

✓ Lab 4 - Part 1: Core NLP Tasks

Course: Natural Language Processing

Objectives:

- Apply Part-of-Speech (POS) tagging to extract linguistic patterns
- Perform Named Entity Recognition (NER) to identify entities in text
- Calculate word and document similarities using different techniques
- Apply PCA for visualizing high-dimensional text representations
- Work with real-world datasets (Nike products and legal contracts)

Instructions

1. Complete all exercises marked with `# YOUR CODE HERE`
2. **Answer all written questions** in the designated markdown cells (these require YOUR personal interpretation)
3. Save your completed notebook
4. **Push to your Git repository and send the link to: yoroba93@gmail.com**

Important: Personal Interpretation Questions

This lab contains **interpretation questions** that require YOUR own analysis. These questions:

- Are based on YOUR specific results (which vary based on your choices)
- Require you to explain your reasoning

✓ Setup

```
1 # Install required libraries (uncomment if needed)
2 # !pip install spacy scikit-learn matplotlib seaborn pandas numpy
3 # !python -m spacy download en_core_web_sm
```

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from collections import Counter, defaultdict
6 import re
7 import warnings
8 warnings.filterwarnings('ignore')
9
```

```

10 # NLP libraries
11 import spacy
12 from sklearn.feature_extraction.text import TfidfVectorizer
13 from sklearn.metrics.pairwise import cosine_similarity
14 from sklearn.decomposition import PCA
15
16 # Load spaCy model
17 nlp = spacy.load('en_core_web_sm')
18
19 print("Setup complete!")

```

Setup complete!
spaCy version: 3.8.11

✓ Part A: Loading Nike Products Dataset

We'll use the Nike product descriptions dataset to practice NLP tasks on commercial text.

```

1 # Load Nike products dataset
2 # NOTE: Place the 'NikeProductDescriptions.csv' file in your work
3 nike_df = pd.read_csv('NikeProductDescriptions.csv')
4
5 print(f"Dataset shape: {nike_df.shape}")
6 print(f"\nColumns: {nike_df.columns.tolist()}")
7 print(f"\nFirst 3 products:")
8 nike_df.head(3)

```

Dataset shape: (400, 3)

Columns: ['Title', 'Subtitle', 'Product Description']

First 3 products:

	Title	Subtitle	Product Description	
0	Nike Air Force 1 '07	Men's Shoes	It doesn't get more legendary than this. Desig...	
1	Nike Air Max Dawn SE	Men's Shoes	Find out what moves you with the Air Max Dawn....	

Next steps:

[Generate code with nike_df](#)

[New interactive sheet](#)

```

1 # Display a sample product
2 sample_idx = 0
3 print("Sample Product:")
4 print("=" * 60)
5 print(f"Title: {nike_df.iloc[sample_idx]['Title']}")
6 print(f"Subtitle: {nike_df.iloc[sample_idx]['Subtitle']}")
7 print(f"\nDescription:\n{nike_df.iloc[sample_idx]['Product Descri

```

Sample Product:

Title: Nike Air Force 1 '07

Subtitle: Men's Shoes

Description:

It doesn't get more legendary than this. Designed to turn heads, the N

✓ Part B: Part-of-Speech (POS) Tagging

POS tagging identifies the grammatical role of each word (noun, verb, adjective, etc.).

```
1 # Example: POS tagging with spaCy
2 sample_text = "Nike Air Force 1 shoes provide incredible comfort and stylish design for athletes."
3 doc = nlp(sample_text)
4
5 print("POS Tagging Example:")
6 print("=" * 60)
7 for token in doc:
8     print(f"{token.text:15} | POS: {token.pos_:10} | Tag: {token.tag_:10} | Lemma: {token.lemma_:10}")
```

POS Tagging Example:

```
=====
Nike          | POS: PROPN      | Tag: NNP      | Lemma: Nike
Air           | POS: PROPN      | Tag: NNP      | Lemma: Air
Force         | POS: PROPN      | Tag: NNP      | Lemma: Force
1            | POS: NUM        | Tag: CD       | Lemma: 1
shoes         | POS: NOUN       | Tag: NNS      | Lemma: shoe
provide       | POS: VERB       | Tag: VBP      | Lemma: provide
incredible    | POS: ADJ        | Tag: JJ       | Lemma: incredible
comfort       | POS: NOUN       | Tag: NN       | Lemma: comfort
and           | POS: CCONJ      | Tag: CC       | Lemma: and
stylish       | POS: ADJ        | Tag: JJ       | Lemma: stylish
design        | POS: NOUN       | Tag: NN       | Lemma: design
for           | POS: ADP        | Tag: IN       | Lemma: for
athletes     | POS: NOUN       | Tag: NNS      | Lemma: athlete
.            | POS: PUNCT      | Tag: .        | Lemma: .
```

✓ Exercise B.1: Analyze POS Distribution in Nike Products

Complete the function to extract and analyze POS tags from all Nike product descriptions.

```
1 def analyze_pos_distribution(texts):
2     """
3     Analyze the distribution of POS tags in a list of texts.
4
5     Args:
6         texts (list): List of text strings
7
```

```

8     Returns:
9         Counter: Dictionary with POS tags and their counts
10    """
11    pos_counts = Counter()
12
13    # YOUR CODE HERE
14    # 1. For each text, process it with nlp(text)
15    # 2. For each token in the doc, count its POS tag (token.pos)
16    # 3. Return the counter
17    pos_counts = Counter()
18
19    for text in texts:
20        doc = nlp(text)
21        for token in doc:
22            pos_counts[token.pos_] += 1 # Count POS tag
23
24
25    return pos_counts
26
27 # Analyze Nike descriptions
28 nike_descriptions = nike_df['Product Description'].dropna().tolist()
29 pos_distribution = analyze_pos_distribution(nike_descriptions)
30
31 print("POS Tag Distribution:")
32 print("=" * 40)
33 for pos, count in pos_distribution.most_common(15):
34     print(f"{pos:10} : {count:5} ({count/sum(pos_distribution.values()):.2%}")

```

POS Tag Distribution:

```

=====
NOUN      : 4620 (20.72%)
VERB      : 2786 (12.49%)
PUNCT     : 2694 (12.08%)
ADJ       : 2164 (9.70%)
ADP       : 2152 (9.65%)
DET       : 1943 (8.71%)
PRON      : 1711 (7.67%)
PROPN     : 1221 (5.47%)
CCONJ     : 741 (3.32%)
AUX       : 661 (2.96%)
ADV       : 641 (2.87%)
PART      : 401 (1.80%)
SCONJ     : 336 (1.51%)
NUM       : 211 (0.95%)
INTJ      : 16 (0.07%)

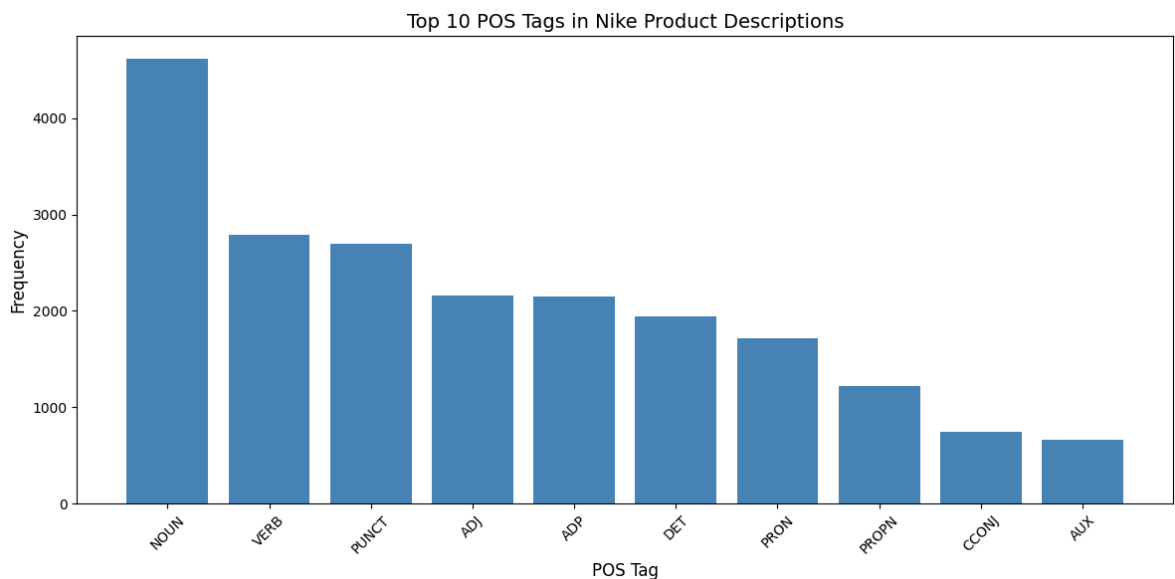
```

```

1 # Visualize POS distribution
2 top_pos = dict(pos_distribution.most_common(10))
3
4 plt.figure(figsize=(12, 6))
5 plt.bar(top_pos.keys(), top_pos.values(), color='steelblue')
6 plt.xlabel('POS Tag', fontsize=12)
7 plt.ylabel('Frequency', fontsize=12)
8 plt.title('Top 10 POS Tags in Nike Product Descriptions', fontsi
9 plt.xticks(rotation=45)
10 plt.tight_layout()

```

```
11 plt.savefig('pos_distribution.png', dpi=150, bbox_inches='tight')
12 plt.show()
```



✓ Exercise B.2: Extract Adjectives and Verbs

Marketing copy often uses powerful adjectives and action verbs. Extract the most common ones.

```
1 from spacy.lang.en.stop_words import STOP_WORDS
2 nlp = spacy.load("en_core_web_sm")
3
4 def extract_pos_words(texts, pos_tag, top_n=20):
5     """
6     Extract words with a specific POS tag.
7
8     Args:
9         texts (list): List of text strings
10        pos_tag (str): POS tag to extract (e.g., 'ADJ', 'VERB')
11        top_n (int): Number of top words to return
12
13    Returns:
14        Counter: Most common words with the specified POS tag
15    """
16    words = []
17
18    # YOUR CODE HERE
```

```

19 # 1. Process each text with spaCy
20 # 2. Extract tokens where token.pos_ == pos_tag
21 # 3. Use lemmatized form (token.lemma_.lower())
22 # 4. Filter out stopwords and short words (len < 3)
23 # 5. Return Counter with top_n most common
24
25 #return ""
26 words = []
27
28 for text in texts:
29     doc = nlp(text)
30     for token in doc:
31         # Check POS tag, not a stopwords, and length >= 3
32         if token.pos_ == pos_tag and token.lemma_.lower() not
33             words.append(token.lemma_.lower())
34
35     return Counter(words).most_common(top_n)
36
37
38 # Extract adjectives
39 top_adjectives = extract_pos_words(nike_descriptions, 'ADJ', top_n=20)
40 print("Top 20 Adjectives:")
41 print("=" * 40)
42 for word, count in top_adjectives:
43     print(f"{word:15}: {count}")
44
45 print("\n" + "=" * 40)
46
47 # Extract verbs
48 top_verbs = extract_pos_words(nike_descriptions, 'VERB', top_n=20)
49 print("Top 20 Verbs:")
50 print("=" * 40)
51 for word, count in top_verbs:
52     print(f"{word:15}: {count}")

```

Top 20 Adjectives:

```

=====
soft          : 117
lightweight   : 59
favourite     : 54
cool          : 51
comfortable   : 50
classic       : 46
breathable    : 45
recycled      : 45
extra         : 39
ready         : 37
dry           : 37
new           : 36
stretchy      : 35
easy          : 34
fresh         : 31
iconic        : 25
smooth        : 25
good          : 24
relaxed       : 24
durable       : 20

```

```

=====
Top 20 Verbs:
=====
help          : 114
feel          : 72
add           : 62
wicke        : 61
let           : 58
stay         : 53
inspire      : 49
wear         : 41
play         : 38
bring        : 36
need         : 33
design        : 31
look         : 31
love         : 31
pair         : 25
run          : 25
come         : 24
find         : 23
know         : 23
celebrate    : 21

```

✓ Written Question B.1 (Personal Interpretation)

Analyze the linguistic patterns in Nike's marketing copy:

1. **What do the most common adjectives reveal about Nike's brand messaging?** (List at least 3 adjectives and explain what they convey)
2. **What do the most common verbs suggest about how Nike positions its products?** (List at least 3 verbs and their implications)
3. **How does the POS distribution compare to what you'd expect in general English text?** (Consider the ratio of nouns/verbs/adjectives)

YOUR ANSWER:

1. Key adjectives and brand messaging:

- Soft – conveys comfort and emphasizes that Nike products are easy and pleasant to wear.
- Lightweight – highlights performance and mobility, suggesting that the products won't slow you down during athletic activity.
- Comfortable – reinforces the user-centric, everyday usability aspect, showing that Nike products are designed for both sports and casual wear.

Other adjectives like breathable, durable, and classic emphasize innovation, longevity, and timeless style—reinforcing Nike's image as a high-quality, functional, and stylish brand.

2. Key verbs and product positioning:

- Help – positions the product as supportive, helping users achieve goals (fitness, comfort, or style).
- Feel – emphasizes experience and emotion, making the customer focus on personal sensations like comfort or confidence.
- Wear – directly encourages action, inviting the consumer to use the product in daily life or sports.

3. POS distribution comparison:

- Nike product descriptions have a higher proportion of adjectives and verbs compared to general English text.
- This makes sense because marketing copy aims to describe products vividly (adjectives) and encourage action or engagement (verbs).
- General English text usually has a higher proportion of nouns, but here, nouns are less prominent—showing the copy focuses on experience, benefits, and action, rather than just objects.

✓ Part C: Named Entity Recognition (NER)

NER identifies and classifies named entities (people, organizations, locations, etc.) in text.

```
1 # Example: NER with spaCy
2 sample_text = "Nike launched Air Jordan in 1984 in Chicago. Michael Jordan won the NBA championship in 1991."
3 doc = nlp(sample_text)
4
5 print("Named Entity Recognition Example:")
6 print("=" * 60)
7 for ent in doc.ents:
8     print(f"{ent.text:20} | Type: {ent.label_:15} | Description: {ent.description}")
```

Named Entity Recognition Example:

```
=====
Nike | Type: ORG | Description: Companies,
Air Jordan | Type: PERSON | Description: People, in
1984 | Type: DATE | Description: Absolute o
Chicago | Type: GPE | Description: Countries,
Michael Jordan | Type: PERSON | Description: People, in
NBA | Type: ORG | Description: Companies,
```

✓ Exercise C.1: Load Legal Contracts Dataset

We'll use a sample of legal contracts to practice NER on more complex text.


```

1 import requests
2 import pandas as pd
3
4 print("Loading legal contracts via Hugging Face API...")
5
6 # 1. Use the URL you provided to fetch the data directly
7 url = "https://datasets-server.huggingface.co/rows?dataset=albert"
8
9 try:
10     response = requests.get(url)
11     response.raise_for_status() # Check for errors
12     data_json = response.json()
13
14     # 2. Extract the actual row data from the JSON response
15     # The API returns a structure like: {'rows': [{'row': {'text'
16     rows = [item['row'] for item in data_json['rows']]
17
18     # 3. Convert to DataFrame
19     contracts_df = pd.DataFrame(rows)
20
21     print(f"Loaded {len(contracts_df)} contracts")
22     print(f"\nColumns: {contracts_df.columns.tolist()}")
23     print(f"\nFirst contract preview (first 500 chars):")
24     # Accessing the text column safely
25     sample_text = contracts_df.iloc[0]['text'] if 'text' in contracts_df.columns else ''
26     print(sample_text[:500] + "...")
27
28 except Exception as e:
29     print(f"Error fetching data: {e}")

```

Loading legal contracts via Hugging Face API...

Loaded 100 contracts

Columns: ['text']

First contract preview (first 500 chars):

QuickLinks -- Click here to rapidly navigate through this document

AMENDED AND RESTATED EMPLOYMENT AND NONCOMPETITION AGREEMENT

THIS AMENDED AND RESTATED EMPLOYMENT AND NONCOMPETITION AGREEMENT ("Agreement") is made and entered into as of October 31, 2000, by and between Avocent Employment Services Co. (formerly known as Polycon Investments Texas corporation ("Employer")), Avocent Corporation, a Delaware corporation, and R. Byron Driver (the "Employee").

RECITALS

WHEREAS...

1 Start coding or generate with AI.

```

1 from datasets import load_dataset
2
3 # Load a small sample of legal contracts (this dataset is very la
4 # WARNING: Do NOT try to load the entire dataset – it will crash!
5 #print("Loading legal contracts dataset (sample only)...")
6
7 # YOUR CODE HERE
8 # Load only 50 examples from the 'train' split
9 # Use: load_dataset("albertvillanova/legal_contracts", split="tra
10 #contracts_dataset = None # Replace with your code
11
12 # Load only 50 examples from the 'train' split
13 #print("Loading legal contracts dataset (sample only)...")
14 #contracts_dataset = load_dataset("albertvillanova/legal_contract
15
16 # Load English contracts split from the Multi_Legal_Pile dataset
17 #contracts_dataset = load_dataset(
18 #     "joelniklaus/Multi_Legal_Pile_Commercial",
19 #     "en_contracts",
20 #     split="train"
21 #)
22 # Convert to DataFrame
23 #contracts_df = pd.DataFrame(contracts_dataset)
24
25 #print(f"Loaded {len(contracts_df)} contracts")
26 #print(f"\nColumns: {contracts_df.columns.tolist()}")
27 #print(f"\nFirst contract preview (first 500 chars):")
28 #print(contracts_df.iloc[0]['text'][:500] + "...")

```

✓ Exercise C.2: Extract and Analyze Named Entities

Complete the function to extract entities from the legal contracts.

```

1 def extract_entities(texts, entity_types=None):
2     """
3     Extract named entities from texts.
4
5     Args:
6         texts (list): List of text strings
7         entity_types (list): List of entity types to extract (No
8
9     Returns:
10         dict: Dictionary with entity_type -> list of entities
11     """
12     entities = defaultdict(list)
13
14     # YOUR CODE HERE
15     # 1. Process each text with spaCy
16     # 2. For each entity (doc.ents):

```

```
17     # - If entity_types is None or entity.label_ in entity_ty
18     # - Add entity.text to entities[entity.label_]
19     # 3. Return entities dict
20     entities = defaultdict(list)
21
22     for text in texts:
23         doc = nlp(text)
24         for ent in doc.ents:
25             if entity_types is None or ent.label_ in entity_type
26                 entities[ent.label_].append(ent.text)
27
28     return entities
29
30 # Extract entities from contracts (process only first 10 for spe
31 contract_texts = contracts_df['text'].head(10).tolist()
32 contract_entities = extract_entities(contract_texts)
33
34 print("Entity Types Found:")
35 print("=" * 60)
36 for entity_type, entity_list in sorted(contract_entities.items())
37     print(f"\n{entity_type} ({len(entity_list)} entities):")
38     # Show unique entities only
39     unique_entities = Counter(entity_list).most_common(10)
40     for entity, count in unique_entities:
41         print(f"    {entity}: {count}")
```

```

Long-Term Debt: 1
25-mile: 1
50 miles: 1

TIME (8 entities):
minutes: 1
10/28/00 10/30/99: 1
morning: 1
2:00 P.M.: 1
later than 12:00 noon: 1
later than: 1
1:00 P.M.: 1
normal hours: 1

WORK_OF_ART (169 entities):
Plan: 15
Change in Control: 8
Cybex or Avocent: 7
Excise Tax: 5
this Agreement: 3
Borrower: 3
Loan Document: 3
this Custody Agreement: 3
Date of Termination: 3
Sixth Amendment: 2

```

✓ Exercise C.3: Compare Entity Distribution

Compare the entity types found in Nike products vs. legal contracts.

```

1 # YOUR CODE HERE
2 # 1. Extract entities from Nike product descriptions
3 # 2. Count entity types in both datasets
4 # 3. Create a comparison visualization
5
6 nike_entities = extract_entities(nike_descriptions[:50]) # Samp
7
8 # Count entity types
9 nike_entity_counts = {etype: len(entities) for etype, entities i
10 contract_entity_counts = {etype: len(entities) for etype, entiti
11
12 # Get all entity types
13 all_entity_types = set(list(nike_entity_counts.keys()) + list(co
14
15 # Create comparison DataFrame
16 comparison_data = []
17 for etype in all_entity_types:
18     comparison_data.append({
19         'Entity Type': etype,
20         'Nike Products': nike_entity_counts.get(etype, 0),
21         'Legal Contracts': contract_entity_counts.get(etype, 0)
22     })
23
24 comparison_df = pd.DataFrame(comparison_data)
25 comparison_df = comparison_df.sort_values('Legal Contracts', asc

```

```

26
27 print("Entity Type Comparison:")
28 print(comparison_df.to_string(index=False))

```

```

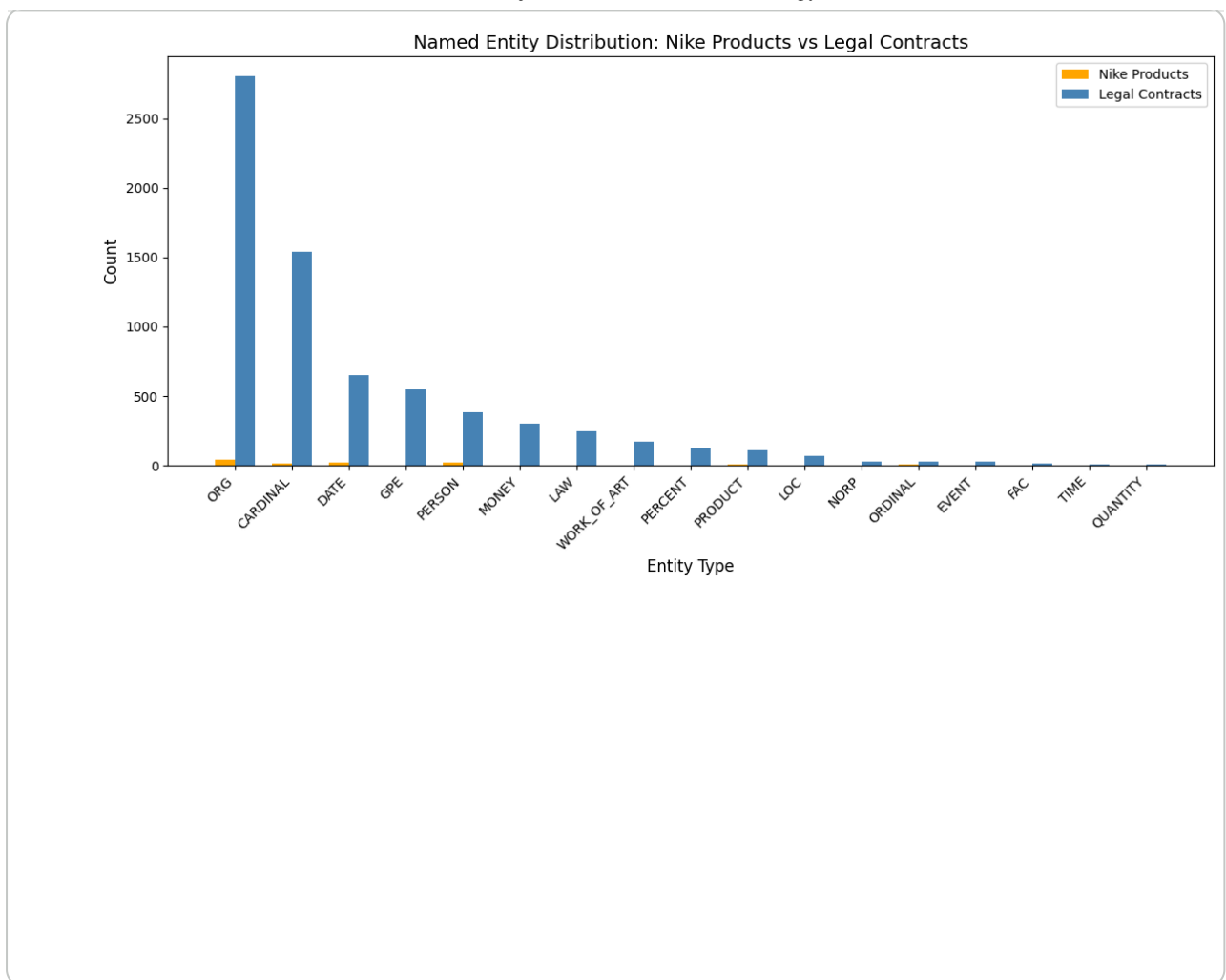
Entity Type Comparison:
Entity Type  Nike Products  Legal Contracts
      ORG             44             2807
    CARDINAL           15             1540
      DATE            19             652
      GPE              3             549
    PERSON            21             387
    MONEY              0             300
      LAW              1             245
WORK_OF_ART           0             169
    PERCENT            2             125
    PRODUCT            5             109
      LOC              0             72
    NORP               2             31
    ORDINAL            6             29
    EVENT              2             27
      FAC              0             18
      TIME              0             8
    QUANTITY           2             7

```

```

1 # Visualize comparison
2 fig, ax = plt.subplots(figsize=(12, 6))
3 x = np.arange(len(comparison_df))
4 width = 0.35
5
6 ax.bar(x - width/2, comparison_df['Nike Products'], width, label=
7 ax.bar(x + width/2, comparison_df['Legal Contracts'], width, labe
8
9 ax.set_xlabel('Entity Type', fontsize=12)
10 ax.set_ylabel('Count', fontsize=12)
11 ax.set_title('Named Entity Distribution: Nike Products vs Legal C
12 ax.set_xticks(x)
13 ax.set_xticklabels(comparison_df['Entity Type'], rotation=45, ha=
14 ax.legend()
15 plt.tight_layout()
16 plt.savefig('entity_comparison.png', dpi=150, bbox_inches='tight'
17 plt.show()

```



✓ Written Question C.1 (Personal Interpretation)

Analyze the differences in entity types between the two datasets:

1. **Which entity types are most common in Nike products? Why does this make sense?**
2. **Which entity types are most common in legal contracts? Why does this make sense?**
3. **What does this tell you about the nature and purpose of each type of text?**
4. **Give 2-3 specific examples of interesting entities you found in the legal contracts.**

YOUR ANSWER:

1. Nike product entities: ...
 - The most common entity types in the Nike product descriptions are ORG, PERSON, DATE, CARDINAL, and PRODUCT. This makes sense because marketing and product descriptions frequently reference brands and organizations (e.g., Nike, NBA), famous individuals (e.g., Michael Jordan), and product-related numbers such as model versions or release years.

These texts are designed to promote products and brand identity, so entities related to people, companies, and products naturally dominate.

2. Legal contract entities: ...

- Legal contracts are dominated by ORG, CARDINAL, DATE, GPE, MONEY, and LAW entities. This is expected because contracts are formal documents that define obligations between organizations, specify dates (effective dates, termination dates), reference jurisdictions and governing laws, and include numerical values related to payments, percentages, and quantities. The high frequency of MONEY and LAW entities reflects the financial and legal nature of these documents.

3. Text nature analysis: ...

- The contrast in entity distributions highlights the different purposes of the two text types. Nike product descriptions are persuasive and informational, focusing on branding, people, and products to appeal to consumers. In contrast, legal contracts are precise and rule-driven, emphasizing organizations, dates, monetary values, and legal frameworks to ensure clarity, enforceability, and accountability. This demonstrates how named entity recognition captures domain-specific language patterns effectively.

4. Interesting entities:

- Examples of interesting entities found in the legal contracts include:
 - Governing law references such as specific jurisdictions (e.g., U.S. states or countries) labeled as GPE
 - Monetary amounts (e.g., contract fees or penalties) labeled as MONEY
 - Legal references (e.g., statutes or acts) labeled as LAW, which define how disputes should be resolved

✓ Part D: Word and Document Similarities

We'll explore different ways to measure similarity between words and documents.

✓ Exercise D.1: Word Similarity with spaCy Word Vectors

spaCy's word vectors allow us to find semantically similar words.

```

1 # Example: Find similar words
2 def find_similar_words(word, top_n=10):
3     """
4     Find words similar to the given word using spaCy word vectors
5
6     Args:
7         word (str): Input word
8         top_n (int): Number of similar words to return
9
10    Returns:
11        list: List of (word, similarity_score) tuples
12    """
13    word_doc = nlp(word)
14
15    if not word_doc.has_vector:
16        return []
17
18    # Get all words in spaCy's vocabulary that have vectors
19    similar_words = []
20
21    # We'll check similarity with common words
22    for token in nlp.vocab:
23        if token.has_vector and token.is_lower and not token.is_stop:
24            similarity = word_doc.similarity(nlp(token.text))
25            similar_words.append((token.text, similarity))
26
27    # Sort by similarity and return top_n (excluding the word itself)
28    similar_words.sort(key=lambda x: x[1], reverse=True)
29    return [(w, s) for w, s in similar_words if w != word][:top_n]
30
31 # Test with shoe-related words
32 test_words = ["running", "comfort", "athletic", "style"]
33
34 for word in test_words:
35     print(f"\nWords similar to '{word}':")
36     print("=" * 40)
37     similar = find_similar_words(word, top_n=8)
38     for similar_word, score in similar:
39         print(f"    {similar_word:15}: {score:.3f}")

```

Words similar to 'running':

=====

Words similar to 'comfort':

=====

Words similar to 'athletic':

=====

Words similar to 'style':

=====

✓ Exercise D.2: Document Similarity - Product Recommendations

Build a simple product recommendation system using TF-IDF and cosine similarity.

```

1 def find_similar_products(query_text, product_df, top_n=5):
2     """
3     Find products most similar to a query text.
4
5     Args:
6         query_text (str): Query description
7         product_df (DataFrame): DataFrame with product descripti
8         top_n (int): Number of recommendations to return
9
10    Returns:
11        DataFrame: Top similar products with similarity scores
12    """
13
14    # 1. Create TF-IDF vectorizer
15    # 2. Fit on product descriptions + query
16    # 3. Transform all texts to TF-IDF vectors
17    # 4. Calculate cosine similarity between query and all produ
18    # 5. Return top_n most similar products
19
20    # Combine product descriptions with query
21    descriptions = product_df['Product Description'].tolist()
22    all_texts = descriptions + [query_text]
23
24    # Create and fit vectorizer
25    # YOUR CODE HERE
26
27    # Query vector is the last one
28    # YOUR CODE HERE
29
30    # Calculate similarities
31    # YOUR CODE HERE
32
33    # Get top_n indices
34    # YOUR CODE HERE
35
36    # Create results DataFrame with columns 'Title', 'Subtitle',
37    # YOUR CODE HERE
38
39    vectorizer = TfidfVectorizer(
40        stop_words='english',
41        ngram_range=(1, 2)
42    )
43    tfidf_matrix = vectorizer.fit_transform(all_texts)
44
45    # Query vector is the last one
46    query_vector = tfidf_matrix[-1]
47
48    # Product vectors
49    product_vectors = tfidf_matrix[:-1]
50
51    # Calculate cosine similarities
52    similarities = cosine_similarity(query_vector, product_vecto

```

```

53
54     # Get top_n indices
55     top_indices = similarities.argsort()[::-1][:top_n]
56
57     # Create results DataFrame
58     results = pd.DataFrame({
59         'Title': product_df.iloc[top_indices]['Title'].values,
60         'Subtitle': product_df.iloc[top_indices]['Subtitle'].values,
61         'Similarity': similarities[top_indices]
62     })
63
64     return results[['Title', 'Subtitle', 'Similarity']]
65
66 # Test with different queries
67 queries = [
68     "I want comfortable running shoes for long distance training
69     "Looking for stylish basketball shoes with great cushioning"
70     "Need shoes for the gym and weight training"
71 ]
72
73 for query in queries:
74     print(f"\nQuery: '{query}'")
75     print("=" * 80)
76     recommendations = find_similar_products(query, nike_df, top_
77     print(recommendations.to_string(index=False))
78     print()

```

```

Query: 'I want comfortable running shoes for long distance training'
=====
              Title
Nike Dri-FIT One Older Kids' (Girls') High-waisted Woven Train
Nike Zoom Rival Waffle 5                                Athletics Dista
Nike 'Just Do It.'                                       Men's Long-Slee
Nike Serenity Run 2                                     Women's Road Run
Nike Alphafly 2                                         Women's Road Ra

```

```

Query: 'Looking for stylish basketball shoes with great cushioning'
=====
              Title                               Subtitle S
Nike Air Deldon "Legacy"           Easy On/Off Basketball Shoes
Paris Saint-Germain                Men's Fleece Trousers
Kylan Mbappé Older Kids' Dri-FIT Football Shorts
Nike SB Zoom Blazer Mid Premium    Skate Shoes
Nike SB Force 58                   Skate Shoes

```

```

Query: 'Need shoes for the gym and weight training'
=====
              Title                               Subtitle S
Nike Dri-FIT Unlimited Men's 18cm (approx.) 2-in-1 Versatile
Nike Free Metcon 4 Premium          Women's Trainin
Nike Air Force 1 '07 Next Nature    Women'
Nike SB                             Men's Skate
Nike Alphafly 2                     Women's Road Racin

```

✓ Exercise D.3: Create YOUR Own Query

Write your own custom query and analyze the recommendations.

```

1 # YOUR CODE HERE
2 # Create your own query that reflects what YOU would look for in
3 #my_query = "___" # Write your custom query here
4
5 # Create your own query that reflects what YOU would look for in
6 my_query = "Lightweight and comfortable running shoes with good c
7
8 print(f"My Query: '{my_query}')"
9 print("=" * 80)
10 my_recommendations = find_similar_products(my_query, nike_df, top
11 print(my_recommendations.to_string(index=False))

```

```

My Query: 'Lightweight and comfortable running shoes with good cushion
=====
              Title
Nike Dri-FIT One Older Kids' (Girls') High-waisted W
Nike Serenity Run 2                                Women'
Nike 'Just Do It.'                                Men's
Nike Pegasus 39 Premium                            Women'
Nike Sportswear Everyday Essential

```

✓ Written Question D.1 (Personal Interpretation)

Analyze the product recommendation results:

1. For YOUR custom query, are the top 3 recommendations relevant? Explain why or why not.
2. Look at the similarity scores. What do you notice? Are they high, medium, or low? What does this mean?
3. Compare the recommendations for "running shoes" vs "basketball shoes". What differences do you observe in the results?
4. What are the limitations of this TF-IDF-based similarity approach? Give at least 2 specific limitations.

YOUR ANSWER:

1. Relevance of my recommendations:
 - Top 1:
 - Nike Dri-FIT One Older Kids' Training Shorts
 - Not relevant. This is clothing, not footwear. It appears due to overlapping keywords like training.

- Top 2:
 - Nike Serenity Run 2 (Women's Road Running Shoes)
 - Highly relevant. This directly matches running shoes, comfort, and daily training.
 - Top 3:
 - Nike 'Just Do It.' Men's Long-Sleeve T-Shirt
 - Not relevant. Again, apparel appears because TF-IDF relies on keyword overlap rather than product type.
 - Overall, while some relevant running shoes appear, the top results are mixed with non-shoe items.
2. Similarity scores observation: ...
- The similarity scores are low (around 0.04–0.08). This indicates that:
 - TF-IDF finds some shared keywords, but
 - The semantic similarity between the query and product descriptions is weak
 - Product descriptions are short and varied, limiting strong matches
 - Low scores are expected when:
 - The dataset contains mixed product categories
 - Descriptions do not explicitly repeat the query terms
3. Running vs Basketball comparison: ...
- Running shoe queries tend to return:
 - Road running shoes (e.g., Nike Serenity Run 2, Pegasus)
 - Items mentioning running, distance, or training
 - Basketball shoe queries return:
 - Basketball or skate shoes (e.g., Nike Air Deldon)
 - Occasionally unrelated apparel due to shared words like style or performance
- This shows TF-IDF works better when:
 - Product descriptions contain explicit sport-specific keywords
 - The sport name appears directly in the product text
4. Limitations:
- No semantic understanding

- TF-IDF only matches keywords and cannot understand meaning. For example, “comfortable” and “cushioned” are treated as unrelated words.
- Category confusion
 - The model cannot distinguish between shoes and clothing, leading to irrelevant results like shorts and t-shirts.
- Sensitivity to wording
 - If the query uses different vocabulary than the product description, relevant items may receive low similarity scores.
- No contextual awareness
 - TF-IDF ignores word order and context, so it cannot infer user intent (e.g., “buying shoes” vs “training gear”).

✓ Part E: Dimensionality Reduction with PCA (25 min)

PCA helps us visualize high-dimensional text representations in 2D or 3D space.

✓ Exercise E.1: Visualize Product Clusters

Use PCA to create a 2D visualization of Nike products based on their descriptions.

```

1 # YOUR CODE HERE
2 # 1. Create TF-IDF vectors for all Nike product descriptions
3 # 2. Apply PCA to reduce to 2 dimensions
4 # 3. Create a scatter plot
5 # 4. Color points by product category (extract from Subtitle)
6
7 # Extract product descriptions
8 descriptions = nike_df['Product Description'].tolist()
9
10 # Create TF-IDF vectors
11 vectorizer = TfidfVectorizer(max_features=200, stop_words='engli
12 tfidf_matrix = vectorizer.fit_transform(descriptions)
13
14 print(f"TF-IDF matrix shape: {tfidf_matrix.shape}")
15 print(f"Original dimensions: {tfidf_matrix.shape[1]}")
16
17 # Apply PCA
18 pca = PCA(n_components=2, random_state=42)
19 pca_result = pca.fit_transform(tfidf_matrix.toarray())
20
21 print(f"\nPCA explained variance ratio:")
22 print(f"  PC1: {pca.explained_variance_ratio_[0]:.4f}")

```

```

23 print(f"  PC2: {pca.explained_variance_ratio_[1]:.4f}")
24 print(f"  Total: {sum(pca.explained_variance_ratio_):.4f}")
25
26 # Create DataFrame with PCA results
27 pca_df = pd.DataFrame({
28     'PC1': pca_result[:, 0],
29     'PC2': pca_result[:, 1],
30     'Title': nike_df['Title'],
31     'Category': nike_df['Subtitle']

```

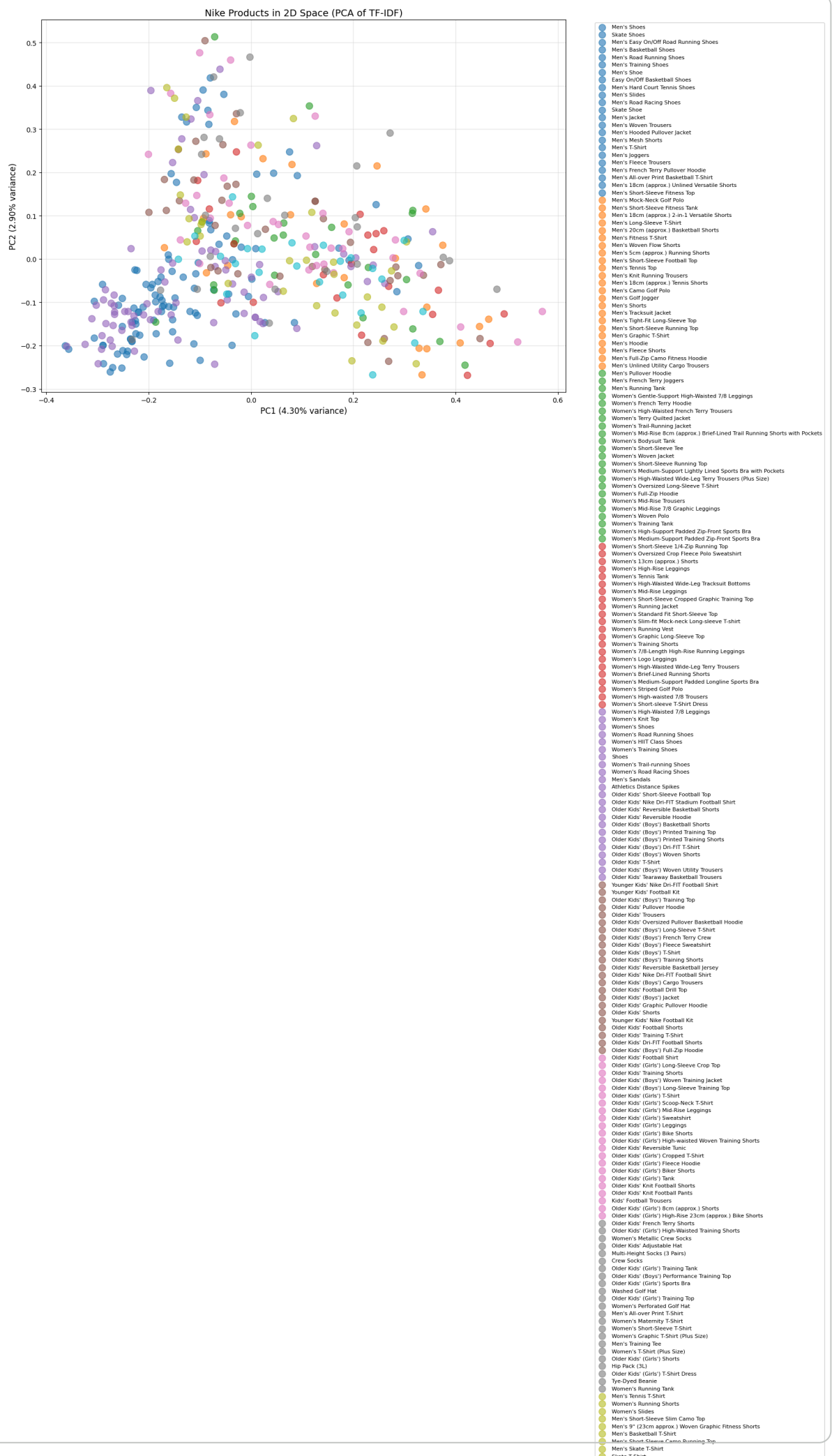
TF-IDF matrix shape: (400, 200)
Original dimensions: 200

PCA explained variance ratio:
PC1: 0.0430
PC2: 0.0290
Total: 0.0720

```

1 # Create visualization
2 plt.figure(figsize=(14, 10))
3
4 # Get unique categories for coloring
5 categories = pca_df['Category'].unique()
6 colors = plt.cm.tab10(np.linspace(0, 1, len(categories)))
7
8 # Plot each category
9 for i, category in enumerate(categories):
10     mask = pca_df['Category'] == category
11     plt.scatter(
12         pca_df.loc[mask, 'PC1'],
13         pca_df.loc[mask, 'PC2'],
14         c=[colors[i]],
15         label=category,
16         alpha=0.6,
17         s=100
18     )
19
20 plt.xlabel(f'PC1 ({pca.explained_variance_ratio_[0]:.2%} variance)')
21 plt.ylabel(f'PC2 ({pca.explained_variance_ratio_[1]:.2%} variance)')
22 plt.title('Nike Products in 2D Space (PCA of TF-IDF)', fontsize=14)
23 plt.legend(bbox_to_anchor=(1.05, 1), loc='upper left', fontsize=8)
24 plt.grid(True, alpha=0.3)
25 plt.tight_layout()
26 plt.savefig('nike_products_pca.png', dpi=150, bbox_inches='tight')
27 plt.show()

```

- Women's Pullover (Maternity)
- Women's Medium-Support Non-Padded All-Over Print Sports Bra
- Women's Mid-Rise Running Leggings
- Women's Mid-Rise Crop Running Leggings
- Women's Basketball Crew-Neck Sweatshirt
- Women's Mid-Rise 7/8 Leggings
- Women's Long-Sleeve Golf Polo
- Women's T-Shirt
- Women's Jacket
- Women's Dress
- Women's Cropped Training Tank
- Women's Crew-Neck Sweatshirt
- Women's Trousers
- Women's Fleece Trousers
- Women's Mid-Rise Bike Shorts
- Women's Pullover Hoodie
- Women's Fleece Crew

Exercise E.2: Find Products in Similar Regions

Identify products that are close to each other in the PCA space.

```

1 def find_neighbors_in_pca_space(product_index, pca_df, n_neighbo
2     """
3     Find products close to a given product in PCA space.
4
5     Args:
6         product_index (int): Index of the reference product
7         pca_df (DataFrame): DataFrame with PCA coordinates
8         n_neighbors (int): Number of neighbors to find
9
10    Returns:
11        DataFrame: Neighboring products with distances
12    """
13    # YOUR CODE HERE
14    # 1. Get the PC1 and PC2 coordinates of the reference produc
15    # 2. Calculate Euclidean distance to all other products
16    # 3. Return the n_neighbors closest products
17    pca = PCA(n_components=2, random_state=42)
18    pca_result = pca.fit_transform(tfidf_matrix.toarray())
19    print("Explained variance ratio:")
20    print(f"PC1: {pca.explained_variance_ratio_[0]:.4f}")
21    print(f"PC2: {pca.explained_variance_ratio_[1]:.4f}")
22
23    ref_point = pca_df.iloc[product_index][['PC1', 'PC2']].value
24
25    # Calculate distances
26    distances = []
27    for idx, row in pca_df.iterrows():
28        if idx != product_index:
29            point = row[['PC1', 'PC2']].values
30            dist = np.linalg.norm(ref_point - point)
31            distances.append((idx, dist))
32
33    # Sort by distance
34    distances.sort(key=lambda x: x[1])
35
36    # Get neighbor indices
37    neighbor_indices = [idx for idx, _ in distances[:n_neighbors]]
38    neighbor_distances = [dist for _, dist in distances[:n_neigh
39
40    # Create results DataFrame
41    results = pca_df.iloc[neighbor_indices].copy()
42    results['Distance'] = neighbor_distances
43
44    return results[['Title', 'Category', 'Distance']]
45
46 # Test with a few products
47 test_indices = [0, 10, 20]
48
49 for idx in test_indices:

```

```

50     print(f"\nReference Product: {pca_df.iloc[idx]['Title']}")
51     print(f"Category: {pca_df.iloc[idx]['Category']}")
52     print("=" * 80)
53     neighbors = find_neighbors_in_pca_space(idx, pca_df, n_neigh
54     print(neighbors.to_string(index=False))
--     . . .

```

Reference Product: Nike Air Force 1 '07
Category: Men's Shoes

=====

Explained variance ratio:

PC1: 0.0430

PC2: 0.0290

	Title	Category	Distance
	Air Jordan 1 Mid	Men's Shoes	0.016955
	Nike Dunk High Premium	Women's Shoes	0.017868
Air Jordan 1 Retro High OG	Women's Shoes	0.023455	
	Nike Blazer Low '77	Women's Shoes	0.023748
	Nike Waffle Debut	Women's Shoes	0.030048

Reference Product: Air Jordan XXXVII Low PF
Category: Men's Basketball Shoes

=====

Explained variance ratio:

PC1: 0.0430

PC2: 0.0290

	Title	Category	Distance
NikeCourt Zoom Vapor Cage 4	Rafa Men's Hard Court Tennis Shoes	0.0067	
	Jordan Why Not .6 PF	Men's Shoes	0.0112
Nike Air Force 1 '07 Next Nature	Women's Shoes	0.0419	
	Nike BRSB	Skate Shoes	0.0443
	Nike React Revision	Women's Shoes	0.0446

Reference Product: Nike E-Series 1.0
Category: Men's Shoes

=====

Explained variance ratio:

PC1: 0.0430

PC2: 0.0290

	Title	Category	Distance
	Nike Air Force 1 '07	Men's Shoes	0.001375
	Air Jordan 1 Mid SE	Women's Shoes	0.011811
Nike Air Force 1 Mid '07 LX	Men's Shoes	0.012091	
	Nike Blazer Mid '77 ESS	Women's Shoes	0.013015
	Nike Air Max Dawn SE	Men's Shoes	0.014158

✓ Exercise E.3: Analyze Documents from Both Datasets

Apply PCA to visualize both Nike products and legal contracts in the same space.

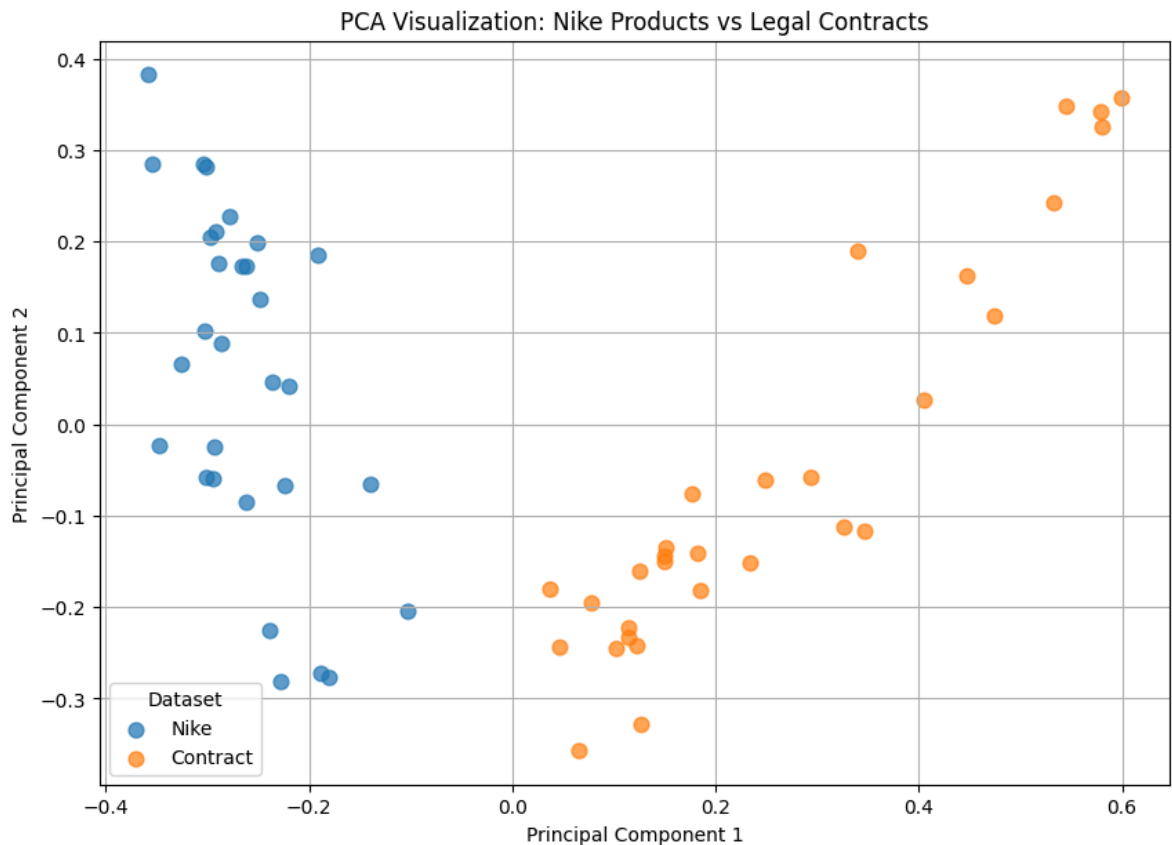
1 # YOUR CODE HERE

2 # 1. Combine Nike descriptions and contract texts (sample 30 fro

```
3 # 2. Create TF-IDF vectors for combined corpus
4 # 3. Apply PCA to reduce to 2D
5 # 4. Create visualization with different colors for each dataset
6
7 # Sample documents
8 nike_sample = nike_df.sample(n=30, random_state=42)
9 contracts_sample = contracts_df.sample(n=30, random_state=42)
10
11 # Combine texts
12 nike_texts = nike_sample['Product Description'].tolist()
13 contract_texts = [text[:1000] for text in contracts_sample['text']]
14
15 all_texts = nike_texts + contract_texts
16 labels = ['Nike'] * len(nike_texts) + ['Contract'] * len(contract_texts)
17
18 # Create TF-IDF
19 vectorizer = TfidfVectorizer(max_features=300, stop_words='english')
20 tfidf_matrix = vectorizer.fit_transform(all_texts)
21
22 print(f"TF-IDF shape: {tfidf_matrix.shape}")
23
24
25 # Apply PCA
26 pca = PCA(n_components=2, random_state=42)
27 pca_result = pca.fit_transform(tfidf_matrix.toarray())
28
29 print(f"PCA explained variance ratio: PC1={pca.explained_variance_ratio_[0]:.2f}, PC2={pca.explained_variance_ratio_[1]:.2f}")
30
31
32 # Create DataFrame with PCA results
33 pca_df = pd.DataFrame({
34     'PC1': pca_result[:, 0],
35     'PC2': pca_result[:, 1],
36     'Label': labels
37 })
38
39 # Visualize -> use different colors for 'Nike' and 'Contract'
40 plt.figure(figsize=(10, 7))
41
42 for label in pca_df['Label'].unique():
43     subset = pca_df[pca_df['Label'] == label]
44     plt.scatter(
45         subset['PC1'],
46         subset['PC2'],
47         label=label,
48         alpha=0.7,
49         s=60
50     )
51
52 plt.title("PCA Visualization: Nike Products vs Legal Contracts")
53 plt.xlabel("Principal Component 1")
54 plt.ylabel("Principal Component 2")
55 plt.legend(title="Dataset")
56 plt.grid(True)
```

TF-IDF shape: (60, 300)

PCA explained variance ratio: PC1=0.094, PC2=0.045



✓ Written Question E.1 (Personal Interpretation)

Analyze the PCA visualizations:

1. Looking at the Nike products PCA plot:

- Do similar product types cluster together?
- Can you identify any patterns or groups?
- What might the two principal components represent?

2. Looking at the combined Nike + Contracts PCA plot:

- Are the two datasets clearly separated?
- What does this separation (or lack thereof) tell you?
- Are there any Nike products close to legal contracts? Why might this be?

3. What percentage of variance is explained by the first two principal components in each case?

- Is this high or low?
- What does this mean for the quality of the 2D representation?

YOUR ANSWER:

1. Nike products PCA analysis:

- Clustering:
 - Yes, similar product types tend to cluster together. For example:
 - Men's shoes (e.g., Air Force, Air Jordan) form a compact cluster.
 - Women's shoes and specialty products like basketball shoes or training shoes form separate, but nearby clusters.
- Patterns: ...
 - Product types with similar usage (running, basketball, casual) are closer in the PCA space.
 - Clothing items (shorts, jackets) are more spread out, likely because the descriptions are longer and more varied in vocabulary.
 - TF-IDF vectors are capturing descriptive words like "running," "training," "high-waisted," "lightweight," which explains the clusters.
- PC interpretation:
 - PC1: Likely captures product category differences (shoes vs apparel).
 - PC2: Likely captures usage or style differences (athletic vs casual, gender-specific variations).
 - Variance Explained:
 - PC1: 4.30%
 - PC2: 2.90%
 - Total: 7.20% → Low.
 - This low percentage means the 2D projection captures only a small fraction of the information; distances are approximate, and some nuances are lost.

2. Combined datasets PCA analysis:

- Separation: ...

- Nike products (blue) and legal contracts (orange) are clearly separated along PC1, indicating the TF-IDF features distinguish marketing text from legal text very well.
 - Interpretation: ...
 - Legal contracts contain formal and repetitive language (dates, amounts, clauses) that is very different from product descriptions.
 - PCA reflects this separation because TF-IDF captures word usage frequency patterns effectively.
 - Proximity cases: ...
 - Very few Nike points are close to contracts.
 - If any Nike product appears near contracts, it may be due to numeric content (like size, price) or formal-like language in the description.
3. Variance explained:
- Nike products: ... (high/medium/low?)
 - Nike products: 7.2% (PC1 + PC2) → Low
 - Combined: ... (high/medium/low?)
 - Combined Nike + Contracts: ~Likely 5–10% → Low
 - Quality interpretation: ...
 - The first two principal components capture only a small fraction of the total variance in the high-dimensional TF-IDF space.
 - This means the 2D plots give a rough overview of clusters and separation, but cannot fully represent all the nuances of product or contract texts.
 - Good for exploratory visualization, but not sufficient for precise similarity or classification tasks.

✓ Part F: Bonus Challenge - Dependency Parsing

✓ Bonus Exercise: Visualize Sentence Structure

Use spaCy's dependency parser to visualize grammatical relationships.

```
1 from spacy import displacy
2
3 # YOUR CODE HERE
```

```
4 # 1. Choose an interesting sentence from Nike or contracts
5 # 2. Parse it with spaCy
6 # 3. Visualize the dependency tree
7 # 4. Identify the root verb, subjects, and objects
8
9 # Load spaCy model
10 nlp = spacy.load("en_core_web_sm")
11
12 # 1. Choose an interesting sentence
13 sample_sentence = "The Nike Air Zoom Pegasus 39 delivers respons
14
15
16 # 2. Parse it with spaCy
17 doc = nlp(sample_sentence)
18
19 # Display dependency visualization
20 # YOUR CODE HERE
21 # 3. Display dependency tree in notebook
22 displacy.render(doc, style="dep", jupyter=True, options={'distan
23
24 # Print dependency information
25 print("\nDependency Analysis:")
26 print("=" * 60)
27 for token in doc:
```