



**CHRIST**  
(DEEMED TO BE UNIVERSITY)  
BANGALORE · INDIA

# TEXT ANALYTICS

## Unit-2

### MISSION

CHRIST is a nurturing ground for an individual's holistic development to make effective contribution to the society in a dynamic environment

### VISION

Excellence and Service

### CORE VALUES

Faith in God | Moral Uprightness  
Love of Fellow Beings  
Social Responsibility | Pursuit of Excellence

## Unit-2

- Text Representation- tokenization, stemming, stop words, TF-IDF, NER, N-gram modeling. Mining Textual Data: Text Clustering, Text Classification,
- **LabExercises:**
  3. Implementation of tokenization, stemming, stop words
  4. Implementation of text classification and clustering with TF-IDF and N-gram.

# Text analysis

- Text mining or natural language processing (NLP), is a field of study focuses on extracting meaningful insights and patterns from unstructured text data
- With the proliferation of digital content, including social media posts, customer reviews, news articles, and scientific literature, and etc.



# Key Concepts in Text Analysis

- **Text Preprocessing**

- tokenization (splitting text into words or sentences),
- removing stopwords (common words like "the", "and", "is"),
- stemming or lemmatization (reducing words to their root form)
- handling special characters and punctuation.

- **Text Representation:** Text data converted into a numerical format for analysis.

- Bag of Words (BoW)
- Term Frequency-Inverse Document Frequency (TF-IDF): A numerical statistic that reflects the importance of a word in a document relative to a corpus.
- Word Embedding: Dense vector representations of words in a high-dimensional space, capturing semantic relationships between words.

- **Text Analysis Techniques:**

- Sentiment Analysis: Determining the sentiment or opinion expressed in text (e.g., positive, negative, neutral).
- Topic Modeling: Discovering latent topics or themes present in a collection of documents.
- Named Entity Recognition (NER): Identifying and classifying named entities such as persons, organizations, locations, and dates mentioned in text.
- Text Classification: Categorizing text documents into predefined classes or categories (e.g., spam detection, topic classification).
- Information Extraction: Extracting structured information from unstructured text (e.g., extracting product names and prices from customer reviews).

- **Tools and Libraries:** Several libraries and tools are available for text analysis in various programming languages. Some popular ones include:

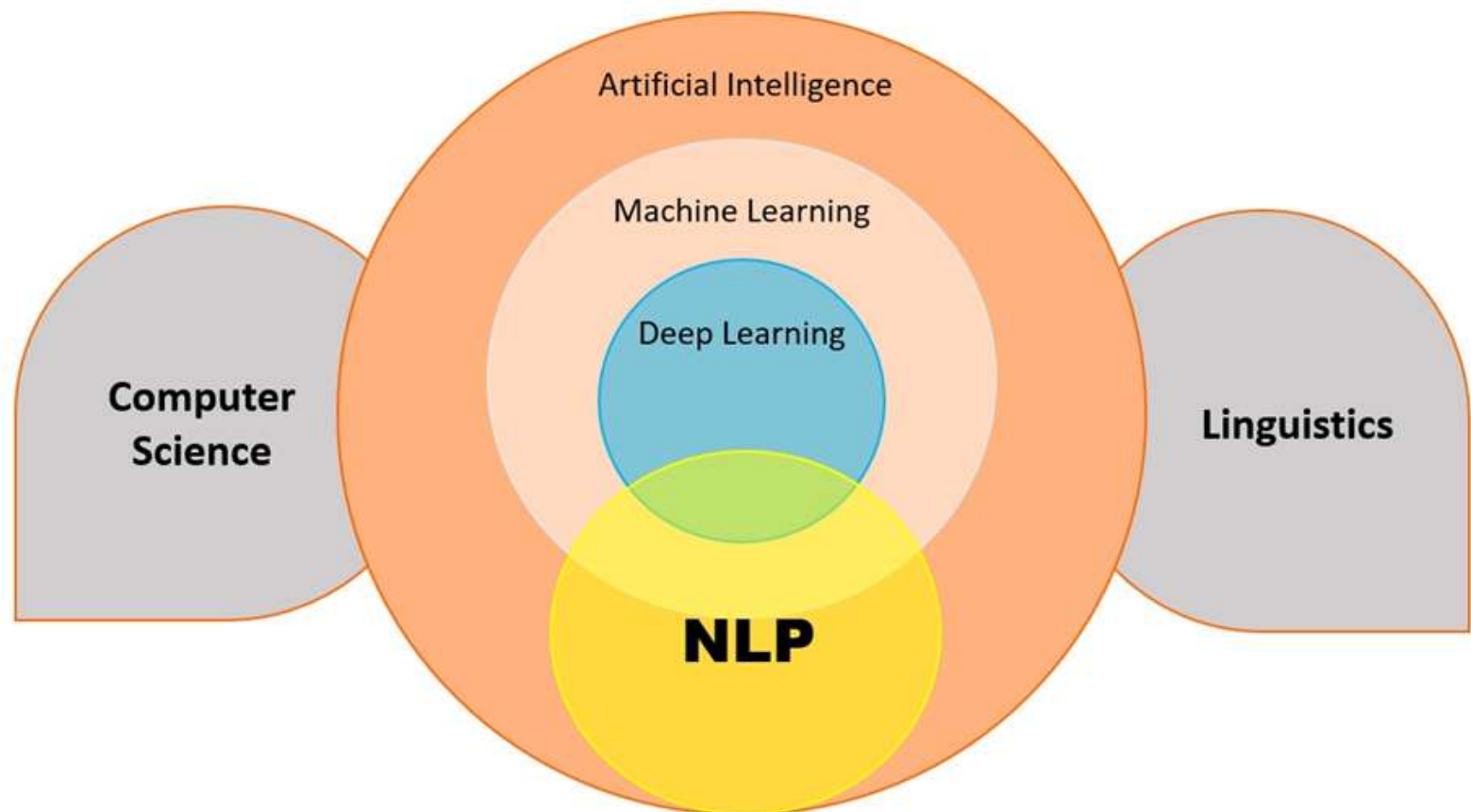
- Python: NLTK (Natural Language Toolkit), spaCy, scikit-learn, gensim.
- R: tm (Text Mining), quanteda, tidytext.
- Java: Apache OpenNLP, Stanford NLP.
- Commercial platforms: IBM Watson Natural Language Understanding, Google Cloud Natural Language API, Microsoft Azure Text Analytics.

## Applications of Text Analysis:

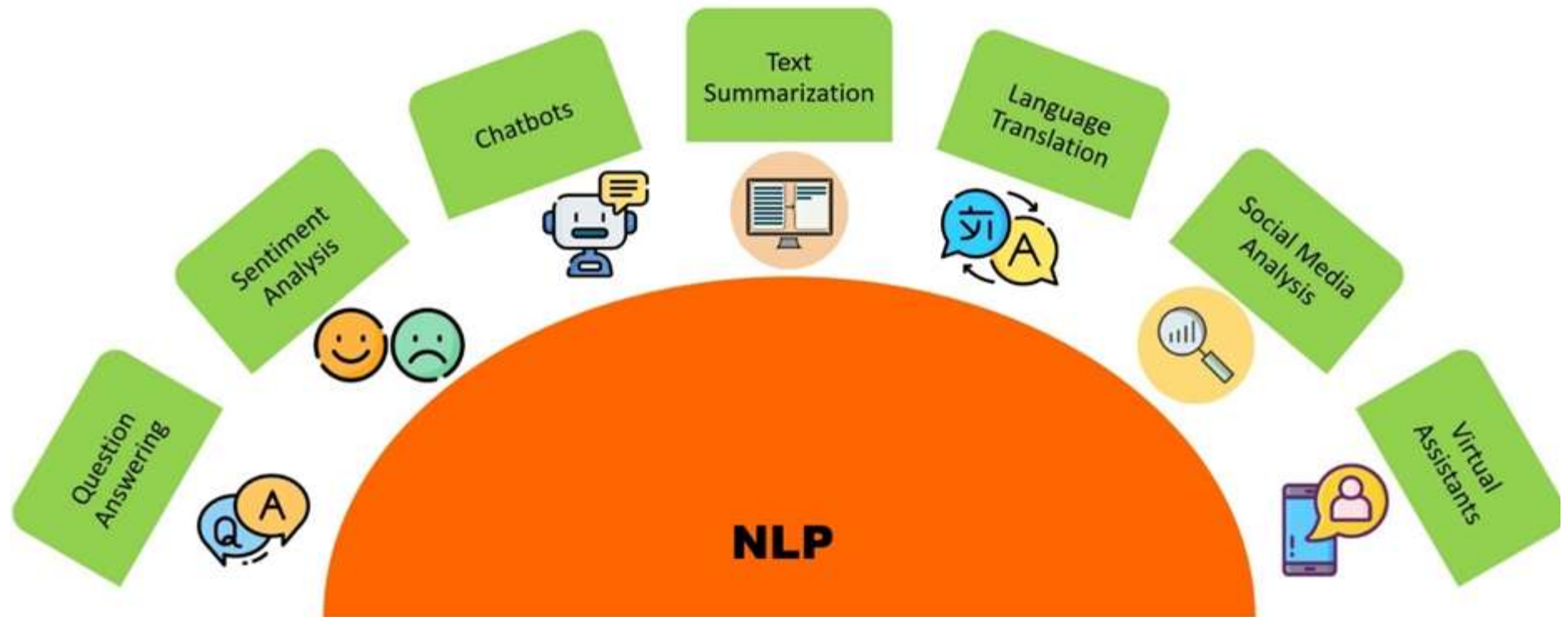
- **Business Intelligence:** Analyzing customer feedback, social media posts, and product reviews to gain insights into customer preferences, sentiment, and market trends.
- **Information Retrieval:** Building search engines that can understand user queries and retrieve relevant documents from a large corpus of text.
- **Healthcare:** Analyzing electronic health records, clinical notes, and biomedical literature for disease diagnosis, drug discovery, and personalized medicine.
- **Social Media Monitoring:** Tracking public opinion, trends, and events by analyzing social media conversations and posts.
- **Legal and Regulatory Compliance:** Analyzing legal documents, contracts, and regulatory texts to extract relevant information and ensure compliance.

# NLP-Natural Language Processing

- NLP first appeared as machine translation in the 1950s as decoding messages during World War II for Russian into English
- Prior 1980s, the main driving force behind NLP was a convoluted system of manual set of instructions



# Most widely used real-word applications of NLP





# Most popular Python libraries for Text Processing



**NLTK**

Natural Language Toolkit

Key technologies:

- Tokenization
- Text Classification
- Stemming
- Tagging
- Parsing
- Semantic Reasoning



**spaCy**

Support custom models in  
PyTorch, TensorFlow

Key technologies:

- Tokenization
- Text Classification
- Lemmatization
- Sentence segmentation, etc



**Gensim**

Topic Modeling for humans

Represents documents as  
semantic vectors

Key technologies:

- Word2Vec
- Latent Semantic Indexing (LSI Model)
- Latent Dirichlet Allocation (LDA Model)



**TensorFlow**

Training & inference of  
neural networks

Key features:

- Pre-trained models and datasets
- Tools to process data
- Deployment options
- Implementing options



**PyTorch**

Key features:

- Production ready
- Distributed training
- Deployment options
- Robust ecosystem
- Native onnx support
- Cloud support

# Named Entity Recognition (NER)

- The ISIS has claimed responsibility for a suicide bomb blast in the Tunisian capital earlier this week, the militant group's Amaq news agency said on Thursday. A militant wearing an explosive belt blew himself up in Tunis.

ORGANISATION

LOCATION

DATE

PERSON

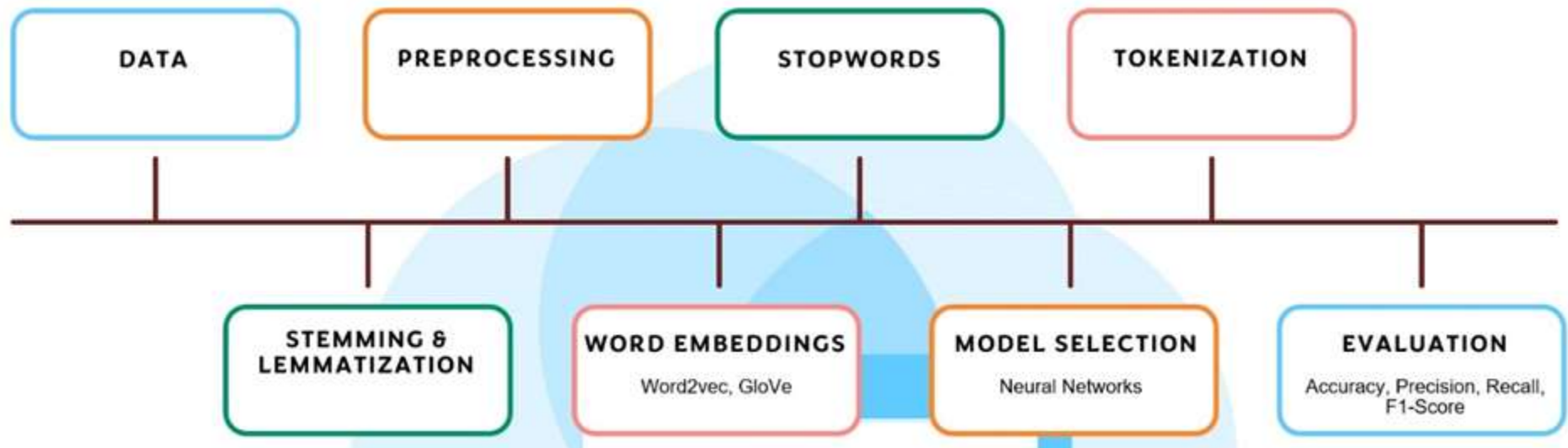
WEAPON

The **ISIS** ORG has claimed responsibility for a suicide bomb blast in the **Tunisian** LOC capital **earlier this week** DATE, the **militant group** ORG's **Amaq news agency** ORG said on **Thursday** DATE. A **militant** PER wearing an **explosives belt** WEAPON blew himself up in **Tunis** LOC.

## NLP Projects Ideas - Others

- Text-to-Speech (TTS) and Speech-to-Text (STT)
- Search Autocorrect and Autocomplete
- Language Translator
- Hiring and Recruitment
- Targeted Advertising
- Survey Analysis
- Social Media Monitoring
- Emotion Detection
- Inspiring Quote Generator
- Email Filtering

# Building blocks of an NLP application



# Text Preprocessing

- **Tokenization:** Splitting text into individual words, phrases, symbols, or other meaningful elements called tokens. This can be done at the word level or sentence level.
- **Lowercasing:** Converting all characters in the text to lowercase to ensure uniformity, as "Word" and "word" should be treated the same.
- **Removing Punctuation:** Eliminating punctuation marks such as periods, commas, and exclamation points to focus on the actual words.
- **Removing Stop Words:** Removing common words that do not add significant meaning to the text, such as "and," "the," "is," etc., to reduce noise in the data.
- **Stemming:** Reducing words to their base or root form. For example, "running," "runner," and "ran" are all stemmed to "run."
- **Lemmatization:** Similar to stemming, but more sophisticated, lemmatization reduces words to their base or dictionary form, considering the context. For example, "better" is lemmatized to "good."

## Text Preprocessing(Cont..)

- **Removing Special Characters:** Eliminating special characters, symbols, or numbers that are not relevant to the analysis.
- **Removing Whitespace:** Trimming leading, trailing, and excessive whitespace within the text to maintain consistency.
- **Normalizing Accents:** Converting accented characters to their unaccented counterparts, e.g., "café" to "cafe."
- **Text Correction:** Correcting spelling and grammatical errors to improve the quality of the text.
- **Removing HTML Tags:** Stripping HTML tags from web-scraped text data to retain only the meaningful content.
- **Removing or Replacing URLs and Email Addresses:** Eliminating or substituting URLs and email addresses with placeholders to avoid irrelevant data.
- **Handling Emojis and Emoticons:** Converting emojis and emoticons to text descriptions if they carry significant meaning for the analysis

## Text Preprocessing(Cont..)

- **Expanding Contractions:** Converting contractions into their full forms, e.g., "don't" to "do not." **Text Segmentation:** Splitting text into meaningful segments or chunks, such as sentences or paragraphs, for more granular analysis
- **Text Segmentation:** Splitting text into meaningful segments or chunks, such as sentences or paragraphs, for more granular analysis.
- **Noise Removal:** Filtering out irrelevant data such as boilerplate text, headers, footers, and advertisements from the text.
- **Feature Extraction:** Converting text into numerical features for machine learning models. Common methods include Bag of Words (BoW), Term Frequency-Inverse Document Frequency (TF-IDF), and word embeddings (e.g., Word2Vec, GloVe, BERT). **Handling Emojis and Emoticons:** Converting emojis and emoticons to text descriptions if they carry significant meaning for the analysis

# Preprocessing Example-1

- Proper nouns and names (e.g., "John," "Jane," "New York")
- Punctuation (e.g., commas, exclamation points)
- Emojis (e.g., 😊, 🏠)
- URLs (e.g., "www.example.com")
- Dates (e.g., "10/12/2023")
- Email addresses (e.g., "john.doe@example.com")
- Contractions (e.g., "you're," "don't")
- Common stop words (e.g., "the," "it's")
- Mixed case text (e.g., "I hope you're doing well!")
- Special characters (e.g., "😊," "🏠")
- Whitespace (leading, trailing, and within the text)
- Numbers (e.g., "10/12/2023," "1")
- Accented characters (e.g., "café's")
- Multiple sentences and phrases that can be tokenized
- Text suitable for stemming and lemmatization (e.g., "visiting," "meeting")
- A mention of a date for segmentation
- Examples of contractions for expansion (e.g., "you're," "I'll")



# Text Data Processing - Tokenization

Text = NLP has made a significant breakthrough



NLP has made a significant breakthrough

## 1. Tokenization

NLP has made a significant breakthrough



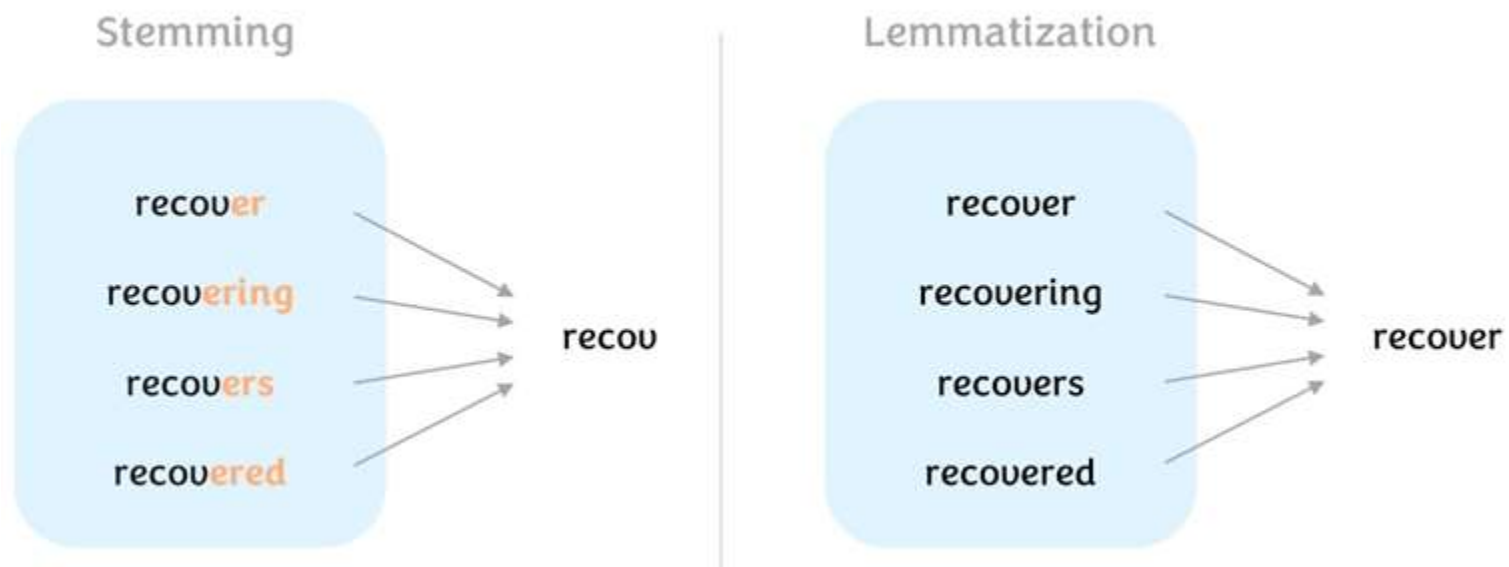
After removing stop words

NLP, significant, breakthrough

## 2. Stop Words Removal

## Text Data Processing Cont..

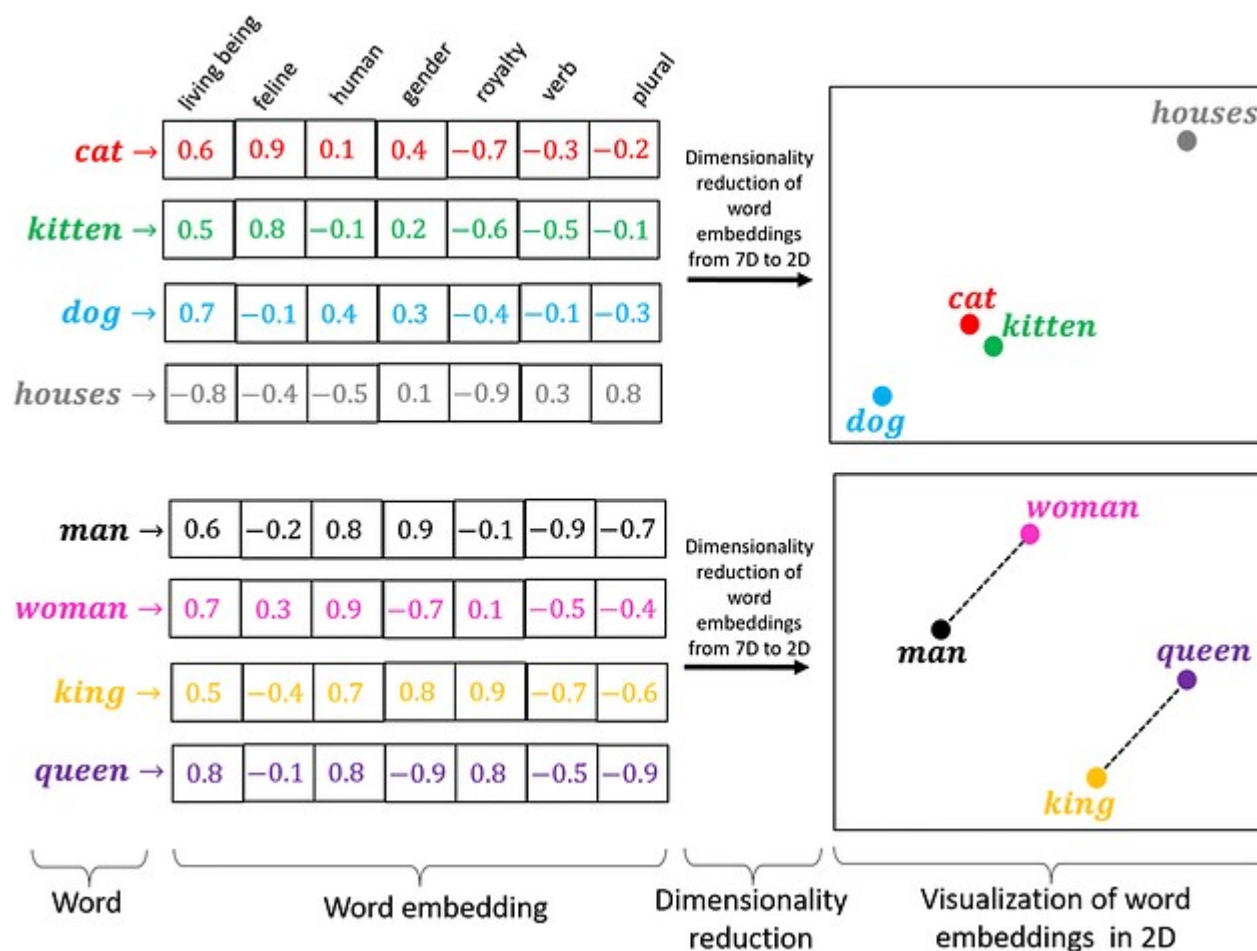
### 3. Stemming & Lemmatization



# Feature Extraction

- Converting text into numerical features for machine learning models.  
Common methods include
  - Bag of Words (BoW)
  - Term Frequency-Inverse Document Frequency (TF-IDF)
  - word embeddings (e.g., Word2Vec, GloVe, BERT).

# Word Embedding



# Vectorization or Word Embedding in NLP

- Word embedding methodologies
  - TF-IDF (Term Frequency Inverse Document Frequency)
  - BOW (Bag-of-Words), Count Vectorization, N-grams Vectorization, Word2Vec, GloVe.
- Word2Vec: The process of 'word embeddings' involves turning each word into a numerical representation of the word (a vector).
  - Each word is converted into a single vector, which is then trained in a manner resembling a neural network.
  - Based on a word's usage in the text, Word2Vec can infer a word's meaning with a high degree of accuracy.
- GloVe : represents a global vector, which is an unsupervised learning technique that generates word vector representations.
  - The advantage of GloVe is that it integrates global statistics to build word vectors, while Word2Vec depends on local statistics (local context knowledge about words) for generation of word vectors.
  - It is based on word-context matrix factorization algorithms. In GloVe, a co-occurrence matrix is used to determine the semantic relationship between words.

- Sentiment Analysis: Perform sentiment analysis on social media data.
- Text Summarization: Implement extractive and abstractive summarization techniques.
- Information Retrieval: Develop a basic search engine using TF-IDF or embeddings.
- Text Generation: Build a simple text generation model using transformers.
- Advanced NLP Application: Develop a question-answering system or a named entity linking application.

## Supervised Learning Approaches & Heuristic Methods

Supervised learning methods involve training a model to identify important sentences based on labeled training data:

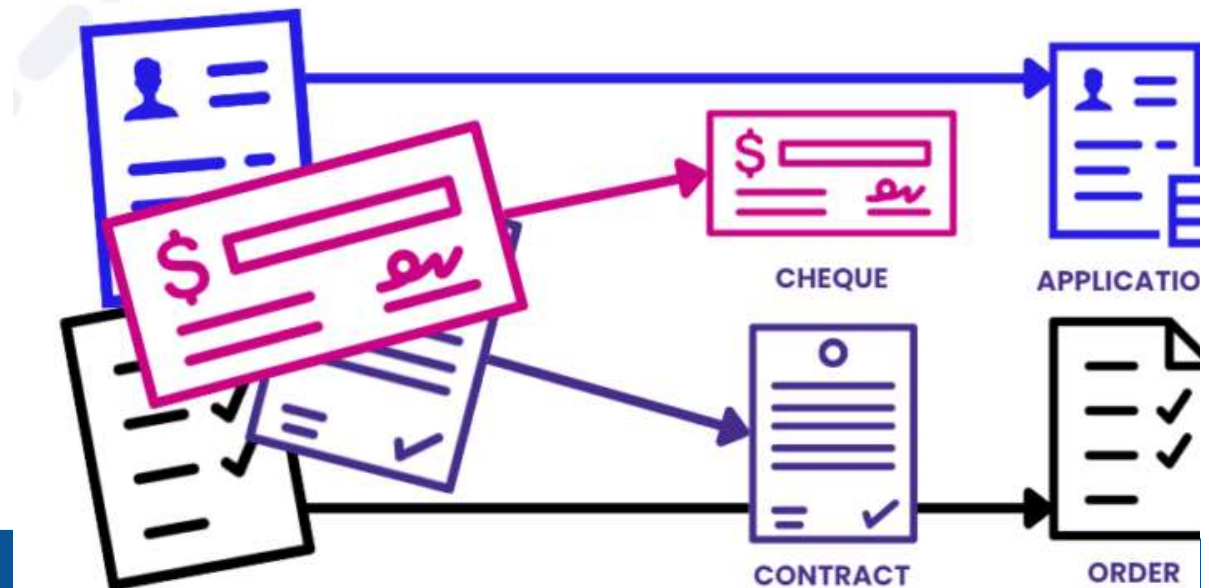
- **Feature Extraction:** Extract features such as sentence position, length, presence of keywords, etc.
- **Model Training:** Train a classifier (e.g., logistic regression, SVM, neural networks) to predict sentence importance.
- **Sentence Scoring:** Score sentences using the trained model.
- **Sentence Selection:** Select the top-scoring sentences for the summary.

Simple heuristic methods rely on predefined rules to select important sentences:

- **Lead-Based Summarization:** Select the first few sentences of the document as the summary (effective for news articles).
- **Title and Heading Matching:** Select sentences that contain keywords from the title or headings.

# Document Classification & Clustering

- **Document classification (document categorization)** refers to recognizing a document category based on its content, visual appearance, and other factors.
- **You can classify documents into folders based on labels you create, for example:**
  - **Level of confidentiality:** public, confidential, top secret.
  - **Type of document:** invoice, corrected invoice, receipt.
  - **Language:** English, French, Spanish.





## Document Processing by Industries

### ACCOUNTS PAYABLE



- invoices
- corrected invoices
- receipts
- purchase orders

### HEALTHCARE

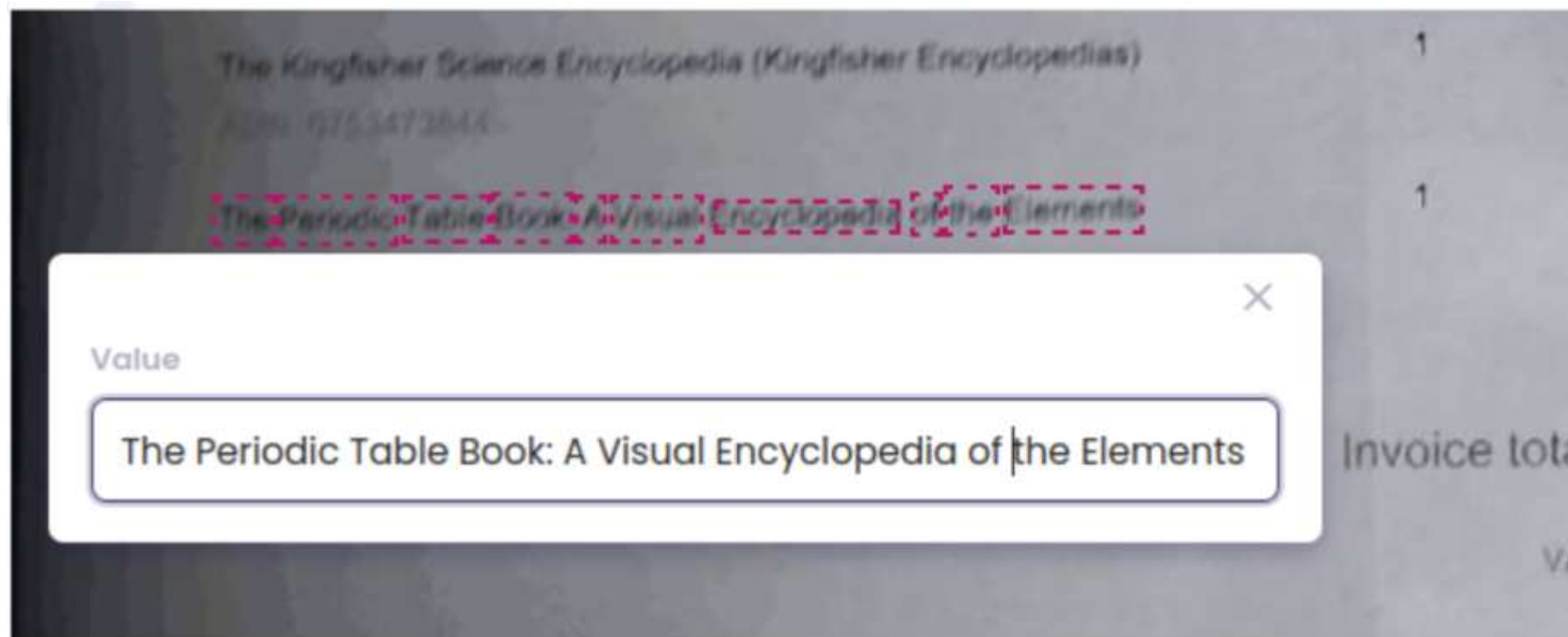


- health records
- ID scans
- medical certificates

### LEGAL



- agreements
- contracts
- notarial deeds
- forms



**prime video**

# INVOICE

Attn: Olivia James  
Amazon Digital UK Limited  
1 Principal Place, Worship Street  
London, EC2A 2FA  
United Kingdom  
VAT #: GB982596668

Order date: 27.12.2020  
Date of Supply: 27.12.2020  
Order #: 001-6134136-6607068

Invoice date / Delivery date: 27.12.2020  
Invoice #: A/VGB-INV-GB-2020-97911707

Issued to:  
Mathew James  
Bourville Garden Centre  
Science Park, Milton Road  
CAMBRIDGE, CAMBS, CB4 0FY  
GB

Qty	Description	Unit price (excl. VAT)	VAT rate	Unit price (incl. VAT)	Total price (incl. VAT)
1	The Chronicles of Narnia - The Lion, the Witch and the Wardrobe	£4.99	20%	£5.99	£5.99
<b>TOTAL:</b>					<b>£5.99</b>

VAT details:

Subtotal (excl. VAT)	VAT rate	VAT amount
£4.99	20%	£0.99
<b>TOTAL:</b>		<b>£5.99</b>

Correction invoice: CI-CK-20-00001  
Reference number: 282000001

Supplier: Release agency  
2020 Milton Road  
11716 New York  
United States  
VAT ID: 080021, Company ID: 2021

Client: Yummy & Tasty  
1410 Flinderston Road  
60605 Chicago  
United States, Illinois  
Company ID: , VAT ID:

Bank: National Bank  
Account no.: 27128049  
Branch no.: 0002  
Routing no.: 0002  
Payment condition: Bank transfer

Invoice date: 17/12/2020  
Supply date: 17/12/2020  
Due date: 01/01/2021  
Related invoice number: 98-20-00011

Before correction

Subject	Count	Unit price	Amount VAT incl.	EUR	VAT
	100.00	2.00	200.00	EUR	21%
	500.00	1.00	500.00	EUR	21%
	10.00	300.00	3,000.00	EUR	21%
			<b>3,700.00</b>	<b>EUR</b>	

Correction

Subject

Advertising materials	Count	Unit price	Amount VAT incl.	EUR	VAT
Leaflets	0.00	2.00	0.00	EUR	21%
Billboards	1000.00	0.00	1,000.00	EUR	21%
	0.00	0.00	0.00	EUR	21%
<b>Total</b>			<b>400.00</b>	<b>EUR</b>	

	Count	Unit price	Amount VAT incl.	EUR	VAT
	100.00	4.00	400.00	EUR	21%
	400.00	1.00	400.00	EUR	21%
Billboards	10.00	300.00	3,000.00	EUR	21%
<b>Total</b>			<b>4,200.00</b>	<b>EUR</b>	

VAT recapitulation

Base	VAT	Total VAT incl.	EUR
21%	21.00	400.00	544.50
<b>Total</b>		<b>400.00</b>	<b>544.50</b>

Payable

Total VAT incl.	Received payment	Total due	EUR
544.50	0.00	544.50	EUR
		<b>544.50</b>	<b>EUR</b>

## How Document Classification Works

