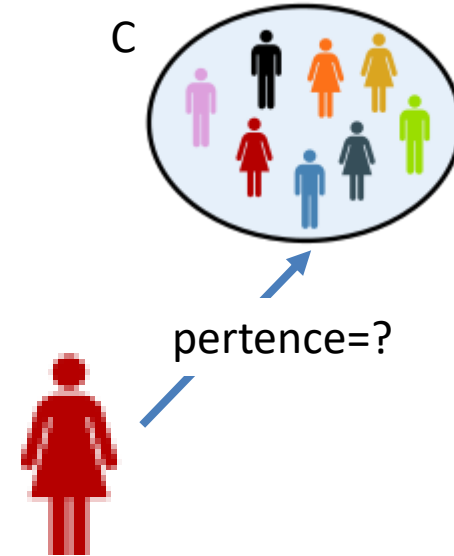


MPEI 2023-2024

Solução Probabilística
para o
Problema da Pertença a um Conjunto

Definição do problema

- Em termos gerais o problema pode ser colocado da seguinte forma:
- Dado um **elemento e** e um **conjunto C** ,
 e pertence a C ?
- O elemento pode ser, por exemplo, uma String



Problema para conjuntos de grandes dimensões

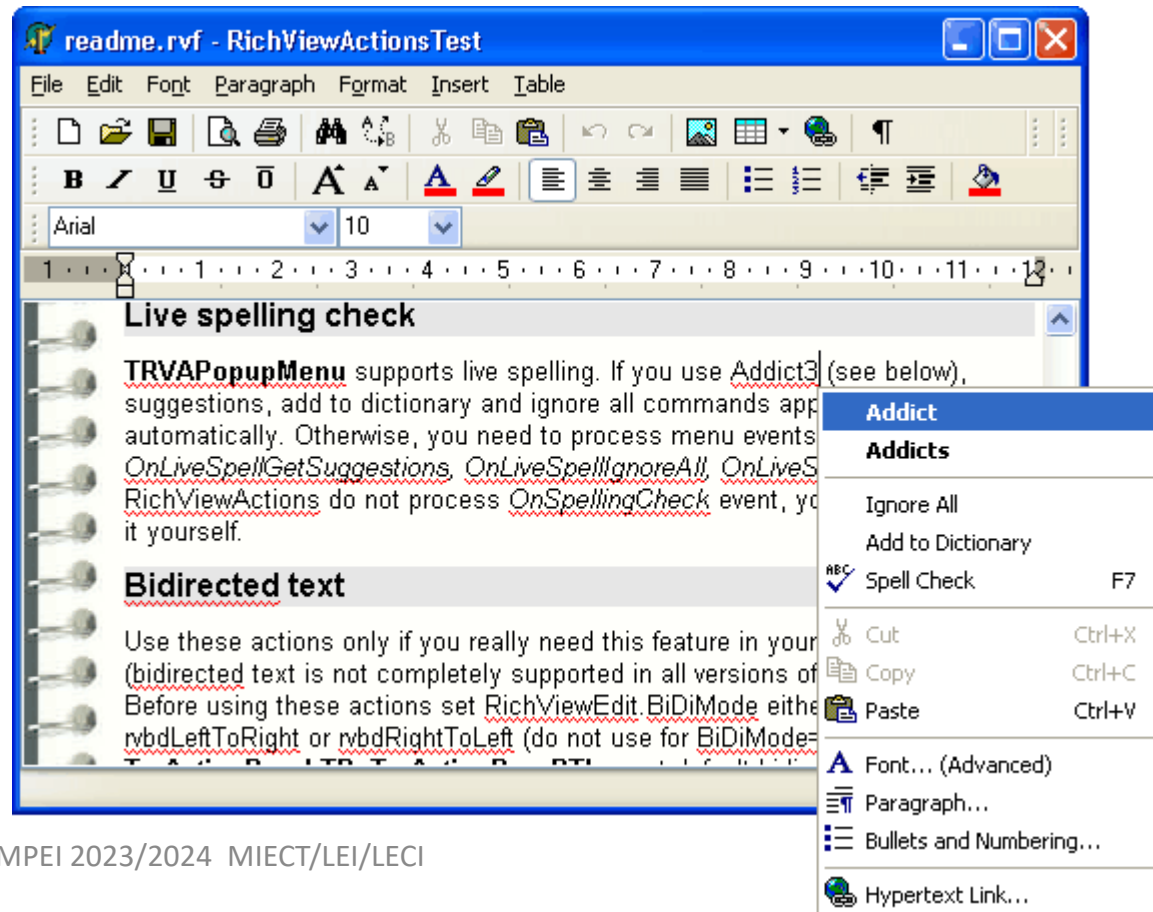
- A resolução deste tipo de problemas para pequenos conjuntos é fácil
 - usando, por exemplo, *hash tables*
- Para **conjuntos de dimensão muito grande** (e apenas passíveis de serem definidos em extensão) não é assim tão simples
 - pode mesmo não haver memória suficiente para armazenar todos os elementos de C
- **Soluções probabilísticas** são necessárias e interessantes
 - cada vez mais relevantes com o aumento de Big Data

Problema 1 - Verificação ortográfica

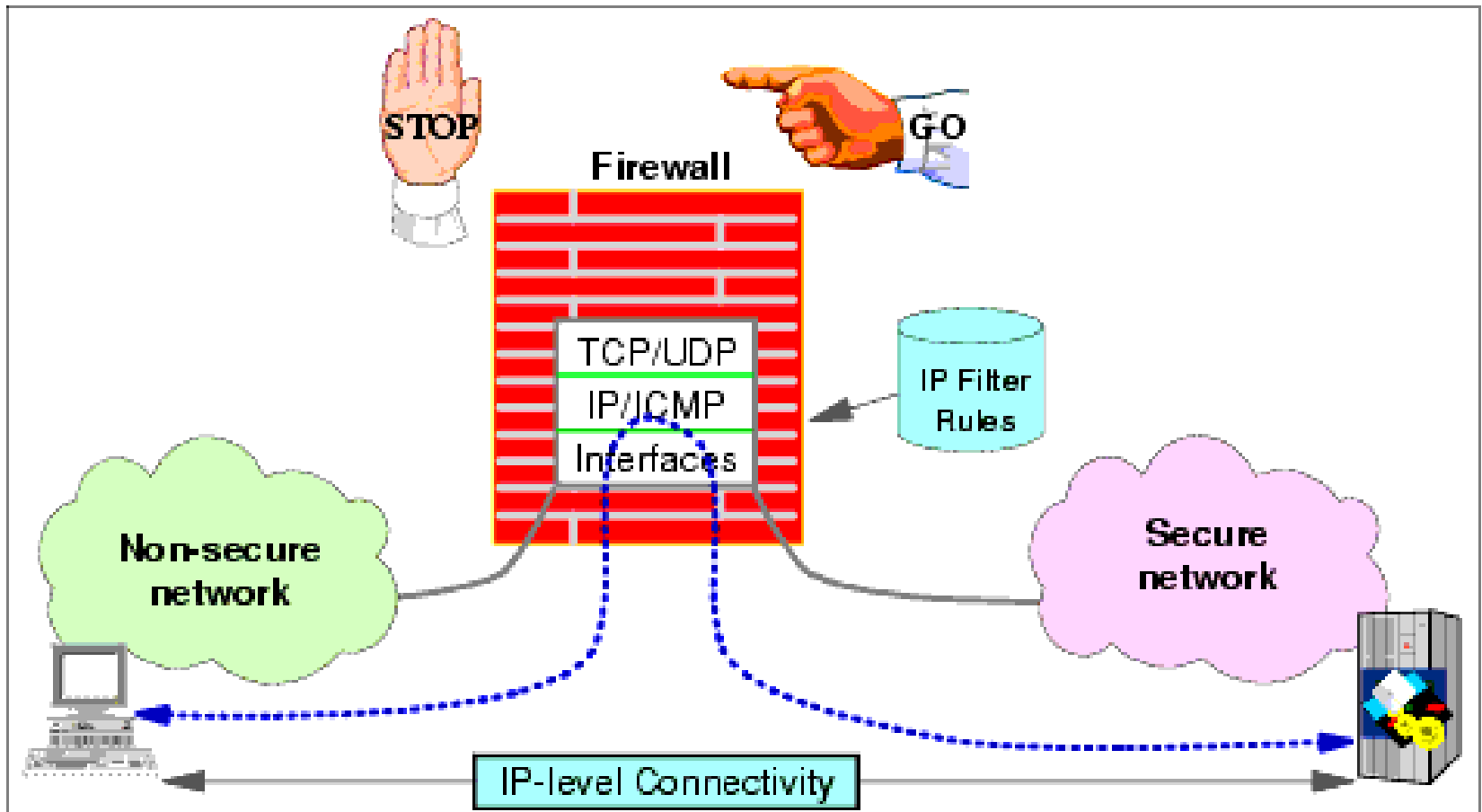


pertence ao dicionário?

- Utiliza dicionário(s)
 - Muitas vezes complementados por regras
- Ortografia validada através da pertença a dicionário
 - Que deve ser de confiança



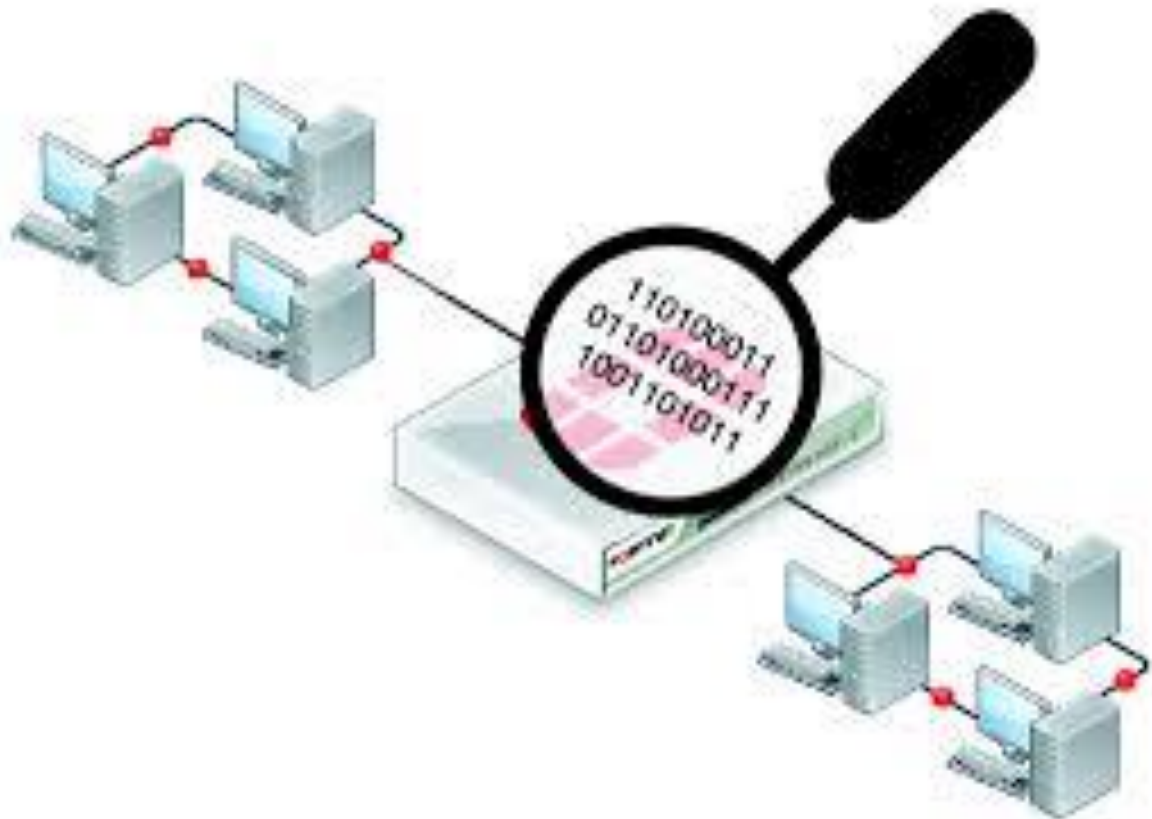
Problema 2 - Firewall



Problema 3

- Detetar sequências de bits (ex: Strings)

Detetar sequências de bits...



Generalizando ...

- Na área da Informática temos muitos outros exemplos
 - Combate ao **SPAM**: pertença a uma lista de emails seguros
 - **Browser**: pertença à lista de endereços que já visitámos
 - ...
- Em muitas aplicações e situações temos de ter uma **forma eficiente** de **determinar se um determinado item pertence ou não a um conjunto**

↳ Mesmo com milhões de elementos...

Ideia base

- Em muitos problemas **apenas pretendemos saber se um elemento pertence ou não ao conjunto**
 - sem necessitarmos de acesso ao conjunto ou informação associada aos elementos
- Nestas situações **podemos eliminar a parte de armazenamento dos elementos**
 - guardando apenas informação de que existe

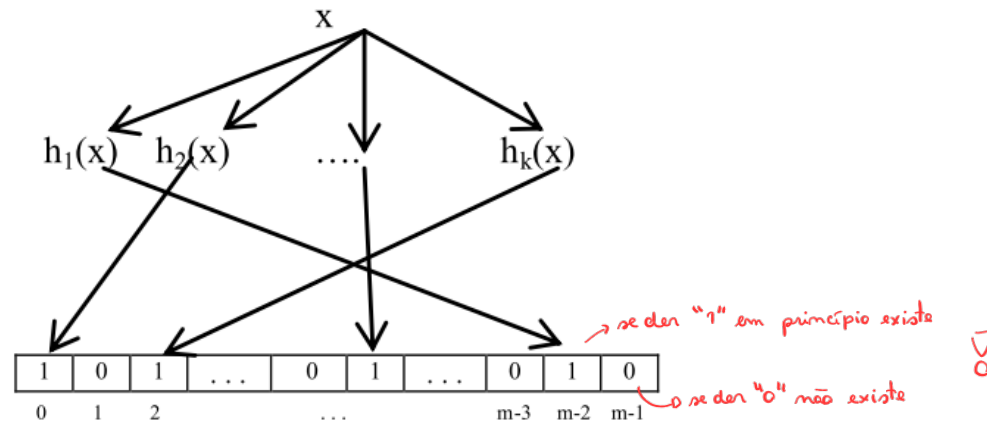
Filtros de Bloom

(Bloom Filters)

Os Filtros de Bloom são uma forma de usar funções de dispersão para determinar se um elemento pertence a um conjunto

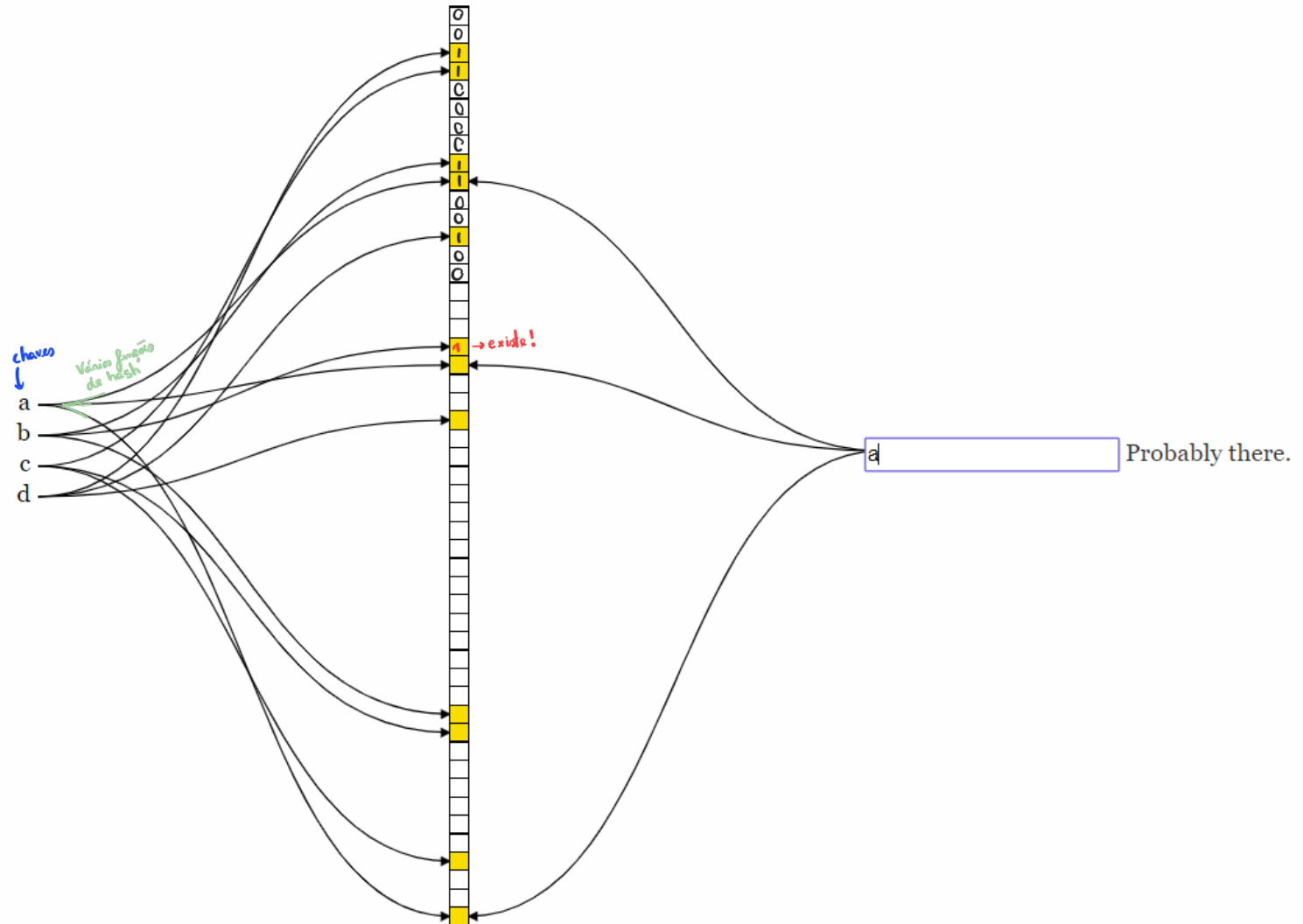
Filtros de Bloom

- Os Filtros de Bloom **usam funções de dispersão para calcular um vetor (o filtro)**
 - que é representativo do conjunto



- A **pertença** ao conjunto é testada através da comparação dos resultados da aplicação das mesmas funções de dispersão aos potenciais membros com o conteúdo desse vetor

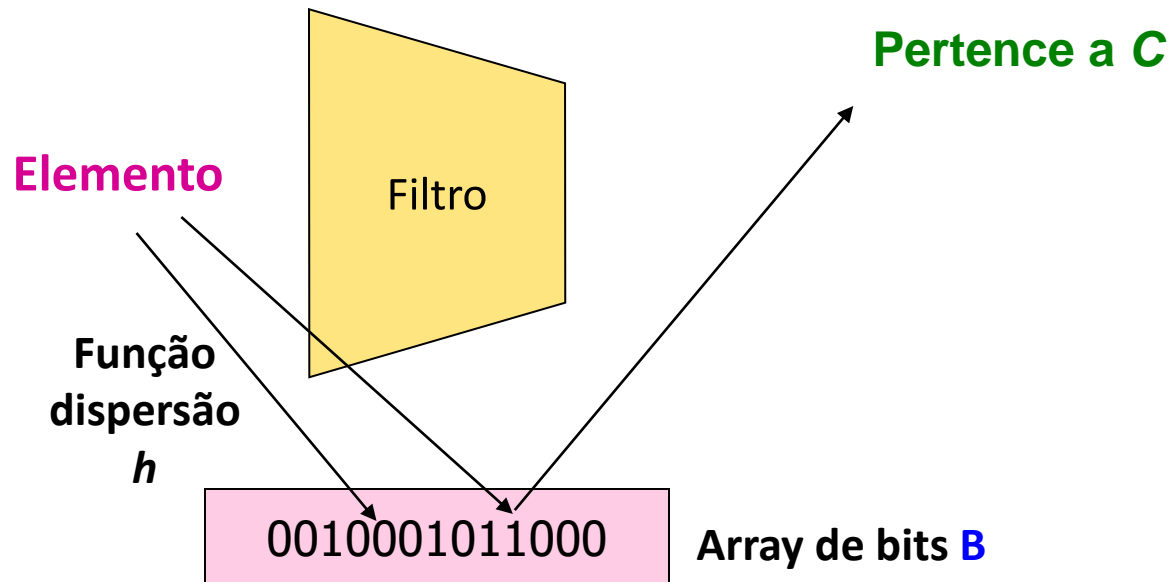
Exemplo de Filtro de Bloom



Filtros de Bloom

- Na sua **forma mais simples** o vetor (filtro) é composto por n posições
 - cada uma de apenas 1 bit
- O bit é apenas colocado a 1 se a função de dispersão mapear nessa posição algum dos elementos do conjunto
- São rápidos, de complexidade temporal constante
- não incorporam qualquer tentativa de resolução de colisões

Filtro



Descartar o elemento

Se é mapeado pela função de dispersão para uma posição contendo 0

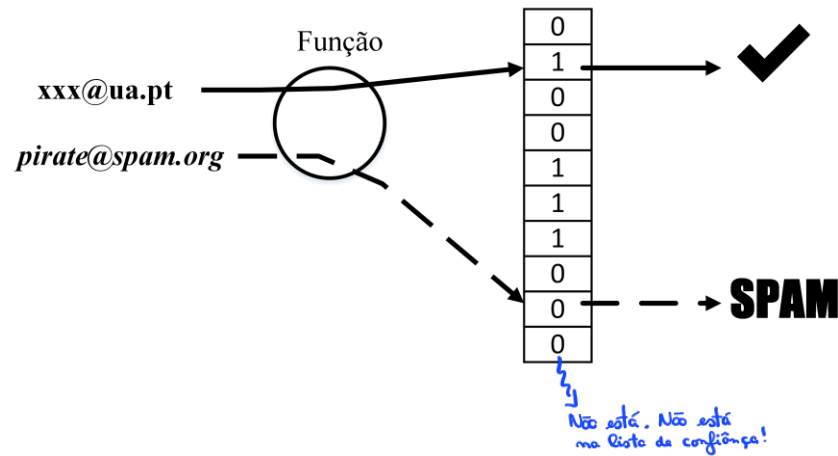
Exemplo – Detecção de SPAM

- Conhecem-se mil milhões (10^9) **de endereços de email de confiança** *→ procurar se um email fog parte de uma lista,,*
 - que constituem o nosso conjunto C
- Se uma mensagem é proveniente de um destes endereços não é SPAM
- Pretende-se filtrar um conjunto de emails recebidos

Exemplo – Detecção de SPAM

- Uma solução consiste em:

1. Criar um vetor de n bits B bastante grande e inicializando todos esses bits com 0
2. Utilizar uma função de dispersão para mapear cada endereço de email numa posição desse vetor
3. Colocar a 1 os bits correspondentes à aplicação da referida função de dispersão a toda a lista de endereços de email “bons” (todos os elementos do conjunto C)
4. Aplicar a função de dispersão ao endereço de cada uma das mensagens que se pretende verificar, considerando SPAM todas as em que o vetor tem 0 na posição correspondente



Ausência de falsos negativos

- Se um endereço que verificamos pertence a C então ele vai ser certamente mapeado pela função de dispersão numa posição do vetor que contém 1
 - Será sempre considerado pertencente a C

Podem haver uma $K_2 : h_{1...K}(K_1) = h_{1...K}(K_2), K_1 \neq K_2$

- Todos os emails de confiança (que integram a lista) serão sempre considerados de confiança
nunca havendo falsos negativos

↳ Se o filtro disser que está na lista, ele não falha!

Generalização

→ ter demais também não é bom, (muitos uns),
precisamos de um equilíbrio

- A solução adotada para o exemplo anterior pode ser generalizada pela utilização de um conjunto de funções de dispersão
 - com algumas vantagens, como veremos mais adiante
- No caso geral temos
 - C como sendo o Conjunto
 - com m elementos (membros) $(\#C = m)$
 - B o vetor (filtro) de Bloom
 - de dimensão n $(\#B = n)$
 - k funções de dispersão independentes h_1, \dots, h_k

Inicialização de um Filtro de Bloom

- Inicializar todas as posições de B com 0
- Aplicar k funções de dispersão a cada um dos elementos $de C$
 - colocando a um todas as posições devolvidas pelas funções de dispersão
 - Ou seja $B[h_i(\text{elemento})] = 1$.

Utilização do Filtro de Bloom

- Para testar um valor x :
 - aplicar as k funções de dispersão
 - e analisar o conteúdo das posições resultantes
- se $B[h_i(x)] == 1$ para todos os valores de $i = 1, \dots, k$ então x provavelmente pertence a C
 - caso contrário não pertence C

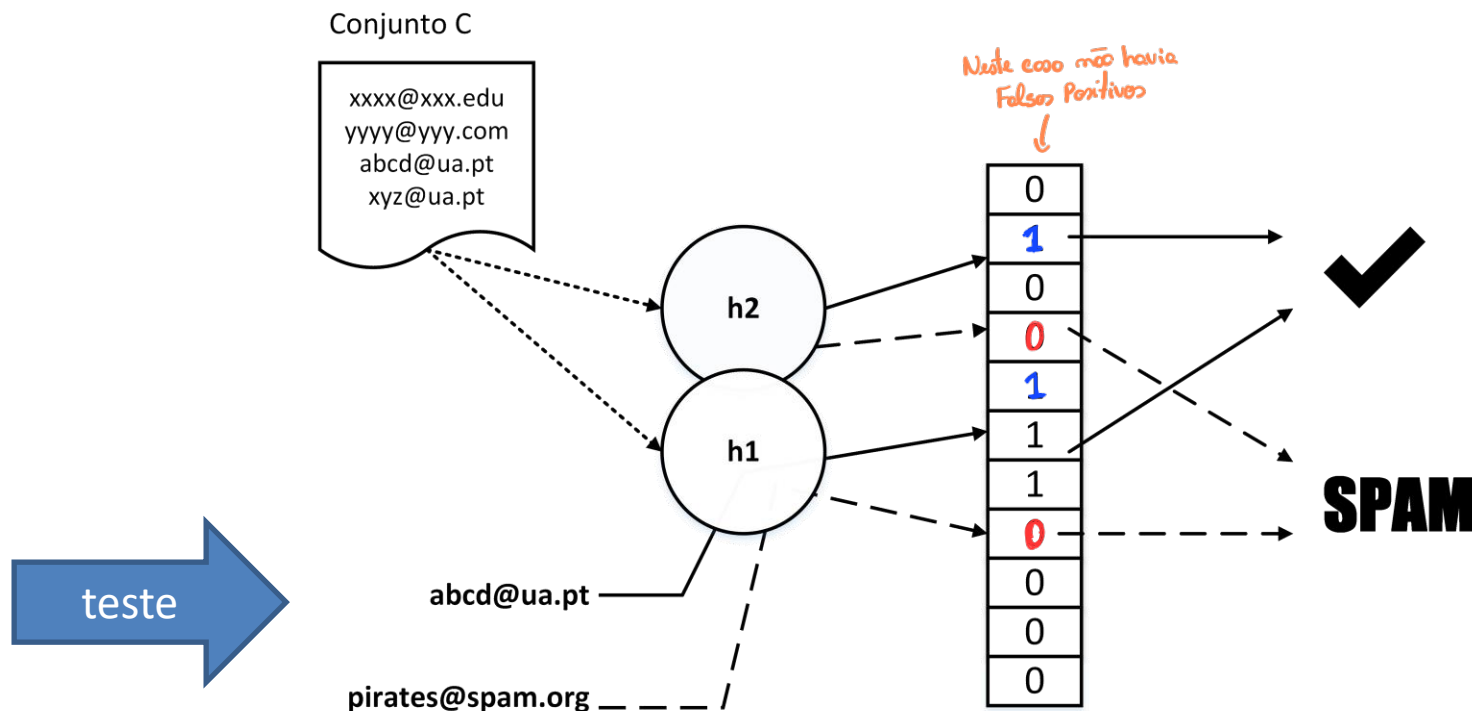
→ Existem uns que correspondem a mais de uma chave

→ Podem haver FALSOS POSITIVOS

↓
Se a função for bem dispersa e tivermos muitos zeros melhoramos os falsos positivos

Retomando o nosso exemplo ...

- Adotando valores pequenos para o número de funções, tamanho do conjunto e do filtro
 - $n = 11$ bits
 - $k = 2$ funções de dispersão.



Erros

Temos de pensar no que vamos considerar "lista negra"
→ ex: garantir que temos SPAM sempre,
e por vezes emails normais eram considerados spam ...

- O teste de associação para um elemento x funciona verificando os elementos que seriam 1 se a chave tivesse sido inserida no vetor
- Se todos os bits apropriados foram colocados a 1 por outras chaves, então x será reportado erradamente como um membro do conjunto
- Temos neste caso o que se designa habitualmente por **falso positivo**.

Exemplo de falso positivo

- Consideremos um conjunto de vegetais contendo batata e couve mas não tomate

Uma parte do filtro:

...	B(1)	C(1)	B(1)	B(1)	C(1)	-(0)	C(1)	..
-----	------	------	------	------	------	------	------	----

- Testemos com Tomate

Resultado de aplicação das funções de dispersão a Tomate:

...	-(0)	T(1)	-(0)	T(1)	T(1)	-(0)	-(0)	..
-----	------	------	------	------	------	------	------	----

→ temos um Falso Positivo

- O filtro da figura identifica incorretamente tomate como vegetal

Parâmetros do Filtro de Bloom

- m : Dimensão do conjunto
 - número de membros do conjunto
- n : Número de posições ou células do filtro
- k : Número de funções de dispersão utilizadas
- Adicionalmente, pode definir-se a fração das posições do filtro com o valor igual a 1 (f)



Implementação

- Em geral, a implementação destes filtros consiste em **3 operações**
 - inicialização
 - adição de elementos
 - teste de pertença de um elemento ao conjunto
- A implementação destas operações é simples

Tipo de dados abstrato

FiltroBloom

n	% número de bits do filtro
m	% número de elementos do conjunto
k	% número de funções de dispersão
+ inicializar (n) % Inicializar filtro com 0s	
+ adicionarElemento (elemento) % Inserir elemento no filtro	
+ membro (elemento) % Testa se elemento existe no filtro	

Implementação

- **inicializar()**
 - Simplesmente o preenchimento com zeros de todo o vetor B
- **adicionarElemento()**
 - calcula os valores das k funções de dispersão do elemento a adicionar
e atualiza as posições apropriadas do vetor B
 - No caso mais simples, coloca a 1 as posições devolvidas pelas funções de dispersão
 - o que requer tempo proporcional ao número de funções

adicionarElemento()

adicionarElemento(B, elemento, k)

for i=1:k do

- Aplicar a função de dispersão $hi()$ a elemento
 - Obtendo o respetivo hash code h
- Colocar a 1 o Vetor B na posição h
 - $B[h] = 1$

endfor

membro()

- Aplica as k funções como adicionarElemento(), mas apenas verifica se as posições contêm o valor 1
 - Se alguma das posições contém 0 não é um membro do conjunto
- A pior situação em termos de tempo de processamento ocorre para membros e para falsos positivos
 - Ambos obrigam a calcular todas as k funções de dispersão

membro()

membro(B, elemento) -> Boolean

i=0

repeat

 i=i+1

 hi é a função de dispersão i ($1 < i \leq k$)

 h é o hash code devolvido por hi (elemento)

until ((i==k) | (B [h] == 0))

if i==k then

 return (B [h] == 1) % True se posição contém 1

else

 return(False)

end.

Complexidade das operações

Operação	Parâmetros	Complexidade temporal
Inicializar()	n (tamanho do vetor)	$O(n)$
adicionarElemento()	Vetor, elemento, k funções de dispersão	$O(k)$
membro()	Vetor, elemento, k funções de dispersão	$O(k)$

Exemplos de aplicações de Filtros de Bloom

Mais Informação em:

BLOOM FILTERS & THEIR APPLICATIONS, International
Journal of Computer Applications and Technology (2278 -
8298) Volume 1– Issue 1, 2012, 25-29

[https://pdfs.semanticscholar.org/d899/05bdf1ff791bdddc7
c471070f34f4da18844.pdf](https://pdfs.semanticscholar.org/d899/05bdf1ff791bdddc7c471070f34f4da18844.pdf)

Aplicações gerais – *Spell Checkers*

- Os filtros Bloom são particularmente úteis na **verificação de ortografia**
 - São usados para determinar se uma palavra é válida numa determinada língua
- Verificação é feita **criando um Filtro Bloom com todas as palavras possíveis dessa linguagem** e verificando uma palavra contra esse filtro
- As **correções sugeridas** são geradas fazendo todas as substituições únicas em palavras rejeitadas e, em seguida, verificando se esses resultados são membros do conjunto

Redes de dados

- Muitas aplicações na área de redes
- São usados, por exemplo, em **caches de servidores proxy** na World Wide Web (WWW)
 - para **determinar eficientemente a existência de um objeto em cache**
 - Servidores proxy interceptam solicitações de clientes e respondem caso tenham a informação solicitada
- O uso de caches ajuda a **reduzir o tráfego da rede**
- Também **melhora o desempenho** quando os clientes obtêm cópias de arquivos de servidores vizinhos em vez do servidor original
 - que pode ser vários links de rede lentos de distância

Segurança e Privacidade

- Sistemas de detecção e prevenção de intrusão (IDS/IPS) usam matching de strings do conteúdo dos pacotes para deteção de conteúdo malicioso
- Os filtros Bloom são particularmente úteis para pesquisar um grande número de strings de forma eficiente
- A ideia básica é encontrar (sub)strings (comumente conhecidas como assinaturas) a alta velocidade

Demonstração Online

- Bloom Filters
 - <https://www.jasondavies.com/bloomfilter/>

Questões?

- Como obter n ?
- Como determinar o melhor valor para k (número de funções de dispersão) ?

Obtenção dos parâmetros de um Filtro de Bloom

Para que se possa definir de forma adequada um filtro tem de primeiro se perceber **como se relacionam com os indicadores de desempenho como a probabilidade de falsos positivos**

Falsos positivos e falsos negativos

Elemento pertence ao conjunto ?	Resultado do teste de pertença	Correto ?	Tipo de erro
SIM	Sim (pertence)	Sim	
	Não	ERRO	Falso negativo
NÃO	Sim	ERRO	Falso positivo
	Não	Sim	<i>Logo que pertence e não pertence !</i>

Obtenção dos parâmetros de um Filtro de Bloom

- Qual a relação dos falsos positivos com os parâmetros do filtro (n e k) e do conjunto (m)?

Lançamento de m dardos para n alvos

- Caso similar ao do lançamento de m dardos para n alvos igualmente prováveis
 - por sua vez similar ao problema dos aniversários

- Se atirmos m dados para n alvos igualmente prováveis, **qual a probabilidade de um alvo ser atingido por pelo menos um dardo ?**

↳ Como é que se relaciona com o (n, K, m) ?

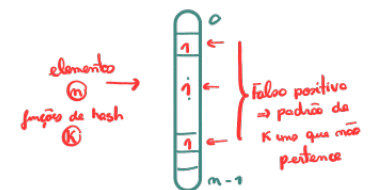
$P_{\text{falso positivo}} = 0.001$ // como?

- No caso dos filtros de Bloom

- os alvos são os vários bits do filtro
- os dardos são os valores assumidos pela função de dispersão

Resultados da função de hash

Bom função de hash: $P = \frac{1}{n}$
↳ tamanho do array!



Probabilidade de falsos positivos ?

- Para termos falso positivo teremos de ter as k posições determinadas pelas funções de dispersão com o valor 1
 - Para elementos não pertencentes ao conjunto
- Como obter essa probabilidade?
- Começamos pela probabilidade de um bit estar a 1 ...

Probabilidades para um bit

- Inicialmente todos os bits estão a zero
- Designando o bit na posição i por b_i ...
- Qual a probabilidade do bit $b_i = 1$ depois de aplicar a primeira função de dispersão na inserção de um elemento?
- Assumindo que a função de dispersão seleciona cada uma das posições do vetor com igual probabilidade ...
- A probabilidade é simplesmente
$$P[b_i = 1] = \frac{1}{n}$$
 - pois os n alvos são equiprováveis
- Em consequência: $P[b_i = 0] = 1 - \frac{1}{n}$

Probabilidade após aplicar k funções de dispersão

- Probabilidade de um bit se manter a zero após a inserção de um elemento ?
- É a probabilidade do bit continuar a ser zero depois de aplicadas as k funções de dispersão
- Se assumirmos que os resultados das funções de dispersão são independentes, será:

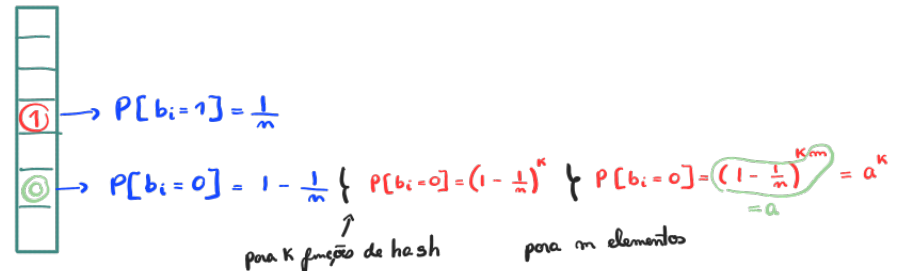
$$P[b_i = 0] = \left(1 - \frac{1}{n}\right)^k$$

Probabilidade após inserir m elementos

- Após a inserção de m elementos, assumindo independência, temos:

- $P[b_i = 0] = \left(1 - \frac{1}{n}\right)^{k m}$

- Fazendo $\left(1 - \frac{1}{n}\right)^m = a$
depois de colocar m elementos



- $P[b_i = 0] = a^k$

- $P[b_i = 1] = 1 - a^k$

- Relacionada com a **probabilidade de um alvo ser atingido por pelo menos um dardo**

Probabilidade de falsos positivos

- Temos um falso positivo **quando temos os k bits iguais a 1** para um elemento não pertencente a C
- A probabilidade de um falso positivo, p_{fp} , após inserirmos m elementos é

$$p_{fp} = \left[1 - \left(1 - \frac{1}{n} \right)^{km} \right]^k = \left(1 - a^k \right)^k$$

- Aplicando $\lim_{n \rightarrow \infty} (1 - 1/n)^n = e^{-1}$
- Pode ser aproximada por:

Simplificada...

$$p_{fp} \approx \left(1 - e^{-km/n} \right)^k$$

Isto não é estritamente correto, pois pressupõe independência para as probabilidades de cada bit ser definido.

Uma análise alternativa, chegando à mesma aproximação sem a hipótese de independência.

- A probabilidade de falsos positivos depende de k, m e n

Efeito de k na p_{fp}

- Exemplo:

$m = 10^9$ (mil milhões)

$n = 8 \times 10^9$ (8 mil milhões)

$$m/n = 1/8$$

- $k=1$

$$p_{fp} = 1 - e^{-1/8} = 0.1175$$

- $k=2$

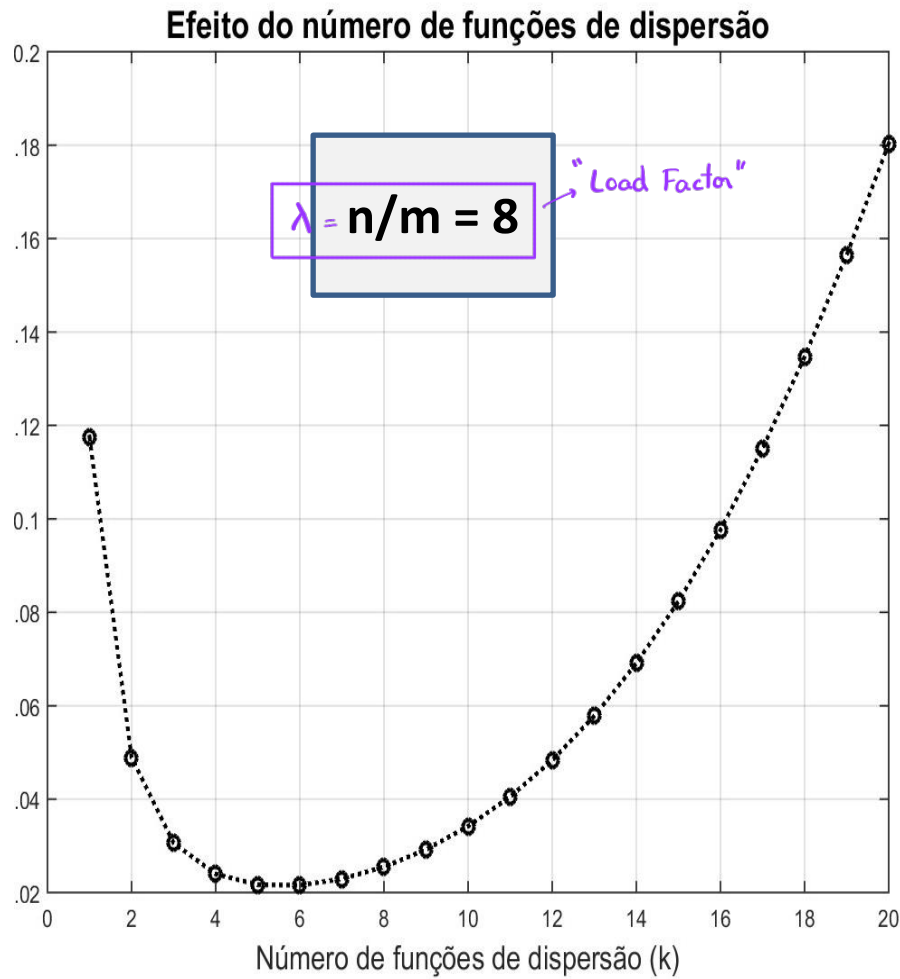
+ Função de hash ?
- Falsos positivos

$$p_{fp} = (1 - e^{-2/8})^2 = 0.0489$$

- Probabilidade de erro diminui com aumento de k (acontece sempre ?)
- O que acontece se formos aumentando k ?

Efeito de k na prob. Falsos Positivos

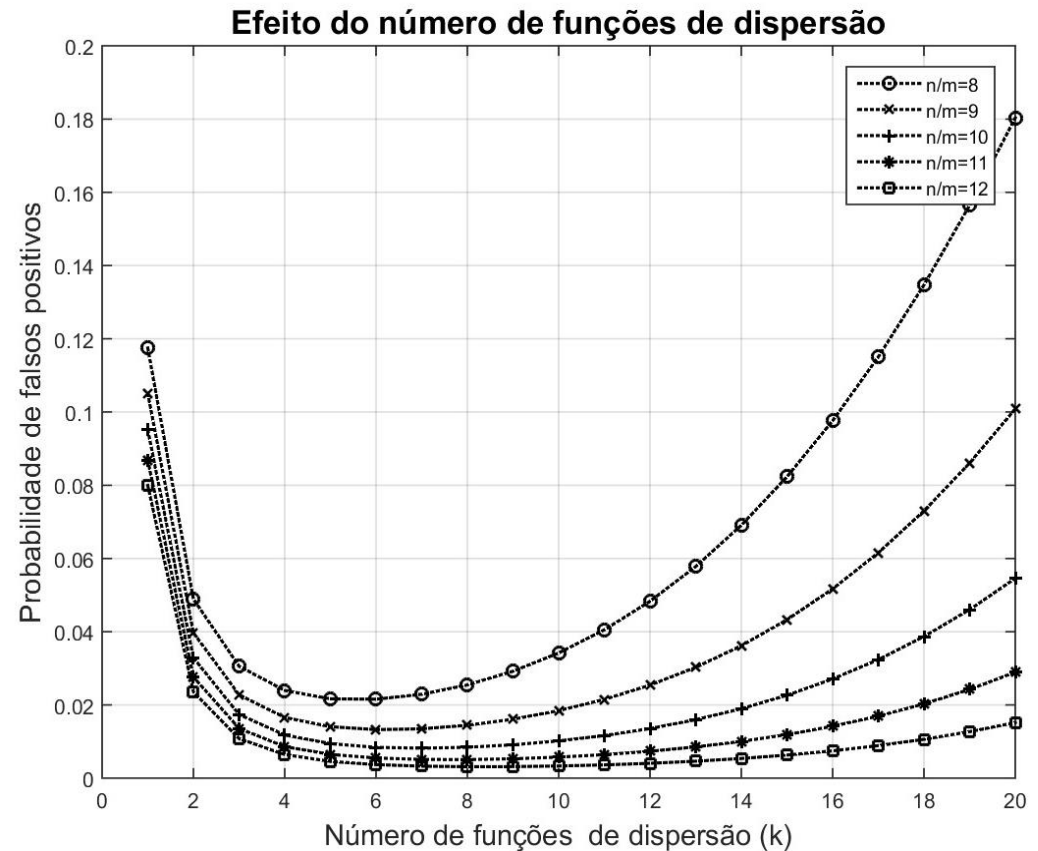
• Com menos funções de hash podemos ter mais hipóteses de cometer erros



- Adicionando mais funções não diminui obrigatoriamente essa probabilidade
- Decresce até um certo valor de k
- Depois aumenta
 - Motivo ?

Redução de falsos positivos

- FP podem ser reduzidos aumentando o tamanho do vetor B
 - neste caso à custa de mais memória.
 - O efeito da relação n/m é ilustrado na figura
- Também podem ser reduzidos aumentando o número de funções de dispersão
 - mas apenas até um certo valor



Então qual o valor ótimo de k ?

- É possível determinar o número de funções de dispersão que minimiza a probabilidade de falsos positivos, p_{fp}
- Para facilitar os cálculos minimiza-se $\ln(p_{fp})$
• Seria derivar e igualar a zero!
- Aplicando $\ln()$ a $p_{pf} = (1 - a^k)^k$
- temos: $\ln(p_{pf}) = k \ln(1 - a^k)$

Minimização

- Como obter k ótimo (que minimiza p) ?
- A solução usual de calcular zeros da derivada

- Derivando (em ordem a k) e igualando a zero
$$(1 - a^k) \ln(1 - a^k) - a^k \ln(a^k) = 0$$

- Que tem por solução

$$a^k = 1/2$$

[usar, por exemplo, fsolve no Matlab]

k ótimo

- O valor ótimo de k pode ser obtido aplicando logaritmos e resolvendo em ordem a k :

$$k_{\text{ótimo}} = \frac{\ln(1/2)}{\ln(a)}$$

- Substituindo o valor de a temos

$$k_{\text{ótimo}} = \frac{\ln(1/2)}{\underbrace{m}_{a^m} \times \ln(1 - 1/n)}$$

- Aproximando $\ln(1 - 1/n)$ pelo primeiro termo da série de Taylor ($-1/n$)

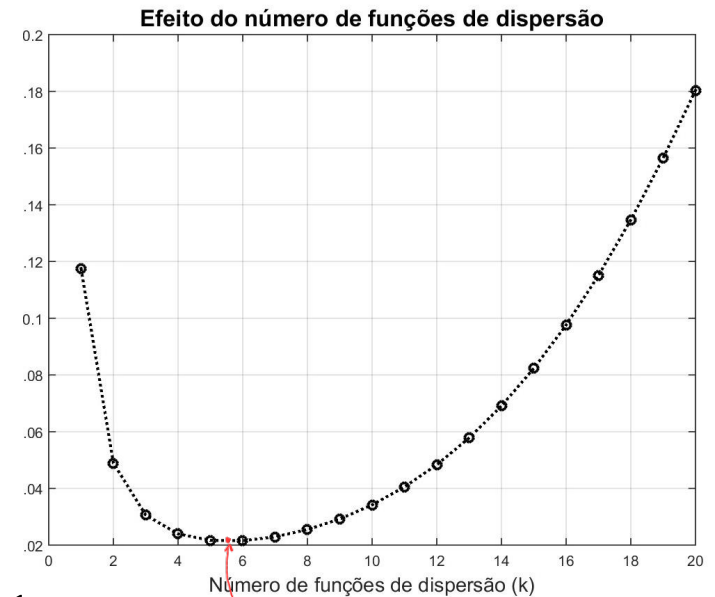
$$k_{\text{ótimo}} \approx \frac{n \ln(2)}{m} = \frac{0.693 n}{m}$$

→ Quantos valores estamos dispostos a ter no array para cada elemento

- Na prática **utiliza-se o inteiro mais próximo**

Exemplo – k ótimo

- $m = 10^{12}$
- $n = 8 \times 10^{12}$
- k ótimo ?



- $k_{\text{ótimo}} \approx \frac{n \ln(2)}{m} = \frac{0.693 \times 8 \times 10^{12}}{10^{12}} = 5.54$
– Aproximadamente 6
- Qual a relação com um dos gráficos que já vimos?

Determinação de n

- O valor de n pode ser calculado substituindo o valor ótimo de k na expressão de probabilidade

$$p_{fp} \approx \left(1 - e^{-km/n}\right)^k$$

- Sendo dados m e um objetivo em termos de probabilidade de falsos positivos p_{fp}
- e assumindo que o valor de k ótimo é adotado

Limite inferior para a probabilidade de erro (FP)

- Se tivéssemos $a^k = 1/2$ para um inteiro $k_{ótimo}$, a equação $(1 - a^k)^k$ resultaria em:

$$\begin{aligned} p_{ótima} &= \left(1 - \frac{1}{2}\right)^{k_{ótimo}} \\ &= \left(\frac{1}{2}\right)^{k_{ótimo}} \\ &= 2^{-k_{ótimo}} \end{aligned}$$

- Que pode ser considerado o limite inferior para a probabilidade de erro
 - Obviamente falsos positivos

Filtros de Bloom – Aspectos positivos

- Os filtros de Bloom garantem **não existência de falsos negativos** ...
- e usam uma **quantidade de memória limitada**
 - Ótimo para pré-processamento antes de processos mais exigentes
- Adequados para implementação em hardware
 - As funções de dispersão podem ser **paralelizadas**

Filtros de Bloom - Limitações

- Como a tabela não pode ser expandida
 - o máximo número de elementos a armazenar no filtro tem de ser conhecido previamente
- Assim que se excede a capacidade para a qual foi projetado, os falsos positivos aumentam rapidamente ao serem inseridos mais elementos

Filtros de Bloom - Compromissos

- Os falsos positivos podem ser diminuídos através de:
 - Aumento do número de funções de dispersão (até ao $k_{ótimo}$)
 - E do espaço alocado para armazenar o vetor

Comentários finais

- Os filtros de Bloom **devem ser considerados para programas em que um teste de pertença imperfeito pode ser aplicado** a um conjunto de dados (muito) grande
- As **grandes vantagens** de um filtro de Bloom são a rapidez e taxa de erro
- Apesar de poder ser aplicado a conjuntos de qualquer dimensão, **árvores e heaps são melhores soluções para conjuntos pequenos**

Filtros Bloom de contagem / Filtros de Contagem

Counting Bloom Filters

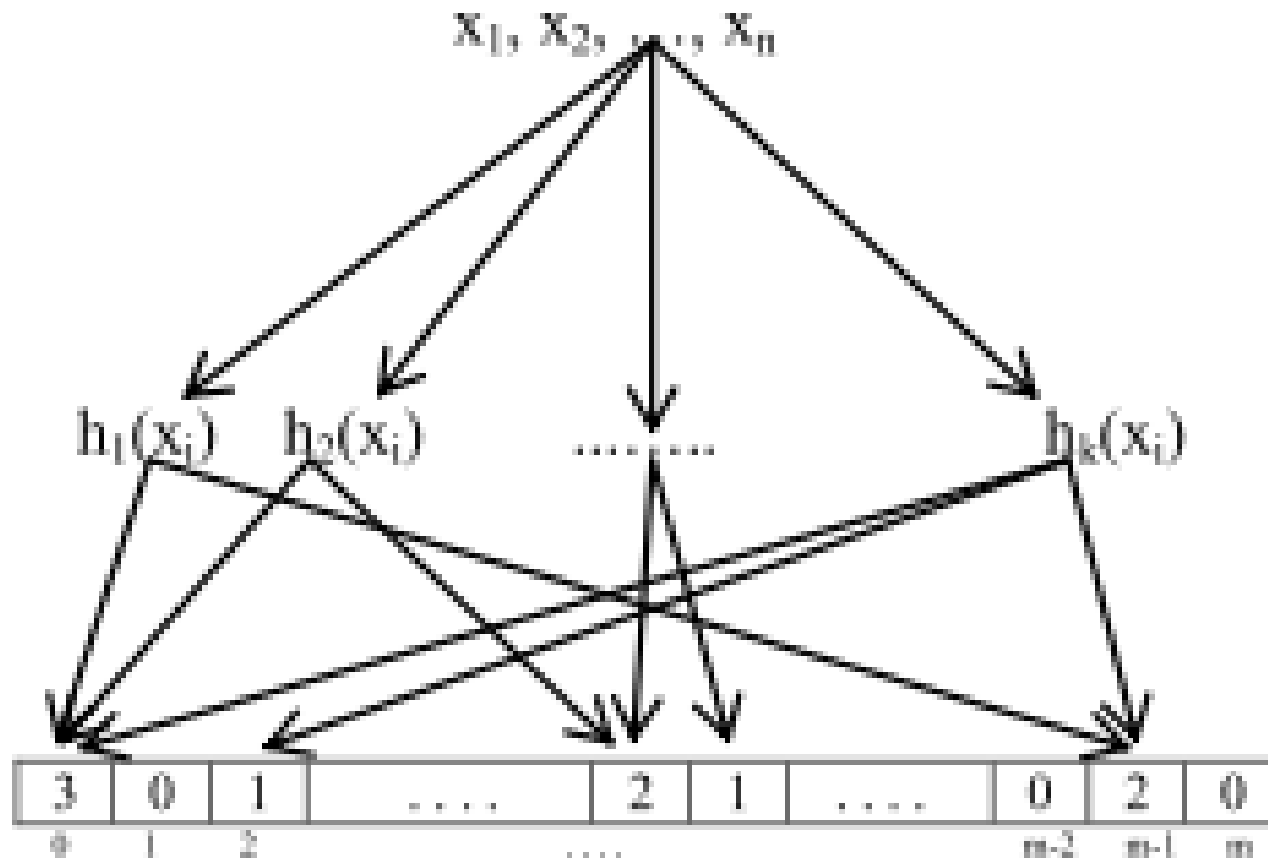
Filtros de Contagem

- Um filtro Bloom básico representa um conjunto, mas:
 - não permite a consulta da **multiplicidade**
(número de vezes que elemento foi inserido)
 - nem suporta a **remoção** de elementos
- Um filtro Bloom de contagem estende um filtro Bloom básico para dar resposta a estas limitações

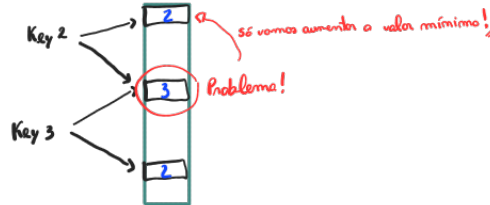
Filtros de Contagem

- As posições do vetor são estendidas de um único bit para um **contador de b bits**
[a versão original utilizava 4 bits]
- Na **inserção** de um elemento o **contador é incrementado**
- na **remoção** é **decrementado**
- Ocupa mais espaço
 - tipicamente 3 a 4 vezes mais que o de Bloom

Filtro de Contagem



Obtenção da multiplicidade



- Problema:
 - Os contadores correspondentes a um elemento podem também ser alterados por outros elementos
- Como ter boa estimativa do número de inserções de um elemento ?
- Solução:
 - Usar o **valor mínimo** entre os vários contadores correspondentes ao elemento

Implementação

- A implementação deste tipo de filtros em Matlab é simples
- Bastam ligeiras alterações a:
 - `adicionar()` : passa a incrementar
 - `membro()` : requer que todos sejam não nulos
valor ≥ 1
- A criação da função adicional `contagem()`
valor mínimo de todas as funções
- E, se necessário, a função `remove()`
retorna 1 a todos

contagem()

- Operação nova
- Para obter a contagem (multiplicidade) associada a um elemento do conjunto:
 1. Determina-se o conjunto de contadores que lhe correspondem
 - Através das k funções de dispersão
 2. Calcula-se o **valor mínimo** armazenado nesses contadores
 - Algoritmo como *minimum selection* (MS)

remove()

- Implementação similar a adicionar()
- **Esta operação introduz a possibilidade de falsos negativos**
 - Quando se coloca um bit a zero que era parte dos k bits de outro elemento, o filtro deixará de o considerar como pertencendo ao conjunto

Problemas ...

- *Overflow do contador*
 - Quando a contagem chega a $2^b - 1$
 - Tipicamente as implementações param de contar
 - Preferível a recomençar em 0 (introduziria falsos negativos)
 - Mas introduz mais erro na estimativa da contagem
- *Escolha de b* (número de bits dos contadores)
 - Um valor grande reduz a poupança de espaço
 - Um valor pequeno rapidamente leva a overflow
 - Escolha do valor é um compromisso e depende dos dados



Note to other teachers and users of these slides: We would be delighted if you found this our material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. If you make use of a significant portion of these slides in your own lecture, please include this message, or a link to our web site: <http://www.mmds.org>

PPT baseado em: Mining Data Streams

Mining of Massive Datasets

Jure Leskovec, Anand Rajaraman, Jeff Ullman

Stanford University

<http://www.mmds.org>



Outras fontes utilizadas

- <http://matthias.vallentin.net/blog/2011/06/a-garden-variety-of-bloom-filters/>
- https://en.wikipedia.org/wiki/Bloom_filter#Counting_filters