
UNIVERSIDADE FEDERAL DE ALAGOAS
INSTITUTO DE COMPUTAÇÃO

Processamento de Linguagem Natural
Professor: Thales Vieira

3a lista de exercícios

27 de setembro de 2024

Instruções:

A lista deve ser respondida por grupos de até 2 pessoas (graduação) e individualmente (mestrado).

Resoluções idênticas de grupos distintos serão desconsideradas.

O código e demais dados devem ser anexados a cada questão.

Data limite para entrega: 20/10/2024.

Usando sua base de textos após os pré-processamento realizados na lista 1, realize as seguintes tarefas:

1. O objetivo dessa questão é desenvolver buscadores de palavras e documentos.

- a) escolha e aplique um modelo do tipo word2vec a seus textos, compatível com o idioma de seus textos (inglês ou português).
- b) escolha 5 palavras de consulta que não estejam em nenhum dos textos. Para cada palavra de consulta, encontre as 3 palavras **de seu conjunto de textos** mais parecidas com cada uma das palavras de consulta e exiba os documentos onde estas palavras aparecem.
- c) Seja d um documento da base e w uma palavra de consulta. Implemente o seguinte algoritmo para buscar documentos:
 1. Encontre $d_{10}(w)$: a lista com as 10 palavras mais parecidas com w em um certo documento d .
 2. Para cada documento d , calcule a distância média $DM_{10}(w)$ entre w e as palavras de $d_{10}(w)$.
 3. Recupere os 3 documentos da base cuja $DM_{10}(w)$ é menor.
- d) aplique o algoritmo da letra c) para buscar documentos em 5 palavras distintas, e exiba os 3 documentos mais próximos de cada um.

2. Resolva novamente a primeira questão da 2a lista e compare com os resultados obtidos anteriormente:

- a) Aplicando a representação vetorial Doc2Vec combinado com os classificadores usados anteriormente.
 - b) Usando uma arquitetura de rede neural que utilize camadas de Embedding e LSTM.
- 3.** Usando sua base de textos e a biblioteca spaCy, realize as seguintes tarefas:
- a) Extraia as etiquetas gramaticais (POS) de cada token do seu textos.
 - b) Calcule e plote um gráfico com as frequências de cada tipo gramatical.
 - c) Extraia entidades do tipo pessoa e lugar dos seus textos.
 - d) Identifique e liste as pessoas mais frequentes nos seus textos. Você só deve contar cada entidade 1 vez por documento.
- 4.** Estude o tutorial *Character-level recurrent sequence-to-sequence model* disponível em https://keras.io/examples/nlp/lstm_seq2seq/.
- a) Treine um outro modelo de tradução entre línguas distintas. Você pode encontrar conjuntos de treinamento em <http://www.manythings.org/anki/>.
 - b) Exiba 5 exemplos de tradução de frases curtas.