

CIS530 HW5 Submission Example

2 SimLex-999 Dataset Revisited (10 points)

Note **5 Extra points** will be awarded for creativity and a more thorough qualitative analysis.)

2.1 What is the least similar 2 pairs of words based on human judgement scores and vector similarity? Do the pairs match? [3 points]

Results

- least similar 2 pairs based on human judgement scores
 1. new
 2. ancient
- least similar 2 pairs based on vector similarity
 1. house
 2. key
- Do the pairs match?
 - No

2.2 What is the most similar 2 pairs of words based on human judgement scores and vector similarity? Do the pairs match? [3 points]

Results

- most similar 2 pairs based on human judgement scores
 3. vanish
 4. disappear
- most similar 2 pairs based on vector similarity
 3. south
 4. north
- Do the pairs match?
 - no

2.3 Provide correlation scores and p values for the following models:

How do those correlation value compare to each other? [4 points]

Results

- glove.6B.50d.magnitude
 - Correlation = 0.18100126067449063,
 - P Value = 1.2242211264976856e-17
- glove.6B.100d.magnitude
 - Correlation = 0.20506409092608713,
 - P Value = 3.41228663395174e-22
- glove.6B.200d.magnitude
 - Correlation = 0.23670323199262908,
 - P Value = 4.9936324557833286e-29
- glove.6B.300d.magnitude
 - Correlation = 0.25894302181101986,
 - P Value = 2.080389068003349e-34
- glove.840B.300d.magnitude
 - Correlation = 0.2860664813618063,
 - P Value = 1.293335613361039e-41
- How do those correlation value compare to each other?
 - When the dimensions of the magnitude increases, the p-value goes down significantly and the correlation increases. Which means the similarity is getting closer to human's judgement.

3.1 Cluster Randomly (1 points)

3.1 Run clustering on dev data, report the f_scores from the dev data [1 point]

- Dev data

Target	k	Paired F-Score
smell.v	4	0.2857
image.n	9	0.1204
express.v	7	0.1562
talk.v	6	0.2428
play.v	34	0.0308
miss.v	8	0.1972
produce.v	7	0.1968
write.v	9	0.1507
provide.v	7	0.2204

party.n	5	0.2017
bank.n	9	0.0741
plan.n	3	0.3693
shelter.n	5	0.2500
difference.n	5	0.2564
eat.v	6	0.2186
mean.v	6	0.1772
treat.v	8	0.1610
use.v	6	0.2407
suspend.v	6	0.1304
judgment.n	7	0.1480
organization.n	7	0.1734
interest.n	5	0.1704
paper.n	7	0.2044
operate.v	7	0.1237
receive.v	13	0.0903
watch.v	5	0.2041
rule.v	7	0.1267
simple.a	5	0.2000
atmosphere.n	6	0.1812
expect.v	6	0.2105
different.a	1	1.0000
begin.v	8	0.0897
note.v	3	0.3684
win.v	4	0.2756
source.n	9	0.1260
performance.n	5	0.2347
wash.v	13	0.0849
hear.v	5	0.1835
climb.v	6	0.1963
degree.n	7	0.1985

=> Average Paired F-Score: 0.1592

3.2 Cluster with Sparse Representations (6 points)

3.2.1 Run clustering on dev data, report the f_scores from the dev data [1 point]

- Dev data

Target	k	Paired F-Score
degree.n	7	0.3273
atmosphere.n	6	0.3058
play.v	34	0.0845
smell.v	4	0.2947
talk.v	6	0.3223
plan.n	3	0.5772
party.n	5	0.2322
express.v	7	0.2340

win.v	4	0.4051
paper.n	7	0.4899
rule.v	7	0.2174
produce.v	7	0.2159
note.v	3	0.5333
performance.n	5	0.2684
bank.n	9	0.3373
miss.v	8	0.2182
operate.v	7	0.2283
write.v	9	0.1660
source.n	9	0.2337
mean.v	6	0.3804
expect.v	6	0.3661
difference.n	5	0.3835
image.n	9	0.1634
shelter.n	5	0.2919
use.v	6	0.3369
wash.v	13	0.1467
interest.n	5	0.2333
organization.n	7	0.2422
treat.v	8	0.2145
begin.v	8	0.2775
climb.v	6	0.3608
different.a	1	1.0000
provide.v	7	0.3276
suspend.v	6	0.2222
simple.a	5	0.1429
hear.v	5	0.3368
eat.v	6	0.3096
judgment.n	7	0.2086
receive.v	13	0.2013
watch.v	5	0.2857

=> Average Paired F-Score: 0.2526

3.2.2 Provide a brief description of your method in the report, making sure to describe the vector space model you chose, the clustering algorithm you used, and the results of any preliminary experiments you might have run on the dev set. [5 points]

- Brief Description of your method
 - [Description / Background]

I have used a different clustering algorithm which is Agglomerative Clustering. The method is the same as Cluster with Sparse Representation. The cooccurrence matrix is the same, but instead of using K-Means algorithm, I switched to use Agglomerative Clustering. Agglomerative Clustering starts the cluster from each point and K-Means assign k random centroids at the beginning. That's the difference. Agglomerative clustering is slower than K-means but it's more stable.

- [the vector space model you chose]

The vector space I used is the same as Sparse Representation which is “coocvec-500mostfreq-window-3.filter.magnitude”

- [the clustering algorithm you used]

I have used a different clustering algorithm which is Agglomerative Clustering, which is a type of hierarchical clustering algorithm. It supports 4 different linkage methods

Linkage	Explanation of distance between clusters
ward	Minimize total within-cluster variance
complete	Max of distances between points in each cluster
average	Average of distances between all pairs
single	Min of distances between points

- [the results of any preliminary experiments]

Linkage	F-score	
ward	0.2366	
complete	0.2832	
average	0.3430	
single	0.3599	Using “single” achieved highest f-score

Using “single” achieves best performance

3.3 Cluster with Dense Representations (8 points)

3.3.1 Run clustering on dev data, report the f_scores from the dev data [1 point]

- Dev data

Target	k	Paired F-Score
write.v	9	0.1633
treat.v	8	0.2658
note.v	3	0.6190
paper.n	7	0.3784
talk.v	6	0.2812
wash.v	13	0.1789
source.n	9	0.3111
degree.n	7	0.3531
operate.v	7	0.2491
bank.n	9	0.7333
smell.v	4	0.3500
use.v	6	0.2926
provide.v	7	0.2495
mean.v	6	0.3145
hear.v	5	0.2712
suspend.v	6	0.3860
party.n	5	0.4209
watch.v	5	0.4964
climb.v	6	0.2798
image.n	9	0.3038
performance.n	5	0.3798
organization.n	7	0.2224
eat.v	6	0.3386

expect.v	6	0.3512
different.a	1	1.0000
shelter.n	5	0.3805
simple.a	5	0.2857
play.v	34	0.1218
receive.v	13	0.1628
interest.n	5	0.3395
begin.v	8	0.3580
atmosphere.n	6	0.3123
difference.n	5	0.4235
plan.n	3	0.4451
rule.v	7	0.2846
miss.v	8	0.2703
produce.v	7	0.2705
express.v	7	0.2698
win.v	4	0.4737
judgment.n	7	0.3219

=> Average Paired F-Score: 0.2984

3.3.2 Provide a brief description of your method in the report that includes the vectors you used, and any experimental results you have from running your model on the dev set. [5 points]

- Brief Description of your method
 - [Description / Background]

Used K-means clustering algorithm to split the words. When the word doesn't exist in the vectors, random select an embedding.

- [the vector space model you chose]

Used GoogleNews-vectors-negative300.magnitude

- [the clustering algorithm you used]

Used K-means clustering algorithm

- [the results of any preliminary experiments]

The f-score is 0.2984. It's more than the f-score from the sparse embedding experiment.

3.3.3 In addition, for Task 3.2 and 3.3, do an analysis of different errors made by each system – i.e. look at instances that the word-context matrix representation gets wrong and dense gets right, and vice versa, and see if there are any interesting patterns. There is no right for this. [2 points]

- To compare, make a dataframe Target word | F score_ sparse | F score_ dense | Difference

K-Means on both sparse and dense embeddings

Target Word	k	F Score Sparse	F Score Dense	Difference
write.v	9	0.2205	0.1633	0.0572

treat.v	8	0.2505	0.2658	-0.0153
note.v	3	0.5333	0.619	-0.0857
paper.n	7	0.3783	0.3784	-1E-04
talk.v	6	0.3507	0.2812	0.0695
wash.v	13	0.12	0.1789	-0.0589
source.n	9	0.1743	0.3111	-0.1368
degree.n	7	0.3279	0.3531	-0.0252
operate.v	7	0.231	0.2491	-0.0181
bank.n	9	0.3373	0.7333	-0.396
smell.v	4	0.2947	0.35	-0.0553
use.v	6	0.4184	0.2926	0.1258
provide.v	7	0.3312	0.2495	0.0817
mean.v	6	0.3804	0.3145	0.0659
hear.v	5	0.3368	0.2712	0.0656
suspend.v	6	0.1905	0.386	-0.1955
party.n	5	0.2762	0.4209	-0.1447
watch.v	5	0.3361	0.4964	-0.1603
climb.v	6	0.35	0.2798	0.0702
image.n	9	0.2062	0.3038	-0.0976
performance.n	5	0.2544	0.3798	-0.1254
organization.n	7	0.2523	0.2224	0.0299
eat.v	6	0.2574	0.3386	-0.0812
expect.v	6	0.3654	0.3512	0.0142
different.a	1	1	1	0
shelter.n	5	0.3129	0.3805	-0.0676
simple.a	5	0.1538	0.2857	-0.1319

play.v	34	0.0903	0.1218	-0.0315
receive.v	13	0.2018	0.1628	0.039
interest.n	5	0.2145	0.3395	-0.125
begin.v	8	0.2022	0.358	-0.1558
atmosphere.n	6	0.3489	0.3123	0.0366
difference.n	5	0.3088	0.4235	-0.1147
plan.n	3	0.572	0.4451	0.1269
rule.v	7	0.2077	0.2846	-0.0769
miss.v	8	0.25	0.2703	-0.0203
produce.v	7	0.2721	0.2705	0.0016
express.v	7	0.3106	0.2698	0.0408
win.v	4	0.4497	0.4737	-0.024
judgment.n	7	0.2245	0.3219	-0.0974

- Now sort based on Difference to identify words where one works better than other
 - [Target word | F score _ sparse | F score _ dense | Difference table]
 - [Error analysis]

When the word has K as 1, there's no difference.

When the word has a very large K, for example “play” has k=34, it’s f-score is very low.

Overall, the F-Score Dense performs better than F-Score Sparse.

Target Word	k	F Score Sparse	F Score Dense	Difference =F Score Dense – F Score Sparse	Comments
bank.n	9	0.3373	0.7333	0.396	This work has highest f-score under densed embedding
suspend.v	6	0.1905	0.386	0.1955	
watch.v	5	0.3361	0.4964	0.1603	
begin.v	8	0.2022	0.358	0.1558	
party.n	5	0.2762	0.4209	0.1447	

source.n	9	0.1743	0.3111	0.1368	
simple.a	5	0.1538	0.2857	0.1319	
performance.n	5	0.2544	0.3798	0.1254	
interest.n	5	0.2145	0.3395	0.125	
difference.n	5	0.3088	0.4235	0.1147	
image.n	9	0.2062	0.3038	0.0976	
judgment.n	7	0.2245	0.3219	0.0974	
note.v	3	0.5333	0.619	0.0857	
eat.v	6	0.2574	0.3386	0.0812	
rule.v	7	0.2077	0.2846	0.0769	
shelter.n	5	0.3129	0.3805	0.0676	
wash.v	13	0.12	0.1789	0.0589	
smell.v	4	0.2947	0.35	0.0553	
play.v	34	0.0903	0.1218	0.0315	
degree.n	7	0.3279	0.3531	0.0252	
win.v	4	0.4497	0.4737	0.024	
miss.v	8	0.25	0.2703	0.0203	
operate.v	7	0.231	0.2491	0.0181	
treat.v	8	0.2505	0.2658	0.0153	
paper.n	7	0.3783	0.3784	1E-04	
different.a	1	1	1	0	
produce.v	7	0.2721	0.2705	-0.0016	
expect.v	6	0.3654	0.3512	-0.0142	
organization.n	7	0.2523	0.2224	-0.0299	
atmosphere.n	6	0.3489	0.3123	-0.0366	
receive.v	13	0.2018	0.1628	-0.039	
express.v	7	0.3106	0.2698	-0.0408	

write.v	9	0.2205	0.1633	-0.0572	
hear.v	5	0.3368	0.2712	-0.0656	
mean.v	6	0.3804	0.3145	-0.0659	
talk.v	6	0.3507	0.2812	-0.0695	
climb.v	6	0.35	0.2798	-0.0702	
provide.v	7	0.3312	0.2495	-0.0817	
use.v	6	0.4184	0.2926	-0.1258	
plan.n	3	0.572	0.4451	-0.1269	

3.4 Cluster without K (6 points)

3.4.1 Run clustering on dev data, report the f_scores from the dev data [1 point]

- Dev data

Target	k	Paired F-Score
write.v	9	0.2859
treat.v	8	0.3296
note.v	3	0.1875
paper.n	7	0.5691
talk.v	6	0.4923
wash.v	13	0.2747
source.n	9	0.2827
degree.n	7	0.4367
operate.v	7	0.3035
bank.n	9	0.4231
smell.v	4	0.4340
use.v	6	0.6310
provide.v	7	0.5948
mean.v	6	0.3425
hear.v	5	0.2222
suspend.v	6	0.4878
party.n	5	0.4551
watch.v	5	0.2927
climb.v	6	0.2400
image.n	9	0.3075
performance.n	5	0.3264
organization.n	7	0.3391
eat.v	6	0.3667
expect.v	6	0.4704
different.a	1	1.0000
shelter.n	5	0.2904

simple.a	5	0.3333	
play.v	34	0.1902	
receive.v	13	0.1961	
interest.n	5	0.3874	
begin.v	8	0.4072	
atmosphere.n	6	0.2947	
difference.n	5	0.4958	
plan.n	3	0.2652	
rule.v	7	0.3906	
miss.v	8	0.3684	
produce.v	7	0.3903	
express.v	7	0.4175	
win.v	4	0.2831	
judgment.n	7	0.3071	

=> Average Paired F-Score: 0.3446

3.4.2 Provide a brief description of your method in the report that includes the vectors you used, and any experimental results you have from running your model on the dev set. [5 points]

- Brief Description of your method
 - [Description / Background]

For each target word's paraphrases, I run the K-Means clustering with different k from 2 to the number of paraphrases. When the paraphrases contain less than 2 words, no need to do clustering.

- [the vector space model you chose]

The vector I used is GoogleNews-vectors-negative300.magnitude

- [the clustering algorithm you used]

The clustering method used is K-Means

- [the results of any preliminary experiments]

The f-score is 0.3446. It's highest comparing to sparse representation and dense representation. Which means that each word might have a best k which can device the paraphrases.

Leaderboards [2 points + 3 bonus]

- No writing is required here, but be sure to submit a valid result (not -1) to Gradescope. [2 points]