

Introdução à Recuperação de Informações

<https://github.com/fccoelho/curso-IRI>

IRI 2: Vocabulário de termos e lista de “postings”

Flávio Codeço Coelho

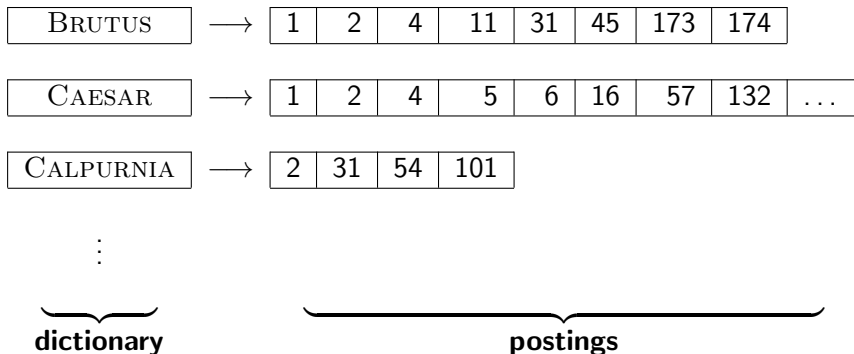
Escola de Matemática Aplicada, Fundação Getúlio Vargas

Sumário da Aula

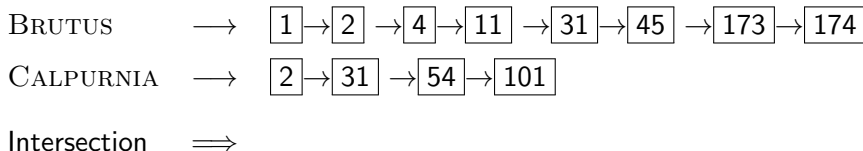
- 1 Recapitulação
- 2 Documentos
- 3 Termos
 - Genericos + outras línguas
 - English
- 4 Consultas de frases

Índice invertido

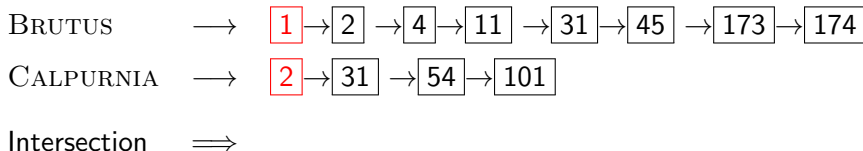
Para cada termo t , armazenamos uma lista de todos os documentos que contém t .



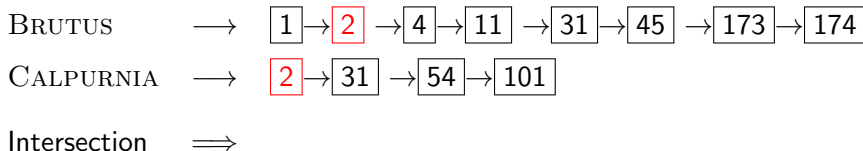
Interseção de duas listas de “postings”



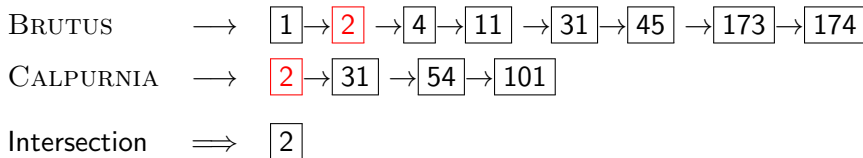
Interseção de duas listas de “postings”



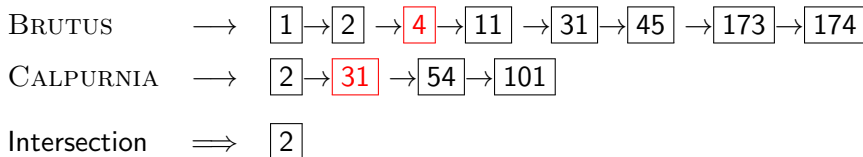
Interseção de duas listas de “postings”



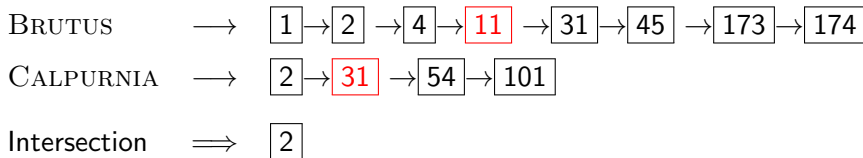
Interseção de duas listas de “postings”



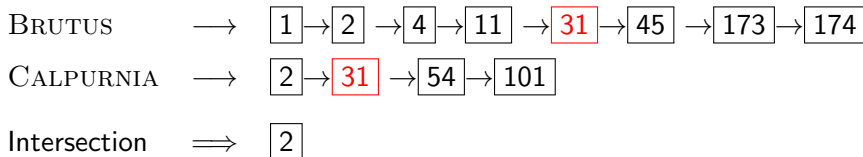
Interseção de duas listas de “postings”



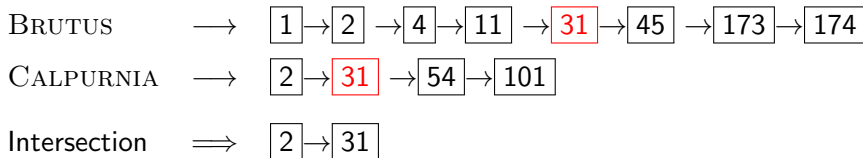
Interseção de duas listas de “postings”



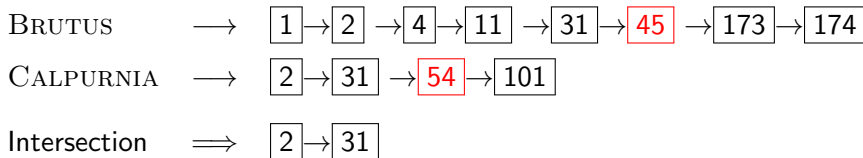
Interseção de duas listas de “postings”



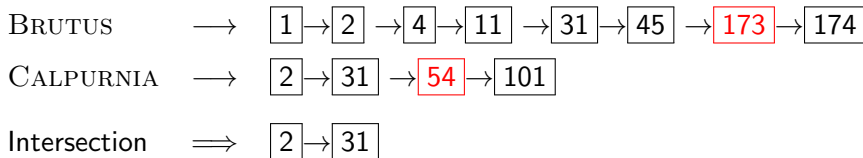
Interseção de duas listas de “postings”



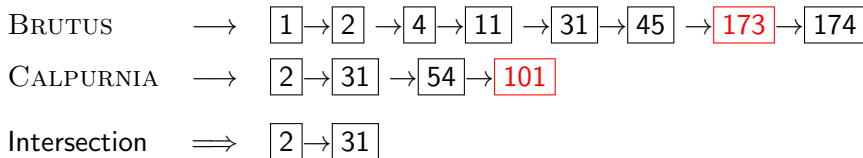
Interseção de duas listas de “postings”



Interseção de duas listas de “postings”



Interseção de duas listas de “postings”



Construindo o índice invertido: ordenando os “postings”

term	docID		term	docID
I	1		ambitious	2
did	1		be	2
enact	1		brutus	1
julius	1		brutus	2
caesar	1		capitol	1
I	1		caesar	1
was	1		caesar	2
killed	1		caesar	2
i'	1		did	1
the	1		enact	1
capitol	1		hath	1
brutus	1		I	1
killed	1		I	1
me	1	⇒	i'	1
so	2		it	2
let	2		julius	1
it	2		killed	1
be	2		killed	1
with	2		let	2
caesar	2		me	1
the	2		noble	2
noble	2		so	2
brutus	2		the	1
hath	2		the	2
told	2		told	2
you	2		you	2
caesar	2		was	1
was	2		was	2
ambitious	2		with	2

O Google usa busca Booleana?

- No Google, a interpretação default de uma consulta $[w_1 w_2 \dots w_n]$ é $w_1 \text{ AND } w_2 \text{ AND } \dots \text{ AND } w_n$
- Casos em que você recebe resultados que não contém uma das w_i :
 - âncora de texto
 - página contém variantes da w_i (morfologia, correção ortográfica, sinônimo)
 - consultas longas (n é grande)
 - expressões booleanas geram poucos resultados.
- Booleana simples vs. Ordenação do conjunto de resultados
 - A busca booleana simples não gera ordenação dos documentos.
 - O Google (e a maioria dos buscadores booleanos bem feitos) ordena o conjunto de resultados – com os melhores resultados (de acordo com um estimador de relevância) no topo.

Documentos

- Última aula: Sistema de recuperação Booleana simples

Documentos

- Última aula: Sistema de recuperação Booleana simples
- Nossos pressupostos eram:

Documentos

- Última aula: Sistema de recuperação Booleana simples
- Nossos pressupostos eram:
 - Nós sabemos o que é um documento.

Documentos

- Última aula: Sistema de recuperação Booleana simples
- Nossos pressupostos eram:
 - Nós sabemos o que é um documento.
 - Documentos são “legíveis por máquina”.

Documentos

- Última aula: Sistema de recuperação Booleana simples
- Nossos pressupostos eram:
 - Nós sabemos o que é um documento.
 - Documentos são “legíveis por máquina”.
- Na prática isto pode ser bem complicado.

Parseando um Documento

- Precisamos lidar com formato e lingua de cada documento.

Parseando um Documento

- Precisamos lidar com formato e lingua de cada documento.
- Em que formato ele está? pdf, word, excel, html etc.

Parseando um Documento

- Precisamos lidar com formato e lingua de cada documento.
- Em que formato ele está? pdf, word, excel, html etc.
- Em que língua ele está?

Parseando um Documento

- Precisamos lidar com formato e lingua de cada documento.
- Em que formato ele está? pdf, word, excel, html etc.
- Em que língua ele está?
- Em que codificação está?

Parseando um Documento

- Precisamos lidar com formato e lingua de cada documento.
- Em que formato ele está? pdf, word, excel, html etc.
- Em que língua ele está?
- Em que codificação está?
- Cada uma destas perguntas implica um problema de classificação.

Parseando um Documento

- Precisamos lidar com formato e lingua de cada documento.
- Em que formato ele está? pdf, word, excel, html etc.
- Em que língua ele está?
- Em que codificação está?
- Cada uma destas perguntas implica um problema de classificação.
- Alternativa: usar heurísticas.

Complicações de formato e Língua

- Um mesmo índice normalmente contém termos de diversas línguas.

Complicações de formato e Língua

- Um mesmo índice normalmente contém termos de diversas línguas.
 - Algumas vezes um documento ou seus componentes contém múltiplas línguas e formatos.

Complicações de formato e Língua

- Um mesmo índice normalmente contém termos de diversas línguas.
 - Algumas vezes um documento ou seus componentes contém múltiplas línguas e formatos.
 - Email em francês com PDF em espanhol anexado

Complicações de formato e Língua

- Um mesmo índice normalmente contém termos de diversas línguas.
 - Algumas vezes um documento ou seus componentes contém múltiplas línguas e formatos.
 - Email em francês com PDF em espanhol anexado
- Qual a unidade do documento para fins de indexação?

Complicações de formato e Língua

- Um mesmo índice normalmente contém termos de diversas línguas.
 - Algumas vezes um documento ou seus componentes contém múltiplas línguas e formatos.
 - Email em francês com PDF em espanhol anexado
- Qual a unidade do documento para fins de indexação?
- Um arquivo?

Complicações de formato e Língua

- Um mesmo índice normalmente contém termos de diversas línguas.
 - Algumas vezes um documento ou seus componentes contém múltiplas línguas e formatos.
 - Email em francês com PDF em espanhol anexado
- Qual a unidade do documento para fins de indexação?
- Um arquivo?
- Um email?

Complicações de formato e Língua

- Um mesmo índice normalmente contém termos de diversas línguas.
 - Algumas vezes um documento ou seus componentes contém múltiplas línguas e formatos.
 - Email em francês com PDF em espanhol anexado
- Qual a unidade do documento para fins de indexação?
- Um arquivo?
- Um email?
- Um email com 5 anexos?

Complicações de formato e Língua

- Um mesmo índice normalmente contém termos de diversas línguas.
 - Algumas vezes um documento ou seus componentes contém múltiplas línguas e formatos.
 - Email em francês com PDF em espanhol anexado
- Qual a unidade do documento para fins de indexação?
- Um arquivo?
- Um email?
- Um email com 5 anexos?
- Um grupo de arquivos(ppt ou latex in HTML)?

Complicações de formato e Língua

- Um mesmo índice normalmente contém termos de diversas línguas.
 - Algumas vezes um documento ou seus componentes contém múltiplas línguas e formatos.
 - Email em francês com PDF em espanhol anexado
- Qual a unidade do documento para fins de indexação?
- Um arquivo?
- Um email?
- Um email com 5 anexos?
- Um grupo de arquivos(ppt ou latex in HTML)?
- Em suma: Responder à pergunta “O que é um documento?” não é trivial e requer algumas decisões.

Definições

- **Palavra** – Uma string delimitada de caracteres como aparece no texto.

Definições

- **Palavra** – Uma string delimitada de caracteres como aparece no texto.
- **Termo** – Uma palavra “normalizada” (capitalização, morfologia, ortografia, etc); Uma classe equivalente de palavras.

Definições

- **Palavra** – Uma string delimitada de caracteres como aparece no texto.
- **Termo** – Uma palavra “normalizada” (capitalização, morfologia, ortografia, etc); Uma classe equivalente de palavras.
- **Token** – Uma instancia de uma palavra ou termo ocorrendo em um documento.

Definições

- **Palavra** – Uma string delimitada de caracteres como aparece no texto.
- **Termo** – Uma palavra “normalizada” (capitalização, morfologia, ortografia, etc); Uma classe equivalente de palavras.
- **Token** – Uma instancia de uma palavra ou termo ocorrendo em um documento.
- **Tipo** – O mesmo que termo na maioria dos casos: Uma classe de equivalência de tokens.

Normalização

- Precisamos “normalizar” os termos no texto indexado, assim como o termos da consulta da mesma maneira.

Normalização

- Precisamos “normalizar” os termos no texto indexado, assim como o termos da consulta da mesma maneira.
- Exemplo: Queremos associar *U.S.A.* com *USA*

Normalização

- Precisamos “normalizar” os termos no texto indexado, assim como o termos da consulta da mesma maneira.
- Exemplo: Queremos associar *U.S.A.* com *USA*
- A maneira mais comum é criar classes de equivalência de termos.

Normalização

- Precisamos “normalizar” os termos no texto indexado, assim como o termos da consulta da mesma maneira.
- Exemplo: Queremos associar *U.S.A.* com *USA*
- A maneira mais comum é criar classes de equivalência de termos.
- Alternativamente: fazer expansões assimétricas

Normalização

- Precisamos “normalizar” os termos no texto indexado, assim como o termos da consulta da mesma maneira.
- Exemplo: Queremos associar *U.S.A.* com *USA*
- A maneira mais comum é criar classes de equivalência de termos.
- Alternativamente: fazer expansões assimétricas
 - window → window, windows

Normalização

- Precisamos “normalizar” os termos no texto indexado, assim como o termos da consulta da mesma maneira.
- Exemplo: Queremos associar *U.S.A.* com *USA*
- A maneira mais comum é criar classes de equivalência de termos.
- Alternativamente: fazer expansões assimétricas
 - window → window, windows
 - windows → Windows, windows

Normalização

- Precisamos “normalizar” os termos no texto indexado, assim como o termos da consulta da mesma maneira.
- Exemplo: Queremos associar *U.S.A.* com *USA*
- A maneira mais comum é criar classes de equivalência de termos.
- Alternativamente: fazer expansões assimétricas
 - window → window, windows
 - windows → Windows, windows
 - Windows (sem expansão)

Normalização

- Precisamos “normalizar” os termos no texto indexado, assim como o termos da consulta da mesma maneira.
- Exemplo: Queremos associar *U.S.A.* com *USA*
- A maneira mais comum é criar classes de equivalência de termos.
- Alternativamente: fazer expansões assimétricas
 - window → window, windows
 - windows → Windows, windows
 - Windows (sem expansão)
- Mais poderoso, porém menos eficiente

Normalização

- Precisamos “normalizar” os termos no texto indexado, assim como o termos da consulta da mesma maneira.
- Exemplo: Queremos associar *U.S.A.* com *USA*
- A maneira mais comum é criar classes de equivalência de termos.
- Alternativamente: fazer expansões assimétricas
 - window → window, windows
 - windows → Windows, windows
 - Windows (sem expansão)
- Mais poderoso, porém menos eficiente
- Porque não colocamos *window*, *Window*, *windows*, e *Windows* na mesma classe de equivalência?

Normalização: outras línguas

- Normalização e detecção de linguagem interagem.

Normalização: outras línguas

- Normalização e detecção de linguagem interagem.
- *PETER WILL NICHT MIT.* → MIT = mit

Normalização: outras línguas

- Normalização e detecção de linguagem interagem.
- *PETER WILL NICHT MIT.* → MIT = mit
- *He got his PhD from MIT.* → MIT \neq mit

Lembre-se: Construção de índices invertidos

- Entrada:

Friends, Romans, countrymen.

So let it be with Caesar ...

Lembre-se: Construção de índices invertidos

- Entrada:

Friends, Romans, countrymen.

So let it be with Caesar ...

- Saída:

friend

roman

countryman

so ...

Lembre-se: Construção de índices invertidos

- Entrada:

Friends, Romans, countrymen. So let it be with Caesar ...

- Saída:

friend roman countryman so ...

- Cada token é candidato a uma entrada de posting.

Lembre-se: Construção de índices invertidos

- Entrada:

Friends, Romans, countrymen. So let it be with Caesar ...

- Saída:

friend roman countryman so ...

- Cada token é candidato a uma entrada de posting.
- Quais são os tokens válidos?

Exercícios

In June, the dog likes to chase the cat in the barn. – Quantos tokens? Quantos tipos de Palavras?

Porque a tokenização é difícil – mesmo em inglês. **Tokenize:** *Mr. O'Neill thinks that the boys' stories about Chile's capital aren't amusing.*

Problemas de tokenização: Uma ou duas palavras? (ou várias)

- Hewlett-Packard

Problemas de tokenização: Uma ou duas palavras? (ou várias)

- Hewlett-Packard
- State-of-the-art

Problemas de tokenização: Uma ou duas palavras? (ou várias)

- Hewlett-Packard
- State-of-the-art
- co-education

Problemas de tokenização: Uma ou duas palavras? (ou várias)

- Hewlett-Packard
- State-of-the-art
- co-education
- the hold-him-back-and-drag-him-away maneuver

Problemas de tokenização: Uma ou duas palavras? (ou várias)

- Hewlett-Packard
- State-of-the-art
- co-education
- the hold-him-back-and-drag-him-away maneuver
- data base

Problemas de tokenização: Uma ou duas palavras? (ou várias)

- Hewlett-Packard
- State-of-the-art
- co-education
- the hold-him-back-and-drag-him-away maneuver
- data base
- San Francisco

Problemas de tokenização: Uma ou duas palavras? (ou várias)

- Hewlett-Packard
- State-of-the-art
- co-education
- the hold-him-back-and-drag-him-away maneuver
- data base
- San Francisco
- Los Angeles-based company

Problemas de tokenização: Uma ou duas palavras? (ou várias)

- Hewlett-Packard
- State-of-the-art
- co-education
- the hold-him-back-and-drag-him-away maneuver
- data base
- San Francisco
- Los Angeles-based company
- cheap San Francisco-Los Angeles fares

Problemas de tokenização: Uma ou duas palavras? (ou várias)

- Hewlett-Packard
- State-of-the-art
- co-education
- the hold-him-back-and-drag-him-away maneuver
- data base
- San Francisco
- Los Angeles-based company
- cheap San Francisco-Los Angeles fares
- York University vs. New York University

Números

- 3/20/91

Números

- 3/20/91
- 20/3/91

Números

- 3/20/91
- 20/3/91
- Mar 20, 1991

Números

- 3/20/91
- 20/3/91
- Mar 20, 1991
- B-52

Números

- 3/20/91
- 20/3/91
- Mar 20, 1991
- B-52
- 100.2.86.144

Números

- 3/20/91
- 20/3/91
- Mar 20, 1991
- B-52
- 100.2.86.144
- (800) 234-2333

Números

- 3/20/91
- 20/3/91
- Mar 20, 1991
- B-52
- 100.2.86.144
- (800) 234-2333
- 800.234.2333

Números

- 3/20/91
- 20/3/91
- Mar 20, 1991
- B-52
- 100.2.86.144
- (800) 234-2333
- 800.234.2333
- Sistemas mais antigos de RI costumavam não indexar números. . .

Números

- 3/20/91
- 20/3/91
- Mar 20, 1991
- B-52
- 100.2.86.144
- (800) 234-2333
- 800.234.2333
- Sistemas mais antigos de RI costumavam não indexar números. . .
- . . . mas as vezes isto pode ser útil.

Chinês: Ná há espaços

莎拉波娃现在居住在美国东南部的佛罗里达。今年4月9日，莎拉波娃在美国第一大城市纽约度过了18岁生日。生日派对上，莎拉波娃露出了甜美的微笑。

Segmentação ambígua em chinês

和尚

Os dois caracteres podem ser tratados como uma palavra significando ‘monge’ ou como uma sequência de duas palavras significando “e” e “parado”.

Outros casos de ausência de espaços

- Palavras compostas em Holandês, Alemão, Sueco in Dutch, German, Swedish

Outros casos de ausência de espaços

- Palavras compostas em Holandês, Alemão, Sueco in Dutch, German, Swedish
- Computerlinguistik → Computer + Linguistik

Outros casos de ausência de espaços

- Palavras compostas em Holandês, Alemão, Sueco in Dutch, German, Swedish
- Computerlinguistik → Computer + Linguistik
- Lebensversicherungsgesellschaftsangestellter

Outros casos de ausência de espaços

- Palavras compostas em Holandês, Alemão, Sueco in Dutch, German, Swedish
- Computerlinguistik → Computer + Linguistik
- Lebensversicherungsgesellschaftsangestellter
- → leben + versicherung + gesellschaft + angestellter

Outros casos de ausência de espaços

- Palavras compostas em Holandês, Alemão, Sueco in Dutch, German, Swedish
- Computerlinguistik → Computer + Linguistik
- Lebensversicherungsgesellschaftsangestellter
- → leben + versicherung + gesellschaft + angestellter
- Inuit: tusaatsiarunnangittualuujunga (Não consigo escutar direito.)

Outros casos de ausência de espaços

- Palavras compostas em Holandês, Alemão, Sueco in Dutch, German, Swedish
- Computerlinguistik → Computer + Linguistik
- Lebensversicherungsgesellschaftsangestellter
- → leben + versicherung + gesellschaft + angestellter
- Inuit: tusaatsiarunnangittualuujunga (Não consigo escutar direito.)
- Muitas outras línguas apresentam dificuldades de segmentação: Finlandês, Urdu, ...

Japonês

ノーベル平和賞を受賞したワンガリ・マータイさんが名誉会長を務めるMOTTA INAIキャンペーンの一環として、毎日新聞社とマガジンハウスは「私の、もったいない」を募集します。皆様が日ごろ「もったいない」と感じて実践していることや、それにまつわるエピソードを800字以内の文章にまとめ、簡単な写真、イラスト、図などを添えて10月20日までにお送りください。大賞受賞者には、50万円相当の旅行券とエコ製品2点の副賞が贈られます。

Japonês

ノーベル平和賞を受賞したワンガリ・マータイさんが名誉会長を務めるMOTTAI NA Iキャンペーンの一環として、毎日新聞社とマガジンハウスは「私の、もったいない」を募集します。皆様が日ごろ「もったいない」と感じて実践していることや、それにまつわるエピソードを800字以内の文章にまとめ、簡単な写真、イラスト、図などを添えて10月20日までにお送りください。大賞受賞者には、50万円相当の旅行券とエコ製品2点の副賞が贈られます。

4 “alfabetos” diferentes: Caracteres chineses, hiragana silabário para finalizações infleccionais e palavras funcionais, katakana silabário para transcrições de palavras estrangeiras. Sem espaços (como no chinês).

Japonês

ノーベル平和賞を受賞したワンガリ・マータイさんが名誉会長を務めるMOTTAINAIキャンペーンの一環として、毎日新聞社とマガジンハウスは「私の、もったいない」を募集します。皆様が日ごろ「もったいない」と感じて実践していることや、それにまつわるエピソードを800字以内の文章にまとめ、簡単な写真、イラスト、図などを添えて10月20日までにお送りください。大賞受賞者には、50万円相当の旅行券とエコ製品2点の副賞が贈られます。

4 “alfabetos” diferentes: Caracteres chineses, hiragana silabário para finalizações infleccionais e palavras funcionais, katakana silabário para transcrições de palavras estrangeiras. Sem espaços (como no chinês).

Um usuário pode fazer uma consulta inteiramente em hiragana!

Árabe

كِتَابٌ ← كِتَابٌ
un b ā t i k
/kitābun/ 'a book'

Árabe: Bidirecionalidade

استقلت الجزائر في سنة 1962 بعد 132 عاما من الاحتلال الفرنسي.

← → ← →

← START

‘Algeria achieved its independence in 1962 after 132 years of French occupation.’

Árabe: Bidirecionalidade

استقلت الجزائر في سنة 1962 بعد 132 عاما من الاحتلال الفرنسي.

← → ← →

← START

‘Algeria achieved its independence in 1962 after 132 years of French occupation.’

Bidirecionalidade não é um problema se o texto estiver codificado em Unicode.

Acentos e diacríticos

- Acentos: Árvore vs. Arvore (Simples omissão de acento)

Acentos e diacríticos

- Acentos: Árvore vs. Arvore (Simples omissão de acento)
- Tremas: Universität vs. Universitaet (substituição por sequencia especial de caracteres: “ae”)

Acentos e diacríticos

- Acentos: Árvore vs. Arvore (Simples omissão de acento)
- Tremas: Universität vs. Universitaet (substituição por sequencia especial de caracteres: "ae")
- Critério mais importante: Qual a forma mais comumente adotada por usuários em suas consultas?

Acentos e diacríticos

- Acentos: Árvore vs. Arvore (Simples omissão de acento)
- Tremas: Universität vs. Universitaet (substituição por sequencia especial de caracteres: "ae")
- Critério mais importante: Qual a forma mais comumente adotada por usuários em suas consultas?
- Mesmo em línguas que possuem acentos por padrão, frequentemente os usuários não os digitam.

Capitalização

- Reduzir todas as letras para minúsculas

Capitalização

- Reduzir todas as letras para minúsculas
- Possíveis exceções: Palavras capitalizadas no meio da frase

Capitalização

- Reduzir todas as letras para minúsculas
- Possíveis exceções: Palavras capitalizadas no meio da frase
- MIT vs. mit

Capitalização

- Reduzir todas as letras para minúsculas
- Possíveis exceções: Palavras capitalizadas no meio da frase
- MIT vs. mit
- Fed vs. fed

Capitalização

- Reduzir todas as letras para minúsculas
- Possíveis exceções: Palavras capitalizadas no meio da frase
- MIT vs. mit
- Fed vs. fed
- Frequentemente é melhor converter tudo para minúsculas pois os usuários normalmente fazem suas consultas sem utilizar a capitalização correta.

“Stop words”

- “stop words” = Palavras extremamente comuns que podem parecer ser de pouco valor na seleção de documentos correspondentes à necessidade do usuário.

“Stop words”

- “stop words” = Palavras extremamente comuns que podem parecer ser de pouco valor na seleção de documentos correspondentes à necessidade do usuário.
- Exemplos: artigos, preposições, etc.

“Stop words”

- “stop words” = Palavras extremamente comuns que podem parecer ser de pouco valor na seleção de documentos correspondentes à necessidade do usuário.
- Exemplos: artigos, preposições, etc.
- Eliminação de “Stop word” costumava ser padrão em sistemas de RI mais antigos.

“Stop words”

- “stop words” = Palavras extremamente comuns que podem parecer ser de pouco valor na seleção de documentos correspondentes à necessidade do usuário.
- Exemplos: artigos, preposições, etc.
- Eliminação de “Stop word” costumava ser padrão em sistemas de RI mais antigos.
- Mas estas palavras são necessárias para consultas frasais, p.ex. “Reino da Dinamarca”

“Stop words”

- “stop words” = Palavras extremamente comuns que podem parecer ser de pouco valor na seleção de documentos correspondentes à necessidade do usuário.
- Exemplos: artigos, preposições, etc.
- Eliminação de “Stop word” costumava ser padrão em sistemas de RI mais antigos.
- Mas estas palavras são necessárias para consultas frasais, p.ex. “Reino da Dinamarca”
- A maioria dos engines de busca na web, indexa “stop words”.

Mais classes de equivalência

- Soundex: equivalencia fonética, Muller = Mueller)

Mais classes de equivalência

- Soundex: equivalencia fonética, Muller = Mueller)
- Thesauri: equivalencia semântica, carro = automóvel)

Lematização

- Reduz variações e inflexões a uma forma básica

Lematização

- Reduz variações e inflexões a uma forma básica
- Exemplo: *sou, são, é* → *ser*

Lematização

- Reduz variações e inflexões a uma forma básica
- Exemplo: *sou, são, é* → *ser*
- Exemplo: *carro, carros* → *carros*

Lematização

- Reduz variações e inflexões a uma forma básica
- Exemplo: *sou, são, é* → *ser*
- Exemplo: *carro, carros* → *carros*
- Exemple: *Os carros do rapaz são de diferentes cores* → *o carro do rapaz ser de diferente cor*

Lematização

- Reduz variações e inflexões a uma forma básica
- Exemplo: *sou, são, é* → *ser*
- Exemplo: *carro, carros* → *carros*
- Exemple: *Os carros do rapaz são de diferentes cores* → *o carro do rapaz ser de diferente cor*
- Lematização implica fazer a redução “correta” às forma mais fundamental segundo o dicionário o [lemma](#)).

Lematização

- Reduz variações e inflexões a uma forma básica
- Exemplo: *sou, são, é* → *ser*
- Exemplo: *carro, carros* → *carros*
- Exemple: *Os carros do rapaz são de diferentes cores* → *o carro do rapaz ser de diferente cor*
- Lematização implica fazer a redução “correta” às forma mais fundamental segundo o dicionário o [lemma](#)).
- Morfologia inflexional (*cortando* → *cortar*) vs. morfologia derivacional (*destruição* → *destruir*)

Truncagem (“Stemming”)

- Definição de truncagem: Heurística bruta que **corta fora as terminações das palavras** na esperança de alcançar o mesmo objetivo da lematização sem a necessidade de conhecimento linguístico.

Truncagem (“Stemming”)

- Definição de truncagem: Heurística bruta que **corta fora as terminações das palavras** na esperança de alcançar o mesmo objetivo da lematização sem a necessidade de conhecimento linguístico.
- Depende da língua

Truncagem (“Stemming”)

- Definição de truncagem: Heurística bruta que **corta fora as terminações das palavras** na esperança de alcançar o mesmo objetivo da lematização sem a necessidade de conhecimento linguístico.
- Depende da língua
- Frequentemente inflexional e derivacional

Truncagem (“Stemming”)

- Definição de truncagem: Heurística bruta que **corta fora as terminações das palavras** na esperança de alcançar o mesmo objetivo da lematização sem a necessidade de conhecimento linguístico.
- Depende da língua
- Frequentemente inflexional e derivacional
- Exemplo de truncagem derivacional: *automatizar*, *automático*, *automatização* se reduzem a *automat*

Algoritmo de Porter

- Algoritmo mais comum para truncagem do inglês

Algoritmo de Porter

- Algoritmo mais comum para truncagem do inglês
- Resultados sugerem que é tão bom quanto outras opções de truncagem

Algoritmo de Porter

- Algoritmo mais comum para truncagem do inglês
- Resultados sugerem que é tão bom quanto outras opções de truncagem
- Convenções + 5 fases de reduções

Algoritmo de Porter

- Algoritmo mais comum para truncagem do inglês
- Resultados sugerem que é tão bom quanto outras opções de truncagem
- Convenções + 5 fases de reduções
- Fases são aplicadas sequencialmente

Algoritmo de Porter

- Algoritmo mais comum para truncagem do inglês
- Resultados sugerem que é tão bom quanto outras opções de truncagem
- Convenções + 5 fases de reduções
- Fases são aplicadas sequencialmente
- Cada fase consiste de um conjunto de comandos.

Algoritmo de Porter

- Algoritmo mais comum para truncagem do inglês
- Resultados sugerem que é tão bom quanto outras opções de truncagem
- Convenções + 5 fases de reduções
- Fases são aplicadas sequencialmente
- Cada fase consiste de um conjunto de comandos.
 - Exemplo de comando: Deleta o *ement* final se o que sobrar for mais longo do que um caractere

Algoritmo de Porter

- Algoritmo mais comum para truncagem do inglês
- Resultados sugerem que é tão bom quanto outras opções de truncagem
- Convenções + 5 fases de reduções
- Fases são aplicadas sequencialmente
- Cada fase consiste de um conjunto de comandos.
 - Exemplo de comando: Deleta o *ement* final se o que sobrar for mais longo do que um caractere
 - replacement → replac

Algoritmo de Porter

- Algoritmo mais comum para truncagem do inglês
- Resultados sugerem que é tão bom quanto outras opções de truncagem
- Convenções + 5 fases de reduções
- Fases são aplicadas sequencialmente
- Cada fase consiste de um conjunto de comandos.
 - Exemplo de comando: Deleta o *ement* final se o que sobrar for mais longo do que um caractere
 - replacement → replac
 - cement → cement

Algoritmo de Porter

- Algoritmo mais comum para truncagem do inglês
- Resultados sugerem que é tão bom quanto outras opções de truncagem
- Convenções + 5 fases de reduções
- Fases são aplicadas sequencialmente
- Cada fase consiste de um conjunto de comandos.
 - Exemplo de comando: Deleta o *ement* final se o que sobrar for mais longo do que um caractere
 - replacement → replac
 - cement → cement
- Exemplo de convenção: Das regras em comando composto, selecione aquela que se aplica ao sufixo mais longo.

Trucador Porter: Algumas regras

Regra

SSES → SS

IES → I

SS → SS

S →

Exemplo

caresses → caress

ponies → poni

caress → caress

cats → cat

Três Truncadores: Uma comparação

Texto exemplo: Such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

Truncator Porter: such an analysis can reveal features that are not easily visible from the variation in the individual gene and can lead to a picture of expression that is more biologically transparent and accessible to interpretation.

Truncador Lovins: such an analysis can reveal features that are not easily visible from the variation in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

Truncador Paice: such an analysis can reveal features that are not easily visible from the variation in the individual genome and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

Truncagem melhora efetividade?

- Em geral, a truncagem aumenta a efetividade para algumas consultas e diminui para outras.

Truncagem melhora efetividade?

- Em geral, a truncagem aumenta a efetividade para algumas consultas e diminui para outras.
- Consultas em que a truncagem pode ajudar: [tartan sweaters], [sightseeing tour san francisco]

Truncagem melhora efetividade?

- Em geral, a truncagem aumenta a efetividade para algumas consultas e diminui para outras.
- Consultas em que a truncagem pode ajudar: [tartan sweaters], [sightseeing tour san francisco]
- (classes de equivalência: {sweater,sweaters}, {tour,tours})

Truncagem melhora efetividade?

- Em geral, a truncagem aumenta a efetividade para algumas consultas e diminui para outras.
- Consultas em que a truncagem pode ajudar: [tartan sweaters], [sightseeing tour san francisco]
- (classes de equivalência: {sweater,sweaters}, {tour,tours})
- Classes de equivalência para *oper* do truncador Porter contém todas estas formas: *operate operating operates operation operative operatives operational*.

Truncagem melhora efetividade?

- Em geral, a truncagem aumenta a efetividade para algumas consultas e diminui para outras.
- Consultas em que a truncagem pode ajudar: [tartan sweaters], [sightseeing tour san francisco]
- (classes de equivalência: {sweater,sweaters}, {tour,tours})
- Classes de equivalência para *oper* do truncador Porter contém todas estas formas: *operate operating operates operation operative operatives operational*.
- Consultas em que a truncagem prejudica: [operational AND research], [operating AND system], [operative AND dentistry]

Exercise: O que o Google faz?

- “Stop words”
- Normalização
- Tokenização
- minúsculização
- Truncagem
- Alfabetos não latinos
- Tremas
- Palavras compostas
- Números

Consultas de frases

- Queremos responder a uma consulta como [stanford university] – como uma frase.

Consultas de frases

- Queremos responder a uma consulta como [stanford university] – como uma frase.
- Logo *The inventor Stanford Ovshinsky never went to university* não deve ser uma resposta.

Consultas de frases

- Queremos responder a uma consulta como [stanford university] – como uma frase.
- Logo *The inventor Stanford Ovshinsky never went to university* não deve ser uma resposta.
- O conceito de consulta de frase é facilmente compreensível para usuários.

Consultas de frases

- Queremos responder a uma consulta como [stanford university] – como uma frase.
- Logo *The inventor Stanford Ovshinsky never went to university* não deve ser uma resposta.
- O conceito de consulta de frase é facilmente compreensível para usuários.
- Cerca de 10% das consultas na web são consultas de frases.

Consultas de frases

- Queremos responder a uma consulta como [stanford university] – como uma frase.
- Logo *The inventor Stanford Ovshinsky never went to university* não deve ser uma resposta.
- O conceito de consulta de frase é facilmente compreensível para usuários.
- Cerca de 10% das consultas na web são consultas de frases.
- CONsequências para o índice invertido: Não é mais suficiente guardar docIDs em listas de postings.

Consultas de frases

- Queremos responder a uma consulta como [stanford university] – como uma frase.
- Logo *The inventor Stanford Ovshinsky never went to university* não deve ser uma resposta.
- O conceito de consulta de frase é facilmente compreensível para usuários.
- Cerca de 10% das consultas na web são consultas de frases.
- CONsequências para o índice invertido: Não é mais suficiente guardar docIDs em listas de postings.
- Duas maneiras de estender os índice invertido:

Consultas de frases

- Queremos responder a uma consulta como [stanford university] – como uma frase.
- Logo *The inventor Stanford Ovshinsky never went to university* não deve ser uma resposta.
- O conceito de consulta de frase é facilmente compreensível para usuários.
- Cerca de 10% das consultas na web são consultas de frases.
- Consequências para o índice invertido: Não é mais suficiente guardar docIDs em listas de postings.
- Duas maneiras de estender os índice invertido:
 - índice de bigramas

Consultas de frases

- Queremos responder a uma consulta como [stanford university] – como uma frase.
- Logo *The inventor Stanford Ovshinsky never went to university* não deve ser uma resposta.
- O conceito de consulta de frase é facilmente compreensível para usuários.
- Cerca de 10% das consultas na web são consultas de frases.
- Consequências para o índice invertido: Não é mais suficiente guardar docIDs em listas de postings.
- Duas maneiras de estender os índice invertido:
 - índice de bigramas
 - Índice posicional

Índices de bigramas

- Indexa cada par consecutivo de termos no texto como uma frase.

Índices de bigramas

- Indexa cada par consecutivo de termos no texto como uma frase.
- Por exemplo, *Friends, Romans, Countrymen* vai gerar dois bigramas: “*friends romans*” e “*romans countrymen*”

Índices de bigramas

- Indexa cada par consecutivo de termos no texto como uma frase.
- Por exemplo, *Friends, Romans, Countrymen* vai gerar dois bigramas: “*friends romans*” e “*romans countrymen*”
- Cada um destes bigramas é agora um termo do vocabulário.

Índices de bigramas

- Indexa cada par consecutivo de termos no texto como uma frase.
- Por exemplo, *Friends, Romans, Countrymen* vai gerar dois bigramas: “*friends romans*” e “*romans countrymen*”
- Cada um destes bigramas é agora um termo do vocabulário.
- Agora consultas de frases de duas palavras podem ser facilmente respondidas.

Consultas de frases mais longas

- Uma frase longa como “*stanford university palo alto*” pode ser representada como a consulta Booleana “STANFORD UNIVERSITY” AND “UNIVERSITY PALO” AND “PALO ALTO”

Consultas de frases mais longas

- Uma frase longa como “*stanford university palo alto*” pode ser representada como a consulta Booleana “STANFORD UNIVERSITY” AND “UNIVERSITY PALO” AND “PALO ALTO”
- É necessário pós-filtrar o conjunto de respostas para identificar o subconjunto que contém a frase de quatro palavras.

Problemas com índices de bigramas

- Porque estes índices são raramente usados?

Problemas com índices de bigramas

- Porque estes índices são raramente usados?
- Falsos Positivos, como indicado acima

Problemas com índices de bigramas

- Porque estes índices são raramente usados?
- Falsos Positivos, como indicado acima
- Índice explode devido ao tamanho exagerado do vocabulário

Índices posicionais

Índices posicionais

- Índices posicionais são mais eficientes que índices de bigramas.

Índices posicionais

- Índices posicionais são mais eficientes que índices de bigramas.
- Listas de Postings em um índice **não posicional**: Cada posting é apenas um docID

Índices posicionais

- Índices posicionais são mais eficientes que índices de bigramas.
- Listas de Postings em um índice **não posicional**: Cada posting é apenas um docID
- Listas de Postings em um índice **posicional**: Cada posting é um docID e **uma lista de posições**

Índices posicionais: Exemplo

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

⟨ 1: ⟨7, 18, 33, 72, 86, 231⟩;
2: ⟨1, 17, 74, 222, 255⟩;
4: ⟨8, 16, 190, 429, 433⟩;
5: ⟨363, 367⟩;
7: ⟨13, 23, 191⟩; ... ⟩

BE, 178239:

⟨ 1: ⟨17, 25⟩;
4: ⟨17, 191, 291, 430, 434⟩;
5: ⟨14, 19, 101⟩; ... ⟩

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

⟨ 1: ⟨7, 18, 33, 72, 86, 231⟩;
2: ⟨1, 17, 74, 222, 255⟩;
4: ⟨8, 16, 190, 429, 433⟩;
5: ⟨363, 367⟩;
7: ⟨13, 23, 191⟩; ... ⟩

BE, 178239:

⟨ 1: ⟨17, 25⟩;
4: ⟨17, 191, 291, 430, 434⟩;
5: ⟨14, 19, 101⟩; ... ⟩

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

⟨ 1: ⟨7, 18, 33, 72, 86, 231⟩;
2: ⟨1, 17, 74, 222, 255⟩;
4: ⟨8, 16, 190, 429, 433⟩;
5: ⟨363, 367⟩;
7: ⟨13, 23, 191⟩; ... ⟩

BE, 178239:

⟨ 1: ⟨17, 25⟩;
4: ⟨17, 191, 291, 430, 434⟩;
5: ⟨14, 19, 101⟩; ... ⟩

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- ⟨ 1: ⟨7, 18, 33, 72, 86, 231⟩;
- 2: ⟨1, 17, 74, 222, 255⟩;
- 4: ⟨8, 16, 190, 429, 433⟩;
- 5: ⟨363, 367⟩;
- 7: ⟨13, 23, 191⟩; ... ⟩

BE, 178239:

- ⟨ 1: ⟨17, 25⟩;
- 4: ⟨17, 191, 291, 430, 434⟩;
- 5: ⟨14, 19, 101⟩; ... ⟩

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;

2: $\langle 1, 17, 74, 222, 255 \rangle$;

4: $\langle 8, 16, 190, 429, 433 \rangle$;

5: $\langle 363, 367 \rangle$;

7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

1: $\langle 17, 25 \rangle$;

4: $\langle 17, 191, 291, 430, 434 \rangle$;

5: $\langle 14, 19, 101 \rangle$; ...

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- 1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
- 2: $\langle 1, 17, 74, 222, 255 \rangle$;
- 4: $\langle 8, 16, 190, 429, 433 \rangle$;
- 5: $\langle 363, 367 \rangle$;
- 7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

- 1: $\langle 17, 25 \rangle$;
- 4: $\langle 17, 191, 291, 430, 434 \rangle$;
- 5: $\langle 14, 19, 101 \rangle$; ...

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- 1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
- 2: $\langle 1, 17, 74, 222, 255 \rangle$;
- 4: $\langle 8, 16, 190, 429, 433 \rangle$;
- 5: $\langle 363, 367 \rangle$;
- 7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

- 1: $\langle 17, 25 \rangle$;
- 4: $\langle 17, 191, 291, 430, 434 \rangle$;
- 5: $\langle 14, 19, 101 \rangle$; ...

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- 1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
- 2: $\langle 1, 17, 74, 222, 255 \rangle$;
- 4: $\langle 8, 16, 190, 429, 433 \rangle$;
- 5: $\langle 363, 367 \rangle$;
- 7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

- 1: $\langle 17, 25 \rangle$;
- 4: $\langle 17, 191, 291, 430, 434 \rangle$;
- 5: $\langle 14, 19, 101 \rangle$; ...

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- 1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
- 2: $\langle 1, 17, 74, 222, 255 \rangle$;
- 4: $\langle 8, 16, 190, 429, 433 \rangle$;
- 5: $\langle 363, 367 \rangle$;
- 7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

- 1: $\langle 17, 25 \rangle$;
- 4: $\langle 17, 191, 291, 430, 434 \rangle$;
- 5: $\langle 14, 19, 101 \rangle$; ...

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- 1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
- 2: $\langle 1, 17, 74, 222, 255 \rangle$;
- 4: $\langle 8, 16, 190, 429, 433 \rangle$;
- 5: $\langle 363, 367 \rangle$;
- 7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

- 1: $\langle 17, 25 \rangle$;
- 4: $\langle 17, 191, 291, 430, 434 \rangle$;
- 5: $\langle 14, 19, 101 \rangle$; ...

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- 1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
- 2: $\langle 1, 17, 74, 222, 255 \rangle$;
- 4: $\langle 8, 16, 190, 429, 433 \rangle$;
- 5: $\langle 363, 367 \rangle$;
- 7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

- 1: $\langle 17, 25 \rangle$;
- 4: $\langle 17, 191, 291, 430, 434 \rangle$;
- 5: $\langle 14, 19, 101 \rangle$; ...

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- 1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
- 2: $\langle 1, 17, 74, 222, 255 \rangle$;
- 4: $\langle 8, 16, 190, 429, 433 \rangle$;
- 5: $\langle 363, 367 \rangle$;
- 7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

- 1: $\langle 17, 25 \rangle$;
- 4: $\langle 17, 191, 291, 430, 434 \rangle$;
- 5: $\langle 14, 19, 101 \rangle$; ...

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- 1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
- 2: $\langle 1, 17, 74, 222, 255 \rangle$;
- 4: $\langle 8, 16, 190, 429, 433 \rangle$;
- 5: $\langle 363, 367 \rangle$;
- 7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

- 1: $\langle 17, 25 \rangle$;
- 4: $\langle 17, 191, 291, 430, 434 \rangle$;
- 5: $\langle 14, 19, 101 \rangle$; ...

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- 1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
- 2: $\langle 1, 17, 74, 222, 255 \rangle$;
- 4: $\langle 8, 16, 190, 429, 433 \rangle$;
- 5: $\langle 363, 367 \rangle$;
- 7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

- 1: $\langle 17, 25 \rangle$;
- 4: $\langle 17, 191, 291, 430, 434 \rangle$;
- 5: $\langle 14, 19, 101 \rangle$; ...

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- 1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
- 2: $\langle 1, 17, 74, 222, 255 \rangle$;
- 4: $\langle 8, 16, 190, 429, 433 \rangle$;
- 5: $\langle 363, 367 \rangle$;
- 7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

- 1: $\langle 17, 25 \rangle$;
- 4: $\langle 17, 191, 291, 430, 434 \rangle$;
- 5: $\langle 14, 19, 101 \rangle$; ...

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- 1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
- 2: $\langle 1, 17, 74, 222, 255 \rangle$;
- 4: $\langle 8, 16, 190, 429, 433 \rangle$;
- 5: $\langle 363, 367 \rangle$;
- 7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

- 1: $\langle 17, 25 \rangle$;
- 4: $\langle 17, 191, 291, 430, 434 \rangle$;
- 5: $\langle 14, 19, 101 \rangle$; ...

Índices posicionais: Exemplo

Query: *"to₁ be₂ or₃ not₄ to₅ be₆"*

TO, 993427:

- 1: $\langle 7, 18, 33, 72, 86, 231 \rangle$;
- 2: $\langle 1, 17, 74, 222, 255 \rangle$;
- 4: $\langle 8, 16, 190, 429, 433 \rangle$;
- 5: $\langle 363, 367 \rangle$;
- 7: $\langle 13, 23, 191 \rangle$; ...

BE, 178239:

- 1: $\langle 17, 25 \rangle$;
- 4: $\langle 17, 191, 291, 430, 434 \rangle$;
- 5: $\langle 14, 19, 101 \rangle$; ...

Documento 4 é a resposta!

Busca por proximidade

- Acabamos de ver como usar um índice posicional para busca de frases.

Busca por proximidade

- Acabamos de ver como usar um índice posicional para busca de frases.
- Também podemos usá-los para busca por proximidade.

Busca por proximidade

- Acabamos de ver como usar um índice posicional para busca de frases.
- Também podemos usá-los para busca por proximidade.
- Por exemplo: employment /4 place

Busca por proximidade

- Acabamos de ver como usar um índice posicional para busca de frases.
- Também podemos usá-los para busca por proximidade.
- Por exemplo: employment /4 place
- Encontrar todos os documentos que contêm EMPLOYMENT e PLACE a até 4 palavras words de distância entre si.

Busca por proximidade

- Acabamos de ver como usar um índice posicional para busca de frases.
- Também podemos usá-los para busca por proximidade.
- Por exemplo: employment /4 place
- Encontrar todos os documentos que contêm EMPLOYMENT e PLACE a até 4 palavras words de distância entre si.
- *Employment agencies that place healthcare workers are seeing growth* é uma resposta.

Busca por proximidade

- Acabamos de ver como usar um índice posicional para busca de frases.
- Também podemos usá-los para busca por proximidade.
- Por exemplo: `employment /4 place`
- Encontrar todos os documentos que contêm `EMPLOYMENT` e `PLACE` a até 4 palavras words de distância entre si.
- *Employment agencies that place healthcare workers are seeing growth* é uma resposta.
- *Employment agencies that have learned to adapt now place healthcare workers* não é uma resposta.

Busca por proximidade

- Use o índice posicional

Busca por proximidade

- Use o índice posicional
- Algoritmo mais simples: Olhe para o Produto vetorial das posições de (i) EMPLOYMENT no documento e (ii) PLACE no documento

Busca por proximidade

- Use o índice posicional
- Algoritmo mais simples: Olhe para o Produto vetorial das posições de (i) EMPLOYMENT no documento e (ii) PLACE no documento
- Muito ineficiente para palavras frequentes, especialmente “stop words”

Busca por proximidade

- Use o índice posicional
- Algoritmo mais simples: Olhe para o Produto vetorial das posições de (i) EMPLOYMENT no documento e (ii) PLACE no documento
- Muito ineficiente para palavras frequentes, especialmente “stop words”
- Note que queremos retornar as posições encontradas e não apenas uma lista de documentos.

Busca por proximidade

- Use o índice posicional
- Algoritmo mais simples: Olhe para o Produto vetorial das posições de (i) EMPLOYMENT no documento e (ii) PLACE no documento
- Muito ineficiente para palavras frequentes, especialmente “stop words”
- Note que queremos retornar as posições encontradas e não apenas uma lista de documentos.
- Isto é importante para sumarizações dinâmicas etc.

Interseção de “proximidade”

POSITIONALINTERSECT(p_1, p_2, k)

```

1  answer  $\leftarrow \langle \rangle$ 
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4      then  $I \leftarrow \langle \rangle$ 
5           $pp_1 \leftarrow \text{positions}(p_1)$ 
6           $pp_2 \leftarrow \text{positions}(p_2)$ 
7          while  $pp_1 \neq \text{NIL}$ 
8          do while  $pp_2 \neq \text{NIL}$ 
9              do if  $|\text{pos}(pp_1) - \text{pos}(pp_2)| \leq k$ 
10                 then  $\text{ADD}(I, \text{pos}(pp_2))$ 
11                 else if  $\text{pos}(pp_2) > \text{pos}(pp_1)$ 
12                     then break
13                  $pp_2 \leftarrow \text{next}(pp_2)$ 
14                 while  $I \neq \langle \rangle$  and  $|I[0] - \text{pos}(pp_1)| > k$ 
15                     do  $\text{DELETE}(I[0])$ 
16                 for each  $ps \in I$ 
17                     do  $\text{ADD}(\text{answer}, \langle \text{docID}(p_1), \text{pos}(pp_1), ps \rangle)$ 
18                  $pp_1 \leftarrow \text{next}(pp_1)$ 
19              $p_1 \leftarrow \text{next}(p_1)$ 
20              $p_2 \leftarrow \text{next}(p_2)$ 
21         else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
22             then  $p_1 \leftarrow \text{next}(p_1)$ 
23             else  $p_2 \leftarrow \text{next}(p_2)$ 
24 return answer

```


Esquema Combinado

- Índices de bigramas e posicionais podem ser combinados proveitosamente.

Esquema Combinado

- Índices de bigramas e posicionais podem ser combinados proveitosamente.
- Muitos bigramas são extremamente frequentes: Michael Jackson, Britney Spears, etc.

Esquema Combinado

- Índices de bigramas e posicionais podem ser combinados proveitosamente.
- Muitos bigramas são extremamente frequentes: Michael Jackson, Britney Spears, etc.
- Para estes bigramas, o aumento de velocidade proporcionado pela indexação de bigramas, é substancial.

Esquema Combinado

- Índices de bigramas e posicionais podem ser combinados proveitosamente.
- Muitos bigramas são extremamente frequentes: Michael Jackson, Britney Spears, etc.
- Para estes bigramas, o aumento de velocidade proporcionado pela indexação de bigramas, é substancial.
- Esquema de combinação: Incluir bigramas frequentes como vocabulário no índice. Todas as outras frases são recuperadas por interseção posicional.

Consultas “Posicionais” no Google

- Para buscas na Web, consultas posicionais são muito mais caras do que consultas Booleanas regulares.

Consultas “Posicionais” no Google

- Para buscas na Web, consultas posicionais são muito mais caras do que consultas Booleanas regulares.
- Vejamos o exemplo de consultas de frases.

Consultas “Posicionais” no Google

- Para buscas na Web, consultas posicionais são muito mais caras do que consultas Booleanas regulares.
- Vejamos o exemplo de consultas de frases.
- Porque são mais custosas?

Consultas “Posicionais” no Google

- Para buscas na Web, consultas posicionais são muito mais caras do que consultas Booleanas regulares.
- Vejamos o exemplo de consultas de frases.
- Porque são mais custosas?
- Você consegue demonstrar no Google que as consultas de frases são mais custosas que consultas Booleanas?