

```
In [1]: 1 # Put these at the top of every notebook, to get automatic reloading and inline
2 from IPython.core.display import display, HTML
3 import pandas as pd
4 import warnings
5 warnings.filterwarnings('ignore')
6
7 %reload_ext autoreload
8 %autoreload 1
9 %matplotlib inline
10
11 pd.set_option('display.max_rows', 500)
12 pd.set_option('display.max_columns', 500)
13 pd.set_option('display.width', 1000)
14
15 display(HTML("<style>.container { width:100% !important; }</style>"))
```

```
In [2]: 1 import os
2 import seaborn as sns
3 import pandas as pd
4 import math
5 import numpy as np
6
7 from Utils.UtilsGeoViz import *
8 from Utils.UtilsViz import *
9 from Utils.DataUtils import *
```

```
In [3]: 1 US_coord = [37.0902, -102]
2 NY_COORD = [40.7128, -74.0060]
3
4 # data_path = "C:\\Users\\sriharis\\OneDrive\\UChicago\\DataMining\\project\\
5 data_path = os.path.join(os.getcwd(), "../data/listings.csv")
6 listings = pd.read_csv(data_path, index_col="id")
7 PERCENTILE_CROP = [1,98]
8 display(listings.shape)

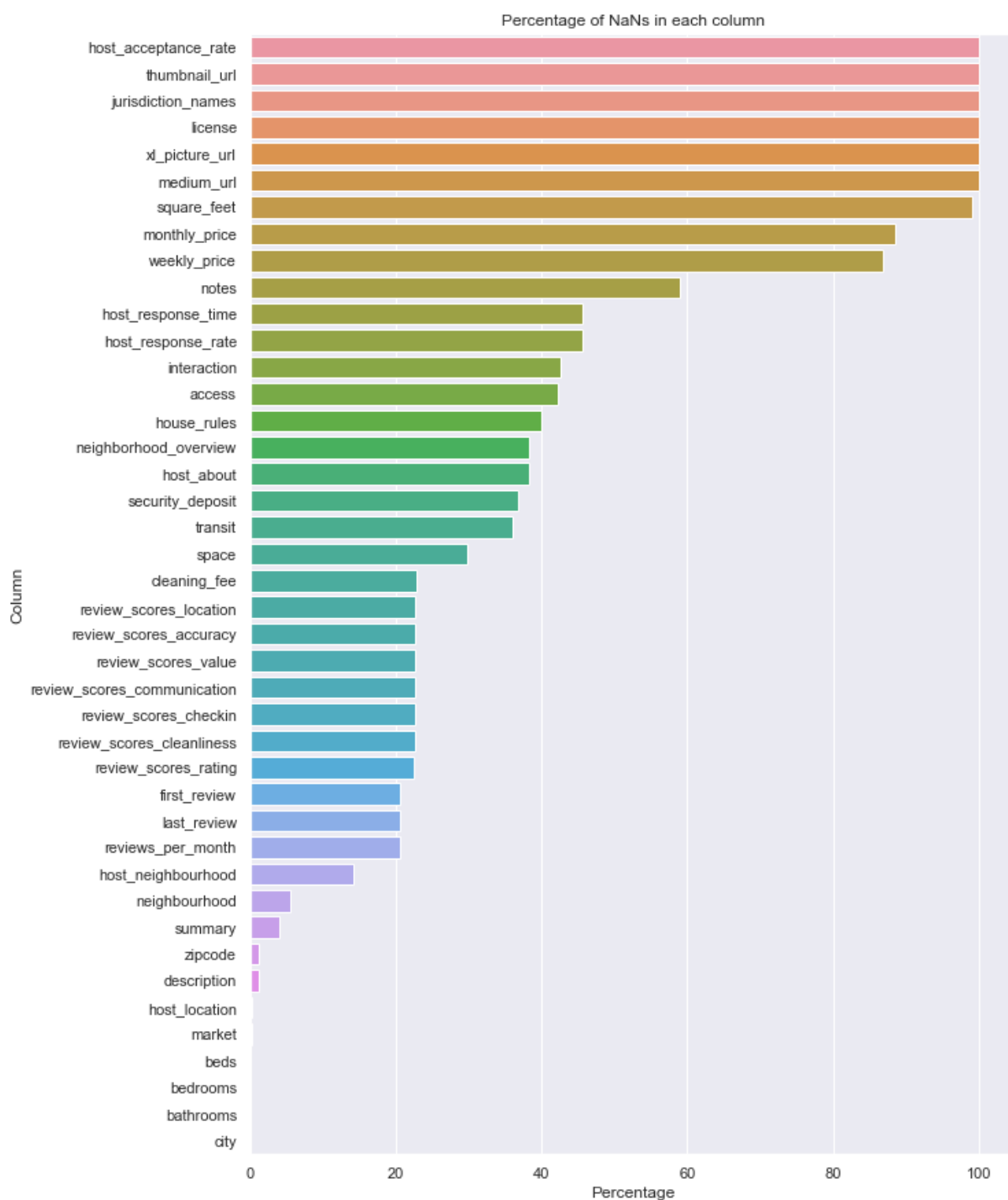
(50228, 105)
```

```
In [4]: 1 def remove_childlist_from_parentlist(parent, child):
2     return [x for x in parent if x not in child]
```

```
In [5]: 1 def plot_nans(df, ax, annot=True, filter_out_zeros=True):
2         nan_df = analyse_nans(df)
3         if filter_out_zeros:
4             cols_to_keep = []
5             for col in nan_df.columns:
6                 if nan_df[col].iloc[1] > 0:
7                     cols_to_keep.append(col)
8             nan_df = nan_df[cols_to_keep]
9         nan_df_transpose = nan_df.T
10        nan_df_transpose.sort_values(by="percentage", ascending=False, axis=0, i
11        sns.barplot(data=nan_df_transpose, y=nan_df_transpose.index, x="percenta
12        ax.set(ylabel="Column", xlabel="Percentage", title="Percentage of NaNs i
13        # plot_bar(data=nan_df_transpose, x=nan_df_transpose.index, y="percentag
14        # ax=ax, annot=annot, highlight_max_min=False,
15        # xlabel="Column", ylabel="Percentage", title="Percentage of Na
16        # xrot=90, plot_mean=False)
```

```
In [6]: 1 rows_to_drop = listings[listings['host_listings_count'] != listings['host_to
2        listings.drop(index=rows_to_drop,axis=0,inplace=True)
```

```
In [43]: 1 f, ax = plt.subplots(1,1, figsize=(10,15))
2         sns.set()
3         plot_nans(listings, ax=ax, annot=False, filter_out_zeros=True)
```



```

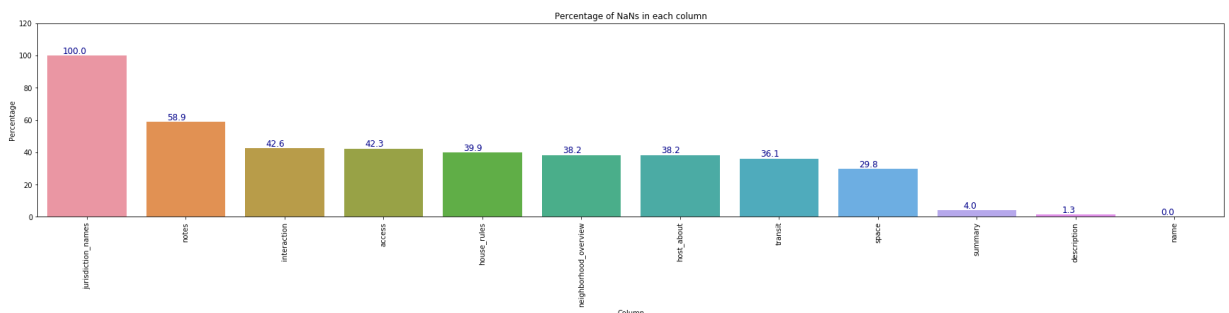
In [14]: 1 monu_cols_to_clean = ['id', 'host_response_time', 'host_response_rate',
2     'host_is_superhost',
3     'host_listings_count',
4     'host_total_listings_count', 'host_verifications',
5     'host_has_profile_pic', 'host_identity_verified',
6     'is_location_exact', 'property_type', 'room_type', 'accommodates',
7     'bathrooms', 'bedrooms', 'beds', 'bed_type', 'amenities', 'square_fee',
8     'price', 'weekly_price', 'monthly_price', 'security_deposit',
9     'cleaning_fee', 'guests_included', 'extra_people', 'instant_bookable',
10    'is_business_travel_ready', 'cancellation_policy',
11    'require_guest_profile_picture', 'require_guest_phone_verification',
12    'calculated_host_listings_count', 'reviews_per_month']
13
14 cols_to_drop = ['listing_url', 'scrape_id', 'last_scraped', 'experiences_offer',
15    'host_url', 'host_name', 'host_acceptance_rate', 'host_thumbnail_url',
16    'license']
17
18 col_list_1 = find_unique_elems([listings.columns, monu_cols_to_clean])
19
20 nlp_cols = \
21 ['jurisdiction_names',
22  'notes',
23  'interaction',
24  'access',
25  'house_rules',
26  'neighborhood_overview',
27  'host_about',
28  'transit',
29  'space',
30  'summary',
31  'name',
32  'description']
33
34 ssh_cols_to_clean = remove_childlist_from_parentlist(col_list_1, cols_to_drop)
35 ssh_cols_to_clean = remove_childlist_from_parentlist(ssh_cols_to_clean, nlp_cols)

```

```

In [17]: 1 nlp_listings = listings[nlp_cols]
2 plot_nans(nlp_listings)

```



```
In [18]: 1 nlp_lencols = ["description", "host_about"]
2 nlp_listings["description"].fillna("", inplace=True)
3 nlp_listings["host_about"].fillna("", inplace=True)
4
5 nlp_listings["description"] = nlp_listings["description"].astype(str)
6 nlp_listings["host_about"] = nlp_listings["host_about"].astype(str)
7
8 nlp_listings["desc_len"] = nlp_listings["description"].apply(len)
9 nlp_listings["host_about_len"] = nlp_listings["host_about"].apply(len)
```

```
In [19]: 1 nlp_len_cols = ["description", "host_about", "desc_len", "host_about_len"]
2 nlp_listings[nlp_len_cols].head()
```

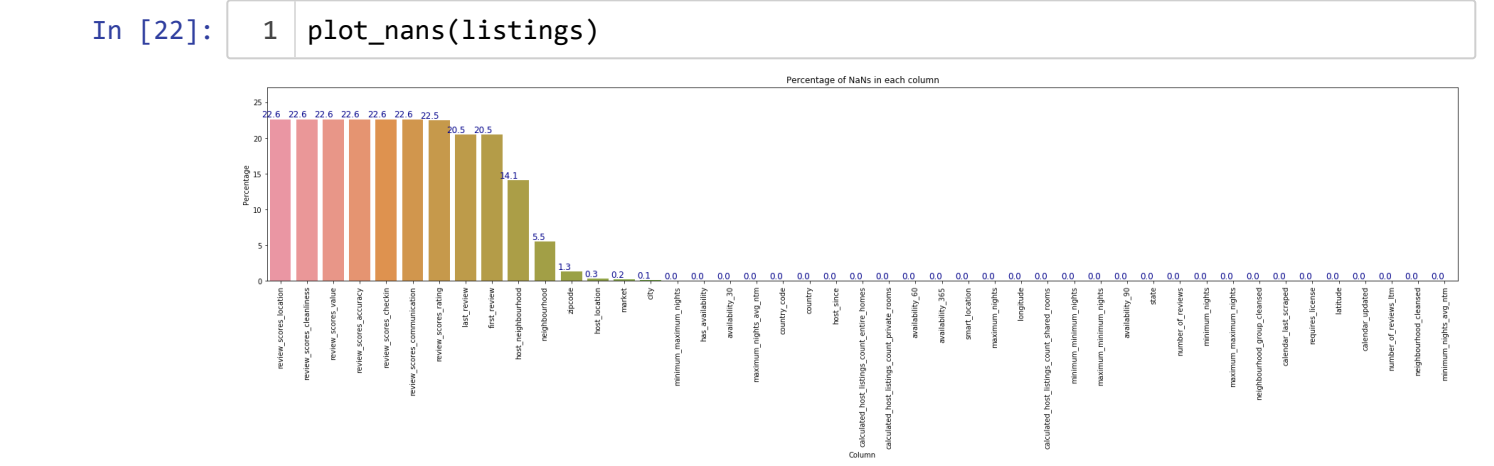
```
Out[19]:
```

	description	host_about	desc_len	host_about_len
0	Great light, exposed brick and 10 feet high ce...		412	0
1	Renovated apt home in elevator building. Spaci...	Educated professional living in Brooklyn. I l...	392	431
2	Find your romantic getaway to this beautiful, ...	A New Yorker since 2000! My passion is creatin...	1000	427
3	This is a spacious, clean, furnished master be...	From Brooklyn with love.	1000	25
4	WELCOME TO OUR INTERNATIONAL URBAN COMMUNITY T...	Make Up Artist National/ (Website hidden by Ai...	1000	354

```
In [20]: 1 cols_to_drop
```

```
Out[20]: ['listing_url',
'scrape_id',
'last_scraped',
'experiences_offered',
'thumbnail_url',
'medium_url',
'picture_url',
'xl_picture_url',
'host_id',
'host_url',
'host_name',
'host_acceptance_rate',
'host_thumbnail_url',
'host_picture_url',
'street',
'license']
```

```
In [21]: 1 listings = listings[ssh_cols_to_clean]
```



Make placeholders

In [23]: 1 more_cols_to_drop = []

In [24]: 1 listings.head()

Out[24]:

	availability_365	neighbourhood_cleaned	longitude	review_scores_value	calculated_host_listing
0	65	Midtown	-73.967679		NaN
1	365	Kensington	-73.972370	10.0	
2	365	Midtown	-73.983774	9.0	
3	290	Williamsburg	-73.942362	10.0	
4	365	Harlem	-73.941902		NaN

In [25]: 1 listings.shape

Out[25]: (50220, 45)

Country

```
In [26]: 1 listings["country"].value_counts()
```

```
Out[26]: United States    50220
         Name: country, dtype: int64
```

```
In [27]: 1 listings["Peru"]==listings["country"]
```

```
Out[27]: availability_365  neighbourhood_cleaned  longitude  review_scores_value  calculated_host_listings_
```

Drop column? Peru seems to be an outlier. listings url is also invalid when checked.

Review scores

Hypothesis: Most rows in the **review_scores_** columns that have NaNs, have them for all the **review_scores_** columns in that particular row.

i.e., if *review_scores_value* has a NaN, then so will *review_scores_location* and the other *review_scores_* columns.

Let's verify this by finding the columns that are NOT NaNs throughout all the *review_scores_** columns.

```
In [7]: 1
2 review_scores_cols = ['review_scores_value', 'review_scores_location', 'review_scores_checkin', 'review_scores_cleanliness', 'review_scores_communication', 'review_scores_accuracy', 'review_scores_rating']
3 nan_df = analyse_nans(listings[review_scores_cols])
4
5
6 # Verify that these columns have missing values in the exact same rows.
7 print("Number of rows between pairs of columns that have NaNs only in one of")
8 for i in range(1, len(review_scores_cols)):
9     print(review_scores_cols[i], "\t", review_scores_cols[i-1], end=" -->\t")
10     idx1 = nan_df[review_scores_cols[i-1]].iloc[2]
11     idx2 = nan_df[review_scores_cols[i]].iloc[2]
12     l = find_unique_elems([idx1, idx2])
13     print(len(l), "rows, \t i.e.", round(100*len(l)/listings.shape[0], 2), "%
```

```
Number of rows between pairs of columns that have NaNs only in one of them -
review_scores_location    review_scores_value -->          17 rows,          i.e. 0.
03 % of total listings rows
review_scores_checkin     review_scores_location -->         19 rows,          i.e. 0.
04 % of total listings rows
review_scores_cleanliness  review_scores_checkin -->         43 rows,
i.e. 0.09 % of total listings rows
review_scores_communication review_scores_cleanliness -->     28 rows,
i.e. 0.06 % of total listings rows
review_scores_accuracy    review_scores_communication -->     31 rows,
i.e. 0.06 % of total listings rows
review_scores_rating      review_scores_accuracy -->         47 rows,          i.e. 0.
09 % of total listings rows
```

Majority of rows have NaNs throughout all the review_scores_* columns. However, there are a few that don't as well.

Plot a distribution of Review Scores

In [8]: 1 listings[review_scores_cols].describe()

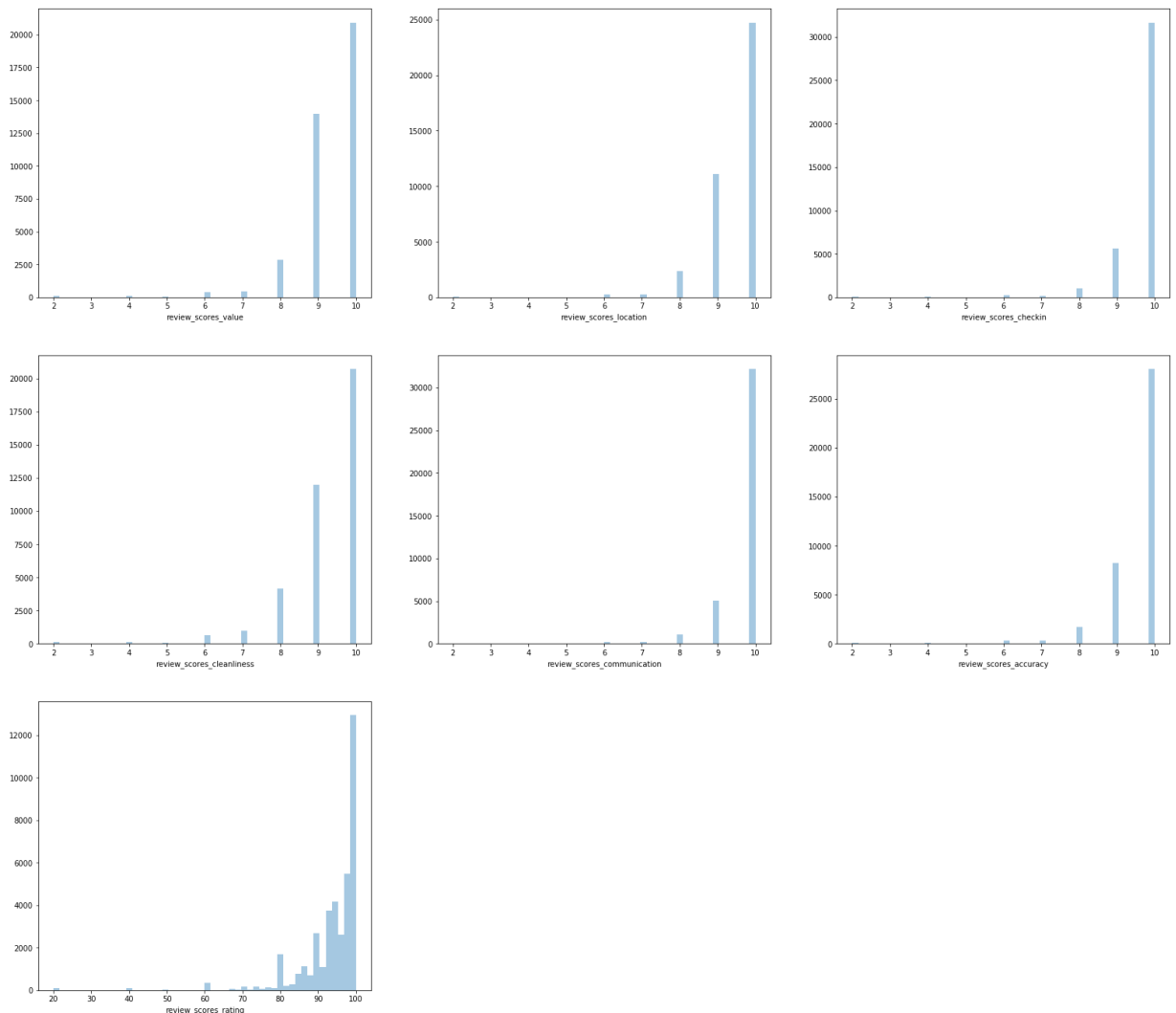
Out[8]:

	review_scores_value	review_scores_location	review_scores_checkin	review_scores_cleanline
count	38848.000000	38847.000000	38852.000000	38889.0000
mean	9.376184	9.527016	9.738315	9.2653
std	0.902918	0.776485	0.711515	1.0786
min	2.000000	2.000000	2.000000	2.0000
25%	9.000000	9.000000	10.000000	9.0000
50%	10.000000	10.000000	10.000000	10.0000
75%	10.000000	10.000000	10.000000	10.0000
max	10.000000	10.000000	10.000000	10.0000


```

In [9]: 1 # Uncomment these lines to see histograms of the review scores columns
2 f, ax = plt.subplots(3,3, figsize=(27,24))
3 col_counter = 0
4 g = sns.distplot(a=listings[~listings[review_scores_cols[col_counter]].isna(
5 col_counter+=1
6 g = sns.distplot(a=listings[~listings[review_scores_cols[col_counter]].isna(
7 col_counter+=1
8 g = sns.distplot(a=listings[~listings[review_scores_cols[col_counter]].isna(
9 col_counter+=1
10 g = sns.distplot(a=listings[~listings[review_scores_cols[col_counter]].isna(
11 col_counter+=1
12 g = sns.distplot(a=listings[~listings[review_scores_cols[col_counter]].isna(
13 col_counter+=1
14 g = sns.distplot(a=listings[~listings[review_scores_cols[col_counter]].isna(
15 col_counter+=1
16 g = sns.distplot(a=listings[~listings[review_scores_cols[col_counter]].isna(
17 col_counter+=1
18 ax[2][1].set_visible(not ax[2][1].get_visible())
19 ax[2][2].set_visible(not ax[2][2].get_visible())

```

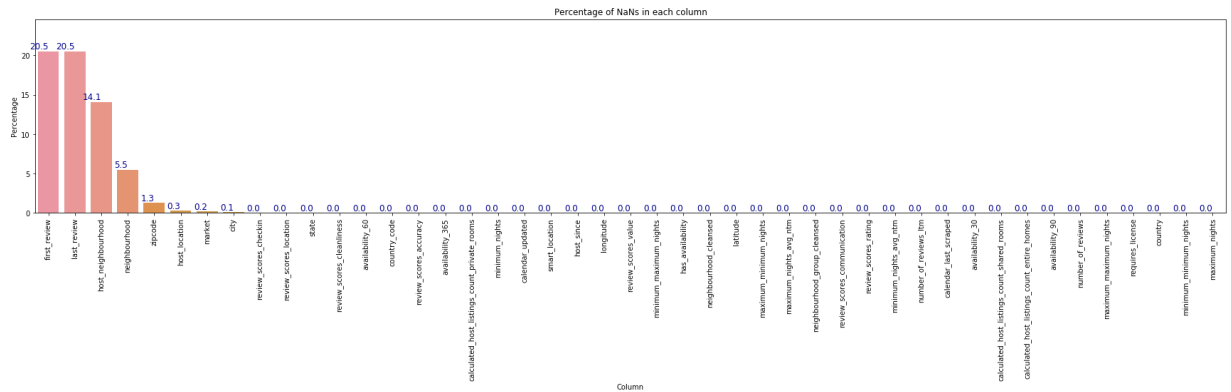


Well, skewed distributions throughout as expected. Transformations might prove useful here to compensate for the large left skew.

We could fill the missing values of each columns with the integer part of its median.

```
In [26]: 1 for col in review_scores_cols:
2         listings[col].fillna(listings[col].median(), inplace=True)
```

```
In [27]: 1 plot_nans(listings)
```



**city, market, zipcode, neighbourhood,
host_neighbourhood**

```
In [28]: 1 geo_cols = ["city", "neighbourhood", "neighbourhood_cleanse", "neighbourhood_group_cleanse"]
2         listings[listings["city"].isna()][geo_cols].head()
```

```
Out[28]:
```

	city	neighbourhood	neighbourhood_cleanse	neighbourhood_group_cleanse	market	hos
--	------	---------------	-----------------------	-----------------------------	--------	-----

8208	NaN	Sunset Park	Sunset Park	Brooklyn	New York
10646	NaN	Astoria	Astoria	Queens	New York
10857	NaN	NaN	Astoria	Queens	New York
10880	NaN	Astoria	Astoria	Queens	New York
10882	NaN	NaN	Astoria	Queens	New York

city

```
In [29]: 1 listings["city"].fillna("NYC", inplace=True)
```

market

```
In [30]: 1 listings["market"].value_counts()
```

```
Out[30]: New York                50066
Other (Domestic)             19
Adirondacks                  2
Catskills and Hudson Valley  1
Los Angeles                  1
Cuba                        1
Agra                        1
Atlanta                     1
Paris                       1
San Francisco                1
D.C.                        1
Boston                      1
Jamaica South Coast          1
South Bay, CA                1
New Orleans                  1
Kyoto                       1
Lagos, NG                    1
Name: market, dtype: int64
```

```
In [31]: 1 no_geonan_listings = listings[geo_cols].dropna()
```

```
In [32]: 1 for market in no_geonan_listings["market"].unique():
2         print(market, " --> ", no_geonan_listings[no_geonan_listings["market"]])
```

```
New York --> ['Manhattan' 'Brooklyn' 'Queens' 'Staten Island' 'Bronx']
Other (Domestic) --> ['Queens' 'Brooklyn' 'Manhattan']
Boston --> ['Brooklyn']
Los Angeles --> ['Manhattan']
Atlanta --> ['Brooklyn']
Paris --> ['Brooklyn']
New Orleans --> ['Brooklyn']
San Francisco --> ['Manhattan']
Cuba --> ['Queens']
D.C. --> ['Brooklyn']
Lagos, NG --> ['Brooklyn']
Agra --> ['Manhattan']
Catskills and Hudson Valley --> ['Brooklyn']
Kyoto --> ['Brooklyn']
Adirondacks --> ['Brooklyn']
Jamaica South Coast --> ['Queens']
```

```
In [33]: 1 map_nbc_market = {}
2         for nbc in no_geonan_listings["neighbourhood_cleansed"].unique():
3             map_nbc_market[nbc] = listings[listings["neighbourhood_cleansed"]==nbc][
```

```
In [34]: 1 map_nbc_market
```

```
Out[34]: {'Midtown': array(['New York', nan], dtype=object),
          'Kensington': array(['New York', 'Other (Domestic)'], dtype=object),
          'Williamsburg': array(['New York', 'D.C.', 'Kyoto', nan], dtype=object),
          'Harlem': array(['New York', 'Agra', nan], dtype=object),
          'Clinton Hill': array(['New York', nan], dtype=object),
          "Hell's Kitchen": array(['New York', nan], dtype=object),
          'Upper West Side': array(['New York', nan], dtype=object),
          'Flatiron District': array(['New York', 'Other (Domestic)', nan], dtype=object),
          'Chinatown': array(['New York'], dtype=object),
          'Upper East Side': array(['New York', nan, 'Other (Domestic)'], dtype=object),
          'South Slope': array(['New York'], dtype=object),
          'West Village': array(['New York', 'San Francisco', nan], dtype=object),
          'East Harlem': array(['New York', nan], dtype=object),
          'Fort Greene': array(['New York'], dtype=object),
          'Chelsea': array(['New York', nan], dtype=object),
          'Crown Heights': array(['New York', 'Paris', nan], dtype=object),
          'Park Slope': array(['New York', nan], dtype=object),
          'East Village': array(['New York', nan], dtype=object),
          'Greenwich Village': array(['New York', nan], dtype=object),
          'SoHo': array(['New York', nan], dtype=object),
          'Lower East Side': array(['New York', nan], dtype=object),
          'Noho': array(['New York', nan], dtype=object),
          'Hudson Yards': array(['New York', nan], dtype=object),
          'Lincoln Square': array(['New York', nan], dtype=object),
          'Lincoln Center': array(['New York', nan], dtype=object),
          'Rockefeller Center': array(['New York', nan], dtype=object),
          'Times Square': array(['New York', nan], dtype=object),
          'Port Authority Bus Terminal': array(['New York', nan], dtype=object),
          'Grand Central Terminal': array(['New York', nan], dtype=object),
          'Penn Station': array(['New York', nan], dtype=object),
          'Madison Avenue': array(['New York', nan], dtype=object),
          'Fifth Avenue': array(['New York', nan], dtype=object),
          'Broadway': array(['New York', nan], dtype=object),
          'Manhattan': array(['New York', nan], dtype=object),
          'New York City': array(['New York', nan], dtype=object),
          'United States': array(['New York', nan], dtype=object),
          'World': array(['New York', nan], dtype=object)}
```

Fill the market using the neighbourhood cleansed column

```
In [35]: 1 def impute_market(row):
2         if type(row["market"]) != type("a"):
3             nb_cleansed = row["neighbourhood_cleansed"]
4             val = map_nbc_market[nb_cleansed][0]
5             return val
6         else:
7             return row
8         return row
9
10 listings["market"] = listings[["market", "neighbourhood_cleansed"]].apply(im
```

```
In [39]: 1 len(listings["zipcode"].unique())
```

```
Out[39]: 378
```

```
In [40]: 1 listings["zipcode"].fillna("XX", inplace=True)
```

First review and last review

Date columns. Why are there missing values?

```
In [10]: 1 review_date_cols = ["first_review", "last_review"]
2 nan_df = analyse_nans(listings)
3
4 # Verify that these columns have missing values in the exact same rows.
5 print("Number of rows between pairs of columns that have NaNs only in one of")
6 for i in range(1, len(review_date_cols)):
7     print(review_date_cols[i], "\t", review_date_cols[i-1], end=" -->\t ")
8     idx1 = nan_df[review_date_cols[i-1]].iloc[2]
9     idx2 = nan_df[review_date_cols[i]].iloc[2]
10    l = find_unique_elems([idx1, idx2])
11    print(len(l), "rows, \ti.e.", round(100*len(l)/listings.shape[0], 2), "%")
```

Number of rows between pairs of columns that have NaNs only in one of them -
last_review first_review --> 0 rows, i.e. 0.0 % of total listings rows

```
In [11]: 1 listings[(~listings["last_review"].isna())&(listings["first_review"].isna())
2          (~listings["first_review"].isna())&(listings["last_review"].isna())]
```

```
Out[11]: listing_url  scrape_id  last_scraped  name  summary  space  description  experiences_offered  id
```

```
In [12]: 1 listings['last_review'] = pd.to_datetime(listings['last_review'])
2 listings['lreview_year'] = listings['last_review'].dt.year
3 listings['lreview_month'] = listings['last_review'].dt.month
4 listings['lreview_day'] = listings['last_review'].dt.day
5
6 listings['first_review'] = pd.to_datetime(listings['first_review'])
7 listings['freview_year'] = listings['first_review'].dt.year
8 listings['freview_month'] = listings['first_review'].dt.month
9 listings['freview_day'] = listings['first_review'].dt.day
```

```
In [13]: 1 cols = ["last_review", 'lreview_year', 'lreview_month', 'lreview_day']  
        2 listings[cols].describe()
```

```
Out[13]:
```

	lreview_year	lreview_month	lreview_day
count	39931.000000	39931.000000	39931.000000
mean	2017.998247	5.644687	14.929178
std	1.133024	4.536616	10.082954
min	2010.000000	1.000000	1.000000
25%	2018.000000	1.000000	5.000000
50%	2018.000000	5.000000	15.000000
75%	2019.000000	10.000000	24.000000
max	2019.000000	12.000000	31.000000

```
In [14]: 1 sns.set(font_scale=2, style="whitegrid")
```

```

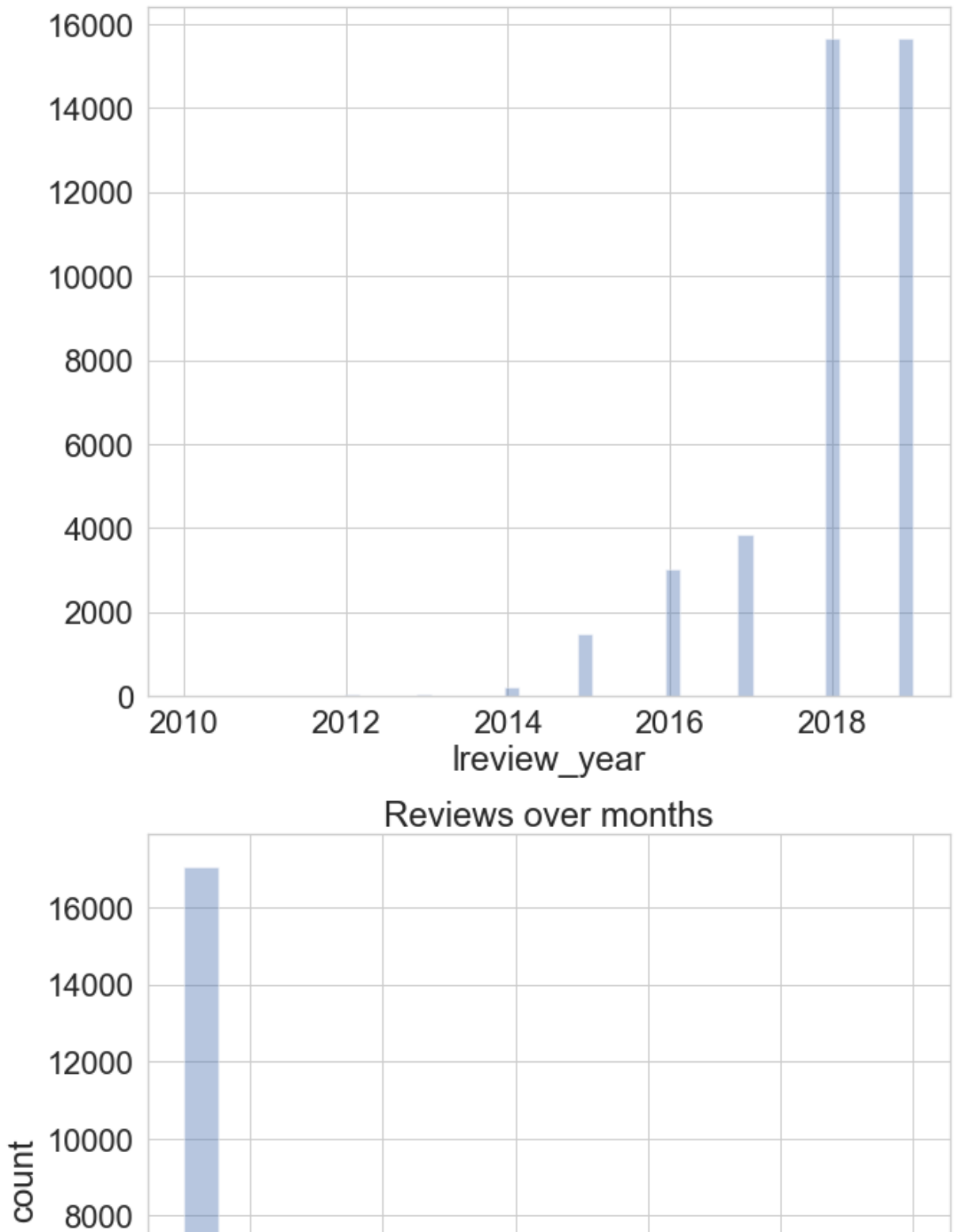
In [16]: 1 f, ax = plt.subplots(3, 1, figsize=(10, 30))
          2 g = sns.distplot(a=listings[~listings["lreview_year"].isna()][ "lreview_year"
          3 g = sns.distplot(a=listings[~listings["lreview_month"].isna()][ "lreview_mont
          4 g.set(title="Reviews over months", xlabel="Month of year", ylabel="count")
          5 g = sns.distplot(a=listings[~listings["lreview_day"].isna()][ "lreview_day"],
          6 g.set(title="Reviews over days", xlabel="Day of month", ylabel="count")

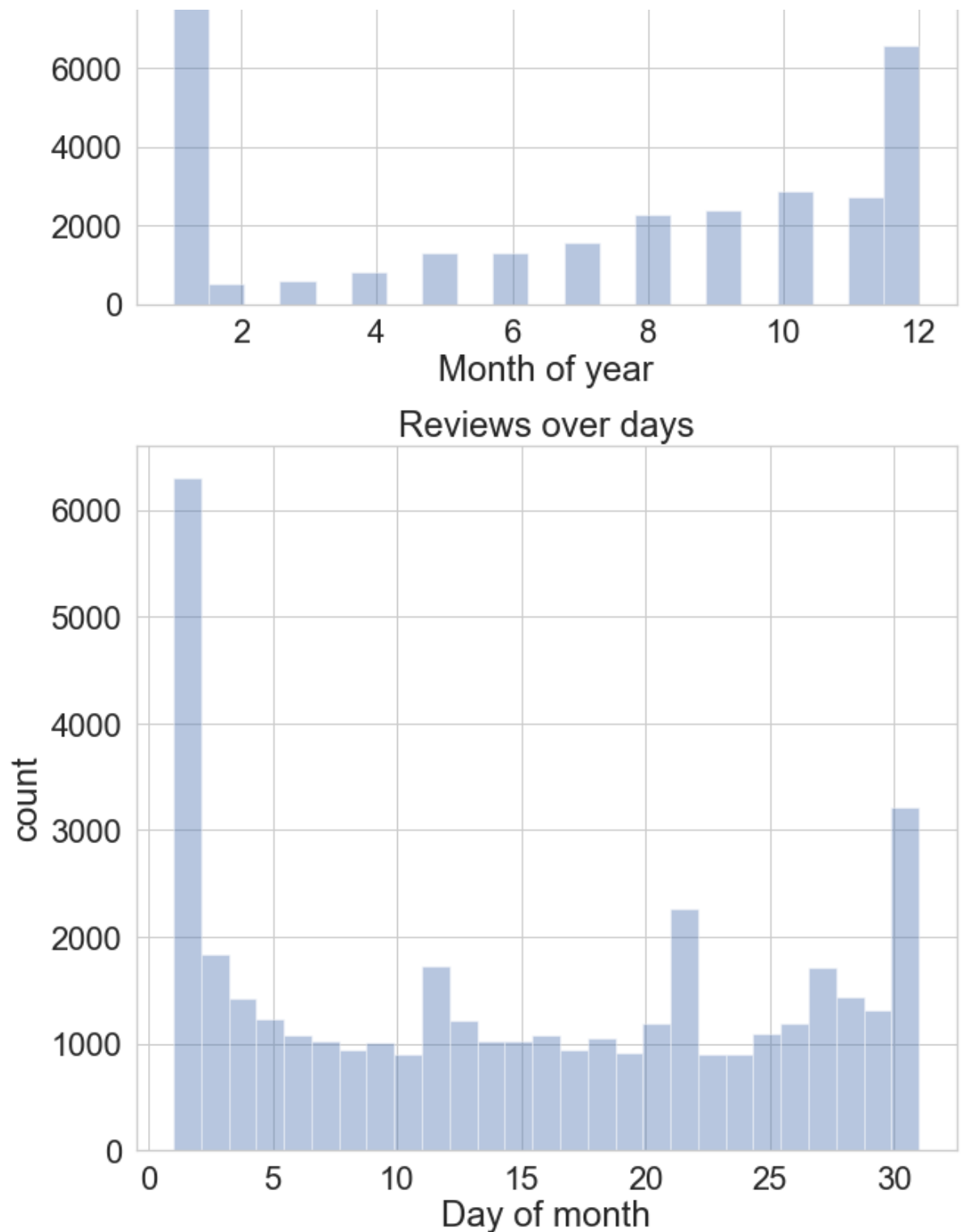
```

```

Out[16]: [Text(0, 0.5, 'count'),
          Text(0.5, 0, 'Day of month'),
          Text(0.5, 1.0, 'Reviews over days')]

```





If we use median, we can estimate that the missing year values is (most probably) 2018, the missing month values is (most probably) 10/11 and for date would be around 16.

Or we can create a separate category called X for each column, to estimate that there has been no review written for those dates yet. Upon checking online using the URLs provided, we can verify that no reviews are provided for those listings yet.

The latter method is used for imputation here.


```
In [47]: 1 listings["lreview_year"].apply(type)
```

```
Out[47]: 0      <class 'float'>
         1      <class 'float'>
         2      <class 'float'>
         3      <class 'float'>
         4      <class 'float'>
         5      <class 'float'>
         6      <class 'float'>
         7      <class 'float'>
         8      <class 'float'>
         9      <class 'float'>
        10      <class 'float'>
        11      <class 'float'>
        12      <class 'float'>
        13      <class 'float'>
        14      <class 'float'>
        15      <class 'float'>
        16      <class 'float'>
        17      <class 'float'>
        18      <class 'float'>
```

```
In [49]: 1 listings["first_review"].fillna("XX", inplace=True)
         2 listings["last_review"].fillna("XX", inplace=True)
         3 listings['lreview_year'].fillna(0, inplace=True)
         4 listings['lreview_month'].fillna(0, inplace=True)
         5 listings['lreview_day'].fillna(0, inplace=True)
         6 listings['freview_year'].fillna(0, inplace=True)
         7 listings['freview_month'].fillna(0, inplace=True)
         8 listings['freview_day'].fillna(0, inplace=True)
```

```
In [54]: 1 listings['ndays_between_f_l_reviews'] = abs(listings['lreview_day'] - listings['freview_day'])
```

```
In [ ]: 1 listings
```

```
In [55]: 1 listings['ndays_between_f_l_reviews'].head()
```

```
Out[55]: 0      0.0
         1     15.0
         2     19.0
         3      3.0
         4      0.0
         Name: ndays_between_f_l_reviews, dtype: float64
```

In [51]:

1 plot_nans(listings)

