

Introdução à Recuperação de Informações

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

IRI 3: Dicionários e Recuperação Tolerante

Flávio Codeço Coelho

Escola de Matemática Aplicada, Fundação Getúlio Vargas

Sumário

- 1 Recapitulação
- 2 Dicionários
- 3 Consultas Coringa
- 4 Distância de Edição

Distinguindo entre Tipo e token

- **Token** – Uma instância de uma palavra ou termo ocorrendo em um documento
- **Tipo** – Uma classe de equivalência de tokens
- *In June, the dog likes to chase the cat in the barn.*
- 12 tokens, 9 tipos de palavras

Problemas na tokenização

- Quais são os delimitadores? Espaços? Apóstrofes? Hífen?
- Para cada um destes: às vezes eles delimitam, às vezes não.
- Muitas línguas não possuem espaços! (P.ex., Chinês)
- Não Há espaços em palavras compostas em Holandês, Alemão e Sueco (*Lebensversicherungsgesellschaftsangestellter*)

Problemas com classes de equivalência

- Um termo é uma classe de equivalência de tokens.
- Como definir Classes de equivalência?
- Números: (3/20/91 vs. 20/3/91)
- Capitalização
- Truncagem, Truncador de Porter
- Análise Morfológica : infleccional vs. derivacional
- Problemas de classes de equivalências em outras línguas
 - Morfologias mais complexas do que o inglês
 - Finlandês: Um único verbo pode ter 12000 formas diferentes
 - Acentos, tremas, etc.

Principais conclusões de hoje

Principais conclusões de hoje

- Recuperação Tolerante: O que fazer se não há correspondência exata entre o termo de consulta e os termos do documento.

Principais conclusões de hoje

- Recuperação Tolerante: O que fazer se não há correspondência exata entre o termo de consulta e os termos do documento.
- Consultas coringas

Principais conclusões de hoje

- Recuperação Tolerante: O que fazer se não há correspondência exata entre o termo de consulta e os termos do documento.
- Consultas coringas
- Correção ortográficas

Índice invertido

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

| | | | | | | | | | |
|--------|---|---|---|---|----|----|----|-----|-----|
| BRUTUS | → | 1 | 2 | 4 | 11 | 31 | 45 | 173 | 174 |
|--------|---|---|---|---|----|----|----|-----|-----|

| | | | | | | | | | | |
|--------|---|---|---|---|---|---|----|----|-----|-----|
| CAESAR | → | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 | ... |
|--------|---|---|---|---|---|---|----|----|-----|-----|

| | | | | | |
|-----------|---|---|----|----|-----|
| CALPURNIA | → | 2 | 31 | 54 | 101 |
|-----------|---|---|----|----|-----|

⋮

dicionário

postings

Índice invertido

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

| | | | | | | | | | |
|--------|---|---|---|---|----|----|----|-----|-----|
| BRUTUS | → | 1 | 2 | 4 | 11 | 31 | 45 | 173 | 174 |
|--------|---|---|---|---|----|----|----|-----|-----|

| | | | | | | | | | | |
|--------|---|---|---|---|---|---|----|----|-----|-----|
| CAESAR | → | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 | ... |
|--------|---|---|---|---|---|---|----|----|-----|-----|

| | | | | | |
|-----------|---|---|----|----|-----|
| CALPURNIA | → | 2 | 31 | 54 | 101 |
|-----------|---|---|----|----|-----|

⋮

dicionário

postings

Dicionários

- O dicionário é a estrutura de dados que é usada para armazenar o vocabulário de termos.

Dicionários

- O dicionário é a estrutura de dados que é usada para armazenar o vocabulário de termos.
- **vocabulário de termos**: os **dados**

Dicionários

- O dicionário é a estrutura de dados que é usada para armazenar o vocabulário de termos.
- **vocabulário de termos**: os **dados**
- **Dicionário**: A **estrutura de dados** para armazenamento do vocabulário

Dicionário como uma matriz com entradas de tamanho fixo

- Para cada termo, precisamos armazenar um par de ítems:

Dicionário como uma matriz com entradas de tamanho fixo

- Para cada termo, precisamos armazenar um par de itens:
 - Frequência de documentos

Dicionário como uma matriz com entradas de tamanho fixo

- Para cada termo, precisamos armazenar um par de itens:
 - Frequência de documentos
 - Ponteiro para a lista de postings

Dicionário como uma matriz com entradas de tamanho fixo

- Para cada termo, precisamos armazenar um par de itens:
 - Frequência de documentos
 - Ponteiro para a lista de postings
 - ...

Dicionário como uma matriz com entradas de tamanho fixo

- Para cada termo, precisamos armazenar um par de itens:
 - Frequência de documentos
 - Ponteiro para a lista de postings
 - ...
- Assuma por ora que podemos armazenar esta informação em uma entrada de tamanho fixo.

Dicionário como uma matriz com entradas de tamanho fixo

- Para cada termo, precisamos armazenar um par de itens:
 - Frequência de documentos
 - Ponteiro para a lista de postings
 - ...
- Assuma por ora que podemos armazenar esta informação em uma entrada de tamanho fixo.
- Assuma que armazenamos estas entradas em uma matriz.

Dicionário como uma matriz com entradas de tamanho fixo

| termo | documento frequência | ponteiro para lista de postings |
|--------|-------------------------|---------------------------------------|
| a | 656.265 | → |
| aachen | 65 | → |
| ... | ... | ... |
| zulu | 221 | → |

espaço necessário: 20 bytes 4 bytes 4 bytes

como acessamos um termo de consulta q_i nesta matriz em tempo de consulta? Ou seja: Que estrutura de dados usamos para localizar a entrada (linha) na matriz onde q_i está armazenado?

Estruturas de dados para acessar os termos

- Duas Classes principais de estruturas de dados: hashes e árvores

Estruturas de dados para acessar os termos

- Duas Classes principais de estruturas de dados: hashes e árvores
- Alguns sistemas de RI usam hashes, outros usam árvores.

Estruturas de dados para acessar os termos

- Duas Classes principais de estruturas de dados: hashes e árvores
- Alguns sistemas de RI usam hashes, outros usam árvores.
- Critérios de escolha:

Estruturas de dados para acessar os termos

- Duas Classes principais de estruturas de dados: hashes e árvores
- Alguns sistemas de RI usam hashes, outros usam árvores.
- Critérios de escolha:
 - Existe um número fixo de termos ou ele crescerá indefinidamente?

Estruturas de dados para acessar os termos

- Duas Classes principais de estruturas de dados: hashes e árvores
- Alguns sistemas de RI usam hashes, outros usam árvores.
- Critérios de escolha:
 - Existe um número fixo de termos ou ele crescerá indefinidamente?
 - Quais as frequências relativas com que as várias chaves serão acessadas?

Estruturas de dados para acessar os termos

- Duas Classes principais de estruturas de dados: hashes e árvores
- Alguns sistemas de RI usam hashes, outros usam árvores.
- Critérios de escolha:
 - Existe um número fixo de termos ou ele crescerá indefinidamente?
 - Quais as frequências relativas com que as várias chaves serão acessadas?
 - Quantos termos teremos?

Hashes

- Cada termo do vocabulário é “hasheado” para um inteiro.

Hashes

- Cada termo do vocabulário é “hasheado” para um inteiro.
- Busca-se evitar colisões

Hashes

- Cada termo do vocabulário é “hasheado” para um inteiro.
- Busca-se evitar colisões
- No momento da consulta, faz-se o seguinte: Hasheia o termo de consulta, resolve as colisões, Localiza entrada em uma matriz de elementos com tamanho constante

Hashes

- Cada termo do vocabulário é “hasheado” para um inteiro.
- Busca-se evitar colisões
- No momento da consulta, faz-se o seguinte: Hasheia o termo de consulta, resolve as colisões, Localiza entrada em uma matriz de elementos com tamanho constante
- Prós: Busca em um hash é mais rápida do que em uma árvore.

Hashes

- Cada termo do vocabulário é “hasheado” para um inteiro.
- Busca-se evitar colisões
- No momento da consulta, faz-se o seguinte: Hasheia o termo de consulta, resolve as colisões, Localiza entrada em uma matriz de elementos com tamanho constante
- Prós: Busca em um hash é mais rápida do que em uma árvore.
 - Tempo de consulta é constante.

Hashes

- Cada termo do vocabulário é “hasheado” para um inteiro.
- Busca-se evitar colisões
- No momento da consulta, faz-se o seguinte: Hasheia o termo de consulta, resolve as colisões, Localiza entrada em uma matriz de elementos com tamanho constante
- Prós: Busca em um hash é mais rápida do que em uma árvore.
 - Tempo de consulta é constante.
- Contras

Hashes

- Cada termo do vocabulário é “hasheado” para um inteiro.
- Busca-se evitar colisões
- No momento da consulta, faz-se o seguinte: Hasheia o termo de consulta, resolve as colisões, Localiza entrada em uma matriz de elementos com tamanho constante
- Prós: Busca em um hash é mais rápida do que em uma árvore.
 - Tempo de consulta é constante.
- Contras
 - Não há forma de encontrar pequenas variações (*resume* vs. *résumé*)

Hashes

- Cada termo do vocabulário é “hasheado” para um inteiro.
- Busca-se evitar colisões
- No momento da consulta, faz-se o seguinte: Hasheia o termo de consulta, resolve as colisões, Localiza entrada em uma matriz de elementos com tamanho constante
- Prós: Busca em um hash é mais rápida do que em uma árvore.
 - Tempo de consulta é constante.
- Contras
 - Não há forma de encontrar pequenas variações (*resume* vs. *résumé*)
 - Não permite busca de prefixos (Todos os termos que começam com *automat*)

Hashes

- Cada termo do vocabulário é “hasheado” para um inteiro.
- Busca-se evitar colisões
- No momento da consulta, faz-se o seguinte: Hasheia o termo de consulta, resolve as colisões, Localiza entrada em uma matriz de elementos com tamanho constante
- Prós: Busca em um hash é mais rápida do que em uma árvore.
 - Tempo de consulta é constante.
- Contras
 - Não há forma de encontrar pequenas variações (*resume* vs. *résumé*)
 - Não permite busca de prefixos (Todos os termos que começam com *automat*)
 - é necessário “rehashar” tudo periodicamente se o vocabulário continua crescendo.

Árvores

- Árvores resolvem o problema do prefixo (Encontrar todos os termos começando com *automat*).

Árvores

- Árvores resolvem o problema do prefixo (Encontrar todos os termos começando com *automat*).
- Árvore mais simples: árvore binária

Árvores

- Árvores resolvem o problema do prefixo (Encontrar todos os termos começando com *automat*).
- Árvore mais simples: árvore binária
- Busca é ligeiramente mais lenta que em hashes: $O(\log M)$, onde M é o tamanho do vocabulário.

Árvores

- Árvores resolvem o problema do prefixo (Encontrar todos os termos começando com *automat*).
- Árvore mais simples: árvore binária
- Busca é ligeiramente mais lenta que em hashes: $O(\log M)$, onde M é o tamanho do vocabulário.
- $O(\log M)$ vale apenas para árvores **balanceadas**.

Árvores

- Árvores resolvem o problema do prefixo (Encontrar todos os termos começando com *automat*).
- Árvore mais simples: árvore binária
- Busca é ligeiramente mais lenta que em hashes: $O(\log M)$, onde M é o tamanho do vocabulário.
- $O(\log M)$ vale apenas para árvores **balanceadas**.
- Rebalancear árvores binárias é caro.

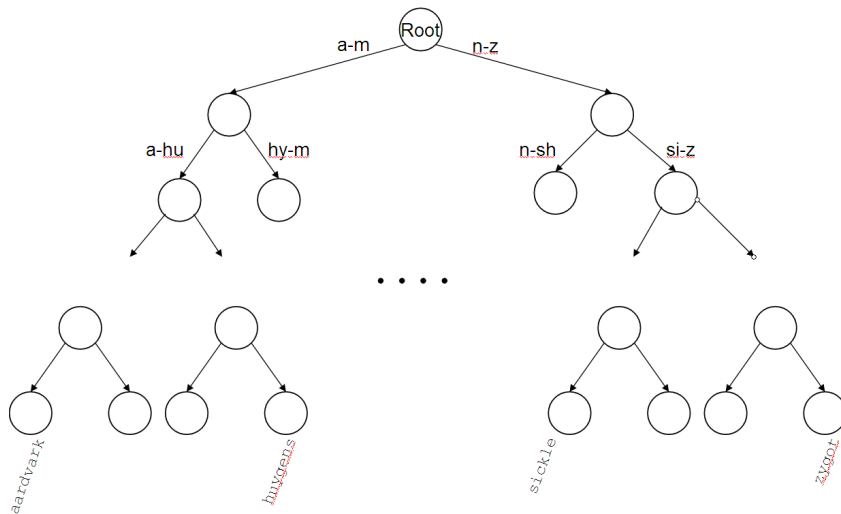
Árvores

- Árvores resolvem o problema do prefixo (Encontrar todos os termos começando com *automat*).
- Árvore mais simples: árvore binária
- Busca é ligeiramente mais lenta que em hashes: $O(\log M)$, onde M é o tamanho do vocabulário.
- $O(\log M)$ vale apenas para árvores **balanceadas**.
- Rebalancear árvores binárias é caro.
- **Arvores-B** resolvem o problema do balanceamento.

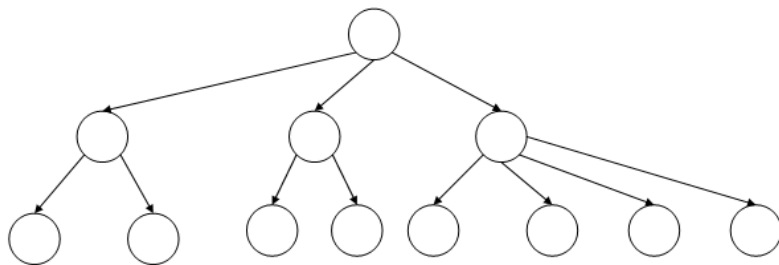
Árvores

- Árvores resolvem o problema do prefixo (Encontrar todos os termos começando com *automat*).
- Árvore mais simples: árvore binária
- Busca é ligeiramente mais lenta que em hashes: $O(\log M)$, onde M é o tamanho do vocabulário.
- $O(\log M)$ vale apenas para árvores **balanceadas**.
- Rebalancear árvores binárias é caro.
- **Arvores-B** resolvem o problema do balanceamento.
- Definição de árvore-B: cada nó interno tem um número de filhos no intervalo $[a, b]$ onde a, b são inteiros positivos apropriados , p.ex., $[2, 4]$.

Árvore Binária



Árvore B



Consultas coringa

- `mon*`: Encontre todos os documentos contendo termos começados por *mon*

Consultas coringa

- mon^* : Encontre todos os documentos contendo termos começados por *mon*
- Fácil com dicionários baseados em árvore B: recupera todos os termos t no intervalo: $mon \leq t < moo$

Consultas coringa

- mon^* : Encontre todos os documentos contendo termos começados por *mon*
- Fácil com dicionários baseados em árvore B: recupera todos os termos t no intervalo: $mon \leq t < moo$
- $*mon$: Encontre todos os documentos contendo termos que terminam com *mon*

Consultas coringa

- mon^* : Encontre todos os documentos contendo termos começados por *mon*
- Fácil com dicionários baseados em árvore B: recupera todos os termos t no intervalo: $mon \leq t < moo$
- $*mon$: Encontre todos os documentos contendo termos que terminam com *mon*
 - Mantém uma árvore adicional para termos *ao contrário*

Consultas coringa

- mon^* : Encontre todos os documentos contendo termos começados por *mon*
- Fácil com dicionários baseados em árvore B: recupera todos os termos t no intervalo: $\text{mon} \leq t < \text{moo}$
- $^*\text{mon}$: Encontre todos os documentos contendo termos que terminam com *mon*
 - Mantém uma árvore adicional para termos *ao contrário*
 - Então recupera todos os termos t no intervalo: $\text{nom} \leq t < \text{non}$

Consultas coringa

- mon^* : Encontre todos os documentos contendo termos começados por *mon*
- Fácil com dicionários baseados em árvore B: recupera todos os termos t no intervalo: $mon \leq t < moo$
- $*mon$: Encontre todos os documentos contendo termos que terminam com *mon*
 - Mantém uma árvore adicional para termos *ao contrário*
 - Então recupera todos os termos t no intervalo: $nom \leq t < non$
- Resultado: Um conjunto de termos que correspondem à consulta coringa

Consultas coringa

- mon^* : Encontre todos os documentos contendo termos começados por *mon*
- Fácil com dicionários baseados em árvore B: recupera todos os termos t no intervalo: $mon \leq t < moo$
- $*mon$: Encontre todos os documentos contendo termos que terminam com *mon*
 - Mantém uma árvore adicional para termos *ao contrário*
 - Então recupera todos os termos t no intervalo: $nom \leq t < non$
- Resultado: Um conjunto de termos que correspondem à consulta coringa
- Então recupera todos os documentos que contenham estes termos

Como lidar com um * no meio de um termo

- Exemplo: m*nchen

Como lidar com um * no meio de um termo

- Exemplo: m^*nchen
- Poderíamos buscar todos os termos que satisfazem m^* e $*nchen$ Na árvore B e reter a interseção dos dois conjuntos.

Como lidar com um * no meio de um termo

- Exemplo: m^*nchen
- Poderíamos buscar todos os termos que satisfazem m^* e $*nchen$ Na árvore B e reter a interseção dos dois conjuntos.
- Mas sai caro

Como lidar com um * no meio de um termo

- Exemplo: m^*nchen
- Poderíamos buscar todos os termos que satisfazem m^* e $*nchen$ Na árvore B e reter a interseção dos dois conjuntos.
- Mas sai caro
- Alternativa: índice [permuterm](#)

Como lidar com um * no meio de um termo

- Exemplo: m^*nchen
- Poderíamos buscar todos os termos que satisfazem m^* e $*nchen$ Na árvore B e reter a interseção dos dois conjuntos.
- Mas sai caro
- Alternativa: índice [permuterm](#)
- Idéia básica: Rotaciona cada consulta coringa, de forma que o * ocorra no final.

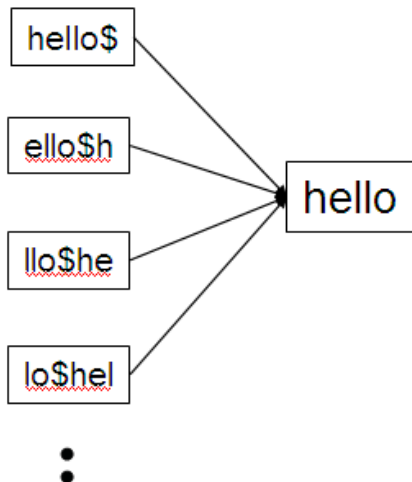
Como lidar com um * no meio de um termo

- Exemplo: m^*nchen
- Poderíamos buscar todos os termos que satisfazem m^* e $*nchen$ Na árvore B e reter a interseção dos dois conjuntos.
- Mas sai caro
- Alternativa: índice [permuterm](#)
- Idéia básica: Rotaciona cada consulta coringa, de forma que o * ocorra no final.
- Armazena cada uma destas rotações no dicionário, por exemplo, em uma árvore B

Índice Permuterm

- Para o termo HELLO: adicione *hello\$*, *ello\$h*, *llo\$he*, *lo\$hel*, *o\$hell*, e *\$hello* à árvore B onde \$ é um símbolo especial

Permuterm → mapeamento de termos



Índice Permuterm

- Para HELLO, adicionamos: *hello\$*, *ello\$h*, *llo\$he*, *lo\$hel*, e *o\$hell*

Índice Permuterm

- Para HELLO, adicionamos: *hello\$, ello\$h, llo\$he, lo\$hel, e o\$hell*
- Consultas

Índice Permuterm

- Para HELLO, adicionamos: *hello\$, ello\$h, llo\$he, lo\$hel, e o\$hell*
- Consultas
 - Para X, acesse X\$

Índice Permuterm

- Para HELLO, adicionamos: *hello\$*, *ello\$h*, *llo\$he*, *lo\$hel*, e *o\$hell*
- Consultas
 - Para X, acesse X\$
 - Para X*, acesse \$X*

Índice Permuterm

- Para HELLO, adicionamos: *hello\$, ello\$h, llo\$he, lo\$hel, e o\$hell*
- Consultas
 - Para X, acesse X\$
 - Para X*, acesse \$X*
 - Para *X, acesse X\$*

Índice Permuterm

- Para HELLO, adicionamos: *hello\$, ello\$h, llo\$he, lo\$hel, e o\$hell*
- Consultas
 - Para X, acesse X\$
 - Para X*, acesse \$X*
 - Para *X, acesse X\$*
 - Para *X*, acesse X*

Índice Permuterm

- Para HELLO, adicionamos: *hello\$, ello\$h, llo\$he, lo\$hel, e o\$hell*
- Consultas
 - Para X, acesse X\$
 - Para X*, acesse \$X*
 - Para *X, acesse X\$*
 - Para *X*, acesse X*
 - Para X*Y, acesse Y\$X*

Índice Permuterm

- Para HELLO, adicionamos: *hello\$, ello\$h, llo\$he, lo\$hel, e o\$hell*
- Consultas
 - Para X, acesse X\$
 - Para X*, acesse \$X*
 - Para *X, acesse X\$*
 - Para *X*, acesse X*
 - Para X*Y, acesse Y\$X*
 - Example: Para hel*o, acesse o\$hel*

Índice Permuterm

- Para HELLO, adicionamos: *hello\$, ello\$h, llo\$he, lo\$hel, e o\$hell*
- Consultas
 - Para X, acesse X\$
 - Para X*, acesse \$X*
 - Para *X, acesse X\$*
 - Para *X*, acesse X*
 - Para X*Y, acesse Y\$X*
 - Example: Para hel*o, acesse o\$hel*
- Um nome mais adequado para um índice Permuterm seria uma árvore permuterm tree.

Índice Permuterm

- Para HELLO, adicionamos: *hello\$, ello\$h, llo\$he, lo\$hel, e o\$hell*
- Consultas
 - Para X, acesse X\$
 - Para X*, acesse \$X*
 - Para *X, acesse X\$*
 - Para *X*, acesse X*
 - Para X*Y, acesse Y\$X*
 - Example: Para hel*o, acesse o\$hel*
- Um nome mais adequado para um índice Permuterm seria uma árvore permuterm tree.
- Mas índice permuterm é o nome mais comum.

Processando um acesso ao índice permuterm

- Rotacione a consulta coringa para a direita

Processando um acesso ao índice permuterm

- Rotacione a consulta coringa para a direita
- Use o acesso à árvore B como descrito anteriormente

Processando um acesso ao índice permuterm

- Rotacione a consulta coringa para a direita
- Use o acesso à árvore B como descrito anteriormente
- Problema: O Permuterm mais do que **quaduplica** o tamanho do dicionário quando comparado a uma árvore B regular. (observação empírica)

Correção ortográfica

Correção ortográfica

- Dois usos principais

Correção ortográfica

- Dois usos principais
 - Correção de documentos a serem indexados

Correção ortográfica

- Dois usos principais
 - Correção de documentos a serem indexados
 - Correção de consultas

Correção ortográfica

- Dois usos principais
 - Correção de documentos a serem indexados
 - Correção de consultas
- Dois métodos diferentes para correção ortográfica

Correção ortográfica

- Dois usos principais
 - Correção de documentos a serem indexados
 - Correção de consultas
- Dois métodos diferentes para correção ortográfica
- Correção de **Palavra isolada**

Correção ortográfica

- Dois usos principais
 - Correção de documentos a serem indexados
 - Correção de consultas
- Dois métodos diferentes para correção ortográfica
- Correção de **Palavra isolada**
 - Verifica cada palavra isoladamente quanto à correção ortográfica

Correção ortográfica

- Dois usos principais
 - Correção de documentos a serem indexados
 - Correção de consultas
- Dois métodos diferentes para correção ortográfica
- Correção de **Palavra isolada**
 - Verifica cada palavra isoladamente quanto à correção ortográfica
 - Não “pega” erros que resultam em palavras corretas, p.ex., *an asteroid that fell **form** the sky*

Correção ortográfica

- Dois usos principais
 - Correção de documentos a serem indexados
 - Correção de consultas
- Dois métodos diferentes para correção ortográfica
- Correção de **Palavra isolada**
 - Verifica cada palavra isoladamente quanto à correção ortográfica
 - Não “pega” erros que resultam em palavras corretas, p.ex., *an asteroid that fell **form** the sky*
- Correção ortográfica **contextual**

Correção ortográfica

- Dois usos principais
 - Correção de documentos a serem indexados
 - Correção de consultas
- Dois métodos diferentes para correção ortográfica
- Correção de **Palavra isolada**
 - Verifica cada palavra isoladamente quanto à correção ortográfica
 - Não “pega” erros que resultam em palavras corretas, p.ex., *an asteroid that fell **form** the sky*
- Correção ortográfica **contextual**
 - Olha para as palavras vizinhas

Correção ortográfica

- Dois usos principais
 - Correção de documentos a serem indexados
 - Correção de consultas
- Dois métodos diferentes para correção ortográfica
- Correção de **Palavra isolada**
 - Verifica cada palavra isoladamente quanto à correção ortográfica
 - Não “pega” erros que resultam em palavras corretas, p.ex., *an asteroid that fell **form** the sky*
- Correção ortográfica **contextual**
 - Olha para as palavras vizinhas
 - Pode corrigir o erro *form/from* acima

Corrigindo documentos

Corrigindo documentos

- Não estamos interessados em correção interativa de documentos (p.ex., MS Word) neste curso.

Corrigindo documentos

- Não estamos interessados em correção interativa de documentos (p.ex., MS Word) neste curso.
- Em RI, Usamos a correção de documentos primariamente para documentos OCR-izados. (OCR = optical character recognition)

Corrigindo documentos

- Não estamos interessados em correção interativa de documentos (p.ex., MS Word) neste curso.
- Em RI, Usamos a correção de documentos primariamente para documentos OCR-izados. (OCR = optical character recognition)
- A filosofia geral em RI é: Não altere os documentos.

Corrigindo consultas

Corrigindo consultas

- Primeiro: correção ortográfica de palavras isoladas

Corrigindo consultas

- Primeiro: correção ortográfica de palavras isoladas
- Premissa 1: Há uma lista de palavras “corretas” a partir da qual as correções podem ser obtidas.

Corrigindo consultas

- Primeiro: correção ortográfica de palavras isoladas
- Premissa 1: Há uma lista de palavras “corretas” a partir da qual as correções podem ser obtidas.
- Premissa 2: Há uma maneira de calcular a **distância** entre uma palavra errada e uma correta.

Corrigindo consultas

- Primeiro: correção ortográfica de palavras isoladas
- Premissa 1: Há uma lista de palavras “corretas” a partir da qual as correções podem ser obtidas.
- Premissa 2: Há uma maneira de calcular a **distância** entre uma palavra errada e uma correta.
- Algoritmo simplificado: retorne a palavra “correta” com a menor distância da palavra errada.

Corrigindo consultas

- Primeiro: correção ortográfica de palavras isoladas
- Premissa 1: Há uma lista de palavras “corretas” a partir da qual as correções podem ser obtidas.
- Premissa 2: Há uma maneira de calcular a **distância** entre uma palavra errada e uma correta.
- Algoritmo simplificado: retorne a palavra “correta” com a menor distância da palavra errada.
- Exemplo: *informaton* → *information*

Corrigindo consultas

- Primeiro: correção ortográfica de palavras isoladas
- Premissa 1: Há uma lista de palavras “corretas” a partir da qual as correções podem ser obtidas.
- Premissa 2: Há uma maneira de calcular a **distância** entre uma palavra errada e uma correta.
- Algoritmo simplificado: retorne a palavra “correta” com a menor distância da palavra errada.
- Exemplo: *informaton* → *information*
- Como lista de palavras corretas, podemos usar o vocabulário de todas as palavras que ocorrem em nossa coleção.

Corrigindo consultas

- Primeiro: correção ortográfica de palavras isoladas
- Premissa 1: Há uma lista de palavras “corretas” a partir da qual as correções podem ser obtidas.
- Premissa 2: Há uma maneira de calcular a **distância** entre uma palavra errada e uma correta.
- Algoritmo simplificado: retorne a palavra “correta” com a menor distância da palavra errada.
- Exemplo: *informaton* → *information*
- Como lista de palavras corretas, podemos usar o vocabulário de todas as palavras que ocorrem em nossa coleção.
- **Porque isto é problemático?**

Alternativas ao uso do vocabulário de termos

Alternativas ao uso do vocabulário de termos

- Um dicionário padrão (Webster's, Aurélio etc.)

Alternativas ao uso do vocabulário de termos

- Um dicionário padrão (Webster's, Aurélio etc.)
- Um dicionário de um domínio específico (Para sistemas de RI especializados)

Alternativas ao uso do vocabulário de termos

- Um dicionário padrão (Webster's, Aurélio etc.)
- Um dicionário de um domínio específico (Para sistemas de RI especializados)
- O vocabulário de termos ponderado, ponderado de forma adequada

Distância entre a palavra errada e a “correta”

Distância entre a palavra errada e a “correta”

- Estudaremos várias alternativas.

Distância entre a palavra errada e a “correta”

- Estudaremos várias alternativas.
- Distância de edição e a distância de Levenshtein

Distância entre a palavra errada e a “correta”

- Estudaremos várias alternativas.
- Distância de edição e a distância de Levenshtein
- Distância de edição ponderada

Distância entre a palavra errada e a “correta”

- Estudaremos várias alternativas.
- Distância de edição e a distância de Levenshtein
- Distância de edição ponderada
- sobreposição de k -grams

Distância de edição

Distância de edição

- A distância de edição entre a string s_1 e a string s_2 é o número mínimo de operações básicas que converte s_1 em s_2 .

Distância de edição

- A distância de edição entre a string s_1 e a string s_2 é o número mínimo de operações básicas que converte s_1 em s_2 .
- Distância de Levenshtein: Operações válidas: inserção, deleção, e substituição

Distância de edição

- A distância de edição entre a string s_1 e a string s_2 é o número mínimo de operações básicas que converte s_1 em s_2 .
- Distância de Levenshtein: Operações válidas: inserção, deleção, e substituição
- Distância de Levenshtein *dog-do*: 1

Distância de edição

- A distância de edição entre a string s_1 e a string s_2 é o número mínimo de operações básicas que converte s_1 em s_2 .
- Distância de Levenshtein: Operações válidas: inserção, deleção, e substituição
- Distância de Levenshtein *dog-do*: 1
- Distância de Levenshtein *cat-cart*: 1

Distância de edição

- A distância de edição entre a string s_1 e a string s_2 é o número mínimo de operações básicas que converte s_1 em s_2 .
- Distância de Levenshtein: Operações válidas: inserção, deleção, e substituição
- Distância de Levenshtein *dog-do*: 1
- Distância de Levenshtein *cat-cart*: 1
- Distância de Levenshtein *cat-cut*: 1

Distância de edição

- A distância de edição entre a string s_1 e a string s_2 é o número mínimo de operações básicas que converte s_1 em s_2 .
- Distância de Levenshtein: Operações válidas: inserção, deleção, e substituição
- Distância de Levenshtein *dog-do*: 1
- Distância de Levenshtein *cat-cart*: 1
- Distância de Levenshtein *cat-cut*: 1
- Distância de Levenshtein *cat-act*: 2

Distância de edição

- A distância de edição entre a string s_1 e a string s_2 é o número mínimo de operações básicas que converte s_1 em s_2 .
- Distância de Levenshtein: Operações válidas: inserção, deleção, e substituição
- Distância de Levenshtein *dog-do*: 1
- Distância de Levenshtein *cat-cart*: 1
- Distância de Levenshtein *cat-cut*: 1
- Distância de Levenshtein *cat-act*: 2
- Distância de Damerau-Levenshtein *cat-act*: 1

Distância de edição

- A distância de edição entre a string s_1 e a string s_2 é o número mínimo de operações básicas que converte s_1 em s_2 .
- Distância de Levenshtein: Operações válidas: inserção, deleção, e substituição
- Distância de Levenshtein *dog-do*: 1
- Distância de Levenshtein *cat-cart*: 1
- Distância de Levenshtein *cat-cut*: 1
- Distância de Levenshtein *cat-act*: 2
- Distância de Damerau-Levenshtein *cat-act*: 1
- Distância de Damerau-Levenshtein inclui a transposição como uma quarta operação possível.

Distância de Levenshtein: Computação

| | | | | | |
|---|---|---|---|---|---|
| | | f | a | s | t |
| | 0 | 1 | 2 | 3 | 4 |
| c | 1 | 1 | 2 | 3 | 4 |
| a | 2 | 2 | 1 | 2 | 3 |
| t | 3 | 3 | 2 | 2 | 2 |
| s | 4 | 4 | 3 | 2 | 3 |

Distância de Levenshtein: Algoritmo

LEVENSHTEINDISTANCE(s_1, s_2)

```

1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 

```

Operações: inserção (custo 1), deleção (custo 1), substituição (custo 1), cópia (custo 0)

Distância de Levenshtein: Algoritmo

LEVENSHTEINDISTANCE(s_1, s_2)

```

1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 

```

Operações: **inserção (custo 1)**, **deleção (custo 1)**, **substituição (custo 1)**, **cópia (custo 0)**

Distância de Levenshtein: Algoritmo

LEVENSHTEINDISTANCE(s_1, s_2)

```

1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 

```

Operações: inserção (custo 1), **deleção (custo 1)**, substituição (custo 1), cópia (custo 0)

Distância de Levenshtein: Algorithm

LEVENSHTEINDISTANCE(s_1, s_2)

```

1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 

```

Operações: inserção (custo 1), deleção (custo 1), **substituição (custo 1)**, cópia (custo 0)

Distância de Levenshtein: Algoritmo

LEVENSHTEINDISTANCE(s_1, s_2)

```

1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 

```

Operações: inserção (custo 1), deleção (custo 1), substituição (custo 1), **cópia (custo 0)**

Distância de Levenshtein: Exemplo

| | | f | | a | | s | | t | | |
|---|--|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| | | <u>0</u> | <u>1</u> | <u>1</u> | <u>2</u> | <u>2</u> | <u>3</u> | <u>3</u> | <u>4</u> | <u>4</u> |
| c | | <u>1</u> | <u>1</u> | <u>2</u> | <u>2</u> | <u>3</u> | <u>3</u> | <u>4</u> | <u>4</u> | <u>5</u> |
| | | <u>1</u> | <u>2</u> | <u>1</u> | <u>2</u> | <u>2</u> | <u>3</u> | <u>3</u> | <u>4</u> | <u>4</u> |
| a | | <u>2</u> | <u>2</u> | <u>2</u> | <u>1</u> | <u>3</u> | <u>3</u> | <u>4</u> | <u>4</u> | <u>5</u> |
| | | <u>2</u> | <u>3</u> | <u>2</u> | <u>3</u> | <u>1</u> | <u>2</u> | <u>2</u> | <u>3</u> | <u>3</u> |
| t | | <u>3</u> | <u>3</u> | <u>3</u> | <u>3</u> | <u>2</u> | <u>2</u> | <u>3</u> | <u>2</u> | <u>4</u> |
| | | <u>3</u> | <u>4</u> | <u>3</u> | <u>4</u> | <u>2</u> | <u>3</u> | <u>2</u> | <u>3</u> | <u>2</u> |
| s | | <u>4</u> | <u>4</u> | <u>4</u> | <u>4</u> | <u>3</u> | <u>2</u> | <u>3</u> | <u>3</u> | <u>3</u> |
| | | <u>4</u> | <u>5</u> | <u>4</u> | <u>5</u> | <u>3</u> | <u>4</u> | <u>2</u> | <u>3</u> | <u>3</u> |

Cada célula da matriz de Levenshtein

| | |
|--|--|
| Custo de chegar aqui a partir do meu vizinho superior esquerdo (cópia or substituição) | Custo de chegar aqui a partir do meu vizinho superior (deleção) |
| Custo de chegar aqui a partir do meu vizinho esquerdo (inserção) | mínimo dos três “movimentos” possíveis; a forma mais barata de chegar aqui |

Distância de Levenshtein: Exemplo

| | | f | | a | | s | | t | | |
|---|--|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| | | <u>0</u> | <u>1</u> | <u>1</u> | <u>2</u> | <u>2</u> | <u>3</u> | <u>3</u> | <u>4</u> | <u>4</u> |
| c | | <u>1</u> | <u>1</u> | <u>2</u> | <u>2</u> | <u>3</u> | <u>3</u> | <u>4</u> | <u>4</u> | <u>5</u> |
| | | <u>1</u> | <u>2</u> | <u>1</u> | <u>2</u> | <u>2</u> | <u>3</u> | <u>3</u> | <u>4</u> | <u>4</u> |
| a | | <u>2</u> | <u>2</u> | <u>2</u> | <u>1</u> | <u>3</u> | <u>3</u> | <u>4</u> | <u>4</u> | <u>5</u> |
| | | <u>2</u> | <u>3</u> | <u>2</u> | <u>3</u> | <u>1</u> | <u>2</u> | <u>2</u> | <u>3</u> | <u>3</u> |
| t | | <u>3</u> | <u>3</u> | <u>3</u> | <u>3</u> | <u>2</u> | <u>2</u> | <u>3</u> | <u>2</u> | <u>4</u> |
| | | <u>3</u> | <u>4</u> | <u>3</u> | <u>4</u> | <u>2</u> | <u>3</u> | <u>2</u> | <u>3</u> | <u>2</u> |
| s | | <u>4</u> | <u>4</u> | <u>4</u> | <u>4</u> | <u>3</u> | <u>2</u> | <u>3</u> | <u>3</u> | <u>3</u> |
| | | <u>4</u> | <u>5</u> | <u>4</u> | <u>5</u> | <u>3</u> | <u>4</u> | <u>2</u> | <u>3</u> | <u>3</u> |