



AULA 11

Processamento de Linguagem
Natural – Parte II

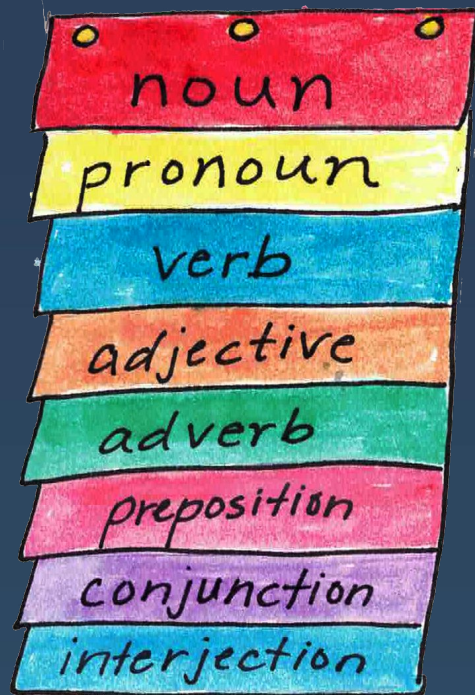
ROTEIRO DA AULA

1. Parts of Speech – POS
2. Named Entity Recognition – NER
3. Extração de Features: Vetorização

-
- A word cloud titled "Natural Language Processing" in a large, bold, black font. The words are arranged in a circular pattern around the title. The words are in various colors (blue, green, yellow, orange, red, pink, purple) and sizes. The words include: DIFFICULT, MACHINE, INFORMATION, MORPHOLOGY, INTERPRET, UNDERSTAND, SUCCESSFUL, GRAMMAR, COMPUTATION, TEXTUAL, NOUN, VERBALIZATION, STARTINGLY, EVALUATION, TAGGERS, COMPLEX, WORDS, ANALYSIS, DEVOYED, PSYCHOLOGICAL, LINGUISTICS, RESOLUTION, AMBIGUOUS, TRANSLATION, SUMMARIZATION, OUTPUT, NATURAL, SYNTACTIC, PROCESSING, DISCOURSE, INPUT, RECOGNITION, WORD, PROJECT, PROBABILITY, CITY, ENTITY, PARSES, RUBRICS, AMBIGUITY, GOLD, MULTIPLE, COMPLEXITY, AUTOMATIC, SUBFIELD, ONTOLOGY, and many others.

Part of Speech – POS

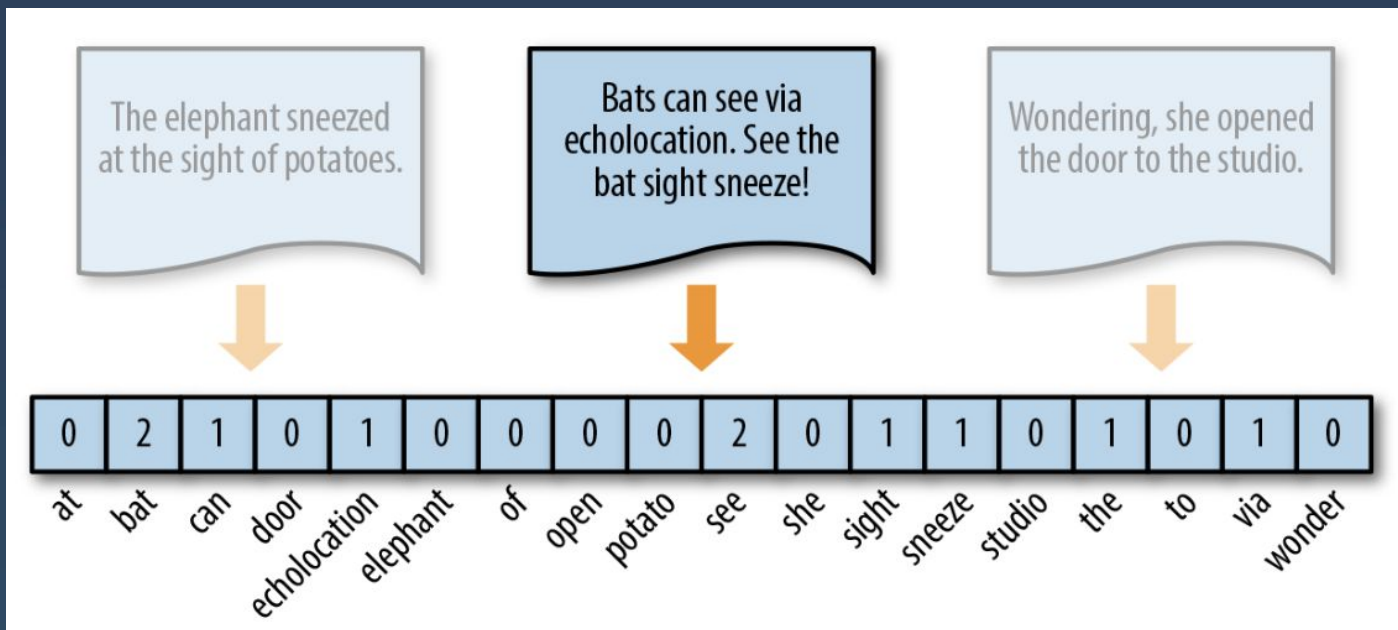
- Técnica que lê um texto, em algum idioma, e assinala, para cada palavra, a classe gramatical a qual ela pertence.
- Pode ser classificada em:
 - **Coarse-grained POS tags:** são tags mais genéricas, representando as grandes classes gramaticais (substantivos, verbos, adjetivos).
 - **Fine-grained POS tags:** são tags mais específicas, contendo detalhes morfológicos do token (substantivo no plural, verbo no futuro do pretérito, adjetivo superlativo).



Extração de Features

- Transformação do texto em um conjunto de dados estruturados (vetores de números) que pode ser utilizado em técnicas analíticas e aprendizado de máquina. Também conhecido como **Vetorização**. Tipicamente a saída da modelagem recebe o nome de Word Embedding.
- Bag of Words (Counter Vectorization)
- TF-IDF

Bag of Words



TF-IDF

- ***Term Frequency – Inverse Document Frequency***
- **Term Frequency:** Frequência da palavra no documento atual
 - $TF = (\# \text{ vezes que o termo } t \text{ aparece no documento}) / (\# \text{ de termos no documento})$
- **Inverse document frequency (IDF):** o quão raro é o termo no documento
 - $IDF = \log(N/n)$, N é o # de documentos e n é o # de documentos que o termo t aparece.
- **TF-IDF:** Importância de um termo para um documento para uma coleção.

TF-IDF

- **Exemplo:**

1. **Documento com 100 termos, no qual o termo “cachorro” aparece 5 vezes**

$$TF = 5 / 100 = 0.05$$

2. **Temos 100 documentos (N) e o termo “cachorro” aparece em 20 desses documentos (n)**

$$IDF = \log(100/20) = 0.69$$

$$TF-IDF = 0.05 * 0.69 = 0.034$$