

In [34]:



```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import matplotlib
5 from sklearn.model_selection import train_test_split
6 from sklearn.compose import ColumnTransformer
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn.model_selection import ParameterGrid
10 from sklearn.metrics import accuracy_score
11 from sklearn.metrics import mean_squared_error
12 from sklearn.model_selection import KFold
13 from sklearn.model_selection import GridSearchCV
14 from sklearn.pipeline import make_pipeline
15 from sklearn.pipeline import Pipeline
16 from sklearn.neighbors import KNeighborsRegressor
17 from sklearn.ensemble import RandomForestRegressor
18 from sklearn.linear_model import LinearRegression
19 from sklearn.linear_model import Lasso
20 from sklearn.linear_model import Ridge
21 from sklearn.linear_model import ElasticNet
22 from sklearn.svm import SVR
```

dataset https://www.kaggle.com/datasets/ramjasmaurya/oyo-rental-price-prediction-in-china/code?select=rental_price.csv (https://www.kaggle.com/datasets/ramjasmaurya/oyo-rental-price-prediction-in-china/code?select=rental_price.csv).

In [2]:



```
1 df = pd.read_csv("data/rentoyo.csv")
2 y=df['price']
3 X=df.loc[:,df.columns!='price']
4 print("feature matrix shape:", X.shape)
5 print("target variable shape:", y.shape)
```

feature matrix shape: (5834, 25)
target variable shape: (5834,)

missing values

In [3]:



```
1 #count nan
2 # print(df.isnull().sum())
3 miss=df.isnull().sum(axis=0)/df.shape[0]
4 print("fraction of missing values in features:")
5 print(miss[miss>0])
6
7 frac_missing = sum(df.isnull().sum(axis=1)!=0)/df.shape[0]
8 print('fraction of points with missing values:', frac_missing)
```

```
fraction of missing values in features:
bathrooms          0.007885
bedrooms           0.001028
beds               0.003942
host_is_superhost  0.002571
host_listings_count 0.002571
review_scores_checkin 0.352588
review_scores_communication 0.352588
review_scores_location 0.352417
review_scores_rating 0.350703
review_scores_value 0.352588
dtype: float64
fraction of points with missing values: 0.35978745286253
```

target variable

In [4]:



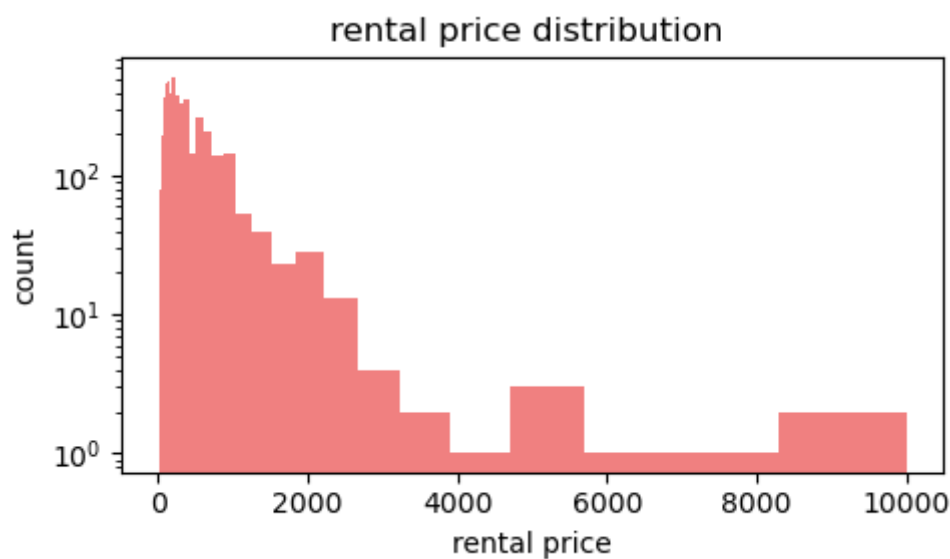
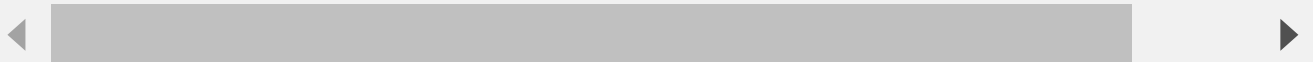
```
1 print(y.describe())
```

```
count    5834.000000
mean      286.219918
std       403.256199
min        0.000000
25%       95.000000
50%      175.000000
75%      325.000000
max     10000.000000
Name: price, dtype: float64
```

In [5]:



```
1 #log form
2 plt.figure(figsize=(5,3))
3 y.plot.hist(log=True,bins = np.logspace(np.log10(1),np.log10(np.max(df['price'])),50),color="red")
4
5 plt.xlabel('rental price')
6 plt.ylabel('count')
7 plt.title("rental price distribution ")
8
9 fig=plt.gcf()
10 plt.tight_layout()
11 plt.savefig("price.png")
```



Exploratory Data Analysis

In [3]:



```
1 # drop worthless columns
2 X_new=X.loc[:,(X.columns!='amenities')&(X.columns!='has_availability')]
```

In [13]:

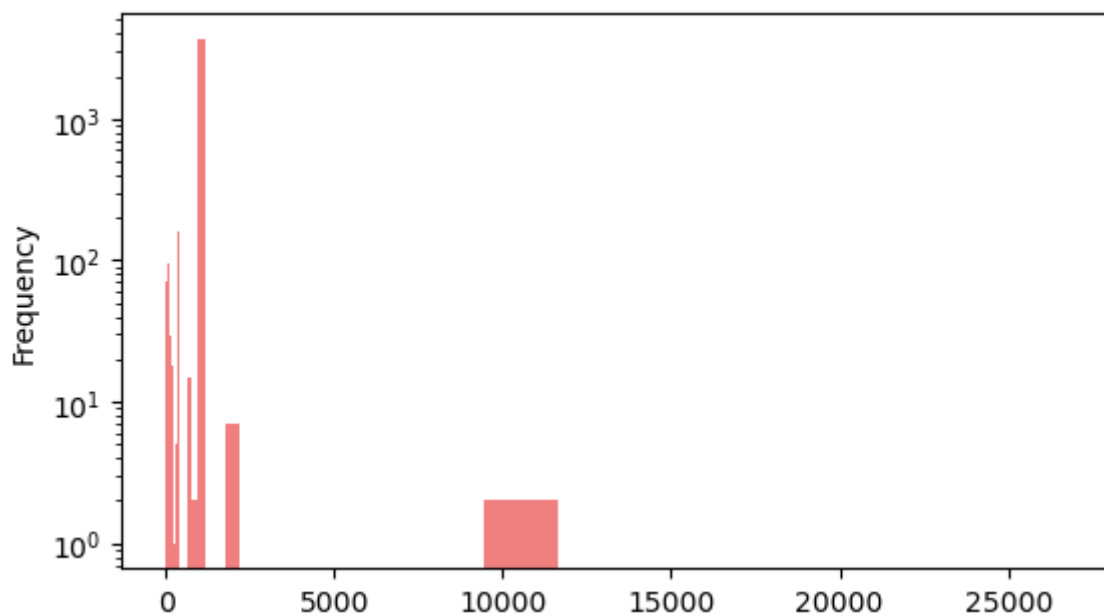


```
1 print(df['maximum_nights'].describe())
2 plt.figure(figsize=(6.4, 3.6))
3 df['maximum_nights'].plot.hist(log=True, bins = np.logspace(np.log10(1), np.log10(np.max(df['max
```

```
count      5834.000000
mean        746.705862
std         641.901800
min           1.000000
25%         30.000000
50%        1125.000000
75%        1125.000000
max       26801.000000
Name: maximum_nights, dtype: float64
```

Out[13]:

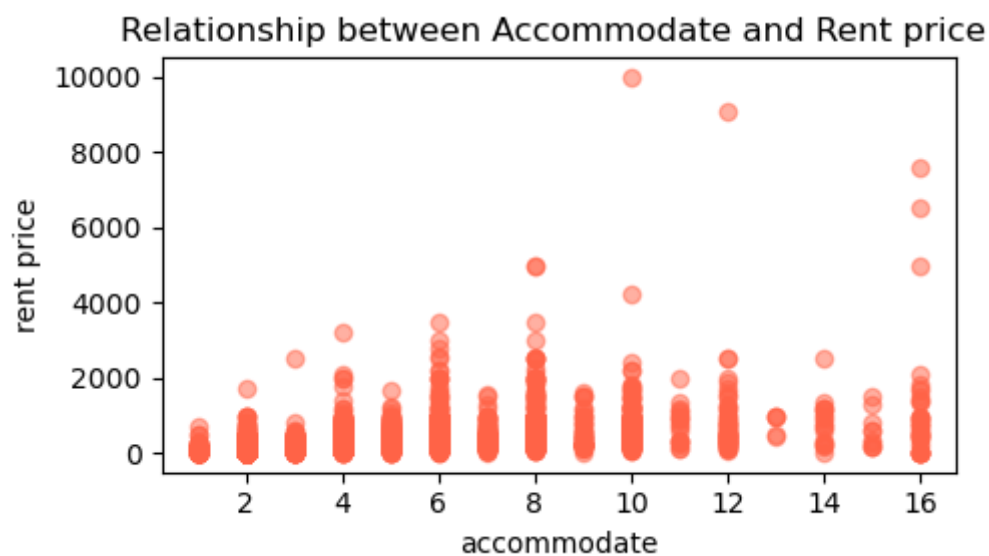
<AxesSubplot:ylabel='Frequency'>



In [16]:



```
1 # numerical feature scatter
2 plt.figure(figsize=(5,3))
3 X1=df['accommodates']
4 Y1=y
5
6 plt.xlabel("accommodate")
7 plt.ylabel("rent price")
8 plt.title("Relationship between Accommodate and Rent price ")
9
10 plt.scatter(X1,Y1,alpha=0.5,color="tomato")
11
12 fig=plt.gcf()
13 plt.tight_layout()
14 plt.savefig("acoomodate.png")
```



In [15]:



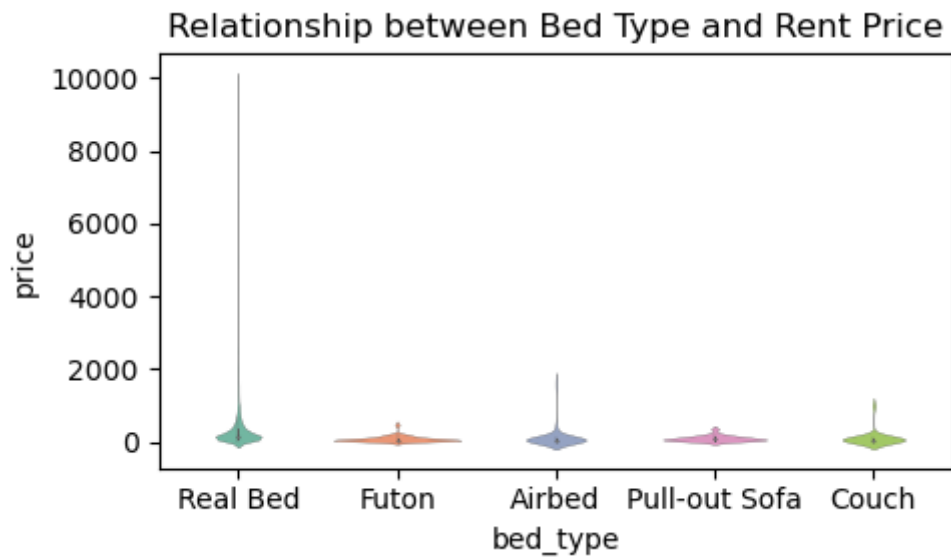
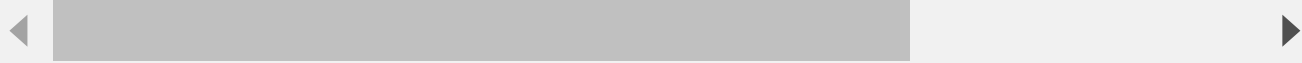
```
1 #property type & AVG price
2 plt.figure(figsize = (15,6))
3 df1=df[['property_type', 'price']]
4 df_group=df1.groupby('property_type')
5 property_type=df_group['price'].mean().reset_index().rename(columns={"property_type":"property_
6
7 ax = sns.barplot(data=property_type, x='property_type', y='avg price')
8 plt.xlabel("Property Type")
9 plt.ylabel("AVG Price")
10 ax.set(title='Relationship between Property Type and AVG Price')
11 g = ax.set_xticklabels(ax.get_xticklabels(),rotation = 90)
12
13 fig=plt.gcf()
14 plt.tight_layout()
15 plt.savefig("property.png")
```



In [18]:



```
1 #categorical violin
2 import seaborn as sns
3 plt.figure(figsize=(5,3))
4 sns.violinplot(x=df['bed_type'],y=df['price'],linewidth=0.2,palette="Set2").set(title='Relation
5
6 fig=plt.gcf()
7 plt.tight_layout()
8 plt.savefig("bedtype.png")
```



In [3]:



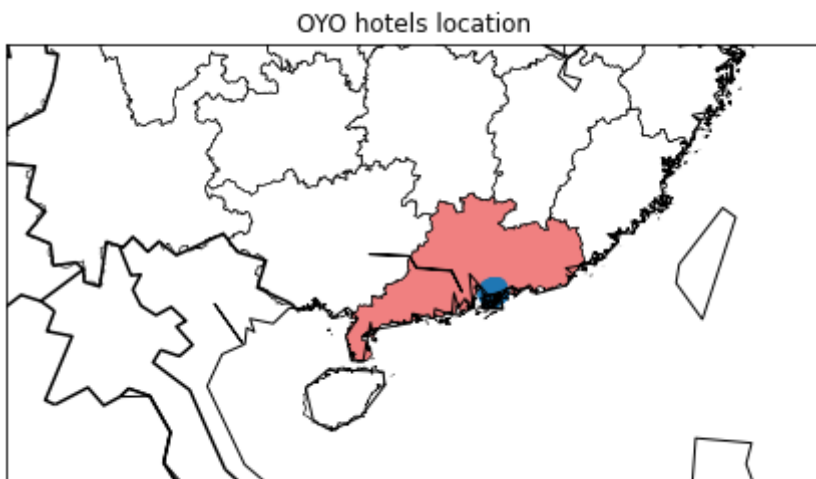
```
1 #latitude longitude basemap
2 #china
3 from mpl_toolkits.basemap import Basemap
4 import matplotlib.pyplot as plt
5 from matplotlib.patches import Polygon
6 ax = plt.gca()
7
8 # plot map
9 m = Basemap(llcrnrlon=82.33,
10             llcrnrlat=3.01,
11             urcrnrlon=138.16,
12             urcrnrlat=53.123,
13             projection='lcc', lat_0 = 42.5, lon_0=120)
14 m.drawcoastlines() # coastline
15 m.drawcountries(linewidth=1.5) # country border
16 # CHN_adml province data of China
17 m.readshapefile(shapefile='data/gadm41_CHN_shp/gadm41_CHN_1',
18                 name='states',
19                 drawbounds=True,color='r')
20
21 for info, shp in zip(m.states_info, m.states):
22     proid = info['NAME_1']
23     if proid == 'Guangdong':
24         poly = Polygon(shp, facecolor='orange', lw=3)
25         ax.add_patch(poly)
26
27 fig=plt.gcf()
28 plt.tight_layout()
29 plt.savefig("china.png")
```



In [4]:



```
1 #latitude longitude basemap
2 #guangdong & oyo's catter
3 from mpl_toolkits.basemap import Basemap
4 import matplotlib.pyplot as plt
5 from matplotlib.patches import Polygon
6 ax = plt.gca()
7
8 # plot map
9 m = Basemap(llcrnrlon=100.33,
10             llcrnrlat=15.01,
11             urcrnrlon=125.16,
12             urcrnrlat=30.12,
13             projection='lcc', lat_0 = 42.5, lon_0=120)
14 m.drawcoastlines() # coastline
15 m.drawcountries(linewidth=1.5) # country border
16 # CHN_adm1 province data of China
17 m.readshapefile(shapefile='D:\BrownUnivercity\DATA1030\midproject\data\gadm41_CHN_shp/gadm41_CH
18                  name='states',
19                  drawbounds=True)
20
21 lon=df['longitude(East)']
22 lat=df['latitude(North)']
23
24 for info, shp in zip(m.states_info, m.states):
25     proid = info['NAME_1']
26     if proid == 'Guangdong':
27         poly = Polygon(shp, facecolor='lightcoral', lw=3)
28         ax.add_patch(poly)
29
30 lon, lat = m(lon, lat)
31 m.scatter(lon, lat, s=60)
32
33 plt.title('OYO hotels location')
34
35 fig=plt.gcf()
36 plt.tight_layout()
37 plt.savefig("guangdong.png")
```



Method ML_pipe_KFold_RMSE

In [33]:



```
1 from sklearn.preprocessing import OneHotEncoder
2 from sklearn.preprocessing import OrdinalEncoder
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.preprocessing import MinMaxScaler
5 from sklearn.impute import SimpleImputer
6 from sklearn.experimental import enable_iterative_imputer
7 from sklearn.impute import IterativeImputer
```

In [26]:



```
1 #regression problem
2 nr_states = 10 #loop 10 times
3
4
5 final_models = []
6
7 def MLpipe_KFold_RMSE(X, y, preprocessor, reg, param_grid):
8     test_scores = np.zeros(nr_states)
9     for i in range(nr_states):
10
11         X_other, X_test, y_other, y_test = train_test_split(X, y, test_size = 0.2, random_state=42*i)
12         kf = KFold(n_splits=4, shuffle=True, random_state=42*i)
13         pipe = make_pipeline(preprocessor, reg)
14         grid = GridSearchCV(pipe, param_grid=param_grid, scoring = 'neg_mean_squared_error',
15                             cv=kf, return_train_score = True, n_jobs=-1, verbose=True)
16         grid.fit(X_other, y_other)
17         results = pd.DataFrame(grid.cv_results_)
18
19         print('best model parameters:', grid.best_params_)
20         print('validation score:', grid.best_score_)
21     # save the model
22     final_models.append(grid)
23     # calculate and save the test score
24     y_test_pred = final_models[-1].predict(X_test)
25     test_scores[i] = np.sqrt(mean_squared_error(y_test, y_test_pred))
26     print('RMSE test score:', test_scores[i])
27     return test_scores, y_test
```

In [10]:



```
1 #read data
2 #preprocessor
3 # collect the various features
4 cat_ftrs = ['bed_type', 'host_is_superhost', 'instant_bookable', 'room_type']
5
6 ordinal_ftrs = ['cancellation_policy', 'property_type']
7 ordinal_cats = [['no_refunds', 'super_strict_30', 'strict', 'moderate', 'flexible'],
8                 ['Hut', 'Condominium', 'Apartment', 'Cabin', 'Villa', 'Boat', 'Tipi', 'Townhouse']]
9
10 num_ftrs1 = ['accommodates', 'availability_30', 'calculated_host_listings_count',
11             'guests_included', 'number_of_reviews', 'bathrooms', 'bedrooms', 'beds', 'host_listings_c',
12             'review_scores_checkin', 'review_scores_communication', 'review_scores_location',
13             'review_scores_rating', 'review_scores_value']
14 num_ftrs2 = ['maximum_nights']
15
16 # one-hot encoder
17 categorical_transformer = Pipeline(steps=[
18     ('imputer', SimpleImputer(strategy='constant', fill_value='missing')),
19     ('onehot', OneHotEncoder(sparse=False, handle_unknown='ignore'))])
20
21 # ordinal encoder
22 ordinal_transformer = Pipeline(steps=[
23     ('imputer2', SimpleImputer(strategy='constant', fill_value='NA')),
24     ('ordinal', OrdinalEncoder(categories = ordinal_cats))])
25
26 # MinMax Scaler
27 numeric_transformer1 = Pipeline(steps=[
28     ('imputer3', SimpleImputer(strategy='mean')),
29     ('scaler', MinMaxScaler(feature_range=(0, 100)))]])
30
31 # Standard scaler
32 numeric_transformer2 = Pipeline(steps=[
33     ('imputer4', SimpleImputer(strategy='mean')),
34     ('scaler', StandardScaler())])
35
36 # collect the transformers
37 preprocessor = ColumnTransformer(
38     transformers=[
39         ('num1', numeric_transformer1, num_ftrs1),
40         ('cat', categorical_transformer, cat_ftrs),
41         ('ord', ordinal_transformer, ordinal_ftrs),
42         ('num2', numeric_transformer2, num_ftrs2)])
```

In [27]:

```
1 # L1 regularized linear regression
2 reg = Lasso(random_state=42)
3 param_grid = {'lasso__alpha': [0.25, 2.5, 25] }
4 L1_c, y_test = MLpipe_KFold_RMSE(X_new, y, preprocessor, reg, param_grid)
```

Fitting 4 folds for each of 3 candidates, totalling 12 fits
best model parameters: {'lasso__alpha': 0.25}
validation score: -95544.89204288769
RMSE test score: 332.5048195343838
Fitting 4 folds for each of 3 candidates, totalling 12 fits
best model parameters: {'lasso__alpha': 0.25}
validation score: -99120.16158444743
RMSE test score: 321.14850380149824
Fitting 4 folds for each of 3 candidates, totalling 12 fits
best model parameters: {'lasso__alpha': 0.25}
validation score: -108367.86935505274
RMSE test score: 249.33466630590902
Fitting 4 folds for each of 3 candidates, totalling 12 fits
best model parameters: {'lasso__alpha': 0.25}
validation score: -90777.64584320017
RMSE test score: 369.3208920924873
Fitting 4 folds for each of 3 candidates, totalling 12 fits
best model parameters: {'lasso__alpha': 2.5}
validation score: -105068.33369778495

In [12]:

```
1 L1_c
```

Out[12]:

```
array([332.50481953, 321.1485038 , 249.33466631, 369.32089209,
       298.27287924, 349.18486624, 364.13289666, 258.50704218,
       229.58742343, 337.10134802])
```

In [29]:

```
1 # y_test
```

In [14]:

```
1 print('test scores mean', np.mean(L1_c))
2 print('test scores standard deviation', np.std(L1_c))
```

test scores mean 310.9095337515326
test scores standard deviation 47.198807526529805

In [17]:

```
1 #RandomForestRegressor
2
3 reg = RandomForestRegressor(random_state=42)
4 param_grid = {'randomforestregressor__max_features': [3,5,7,9] ,
5               'randomforestregressor__max_depth': [3,5,7,9]}
6 ran_c=MLpipe_KFold_RMSE(X_new, y, preprocessor, reg, param_grid)
7
```

Fitting 4 folds for each of 16 candidates, totalling 64 fits
best model parameters: {'randomforestregressor__max_depth': 9, 'randomforestregressor__max_features': 7}
validation score: -81654.734316264
RMSE test score: 307.10369092624916
Fitting 4 folds for each of 16 candidates, totalling 64 fits
best model parameters: {'randomforestregressor__max_depth': 9, 'randomforestregressor__max_features': 3}
validation score: -93287.08835230021
RMSE test score: 305.09515619169105
Fitting 4 folds for each of 16 candidates, totalling 64 fits
best model parameters: {'randomforestregressor__max_depth': 9, 'randomforestregressor__max_features': 5}
validation score: -94111.43843602226
RMSE test score: 215.86303582933243
Fitting 4 folds for each of 16 candidates, totalling 64 fits
best model parameters: {'randomforestregressor__max_depth': 9, 'randomforestregressor__max_features': 7}
validation score: -79966.11391162813

In [18]:

```
1 ran_c
```

Out[18]:

```
array([307.10369093, 305.09515619, 215.86303583, 344.76157886,
       292.34805259, 336.76020741, 360.37117433, 239.62375777,
       218.45557428, 299.79961715])
```

In [107]:

```
1 print('test scores mean', np.mean(ran_c))
2 print('test scores standard deviation', np.std(ran_c))
```

test scores mean 292.01818453333726
test scores standard deviation 48.858289905085414

In [19]:

```
1 # SVR
2 reg = SVR(kernel='rbf')
3 param_grid = {'svr__gamma': [0.1, 10, 100],
4               'svr__C': [0.1, 1, 10]}
5 SVR_c=MLpipe_KFold_RMSE(X, y, preprocessor, reg, param_grid)
```

Fitting 4 folds for each of 9 candidates, totalling 36 fits
best model parameters: {'svr__C': 10, 'svr__gamma': 0.1}
validation score: -166525.97989233583
RMSE test score: 457.0633809705048
Fitting 4 folds for each of 9 candidates, totalling 36 fits
best model parameters: {'svr__C': 10, 'svr__gamma': 0.1}
validation score: -170334.21140706414
RMSE test score: 439.1388000674527
Fitting 4 folds for each of 9 candidates, totalling 36 fits
best model parameters: {'svr__C': 10, 'svr__gamma': 0.1}
validation score: -186566.7821397371
RMSE test score: 352.79752267955473
Fitting 4 folds for each of 9 candidates, totalling 36 fits
best model parameters: {'svr__C': 10, 'svr__gamma': 0.1}
validation score: -155341.4636547502
RMSE test score: 502.45544361279406
Fitting 4 folds for each of 9 candidates, totalling 36 fits
best model parameters: {'svr__C': 10, 'svr__gamma': 0.1}
validation score: -177387.0667796479

In [20]:

```
1 SVR_c
```

Out[20]:

```
array([457.06338097, 439.13880007, 352.79752268, 502.45544361,
       409.21875057, 442.50978551, 495.11179577, 357.68006883,
       299.07785263, 460.71044277])
```

In [109]:

```
1 print('test scores mean', np.mean(SVR_c))
2 print('test scores standard deviation', np.std(SVR_c))
```

test scores mean 421.57638434199333
test scores standard deviation 62.81576114486148

In [21]:

```
1 #KNN
2 reg = KNeighborsRegressor(weights='uniform', n_jobs=-1)
3 param_grid = {'kneighborsregressor__n_neighbors': [5, 25, 50]}
4 KNN_c=MLpipe_KFold_RMSE(X, y, preprocessor, reg, param_grid)
```

best model parameters: {'kneighborsregressor__n_neighbors': 25}
validation score: -94949.25161277174
RMSE test score: 335.33792052276357
Fitting 4 folds for each of 3 candidates, totalling 12 fits
best model parameters: {'kneighborsregressor__n_neighbors': 25}
validation score: -91057.81310763108
RMSE test score: 384.5499262215294
Fitting 4 folds for each of 3 candidates, totalling 12 fits
best model parameters: {'kneighborsregressor__n_neighbors': 25}
validation score: -105365.10525826583
RMSE test score: 252.86243300463084
Fitting 4 folds for each of 3 candidates, totalling 12 fits
best model parameters: {'kneighborsregressor__n_neighbors': 25}
validation score: -108226.16542301633
RMSE test score: 222.87348306603138
Fitting 4 folds for each of 3 candidates, totalling 12 fits
best model parameters: {'kneighborsregressor__n_neighbors': 25}
validation score: -94902.09452370816
RMSE test score: 338.9797156249049

In [22]:

```
1 KNN_c
```

Out[22]:

```
array([318.20838074, 320.60138962, 265.89838554, 371.5695846 ,
       289.94381831, 335.33792052, 384.54992622, 252.862433 ,
       222.87348307, 338.97971562])
```

In [111]:

```
1 print('test scores mean', np.mean(KNN_c))
2 print('test scores standard deviation', np.std(KNN_c))
```

test scores mean 310.0825037242632
test scores standard deviation 49.20667904497913

Results

In [23]:



```
1 allscores=np.r_[L1_c, ran_c, SVR_c, KNN_c]#array 行合并 （列合并 np.c_）
2 print(allscores)
```

```
[332.50481953 321.1485038 249.33466631 369.32089209 298.27287924
349.18486624 364.13289666 258.50704218 229.58742343 337.10134802
307.10369093 305.09515619 215.86303583 344.76157886 292.34805259
336.76020741 360.37117433 239.62375777 218.45557428 299.79961715
457.06338097 439.13880007 352.79752268 502.45544361 409.21875057
442.50978551 495.11179577 357.68006883 299.07785263 460.71044277
318.20838074 320.60138962 265.89838554 371.5695846 289.94381831
335.33792052 384.54992622 252.862433 222.87348307 338.97971562]
```

In [24]:



```
1 print('test scores mean', np.mean(allscores))
2 print('test scores standard deviation', np.std(allscores))
```

```
test scores mean 333.64665158778155
test scores standard deviation 73.34642083766414
```

In [30]:



```
1 # baseline
2 from sklearn.metrics import mean_squared_error
3 y_mean=np.mean(y_test)
4 y_std=np.std(y_test, ddof=0)
5 #create full y_mean ndarray
6 array=np.full((len(y_test),1), y_mean)
7 y_pred1=pd.DataFrame(array)
8
9 print("baseline RMSE:", np.sqrt(mean_squared_error(y_test, y_pred1)))
```

```
baseline RMSE: 445.7852569807543
```

In [31]:



```
1 # test scores of different models
2 test_temp = {'model': ['Lasso', 'RandomForest', 'SVR', 'KNeighbor'],
3              'test_score': [np.mean(L1_c), np.mean(ran_c), np.mean(SVR_c), np.mean(KNN_c)]}
4 test_data = pd.DataFrame.from_dict(test_temp)
5 # test_data
```

In [74]:

```
1 base=np.full((len(test_data)), 445.78)
2 error=(test_data['test_score'].values-base)/73.34
3 error
4 # error=(test_data['test_score'].values-base1)/73.34
5 # error #Series
```

Out[74]:

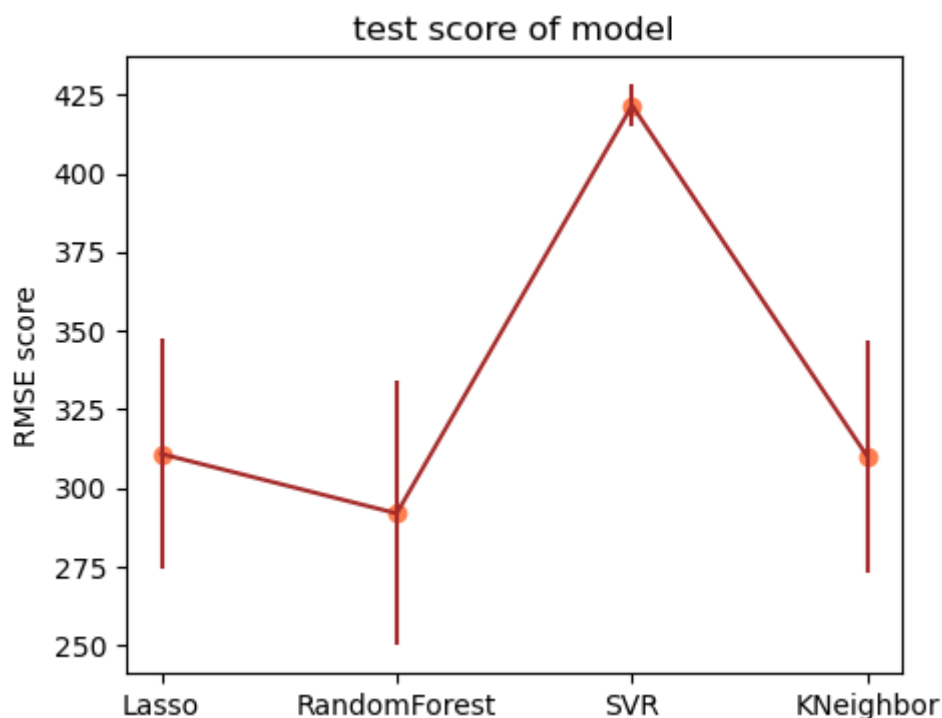
```
array([-1.83897554, -2.09656143, -0.3300193 , -1.8502522 ])
```

In [82]:

```
1 plt.rcParams.update({'font.size': 10})
2 plt.figure(figsize=(5,4))
3 # plt.plot(test_data['model'], test_data['test_score'], c='brown')
4 plt.scatter(test_data['model'], test_data['test_score'], c='coral')
5 plt.errorbar(test_data['model'], test_data['test_score'], yerr=-error*20, c='brown')
6 plt.ylabel('RMSE score')
7 plt.title('test score of model')
8 plt.plot()
```

Out[82]:

[]



errorbar <https://vimsky.com/examples/usage/matplotlib-pyplot-errorbar-in-python.html>
(<https://vimsky.com/examples/usage/matplotlib-pyplot-errorbar-in-python.html>)

In [84]:



```
1 # the best model is RandomForestRegressor
2 '''
3 best model parameters: {'randomforestregressor__max_depth': 9, 'randomforestregressor__max_feat
4 validation score: -94111.43843602226
5 RMSE test score: 215.8630358293324
6 '''
7 import time
8 start=time.time()
9 random_state=30
10 #split
11 X_train, X_other, y_train, y_other = train_test_split(X_new,y,test_size = 0.2,random_state=rand
12 X_val, X_test, y_val, y_test = train_test_split(X_other,y_other,test_size = 0.5,random_state=ra
13
14
15 # fit_transform the training set
16 X_prep = preprocessor.fit_transform(X_train)
17
18 #the feature names after fit
19 feature_names = preprocessor.get_feature_names_out()
20
21 #transform the train
22 df_train = pd.DataFrame(data=X_prep, columns=feature_names)
23 print(df_train.shape)
24
25 #transform the val
26 df_val = preprocessor.transform(X_val)
27 df_val = pd.DataFrame(data=df_val, columns = feature_names)
28 print(df_val.shape)
29
30 #transform the test
31 df_test = preprocessor.transform(X_test)
32 df_test = pd.DataFrame(data=df_test, columns = feature_names)
33
34 #fit model
35 reg = RandomForestRegressor(random_state=42,max_features=9,max_depth=5)
36 reg.fit(df_train,y_train)
37
38 y_test_pred = reg.predict(df_test)
39 test_scores = np.sqrt(mean_squared_error(y_test,y_test_pred))
40
41 end=time.time()
42
43 print('running time',end-start)
44 print('test scores mean',np.mean(test_scores))
45 print('test scores standard deviation',np.std(test_scores))
```

(4667, 30)

(583, 30)

running time 0.26233458518981934

test scores mean 202.72660072854498

test scores standard deviation 0.0

• Global

In [69]:



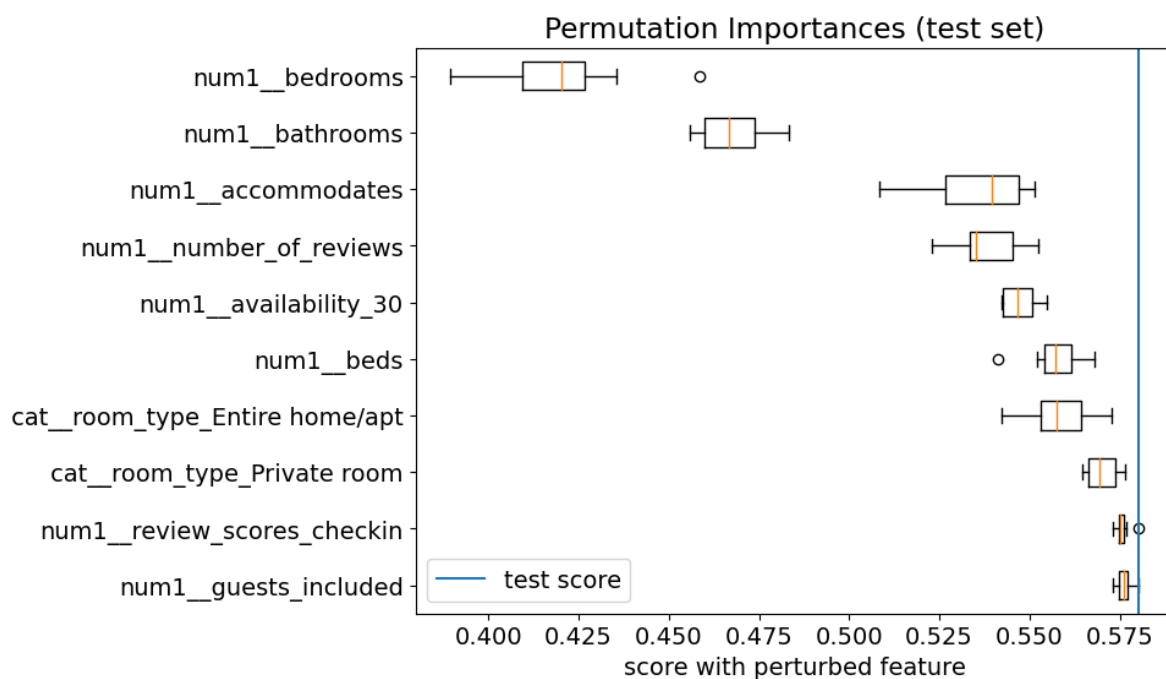
```
1 # G1 shuffling
2 from sklearn.metrics import r2_score
3 import warnings
4 warnings.filterwarnings("ignore")
5
6 test_score=r2_score(y_test,y_test_pred)
7 nr_runs = 10
8 ftr_names=df_test.columns
9 scores = np.zeros([len(ftr_names),nr_runs])
10
11
12 for i in range(len(ftr_names)):
13     print('shuffling '+str(ftr_names[i]))
14     r2_scores = []
15     for j in range(nr_runs):
16         X_test_shuffled = df_test.copy()
17         X_test_shuffled[ftr_names[i]] = np.random.permutation(df_test[ftr_names[i]].values)
18         y_test_pred=reg.predict(X_test_shuffled)
19         r2_scores.append(r2_score(y_test,y_test_pred))
20     print('    shuffled test score:', np. around(np. mean(r2_scores), 3), '+/-' , np. around(np. std(r2_s
21     scores[i] = r2_scores
22
```

```
shuffling num1__accommodates
    shuffled test score: 0.536 +/- 0.013
shuffling num1__availability_30
    shuffled test score: 0.547 +/- 0.005
shuffling num1__calculated_host_listings_count
    shuffled test score: 0.579 +/- 0.001
shuffling num1__guests_included
    shuffled test score: 0.576 +/- 0.002
shuffling num1__number_of_reviews
    shuffled test score: 0.538 +/- 0.009
shuffling num1__bathrooms
    shuffled test score: 0.468 +/- 0.009
shuffling num1__bedrooms
    shuffled test score: 0.42 +/- 0.018
shuffling num1__beds
    shuffled test score: 0.557 +/- 0.007
shuffling num1__host_listings_count
    shuffled test score: 0.58 +/- 0.001
shuffling num1__review_scores_checkin
```

In [70]:



```
1 #plot
2 #The smaller the R2, the more important the features are
3 import matplotlib.pyplot as plt
4 plt.rcParams.update({'font.size': 14})
5 plt.figure(figsize=(8,6))
6
7 sorted_indices = np.argsort(np.mean(scores,axis=1))[:,1]
8 sorted_indices = np.array(sorted_indices)[9::-1] #top 10
9
10 plt.boxplot(scores[sorted_indices].T, labels=ftr_names[sorted_indices], vert=False)
11 plt.axvline(test_score, label='test score')
12
13 plt.title("Permutation Importances (test set)")
14 plt.xlabel('score with perturbed feature')
15 plt.legend()
16 plt.show()
17
18 # fig=plt.gcf()
19 # plt.tight_layout()
20 # plt.savefig("Permutation Importances ")
```



In [86]:

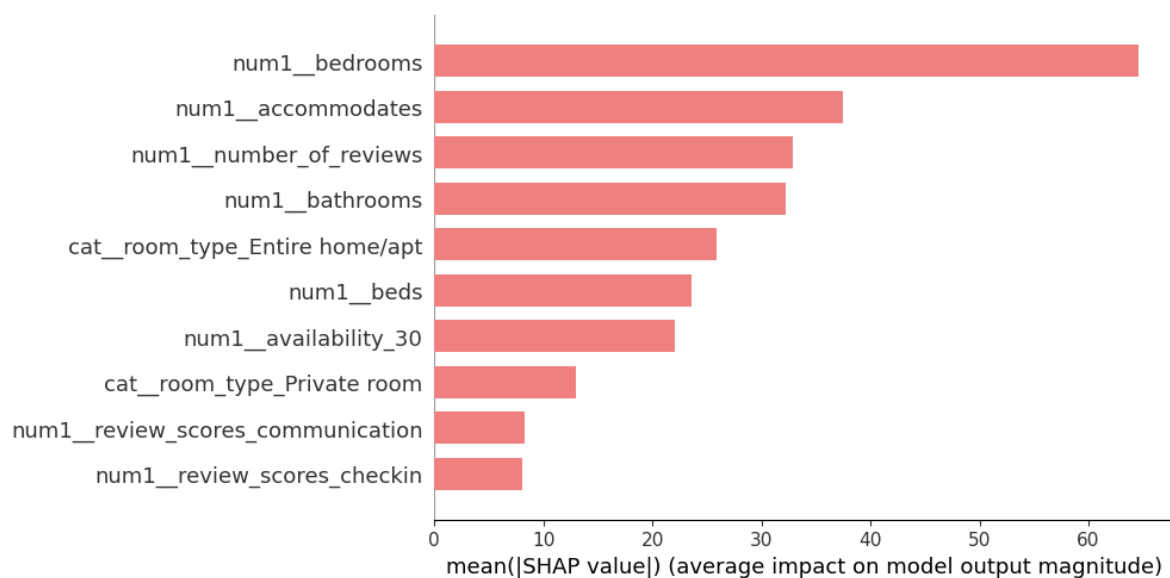
```
1 #G2 shap.global
2 import shap
3 shap.initjs()
4
5 # create the explainer object with RandomForestRegressor
6 # model = grid.best_estimator_
7 explainer = shap.TreeExplainer(reg)
```



In [88]:

```
1 shap_values = explainer.shap_values(df_test)
2 print(np.shape(shap_values))
3 #Global features importances
4 shap.summary_plot(shap_values, df_test, max_display=10, plot_type='bar', color='lightcoral')
```

(584, 30)



In [101]:

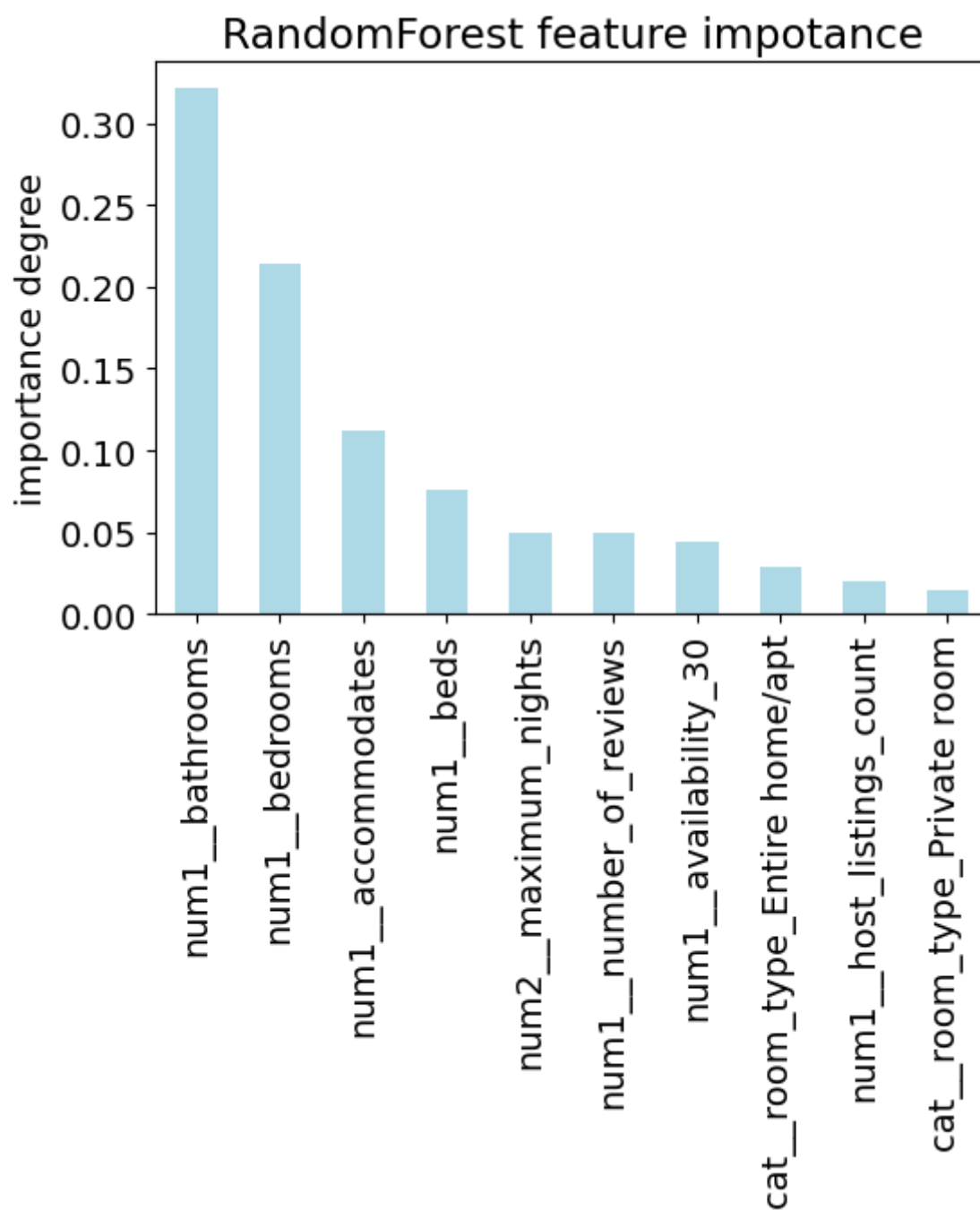
```
1 #G3 Randomforest SKlearn important features
2 # importances = reg.feature_importances_
3 importances = np.sort(reg.feature_importances_)[:-1]
4 loc=np.argsort(reg.feature_importances_)[:-1]
5 index_name=[]
6 for i in loc:
7     index_name.append(ftr_names[i])
```

In [110]:

```
1 plt.figure(figsize=(6,4))
2 forest_importances = pd.Series(importances[:10], index=index_name[:10])
3 forest_importances.plot.bar(color='lightblue')
4 plt.ylabel('importance degree')
5 plt.title('RandomForest feature impotence')
```

Out[110]:

Text(0.5, 1.0, 'RandomForest feature impotence')

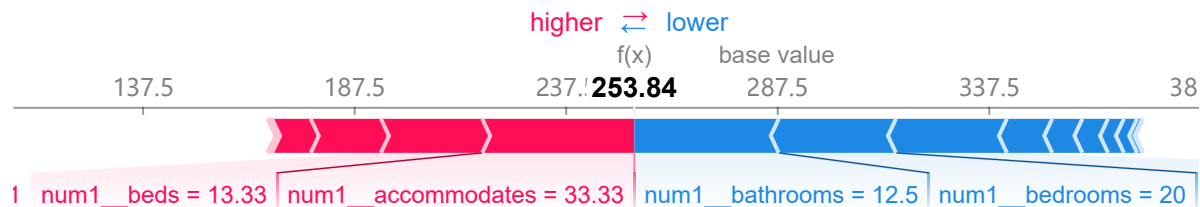


• Local

In [89]:

```
1 #point 0
2 shap.force_plot(explainer.expected_value, shap_values[0,:], df_test.iloc[0,:])
```

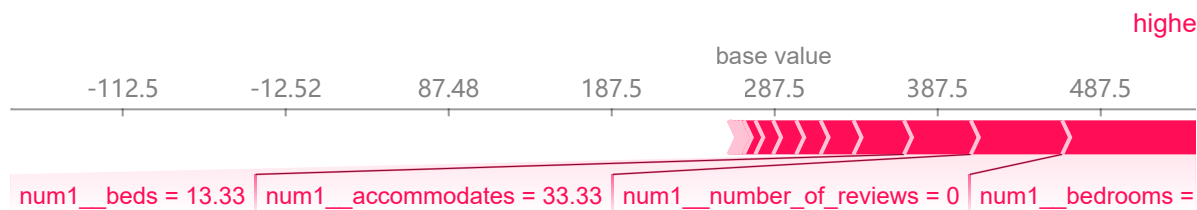
Out[89]:



In [90]:

```
1 #point
2 shap.force_plot(explainer.expected_value, shap_values[20,:], df_test.iloc[20,:])
```

Out[90]:



In []:

```
1
```