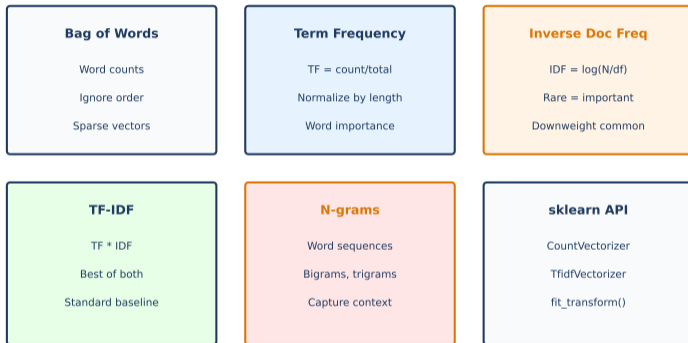


Lesson 38 Summary: BOW & TF-IDF

Data Science with Python – Key Concepts

Data Science Program

Bag of Words & TF-IDF



Convert text to numerical vectors for ML

Count-based representation:

- **Vocabulary:** All unique words
- **Vector:** Word counts per document
- **Sparse:** Most entries are zero

Ignores word order, captures frequency

Weight by importance:

- **TF:** Term frequency in document
- **IDF:** $\log\left(\frac{N}{df}\right)$ penalizes common words
- **TF-IDF:** $TF \times IDF$

Rare words get higher weights

Capture word sequences:

- **Unigrams:** Single words
- **Bigrams:** Word pairs (“not good”)
- **Trigrams:** Word triples

`ngram_range=(1,2)` includes uni and bigrams

Two main classes:

- **CountVectorizer:** Raw counts (BOW)
- **TfidfVectorizer:** TF-IDF weights
- **API:** `fit_transform(docs)`

Same interface, different weighting

Essential Commands:

Task	Code
BOW	<code>CountVectorizer().fit_transform(docs)</code>
TF-IDF	<code>TfidfVectorizer().fit_transform(docs)</code>
Bigrams	<code>TfidfVectorizer(ngram_range=(1,2))</code>
Max features	<code>TfidfVectorizer(max_features=1000)</code>
Get names	<code>vec.get_feature_names_out()</code>

TF-IDF is the standard baseline for text classification