

Table: Advantages and Disadvantages of One-Hot Encoding, Bag-of-Words (BoW), N-Grams, and TF-IDF

Method	Advantages	Disadvantages
One-Hot Encoding	- Simple and easy to understand.	- High-dimensional representation, leading to the curse of dimensionality.
	- Preserves the categorical nature of data.	- Does not capture semantic relationships between categories.
	- Useful for algorithms that require numerical input.	- Inefficient for large categorical variables.
	- Does not require scaling of features.	- Cannot handle out-of-vocabulary terms.
	- Can handle missing values effectively by representing them as all-zero vectors.	- Increases computational complexity for large datasets.
Bag-of-Words (BoW)	- Simple and straightforward representation of text data.	- Ignore word order and grammar, leading to loss of context.
	- Efficient for large datasets.	- Treats all words as equally important, ignoring word semantics.
	- Can be easily interpreted and visualized.	- Requires extensive preprocessing for better results.
	- Captures basic word frequency information.	- Doesn't handle typos or misspellings well.

	<ul style="list-style-type: none"> <li>- Can be combined with various algorithms for text classification and clustering.</li> </ul>	<ul style="list-style-type: none"> <li>- Sparse representation can be memory-intensive for large vocabularies.</li> </ul>
N-Grams	<ul style="list-style-type: none"> <li>- Captures local word order information.</li> </ul>	<ul style="list-style-type: none"> <li>- Increases feature dimensionality, especially with higher n-values.</li> </ul>
	<ul style="list-style-type: none"> <li>- Preserves some context around words.</li> </ul>	<ul style="list-style-type: none"> <li>- Limited context windows may fail to capture long-range dependencies.</li> </ul>
	<ul style="list-style-type: none"> <li>- Useful for capturing short phrases and idiomatic expressions.</li> </ul>	<ul style="list-style-type: none"> <li>- Sensitive to noise and variability in language usage.</li> </ul>
	<ul style="list-style-type: none"> <li>- Can be used to model syntactic and semantic relationships.</li> </ul>	<ul style="list-style-type: none"> <li>- Requires careful selection of n-value.</li> </ul>
	<ul style="list-style-type: none"> <li>- Provides a compromise between Bag-of-Words and full sequence models.</li> </ul>	<ul style="list-style-type: none"> <li>- Memory and computation-intensive for large n-values and datasets.</li> </ul>
TF-IDF	<ul style="list-style-type: none"> <li>- Highlight important terms while downweighting common terms.</li> </ul>	<ul style="list-style-type: none"> <li>- Does not capture semantic relationships between words.</li> </ul>
	<ul style="list-style-type: none"> <li>- Effective in reducing the impact of noise words (stop-words).</li> </ul>	<ul style="list-style-type: none"> <li>- Requires a representative corpus for accurate IDF calculation.</li> </ul>
	<ul style="list-style-type: none"> <li>- Handles out-of-vocabulary terms gracefully.</li> </ul>	<ul style="list-style-type: none"> <li>- Sensitive to term variations and synonyms.</li> </ul>
	<ul style="list-style-type: none"> <li>- Captures the importance of a term within a document and across a corpus.</li> </ul>	<ul style="list-style-type: none"> <li>- May not perform well with short documents or small corpora.</li> </ul>

	<ul style="list-style-type: none"><li>- Widely used in information retrieval, text mining, and document classification tasks.</li></ul>	<ul style="list-style-type: none"><li>- Computational overhead in calculating TF-IDF scores for large datasets.</li></ul>
--	---	---