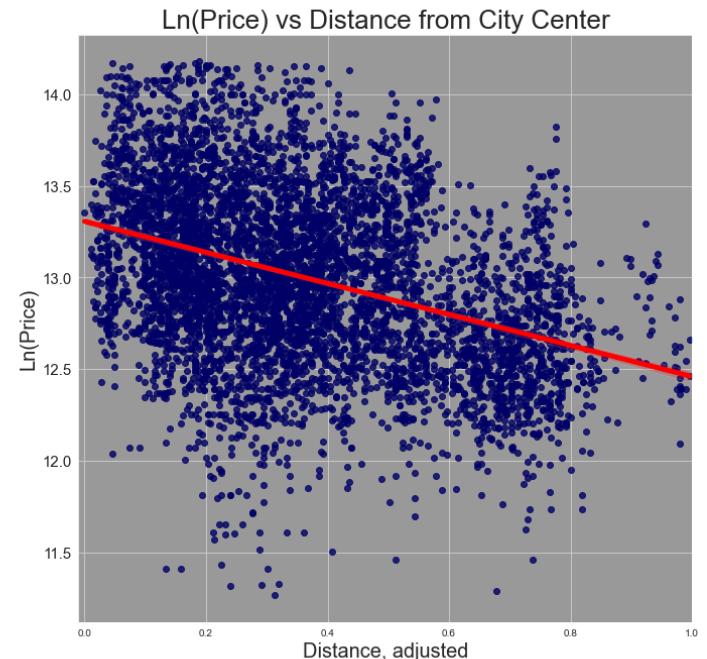
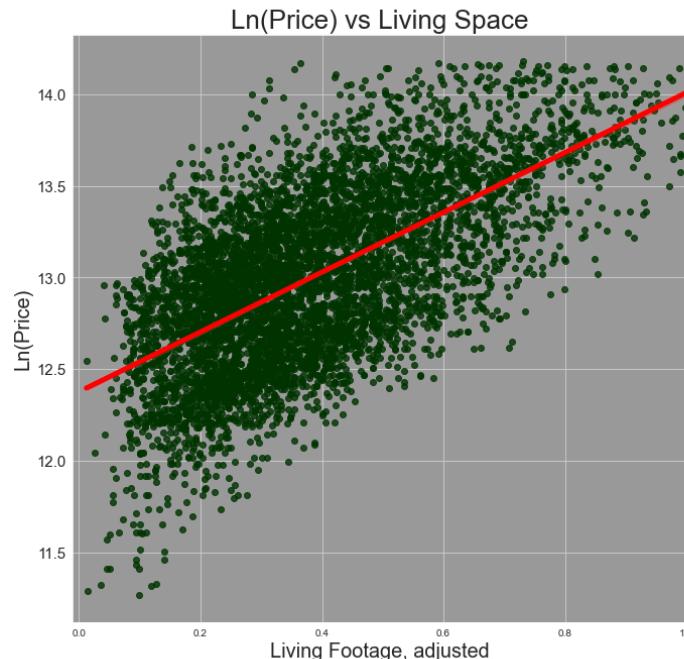
## Regression plots of Price vs Four Independent Variables



In [7]:

```
1 # Regplots for all four variables
2
3 sns.set_style("darkgrid", {"axes.facecolor": ".6"})
4 fig, axes = plt.subplots(figsize=(20,20), ncols=2, nrows=2)
5
6 g1=sns.regplot(data=df_test, x="sqft_living_st", y="log_price", color="#003300", fit_reg=True,
7                  ax=axes[0,0], line_kws={"color": "red", "lw":5});
8 g2=sns.regplot(data=df_test, x="distance_st", y="log_price", color="#000066",
9                  ax=axes[0,1], line_kws={"color": "red", "lw":5});
10 g3=sns.regplot(data=df_test, x="bathrooms_st", y="log_price", color="#009900",
11                  ax=axes[1,0], line_kws={"color": "red", "lw":5});
12 g4=sns.regplot(data=df_test, x="grade_st", y="log_price", color="#0000ff",
13                  ax=axes[1,1], line_kws={"color": "red", "lw":5});
14
15 axes[0,0].set_title("Ln(Price) vs Living Space", fontsize=26);
16 axes[0,0].set_ylabel('Ln(Price)', fontsize=20)
17 axes[0,0].set_xlabel('Living Footage, adjusted', fontsize=20)
18 axes[0,0].set_xlim(-0.01, 1.0)
19 ylabels = ['{:.1f}'.format(x) for x in g1.get_yticks()]
20 axes[0,0].set_yticklabels(ylabels, size=15)
21 axes[0,0].grid(color='lightgrey')
22
23 axes[0,1].set_title("Ln(Price) vs Distance from City Center", fontsize=26);
24 axes[0,1].set_ylabel('Ln(Price)', fontsize=20)
25 axes[0,1].set_xlabel('Distance, adjusted', fontsize=20)
26 axes[0,1].set_xlim(-0.01, 1.0)
27 ylabels = ['{:.1f}'.format(x) for x in g2.get_yticks()]
28 axes[0,1].set_yticklabels(ylabels, size=15)
29 axes[0,1].grid(color='lightgrey')
30
31 axes[1,0].set_title("Ln(Price) vs Number of Bathrooms", fontsize=26);
32 axes[1,0].set_ylabel('Ln(Price)', fontsize=20)
33 axes[1,0].set_xlabel('Number of Bathrooms, adjusted', fontsize=20)
34 axes[1,0].set_xlim(-0.01, 1.0)
35 ylabels = ['{:.1f}'.format(x) for x in g3.get_yticks()]
36 axes[1,0].set_yticklabels(ylabels, size=15)
37 axes[1,0].grid(color='lightgrey')
38
39 axes[1,1].set_title("Ln(Price) vs Building Grade", fontsize=26);
40 axes[1,1].set_ylabel('Ln(Price)', fontsize=20)
41 axes[1,1].set_xlabel('Building Grade, adjusted', fontsize=20)
42 axes[1,1].set_xlim(-0.01, 1.0)
43 ylabels = ['{:.1f}'.format(x) for x in g4.get_yticks()]
44 axes[1,1].set_yticklabels(ylabels, size=15)
45 axes[1,1].grid(color='lightgrey')
46
47 plt.suptitle("Regression plots of Ln(Price) vs Four Independent Variables", size=30, c="Blue")
48 plt.tight_layout(pad=3)
```

## Regression plots of Ln(Price) vs Four Independent Variables

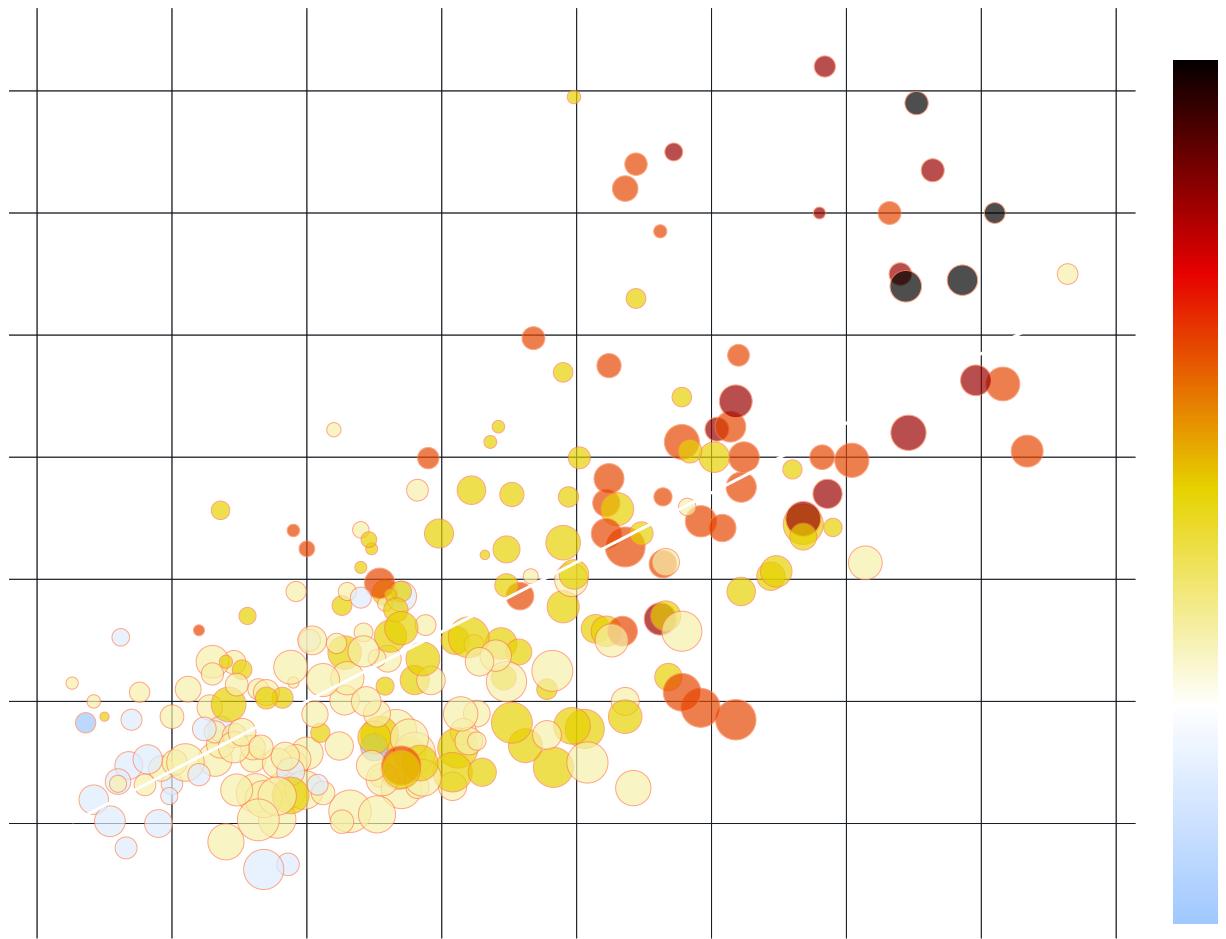


```
In [8]: 1 df_test_sample=df_test.sample(n=250, random_state=123)
```

```

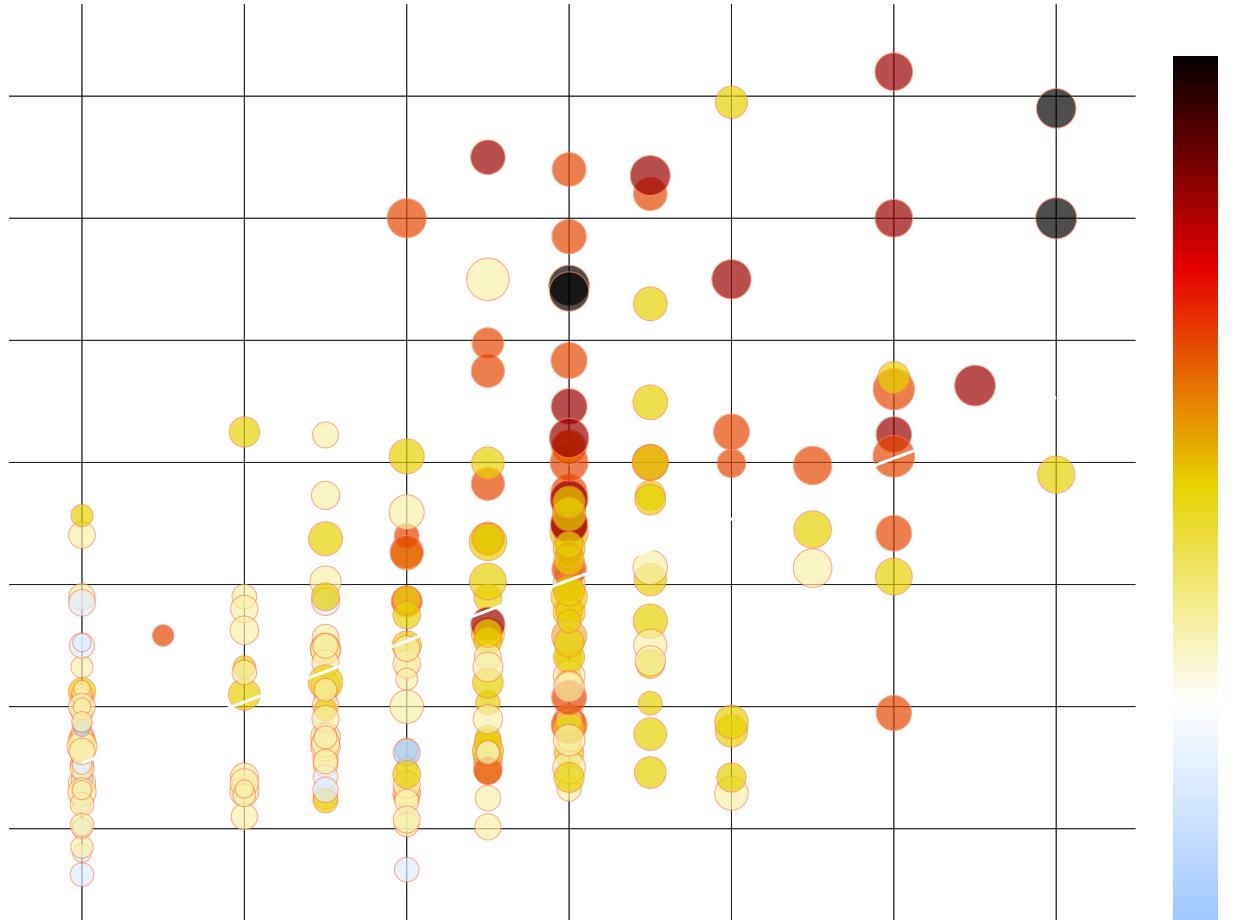
In [9]: 1 fig = px.scatter(df_test_sample, x='sqft_living', y='price', trendline='ols', trendline_color_override='white',
2                     color='grade', size='distance', width=1000, height=800, size_max=20,
3                     color_continuous_scale=px.colors.sequential.Blackbody_r,
4                     labels={
5                         "price": "Price",
6                         "sqft_living": "Living Space (sq ft)",
7                         "grade": "Building Grade"
8                     },
9                     title="Correlation: Property Price vs Living Space Footage  
Sized by Distance", template
10
11 fig.update_traces(marker=dict(
12     line=dict(
13         color='coral',
14         width=0.5
15     )))
16 fig.update_layout(
17     font_family="Arial",
18     font_size=18,
19     font_color="white",
20     title_font_family="Arial",
21     title_font_color="white",
22 )
23 )
24 fig.update_layout(
25     title={
26         'y':0.95,
27         'x':0.5,
28         'xanchor': 'center',
29         'yanchor': 'top',
30     })
31
32
33 fig.show()

```



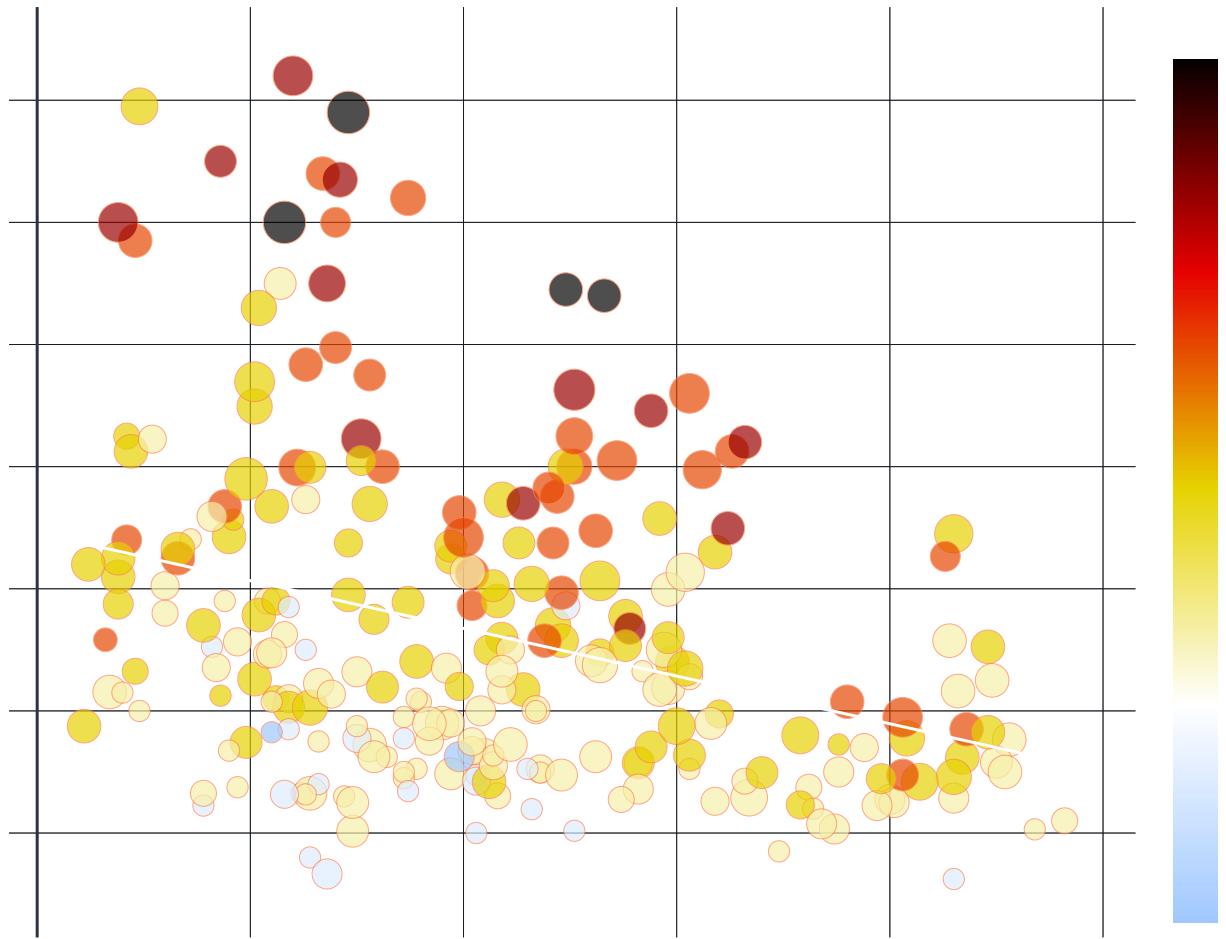


```
In [10]: 1 fig = px.scatter(df_test_sample, x='bathrooms', y='price', color_continuous_scale=px.colors.sequential.Blackbody_r,
2                     size='sqft_living', size_max=20,
3                     trendline='ols', trendline_color_override='white', color='grade',
4                     width=1000, height=800, labels={
5                         "price": "Price",
6                         "bathrooms": "Number of bathrooms",
7                         "grade": "Building Grade"
8                     },
9                     title="Correlation: Property Price vs Number of Bathrooms  
Sized by Living Space",
10                    template = "plotly_dark")
11
12
13 fig.update_traces(marker=dict(
14     line=dict(
15         color='coral',
16         width=0.5
17     )))
18 fig.update_layout(
19     font_family="Arial",
20     font_size=18,
21     font_color="white",
22     title_font_family="Arial",
23     title_font_color="white",
24 )
25
26
27 fig.update_layout(
28     title={
29         'y':0.95,
30         'x':0.5,
31         'xanchor': 'center',
32         'yanchor': 'top',
33     })
34
35
36 fig.show()
```





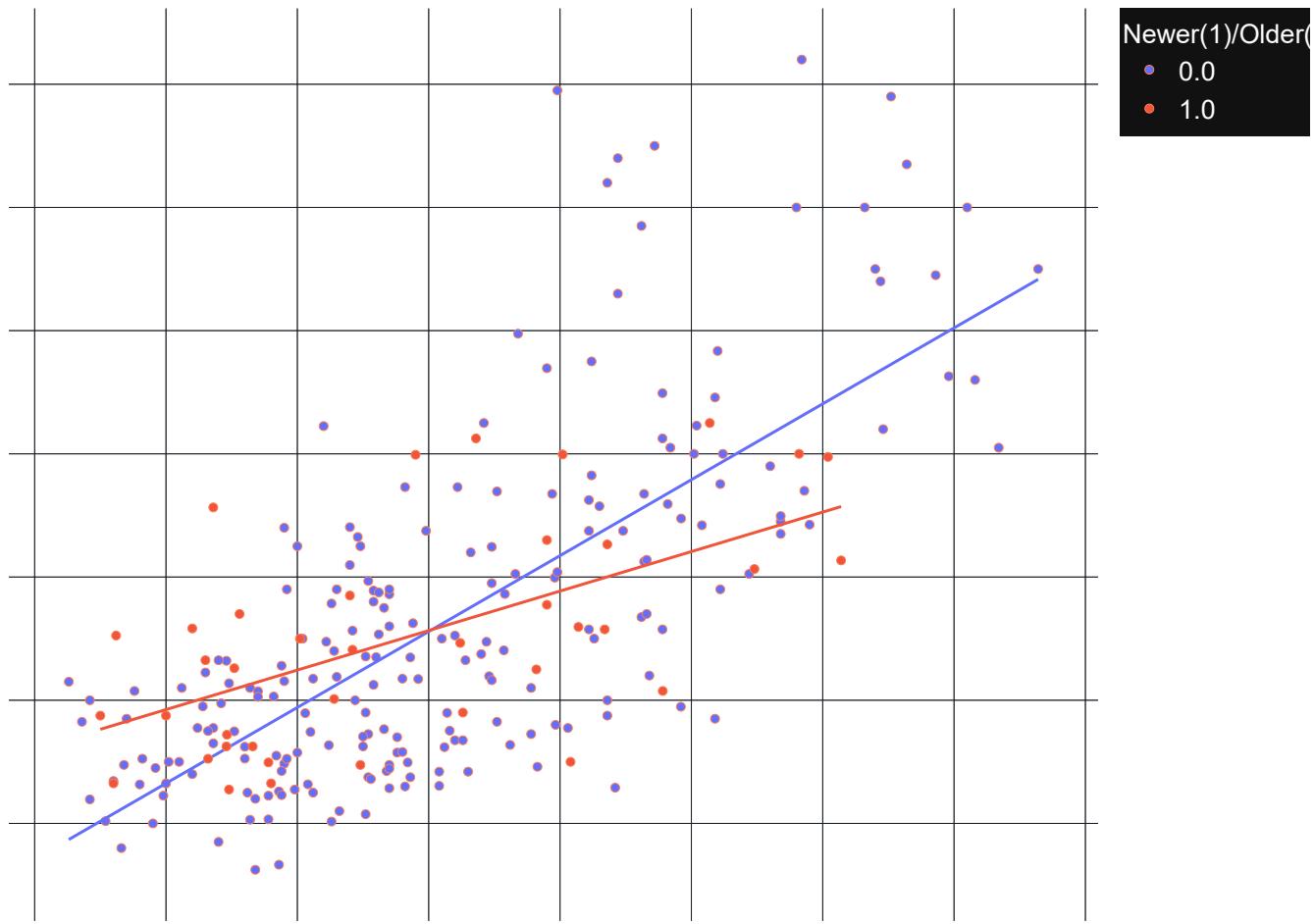
```
In [11]: 1 fig = px.scatter(df_test_sample, x='distance', y='price', trendline='ols', trendline_color_override='white',
2                     color='grade', size='bathrooms', width=1000, height=800, size_max=20,
3                     color_continuous_scale=px.colors.sequential.Blackbody_r,
4                     labels={
5                         "price": "Price",
6                         "distance": "Distance",
7                         "grade": "Building Grade"
8                     },
9                     title="Correlation: Property Price vs Distance from the City Center  
Sized by Number of Bathrooms")
10
11 fig.update_traces(marker=dict(
12     line=dict(
13         color='coral',
14         width=0.5
15     )))
16 fig.update_layout(
17     font_family="Arial",
18     font_size=18,
19     font_color="white",
20     title_font_family="Arial",
21     title_font_color="white",
22 )
23 )
24 fig.update_layout(
25     title={
26         'y':0.95,
27         'x':0.5,
28         'xanchor': 'center',
29         'yanchor': 'top',
30     })
31
32
33 fig.show()
```





In [12]:

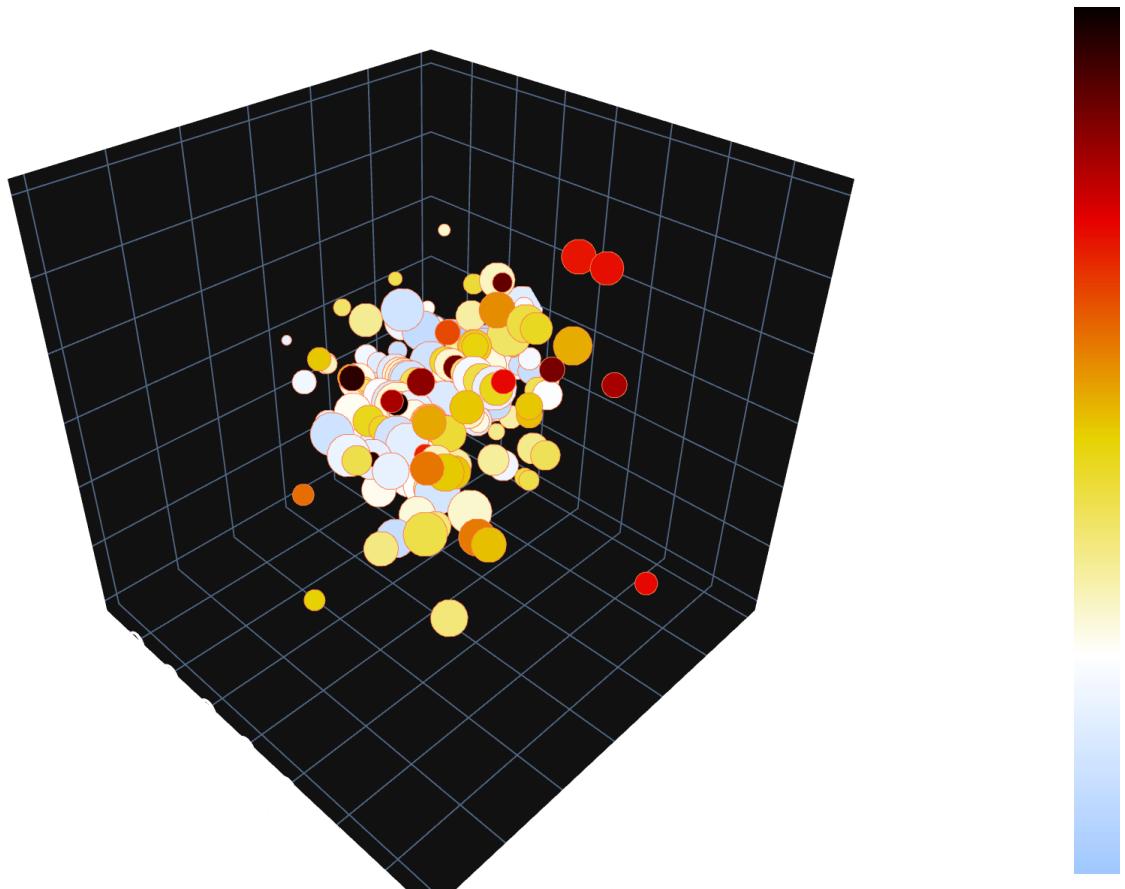
```
1 fig = px.scatter(df_test_sample, x='sqft_living', y='price', trendline='ols',
2                   color='recent_renovation_new_str', width=1000, height=800, size_max=20,
3                   labels={
4                     "price": "Price",
5                     "sqft_living": "Living Space (sq ft)",
6                     "recent_renovation_new_str": "Newer(1)/Older(0)"
7                   },
8                   title="Correlation: Property Price vs Living Space Footage of Newer vs Older Properties",
9                   template = "plotly_dark")
10
11 fig.update_traces(marker=dict(
12   line=dict(
13     color='coral',
14     width=0.5
15   )))
16 fig.update_layout(
17   font_family="Arial",
18   font_size=18,
19   font_color="white",
20   title_font_family="Arial",
21   title_font_color="white",
22 )
23 )
24 fig.update_layout(
25   title={
26     'y':0.95,
27     'x':0.5,
28     'xanchor': 'center',
29     'yanchor': 'top',
30   })
31
32
33 fig.show()
```





In [13]:

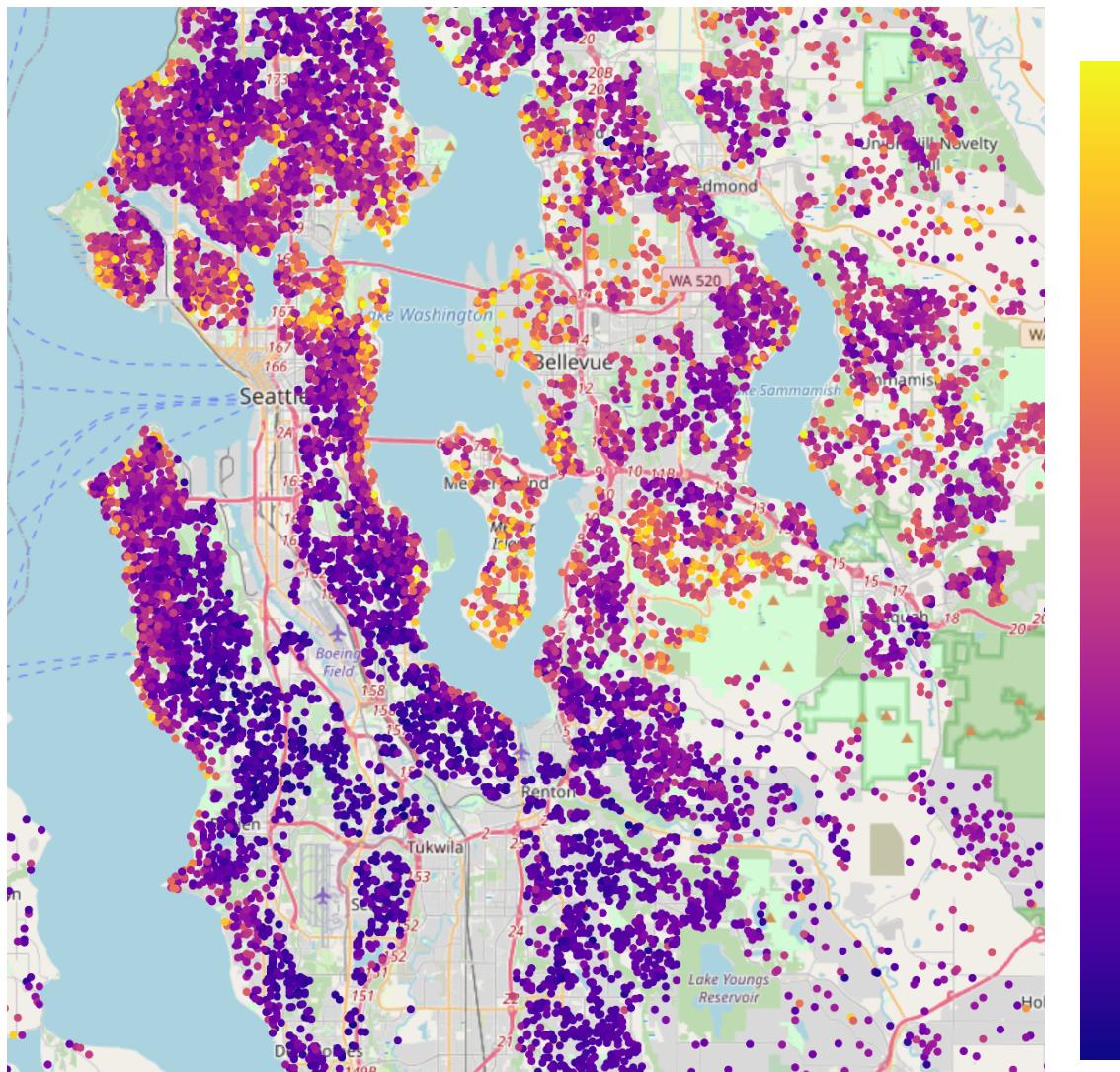
```
1 fig = px.scatter_3d(df_test_sample, x='bathrooms', z='grade', y='sqft_living',
2                     color='price', size='distance', size_max=50, opacity=1, width=1000, height=800,
3                     color_continuous_scale=px.colors.sequential.Blackbody_r,
4                     labels={
5                         "bathrooms": "Number of Bathrooms",
6                         "sqft_living": "Living Space (sq ft)",
7                         "grade": "Grade",
8                         "price": "Price"
9                     },
10                    title="3D plot: Living Space Footage, Number of Bathrooms and Grade of Sold Properties",
11                    template = "plotly_dark")
12
13 fig.update_traces(marker=dict(
14     line=dict(
15         color='coral',
16         width=0.5
17     )))
18 fig.update_layout(
19     font_family="Arial",
20     font_size=16,
21     font_color="white",
22     title_font_family="Arial",
23     title_font_color="white"
24 )
25 fig.update_layout(
26     title={
27         'y':0.9,
28         'x':0.5,
29         'xanchor': 'center',
30         'yanchor': 'top',
31     })
32
33
34 fig.show()
```



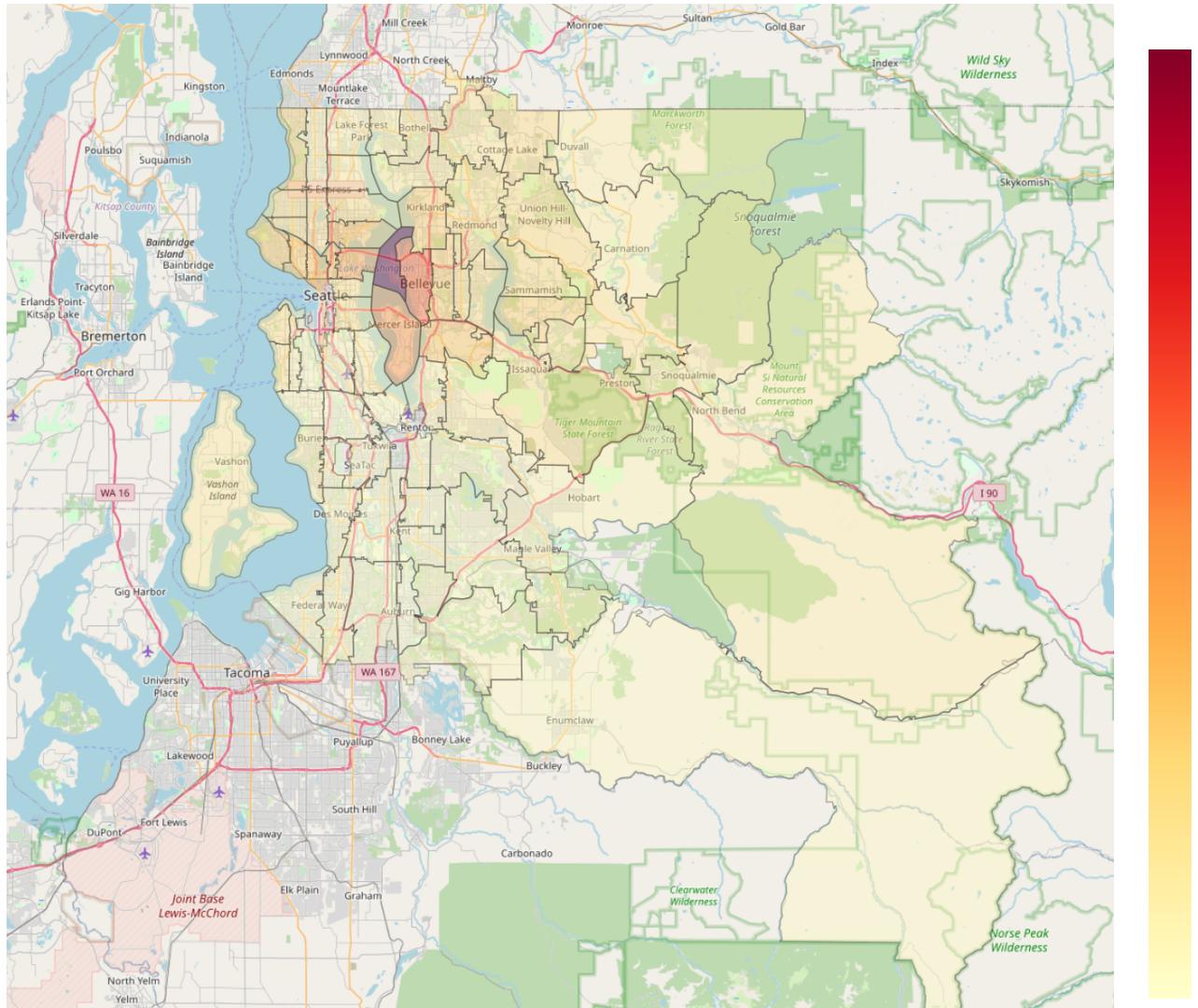


In [14]:

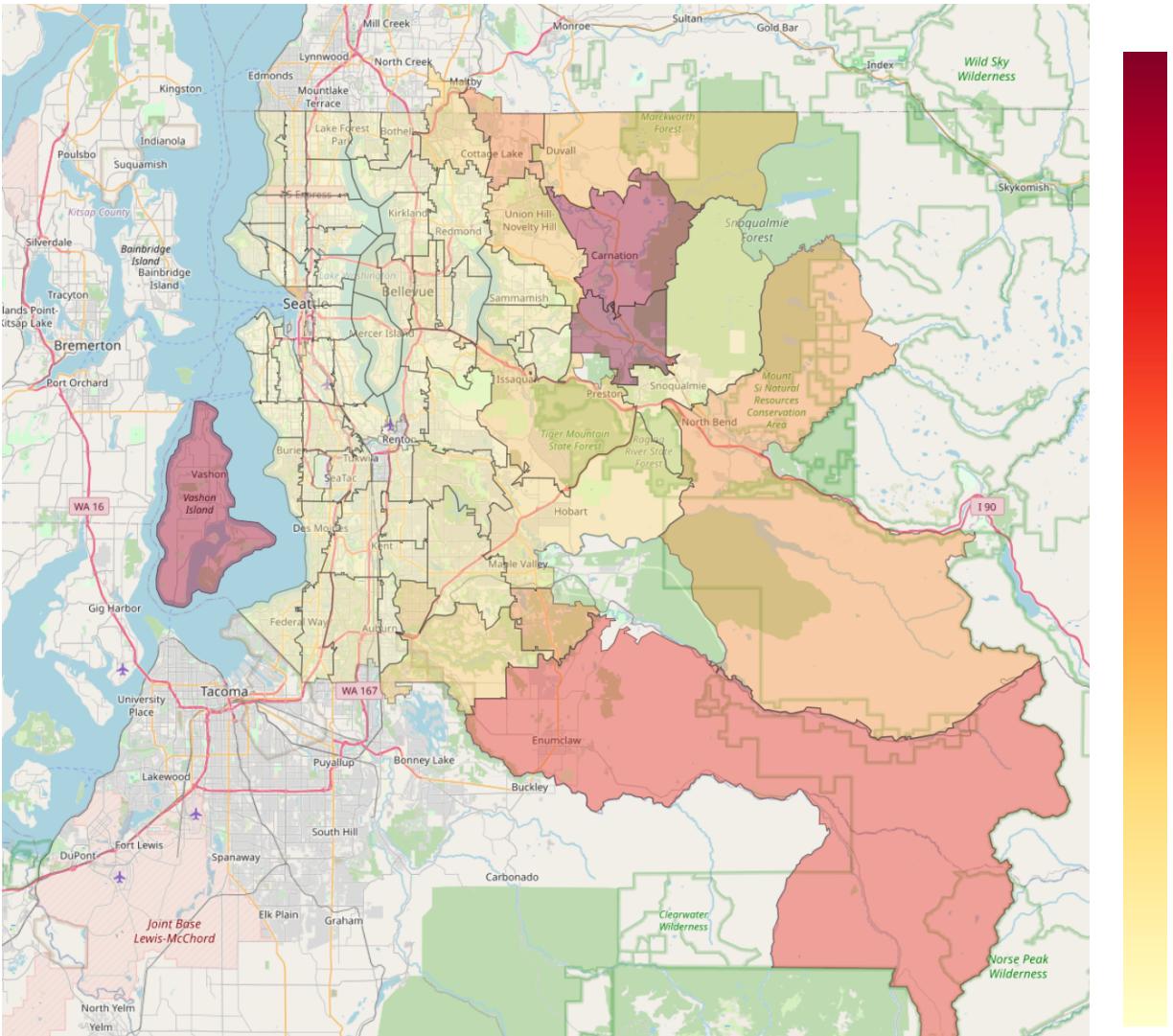
```
1 df_cut_off=df.copy()
2 df_cut_off=df_cut_off[(df_cut_off['price']<1500000)]
3
4 fig = px.scatter_mapbox(df_cut_off, lat="lat", lon="long", color="price",
5                         color_continuous_scale=px.colors.sequential.Plasma, zoom=10,
6                         mapbox_style='open-street-map', width=900, height=900,
7                         title="Properties Sold in King County in 2014-2015",
8                         template = "plotly_dark")
9
10 fig.update_layout(
11     font_family="Arial",
12     font_size=20,
13     font_color="white",
14     title_font_family="Arial",
15     title_font_color="white"
16 )
17 fig.update_layout(
18     title={
19         'y':0.98,
20         'x':0.5,
21         'xanchor': 'center',
22         'yanchor': 'top',
23     })
24
25
26 fig.show()
```



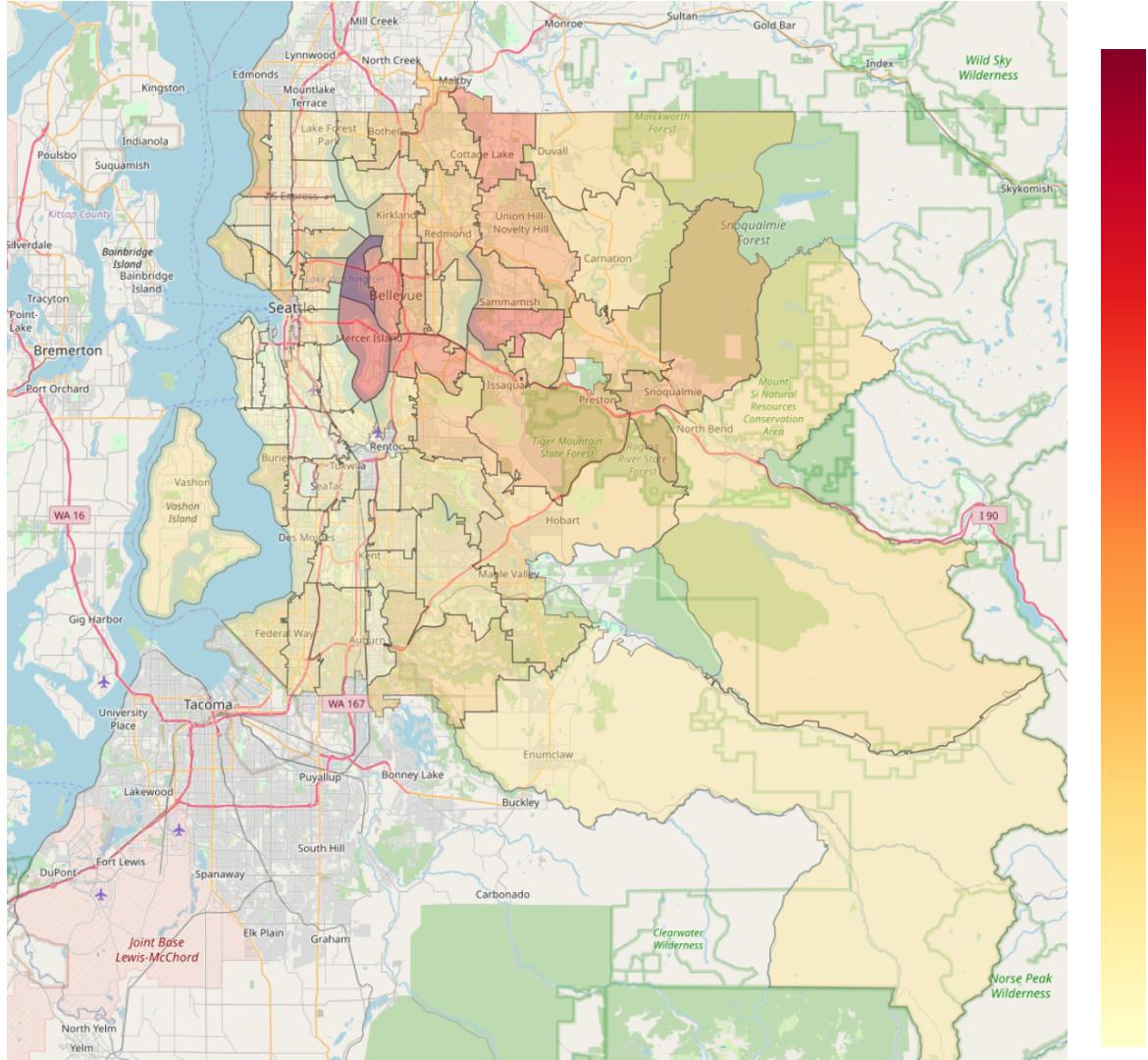
```
In [15]: 1 map_choropleth_zip(df_zipcode_viz, 'price', "Average Prices of Sold Properties per Zipcode (King County, 2014-2015)",  
2 "Price")
```



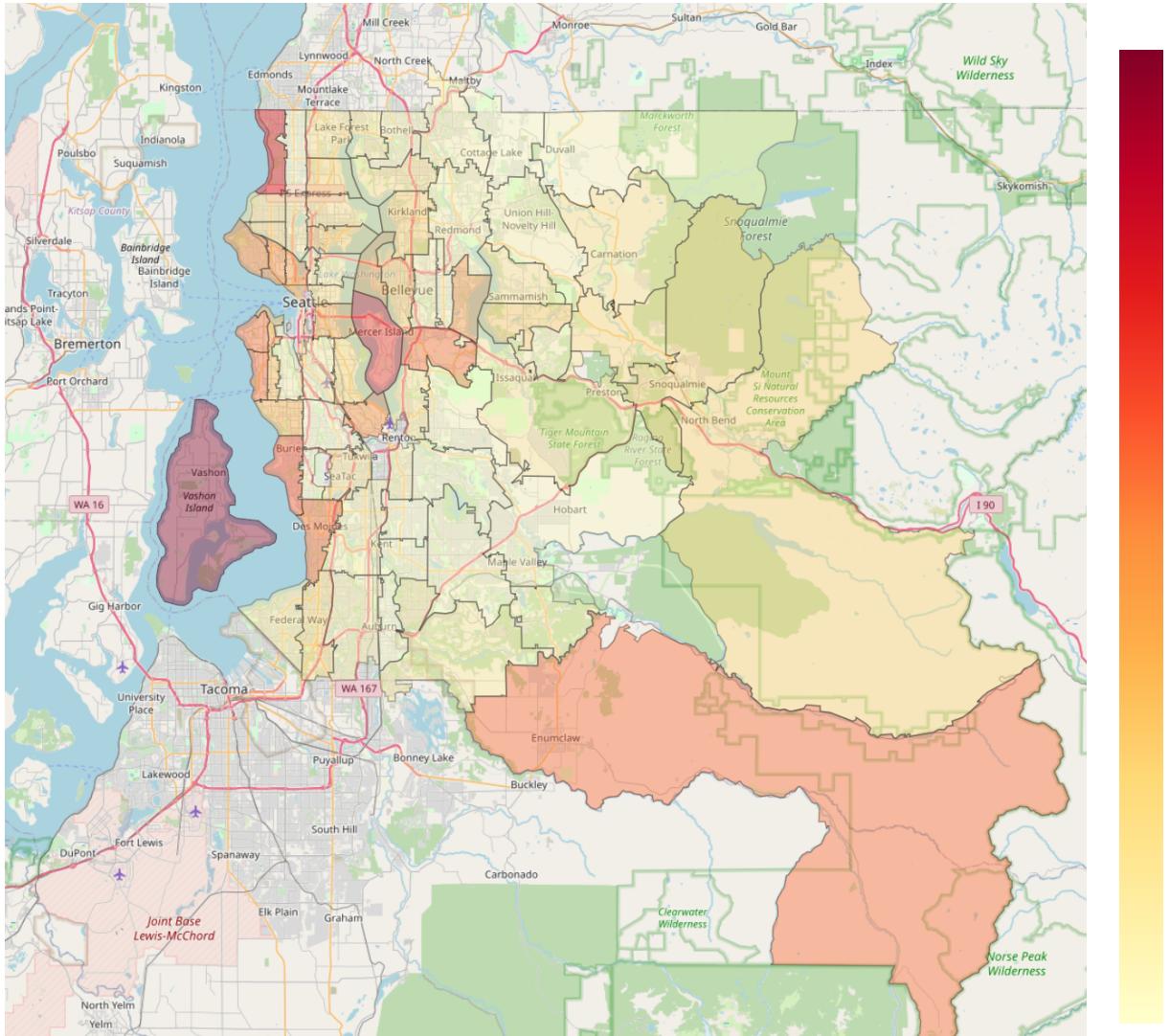
```
In [16]: 1 map_choropleth_zip(df_zipcode_viz, 'sqft_lot', "Average Lot Size of Sold Properties per Zipcode (King County, 2014-2015)"  
2 "Lot size (sq ft)")
```



```
In [17]: 1 map_choropleth_zip(df_zipcode_viz, 'sqft_living', "Average Living Space of Sold Properties per Zipcode (King County, 2014  
2           "Living Space (sq ft) ")
```



```
In [18]: 1 map_choropleth_zip(df_zipcode_viz, 'view', "Average View Category of Sold Properties per Zipcode (King County, 2014-2015)  
2      'View Category")
```



```
In [19]: 1 map_choropleth_zip(df_zipcode_viz, 'yr_built', "Average Year Built of Sold Properties per Zipcode (King County, 2014-2015  
2 "Year Built")
```

