In [91]:
```python
# Put these at the top of every notebook, to get automatic reloading and inl
from IPython.core.display import display, HTML
import pandas as pd
import warnings
warnings.filterwarnings('ignore')

%reload_ext autoreload
%autoreload 1
%matplotlib inline

pd.set_option('display.max_rows', 500)
pd.set_option('display.max_columns', 500)
pd.set_option('display.width', 1000)

display(HTML("<style>.container { width:100% !important; }</style>"))
```

In [92]:
```python
import os
import seaborn as sns
sns.set(font_scale=2, style="whitegrid")
import pandas as pd
import math

from Utils.UtilsViz import *
from Utils.DataUtils import *
```

In [93]:
```python
US_coord = [37.0902, -102]
NY_COORD = [40.7128, -74.0060]

# ny_data_path = os.getcwd()
ny_datapath = "C:\\Users\\sriharis\\OneDrive\\UChicago\\DataMining\\project\
# ny_datapath = "C:\\Users\\Ssrih\\OneDrive\\UChicago\\DataMining\\project\\
```

In [94]:
```python
listings = pd.read_csv(os.path.join(ny_datapath, "listings.csv"))
# print(os.getcwd())
# ny_datapath = os.path.join(os.getcwd(), "../data/listings_no_nlp.csv")
# listings = pd.read_csv(ny_datapath)
# listings = pd.read_csv(ny_datapath, index_col="Unnamed: 0")
```
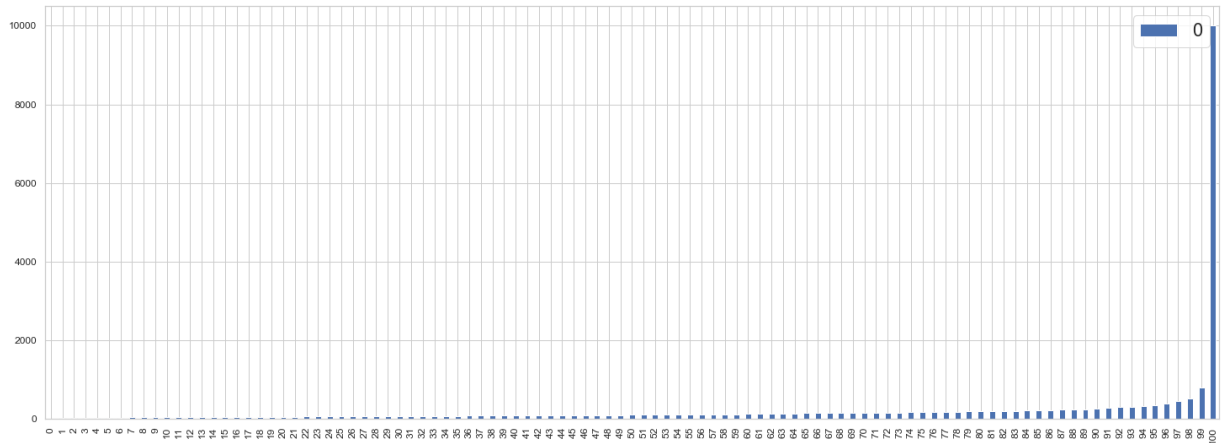
In [95]:    `1  listings.head()`

Out[95]:

| | id | listing_url | scrape_id | last_scraped | name | summary |
|---|---|---|---|---|---|---|
| 0 | 2454 | https://www.airbnb.com/rooms/2454 | 20190201155637 | 2019-02-01 | superCondo | Great light, exposed brick and 10 feet high ce... |
| 1 | 2539 | https://www.airbnb.com/rooms/2539 | 20190201155637 | 2019-02-02 | Clean & quiet apt home by the park | Renovated apt home in elevator building. |
| 2 | 2595 | https://www.airbnb.com/rooms/2595 | 20190201155637 | 2019-02-02 | Skylit Midtown Castle | Find your romantic getaway to this beautiful, ... |
| 3 | 3330 | https://www.airbnb.com/rooms/3330 | 20190201155637 | 2019-02-02 | ++ Brooklyn Penthouse Guestroom ++ | This is a spacious, clean, furnished master be... |
| 4 | 3647 | https://www.airbnb.com/rooms/3647 | 20190201155637 | 2019-02-02 | THE VILLAGE OF HARLEM....NEW YORK ! | NaN |

# Price preprocessing

In [11]:    `1  listings['price'] = listings['price'].str.strip('').str.strip('$').str.repla`
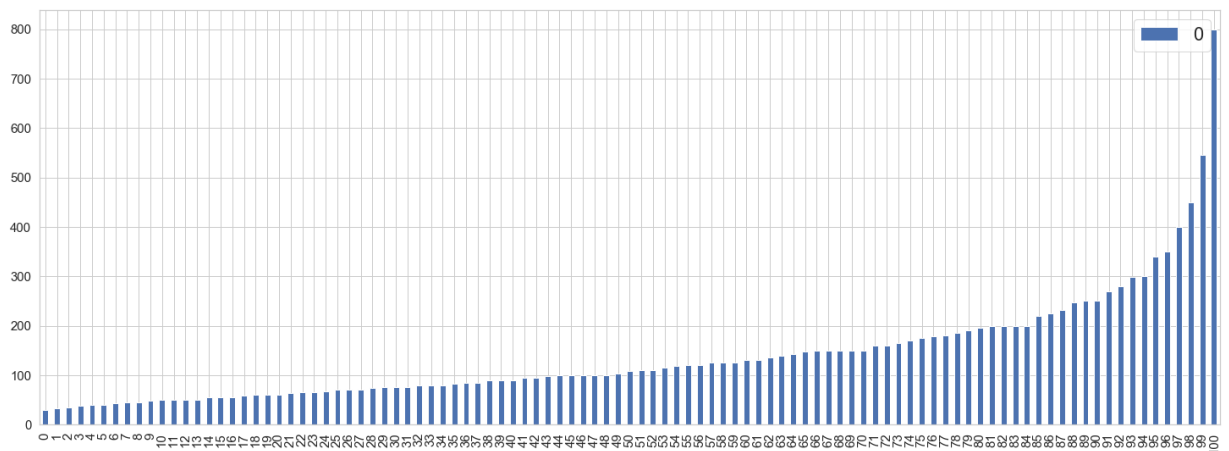
```
In [12]:    1  percentiles = list(range(0,101))
            2  price_percentile = {}
            3  for p in percentiles:
            4      price_percentile[p] = np.percentile(listings['price'].values, p)
            5
            6  price_percentile = pd.DataFrame.from_dict(price_percentile, orient='index')
            7  g = price_percentile.plot(kind='bar', figsize=(25,9), grid=True)
            8  t = g.tick_params(labelsize=12)
            9
```



```
In [13]:    1  listings = listings[listings["price"] <= price_percentile.iloc[99,:].values[
            2  listings = listings[listings["price"] >= price_percentile.iloc[1,:].values[0
            3  listings["price"].describe()
```

```
Out[13]:  count    49251.000000
          mean       137.606262
          std        101.958580
          min         30.000000
          25%         70.000000
          50%        108.000000
          75%        175.000000
          max        799.000000
          Name: price, dtype: float64
```
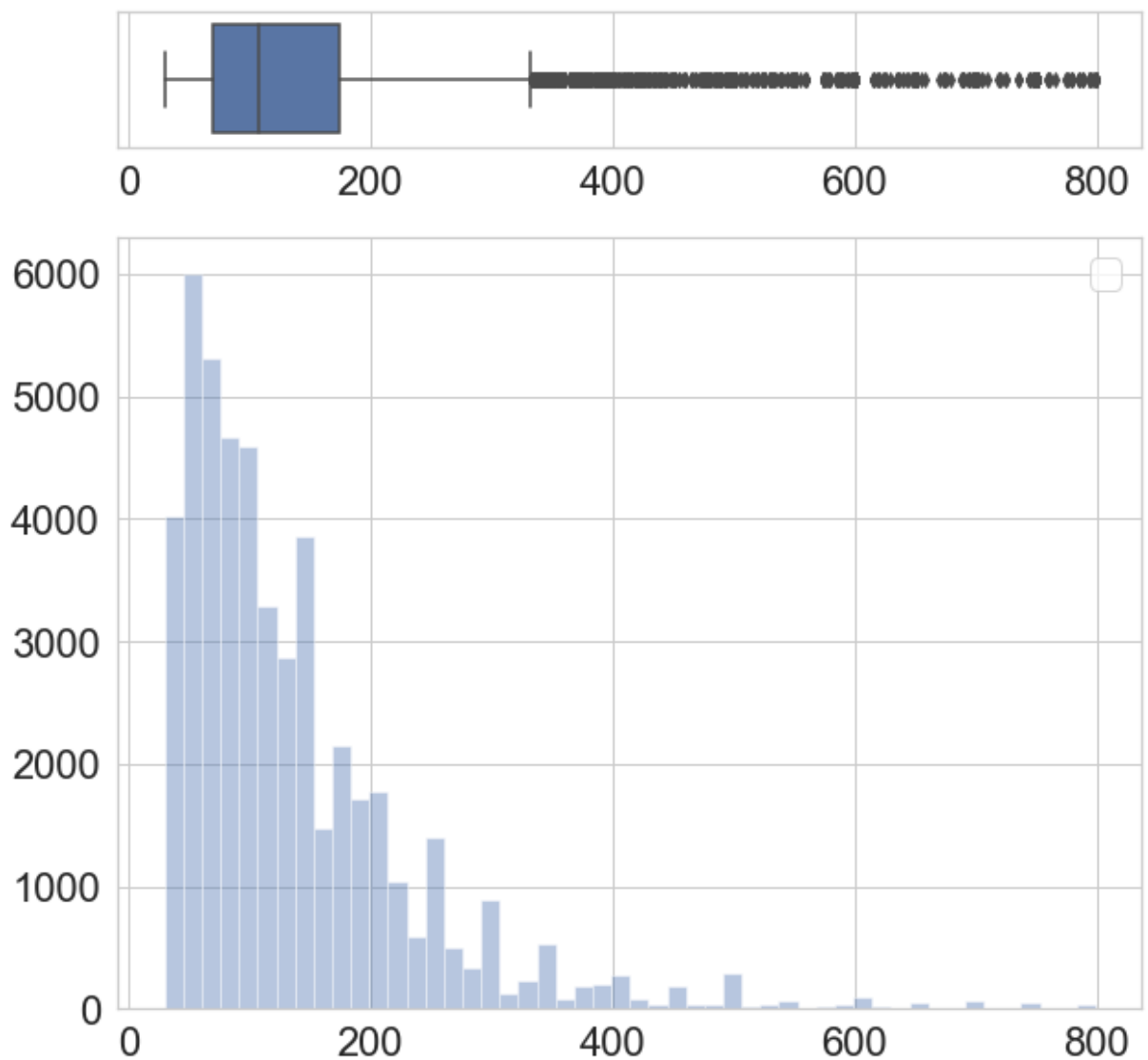
In [14]:

```python
percentiles = list(range(0,101))
price_percentile = {}
for p in percentiles:
    price_percentile[p] = np.percentile(listings['price'].values, p)

price_percentile = pd.DataFrame.from_dict(price_percentile, orient='index')
g = price_percentile.plot(kind='bar', figsize=(25,9), grid=True)
t = g.tick_params(labelsize=15)
```
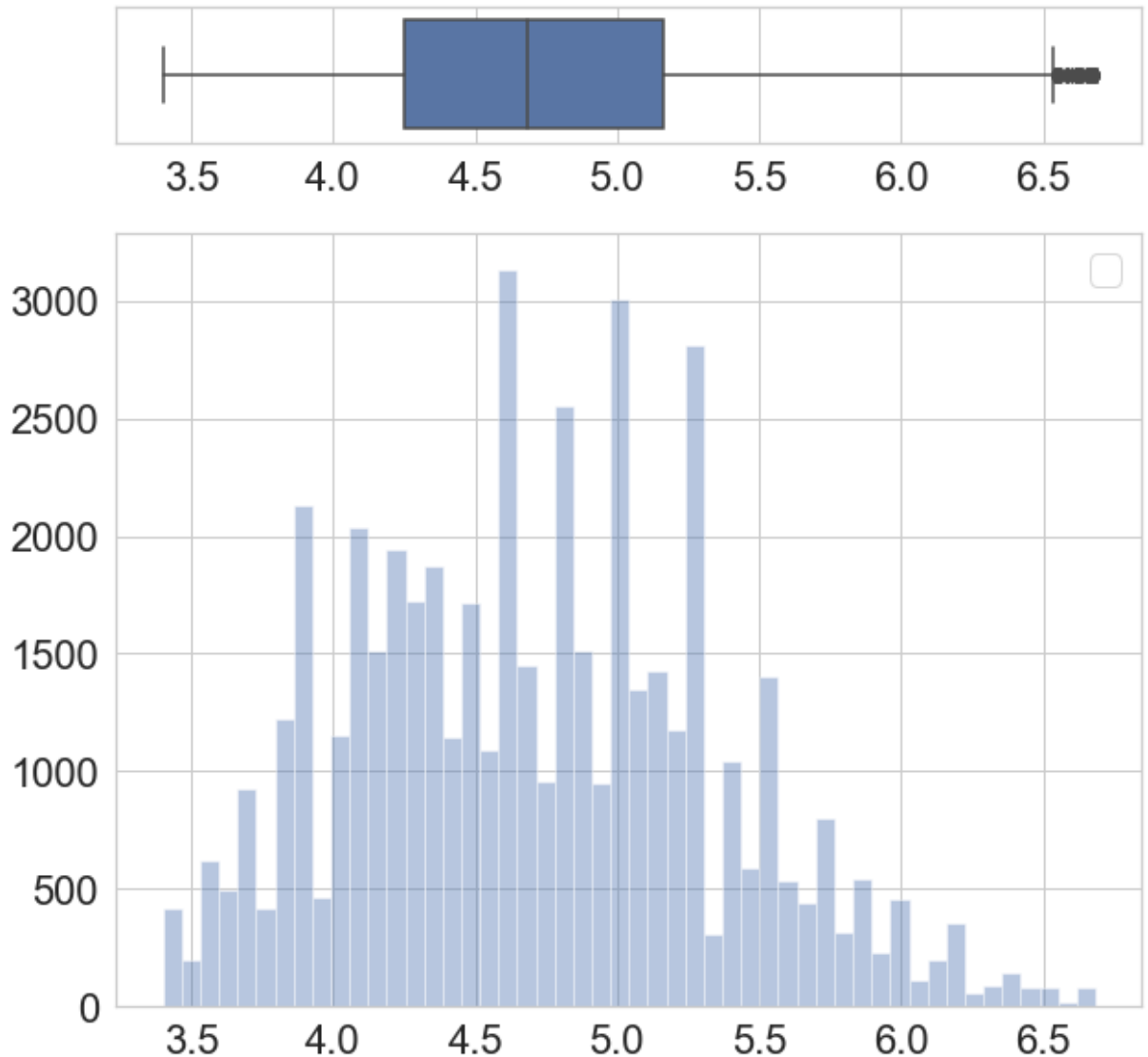
In [15]:     1  plot_dist(data=listings, colname="price", kde=False)

No handles with labels found to put in legend.

```
In [16]:   1  def get_logprice(price):
           2      if price <= 0.0:
           3          return 0.0
           4      else:
           5          return np.log(price)
           6  listings['price_log'] = listings['price'].apply(get_logprice)
           7  plot_dist(data=listings, colname="price_log", kde=False)
```

No handles with labels found to put in legend.



```
In [17]:   1  listings.drop("price_log", axis=1, inplace=True)
```

# Host response rate

```
In [ ]:    1  def get_hrr_fillval(listings):
           2      non_null = listings['host_response_rate'].dropna(axis=0)
           3      fv = non_null.str.strip('%').astype('int').median()
           4      return fv
           5  listings['host_response_rate'] = listings['host_response_rate'].fillna(str(i
           6  listings['host_response_rate'] = listings['host_response_rate'].str.strip('%
```

```
In [50]:   1  def roundto(row, base=5):
           2      return int(base * round(float(row) / base))
           3  listings["price_rounded"] = listings["price"].apply(roundto)
```

```
In [88]:   1  f, ax = plt.subplots(1,1,figsize=(30,9))
           2  g = sns.pointplot(x=listings["price_rounded"].values[::20], y=listings["host
           3  t = g.set(title="Host response rate over price bins", ylabel="response rate"
           4  t = g.tick_params(labelsize=10)
```



Host response rate over price bins

```
In [21]:   1  listings.drop("price_rounded", axis=1, inplace=True)
```

# Host count of listings

```
In [ ]:    1  listings["host_listings_count"]
```

```
In [90]:   1  tmp = listings[["host_id", "host_listings_count"]].groupby("host_id", as_ind
           2  tmp.head()
```
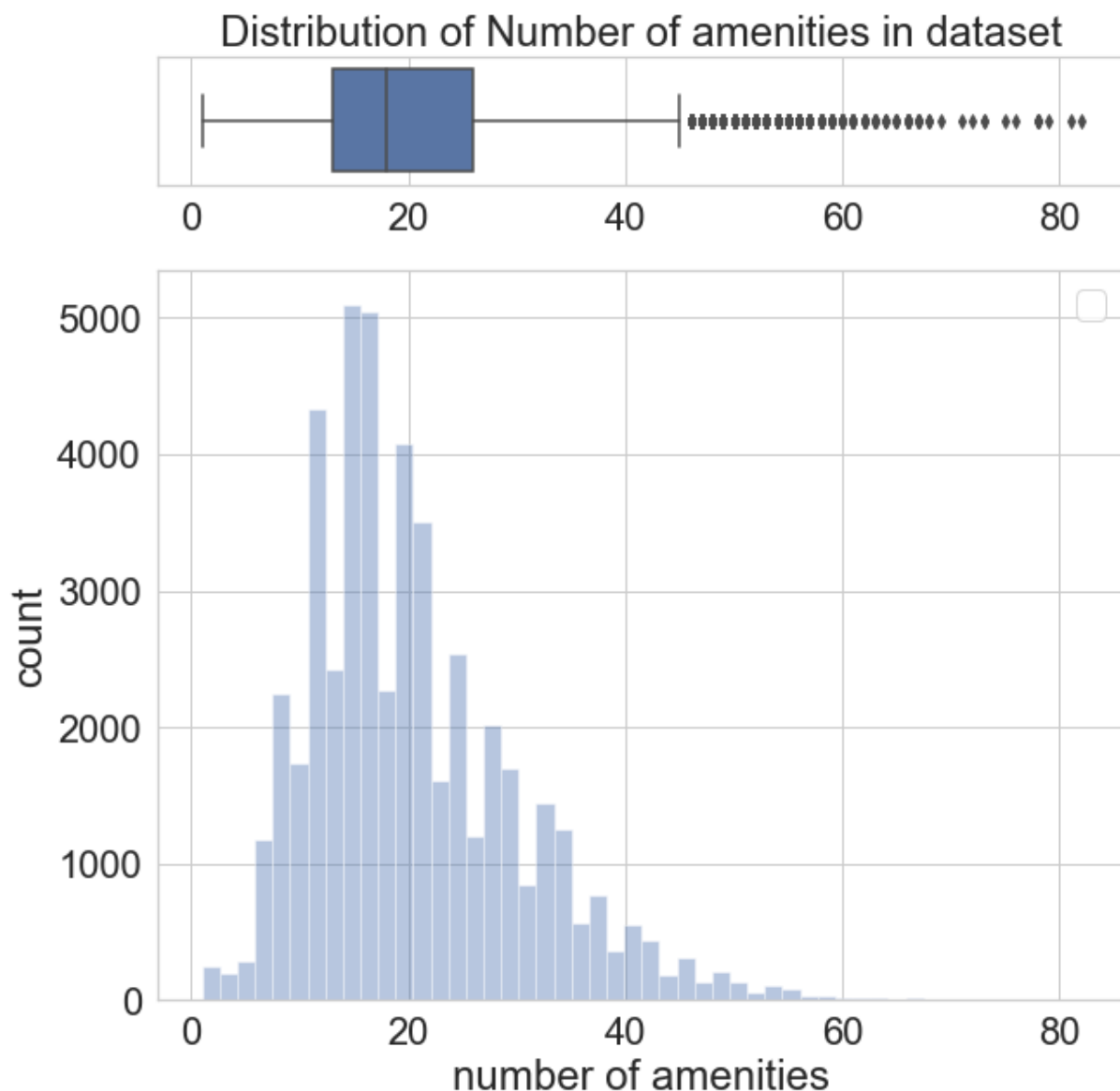
Out[90]:

|   | host_id | host_listings_count |
|---|---------|---------------------|
| 0 | 2571    | 1                   |
| 1 | 2688    | 1                   |
| 2 | 2787    | 8                   |
| 3 | 2845    | 2                   |
| 4 | 2881    | 2                   |

# Amenities

Let's have a look at the distribution of amenities

```
In [22]:  1  # Amenities
          2  def get_num_amenities(row):
          3      a = row[1:-1].split(",")
          4      return len(a)
          5
          6  listings["num_amenities"] = listings["amenities"].apply(get_num_amenities)
          7  g = plot_dist(data=listings, colname="num_amenities", xlabel="number of amen
```
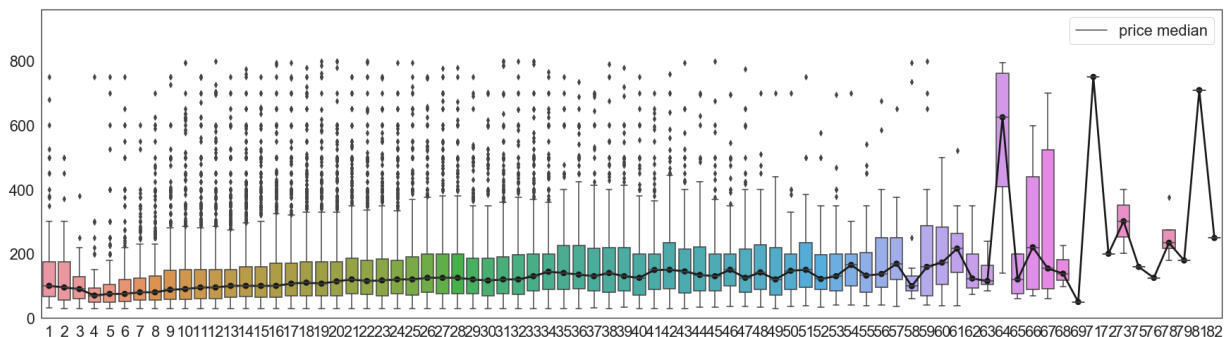
No handles with labels found to put in legend.



Distribution of Number of amenities in dataset

How does price behave based on number of amenities?
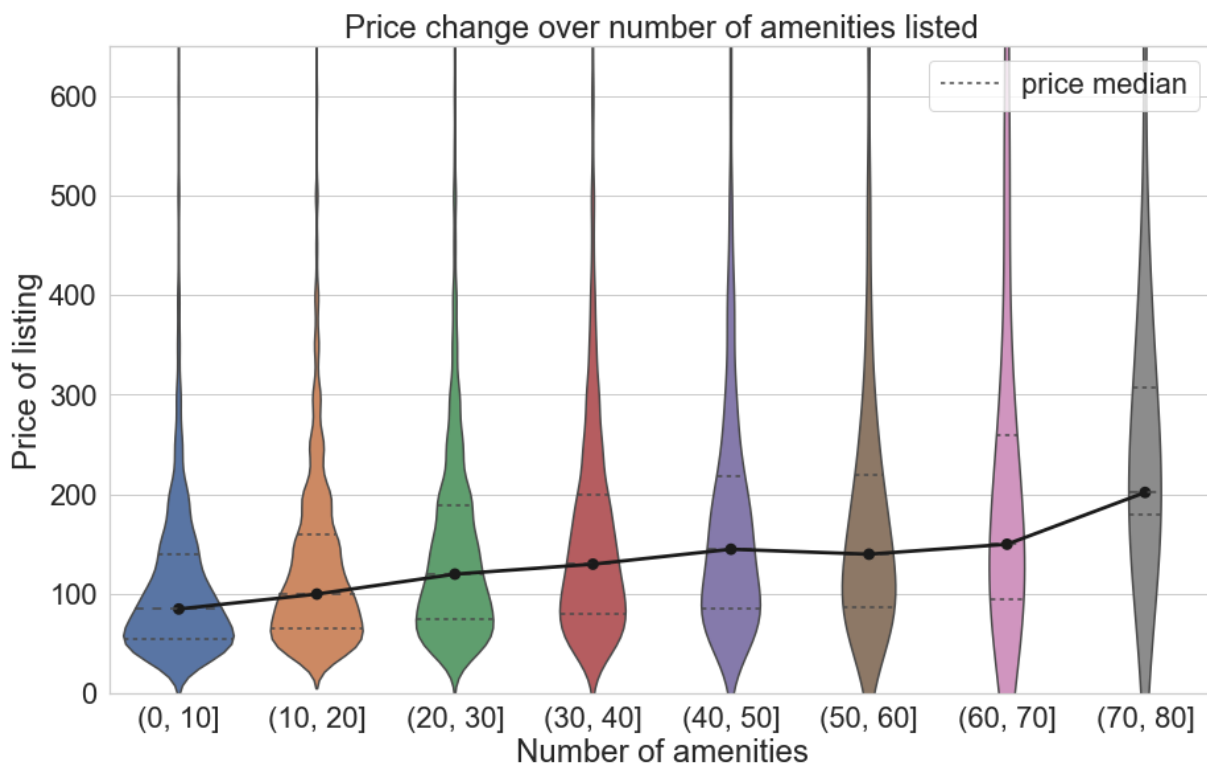
Let's filter out all the prices above 500$

In [53]:
```python
f, ax = plt.subplots(1,1,figsize=(30, 8))
g = plot_box(x="num_amenities", y="price", data=listings, ax=ax, agg_rule="m
```



Let's group price into bins and see how behaviour is

In [24]:
```python
bins = list(range(0,90, 10))
listings['amenities_binned'] = pd.cut(listings['num_amenities'], bins)
```

In [25]:
```python
f, ax = plt.subplots(1,1,figsize=(15, 9))
g = plot_violin(x="amenities_binned", y="price", data=listings, ax=ax, agg_r
            title="Price change over number of amenities listed",
            xlabel="Number of amenities", ylabel="Price of listing", ylim=65
```



# Accommodates

Does the price vary significantly as accommodation increases?

In [26]:
```python
def group_accommodates(row):
    if row < 10:
        return str(int(row))
    elif row >= 10:
        return "10+"
    else:
        return row

listings["acc_group"] = listings["accommodates"].apply(group_accommodates)
# listings[["bedrooms", "bedroom_group"]]
```

In [27]:

```python
# f, ax = plt.subplots(1,1,figsize=(15, 12))
# g = plot_box(x="accommodates", y="price", data=listings, ax=ax, agg_rule="

f, ax = plt.subplots(1,1,figsize=(15, 12))

listings['acc_group'] = pd.Categorical(
    listings['acc_group'],
    categories=['0','1','2','3','4','5','6', '7', '8', '9', '10+'],
    ordered=True
)
listings.sort_values(by="acc_group", inplace=True)

g = sns.boxplot(x="acc_group", y="price", data=listings, ax=ax)

agg_data = listings[["price", "acc_group"]].groupby(by=["acc_group"], as_ind
g = sns.pointplot(x="acc_group", y="price", data=agg_data, ax=ax, color="k")
t = g.set(title="Price variation over number of people who can be accommodat
```



Price variation over number of people who can be accommodated

## Bathrooms, Bedrooms and Beds

Manhattan is definitely the most expensive.

In [28]:
```python
1  cols = ["bathrooms", "bedrooms", "beds"]
2  listings[cols].describe()
```

Out[28]:

|        | bathrooms    | bedrooms     | beds         |
|--------|--------------|--------------|--------------|
| count  | 49185.000000 | 49208.000000 | 49210.000000 |
| mean   | 1.135814     | 1.173529     | 1.536862     |
| std    | 0.404808     | 0.737600     | 1.059946     |
| min    | 0.000000     | 0.000000     | 0.000000     |
| 25%    | 1.000000     | 1.000000     | 1.000000     |
| 50%    | 1.000000     | 1.000000     | 1.000000     |
| 75%    | 1.000000     | 1.000000     | 2.000000     |
| max    | 7.500000     | 11.000000    | 21.000000    |

Bathrooms and bedrooms may not add as much value to the dataset when over 75% of the column is just 1.

Beds can be filled with the median, 1.0

# Security Deposit

In [29]:
```python
1  listings['security_deposit'] = listings['security_deposit'].dropna(axis=0).s
```
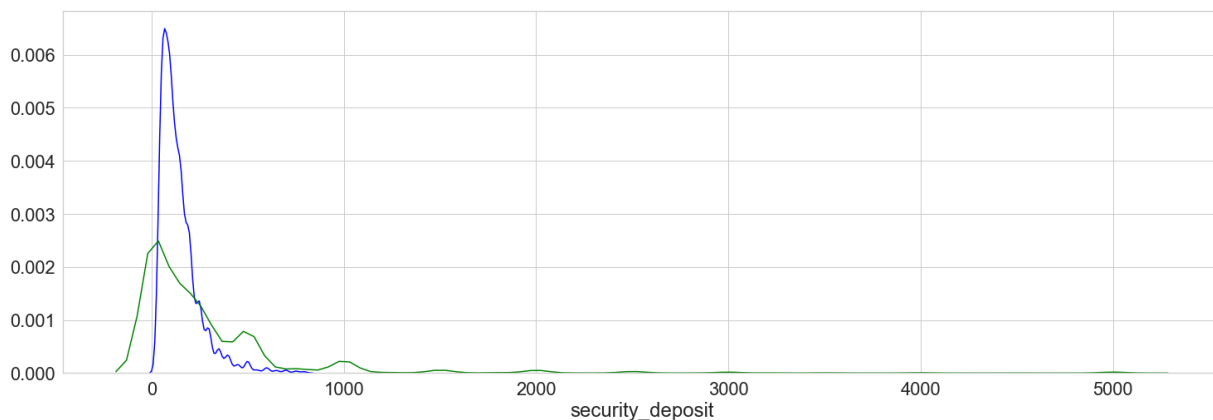
In [30]:
```python
1  tmp = listings[["price", "security_deposit"]]
2  tmp = tmp.dropna()
3  tmp.corr()
```

Out[30]:

|                  | price    | security_deposit |
|------------------|----------|------------------|
| price            | 1.000000 | 0.245089         |
| security_deposit | 0.245089 | 1.000000         |

```
In [31]:    1  f, ax = plt.subplots(1, 1, figsize=(25, 8))
            2  sns.distplot(tmp["price"], ax=ax, hist=False, kde=True, color="blue")
            3  sns.distplot(tmp["security_deposit"], ax=ax, hist=False, kde=True, color="gr
```

Out[31]:  <matplotlib.axes._subplots.AxesSubplot at 0xc1b7f0>



Definitely have to fill with median. Nothing else can really help over here.

# Cleaning Fee

```
In [32]:    1  listings['cleaning_fee'] = listings['cleaning_fee'].str.strip('$').str.repla
```
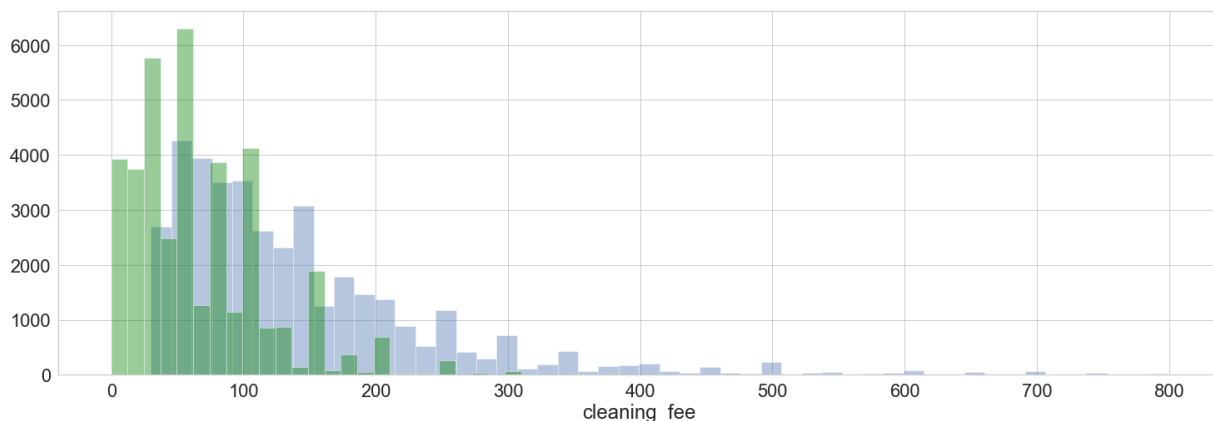
```
In [33]:    1  tmp = listings[["price", "cleaning_fee"]]
            2  tmp = tmp.dropna()
            3  tmp.corr()
```

Out[33]:

|              | price    | cleaning_fee |
|--------------|----------|--------------|
| **price**        | 1.000000 | 0.554281     |
| **cleaning_fee** | 0.554281 | 1.000000     |

```
In [34]:    1  f, ax = plt.subplots(1, 1, figsize=(25, 8))
            2  sns.distplot(tmp["price"], ax=ax, hist=True, kde=False)
            3  sns.distplot(tmp["cleaning_fee"], ax=ax, hist=True, kde=False, color="green"
```
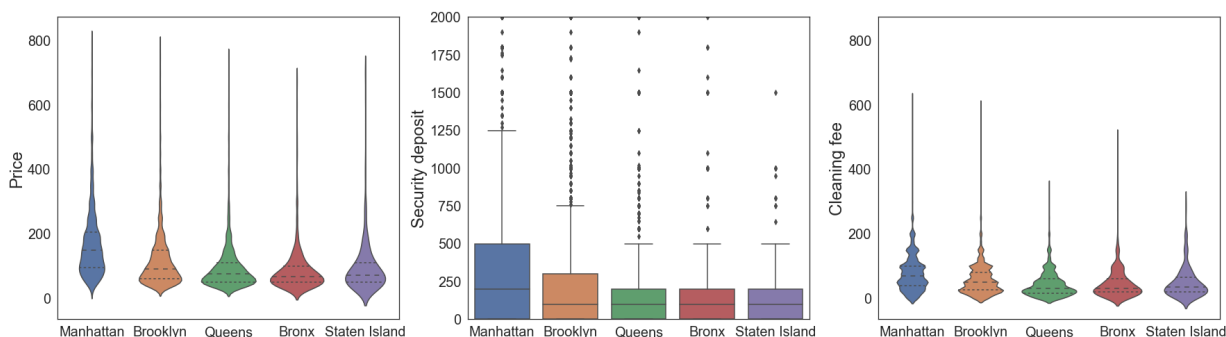
Out[34]:  <matplotlib.axes._subplots.AxesSubplot at 0x10178d0>

Definitely have to fill with median. Nothing else can really help over here.

# Neighbourhood Group

```
In [35]:  1  listings["neighbourhood_group_cleansed"].unique()
```

```
Out[35]:  array(['Manhattan', 'Brooklyn', 'Queens', 'Bronx', 'Staten Island'],
                dtype=object)
```

```
In [36]:   1  sns.set_style("white")
           2
           3  f, ax = plt.subplots(1,3,figsize=(30, 8))
           4  g = sns.violinplot(x="neighbourhood_group_cleansed", y="price", data=listing
           5  # t = g.set_title("Price vs.Neighbourhood Group")
           6  t = g.set_ylabel("Price")
           7  t = g.set_xlabel("")
           8  t = g.tick_params(labelsize=20)
           9  matching_ylim = g.get_ylim()
          10
          11  g = sns.boxplot(x="neighbourhood_group_cleansed", y="security_deposit", data
          12  t = g.set_title("")
          13  t = g.set_ylim([0, 2000])
          14  t = g.set_xlabel("")
          15  t = g.set_ylabel("Security deposit")
          16  t = g.tick_params(labelsize=20)
          17
          18  g = sns.violinplot(x="neighbourhood_group_cleansed", y="cleaning_fee", data=
          19  t = g.set_title("")
          20  t = g.set_ylim(matching_ylim)
          21  t = g.set_xlabel("")
          22  t = g.set_ylabel("Cleaning fee")
          23  t = g.tick_params(labelsize=20)
          24
```

# Property Type and Bedrooms

Club the lower frequency elements together

In [37]:
```python
strings = ("Apartment", "House", "Townhouse", "Loft", "Condominium", "Servic
apartment_list = list([])
for line in listings['property_type']:
    if any(s in line for s in strings):
        apartment_list.append('yes')
    else:
        apartment_list.append('no')

listings['prop'] = apartment_list
listings.loc[listings['prop'] == 'no', 'property_type'] = 'Other'
listings.loc[listings['property_type'] == 'Houseboat', 'property_type'] = 'O
listings.drop(['prop'], axis=1, inplace=True)
```

In [38]:
```python
listings["property_type"].value_counts()
```

Out[38]:
```
Apartment              39503
House                   3603
Townhouse               1557
Loft                    1447
Condominium             1370
Other                   1066
Serviced apartment       705
Name: property_type, dtype: int64
```
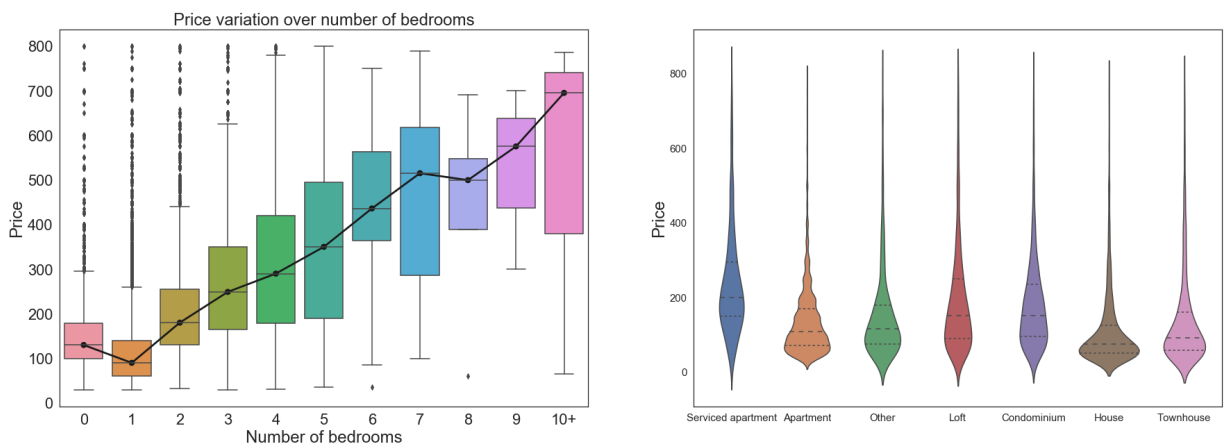
Group all bedrooms 10 and above into one category 10+

In [39]:
```python
def group_bedrooms(row):
    if row < 10:
        return str(int(row))
    elif row >= 10:
        return "10+"
    else:
        return row

listings["bedroom_group"] = listings["bedrooms"].apply(group_bedrooms)
# listings[["bedrooms", "bedroom_group"]]
```

In [40]:
```python
f, ax = plt.subplots(1,2,figsize=(30, 10))
# plot_box(x="bedroom_group", y="price", data=listings, agg_rule="median", a

listings['bedroom_group'] = pd.Categorical(
    listings['bedroom_group'],
    categories=['0','1','2','3','4','5','6', '7', '8', '9', '10+'],
    ordered=True
)
listings.sort_values(by="bedroom_group", inplace=True)

g = sns.boxplot(x="bedroom_group", y="price", data=listings, ax=ax[0])

agg_data = listings[["price", "bedroom_group"]].groupby(by=["bedroom_group"]
g = sns.pointplot(x="bedroom_group", y="price", data=agg_data, ax=ax[0], col
t = g.set(title="Price variation over number of bedrooms", xlabel="Number of

g = sns.violinplot(x="property_type", y="price", data=listings, ax=ax[1], in
ty = g.set(title="")
t = g.set_ylabel("Price")
t = g.set_xlabel("")
t = g.tick_params(labelsize=15)
```



In [56]:
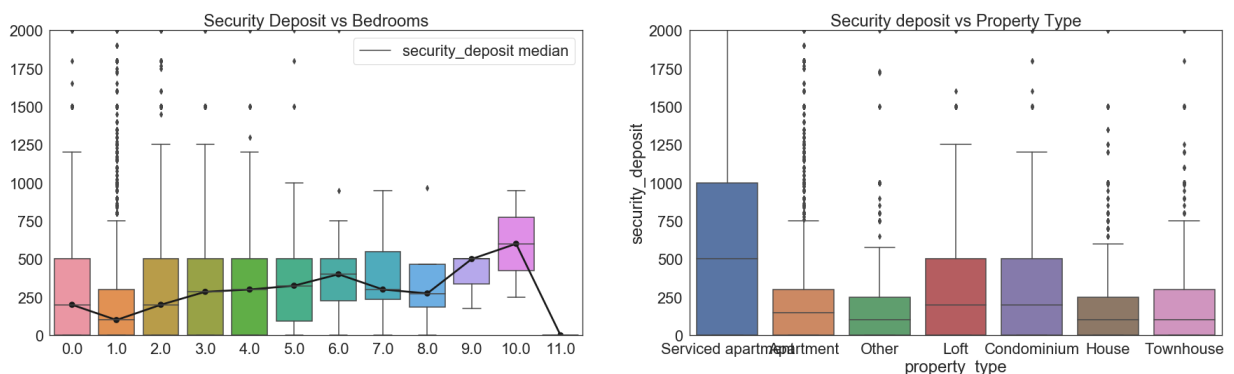```python
f, ax = plt.subplots(1,2,figsize=(30, 8))
plot_box(x="bedrooms", y="security_deposit", data=listings, agg_rule="median
g = sns.boxplot(x="property_type", y="security_deposit", data=listings, ax=a
g.set(title="Security deposit vs Property Type")
yl = g.set_ylim(0,2000)
```
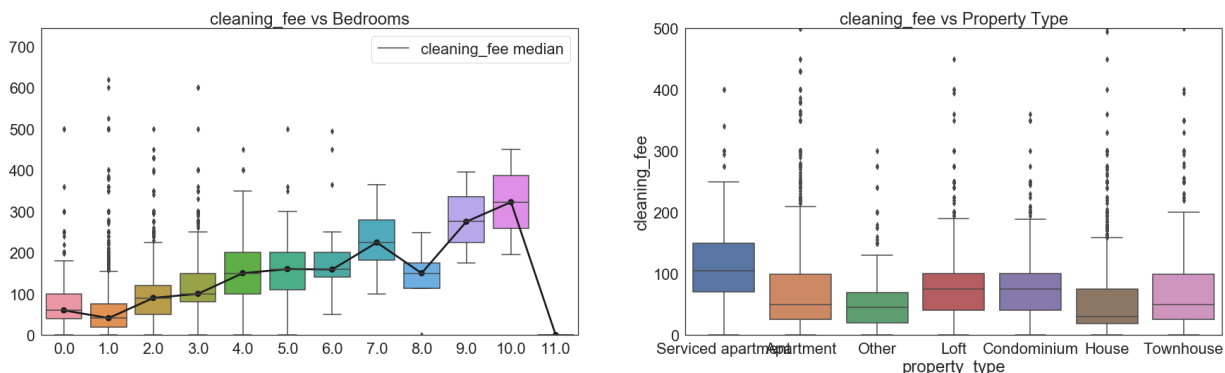
```
In [42]:    1  f, ax = plt.subplots(1,2,figsize=(30, 8))
            2  plot_box(x="bedrooms", y="cleaning_fee", data=listings, agg_rule="median", a
            3  g = sns.boxplot(x="property_type", y="cleaning_fee", data=listings, ax=ax[1]
            4  g.set(title="cleaning_fee vs Property Type")
            5  yl = g.set_ylim(0,500)
```

How do amenities behave for each property type?
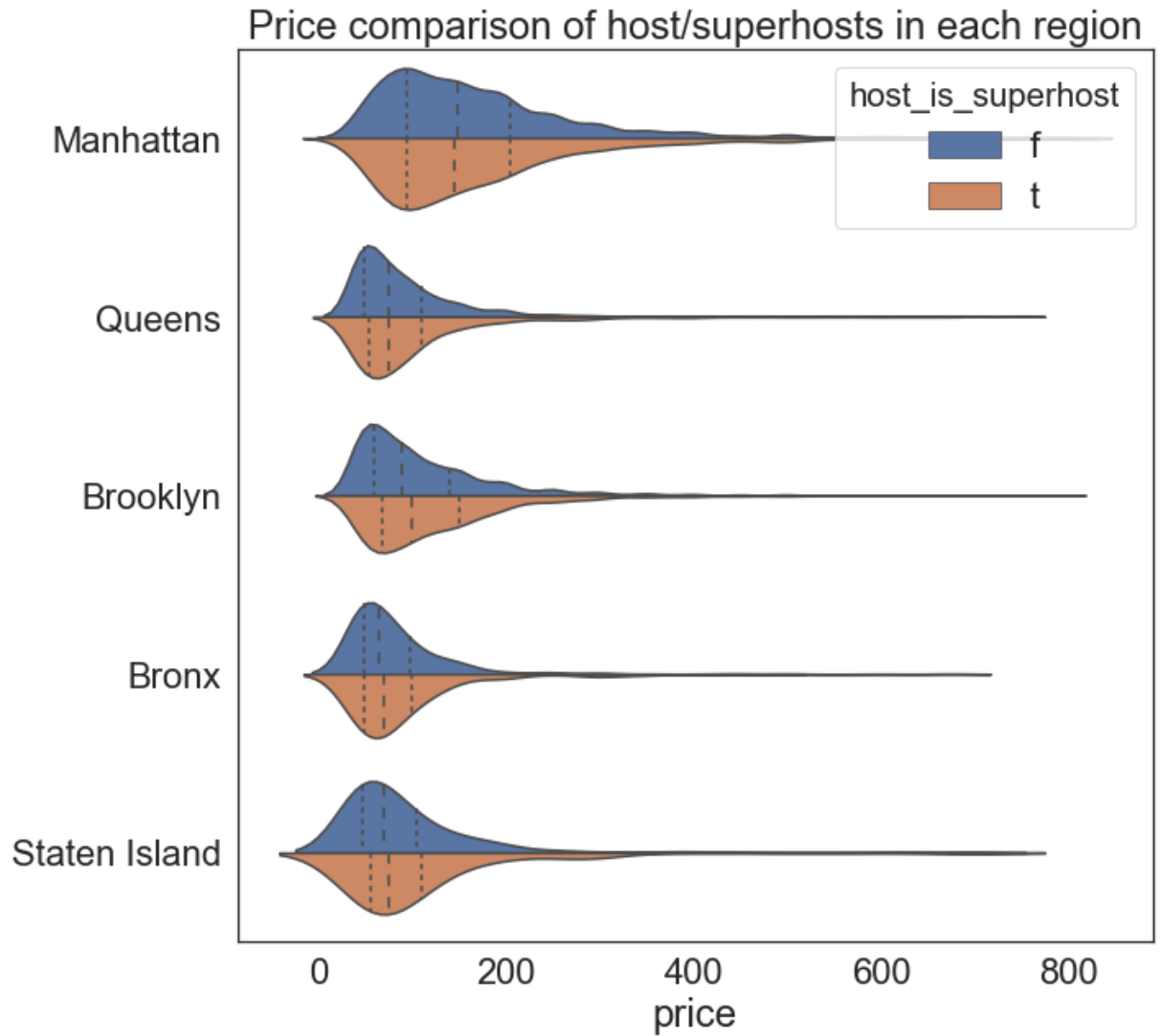
```
In [43]:    1  # listings[["num_amenities", "property_type"]].head()
            2  f, ax = plt.subplots(1, 2, figsize=(25, 8))
            3  g = sns.violinplot(x="property_type", y="num_amenities", data=listings, ax=a
            4  t = g.set(title="Distribution of amenities offered over Property types", xla
            5  g = sns.violinplot(x="neighbourhood_group_cleansed", y="num_amenities", data
            6  t = g.set(title="Distribution of amenities offered over Neighbourhood Group"
```

# Super Hosts

```
In [44]:    1  shost_subset = listings[["host_is_superhost", "neighbourhood_group_cleansed"
```

In [45]:
```python
import seaborn as sns
# sns.set(style="whitegrid")
f, ax = plt.subplots(1,1,figsize=(10, 10))
# g = sns.boxplot(y="neighbourhood_group_cleansed", x="price", data=shost_su
g = sns.violinplot(y="neighbourhood_group_cleansed", x="price", data=shost_s
t = g.set(title="Price comparison of host/superhosts in each region", ylabel
```



Price comparison of host/superhosts in each region
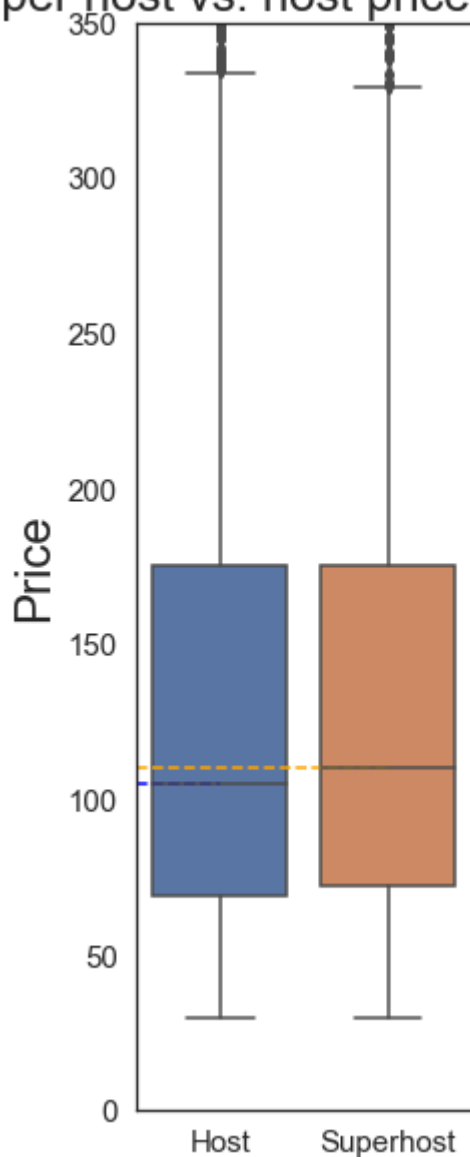
```
In [62]:    1  f, ax = plt.subplots(1, 1, figsize=(3, 10))
            2  shost_subset["dummy"] = 1
            3  g = sns.boxplot(x="host_is_superhost", y="price", data=shost_subset, ax=ax)
            4  yl = g.set(ylim=[0, 350], title="Super host vs. host price variation", xlabe
            5  g.set_xticklabels(["Host", "Superhost"])
            6  t = g.tick_params(labelsize=15)
            7  shost_median = shost_subset[shost_subset["host_is_superhost"]=="t"]["price"]
            8  host_median = shost_subset[shost_subset["host_is_superhost"]=="f"]["price"].
            9  l = plt.axhline(y=shost_median, xmax=0.75, color="orange", linestyle="--")
           10  l = plt.axhline(y=host_median, xmax=0.25, color="blue", linestyle="--")
           11  display(shost_median, host_median)
```

110.0

105.0



## Reviews Per month

How many listings have reviews and how many do not?

In [71]:
```
1  nan_df = analyse_nans(listings[["reviews_per_month"]])
2  nan_df.head()
```

Out[71]:

|  | reviews_per_month |
|---|---|
| total | 9933 |
| percentage | 20.2 |
| idx_list | [37563, 21942, 46229, 21527, 39753, 6802, 2230... |

In [85]:
```
1  num_nans= nan_df.iloc[1,:].values[0]
2  num_revs = 100 - num_nans
3  x = ['None', 'Available']
4  y = [num_nans, num_revs]
5  g = sns.barplot(x=x, y=y)
6  t = g.set(title="Listings with and without reviews")
```