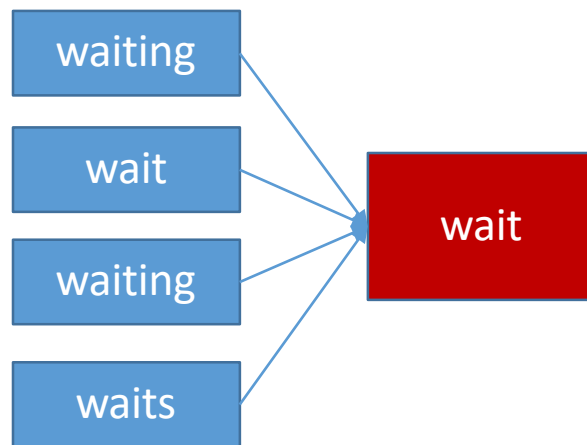



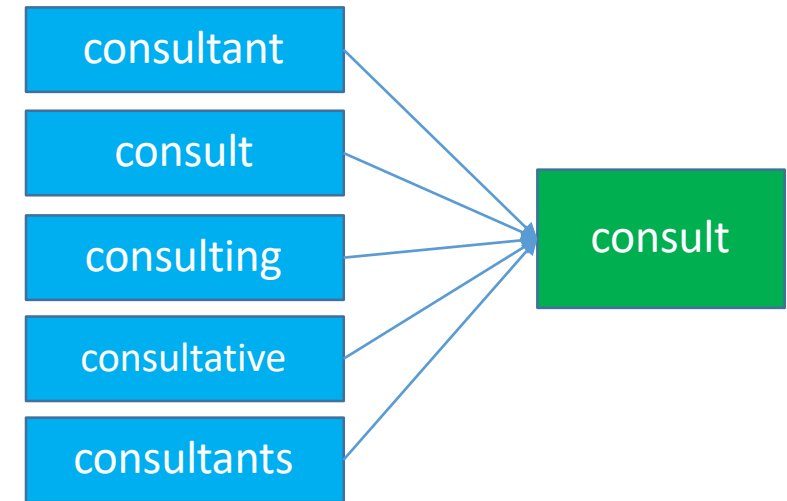
1. Word tokens
2. Character tokens
3. Sentence tokens
4. Named entity tokens
5. Part-of-speech (POS) tags
6. Sub-word tokens

Stemming and Lemmatization are Text Normalization (or sometimes called Word Normalization) techniques in the field of Natural Language Processing that are used to prepare text, words, and documents for further processing. Stemming and Lemmatization have been studied, and algorithms have been developed in Computer Science since the 1960s. We will learn about Stemming and Lemmatization in a practical approach covering the background, some famous algorithms, applications of Stemming and Lemmatization, and how to stem and lemmatize words, sentences, and documents using the Python **nltk package** which is the Natural Language Tool Kit package provided by Python for Natural Language Processing tasks.

Stemming is the process of reducing a word to its word stem that affixes to suffixes and prefixes or to the roots of words known as a lemma. Stemming is important in natural language understanding (NLU) and natural language processing (NLP). When a new word is found, it can present new research opportunities. For example -



change
changing
changes
changed
changer

A diagram illustrating the stemming process for the word 'change'. On the left, five lines of text are listed: 'change', 'changing', 'changes', 'changed', and 'changer'. Lines from the words 'changing', 'changes', 'changed', and 'changer' converge on a single point to the right, representing the root 'change'.

- **Porter Stemmer():** The Porter stemming algorithm (or 'Porter stemmer') removes the commoner morphological endings from words in English.
- **Lovins Stemmer**
- **Dawson Stemmer**
- **Krovetz Stemmer**
- **Xerox Stemmer**
- **N-Gram Stemmer**
- **Snowball Stemmer**
- **Lancaster Stemmer**

Lemmatization usually refers to doing things properly with the use of a vocabulary and morphological analysis of words, normally aiming to remove inflectional endings only and to return the base or dictionary form of a word, which is known as the lemma.

- Word Net Lemmatizer
- Spacy Lemmatizer
- TextBlob
- Gensim Lemmatizer
- TreeTagger

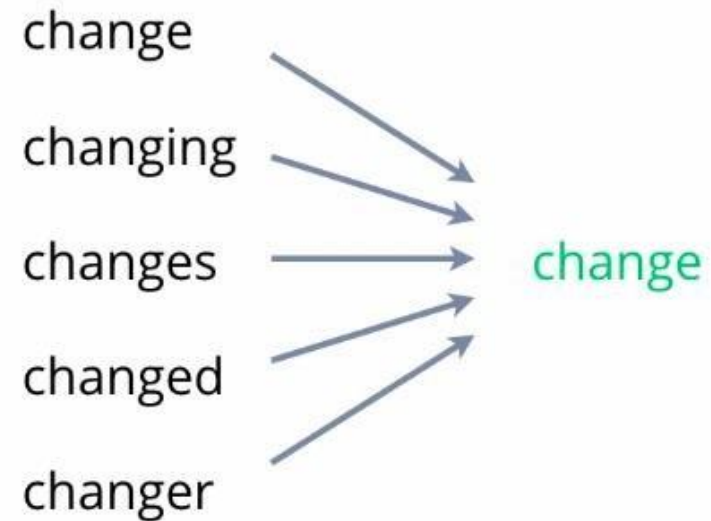


| | original_word | lemmatized_word |
|---|---------------|-----------------|
| 0 | trouble | trouble |
| 1 | troubling | trouble |
| 2 | troubled | trouble |
| 3 | troubles | trouble |

Stemming



Lemmatization



Vectorizer!

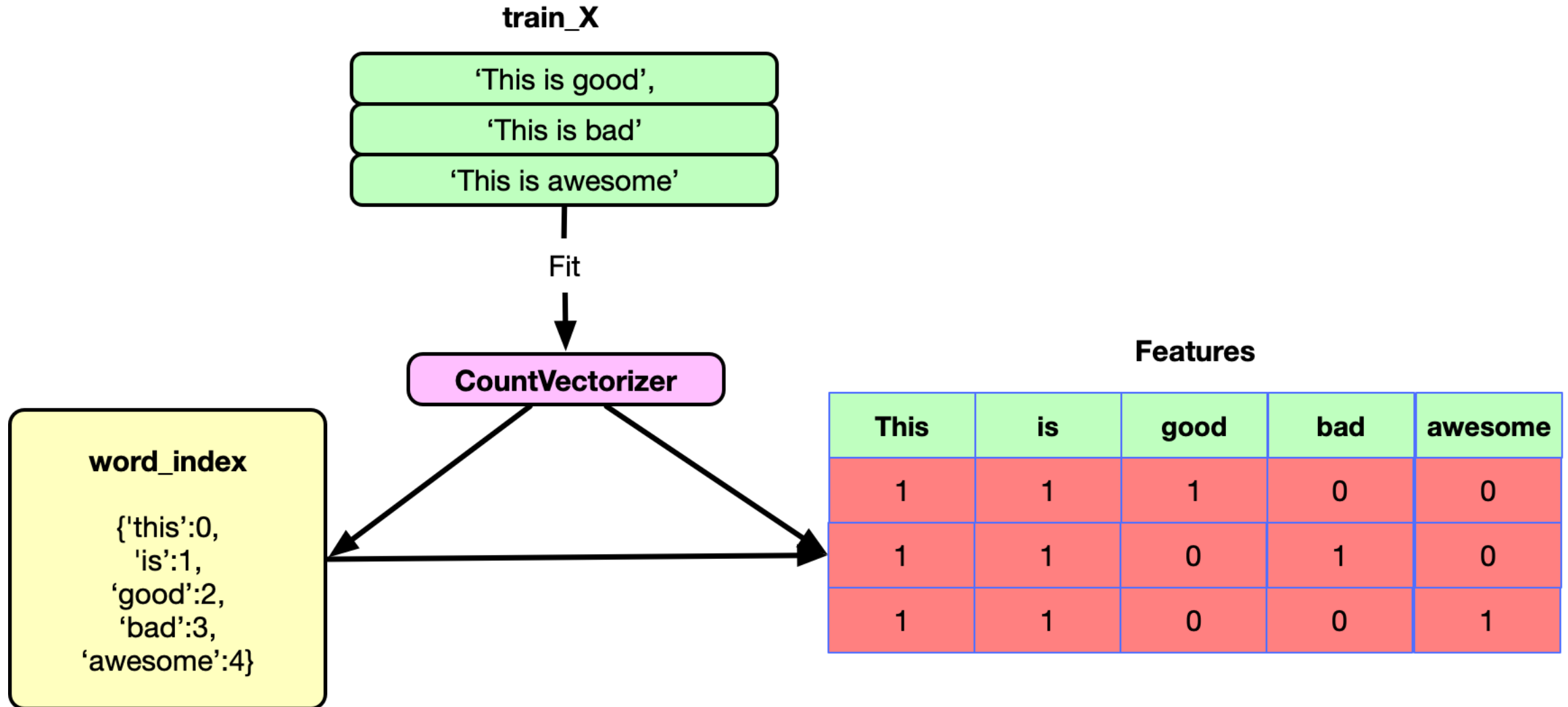
www.aiquest.org

Feature transformation: Transformation of data to improve the accuracy of the algorithm.

Feature selection: Removing unnecessary features.

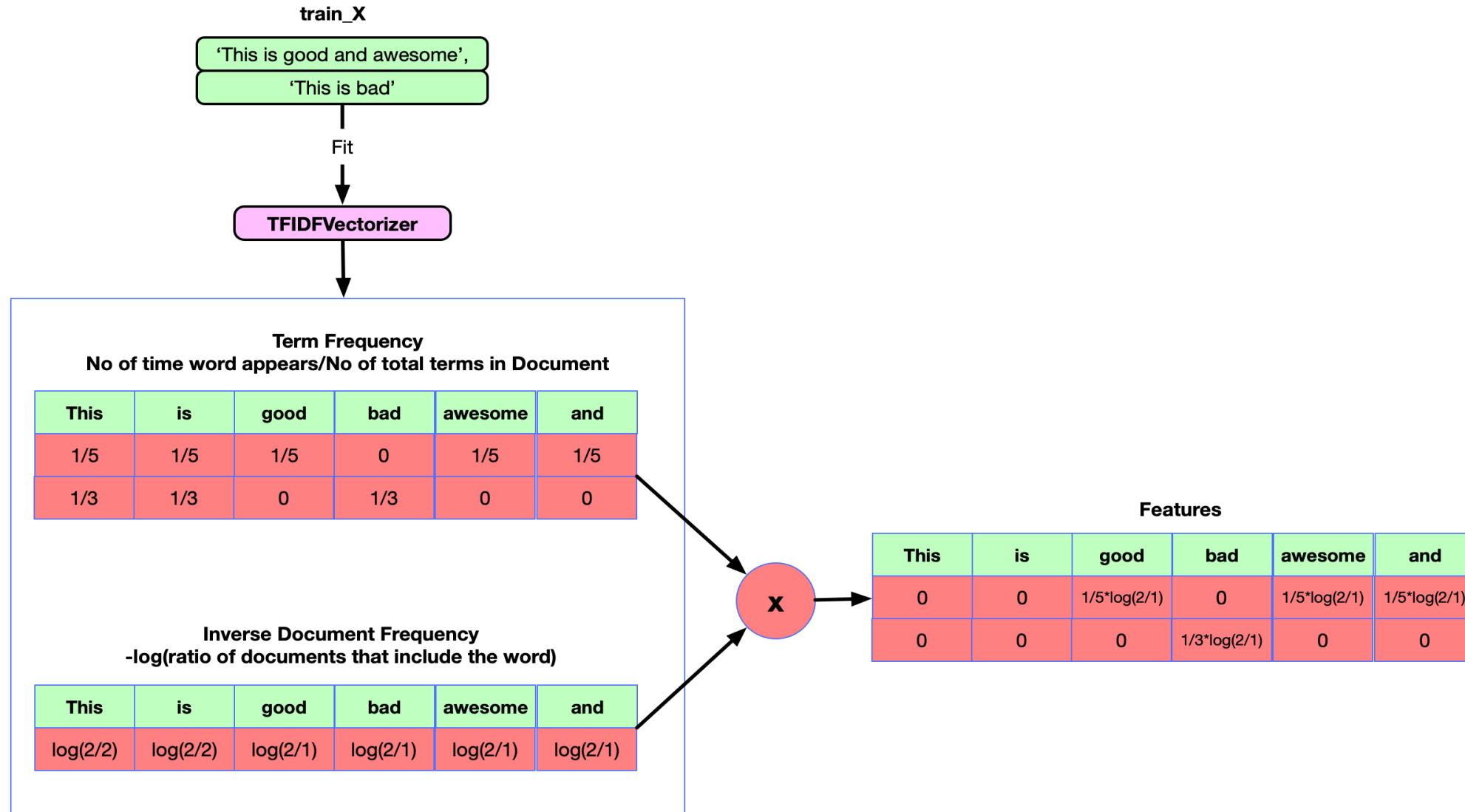
Feature extraction: The transformation of raw data into features suitable for modeling.

- Bag of Words: Count Vectorizer
- TF-IDF Vectorizer
- Word2Vec
- Global Vectors for Word Representation



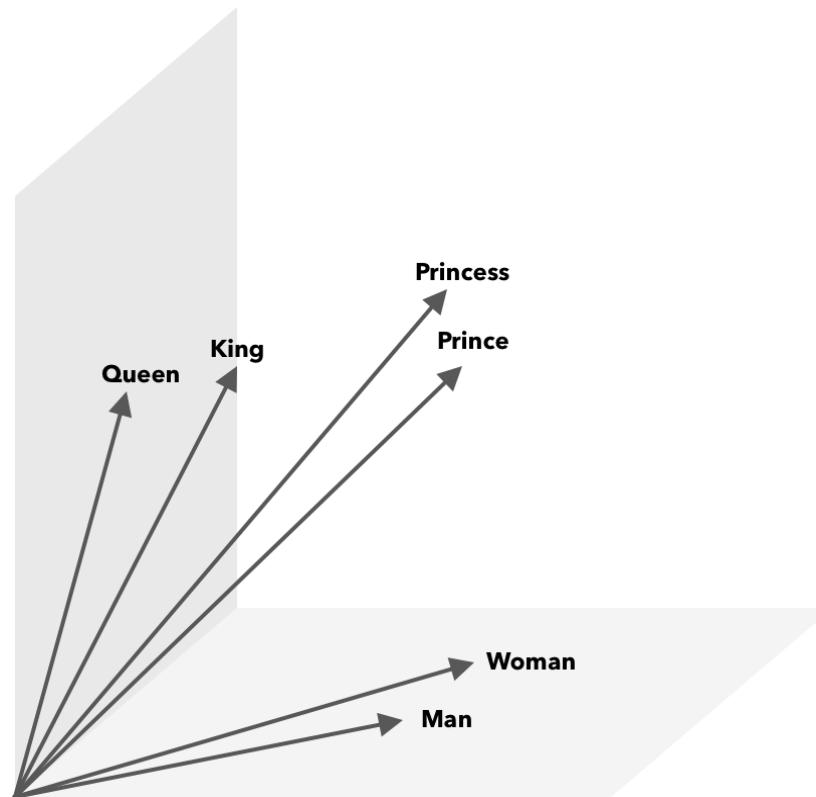
TF-IDF Vectorizer

Mathematical architecture



- TF-IDF is better than Count Vectorizers because it not only focuses on the frequency of words present in the corpus but also provides the importance of the words. The term "df" is called document frequency which means in how many documents the word “subfield” is present within the corpus.
 - ☐ Can TF IDF Be Negative?
 - ☐ No. The lowest value is 0. Both term frequency and inverse document frequency are positive numbers.
- In AI inference and machine learning, sparsity refers to a matrix of numbers that includes many zeros or values that
- will not significantly impact a calculation.

Word2vec is a Neural Network that processes text by “vectorizing” words. Its input is a text corpus and its output is a set of vectors: feature vectors that represent words in that corpus.



- $\text{King} - \text{Man} + \text{Women} = \text{Queen}$
- $\text{Prince} + \text{mom} = \text{Queen}$
- $\text{vec}(\text{"king"}) - \text{vec}(\text{"man"}) + \text{vec}(\text{"woman"}) \approx \text{vec}(\text{"queen"})$