

# NATURAL LANGUAGE PROCESSING WITH MACHINE LEARNING

Yogesh Haribhau Kulkarni



# About Me

# Yogesh Haribhau Kulkarni

## Bio:

- ▶ 20+ years in CAD/Engineering software development
- ▶ Got Bachelors, Masters and Doctoral degrees in Mechanical Engineering (specialization: Geometric Modeling Algorithms).
- ▶ Currently doing Coaching in fields such as Data Science, Artificial Intelligence Machine-Deep Learning (ML/DL) and Natural Language Processing (NLP).
- ▶ Feel free to follow me at:
  - ▶ Github ([github.com/yogeshhk](https://github.com/yogeshhk))
  - ▶ LinkedIn ([www.linkedin.com/in/yogeshkulkarni/](https://www.linkedin.com/in/yogeshkulkarni/))
  - ▶ Medium ([yogeshharibhaukulkarni.medium.com](https://yogeshharibhaukulkarni.medium.com))
  - ▶ Send email to [yogeshkulkarni at yahoo dot com](mailto:yogeshkulkarni@yahoo.com)



Office Hours:  
Saturdays, 2 to 5pm  
(IST); Free-Open to all;  
email for appointment.

# Introduction

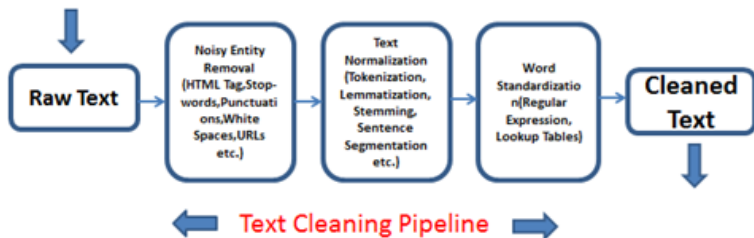
- ▶ Natural Language Processing (NLP) and Machine Learning (ML) form the foundation for modern intelligent systems.
- ▶ Retrieval Augmented Generation (RAG) leverages NLP pipelines, embeddings, and ML for knowledge-grounded responses.
- ▶ This presentation explores:
  - ▶ Text preprocessing and representation
  - ▶ spaCy-based NLP components
  - ▶ ML techniques for classification and clustering
  - ▶ Embedding-based search and semantic understanding
  - ▶ Sentiment analysis applications



(Ref: Stanford NLP Group)

# Text Preprocessing

- ▶ Clean and normalize raw text data for consistent downstream analysis.
- ▶ Common steps include lowercasing, punctuation & digit removal, and whitespace normalization.
- ▶ Ensures reliable tokenization and model training.
- ▶ Helps retain semantic integrity by focusing on meaningful linguistic content.



(Ref: Text Cleaning Steps ResearchGate )

# Tokenization

- ▶ Tokenization splits text into smaller linguistic units (words, subwords, or sentences).
- ▶ spaCy provides efficient rule-based and statistical tokenizers.
- ▶ Basis for vectorization, tagging, and parsing.

```
1 import spacy
  nlp = spacy.load("en_core_web_sm")
3 doc = nlp("Retrieval-Augmented Generation is powerful!")
  tokens = [token.text for token in doc]
5 print(tokens)
  # ['Retrieval', '-', 'Augmented', 'Generation', 'is', 'powerful', '!']
7
```

## Lowercasing & Cleaning

- ▶ Lowercasing removes case sensitivity in NLP models.
- ▶ Cleaning eliminates noise: URLs, HTML tags, punctuation, numbers.
- ▶ Helps search and RAG systems achieve higher recall and consistency.

```
import re
2 text = "Visit https://openai.com for <b>AI Research!</b>"
cleaned = re.sub(r"http\S+|<.*?>", "", text).lower()
4 print(cleaned)
# 'visit for ai research!'
6
```



# Stop Word Removal

- ▶ Removes high-frequency words that add little meaning.
- ▶ Reduces dimensionality and improves efficiency of embeddings.
- ▶ Example stop words: “the”, “is”, “in”, “on”.

```
1 from spacy.lang.en.stop_words import STOP_WORDS
  tokens = ["Retrieval", "Augmented", "Generation", "is", "powerful"]
3 filtered = [w for w in tokens if w.lower() not in STOP_WORDS]
  print(filtered)
5 # ['Retrieval', 'Augmented', 'Generation', 'powerful']
```

# Stemming & Lemmatization

- ▶ Stemming removes suffixes using heuristic rules.
- ▶ Lemmatization uses grammar-based transformations to return dictionary form.
- ▶ Reduces vocabulary size and enhances text matching in retrieval.

```
1 from nltk.stem import WordNetLemmatizer
  lem = WordNetLemmatizer()
3 print(lem.lemmatize("running", "v")) # 'run'
  print(lem.lemmatize("better", "a")) # 'good'
5
```

# Regular Expressions

- ▶ Enable pattern-based extraction from text.
- ▶ Examples: extract emails, phone numbers, or URLs.
- ▶ Essential for data cleaning, anonymization, and text mining.

```
import re
2 text = "Contact us at info@company.com or sales@domain.org"
emails = re.findall(r"\b[A-Za-z0-9._%+-]+@[A-Za-z]+\.[A-Z|a-z]{2,}\b", text)
4 print(emails)
# ['info@company.com', 'sales@domain.org']
6
```

# Project: Text Preprocessing Pipeline

- ▶ Modular preprocessing pipeline for flexible text workflows.
- ▶ Each function handles one cleaning step.
- ▶ Easily extendable for lemmatization, stemming, or entity masking.

```
1 class TextPipeline:
2     def __init__(self, steps):
3         self.steps = steps # list of functions
4
5     def run(self, text):
6         for fn in self.steps:
7             text = fn(text)
8         return text
9
10 pipeline = TextPipeline([basic_clean])
11 print(pipeline.run("Hello World! This is 2025."))
```

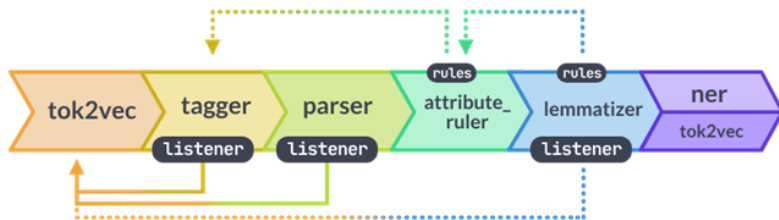
# Project: Document Cleaner

- ▶ Extracts text from PDFs, HTML, and Word documents.
- ▶ Normalizes whitespace, removes headers/footers, handles encoding.
- ▶ Outputs clean text ready for embedding or indexing.

```
1 from bs4 import BeautifulSoup
2 from PyPDF2 import PdfReader
3
4 def clean_document(path):
5     if path.endswith(".pdf"):
6         text = "".join([page.extract_text() for page in PdfReader(path).pages])
7     elif path.endswith(".html"):
8         html = open(path).read()
9         text = BeautifulSoup(html, "html.parser").get_text()
10    else:
11        text = open(path).read()
12    return re.sub(r"\s+", " ", text.strip())
13
```

# spaCy Pipeline Basics

- ▶ spaCy's pipeline includes tokenization, tagging, parsing, and entity recognition.
- ▶ Designed for efficiency and extensibility.
- ▶ Supports custom components for preprocessing and postprocessing.



(Ref: spaCy Docs)

# Part-of-Speech Tagging

- ▶ Assigns syntactic roles: noun, verb, adjective, etc.
- ▶ Enables grammar-aware retrieval, information extraction, and text summarization.

```
doc = nlp("NLP models enhance communication.")
2 for token in doc:
    print(token.text, token.pos_)
4 # Output: NLP NOUN, models NOUN, enhance VERB, communication NOUN
```

# Named Entity Recognition (NER)

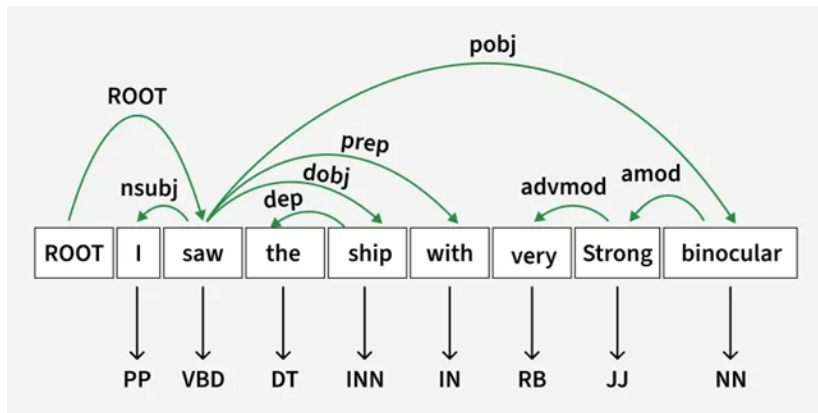
- Identifies named entities: PERSON, ORG, LOC, DATE, MONEY, etc.
- Crucial for knowledge extraction and linking to external databases.

```
doc = nlp("OpenAI was founded in San Francisco in 2015.")
2 for ent in doc.ents:
    print(ent.text, ent.label_)
4 # OpenAI ORG, San Francisco GPE, 2015 DATE
```



# Dependency Parsing

- Analyzes grammatical structure of sentences.
- Reveals subject-object relationships.
- Basis for relation extraction and semantic role labeling.



(Ref: Dependency Parsing with NLTK - GeeksforGeeks)

# Text Similarity

- ▶ Measures semantic closeness between two texts.
- ▶ Useful in RAG for ranking documents by contextual relevance.

```
1 doc1 = nlp("AI improves healthcare")
2 doc2 = nlp("Artificial intelligence helps medicine")
3 print(doc1.similarity(doc2))
4 # Output: similarity score ~0.85
```

## Project: Named Entity Extractor

- ▶ Extracts and exports named entities from text.
- ▶ Can be integrated with visualization tools or stored in databases.

```
def extract_entities(text):  
2     nlp = spacy.load("en_core_web_sm")  
    doc = nlp(text)  
4     return {ent.text: ent.label_ for ent in doc.ents}  
  
6 print(extract_entities("Elon Musk leads SpaceX and Tesla."))
```

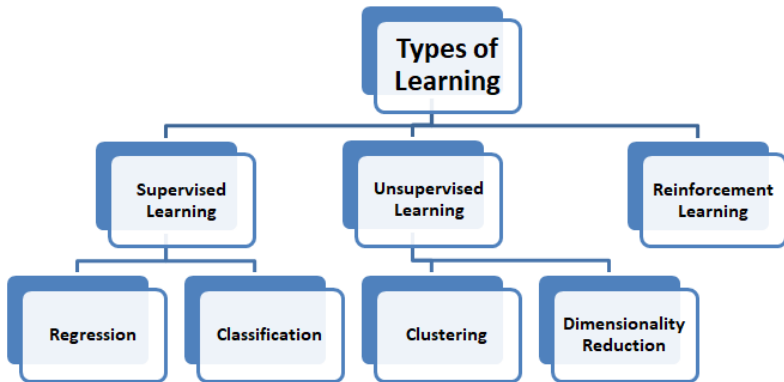
# Project: Text Similarity Engine

- Uses embeddings to compute cosine similarity between texts.
- Enables semantic document search and clustering.

```
2 from sklearn.metrics.pairwise import cosine_similarity
4
6 def rank_docs(query, docs, model):
8     q_vec = model.encode([query])
9     d_vecs = model.encode(docs)
10    scores = cosine_similarity(q_vec, d_vecs)[0]
11    ranked = sorted(zip(docs, scores), key=lambda x: x[1], reverse=True)
12    return ranked[:3]
```

# Types of Machine Learning

- ▶ **Supervised:** learns from labeled data (e.g., text classification).
- ▶ **Unsupervised:** discovers hidden structure (e.g., clustering).
- ▶ **Reinforcement:** learns from feedback/reward.
- ▶ NLP tasks map naturally to these paradigms.



(Ref: <https://www.studytrigger.com/article/types-of-learning-in-machine-learning/>)

# Classification

- ▶ Predicts categorical labels (spam, positive/negative sentiment).
- ▶ Logistic Regression is a strong baseline for text classification.

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.linear_model import LogisticRegression
3
4 vec = TfidfVectorizer()
5 X_train = vec.fit_transform(train_texts)
6 model = LogisticRegression(max_iter=200)
7 model.fit(X_train, y_train)
8 preds = model.predict(vec.transform(test_texts))
```

# Clustering

- ▶ Groups similar documents based on content.
- ▶ Common methods: KMeans, DBSCAN, Agglomerative.
- ▶ Supports topic discovery and unsupervised organization.

```
1 from sklearn.cluster import KMeans
2 from sklearn.decomposition import PCA
3 import matplotlib.pyplot as plt
4
5 km = KMeans(n_clusters=5, random_state=42)
6 km.fit(embeddings)
7 labels = km.labels_
8
9 pca = PCA(n_components=2)
10 reduced = pca.fit_transform(embeddings)
11 plt.scatter(reduced[:,0], reduced[:,1], c=labels)
12 plt.show()
```

# Project: Iris Flower Classifier

- ▶ Classic ML dataset demonstration.
- ▶ Illustrates feature scaling, model training, and evaluation.

```
from sklearn.datasets import load_iris
2 from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import accuracy_score

6 X, y = load_iris(return_X_y=True)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)
8 model = LogisticRegression(max_iter=300)
model.fit(X_train, y_train)
10 print("Accuracy:", accuracy_score(y_test, model.predict(X_test)))
```



## Project: Document Clusterer

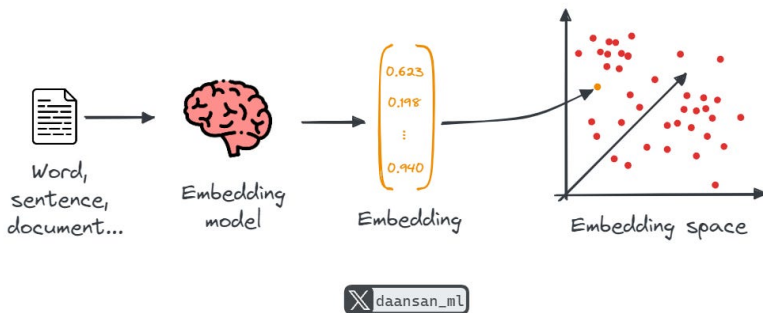
- ▶ Embeds documents using spaCy or sentence-transformers.
- ▶ Applies KMeans for unsupervised grouping.
- ▶ Visualizes using t-SNE or PCA for semantic clusters.

```
import spacy
2 from sklearn.cluster import KMeans
nlp = spacy.load("en_core_web_md")
4
docs = ["AI in healthcare", "Quantum computing basics", "AI in finance"]
6 vectors = [nlp(d).vector for d in docs]
labels = KMeans(n_clusters=2).fit_predict(vectors)
8 print(list(zip(docs, labels)))
```

# Word Embeddings Concepts

- ▶ Represent words as dense numerical vectors capturing context.
- ▶ Models: Word2Vec, GloVe, FastText, Transformer embeddings.
- ▶ Enable semantic understanding and transfer learning.

## Embeddings



(Ref: <https://mlpills.substack.com/p/issue-58-embeddings-in-nlp>)

# Project: Semantic Search Engine

- ▶ Retrieve documents semantically, not just by keyword match.
- ▶ Uses embedding models and cosine similarity for ranking.

```
1 from sentence_transformers import SentenceTransformer, util
2 model = SentenceTransformer("all-MiniLM-L6-v2")
3
4 def semantic_search(query, docs):
5     q_emb = model.encode(query, convert_to_tensor=True)
6     d_emb = model.encode(docs, convert_to_tensor=True)
7     scores = util.cos_sim(q_emb, d_emb)
8     return sorted(zip(docs, scores[0]), key=lambda x: float(x[1]),
9                   reverse=True)[:3]
```

## Project: Word Analogy Solver

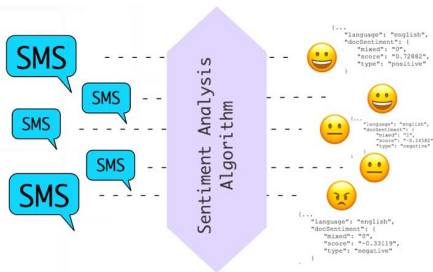
- Demonstrates vector arithmetic:  $king - man + woman \approx queen$
- Captures latent gender, tense, and semantic relations.

```
from gensim.models import KeyedVectors
2 model = KeyedVectors.load_word2vec_format("GoogleNews-vectors.bin",
    binary=True)
result = model.most_similar(positive=["king", "woman"], negative=["man"],
    topn=1)
4 print(result)  # [('queen', 0.75)]
```

# Sentiment Analysis Basics

- ▶ Detects emotional polarity (positive, neutral, negative).
- ▶ Used in reviews, feedback, and social media monitoring.
- ▶ Can be rule-based, ML-based, or transformer-based.

## What is Sentiment Analysis?



(Ref: Sentiment Analysis — Engati)

## Project: Review Sentiment Analyzer

- ▶ Train a sentiment classifier using Naive Bayes.
- ▶ Vectorize text with TF-IDF.
- ▶ Evaluate accuracy using test set.

```
1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.metrics import accuracy_score
3
4 X = vec.fit_transform(reviews)
5 model = MultinomialNB()
6 model.fit(X, labels)
7 preds = model.predict(vec.transform(test_reviews))
8 print("Accuracy:", accuracy_score(test_labels, preds))
```

## Project: Tweet Sentiment Dashboard

- ▶ Streams tweets and classifies their sentiment in real-time.
- ▶ Visualizes sentiment trends dynamically using dashboards.

```
import tweepy, pandas as pd
2 from textblob import TextBlob

4 def classify_sentiment(text):
    polarity = TextBlob(text).sentiment.polarity
    6     return "positive" if polarity > 0 else "negative" if polarity < 0 else
        "neutral"

8 # TODO: integrate with dashboard for visualization
```

# Conclusions

- ▶ Robust NLP preprocessing and embeddings are critical for effective information retrieval.
- ▶ spaCy and scikit-learn streamline NLP + ML workflows.
- ▶ Vector-based understanding enables semantic and contextual search.
- ▶ Combining these components forms the backbone of Retrieval-Augmented Generation.
- ▶ Future directions:
  - ▶ Fine-tuning LLMs for domain tasks
  - ▶ Efficient vector database retrieval (FAISS, Milvus)
  - ▶ Integrating sentiment and entity-level reasoning



# References

- ▶ Explosion AI, *spaCy Documentation*, <https://spacy.io/>
- ▶ Pedregosa et al., *Scikit-learn: Machine Learning in Python*, JMLR, 2011.
- ▶ Mikolov et al., *Efficient Estimation of Word Representations in Vector Space*, arXiv, 2013.
- ▶ Bird et al., *Natural Language Processing with Python*, O'Reilly, 2009.
- ▶ OpenAI, *Retrieval-Augmented Generation (RAG) Overview*, 2024.
- ▶ Manning & Schütze, *Foundations of Statistical Natural Language Processing*, MIT Press.
- ▶ Jurafsky & Martin, *Speech and Language Processing*, 3rd Ed. Draft, 2023.

# Thanks ...

- ▶ Office Hours: Saturdays, 3 to 5 pm (IST);  
Free-Open to all; email for appointment to  
yogeshkulkarni at yahoo dot com
- ▶ Call + 9 1 9 8 9 0 2 5 1 4 0 6



(<https://www.linkedin.com/in/yogeshkulkarni/>)



(<https://medium.com/@yogeshharibhaukulkarni> )



(<https://www.github.com/yogeshhk/> )