

▼ Assignment 1 for FIT5212, Semester 1

Student Name: Tsz Yan CHUNG

Student ID: 32973381

```
1 # !pip3 install torch==1.9.1+cu111 torchvision==0.10.1+cu111 torchaudio==0.9.1 -f https:
```

```
1 # !pip3 install torchtext==0.10.0
```

▼ Import Libraries

```
1 from google.colab import drive
2 drive.mount('/content/drive')
3

4 import copy
5 import random
6 import time
7 import re
8 import numpy as np
9 import pandas as pd
10 import seaborn as sns
11 import matplotlib.pyplot as plt
12
13 import spacy
14 from spacy.lang.en.stop_words import STOP_WORDS
15
16 import nltk
17 from nltk.corpus import stopwords
18 from nltk import word_tokenize
19 from nltk.tokenize import wordpunct_tokenize
20 from nltk.stem import WordNetLemmatizer
21 from nltk.stem import PorterStemmer
22 nltk.download('punkt')
23 nltk.download('omw-1.4')
24 nltk.download('wordnet')
25 nltk.download('stopwords')
26
27 import torch
28 import torchtext
29 import torch.nn as nn
```

```
28 import torch.optim as optim
29 from torchtext.legacy import data
30 from torchtext.legacy.data import Field, LabelField, TabularDataset, Dataset
31
32 from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score, con
33 from sklearn.metrics import precision_recall_curve
34 from sklearn.model_selection import cross_val_score, train_test_split
35 from sklearn.feature_extraction.text import TfidfVectorizer
36 from sklearn.linear_model import LogisticRegression
37 from sklearn.naive_bayes import MultinomialNB, GaussianNB, BernoulliNB
38 from sklearn.svm import LinearSVC, SVC
39 from sklearn.ensemble import RandomForestClassifier
40
[nltk_data] Downloading package punkt to
[nltk_data]      C:\Users\Sunny\AppData\Roaming\nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package omw-1.4 to
[nltk_data]      C:\Users\Sunny\AppData\Roaming\nltk_data...
[nltk_data] Package omw-1.4 is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]      C:\Users\Sunny\AppData\Roaming\nltk_data...
[nltk_data] Package wordnet is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]      C:\Users\Sunny\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

```
1 %matplotlib inline
2 import nltk
3 import torch
4 import torchtext
5 from torchtext.legacy import data
6 from torchtext.legacy.data import Field, LabelField, TabularDataset, Dataset
7
8 import nltk
9 import pandas as pd
10 import numpy as np
11 import seaborn as sns
12 import matplotlib.pyplot as plt
13 import spacy
14
15 from nltk.corpus import stopwords
16 from nltk import word_tokenize
17 from nltk.tokenize import wordpunct_tokenize
18 from nltk.stem import WordNetLemmatizer
19
20 from sklearn.feature_extraction.text import TfidfVectorizer
21 from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score, con
22 from sklearn.linear_model import LogisticRegression
23 from sklearn.naive_bayes import MultinomialNB, GaussianNB, BernoulliNB
24 from sklearn.svm import LinearSVC, SVC
25 from sklearn.ensemble import RandomForestClassifier
```

```
26 from sklearn.model_selection import cross_val_score
27 from sklearn.model_selection import train_test_split
28 from sklearn.metrics import precision_recall_curve
29
30
31 from torchtext.legacy import data
32 from torchtext.legacy.data import Field, LabelField, TabularDataset, Dataset
33
34 import copy
35 import re
36 from nltk.stem import PorterStemmer
37
38 # Download necessary NLTK datasets
39 nltk.download('punkt')
40 nltk.download('omw-1.4')
41 nltk.download('wordnet')
42 nltk.download('stopwords')
43 # Load spacy language model
44 spacy_en = spacy.load("en_core_web_sm")
45
46 # Load NLTK English stopwords
47 stop_words = set(stopwords.words('english'))
48
49 # Define NLTK WordNet lemmatizer
50 lemmatizer = WordNetLemmatizer()
51
52 # Define Porter Stemmer
53 porter_stemmer = PorterStemmer()
54
[nltk_data] Downloading package punkt to
[nltk_data]     C:\Users\Sunny\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package omw-1.4 to
[nltk_data]     C:\Users\Sunny\AppData\Roaming\nltk_data...
[nltk_data]   Package omw-1.4 is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]     C:\Users\Sunny\AppData\Roaming\nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]     C:\Users\Sunny\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

▼ Part 1: Text Classification

In this task, we will first go through the statistical method applied for this classification task, followed by the RNN algorithm.

Some parts of the code were assisted by ChatGPT, and those will be labelled for clarity

▼ Part 1A: Statistical Method

SVC algorithm was applied in this part of the task.

Week 4 tutorial materials were referenced for this part of the code.

Code are put into functions for the sake of accessibility in further iterations.

```
1 # Create empty df to store all evaluation values
2 metric_results_df = pd.DataFrame()
3
4 # Model training, building, and evaluating function
5
6 def evaluate_models(tokenizer, df_train, target_field, df_test, metric_results_df):
7
8     # Preprocess and training data
9     label_names = ["InformationTheory", "ComputationalLinguistics", "ComputerVision"]
10
11    trainDocs = df_train[target_field].tolist()
12    trainLabels = df_train[label_names].to_numpy()
13
14    vectorizer = TfidfVectorizer(analyzer='word', input='content', lowercase=True, min_d
15    x_train = vectorizer.fit_transform(trainDocs)
16
17    # Set up SVC model
18    model = LinearSVC()
19    CV = 5
20    entries = []
21
22    for idx, label_name in enumerate(label_names):
23        accuracies = cross_val_score(model, x_train, trainLabels[:, idx], scoring='accur
24        for fold_idx, accuracy in enumerate(accuracies):
25            entries.append((label_name, model.__class__.__name__, fold_idx, accuracy))
26
27    cv_df = pd.DataFrame(entries, columns=['label_name', 'model_name', 'fold_idx', 'accu
28    sns.boxplot(x='label_name', y='accuracy', data=cv_df)
29    sns.stripplot(x='label_name', y='accuracy', data=cv_df, size=8, jitter=True, edgecol
30    plt.show()
31
32    testDocs = df_test[target_field].tolist()
33    testLabels = df_test[label_names].to_numpy()
34    x_test = vectorizer.transform(testDocs)
35
36    # Create an empty list to store metric results
37    metric_results = []
38
39    for idx, label_name in enumerate(label_names):
40        print(f"Evaluating for {label_name}")
41        y_train, y_test = trainLabels[:, idx], testLabels[:, idx]
42        model.fit(x_train, y_train)
```

```
43     y_predict = model.predict(x_test)
44     y_probs = model.decision_function(x_test) # Get the decision function scores fo
45
46     # Create a dictionary to store the evaluation metrics for the current label
47     label_metric_results = {'label_name': label_name, 'model_name': model.__class__.
48
49     metrics = {
50         'Accuracy': accuracy_score,
51         'Macro Precision': precision_score,
52         'Macro Recall': recall_score,
53         'Macro F1 score': f1_score,
54         'MCC': matthews_corrcoef
55     }
56     print(model.__class__.__name__)
57     print(confusion_matrix(y_test, y_predict))
58     for metric_name, metric_func in metrics.items():
59         score = metric_func(y_test, y_predict, average='macro') if 'Macro' in metric
60         label_metric_results[metric_name] = score
61         print(f'{metric_name}: {score}')
62
63
64     # Get and store precision-recall curve values
65     precision, recall, _ = precision_recall_curve(y_test, y_probs)
66     label_metric_results['Recall'] = recall
67     label_metric_results['Precision'] = precision
68
69     # Append dictionary to list
70     metric_results.append(label_metric_results)
71
72     print('-----\n')
73 return metric_results_df.append(metric_results, ignore_index = True)
```

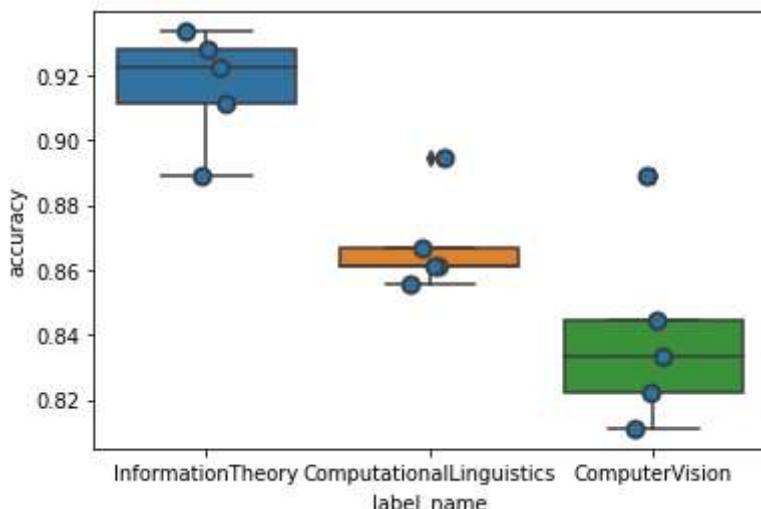
```
1 # Assisted by ChatGPT to convert tokenization methods from the RNN application format to
2
3 # Preprocessing method 1 (P1)
4 # Removing stopwords
5
6 def custom_tokenizer(text):
7
8     return [tok.text for tok in spacy_en.tokenizer(text) if tok.text.lower() not in STOP
9
10 # Preprocessing method 2 (P2)
11 # Stemming tokens
12 def custom_tokenizer_P2(text):
13
14     ps = PorterStemmer()
15
16     return [ps.stem(tok.text) for tok in spacy_en.tokenizer(text)]
```

```
1 # Splitting data
2
3 split_ratio = 0.9
4 SEED = 1
5 train_data = pd.read_csv('train.csv')
6 train_1000 = train_data.head(1000)
7 train_data, validation = train_test_split(train_data, test_size=1-split_ratio, random_st
8 train_900, valid_100 = train_test_split(train_1000, test_size=1-split_ratio, random_stat
9 test_data = pd.read_csv("test.csv")
10
11 len(train_data), len(validation), len(train_900), len(valid_100)
```

(112500, 12500, 900, 100)

▼ P1 Title 1000

```
1 # Defining parameters
2
3 tokenizer = custom_tokenizer
4 df_train = train_900
5 target_field = "title"
6 df_test = test_data
7 metric_results_df = evaluate_models(tokenizer, df_train, target_field, df_test, metric_r
```

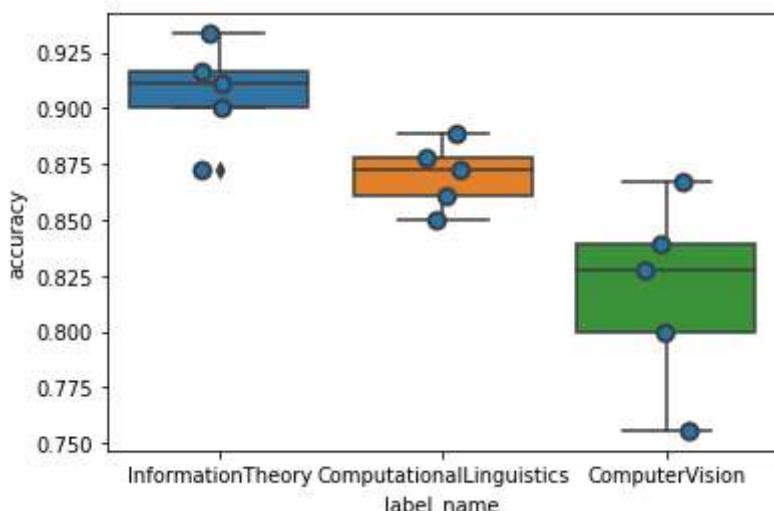


```
Evaluating for InformationTheory
LinearSVC
[[9249 338]
 [2616 5863]]
Accuracy: 0.8364884313074283
Macro Precision: 0.8625061289612967
Macro Recall: 0.8281084875655896
```

▼ P2 Title 1000

Evaluating for ComputationalLinguistics

```
1 # Defining parameters
2
3 tokenizer = custom_tokenizer_P2
4 df_train = train_900
5 target_field = "title"
6 df_test = test_data
7 metric_results_df = evaluate_models(tokenizer, df_train, target_field, df_test, metric_r
```



Evaluating for InformationTheory

LinearSVC

```
[[9245  342]
 [2424 6055]]
```

Accuracy: 0.8468947193623381

Macro Precision: 0.8694037784149105

Macro Recall: 0.8392219616007361

Macro F1 score: 0.8419674781353987

MCC: 0.707982695654716

Evaluating for ComputationalLinguistics

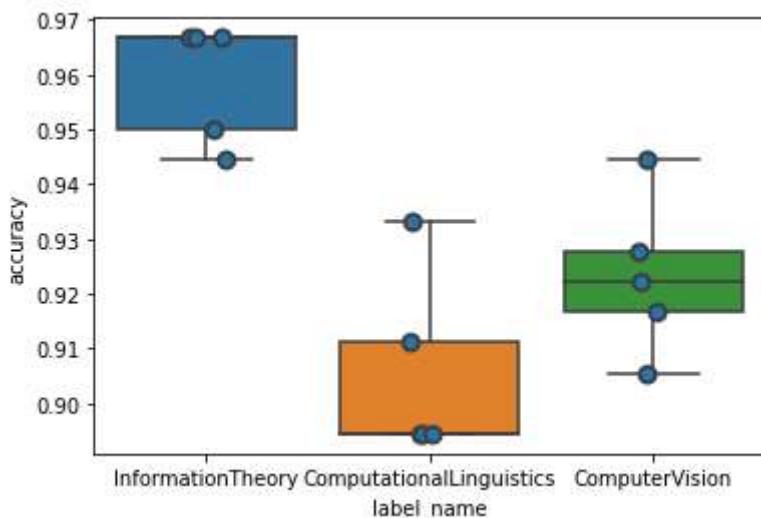
LinearSVC

```
[[14082   616]]
```

▼ P1 Abstract 1000

Macro F1 score: 0.808611702132926

```
1 tokenizer = custom_tokenizer
2 df_train = train_900
3 target_field = "abstract"
4 df_test = test_data
5 metric_results_df = evaluate_models(tokenizer, df_train, target_field, df_test, metric_r
```



Evaluating for InformationTheory

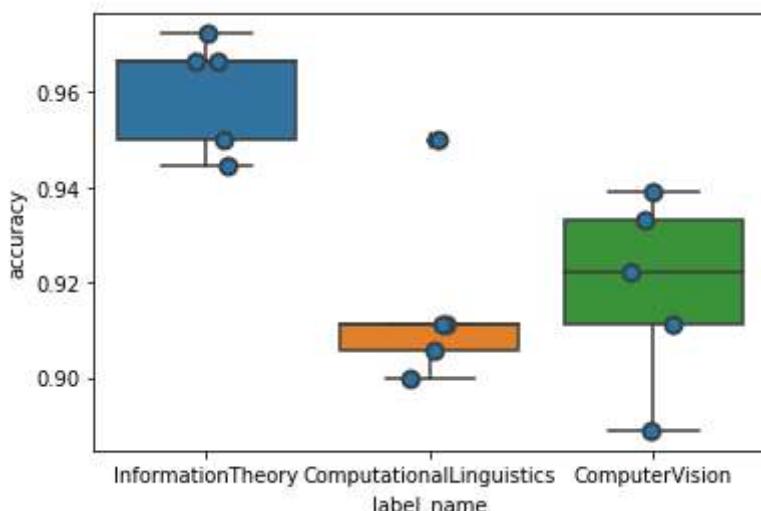
```
LinearSVC
[[9524  63]
 [1297 7182]]
Accuracy: 0.9247204693900144
Macro Precision: 0.9357224077176827
Macro Recall: 0.9202312247810023
Macro F1 score: 0.9234337400519605
MCC: 0.8558134400901379
-----
```

Evaluating for ComputationalLinguistics

```
LinearSVC
[[14418   280]
 [  922 2446]]
Accuracy: 0.9334661795638215
Macro Precision: 0.9185905486880406
Macro Recall: 0.8535984099829053
Macro F1 score: 0.8813704151090274
```

▼ P2 Abstract 1000

```
LinearSVC
1 tokenizer = custom_tokenizer_P2
2 df_train = train_900
3 target_field = "abstract"
4 df_test = test_data
5 metric_results_df = evaluate_models(tokenizer, df_train, target_field, df_test, metric_r
```



Evaluating for InformationTheory

LinearSVC

```
[[9513  74]
 [1287 7192]]
```

Accuracy: 0.9246651167939777

Macro Precision: 0.9353244563721442

Macro Recall: 0.920247223418344

Macro F1 score: 0.9234008159951866

MCC: 0.855438820901866

Evaluating for ComputationalLinguistics

LinearSVC

```
[[14435  263]
 [ 975 2393]]
```

Accuracy: 0.9314734861064984

Macro Precision: 0.9188541560792163

Macro Recall: 0.8463085489356736

Macro F1 score: 0.8766850361275769

MCC: 0.7617158919373275

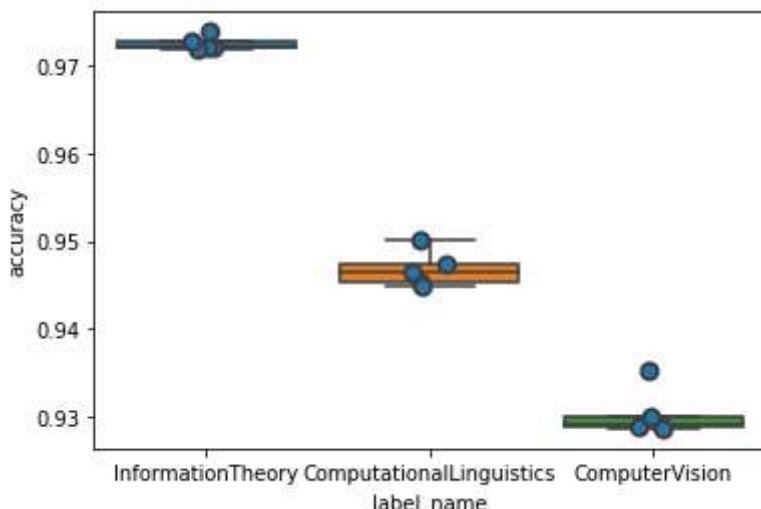
Evaluating for ComputerVision

LinearSVC

▼ P1 Title ALL

Macro Recall: 0.927075682062057

```
1 tokenizer = custom_tokenizer
2 df_train = train_data
3 target_field = "title"
4 df_test = test_data
5 metric_results_df = evaluate_models(tokenizer, df_train, target_field, df_test, metric_r
```



Evaluating for InformationTheory

LinearSVC

```
[[9407 180]
 [ 901 7578]]
```

Accuracy: 0.9401638436842688

Macro Precision: 0.9446951526397626

Macro Recall: 0.9374810219931994

Macro F1 score: 0.9395442009454607

MCC: 0.8821466768112096

Evaluating for ComputationalLinguistics

LinearSVC

```
[[14318 380]
 [ 514 2854]]
```

Accuracy: 0.9505147791431419

Macro Precision: 0.9239218267477629

Macro Recall: 0.9107666578644824

Macro F1 score: 0.917156095808995

MCC: 0.8345848116734086

Evaluating for ComputerVision

LinearSVC

```
[[10957 890]
 [ 406 5813]]
```

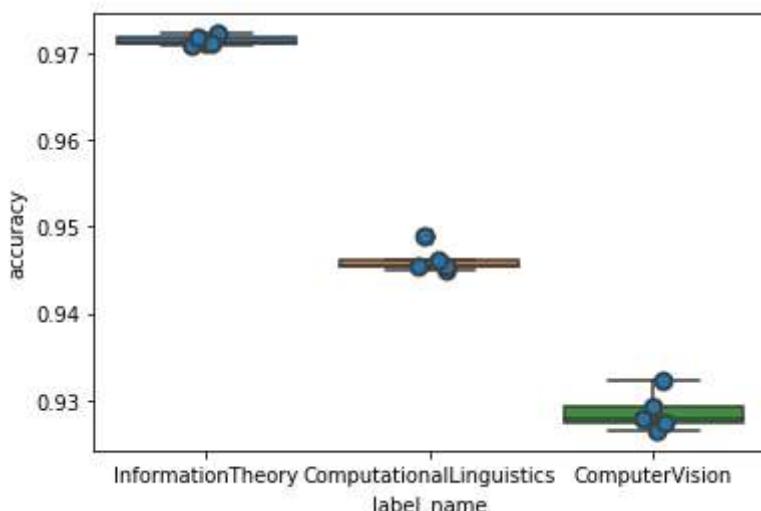
Accuracy: 0.9282630355363667

Macro Precision: 0.9157468151649284

Macro Recall: 0.9297958441100067

▼ P2 Title ALL

```
1 tokenizer = custom_tokenizer_P2
2 df_train = train_data
3 target_field = "title"
4 df_test = test_data
5 metric_results_df = evaluate_models(tokenizer, df_train, target_field, df_test, metric_r
```



Evaluating for InformationTheory

LinearSVC

```
[[9408 179]
 [ 890 7589]]
```

Accuracy: 0.9408280748367098

Macro Precision: 0.9452660985835291

Macro Recall: 0.9381818373504347

Macro F1 score: 0.9402220613050182

MCC: 0.883419531564003

Evaluating for ComputationalLinguistics

LinearSVC

```
[[14302 396]
 [ 520 2848]]
```

Accuracy: 0.9492970220303332

Macro Precision: 0.9214227493007445

Macro Recall: 0.9093316297820668

Macro F1 score: 0.9152170972707914

MCC: 0.8306663850150988

Evaluating for ComputerVision

LinearSVC

```
[[10981 866]
 [ 431 5788]]
```

Accuracy: 0.9282076829403298

Macro Precision: 0.9160427288188252

Macro Recall: 0.928798789323482

Macro F1 score: 0.9217415622356171

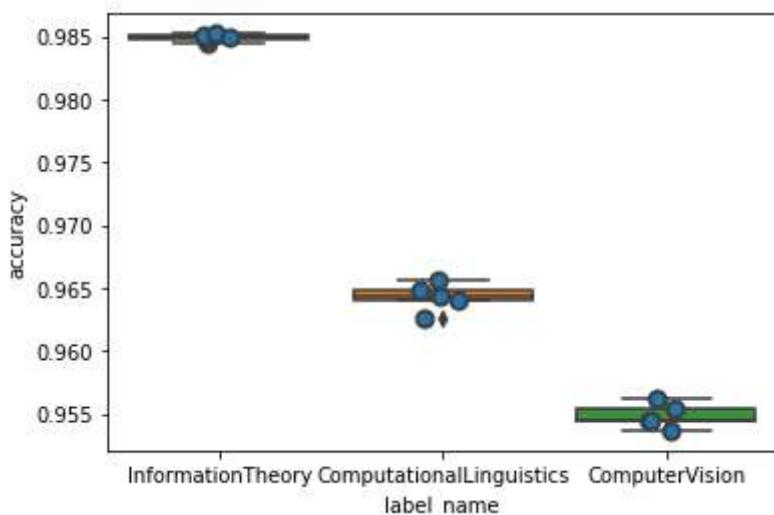
MCC: 0.8447452122962283

▼ P1 Abstract ALL

```
1 tokenizer = custom_tokenizer
2 df_train = train_data
```

https://colab.research.google.com/drive/15_VNNDEW8oOURT2Me9nFdCi-bo3mhAD9#scrollTo=fb95be87

```
3 target_field = "abstract"
4 df_test = test_data
5 metric_results_df = evaluate_models(tokenizer, df_train, target_field, df_test, metric_r
```



Evaluating for InformationTheory

LinearSVC

```
[[9491  96]
 [ 574 7905]]
```

Accuracy: 0.9629137606553747

Macro Precision: 0.9654860951620987

Macro Recall: 0.9611448888142682

Macro F1 score: 0.9626257190651313

MCC: 0.9266208147848048

Evaluating for ComputationalLinguistics

LinearSVC

```
[[14468  230]
 [ 314 3054]]
```

Accuracy: 0.9698881877560057

Macro Precision: 0.9543607040264099

Macro Recall: 0.9455606043319029

Macro F1 score: 0.9498834477913826

MCC: 0.8998782803705756

Evaluating for ComputerVision

LinearSVC

```
[[11296  551]
 [ 228 5991]]
```

Accuracy: 0.9568803276873685

Macro Precision: 0.9479950977057862

Macro Recall: 0.9584142461829718

Macro F1 score: 0.9528113601579944

MCC: 0.9063494580093591

▼ P2 Abstract ALL

```
1 tokenizer = custom_tokenizer_P2
2 df_train = train_data
3 target_field = "abstract"
4 df_test = test_data
5 metric_results_df = evaluate_models(tokenizer, df_train, target_field, df_test, metric_r
```



```
1 # Saving all evaluation metrics into dataframe  
2  
3 metric_results_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 24 entries, 0 to 23  
Data columns (total 12 columns):  
 #   Column           Non-Null Count  Dtype     
---  --  
 0   label_name        24 non-null    object    
 1   model_name        24 non-null    object    
 2   target_field      24 non-null    object    
 3   train_size         24 non-null    int64     
 4   tokenizer          24 non-null    object    
 5   Accuracy          24 non-null    float64  
 6   Macro Precision   24 non-null    float64  
 7   Macro Recall      24 non-null    float64  
 8   Macro F1 score    24 non-null    float64  
 9   MCC                24 non-null    float64  
 10  Recall              24 non-null    object    
 11  Precision           24 non-null    object    
 dtypes: float64(5), int64(1), object(6)  
 memory usage: 2.4+ KB
```

```
1 metric_results_df
```

	label_name	model_name	target_field	train_size	tokenizer	Acc
0	InformationTheory	LinearSVC	title	900	custom_tokenizer	0.81
1	ComputationalLinguistics	LinearSVC	title	900	custom_tokenizer	0.81
2	ComputerVision	LinearSVC	title	900	custom_tokenizer	0.81
3	InformationTheory	LinearSVC	title	900	custom_tokenizer_P2	0.81

Libraries imports should have been integrated to the first part of the code.

This part is left here for error handling purpose.

```
1 # # Import Libraries
2 # import torch
3 # import torchtext
4
5 # from torchtext.legacy import data
6 # import torch.nn as nn
7 # from torchtext.legacy.data import Field, LabelField, TabularDataset, Dataset
8 # import copy
9
10 # import torch.optim as optim
11 # import collections
12 # import random
13 # import time
14 # import pandas as pd
15 # import re
16 # import spacy
17 # from nltk.stem import PorterStemmer
18
19 # from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score, c
20 # import numpy as np
```

```
21 # spacy_en = spacy.load("en_core_web_sm")  
  
1 # import torch  
2 # import torchtext  
3 # import torch.nn as nn  
4 # import torch.optim as optim  
5  
6 # from torchtext.legacy import data  
7 # from torchtext.legacy.data import Field, LabelField, TabularDataset  
8 # from sklearn.metrics import precision_score, recall_score, f1_score, accuracy_score, c  
9 # from sklearn.metrics import precision_recall_curve  
10  
11 # import spacy  
12 # from spacy.lang.en.stop_words import STOP_WORDS  
13 # import numpy as np  
14 # import random  
15 # import time  
16 # import copy  
17 # import collections  
18  
19 # from nltk.stem import PorterStemmer  
20
```

▼ Part 1B: RNN Method

Code was referenced to Week 4 tutorial material on Building RNN models, and converted into functions for the sake of accessibility.

Summarising codes into functions were assisted by ChatGPT.

```
1 # I accidentally re-run this cell. Hopefully it will not affect the output of the later  
2  
3 spacy_en = spacy.load("en_core_web_sm")  
4 SEED = 1  
5 torch.manual_seed(SEED)  
6 torch.backends.cudnn.deterministic = True  
7 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')  
8 device
```

```

1 def load_data(path, train_file, test_file, target_field, tokenizer, MAX_VOCAB_SIZE):
2     TEXT = data.Field(sequential=True, tokenize=tokenizer, lower=True)
3     LABEL = data.LabelField(dtype=torch.float, use_vocab=False, preprocessing=int)
4     if target_field == "title":
5         train_datafield = [
6             (target_field, TEXT),
7             ("abstract", None),
8             ("InformationTheory", LABEL),
9             ("ComputationalLinguistics", LABEL),
10            ("ComputerVision", LABEL)
11        ]
12    else:
13        train_datafield = [
14            ("title", None),
15            (target_field, TEXT),
16            ("InformationTheory", LABEL),
17            ("ComputationalLinguistics", LABEL),
18            ("ComputerVision", LABEL)
19        ]
20
21    train_data, test_data = TabularDataset.splits(
22        path=path,
23        train=train_file, test=test_file, format="csv",
24        skip_header=True, fields=train_datafield
25    )
26    train_data, valid_data = train_data.split(split_ratio=0.9, random_state=random.getstate())
27
28    MAX_VOCAB_SIZE = MAX_VOCAB_SIZE
29
30    TEXT.build_vocab(train_data, max_size=MAX_VOCAB_SIZE)
31
32    return TEXT, LABEL, train_data, valid_data, test_data
33

```

```

1 def generate_label_iterator(train_data, valid_data, test_data, label, target_field):
2
3     label_attr = f"label_{label}"
4     BATCH_SIZE = 64
5     iterators = data.BucketIterator.splits(
6         (train_data, valid_data, test_data),
7         batch_size = BATCH_SIZE,
8         device = device,
9         sort_key = lambda x: len(getattr(x, target_field)),
10        sort_within_batch = False)
11    return iterators[0], iterators[1], iterators[2]

```

```

1
2 class RNN(nn.Module):

```

```

3     def __init__(self, input_dim, embedding_dim, hidden_dim, output_dim, pos_weight = None):
4
5         super(RNN, self).__init__()
6
7         self.embedding = nn.Embedding(input_dim, embedding_dim)
8
9         self.rnn = nn.RNN(embedding_dim, hidden_dim)
10
11        self.fc = nn.Linear(hidden_dim, output_dim)
12
13        self.pos_weight = pos_weight
14
15        if self.pos_weight is not None:
16            self.criterion = nn.BCEWithLogitsLoss(pos_weight=self.pos_weight)
17        else:
18            self.criterion = nn.BCEWithLogitsLoss()
19
20    def forward(self, text):
21
22        embedded = self.embedding(text)
23
24        output, hidden = self.rnn(embedded)
25
26        assert torch.equal(output[-1,:,:], hidden.squeeze(0))
27
28        return self.fc(hidden.squeeze(0))

```

```

1 #·Assisted·by·ChatGPT
2
3 def get_pos_weight(train_data):
4     pos_weights = []
5     for field in ["InformationTheory", "ComputationalLinguistics", "ComputerVision"]:
6         negative_count, positive_count = 0, 0
7         for example in train_data.examples:
8             if getattr(example, field) == 0:
9                 negative_count += 1
10            else:
11                positive_count += 1
12
13    pos_weights.append(torch.tensor([negative_count / positive_count]).to(device))
14
15    return pos_weights

```

```

1 def generate_model_and_optimizer(pos_weight, embedding_dim=100, hidden_dim=256, output_d
2 INPUT_DIM = len(TEXT.vocab)
3
4 model = RNN(INPUT_DIM, embedding_dim, hidden_dim, output_dim, pos_weight)
5
6 optimizer = optim.SGD(model.parameters(), lr=lr)
7

```

```
8     model = model.to(device)
9
10    return model, optimizer

1 # Evaluation functions
2
3 def binary_accuracy(preds, y):
4
5     #round predictions to the closest integer
6     rounded_preds = torch.round(torch.sigmoid(preds))
7     correct = (rounded_preds == y).float() #convert into float for division
8     acc = correct.sum() / len(correct)
9     return acc
10
11 def train(model, iterator, optimizer, criterion, label_field, target_field):
12
13     epoch_loss = 0
14     epoch_acc = 0
15
16     model.train()
17
18     for batch in iterator:
19
20         optimizer.zero_grad()
21
22         predictions = model(getattr(batch, target_field)).squeeze(1)
23
24         # Use the specific label field
25         loss = criterion(predictions, getattr(batch, label_field))
26
27         acc = binary_accuracy(predictions, getattr(batch, label_field))
28
29         loss.backward()
30
31         optimizer.step()
32
33         epoch_loss += loss.item()
34         epoch_acc += acc.item()
35
36     return epoch_loss / len(iterator), epoch_acc / len(iterator)
37
38
39 def evaluate(model, iterator, criterion, label_field, target_field):
40
41     epoch_loss = 0
42     epoch_acc = 0
43
44     model.eval()
45
46     with torch.no_grad():
47
```

```
48     for batch in iterator:
49
50         predictions = model(getattr(batch, target_field)).squeeze(1)
51
52         # Use the specific label field
53         loss = criterion(predictions, getattr(batch, label_field))
54
55         acc = binary_accuracy(predictions, getattr(batch, label_field))
56
57         epoch_loss += loss.item()
58         epoch_acc += acc.item()
59
60     return epoch_loss / len(iterator), epoch_acc / len(iterator)
61
62 def epoch_time(start_time, end_time):
63     elapsed_time = end_time - start_time
64     elapsed_mins = int(elapsed_time / 60)
65     elapsed_secs = int(elapsed_time - (elapsed_mins * 60))
66     return elapsed_mins, elapsed_secs
67
68
69 # Assisted by ChatGPT to build up the train_loop, creating early stopping and error hand
70
71
72 rnn_results_df = pd.DataFrame()
73
74
75 def train_loop(ModelsList, OptimizersList, IteratorsList, N_EPOCHS, label_names, target_
76   result = []
77   best_model_states = []
78   best_epochs = [0, 0, 0]
79   for idx, (label_name, model, optimizer, iterators_) in enumerate(zip(label_names, Mo
80       print(f"Training model for {preprocess}_{label_name}, using {train_size} data of
81
82       bad_epochs = 0
83       best_model_state = None
84       best_epoch = 0
85
86       for epoch in range(N_EPOCHS):
87           start_time = time.time()
88
89           train_iterator = iterators_[0]
90           test_iterator = iterators_[2] if len(iterators_) > 2 else iterators_[1]
91           valid_iterator = iterators_[1] if len(iterators_) > 2 else None
92
93           train_loss, train_acc = train(model, train_iterator, optimizer, model.criter
94
95           if valid_iterator:
96               valid_loss, valid_acc = evaluate(model, valid_iterator, model.criterion,
97
98               test_loss, test_acc = evaluate(model, test_iterator, model.criterion, label_
99
100              end_time = time.time()
101              epoch_mins, epoch_secs = epoch_time(start_time, end_time)
102
103      result.append({
104          "label_name": label_name,
105          "model": model,
106          "optimizer": optimizer,
107          "train_loss": train_loss,
108          "train_acc": train_acc,
109          "valid_loss": valid_loss,
110          "valid_acc": valid_acc,
111          "test_loss": test_loss,
112          "test_acc": test_acc,
113          "epoch_time": f"{epoch_mins}m {epoch_secs}s"
114      })
115
116  best_model_state = result[best_epochs[0]]["model"]
117  best_epochs[0] = result.index(best_model_state)
118
119  best_model_states.append(best_model_state)
120  best_epochs[1] = result.index(best_model_states[-1])
121
122  best_model_state = result[best_epochs[1]]["model"]
123  best_epochs[2] = result.index(best_model_state)
124
125  best_model_states.append(best_model_state)
```

```
32
33     if test_loss < bestLossesList[idx]:
34         bestLossesList[idx] = test_loss
35         best_model_state = copy.deepcopy(model.state_dict())
36 #             best_model_state = model.state_dict()
37         bad_epochs = 0
38         best_epochs[idx] = epoch + 1
39
40         print(f'Epoch: {epoch+1:02} | Epoch Time: {epoch_mins}m {epoch_secs}s')
41         print(f'{label_name} Train Loss: {train_loss:.3f} | Train Acc: {train_ac}
42         print(f'{label_name} Valid Loss: {valid_loss:.3f} | Valid Acc: {valid_ac
43         print(f'{label_name} Test Loss: {test_loss:.3f} | Test Acc: {test_acc*10
44 else:
45     bad_epochs += 1
46
47 if bad_epochs > patience:
48     print(f"Early stopping at epoch {epoch+1} for {label_name} model.")
49 #
50     model.load_state_dict(best_model_state)
51     break
52
53     # For predictions and ground truth collection
54     model.eval()
55     y_predict, y_test = [], []
56     with torch.no_grad():
57         for batch in test_iterator:
58             predictions = model(getattr(batch, target_field)).squeeze(1)
59             rounded_preds = torch.round(torch.sigmoid(predictions))
60             y_predict += rounded_preds.tolist()
61             y_test += getattr(batch, label_name).tolist()
62     print(f"Best performing epoch for {label_name} model: {best_epochs[idx]}")
63     best_model_states.append(best_model_state)
64
65     model.load_state_dict(best_model_state)
66     # Call 'evaluate_model' and store its return value
67     rnn_eval_metrics = evaluate_model(model, test_iterator, label_name, target_field
68     print(rnn_eval_metrics)
69     # Append the evaluation metrics to the 'rnn_metric_results' list
70     rnn_metric_results.append({
71         "accuracy": rnn_eval_metrics["accuracy"],
72         'label_name': label_name,
73         'model_name': model.__class__.__name__,
74         'target_field': target_field,
75         'train_size': train_size,
76         'tokenizer': preprocess,
77         **rnn_eval_metrics
78     })
79     result.append((y_predict, y_test))
80
81     return rnn_metric_results, best_model_states
82
```

```
1 # Assisted by ChatGPT to return evaluation metrics to df
2
3
4 def evaluate_model(model, iterator, label_field, target_field):
5
6
7     y_probs = []
8     y_predict = []
9     y_test = []
10
11    model.eval()
12    with torch.no_grad():
13        for batch in iterator:
14            predictions = model(getattr(batch, target_field)).squeeze(1)
15            probs = torch.sigmoid(predictions)
16            rounded_preds = torch.round(torch.sigmoid(predictions))
17
18            y_predict += rounded_preds.tolist()
19            y_test += getattr(batch, label_field).tolist()
20            y_probs += probs.tolist()
21
22    y_predict = np.asarray(y_predict)
23    y_test = np.asarray(y_test)
24    y_probs = np.asarray(y_probs)
25    # Compute metrics
26    recall_macro = recall_score(y_test, y_predict, average='macro')
27    precision_macro = precision_score(y_test, y_predict, average='macro')
28    #    recall = recall_score(y_test, y_predict, average=None)
29    #    precision = precision_score(y_test, y_predict, average=None)
30    f1score = f1_score(y_test, y_predict, average='macro')
31    accuracy = accuracy_score(y_test, y_predict)
32    matthews = matthews_corrcoef(y_test, y_predict)
33    #    precision_curve, recall_curve, _ = precision_recall_curve(y_test, y_probs)
34    precision, recall, _ = precision_recall_curve(y_test, y_probs)
35
36    # Print metrics
37
38    print(f"{label_field}:")
39    print(confusion_matrix(y_test, y_predict))
40    print('Accuracy:', accuracy)
41    print('Macro Precision:', precision_macro)
42    print('Macro Recall:', recall_macro)
43    print('Macro F1 score:', f1score)
44    print('MCC:', matthews)
45    print("\n")
46    return {
47        "accuracy": accuracy,
48        'Macro Recall': recall_macro,
49        'Macro Precision': precision_macro,
50        'Macro F1 score': f1score,
```

```

51     'MCC': matthews,
52     'Recall': recall,
53     'Precision': precision
54 }
55

```

```
1
```

P1 Title 1000

```

1 def custom_tokenizer(text):
2     return [tok.text for tok in spacy_en.tokenizer(text) if tok.text.lower() not in STOP
3
4 path = './'
5 train_file = "train_1000.csv"
6 test_file = "test.csv"
7 target_field = "title"
8 tokenizer = custom_tokenizer
9 MAX_VOCAB_SIZE = 5000
10
11 TEXT, LABEL, train_data, valid_data, test_data = load_data(path, train_file, test_file,
12
13
14 train_iterator_IT, validation_IT, test_iterator_IT = generate_label_iterator(train_data,
15 train_iterator_CL, validation_CL, test_iterator_CL = generate_label_iterator(train_data,
16 train_iterator_CV, validation_CV, test_iterator_CV = generate_label_iterator(train_data,
17
18 pos_weights = get_pos_weight(train_data)
19
20 model_IT, optimizer_IT = generate_model_and_optimizer(pos_weight=pos_weights[0])
21 model_CL, optimizer_CL = generate_model_and_optimizer(pos_weight=pos_weights[1])
22 model_CV, optimizer_CV = generate_model_and_optimizer(pos_weight=pos_weights[2])
23
24 N_EPOCHS = 2
25 label_names = ["InformationTheory", "ComputationalLinguistics", "ComputerVision"]
26 models = [model_IT, model_CL, model_CV]
27 optimizers = [optimizer_IT, optimizer_CL, optimizer_CV]
28 iterators = [(train_iterator_IT, validation_IT, test_iterator_IT), (train_iterator_CL, v
29 model_file_names = ["RNN_model_IT_T_P1_1000", "RNN_model_CL_T_P1_1000.pt", "RNN_model_CV
30 best_valid_losses = [float("inf"), float("inf"), float("inf")]

```



```

1 # Making sure data are loaded correctly
2
3 # print(TEXT.vocab.freqs.most_common(200))
4 len(train_data), len(valid_data), len(test_data), len(TEXT.vocab)
5
6

```

```
(900, 100, 18066, 2497)
```

```
1 # Run model and store output
2
3 rnn_results, best_model_states = train_loop(models, optimizers, iterators, N_EPOCHS, lab
4 rnn_results_df = rnn_results_df.append(pd.DataFrame(rnn_results), ignore_index = True)
5
```

Training model for P1_InformationTheory, using 900 data of title...

Epoch: 01 | Epoch Time: 0m 2s
InformationTheory Train Loss: 1.057 | Train Acc: 40.73%
InformationTheory Valid Loss: 0.978 | Valid Acc: 22.05%
InformationTheory Test Loss: 1.323 | Test Acc: 51.96%
Best performing epoch for InformationTheory model: 1
InformationTheory:
[[6510 3077]
 [5607 2872]]
Accuracy: 0.5193180560168272
Macro Precision: 0.5100159559607327
Macro Recall: 0.5088818640320529
Macro F1 score: 0.4990020990638365
MCC: 0.01886375985791359

Training model for P1_ComputationalLinguistics, using 900 data of title...

Epoch: 01 | Epoch Time: 0m 0s
ComputationalLinguistics Train Loss: 1.085 | Train Acc: 49.38%
ComputationalLinguistics Valid Loss: 1.272 | Valid Acc: 67.36%
ComputationalLinguistics Test Loss: 1.002 | Test Acc: 54.41%
Best performing epoch for ComputationalLinguistics model: 1
ComputationalLinguistics:
[[8274 6424]
 [1828 1540]]
Accuracy: 0.5432303775047049
Macro Precision: 0.5062079496319858
Macro Recall: 0.5100891940312787
Macro F1 score: 0.46952737323935645
MCC: 0.015828229006873723

Training model for P1_ComputerVision, using 900 data of title...

Epoch: 01 | Epoch Time: 0m 0s
ComputerVision Train Loss: 0.678 | Train Acc: 48.44%
ComputerVision Valid Loss: 0.678 | Valid Acc: 49.05%
ComputerVision Test Loss: 0.693 | Test Acc: 50.67%
Best performing epoch for ComputerVision model: 1
ComputerVision:
[[5748 6099]
 [2802 3417]]
Accuracy: 0.5073065426768515
Macro Precision: 0.5156800734497025
Macro Recall: 0.5173156857506777
Macro F1 score: 0.4989653472003046
MCC: 0.032955195305298446

```

1 rnn_results_df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3 entries, 0 to 2
Data columns (total 12 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   accuracy        3 non-null      float64
 1   label_name      3 non-null      object 
 2   model_name      3 non-null      object 
 3   target_field    3 non-null      object 
 4   train_size      3 non-null      int64  
 5   tokenizer       3 non-null      object 
 6   Macro Recall    3 non-null      float64
 7   Macro Precision 3 non-null      float64
 8   Macro F1 score  3 non-null      float64
 9   MCC              3 non-null      float64
 10  Recall          3 non-null      object 
 11  Precision        3 non-null      object 

dtypes: float64(5), int64(1), object(6)
memory usage: 416.0+ bytes

```

P2 Title 1000

```

1 def custom_tokenizer_P2(text):
2     ps = PorterStemmer()
3     return [ps.stem(tok.text) for tok in spacy_en.tokenizer(text)]
4
5 path = './'
6 train_file = "train_1000.csv"
7 test_file = "test.csv"
8 target_field = "title"
9 tokenizer = custom_tokenizer_P2
10 MAX_VOCAB_SIZE = 5000
11
12 TEXT, LABEL, train_data, valid_data, test_data = load_data(path, train_file, test_file,
13
14
15 train_iterator_IT, validation_IT, test_iterator_IT = generate_label_iterator(train_data,
16 train_iterator_CL, validation_CL, test_iterator_CL = generate_label_iterator(train_data,
17 train_iterator_CV, validation_CV, test_iterator_CV = generate_label_iterator(train_data,
18
19 pos_weights = get_pos_weight(train_data)
20
21 model_IT, optimizer_IT = generate_model_and_optimizer(pos_weight=pos_weights[0])
22 model_CL, optimizer_CL = generate_model_and_optimizer(pos_weight=pos_weights[1])
23 model_CV, optimizer_CV = generate_model_and_optimizer(pos_weight=pos_weights[2])

```

24

```
25 N_EPOCHS = 30
26 label_names = ["InformationTheory", "ComputationalLinguistics", "ComputerVision"]
27 models = [model_IT, model_CL, model_CV]
28 optimizers = [optimizer_IT, optimizer_CL, optimizer_CV]
29 iterators = [(train_iterator_IT, validation_IT, test_iterator_IT), (train_iterator_CL, v
30 model_file_names = ["RNN_model_IT_T_P2_1000", "RNN_model_CL_T_P2_1000.pt", "RNN_model_CV
31 best_valid_losses = [float("inf"), float("inf"), float("inf")]
```

```
1 print(TEXT.vocab.freqs.most_common(200))
2 len(train_data), len(valid_data), len(test_data), len(TEXT.vocab)
3
4
```

```
[('-', 602), ('\n ', 444), ('for', 353), ('of', 247), (':', 221), ('and', 202), ('a',
(900, 100, 18066, 2072)
```



```
1 rnn_results, best_model_states = train_loop(models, optimizers, iterators, N_EPOCHS, lab
2 rnn_results_df = rnn_results_df.append(pd.DataFrame(rnn_results), ignore_index = True)
3
```



```

ComputerVision Test Loss: 0.682 | Test Acc: 52.54%
Epoch: 17 | Epoch Time: 0m 0s
ComputerVision Train Loss: 0.680 | Train Acc: 48.44%
ComputerVision Valid Loss: 0.684 | Valid Acc: 51.04%
ComputerVision Test Loss: 0.682 | Test Acc: 52.75%
Epoch: 19 | Epoch Time: 0m 0s
ComputerVision Train Loss: 0.677 | Train Acc: 50.00%
ComputerVision Valid Loss: 0.684 | Valid Acc: 51.04%
ComputerVision Test Loss: 0.682 | Test Acc: 52.52%
Epoch: 21 | Epoch Time: 0m 0s
ComputerVision Train Loss: 0.677 | Train Acc: 52.19%
ComputerVision Valid Loss: 0.686 | Valid Acc: 35.50%
ComputerVision Test Loss: 0.682 | Test Acc: 52.27%
Epoch: 25 | Epoch Time: 0m 0s
ComputerVision Train Loss: 0.673 | Train Acc: 48.54%
ComputerVision Valid Loss: 0.680 | Valid Acc: 51.82%
ComputerVision Test Loss: 0.678 | Test Acc: 53.37%
Epoch: 27 | Epoch Time: 0m 0s
ComputerVision Train Loss: 0.677 | Train Acc: 52.60%
ComputerVision Valid Loss: 0.685 | Valid Acc: 51.82%
ComputerVision Test Loss: 0.678 | Test Acc: 53.41%
Best performing epoch for ComputerVision model: 27
ComputerVision:
[[7128 4719]
 [3706 2513]]
Accuracy: 0.5336543783903465
Macro Precision: 0.50270607496312

```

P1 Title ALL

```

1 path = './'
2 train_file = "train.csv"
3 test_file = "test.csv"
4 target_field = "title"
5 tokenizer = custom_tokenizer
6 MAX_VOCAB_SIZE = 5000
7
8 TEXT, LABEL, train_data, valid_data, test_data = load_data(path, train_file, test_file,
9
10
11 train_iterator_IT, validation_IT, test_iterator_IT = generate_label_iterator(train_data,
12 train_iterator_CL, validation_CL, test_iterator_CL = generate_label_iterator(train_data,
13 train_iterator_CV, validation_CV, test_iterator_CV = generate_label_iterator(train_data,
14
15 pos_weights = get_pos_weight(train_data)
16
17 model_IT, optimizer_IT = generate_model_and_optimizer(pos_weight=pos_weights[0])
18 model_CL, optimizer_CL = generate_model_and_optimizer(pos_weight=pos_weights[1])
19 model_CV, optimizer_CV = generate_model_and_optimizer(pos_weight=pos_weights[2])
20
21 N_EPOCHS = 30

```

```
22 label_names = ["InformationTheory", "ComputationalLinguistics", "ComputerVision"]
23 models = [model_IT, model_CL, model_CV]
24 optimizers = [optimizer_IT, optimizer_CL, optimizer_CV]
25 iterators = [(train_iterator_IT, validation_IT, test_iterator_IT), (train_iterator_CL, v
26 model_file_names = ["RNN_model_IT_T_P1_ALL", "RNN_model_CL_T_P1_ALL.pt", "RNN_model_CV_T
27 best valid losses = [float("inf"), float("inf"), float("inf")]

1 print(TEXT.vocab.freqs.most_common(50))
2 len(train_data), len(valid_data), len(test_data), len(TEXT.vocab)

[('-', 75468), ('\n ', 55134), (':', 27878), ('learning', 15463), ('based', 10302), ('112500, 12500, 18066, 5002)
```



```
1 rnn_results, best_model_states = train_loop(models, optimizers, iterators, N_EPOCHS, lab
2 rnn_results_df = rnn_results_df.append(pd.DataFrame(rnn_results), ignore_index = True)
3
```



```
Epoch: 13 | Epoch Time: 0m 7s
ComputerVision Train Loss: 0.242 | Train Acc: 89.64%
ComputerVision Valid Loss: 0.318 | Valid Acc: 86.01%
ComputerVision Test Loss: 0.345 | Test Acc: 84.18%
Epoch: 14 | Epoch Time: 0m 7s
ComputerVision Train Loss: 0.234 | Train Acc: 90.07%
ComputerVision Valid Loss: 0.314 | Valid Acc: 86.07%
ComputerVision Test Loss: 0.336 | Test Acc: 84.60%
Epoch: 15 | Epoch Time: 0m 7s
ComputerVision Train Loss: 0.222 | Train Acc: 90.65%
ComputerVision Valid Loss: 0.304 | Valid Acc: 86.82%
ComputerVision Test Loss: 0.330 | Test Acc: 85.10%
Epoch: 18 | Epoch Time: 0m 7s
ComputerVision Train Loss: 0.202 | Train Acc: 91.54%
ComputerVision Valid Loss: 0.295 | Valid Acc: 87.06%
ComputerVision Test Loss: 0.305 | Test Acc: 86.31%
Epoch: 26 | Epoch Time: 0m 7s
ComputerVision Train Loss: 0.164 | Train Acc: 93.27%
ComputerVision Valid Loss: 0.275 | Valid Acc: 87.83%
ComputerVision Test Loss: 0.298 | Test Acc: 86.78%
Best performing epoch for ComputerVision model: 26
ComputerVision:
```

P2 Title ALL

```
1 path = './'
2 train_file = "train.csv"
3 test_file = "test.csv"
4 target_field = "title"
5 tokenizer = custom_tokenizer_P2
6 MAX_VOCAB_SIZE = 5000
7
8 TEXT, LABEL, train_data, valid_data, test_data = load_data(path, train_file, test_file,
9
10
11 train_iterator_IT, validation_IT, test_iterator_IT = generate_label_iterator(train_data,
12 train_iterator_CL, validation_CL, test_iterator_CL = generate_label_iterator(train_data,
13 train_iterator_CV, validation_CV, test_iterator_CV = generate_label_iterator(train_data,
14
15 pos_weights = get_pos_weight(train_data)
16
17 model_IT, optimizer_IT = generate_model_and_optimizer(pos_weight=pos_weights[0])
18 model_CL, optimizer_CL = generate_model_and_optimizer(pos_weight=pos_weights[1])
19 model_CV, optimizer_CV = generate_model_and_optimizer(pos_weight=pos_weights[2])
20
21 N_EPOCHS = 30
22 label_names = ["InformationTheory", "ComputationalLinguistics", "ComputerVision"]
23 models = [model_IT, model_CL, model_CV]
24 optimizers = [optimizer_IT, optimizer_CL, optimizer_CV]
25 iterators = [(train_iterator_IT, validation_IT, test_iterator_IT), (train_iterator_CL, v
```

```
26 model_file_names = ["RNN_model_IT_T_P2_ALL", "RNN_model_CL_T_P2_ALL.pt", "RNN_model_CV_T
```

```
1 print(TEXT.vocab.freqs.most_common(50))
2 len(train_data), len(valid_data), len(test_data), len(TEXT.vocab)
```

```
[('-', 75468), ('\n ', 55134), ('for', 44553), ('of', 28257), (':', 27878), ('and', 27
112500, 12500, 18066, 5002)
```



```
1 rnn_results, best_model_states = train_loop(models, optimizers, iterators, N_EPOCHS, lab
2 rnn_results_df = rnn_results_df.append(pd.DataFrame(rnn_results), ignore_index = True)
3
```



```
ComputerVision Test Loss: 0.299 | Test Acc: 87.24%
Epoch: 24 | Epoch Time: 0m 8s
ComputerVision Train Loss: 0.179 | Train Acc: 92.69%
ComputerVision Valid Loss: 0.272 | Valid Acc: 88.37%
ComputerVision Test Loss: 0.294 | Test Acc: 87.40%
Epoch: 25 | Epoch Time: 0m 8s
ComputerVision Train Loss: 0.177 | Train Acc: 92.85%
ComputerVision Valid Loss: 0.265 | Valid Acc: 88.32%
ComputerVision Test Loss: 0.281 | Test Acc: 88.20%
Epoch: 26 | Epoch Time: 0m 8s
ComputerVision Train Loss: 0.170 | Train Acc: 93.10%
ComputerVision Valid Loss: 0.267 | Valid Acc: 88.52%
ComputerVision Test Loss: 0.273 | Test Acc: 88.56%
Best performing epoch for ComputerVision model: 26
ComputerVision:
[[10595 1252]
 [ 818 5401]]
```

1

P1 Abstract 1000

```
1 path = './'
2 train_file = "train_1000.csv"
3 test_file = "test.csv"
4 target_field = "abstract"
5 tokenizer = custom_tokenizer
6 MAX_VOCAB_SIZE = 10000
7
8 TEXT, LABEL, train_data, valid_data, test_data = load_data(path, train_file, test_file,
9
10
11 train_iterator_IT, validation_IT, test_iterator_IT = generate_label_iterator(train_data,
12 train_iterator_CL, validation_CL, test_iterator_CL = generate_label_iterator(train_data,
13 train_iterator_CV, validation_CV, test_iterator_CV = generate_label_iterator(train_data,
14
15 pos_weights = get_pos_weight(train_data)
16
17 model_IT, optimizer_IT = generate_model_and_optimizer(pos_weight=pos_weights[0])
18 model_CL, optimizer_CL = generate_model_and_optimizer(pos_weight=pos_weights[1])
19 model_CV, optimizer_CV = generate_model_and_optimizer(pos_weight=pos_weights[2])
20
21 N_EPOCHS = 30
22 label_names = ["InformationTheory", "ComputationalLinguistics", "ComputerVision"]
23 models = [model_IT, model_CL, model_CV]
24 optimizers = [optimizer_IT, optimizer_CL, optimizer_CV]
25 iterators = [(train_iterator_IT, validation_IT, test_iterator_IT), (train_iterator_CL, v
26 model_file_names = ["RNN_model_IT_A_P1_1000", "RNN_model_CL_A_P1_1000.pt", "RNN_model_CV
27 best_valid_losses = [float("inf"), float("inf"), float("inf")]
```

```
1 print(TEXT.vocab.freqs.most_common(50))
2 len(train_data), len(valid_data), len(test_data), len(TEXT.vocab)
[( '\n', 12514), ('.', 6196), ('-', 4941), (')', 1610), ('(', 1537), (' ', 900), ('mode', 900, 100, 18066, 10002)
```



```
1 rnn_results, best_model_states = train_loop(models, optimizers, iterators, N_EPOCHS, lab
2 rnn_results_df = rnn_results_df.append(pd.DataFrame(rnn_results), ignore_index = True)
3
```

```
Accuracy: 0.4754234473596812
Macro Precision: 0.5606306227183617
Macro Recall: 0.5048577669472286
Macro F1 score: 0.3429903488800032
MCC: 0.03432371978857447
```



```
Training model for P1_ComputationalLinguistics, using 900 data of abstract...
```

```
Epoch: 01 | Epoch Time: 0m 2s
```

```
ComputationalLinguistics Train Loss: 1.078 | Train Acc: 58.54%
```

```
ComputationalLinguistics Valid Loss: 1.229 | Valid Acc: 34.64%
```

```
ComputationalLinguistics Test Loss: 1.007 | Test Acc: 59.05%
```

```
Early stopping at epoch 12 for ComputationalLinguistics model.
```

```
Best performing epoch for ComputationalLinguistics model: 1
```

```
ComputationalLinguistics:
```

```
[[9462 5236]
 [2143 1225]]
```

```
Accuracy: 0.5915531938447913
```

```
ComputerVision Test Loss: 0.649 | Test Acc: 65.49%
Epoch: 25 | Epoch Time: 0m 2s
ComputerVision Train Loss: 0.678 | Train Acc: 46.77%
ComputerVision Valid Loss: 0.679 | Valid Acc: 51.30%
ComputerVision Test Loss: 0.641 | Test Acc: 65.53%
Best performing epoch for ComputerVision model: 25
ComputerVision:
[[11820    27]
 [ 6204    15]]
Accuracy: 0.655097974094985
Macro Precision: 0.5064675670534525
Macro Recall: 0.5000664526743964
Macro F1 score: 0.39809729994221926
MCC: 0.0013111630372156966
```

1

P2 Abstract 1000

```
1 path = './'
2 train_file = "train_1000.csv"
3 test_file = "test.csv"
4 target_field = "abstract"
5 tokenizer = custom_tokenizer_P2
6 MAX_VOCAB_SIZE = 10000
7
8 TEXT, LABEL, train_data, valid_data, test_data = load_data(path, train_file, test_file,
9
10
11 train_iterator_IT, validation_IT, test_iterator_IT = generate_label_iterator(train_data,
12 train_iterator_CL, validation_CL, test_iterator_CL = generate_label_iterator(train_data,
13 train_iterator_CV, validation_CV, test_iterator_CV = generate_label_iterator(train_data,
14
15 pos_weights = get_pos_weight(train_data)
16
17 model_IT, optimizer_IT = generate_model_and_optimizer(pos_weight=pos_weights[0])
18 model_CL, optimizer_CL = generate_model_and_optimizer(pos_weight=pos_weights[1])
19 model_CV, optimizer_CV = generate_model_and_optimizer(pos_weight=pos_weights[2])
20
21 N_EPOCHS = 30
22 label_names = ["InformationTheory", "ComputationalLinguistics", "ComputerVision"]
23 models = [model_IT, model_CL, model_CV]
24 optimizers = [optimizer_IT, optimizer_CL, optimizer_CV]
25 iterators = [(train_iterator_IT, validation_IT, test_iterator_IT), (train_iterator_CL, v
26 model_file_names = ["RNN_model_IT_A_P2_1000", "RNN_model_CL_A_P2_1000.pt", "RNN_model_CV
27 best_valid_losses = [float("inf"), float("inf"), float("inf")]]
```

```
1 print(TEXT.vocab.freqs.most_common(50))
2 len(train_data), len(valid_data), len(test_data), len(TEXT.vocab)

[('\'\n', 12514), ('the', 8593), ('.', 6196), ('of', 4993), ('-', 4941), ('and', 4018), (
(900, 100, 18066, 7054)
```



```
1 rnn_results, best_model_states = train_loop(models, optimizers, iterators, N_EPOCHS, lab
2 rnn_results_df = rnn_results_df.append(pd.DataFrame(rnn_results), ignore_index = True)
3

Initiating model for PL_CoMPuLATIONaLlInGuIStICs, using GPU 0 of available...
Epoch: 01 | Epoch Time: 0m 3s
ComputationalLinguistics Train Loss: 1.080 | Train Acc: 34.38%
ComputationalLinguistics Valid Loss: 1.228 | Valid Acc: 32.47%
ComputationalLinguistics Test Loss: 0.999 | Test Acc: 58.32%
Epoch: 02 | Epoch Time: 0m 3s
ComputationalLinguistics Train Loss: 1.109 | Train Acc: 33.54%
```



```

ComputerVision Valid Loss: 0.677 | Valid Acc: 51.91%
ComputerVision Test Loss: 0.667 | Test Acc: 63.32%
Epoch: 25 | Epoch Time: 0m 3s
ComputerVision Train Loss: 0.678 | Train Acc: 46.98%
ComputerVision Valid Loss: 0.684 | Valid Acc: 51.91%
ComputerVision Test Loss: 0.654 | Test Acc: 64.88%
Best performing epoch for ComputerVision model: 25
ComputerVision:
[[11645 202]
 [ 6154   65]]
Accuracy: 0.6481788995903908
Macro Precision: 0.44884796583067726
Macro Recall: 0.49670055549468134
Macro F1 score: 0.4028233119980119
----
```

P1 Abstract ALL

```

1 path = './'
2 train_file = "train.csv"
3 test_file = "test.csv"
4 target_field = "abstract"
5 tokenizer = custom_tokenizer
6 MAX_VOCAB_SIZE = 10000
7
8 TEXT, LABEL, train_data, valid_data, test_data = load_data(path, train_file, test_file,
9
10
11 train_iterator_IT, validation_IT, test_iterator_IT = generate_label_iterator(train_data,
12 train_iterator_CL, validation_CL, test_iterator_CL = generate_label_iterator(train_data,
13 train_iterator_CV, validation_CV, test_iterator_CV = generate_label_iterator(train_data,
14
15 pos_weights = get_pos_weight(train_data)
16
17 model_IT, optimizer_IT = generate_model_and_optimizer(pos_weight=pos_weights[0])
18 model_CL, optimizer_CL = generate_model_and_optimizer(pos_weight=pos_weights[1])
19 model_CV, optimizer_CV = generate_model_and_optimizer(pos_weight=pos_weights[2])
20
21 N_EPOCHS = 30
22 label_names = ["InformationTheory", "ComputationalLinguistics", "ComputerVision"]
23 models = [model_IT, model_CL, model_CV]
24 optimizers = [optimizer_IT, optimizer_CL, optimizer_CV]
25 iterators = [(train_iterator_IT, validation_IT, test_iterator_IT), (train_iterator_CL, v
26 model_file_names = ["RNN_model_IT_A_P2_1000", "RNN_model_CL_A_P2_1000.pt", "RNN_model_CV
27 best_valid_losses = [float("inf"), float("inf"), float("inf")]

1 # print(TEXT.vocab.freqs.most_common(50))
2 len(train_data), len(valid_data), len(test_data), len(TEXT.vocab)

(112500, 12500, 18066, 10002)
```

```
1 rnn_results, best_model_states = train_loop(models, optimizers, iterators, N_EPOCHS, lab
2 rnn_results_df = rnn_results_df.append(pd.DataFrame(rnn_results), ignore_index = True)
3
```

```
Best performing epoch for ComputerVision model: 24
ComputerVision:
[[ 719 11128]
 [ 446 5773]]
Accuracy: 0.35934905347060775
Macro Precision: 0.4793724017143004
Macro Recall: 0.4944873801200065
```

P2 Abstract ALL

```
1 path = './'
2 train_file = "train.csv"
3 test_file = "test.csv"
4 target_field = "abstract"
5 tokenizer = custom_tokenizer_P2
6 MAX_VOCAB_SIZE = 10000
7
8 TEXT, LABEL, train_data, valid_data, test_data = load_data(path, train_file, test_file,
9
10 train_iterator_IT, validation_IT, test_iterator_IT = generate_label_iterator(train_data,
11 train_iterator_CL, validation_CL, test_iterator_CL = generate_label_iterator(train_data,
12 train_iterator_CV, validation_CV, test_iterator_CV = generate_label_iterator(train_data,
13
14 pos_weights = get_pos_weight(train_data)
15
16 model_IT, optimizer_IT = generate_model_and_optimizer(pos_weight=pos_weights[0])
17 model_CL, optimizer_CL = generate_model_and_optimizer(pos_weight=pos_weights[1])
18 model_CV, optimizer_CV = generate_model_and_optimizer(pos_weight=pos_weights[2])
19
20 N_EPOCHS = 30
21 label_names = ["InformationTheory", "ComputationalLinguistics", "ComputerVision"]
22 models = [model_IT, model_CL, model_CV]
23 optimizers = [optimizer_IT, optimizer_CL, optimizer_CV]
24 iterators = [(train_iterator_IT, validation_IT, test_iterator_IT), (train_iterator_CL, v
25 model_file_names = ["RNN_model_IT_A_P2_ALL", "RNN_model_CL_A_P2_ALL.pt", "RNN_model_CV_A
26 best_valid_losses = [float("inf"), float("inf"), float("inf")]

1 # print(TEXT.vocab.freqs.most_common(50))
2 len(train_data), len(valid_data), len(test_data), len(TEXT.vocab)

(112500, 12500, 18066, 10002)

1 rnn_results, best_model_states = train_loop(models, optimizers, iterators, N_EPOCHS, lab
2 rnn_results_df = rnn_results_df.append(pd.DataFrame(rnn_results), ignore_index = True)
```

```
ComputationalLinguistics Test Loss: 0.988 | Test Acc: 80.77%
Epoch: 03 | Epoch Time: 0m 58s
ComputationalLinguistics Train Loss: 1.053 | Train Acc: 46.70%
ComputationalLinguistics Valid Loss: 1.112 | Valid Acc: 75.64%
ComputationalLinguistics Test Loss: 0.968 | Test Acc: 80.78%
Early stopping at epoch 9 for ComputationalLinguistics model.
Best performing epoch for ComputationalLinguistics model: 3
ComputationalLinguistics:
[[14573 125]
 [ 3346 22]]
Accuracy: 0.8078711391564264
Macro Precision: 0.4814653468954968
Macro Recall: 0.49901375403249393
Macro F1 score: 0.4530504408879102
MCC: -0.008550959451065022
```

Training model for P2_ComputerVision, using 112500 data of abstract...

```
Epoch: 01 | Epoch Time: 0m 58s
ComputerVision Train Loss: 0.666 | Train Acc: 50.03%
ComputerVision Valid Loss: 0.667 | Valid Acc: 50.67%
ComputerVision Test Loss: 0.688 | Test Acc: 43.47%
Epoch: 02 | Epoch Time: 0m 58s
ComputerVision Train Loss: 0.665 | Train Acc: 50.25%
ComputerVision Valid Loss: 0.665 | Valid Acc: 51.70%
ComputerVision Test Loss: 0.684 | Test Acc: 43.22%
Epoch: 04 | Epoch Time: 0m 57s
ComputerVision Train Loss: 0.665 | Train Acc: 50.01%
ComputerVision Valid Loss: 0.664 | Valid Acc: 52.84%
ComputerVision Test Loss: 0.674 | Test Acc: 52.54%
Epoch: 06 | Epoch Time: 0m 59s
ComputerVision Train Loss: 0.665 | Train Acc: 50.11%
ComputerVision Valid Loss: 0.662 | Valid Acc: 54.05%
ComputerVision Test Loss: 0.670 | Test Acc: 52.59%
Epoch: 08 | Epoch Time: 0m 59s
ComputerVision Train Loss: 0.665 | Train Acc: 50.13%
ComputerVision Valid Loss: 0.662 | Valid Acc: 54.44%
ComputerVision Test Loss: 0.668 | Test Acc: 53.50%
Epoch: 14 | Epoch Time: 0m 59s
ComputerVision Train Loss: 0.665 | Train Acc: 50.18%
ComputerVision Valid Loss: 0.660 | Valid Acc: 54.11%
ComputerVision Test Loss: 0.659 | Test Acc: 56.98%
Epoch: 16 | Epoch Time: 0m 58s
ComputerVision Train Loss: 0.665 | Train Acc: 50.07%
ComputerVision Valid Loss: 0.659 | Valid Acc: 53.48%
ComputerVision Test Loss: 0.657 | Test Acc: 59.57%
Early stopping at epoch 22 for ComputerVision model.
Best performing epoch for ComputerVision model: 16
ComputerVision:
[[8012 3835]
 [3480 2739]]
Accuracy: 0.5950957599911436
Macro Precision: 0.5569109808376108
Macro Recall: 0.5583569307513049
Macro F1 score: 0.5572905202512220
```

```
1 rnn_results_df.columns  
  
Index(['accuracy', 'label_name', 'model_name', 'target_field', 'train_size',  
       'tokenizer', 'Macro Recall', 'Macro Precision', 'Macro F1 score', 'MCC',  
       'Recall', 'Precision'],  
      dtype='object')
```

▼ Storing output

```
1 metric_results_df.to_csv("svc_results.csv", index = False)
```

```
1 rnn_results_df.to_csv("rnn_results.csv", index = False)
```

▼ Part 1C: Results for Methods

F1, precision, etc.

```
1 from google.colab import drive  
2 drive.mount('/content/drive')  
3
```

Mounted at /content/drive

```
1 import pandas as pd  
2 import os  
3  
4 os.chdir("/content/drive/MyDrive/")  
5  
6 rnn_results = pd.read_csv("rnn_results.csv")  
7 svc_results = pd.read_csv("svc_results.csv")  
8
```

```
1 combined_df = pd.concat([rnn_results, svc_results], ignore_index = True)
```

```
1 combined_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 48 entries, 0 to 47  
Data columns (total 13 columns):  
 #   Column           Non-Null Count  Dtype     
 ---  --     
 0   accuracy        24 non-null     float64  
 1   label_name      48 non-null     object  
 2   model_name      48 non-null     object  
 3   target_field    48 non-null     object
```

```
4  train_size      48 non-null    int64
5  tokenizer       48 non-null    object
6  Macro Recall   48 non-null    float64
7  Macro Precision 48 non-null    float64
8  Macro F1 score 48 non-null    float64
9  MCC             48 non-null    float64
10 Recall          48 non-null    object
11 Precision       48 non-null    object
12 Accuracy        24 non-null    float64
dtypes: float64(6), int64(1), object(6)
memory usage: 5.0+ KB
```

▼ RNN models result

```
1 rnn_results = rnn_results[['label_name', 'model_name', 'target_field', 'train_size',
2                               'tokenizer', 'accuracy', 'Macro Recall', 'Macro Precision',
3                               'Recall', 'Precision']]
4 rnn_results = rnn_results.rename(columns = {"accuracy": "Accuracy"})
5 rnn_results
```

	label_name	model_name	target_field	train_size	tokenizer	Accuracy	
0	InformationTheory	RNN	title	900	P1	0.519318	0.519318
1	ComputationalLinguistics	RNN	title	900	P1	0.543230	0.543230

▼ SVC models result

```
1 svc_results
```

	label_name	model_name	target_field	train_size	tokenizer	Acc
0	InformationTheory	LinearSVC	title	900	custom_tokenizer	0.81
1	ComputationalLinauistics	LinearSVC	title	900	custom tokenizer	0.81
1						

▼ Part 1D: Plots for Methods

F1, precision, etc.

```

1 %matplotlib inline
2
3 import matplotlib.pyplot as plt
4 from sklearn.metrics import auc
5
6
7 def plot_precision_recall_curve(df):
8     plt.figure(figsize=(10, 6))
9
10    for index, row in df.iterrows():
11        Model = "RNN" if row["model_name"] == "RNN" else "SVC"
12        tag = ''
13        target_field = 'T' if row['target_field'] == 'title' else 'A'
14        size = str(1000) if row['train_size'] == 900 else 'All'
15        tokenizer = 'P1' if row['tokenizer'] == 'P1' else 'P2'
16        if row["label_name"]== "InformationTheory":
17            tag = "IT"
18        elif row["label_name"] == "ComputationalLinguistics":
19            tag = "CL"
20        else:
21            tag = "CV"
22        label = f"{Model}{tokenizer}{size}{target_field}_{tag}"
23 #
24         plt.plot(row['Recall'], row['Precision'], label=f"{label} (AUC: {auc(row['Recall'], row['Precision'])})")
25
26    plt.xlabel("Recall")
27    plt.ylabel("Precision")
28    plt.title("Precision-Recall Curve")
29    plt.legend(loc='best')
30    plt.show()
31
32

```

```

1 # I accidentally run this part of code again, hopefully it wont affect the following out
2
3 rnn_results_df = rnn_results_df[['label_name', 'model_name', 'target_field', 'train_size',
4                                     'tokenizer', 'accuracy', 'Macro Recall', 'Macro Precision',
5                                     'Recall', 'Precision']]
6 rnn_results_df = rnn_results_df.rename(columns = {"accuracy": "Accuracy"})
7
8 combined_df = pd.concat([rnn_results_df, metrics_results_df], axis = 0, ignore_index = T

```

Traceback (most recent call last):

```
Cell In[106], line 1
  rnn_results_df = rnn_results_df[['label_name', 'model_name', 'target_field',
'train_size',

  File ~\anaconda3\envs\5212A1\lib\site-packages\pandas\core\frame.py:3813 in
    __getitem__
        indexer = self.columns.get_indexer.strict_map["columns"][[1]

1 combined_df = pd.concat([rnn_results_df, metric_results_df], axis = 0, ignore_index = Tr
2 combined_df.info()
3
4 combined_df.loc[combined_df["tokenizer"] == "custom_tokenizer", "tokenizer"] = "P1"
5 combined_df.loc[combined_df["tokenizer"] == "custom_tokenizer_P2", "tokenizer"] = "P2"
6
7 combined_df
```

			RNN.ipynb - Colaboratory				
16	ComputationalLinguistics	RNN	abstract	900	P2	0.795749	▲
17	ComputerVision	RNN	abstract	900	P2	0.648179	
18	InformationTheory	RNN	abstract	112500	P1	0.517547	
19	ComputationalLinguistics	RNN	abstract	112500	P1	0.807041	
20	ComputerVision	RNN	abstract	112500	P1	0.359349	
21	InformationTheory	RNN	abstract	112500	P2	0.521255	▼
22	ComputationalLinguistics	RNN	abstract	112500	P2	0.807871	
23	ComputerVision	RNN	abstract	112500	P2	0.595096	
24	InformationTheory	LinearSVC	title	900	P1	0.836488	
25	ComputationalLinguistics	LinearSVC	title	900	P1	0.881103	
26	ComputerVision	LinearSVC	title	900	P1	0.834883	
27	InformationTheory	LinearSVC	title	900	P2	0.846895	▼

```
1 df_abstract = combined_df[combined_df["target_field"] == "abstract"]
2 df_title = combined_df[combined_df["target_field"] == "title"]

1 df_abstract.info()

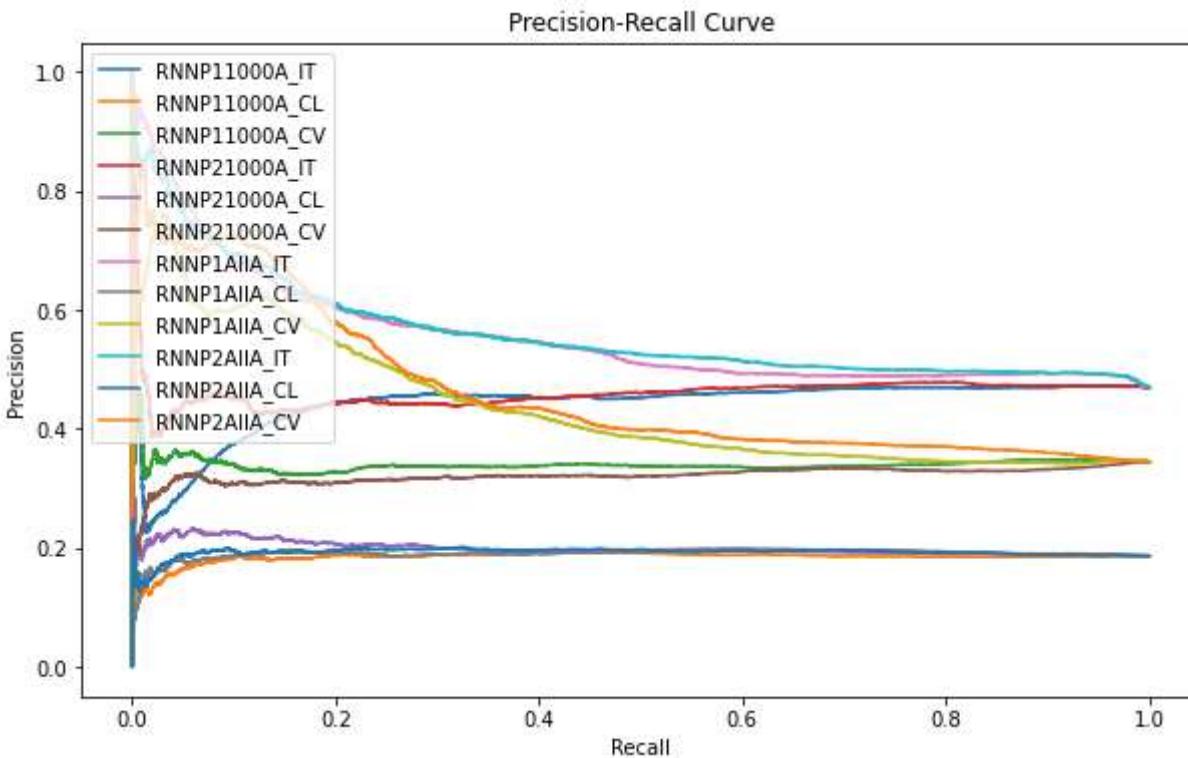
<class 'pandas.core.frame.DataFrame'>
Int64Index: 24 entries, 12 to 47
Data columns (total 12 columns):
 #   Column           Non-Null Count  Dtype  
---  --  
 0   label_name       24 non-null     object  
 1   model_name       24 non-null     object  
 2   target_field     24 non-null     object  
 3   train_size       24 non-null     int64  
 4   tokenizer        24 non-null     object  
 5   Accuracy         24 non-null     float64 
 6   Macro Recall    24 non-null     float64 
 7   Macro Precision  24 non-null     float64 
 8   Macro F1 score  24 non-null     float64 
 9   MCC              24 non-null     float64 
 10  Recall            24 non-null     object  
 11  Precision         24 non-null     object  
dtypes: float64(5), int64(1), object(6)
memory usage: 2.4+ KB

1 df_title.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 24 entries, 0 to 41
Data columns (total 12 columns):
 #   Column           Non-Null Count  Dtype  
---  --  
 0   label_name       24 non-null     object  
 1   model_name       24 non-null     object  
 2   target_field     24 non-null     object  
 3   train_size       24 non-null     int64  
 4   tokenizer        24 non-null     object  
 5   Accuracy         24 non-null     float64 
 6   Macro Recall    24 non-null     float64 
 7   Macro Precision  24 non-null     float64 
 8   Macro F1 score  24 non-null     float64 
 9   MCC              24 non-null     float64 
 10  Recall            24 non-null     object  
 11  Precision         24 non-null     object  
dtypes: float64(5), int64(1), object(6)
memory usage: 2.4+ KB
```

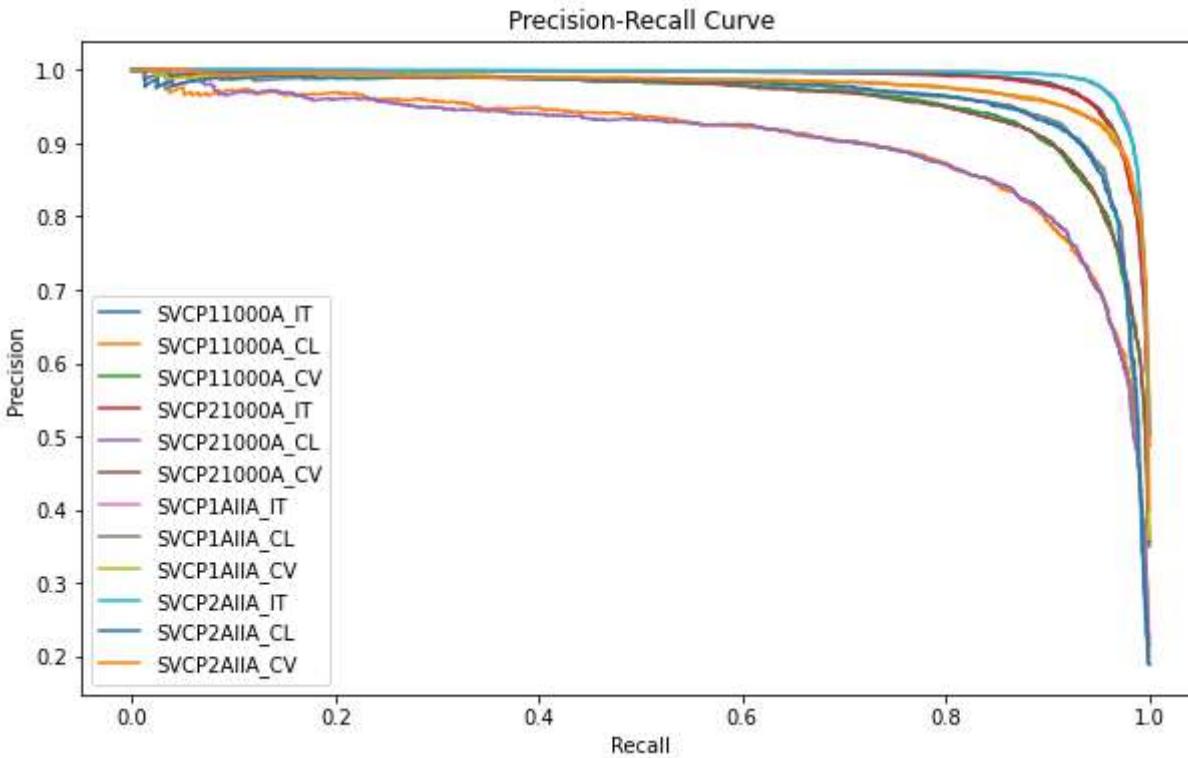
▼ ALL RNN models training on the "abstract" column

```
1 plot_precision_recall_curve(df_abstract[df_abstract["model_name"] == "RNN"])
```



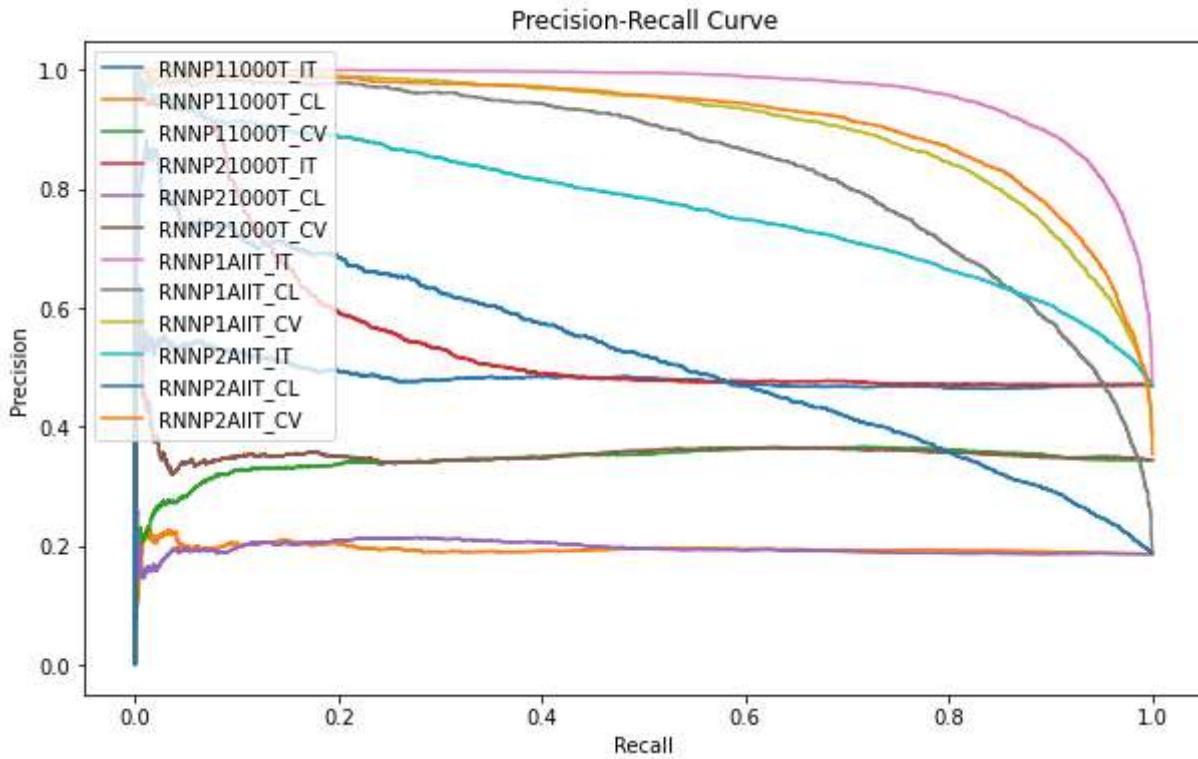
▼ ALL SVC models training on the "abstract" column

```
1 plot_precision_recall_curve(df_abstract[df_abstract["model_name"] != "RNN"])
```



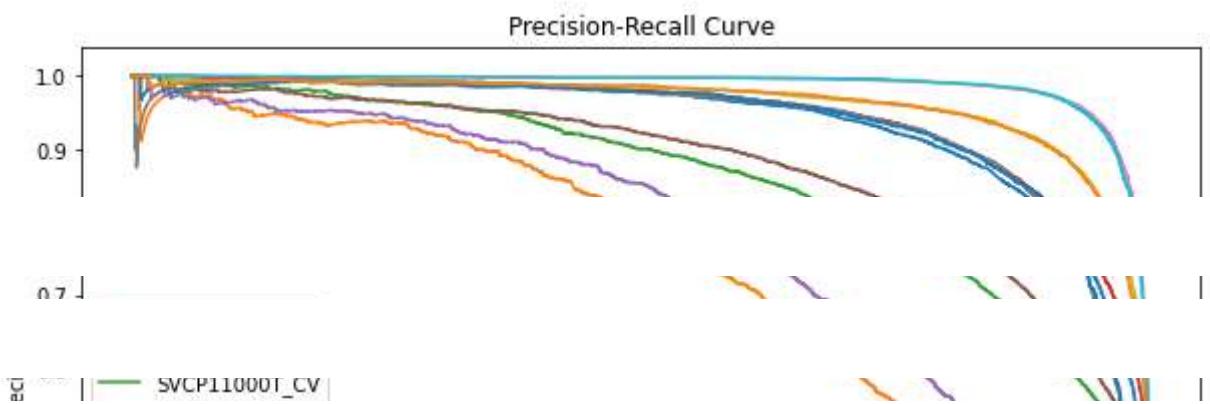
▼ ALL RNN models training on the "title" column

```
1 plot_precision_recall_curve(df_title[df_title["model_name"] == "RNN"])
```



▼ ALL SVC models training on the "abstract" column

```
1 plot_precision_recall_curve(df_title[df_title["model_name"] != "RNN"])
```



▼ Part 2: Topic Modelling

```
||— SVCP1AIIT CL ||  
1 # !pip3 install pyldavis==2.1.2  
2 # !pip3 install scikit-learn==0.24.2  
3 # !pip install gensim
```

```
    self.initialize()
File "C:\Users\Sunny\AppData\Local\Temp\pip-build-env-ocn1k1ah\overlay\Lib\site-
  vc_env = _get_vc_env(plat_spec)
File "C:\Users\Sunny\AppData\Local\Temp\pip-build-env-ocn1k1ah\overlay\Lib\site-
  return _msvc14_get_vc_env(plat_spec)
File "C:\Users\Sunny\AppData\Local\Temp\pip-build-env-ocn1k1ah\overlay\Lib\site-
  raise distutils.errors.DistutilsPlatformError(
distutils.errors.DistutilsPlatformError: Microsoft Visual C++ 14.0 or greater is r
[end of output]
```

note: This error originates from a subprocess, and is likely not a problem with pip.
error: metadata-generation-failed

Encountered error while generating package metadata.

See above for output.

note: This is an issue with the package mentioned above, not pip.

hint: See above for details.

Requirement already satisfied: gensim in c:\users\sunny\appdata\local\programs\python
Requirement already satisfied: gensim in c:\users\sunny\appdata\local\programs\python

```
1 import warnings
2 warnings.filterwarnings('ignore')
```

```
1 import re, nltk, spacy, string
2
3 from sklearn.decomposition import LatentDirichletAllocation
4 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
5 from pprint import pprint
6
7 import pyLDAvis
8 import pyLDAvis.sklearn
9 import matplotlib.pyplot as plt
10 %matplotlib inline
11
12 from plotly.offline import plot
13 import plotly.graph_objects as go
14 import plotly.express as px
15
16 import pandas as pd
17 import seaborn as sns
18 from nltk.tokenize import RegexpTokenizer
19 import nltk
20 from nltk.stem.wordnet import WordNetLemmatizer
21 from gensim.models import Phrases
22 from gensim.corpora import Dictionary
23 from gensim.models.coherencemodel import CoherenceModel
24 from gensim.models import LdaModel
25 from spacy.lang.en.stop_words import STOP_WORDS
26
```

```
1 import random
2 SEED = 1
3 random.seed(SEED)

1 def prepare_data(df_train, tokenizer_func, tokenizer, *args, **kwargs):
2     docs = df_train['abstract'].tolist()
3     raw_docs = docs.copy()
4
5     for idx in range(len(docs)):
6         docs[idx] = docs[idx].lower()
7         docs[idx] = tokenizer_func(docs[idx], tokenizer, *args, **kwargs)
8
9     bigram = Phrases(docs, min_count=20)
10    for idx in range(len(docs)):
11        for token in bigram[docs[idx]]:
12            if '_' in token:
13                docs[idx].append(token)
14
15    dictionary = Dictionary(docs)
16
17    dictionary.filter_extremes(no_below=20, no_above=0.5)
18
19    corpus = [dictionary.doc2bow(doc) for doc in docs]
20
21    return corpus, dictionary
22

1 def train_lda(corpus, dictionary, num_topics_range):
2     best_avg_topic_coherence = -100
3     best_num_topics = 0
4     avg_topic_coherence_values = []
5     chunksize = 2000
6     passes = 20
7     iterations = 400
8     eval_every = None # Don't evaluate model perplexity, takes too much time.
9
10    # Initialize variables
11    best_avg_topic_coherence = -100
12    best_num_topics = 0
13    avg_topic_coherence_values = []
14    for num_topics in num_topics_range:
15        print(f"Now processing LDA model with {num_topics} number of topics")
16        temp = dictionary[0]
17        id2word = dictionary.id2token
18        model = LdaModel(
19            corpus=corpus,
20            id2word=id2word,
21            chunksize=chunksize,
22            alpha='auto',
```

```
23         eta='auto',
24         iterations=iterations,
25         num_topics=num_topics,
26         passes=passes,
27         eval_every=eval_every,
28     )
29
30     top_topics = model.top_topics(corpus)
31     avg_topic_coherence = sum([t[1] for t in top_topics]) / num_topics
32     avg_topic_coherence_values.append(avg_topic_coherence)
33
34     if avg_topic_coherence > best_avg_topic_coherence:
35         best_avg_topic_coherence = avg_topic_coherence
36         best_num_topics = num_topics
37         best_model = model
38     print(f"Best number of topics: {best_num_topics}, average topic coherence: {best_avg
39
40 # Plot the average topic coherence values for different numbers of topics
41 plt.plot(num_topics_range, avg_topic_coherence_values)
42 plt.xlabel("Number of Topics")
43 plt.ylabel("Average Topic Coherence")
44 plt.show()
45
46 return best_num_topics, best_avg_topic_coherence, avg_topic_coherence_values, best_m
47
```

```
1 def plot_topic_coherence(num_topics_range, avg_topic_coherence_values):
2     plt.plot(num_topics_range, avg_topic_coherence_values)
3     plt.xlabel("Number of Topics")
4     plt.ylabel("Average Topic Coherence")
5     plt.show()

1 def get_document_topics(ldamodel, corpus, texts):
2     # Init output
3     document_topics_df = pd.DataFrame()
4
5     # Get main topic in each document
6     for i, row in enumerate(ldamodel[corpus]):
7         row = sorted(row, key=lambda x: (x[1]), reverse=True)
8         # Get the Dominant topic, Perc Contribution and Keywords for each document
9         for j, (topic_num, prop_topic) in enumerate(row):
10             if j == 0: # => dominant topic
11                 wp = ldamodel.show_topic(topic_num)
12                 topic_keywords = ", ".join([word for word, prop in wp])
13                 document_topics_df = document_topics_df.append(pd.Series([int(topic_num),
14
15             else:
16                 break
17
18             # Add original text to the end of the output
```

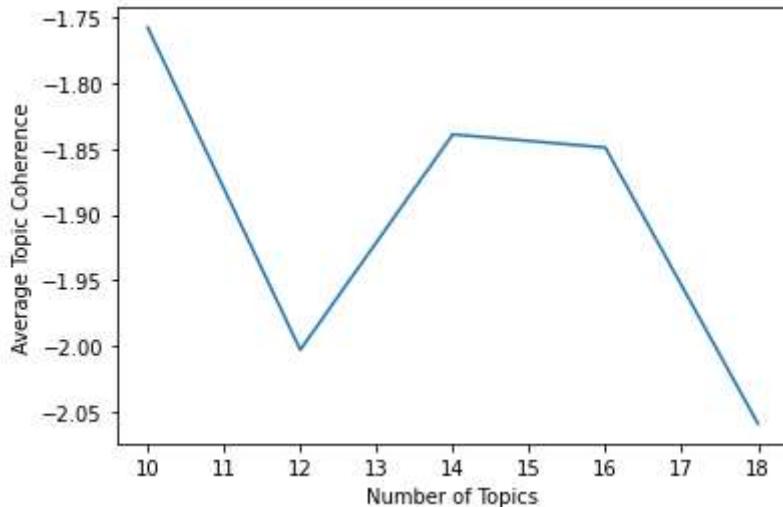
```
19     contents = pd.Series(texts)
20     document_topics_df = pd.concat([document_topics_df, contents], axis=1)
21
22     document_topics_df.columns = ['Dominant_Topic', 'Perc_Contribution', 'Topic_Keywords'
23
24     return document_topics_df
25
26
27 def find_top_k_doc(doc_topic_df, k=5):
28
29     doc_topics_sorted_df = pd.DataFrame()
30
31     doc_topic_df_grpd = doc_topic_df.groupby('Dominant_Topic')
32
33     for i, grp in doc_topic_df_grpd:
34         doc_topics_sorted_df = pd.concat([doc_topics_sorted_df,
35                                         grp.sort_values(['Perc_Contribution'], ascending=False)
36                                         .head(k)])
37
38     doc_topics_sorted_df.reset_index(drop=True, inplace=True)
39     doc_topics_sorted_df.columns = ['Topic_Num', "Topic_Perc_Contrib", "Keywords", "Text"]
40
41     return doc_topics_sorted_df
42
43
44 stop_words = set(STOP_WORDS)
45
46 def custom_tokenizer(text, tokenizer, stop_words):
47     return [tok for tok in tokenizer.tokenize(text) if tok.lower() not in stop_words]
48
49 def custom_tokenizer_P2(text, tokenizer, stemmer):
50     return [stemmer.stem(tok) for tok in tokenizer.tokenize(text)]
51
52
53 import pyLDAvis.gensim
54 from pprint import pprint
55
56 def summarize_lda_results(ldamodel, corpus, dictionary, raw_docs, num_topics, top_k=5):
57     # Display LDA visualization
58     # lda_display = pyLDAvis.gensim.prepare(ldamodel, corpus, dictionary, sort_topics=False)
59     # pyLDAvis.display(lda_display)
60
61     # Get document topics
62     doc_topic_df = get_document_topics(ldamodel, corpus, raw_docs)
63
64     # Find top k documents for each topic
65     top_k_df = find_top_k_doc(doc_topic_df, k=top_k)
66
67
68     return doc_topic_df, top_k_df
```

```
1 train_data = pd.read_csv("train.csv")
2 train_1000 = train_data.head(1000)
3 train_data = train_data.head(20000)
4 test_data = pd.read_csv("test.csv")
```

▼ P1 1000

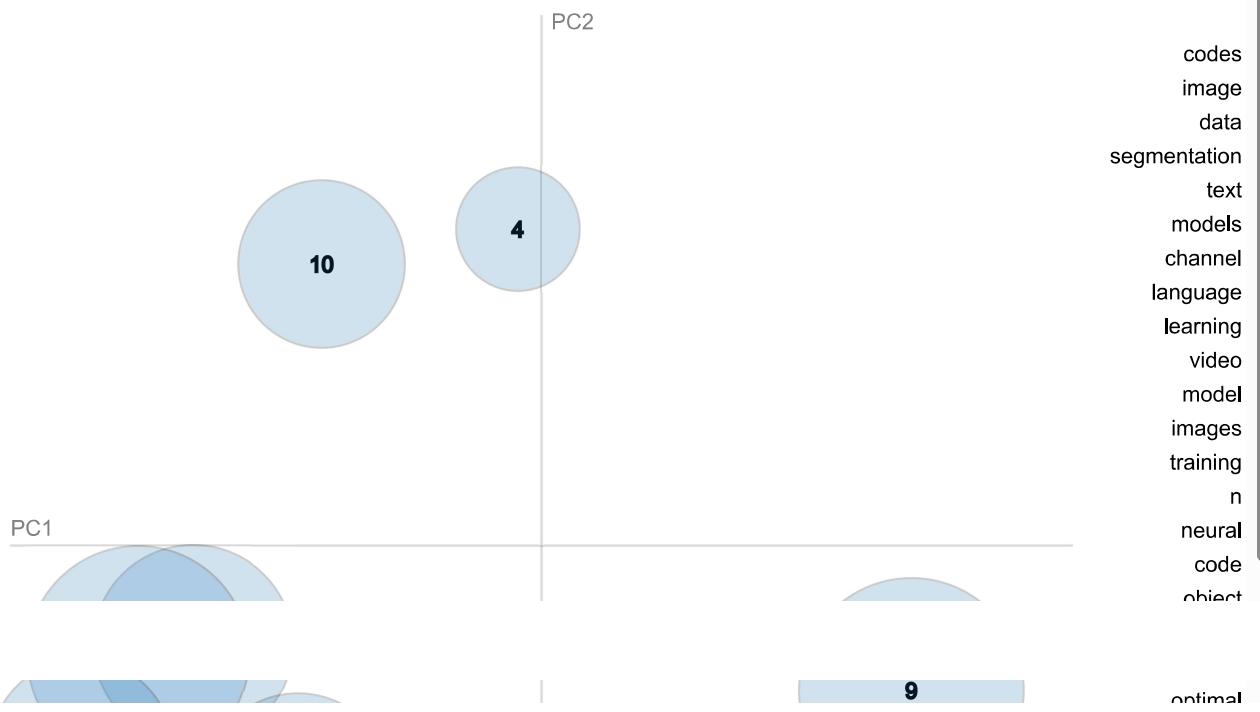
```
1 df_train = train_1000
2
3 raw_docs = df_train['abstract'].tolist().copy
4
5 tokenizer = RegexpTokenizer(r'\w+')
6 corpus, dictionary = prepare_data(df_train, custom_tokenizer, tokenizer, stop_words=stop
7
8
9 # Train LDA model
10 num_topics_range = range(10, 15, 2)
11 best_num_topics, best_avg_topic_coherence, avg_topic_coherence_values, best_model = trai
12
13 doc_topic_df_P1_1000, top_k_df_P1_1000 = summarize_lda_results(ldamodel=best_model, corp
14
15 lda_display = pyLDAvis.gensim.prepare(best_model, corpus, dictionary, sort_topics=False)
16
17 pyLDAvis.display(lda_display)
```

Now processing LDA model with 10 number of topics
 Now processing LDA model with 12 number of topics
 Now processing LDA model with 14 number of topics
 Now processing LDA model with 16 number of topics
 Now processing LDA model with 18 number of topics
 Best number of topics: 10, average topic coherence: -1.757952894995696



Selected Topic: [Previous Topic](#) [Next Topic](#) [Clear Topic](#)

Intertopic Distance Map (via multidimensional scaling)

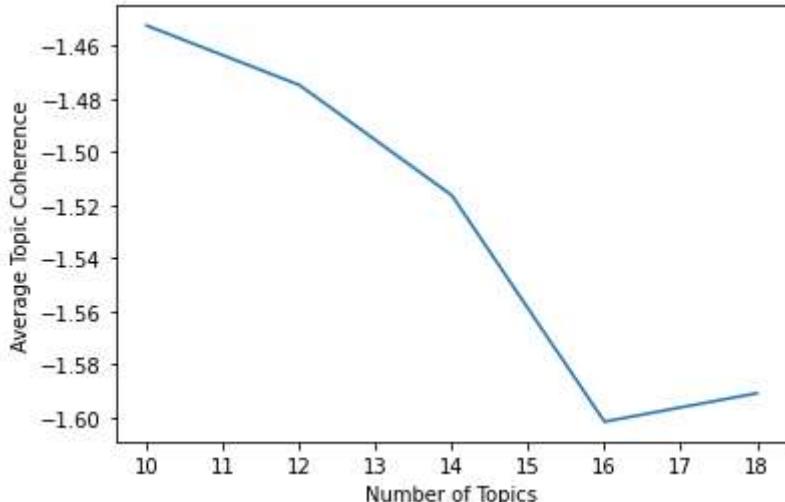


P2 1000

```
1 df_train = train_1000
2
```

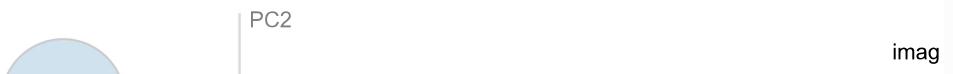
```
3 stemmer = PorterStemmer()
4
5 tokenizer = RegexpTokenizer(r'\w+')
6 corpus, dictionary = prepare_data(df_train, custom_tokenizer_P2, tokenizer, stemmer=stem
7
8 # Train LDA model
9 num_topics_range = range(10, 15, 2)
10 best_num_topics, best_avg_topic_coherence, avg_topic_coherence_values, best_model = train
11
12 doc_topic_df_P2_1000, top_k_df_P2_1000 = summarize_lda_results(ldamodel=best_model, corp
13
14 lda_display = pyLDAvis.gensim.prepare(best_model, corpus, dictionary, sort_topics=False)
15
16 pyLDAvis.display(lda_display)
```

Now processing LDA model with 10 number of topics
Now processing LDA model with 12 number of topics
Now processing LDA model with 14 number of topics
Now processing LDA model with 16 number of topics
Now processing LDA model with 18 number of topics
Best number of topics: 10, average topic coherence: -1.4523068232507323



Selected Topic: [Previous Topic](#) [Next Topic](#) [Clear Topic](#)

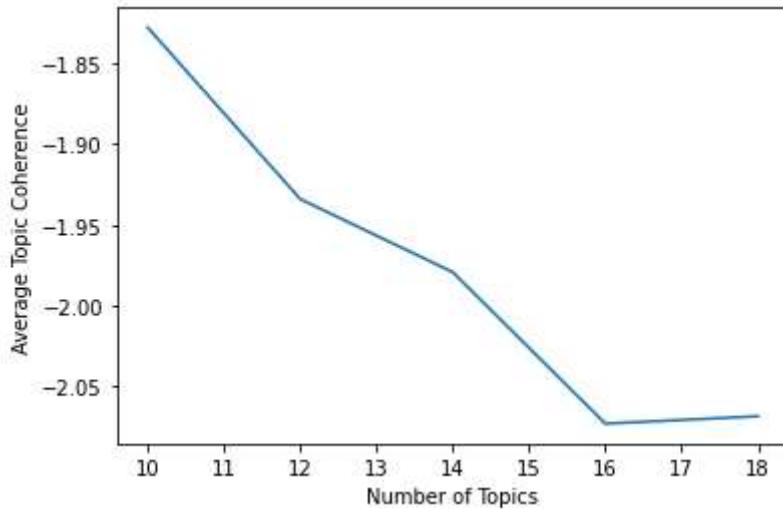
Intertopic Distance Map (via multidimensional scaling)



▼ P1 ALL

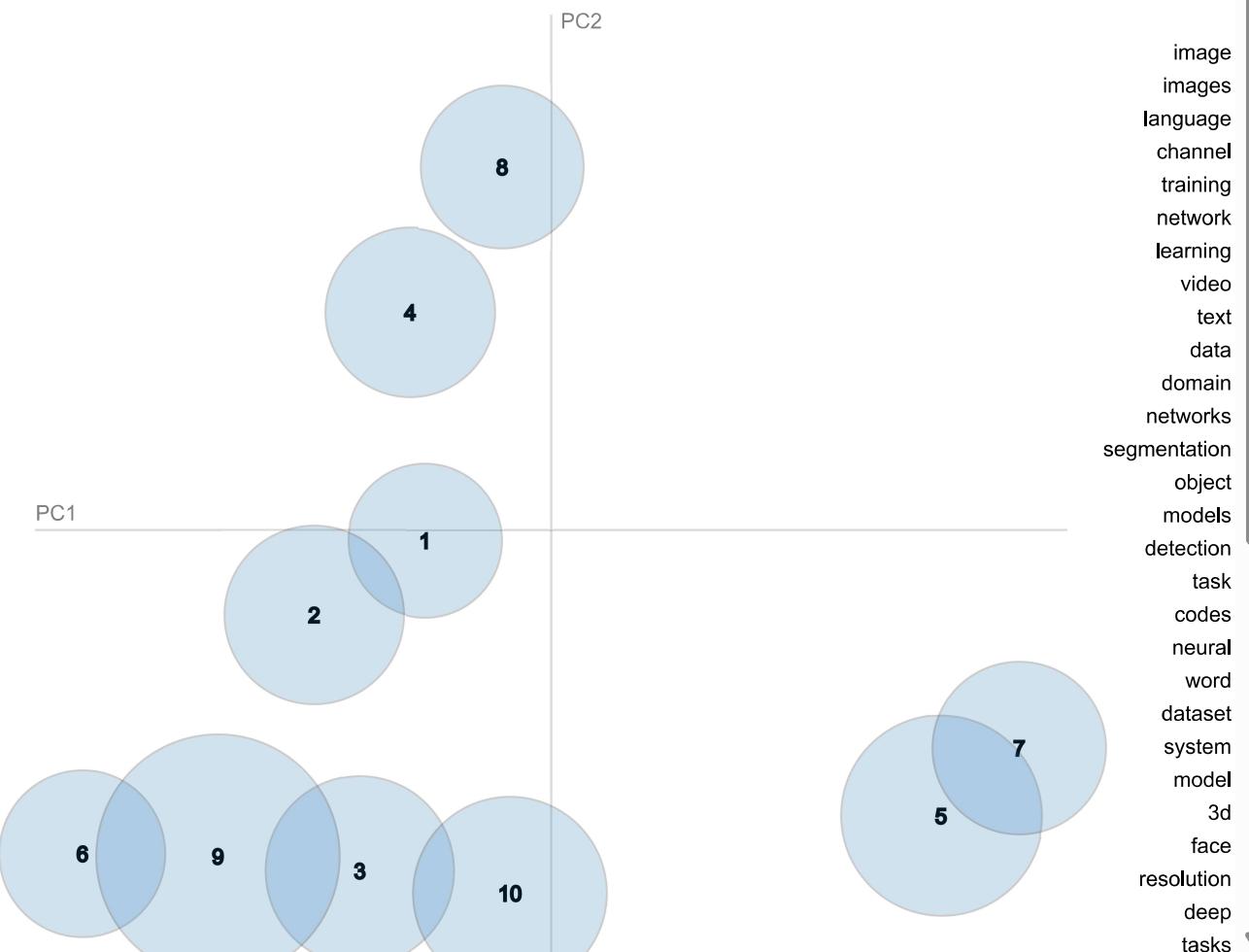
```
1 df_train = train_data
2
3 raw_docs = df_train['abstract'].tolist().copy
4
5 tokenizer = RegexpTokenizer(r'\w+')
6 corpus, dictionary = prepare_data(df_train, custom_tokenizer, tokenizer, stop_words=stop
7
8
9 # Train LDA model
10 num_topics_range = range(10, 15, 2)
11 best_num_topics, best_avg_topic_coherence, avg_topic_coherence_values, best_model = trai
12
13 doc_topic_df_P1_ALL, top_k_df_P1_ALL = summarize_lda_results(ldamodel=best_model, corpus
14
15 lda_display = pyLDAvis.gensim.prepare(best_model, corpus, dictionary, sort_topics=False)
16
17 pyLDAvis.display(lda_display)
```

Now processing LDA model with 10 number of topics
Now processing LDA model with 12 number of topics
Now processing LDA model with 14 number of topics
Now processing LDA model with 16 number of topics
Now processing LDA model with 18 number of topics
Best number of topics: 10, average topic coherence: -1.8275886401598131



Selected Topic: 0

Intertopic Distance Map (via multidimensional scaling)



▼ P2 ALL

```
1 df_train = train_data
2
3 stemmer = PorterStemmer()
4
5 tokenizer = RegexpTokenizer(r'\w+')
6 corpus, dictionary = prepare_data(df_train, custom_tokenizer_P2, tokenizer, stemmer=stem
7
8 # Train LDA model
9 num_topics_range = range(10, 15, 2)
10 best_num_topics, best_avg_topic_coherence, avg_topic_coherence_values, best_model = trai
11
12 doc_topic_df_P2_ALL, top_k_df_P2_ALL = summarize_lda_results(ldamodel=best_model, corpus
13
14 lda_display = pyLDAvis.gensim.prepare(best_model, corpus, dictionary, sort_topics=False)
15
16 pyLDAvis.display(lda_display)
```