

Breaking Even With



[REDACTED] & Victoria Sharam

Problem
Overview

Big Data
Issue:
Volume

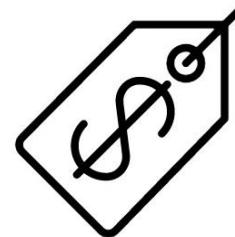
Analyzing
Text Data

Machine
Learning

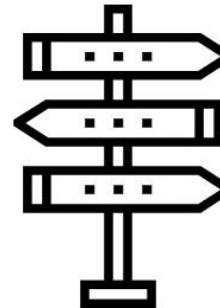
Solution

How can an NYC Airbnb host optimize their experience and aim to break even on average monthly rent

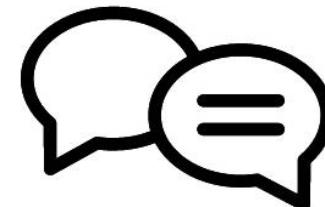
Price



Borough



Number of
Reviews



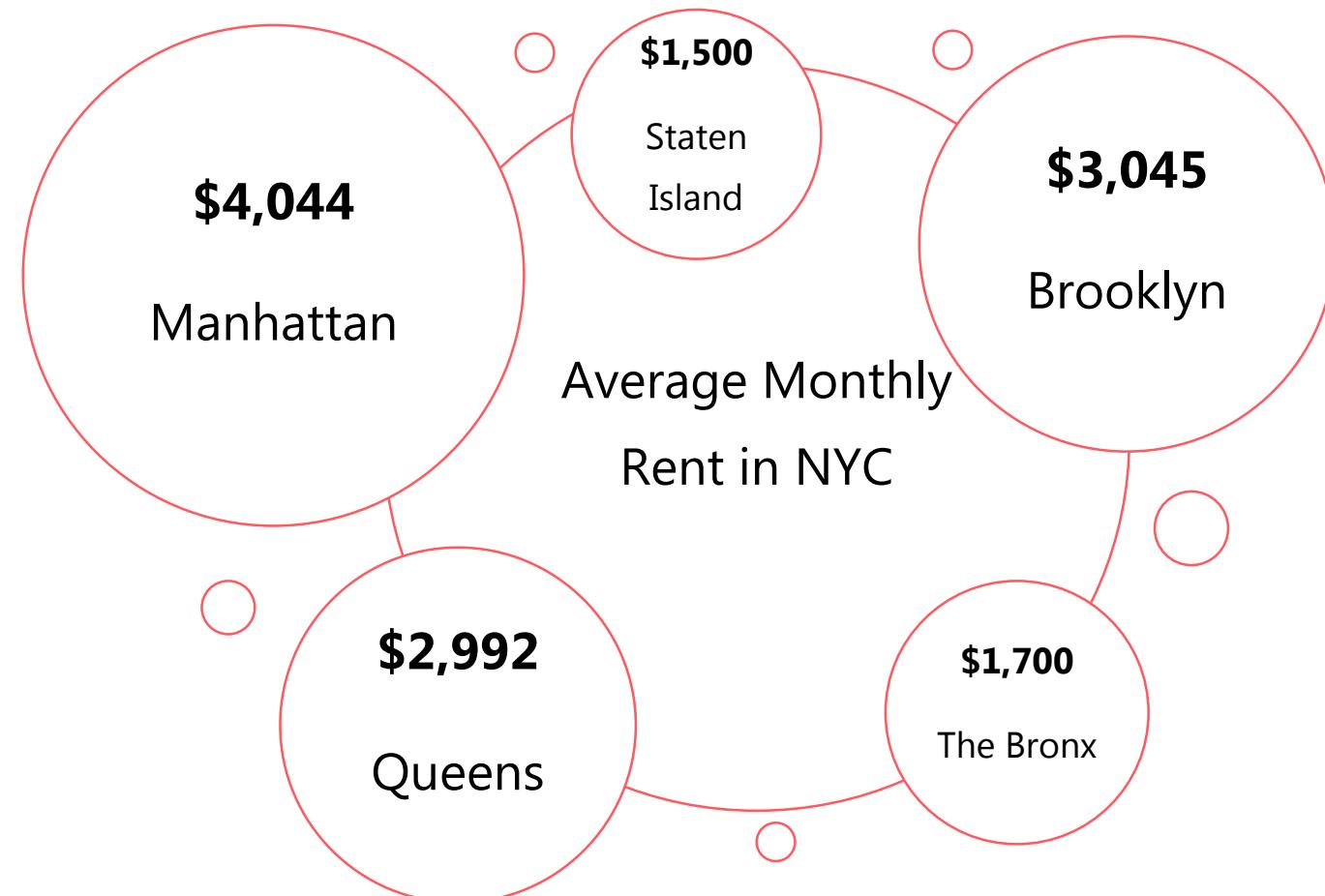
Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution



Problem Overview

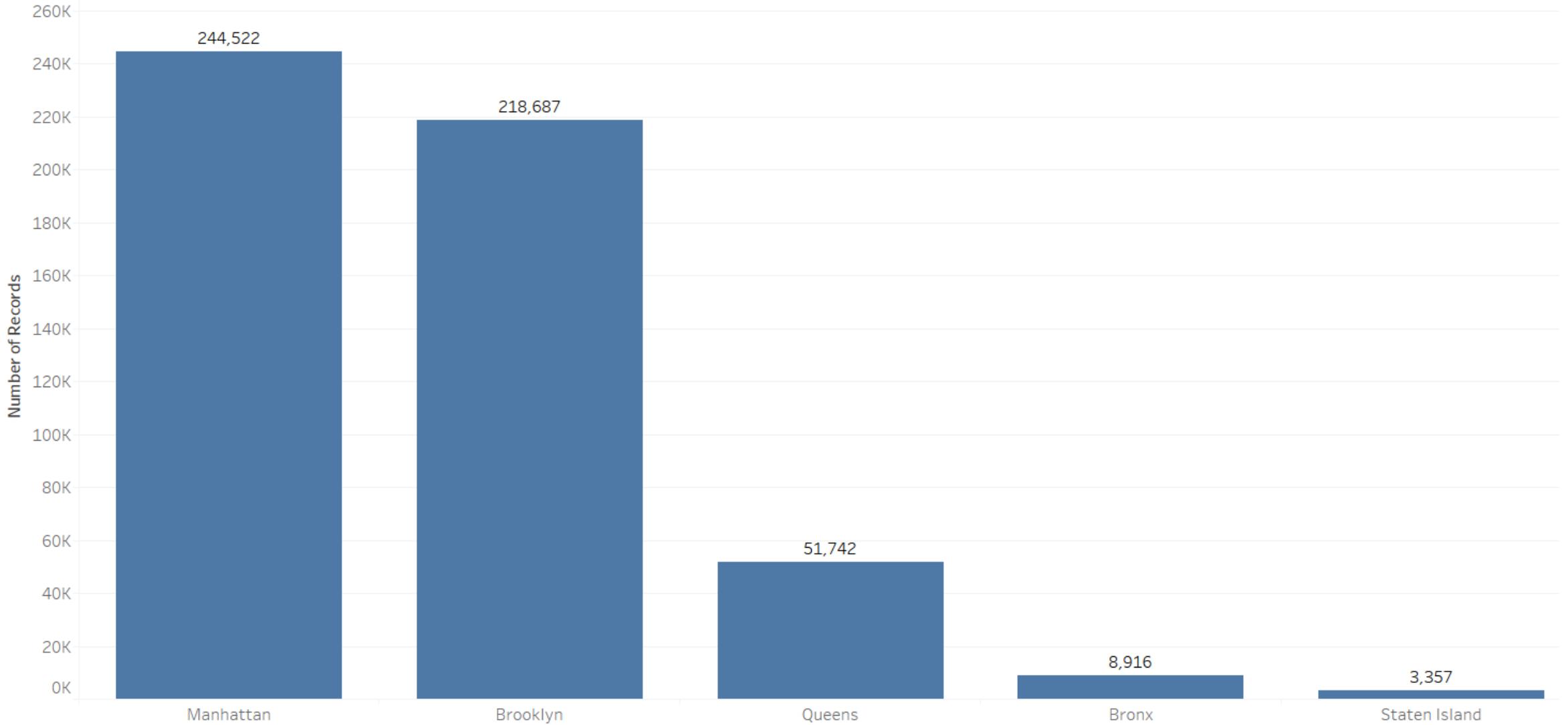
Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Number of Listings by Borough



Problem Overview

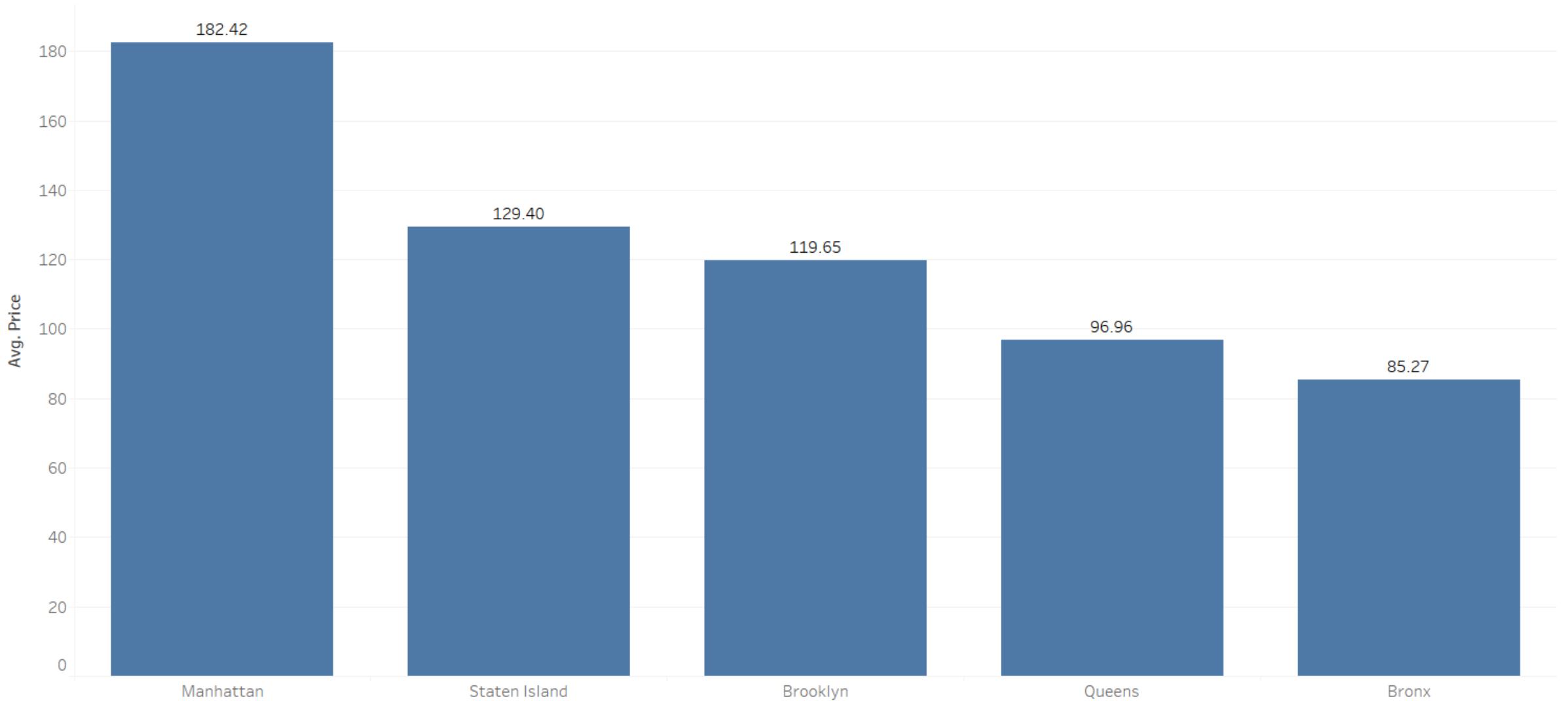
Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Average Daily Price by Borough



Problem Overview

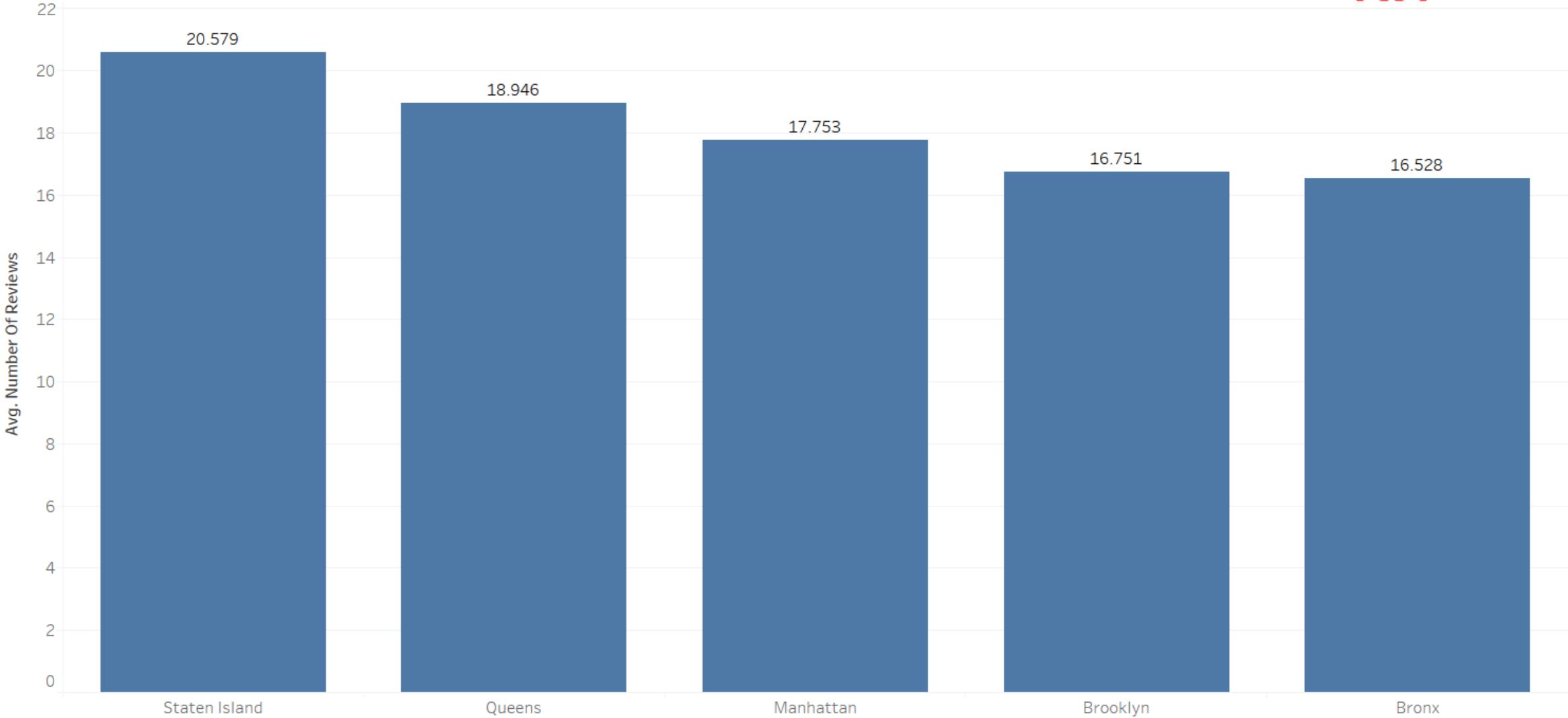
Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Average Number of Reviews by Borough



Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Dataset Structure



	V1	V2	V3	V4	V5	...	n variables
3/17							
4/17							
5/17							
6/17							
7/17							
8/17							
9/17							
10/17							
11/17							
12/17							
01/18							
02/19							

Problem Overview

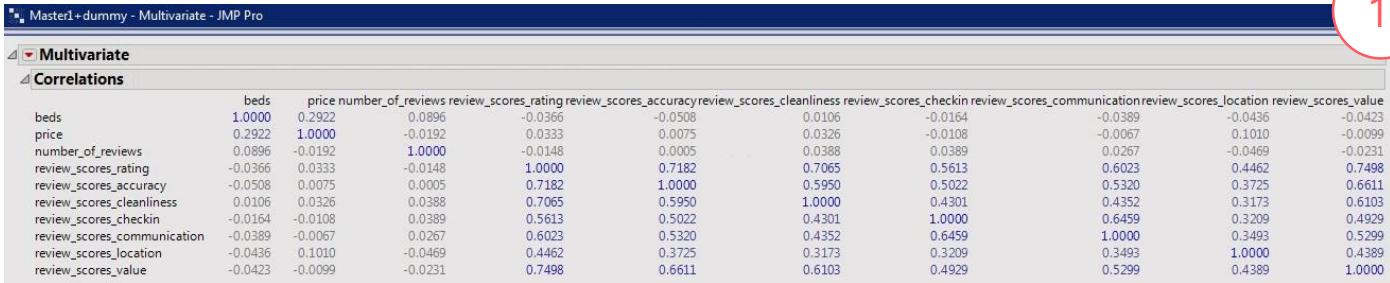
Big Data Issue:
Volume

Analyzing Text Data

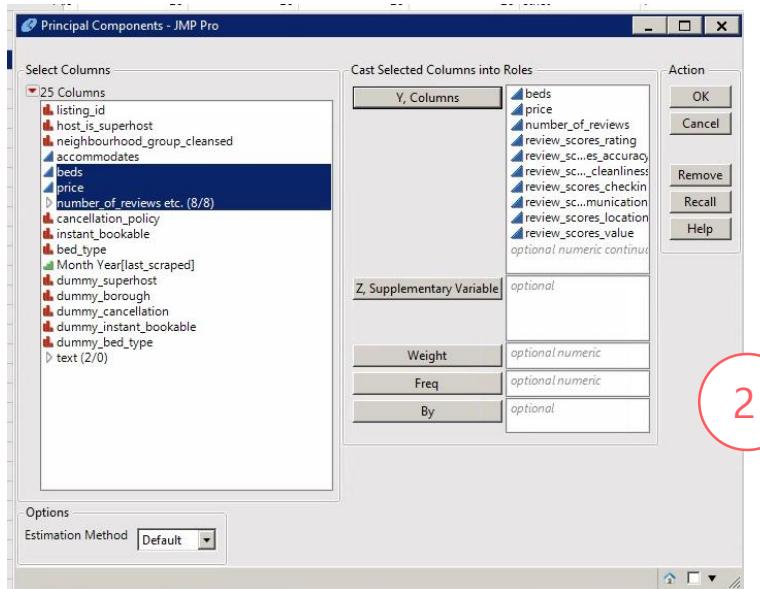
Machine Learning

Solution

Principal Component Analysis

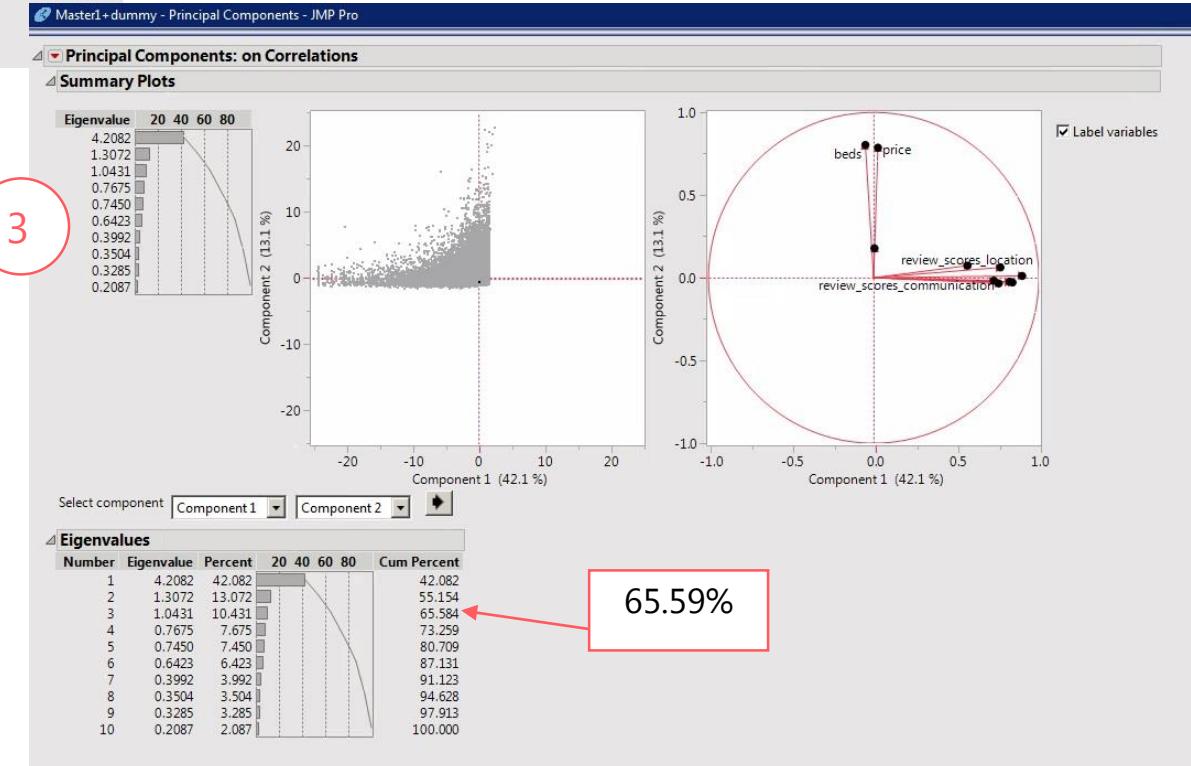


There are 129368 missing values. The correlations are estimated by Pairwise method.



Problem Overview

Big Data Issue: Volume

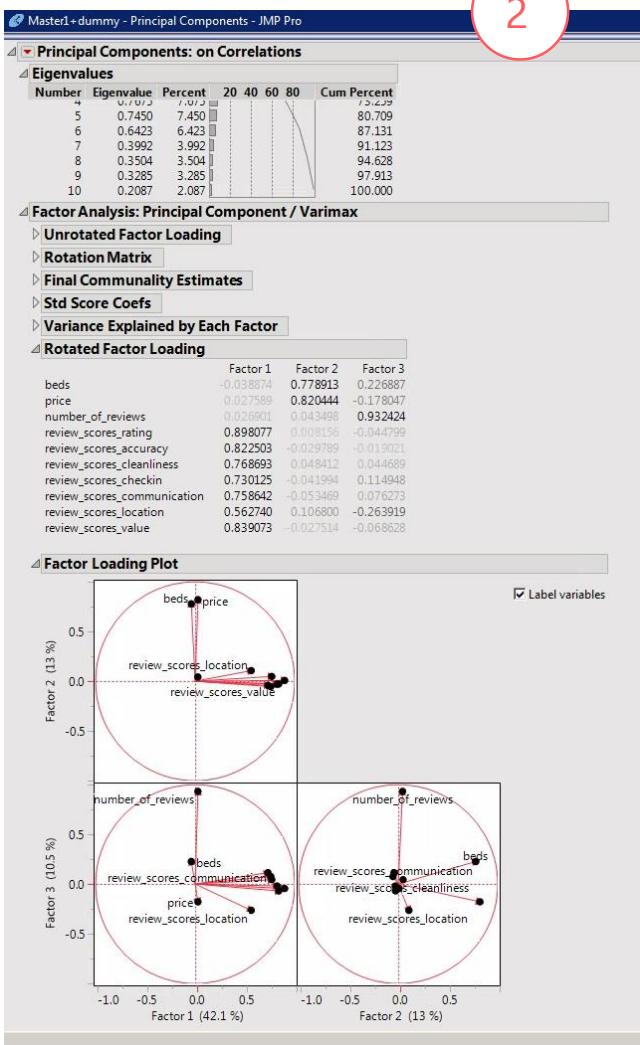
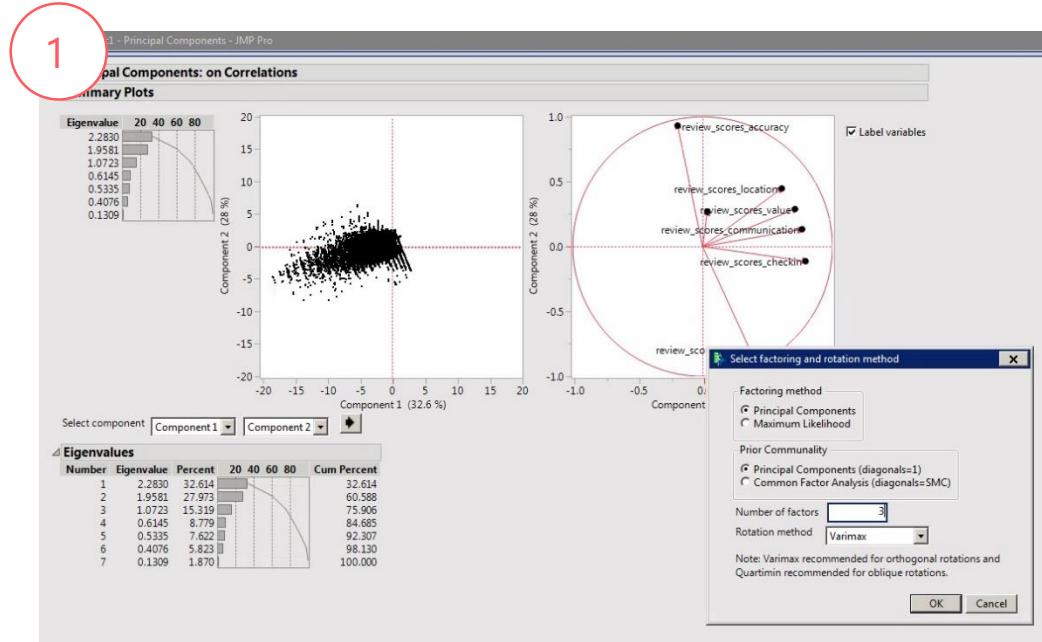


Analyzing Text Data

Machine Learning

Solution

Factor Analysis



3

	Ratings	Price	Demand
1t...	0.7688540069	0.1276002768	1.7975116064
1s...	0.7506345094	-0.48005842	1.825933107
r...	*	*	
1t...	*	*	
e...	0.7226102519	-0.320597141	0.5840721271
e...	0.7846407826	-0.400081871	0.0718560692
r...	0.6245799608	-0.331866173	-0.291992608
e...	0.1317264116	-1.003301431	0.6235131406
e...	0.476683409	-0.690226063	0.6095977309
1t...	0.0599015674	-0.57740105	0.6093752716
r...	-0.69222349	-0.20961703	-1.033723603
1t...	-0.172765094	-0.455334356	0.1297546837
1t...	0.3736956841	-0.662041801	0.0602820684
r...	0.5366883785	-0.486688509	-0.182506436
1t...	-1.055139993	-0.001345066	1.8603857767
r...	-0.082958683	0.7271330068	5.84782703
r...	-0.20265644	-0.949411994	0.4870919543
r...	-0.262830801	-0.671577487	1.2717512605
r...	-0.005128513	-0.599494862	-0.157190113
e...	0.7662425137	-0.590469485	-0.536800611
e...	0.2832513796	-0.651123952	-0.094352287
1t...	-0.912705002	-0.09579828	2.7203849372
1t...	0.1008592407	-0.687577331	0.6405033561
r...	-0.132641127	1.3240136008	1.0421945062
r...	-0.952935959	0.0143658987	2.9016628482
r...	0.168575467	-0.172712373	0.3161820943
e...	0.4309000775	-0.564071405	-0.69266885
1s...	-0.415845518	-0.223646457	2.0878035698
e...	0.3816894332	-0.706253075	1.4367930193

Problem Overview

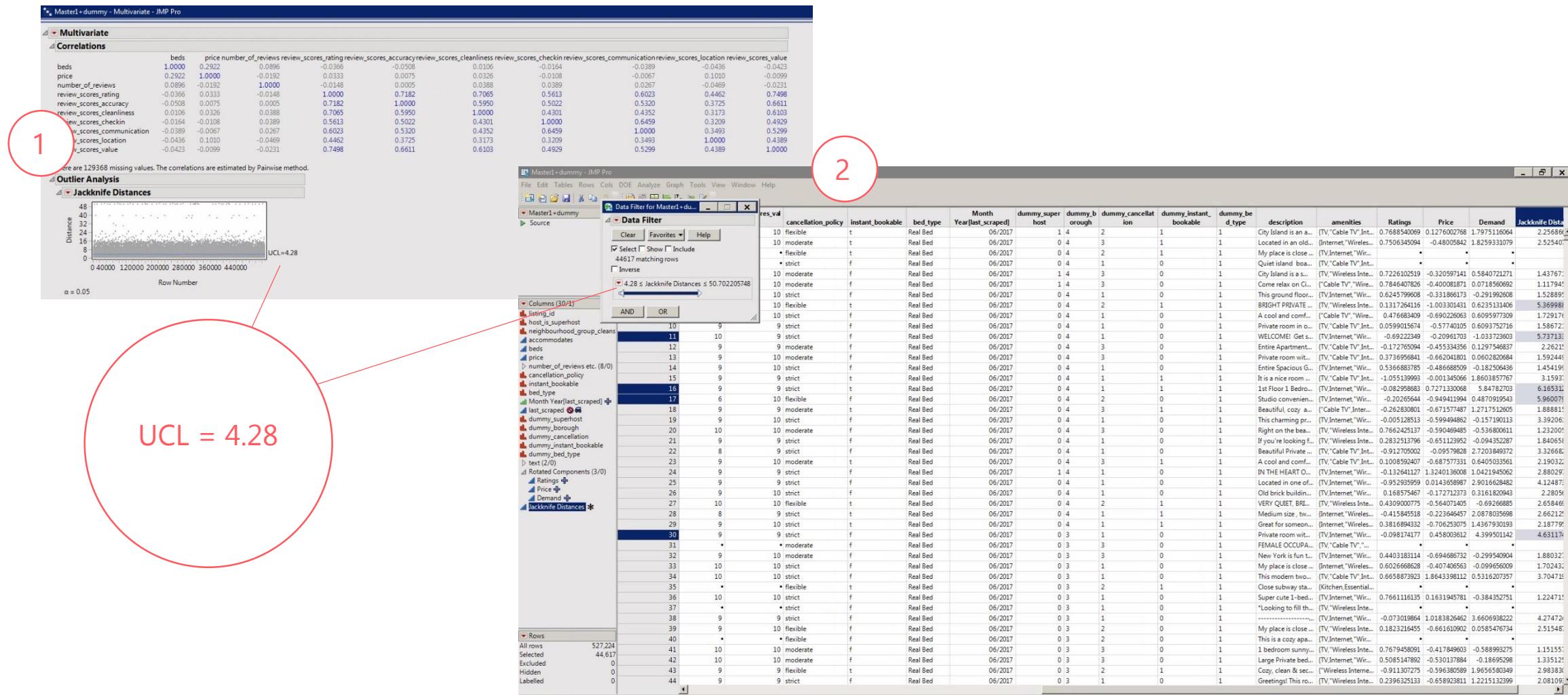
Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Clustering: Outlier Analysis



Problem Overview

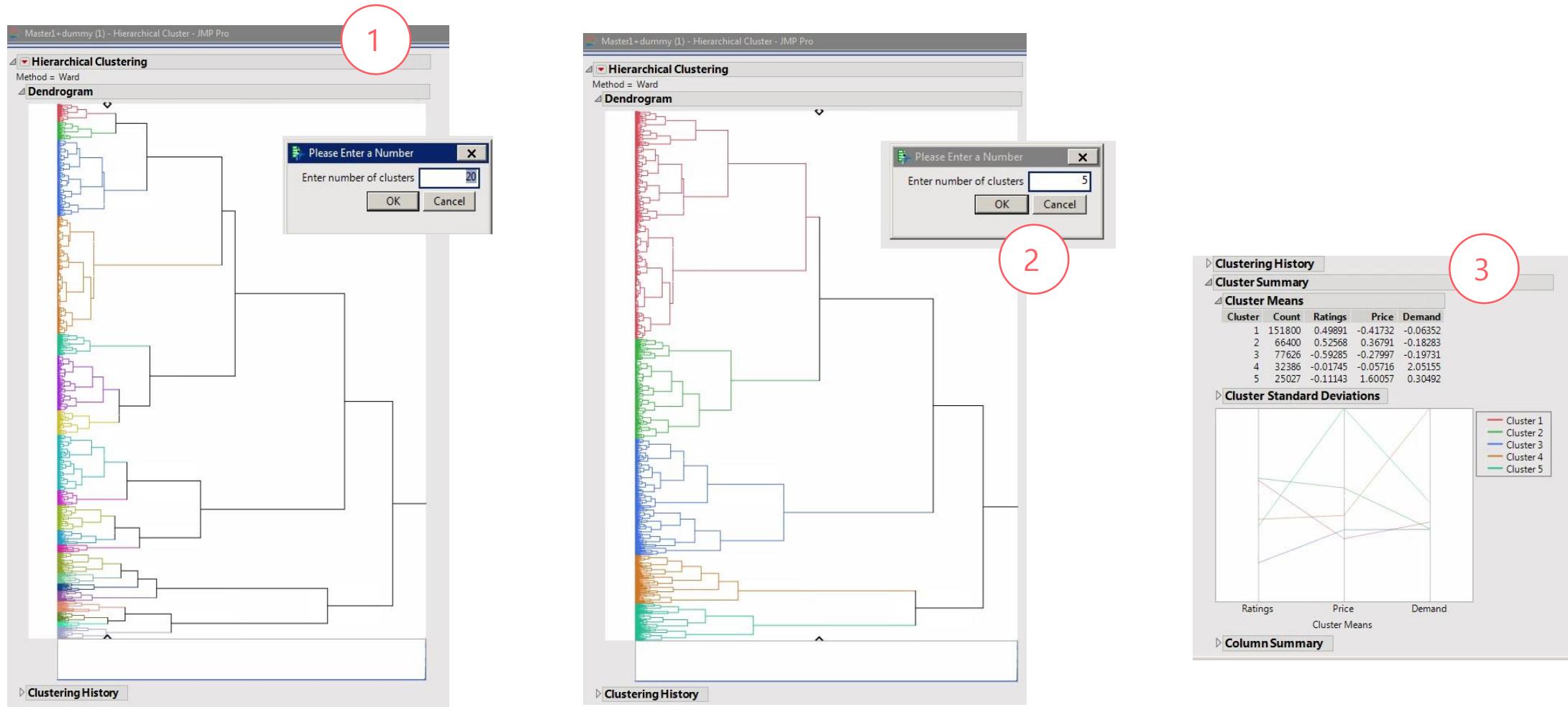
Big Data Issue: Volume

Analyzing Text Data

Machine Learning

Solution

Hierarchical Clustering by Rotated Components



Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Clusters

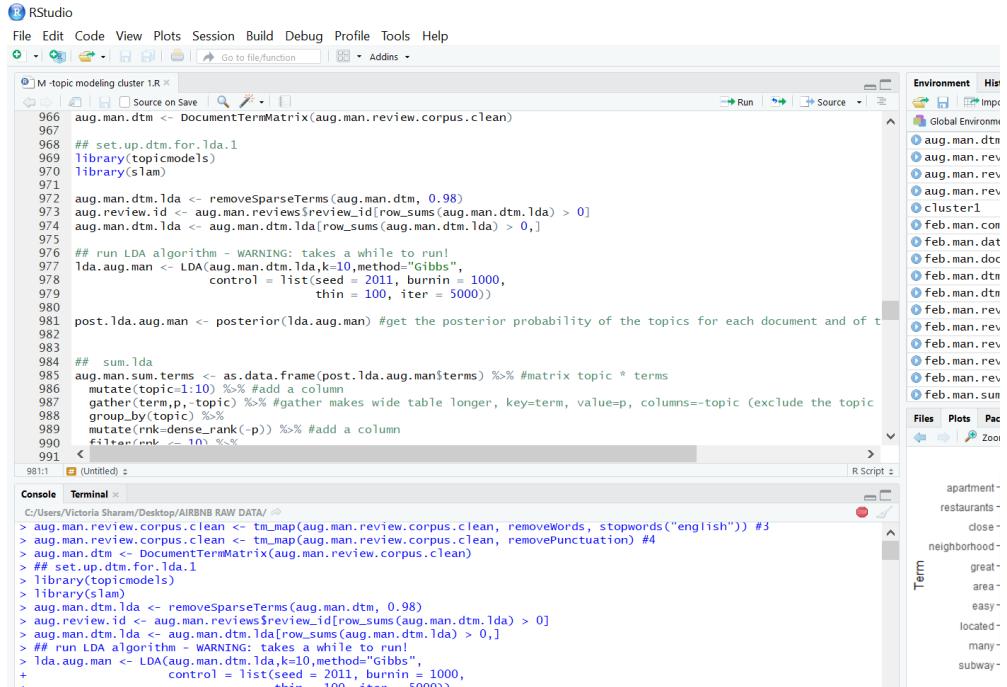


Text Mining: Topic Modeling

1

Joined clusters with review data on JMP

2



The screenshot shows an RStudio interface with the following details:

- Code Editor:** An R script titled "M-topic modeling cluster 1.R" containing code for Topic Modeling. The code includes steps like cleaning the corpus, removing sparse terms, calculating document-term matrices, and running the LDA algorithm with Gibbs sampling.
- Environment View:** Shows various objects in the global environment, including "aug.man.dtm", "aug.man.rev", and several "feb.man" objects.
- Term Frequency Matrix:** A large matrix titled "Term" showing term frequencies across documents. The visible terms include "apartment", "restaurants", "close", "neighborhood", "great", "area", "easy", "located", "many", and "subway".
- Console:** Displays the R commands run in the terminal, corresponding to the code in the editor.

Problem Overview

Big Data Issue:
Volume

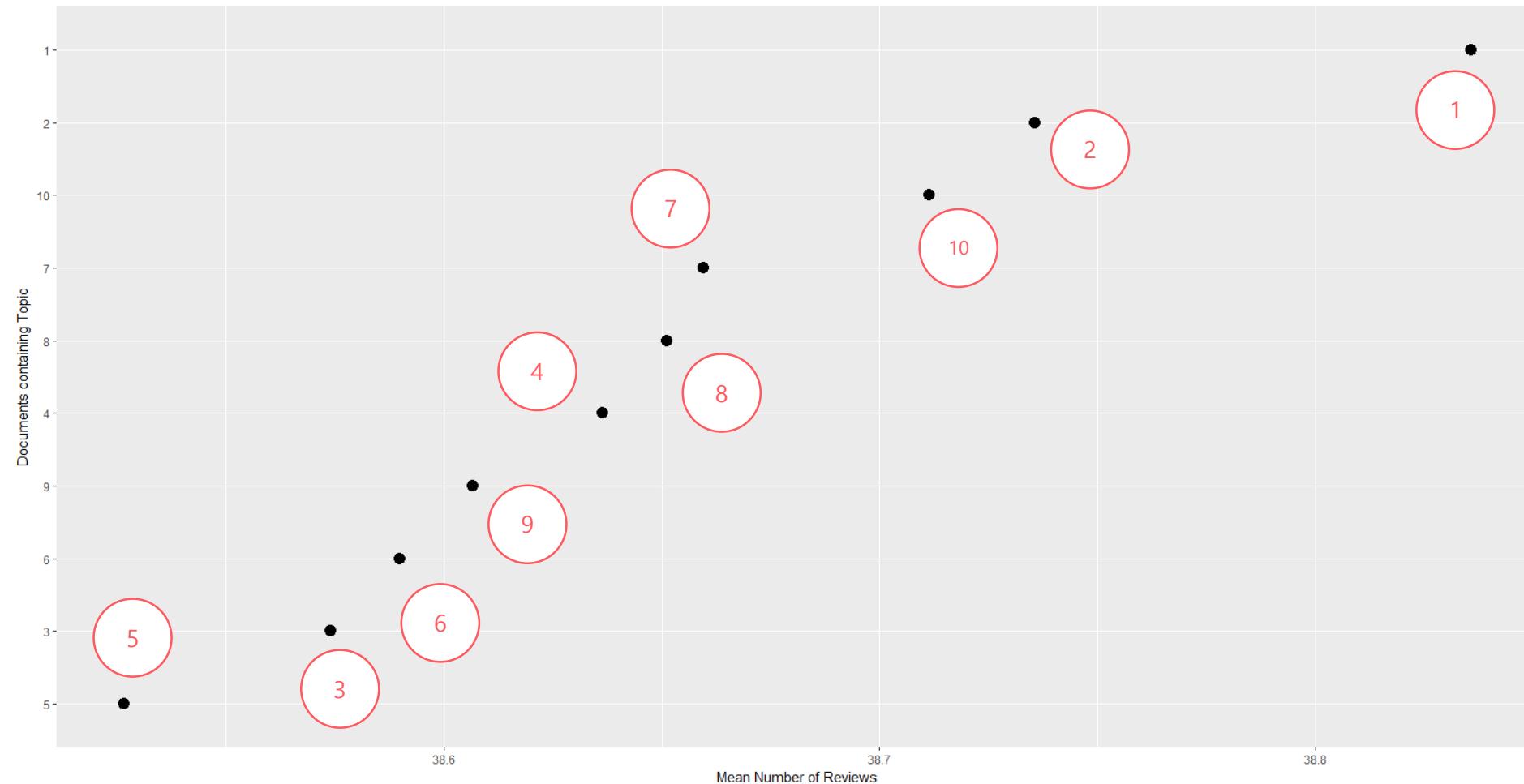
Analyzing Text Data

Machine Learning

Solution

Text Mining: Topic Modeling

Cluster One, Manhattan, February



Problem Overview

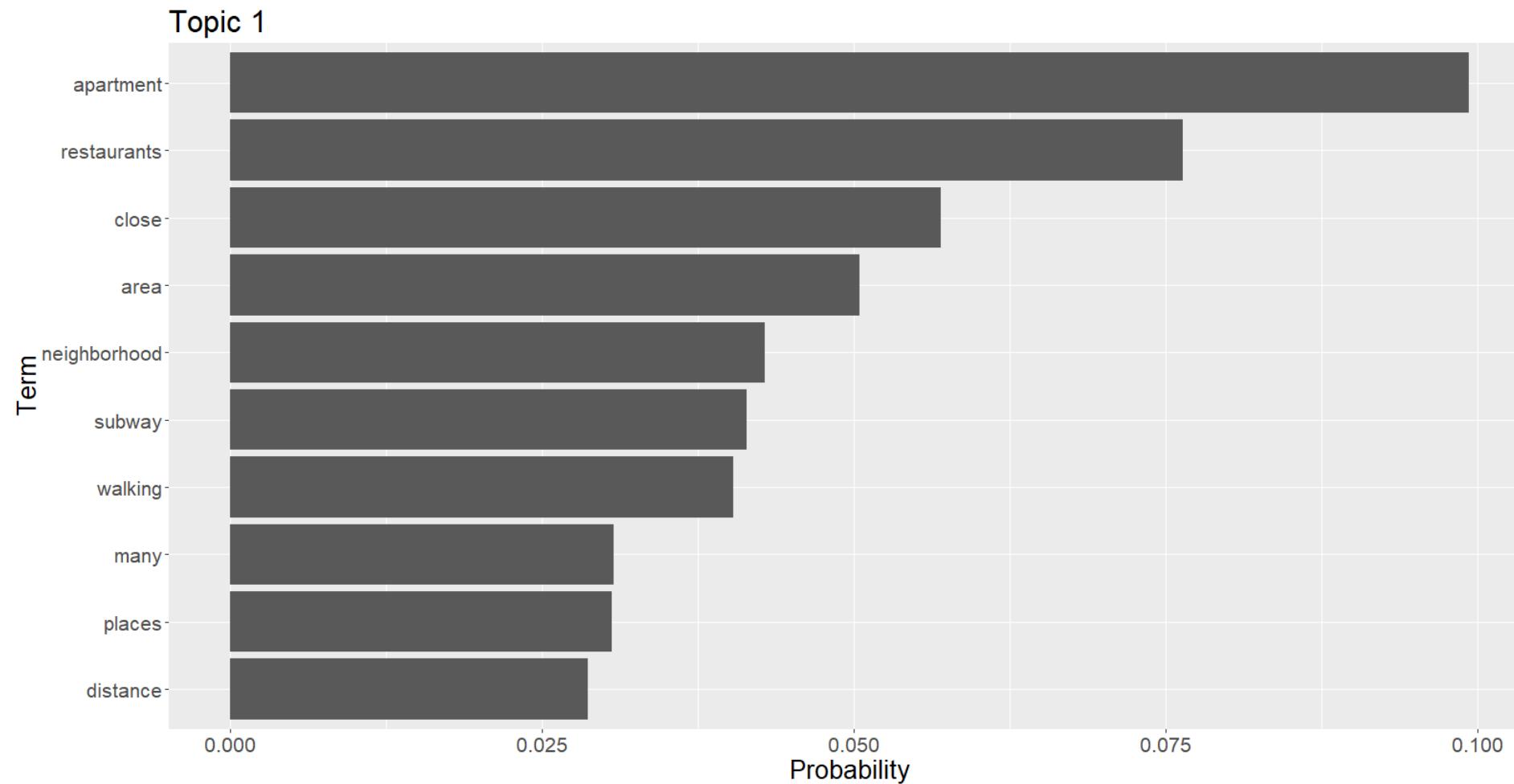
Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Text Mining: Locality



Problem Overview

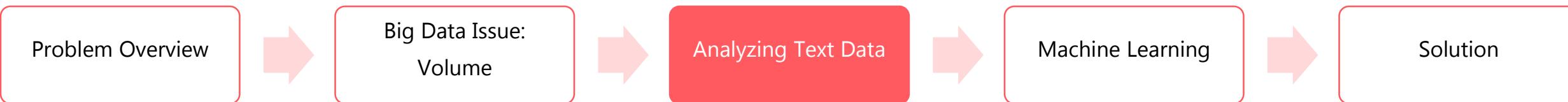
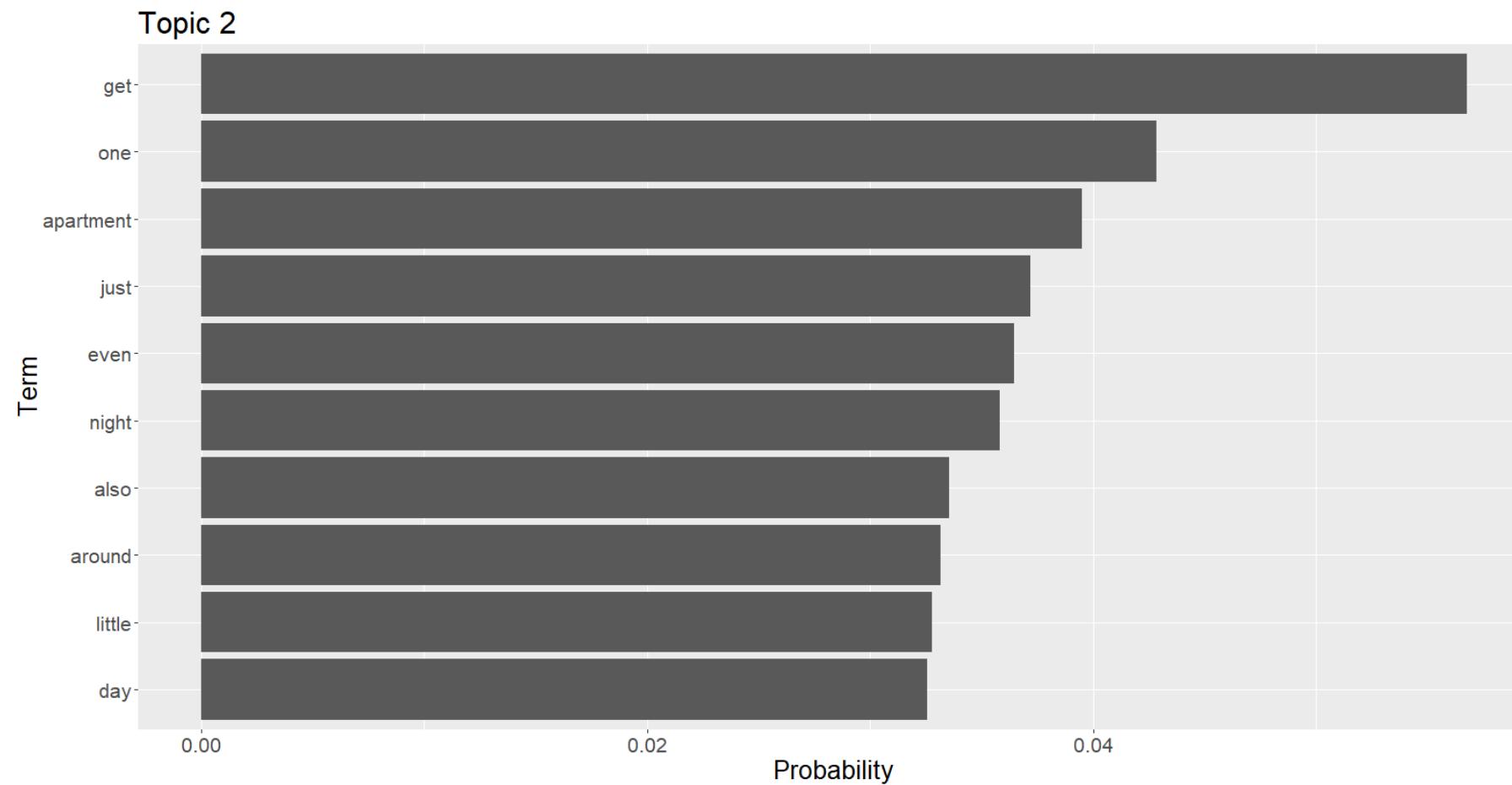
Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Text Mining: Ephemeral



Dataset Structure



	Price	Demand	Superhost	Bed Type	Accommodates	Beds	Cancellation Policy	Review Score Rating	Instant Bookable	Topic Modeling
3/17										
4/17										
5/17										
6/17										
7/17										
8/17										
9/17										
10/17										
11/17										
12/17										
01/18										
02/19										

Problem Overview

Big Data Issue:
Volume

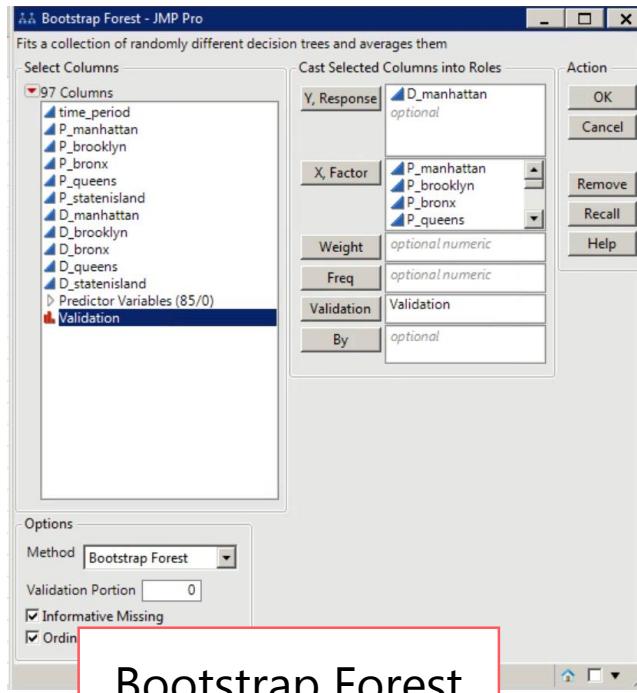
Analyzing Text Data

Machine Learning

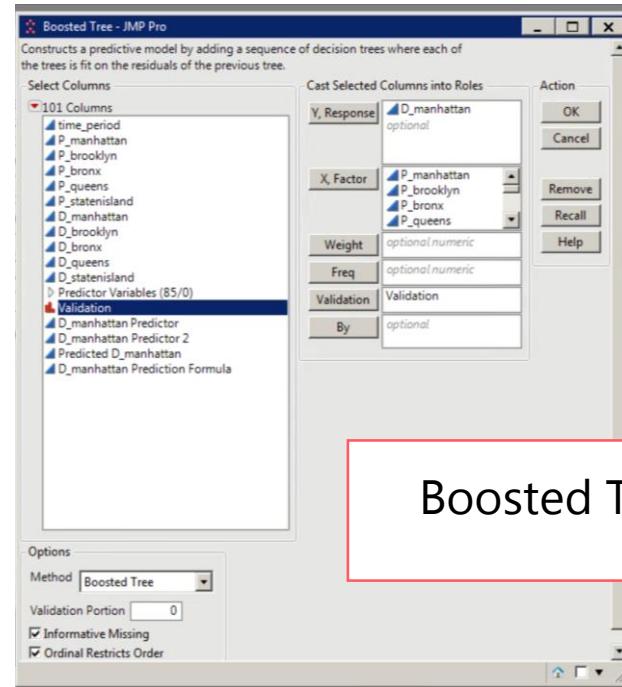
Solution

Machine Learning

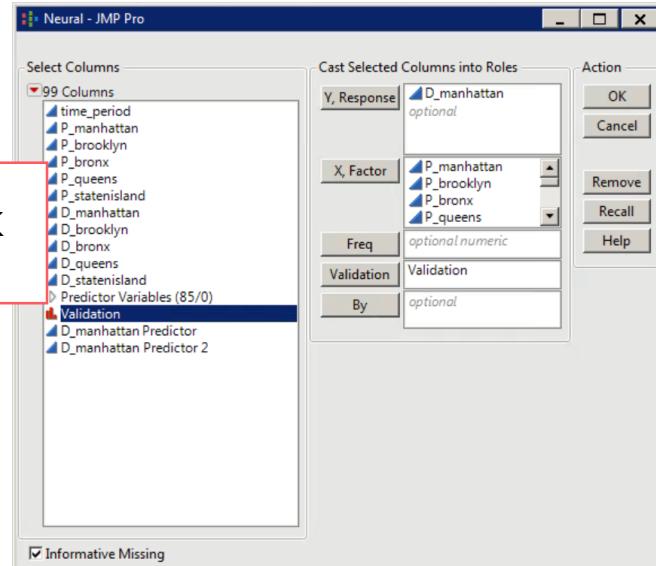
Cluster One, Manhattan



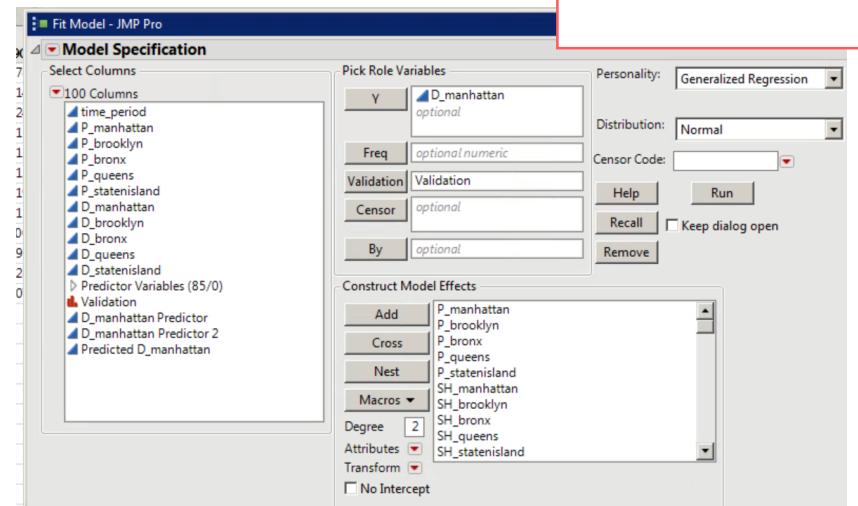
Bootstrap Forest



Boosted Tree



Neural Network



Elastic Net

Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Machine Learning with Tuned Parameters

Bootstrap Forest

Bootstrap Forest Specification

Number of Rows: 12
Number of Terms: 90

Forest

- Number of Trees in the Forest: **300**
- Number of Terms Sampled per Split: **10**
- Bootstrap Sample Rate: 1
- Minimum Splits per Tree: 10
- Maximum Splits per Tree: 2000
- Minimum Size Split: 5
- Early Stopping

OK **Cancel**

Boosted Tree

Gradient-Boosted Trees Specification

Boosting

- Number of Layers: **80**
- Splits per Tree: **2**
- Learning Rate: 0.1
- Minimum Size Split: 6

Multiple Fits

- Multiple Fits over Splits and Learning Rate
- Max Splits Per Tree: 4
- Max Learning Rate: 0.1
- Use Tuning Design Table

Stochastic Boosting

- Row Sampling Rate: 1.0000
- Column Sampling Rate: 1.0000

Reproducibility

- Suppress Multithreading
- Random Seed: 0

Early Stopping

OK **Cancel**

Sheet1 - Neural of D_manhattan - JMP Pro

Neural

Validation Column: Validation
Informative Missing

Model Launch

Hidden Layer Structure

Number of nodes of each activation type
Activation Sigmoid Identity Radial

Layer	Tanh	Linear	Gaussian
First	1	0	0
Second	0	0	0

Second layer is closer to X's in two layer models.

Boosting

Fit an additive sequence of models scaled by the learning rate.

Number of Models: 0
Learning Rate: 0.1

Fitting Options

- Transform Covariates
- Robust Fit
- Penalty Method: Squared
- Number of Tours: 1

Go

Model Comparison

Model Comparison - JMP Pro

Comparing predictors to see which performs better.

Select Columns

104 Columns

- D_manhattan_P..dition Formula
- D_manhattan Predictor 3
- D_manhattan Predictor 4
- Predicted D_manhattan 2
- D_manhattan Predictor
- D_manhattan Predictor 2
- D_manhattan Predictor 3
- D_manhattan Predictor 4
- Predicted D_manhattan 2

Cast Selected Columns into Roles

Y, Predictors: D_manhattan Predictor 3
Group: optional
Weight: optional numeric
Freq: optional numeric
By: Validation
Validation: optional

If you choose no Predictor columns, it will find and analyze all predictors.

Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Machine Learning Model Comparison: RASE

Without Tuning

Sheet1 - Model Comparison - JMP Pro						
Model Comparison Validation=Training						
Predictors						
Measures of Fit for D_manhattan						
Predictor	Creator	.2 .4 .6 .8	RSquare	RASE	AAE	Freq
D_manhattan Predictor	Bootstrap Forest		-0.017	0.4847	0.4079	8
D_manhattan Predictor 2	Boosted Tree		0.0000	0.4806	0.4079	8
Predicted D_manhattan	Neural		0.8277	0.1995	0.1655	8
D_manhattan Prediction Formula	Fit Generalized Adaptive Elastic Net		0.9237	0.1327	0.1078	8
Model Comparison Validation=Validation						
Predictor	Creator	.2 .4 .6 .8	RSquare	RASE	AAE	Freq
D_manhattan Predictor	Bootstrap Forest		-0.023	0.4833	0.4253	4
D_manhattan Predictor 2	Boosted Tree		-0.080	0.4965	0.4566	4
Predicted D_manhattan	Neural		0.9575	0.0985	0.0951	4
D_manhattan Prediction Formula	Fit Generalized Adaptive Elastic Net		0.8082	0.2092	0.1740	4

With Tuning

Sheet1 - Model Comparison - JMP Pro						
Model Comparison Validation=Training						
Predictors						
Measures of Fit for D_manhattan						
Predictor	Creator	.2 .4 .6 .8	RSquare	RASE	AAE	Freq
D_manhattan Prediction Formula	Fit Generalized Adaptive Elastic Net		0.9237	0.1327	0.1078	8
D_manhattan Predictor 3	Bootstrap Forest		-0.091	0.5020	0.4079	8
D_manhattan Predictor 4	Boosted Tree		0.0000	0.4806	0.4079	8
Predicted D_manhattan 2	Neural		0.9424	0.1153	0.0897	8
Model Comparison Validation=Validation						
Predictor	Creator	.2 .4 .6 .8	RSquare	RASE	AAE	Freq
D_manhattan Prediction Formula	Fit Generalized Adaptive Elastic Net		0.8082	0.2092	0.1740	4
D_manhattan Predictor 3	Bootstrap Forest		-0.000	0.4779	0.3841	4
D_manhattan Predictor 4	Boosted Tree		-0.080	0.4965	0.4566	4
Predicted D_manhattan 2	Neural		0.9436	0.1135	0.0845	4

Neural Network

Problem Overview

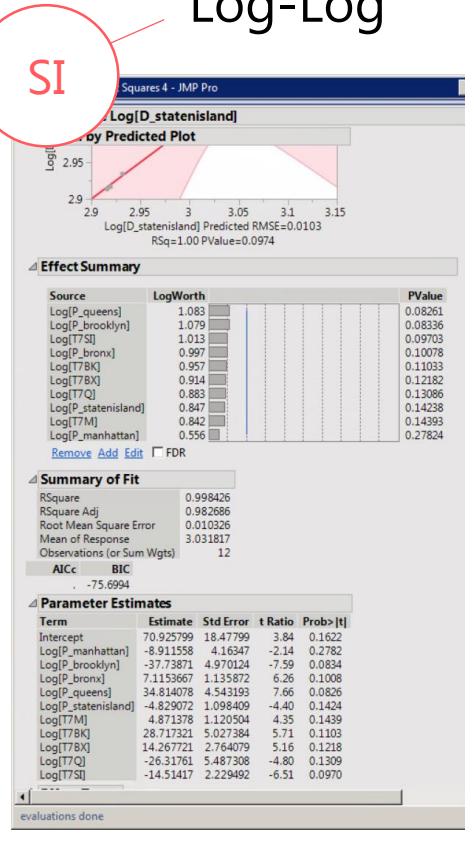
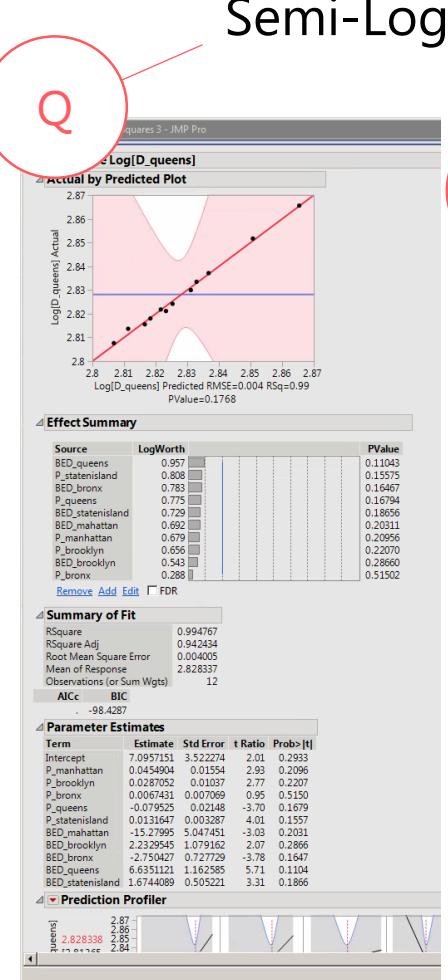
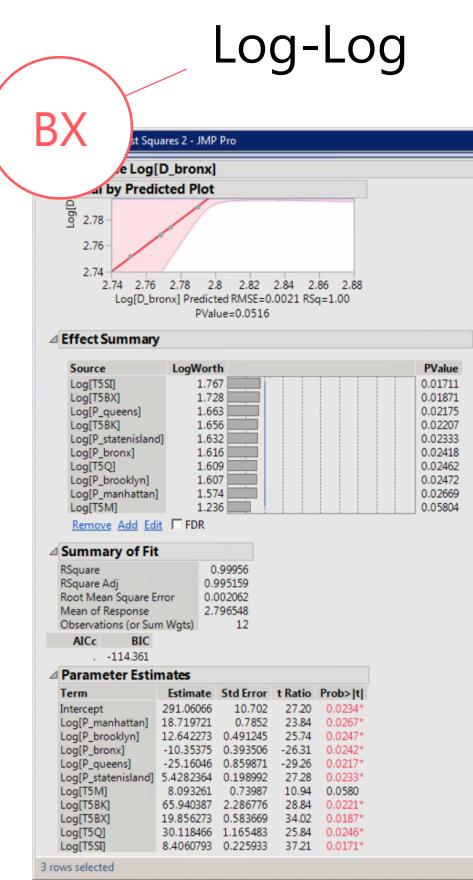
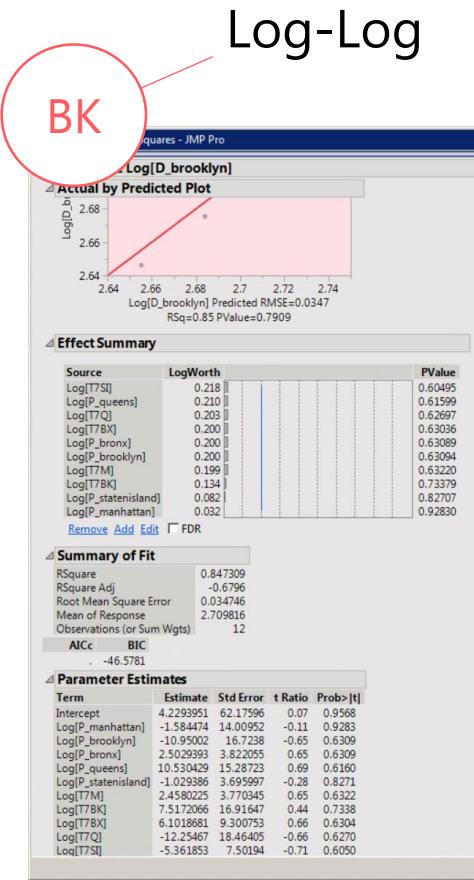
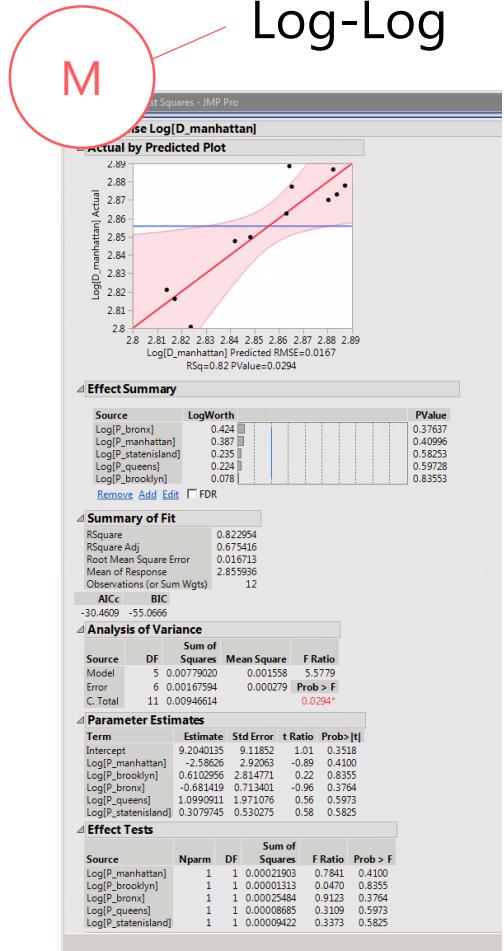
Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Optimal Pricing for Airbnb Hosts



Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

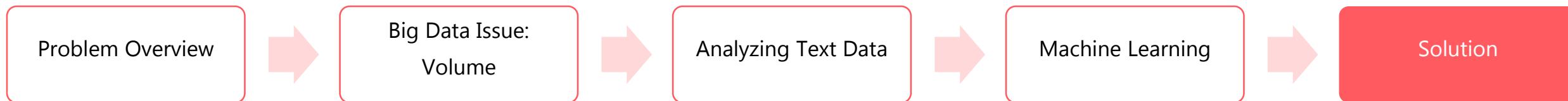
Solution

Manhattan Optimal Daily Price



\$220 / day

Based on this price, a Manhattan host would have to host a guest for **19 days** to make their monthly rent.



Brooklyn Optimal Daily Price



\$112 / day

Based on this price, a Brooklyn host would have to host a guest for **28 days** to make their monthly rent.



Bronx Optimal Daily Price



\$63 / day

Based on this price, a Bronx host would have to host a guest for **27 days** to make their monthly rent.



Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Queens Optimal Daily Price

\$120 / day

Based on this price, a Queens host would have to host a guest for **25 days** to make their monthly rent.



Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

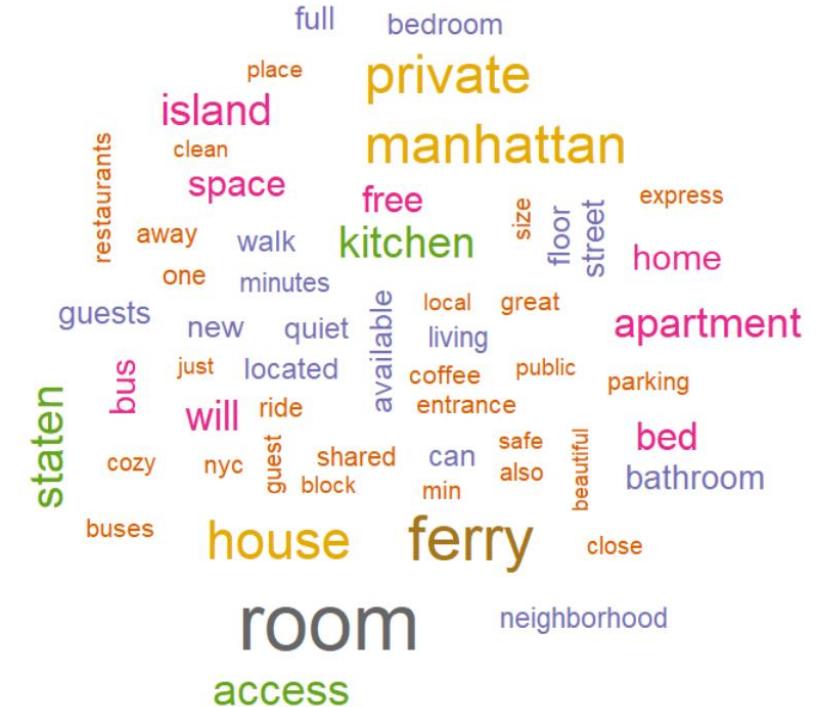
Solution

Staten Island Optimal Daily Price



\$63 / day

Based on this price, a Staten Island host would have to host a guest for **24 days** to make their monthly rent.



Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution



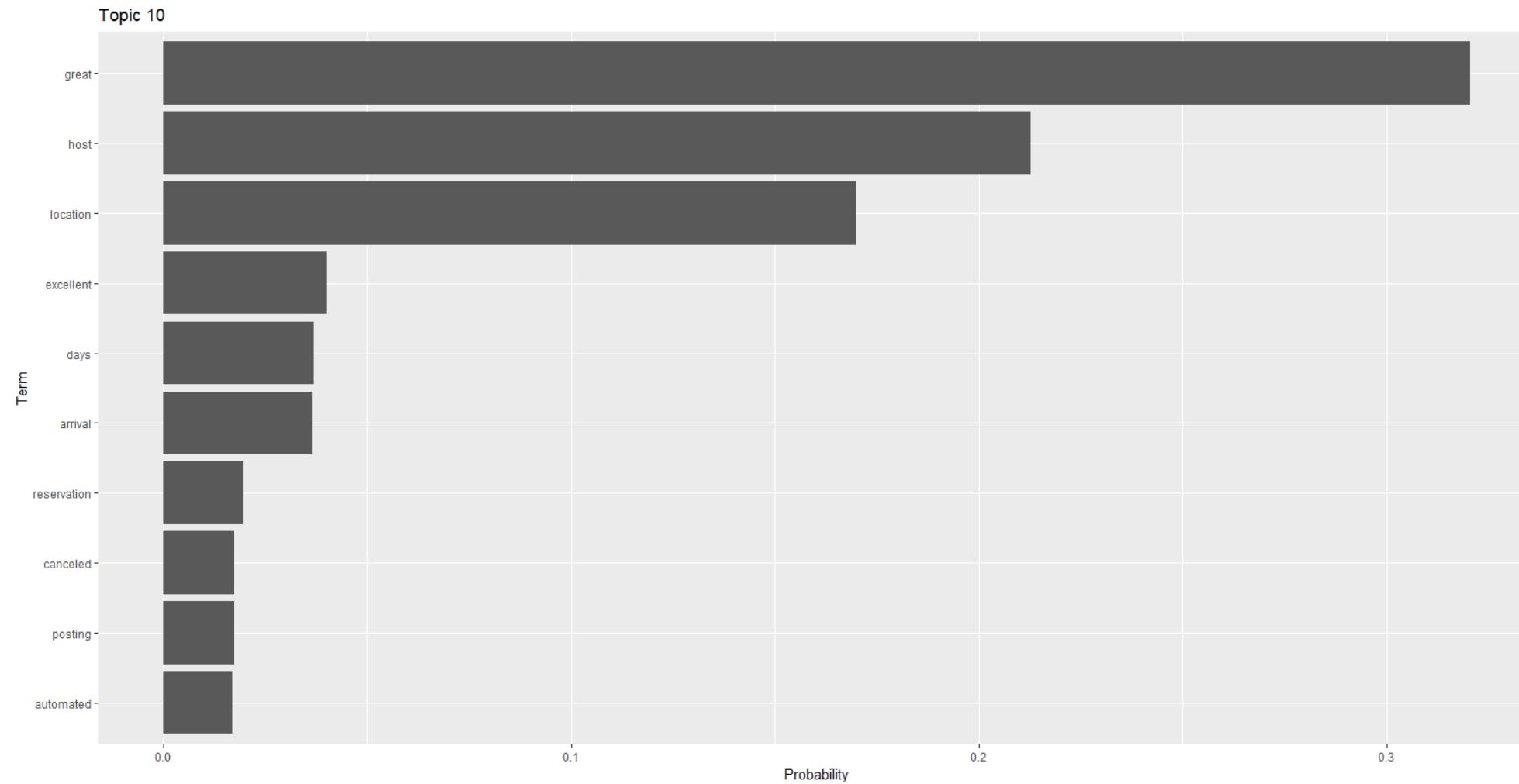
Thank you!

Any questions?



Appendix

Text Mining: Ease of Booking



Problem Overview

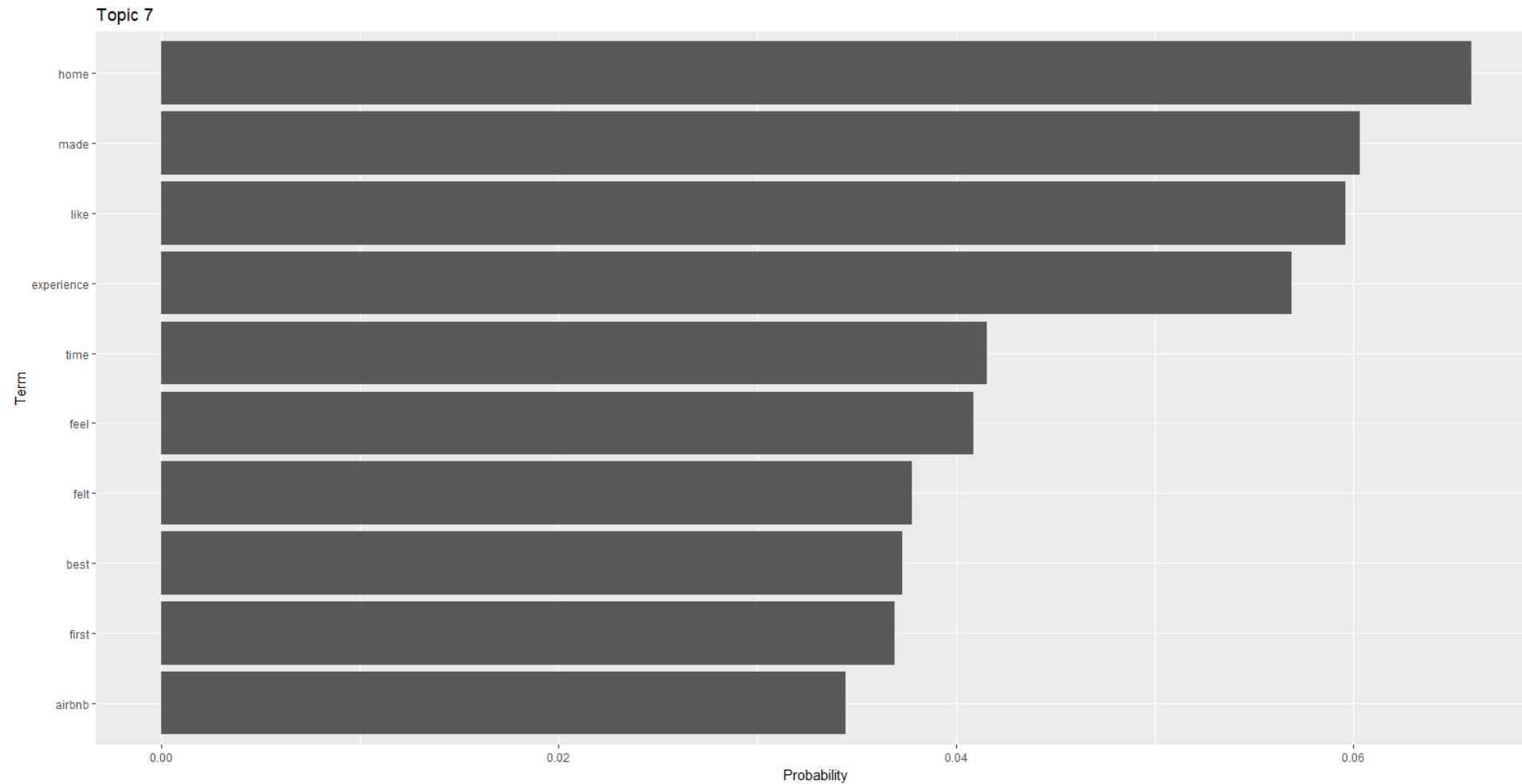
Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Text Mining: Satisfaction



Problem Overview

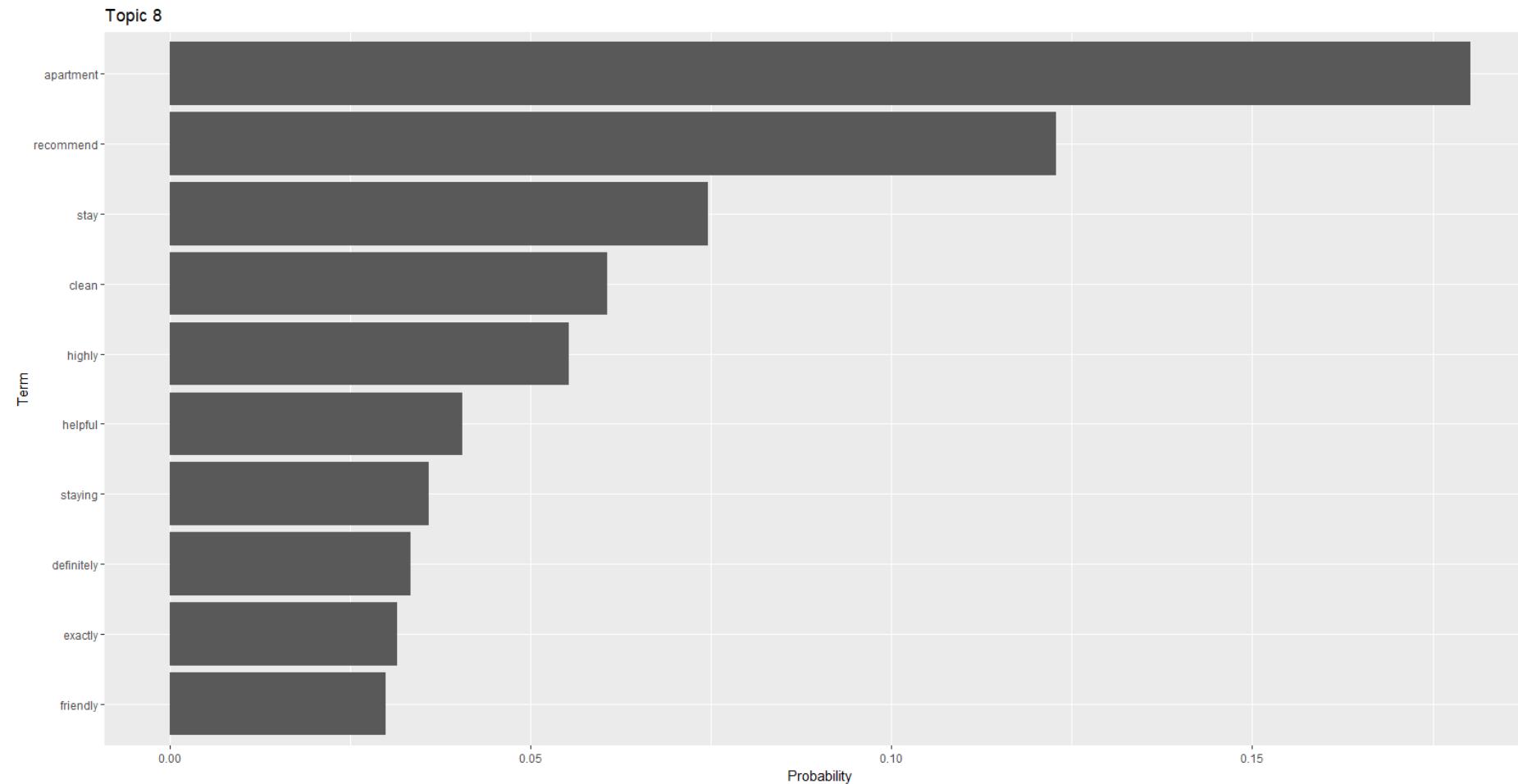
Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Text Mining: Would Recommend



Problem Overview

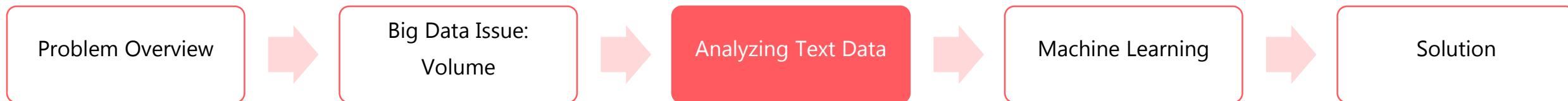
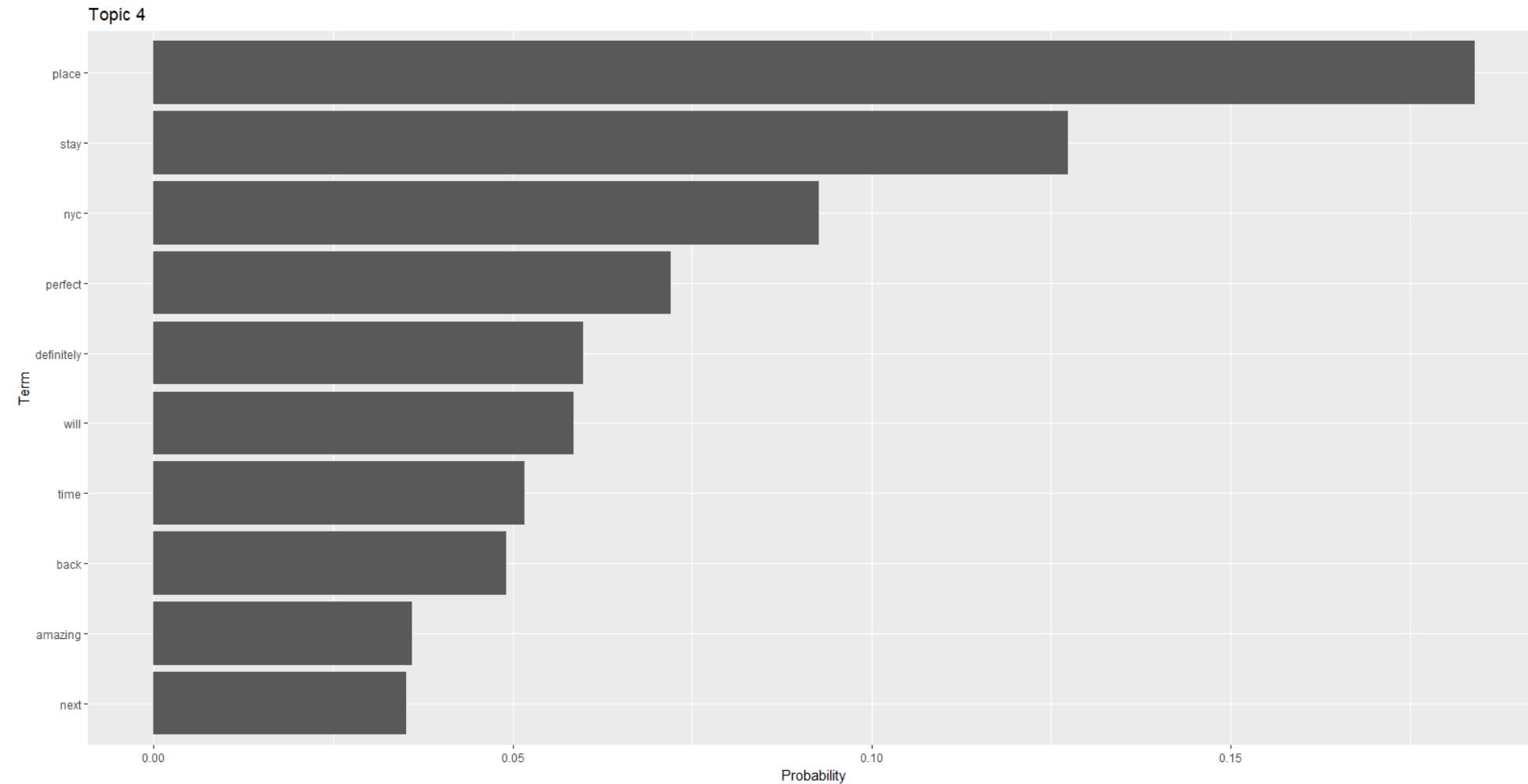
Big Data Issue:
Volume

Analyzing Text Data

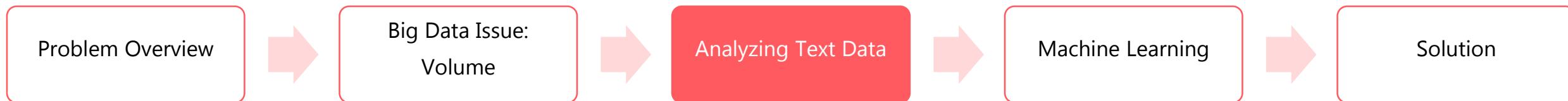
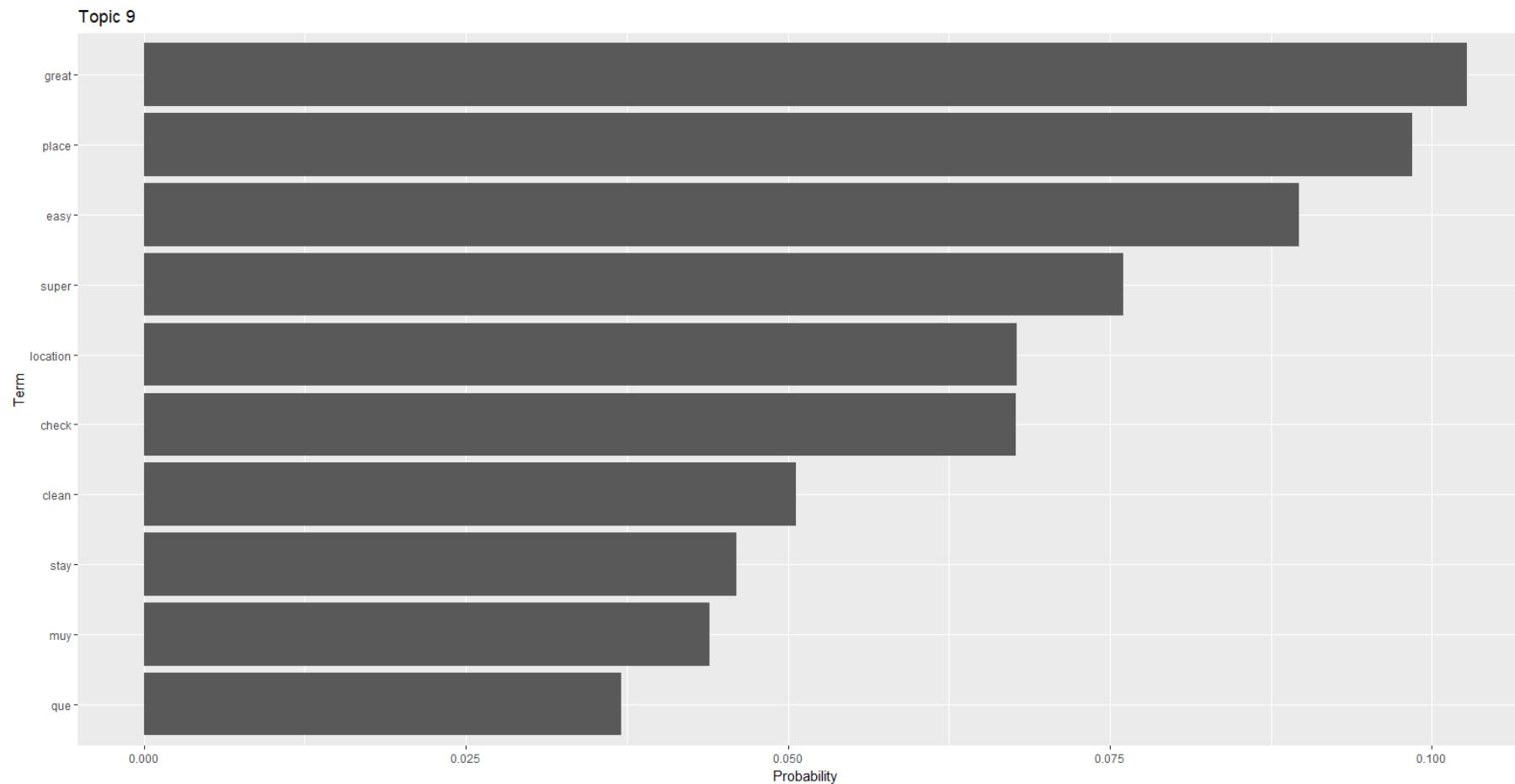
Machine Learning

Solution

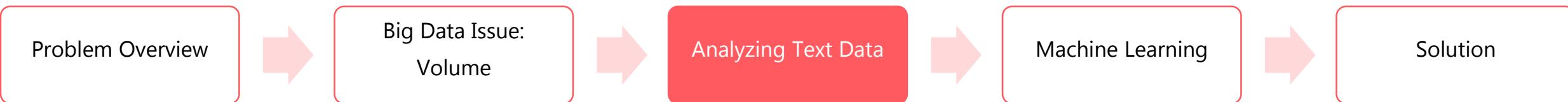
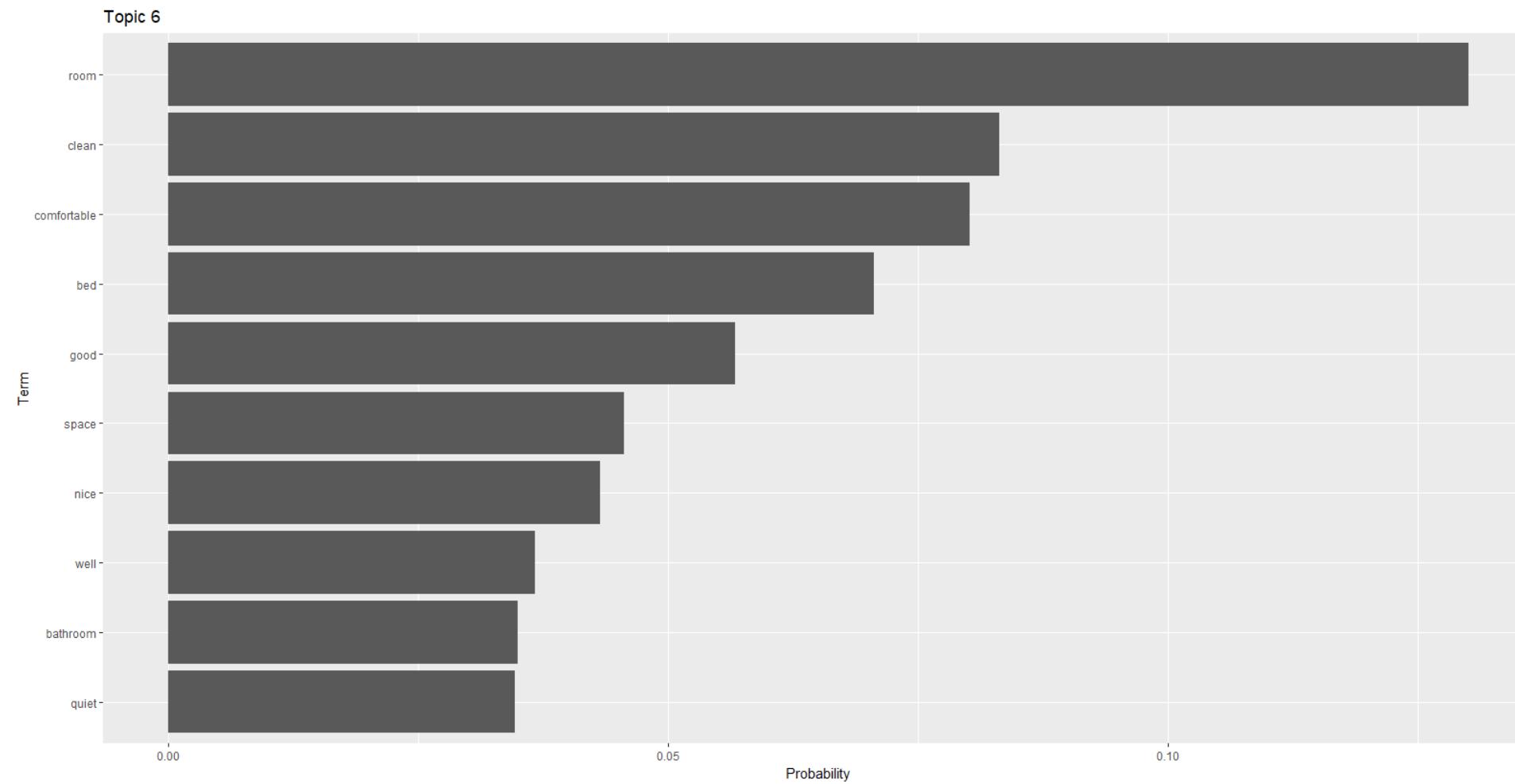
Text Mining: Satisfaction II



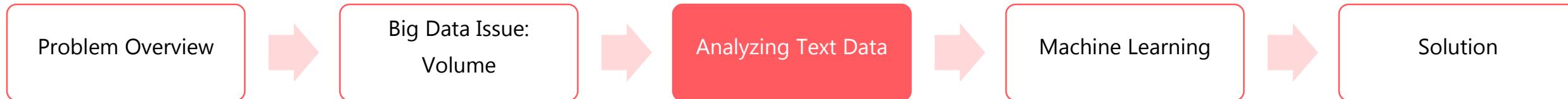
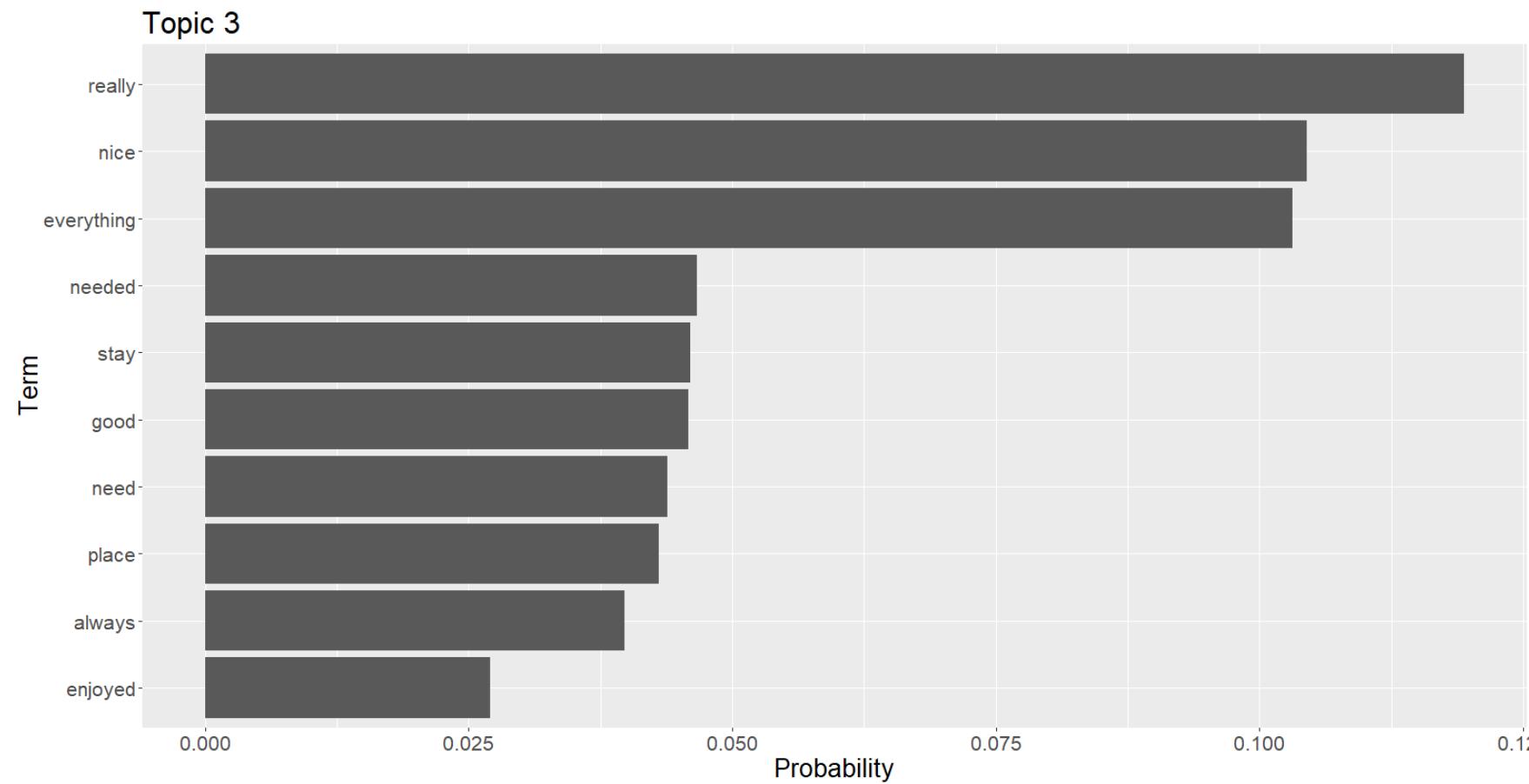
Text Mining: Overall Experience



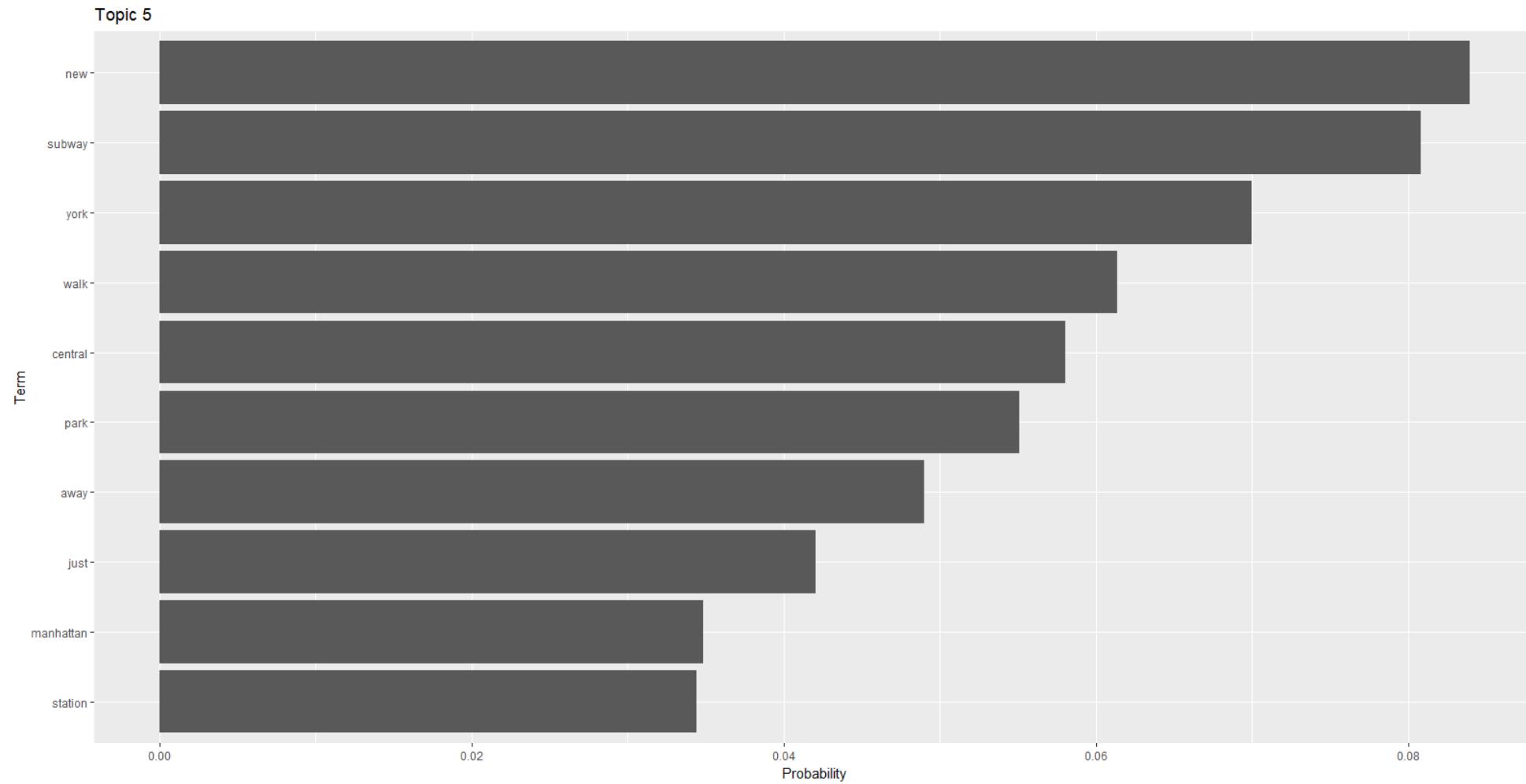
Text Mining: Quality



Text Mining: Enjoyment



Text Mining: Location Perks



Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Machine Learning

Cluster One, Brooklyn

Sheet1 - Model Comparison - JMP Pro						
Model Comparison Validation=Training						
Predictors						
Measures of Fit for D_brooklyn						
Predictor	Creator	.2	.4	.6	.8	RSquare
D_brooklyn Predictor	Bootstrap Forest					-0.010
D_brooklyn Predictor 2	Boosted Tree					0.0000
Predicted D_brooklyn	Neural					1.0000
D_brooklyn Prediction Formula	Fit Generalized Adaptive Elastic Net					0.7925
Model Comparison Validation=Validation						
Predictors						
Measures of Fit for D_brooklyn						
Predictor	Creator	.2	.4	.6	.8	RSquare
D_brooklyn Predictor	Bootstrap Forest					-0.000
D_brooklyn Predictor 2	Boosted Tree					-0.017
Predicted D_brooklyn	Neural					0.9622
D_brooklyn Prediction Formula	Fit Generalized Adaptive Elastic Net					0.5924
RASE	AAE	Freq	RASE	AAE	Freq	RASE
0.4073	0.2803	8	0.4053	0.2905	8	0.1846
0.0000	0.0000	8	0.0626	0.0594	4	0.2055
						0.1850

Neural Network

Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Machine Learning

Cluster One, Bronx

Sheet1 - Model Comparison - JMP Pro

Model Comparison Validation=Training

Predictors

Measures of Fit for D_bronx

Predictor	Creator	.2	.4	.6	.8	RSquare	RASE	AAE	Freq
D_bronx Predictor	Bootstrap Forest					-0.064	0.5390	0.4019	8
D_bronx Predictor 2	Boosted Tree					0.0000	0.5226	0.3981	8
Predicted D_bronx	Neural					1.0000	0.0018	0.0015	8
D_bronx Prediction Formula	Fit Generalized Adaptive Elastic Net					0.7718	0.2496	0.1890	8

Model Comparison Validation=Validation

Predictors

Measures of Fit for D_bronx

Predictor	Creator	.2	.4	.6	.8	RSquare	RASE	AAE	Freq
D_bronx Predictor	Bootstrap Forest					-0.001	0.3176	0.2386	4
D_bronx Predictor 2	Boosted Tree					-0.200	0.3478	0.3047	4
Predicted D_bronx	Neural					0.1382	0.2947	0.2529	4
D_bronx Prediction Formula	Fit Generalized Adaptive Elastic Net					0.7188	0.1683	0.1490	4

Elastic Net

Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Machine Learning

Cluster One, Queens

Sheet1 - Model Comparison - JMP Pro

Model Comparison Validation=Training

Predictors

Measures of Fit for D_queens

Predictor	Creator	.2	.4	.6	.8	RSquare	RASE	AAE	Freq
D_queens Predictor	Bootstrap Forest					-0.014	0.2855	0.2140	8
D_queens Predictor 2	Boosted Tree					0.0000	0.2836	0.2140	8
Predicted D_queens	Neural					0.2751	0.2414	0.2159	8
D_queens Prediction Formula	Fit Generalized Adaptive Elastic Net					0.7565	0.1399	0.1039	8

Model Comparison Validation=Validation

Predictors

Measures of Fit for D_queens

Predictor	Creator	.2	.4	.6	.8	RSquare	RASE	AAE	Freq
D_queens Predictor	Bootstrap Forest					-0.000	0.2479	0.2091	4
D_queens Predictor 2	Boosted Tree					-0.017	0.2500	0.2257	4
Predicted D_queens	Neural					0.8544	0.0946	0.0730	4
D_queens Prediction Formula	Fit Generalized Adaptive Elastic Net					0.6258	0.1517	0.1173	4

Neural Network

Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution

Machine Learning

Cluster One, Staten Island

Sheet1 - Model Comparison - JMP Pro

Model Comparison Validation=Training

Predictors

Measures of Fit for D_statenisland

Predictor	Creator	.2	.4	.6	.8	RSquare	RASE	AAE	Freq
D_statenisland Predictor	Bootstrap Forest					-0.180	1.5567	1.2299	8
D_statenisland Predictor 2	Boosted Tree					0.0000	1.4328	1.1686	8
Predicted D_statenisland	Neural					0.6677	0.8259	0.6441	8
D_statenisland Prediction Formula	Fit Generalized Adaptive Elastic Net					0.7760	0.6782	0.5218	8

Model Comparison Validation=Validation

Predictors

Measures of Fit for D_statenisland

Predictor	Creator	.2	.4	.6	.8	RSquare	RASE	AAE	Freq
D_statenisland Predictor	Bootstrap Forest					-0.000	1.6603	1.4152	4
D_statenisland Predictor 2	Boosted Tree					-0.148	1.7789	1.7195	4
Predicted D_statenisland	Neural					0.9322	0.4321	0.3823	4
D_statenisland Prediction Formula	Fit Generalized Adaptive Elastic Net					0.5948	1.0567	1.0341	4

Neural Network

Problem Overview

Big Data Issue:
Volume

Analyzing Text Data

Machine Learning

Solution