

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

↗ Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

```
1 import pandas as pd
2
3 # Define file path
4 file_path = '/content/drive/My Drive/NLPHW4/hw4.pk'
5
6 # Load the pickle file into a DataFrame
7 data = pd.read_pickle(file_path)
8
9 # Display the first few rows to confirm the data
10 print(data.head())
11
```

```
↗
      body                                label
0  We use essential cookies to make Venngage wor...  legal_contract_examples
1  A legal contract is a written document that is...  legal_contract_examples
2  November 27 2023 14 min Author Olga Asheyichik...  legal_contract_examples
3  Accelerate contracts with AI native workflows ...  legal_contract_examples
4  Create smarter agreements commit to them more ...  legal_contract_examples
```

```
1 import pickle
2 import nltk
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.feature_extraction.text import TfidfVectorizer
6 from sklearn.ensemble import RandomForestClassifier
7 from sklearn.metrics import classification_report, accuracy_score
8 from nltk.tokenize import word_tokenize
9 from nltk.corpus import stopwords
10 import string
11
12 # Download NLTK resources
13 nltk.download('punkt_tab')
14 nltk.download('stopwords')
15
16 # Define file path for pickle file
17 file_path = '/content/drive/My Drive/NLPHW4/hw4.pk'
18
19 # Load the pickle file
20 with open(file_path, 'rb') as f:
21     data = pickle.load(f)
22
23 # Extract body (text of document) and labels
24 documents = data['body']
25 labels = data['label']
26
27 # Preprocess text: tokenization, lowercasing, removing stopwords and punctuation
28 stop_words = set(stopwords.words('english'))
29
30 def preprocess_text(text):
31     # Tokenize the text
32     tokens = word_tokenize(text.lower()) # Convert to lowercase and tokenize
33     # Remove stopwords and punctuation
34     tokens = [word for word in tokens if word.isalpha() and word not in stop_words]
35     return " ".join(tokens)
36
37 # Preprocess all documents
38 preprocessed_documents = [preprocess_text(doc) for doc in documents]
39
40 # Split the data into training and test sets (70% train, 30% test)
41 X_train, X_test, y_train, y_test = train_test_split(preprocessed_documents, labels, test_size=0.3, random_state=42)
42
43 # Convert text data into TF-IDF features
44 tfidf_vectorizer = TfidfVectorizer()
45 X_train_tfidf = tfidf_vectorizer.fit_transform(X_train)
46 X_test_tfidf = tfidf_vectorizer.transform(X_test)
47
48 # Train a Random Forest classifier
49 classifier = RandomForestClassifier(n_estimators=100, random_state=42)
50 classifier.fit(X_train_tfidf, y_train)
51
```

```

52 # Make predictions on the test set
53 y_pred = classifier.predict(X_test_tfidf)
54
55 # Evaluate the model
56 accuracy = accuracy_score(y_test, y_pred)
57 print(f"Accuracy: {accuracy}")
58 print("\nClassification Report:")
59 print(classification_report(y_test, y_pred))
60
61 # Classify all documents in the dataset (train and test combined)
62 all_tfidf = tfidf_vectorizer.transform(preprocessed_documents)
63 all_predictions = classifier.predict(all_tfidf)
64
65 # Print the predictions for all documents (first 10 for brevity)
66 for doc, label in zip(documents[:10], all_predictions[:10]): # Displaying the first 10 for brevity
67     print(f"Document: {doc[:100]}... => Predicted label: {label}")
68
69 # Optionally: Save the predictions with document texts
70 classified_documents = pd.DataFrame({
71     'Document': documents,
72     'Predicted Label': all_predictions
73 })
74
75 # Save to a CSV file
76 classified_documents.to_csv('/content/drive/My Drive/NLPHW4/classified_documents.csv', index=False)
77

```

```

[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data] Package punkt_tab is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
Accuracy: 0.7794117647058824

```

Classification Report:

	precision	recall	f1-score	support
engineering_specification_examples	0.62	0.95	0.75	22
legal_contract_examples	0.89	0.80	0.84	20
marketing_material_examples	1.00	0.62	0.76	26
accuracy			0.78	68
macro avg	0.84	0.79	0.78	68
weighted avg	0.84	0.78	0.78	68

Document: We use essential cookies to make Venngage work By clicking Accept All Cookies you agree to the stor... => Predicted label: le

Document: A legal contract is a written document that is drawn up by a party and is agreed upon by all parties... => Predicted label: le

Document: November 27 2023 14 min Author Olga Asheychik Senior Web Analytics Manager at PandaDoc Choosing the... => Predicted label: le

Document: Accelerate contracts with AI native workflows Advanced electronic signature on any device Create con... => Predicted label: le

Document: Create smarter agreements commit to them more efficiently and manage them to realize their full valu... => Predicted label: le

Document: A contract is an agreement between two parties that creates an obligation to perform or not perform ... => Predicted label: le

Document: Please enable JS and disable any ad blocker... => Predicted label: legal_contract_examples

Document: WEBINAR NOV 14th Use Contract Data to Drive New Insights REGISTER NOW In the realm of legal agreeme... => Predicted label: le

Document: luctus etiam leo nulla etiam convallis tincidunt integer Pellentesque suscipit adipiscing nullam lu... => Predicted label: le

Document: Entering into contracts is part of running a small business and it s important to manage your contra... => Predicted label: le