

Introduction to Dataset and Algorithm

Q1: Sets of Amenities w/ High Review Scores

Across Entire Dataset vs. Within Binned Dataset

Apriori Algorithm, FP-Growth Algorithm, SparkML FP-Growth Algorithm

Q2: Sets of Amenities w/ High Review Scores Based on Location

Across Entire Dataset vs. Within Binned Dataset

Apriori Algorithm, FP-Growth Algorithm

Recommendation

2

2.1

2.1.1

3

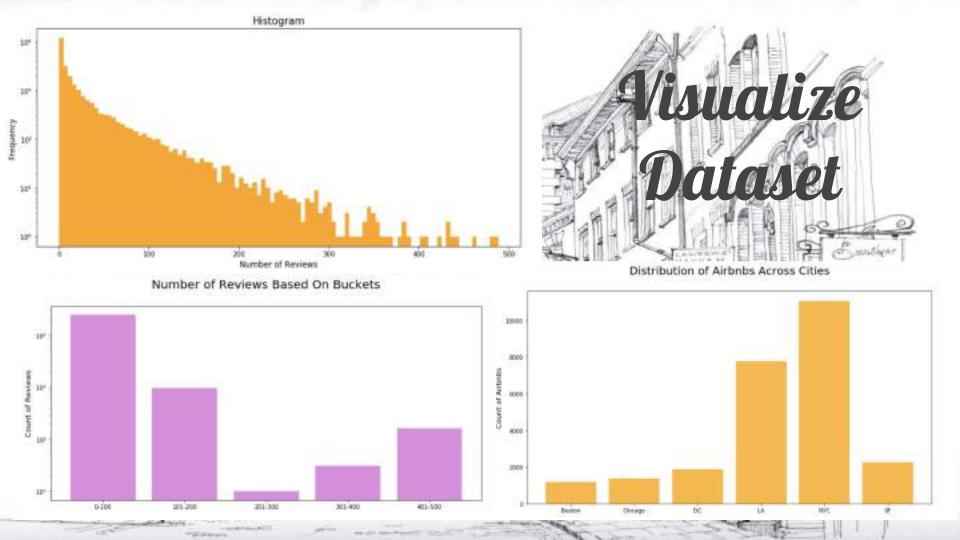
3.1

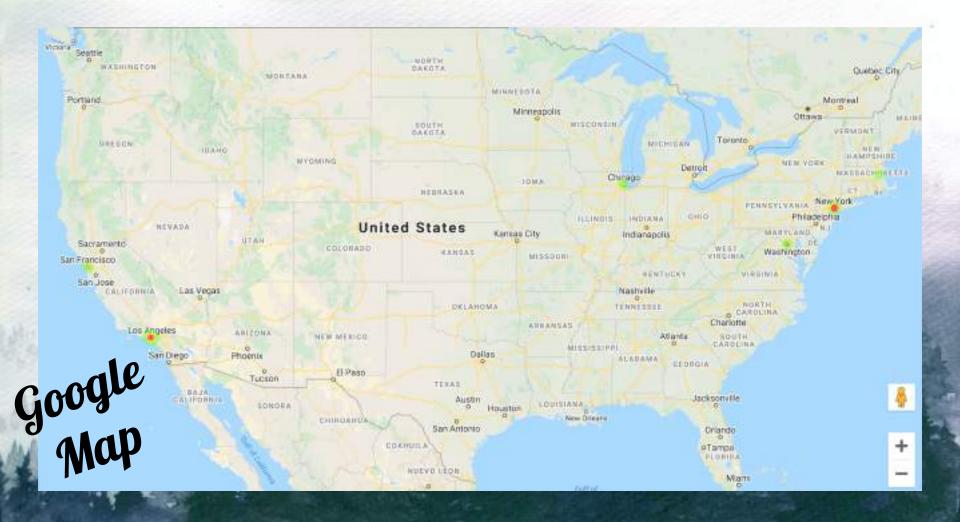
3.1.1

4

Introduction to Dataset

- 25,458 rows of data, 213 columns
- All columns are of type double and int
- There are no null values in the dataset
- Columns are mostly in binary form
- Example of columns are:
 - Property_Type_Condominium (0 for no, 1 for yes)
 - City_SF (0 for no, 1 for yes)
 - Firm mattress (0 for no, 1 for yes)
 - Washer (0 for no, 1 for yes)





accommod	dates bath	rooms	bedroo	ms I	beds cl	eaning_	fee hos	t_resp	ponse_rate	ld	latit	ude log	_price	longitude	numbe	er_of_r	reviews
	2	1.0		1.0	1.0		1		100	3895911	34.028	372	0.0 -	118.494449			6
	3	1.0		1.0	1.0		1		100	9710289	40.720	380	0.0	-73.942329			2
	1	1.0		1.0	1.0		1		100	9051635	37.785	434	0.0 -	122.470284			2
	1	1.0		0.0	1.0		1		100	708374	33.976	026	0.0 -	118.463471			7
	2	1.0		1.0	1,0		1		0	626296	40.735	573	0.0	-74.005996			0
Family/kid friendly	Fir extinguishe			Firm atress	Firm mattres	airl	Fixed grab bars for shower & toilet	Flat	Flat smooth pathway to front door	Free parking on premises	Free parking on street	Game console	Garde o backyan	shower	Ground floor access	Hair dryer	Hand or paper towel
0		1	0	0		0 1	0	0	0	1	0	0		0 0	0	1	0
0	9	0	0	0	(6)	0 0	0	0	0	0	0	0		0 0	0	1	c

When is Frequent Pattern Mining Algorithm Used?

- For companies to increase profits by understanding purchasing patterns
- In medical diagnosis to understand symptoms that coexist
- By studying behavior of users while navigating through websites

These information helps companies make decisions like, the kind of group discounts can be offered to customers, how to orient products at the stores or online, how to effectively advertise to customers, and what products/services/medications/websites to recommend to their customers.

Introduction to Association Rules

Association rules are most often used in businesses to help identify what items are commonly bought together. It determines how items are associated to each other.

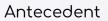
Disadvantage: To identify the associations between lists of itemsets, the support values for each possible configuration of items need to be calculated, and then shortlist itemsets that meet minimum support threshold.

$$Support(\{X\} \rightarrow \{Y\}) = \frac{Transactions\ containing\ both\ X\ and\ Y}{Total\ number\ of\ transactions}$$

$$Confidence(\{X\} \rightarrow \{Y\}) = \frac{Transactions\ containing\ both\ X\ and\ Y}{Transactions\ containing\ X}$$

$$Lift(\{X\} \rightarrow \{Y\}) = \frac{(Transactions\ containing\ both\ X\ and\ Y)/(Transactions\ containing\ X)}{Fraction\ of\ transactions\ containing\ Y}$$

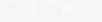
Association Rules











Consequent



Support Measur e













Confidence Measure







Lift Measur

е







$$\equiv \frac{7}{6}$$









Introduction to Apriori Algorithm

Compared to Association Rules, the number of itemsets the Apriori Algorithm examines is reduced.

Apriori algorithm uses frequent itemsets to generate association rules. Frequent itemset is an itemset whose support value is greater than a threshold value.

Disadvantages: Computationally expensive for large datasets (or when support value is low)

	ID	"Bed linens"	"Buzzer/wireless intercom"	"Cable TV"	"Carbon monoxide detector"	"Hair dryer"	
	81	0			0	1	
	82	0	0	0	0	1	
L	83	0	1	1	1	1	
2	84	0	1	1	1	0	ø
ą	85	0		0	1	0	ě
3	86	0			1	0	k
ŝ	87	0	0	0	1	1	8
ĕ	88	0	1	1	1	1	E
e	89	1	0	1	0	1	
	810	0		0	0	1	9

B3, B8

k = 1, min	nimum support = 2
Item	IDSet
Buzzer	B3, B4, B8
Cable TV	B3, B4, B8, B9
CO Detector	83, 84, 85, 86, 87, 88
Hair Dryer	B1, B2, B3, B7, B8, B9, B10

k = 3, minimum suppo	ort = 2
Item	IDSet
Buzzer, Cable TV, CO Detector	B3, B4, B8
Buzzer, Cable TV, Hair Dryer	B3, B8
k = 4, minimum suppo	ort = 2
Item	IDSet

Buzzer, Cable TV, CO

Detector, Hair Dryer

k = 2, minimum si	upport = 2
Item	IDSet
Buzzer, Cable TV	B3, B4, B8
Buzzer, CO Detector	B3, B4, B8
Buzzer, Hair Dryer	B3, B8

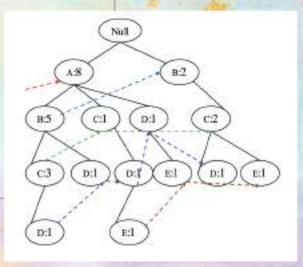
Items Offered	Recommended Items
Buzzer	Cable TV
Buzzer	CO Detector
Buzzer	Hair Dryer
Cable TV	CO Detector

Introduction to FP-Growth Algorithm

FP growth technique uses pattern fragment growth to mine the frequent patterns from large database.

An extended prefix tree structure is used for storing crucial and compressed information about frequent patterns. FP growth discovers the frequent itemsets without candidate itemset generation.

Advantages: Reduced computation time and cost

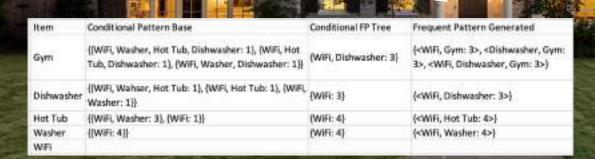


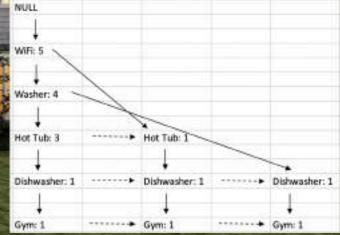


ID	Items
B1	Dryer, Washer, WiFi, Hot Tub, Patio
B2	WiFi, Heater, Hot Tub, Dishwasher, Gym
83	WiFi, Washer, Hot Tub
B4	Dishwasher, Washer, WiFi, Gym
85	Washer, WiFi, Hot Tub, Dishwasher, Gym

	Item	Frequence
į	Dryer	1
	Washer	4
	WiFi	5
	Hot Tub	4
	Patio	1
j	Heater	1
	Dishwasher	3
	Gym	3

		W10.486
Air.		L = { WiFi, Washer, Hot Tub, Dishwasher, Gym }
S. Salak	ID	Ordered-Item Set
	B1	WiFi, Washer, Hot Tub
W -	B2	WiFi, Hot Tub, Dishwasher, Gym
	B3	WiFi, Washer, Hot Tub
3	B4	WiFi, Washer, Dishwasher, Gym
4	B5	WiFi, Washer, Hot Tub, Dishwasher, Gym
	11.5	





Find Amenities w/ High Review Scores



- 1. Divided into two parts:
 - a. Across entire dataset
 - b. Within selected bins
- 2. Using different algorithms within FPM

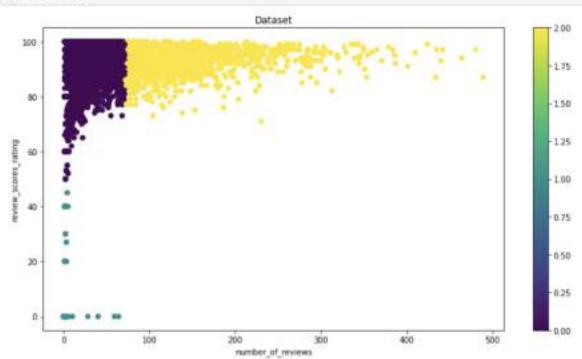


```
cols = [c for c in df_amenities.columns if c.lower()[:14] != 'property type ']
  cols = [c for c in df amenities.columns if c.lower()[:19] != 'translation missing']
  df_amenities = df_amenities cols
  df amenities = df amenities.drop(df amenities.loc[:, 'accommodates':'instant bookable t'l.columns, axis = 1)
  df_amenities = df_amenities.drop(df_amenities.loc[:, 'first_review_Year':'numReviews_buckets2'].columns,
                                     axis = 1)
  df_amenities = df_amenities.drop(['Internet', 'Shampoo', 'Hangers', 'Bath towel', 'Bed linens',
                                    'Body soap', 'Firm matress', 'Grab-rails for shower and toilet',
                                   'Hand soap', 'Washer / Dryer', 'Wide clearance to shower and toilet'], axis=1]
  df amenities.head(4)
                                                                                       Bathtub
    smooth
                            Accessible-
                                                                                                       Buzzer/wireless
           check-
      door
                                                                                         chair
        0
               0
                         D
                                   0
                                              0
                                                              0
                                                                      0
                                                                                    0
                                                                                                     0
                                                                                                                  0
        0
               0
                         D
                                   0
 2
                                   0
                                        len(df amenities (df amenities ('Wireless Internet') == 1)])
                         Ü
                                   0
 3
               0
                                      11793
  len(df amenities.columns)
                                        len(df_amenities[(df_amenities['Internet'] == 1)])
117
                                      7828
                                        len(df_amenities['df_amenities['wireless Internet'] == 1) & (df_amenities['Internet'] == 1)])
                                      7798
```

```
df_bin2.count()
25458
 df_bin2.filter(df_bin2['numReviews_buckets2'] == 0).count()
11887
 df_bin2.filter(df_bin2['review_scores_rating'] > 0).select('review_scores_rating').describe().show()
 summary | review_scores_rating
                        19741
   count
            94.011549566891251
    mean !
            7.978032677622043
  stddevl
     minl
                         20.01
                        100.0
 df.createOrReplaceTempView('df')
 spark.sql('select percentile(cast(review_scores_rating as BIGINT), 0.5) from df').collect()
[Row(percentile(CAST(review_scores_rating AS BIGINT), CAST(0.5 AS DOUBLE), (1)=94.0)]
 df amenities = df amenities.filter(df amenities['review scores rating'] > 94)
 df amenities.count()
12053
```

```
x1 = np.array(pd_transformed['number_of_reviews'])
x2 = np.array(pd_transformed['review_scores_rating'])
f, ax = plt.subplots(figsize=(13,7))
plt.title('Dataset')
plt.scatter(x1, x2, c=c)
plt.colorbar()
plt.xlabel('number_of_reviews')
plt.ylabel('review_scores_rating')
plt.show();
```

K-means?



Apriori Algorithm

from mlxtend.frequent_patterns import apriori, association_rules

frq_items = apriori(df_amenities, min_support=0.7, use_colnames = True)

frq_items['length'] = frq_items['itemsets'].apply(lambda x: len(x))
frq_items = frq_items.sort_values(by='support', ascending = False)
frq_items[(frq_items['length'] == 1)].head(20)

2 0.978429 (Wireless Internet) 4 0.933544 (Heating) 5 0.917116 (Kitchen) 3 0.899693 (Essentials)	1
5 0.917116 (Kitchen)	
A MANAGER AND	1
3 0.899693 (Essentials)	1
	1
1 0.873475 (Smoke detector)	1
0 0.753422 (Air conditioning)	1
6 0.743217 (TV)	1

frq_items[(frq_items['length'] == 2)].head(38)

length	itemsets	support	
2	(Wireless Internet, Heating)	0.922675	15
2	(Kitchen, Wireless Internet)	0.906164	16
2	(Wireless Internet, Essentials)	0.889073	14
2	(Kitchen, Heating)	0.867502	20
2	(Smoke detector, Wireless Internet)	0.862607	10
2	(Heating, Essentials)	0.852983	18
2	(Kitchen, Essentials)	0.833900	19
2	(Smoke detector, Heating)	0.830913	12
2	(Smoke detector, Essentials)	0.814403	11
2	(Kitchen, Smoke detector)	0.811250	13
2	(Wireless Internet, Air conditioning)	0.745209	7
2	(TV, Wireless Internet)	0.737161	17
2	(Heating, Air conditioning)	0.727122	В
2	(TV, Heating)	0.711441	21
2	(Kitchen, Air conditioning)	0.704638	9

Apriori Algorithm

frq_items[(frq_items['length'] == 3)].head(30)

	support	itemsets	length
31	0.858790	(Kitchen, Wireless Internet, Heating)	3
29	0.844769	(Wireless Internet, Essentials, Heating)	3
30	0.825438	(Kitchen, Wireless Internet, Essentials)	3
24	0.822368	(Smoke detector, Wireless Internet, Heating)	3
23	0.806355	(Smoke detector, Wireless Internet, Essentials)	3
25	0.802788	(Kitchen, Smoke detector, Wireless Internet)	3
33	0.793578	(Kitchen, Heating, Essentials)	3
26	0.779142	(Smoke detector, Heating, Essentials)	3
28	0.773915	(Kitchen, Heating, Smoke detector)	3
27	0.757654	(Kitchen, Essentials, Smoke detector)	3
22	0.719489	(Wireless Internet, Heating, Air conditioning)	3
32	0.706214	(TV, Wireless Internet, Heating)	3

frq_items[(frq_items['length'] == 4)].head(30)

	support	itemsets	length
38	0.786858	(Kitchen, Wireless Internet, Essentials, Heating)	4
34	0.772339	(Smoke detector, Wireless Internet, Essentials	4
36	0.766780	(Kitchen, Smoke detector, Wireless Internet, H	4
35	0.750850	(Kitchen, Smoke detector, Wireless Internet, E	4
37	0.725877	(Kitchen, Heating, Essentials, Smoke detector)	4
f	rq_item:	s[(frq_items['length'] == 5)].hea	ad (30)
	support	itemsets	length
39	0.720153	(Kitchen, Smoke detector, Wireless Internet, E	5

Association Rules

rules = association_rules(frq_items, metric ="lift", min_threshold = 1.0)

rules = rules.sort_values(['confidence', 'lift'], ascending =[False, False])

rules.head()

	antecedents	consequents	antecedent support	consequent support	support	confidence	Hft	leverage	conviction
195	(TV, Heating)	(Wireless Internet)	0.711441	0.978429	0.706214	0.992653	1.014538	0.010120	2.936106
158	(Kitchen, Heating, Essentials, Smoke detector)	(Wireless Internet)	0.725877	0.978429	0.720153	0.992113	1.013986	0.009933	2.735190
138	(IV)	(Wireless Internet)	0.743217	0.976429	0.737161	0.991851	1.013718	0.009976	2.647076
64	(Kitchen, Heating, Essentials)	(Wineless Internet)	0.793578	0.978429	0.786858	0.991532	1.013392	0.010398	2.547289
90	(Smoke detector, Heating, Essentials)	(Wireless Internet)	0.779142	0.978429	0.772339	0.991268	1.013123	0.010004	2.470451













	antecedents	consequents	antecedent	consequent	support	confidence	lift	leverage	conviction	antecedent len
rules[(rules['an & (rules['		len'] == 4) 1.0)].sort_v					10)			

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	leverage	conviction	antecedent_len
158	(Kitchen, Heating, Essentials, Smoke detector)	(Wireless Internet)	0.725877	0.978429	0.720153	0.992113	1.013986	0.009933	2.735190	4
156	(Kitchen, Wireless Internet, Essentials, Smoke	(Heating)	0.750850	0.933544	0.720153	0.959116	1.027393	0.019201	1.625490	4
157	(Kitchen, Wireless Internet, Heating, Smoke de	(Essentials)	0.766780	0.899693	0.720153	0.939191	1.043901	0.030286	1.649532	4
160	(Smoke detector, Wireless Internet, Essentials	(Kitchen)	0.772339	0.917116	0.720153	0.932431	1.016699	0.011828	1.226656	4
159	(Kitchen, Wireless Internet, Essentials, Heating)	(Smoke detector)	0.786858	0.873475	0.720153	0.915226	1.047798	0.032851	1.492486	4
	Essentials)	Internet)	******	MIN. M. 184	****	****		w.w.uww.	*:States	₹4
104	(Kitchen, Wireless Internet, Heating)	(Smoke detector)	0.858790	0.873475	0.766780	0.892861	1.022193	0.016648	1.180934	3
103	(Kitchen, Heating, Smoke detector)	(Wireless Internet)	0.773915	0.978429	0.766780	0.990780	1.012624	0.009559	2.339744	3

from mlxtend.frequent_patterns import fpgrowth

fp_frqitems = fpgrowth(df_amenities, min_support=0.7, use_colnames = True)

fp_frqitems['length'] = fp_frqitems['itemsets'].apply(lambda x: len(x))
fp_frqitems[(fp_frqitems['length'] == 1)].sort_values('support', ascending = False).head(20)

	support	Itemsets	length
5	0.978429	(Wireless Internet)	3
0	0.933544	(Heating)	1
1	0.917116	(Kitchen)	1
2	0.899693	(Essentials)	1
3	0.873475	(Smoke detector)	1
6	0.753422	(Air conditioning)	1
4	0.743217	(TV)	1



fp_frqitems[(fp_frqitems['length'] == 2)].head(20)

length	itemsets	support	
2	(Wireless Internet, Heating)	0.922675	7
2	(Kitchen, Heating)	0.867502	8
2	(Kitchen, Wireless Internet)	0.906164	9
2	(Heating, Essentials)	0.852983	11
2	(Kitchen, Essentials)	0.833900	12
2	(Wireless Internet, Essentials)	0.889073	13
2	(Smoke detector, Heating)	0.830913	18
2	(Smoke detector, Essentials)	0.814403	19
2	(Kitchen, Smoke detector)	0.811250	20
2	(Smoke detector, Wireless Internet)	0.862607	21
2	(TV, Heating)	0.711441	33
2	(TV, Wireless Internet)	0.737161	34
2	(Wireless Internet, Air conditioning)	0.745209	36
2	(Heating, Air conditioning)	0.727122	37
2	(Kitchen, Air conditioning)	0.704638	38



```
fp_frqitems[(fp_frqitems['length'] == 3)].\
sort_values('support', ascending = False).head(20)
```

length	itemsets	support	
3	(Kitchen, Wireless Internet, Heating)	0.858790	10
3	(Wireless Internet, Essentials, Heating)	0.844769	14
3	(Kitchen, Wireless Internet, Essentials)	0.825438	16
3	(Smoke detector, Wireless Internet, Heating)	0.822368	22
3	(Smoke detector, Wireless Internet, Essentials)	0.806355	24
3	(Kitchen, Wireless Internet, Smoke detector)	0.802788	28
3	(Kitchen, Heating, Essentials)	0.793578	15
3	(Smoke detector, Heating, Essentials)	0.779142	23
3	(Kitchen, Heating, Smoke detector)	0.773915	26
3	(Kitchen, Essentials, Smoke detector)	0.757654	27
3	(Wireless Internet, Heating, Air conditioning)	0.719489	39
3	(TV, Wireless Internet, Heating)	0.706214	35

```
fp_frqitems[(fp_frqitems['length'] == 4)].\
sort_values('support', ascending = False).head(20)
```

	support	itemsets	length
17	0.786858	(Kitchen, Wireless Internet, Essentials, Heating)	4
25	0.772339	(Smoke detector, Wireless Internet, Essentials	4
29	0.766780	(Kitchen, Wireless Internet, Heating, Smoke de	4
31	0.750850	(Kitchen, Wireless Internet, Essentials, Smoke	4
30	0.725877	(Kitchen, Heating, Essentials, Smoke detector)	4
		ems[(fp_frqitems['length'] == 5)] ues('support', ascending = False)	

	support	itemsets	length
32	0.720153	(Kitchen, Smoke detector, Wireless Internet, E	5

```
spark_test = spark.read.format('com.databricks.spark.csv').options(header='true',
                                inferschema='true').load('df amenities.csv')
 spark test.select('amenities list').show(3)
       amenities list!
                                                                                         Spark MI
 , Cable TV, Eleva...
 , Carbon monoxide...
 , Air conditionin...
only showing top 3 rows
 spark test = spark test.toPandas()
  spark test['amenities list'] = spark test['amenities list'].apply(lambda x; str(x))\
  .apply(lambda x: x.split(',')).apply(lambda x: list(x))
  spark test['amenities list'] = spark test['amenities list'].\
 apply(lambda row: [val for val in row if val != "'])
  spark_test = spark.createDataFrame(spark_test)
 spark test['amenities list'].head(5)
    [ Cable TV, Elevator in building, Fire extinguisher, First aid kit, Free parking on premises, Hair dr
yer, Laptop friendly workspace, Safety...
    [ Carbon monoxide detector, Family/kid friendly, First aid kit, Hair dryer, Laptop friendly workspace,
Smoke detector, Wireless Internet, E...
    [ Air conditioning, Carbon monoxide detector, Family/kid friendly, Fire extinguisher, Hair dryer, Hot
tub, Laptop friendly workspace, Smoke...
    [ Air conditioning, Carbon monoxide detector, Family/kid friendly, Hair dryer, Pets allowed, Private
entrance, Smoke detector, Wireless Int...
    [ Air conditioning, Carbon monoxide detector, Fire extinguisher, First aid kit, Free parking on premis
es, Hair dryer, Laptop friendly worksp...
Name: amenities_list, dtype: object
```

Spark ML

spark_frq[(spark_frq['length'] - 1)].sort_values('freq', ascending - False).head(20)

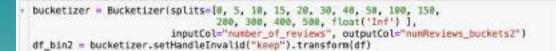
dt.	amenit	tio	e he	sad(2)																items	freq	length
	and the same	-10		1801 < 1															16	[Wireless Internet]	11793	1
														_					22	[Heating]	11252	1
Gym	Heating	l li	ron	Keypad	Kitchen	Lockbox	Microwave	Other	Over	Pool	Refrigerator	Bmartlock	Blove	TV	Washer	Waterfront	amenties	lest	24	[Kitchen]	11054	1
			121														, Cable		32	[Essentials]	10844	1
0	- 3		-1	ū	- 1	٥	0	0	0	0	6	0	à	1	- 3	0	Elevak building, eating	Fire	0	[Smoke detector]	10528	1
																	, Car mono	xide	17	[Air conditioning]	9081	1
0			0	0	23	۰	0	0		0	0		0			0	Family frien	/kid	28	[TV]	8958	- 1

antecedent	consequent	confidence	llift
[TV, Heating]	[Wireless Internet] 0.9926530612244898	1.014538060454403
[Smoke detector, Essentials, Kitchen, H	Heating] [Wireless Internet	1 0.9921133843867871	11.0139864853738612
[VT]	[Wireless Internet	1 0.9918508595668676	1.0137181726752695
[Essentials, Kitchen, Heating]	[Wireless Internet] [0.9915316257187663	11.0133919006858552
[Smoke detector, Essentials, Heating]	[[Wireless Internet] 0.9912682355446705	1.0131227035546437
[Smoke detector, Essentials, Kitchen]	Wireless Internet] 0.9910205869469996	1.0128695950540307
[Smoke detector, Kitchen, Heating]	[Wireless Internet] 0.9907804459691252	11.0126241596935357
[Essentials, Heating]	[Wireless Internet] [0.9903705865188211	11.0122052640813493
[Smoke detector, Essentials]	[[Wireless Internet	1 0.990118174409128	1.011947287047674
[Kitchen, Heating]			11.0117834983872203

Kitchen, Wireless Internet, Heating, Smoke Detector, Essentials



Data Selection

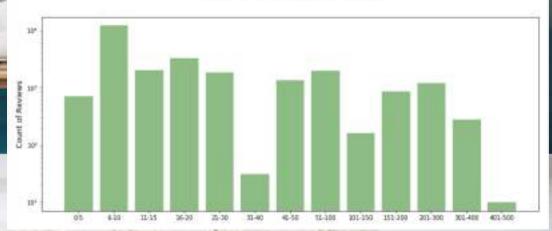


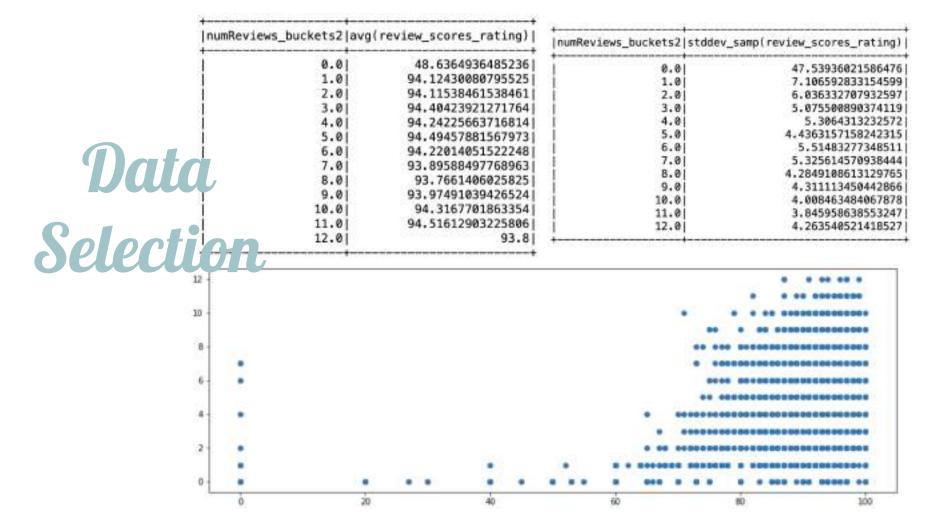
df_bin2.select('number of reviews', 'review scores rating', 'numReviews buckets2').show(18)

Reviews_buckets2	iew_scores_rating numRevi	number_of_reviews
1.0	97.01	61
0.0	80.8	21
0.0	180.8	2
1.0	94.0	7
0.8	0.8	0
1.0	98.81	8
1.0	97.8	6
0.0	93.81	3 3
3.0	92.8	15
0.0	0.0	0

only showing top 18 raws

Number of Reviews Based On Buckets





Data Selection

```
df_bin_amenities = df_bin2.filter(df_bin2|'nunReviews_buckets2"| > 2)
 df_bin_amenities = df_norm.toPandas()

    df_bin_amenities = df_bin_amenities[((df_bin_amenities['mumReviews_buckets2'] == 3.0) 6

                                      (df bin amenities['review scores rating'] > 99:47))
                                      [[df_bin_amenities['numfleviews_buckets2'] == 4.0) &
                                      (df_bin_amenities['review_scures_rating'] > 99.54)) |
                                      ((df bin amenities['numReviews buckets2'] == 5.0) 5
                                      (df bin amenities['review scores rating'] > 98.92)) |
                                      ((df_bin_amenities['munReviews_buckets2'] - 6.0) 4
                                      (df_bin_amenities['review_scores_rating'] > 99.73)) |
                                      ((df bin amenities['numReviews buckets2'] = 7.0) 4
                                      (df bin amenities['review scores rating'] > 99.21)) |
                                      ((df_bin_amenities['musReviews_buckets2'] - 8.0) &
                                      (df bin amenities['review scores rating'] > 98.04)) [
                                      ([df_bin_amenities['munReviews_buckets2'] == 9.8) &
                                      (df_bin_amenities['review_scores_rating'] > 98.28)) |
                                      ((df_bin_amenities['munReviews_buckets2'] == 18.0) &
                                      (df bin amenities['review scores rating'] > 90.31)) [
                                      ([df bin amenities['numberiess buckets2'] == 11.0) 4
                                      (df_bin_amenities['review_scores_rating'] > 98.35)) |
                                      ([df bin amenities['numReviews buckets2'] = 12.0) &
                                      (df bin amenities['review scores rating'] > 40.06))]
```

df bin amenities.count()

numReviews_buckets2	654
review_scores_rating	656
r-score	654
dtype: int64	





Binned Unnormalized

Binned Unnormalized (High Buckets)

```
df_high_bucket = df_high_bucket[(idf_high_bucket['numfleviews_buckets2'] == 8.8) 6

(idf_high_bucket['review_scores_rating'] > 98.84)) |

(idf_high_bucket['review_scores_rating'] > 98.25)) |

(idf_high_bucket['review_scores_rating'] > 98.25)) |

(idf_high_bucket['review_scores_rating'] > 98.31)) |

(idf_high_bucket['review_scores_rating'] > 98.35)) |

(idf_high_bucket['review_scores_rating'] > 98.35)) |

(idf_high_bucket['review_scores_rating'] > 98.35)) |

(idf_high_bucket['review_scores_rating'] > 98.86))|
```

lenidf_high_bucket

87









```
df_norm.groupBy('numReviews_buckets2').agg(mean("z-score"),
stddev("z-score")).sort('numReviews_buckets2', ascending = True).show(10)
numReviews_buckets2|
                            avg(z-score)|stddev_samp(z-score)
                                            1.0000420654117332
                0.0|5.929657557617728...
                1.0 | -8.83210361417153...
                                            1.0001554122318808
                2.0|3.200319002158346...
                                            1.0002531325189254
                3.0 |-6.19909540472002...
                                            1.0003787161658209
                4.0|3.035919598311269...
                                            1.0002766634442195
                5.015.540929590789772...
                                            1.0004172752111498
                6.0 | 8.320172549884078...
                                            1.0005859947763391
                7.0 | 9.564320910206604...
                                            1.0002479851247048
                8.0 4.077720147834147...
                                            1.0007181329471326
                9.0 |-9.23196206856761...
                                            1.0017969466424617
```

only showing top 10 rows

```
df_norm.filter(df_norm['z-score'] > 1.0).count()
```

Binned Normalized Z-Score

4581

Apriori Algorithm (Binned Normalized)

length	itemsets	support	
3	(Kitchen, Heating, Wireless Internet)	0.831041	28
3	(Essentials, Wireless Internet, Kitchen)	0.810303	27
3	(Essentials, Heating, Wireless Internet)	0.800262	26
3	(Kitchen, Smoke detector, Wireless Internet)	0.782580	22
3	(Smoke detector, Heating, Wireless Internet)	0.771011	21
3	(Essentials, Smoke detector, Wireless Internet)	0.765553	20
3	(Kitchen, Heating, Essentials)	0.761842	29
3	(Kitchen, Heating, Smoke detector)	0.738048	25
3	(Kitchen, Smoke detector, Essentials)	0.731281	24
3	(Essentials, Heating, Smoke detector)	0.726916	23

	support	itemsets	length
2	0.964637	(Wireless Internet)	1
5	0.917704	(Kitchen)	1
4	0.898494	(Heating)	1
3	0.877538	(Essentials)	1
1	0.847632	(Smoke detector)	- 1
0	0.742851	(Air conditioning)	1
6	0.714473	(TV)	1

		support	itemsets	length
	15	0.901986	(Kitchen, Wireless Internet)	2
	14	0.883650	(Heating, Wireless Internet)	2
	13	0.863130	(Essentials, Wireless Internet)	2
	19	0.843047	(Kitchen, Heating)	2
	9	0.830605	(Smoke detector, Wireless Internet)	2
	18	0.821436	(Kitchen, Essentials)	2
	17	0.810968	(Essentials, Heating)	2
	12	0.795678	(Kitchen, Smoke detector)	2
	11	0.783890	(Heating, Smoke detector)	2
	10	0.777123	(Essentials, Smoke detector)	2
	7	0.731718	(Air conditioning, Wireless Internet)	2
. 7	16	0.704868	(TV; Wireless Internet)	2
157	8	0.701594	(Air conditioning, Heating)	2

N.		support	itemsets	length
	33	0.753329	(Essentials, Heating, Wireless Internet, Kitchen)	4
	32	0.727352	(Kitchen, Smoke detector, Heating, Wireless In	4
	31	0.721676	(Essentials, Smoke detector, Wireless Internet	4
	30	0.717311	(Essentials, Smoke detector, Heating, Wireless	4

Binned Unnormalized

١		support	itemsets	length
ľ	5	0.992308	(Wireless Internet)	- 1
ı	7	0.969231	(Heating)	- 1
	6	0.936923	(Essentials)	- 1
ĺ	4	0.930769	(Smoke detector)	- 1
Ī	9	0.886154	(Kitchen)	1
ı	2	0.815385	(Hair dryer)	1
l	1	0.803077	(Carbon monoxide detector)	1
ļ	10	0.792308	(TV)	- 1
4	8	0.769231	(Iron)	-1
ı	3	0.767692	(Laptop friendly workspace)	1
	0	0.760000	(Air conditioning)	1

priori

Algorithm

Unnormalized (High Buckets)

			And the second	
N		support	itemsets	length
	5	1.000000	(Wireless Internet)	1
Ц	7	0.965517	(Heating)	1
3	6	0.919540	(Essentials)	1
	4	0.908046	(Smoke detector)	1
	2	0.885057	(Hair dryer)	1
	8	0.873563	(Iron)	1
	9	0.873563	(Kitchen)	1
4	3	0.839080	(Laptop friendly workspace)	1
1	0	0.816092	(Carbon monoxide detector)	1
170		7 4 5	THE PERSON NAMED IN COLUMN TWO IS NOT THE PERSON NAMED IN COLUMN TWO IS NAMED IN CO	The second second

Binned Normalized bin_frqitems[(bin_frqitems['length'] == 4)].head(20)

support	itemsets	length
.753329	(Essentials, Heating, Wireless Internet, Kitchen)	- 4
727352	(Kitchen, Smoke detector, Heating, Wireless In	4
721676	(Essentials, Smoke detector, Wireless Internet	4
717311	(Essentials, Smoke detector, Heating, Wireless	- 4
	753329 727352 721676	.753329 (Essentials, Heating, Wireless Internet, Kitchen) .727352 (Kitchen, Smoke detector, Heating, Wireless In

Association Rules

Binned
Normalized

	antecedents	consequents	antocedent support	consequent support	support	confidence	in	leverage	conviction	antecedent_len
(Es	sentials, Heating, Kitcheni	(Wireless Internet)	0.761842	0.964637	0.753329	0.988825	1.025075	0.018428	3.164576	3
(Esc	entials, Heating, Wireless Internetj	(Kitchen)	0.800262	0.917704	0.753329	0.941353	1.025770	0.018926	1.403250	3
(fix	ventials, Wireless Internet, Kitcheni	(Heating)	0.810303	0.898464	0.753329	0.929688	1,034718	0.025276	1.443644	9
. 0	Otchen, Heating, Wireless Interneti	(Essertish)	0.831041	0.877538	0.753329	0.906488	1.032990	0.024069	1.309590	3
	(Kitchen, Heating, Smoke detector)	(Wireless Internet)	0.738048	0.964637	0.727352	0.985507	1.021636	0.015404	2.440079	0.
0	leating, Wireless Internet, Smoke detector)	(Kitcheri)	0.771011	0.917704	0.727352	0.943375	1.027973	0.019793	1.453055	3

Binned Unnormalized (High Buckets)

	antecedents	consequents	antecedent support	consequent support	support	confidence	Int	leverage	conviction	antecedent_len
43	(Essentials, Wireless Internet, Smoke detector)	(Heating)	0.896652	0.965517	0.873563	0.974359	1.009158	0.007927	1.344828	3
44	Essentials, Heating, Wireless Internet)	(Smoke detector)	0.896552	0.908048	0.873563	0.974359	1.073028	0.059453	3.586207	3
45	(Heating, Wireless Internet, Smoke detector)	(Essentials)	0.885057	0.919540	0.873563	0.987013	1.073377	0.059717	6.195402	3
64	(Essentials, Wireless internet, Smoke detector)	(Hair dryer)	0.896552	0.885067	0.862068	0.961538	1.086414	0.068669	2.988508	3
88	(Essentials, Wireless Internet, Hair dryer)	(Smoke detector)	0.873563	0.908046	0.852069	0.986842	1.086775	0.068833	6.988506	3

FP-Growth

Binned Normalized

	support	itemsets	length
14	0.753329	(Kitchen, Heating, Wireless Internet, Essentials)	4
27	0.727352	(Kitchen, Smoke detector, Heating, Wireless Internet)	4
23	0.721676	(Kitchen, Smoke detector, Wireless Internet, Essentials)	4
24	0.717311	(Essentials, Smoke detector, Heating, Wireless Internet)	4

Binned Unnormalized (High Buckets)

	support	itemsets	length
36	0.839080	(Wireless Internet, Hair dryer, Essentials, Smoke detector, Heating)	5
84	0.839080	(Iron, Wireless Internet, Hair dryer, Essentials, Smoke detector)	5
92	0.827586	(Iron, Wireless Internet, Hair dryer, Essentials, Heating)	5
96	0.827586	(Iron, Wireless Internet, Essentials, Smoke detector, Heating)	5
148	0.816092	(Laptop friendly workspace, Wireless Internet, Hair dryer, Essentials, Smoke detector)	5
97	0.816092	(Iron, Hair dryer, Essentials, Smoke detector, Heating)	5
98	0.816092	(Iron, Wireless Internet, Hair dryer, Smoke detector, Heating)	5

Binned Normalized

Binned Overall & Binned (High Buckets)

Internet,

Essentials,

Heating, Kitchen,

Smoke Detector

Internet, Heating,

Essentials,

Smoke Detector, Hair

Dryer, Iron, Laptop

Friendly Workspace

Hmmmm.....

df_low = df_bin2[df_bin2['review_scores_rating'] < 10]
df_low.count()</pre>

5717

	support	itemsets	length
0	0.910093	(Wireless Internet)	1
1	0.886304	(Kitchen)	1
2	0.823159	(Heating)	1
3	0.794123	(Essentials)	1
4	0.757740	(Smoke detector)	1
5	0.700542	(Air conditioning)	1

	support	itemsets	length
9	0.746545	(Kitchen, Heating, Wireless Internet)	3
12	0.718384	(Kitchen, Wireless Internet, Essentials)	3

	support	itemsets	length	7	support	itemsets l	ength
2	0.978429	(Wireless Internet)	1	5	0.978429	(Wireless Internet)	1
4	0.933544	(Heating)	1	0	0.933544	(Heating)	1:
5	0.917116	(Kitchen)	1	1	0.917116	(Kitchen)	1
3	0.899693	(Essentials)	1	2	0.899693	(Essentials)	1
1	0.873475	(Smoke detector)	1	3	0.873475	(Smoke detector)	1
0	0.753422	(Air conditioning)	1	6	0.753422	(Air conditioning)	1
6	0.743217	(TV)	- 1	4	0.743217	(TV)	1
	support	itemsets	length	→	support	itemsets	length
	0.36380000	version and a second second		5	1,000000	(Wireless Internet)	. 1
2	0.964637	(Wireless Internet)	1	7	0.966517	(Heating	. 1
5	0.917704	(Kitchen)	1	6	0.919540	(Essentials)	1
4	0.898494	(Heating)	1	4	0.908046	(Smoke detector)	. 1
3	0.877538	(Essentials)	1	2	0.886057	(Hair dryer)	1
1	0.847632	(Smoke detector)	1	8	0.873563	(fron)	()i 1
	0.047002	(Sillioke detector)		9	0.873563	(Kitcher)	00 B
0	0.742851	(Air conditioning)	16	3	0.839080	(Laptop friendly workspace)	
6	0.714473	(TV)	1	0	0.816092	(Carbon monoxide detector)	





Apriori Algorithm

length	itemsets	support	
3	(Kitchen, Heating, Wireless Internet)	0.882174	40
3	(Essentials, Heating, Wireless Internet)	0.850520	38
3	(Air conditioning, Heating, Wireless Internet)	0.841886	27
3	(Essentials, Wireless Internet, Kitchen)	0.836611	39
3	(Kitchen, Air conditioning, Wireless Internet)	0.823341	28
3	(Smoke detector, Heating, Wireless Internet)	0.813749	33

The same of the same of	The state of the s		
length	itemsets	support	
3	(Essentials, Heating, Wireless Internet)	0.826313	37
3	(Smoke detector, Heating, Wireless Internet)	0.818127	30
3	(Kitchen, Heating, Wireless Internet)	0.817528	10
3	(Essentials, Smoke detector, Wireless Internet)	0.816530	29
3	(Essentials, Wireless Internet, Kitchen)	0.802156	38

East Coast



West Coast



FP-Growth Binned Unnormalized







East Coast

	support	itemsets	length
37	0.845361	(Wireless Internet, Essentials, Kitchen, Air conditioning, Heating)	5
53	0.835052	(Wireless Internet, Hair dryer, Essentials, Air conditioning, Heating)	5
64	0.814433	(Wireless Internet, Hair dryer, Kitchen, Air conditioning, Heating)	5
85	0.814433	(Wireless Internet, Essentials, Smoke detector, Air conditioning, Heating)	5

West Coast

itemsets	support	
(Wireless Internet, Hair dryer, Essentials Smoke detector, Heating)	0.848168	37
aptop friendly workspace Wireless Internet, Hair dryer, Essentials, Smoke detector)	0.837696	85
(Iron, Wireless Internet, Hair dryer, Essentials, Smoke detector)	0.832461	53
(Iron, Wireless Internet, Hair dryer, Essentials, Heating)	0.821990	61
(Iron, Wireless Internet, Essentials, Smoke detector, Heating)	0.821990	65
	(Wireless Internet, Hair dryer, Essentials Smoke detector, Heating) Laptop friendly workspace. Wireless Internet, Hair dryer, Essentials, Smoke detector) (Iron, Wireless Internet, Hair dryer, Essentials, Smoke detector) (Iron, Wireless Internet, Hair dryer, Essentials, Heating)	0.848168 (Wireless Internet, Hair dryer, Essentials, Smoke detector, Heating) 0.837696 (Laptop friendly workspace) Wireless Internet, Hair dryer, Essentials, Smoke detector) 0.832461 (Iron, Wireless Internet, Hair dryer, Essentials, Smoke detector) 0.821990 (Iron, Wireless Internet, Hair dryer, Essentials, Heating)

A retire or

	support	itemsets	length
9	0.854575	(Kitchen, Heating, Wireless Internet)	3
13	0.821015	(Kitchen, Wireless Internet, Essentials)	3
14	0.808230	(Essentials, Heating, Wireless Internet)	3

		support	itemsets	length
-	11	0.786566	(Kitchen, Wireless Internet, Essentials)	3
-	36	0.786024	(Kitchen, Heating, Wireless Internet)	3
	37	0.777898	(Essentials, Heating, Wireless Internet)	3
	16	0.774648	(Kitchen, Smoke detector, Wireless Internet)	3

FP-Growth Binned Normalized

East Coast

.

- West Coast

Conclusion

- The number of reviews do not have a significant impact on review score rating.
- 2. Sets of amenities do not provide a clear indication of what amenity, or amenities are needed to guarantee a high review score.
- 3. For Airbnb hosts looking to increase their review scores:
 - a. There seems to be a mandatory set of amenities that are needed.

Conclusion

- 4. The seemingly "mandatory" set of amenities show up in "poor" review scored AirBnB's.
- 5. We now know what information we need to come to more definitive conclusion.
 - a. Specific information about why a review score is bad.
 - i. Business trip vs. vacation vs. personal travel
 - ii. Bad neighborhood (factors indicating why)

References

- 1. https://towardsdatascience.com/plotting-business-locations-on-maps-using-multiple-plotting-libraries-in-python-45a00ea770af
- 2. https://www.pinterest.com/pin/163044448983653898/
- 3. https://www.pinterest.com/pin/87820261463065433/
- 4. https://paintingvalley.com/architecture-sketch-wallpaper
- 5. https://www.elledecor.com/design-decorate/g9096337/beautiful-views/
- 6. https://pdfs.semanticscholar.org/3bc3/6381f77ec371b0931d3276a3f4f58d186427.pdf
- 7. https://www.thesleepjudge.com/white-bedroom-ideas/
- 8. https://www.youtube.com/watch?v=quVvtZ7ZClw
- 9. https://www.kdnuggets.com/2016/04/association-rules-apriori-algorithm-tutorial.html/2
- 10. https://www.muralswallpaper.com/shop-murals/black-brick-wallpaper-mural/
- 11. https://bluemarbleprivate.com/insight/arijiju-sirai-luxury-safari-lodges
- 12. https://beyond-universe.fandom.com/wiki/Corridor of All-Creator
- 13. https://brickhunter-store.myshopify.com/collections/wall-thin-brick/products/old-chicago-blanc-flat
- 14. https://www.pinterest.com/pin/196188127497314831/
- 15. https://waliicorners.com/product/custom-fashion-wallpaper-geometric-gold-lines-modern-minimalistic-luxury-abstract-background-wall-papers-home-decor/
- 16. https://www.pinterest.com/pin/181340322471070745/
- 17. http://rasbt.github.io/mlxtend/user_guide/frequent_patterns/association_rules/
- 18. http://rasbt.github.io/mlxtend/api_subpackages/mlxtend.frequent_patterns/#fpgrowth