

Manipulando textos e imagens

Prof.: Eduardo Vargas Ferreira

- Em machine learning são comuns as aplicações em que x representa objetos não convencionais, como imagens e documentos de texto.
- Entretanto, computadores não entendem tais documentos diretamente como nós;
- Ou seja, utilizamos conhecimentos “não lógicos” para reconhecer imagens e textos. Por exemplo:
 - ★ Quebramos automaticamente sentenças em unidades de significado;
 - ★ Reconhecemos padrões em imagens com certa facilidade.

- Em machine learning são comuns as aplicações em que x representa objetos não convencionais, como imagens e documentos de texto.
- Entretanto, computadores não entendem tais documentos diretamente como nós;
- Ou seja, utilizamos conhecimentos “não lógicos” para reconhecer imagens e textos. Por exemplo:
 - ★ Quebramos automaticamente sentenças em unidades de significado;
 - ★ Reconhecemos padrões em imagens com certa facilidade.

- Em machine learning são comuns as aplicações em que x representa objetos não convencionais, como imagens e documentos de texto.
- Entretanto, computadores não entendem tais documentos diretamente como nós;
- Ou seja, utilizamos conhecimentos “não lógicos” para reconhecer imagens e textos. Por exemplo:
 - ★ Quebramos automaticamente sentenças em unidades de significado;
 - ★ Reconhecemos padrões em imagens com certa facilidade.

- Em machine learning são comuns as aplicações em que x representa objetos não convencionais, como imagens e documentos de texto.
- Entretanto, computadores não entendem tais documentos diretamente como nós;
- Ou seja, utilizamos conhecimentos “não lógicos” para reconhecer imagens e textos. Por exemplo:
 - ★ Quebramos automaticamente sentenças em unidades de significado;
 - ★ Reconhecemos padrões em imagens com certa facilidade.

- Em machine learning são comuns as aplicações em que x representa objetos não convencionais, como imagens e documentos de texto.
- Entretanto, computadores não entendem tais documentos diretamente como nós;
- Ou seja, utilizamos conhecimentos “não lógicos” para reconhecer imagens e textos. Por exemplo:

★ Quebramos automaticamente sentenças em unidades de significado;

★ Reconhecemos padrões em imagens com certa facilidade.



0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9

- Vamos começar falando sobre mineração de texto (*text mining*);
- Sua apresentação (assim como as imagens) são **não estruturados**, ou seja, os dados são desorganizados e difíceis de trabalhar;
 - ★ P. ex., artigos de jornais, *social media*, vídeo, e-mail etc.
- Os **estruturados** seriam dados organizados de forma gerenciável.
 - ★ P. ex., OLAP (*Online Analytical Processing*), XML (*eXtensible Markup Language*) etc.
- Merrill Lynch projeta que em torno de 80-90% de toda informação potencialmente útil está na forma não estruturada;
- Em 2010, Computer World estimaram que a informação não estruturada representa 70-80% dos dados de uma empresa;

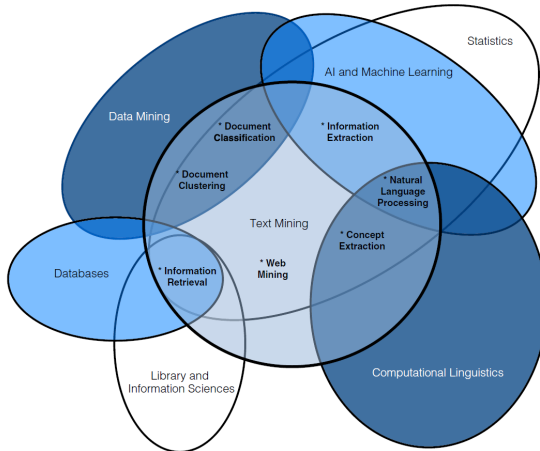
- Vamos começar falando sobre mineração de texto (*text mining*);
- Sua apresentação (assim como as imagens) são **não estruturados**, ou seja, os dados são desorganizados e difíceis de trabalhar;
 - ★ P. ex., artigos de jornais, *social media*, vídeo, e-mail etc.
- Os **estruturados** seriam dados organizados de forma gerenciável.
 - ★ P. ex., OLAP (*Online Analytical Processing*), XML (*eXtensible Markup Language*) etc.
- Merrill Lynch projeta que em torno de 80-90% de toda informação potencialmente útil está na forma não estruturada;
- Em 2010, Computer World estimaram que a informação não estruturada representa 70-80% dos dados de uma empresa;

- Vamos começar falando sobre mineração de texto (*text mining*);
- Sua apresentação (assim como as imagens) são **não estruturados**, ou seja, os dados são desorganizados e difíceis de trabalhar;
 - ★ P. ex., artigos de jornais, *social media*, vídeo, e-mail etc.
- Os **estruturados** seriam dados organizados de forma gerenciável.
 - ★ P. ex., OLAP (*Online Analytical Processing*), XML (*eXtensible Markup Language*) etc.
- Merrill Lynch projeta que em torno de 80-90% de toda informação potencialmente útil está na forma não estruturada;
- Em 2010, Computer World estimaram que a informação não estruturada representa 70-80% dos dados de uma empresa;

- Vamos começar falando sobre mineração de texto (*text mining*);
- Sua apresentação (assim como as imagens) são **não estruturados**, ou seja, os dados são desorganizados e difíceis de trabalhar;
 - ★ P. ex., artigos de jornais, *social media*, vídeo, e-mail etc.
- Os **estruturados** seriam dados organizados de forma gerenciável.
 - ★ P. ex., OLAP (*Online Analytical Processing*), XML (*eXtensible Markup Language*) etc.
- Merrill Lynch projeta que em torno de 80-90% de toda informação potencialmente útil está na forma não estruturada;
- Em 2010, Computer World estimaram que a informação não estruturada representa 70-80% dos dados de uma empresa;

- Vamos começar falando sobre mineração de texto (*text mining*);
- Sua apresentação (assim como as imagens) são **não estruturados**, ou seja, os dados são desorganizados e difíceis de trabalhar;
 - ★ P. ex., artigos de jornais, *social media*, vídeo, e-mail etc.
- Os **estruturados** seriam dados organizados de forma gerenciável.
 - ★ P. ex., OLAP (*Online Analytical Processing*), XML (*eXtensible Markup Language*) etc.
- Merrill Lynch projeta que em torno de 80-90% de toda informação potencialmente útil está na forma não estruturada;
- Em 2010, Computer World estimaram que a informação não estruturada representa 70-80% dos dados de uma empresa;

Mineração de texto no espaço da TI



- Existem pelo menos 5 principais áreas práticas na mineração de texto:
 - ★ **Extração de informação:** Identificação e extração de informação de fatores relevantes e relações entre textos não estruturados;
 - ★ **Clusterização de documentos:** agrupa e categoriza termos, fragmentos, parágrafos ou documentos utilizando métodos de agrupamentos;
 - ★ **Classificação de documentos:** classifica termos, fragmentos, parágrafos ou documentos a partir de exemplos já classificados (dados de treinamento);
 - ★ **Mineração web:** minera dados da internet com foco em interconexões da web;
 - ★ **Natural language processing (NLP):** se preocupa com a interação entre o computador e a linguagem (natural) humana. P ex., reconhecimento de voz, face etc.

- Existem pelo menos 5 principais áreas práticas na mineração de texto:
 - ★ **Extração de informação:** Identificação e extração de informação de fatores relevantes e relações entre textos não estruturados;
 - ★ **Clusterização de documentos:** agrupa e categoriza termos, fragmentos, parágrafos ou documentos utilizando métodos de agrupamentos;
 - ★ **Classificação de documentos:** classifica termos, fragmentos, parágrafos ou documentos a partir de exemplos já classificados (dados de treinamento);
 - ★ **Mineração web:** minera dados da internet com foco em interconexões da web;
 - ★ **Natural language processing (NLP):** se preocupa com a interação entre o computador e a linguagem (natural) humana. P ex., reconhecimento de voz, face etc.

- Existem pelo menos 5 principais áreas práticas na mineração de texto:
 - ★ **Extração de informação:** Identificação e extração de informação de fatores relevantes e relações entre textos não estruturados;
 - ★ **Clusterização de documentos:** agrupa e categoriza termos, fragmentos, parágrafos ou documentos utilizando métodos de agrupamentos;
 - ★ **Classificação de documentos:** classifica termos, fragmentos, parágrafos ou documentos a partir de exemplos já classificados (dados de treinamento);
 - ★ **Mineração web:** minera dados da internet com foco em interconexões da web;
 - ★ **Natural language processing (NLP):** se preocupa com a interação entre o computador e a linguagem (natural) humana. P ex., reconhecimento de voz, face etc.

- Existem pelo menos 5 principais áreas práticas na mineração de texto:
 - ★ **Extração de informação:** Identificação e extração de informação de fatores relevantes e relações entre textos não estruturados;
 - ★ **Clusterização de documentos:** agrupa e categoriza termos, fragmentos, parágrafos ou documentos utilizando métodos de agrupamentos;
 - ★ **Classificação de documentos:** classifica termos, fragmentos, parágrafos ou documentos a partir de exemplos já classificados (dados de treinamento);
 - ★ **Mineração web:** minera dados da internet com foco em interconexões da web;
 - ★ **Natural language processing (NLP):** se preocupa com a interação entre o computador e a linguagem (natural) humana. P ex., reconhecimento de voz, face etc.

- Existem pelo menos 5 principais áreas práticas na mineração de texto:
 - ★ **Extração de informação:** Identificação e extração de informação de fatores relevantes e relações entre textos não estruturados;
 - ★ **Clusterização de documentos:** agrupa e categoriza termos, fragmentos, parágrafos ou documentos utilizando métodos de agrupamentos;
 - ★ **Classificação de documentos:** classifica termos, fragmentos, parágrafos ou documentos a partir de exemplos já classificados (dados de treinamento);
 - ★ **Mineração web:** minera dados da internet com foco em interconexões da web;
 - ★ **Natural language processing (NLP):** se preocupa com a interação entre o computador e a linguagem (natural) humana. P ex., reconhecimento de voz, face etc.

- Existem pelo menos 5 principais áreas práticas na mineração de texto:
 - ★ **Extração de informação:** Identificação e extração de informação de fatores relevantes e relações entre textos não estruturados;
 - ★ **Clusterização de documentos:** agrupa e categoriza termos, fragmentos, parágrafos ou documentos utilizando métodos de agrupamentos;
 - ★ **Classificação de documentos:** classifica termos, fragmentos, parágrafos ou documentos a partir de exemplos já classificados (dados de treinamento);
 - ★ **Mineração web:** minera dados da internet com foco em interconexões da web;
 - ★ **Natural language processing (NLP):** se preocupa com a interação entre o computador e a linguagem (natural) humana. P ex., reconhecimento de voz, face etc.

- Para transformar dados não estruturados em estruturados (na forma numérica) devemos empregar algumas técnicas;
- Este processo deve ser ao mesmo tempo
 - (i) Rápido;
 - (ii) Informativo.
- Primeiramente, separamos os elementos constituintes do texto, identificando cada palavra através de um processo chamado **tokenização**;
- No R, isto pode ser feito facilmente

```
string2 <- "Olá professor, sou aluna de Estatística"  
strsplit(string2, " ")[[1]]  
# [1] "Olá" "professor," "sou" "aluna" "de" "Estatística"
```

- Para transformar dados não estruturados em estruturados (na forma numérica) devemos empregar algumas técnicas;
- Este processo deve ser ao mesmo tempo
 - (i) Rápido;
 - (ii) Informativo.
- Primeiramente, separamos os elementos constituintes do texto, identificando cada palavra através de um processo chamado **tokenização**;
- No R, isto pode ser feito facilmente

```
string2 <- "Olá professor, sou aluna de Estatística"  
strsplit(string2, " ")[[1]]  
# [1] "Olá" "professor," "sou" "aluna" "de" "Estatística"
```

- Para transformar dados não estruturados em estruturados (na forma numérica) devemos empregar algumas técnicas;
- Este processo deve ser ao mesmo tempo
 - (i) Rápido;
 - (ii) Informativo.
- Primeiramente, separamos os elementos constituintes do texto, identificando cada palavra através de um processo chamado **tokenização**;
- No R, isto pode ser feito facilmente

```
string2 <- "Olá professor, sou aluna de Estatística"  
strsplit(string2, " ")[[1]]  
# [1] "Olá" "professor," "sou" "aluna" "de" "Estatística"
```

- Para transformar dados não estruturados em estruturados (na forma numérica) devemos empregar algumas técnicas;
- Este processo deve ser ao mesmo tempo
 - (i) Rápido;
 - (ii) Informativo.
- Primeiramente, separamos os elementos constituintes do texto, identificando cada palavra através de um processo chamado **tokenização**;
- No R, isto pode ser feito facilmente

```
string2 <- "Olá professor, sou aluna de Estatística"  
strsplit(string2, " ")[[1]]  
# [1] "Olá" "professor," "sou" "aluna" "de" "Estatística"
```

- **Stemming** é a redução de um termo ao seu radical, removendo as desinências, afixos e vogais temáticas;
- Tal abordagem auxilia na filtragem e classificação do documento;
- Por exemplo, considere o conjunto de palavras: {prática, praticada, praticados, praticando, praticante, praticar, praticaram, praticidade};
- Apesar de terem características diferentes preservam o mesmo radical **PRATIC**;
- O processo ocorre em etapas, e em cada uma delas uma decisão.

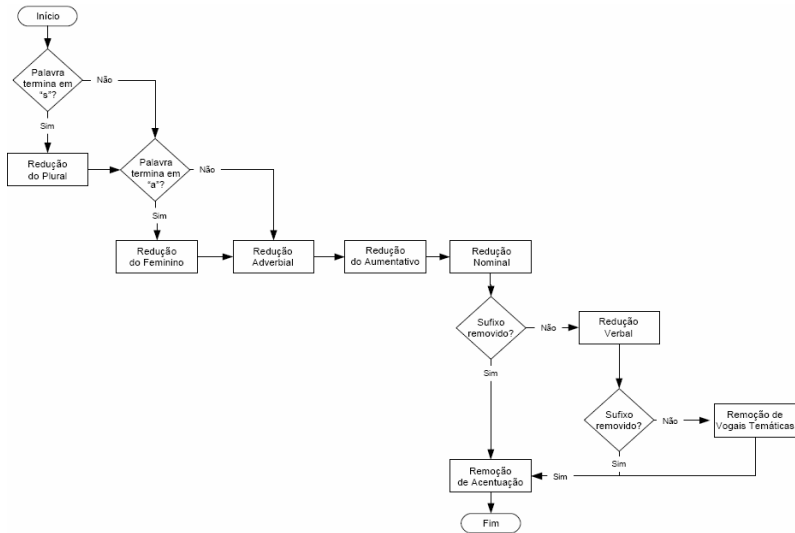
- **Stemming** é a redução de um termo ao seu radical, removendo as desinências, afixos e vogais temáticas;
- Tal abordagem auxilia na filtragem e classificação do documento;
- Por exemplo, considere o conjunto de palavras: {prática, praticada, praticados, praticando, praticante, praticar, praticaram, praticidade};
- Apesar de terem características diferentes preservam o mesmo radical **PRATIC**;
- O processo ocorre em etapas, e em cada uma delas uma decisão.

- **Stemming** é a redução de um termo ao seu radical, removendo as desinências, afixos e vogais temáticas;
- Tal abordagem auxilia na filtragem e classificação do documento;
- Por exemplo, considere o conjunto de palavras: {prática, praticada, praticados, praticando, praticante, praticar, praticaram, praticidade};
- Apesar de terem características diferentes preservam o mesmo radical **PRATIC**;
- O processo ocorre em etapas, e em cada uma delas uma decisão.

- **Stemming** é a redução de um termo ao seu radical, removendo as desinências, afixos e vogais temáticas;
- Tal abordagem auxilia na filtragem e classificação do documento;
- Por exemplo, considere o conjunto de palavras: {prática, praticada, praticados, praticando, praticante, praticar, praticaram, praticidade};
- Apesar de terem características diferentes preservam o mesmo radical **PRATIC**;
- O processo ocorre em etapas, e em cada uma delas uma decisão.

- **Stemming** é a redução de um termo ao seu radical, removendo as desinências, afixos e vogais temáticas;
- Tal abordagem auxilia na filtragem e classificação do documento;
- Por exemplo, considere o conjunto de palavras: {prática, praticada, praticados, praticando, praticante, praticar, praticaram, praticidade};
- Apesar de terem características diferentes preservam o mesmo radical **PRATIC**;
- O processo ocorre em etapas, e em cada uma delas uma decisão.

Técnicas para estruturar os dados: stemming



Técnicas para estruturar os dados: stopwords

- Outra tarefa na preparação dos dados é a identificação das palavras que podem ser desconsideradas nos passos posteriores da análise;
- Nesta fase, tenta-se retirar tudo que não constitui conhecimento no texto;
- Palavras que são muito comuns muitas vezes não são informativas (e.g., “a”, “esse”, ...);
- O resultado é uma lista com as palavras a serem descartadas conhecido como **stopwords** ou **stoplist**.

Técnicas para estruturar os dados: stopwords

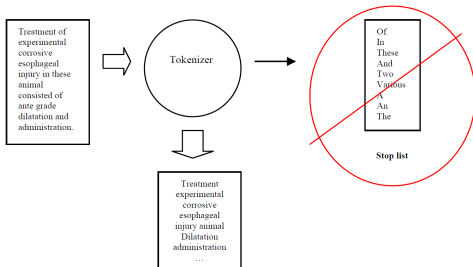
- Outra tarefa na preparação dos dados é a identificação das palavras que podem ser desconsideradas nos passos posteriores da análise;
- Nesta fase, tenta-se retirar tudo que não constitui conhecimento no texto;
- Palavras que são muito comuns muitas vezes não são informativas (e.g., “a”, “esse”, ...);
- O resultado é uma lista com as palavras a serem descartadas conhecido como **stopwords** ou **stoplist**.

Técnicas para estruturar os dados: stopwords

- Outra tarefa na preparação dos dados é a identificação das palavras que podem ser desconsideradas nos passos posteriores da análise;
- Nesta fase, tenta-se retirar tudo que não constitui conhecimento no texto;
- Palavras que são muito comuns muitas vezes não são informativas (e.g., “a”, “esse”, ...);
- O resultado é uma lista com as palavras a serem descartadas conhecido como **stopwords** ou **stoplist**.

Técnicas para estruturar os dados: stopwords

- Outra tarefa na preparação dos dados é a identificação das palavras que podem ser desconsideradas nos passos posteriores da análise;
- Nesta fase, tenta-se retirar tudo que não constitui conhecimento no texto;
- Palavras que são muito comuns muitas vezes não são informativas (e.g., “a”, “esse”, ...);
- O resultado é uma lista com as palavras a serem descartadas conhecido como **stopwords** ou **stoplist**.



Inverse document frequency (IDF)

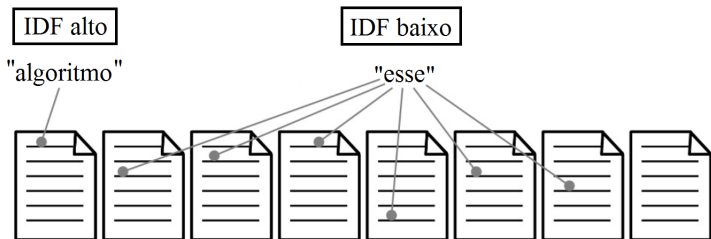


- Uma maneira mais eficiente de resolver a questão de palavras comuns é a chamada **Inverse document frequency (IDF)**;
- A ideia é simples, palavras muito comuns (mais frequentes) recebem menores pesos, pois discriminam menos os documentos;

Inverse document frequency (IDF)



- Uma maneira mais eficiente de resolver a questão de palavras comuns é a chamada **Inverse document frequency (IDF)**;
- A ideia é simples, palavras muito comuns (mais frequentes) recebem menores pesos, pois discriminam menos os documentos;



- Assim, para o i -ésimo termo do k -ésimo documento temos a seguinte formulação

$$a_{ik} = f_{ik} \log \left(\frac{D}{n_i} \right),$$

em que

- ★ a_{ik} é o peso atribuído ao termo i no documento k ;
 - ★ f_{ik} a frequência do termo no documento;
 - ★ D é o número total de documentos;
 - ★ n_i o número de documentos que contém o termo i .
- Obs.: existem outras formas de se ponderar as frequências (e.g. pela raiz quadrada ao invés do logaritmo);

- Queremos minimizar distâncias entre vetores com características similares;
- Considere o exemplo abaixo:

"Eduardo quando crescer será Estatístico";

"eduardo qundo crecser será estatístico".

- Gostaríamos que nosso algoritmo pontuasse as duas respostas da mesma forma;
- Entretanto, simples tokenização não captará tal semelhança;
- Assim, além de incluirmos apenas entradas minúsculas, devemos corrigir a ortografia (e.g. [Peter Norvig's method](#)).
- **Atenção:** sempre que fazemos isso perdemos informação!

- Queremos minimizar distâncias entre vetores com características similares;
- Considere o exemplo abaixo:

“Eduardo quando crescer será Estatístico”;

“eduardo qundo crecser será estatístico”.

- Gostaríamos que nosso algoritmo pontuasse as duas respostas da mesma forma;
- Entretanto, simples tokenização não captará tal semelhança;
- Assim, além de incluirmos apenas entradas minúsculas, devemos corrigir a ortografia (e.g. [Peter Norvig's method](#)).
- **Atenção:** sempre que fazemos isso perdemos informação!

Minimizando distâncias entre vetores



- Queremos minimizar distâncias entre vetores com características similares;
- Considere o exemplo abaixo:

“Eduardo quando crescer será Estatístico”;

“eduardo qundo crecser será estatístico”.

- Gostaríamos que nosso algoritmo pontuasse as duas respostas da mesma forma;
- Entretanto, simples tokenização não captará tal semelhança;
- Assim, além de incluirmos apenas entradas minúsculas, devemos corrigir a ortografia (e.g. [Peter Norvig's method](#)).
- **Atenção:** sempre que fazemos isso perdemos informação!

Minimizando distâncias entre vetores



- Queremos minimizar distâncias entre vetores com características similares;
- Considere o exemplo abaixo:

“Eduardo quando crescer será Estatístico”;

“eduardo qundo crecser será estatístico”.

- Gostaríamos que nosso algoritmo pontuasse as duas respostas da mesma forma;
- Entretanto, simples tokenização não captará tal semelhança;
- Assim, além de incluirmos apenas entradas minúsculas, devemos corrigir a ortografia (e.g. [Peter Norvig's method](#)).
- **Atenção:** sempre que fazemos isso perdemos informação!

- Queremos minimizar distâncias entre vetores com características similares;
- Considere o exemplo abaixo:

“Eduardo quando crescer será Estatístico”;

“eduardo qundo crecser será estatístico”.

- Gostaríamos que nosso algoritmo pontuasse as duas respostas da mesma forma;
- Entretanto, simples tokenização não captará tal semelhança;
- Assim, além de incluirmos apenas entradas minúsculas, devemos corrigir a ortografia (e.g. [Peter Norvig's method](#)).
- **Atenção:** sempre que fazemos isso perdemos informação!

- Queremos minimizar distâncias entre vetores com características similares;
- Considere o exemplo abaixo:

“Eduardo quando crescer será Estatístico” ;

“eduardo qundo crecser será estatístico” .

- Gostaríamos que nosso algoritmo pontuasse as duas respostas da mesma forma;
- Entretanto, simples tokenização não captará tal semelhança;
- Assim, além de incluirmos apenas entradas minúsculas, devemos corrigir a ortografia (e.g. [Peter Norvig's method](#)).
- **Atenção:** sempre que fazemos isso perdemos informação!

Exemplo: Motores de busca (*Search Engine*)

Motores de busca (Search Engine)



- Como encontramos informações que procuramos na internet?
- Como é determinado o ranking dos sites?

The image shows a Google search results page for the query "san francisco dentist". The search bar at the top contains the text "san francisco dentist" and a magnifying glass icon. Below the search bar, the text "Search About 13,500,000 results (0.35 seconds)" is displayed. On the left side, there is a vertical menu with categories: "Everything", "Images", "Maps", "Videos", "News", "Shopping", "More", "Sandwich, MA", "Change location", and "Show search tools". The main content area displays several search results. The first result is an advertisement for "San Francisco Dentist - 'Best of the Bay' in Dentistry" from drzabek.com, with a price of "\$49 Exam & Xrays. Call Today!". Below it is a result for "Blende Dental Group | DrBlende.com" with the address "390 Laurel St # 310, San Francisco, CA". The third result is for "Dr Mike Hack, SF Dentist - 30 Yrs in SF Financial District" from franciadistrictdental.com. The fourth result is for "Cosmetic Dentist" from cosmeticimplantdentistryca.com. The fifth result is for "Need A Good Dentist?" from drvaksmans.com, with a BBB rating and address "450 Sutter Street, San Francisco". The sixth result is for "Rincon Dental SF" from rincondental.com. A map on the right side shows the location of the search results in San Francisco. Red arrows and boxes highlight specific elements: a red box around the search bar, a red box around the search term "san francisco dentist", a red box around the text "Paid ads here", a red box around the text "#1 ranking for 'san francisco dentist'", and a red box around the first search result.

Google

san francisco dentist

Search

About 13,500,000 results (0.35 seconds)

Search term "san francisco dentist"

Everything

Images

Maps

Videos

News

Shopping

More

Sandwich, MA

Change location

Show search tools

San Francisco Dentist - "Best of the Bay" in Dentistry

www.drzabek.com/

\$49 Exam & Xrays. Call Today!

1 Suite 404, 490 Post Street, San Francisco

Directions

Blende Dental Group | DrBlende.com

www.drblende.com/

Painless Sedation & Sleep Dentistry San Francisco Dentist

2 390 Laurel St # 310, San Francisco, CA

(800) 575-3375 - Directions

Dr Mike Hack, SF Dentist - 30 Yrs in SF Financial District

www.franciadistrictdental.com/

Experienced, Friendly & Convenient

→ General Dentistry - Ongoing Care & Cleaning - Experienced Cosmetic Dentistry

Cosmetic Dentist

www.cosmeticimplantdentistryca.com/

Advanced Cosmetic Dentistry For A Beautiful & Healthy Smile. Call Us!

Need A Good Dentist?

www.drvaksmans.com/

A+ BBB Rated Office

Voted Best of the Bay Dentist

450 Sutter Street, San Francisco

(415) 404-6644

★★★★★ 7 reviews

Rincon Dental SF

www.rincondental.com/

San Francisco Dentist, Cosmetic Dentist, Call Us (415) 944-1447 ...

www.aesthetika.net/san-francisco-dentist.html

San Francisco Dentist offering patients dental implants, cosmetic, sedation, TMJ & Invisalign. Call us today at (415) 944-1447 to make your appointment.

San Francisco Dentist

www.drvaksmans.com/

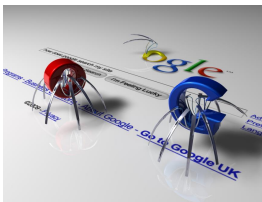
Irena Vaksmans, DDS is a general and cosmetic dental office. Using state of the art

Map for san francisco dentist

Paid ads here

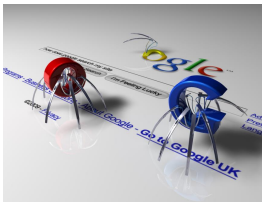
#1 ranking for "san francisco dentist"

- Motores de busca são mecanismos para encontrar informações de texto a partir de palavras-chave indicadas pelo utilizador;
- Eles percorrem “toda” a internet em busca da informação que se pretende (documentos ou endereços de páginas web).
- A forma como a informação é indexada depende de cada motor de busca



- ★ Por palavras, títulos e URL's (como é o caso do Google);
- ★ Ou diretorias (como o Yahoo).

- Motores de busca são mecanismos para encontrar informações de texto a partir de palavras-chave indicadas pelo utilizador;
- Eles percorrem “toda” a internet em busca da informação que se pretende (documentos ou endereços de páginas web).
- A forma como a informação é indexada depende de cada motor de busca

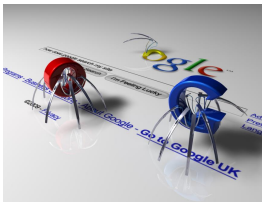


- ★ Por palavras, títulos e URL's (como é o caso do Google);
- ★ Ou diretorias (como o Yahoo).

Motores de busca (Search Engine)



- Motores de busca são mecanismos para encontrar informações de texto a partir de palavras-chave indicadas pelo utilizador;
- Eles percorrem “toda” a internet em busca da informação que se pretende (documentos ou endereços de páginas web).
- A forma como a informação é indexada depende de cada motor de busca



- ★ Por palavras, títulos e URL's (como é o caso do Google);
- ★ Ou diretorias (como o Yahoo).

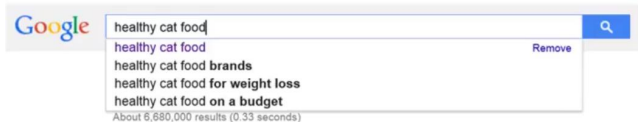


- Vamos construir nosso motor de busca para encontrar em um grupo de 7 websites (o Google utiliza mais de mil milhões) o melhor ranking;
- Nossa pesquisa será
- A visualização dos vetores do espaço fica então dessa forma.

Motores de busca (Search Engine)



- Vamos construir nosso motor de busca para encontrar em um grupo de 7 websites (o Google utiliza mais de mil milhões) o melhor ranking;
- Nossa pesquisa será

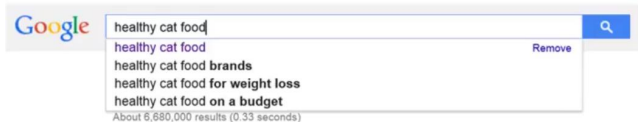


- A visualização dos vetores do espaço fica então dessa forma.

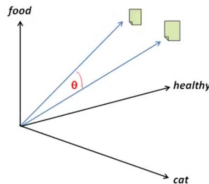
Motores de busca (Search Engine)



- Vamos construir nosso motor de busca para encontrar em um grupo de 7 websites (o Google utiliza mais de mil milhões) o melhor ranking;
- Nossa pesquisa será



- A visualização dos vetores do espaço fica então dessa forma.



- Primeiramente, devemos construir o **Corpus** (que é uma coleção de documentos de texto);
- Abaixo um caso de texto não estruturado para o nosso exemplo;

Web Page	Text Field
1	"Stray cats are running all over the place. I see 10 a day!"
2	"Cats are killers. They kill billions of animals a year."
3	"The best food in Columbus, OH is the North Market."
4	"Brand A is the best tasting cat food around. Your cat will love it."
5	"Buy Brand C cat food for your cat. Brand C makes healthy and happy cats."
6	"The Arnold Classic came to town this weekend. It reminds us to be healthy."
7	"I have nothing to say. In summary, I have told you nothing."

- Note que a maioria dos documentos contém alguma referência sobre **cat**, **healthy** ou **food**;

- Primeiramente, devemos construir o **Corpus** (que é uma coleção de documentos de texto);
- Abaixo um caso de texto não estruturado para o nosso exemplo;

Web Page	Text Field
1	"Stray cats are running all over the place. I see 10 a day!"
2	"Cats are killers. They kill billions of animals a year."
3	"The best food in Columbus, OH is the North Market."
4	"Brand A is the best tasting cat food around. Your cat will love it."
5	"Buy Brand C cat food for your cat. Brand C makes healthy and happy cats."
6	"The Arnold Classic came to town this weekend. It reminds us to be healthy."
7	"I have nothing to say. In summary, I have told you nothing."

- Note que a maioria dos documentos contém alguma referência sobre **cat**, **healthy** ou **food**;

- Como já dito, a fim de melhorar a qualidade da busca, precisamos preparar os dados antes da análise. Por exemplo em:

Stray cats are running all over the place. I see 10 a day!

- Devemos realizar os seguintes passos:
 - ★ Remover pontuação:
Stray cats are running all over the place. I see 10 a day!
 - ★ Stemming:
Stray cats are running all over the place I see 10 a day
 - ★ Trocar os termos em letra maiúscula:
Stray cat are run all over the place I see 10 a day
 - ★ Remover os números:
stray cat are run all over the place I see 10 a day
 - ★ Eliminar os espaços desnecessários:
stray cat are run all over the place I see a day!

- Como já dito, a fim de melhorar a qualidade da busca, precisamos preparar os dados antes da análise. Por exemplo em:

Stray cats are running all over the place. I see 10 a day!

- Devemos realizar os seguintes passos:
 - ★ Remover pontuação:
Stray cats are running all over the place. I see 10 a day!
 - ★ Stemming:
Stray cats are running all over the place I see 10 a day
 - ★ Trocar os termos em letra maiúscula:
Stray cat are run all over the place I see 10 a day
 - ★ Remover os números:
stray cat are run all over the place I see 10 a day
 - ★ Eliminar os espaços desnecessários:
stray cat are run all over the place I see a day!

- Como já dito, a fim de melhorar a qualidade da busca, precisamos preparar os dados antes da análise. Por exemplo em:

Stray cats are running all over the place. I see 10 a day!

- Devemos realizar os seguintes passos:
 - ★ Remover pontuação:
Stray cats are running all over the place. I see 10 a day!
 - ★ Stemming:
Stray cats are running all over the place I see 10 a day
 - ★ Trocar os termos em letra maiúscula:
Stray cat are run all over the place I see 10 a day
 - ★ Remover os números:
stray cat are run all over the place I see 10 a day
 - ★ Eliminar os espaços desnecessários:
stray cat are run all over the place I see a day!

- Como já dito, a fim de melhorar a qualidade da busca, precisamos preparar os dados antes da análise. Por exemplo em:

Stray cats are running all over the place. I see 10 a day!

- Devemos realizar os seguintes passos:
 - ★ Remover pontuação:
Stray cats are running all over the place. I see 10 a day!
 - ★ Stemming:
Stray cats are running all over the place I see 10 a day
 - ★ Trocar os termos em letra maiúscula:
Stray cat are run all over the place I see 10 a day
 - ★ Remover os números:
stray cat are run all over the place I see 10 a day
 - ★ Eliminar os espaços desnecessários:
stray cat are run all over the place I see a day!

- Como já dito, a fim de melhorar a qualidade da busca, precisamos preparar os dados antes da análise. Por exemplo em:

Stray cats are running all over the place. I see 10 a day!

- Devemos realizar os seguintes passos:
 - ★ Remover pontuação:
Stray cats are running all over the place. I see 10 a day!
 - ★ Stemming:
Stray cats are running all over the place I see 10 a day
 - ★ Trocar os termos em letra maiúscula:
Stray cat are run all over the place I see 10 a day
 - ★ Remover os números:
stray cat are run all over the place I see 10 a day
 - ★ Eliminar os espaços desnecessários:
stray cat are run all over the place I see a day!

- Como já dito, a fim de melhorar a qualidade da busca, precisamos preparar os dados antes da análise. Por exemplo em:

Stray cats are running all over the place. I see 10 a day!

- Devemos realizar os seguintes passos:
 - ★ Remover pontuação:
Stray cats are running all over the place. I see 10 a day!
 - ★ Stemming:
Stray cats are running all over the place I see 10 a day
 - ★ Trocar os termos em letra maiúscula:
Stray cat are run all over the place I see 10 a day
 - ★ Remover os números:
stray cat are run all over the place I see 10 a day
 - ★ Eliminar os espaços desnecessários:
stray cat are run all over the place I see a day!

- Como já dito, a fim de melhorar a qualidade da busca, precisamos preparar os dados antes da análise. Por exemplo em:

Stray cats are running all over the place. I see 10 a day!

- Devemos realizar os seguintes passos:
 - ★ Remover pontuação:
Stray cats are running all over the place. I see 10 a day!
 - ★ Stemming:
Stray cats are running all over the place I see 10 a day
 - ★ Trocar os termos em letra maiúscula:
Stray cat are run all over the place I see 10 a day
 - ★ Remover os números:
stray cat are run all over the place I see 10 a day
 - ★ Eliminar os espaços desnecessários:
stray cat are run all over the place I see a day!

Transformando o texto em matriz



- Neste caso, as linhas são os termos e as colunas são os documentos.

Terms	Web Page 1	Web Page 2	Web Page 3	Web Page 4	Web Page 5	Web Page 6	Web Page 7	Query
all	1	0	0	0	0	0	0	0
and	0	0	0	0	1	0	0	0
anim	0	1	0	0	0	0	0	0
are	1	1	0	0	0	0	0	0
arnold	0	0	0	0	0	1	0	0
around	0	0	0	1	0	0	0	0
best	0	0	1	1	0	0	0	0
billion	0	1	0	0	0	0	0	0
brand	0	0	0	1	2	0	0	0
buy	0	0	0	0	1	0	0	0
came	0	0	0	0	0	1	0	0
cat	1	1	0	2	3	0	0	1
classic	0	0	0	0	0	1	0	0
columbus	0	0	1	0	0	0	0	0
classic	0	0	0	0	0	1	0	0
columbus	0	0	1	0	0	0	0	0



Esta linha contém os valores do termo de consulta

- Na forma matricial fica

$$idf.matrix = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \dots & 0 \\ 1 & 1 & 0 & 2 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots \end{bmatrix}$$

Transformando o texto em matriz



- Neste caso, as linhas são os termos e as colunas são os documentos.

Terms	Web Page 1	Web Page 2	Web Page 3	Web Page 4	Web Page 5	Web Page 6	Web Page 7	Query
all	1	0	0	0	0	0	0	0
and	0	0	0	0	1	0	0	0
anim	0	1	0	0	0	0	0	0
are	1	1	0	0	0	0	0	0
arnold	0	0	0	0	0	1	0	0
around	0	0	0	1	0	0	0	0
best	0	0	1	1	0	0	0	0
billion	0	1	0	0	0	0	0	0
brand	0	0	0	1	2	0	0	0
buy	0	0	0	0	1	0	0	0
came	0	0	0	0	0	1	0	0
cat	1	1	0	2	3	0	0	1
classic	0	0	0	0	0	1	0	0
columbus	0	0	1	0	0	0	0	0
classic	0	0	0	0	0	1	0	0
columbus	0	0	1	0	0	0	0	0



Esta linha contém os valores do termo de consulta

- Na forma matricial fica

$$idf.matrix = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \dots & 0 \\ 1 & 1 & 0 & 2 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots \end{bmatrix}$$

- Note que os valores na nossa matriz são simples frequências observadas;
- Mas parece razoável supor que palavras raras pode impulsionar nosso algoritmo;
- P. ex., a palavra **healthy** aparece em somente um documento, enquanto **cat** aparece em 4;
- Então vamos ponderar as palavras pelo inverso da sua frequência (através do método IDF);

- Note que os valores na nossa matriz são simples frequências observadas;
- Mas parece razoável supor que palavras raras pode impulsionar nosso algoritmo;
- P. ex., a palavra **healthy** aparece em somente um documento, enquanto **cat** aparece em 4;
- Então vamos ponderar as palavras pelo inverso da sua frequência (através do método IDF);

- Note que os valores na nossa matriz são simples frequências observadas;
- Mas parece razoável supor que palavras raras pode impulsionar nosso algoritmo;
- P. ex., a palavra **healthy** aparece em somente um documento, enquanto **cat** aparece em 4;
- Então vamos ponderar as palavras pelo inverso da sua frequência (através do método IDF);

Ponderando os termos nos documentos



- Note que os valores na nossa matriz são simples frequências observadas;
- Mas parece razoável supor que palavras raras pode impulsionar nosso algoritmo;
- P. ex., a palavra **healthy** aparece em somente um documento, enquanto **cat** aparece em 4;
- Então vamos ponderar as palavras pelo inverso da sua frequência (através do método IDF);

Terms	Web Page 1	Web Page 2	Web Page 3	Web Page 4	Web Page 5	Web Page 6	Web Page 7	Query
cat	1	1	0	2	3	0	0	1



Terms	Weighting 1	Weighting 2	Weighting 3	Weighting 4	Weighting 5	Weighting 6	Weighting 7	Query
cat	0.8073549	0.8073549	0	1.61471	2.086982	0	0	0.807355

- Uma das vantagens de se trabalhar com vetores no espaço é poder calcular correlações através da geometria;
- Assim, a partir do produto interno entre os vetores (normalizados) temos uma medida sobre o grau de similaridade entre eles ($\cos(\theta)$);

- Para tanto, considere `query.vector` como última coluna da matriz `idf.matrix`

```
query.vector <- idf.matrix[, (N.docs+1)]  
idf.matrix <- idf.matrix[, 1:N.docs]
```

- Então, temos

```
doc.scores <- t(query.vector) %*% idf.matrix
```

- Uma das vantagens de se trabalhar com vetores no espaço é poder calcular correlações através da geometria;
- Assim, a partir do produto interno entre os vetores (normalizados) temos uma medida sobre o grau de similaridade entre eles ($\cos(\theta)$);

- Para tanto, considere `query.vector` como última coluna da matriz `idf.matrix`

```
query.vector <- idf.matrix[, (N.docs+1)]  
idf.matrix <- idf.matrix[, 1:N.docs]
```

- Então, temos

```
doc.scores <- t(query.vector) %*% idf.matrix
```

- Uma das vantagens de se trabalhar com vetores no espaço é poder calcular correlações através da geometria;
- Assim, a partir do produto interno entre os vetores (normalizados) temos uma medida sobre o grau de similaridade entre eles ($\cos(\theta)$);
- Para tanto, considere `query.vector` como última coluna da matriz `idf.matrix`

```
query.vector <- idf.matrix[, (N.docs+1)]  
idf.matrix <- idf.matrix[, 1:N.docs]
```

- Então, temos

```
doc.scores <- t(query.vector) %*% idf.matrix
```

- Uma das vantagens de se trabalhar com vetores no espaço é poder calcular correlações através da geometria;
- Assim, a partir do produto interno entre os vetores (normalizados) temos uma medida sobre o grau de similaridade entre eles ($\cos(\theta)$);
- Para tanto, considere `query.vector` como última coluna da matriz `idf.matrix`

```
query.vector <- idf.matrix[, (N.docs+1)]  
idf.matrix <- idf.matrix[, 1:N.docs]
```

- Então, temos

```
doc.scores <- t(query.vector) %*% idf.matrix
```

- O que fizemos foi (não com esses valores!)

$$\begin{matrix} \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \\ t(\text{query.vector}) \end{matrix} \begin{matrix} \begin{bmatrix} 2 & 1 & 3 \\ 3 & 3 & 2 \\ 4 & 1 & 2 \end{bmatrix} \\ \text{idf.matrix} \end{matrix} = \begin{matrix} \begin{bmatrix} 1 \times 2 + 2 \times 3 + 3 \times 4 \\ 1 \times 1 + 2 \times 3 + 3 \times 1 \\ 1 \times 3 + 2 \times 2 + 3 \times 2 \end{bmatrix} \\ \text{doc.scores} \end{matrix}$$

- Com os escores nas mãos, basta ordená-los e descobrir as melhores indicações

- Note que, devido a ponderação, a *web page 6* ficou em segundo lugar, ainda que apresentou somente um termo, porém "raro".

- O que fizemos foi (não com esses valores!)

$$\begin{matrix} \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \\ t(\text{query.vector}) \end{matrix} \begin{matrix} \begin{bmatrix} 2 & 1 & 3 \\ 3 & 3 & 2 \\ 4 & 1 & 2 \end{bmatrix} \\ \text{idf.matrix} \end{matrix} = \begin{matrix} \begin{bmatrix} 1 \times 2 + 2 \times 3 + 3 \times 4 \\ 1 \times 1 + 2 \times 3 + 3 \times 1 \\ 1 \times 3 + 2 \times 2 + 3 \times 2 \end{bmatrix} \\ \text{doc.scores} \end{matrix}$$

- Com os escores nas mãos, basta ordená-los e descobrir as melhores indicações

Web Page	Score	Text Field
5	0.344	Buy Brand C cat food for your cat . Brand C makes healthy and happy cats .
6	0.183	The Arnold Classic came to town this weekend. It reminds us to be healthy .
4	0.177	Brand A is the best tasting cat food around. Your cat will love it.
3	0.115	The best food in Columbus, OH is the North Market.
2	0.039	Cats are killers. They kill billions of animals a year.
1	0.036	Stray cats are running all over the place. I see 10 a day!
7	0.000	I have nothing to say. In summary, I have told you nothing.

- Note que, devido a ponderação, a *web page* 6 ficou em segundo lugar, ainda que apresentou somente um termo, porém "raro".

- O que fizemos foi (não com esses valores!)

$$\begin{matrix} \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \\ t(\text{query.vector}) \end{matrix} \begin{matrix} \begin{bmatrix} 2 & 1 & 3 \\ 3 & 3 & 2 \\ 4 & 1 & 2 \end{bmatrix} \\ \text{idf.matrix} \end{matrix} = \begin{matrix} \begin{bmatrix} 1 \times 2 + 2 \times 3 + 3 \times 4 \\ 1 \times 1 + 2 \times 3 + 3 \times 1 \\ 1 \times 3 + 2 \times 2 + 3 \times 2 \end{bmatrix} \\ \text{doc.scores} \end{matrix}$$

- Com os escores nas mãos, basta ordená-los e descobrir as melhores indicações

Web Page	Score	Text Field
5	0.344	Buy Brand C cat food for your cat . Brand C makes healthy and happy cats .
6	0.183	The Arnold Classic came to town this weekend. It reminds us to be healthy .
4	0.177	Brand A is the best tasting cat food around. Your cat will love it.
3	0.115	The best food in Columbus, OH is the North Market.
2	0.039	Cats are killers. They kill billions of animals a year.
1	0.036	Stray cats are running all over the place. I see 10 a day!
7	0.000	I have nothing to say. In summary, I have told you nothing.

- Note que, devido a ponderação, a *web page* 6 ficou em segundo lugar, ainda que apresentou somente um termo, porém “raro”.

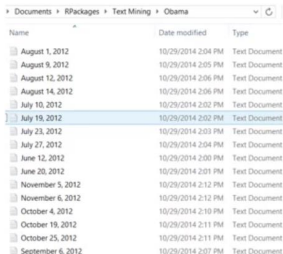
Exemplo: Classificação de textos

- Os discursos de Obama e Romney foram gravados e transcritos;
- A partir desses documentos, queremos encontrar um padrão no discurso;

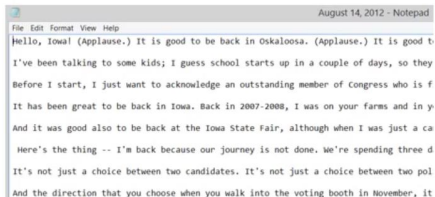


- O algoritmo preditivo deve ser capaz de, a partir de um discurso desconhecido, determinar qual candidato o fez;
- Essa técnica de classificação tem aplicações em detecção de spam, fraude, dentre outras.

- Todos os discursos foram arquivados em pastas para ambos os candidatos. Sem se fazer qualquer “limpeza” no texto.



Name	Date modified	Type
August 1, 2012	10/29/2014 2:04 PM	Text Document
August 9, 2012	10/29/2014 2:05 PM	Text Document
August 12, 2012	10/29/2014 2:06 PM	Text Document
August 14, 2012	10/29/2014 2:06 PM	Text Document
July 10, 2012	10/29/2014 2:02 PM	Text Document
July 19, 2012	10/29/2014 2:02 PM	Text Document
July 23, 2012	10/29/2014 2:03 PM	Text Document
July 27, 2012	10/29/2014 2:04 PM	Text Document
June 12, 2012	10/29/2014 2:00 PM	Text Document
June 20, 2012	10/29/2014 2:01 PM	Text Document
November 5, 2012	10/29/2014 2:12 PM	Text Document
November 6, 2012	10/29/2014 2:12 PM	Text Document
October 4, 2012	10/29/2014 2:10 PM	Text Document
October 19, 2012	10/29/2014 2:11 PM	Text Document
October 25, 2012	10/29/2014 2:11 PM	Text Document
September 6, 2012	10/29/2014 2:07 PM	Text Document



File Edit Format View Help

Hello, Iowa! (Applause.) It is good to be back in Oskaloosa. (Applause.) It is good to be back in Iowa. (Applause.) I've been talking to some kids; I guess school starts up in a couple of days, so they're excited. Before I start, I just want to acknowledge an outstanding member of Congress who is from Iowa. It has been great to be back in Iowa. Back in 2007-2008, I was on your farms and in your communities. And it was good also to be back at the Iowa State Fair, although when I was just a candidate. Here's the thing -- I'm back because our journey is not done. We're spending three days in Iowa. It's not just a choice between two candidates. It's not just a choice between two pols. And the direction that you choose when you walk into the voting booth in November, it

- Inicialmente, precisamos construir o Corpus (coleção de textos);
- Para tanto, vamos aplicar um loop nos discursos removendo pontuações, espaço em branco, *stopwords* etc.;

```
cleanCorpus <- function(corpus){  
  corpus.tmp <- tm_map(corpus, removePunctuation)  
  corpus.tmp <- tm_map(corpus.tmp, stripWhitespace)  
  corpus.tmp <- tm_map(corpus.tmp, tolower)  
  corpus.tmp <- tm_map(corpus.tmp, removeWords, stopwords('english'))  
  return(corpus.tmp)  
}
```

- Inicialmente, precisamos construir o Corpus (coleção de textos);
- Para tanto, vamos aplicar um loop nos discursos removendo pontuações, espaço em branco, *stopwords* etc.;

```
cleanCorpus <- function(corpus){  
  corpus.tmp <- tm_map(corpus, removePunctuation)  
  corpus.tmp <- tm_map(corpus.tmp, stripWhitespace)  
  corpus.tmp <- tm_map(corpus.tmp, tolower)  
  corpus.tmp <- tm_map(corpus.tmp, removeWords, stopwords("english"))  
  return(corpus.tmp)  
}
```

Criando a matriz de documentos



- Neste exemplo temos 1330 termos e nas colunas as frequências de cada um nos discursos.

Document	Candidate	ability	achieve	across	act	advantage	afghanistan	agenda	always	america
August 1, 2012	Romney	1	1	0	1	0	0	1	0	7
August 12, 2012	Romney	0	1	2	0	1	0	3	0	19
August 14, 2012	Romney	3	1	2	1	0	1	0	0	11
August 9, 2012	Romney	0	0	0	0	1	0	0	8	53
July 10, 2012	Romney	0	0	2	0	2	0	1	2	7
July 19, 2012	Romney	2	0	1	0	2	3	0	2	17
July 23, 2012	Romney	1	1	0	0	0	1	0	1	8
July 27, 2012	Romney	0	1	2	1	0	2	0	0	5
June 12, 2012	Romney	1	0	3	0	0	0	1	1	29
June 20, 2012	Romney	0	0	1	1	0	0	0	1	18
November 5, 2012	Romney	0	1	4	0	0	0	0	0	12
November 6, 2012	Romney	0	1	6	5	1	0	1	0	9
October 19, 2012	Romney	0	0	2	0	0	0	0	0	5
August 1, 2012	Obama	0	0	1	1	0	0	0	2	12
August 12, 2012	Obama	0	0	3	3	0	1	0	5	20
August 14, 2012	Obama	0	0	9	0	0	1	0	1	27
August 9, 2012	Obama	0	0	4	2	0	1	0	0	12
July 10, 2012	Obama	0	0	5	1	0	1	0	4	15
July 19, 2012	Obama	0	0	7	3	0	1	0	4	16
July 23, 2012	Obama	0	0	3	2	0	6	0	3	20
July 27, 2012	Obama	0	0	0	2	0	1	0	3	17
June 12, 2012	Obama	0	0	7	0	0	1	0	3	12
June 20, 2012	Obama	0	0	1	0	0	1	0	5	12
November 5, 2012	Obama	0	0	3	1	0	1	0	9	8
November 6, 2012	Obama	0	0	0	0	0	0	0	4	15

Preparação dos dados para predição



- Separamos os dados em treino (70%) e teste (30%);

Dados de treino

Candidate	ability	achieve	across	act	advantage	afghanistan	agenda	always	america
Romney	1	1	0	1	0	0	1	0	7
Romney	0	1	2	0	1	0	3	0	19
Romney	3	1	2	1	0	1	0	0	11
Romney	0	0	0	0	1	0	0	8	53
Romney	0	0	2	0	2	0	1	2	7
Romney	2	0	1	0	2	3	0	2	17
Romney	1	1	0	0	0	1	0	1	8
Romney	0	1	2	1	0	2	0	0	5
Romney	1	0	3	0	0	0	1	1	29
Obama	0	0	1	1	0	0	0	2	12
Obama	0	0	3	3	0	1	0	5	20
Obama	0	0	9	0	0	1	0	1	27
Obama	0	0	4	2	0	1	0	0	12
Obama	0	0	5	1	0	1	0	4	15
Obama	0	0	7	3	0	1	0	4	16
Obama	0	0	3	2	0	6	0	3	20

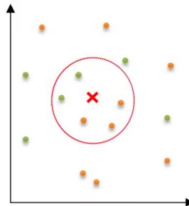
Dados de teste

Speech	ability	achieve	across	act	advantage	afghanistan	agenda	always	america
Speech 1	0	0	1	1	0	0	0	1	18
Speech 2	0	1	4	0	0	0	0	0	12
Speech 3	0	1	6	5	1	0	1	0	9
Speech 4	0	0	2	0	0	0	0	0	5
Speech 5	0	0	0	2	0	1	0	3	17
Speech 6	0	0	7	0	0	1	0	3	12
Speech 7	0	0	1	0	0	1	0	5	12
Speech 8	0	0	3	1	0	1	0	9	8
Speech 9	0	0	0	0	0	0	0	4	15

- Note que a coluna dos candidatos foi removida dos dados de teste, pois será predita a partir dos discursos.

- Utilizando um algoritmo de classificação, chegamos no seguinte resultado

Speech	ability	achieve	across	act	advantage	afghanistan	agenda	always	america
Speech 1	0	0	1	1	0	0	0	1	18
Speech 2	0	1	4	0	0	0	0	0	12
Speech 3	0	1	6	5	1	0	1	0	9
Speech 4	0	0	2	0	0	0	0	0	5
Speech 5	0	0	0	2	0	1	0	3	17
Speech 6	0	0	7	0	0	1	0	3	12
Speech 7	0	0	1	0	0	1	0	5	12
Speech 8	0	0	3	1	0	1	0	9	8
Speech 9	0	0	0	0	0	0	0	4	15



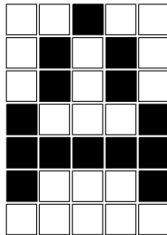
Speech	Actual	Predicted
Speech 1	Romney	Romney
Speech 2	Obama	Obama
Speech 3	Romney	Romney
Speech 4	Romney	Obama
Speech 5	Obama	Obama
Speech 6	Romney	Romney
Speech 7	Obama	Obama
Speech 8	Obama	Obama
Speech 9	Obama	Obama

- [▶ Introduction to the tm Package - Text Mining in R](#)
- [▶ Text Mining Infrastructure in R](#)
- [▶ Text Mining Handbook](#)
- [▶ Distributed Text Mining in R](#)
- [▶ Text mining with Twitter and R](#)

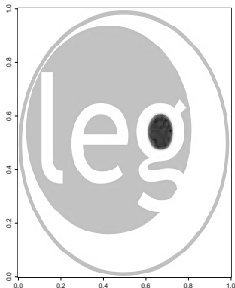
- Inicialmente, vamos entender como as imagens do tipo *raster* são representadas (e.g. JPEG, PNG, ...);
 - ★ *Raster* são imagens que contêm a descrição de cada pixel, em oposição aos gráficos vetoriais.
- Vamos começar com uma ideia simples, utilizando uma matriz binária

- Inicialmente, vamos entender como as imagens do tipo *raster* são representadas (e.g. JPEG, PNG, ...);
 - ★ *Raster* são imagens que contêm a descrição de cada pixel, em oposição aos gráficos vetoriais.
- Vamos começar com uma ideia simples, utilizando uma matriz binária

$$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$



- Podemos ir um passo além



- Ao invés de usar apenas 0 (branco) e 1 (preto), usamos números entre 0 e 1 para denotar a intensidade de cinza.
- Quanto mais pixels, maior a resolução.

```
library(jpeg)
imagem=readJPEG("282px-leg.jpg")
image(t(imagem[282:1, ,3]), col = grey.colors(1000, start = 0, end = 1))
```

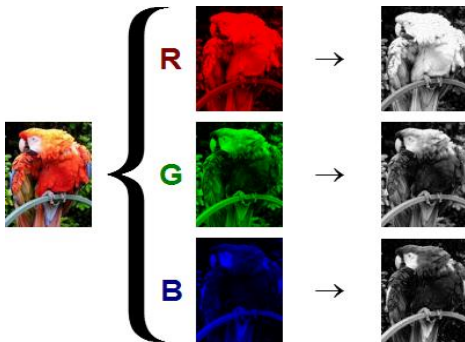

- Usando essa ideia, podemos representar uma imagem a partir de três matrizes simultaneamente (com as cores primárias. Cada elemento é um número entre 0 e 1)



- A primeira indica o quanto de azul em cada pixel;
- A segunda indica o quanto de amarelo;
- A terceira indica o quanto de vermelho.

```
library(jpeg)
imagem=readJPEG("282px-leg.jpg")
rasterImage(imagem, 0, 0, 1, 1)
```

- Podemos, ao invés das cores primárias, utilizar o vermelho, verde e azul (*RGB channels*)

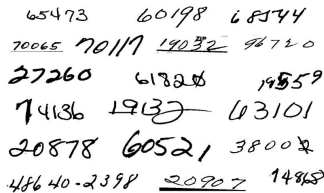


Exemplo: Classificação de dígitos
escritos à mão

- O reconhecimento de imagens é um assunto muito estudado devido a sua variedade de aplicações;
- A classificação de dígitos escritos à mão é um dos assuntos mais discutidos nesta área, e muitos métodos foram desenvolvidos ao longo dos anos;
- A dificuldade deste reconhecimento é causada pela alta variabilidade das imagens;
- Neste exemplo vamos aplicar algumas técnicas vistas nas aulas anteriores para classificar os dígitos 1, 2 e 7 escritos à mão.

65473 60198 68544
70065 70117 19032 96720
27260 61820 19559
74136 19137 63101
20878 60521 38002
48640-2398 20907 14868

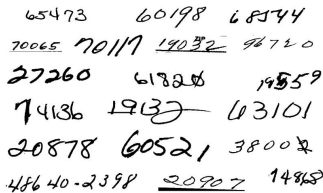
- O reconhecimento de imagens é um assunto muito estudado devido a sua variedade de aplicações;
- A classificação de dígitos escritos à mão é um dos assuntos mais discutidos nesta área, e muitos métodos foram desenvolvidos ao longo dos anos;
- A dificuldade deste reconhecimento é causada pela alta variabilidade das imagens;
- Neste exemplo vamos aplicar algumas técnicas vistas nas aulas anteriores para classificar os dígitos 1, 2 e 7 escritos à mão.



Handwritten digits for classification:

65473	60198	68544	
<u>70065</u>	70117	<u>19032</u>	96720
27260	61820	19559	
74136	19137	63101	
20878	60521	38002	
48640-2398	<u>20907</u>	14868	

- O reconhecimento de imagens é um assunto muito estudado devido a sua variedade de aplicações;
- A classificação de dígitos escritos à mão é um dos assuntos mais discutidos nesta área, e muitos métodos foram desenvolvidos ao longo dos anos;
- A dificuldade deste reconhecimento é causada pela alta variabilidade das imagens;
- Neste exemplo vamos aplicar algumas técnicas vistas nas aulas anteriores para classificar os dígitos 1, 2 e 7 escritos à mão.



Handwritten digits for classification:

65473	60198	68544
<u>70065</u>	70117	<u>19032</u> 96720
27260	61820	19559
74136	1913	63101
20878	60521	38002
48640-2398	<u>20907</u>	14868

- O reconhecimento de imagens é um assunto muito estudado devido a sua variedade de aplicações;
- A classificação de dígitos escritos à mão é um dos assuntos mais discutidos nesta área, e muitos métodos foram desenvolvidos ao longo dos anos;
- A dificuldade deste reconhecimento é causada pela alta variabilidade das imagens;
- Neste exemplo vamos aplicar algumas técnicas vistas nas aulas anteriores para classificar os dígitos 1, 2 e 7 escritos à mão.

65473 60198 68544
70065 70117 19032 96720
27260 61820 19559
74136 19137 63101
20878 60521 38002
48640-2398 20907 14868

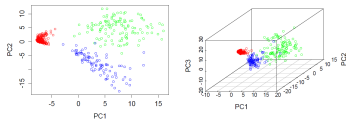
- Os dígitos originais, de diferentes tamanhos e orientações, foram normalizados resultando em imagens 16×16 em escala de cinza;
- E, posteriormente, rearranjados em uma coluna de 256 colunas relativa às cores e uma referente ao dígito em questão;
- Note que a variável resposta é categórica (assume somente três valores: 1,2 ou 7);
- Os métodos em competição são:
 - ★ *k-Nearest Neighbors* (K-NN);
 - ★ *Linear discriminant analysis* (LDA);
 - ★ *Quadratic discriminant analysis* (QDA);
 - ★ *Support vector machine* (SVM);
 - ★ Regressão logística.

- Os dígitos originais, de diferentes tamanhos e orientações, foram normalizados resultando em imagens 16×16 em escala de cinza;
- E, posteriormente, rearranjados em uma coluna de 256 colunas relativa às cores e uma referente ao dígito em questão;
- Note que a variável resposta é categórica (assume somente três valores: 1,2 ou 7);
- Os métodos em competição são:
 - ★ *k-Nearest Neighbors* (K-NN);
 - ★ *Linear discriminant analysis* (LDA);
 - ★ *Quadratic discriminant analysis* (QDA);
 - ★ *Support vector machine* (SVM);
 - ★ Regressão logística.

- Os dígitos originais, de diferentes tamanhos e orientações, foram normalizados resultando em imagens 16×16 em escala de cinza;
- E, posteriormente, rearranjados em uma coluna de 256 colunas relativa às cores e uma referente ao dígito em questão;
- Note que a variável resposta é categórica (assume somente três valores: 1,2 ou 7);
- Os métodos em competição são:
 - ★ *k-Nearest Neighbors* (K-NN);
 - ★ *Linear discriminant analysis* (LDA);
 - ★ *Quadratic discriminant analysis* (QDA);
 - ★ *Support vector machine* (SVM);
 - ★ Regressão logística.

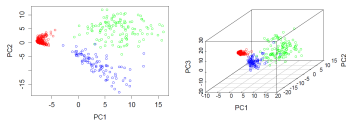
- Os dígitos originais, de diferentes tamanhos e orientações, foram normalizados resultando em imagens 16×16 em escala de cinza;
- E, posteriormente, rearranjados em uma coluna de 256 colunas relativa às cores e uma referente ao dígito em questão;
- Note que a variável resposta é categórica (assume somente três valores: 1,2 ou 7);
- Os métodos em competição são:
 - ★ *k-Nearest Neighbors* (K-NN);
 - ★ *Linear discriminant analysis* (LDA);
 - ★ *Quadratic discriminant analysis* (QDA);
 - ★ *Support vector machine* (SVM);
 - ★ Regressão logística.

- Utilizamos análise das componentes principais a fim de reduzir a dimensão do espaço;



- Aplicamos o método de validação cruzada (k -dobras, com $k = 4$) para avaliar/comparar o desempenho dos modelos ajustados;
- Particionamos a amostra em k grupos de tamanhos iguais;
- Um grupo é separado para validação e modelo é ajustado para os $k - 1$ demais grupos.

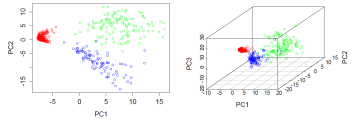
- Utilizamos análise das componentes principais a fim de reduzir a dimensão do espaço;



- Aplicamos o método de validação cruzada (k -dobras, com $k = 4$) para avaliar/comparar o desempenho dos modelos ajustados;
- Particionamos a amostra em k grupos de tamanhos iguais;
- Um grupo é separado para validação e modelo é ajustado para os $k - 1$ demais grupos.

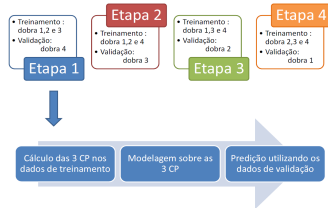
Forma de avaliar o desempenho

- Utilizamos análise das componentes principais a fim de reduzir a dimensão do espaço;



- Aplicamos o método de validação cruzada (k -dobras, com $k = 4$) para avaliar/comparar o desempenho dos modelos ajustados;

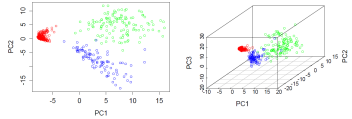
- Particionamos a amostra em k grupos de tamanhos iguais;



- Um grupo é separado para validação e modelo é ajustado para os $k - 1$ demais grupos.

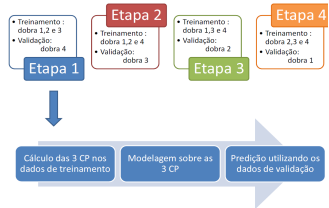
Forma de avaliar o desempenho

- Utilizamos análise das componentes principais a fim de reduzir a dimensão do espaço;



- Aplicamos o método de validação cruzada (k -dobras, com $k = 4$) para avaliar/comparar o desempenho dos modelos ajustados;

- Particionamos a amostra em k grupos de tamanhos iguais;



- Um grupo é separado para validação e modelo é ajustado para os $k - 1$ demais grupos.

- Os modelos de classificação K -NN foram ajustados com valores de $K = 1, 3, 5, 10$ e 15 ;
- Os modelos de classificação SVM foram ajustados utilizando *kernel* linear, polinomial, sigmoide e radial;
- Para medir a qualidade do ajuste, utilizamos o percentual de classificações incorretas;

- Os modelos de classificação K -NN foram ajustados com valores de $K = 1, 3, 5, 10$ e 15 ;
- Os modelos de classificação SVM foram ajustados utilizando *kernel* linear, polinomial, sigmoide e radial;
- Para medir a qualidade do ajuste, utilizamos o percentual de classificações incorretas;

- Os modelos de classificação K -NN foram ajustados com valores de $K = 1, 3, 5, 10$ e 15 ;
- Os modelos de classificação SVM foram ajustados utilizando *kernel* linear, polinomial, sigmoide e radial;
- Para medir a qualidade do ajuste, utilizamos o percentual de **classificações incorretas**;

	Treinamento	Rank	Validação	Rank	Kaggle	Rank
Lógica	0.003	2	0.013	8	0.04	10
LDA	0.023	10	0.023	10	0.04	10
QDA	0.013	9	0.02	9	0.04	10
K-NN (1)	0	1	0.01	7	0.03	7
K-NN (3)	0.005	5	0.008	5	0.01	1
K-NN(5)	0.007	7	0.005	3	0.02	5
K-NN (10)	0.007	7	0.005	3	0.02	5
K-NN(15)	0.009	8	0.008	5	0.03	7
SVM-Linear	0.004	4	0.005	3	0.02	5
SVM-Polinomial	0.048	12	0.055	12	0.06	12
SVM-Sigmoide	0.036	11	0.033	11	0.05	11
SVM-Radial	0.004	4	0.01	7	0.02	5