

```

1 ### md
2 4.1 (8 points) Evaluate the word vectors (EVALWS)
   corresponding to the three ways of computing vectors
   (counts, IDF, and PMI), three values of w (1, 3, and
   6), and two context vocabularies
3 (vocab-15kws.txt and vocab-5k.txt). For all cases,
   use vocab-15kws.txt for V . Report the
4 results (there should be 36 correlations in all) and
   describe your findings. What happens as window
5 size changes for different methods of creating word
   vectors? What happens when context vocabulary
6 changes? Why do you think you observe the trends you
   see? Do you see the same trends for MEN and
7 SimLex or do they differ?
8
9 For Count method
10 ###
11 from collections import defaultdict
12 import math
13 from scipy.stats import rankdata
14
15 # Function to find indices for middle word.
16 def find_indexes(arr, words_set):
17     return [(i, w) for i, w in enumerate(arr) if w in
        words_set]
18
19 # Function to find count for word vectors.
20 def word_vector_count(w, vocab, vocab_context, lines
    ):
21     V = set(vocab)
22     Vc = set(vocab_context)
23
24     my_dict = defaultdict(int)
25
26     for line in lines:
27         words = line.split()
28         length = len(words)
29         word_indexes = find_indexes(words, V)
30
31         for idx, word in word_indexes:
32             ind_l = max(0, idx - w)

```

```

33         ind_h = min(length, idx + w + 1)
34         context_window = words[ind_l:idx] + words
           [idx + 1:ind_h]
35
36         for context_word in context_window:
37             if context_word in Vc:
38                 my_dict[(word, context_word)] +=
           1
39
40     return my_dict
41 ###
42 from collections import defaultdict
43 import math
44 from scipy.stats import rankdata
45 # Function to deduce word pairs and their scores from
   given dataset.
46 def make_word_pairs(path):
47     word_pairs = []
48
49     word_pair_scores = {}
50     with open(path) as f3:
51         next(f3)
52         for line in f3:
53             word1, word2, score = line.split()
54             pair = (word1, word2)
55             word_pairs.append(pair)
56
57             word_pair_scores[pair] = float(score)
58
59     return word_pairs, word_pair_scores
60 ###
61 from collections import defaultdict
62 import math
63 from scipy.stats import rankdata
64
65 # Function to calculate cosine similarity.
66 def cosine_similarity(vec1, vec2):
67     dot_product = sum(vec1[word] * vec2.get(word, 0)
           for word in vec1)
68     norm_vec1 = math.sqrt(sum(value ** 2 for value in
           vec1.values()))

```

```

69     norm_vec2 = math.sqrt(sum(value ** 2 for value
    in vec2.values()))
70
71     if norm_vec1 == 0 or norm_vec2 == 0:
72         return 0.0
73
74     return dot_product / (norm_vec1 * norm_vec2)
75 ###
76 from collections import defaultdict
77 import math
78 from scipy.stats import rankdata
79
80
81 # Function to calculate Spearman correlation.
82 def spearman_rank_correlation(predicted_similarities
    , actual_scores):
83     predicted = [predicted_similarities[pair] for
    pair in actual_scores if pair in
    predicted_similarities]
84     actual = [actual_scores[pair] for pair in
    actual_scores if pair in predicted_similarities]
85
86     if not predicted:
87         return None
88
89     predicted_ranks = rankdata(predicted)
90     actual_ranks = rankdata(actual)
91
92     n = len(predicted_ranks)
93     rank_differences_squared = [(pred_r - act_r) **
    2 for pred_r, act_r in zip(predicted_ranks,
    actual_ranks)]
94
95     return 1 - (6 * sum(rank_differences_squared
    )) / (n * (n**2 - 1))
96 ###
97 from collections import defaultdict
98 import math
99 from scipy.stats import rankdata
100
101 def load_file(path):

```

```

102     with open(path, "r") as file:
103         return file.read().splitlines()
104
105 def count_lines(Vc, lines):
106     total_lines = 0
107     context_word_line_counts = defaultdict(int)
108
109     for line in lines:
110         total_lines += 1
111         words = set(line.split())
112         for word in words:
113             if word in Vc:
114                 context_word_line_counts[word] += 1
115
116     return total_lines, context_word_line_counts
117
118 # Defining context window size, different paths,
119 different context word vocab and different methods.
120 context_window_sizes = [1, 3, 6]
121 word_pairs_paths = ["/men.txt", "/simlex-999.txt"]
122 context_vocab_paths = ["/vocab-15kws.txt", "/vocab-5k.txt"]
123 methods = {"count", "tf-idf", "pmi"}
124
125 wiki_lines = load_file("/wiki-1percent.txt")
126
127 for path in word_pairs_paths:
128     word_pairs, word_pair_scores = make_word_pairs(
129         path)
130
131     for context_vocab in context_vocab_paths:
132         vocab_context = load_file(context_vocab)
133
134         for w in context_window_sizes:
135             word_vectors = defaultdict(dict)
136             vocab = load_file("/vocab-15kws.txt")
137             my_dict = word_vector_count(w, vocab,
138                 vocab_context, wiki_lines)
139             for (word, context_word), count in
140                 my_dict.items():
141                 word_vectors[word][context_word] =

```

```

137 count
138
139
140         cosine_similarities = {pair:
    cosine_similarity(word_vectors.get(pair[0], {}),
    word_vectors.get(pair[1], {}))
141                             for pair in
    word_pairs}
142
143         spearman_rho = spearman_rank_correlation
    (cosine_similarities, word_pair_scores)
144         print(f"Spearman's  $\rho$  for path = {path},
    context_vocab = {context_vocab}, w = {w}, method =
    count: {spearman_rho}")
145
146 ### md
147 4.1) For TF-IDF Method
148 ###
149 for path in word_pairs_paths:
150     word_pairs, word_pair_scores = make_word_pairs(
    path)
151
152     for context_vocab in context_vocab_paths:
153         vocab_context = load_file(context_vocab)
154
155         for w in context_window_sizes:
156             word_vectors = defaultdict(dict)
157             vocab = load_file("/vocab-15kws.txt")
158             my_dict = word_vector_count(w, vocab,
    vocab_context, wiki_lines)
159             total_lines, count_context_word_lines =
    count_lines(set(vocab_context), wiki_lines)
160             for (word, context_word), count in
    my_dict.items():
161                 idf = (total_lines /
    count_context_word_lines[context_word]) if
    count_context_word_lines[context_word] > 0 else 0
162                 word_vectors[word][context_word] =
    count * idf
163
164

```

```

165         cosine_similarities = {pair:
    cosine_similarity(word_vectors.get(pair[0], {}),
    word_vectors.get(pair[1], {}))
166         for pair in
    word_pairs}
167
168         spearman_rho = spearman_rank_correlation
    (cosine_similarities, word_pair_scores)
169         print(f"Spearman's  $\rho$  for path = {path},
    context_vocab = {context_vocab}, w = {w}, method =
    tf-idf: {spearman_rho}")
170 ### md
171 4.1) For PMI method
172 ###
173 for path in word_pairs_paths:
174     word_pairs, word_pair_scores = make_word_pairs(
    path)
175
176     for context_vocab in context_vocab_paths:
177         vocab_context = load_file(context_vocab)
178
179         for w in context_window_sizes:
180             word_vectors = defaultdict(dict)
181             vocab = load_file("/vocab-15kws.txt")
182             my_dict = word_vector_count(w, vocab,
    vocab_context, wiki_lines)
183             total_count_N = sum(my_dict.values())
184             joint_probabilities = {pair: count /
    total_count_N for pair, count in my_dict.items()}
185
186             partial_prob_x = {word: sum(
    joint_probabilities.get((word, cw), 0) for cw in set
    (vocab_context)) for word in set(vocab)}
187             partial_prob_y = {cw: sum(
    joint_probabilities.get((word, cw), 0) for word in
    set(vocab)) for cw in set(vocab_context)}
188
189             for (word, context_word), joint_prob in
    joint_probabilities.items():
190                 pmi = math.log2(joint_prob / (
    partial_prob_x[word] * partial_prob_y[context_word

```

```
190 ))) if joint_prob > 0 else 0
191         word_vectors[word][context_word] = pmi
192
193
194         cosine_similarities = {pair:
    cosine_similarity(word_vectors.get(pair[0], {}),
    word_vectors.get(pair[1], {}))
195         for pair in
    word_pairs}
196
197         spearman_rho = spearman_rank_correlation
    (cosine_similarities, word_pair_scores)
198         print(f"Spearman's ρ for path = {path},
    context_vocab = {context_vocab}, w = {w}, method =
    pmi: {spearman_rho}")
```