```python
 1  import string
 2  import os
 3  import json
 4
 5
 6  # ------------------------------------------------ load file
    ---------------------------------------------- #
 7
 8  def load_vocab(vocab_path):
 9      vocabs = set()
10      with open(vocab_path) as vocab_file:
11          for word in vocab_file:
12              vocabs.add(word.strip())
13      return vocabs
14
15
16  def save_file(file_content, file_path):
17      with open(file_path, 'w') as file:
18          file.write(file_content)
19
20
21  # --------------------------------------------- preprocess files
    --------------------------------------------- #
22
23  def is_punc(char):
24      return char in string.punctuation
25
26
27  def lowercase_sentence(comment: str) -> str:
28      return comment.lower()
29
30
31  def separate_punctuation(comment: str) -> str:
32      res = []
33      for word in comment.split():
34          start = 0 # start index of current separation
35          for i, c in enumerate(word):
36              if is_punc(c):
37                  if start < i:
38                      res.append(word[start: i])
39                  res.append(c)
40                  start = i + 1
41          if start < len(word):
42              res.append(word[start:])
43      return ' '.join(res)
44
45
46  def contains_strong_pos_word(words:list) -> bool:
47      positive_words = {'excellent', 'amazing', 'great', 'fantastic', '
    outstanding', 'terrific', 'phenomenal', 'superb',
48                        'brilliant', 'impressive'}
49      for word in words:
50          if word in positive_words:
51              return True
52      return False
53
54
55  def contains_strong_neg_word(words: list) -> bool:
56      negative_words = {'disappointing', 'terrible', 'awful', 'horrible'
    , 'dreadful', 'abysmal', 'appalling', 'atrocious',
57                        'repulsive', 'disgusting'}
```

```python
 58        for word in words:
 59            if word in negative_words:
 60                return True
 61        return False
 62
 63
 64 def preprocess_comment(comment: str) -> str:
 65     return separate_punctuation(lowercase_sentence(comment))
 66
 67
 68 def preprocess_file(file_path: str) -> str:
 69     with open(file_path) as file:
 70         comment = ' '.join([line.strip() for line in file])
 71         comment = lowercase_sentence(comment)
 72         comment = separate_punctuation(comment)
 73         return comment
 74
 75
 76 def build_bag_of_word_vector(comment: str, vocabs: set):
 77     vector = {}
 78     for word in comment.split():
 79         if word in vocabs:
 80             if word in vector:
 81                 vector[word] += 1
 82             else:
 83                 vector[word] = 1
 84     return vector
 85
 86
 87 def preprocess_folder(folder_path: str, vocabs):
 88     vector_list = []
 89     for filename in os.listdir(folder_path):
 90         comment = preprocess_file(f'{folder_path}/{filename}')
 91         vector_list.append(build_bag_of_word_vector(comment, vocabs))
 92     return vector_list
 93
 94
 95 def preprocess(folder_path1, folder_path2, vocab_path, path1_class,
    path2_class, output_path):
 96     # label######{json format of vector}
 97     # ###### is the separator of column to access easier
 98     vocabs = load_vocab(vocab_path)
 99     folder_path1_vectors = preprocess_folder(folder_path1, vocabs)
100     folder_path2_vectors = preprocess_folder(folder_path2, vocabs)
101
102     res = []
103     for vector in folder_path1_vectors:
104         res.append(f'{path1_class}#####{json.dumps(vector)}#####{int(
    contains_strong_pos_word(list(vector.keys())))}#####{int(
    contains_strong_neg_word(list(vector.keys())))}')
105     for vector in folder_path2_vectors:
106         res.append(f'{path2_class}#####{json.dumps(vector)}#####{int(
    contains_strong_pos_word(list(vector.keys())))}#####{int(
    contains_strong_neg_word(list(vector.keys())))}')
107     save_file('\n'.join(res), output_path)
108
109
110
```

```python
1  import os
2  import json
3  import decimal
4
5  from pre_process import preprocess_comment, preprocess,
   preprocess_file
6
7
8  # ------------------------------------------------ load file
   ----------------------------------------------- #
9
10 def load_vocab(vocab_path) -> set:
11     vocabs = set()
12     with open(vocab_path) as vocab_file:
13         for word in vocab_file:
14             vocabs.add(word.strip())
15     return vocabs
16
17
18 def save_model(file_content, file_path):
19     with open(file_path, 'w') as file:
20         json.dump(file_content, file)
21
22
23 def save_file(file_content, file_path):
24     with open(file_path, 'w') as file:
25         file.write(file_content)
26
27
28 def load_json(file_path: str) -> dict:
29     with open(file_path) as file:
30         return json.load(file)
31
32
33 def load_one_vector(vector_file_path: str) -> dict:
34     return load_json(vector_file_path)
35
36
37 def load_naive_bayes_classifier(classifier_path: str) -> dict:
38     return load_json(classifier_path)
39
40
41 # ---------------------------------------- train naive bayes
   classifier ---------------------------------------- #
42
43 def initialize_counter(vocabs) -> dict:
44     '''
45     load vocab and build  a dictionary that contains all vocab as key
   , and value set to 0
46     :return: a dictionary that contains all vocab as key, and value
   set to 0
47     '''
48
49     return {vocab: 0 for vocab in vocabs}
50
51
52 def preprocessed_file_decoder(training_file_path: str):
53     with open(training_file_path) as file:
54         line = file.readline()
55         while line:
56             line = line.rstrip('\n')
```

```python
57                line = line.split('#####')
58                yield line[0], json.loads(line[1]), bool(line[2]), bool(
     line[3])
59                line = file.readline()
60
61
62  def aggregate_vector_into_counter(counter: dict, vector: dict,
     class_type, contains_strong_pos_word, contains_strong_neg_word):
63      total_token = 0
64      for word, freq in vector.items():
65          counter[word] += freq
66          total_token += freq
67
68      if class_type == 'pos' and contains_strong_pos_word or class_type
     == 'neg' and contains_strong_neg_word:
69          for word, freq in vector.items():
70              counter[word] += 100
71              total_token += 100
72      else:
73          for word, freq in vector.items():
74              if counter[word] > 100:
75                  counter[word] -= 100
76                  total_token -= 100
77              else:
78                  total_token -= counter[word]
79                  counter[word] = 0
80      return total_token
81
82
83  def train_naive_bayes_class_recognizer(counter: dict, total_token:
     int) -> dict:
84      recognizer = {}
85      total_vocab = len(counter)
86      add_one_smoothing_total_token = total_vocab + total_token
87      for word, freq in counter.items():
88          recognizer[word] = (freq + 1) / add_one_smoothing_total_token
89      return recognizer
90
91
92  def naive_bayes(training_file_path, result_model_path, vocab_path=""
     , class_1="pos", class_2="neg"):
93      vocabs = load_vocab(vocab_path)
94      counter = {
95          class_1: {
96              'counter': initialize_counter(vocabs),
97              'total_token': 0,
98              'class_recognizer': None,
99              'prior_prob': 0
100         },
101         class_2: {
102             'counter': initialize_counter(vocabs),
103             'total_token': 0
104         },
105     }
106
107     for class_type, vector, contains_strong_pos_word,
     contains_strong_neg_word in preprocessed_file_decoder(
     training_file_path):
108         counter[class_type]['total_token'] +=
     aggregate_vector_into_counter(counter[class_type]['counter'],
109
```

```
109             vector,
110
                class_type,
111
                contains_strong_pos_word,
112
                contains_strong_neg_word
113
                )
114
115         total_token = counter[class_1]['total_token'] + counter[class_2][
    'total_token']
116         naive_bayes_classifier = {class_1:
    train_naive_bayes_class_recognizer(counter[class_1]['counter'],
117
            counter[class_1]['total_token']
118
            ),
119                                     class_2:
    train_naive_bayes_class_recognizer(counter[class_2]['counter'],
120
            counter[class_2]['total_token']
121
            ),
122                                     f'{class_1}_prior': counter[class_1]['
    total_token'] / total_token,
123                                     f'{class_2}_prior': counter[class_2]['
    total_token'] / total_token,
124                                     "class_1": class_1,
125                                     "class_2": class_2
126                                     }
127
128         save_model(naive_bayes_classifier, result_model_path)
129         return naive_bayes_classifier
130
131
132 # ----------------------------------------------- evaluate test data
    --------------------------------------------- #
133
134 def compute_prob(comment: str | list, class_recognizer: dict,
    prior_prob: float) -> decimal.Decimal:
135     if type(comment) is str:
136         comment = comment.split()
137
138     # the min of a float is around 1e-310, it's very likely to have
    the prob of the sentence less than this
139     prob = decimal.Decimal(prior_prob)
140     for word in comment:
141         if word not in class_recognizer:
142             continue
143         prob = prob * decimal.Decimal(class_recognizer[word])
144     return prob
145
146
147 class NaiveBayesClassifier:
148     def __init__(self, path_to_model="", model=None):
149         self.model = model
150         self.class_1 = None
151         self.class_2 = None
152         if not model and path_to_model:
153             self.load_model(path_to_model)
```

```python
154
155    def load_model(self, path_to_model: str):
156        with open(path_to_model) as model_file:
157            self.model = json.load(model_file)
158            self.class_1 = self.model["class_1"]
159            self.class_2 = self.model["class_2"]
160
161    def classify(self, comment: str):
162        comment = preprocess_comment(comment)
163        word_list = comment.split()
164        class_1_prob = compute_prob(word_list, self.model[self.
    class_1], self.model[f"{self.class_1}_prior"])
165        class_2_prob = compute_prob(word_list, self.model[self.
    class_2], self.model[f"{self.class_2}_prior"])
166
167        # print(self.class_1, "probability is", class_1_prob)
168        # print(self.class_2, "probability is", class_2_prob)
169
170        return self.class_1 if class_1_prob > class_2_prob else self.
    class_2
171
172
173 # ------------------------------------- main (training and
    evaluate) ------------------------------------- #
174
175 # Train a classifier use a small corpus
176 def problem_2b():
177     preprocess(folder_path1="./data/movie_review_small/action",
178               folder_path2="./data/movie_review_small/comedy",
179               vocab_path="./data/movie_review_small/
    movie_review_small.vocab",
180               path1_class="action",
181               path2_class="comedy",
182               output_path="./preprocessed/movie_review_small.txt"
183               )
184
185     naive_bayes(training_file_path="./preprocessed/movie_review_small
    .txt",
186               result_model_path="models/movie_review_small.NB",
187               class_1="action",
188               class_2="comedy",
189               vocab_path="./data/movie_review_small/
    movie_review_small.vocab",
190               )
191
192
193 def problem_2c():
194     comment = "fast, couple, shoot, fly"
195     naive_bayes_classifier = NaiveBayesClassifier(path_to_model='
    models/movie_review_small.NB')
196     class_estimation = naive_bayes_classifier.classify(comment)
197
198     print(f"Class of sentence {comment} is: {class_estimation}")
199
200
201 def problem_2d():
202     # preprocess training data and train model
203     preprocess(folder_path1="./data/train/pos",
204               folder_path2="./data/train/neg",
205               vocab_path="./data/imdb.vocab",
206               path1_class="pos",
```

```python
207                     path2_class="neg",
208                     output_path="./preprocessed/movie_review_BOW.txt"
209                     )
210
211     naive_bayes(training_file_path="./preprocessed/movie_review_BOW.
    txt",
212                 result_model_path="./models/movie_review_BOW.NB",
213                 class_1="pos",
214                 class_2="neg",
215                 vocab_path="./data/imdb.vocab",
216                 )
217
218     naive_bayes_classifier = NaiveBayesClassifier(path_to_model='./
    models/movie_review_BOW.NB')
219     pos_test_folder = './data/test/pos'
220     neg_test_folder = './data/test/neg'
221     pos_test_files = os.listdir(pos_test_folder)
222     neg_test_files = os.listdir(neg_test_folder)
223
224     result = []  # [[estimation, comment],...]
225     incorrect = []
226     total_est = len(neg_test_files) + len(pos_test_files)
227
228     for file in pos_test_files:
229         comment = preprocess_file(file_path=f'{pos_test_folder}/{file
    }')
230         class_est = naive_bayes_classifier.classify(comment)
231         result.append(f'{class_est}    {comment}')
232         if class_est != 'pos':
233             incorrect.append(f'{class_est}    {comment}')
234
235     for file in neg_test_files:
236         comment = preprocess_file(file_path=f'{neg_test_folder}/{file
    }')
237         class_est = naive_bayes_classifier.classify(comment)
238         result.append(f'{class_est}    {comment}')
239         if class_est != 'neg':
240             incorrect.append(f'{class_est}    {comment}')
241
242     accuracy = (total_est - len(incorrect)) / total_est
243     result.append(f'Accuracy: {accuracy}    Total Estimations: {
    total_est}    Incorrect Estimations: {len(incorrect)}')
244     save_file('\n'.join(result), './report.txt')
245     save_file('\n'.join(incorrect), './incorrect.txt')
246
247 problem_2d()
248
249
```