

## TF-IDF

### Step 1 — Load a small 20 Newsgroups subset and inspect it


We will load three concise categories from scikit-learn's 20 Newsgroups dataset for fast experiments.

This cell loads the training split for those categories, prints the number of documents and categories, shows how many documents per category, and prints short snippets from the first few documents so you can see the raw text we'll work with.

```
1 from sklearn.datasets import fetch_20newsgroups
2 import pandas as pd
3
4 # choose categories
5 categories = ['sci.space', 'rec.sport.baseball', 'comp.graphics']
6
7 # load dataset
8 newsgroups = fetch_20newsgroups(
9     subset='train',
10    categories=categories,
11    remove=('headers', 'footers', 'quotes'))
12 )
13
14 texts = newsgroups.data
15 targets = newsgroups.target
16 target_names = newsgroups.target_names
17
18 # build a dataframe properly
19 df = pd.DataFrame({
20     "text": texts,
21     "target_id": targets,
22     "target_name": [target_names[i] for i in targets]
23 })
24
25 print("Loaded categories:", target_names)
26 print("Number of documents:", len(texts))
27
28 # count documents per category
29 counts = df["target_name"].value_counts()
30 print("\nDocuments per category:")
31 print(counts)
32
33 df.head()
34
```

Loaded categories: ['comp.graphics', 'rec.sport.baseball', 'sci.space']  
Number of documents: 1774

Documents per category:  
target\_name  
rec.sport.baseball 597  
sci.space 593  
comp.graphics 584  
Name: count, dtype: int64

	text	target_id	target_name	
0	\nBy '8 grey level images' you mean 8 items of...	0	comp.graphics	
1	\nThere was a very useful article in one of th...	0	comp.graphics	
2	Bunker & McNally were later.\n\nPappas, Estrad...	1	rec.sport.baseball	
3	\nCalled "gas".\n\n\nThe balloons were in suff...	2	sci.space	
4	\nBut you still need the pitching staff to hol...	1	rec.sport.baseball	

Next steps: [Generate code with df](#) [New interactive sheet](#)

### Step 2 — Build TF-IDF, inspect vocabulary, and examine one document's TF-IDF weights

What this step does

We create a `TfidfVectorizer` and fit it to our dataset.

This converts every document into a TF-IDF vector based on word importance.

Why these parameters are chosen

- `lowercase=True`  
Ensures consistent matching of words by converting everything to lowercase.
- `stop_words='english'`  
Removes extremely common English words (like “the”, “and”) that carry little meaning and would dominate the vectors.
- `max_df=0.9`  
Ignores words that appear in more than 90% of documents.  
Such words are too common to be useful for distinguishing documents.
- `min_df=2`  
Keeps only words that appear in at least 2 documents.  
Extremely rare words often add noise and inflate vocabulary size.

## What we output

1. The total vocabulary size after filtering.
2. The TF-IDF weights of one document.
3. The **top 10 highest-weight tokens** for that document, showing what the model thinks is most important.

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 import numpy as np
3
4 vectorizer = TfidfVectorizer(lowercase=True, stop_words= "english",
5                             max_df= 0.9, min_df = 2, )
6 #print(vectorizer.get_stop_words())
7
8 tfidf_matrix = vectorizer.fit_transform(texts)
9 vocab = vectorizer.get_feature_names_out()
10 #print(vocab)
11
12 print("Vocabulary size:", len(vocab))
13 #print(tfidf_matrix)
```

Vocabulary size: 10416

```
1 doc_index = 5 # lets inspect for the sixth document in our dataset
2 doc_vector = tfidf_matrix[doc_index]
3 nonzero_indices = doc_vector.indices
4 nonzero_data = doc_vector.data
5
6
7 doc_df = pd.DataFrame({
8     "token": vocab[nonzero_indices],
9     "tfidf": nonzero_data
10 })
11
12 print(f"{doc_index}'th document with topic {df.target_name[doc_index]} has the following top 5 frequent words:")
13 doc_df.head()
14
```

5'th document with topic rec.sport.baseball has the following top 5 frequent words:

	token	tfidf
0	win	0.382816
1	games	0.214376
2	essentially	0.193526
3	year	0.167878
4	pitching	0.162579

Next steps: [Generate code with doc\\_df](#) [New interactive sheet](#)

## ✓ Step 3 — Compute cosine similarity between a query and all documents

### Goal of this step

We want to see how TF-IDF is actually *used* to find relevant documents.

We will take a short text query, transform it into a TF-IDF vector using the **same vectorizer**, and compute cosine similarity with every document.

This shows how IR systems rank documents based on vector similarity.

### Why cosine similarity?

- TF-IDF vectors vary in length (different document sizes).

- Cosine similarity compares only the **angle** between vectors, not their magnitude.
- This makes similarity robust and widely used in IR systems.

What this cell does

1. Define a **sample query** such as "space shuttle mission".
2. Convert it to a TF-IDF vector.
3. Compute cosine similarity between the query and all documents.
4. Rank the top 5 most relevant documents.
5. Print their category + a snippet so you can see *why* they matched.

After this step, we will visually inspect how well TF-IDF retrieves meaningful documents.

## What is the query vector?

The **query vector** is the TF-IDF representation of the user's search query.

We take the raw text query (e.g., "space shuttle mission") and apply the **same trained TF-IDF vectorizer** that was fitted on the document corpus.

This converts the query into a numerical vector with the same dimensions as the document TF-IDF matrix.

Each dimension represents a word in the vocabulary, and its value reflects how important that word is in the query relative to the whole dataset.

This allows us to directly compare the query to every document using cosine similarity.

```
1 query = ["Space shuttle mission"]
2 query_vector = vectorizer.transform(query)
3 indices = query_vector.nonzero()[1]
4 values = query_vector.data # actual tf-idf values
5
6 for idx, val in zip(indices, values):
7     print(f"Token '{vocab[idx]}' has TF-IDF weight {val}")
8 #print(query_vector.toarray()[0][indices])
9
```

```
Token 'mission' has TF-IDF weight 0.6652845186983247
Token 'shuttle' has TF-IDF weight 0.6100211130597171
Token 'space' has TF-IDF weight 0.43043088969278487
```

```
1 from sklearn.metrics.pairwise import cosine_similarity
2
3 similarities = cosine_similarity(query_vector, tfidf_matrix).flatten()
4 top_indices = np.argsort(similarities)[-1:5] # sorting descending top 5
5
6 print("Query:", query[0])
7 print("\nTop 5 most similar documents:\n")
8
9 for rank, idx in enumerate(top_indices, start=1):
10     snippet = texts[idx].strip().replace("\n", " ")
11     if len(snippet) > 300:
12         snippet = snippet[:300] + "..."
13     print(f"Rank {rank} | Document index {idx} | Similarity: {similarities[idx]:.4f} | Category: {target_names[target_idx]}")
14     print(snippet)
15     print("-"*1000)
16
```

Query: Space shuttle mission

Top 5 most similar documents:

```
Rank 1 | Document index 1760 | Similarity: 0.3882 | Category: sci.space
Ed Campion Headquarters, Washington, D.C. April 23, 1993 (Phone: 202/358-1780) Kyle Herring Johnson Space Center,
-----
Rank 2 | Document index 870 | Similarity: 0.3592 | Category: sci.space
Better idea for use of NASA Shuttle Astronauts and Crew is have them be found lost in space after a accident with a worm hole or other space/tim
-----
Rank 3 | Document index 403 | Similarity: 0.3335 | Category: sci.space
Archive-name: space/controversy Last-modified: $Date: 93/04/01 14:39:06 $ CONTROVERSIAL QUESTIONS These issues periodically come up with m
-----
Rank 4 | Document index 974 | Similarity: 0.3332 | Category: sci.space
Archive-name: space/schedule Last-modified: $Date: 93/04/01 14:39:23 $ SPACE SHUTTLE ANSWERS, LAUNCH SCHEDULES, TV COVERAGE SHUTTLE LAUNCH
-----
Rank 5 | Document index 428 | Similarity: 0.3013 | Category: sci.space
For an essay, I am writing about the space shuttle and a need for a better propulsion system. Through research, I have found that it is rather
-----
```

## Step 4 — Visualize TF-IDF document vectors using PCA

## Goal of this step

We want to see how TF-IDF separates documents from different categories.

Because TF-IDF vectors are high-dimensional (thousands of features), we reduce them to **2 dimensions** using PCA so we can plot them.

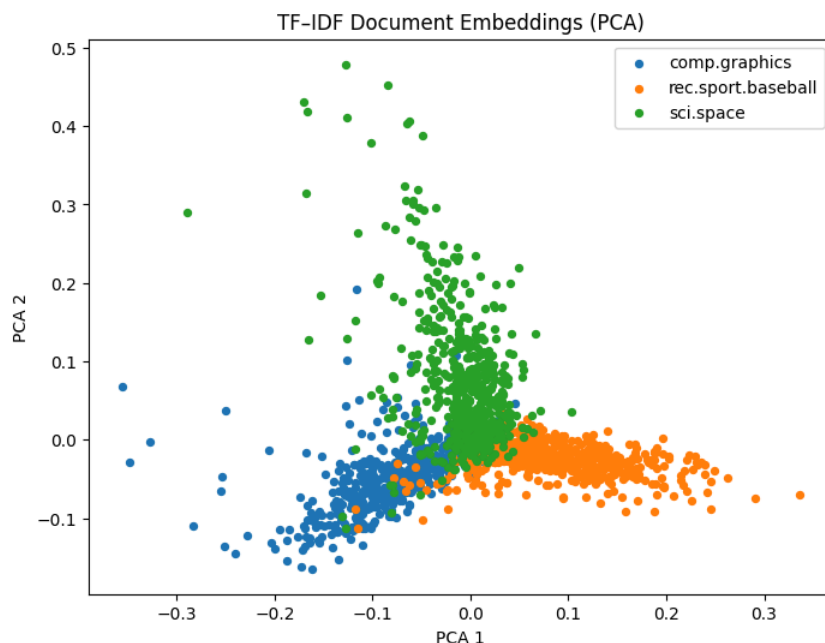
## Why PCA?

- It is a linear dimensionality-reduction method.
- Fast and suitable for sparse TF-IDF matrices.
- Helps us visually inspect whether categories form clusters.

## What this cell does

1. Applies PCA to reduce all TF-IDF vectors to 2D.
2. Plots them as a scatter plot.
3. Colors each point by its category (space, baseball, graphics).
4. Helps you visually understand how TF-IDF groups similar documents.

```
1 from sklearn.decomposition import PCA
2 import matplotlib.pyplot as plt
3 pca = PCA(n_components=2)
4 reduced = pca.fit_transform(tfidf_matrix.toarray())
5
6
7 plt.figure(figsize=(8,6))
8
9 for i, name in enumerate(target_names):
10     points = reduced[np.array(targets) == i]
11     plt.scatter(points[:,0], points[:,1], label=name, s=18)
12
13 plt.xlabel("PCA 1")
14 plt.ylabel("PCA 2")
15 plt.title("TF-IDF Document Embeddings (PCA)")
16 plt.legend()
17 plt.show()
18
```



## Step 5 — Compare `CountVectorizer` (raw counts) vs `TfidfVectorizer` (TF-IDF) for the same query

### What this step does

We build a `CountVectorizer` with the same preprocessing choices as the TF-IDF vectorizer and compare how each method ranks documents for the same query. This highlights the practical effect of IDF weighting: TF-IDF downweights very common terms across the corpus and upweights distinctive terms, which can change ranking order.

### Why compare these two

- `CountVectorizer` represents documents by raw token counts; it treats all terms equally regardless of global rarity.
- `TfidfVectorizer` accounts for a term's rarity across documents (the IDF part) so rare-but-informative words get more weight.
- Observing differences for the same query shows when IDF helps (e.g., when common words would otherwise dominate) and when raw counts might still be useful.

### Parameters chosen and why

- `lowercase=True` and `stop_words='english'` to match preprocessing used earlier and ensure a fair comparison.
- `max_df=0.9` and `min_df=2` to ignore extremely common and extremely rare terms, keeping the vocabulary manageable and comparable to the TF-IDF run.

## Inverted Index

### ✓ Inverted Index — Step 1: Prepare corpus and build a simple inverted index

What this step does

1. Load the same 20 Newsgroups subset (three categories) so experiments stay consistent.
2. Tokenize each document into normalized word tokens and remove English stop words.
3. Build an inverted index mapping `(term -> postings_list)` where each posting is `(doc_id, [positions])`.
4. Compute document frequencies and document-length statistics.
5. Print dataset stats, the top-10 most document-frequent terms, and example postings for inspection.

Why these choices

- Tokenization with `\w+` is simple and interpretable for a learning notebook.
- Lowercasing ensures tokens match regardless of capitalization.
- Removing English stop words keeps the index focused on informative terms.
- Storing positions enables phrase and proximity queries later.
- Document IDs are the Python default (starting at 0) unless you request a different convention.

What you should check in the output

- Number of documents and vocabulary size to confirm the corpus loaded correctly.
- Document length stats to understand corpus variability (useful later for BM25).
- Top terms by document frequency to verify that common tokens were retained/filtered as expected.
- Postings format `(doc_id, [positions])` to ensure positions were recorded correctly.

```

1 from sklearn.datasets import fetch_20newsgroups
2 from sklearn.feature_extraction.text import ENGLISH_STOP_WORDS
3 import re
4 from collections import defaultdict, Counter
5 import pandas as pd
6
7 categories = ['sci.space', 'rec.sport.baseball', 'comp.graphics']
8 raw = fetch_20newsgroups(subset='train', categories=categories, remove=('headers', 'footers', 'quotes'))
9 docs = raw.data
10
11 token_pattern = re.compile(r"\w+")
12 stop_words = set(ENGLISH_STOP_WORDS)
13
14 inverted_index = defaultdict(list)
15 doc_freq = Counter()
16 doc_lengths = {}
17
18 for doc_id, text in enumerate(docs):
19     tokens = [t.lower() for t in token_pattern.findall(text)]
20     tokens = [t for t in tokens if t not in stop_words]
21     doc_lengths[doc_id] = len(tokens)
22     positions_by_term = defaultdict(list)
23     for pos, tok in enumerate(tokens):
24         positions_by_term[tok].append(pos)
25     for term, positions in positions_by_term.items():
26         inverted_index[term].append((doc_id, positions))
27     for term in positions_by_term.keys():
28         doc_freq[term] += 1
29
30 vocab_size = len(inverted_index)
31 num_docs = len(docs)
32
33 print("Number of documents:", num_docs)
34 print("Vocabulary size (after stop-word removal):", vocab_size)
35
36 lengths = list(doc_lengths.values())

```

```

37 print("\nDocument length stats:")
38 print(pd.Series(lengths).describe()[["min", "50%", "mean", "max"]])
39
40 most_common_terms = doc_freq.most_common(10)
41 print("\nTop 10 terms by document frequency:")
42 for term, df in most_common_terms:
43     print(f"{term} - df={df}")
44
45 examples = [most_common_terms[0][0], most_common_terms[3][0]]
46 print("\nExample postings lists (first 10 postings shown for each term):")
47 for term in examples:
48     postings = inverted_index[term][:10]
49     print(f"\nTerm: '{term}' | df={doc_freq[term]}")
50     print(postings)
51

```

Number of documents: 1774

Vocabulary size (after stop-word removal): 22383

Document length stats:

```

min      0.000000
50%     35.000000
mean     97.421646
max     6990.000000
dtype: float64

```

Top 10 terms by document frequency:

```

s - df=670
t - df=620
like - df=379
just - df=338
know - df=320
don - df=311
think - df=299
1 - df=293
m - df=293
2 - df=280

```

Example postings lists (first 10 postings shown for each term):

Term: 's' | df=670

```

[[0, [42]], (6, [1, 9]), (13, [7, 53, 73]), (17, [49, 83]), (18, [0]), (20, [4, 13]), (23, [17]), (25, [14]), (27, [64, 105, 118, 140]), (28, [5

```

Term: 'just' | df=338

```

[[27, [89]], (34, [3159]), (44, [160]), (47, [20]), (48, [104]), (49, [0]), (51, [4]), (52, [2, 19]), (54, [41]), (64, [13])]

```

## ✓ Quick checks and examples

Below you will see:

- The 10 most document-frequent tokens (terms that appear in the most documents) and their document frequencies.
- Example postings lists for two informative tokens to inspect the stored `(doc_id, positions)` structure. This helps confirm the index is built correctly and shows how postings store both which documents contain a term and where that term occurs inside the document.

```

1 most_common_terms = doc_freq.most_common(10)
2 print("Top 10 terms by document frequency:\n")
3 for term, df in most_common_terms:
4     print(f"{term} - df={df}")
5
6 examples = [most_common_terms[0][0], most_common_terms[3][0]]
7 print("\n\nExample postings lists:\n")
8 for term in examples:
9     postings = inverted_index[term][:10]
10    print(f"Term: '{term}' | df={doc_freq[term]} | first 10 postings (doc_id, positions):")
11    print(postings)
12    print()
13

```

Top 10 terms by document frequency:

```

s - df=670
t - df=620
like - df=379
just - df=338
know - df=320
don - df=311
think - df=299
1 - df=293
m - df=293
2 - df=280

```

Example postings lists:

Term: 's' | df=670 | first 10 postings (doc\_id, positions):

```
[(1, [42]), (7, [1, 9]), (14, [7, 53, 73]), (18, [49, 83]), (19, [0]), (21, [4, 13]), (24, [17]), (26, [14]), (28, [64, 105, 118, 140]), (29, [5
Term: 'just' | df=338 | first 10 postings (doc_id, positions):
[(28, [89]), (35, [3159]), (45, [160]), (48, [20]), (49, [104]), (50, [0]), (52, [4]), (53, [2, 19]), (55, [41]), (65, [13])]
```

## How to read and interpret the output of Step 1

When you run Step 1, you'll see several kinds of information.  
Here is how to understand each part of the printed output:

### 1. "Number of documents: X"

This tells you how many documents from the selected 20 Newsgroups categories are loaded.  
These are the documents we build the inverted index over.

### 2. "Vocabulary size: Y"

This is the total number of **unique terms** that survived:

- tokenization
- lowercasing
- stop-word removal
- min\_df=1 (every remaining word appeared at least once somewhere)

This number is usually in the thousands: that's your **index vocabulary**.

### 3. Document length statistics

You will see:

- **min**: shortest document (after removing stop words)
- **50% (median)**: typical document length
- **mean**: average length
- **max**: longest document

These lengths matter later for ranking functions (BM25, etc.).

### 4. Top 10 terms by document frequency

Output looks like:

term — df=42 graphics — df=35 game — df=30 ..

- **df (document frequency)** = number of documents that contain this term.
- High-df terms are common across many documents.
- These are often topic-general words (like "game", "team", "graphics").

This helps verify the index is correct and gives a sense of which concepts dominate the dataset.

### 5. Example postings lists

You'll see something like:

Term: 'game' | df=30 | first 10 postings (doc\_id, positions): [(3, [5, 17]), (7, [8]), (12, [1, 9, 16]), ...]

This is the heart of the inverted index.

Each entry `(doc_id, positions)` means:

- **doc\_id** — which document contains the term
- **positions** — the word positions where that term appears in that document  
(position 0 = first token after preprocessing)

Example interpretation:

`(12, [1, 9, 16])` means

- Term appears **in document 12**
- At token positions **1, 9, and 16** within that document

This is precisely the structure search engines use for:

- Boolean queries
- Phrase search ("space shuttle mission")
- Proximity queries
- Ranking and scoring later

## evaluation

### ✓ Evaluation Step 1 — Load train/test datasets and build the train index & TF-IDF

What this step does

1. Load the 20 Newsgroups dataset using **train split** for building all retrieval models.
2. Load the **test split** which will supply our evaluation queries.
3. From the train set:
  - Tokenize documents
  - Remove English stop words
  - Build an inverted index mapping:  
`term → list of (doc_id_train, [positions])`
  - Save document lengths and document frequencies
4. Build the **TF-IDF vectorizer** and matrix on the train documents only.

Why this is the correct approach

- Retrieval systems are trained/indexed on a fixed corpus (train).
- Queries come from unseen data (test).
- Relevance is defined by category labels:  
documents in the train corpus with the same category as the test query document.

Outputs of this step

- `train_docs`, `test_docs`
- `inverted_index`
- `tfidf_matrix` and `vectorizer`
- Train document frequencies, document lengths
- Ready for query generation and evaluation in the next step.

```

1 from sklearn.datasets import fetch_20newsgroups
2 from sklearn.feature_extraction.text import ENGLISH_STOP_WORDS
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 import re
5 from collections import defaultdict, Counter
6 import pandas as pd
7 import numpy as np
8
9 categories = ['sci.space', 'rec.sport.baseball', 'comp.graphics']
10 train = fetch_20newsgroups(subset='train', categories=categories, remove=('headers', 'footers', 'quotes'))
11 test = fetch_20newsgroups(subset='test', categories=categories, remove=('headers', 'footers', 'quotes'))
12
13 train_docs = train.data
14 train_targets = train.target
15 test_docs = test.data
16 test_targets = test.target
17 target_names = train.target_names
18
19 token_pattern = re.compile(r"\w+")
20 stop_words = set(ENGLISH_STOP_WORDS)
21
22 inverted_index = defaultdict(list)
23 train_doc_freq = Counter()
24 train_doc_lengths = {}
25 num_train_docs = len(train_docs)
26
27 for doc_id, text in enumerate(train_docs):
28     tokens = [t.lower() for t in token_pattern.findall(text)]
29     tokens = [t for t in tokens if t not in stop_words]
30     train_doc_lengths[doc_id] = len(tokens)
31     positions_by_term = defaultdict(list)
32     for pos, tok in enumerate(tokens):
33         positions_by_term[tok].append(pos)
34     for term, positions in positions_by_term.items():
35         inverted_index[term].append((doc_id, positions))
36     for term in positions_by_term.keys():
37         train_doc_freq[term] += 1
38
39 vectorizer = TfidfVectorizer(lowercase=True, stop_words='english', max_df=0.9, min_df=2)
40 tfidf_matrix = vectorizer.fit_transform(train_docs)
41

```



```

42 print("Train docs:", len(train_docs), "Test docs:", len(test_docs))
43 print("Train vocabulary size:", len(vectorizer.get_feature_names_out()))
44 print("Train document length stats:")
45 print(pd.Series(list(train_doc_lengths.values())).describe()[["min", "50%", "mean", "max"]])
46

```

```

Train docs: 1774 Test docs: 1180
Train vocabulary size: 10416
Train document length stats:
min      0.000000
50%     35.000000
mean     97.421646
max     6990.000000
dtype: float64

```

## Step 2: Create test queries from the test split

What this cell does

1. For each selected test document we produce a short query formed by the top-k TF-IDF tokens where TF-IDF is computed by applying the train-fit vectorizer to the test document.
2. We sample a reproducible subset of test documents as queries (fixed random seed).
3. Each query record contains: `qid` (test doc index), `query_text`, `query_terms`, and `relevant_label` (test doc category).

Why this approach

- Using top TF-IDF tokens from test documents yields concise, realistic queries that a retriever should answer.
- Using a random but seeded subset keeps evaluation quick and reproducible.

Outputs

- `queries` list with structured query info for evaluation.

```

1 import random
2 np.random.seed(42)
3 random.seed(42)
4
5 num_queries = 100
6 num_test_docs = len(test_docs)
7 if num_queries > num_test_docs:
8     num_queries = num_test_docs
9
10 query_indices = random.sample(list(range(num_test_docs)), num_queries)
11 top_k_terms = 5
12
13 test_tfidf = vectorizer.transform(test_docs)
14
15 def top_k_terms_from_vector(row_vec, k=top_k_terms):
16     arr = row_vec.toarray().flatten()
17     nz = np.where(arr > 0)[0]
18     if len(nz) == 0:
19         return []
20     top_idx = nz[np.argsort(arr[nz])][::-1][:k]
21     features = vectorizer.get_feature_names_out()
22     return [features[i] for i in top_idx]
23
24 queries = []
25 for qi in query_indices:
26     row = test_tfidf[qi]
27     terms = top_k_terms_from_vector(row, top_k_terms)
28     query_text = " ".join(terms) if terms else " ".join(token_pattern.findall(test_docs[qi])[:top_k_terms])
29     queries.append({
30         "qid": qi,
31         "query_text": query_text,
32         "query_terms": terms,
33         "relevant_label": test_targets[qi]
34     })
35
36 print("Prepared queries:", len(queries))
37 print("Example query:", queries[0])
38

```

```

Prepared queries: 100
Example query: {'qid': 228, 'query_text': 'library windows exhibit viewers connection', 'query_terms': ['library', 'windows', 'exhibit', 'viewer

```

## Step 3: Helper functions for ranking and textbook metrics

What this cell does

1. Implements retrievers:

- TF-IDF cosine ranking using the train TF-IDF matrix and train vectorizer.
- Term-overlap baseline that counts how many query terms occur in each train doc using the inverted index.

## 2. Implements metrics:

- precision@r and recall@r arrays along a full ranked list,
- Average Precision (AP) per textbook (mean of precision at ranks with relevant docs),
- 11-point interpolated precision (recall levels 0.0,0.1,...,1.0) per query.

## Why these implementations

- These functions directly map to the textbook evaluation procedures so results are comparable to chapter descriptions.

```

1 from sklearn.metrics.pairwise import cosine_similarity
2
3 def rank_tfidf_cosine(query_text):
4     qv = vectorizer.transform([query_text])
5     sims = cosine_similarity(qv, tfidf_matrix).flatten()
6     ranked_indices = np.argsort(sims)[::-1]
7     ranked_scores = sims[ranked_indices]
8     return ranked_indices, ranked_scores
9
10 def rank_term_overlap(query_terms):
11     counts = np.zeros(num_train_docs, dtype=int)
12     for t in query_terms:
13         postings = inverted_index.get(t, [])
14         for doc_id, _ in postings:
15             counts[doc_id] += 1
16     ranked_indices = np.argsort(counts)[::-1]
17     ranked_scores = counts[ranked_indices].astype(float)
18     return ranked_indices, ranked_scores
19
20 def precision_recall_at_ranks(ranked_indices, relevant_set):
21     relevant_flags = np.array([1 if doc in relevant_set else 0 for doc in ranked_indices], dtype=int)
22     cum_relevant = np.cumsum(relevant_flags)
23     ranks = np.arange(1, len(ranked_indices) + 1)
24     precision_at_r = cum_relevant / ranks
25     total_relevant = len(relevant_set)
26     if total_relevant == 0:
27         recall_at_r = np.zeros_like(precision_at_r)
28     else:
29         recall_at_r = cum_relevant / total_relevant
30     return precision_at_r, recall_at_r, relevant_flags
31
32 def average_precision(precision_at_r, relevant_flags):
33     relevant_positions = np.where(relevant_flags == 1)[0]
34     if relevant_positions.size == 0:
35         return 0.0
36     precisions_at_relevant = precision_at_r[relevant_positions]
37     ap = precisions_at_relevant.mean()
38     return ap
39
40 def eleven_point_interpolated_precision(precision_at_r, recall_at_r):
41     recall_levels = np.linspace(0.0, 1.0, 11)
42     interp = []
43     for r in recall_levels:
44         mask = recall_at_r >= r - 1e-12
45         if mask.any():
46             interp.append(np.max(precision_at_r[mask]))
47         else:
48             interp.append(0.0)
49     return np.array(interp), recall_levels
50

```

## Step 4: Run evaluation over all test queries (TF-IDF vs term-overlap)

### What this cell does

1. For each test query:
  - Finds the set of relevant train documents (same category label as test query).
  - Computes TF-IDF ranking and term-overlap ranking.
  - Computes precision@r, recall@r, AP, and 11-point interpolated precision for each method.
2. Aggregates:
  - Per-query AP arrays for MAP calculation.
  - Per-query 11-point arrays for averaged interpolated precision curves.

### Notes

- Query document itself is in the test split and is not part of the train corpus, so no need to exclude it from the relevant set.

```

1 tfidf_APs = []
2 overlap_APs = []
3 tfidf_interps = []
4 overlap_interps = []
5
6 for q in queries:
7     qid = q["qid"]
8     qtext = q["query_text"]
9     qterms = q["query_terms"]
10    rel_label = q["relevant_label"]
11    relevant_docs = set([i for i, lbl in enumerate(train_targets) if lbl == rel_label])
12    ranked_tfidf, _ = rank_tfidf_cosine(qtext)
13    p_tfidf, r_tfidf, flags_tfidf = precision_recall_at_ranks(ranked_tfidf, relevant_docs)
14    ap_tfidf = average_precision(p_tfidf, flags_tfidf)
15    interp_tfidf, recall_levels = eleven_point_interpolated_precision(p_tfidf, r_tfidf)
16    tfidf_APs.append(ap_tfidf)
17    tfidf_interps.append(interp_tfidf)
18    ranked_ov, _ = rank_term_overlap(qterms)
19    p_ov, r_ov, flags_ov = precision_recall_at_ranks(ranked_ov, relevant_docs)
20    ap_ov = average_precision(p_ov, flags_ov)
21    interp_ov, _ = eleven_point_interpolated_precision(p_ov, r_ov)
22    overlap_APs.append(ap_ov)
23    overlap_interps.append(interp_ov)
24
25 tfidf_MAP = float(np.mean(tfidf_APs))
26 overlap_MAP = float(np.mean(overlap_APs))
27 tfidf_mean_interp = np.mean(np.vstack(tfidf_interps), axis=0)
28 overlap_mean_interp = np.mean(np.vstack(overlap_interps), axis=0)
29
30 print("Evaluated", len(queries), "queries")
31 print("TF-IDF MAP:", round(tfidf_MAP, 4))
32 print("Term-overlap MAP:", round(overlap_MAP, 4))
33

```

Evaluated 100 queries  
 TF-IDF MAP: 0.3926  
 Term-overlap MAP: 0.3894

## ✓ Step 5: Plot averaged 11-point interpolated precision–recall curves

What this cell does

1. Plots the averaged 11-point interpolated precision values for TF-IDF and term-overlap across the standard recall levels.
2. Adds MAP values to the legend for quick numeric comparison.

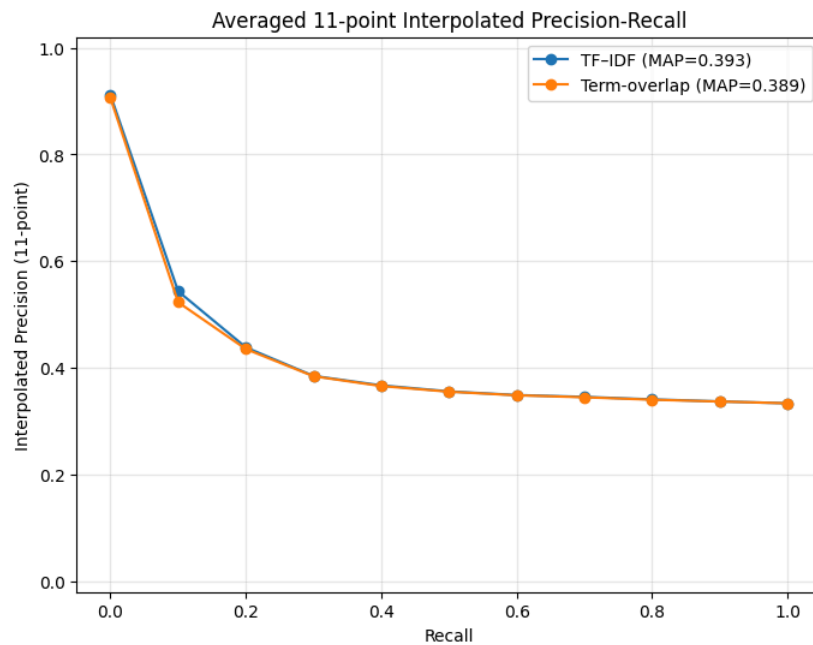
Why this plot

- This plot is the textbook's standard compact way to compare ranked retrieval systems.

```

1 import matplotlib.pyplot as plt
2
3 plt.figure(figsize=(8,6))
4 plt.plot(recall_levels, tfidf_mean_interp, marker='o', label=f"TF-IDF (MAP={tfidf_MAP:.3f})")
5 plt.plot(recall_levels, overlap_mean_interp, marker='o', label=f"Term-overlap (MAP={overlap_MAP:.3f})")
6 plt.xlabel("Recall")
7 plt.ylabel("Interpolated Precision (11-point)")
8 plt.title("Averaged 11-point Interpolated Precision-Recall")
9 plt.ylim(-0.02, 1.02)
10 plt.grid(alpha=0.3)
11 plt.legend()
12 plt.show()
13

```

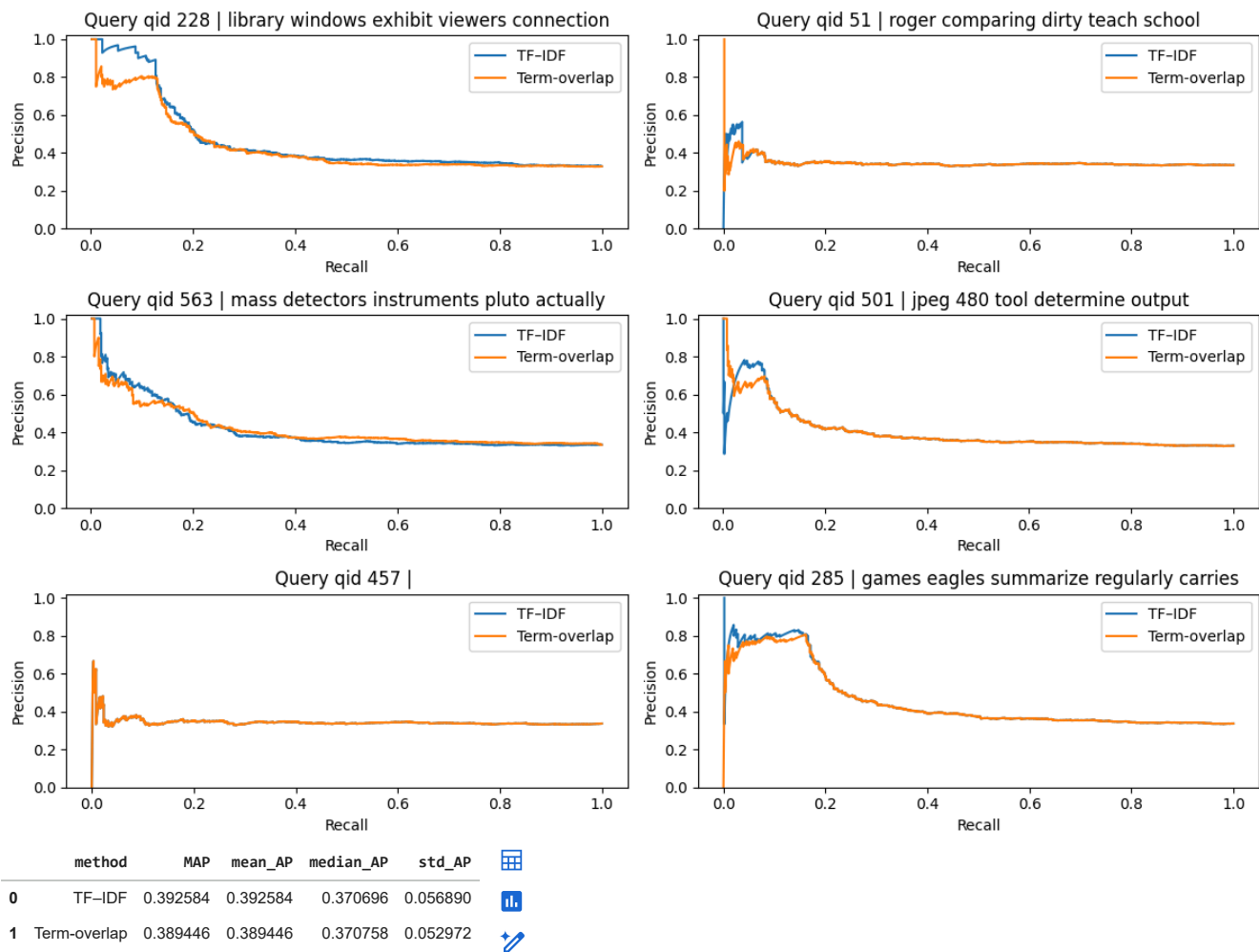


## Step 6 — Per-query precision–recall plots and summary table

What this cell does

1. Plots precision vs recall for the first six queries for quick inspection.
2. Displays a summary table with MAP, mean AP, median AP, and std AP for both methods.

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3
4 example_qs = queries[:6]
5
6 plt.figure(figsize=(12,8))
7 for i, q in enumerate(example_qs, start=1):
8     rel_label = q["relevant_label"]
9     relevant_docs = set([i for i, lbl in enumerate(train_targets) if lbl == rel_label])
10    ranked_tfidf, _ = rank_tfidf_cosine(q["query_text"])
11    p_tfidf, r_tfidf, _ = precision_recall_at_ranks(ranked_tfidf, relevant_docs)
12    ranked_ov, _ = rank_term_overlap(q["query_terms"])
13    p_ov, r_ov, _ = precision_recall_at_ranks(ranked_ov, relevant_docs)
14    plt.subplot(3,2,i)
15    plt.plot(r_tfidf, p_tfidf, label='TF-IDF')
16    plt.plot(r_ov, p_ov, label='Term-overlap')
17    plt.xlabel("Recall")
18    plt.ylabel("Precision")
19    title_txt = q["query_text"] if len(q["query_text"]) < 60 else q["query_text"][:57] + "..."
20    plt.title(f"Query qid {q['qid']} | {title_txt}")
21    plt.legend()
22    plt.ylim(0,1.02)
23 plt.tight_layout()
24 plt.show()
25
26 summary = pd.DataFrame({
27     "method": ["TF-IDF", "Term-overlap"],
28     "MAP": [tfidf_MAP, overlap_MAP],
29     "mean_AP": [float(np.mean(tfidf_APs)), float(np.mean(overlap_APs))],
30     "median_AP": [float(np.median(tfidf_APs)), float(np.median(overlap_APs))],
31     "std_AP": [float(np.std(tfidf_APs)), float(np.std(overlap_APs))]
32 })
33 display(summary)
34
```



Next steps: [Generate code with summary](#) [New interactive sheet](#)

## Step 8 — Precision vs Recall curve for a single query (non-interpolated, with markers and vertical stems)

What this cell does

1. Selects one example test query (the first in `queries`).
2. Computes the TF-IDF ranked list for that query and the precision and recall at every rank using the textbook definitions.
3. Plots the **non-interpolated** precision-recall curve with:
  - connected markers (precision points at each rank),
  - light vertical stems from the x-axis up to each precision point (to match the style you provided),
  - axis limits and grid similar to the example figure.

How to read the plot

- Each marker corresponds to a rank in the ranked list; x = recall after that rank, y = precision after that rank.
- Vertical stems visually emphasize the jumps in precision at the ranks where relevant documents are found.
- This is the exact per-rank precision/recall behavior the textbook shows (Figure 11.3) and matches the non-interpolated example you provided.

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 query_text = "space shuttle"
5 query_label = target_names.index("sci.space")
6 relevant_docs = set([i for i, lbl in enumerate(train_targets) if lbl == query_label])
7
8 ranked, = rank_tfidf_cosine(query_text)
```

```
9 precision_at_r, recall_at_r, _ = precision_recall_at_ranks(ranked, relevant_docs)
10
11 xs = recall_at_r
12 ys = precision_at_r
13
14 plt.figure(figsize=(8,4.2))
15 plt.plot(xs, ys, linestyle='-', marker='o', markersize=6, color='tab:blue')
16 for x, y in zip(xs, ys):
17     plt.vlines(x, 0, y, colors='tab:blue', alpha=0.25, linewidth=1)
18 plt.xlabel("Recall")
19 plt.ylabel("Precision")
20 plt.title('Precision vs Recall (non-interpolated) — query: "space shuttle"')
21 plt.xlim(0.0, 1.0)
22 plt.ylim(0.0, 1.02)
23 plt.grid(alpha=0.3)
24 plt.show()
25
```

