```python
1  import nltk
2  import pandas as pd
3  import numpy as np
4  import random
5  import re
6  import joblib
7  import json
8  from collections import defaultdict
9  from nltk.corpus import stopwords
10 from sklearn.feature_extraction.text import CountVectorizer
11 from sklearn.linear_model import SGDClassifier
12 from sklearn.naive_bayes import MultinomialNB
13 from sklearn.metrics import accuracy_score
14
15
16 def preprocess_data(csvfile):
17     """
18     Read and preprocess CSV file complaint data
19
20     """
21
22     stops = stopwords.words('english')
23
24     df = pd.read_csv(csvfile)
25
26     ids = np.array(df["complaint_id"])
27     product_groups = np.array(df["product_group"])
28     raw_text_list = np.array(df["text"])
29
30     clean_text_list = raw_text_list
31
32     return ids, product_groups, clean_text_list
33
34 def define_group_labels(product_groups):
35     """
36     Map product group text label to numerical labels (random assignment) and return diction
37
38     """
39
40     group_to_label = dict()
41     label_to_group = dict()
42     for idx, g in enumerate(set(product_groups)):
43         group_to_label[g] = idx
44         label_to_group[idx] = g
45
46     return group_to_label, label_to_group
47
48 def create_vectorizer(document_list, vocabulary=None):
49     """
50     Function to create a vectorizer for input
51
52     """
53     vectorizer = CountVectorizer(ngram_range=(1,2), vocabulary=vocabulary, max_features=500
54     fit_vectorizer = vectorizer.fit(document_list)
55
56     return fit_vectorizer
57
58 def vectorize_documents(document_list, vectorizer):
59     """
60     Function to vectorize input text documents
61
62     """
63     X = vectorizer.transform(document_list)
64
65     return X
66
```

```python
67 def create_balanced_binary_sample(document_list, labels, focus_label):
68     """
69     Create equal, binary test data for One-Versus_Rest classification
70
71     """
72     target_documents = []
73     target_labels = [] # Class 1
74     other_documents = []
75     other_labels = [] # All other classes mapped to Class 0
76
77     for i in range(len(document_list)):
78         if labels[i] == focus_label:
79             target_documents.append(document_list[i])
80             target_labels.append(1)
81         else:
82             other_documents.append(document_list[i])
83             other_labels.append(0)
84
85     to_fill = len(target_documents)
86
87     # Sample other classes randomly
88     ziplist = list(zip(other_documents, other_labels))
89     sampled_others = random.sample(ziplist, to_fill)
90     second_documents, second_labels = zip(*sampled_others)
91
92     # Store in lists
93     output_documents = [d for d in target_documents]
94     output_documents.extend([d for d in second_documents])
95     output_labels = [l for l in target_labels]
96     output_labels.extend([l for l in second_labels])
97
98     return output_documents, output_labels
99
100 def create_kfolds(document_list, labels, k):
101     """
102     Randomize and split dataset into K folds (close but not exactly equal size folds)
103     (NOTE: Unused for final model training)
104
105     """
106     x_folds = []
107     y_folds = []
108
109     combined_lists = list(zip(document_list, labels))
110     random.shuffle(combined_lists)
111     rand_documents, rand_labels = zip(*combined_lists)
112
113     slicesize = len(rand_documents) // k
114
115     for num in range(k):
116         if num != k-1:
117             x_folds.append(rand_documents[num*slicesize:(num+1)*slicesize])
118             y_folds.append(rand_documents[num*slicesize:(num+1)*slicesize])
119         else:
120             x_folds.append(rand_documents[num*slicesize:])
121             y_folds.append(rand_documents[num*slicesize:])
122
123     return x_folds, y_folds
124
125 def elasticNet_feature_selection(x_folds, y_folds, vectorizer):
126     """
127     Function to train One-Versus-Rest Elastic Net Logistic Regression to find informative f
128     (NOTE: Unused for final model training)
129
130     """
131
132     clf  = SGDClassifier(loss="log",penalty="elasticnet")
```

```
133
134    for fold_idx in range(len(x_folds)):
135
136        test_documents = x_folds[fold_idx]
137        test_labels = y_folds[fold_idx]
138
139        train_documents = []
140        train_labels = []
141        for idx in range(len(x_folds)):
142            if idx != fold_idx:
143                train_documents.extend(x_folds[idx])
144                train_labels.extend(y_folds[idx])
145
146        X_train = vectorize_documents(train_documents, vectorizer)
147        y_train = np.array(train_labels)
148
149        print(X_train.shape)
150        print(len(y_train))
151
152        clf.fit(X_train, y_train)
153        y_pred = clf.predict(X_new)
154        model_accuracy = accuracy_score(y_pred, y_new)
155
156        print("SGD Model Accuracy on Fold %d: " %fold_idx)
157        print(model_accuracy)
158
159    return best_model
160
161 if __name__ == "__main__":
162
163    # Preprocess data
164    print("\nPreprocessing documents...")
165    ids, product_groups, clean_text_list = preprocess_data("case_study_data.csv")
166
167    # Define numerical target labels (y)
168    group_to_label, label_to_group = define_group_labels(product_groups)
169    group_numbers = [group_to_label[group] for group in product_groups]
170
171    # Create initial vectorizer
172    print("\nCreating Vectorizer...")
173    firstvectorizer = create_vectorizer(clean_text_list)
174
175    for group, label in group_to_label.items():
176
177        binary_docs, binary_labels = create_balanced_binary_sample(clean_text_list, group_n
178
179        print(label, group)
180        print(len(binary_docs))
181
182        #x_folds, y_folds = create_kfolds(binary_docs, binary_labels, 5)
183
184        #elasticNet_feature_selection(x_folds, y_folds, firstvectorizer)
185
186        traintest_split = int(len(binary_docs)*.8)
187
188        X_train = vectorize_documents(binary_docs[0:traintest_split], firstvectorizer)
189        y_train = np.array(binary_labels[0:traintest_split])
190
191        X_test = vectorize_documents(binary_docs[traintest_split:], firstvectorizer)
192        y_test = np.array(binary_labels[traintest_split:])
193
194        clf  = SGDClassifier(loss="log",penalty="elasticnet")
195
196        clf.fit(X_train, y_train)
197        y_pred = clf.predict(X_test)
198
```

```
199        print("SGD Model Accuracy: ")
200        print(accuracy_score(y_pred, y_test))
201
202        idx_to_feature = dict()
203        feature_to_weights = defaultdict(list)
204        selected_vocabulary = set()
205
206        for key, value in firstvectorizer.vocabulary_.items():
207            idx_to_feature[value] = key
208        for idx, c in enumerate(clf.coef_[0]):
209            if abs(c) > 0:
210                selected_vocabulary.add(idx_to_feature[idx])
211
212    print("\nVocabulary Size: %d" %len(selected_vocabulary))
213    print("Creating Fine Tuned Vectorizer...")
214    newvectorizer = create_vectorizer(clean_text_list, vocabulary=selected_vocabulary)
215
216    splitidx = int(len(ids) * .8)
217
218    combined_lists = list(zip(product_groups, clean_text_list))
219    random.shuffle(combined_lists)
220    rand_product_groups, rand_text_list = zip(*combined_lists)
221
222
223    print("\nVectorizing documents...")
224    X = vectorize_documents(rand_text_list[0:splitidx], newvectorizer)
225    y = np.array([group_to_label[group] for group in rand_product_groups[0:splitidx]])
226
227    print("\nVectorizing documents...")
228    X_new = vectorize_documents(rand_text_list[splitidx:], newvectorizer)
229    y_new = np.array([group_to_label[group] for group in rand_product_groups[splitidx:]])
230
231    alphas = [.001, .01, .05, .5, .75, 1]
232
233    best_alpha = .05
234    current_score = 0
235    for a in alphas:
236        clf  = MultinomialNB(alpha=a)
237
238        clf.fit(X, y)
239        y_pred = clf.predict(X_new)
240
241        acc = accuracy_score(y_pred, y_new)
242
243        print("Model With Alpha %f: " %a)
244        print(acc)
245        if acc > current_score:
246            current_score = acc
247            best_alpha = a
248
249    print("\nVectorizing Final documents...")
250    X_final = vectorize_documents(rand_text_list, newvectorizer)
251    y_final = np.array([group_to_label[group] for group in rand_product_groups])
252
253    final_model = MultinomialNB(alpha=best_alpha)
254    final_model.fit(X_final, y_final)
255
256    joblib.dump(newvectorizer, "NWCaseStudyvectorizer.joblib")
257    joblib.dump(final_model, "NWCaseStudyNBModel.joblib")
258
259    with open("group_mapping.json", "w") as file:
260        file.write(json.dumps(group_to_label))
```