

## Assignment 1.0

In this project, we are going to implement a linear regression model to predict the living space from the dataset we have on apartment rental offers in Germany.

This code is written with the Spyder editor 4.2.

This linear model implementation is going to be without using packages with ready to go linear regression models, but it is going to be compared to one.

Linear regression attempts to model the relationship between two variables by fitting a linear equation to observed data. One variable is considered to be an explanatory variable, and the other is considered to be a dependent variable. For example, a modeler might want to relate the weights of individuals to their heights using a linear regression model.

Before attempting to fit a linear model to observed data, a modeler should first determine whether or not there is a relationship between the variables of interest. This does not necessarily imply that one variable causes the other (for example, higher SAT scores do not cause higher college grades), but that there is some significant association between the two variables.

A scatterplot can be a helpful tool in determining the strength of the relationship between two variables. If there appears to be no association between the proposed explanatory and dependent variables (i.e., the scatterplot does not

indicate any increasing or decreasing trends), then fitting a linear regression model to the data probably will not provide a useful model.

A valuable numerical measure of association between two variables is the correlation coefficient, which is a value between -1 and 1 indicating the strength of the association of the observed data for the two variables.

A linear regression line has an equation of the form  $Y = a + bX$ , where  $X$  is the explanatory variable and  $Y$  is the dependent variable. The slope of the line is  $b$ , and  $a$  is the intercept (the value of  $y$  when  $x = 0$ ).

Also a correlation matrix is implemented on the data for a better understanding of features correlation with the LivingSpace.

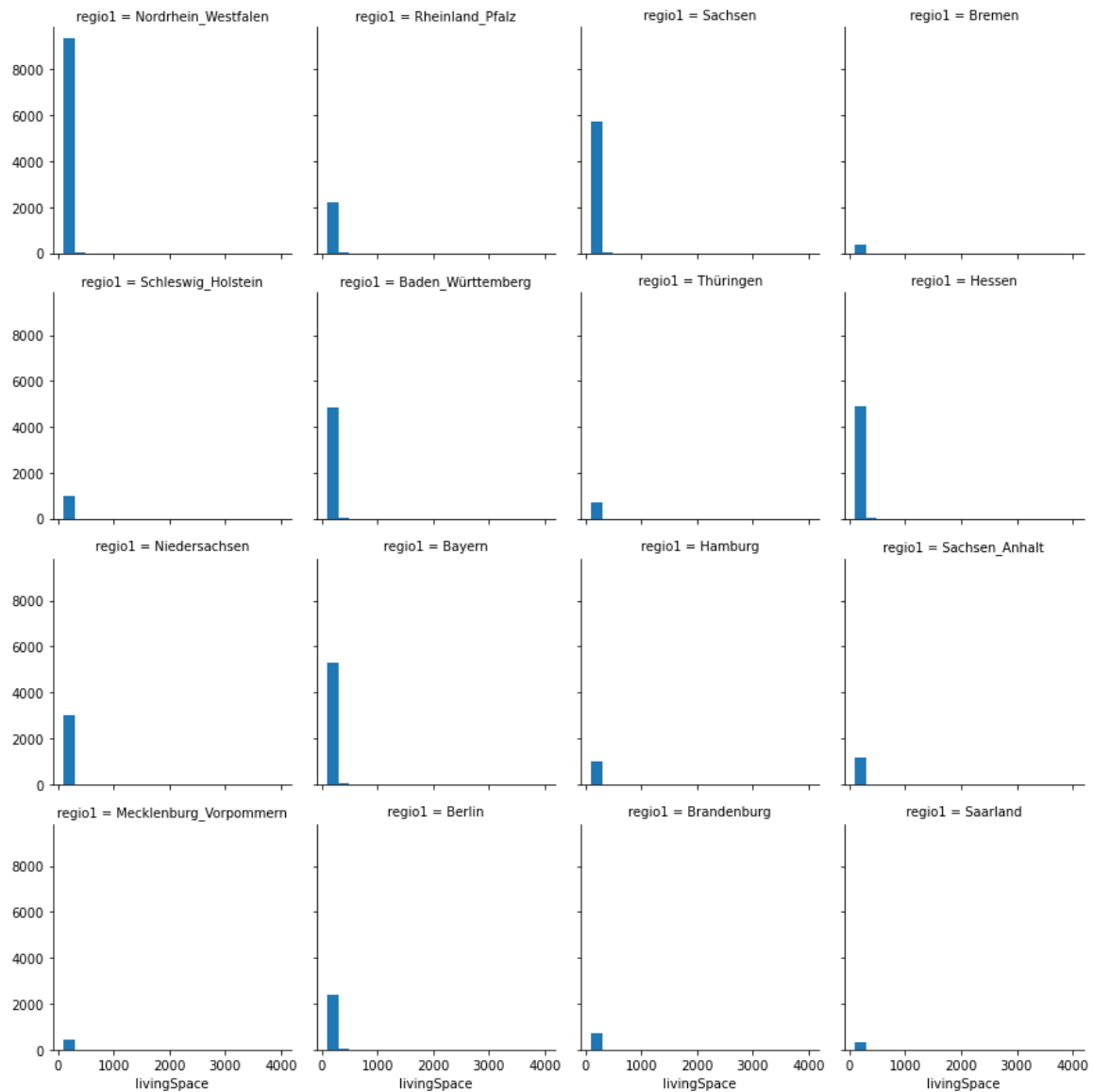
A correlation matrix is a table showing correlation coefficients between variables. Each cell in the table shows the correlation between two variables. A correlation matrix is used to summarize data, as an input into a more advanced analysis, and as a diagnostic for advanced analyses.

Before we begin fitting our model to the data, some preprocessing is necessary.

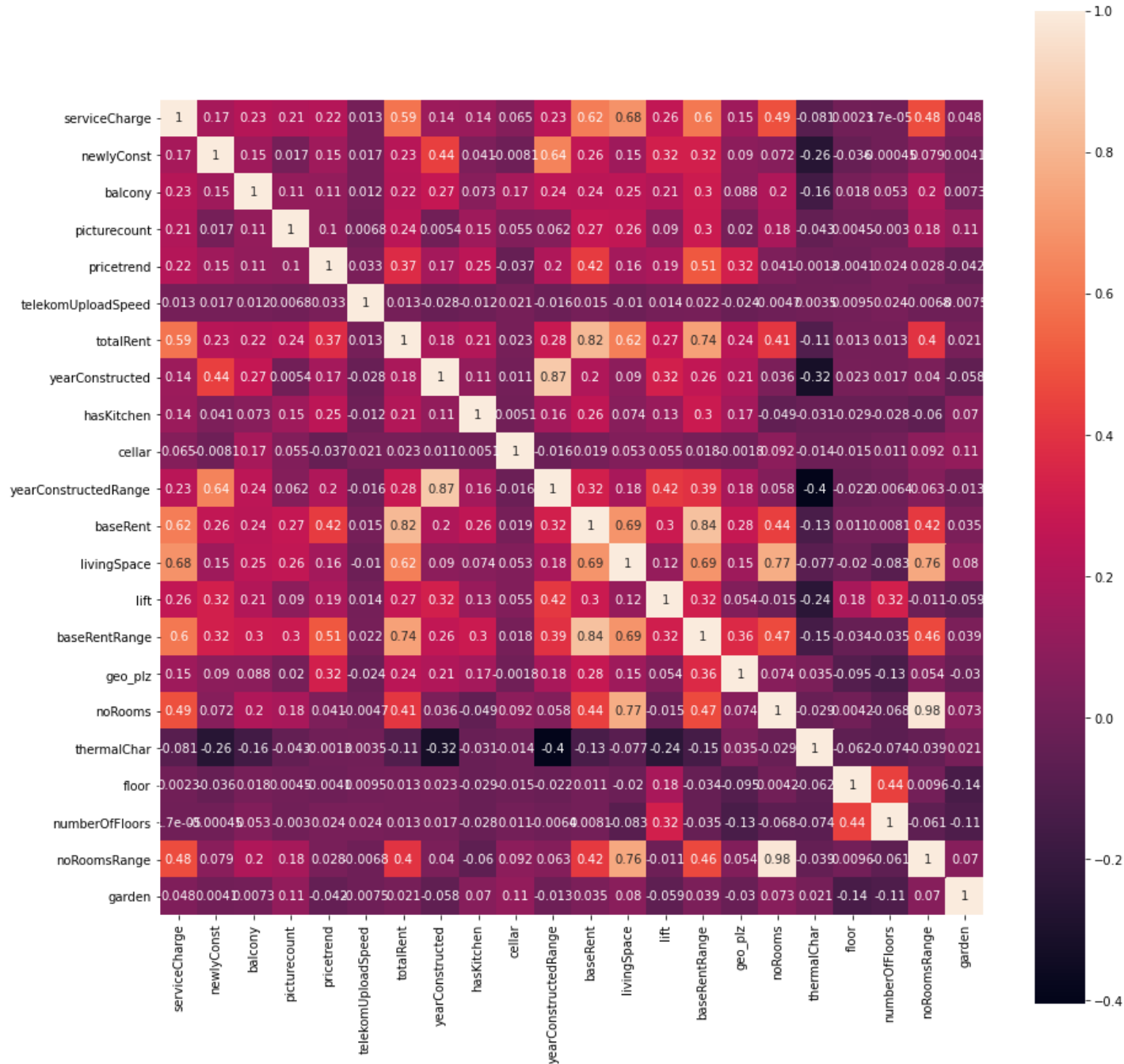
For this step:

First, I tried to clean and process the datasets given. For this purpose, I used the videos and codes that Mr. Sharifi has shared as a reference. For this part, I have removed some or replaced some of the null data, numeric and categorical. I have also removed and cleaned up the categorical features I believed couldn't be of help.

To get a better understanding of how `regio1` feature categories should be handled I used this plot for visualization and merged some of the categories into a single category called others:



From the correlation matrix, some of the features mostly correlated with “Living Space” are number of rooms and the range of number of rooms, base rent range, total range and service charge, which all have correlations over 0.60. Whereas the number of floors has a negative correlation with living space with a score of -0.08.



1) In this step, we should introduce five hypotheses and investigate the them with useful hypothesis test techniques.

Some of the well-known hypothesis test methods are:

Pearson's Correlation Coefficient

Spearman's Rank Correlation

Kendall's Rank Correlation

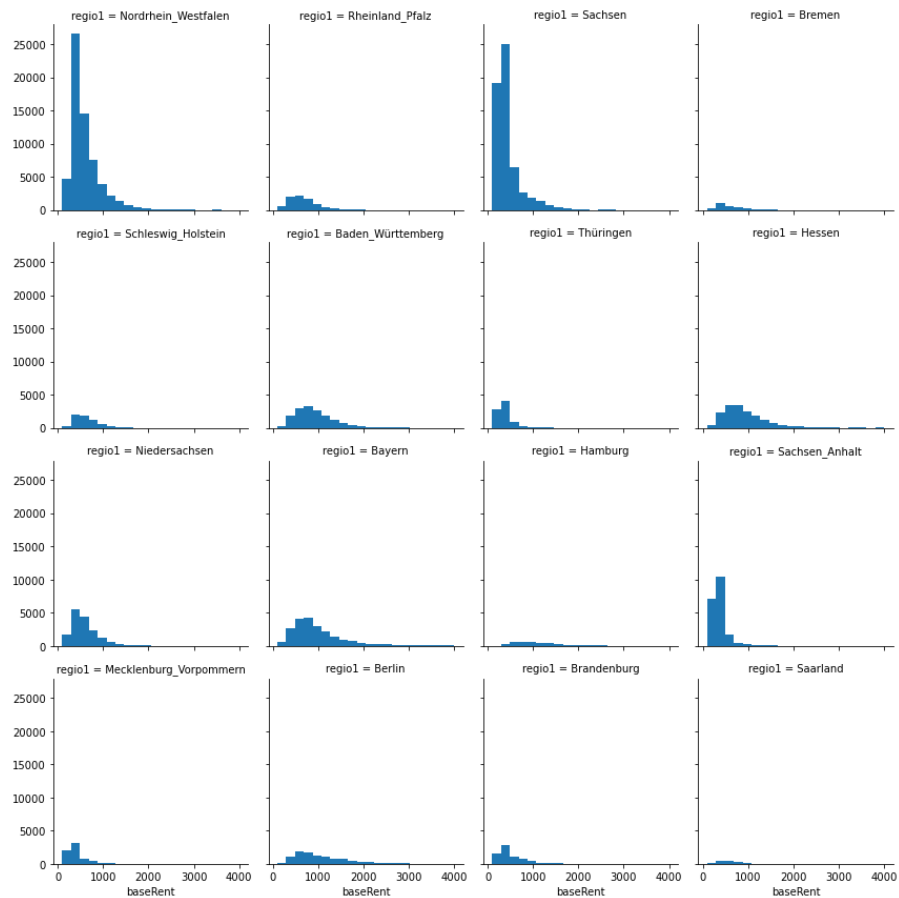
Chi-Squared Test

- a. First, let's investigate the relation of living space and the number of rooms.

The result for this hypothesis is `SpearmanrResult`  
(correlation=0.7815992067597107, pvalue=0.0)

- b. To investigate the relation between `regio1` and base rent, Hypothesis test Anova was run with the results "Oneway Anova totalRent ~ edit region1 F=5430.63, p-value=0.000000E+00"

To better understand this feature, let's see this grid face plot :



c. To investigate the relation of living space and no. of floors

The results show negative relation: SpearmanrResult (correlation=-0.10258377605772608, pvalue=0.0)

d. The relation of having kitchen with living space was tested with the results: PointbiseriialrResult (correlation=0.07421226391198485, pvalue=4.18269128783e-312)

e. Relation of having a kitchen and price of rent was tested and the results are: PointbiseriialrResult (correlation=0.21400248713621559, pvalue=0.0)

Some of the hypothesizes we investigated had relations and were correct.

2) In this step, I implement a linear regression model without using scikitlearn or statsmodels.

Before starting this step, replacing categorical data with dummy variables is necessary. I also split the target from features and divided the dataset into an 80 percent part train and a 20 percent part test after shuffling datapoints with this command

```
df = shuffle(df)
```

Then the regression model implementation was done with no packages other than NumPy. This part of the code can be seen in the #%%linear regr no package cell.

The difference of some predicted values vs the real test data values can be seen in the output below:

```
Out[207]:
```

	test	pred
0	73.61	84.908897
1	59.00	59.868465
2	156.00	101.394185
3	60.88	65.675451
4	117.13	124.502588

4) For this part, the linear regression model is implemented with scikitlearn package and the MSE for this model for all features is

The MSE for this model with all features was 247.4 with a learning rate of 0.001.

The most correlated feature with living space is number of rooms with a correlation score of 0.77, I only run a linear regr with this feature in use and the MSE for this model was 459.0.

The most negatively correlated feature with living space is number of floors with a correlation score of -0.08, I only run a linear regr with this feature in use and the MSE for this model was 1032.3.

As implying PCA to reduce features, this is what we got:

with a pov of 0.99 we got these results for the test and train sets:

Train:(206602, 36)

Test:(51651, 36)

with a pov of 0.95 we got these results for the test and train sets:

Train:(206602, 24)

Test:(51651, 24)

with a pov of 0.90 we got these results for the test and train sets:

Train:(206602, 17)

Test:(51651, 17)

with a pov of 0.70 we got these results for the test and train sets:

Train:(206602, 8)

Test:(51651, 8)

For a margin of 0.1, we test our model's accuracy and the outcome is 246.6 for MSE and 45.1 accuracy which is not a very good score.

## References:

1. Codes that the Teacher Assistants of this course have shared and Mr. Sharifi's videos.
2. [sharpsightlabs.com](https://sharpsightlabs.com)
3. [machinelearningmastery.com](https://machinelearningmastery.com)
4. [medium.com](https://medium.com), [towardsdatascience.com](https://towardsdatascience.com)
5. [scikit-learn.org](https://scikit-learn.org)
6. [kaggle.com](https://kaggle.com)



# Appendix

```
1. # -*- coding: utf-8 -*-
2. """
3. @author: Sayeh Sobhani
4. written in Spyder
5. """
6. ###Importing packages and dataset
7. import numpy as np
8. import pandas as pd
9. import seaborn as sns
10. import matplotlib.pyplot as plt
11. from scipy import stats
12. from sklearn.utils import shuffle
13. from sklearn.decomposition import PCA
14. from sklearn.linear_model import LinearRegression
15. from sklearn.metrics import mean_squared_error
16.
17. dfr = pd.read_csv(r'C:\Users\Asus\Desktop\Data Mining\01_assignment\1.0\immo_data.csv')
18.
19.
20. ###Understanding and cleaning dataset
21.
22. #to get a better understanding of how our data looks we look at the first 10 rows
23. dfr.shape
24. dfr.head(10)
25.
26. #lets see some info about this dataset and its features
27. dfr.info()
28.
29. dfr.describe()
30.
31. #the mean for each numeric feature?
32. #miangin = dfr._get_numeric_data().mean()
33. #print(miangin)
34.
35. #useless columns deletion '',
36. dfr = dfr.drop(columns=['description'])
37. dfr = dfr.drop(columns=['facilities'])
38. dfr = dfr.drop(columns=['scoutId'])
39. dfr = dfr.drop(columns=['geo_bln','date'])
40. dfr = dfr.drop(columns=['houseNumber','livingSpaceRange','firingTypes'])
41. dfr = dfr.drop(columns=['streetPlain','street','geo_krs','regio2','regio3'])
42. dfr = dfr.drop(dfr[dfr['livingSpace'] == 0.0].index)
43. dfr = dfr.drop(dfr[dfr['totalRent'] == 0.0].index)
44.
45. #how null data looks like
46. dfr.isna()
47.
48. dfr.columns
49.
50. #columns with null data
51. dfr.isna().sum()/len(dfr)
52.
53. #which columns hve >50 percent null
54. dfr.columns[((dfr.isna().sum()/len(dfr)) > 0.50)]
55.
56. #removing those columns with >50 null
57. dfr = dfr.drop(columns=dfr.columns[((dfr.isna().sum()/len(dfr)) > 0.50)])
```

```

58. dfr.columns
59. dfr.shape
60.
61. #filling the remaining data with means and max seen data
62.     #numeric
63. dfr._get_numeric_data().mean()
64. dfr.fillna(dfr._get_numeric_data().mean(),inplace = True)
65. dfr.isna().sum()
66.     #Categorical most seems
67. for cols in dfr.columns:
68.     if dfr[cols].dtype == 'object' or dfr[cols].dtype == 'bool':
69.         print('column : ',cols)
70.         print(dfr[cols].value_counts().head())
71.
72.     #Fill categorical nan with most seen data in each feature
73. for cols in dfr.columns:
74.     if dfr[cols].dtype == 'object' or dfr[cols].dtype == 'bool':
75.         print('cols : {} , value : {}'.format(cols ,
76.         dfr[cols].value_counts().head(1).index[0]))
77.         dfr[cols].fillna(dfr[cols].value_counts().head(1).index[0],inplace = True)
78. #to check if nan left
79. dfr.isna().sum()
80. dfr.shape
81. dfr.info()
82. "Some Visualization"
83. #region1 categorical by base rent amounts PLOT visualization
84. #dfr['region1'].value_counts()
85. #g = sns.FacetGrid(dfr, col='region1', col_wrap=4)
86. #g = g.map(plt.hist, 'livingSpace', bins=20, range=(100,4000))
87.
88. #outliers finding and removing
89. for cols in dfr.columns:
90.     if dfr[cols].dtype == 'int64' or dfr[cols].dtype == 'float64':
91.         upper_range = dfr[cols].mean() + 3 * dfr[cols].std()
92.         lower_range = dfr[cols].mean() - 3 * dfr[cols].std()
93.
94.         indexs = dfr[(dfr[cols] > upper_range) | (dfr[cols] < lower_range)].index
95.         dfr = dfr.drop(indexs)
96.
97. dfr.info()
98. dfr.shape
99.
100. #%%categorical feature inspection
101. #the number of categories in each feature
102. for cols in dfr.columns:
103.     if dfr[cols].dtype == 'object' or dfr[cols].dtype == 'bool':
104.         print('cols : {} , unique values : {}'.format(cols,dfr[cols].nunique()))
105.
106.
107. #the categories and their count in each feature
108. for cols in dfr.columns:
109.     if dfr[cols].dtype == 'object' or dfr[cols].dtype == 'bool':
110.         print('cols : {} ,\n {}'.format(cols,dfr[cols].value_counts()))
111.
112. #region1 edition based on plots from base rent amounts in line 75
113. dfr['region1'].value_counts()
114. def edit_region1(x):
115.     if x in ['Hamburg','Bremen','Saarland']:
116.         return 'other'
117.     else:
118.         return x
119. dfr['region1_'] = dfr['region1'].apply(edit_region1)
120. dfr = dfr.drop(columns = ['region1'])
121. dfr['region1_'].value_counts()*100 / len(dfr)

```

```

122.
123.     #visualization
124. plt_regio1_ = dfr['regio1_'].value_counts().plot(kind='bar')
125.
126. #heatingType
127. dfr['heatingType'].value_counts()
128. others = list(dfr['heatingType'].value_counts().tail(12).index)
129. def edit_heatingType(x):
130.     if x in others:
131.         return 'other'
132.     else:
133.         return x
134.
135. dfr['heatingType_'] = dfr['heatingType'].apply(edit_heatingType)
136. dfr = dfr.drop(columns = ['heatingType'])
137. dfr['heatingType_'].value_counts()*100 / len(dfr)
138.
139. #telekomTvOffer
140. dfr['telekomTvOffer'].value_counts()
141. others = list(dfr['telekomTvOffer'].value_counts().tail(2).index)
142. def edit_telekomTvOffer(x):
143.     if x in others:
144.         return 'other'
145.     else:
146.         return x
147.
148. dfr['telekomTvOffer_'] = dfr['telekomTvOffer'].apply(edit_telekomTvOffer)
149. dfr = dfr.drop(columns = ['telekomTvOffer'])
150. dfr['telekomTvOffer_'].value_counts()*100 / len(dfr)
151.
152. #typeOfFlat
153. dfr['typeOfFlat'].value_counts()
154. others = list(dfr['typeOfFlat'].value_counts().tail(6).index)
155. def edit_typeOfFlat(x):
156.     if x in others:
157.         return 'other'
158.     else:
159.         return x
160.
161. dfr['typeOfFlat_'] = dfr['typeOfFlat'].apply(edit_typeOfFlat)
162. dfr = dfr.drop(columns = ['typeOfFlat'])
163. dfr['typeOfFlat_'].value_counts()*100 / len(dfr)
164.
165. #condition
166. dfr['condition'].value_counts()
167. others = list(dfr['condition'].value_counts().tail(4).index)
168. def edit_condition(x):
169.     if x in others:
170.         return 'other'
171.     else:
172.         return x
173.
174. dfr['condition_'] = dfr['condition'].apply(edit_condition)
175. dfr = dfr.drop(columns = ['condition'])
176. dfr['condition_'].value_counts()*100 / len(dfr)
177.
178. %%normalization of numeric data
179. for cols in dfr.columns:
180.     if dfr[cols].dtype == 'int64' or dfr[cols].dtype == 'float64':
181.         if cols != 'livingSpace':
182.             dfr[cols] = ((dfr[cols] - dfr[cols].mean())/(dfr[cols].std()))
183.
184.
185. %%correlation matrix
186. corr = dfr.corr()

```

```

187. corr
188. f, ax = plt.subplots(figsize=(15, 15))
189.
190. sns.heatmap(corr, square = True ,annot = True)
191.
192. ###Hypothesis tests
193. dfr["region1_"]
194.     #relation of noRoomsRange&livingSpace
195. stats.stats.spearmanr(dfr['livingSpace'],dfr['noRoomsRange'])
196.
197.     #relation region1 and rent
198. fstat, pval = stats.f_oneway(*[dfr.baseRent[dfr.region1_ == s]
199. for s in dfr.region1_.unique()])
200. print("Oneway Anova totalRent ~ edit region1 F=%.2f, p-value=%E" % (fstat, pval))
201.
202.     #living space relation with no of floors
203. stats.stats.spearmanr(dfr['livingSpace'],dfr['numberOfFloors'])
204.
205.     #relation of having kitchen vs living space
206. stats.pointbiserialr(dfr['hasKitchen'],dfr["livingSpace"])
207.
208.     #relation between owning a kitchen and paying rent in total
209. stats.pointbiserialr(dfr['totalRent'],dfr["hasKitchen"])
210.
211. ###Dummy variables for categorical
212. columns = []
213. for cols in dfr.columns:
214.     if dfr[cols].dtype == 'object' or dfr[cols].dtype == 'bool':
215.         columns.append(cols)
216. columns
217.
218. dummies_feature = pd.get_dummies(dfr[columns])
219. dummies_feature.head()
220. dummies_feature.shape
221. dfr = pd.concat([dfr, dummies_feature], axis=1)
222. dfr.head()
223. dfr = dfr.drop(columns=columns)
224. dfr.head()
225. dfr.info()
226.
227. ###shuffling data and dividing test and train
228. dfr = shuffle(dfr)
229. y = dfr['livingSpace'].values
230. x = dfr.drop(columns = ['livingSpace']).values
231.
232. print(x.shape)
233. print(y.shape)
234.
235. #splitting test and train
236. train_size = int(0.8 * x.shape[0])
237. train_size
238.
239. x_train = x[:train_size]
240. y_train = y[:train_size]
241.
242. print(x_train.shape)
243. print(y_train.shape)
244.
245. x_test = x[train_size:]
246. y_test = y[train_size:]
247.
248. print(x_test.shape)
249. print(y_test.shape)
250.
251. ###linear regr no package

```

```

252.
253. #y^=w1*x1+w2*x2+...+w51*x51+bias is the line we're lookin for
254. #setting parameters "w" and "b" randomly
255. np.random.seed(42)
256. b = np.random.randn(1)
257. w = np.random.randn(51)
258. n = x_train.shape[0]
259.
260. #Sets learning rate
261. lr = 0.001
262.
263. #number of epochs
264. n_epochs = 500
265.
266. test_error = y_test - np.sum((w * x_test) + b , axis = 1)
267. test_mse = (test_error**2).mean()
268. print('init MSE : ',test_mse)
269.
270. for epoch in range(n_epochs):
271.     error = y_train - np.sum((w * x_train) + b , axis = 1)
272.
273.     if epoch % 100 == 0:
274.         print('epoch {} , MSE : {}'.format(epoch,(error**2).mean()))
275.
276.     w_grad = [0] * x_train.shape[1]
277.     b_grad = 0
278.
279.     for i in range(x_train.shape[1]):
280.         w_grad[i] = -1 * (x_train[:,i] * error).mean()
281.         w[i] = w[i] - (lr * w_grad[i])
282.
283.     b_grad = -1 * error.mean()
284.     b = b - (lr*b_grad)
285.
286.
287. test_error = y_test - np.sum((w * x_test) + b , axis = 1)
288. test_mse = (test_error**2).mean()
289. print('Final MSE : ',test_mse)
290.
291. #comparison of the predicted value and the value in dataset
292. y_pred = np.sum((w * x_test) + b , axis = 1)
293. temp = pd.DataFrame({'test':y_test,'pred':y_pred})
294. temp.head()
295.
296. temp['upper_range'] = temp['test'] * 1.2
297. temp['lower_range'] = temp['test'] * 0.8
298.
299. temp[(temp['upper_range'] >=temp['pred']) & (temp['pred'] >=
temp['lower_range'])].shape[0] * 100/temp.shape[0]
300.
301. #%%only noRooms feature linear regr no package
302.
303. X = dfr['noRooms'].values
304. X.shape
305. X_train = X[:train_size]
306. print(X_train.shape)
307. X_test = X[train_size:]
308. print(X_test.shape)
309.
310. np.random.seed(42)
311. b = np.random.randn(1)
312. w = np.random.randn(1)
313. n = X_train.shape[0]
314.
315. # Sets learning rate

```

```

316. lr = 0.1
317.
318.
319. n_epochs = 500
320. for epoch in range(n_epochs):
321.     error = y_train - ((w * X_train) + b)
322.
323.     if epoch % 100 == 0:
324.         print('epoch {} , MSE : {}'.format(epoch, (error**2).mean()))
325.
326.         # adoptive learning rate
327.         if epoch % 200 == 0:
328.             lr = lr * 0.1
329.
330.
331.         w_grad = 0
332.         b_grad = 0
333.
334.         w_grad = -1 * (X_train * error).mean()
335.         w = w - (lr * w_grad)
336.
337.         b_grad = -1 * error.mean()
338.         b = b - (lr*b_grad)
339.
340.     test_error = y_test - ((w * X_test) + b)
341.     test_mse = (test_error**2).mean()
342.     print('Final MSE : ', test_mse)
343.     print('Final learning rate : ', lr)
344.
345.
346.
347. #%%Linear regr with scikitlearn package all features
348. Linear = LinearRegression()
349. Linear.fit(x_train, y_train)
350. print("the coefficients are:")
351. print(Linear.coef_)
352. print("the intercept is:")
353. print(Linear.intercept_)
354.
355. y_pred = Linear.predict(x_test)
356. mean_squared_error(y_pred, y_test)
357.
358. temp = pd.DataFrame({'test': y_test, 'pred': y_pred})
359. temp.head()
360.
361. temp['upper_range'] = temp['test'] * 1.2
362. temp['lower_range'] = temp['test'] * 0.8
363.
364. temp[(temp['upper_range'] >= temp['pred']) & (temp['pred'] >=
temp['lower_range'])].shape[0] * 100/temp.shape[0]
365.
366. #%%Linear regr package floor feature
367.
368. X = dfr['numberOfFloors'].values
369. X.shape
370. X_train = X[:train_size]
371. print(X_train.shape)
372. X_test = X[train_size:]
373. print(X_test.shape)
374.
375. Linear = LinearRegression()
376. Linear.fit(x_train, y_train)
377. print("the coefficients are:")
378. print(Linear.coef_)
379. print("the intercept is:")

```

```

380. print(Linear.intercept_)
381.
382. y_pred = Linear.predict(x_test)
383. mean_squared_error(y_pred,y_test)
384.
385. temp = pd.DataFrame({'test':y_test,'pred':y_pred})
386. temp.head()
387.
388. temp['upper_range'] = temp['test'] * 1.2
389. temp['lower_range'] = temp['test'] * 0.8
390.
391. temp[(temp['upper_range'] >=temp['pred']) & (temp['pred'] >=
temp['lower_range'])].shape[0] * 100/temp.shape[0]
392.
393. ###number of floors feature linear regr no package
394.
395. X = dfr['numberOfFloors'].values
396. X.shape
397. X_train = X[:train_size]
398. print(X_train.shape)
399. X_test = X[train_size:]
400. print(X_test.shape)
401.
402. np.random.seed(42)
403. b = np.random.randn(1)
404. w = np.random.randn(1)
405. n = X_train.shape[0]
406.
407. # Sets learning rate
408. lr = 0.1
409.
410.
411. n_epochs = 500
412. for epoch in range(n_epochs):
413.     error = y_train - ((w * X_train) + b)
414.
415.     if epoch % 100 == 0:
416.         print('epoch {} , MSE : {}'.format(epoch,(error**2).mean()))
417.
418.     # adoptive learning rate
419.     if epoch % 200 == 0:
420.         lr = lr * 0.1
421.
422.
423.     w_grad = 0
424.     b_grad = 0
425.
426.     w_grad = -1 * (X_train * error).mean()
427.     w = w - (lr * w_grad)
428.
429.     b_grad = -1 * error.mean()
430.     b = b - (lr*b_grad)
431.
432. test_error = y_test - ((w * X_test) + b)
433. test_mse = (test_error**2).mean()
434. print('Final MSE : ',test_mse)
435. print('Final learning rate : ',lr)
436.
437. ###PCA
438.
439. #pov 0.99
440. pca = PCA(0.99)
441. x_pca = pca.fit_transform(x)
442. x_pca.shape
443. x_pca_train = x_pca[:train_size]

```

```

444. print(x_pca_train.shape)
445. x_pca_test = x_pca[train_size:]
446. print(x_pca_test.shape)
447.
448. #pov 0.95
449. pca = PCA(0.95)
450. x_pca = pca.fit_transform(x)
451. x_pca.shape
452. x_pca_train = x_pca[:train_size]
453. print(x_pca_train.shape)
454. x_pca_test = x_pca[train_size:]
455. print(x_pca_test.shape)
456.
457. #pov 0.90
458. pca = PCA(0.90)
459. x_pca = pca.fit_transform(x)
460. x_pca.shape
461. x_pca_train = x_pca[:train_size]
462. print(x_pca_train.shape)
463. x_pca_test = x_pca[train_size:]
464. print(x_pca_test.shape)
465.
466. #pov 0.70
467. pca = PCA(0.70)
468. x_pca = pca.fit_transform(x)
469. x_pca.shape
470. x_pca_train = x_pca[:train_size]
471. print(x_pca_train.shape)
472. x_pca_test = x_pca[train_size:]
473. print(x_pca_test.shape)
474.
475.
476. ## linear regr no package different margin
477. #margin or 0.1
478. np.random.seed(42)
479. b = np.random.randn(1)
480. w = np.random.randn(51)
481. n = x_train.shape[0]
482.
483. #Sets learning rate
484. lr = 0.001
485.
486. #number of epochs
487. n_epochs = 500
488.
489. test_error = y_test - np.sum((w * x_test) + b , axis = 1)
490. test_mse = (test_error**2).mean()
491. print('init MSE : ',test_mse)
492.
493. for epoch in range(n_epochs):
494.     error = y_train - np.sum((w * x_train) + b , axis = 1)
495.
496.     if epoch % 100 == 0:
497.         print('epoch {} , MSE : {}'.format(epoch,(error**2).mean()))
498.
499.     w_grad = [0] * x_train.shape[1]
500.     b_grad = 0
501.
502.     for i in range(x_train.shape[1]):
503.         w_grad[i] = -1 * (x_train[:,i] * error).mean()
504.         w[i] = w[i] - (lr * w_grad[i])
505.
506.     b_grad = -1 * error.mean()
507.     b = b - (lr*b_grad)
508.

```



```
509.
510. test_error = y_test - np.sum((w * x_test) + b , axis = 1)
511. test_mse = (test_error**2).mean()
512. print('Final MSE : ',test_mse)
513.
514. #comparison of the predicted value and the value in dataset
515. y_pred = np.sum((w * x_test) + b , axis = 1)
516. temp = pd.DataFrame({'test':y_test,'pred':y_pred})
517. temp.head()
518.
519. temp['upper_range'] = temp['test'] * 1.1
520. temp['lower_range'] = temp['test'] * 0.9
521.
522. temp[(temp['upper_range'] >=temp['pred']) & (temp['pred'] >=
    temp['lower_range'])].shape[0] * 100/temp.shape[0]
523.
```