

```
1 #%% md
2 # Import Functions
3 #%%
4 import pandas as pd
5 from tqdm import tqdm
6 import random as rd
7 import torch
8 from torch.utils.data import Dataset, DataLoader
9 import numpy as np
10 from sklearn.model_selection import train_test_split
11 from sklearn.linear_model import LinearRegression
12 from sklearn.preprocessing import LabelBinarizer
13 from sklearn.metrics import classification_report
14 #%% md
15 # Initiate lists of data to load into training
16 #%% md
17 ## Product Types list and query funciton
18 #%%
19 def product_list_creation(file_name: str) -> list:
20     products_df = pd.read_excel(file_name, sheet_name='product_type')
21
22     unique_family = products_df['Family'].dropna().unique().tolist()
23     product_types = products_df['ProductType'].dropna().tolist()
24     keywords = [i.split(', ') for i in products_df['Keywords'].dropna().tolist()]
25     keyword_list = list(set([w for word in keywords for w in word]))
26
27     product_types.extend(keyword_list)
28     product_types.extend(unique_family)
29     product_types = list(set(product_types))
30     return product_types
31 #%%
32 def product_query(product_types):
33     queries = []
34     labels = []
35     for word in product_types:
36         labels_product = ['B-Product'] + ['I-Product']
```

```
36 ] * (len(word.split(' ')) - 1)
37         queries.append(word.split(' '))
38         labels.append(labels_product)
39     return queries, labels
40 #%% md
41 ## Material Types list and query function
42 #%%
43 def material_list_creation(file_name: str) -> list:
44     materials_df = pd.read_excel(file_name,
        sheet_name='material_type')
45
46     material_types = materials_df['material_type'].
        dropna().tolist()
47     mat_fam = materials_df['family'].dropna().tolist()
48     mat_k_word = list(set([w for words in [i.split(
        ', ') for i in materials_df['keywords'].dropna().
        tolist()] for w in words]))
49
50     material_types.extend(mat_fam)
51     material_types.extend(mat_k_word)
52     material_types = list(set(material_types))
53     return material_types
54 #%%
55 def material_query(material_types):
56     queries = []
57     labels = []
58     for word in material_types:
59         word = word.lower()
60         labels_material = ['B-Material'] + ['I-
    Material'] * (len(word.split(' ')) - 1)
61         queries.append(word.split(' '))
62         labels.append(labels_material)
63     return queries, labels
64 #%% md
65 ## Building Applications list and query function
66 #%%
67 def building_app_list_creation(file_name: str) ->
    list:
68     building_app_df = pd.read_excel(file_name,
        sheet_name='building_applications')
```

```

69
70     building_app_list = building_app_df['
  building_applications'].dropna().tolist()
71     keywords = list(set([w for words in [i.split(
    ', ') for i in building_app_df['keywords'].dropna().tolist()] for w in words]))
72     building_app_list.extend(keywords)
73     return building_app_list
74 #%%
75 def building_app_query(building_app_list):
76     queries = []
77     labels = []
78     for word in building_app_list:
79         labels_building_application = ['B-
  Application'] + ['I-Application'] * (len(word.split(
    ' ')) - 1)
80         queries.append(word.split(' '))
81         labels.append(labels_building_application)
82     return queries, labels
83 #%% md
84 ## Country list and query function
85 #%%
86 def country_list_creation(file_name: str) -> list:
87     country_df = pd.read_excel(file_name, sheet_name =
  'countries')
88     country_list = country_df['country'].dropna().tolist()
89     return country_list
90 #%%
91 def country_query(country_list):
92     queries = []
93     labels = []
94     for word in country_list:
95         labels_country = ['B-Country'] + ['I-Country
  '] * (len(word.split(' ')) - 1)
96         queries.append(word.split(' '))
97         labels.append(labels_country)
98     return queries, labels
99 #%% md
100 ## Recycled Content list and query
101 #%%

```

```

102 def recycled_list_creation(file_name: str) -> list:
103     recycled_content_df = pd.read_excel(file_name,
104                                         sheet_name='recycled_content')
105     recycled_list = recycled_content_df['
106                                         recycled_content'].dropna().str.lower().tolist()
107     recycle_keywords = list(set([w for words in [i.
108                                         split(', ') for i in recycled_content_df['keywords'
109                                         ].dropna().tolist()] for w in words]))
110     recycled_list.extend(recycle_keywords)
111
112     final_list = []
113     for item in recycled_list:
114         if 'x%' in item:
115             for _ in range(21):
116                 final_list.append(item.replace('x%', str(rd.randint(0, 100)) + '%'))
117         elif '.x' in item:
118             for _ in range(21):
119                 final_list.append(item.replace('.x', '.' + str(rd.randint(0, 100))))
120     else:
121         final_list.append(item)
122
123     return final_list
124 #%%
125 def recycle_query(recycled_list):
126     queries = []
127     labels = []
128     for word in recycled_list:
129         labels_country = ['B-Recycle'] + ['I-Recycle'
130                                         ] * (len(word.split(' ')) - 1)
131         queries.append(word.split(' '))
132         labels.append(labels_country)
133
134     return queries, labels
135 #%%
136 '''
137 EXCEL FILE WITH SHEETS OF UNIQUE PRODUCT/MATERIAL/
138 BUILDINGAPPLICATION TYPE REDACTED
139 PROPRIETARY KNOWLEDGE OF 2050 MATERIALS
140 CONTACT info@2050-materials.com FOR MORE INFORMATION
141 '''
142 def list_initialization(file_name: str = '

```

```
134 YOUR_EXCEL_FILE.xlsx'):
135     product_types = product_list_creation(file_name)
136     material_types = material_list_creation(
137         file_name)
138     building_app_list = building_app_list_creation(
139         file_name)
140     country_list = country_list_creation(file_name)
141     recycled_list = recycled_list_creation(file_name)
142 )
143
144     return product_types, material_types,
145         building_app_list, country_list, recycled_list
146 #%% md
147 ## Creates query with proper label for every word in
148     product types, material types, building
149     applications, countries, and recycled content
150 #%%
151
152 def all_queries(product_types, material_types,
153     building_app_list, country_list, recycled_list):
154     queries = []
155     labels = []
156
157     p_queries, p_labels = product_query(
158         product_types)
159     queries.extend(p_queries)
160     labels.extend(p_labels)
161
162     m_queries, m_labels = material_query(
163         material_types)
164     queries.extend(m_queries)
165     labels.extend(m_labels)
166
167     b_queries, b_labels = building_app_query(
168         building_app_list)
169     queries.extend(b_queries)
170     labels.extend(b_labels)
171
172     c_queries, c_labels = country_query(country_list)
173     queries.extend(c_queries)
174     labels.extend(c_labels)
```

```

164
165     r_queries, r_labels = recycle_query(
166         recycled_list)
167     queries.extend(r_queries)
168     labels.extend(r_labels)
169
170     return queries, labels
171 ## md
172 ## Creation of template sentences, all words in
173 ## database added with "all_words='y'"
174 ##%
175 #%% md
176 def new_sentences(num, product_types, material_types,
177   , building_app_list, country_list, recycled_list,
178   all_words='n'):
179     # template sentences that selects information at
180     # random to be selected for training data
181     templates = [
182       "{product_type} products
183       manufactured {country} for {building_application}",
184       "{material_type} manufactured in {
185       country} for {building_application} with {
186       recycled_content} recycled content",
187       "{product_type} manufactured {
188       country} with {recycled_content} recycled content",
189       "{product_type} made in {country}
190       for {building_application} with {recycled_content}
191       recycle content",
192       "Find {product_type} from {country}
193       } used in {building_application}",
194       "Do you have any {product_type}
195       used in {building_application}?",
196       "I'm looking for {material_type}
197       products from {country} for {building_application}."
198       ,
199       "{product_type} from {country}
200       suitable for {building_application}?",
201       "Search for {material_type}
202       materials from {country} used in {
203       building_application} applications",
204       "Can I find {product_type} suitable
205       for {building_application} applications?",
```

```
186          "What types of {material_type} from
187              {country} are used in {building_application}?",
188          "Show me {product_type} suitable
189              for {building_application}",
190          "List all {material_type} from {
191              country} used in {building_application}.",
192          "Are there any {product_type} from
193              {country} for {building_application}?",
194          "Find {product_type} with {
195              recycled_content} recycled content.",
196          "Find {material_type} with {
197              recycled_content} recycled content.",
198          "List all {material_type} from {
199              country} with recycled content {recycled_content} "
200          ]
201
202
203     queries = []
204     labels = []
205
206     for null in range(num):
207         temp = rd.choice(templates)
208         material_type = rd.choice(material_types)
209         product_type = rd.choice(product_types)
210         building_application = rd.choice(
211             building_app_list)
212         country = rd.choice(country_list)
213         recycled_content = rd.choice(recycled_list)
214         sentence = temp.format(
215             material_type=
216             material_type,
217             product_type=
218             product_type,
219             building_application=
220             building_application,
221             country=country,
222             recycled_content=
223             recycled_content
224         )
225
226         queries.append(sentence.lower())
227
228     # Generate labels
```

```

215     labels_product = ['B-Product'] + ['I-Product']
216         ''] * (len(product_type.split(' ')) - 1)
217     labels_material = ['B-Material'] + ['I-
218 Material'] * (len(material_type.split(' ')) - 1)
219     labels_country = ['B-Country'] + ['I-Country']
220         ''] * (len(country.split(' ')) - 1)
221     labels_building_application = ['B-
222 Application'] + ['I-Application'] * (len(
223 building_application.split(' ')) - 1)
224     labels_recycled_content = ['B-Recycle'] + ['
225 I-Recycle'] * (len(str(recycled_content).split(' '
226 )) - 1)
227
228     # Identify the placeholders in the template
229     and replace them with the corresponding labels
230     template_labels = temp.split()
231     new_template_labels = []
232     for word in template_labels:
233         if 'product_type' in word:
234             new_template_labels.extend(
235             labels_product)
236         elif 'material_type' in word:
237             new_template_labels.extend(
238             labels_material)
239         elif 'country' in word:
240             new_template_labels.extend(
241             labels_country)
242         elif 'building_application' in word:
243             new_template_labels.extend(
244             labels_building_application)
245         elif 'recycled_content' in word:
246             new_template_labels.extend(
247             labels_recycled_content)
248         else:
249             new_template_labels.append('0')
250
251     template_labels = new_template_labels
252     labels.append(template_labels)
253
254     if all_words == 'y':
255         solo_queries, solo_labels = all_queries(

```

```
242 product_types, material_types, building_app_list,
243     country_list, recycled_list)
244     queries.extend(solo_queries)
245     labels.extend(solo_labels)
246     print(f'{len(solo_queries)} words from
247 database and {num} random sentences added to queries
248 and labels data.')
249 else:
250     print(f'{num} random sentences added to
251 queries and labels data.')
252 return queries, labels
253 %% md
254 ## CustomDataset class to process words for
255 tokeninzer
256 %%
257 class CustomDataset(Dataset):
258     def __init__(self, tokenizer, sentences, labels
259 , max_len, pad_token_label_id):
260         self.len = len(sentences)
261         self.sentences = sentences
262         self.labels = labels
263         self.tokenizer = tokenizer
264         self.max_len = max_len
265         self.pad_token_label_id = pad_token_label_id
266
267     def __getitem__(self, index):
268         sentence = str(self.sentences[index])
269         inputs = self.tokenizer.encode_plus(
270             sentence,
271             None,
272             add_special_tokens=True,
273             max_length=self.max_len,
274             padding='max_length',
275             return_token_type_ids=True
276         )
277         ids = inputs['input_ids']
278         mask = inputs['attention_mask']
279
280         label = ['0'] + self.labels[index] + ['0']
281         label = [Token_Classification_Labels.get(l)
282 for l in label]
```

```
276         padding_length = self.max_len - len(label)
277         label.extend([self.pad_token_label_id]*padding_length)
278         label = label[:self.max_len]
279     return {
280         'input_ids': torch.tensor(ids, dtype=torch.long),
281         'attention_mask': torch.tensor(mask, dtype=torch.long),
282         'labels': torch.tensor(label, dtype=torch.long)
283     }
284
285     def __len__(self):
286         return self.len
287 #%% md
288 ## Loads queries and labels into train dataset to
289 ## prepare for training
290 #%% md
291     def train_data_info(queries, labels, tokenizer):
292         # First, split into train+val and test sets
293         sentences_temp, sentences_test, labels_temp,
294         labels_test = train_test_split(queries, labels,
295             test_size=0.1, random_state=42)
296
297         # Then split train+val into separate train and
298         # val sets
299         sentences_train, sentences_val, labels_train,
300         labels_val = train_test_split(sentences_temp,
301             labels_temp, test_size=0.1, random_state=42)
302
303         # Padding label
304         PAD_TOKEN_LABEL_ID = -100
305
306         # Now, prepare the training, validation, and
307         # test data
308         train_data = CustomDataset(tokenizer,
309             sentences_train, labels_train, max_len=50,
310             pad_token_label_id=PAD_TOKEN_LABEL_ID)
311         valid_data = CustomDataset(tokenizer,
312             sentences_val, labels_val, max_len=50,
```

```
302     pad_token_label_id=PAD_TOKEN_LABEL_ID)
303     test_data = CustomDataset(tokenizer,
304         sentences_test, labels_test, max_len=50,
305         pad_token_label_id=PAD_TOKEN_LABEL_ID)
306
307     # Initialize the data loaders
308     train_dataloader = DataLoader(train_data,
309         batch_size=32, shuffle=True)
310     valid_dataloader = DataLoader(valid_data,
311         batch_size=32, shuffle=True)
312     test_dataloader = DataLoader(test_data,
313         batch_size=32, shuffle=False) # we don't need to
314         shuffle the test set
315
316     return train_dataloader, valid_dataloader,
317         test_dataloader
318
319 #%%
320
321 def Initialize_tokenizer():
322     from transformers import DistilBertTokenizerFast
323     # Initialize the tokenizer
324     tokenizer = DistilBertTokenizerFast.
325         from_pretrained('distilbert-base-uncased')
326
327     return tokenizer
328
329 #%%
330
331 def Token_Classification_Initialization():
332     from transformers import
333         DistilBertForTokenClassification
334
335     Token_Classification_Labels = {
336
337         '0': 0,
338         'B-Product': 1,
339         'I-Product': 2,
340         'B-Material': 3,
341         'I-Material': 4,
342         'B-Country': 5,
343         'I-Country': 6,
344         'B-Application': 7
345
346         ,
347         'I-Application': 8
348
349         ,
350         'B-Recycle': 9,
351         'I-Recycle': 10
352     }
```

```
332     device = torch.device('cuda') if torch.cuda.
333         is_available() else torch.device('cpu')
334     Token_Classification_model =
335         DistilBertForTokenClassification.from_pretrained(
336             "distilbert-base-
337             uncased",
338             num_labels=len(
339                 Token_Classification_Labels)).to(device)
340     return Token_Classification_Labels,
341     Token_Classification_model
342 #%%
343 def Token_Classification_Model_Training(
344
345     Token_Classification_Labels: dict,
346
347     Token_Classification_model,
348
349     train_dataloader,
350
351     test_dataloader,
352
353     Model_load_file: None|str = None,
354
355     learning_rate: float = 1e-5,
356
357     epoch_loop
358     : int = 5
359
360     ) -> list:
361
362     optimizer = torch.optim.AdamW(params=
363         Token_Classification_model.parameters(), lr=
364         learning_rate)
365
366     # Loads model if there is provided file name
367     if Model_load_file:
368
369         Token_Classification_model.load_state_dict(
370             torch.load(Model_load_file))
371
372         print(f'{Model_load_file} successfully
373             loaded.')
374
375     # Set the device
376     device = torch.device("cuda" if torch.cuda.
377         is_available() else "cpu")
378
379     Token_Classification_model.to(device)
380
381
```

```
356     Token_Classification_model.train()
357
358     # Training loop
359     training_loss = []
360     raw_results = []
361     print('Training starting.')
362     for null in tqdm(range(epoch_loop)): # loop
363         over the dataset multiple time
363         total_loss = 0
364         for i, batch in enumerate(train_dataloader,
365             0):
365             # zero the parameter gradients
366             optimizer.zero_grad()
367
368             # get the inputs
369             inputs = batch['input_ids'].to(device)
370             labels = batch['labels'].to(device)
371             masks = batch['attention_mask'].to(
371                 device)
372
373             # forward + backward + optimize
374             outputs = Token_Classification_model(
374                 inputs, attention_mask=masks, labels=labels)
375             loss = outputs.loss
376
377             loss.backward()
378             optimizer.step()
379
380             # add the loss to the list
381             total_loss += loss.item()
382
383             Token_Classification_model.eval()
384
385             predictions = []
386             true_labels = []
387
388             for i, batch in enumerate(test_dataloader):
389                 b_input_ids = batch['input_ids'].to(
389                     device)
390                 b_input_mask = batch['attention_mask'].to(
390                     device)
```

```
391         b_labels = batch['labels'].to(device)
392
393         with torch.no_grad():
394             outputs = Token_Classification_model
395             (b_input_ids, attention_mask=b_input_mask)
396
397             logits = outputs[0].detach().cpu().numpy()
398             ()
399             label_ids = b_labels.to('cpu').numpy()
400
401             # Calculate the predictions
402             predictions.extend([list(p) for p in np.
403             argmax(logits, axis=2)])
404             true_labels.extend(label_ids)
405
406
407             # Define mapping from index to labels
408             idx2label = {v: k for k, v in
409             Token_Classification_Labels.items()}
410             idx2label[-100] = 'PAD'
411
412             # Flatten the output tensors
413             predictions_flat = [item for sublist in
414             predictions for item in sublist]
415             true_labels_flat = [item for sublist in
416             true_labels for item in sublist]
417
418             # Map indices back to labels
419             pred_tags_flat = [idx2label[p] for p in
420             predictions_flat]
421             true_tags_flat = [idx2label[t] for t in
422             true_labels_flat]
423
424
425             # Filter out 'PAD' labels
426             pred_tags = [pred for pred, true in zip(
427             pred_tags_flat, true_tags_flat) if true != 'PAD']
428             true_tags = [true for true in true_tags_flat
429             if true != 'PAD']
430
431             # Binarize the labels
432             lb = LabelBinarizer()
433             y_true_combined = lb.fit_transform(true_tags)
```

```
421 )
422         y_pred_combined = lb.transform(pred_tags)
423
424         # appends classification report into list
425         # for each epoch
426         class_report = classification_report(
427             y_true_combined, y_pred_combined, target_names=lb.
428             classes_, digits=4, zero_division=1, output_dict=
429             True)
430         raw_results.append(class_report)
431         Token_Classification_model.train()
432
433         # Save training loss
434         training_loss.append(total_loss/len(
435             train_dataloader))
436
437         raw_results.append(training_loss)
438         print('Training Complete.')
439         if Model_load_file:
440             torch.save(Token_Classification_model.
441                 state_dict(), Model_load_file)
442             print(f'{Model_load_file} successfully saved
443 .')
444         else:
445             Model_load_file = 'Token_Class_Model.pth'
446             torch.save(Token_Classification_model.
447                 state_dict(), Model_load_file)
448             print(f'{Model_load_file} successfully saved
449 .')
450
451     return raw_results
452 #%% md
453 ## Organizes results from training into dictionary
454 #%%
455 def results_dict_creation(raw_results: list) -> dict
456     :
457         training_loss = raw_results.pop()
458         results_dict = {
459             'epoch': [],
460             'precision': [],
461             'recall': [],
462             'f1_score': [],
```

```

452                 'training loss': []
453             }
454         for index, item in enumerate(raw_results):
455             results_dict['epoch'].append(index+1)
456             results_dict['precision'].append(item['
457               weighted avg']['precision'])
458             results_dict['recall'].append(item['weighted
459               avg']['recall'])
460             results_dict['f1_score'].append(item['
461               weighted avg']['f1-score'])
462             results_dict['training loss'].append(
463               training_loss[index])
464         return results_dict
465 %% md
466 ## Takes Epoch df column and column to compare to
467   output graph of progression over each epoch
468 %%%
469 def plot_creation(epochs, y_value):
470     import matplotlib.pyplot as plt
471
472     epoch_nums = epochs.tolist()
473     y_val_list = y_value.tolist()
474
475     x = epochs.values.reshape(-1,1)
476     y = y_value.values.reshape(-1,1)
477
478     linear_r = LinearRegression()
479     linear_r.fit(x, y)
480     y_pred = linear_r.predict(x)
481     plt.plot(epoch_nums, y_pred, color='red',
482               linestyle='dashdot')
483
484     plt.plot(epoch_nums, y_val_list, color='blue')
485     plt.title(y_value.name)
486     plt.xlabel('epoch')
487     plt.legend(['linear regression', y_value.name])
488
489     plt.tight_layout()
490     plt.show()
491 %% md
492 ## Saves results of training into .csv file

```

```
487 #%%
488 def save_results(df):
489     import datetime
490     import pytz
491     '''
492         SAVES ALL TRAINING RESULTS, PRECISION RECALL
493         F1_SCORE TRAINING LOSS
494         INTO .csv FILE IN YEAR-MONTH-DAY-HOUR-MINUTE-
495         SECOND.csv FILE NAME,
496         LINKED WITH YOUR TIMEZONE SO NO DUPLICATE FILES
497         AND KNOWLEDGE OF WHEN TRAINING IS COMPLETE
498         '''
499     current_time = datetime.datetime.now(pytz.
500                                         timezone('REPLACE_WITH_YOUR_TIMEZONE'))
501     save_location = (
502         'FOLDER_LOCATION_OF_TRAINING_DATA' + current_time.
503         strftime('%Y-%m-%d_%H-%M-%S') + '.csv')
504     df.to_csv(save_location, sep=',', index=False)
505     print(f'File {save_location} saved.')
506 #%% md
507 ## Function to take sentence/query from user and
508 process result, output in dictionary {str:list}
509 #%%
510 def Token_Classification_Results(List_of_Sentences:
511     list,
512     Token_Classification_Labels: dict,
513     Token_Classification_model,
514     tokenizer,
515     Token_Class_File:
516     str
517             ) -> dict:
518     # Puts model into eval mode
519     Token_Classification_model.eval()
520
521     # Loads the saved file of Pre-trained model
522     Token_Classification_model.load_state_dict(torch
523         .load(Token_Class_File))
524     device = torch.device('cuda') if torch.cuda.
525     is_available() else torch.device('cpu')
```

```
515
516      # Inverse mapping from label indices to labels
517      inv_label_map = {v: k for k, v in
518          Token_Classification_Labels.items()}
519
520      # Creating Final Dictionary for results
521      final_dict = {
522          'Product Type': [],
523          'Material Type': [],
524          'Building applications': [],
525          'Countries': [],
526          'Recycled Content': []
527      }
528
529      # Cycling through query
530      for sentence in List_of_Sentences:
531          sentence = sentence.lower()
532          # Encode the sentence
533          inputs = tokenizer.encode_plus(
534              sentence,
535              None,
536              add_special_tokens=True,
537              padding='longest',
538              return_token_type_ids=True
539          )
540
541          # Create torch tensors and move them to the
542          # device
543          input_ids = torch.tensor([inputs['input_ids']],
544          dtype=torch.long).to(device)
545          attention_mask = torch.tensor([inputs['
546          attention_mask']], dtype=torch.long).to(device)
547
548          # Run the sentence through the model
549          with torch.no_grad():
550              outputs = Token_Classification_model(
551                  input_ids, attention_mask=attention_mask)
```

```
548     # Get the token-level class probabilities
549     logits = outputs[0]
550
551     # Compute the predicted labels
552     predictions = torch.argmax(logits, dim=-1)
553
554     # Remove padding and special tokens
555     input_ids = input_ids[0].tolist()
556     predictions = predictions[0].tolist()
557
558     # real_input_ids = [id for id, pred in zip(
559     #     input_ids, predictions) if id != 0 and id != 101 and
560     #     id != 102]
561     real_predictions = [pred for id, pred in zip(
562     input_ids, predictions) if id != 0 and id != 101
563     and id != 102]
564
565     # real_input_ids = []
566     # real_predictions = []
567     #
568     # for id, pred in zip(input_ids, predictions):
569     #     if id != 0 and id != 101 and id != 102
570     #: # Ignore [PAD], [CLS] and [SEP] tokens
571     #     real_input_ids.append(id)
572     #     real_predictions.append(pred)
573
574     # Map predicted label indices back to label
575     # strings
576     predicted_labels = [inv_label_map[label] for
577     label in real_predictions]
578
579     # Combine tokens and their predicted labels
580     # into dict
581     results = dict(zip(sentence.split(' '), 
582     predicted_labels))
583
584     # Put results into final_dict, key is word
585     # and value is what it is labeled as
586     for word, label in results.items():
587         print(word, label)
```

```
578         if word in ['for', 'and', 'from']:
579             continue
580         if label != 'O':
581             # Adds word to proper dictionary if
582             # belongs to that type
583             # Checks if there is a B-word before
584             # the I-word so it will add to it
585             # It will make a new word if there
586             # is no B-word
587             if label == 'B-Product':
588                 final_dict['Product Type'].
589                 append(word)
590             elif label == 'I-Product':
591                 if final_dict['Product Type']:
592                     final_dict['Product Type'][-1]
593                     ] += ' ' + word
594                 else:
595                     final_dict['Product Type'].
596                     append(word)
597             elif label == 'B-Material':
598                 final_dict['Material Type'].
599                 append(word)
600             elif label == 'I-Material':
601                 if final_dict['Material Type']:
602                     final_dict['Material Type'][-1]
603                     ] += ' ' + word
604                 else:
605                     final_dict['Material Type'].
606                     append(word)
607             elif label == 'B-Country':
608                 final_dict['Countries'].append(
609                     word)
```

```
609
610             elif label == 'I-Country':
611
612                 if final_dict['Countries']:
613                     final_dict['Countries'][-1] +=
614                         ' ' + word
615
616             else:
617                 final_dict['Countries'].append(
618                     word)
619
620             elif label == 'B-Application':
621
622                 final_dict['Building
623 applications'].append(word)
624
625             elif label == 'I-Application':
626
627                 if final_dict['Building
628 applications']:
629                     final_dict['Building
630 applications'][-1] += ' ' + word
631
632             else:
633                 final_dict['Building
634 applications'].append(word)
635
636             elif label == 'B-Recycle':
637
638                 final_dict['Recycled Content'].
639                     append(word)
640
641             elif label == 'I-Recycle':
642
643                 if final_dict['Recycled Content
644 ']:
645                     final_dict['Recycled Content
646 '][-1] += ' ' + word
647
648             else:
649                 final_dict['Recycled Content
650 '].append(word)
651
652             print('\n')
653
654
655     return final_dict
```

```
640 %% md
641 # Testing
642 %% md
643 ## Model Initialization
644 %%
645 Token_Classification_Labels,
    Token_Classification_model =
    Token_Classification_Initialization()
646 %%
647 '''
648 EXCEL FILE WITH SHEETS OF UNIQUE PRODUCT/MATERIAL/
    BUILDINGAPPLICATION TYPE REDACTED
649 PROPRIETARY KNOWLEDGE OF 2050 MATERIALS
650 CONTACT info@2050-materials.com FOR MORE INFORMATION
651 '''
652 excel_file = 'YOUR_EXCEL_FILE.xlsx'
653
654 product_types, material_types, building_app_list,
    country_list, recycled_list = list_initialization(
        excel_file)
655 %%
656 queries, labels = new_sentences(1000, product_types
    , material_types, building_app_list, country_list,
    recycled_list, all_words='n')
657 %%
658 tokenizer = Initialize_tokenizer()
659 %%
660 train_dataloader, valid_dataloader, test_dataloader
    = train_data_info(queries=queries, labels=labels,
        tokenizer=tokenizer)
661 '''
662 TOKEN CLASSIFICATION MODEL SAVE PATH LOCATION
663 IF FIRST TIME TRAINING, OR NO FILE IS PASSED,
664 IT WILL SAVE 'Token_Class_Model.pth' IN THIS FOLDER
    WHEN COMPLETE WITH TRAINING
665 '''
666 model_file = 'YOUR_TOKEN_CLASSIFICATION_MODEL.pth'
667 %%
668 raw_results = Token_Classification_Model_Training(
    Token_Classification_Labels=
    Token_Classification_Labels,
```

```
669     Token_Classification_model=
670         Token_Classification_model,
671     train_dataloader=train_dataloader,
672     test_dataloader=test_dataloader,
673     Model_load_file=model_file,
674     learning_rate=1e-5,
675     epoch_loop=5)
676 #%%
677 results_dict = results_dict_creation(raw_results)
678 df = pd.DataFrame.from_dict(results_dict)
679 save_results(df)
680 #%%
681
682 for index in range(1, 5):
683     plot_creation(df.iloc[:, 0], df.iloc[:, index])
684 #%%
685 '''
686 CREATE SENTENCE TO TEST MODEL
687 '''
688 list_of_sentences = ['TEST SENTENCE PASSED INTO
689                         MODEL TO DETERMINE OUTPUT']
690 final_output = Token_Classification_Results(
691     List_of_Sentences=list_of_sentences,
692     Token_Classification_Labels=
693         Token_Classification_Labels,
694     Token_Classification_model=
695         Token_Classification_model,
696     tokenizer=tokenizer,
697     Token_Class_File=model_file
698 )
699 
```

```
695 #%%
696 final_output
697 #%%
698
```