

Indexação e Tratamento de Dados Heterogêneos: Variedade

Mecanismos de Busca



Conteúdo do Bloco

Indexação

Recuperação de Informação

Solr

Elasticsearch



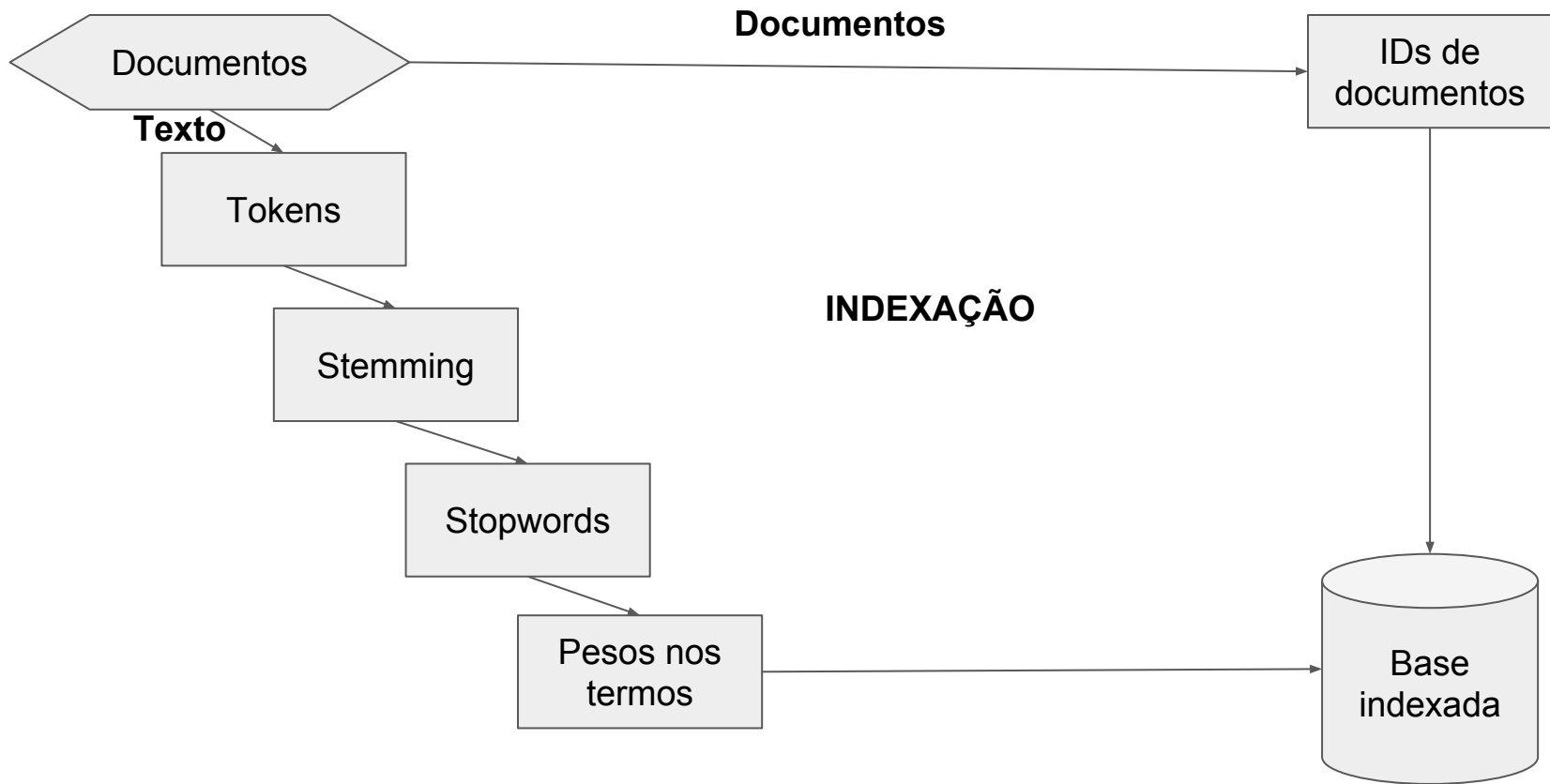
Aula 2

Métricas em information retrieval

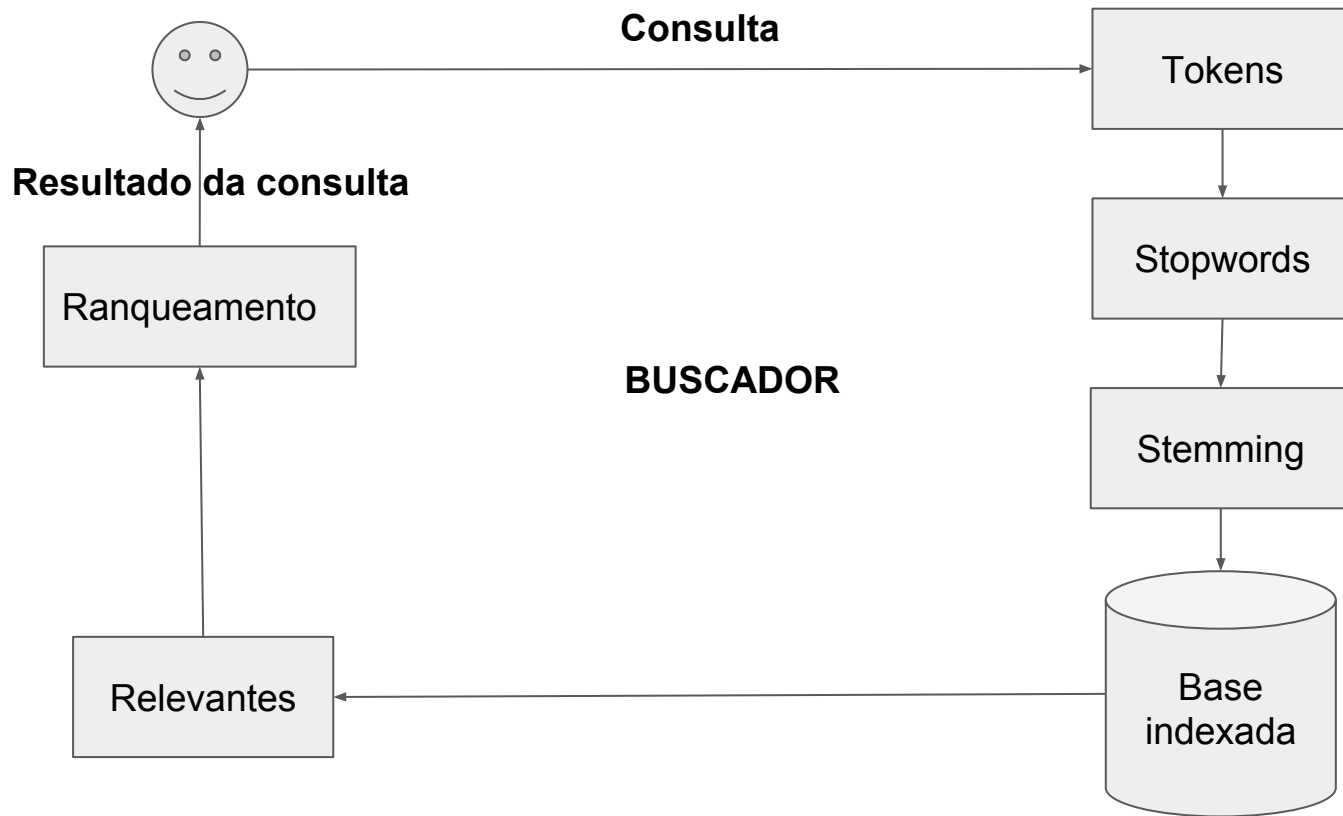
Um pouco mais sobre modelo vetorial

Brincando com NLTK

Relembrando



Relembrando



Relembrando

Similaridade entre consulta e documento no modelo vetorial

- Cosseno

Algoritmos de similaridade diferentes interferem no resultado final

- Existem outras métricas de similaridade?

Outras técnicas de similaridade

Coeficiente de Jaccard $s_{AB} = \frac{|A \cap B|}{|A \cup B|}$

Qual a distância entre dois objetos? Qual a similaridade entre eles?

- D1 = {aula, mba, bigdata, infnet}
- D2 = {aula, analise, sistemas, infnet}
 - $S(d1, d2) = \{aula, infnet\} / \{aula, mba, bigdata, analise, sistemas, infnet\} = 1/3$
 - Coeficiente de Jaccard = $1/3$ (similaridade)
 - Distância de Jaccard = $2/3$ (não similaridade)

Coeficiente de Pearson

Distância Euclidiana

Métricas

Como saber se um sistema de busca e recuperação da informação está se comportando bem?

- Quanto preciso é o meu buscador?
- Quantos documentos relevantes meu buscador exibe como resultado?
- Existe uma hierarquia de relevância dentro dos documentos relevantes, o buscador exibe o resultado na ordem correta?

Métricas

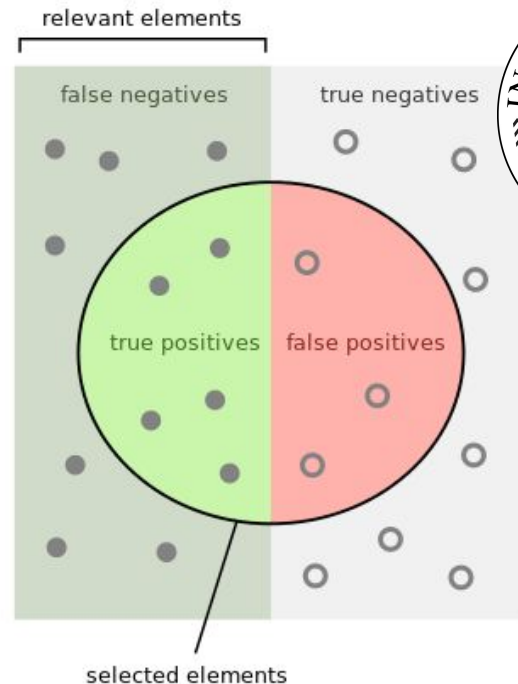
Precisão e revocação (precision and recall)

- Precision = relevantes / todos documentos
 - Quão útil foi o resultado?
- Recall = relevantes / todos relevantes
 - Quão completo foi o resultado?

Precision @ K

- K documentos relevantes na primeira página de resultado da consulta

$$F1 \text{ score} = 2 * ((\text{precision} + \text{recall}) / (\text{precision} * \text{recall}))$$



How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

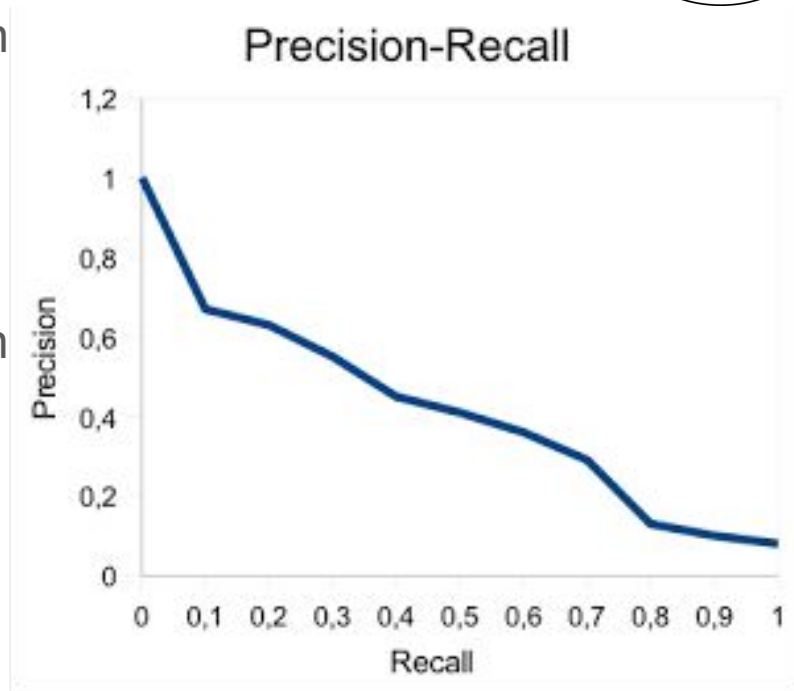
How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

Métricas

Interpolação de 11 pontos de recall e precision

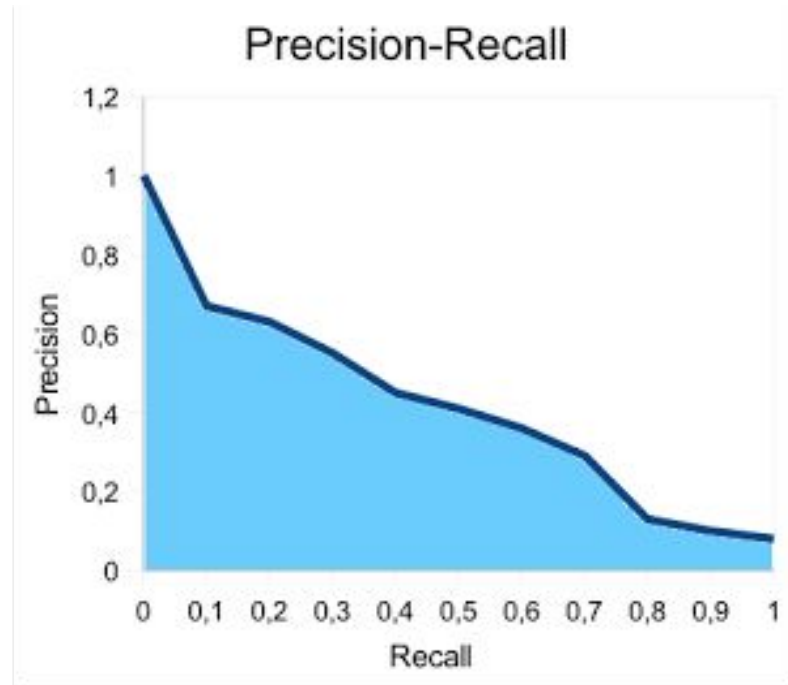
- Precision e recall são inversamente proporcionais
- Quanto mais documentos retornados: Maior fica o recall e menor fica o precision



Métricas

Mean Average Precision (MAP)

- Precisão média de cada consulta / total de consultas



Métricas

DCG (Discounted Cumulative Gain)

- Ganho cumulativo descontado
 - Páginas relevantes que aparecem com ranqueamento baixo são penalizadas

Rank	Relevância	Ganho descontado	DCG
1	2	2/1	2
2	0	0/2	2
3	3	3/3	3
4	2	2/4	3,5

Métricas

NDCG (Normalized DCG)

- $NDCG = DCG / IDC$
 - IDC = Ideal discounted cumulative gain

Rank	Relevância	Ganho descontado	DCG	Ganho ideal descontado	IDC	NDCG
1	2	2/1	2	3/1	3	0,67
2	0	0/2	2	2/2	4	0,5
3	3	3/3	3	2/3	4,67	0,64
4	2	2/4	3,5	0/4	4,67	0,75

Part of speech tag

Serve para diminuir ambiguidades

Melhora pré-processamento

Word

Tag

heat

verb (noun)

water

noun (verb)

in

prep (noun, adv)

a

det (noun)

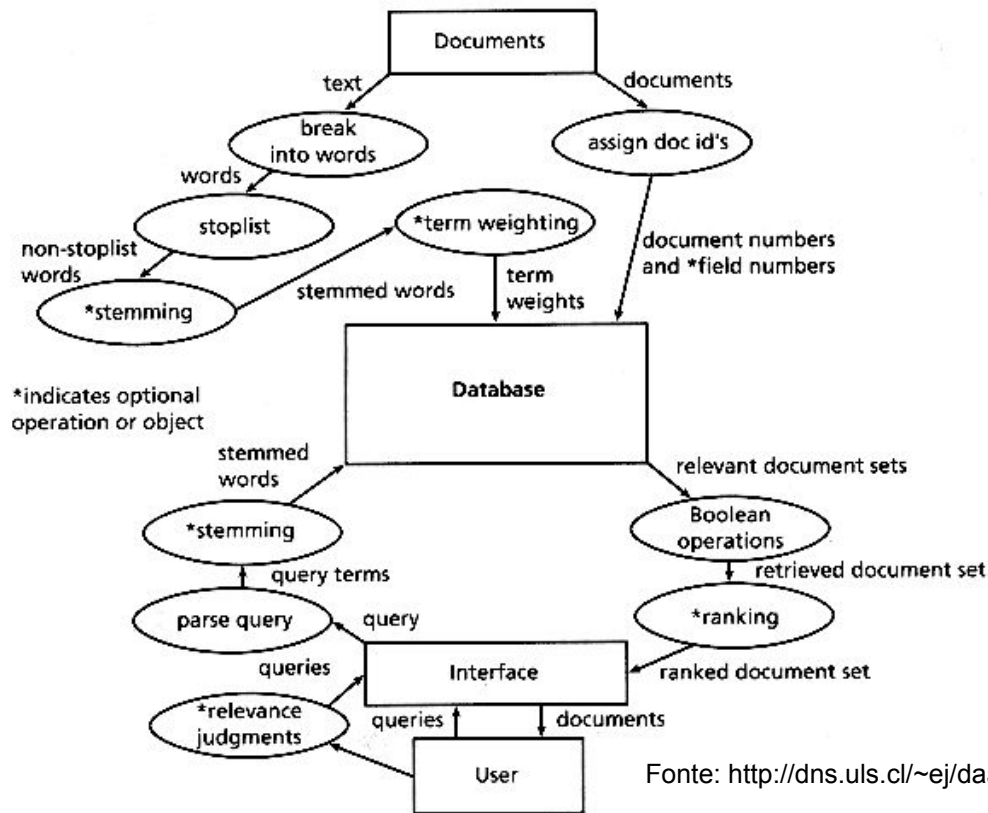
large

adj (noun)

vessel

noun

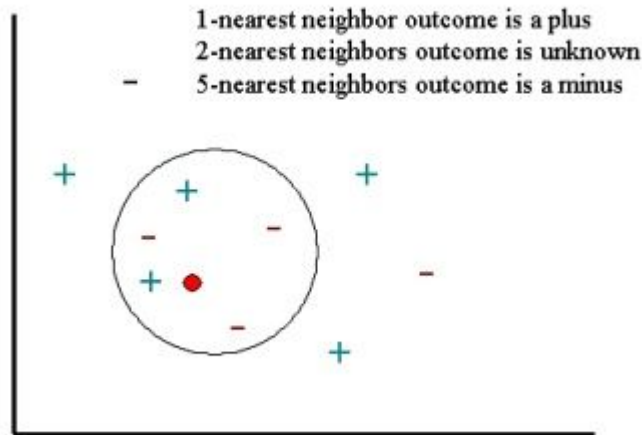
Recuperação da Informação



Recuperação da Informação

Algoritmos comumente usados em Information Retrieval

- K-NN (K nearest neighbor)
 - Por exemplo, usado para classificar documentos dentro de uma categoria
- Outros métodos serão descritos no decorrer do curso



Praticando com o NLTK

Vamos testar o NLTK segundo o livro da ferramenta

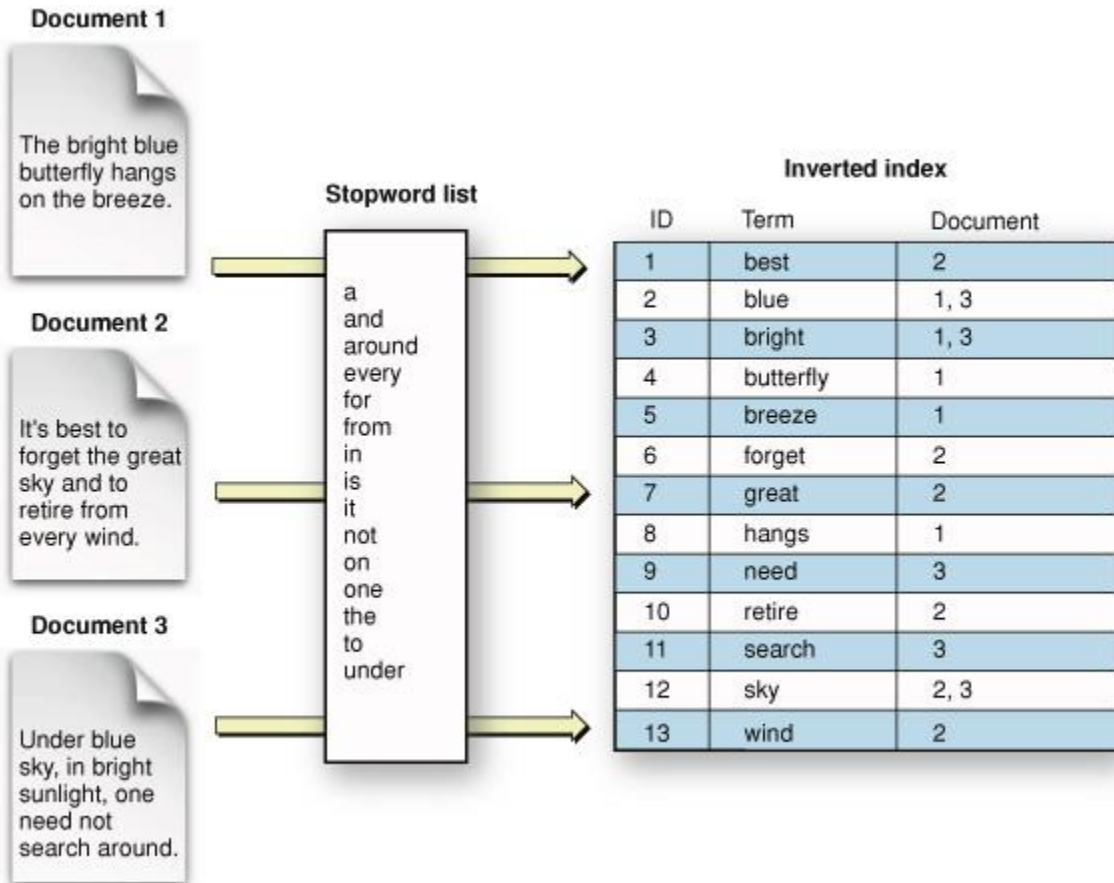
- <http://www.nltk.org/book/>
 - Capítulo 1
 - Capítulo 3
 - Capítulo 5
- Capítulos acima nortearão a implementação do segundo trabalho semanal
- Praticando enquanto construímos nosso próprio sistema de IR

Exercício

Criar uma lista invertida em python:

- Pegar 3 textos da imagem seguinte e inicializar cada um em um arquivo diferente
- Extrair o texto desses arquivos para dentro do python (3 variáveis ou array)
- Aplicar o tratamento aprendido na aula 1 e praticado na aula 2 para criação dos termos e remoção de stopwords
- Criar a lista de ocorrências de termos dentro dos moldes de uma lista invertida e salvar em arquivo
- Usar a figura seguinte como gabarito (aproximado)

Exercício



Exercício

Sugestão para estrutura a ser salva em arquivo

- Chave é o termo, valor é um array de documento/ocorrências desses termos
 - blue;['1', '3']

Enviar arquivo (zipado) com código do exercício

- raul.ferreira@prof.infnet.edu.br
- Colocar no assunto do email “Trabalho 2 - Turma de quinta - Infnet”
- Colocar identificação no corpo do email (nome e sobrenome)

Trabalho individual. Trabalhos copiados = zero



Próxima aula

Recuperação da informação

Mais implementações em sala usando python

Continuar a construção de parte do nosso buscador