

Exploratory Data Analysis

This notebook focuses the EDA on sold condos and the data scraped/processed from Redfin.

```
In [1]: 1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib
5 import matplotlib.pyplot as plt
6
7 color = sns.color_palette()
8 %matplotlib inline
```

```
In [2]: 1 def price_vs_house_feature(df, feature, subtitle, ax):
2     sns.regplot(x=df[feature], y=df.LOG_SOLDPRICE, scatter_kws={'alpha':0.2}, color=color[0], ax=ax)
3     # set titles
4     ax.set_title("Log Sold Price vs. " + feature + "\n" + subtitle, fontsize=16)
5     # label axes
6     ax.set_ylabel("Log Sold Price", fontsize=14)
7     ax.set_xlabel(feature, fontsize=14)
8
9 def dom_vs_house_feature(df, feature, subtitle, ax):
10     sns.regplot(x=df[feature], y=df.LOG_DOM, scatter_kws={'alpha':0.2}, color=color[0], ax=ax)
11     # set titles
12     ax.set_title("Log DOM vs. " + feature + "\n" + subtitle, fontsize=16)
13     # label axes
14     ax.set_ylabel("Log DOM", fontsize=14)
15     ax.set_xlabel(feature, fontsize=14)
16
17 def plot_sold_price_distribution(df, subtitle):
18     fig, ax = plt.subplots(1, 2, figsize=(14, 5))
19     sns.distplot(df.SOLDPRICE, kde=False, ax=ax[0])
20     sns.distplot(df.LOG_SOLDPRICE, kde=False, ax=ax[1])
21
22     ax[0].set_title("Distribution of Sold Price \n" + subtitle, fontsize=16)
23     ax[1].set_title("Distribution of Log Sold Price \n" + subtitle, fontsize=16)
24
25     ax[0].set_xlabel("Sold Price", fontsize=14)
26     ax[1].set_xlabel("Log Sold Price", fontsize=14)
27     ax[0].set_ylabel("Frequency", fontsize=14)
28     ax[1].set_ylabel("Frequency", fontsize=14)
29
30 def plot_dom_distribution(df, subtitle):
31     fig, ax = plt.subplots(1, 2, figsize=(14, 5))
32     sns.distplot(df.DOM, kde=False, ax=ax[0])
33     sns.distplot(df.LOG_DOM, kde=False, ax=ax[1])
34
35     ax[0].set_title("Distribution of DOM \n" + subtitle, fontsize=16)
36     ax[1].set_title("Distribution of Log (Log DOM) \n" + subtitle, fontsize=16)
37
38     ax[0].set_xlabel("DOM", fontsize=14)
39     ax[1].set_xlabel("Log DOM", fontsize=14)
40     ax[0].set_ylabel("Frequency", fontsize=14)
41     ax[1].set_ylabel("Frequency", fontsize=14)
42
```

```
In [3]: 1 df_con = pd.read_csv("data/features/CON_feats_no_img.csv", index_col=0)
2 df_con = df_con.fillna(-1)
3 df_con['LOG_SOLDPRICE'] = np.log(df_con['SOLDPRICE'])
4 df_con['LOG_DOM'] = np.log(np.log(df_con['DOM']+20))
5 df_con['log_sqft_finished'] = np.log(df_con['sqft_finished']+2)
6 df_con['log_sqft_unfinished'] = np.log(df_con['sqft_unfinished']+2)
7 df_con.head()
```

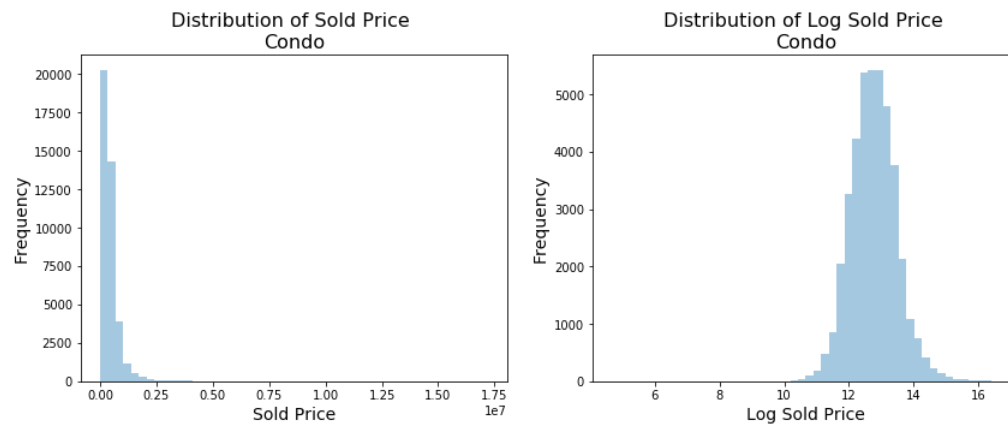
```
Out[3]:
```

	MLSNUM	LISTPRICE	SOLDPRICE	DOM	DTO	AGE	LOTSIZE	GARAGE	LISTMONTH	SOLDMONTH	...	school_distances_high_min	walk_score	transit_score	bike_score	nu
0	71498924	169900.0	177500.0	709	618.0	11.0	-1.0	2.0	3	1	...	0.1	34.0	-1.0	-1.0	
1	71905628	242000.0	235000.0	91	81.0	14.0	-1.0	1.0	9	1	...	0.3	31.0	-1.0	-1.0	
2	71918879	209000.0	209000.0	78	37.0	32.0	-1.0	0.0	10	1	...	0.5	14.0	-1.0	-1.0	
3	71952614	339900.0	350695.0	1	1.0	3.0	-1.0	1.0	1	1	...	0.9	17.0	-1.0	-1.0	
4	71912071	299900.0	280000.0	83	71.0	10.0	-1.0	1.0	9	2	...	0.9	16.0	-1.0	-1.0	

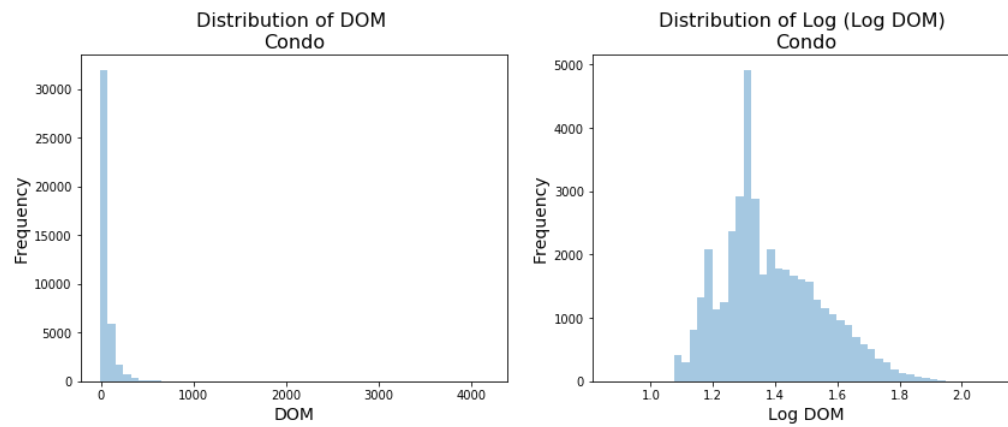
5 rows × 42 columns

Response Distribution

```
In [4]: 1 plot_sold_price_distribution(df_con, "Condo")
```



```
In [5]: 1 plot_dom_distribution(df_con, "Condo")
```



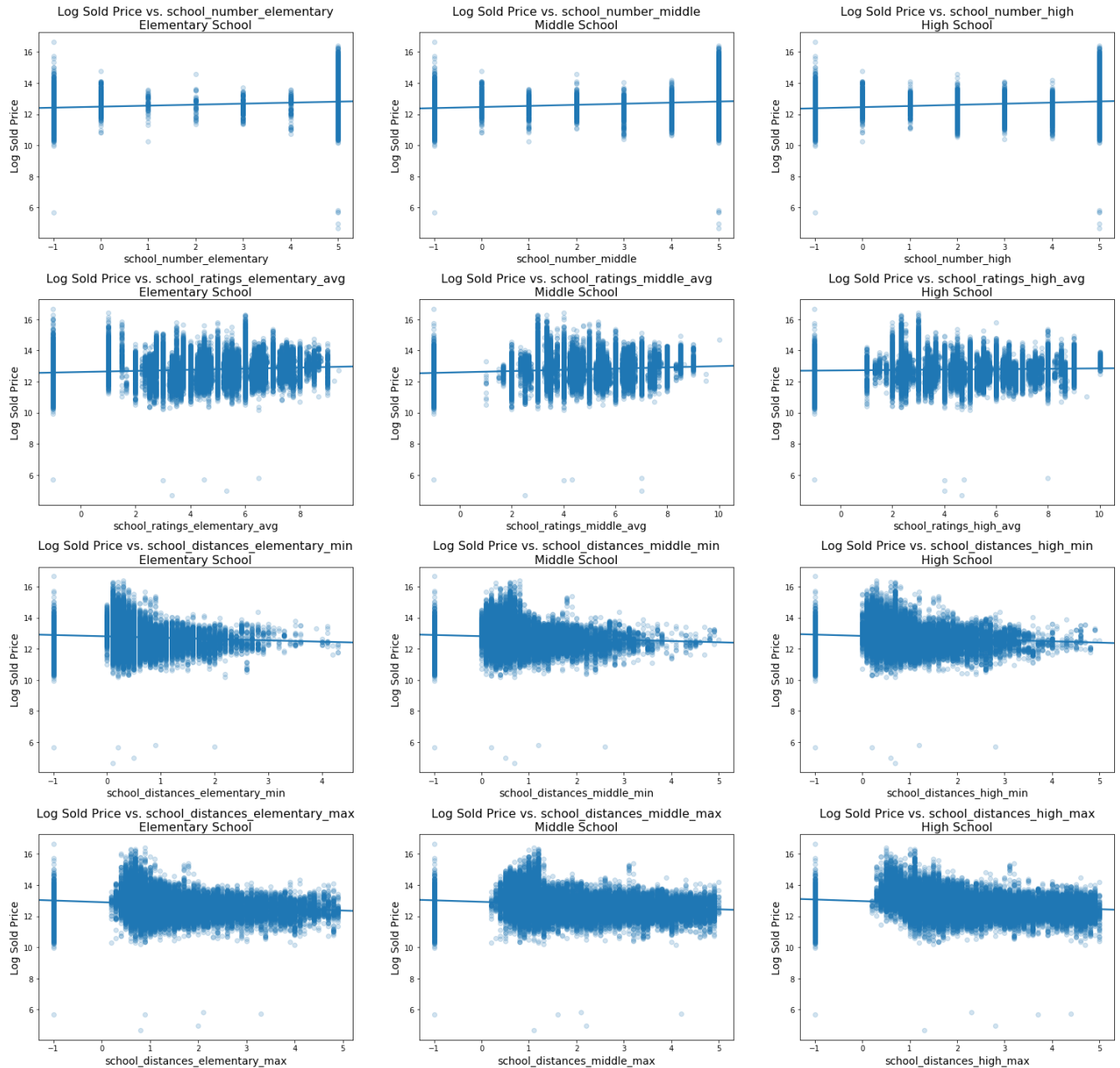
Price Trends

Price vs. School information

```

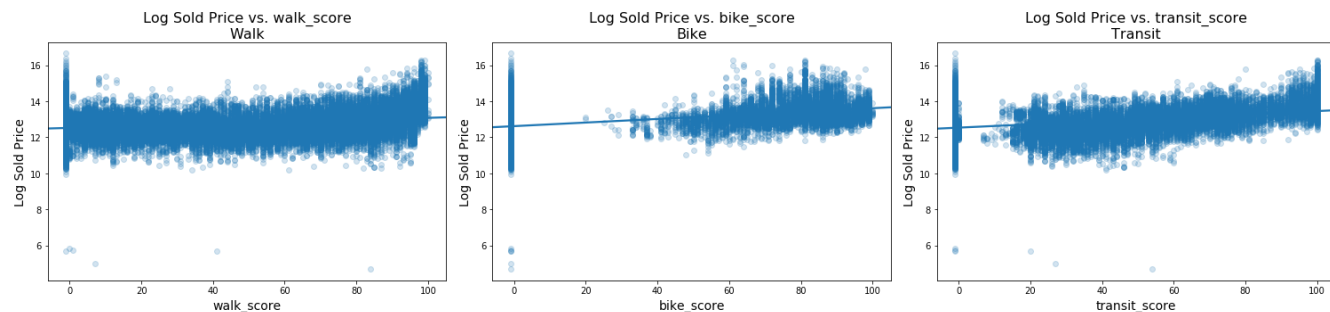
In [6]: 1 # price vs school information
2 fig, ax = plt.subplots(4, 3, figsize=(21, 20))
3 # number of schools
4 price_vs_house_feature(df_con, 'school_number_elementary', 'Elementary School', ax=ax[0][0])
5 price_vs_house_feature(df_con, 'school_number_middle', 'Middle School', ax=ax[0][1])
6 price_vs_house_feature(df_con, 'school_number_high', 'High School', ax=ax[0][2])
7
8 # average school rating
9 price_vs_house_feature(df_con, 'school_ratings_elementary_avg', 'Elementary School', ax=ax[1][0])
10 price_vs_house_feature(df_con, 'school_ratings_middle_avg', 'Middle School', ax=ax[1][1])
11 price_vs_house_feature(df_con, 'school_ratings_high_avg', 'High School', ax=ax[1][2])
12
13 # min school distance
14 price_vs_house_feature(df_con, 'school_distances_elementary_min', 'Elementary School', ax=ax[2][0])
15 price_vs_house_feature(df_con, 'school_distances_middle_min', 'Middle School', ax=ax[2][1])
16 price_vs_house_feature(df_con, 'school_distances_high_min', 'High School', ax=ax[2][2])
17
18 # max school distance
19 price_vs_house_feature(df_con, 'school_distances_elementary_max', 'Elementary School', ax=ax[3][0])
20 price_vs_house_feature(df_con, 'school_distances_middle_max', 'Middle School', ax=ax[3][1])
21 price_vs_house_feature(df_con, 'school_distances_high_max', 'High School', ax=ax[3][2])
22 plt.tight_layout()

```



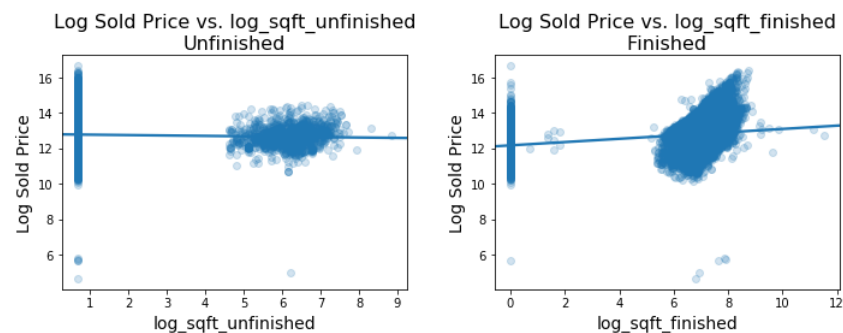
Price vs. Convenience scores

```
In [7]: 1 # price vs convenience scores
2 fig, ax = plt.subplots(1, 3, figsize=(21, 5))
3 # number of schools
4 price_vs_house_feature(df_con, 'walk_score', 'Walk', ax=ax[0])
5 price_vs_house_feature(df_con, 'bike_score', 'Bike', ax=ax[1])
6 price_vs_house_feature(df_con, 'transit_score', 'Transit', ax=ax[2])
7
8 plt.tight_layout()
```



Price vs. Square footage

```
In [8]: 1 # price vs sqft
2 fig, ax = plt.subplots(1, 2, figsize=(10, 4))
3 # number of schools
4 price_vs_house_feature(df_con, 'log_sqft_unfinished', 'Unfinished', ax=ax[0])
5 price_vs_house_feature(df_con, 'log_sqft_finished', 'Finished', ax=ax[1])
6
7 plt.tight_layout()
```



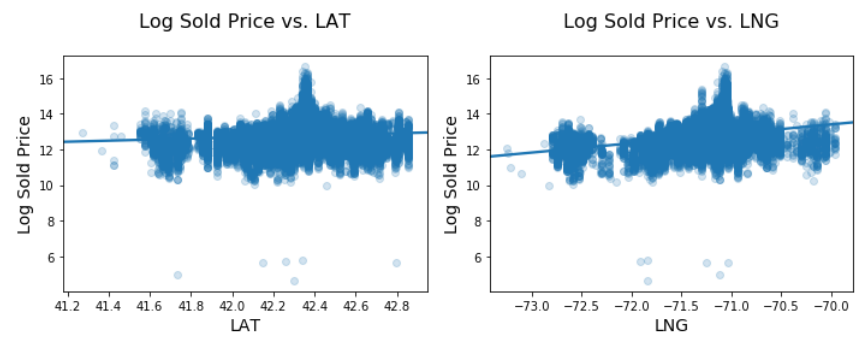
Price vs. number of photos on redfin

```
In [9]: 1 # price vs number of images posted on redfin
2 fig, ax = plt.subplots(1, 1, figsize=(5, 4))
3 # number of schools
4 price_vs_house_feature(df_con, 'num_photo', '', ax=ax)
5
6 plt.tight_layout()
```



Price vs. latitude/longitude

```
In [10]: 1 # price vs latitude/longitude
2 fig, ax = plt.subplots(1, 2, figsize=(10, 4))
3 # number of schools
4 price_vs_house_feature(df_con, 'LAT', '', ax=ax[0])
5 price_vs_house_feature(df_con, 'LNG', '', ax=ax[1])
6
7 plt.tight_layout()
```

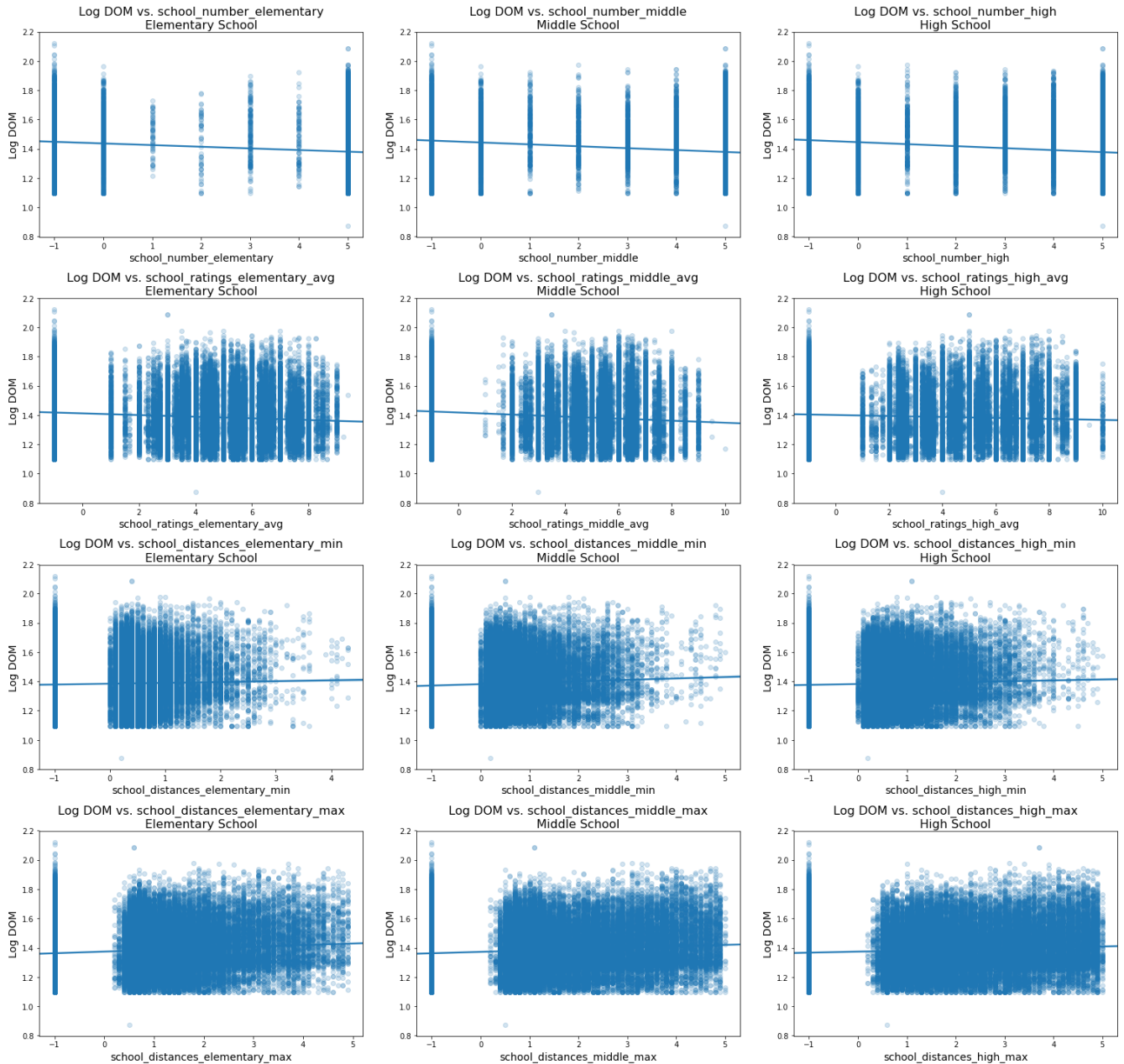


DOM Trend

DOM vs. School information

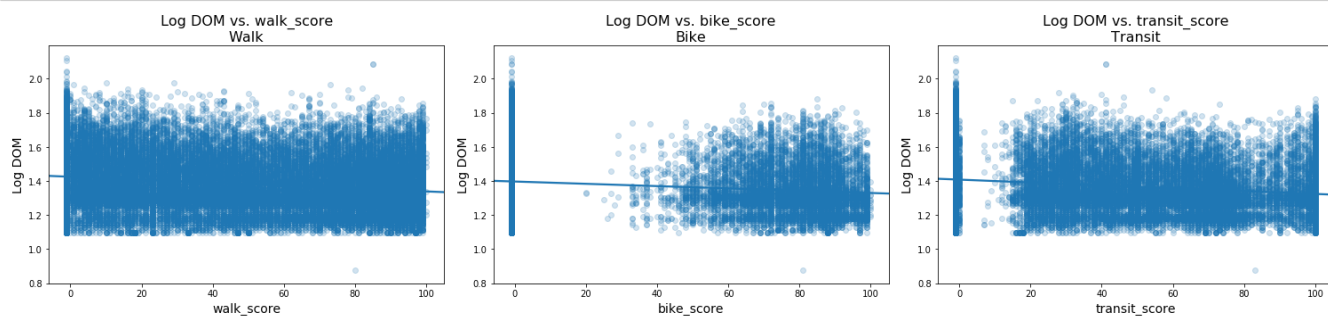
In [11]:

```
1 # dom vs school information
2 fig, ax = plt.subplots(4, 3, figsize=(21, 20))
3 # number of schools
4 dom_vs_house_feature(df_con, 'school_number_elementary', 'Elementary School', ax=ax[0][0])
5 dom_vs_house_feature(df_con, 'school_number_middle', 'Middle School', ax=ax[0][1])
6 dom_vs_house_feature(df_con, 'school_number_high', 'High School', ax=ax[0][2])
7
8 # average school rating
9 dom_vs_house_feature(df_con, 'school_ratings_elementary_avg', 'Elementary School', ax=ax[1][0])
10 dom_vs_house_feature(df_con, 'school_ratings_middle_avg', 'Middle School', ax=ax[1][1])
11 dom_vs_house_feature(df_con, 'school_ratings_high_avg', 'High School', ax=ax[1][2])
12
13 # min school distance
14 dom_vs_house_feature(df_con, 'school_distances_elementary_min', 'Elementary School', ax=ax[2][0])
15 dom_vs_house_feature(df_con, 'school_distances_middle_min', 'Middle School', ax=ax[2][1])
16 dom_vs_house_feature(df_con, 'school_distances_high_min', 'High School', ax=ax[2][2])
17
18 # max school distance
19 dom_vs_house_feature(df_con, 'school_distances_elementary_max', 'Elementary School', ax=ax[3][0])
20 dom_vs_house_feature(df_con, 'school_distances_middle_max', 'Middle School', ax=ax[3][1])
21 dom_vs_house_feature(df_con, 'school_distances_high_max', 'High School', ax=ax[3][2])
22 plt.tight_layout()
```



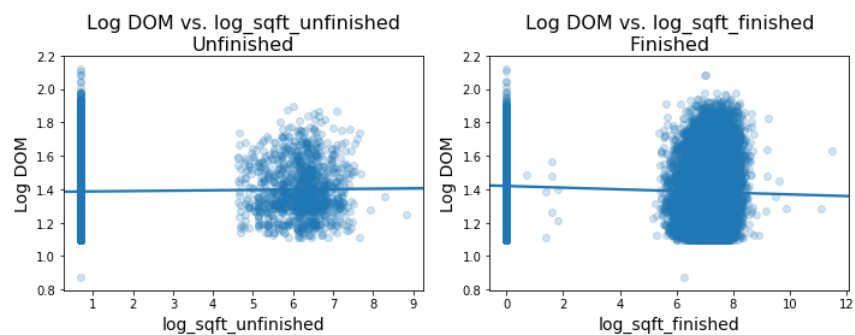
DOM vs. convenience scores

```
In [12]: 1 # dom vs convenience scores
2 fig, ax = plt.subplots(1, 3, figsize=(21, 5))
3 # number of schools
4 dom_vs_house_feature(df_con, 'walk_score', 'Walk', ax=ax[0])
5 dom_vs_house_feature(df_con, 'bike_score', 'Bike', ax=ax[1])
6 dom_vs_house_feature(df_con, 'transit_score', 'Transit', ax=ax[2])
7
8 plt.tight_layout()
```



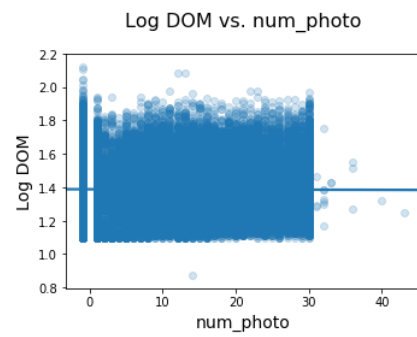
DOM vs. Square Footage

```
In [13]: 1 # dom vs sqft
2 fig, ax = plt.subplots(1, 2, figsize=(10, 4))
3 # number of schools
4 dom_vs_house_feature(df_con, 'log_sqft_unfinished', 'Unfinished', ax=ax[0])
5 dom_vs_house_feature(df_con, 'log_sqft_finished', 'Finished', ax=ax[1])
6
7 plt.tight_layout()
```

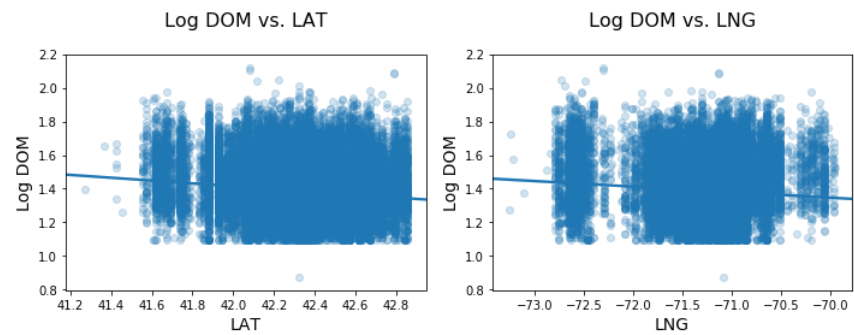


DOM vs. number of photos on redfin

```
In [14]: 1 # dom vs number of images posted on redfin
2 fig, ax = plt.subplots(1, 1, figsize=(5, 4))
3 # number of schools
4 dom_vs_house_feature(df_con, 'num_photo', '', ax=ax)
5
6 plt.tight_layout()
```



```
In [15]: 1 # price vs latitude/longitude
2 fig, ax = plt.subplots(1, 2, figsize=(10, 4))
3 # number of schools
4 dom_vs_house_feature(df_con, 'LAT', '', ax=ax[0])
5 dom_vs_house_feature(df_con, 'LNG', '', ax=ax[1])
6
7 plt.tight_layout()
```



```
In [ ]: 1
```