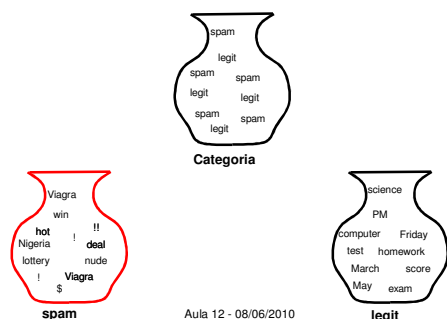




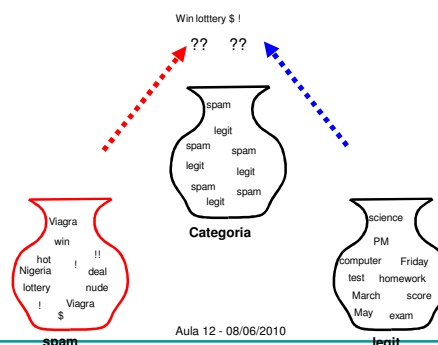
## Modelo Naïve Bayes para textos



Aula 12 - 08/06/2010

7

## Classificação Naïve Bayes



Aula 12 - 08/06/2010

8

## Algoritmo Naive Bayes para Textos (Treinamento)

Seja  $D$  um conjunto de documentos  
 Seja  $V$  o vocabulário de todas as palavras nos documentos de  $D$   
 Para cada classe  $c_i \in C$   
 Seja  $D_i$  o subconjunto de documentos em  $D$  que pertencem à categoria  $c_i$   
 $P(c_i) = |D_i| / |D|$   
 Seja  $T_i$  a concatenação de todos os documentos em  $D_i$   
 Seja  $n_i$  o número total de ocorrências de palavras em  $T_i$   
 Para cada palavra  $w_j \in V$   
 Seja  $n_{ij}$  o número de ocorrências de  $w_j$  em  $T_i$   
 Let  $P(w_{ij} | c_i) = (n_{ij} + 1) / (n_i + |V|)$

Aula 12 - 08/06/2010

9

## Algoritmo Naive Bayes para Textos (Teste)

Dado um documento de teste  $X$   
 Seja  $n$  o número de ocorrências de palavras em  $X$   
 Retorne a classe:

$$\operatorname{argmax}_{c_i \in C} P(c_i) \prod_{i=1}^n P(a_i | c_i)$$

onde  $a_i$  é a palavra que ocorre na  $i$ -ésima posição de  $X$

Aula 12 - 08/06/2010

10

## Prevenção de *Underflow*

- Multiplicar muitas probabilidades, que estão entre 0 e 1, pode resultar num *underflow* de ponto flutuante.
- Como  $\log(xy) = \log(x) + \log(y)$ , é melhor fazer todos os cálculos somando logs de probabilidades ao invés de multiplicar probabilidades.
- Classe com maior valor de log-probabilidade é também a mais provável na escala normal.

Aula 12 - 08/06/2010

11

## Métricas de Similaridade de Texto

- Medir a similaridade de textos é um problema bastante estudado.
- Métricas são baseadas no modelo "*bag of words*".
- Normalmente é feito um pré-processamento: "*stop words*" são removidas e as palavras são reduzidas à sua raiz morfológica.
- Modelo vetorial de Recuperação de Informação (IR) é a abordagem padrão.

Aula 12 - 08/06/2010

12

## O modelo vetorial

- Supõe-se que  $t$  termos distintos restam após o pré-processamento; chamados de termos do vocabulário.
- Estes termos “ortogonais” formam um espaço vetorial.  
Dimensão =  $t = |\text{vocabulário}|$
- Cada termo,  $i$ , num documento ou consulta,  $j$ , tem um peso dado por um número real,  $w_{ji}$ .
- Tanto documentos quanto consultas são representados por vetores  $t$ -dimensionais:  
 $d_j = (w_{1j}, w_{2j}, \dots, w_{tj})$

Aula 12 - 08/06/2010

13

## Representação gráfica

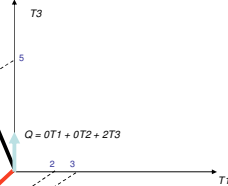
Exemplo:

$$D_1 = 2T_1 + 3T_2 + 5T_3$$

$$D_2 = 3T_1 + 7T_2 + T_3$$

$$Q = 0T_1 + 0T_2 + 2T_3$$

$$D_1 = 2T_1 + 3T_2 + 5T_3$$



- Quem é mais similar a Q? D1 or D2?
- Como medir o grau de similaridade? Distância? Ângulo? Projeção?

Aula 12 - 08/06/2010

14

## Coleção de Documentos

- Uma coleção de  $n$  documentos pode ser representada no modelo vetorial por uma matriz.
- Uma entrada na matriz corresponde ao “**peso**” do termo no documento; zero indica que o termo não é significativo no documento ou simplesmente não existe no documento.

$$\begin{pmatrix} D_1 & T_1 & T_2 & \dots & T_t \\ w_{11} & w_{21} & \dots & w_{t1} \\ D_2 & w_{12} & w_{22} & \dots & w_{t2} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ D_n & w_{1n} & w_{2n} & \dots & w_{tn} \end{pmatrix}$$

Aula 12 - 08/06/2010

15

## Pesos: Frequência dos Termos

- Termos frequentes em um documento são mais importantes, i.e. mais indicativos do tópico do documento.

$$f_{ij} = \text{frequência do termo } i \text{ no documento } j$$

- Podemos obter a *frequência do termo* ( $tf$ ) dividindo  $f$  pela frequência do termo mais comum no documento:

$$tf_{ij} = f_{ij} / \max_k \{f_{ik}\}$$

Aula 12 - 08/06/2010

16

## Pesos: Frequência Inversa dos Documentos

- Termos que aparecem em muitos documentos *diferentes* são *menos* significativos.  
 $df_i$  = frequência em documentos do termo  $i$   
= número de documentos contendo o termo  $i$   
 $idf_i$  = frequência inversa em documentos do termo  $i$ ,  
=  $\log_2 (N / df_i)$   
( $N$ : número total de documentos)
- É uma indicação do *poder de discriminação* do termo.
- Log é usado para diminuir o efeito em relação a  $tf$ .

Aula 12 - 08/06/2010

17

## Ponderação TF-IDF

- Uma ponderação tipicamente utilizada é:  
 $w_{ij} = tf_{ij} idf_i = tf_{ij} \log_2 (N / df_i)$
- Um termo que ocorre com frequência no documento mas raramente no resto da coleção tem peso maior.
- Muitas outras formas de ponderação foram propostas.
- Experimentalmente, determinou-se que a ponderação *tf-idf* funciona bem.

Aula 12 - 08/06/2010

18

## Medida de Similaridade de Cosseno

- Mede o cosseno do ângulo entre dois vetores.
- Produto interno normalizado pelo comprimento dos vetores.

$$\text{CosSim}(\vec{d}, \vec{q}) = \frac{\vec{d} \cdot \vec{q}}{|\vec{d}| \cdot |\vec{q}|} = \frac{\sum_{i=1}^n (w_{ij} \cdot w_{iq})}{\sqrt{\sum_{i=1}^n w_{ij}^2} \cdot \sqrt{\sum_{i=1}^n w_{iq}^2}}$$

$$\begin{aligned} D1 &= 2T1 + 3T2 + 5T3 & \text{CosSim}(D1, Q) &= 10 / \sqrt{(4+9+25)(0+0+4)} = 0.81 \\ D2 &= 3T1 + 7T2 + 1T3 & \text{CosSim}(D2, Q) &= 2 / \sqrt{(9+49+1)(0+0+4)} = 0.13 \\ Q &= 0T1 + 0T2 + 2T3 \end{aligned}$$

*D1* é 6 vezes melhor que *D2* usando similaridade de cosseno mas só 5 vezes melhor usando produto interno.

Aula 12 - 08/06/2010

19

## K-NN para Textos

### Treinamento:

Para cada exemplo de treinamento  $\langle x, c(x) \rangle \in D$

Calcule o vetor TF-IDF correspondente,  $\mathbf{d}_x$ , para o documento  $x$

### Exemplo de teste $y$ :

Calcule o vetor TF-IDF  $\mathbf{d}$  para o documento  $y$

Para cada  $\langle x, c(x) \rangle \in D$

Seja  $s_x = \text{cosSim}(\mathbf{d}, \mathbf{d}_x)$

Ordene os exemplos,  $x$ , em  $D$  por valor decrescente de  $s_x$

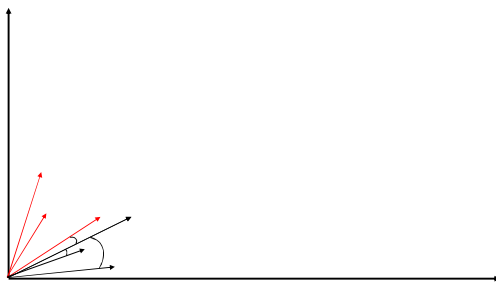
Seja  $N$  o conjunto dos primeiros  $k$  exemplos de  $D$ .

Retorne a classe majoritária dos exemplos em  $N$ .

Aula 12 - 08/06/2010

20

## Exemplo: 3-NN para Textos



Aula 12 - 08/06/2010

21

## Índice Invertido

- Busca linear na base de treinamento não é escalável.
- Índice invertido: estrutura mapeando palavras a documentos.
- Quando as *stopwords* são removidas, as palavras que sobram são raras, então um índice invertido ajuda a eliminar boa parte dos documentos que não tem muitas palavras em comum com o documento de teste.

Aula 12 - 08/06/2010

22

## Conclusões

- Existem muitas aplicações importantes da classificação de textos.
- Requer uma técnica que lide bem com vetores esparsos de muitos atributos, porque tipicamente cada palavra é um atributo e a maioria das palavras é rara.
  - Naïve Bayes
  - kNN com similaridade de cosseno
  - SVMs

Aula 12 - 08/06/2010

23