

---

# Uma Introdução à Mineração de Informações

Fabrício J. Barth

Apontador

<http://www.apontador.com.br>

<http://www.apontador.com.br/institucional/>

[fabricao.barth@lbslocal.com](mailto:fabricao.barth@lbslocal.com)

Outubro de 2010

---

---

# Objetivo

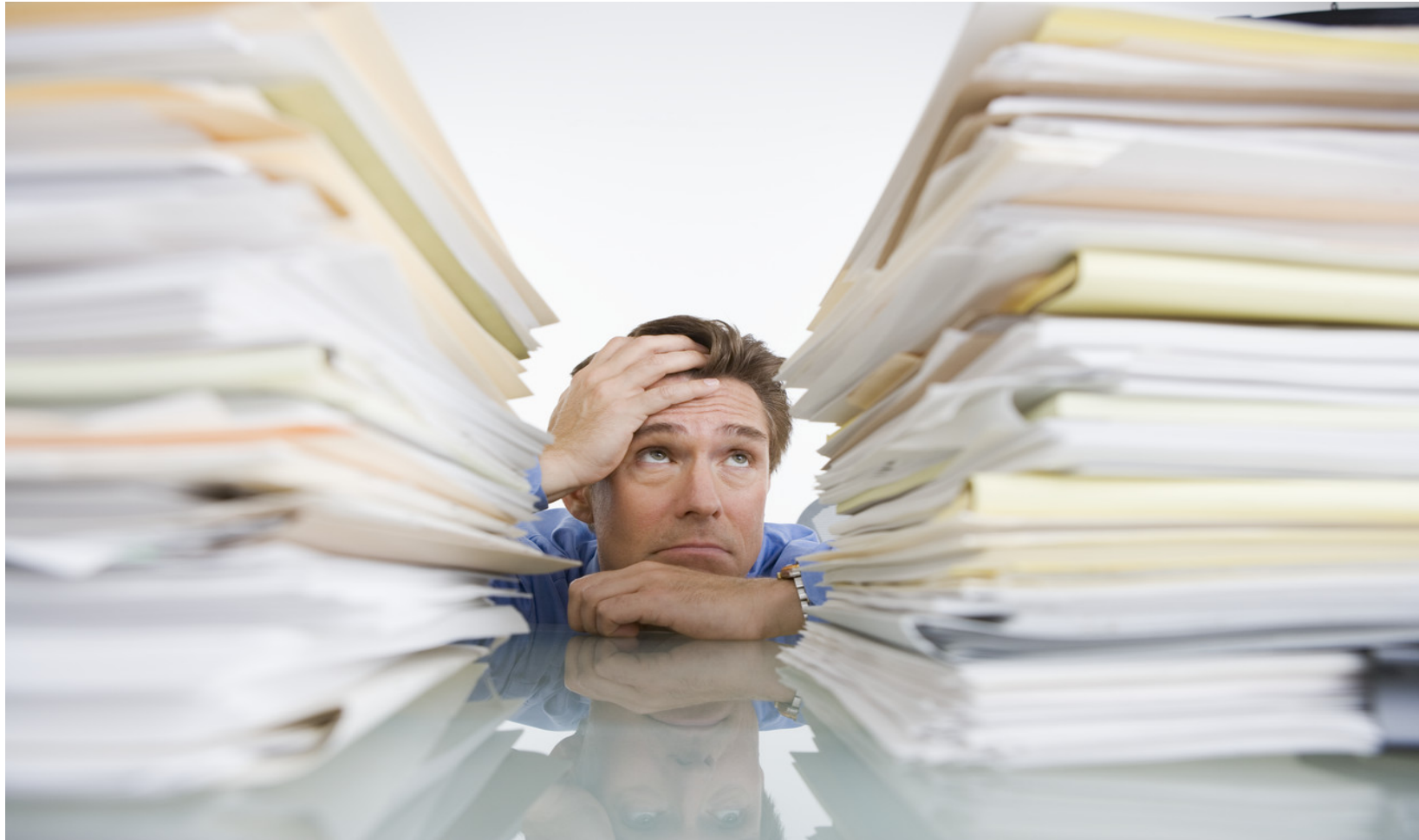
Apresentar a importância do tema, os conceitos relacionados e alguns exemplos de aplicações.

---

# Importância do Tema

---

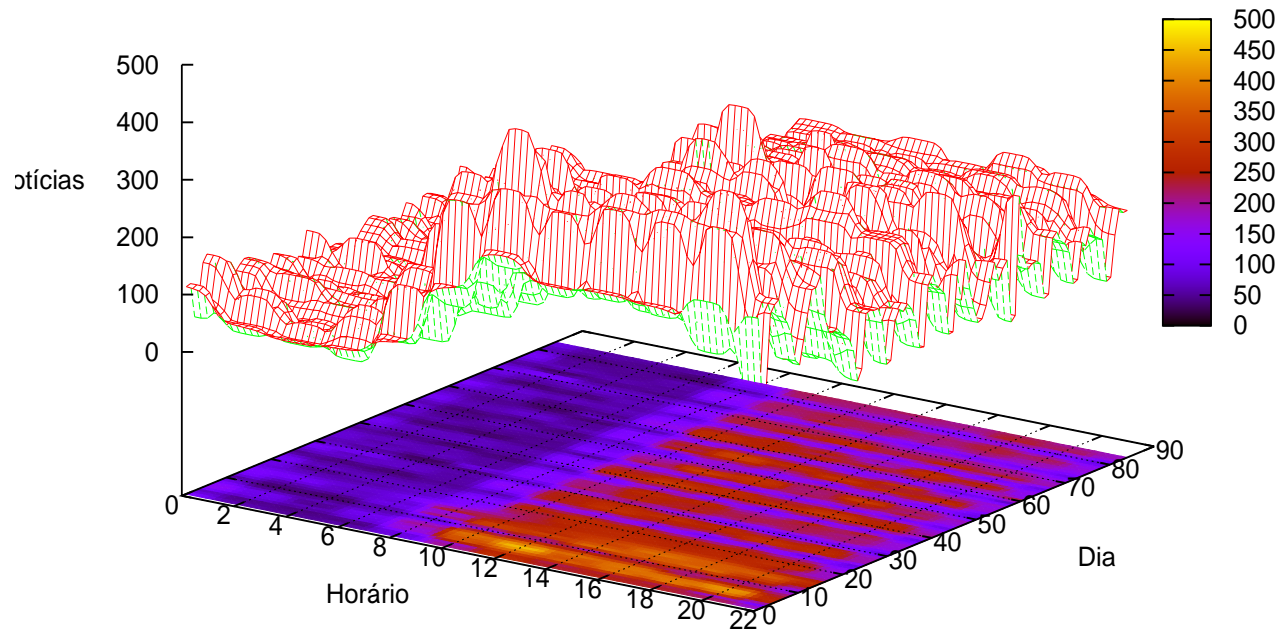
# Problema



<http://investingcaffeine.com/2010/01/07/tmi-the-age-of-information-overload/>

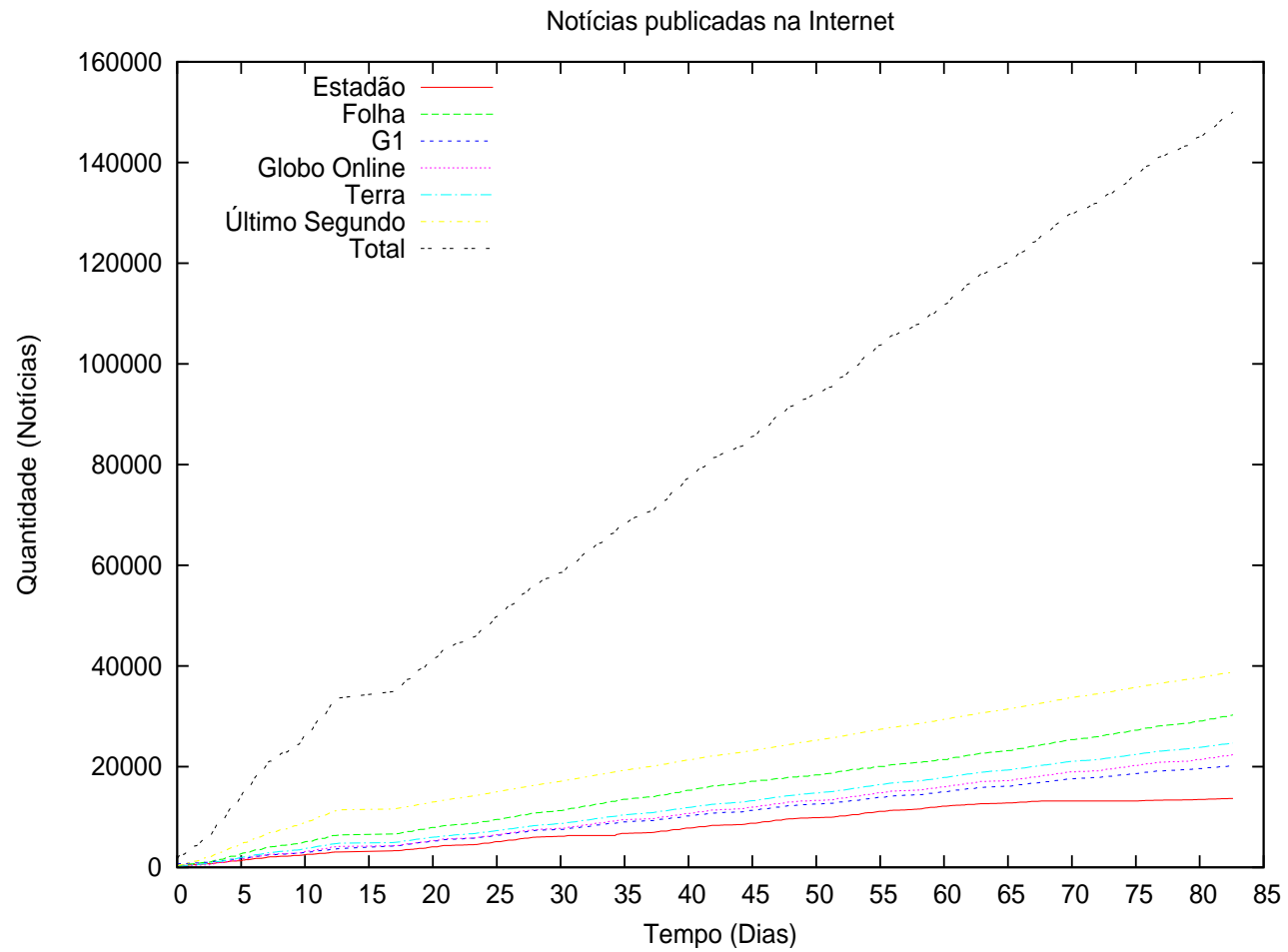
# Alguns dados...

Relação Horário x Dia x Quantidade de Notícias Produzidas



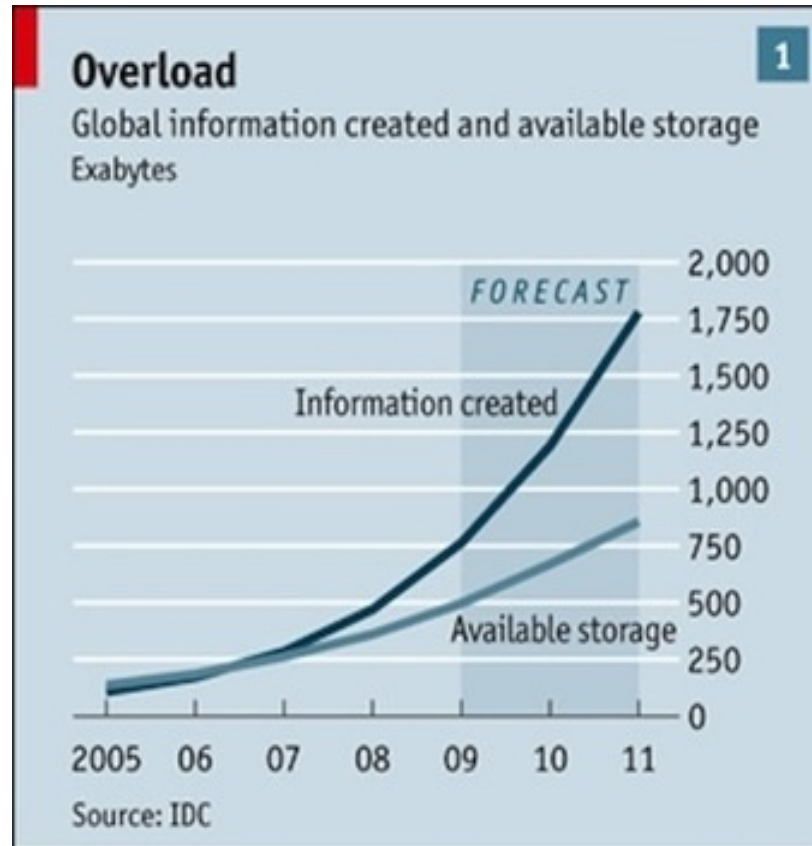
Quantidade de notícias publicadas na Web por apenas seis veículos de notícias (D0 = 17/07/2007)

# Mais dados...



D0 = 17/07/2007

# Big Data



*"We collect an astonishing amount of digital information... ...we've long since surpassed our ability to store and process it all. Big data is here, and it's causing big problems..."* [Data, data everywhere. A special report on managing information 2010]

---

## Por que minerar informações?

- Explicitar conhecimento médico a partir de registros médicos.
- Sumarizar tendências de publicações de artigos e patentes sobre um determinado tema.
- Sumarizar e filtrar notícias relevantes.



- 
- Sumarizar a opinião expressa na Web sobre a sua empresa.
  - Identificar padrões de navegação em sites.
  - Identificar conteúdo impróprio em sites.

---

**Explicitar  
conhecimento médico  
a partir de registros  
médicos**

---

# Diagnóstico para o uso de lentes de contato

O setor de oftalmologia de um hospital da cidade de São Paulo possui, no seu banco de dados, um histórico de pacientes que procuraram o hospital queixando-se de problemas na visão.

A conduta, em alguns casos, realizada pelo corpo clínico de oftalmologistas do hospital é indicar o uso de lentes ao paciente.

**Problema: Extrair do banco de dados do hospital uma hipótese que explica que paciente deve usar ou não lente de contatos.**

---

**Quais são os  
componentes deste  
sistema?**

---

# Componentes

- Que objetos são relevantes para a criação da hipótese?
- Como representá-los?
- Que linguagem de representação de conhecimento deve-se utilizar para representar a hipótese?
- Que algoritmo utilizar para gerar a hipótese?

- 
- Que objetos são relevantes?
    - ★ Depois de um estudo detalhado do problema com especialistas da área...
    - ★ **Idade** do paciente.
    - ★ Se o paciente tem ou não **miopia**.
    - ★ Se o paciente tem ou não **astigmatismo**.
    - ★ Qual é a taxa de **lacrimejamento** dos olhos do paciente.
  - Como representá-los? **Atributo/Valor**

---

# Atributos

- idade (jovem, adulto, idoso)
- miopia (míope, hipermetrópe)
- astigmatismo (não, sim)
- taxa de lacrimejamento (reduzido, normal)
- lentes de contato (forte, fraca, nenhuma)

---

# Dados

Idade	Miopia	Astigmat.	Lacrimenj.	Lentes
jovem	míope	não	reduzido	nenhuma
jovem	míope	não	normal	fraca
jovem	míope	sim	reduzido	nenhuma
jovem	míope	sim	normal	forte
jovem	hiper	não	reduzido	nenhuma
jovem	hiper	não	normal	fraca
jovem	hiper	sim	reduzido	nenhuma
jovem	hiper	sim	normal	forte
adulto	míope	não	reduzido	nenhuma



---

Idade	Miopia	Astigmat.	Lacrimenj.	Lentes
adulto	míope	não	normal	fraca
adulto	míope	sim	reduzido	nenhuma
adulto	míope	sim	normal	forte
adulto	hiper	sim	reduzido	nenhuma
adulto	hiper	não	normal	fraca
adulto	hiper	sim	reduzido	nenhuma
adulto	hiper	sim	normal	nenhuma

---

Idade	Miopia	Astigmat.	Lacrimej.	Lentes
idoso	míope	não	reduzido	nenhuma
idoso	míope	não	normal	nenhuma
idoso	míope	sim	reduzido	nenhuma
idoso	míope	sim	normal	forte
idoso	hiper	não	reduzido	nenhuma
idoso	hiper	não	normal	fraca
idoso	hiper	sim	reduzido	nenhuma
idoso	hiper	sim	normal	nenhuma

---

## Extração de “conhecimento”

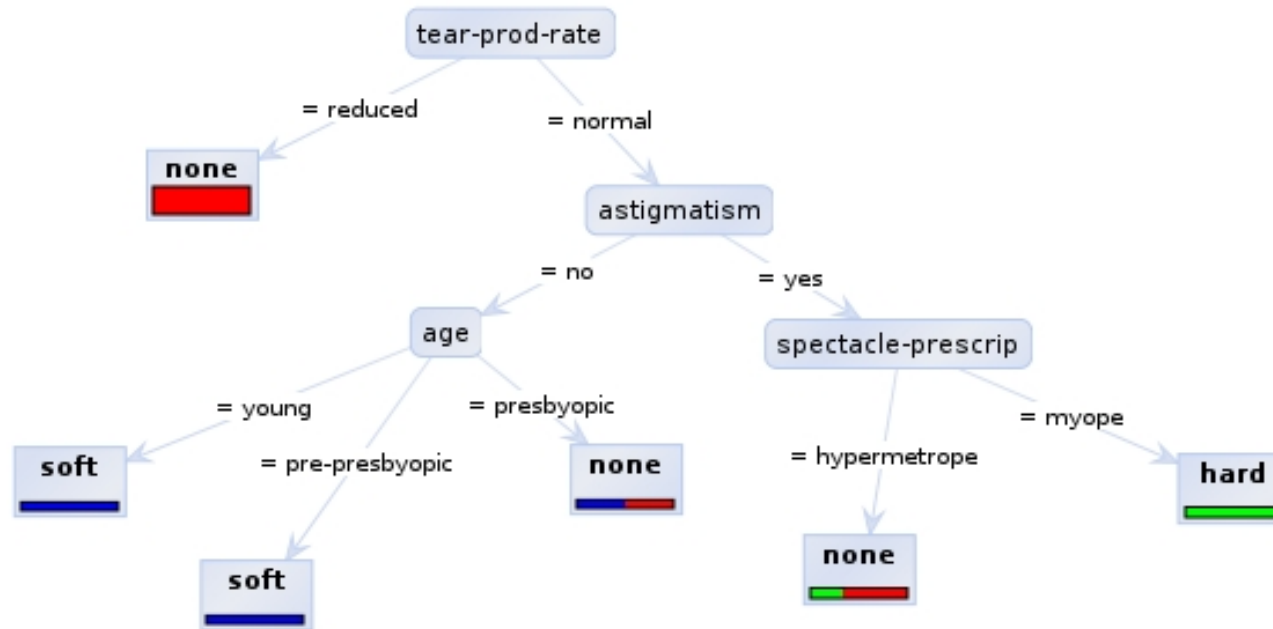
- O que foi apresentado nos slides anteriores pode ser considerado como conhecimento? **Não**
- Pode ser apresentado como uma informação que consegue explicar a tomada de decisão dos especialistas? **Não**
- **O que fazer?**

---

## Extração de “conhecimento”

- Extrair a informação realmente relevante.
- Utilizar uma linguagem de representação **compreensível** ao ser humano.

# Árvore de decisão



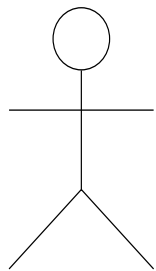
- Cada nodo interno testa um atributo.
- Cada aresta corresponde a um valor de atributo.
- Cada nodo folha retorna uma classificação.

---

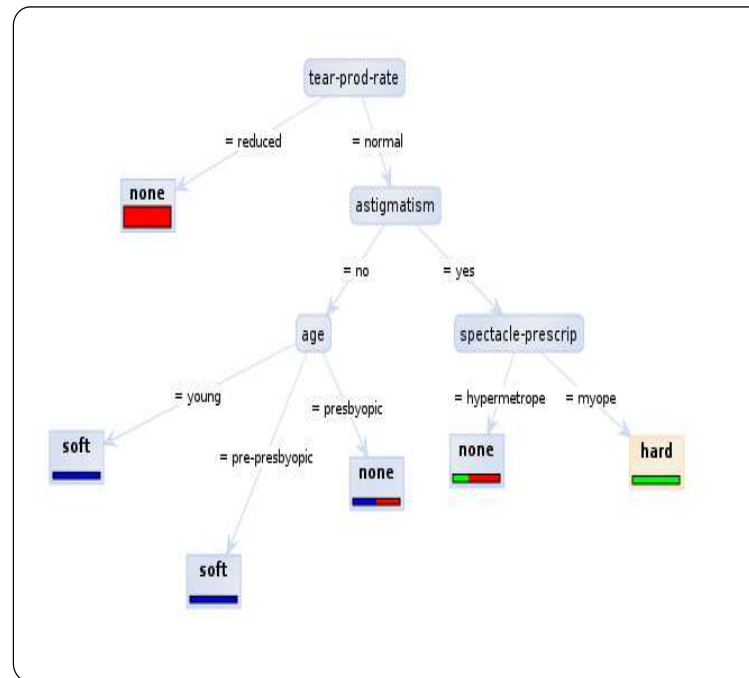
# Algoritmos Indutores de Árvores de Decisão

- **Que algoritmo utilizar para gerar hipóteses na forma de árvores de decisão?**
- ID3, C4.5[[Quinlan 1988](#)]: são algoritmos indutores de árvore de decisão, **top-down**, **recursivos** e que fazem uso do conceito de **entropia** para identificar os melhores atributos que representam o conjunto de dados.

# Resultado: Sistema Especialista



(none, soft, hard)



---

# Sistema Especialista: Regras de Produção

- Baseado na premissa que o processo de tomada de decisão humano pode ser modelado por meio de regras do tipo **SE condições ENTÃO conclusões [FAÇA ações]**
- Convertendo uma árvore de decisão em regras de produção:



- 
- SE lacrimejamento=reduzido ENTÃO lente=nenhuma
  - SE lacrimejamento=normal E astigmatismo=não ENTÃO lente=fraca
  - SE lacrimejamento=normal E astigmatismo=sim E miopia=miope ENTÃO lente=forte
  - SE lacrimejamento=normal E astigmatismo=sim E miopia=hipermetrope ENTÃO lente=nenhuma

---

## Um pouco de código...

- Gerar a árvore de decisão usando o RapidMiner<sup>a</sup>.
- Codificar as regras de produção usando o *Drools Expert*<sup>b</sup>.

---

<sup>a</sup><http://www.rapidminer.com>

<sup>b</sup><http://www.jboss.org/drools/drools-expert.html>

---

# Organizar documentos

---

# O que fazer com grandes quantidades de documentos?

- Notícias, patentes, artigos, mensagens de twitter...
- **Para tirar proveito desta informação é necessário organizá-la de alguma forma:**
  - ★ Agrupamento de notícias, patentes, artigos e mensagens.
  - ★ Classificação, Recomendação e Filtragem de documentos (notícias, relatórios, mensagens do twitter, avaliação de itens).

# Exemplo de classificação/agrupamento

The screenshot shows the Google News Brazil interface. At the top, there are navigation links for Web, Imagens, Vídeos, Mapas, Notícias, Livros, Gmail, and mais. The user's email is fabricio.barth@gmail.com. The main search area includes 'Pesquisar notícias' and 'Pesquisar na web' buttons. Below the search area, there are links for 'Editar esta página' and 'Adicionar uma seção'. The main content area is titled 'Últimas notícias' and is updated 'há 5 minutos'. On the left, a navigation menu is highlighted with a red circle, containing links for 'Últimas notícias', 'Com estrela', 'Mundo', 'Brasil', 'Negócios', 'Ciência/Tecnologia', 'Entretenimento', 'Esportes', 'Saúde', and 'Mais populares'. Below this menu, there are links for 'Qualquer conteúdo', 'Manchetes', and 'Imagens'. A red arrow points to the 'Imagens' link. The main content area displays several news articles. The first article is 'Documentos do IOF e juros chineses ajudam a valorizar o dólar', dated 'há 36 minutos'. The second article is 'Polícia Federal prende 24 pessoas em operação contra golpes em MG', dated 'há 12 minutos'. The third article is 'Greves e manifestações desafiam reforma previdenciária na França', dated 'há 1 hora'. On the right side, there are more articles, including 'Petróleo cai abaixo de US\$ 80 após alta do juro na China', 'Preço torna banda larga proibitiva em países de baixa renda, diz UIT', 'Em nota, diretor de "Tropa de Elite" nega que tenha apoiado Dilma', 'Real Madrid 2x0 Milan – Placar modesto para o passeio no Bernabéu', 'MP-RJ apura morte de aposentada que aguardava leito', 'Margaret Thatcher é internada após apresentar quadro de infecção', and 'Delegado depõe e diz que versões de Mizael não batem'. At the bottom right, there is a section for 'Notícias em destaque' with links to 'Mércia Nakashima', 'Tribunal de Justiça', 'Paul McCartney', 'Nicolas Sarkozy', 'Wesliam Roriz', 'Ércio Quaresma', 'Vox Populi', 'Polícia Federal', 'Grand Prix', and 'Vanderlei Luxemburgo'.

---

## Exemplo de classificação/agrupamento

- E quando não é possível fazer manualmente?
- Qual é o processo para classificar e agrupar documentos de forma automática?

---

# Formato de um documento

... Esta disciplina tem como objetivo apresentar os principais conceitos da área de Inteligência Artificial, caracterizar as principais técnicas e métodos, e implementar alguns problemas clássicos desta área sob um ponto de vista introdutório.

A estratégia de trabalho, o conteúdo ministrado e a forma dependerão dos projetos selecionados pelos alunos.

Inicialmente, os alunos deverão trazer os seus Projetos de Conclusão de Curso, identificar intersecções entre o projeto e a disciplina, e propor atividades para a disciplina. ...

---

## Atributo/Valor usando vetores

Como representar os documentos?

$$\vec{d}_i = (p_{i1}, p_{i2}, \dots, p_{in}) \quad (1)$$

- Os atributos são as palavras que aparecem nos documentos.
- Se todas as palavras que aparecem nos documentos forem utilizadas, o vetor não ficará muito grande?



---

# Diminuindo a dimensionalidade do vetor

- Como filtrar as palavras que devem ser usadas como atributos?
- Em todos os idiomas existem átomos (palavras) que não significam muito. **Stop-words**

Esta disciplina **tem como** objetivo apresentar **os** principais conceitos **da** área **de** Inteligência Artificial, caracterizar **as** principais técnicas **e** métodos, **e** implementar alguns problemas clássicos **desta** área **sob um** ponto **de** vista introdutório.

...

---

## Diminuindo ainda mais a dimensionalidade do vetor

- Algumas palavras podem aparecer no texto de diversas maneiras: **técnica**, **técnicas**, **implementar**, **implementação**...
- **Stemming** - encontrar o radical da palavra e usar apenas o radical.

---

# Atributo/Valor usando vetores

- Já conhecemos os atributos.
- E os valores?
  - ★ **Booleana** - se a palavra aparece ou não no documento (1 ou 0)
  - ★ **Por frequência do termo** - a frequência com que a palavra aparece no documento (normalizada ou não)
  - ★ **Ponderação tf-idf** - o peso é proporcional ao número de ocorrências do termo no documento e inversamente proporcional ao número de documentos onde o termo aparece.

---

## Por freqüência do termo

(apresent,0.33) (form,0.33) (tecnic,0.33) (caracteriz,0.33)  
(projet,1.0) (introdutori,0.33) (objet,0.33) (inteligente,0.33)  
(conclusa,0.33) (selecion,0.33) (intersecco,0.33) (classic,0.33)  
(identific,0.33) (conceit,0.33) (trabalh,0.33) (disciplin,1.0)  
(traz,0.33)

---

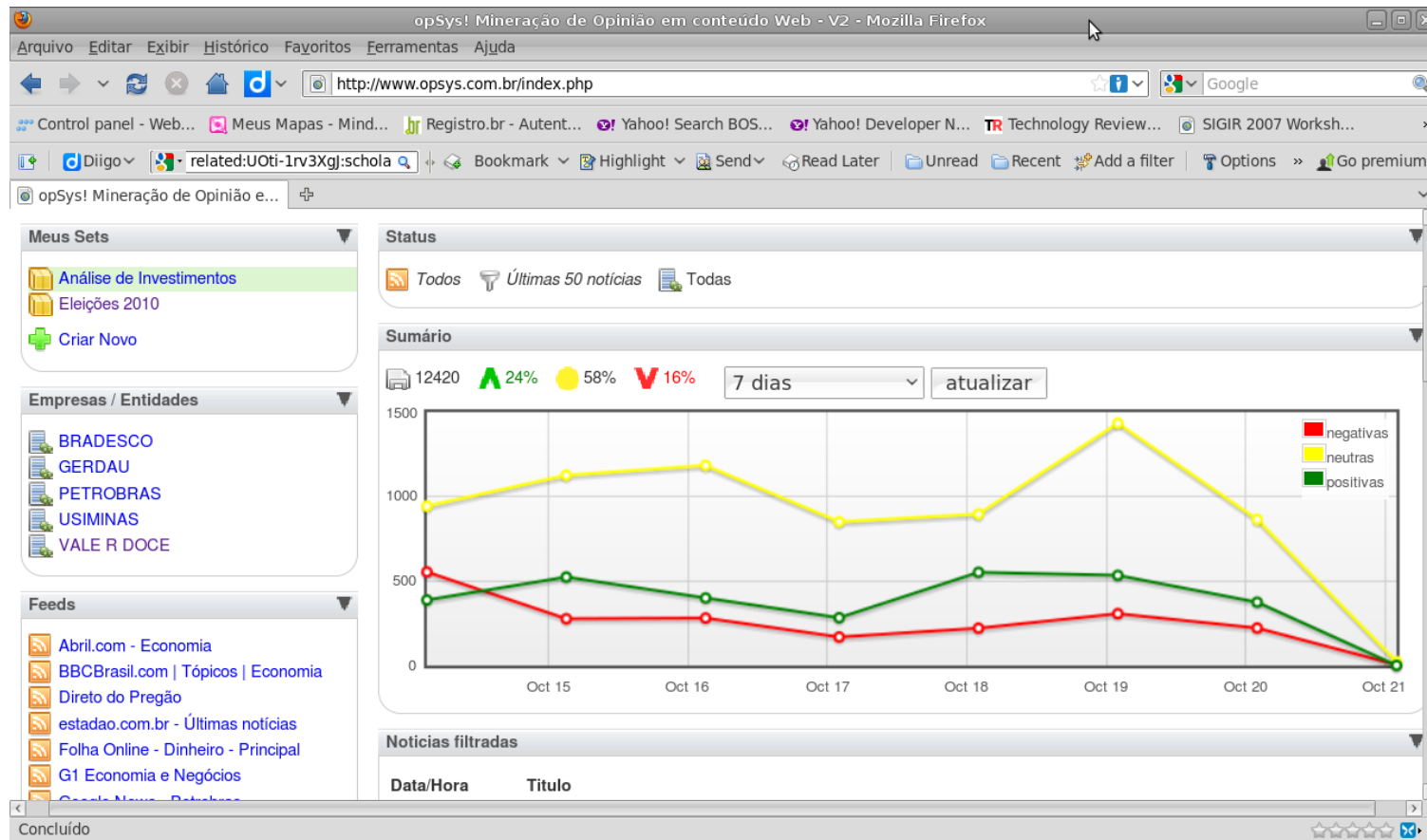
## Conjunto de Exemplos - Atributo/Valor

<b>Doc.</b>	<b>apresent</b>	<b>form</b>	<b>tecnic</b>	<b>caracteriz</b>	<b>...</b>
$d_1$	0.33	0.33	0.33	0.33	...
$d_2$	0	0.5	0.2	0.33	...
$d_3$	1	0.6	0	0	...
$d_4$	0.4	0.3	0.33	0.4	...
$d_5$	1	0.4	0.1	0.1	...
$d_n$	...	...	...	...	...

---

# Classificação de documentos

# Análise de Sentimento em mensagens no Twitter



Teor das mensagens sobre a Vale nos últimos sete dias.

---

# Conjunto de Exemplos Rotulados

<b>Doc.</b>	<b>Mensagem</b>	<b>Classe</b>
$d_1$	A empresa X é uma empresa muito séria	Positivo
$d_2$	O produto Y é uma porcaria	Negativo
$d_3$	Gostei muito da palestra de fulano	Positivo
$d_4$	Aquela praia é muito bonita	Positivo
$d_5$	Gostei daquele restaurante	Positivo
$d_n$		...



---

# Conjunto de Exemplos - Atributo/Valor e Classe

<b>Doc.</b>	<b>restaur</b>	<b>empres</b>	<b>bom</b>	<b>caracteriz</b>	<b>...</b>	<b>Classe</b>
$d_1$	0.33	0.33	0.33	0.33	...	Positivo
$d_2$	0	0.5	0.2	0.33	...	Negativo
$d_3$	1	0.6	0	0	...	Positivo
$d_4$	0.4	0.3	0.33	0.4	...	Positivo
$d_5$	1	0.4	0.1	0.1	...	Positivo
$d_n$	...	...	...	...	...	...

---

# Algoritmo Naïve Bayes

**NaiveBayesLearn(exemplos):**  $P'(v_j)$  e  $P'(a_i|v_j)$

**for all** valor alvo  $v_j$  **do**

$P'(v_j) \leftarrow$  estimar  $P(v_j)$

**for all** valor de atributo  $a_i$  de cada atributo  $a$  **do**

$P'(a_i|v_j) \leftarrow$  estimar  $P(a_i|v_j)$

**end for**

**end for**

**ClassifyNewInstance(x):**  $V_{NB}$

$$V_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_i P(a_i|v_j)$$

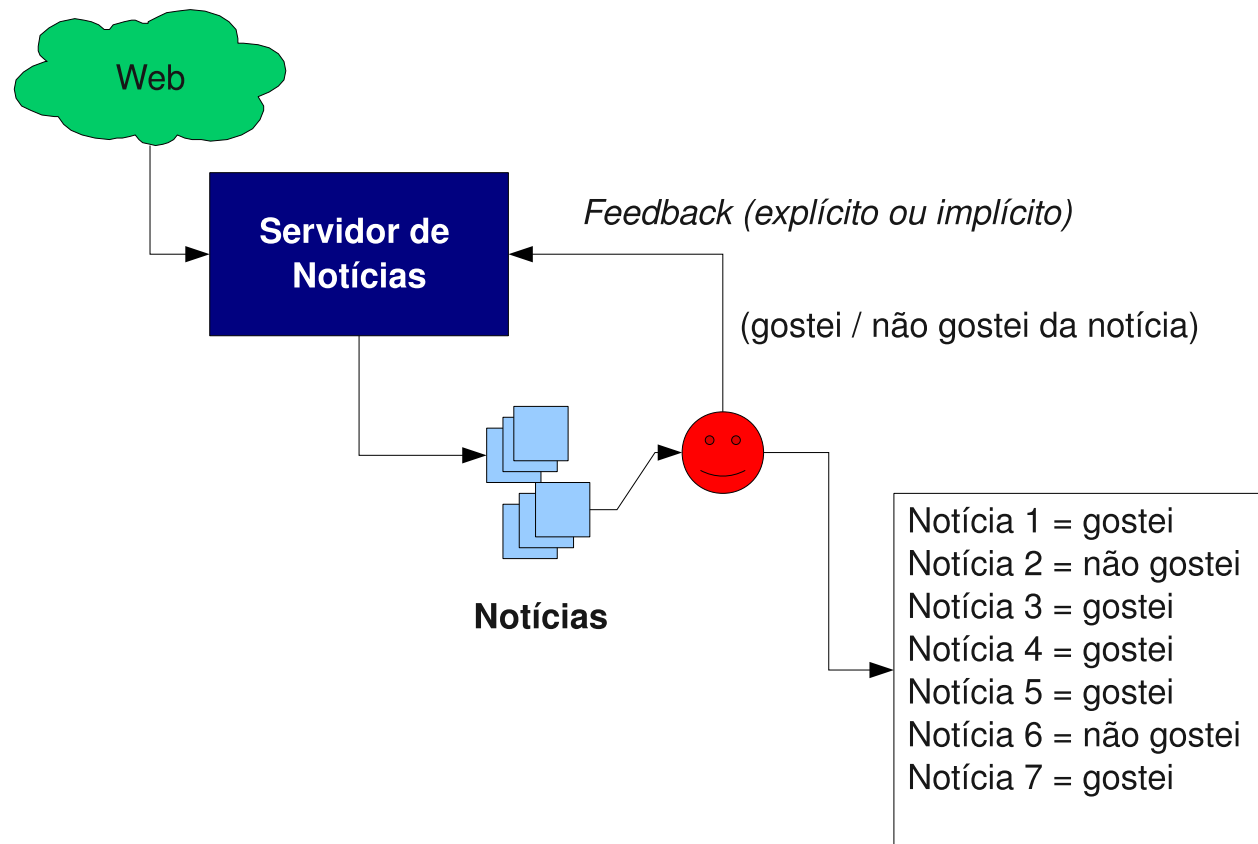
---

# Exemplo

## Execução de um demo usando estes conceitos...

Transformando o conjunto de treinamento em um vetor de palavras  
Criando o modelo  
Aplicando o modelo a novos casos  
Mensagem "Meu\_voto\_e\_para\_X,\_com\_certeza!" e classificada como POSITIVA  
Mensagem "Este\_produto\_e\_muito\_ruim" e classificada como NEGATIVA  
Mensagem "Nunca\_mais\_compro\_naquela\_loja!" e classificada como NEGATIVA  
Mensagem "Fulano\_e\_um\_mentiroso!" e classificada como NEGATIVA  
Mensagem "X\_lidera\_intenções\_de\_voto" e classificada como POSITIVA

# Classificação e Filtragem de Notícias

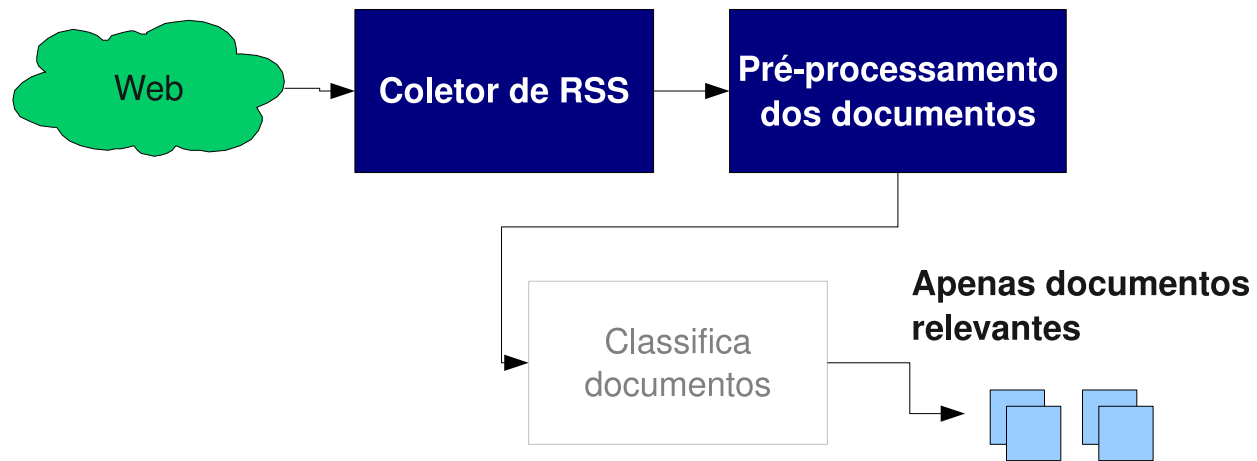


---

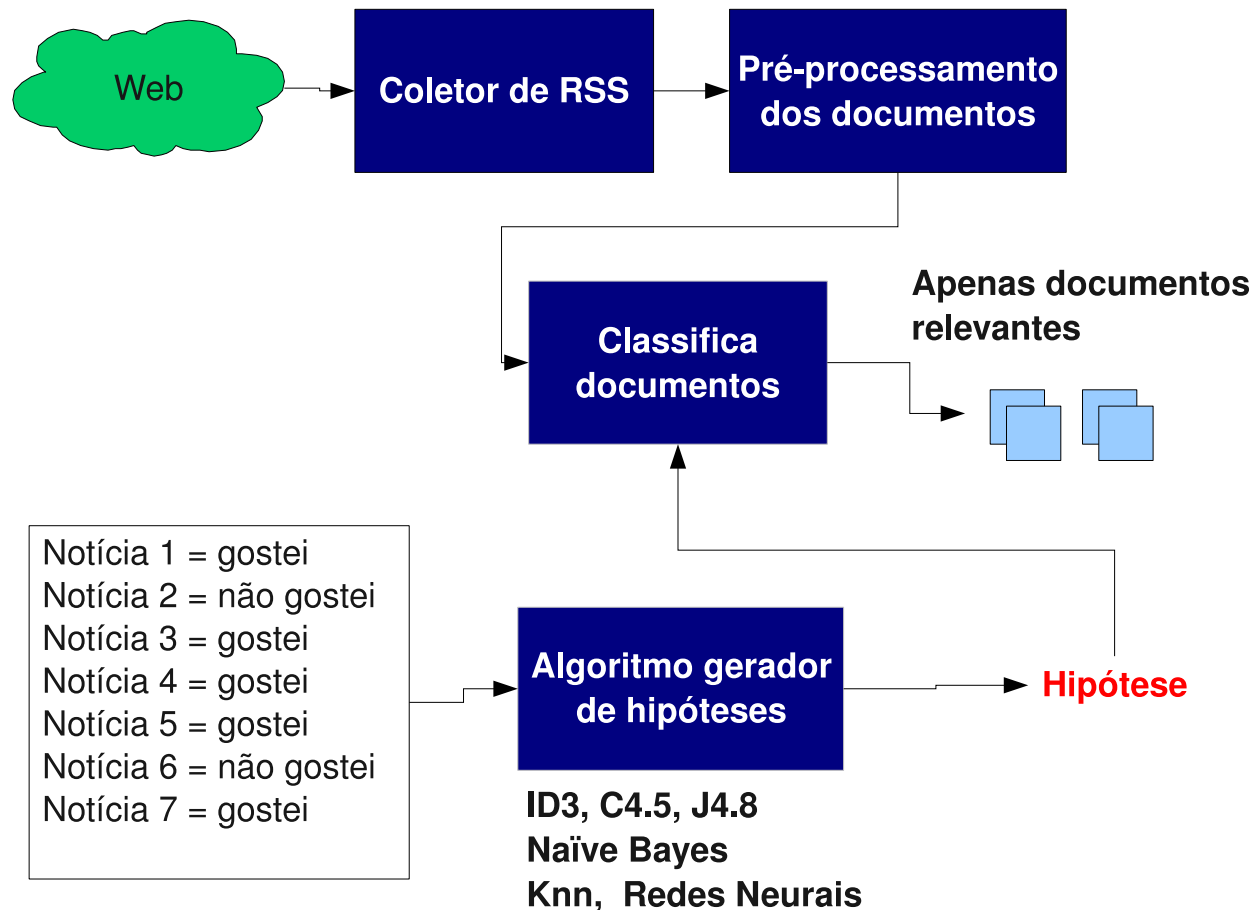
# Conjunto de Exemplos - Atributo/Valor e Classe

<b>Doc.</b>	<b>apresent</b>	<b>form</b>	<b>tecnic</b>	<b>caracteriz</b>	<b>...</b>	<b>Relevante</b>
$d_1$	0.33	0.33	0.33	0.33	...	1
$d_2$	0	0.5	0.2	0.33	...	0
$d_3$	1	0.6	0	0	...	1
$d_4$	0.4	0.3	0.33	0.4	...	1
$d_5$	1	0.4	0.1	0.1	...	1
$d_n$	...	...	...	...	...	...

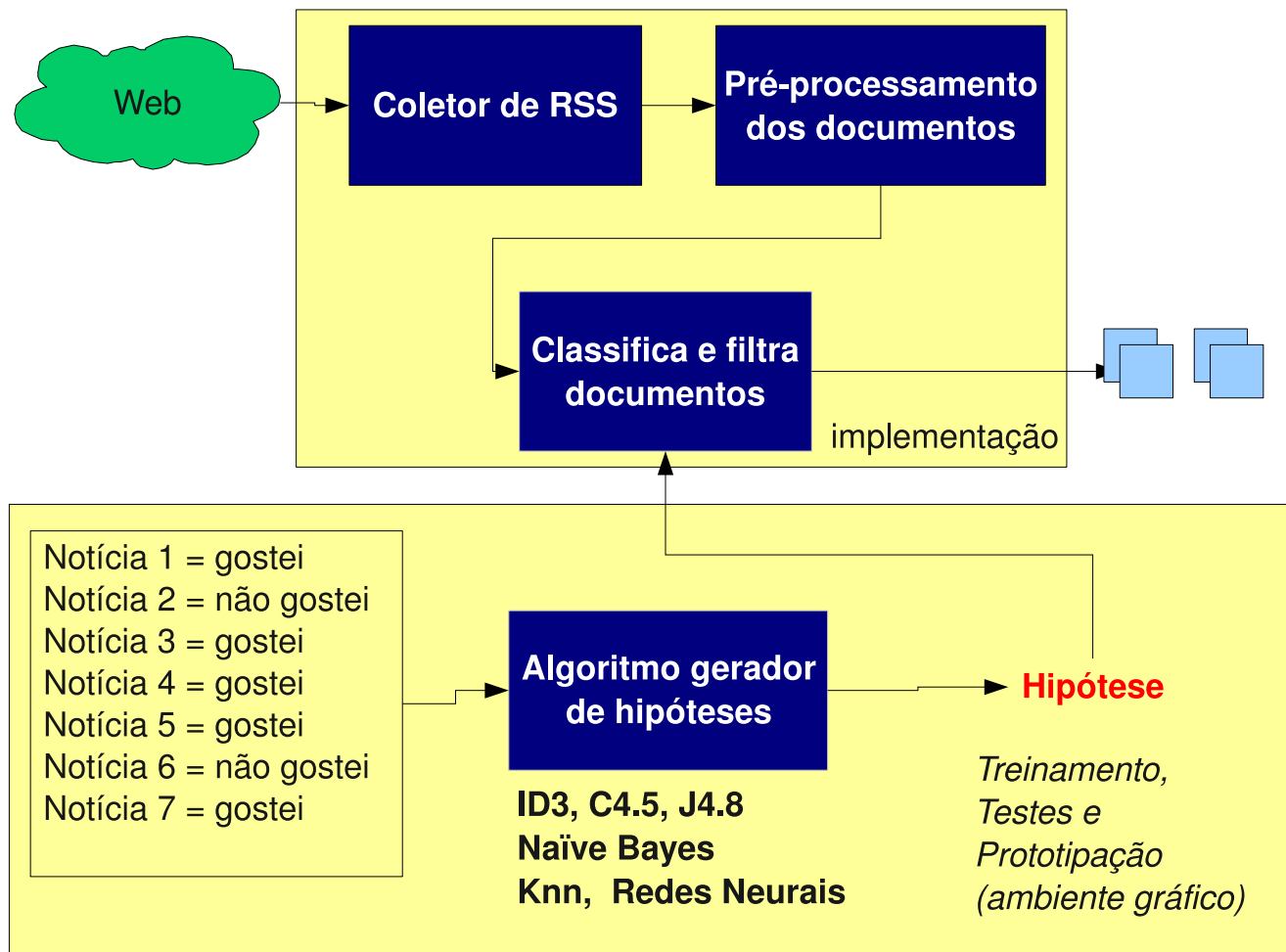
# Qual é o problema?



# Uma solução...



# Processo de trabalho





---

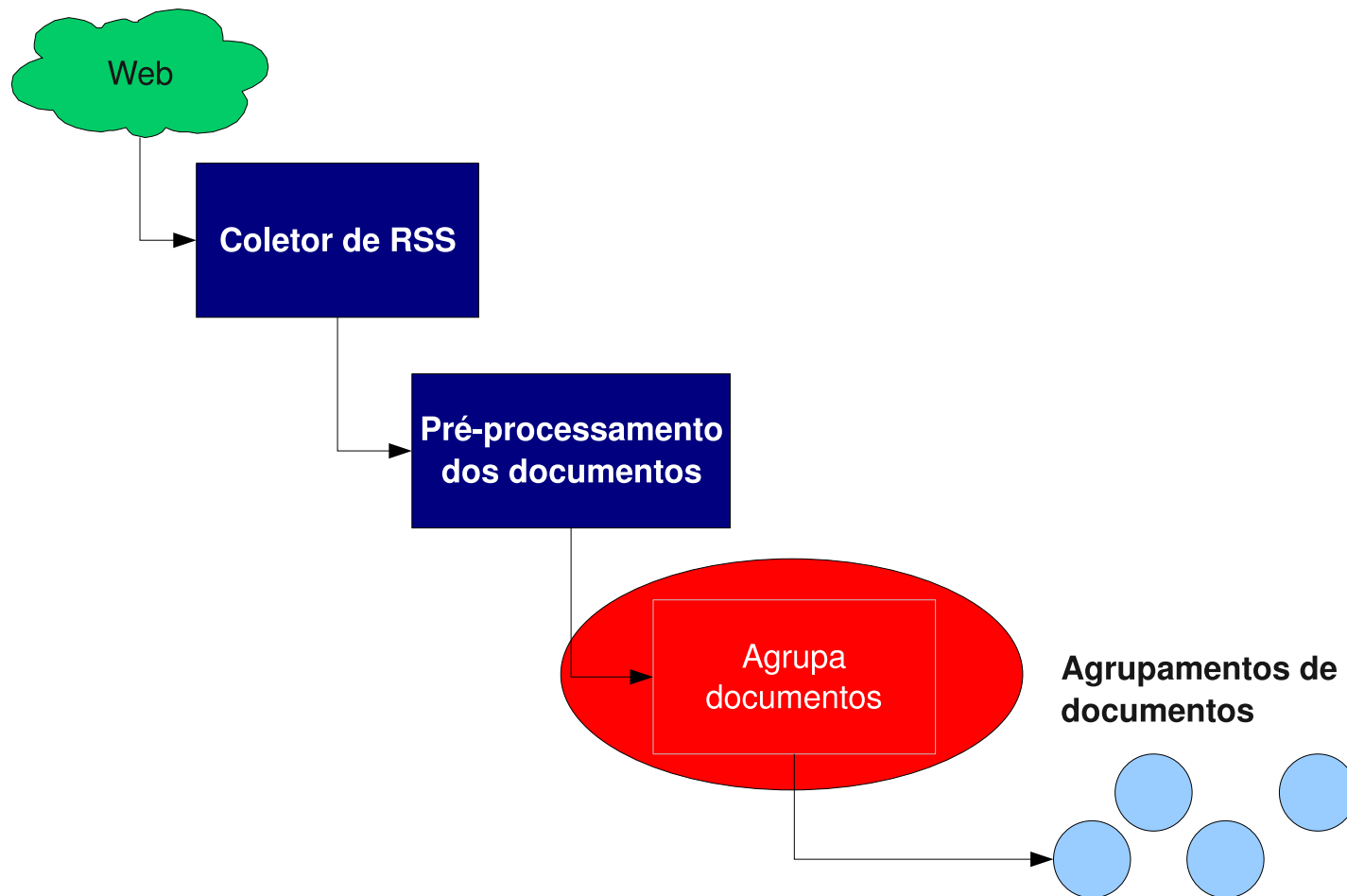
# Agrupamento de documentos

---

## Definições de Algoritmos de Agrupamento

- O objetivo dos algoritmos de agrupamento é colocar os objetos **similares** em um **mesmo grupo** e objetos **não similares** em **grupos diferentes**.
- Normalmente, objetos são descritos e agrupados usando um conjunto de **atributos e valores**.
- Não existe nenhuma informação sobre a classe ou categoria dos objetos.

# Componentes para uma solução...

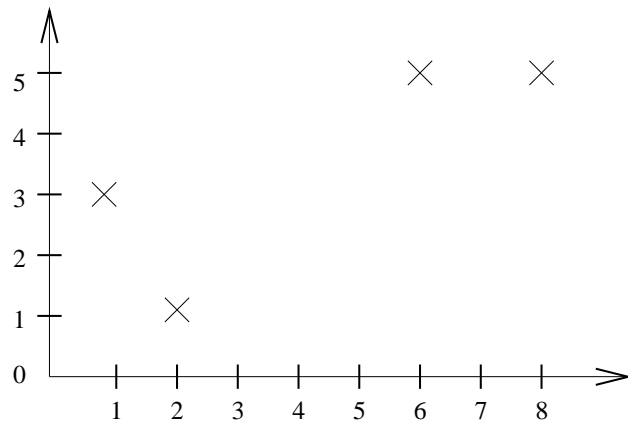


---

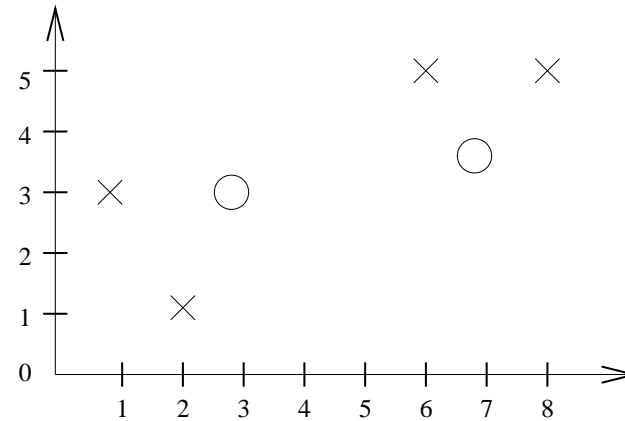
## Algoritmos para Agrupamento - *K-means*

- **K** significa o número de agrupamentos (que deve ser informado à priori).
- Sequência de ações **iterativas**.
- A parada é baseada em algum critério de qualidade dos agrupamentos (por exemplo, similaridade média).

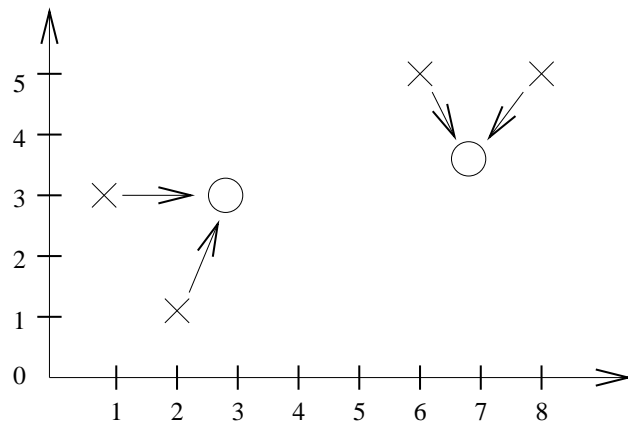
# Algoritmo para Agrupamento - *K-means*



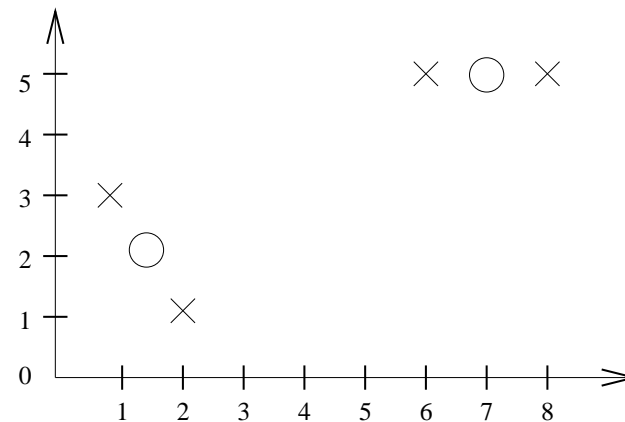
(1) Objetos que devem ser agrupados



(2) Sorteio dos pontos centrais dos agrupamentos



(3) Atribuição dos objetos aos agrupamentos



(4) Definição do centro do agrupamento

---

# Algoritmos para agrupamento dos documentos - WEKA

Execução do *K-means* no WEKA<sup>a</sup>.

---

<sup>a</sup><http://www.cs.waikato.ac.nz/ml/weka/>

---

# Algoritmo para agrupamento dos documentos - Resultados

```
A instância 0.1,0.1,0.1,0.1,0.1 está no cluster 1
A instância 0.1,0.2,0.3,0.1,0.8 está no cluster 1
A instância 0.3,0.4,0.5,0.8,0.9 está no cluster 0
A instância 0.3,0.1,0.1,0.1,0.1 está no cluster 1
A instância 0.3,0.1,0.1,0.1,0.1 está no cluster 1
A instância 0.8,0.7,0.8,0.8,0.8 está no cluster 0
A instância 0.1,0.1,0.1,0.1,0.1 está no cluster 1
A instância 0.1,0.1,0.1,0.1,0.1 está no cluster 1
A instância 0.1,0.1,0.1,0.1,0.1 está no cluster 1
A instância 0.6,0.5,0.6,0.6,0.6 está no cluster 0
A instância 0.6,0.5,0.6,0.6,0.6 está no cluster 0
A instância 0.1,0.1,0.1,0.1,0.1 está no cluster 1
A instância 0.2,0.8,0.8,0.7,0.9 está no cluster 0
A instância 0.1,0.1,0.1,0.1,0.1 está no cluster 1
```

---

# Minerando o log de um servidor Web



# Exemplo típico de log

1	2006-02-01 00:08:43 1.2.3.4 - GET /classes/cs589/papers.html - 200 9221 HTTP/1.1 maya.cs.depaul.edu Mozilla/4.0+(compatible;+MSIE+6.0;+Windows+NT+5.1;+SV1;+.NET+CLR+2.0.50727) http://dataminingresources.blogspot.com/
2	2006-02-01 00:08:46 1.2.3.4 - GET /classes/cs589/papers/cms-tai.pdf - 200 4096 HTTP/1.1 maya.cs.depaul.edu Mozilla/4.0+(compatible;+MSIE+6.0;+Windows+NT+5.1;+SV1;+.NET+CLR+2.0.50727) http://maya.cs.depaul.edu/~classes/cs589/papers.html
3	2006-02-01 08:01:28 2.3.4.5 - GET /classes/ds575/papers/hyperlink.pdf - 200 318814 HTTP/1.1 maya.cs.depaul.edu Mozilla/4.0+(compatible;+MSIE+6.0;+Windows+NT+5.1) http://www.google.com/search?hl=en&lr=&q=hyperlink+analysis+for+the+web+survey
4	2006-02-02 19:34:45 3.4.5.6 - GET /classes/cs480/announce.html - 200 3794 HTTP/1.1 maya.cs.depaul.edu Mozilla/4.0+(compatible;+MSIE+6.0;+Windows+NT+5.1;+SV1) http://maya.cs.depaul.edu/~classes/cs480/
5	2006-02-02 19:34:45 3.4.5.6 - GET /classes/cs480/styles2.css - 200 1636 HTTP/1.1 maya.cs.depaul.edu Mozilla/4.0+(compatible;+MSIE+6.0;+Windows+NT+5.1;+SV1) http://maya.cs.depaul.edu/~classes/cs480/announce.html
6	2006-02-02 19:34:45 3.4.5.6 - GET /classes/cs480/header.gif - 200 6027 HTTP/1.1 maya.cs.depaul.edu Mozilla/4.0+(compatible;+MSIE+6.0;+Windows+NT+5.1;+SV1) http://maya.cs.depaul.edu/~classes/cs480/announce.html

# Pré-processamento do log: identificação de usuários

Time	IP	URL	Ref	Agent
0:01	1.2.3.4	A	-	IE5;Win2k
0:09	1.2.3.4	B	A	IE5;Win2k
0:10	2.3.4.5	C	-	IE6;WinXP;SP1
0:12	2.3.4.5	B	C	IE6;WinXP;SP1
0:15	2.3.4.5	E	C	IE6;WinXP;SP1
0:19	1.2.3.4	C	A	IE5;Win2k
0:22	2.3.4.5	D	B	IE6;WinXP;SP1
0:22	1.2.3.4	A	-	IE6;WinXP;SP2
0:25	1.2.3.4	E	C	IE5;Win2k
0:25	1.2.3.4	C	A	IE6;WinXP;SP2
0:33	1.2.3.4	B	C	IE6;WinXP;SP2
0:58	1.2.3.4	D	B	IE6;WinXP;SP2
1:10	1.2.3.4	E	D	IE6;WinXP;SP2
1:15	1.2.3.4	A	-	IE5;Win2k
1:16	1.2.3.4	C	A	IE5;Win2k
1:17	1.2.3.4	F	C	IE6;WinXP;SP2
1:26	1.2.3.4	F	C	IE5;Win2k
1:30	1.2.3.4	B	A	IE5;Win2k
1:36	1.2.3.4	D	B	IE5;Win2k

User 1

0:01	1.2.3.4	A	-
0:09	1.2.3.4	B	A
0:19	1.2.3.4	C	A
0:25	1.2.3.4	E	C
1:15	1.2.3.4	A	-
1:26	1.2.3.4	F	C
1:30	1.2.3.4	B	A
1:36	1.2.3.4	D	B

User 2

0:10	2.3.4.5	C	-
0:12	2.3.4.5	B	C
0:15	2.3.4.5	E	C
0:22	2.3.4.5	D	B

User 3

0:22	1.2.3.4	A	-
0:25	1.2.3.4	C	A
0:33	1.2.3.4	B	C
0:58	1.2.3.4	D	B
1:10	1.2.3.4	E	D
1:17	1.2.3.4	F	C

---

# Pré-processamento do log: identificação das seções

User 1

Time	IP	URL	Ref
0:01	1.2.3.4	A	-
0:09	1.2.3.4	B	A
0:19	1.2.3.4	C	A
0:25	1.2.3.4	E	C
1:15	1.2.3.4	A	-
1:26	1.2.3.4	F	C
1:30	1.2.3.4	B	A
1:36	1.2.3.4	D	B

Session 1

0:01	1.2.3.4	A	-
0:09	1.2.3.4	B	A
0:19	1.2.3.4	C	A
0:25	1.2.3.4	E	C

Session 2

1:15	1.2.3.4	A	-
1:26	1.2.3.4	F	C
1:30	1.2.3.4	B	A
1:36	1.2.3.4	D	B

# Matriz de transações

**Pageviews**

	<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>	<b>F</b>
<b>user0</b>	15	5	0	0	0	185
<b>user1</b>	0	0	32	4	0	0
<b>user2</b>	12	0	0	56	236	0
<b>user3</b>	9	47	0	0	0	134
<b>user4</b>	0	0	23	15	0	0
<b>user5</b>	17	0	0	157	69	0
<b>user6</b>	24	89	0	0	0	354
<b>user7</b>	0	0	78	27	0	0
<b>user8</b>	7	0	45	20	127	0
<b>user9</b>	0	38	57	0	0	15

---

# Regras de Associação

- **Caso do supermercado** (fralda  $\rightarrow$  cerveja)
- Quem acessa a página sobre futebol também acessa a página de volei em **90%** dos casos (futebol  $\rightarrow$  volei).
- Quem acessa a página de ofertas e a página de material de construção também finaliza a compra em **83%** dos casos (ofertas  $\wedge$  material\_construção  $\rightarrow$  compra)

---

# Considerações Finais

---

# Considerações Finais

- Foram vistos: problemas de classificação, agrupamento e análise de log. Tem muito mais de onde vieram estes...
- **Atenção para o processo!** Pré-processamento, criação dos modelos, avaliação e aplicação.
- Alguns algoritmos para mineração de informação são pesados. Talvez, parte da solução esteja na adoção de *cloud computing*.
- **Muitos dados... Muitas oportunidades...**

---

# Outros exemplos

- Wiki2Group<sup>a</sup> - 2010
- Sistema Folkaliza<sup>b</sup> - 2009
- Sistema opSys<sup>c</sup> - 2008
- Sistema FaroFino - 2005
- Mais informações podem ser encontradas em  
**<http://fbarth.net.br>** e  
**<http://fbarth.net.br/projetos/riInteligente.html>**

---

<sup>a</sup><http://trac.fbarth.net.br/wikiAnalysis>

<sup>b</sup><http://www.jessicacintra.com.br/jeh/folkaliza/Home/Default.aspx>

<sup>c</sup><http://www.opsys.com.br>



---

# Referências

# References

[Data, data everywhere. A special report on managing information 20  
DATA, data everywhere. A special report on managing information.  
mation. **The Economist**, p. 1–16, February 2010.

[Liu 2009]LIU, B. **Web Data Mining: Exploring Hyperlinks, Contents, and Usage Data (Data-Centric Systems and Applications)**. 1st ed. 2007. corr. 2nd printing. ed. Springer, 2009. Hardcover. ISBN 3540378812. Disponível em: <http://www.amazon.com/exec/obidos/redirect?tag=citeulike07-20&path=ASIN/3540378812>.

[Mitchell 1997]MITCHELL, T. M. **Machine Learning**. [S.I.]: McGraw-Hill, 1997.

[Quinlan 1988]QUINLAN, J. R. Knowledge acquisition for knowledge-based systems. In: \_\_\_\_\_. [S.I.]: Academic Press, 1988. cap. Simplifying Decision Trees.

[Russel e Norvig 2003]RUSSEL, S. J.; NORVIG, P. **Artificial intelligence: a modern approach**. 2. ed. [S.I.]: Prentice-Hall, 2003. ISBN 0-13-790395-2.

[Witten e Frank 2005]WITTEN, I. H.; FRANK, E. **Data Mining: Practical Machine Learning Tools and Techniques**. Second. [S.I.]: Elsevier, 2005.