

Introdução

Inteligência na Web e Big Data

Fabricio Olivetti de França e Thiago Ferreira Covões
folivetti@ufabc.edu.br, thiago.covoes@ufabc.edu.br

Centro de Matemática, Computação e Cognição
Universidade Federal do ABC



Introdução

Extração e descoberta de conhecimentos de bases de dados.

Interdisciplinar: aprendizado de máquina, estatística e banco de dados.

Objetivo: conhecimento entregue de forma legível ou modelo computacional.

Agrupamento dos Dados:

- Segmentação
- Sumarização
- Redução

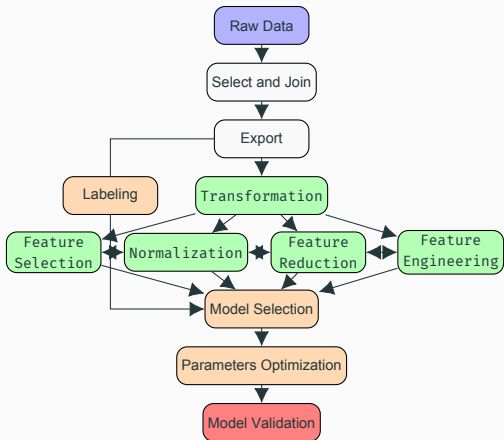
Classificação:

- Rotular dados em classes
- Aplicável em dados ainda não rotulados
- Automatizar tarefa

Regressão:

- Modelo matemático descritivo
- Entender fenômenos

Pipeline



Novas tecnologias:

- Poder de processamento
- Capacidade de armazenamento
- Capacidade de coleta de dados

Poder de Processamento

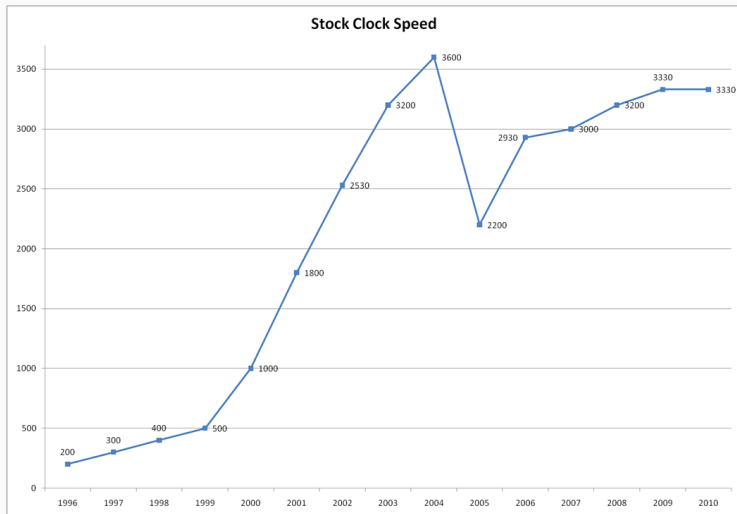


Figura 1: Fonte: maximumpc.com

Não leva em conta apenas o clock:

- Instruções utilizando menos ciclos
- Cache maior
- Pipeline de instruções / Paralelização
- Mais instruções → otimização
- Intrinsics
- Multicore

Armazenamento

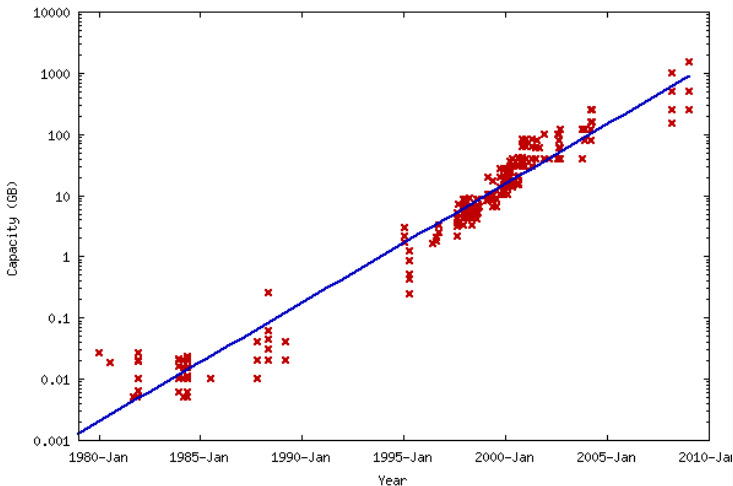


Figura 2: Fonte: wikimedia.org

- Velocidade e durabilidade também aumentaram.
- Custo diminuiu com o tempo.
- Computadores capazes de gerenciar múltiplos discos.

Novas tecnologias, novas fontes:

- Transações de compras na internet
- Dados de navegação (Google, Facebook)
- Dados de satélite
- Experimentos com genes

O crescimento de capacidade de armazenamento e coleta levou a uma imensa quantidade de dados, que devem ser processados.

Em 2012, a Amazon vendia 158 itens por segundo: **13.651.200** itens em um dia!

Assumindo que cada transação de um cliente continha 5 itens, teremos **2.730.240** transações em um dia.

Em um ano teremos **996.537.600** transações para serem analisadas por técnicas de mineração de dados.

Se cada transação for representada por 24 bytes (otimista), teríamos cerca de 24 Gigabytes de dados.

Em 2010, a estimativa era de 3 bilhões de spam enviados por dia.

Sem imagens, um e-mail de spam tem um tamanho médio de 5kb (2007).

Em um dia teremos 15TB de informação para processar!

Durante o modismo do Big Data foram definidos 4 desafios dessa área:

- Volume
- Velocidade
- Variedade
- Veracidade

Os resultados extraídos da base de dados são relevantes?

Possibilidade alta de ruídos aleatórios.

Torture os dados, e confessará qualquer coisa. — *Ronald Coase, economista.*

Uma crítica da área de Big Data é que a grande quantidade de dados pode te ajudar a encontrar padrões que, na verdade, eram esperados.

Se seu modelo estatístico retorna mais casos positivos do que o esperado, então seu modelo (ou hipótese) não funciona.

Exemplo Fictício

Hipótese: grupos criminosos se encontram em um hotel para planejar o crime.

Objetivo: encontrar pessoas sem relação aparente que, por pelo menos duas vezes, se hospedaram no mesmo hotel no mesmo dia.

Exemplo Fictício

Experimento: rastrear 10^9 pessoas hospedadas em hotéis em um período de 1000 dias.

Cada pessoa ficou em um hotel cerca de 1% desse período (10 dias).

Temos 10^5 hotéis e cada hotel comporta 100 pessoas e assumimos que eles estão sempre lotados.

Pergunta: se todos agirem aleatoriamente, nosso modelo estatístico encontrará algo suspeito nos dados coletados?

Duas pessoas p, q escolhem um dia específico:

$$P(p, d1) \times P(q, d2) = \frac{1}{100} \times \frac{1}{100} = 10^{-4}$$

Duas pessoas p, q escolhem um hotel **específico**:

$$P(p, h1) \times P(q, h2) = \frac{1}{10^{-5}} \times \frac{1}{10^{-5}} = 10^{-10}$$

Essas duas pessoas escolhem o mesmo hotel:

$$P(p, q, h) = 10^{-10} \times 10^{-5} = 10^{-5}$$

Elas escolhem o mesmo dia:

$$P(p, q, h, d) = 10^{-4} \times 10^{-5} = 10^{-9}$$

Isso vai ocorrer por duas vezes:

$$P(p, q, h, d1, d2) = 10^{-9} \times 10^{-9} = 10^{-18}$$

Como temos 5×10^5 pares de dias, a probabilidade de observar tal situação em um deles é:

$$5 \times 10^5 \times 10^{-18} = 5 \times 10^{-13}$$

Como temos um total de 5×10^{17} pares de pessoas para investigar, a quantidade esperada de observações que faremos em nossa base de dados será:

$$5 \times 10^{17} \times 5 \times 10^{-13} = 2.5 \times 10^5$$

Exemplo Fictício

Se temos 10 pares de criminosos que apresentou esse comportamento, teremos que procurá-los dentro das 250.000 pares de pessoas que fizeram o mesmo...

Tenha certeza que a hipótese não permita um número de eventos aleatórios maior do que o que se busca.

Exercício

Em uma base de dados de avaliações de produtos, queremos detectar usuários que são pagos para avaliar uma série de produtos positivamente. É razoável imaginar que se dois usuários avaliaram positivamente em comum 90% dos produtos, então esses foram pagos para isso?

Como armazenar e tratar a grande quantidade de dados coletada?

- Dividir um arquivo muito grande em várias máquinas conectadas em rede
- Redundância dos dados
- Transparência no uso

Arquivos seguem a filosofia *write-once-read-many-times*.

Ou seja, os arquivos armazenados se tornam somente leitura.

Isso não é problema, pois queremos extrair informações, não alterar os dados.

Tarefas simples e comuns como cálculo de estatística básica, procedimentos de seleção, filtragem, junção não são triviais.

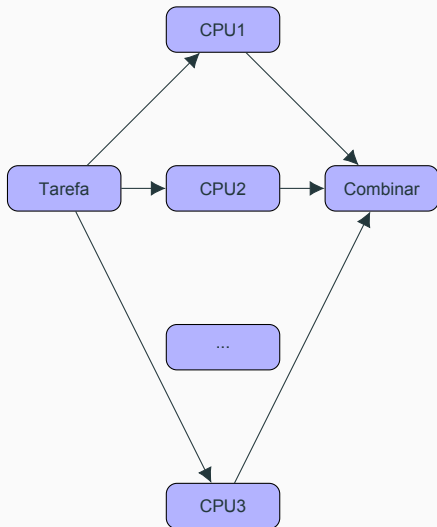
Demandam algoritmos que podem ser aplicados localmente em cada uma das máquinas e combinados ao final.

A imutabilidade dos dados evita o problema de Race Condition.

Não existe alteração, não existe mudança de estado, não existe conflito!

Particionando a Tarefa

Possibilidade de executar tarefas localmente.



Paradigma Funcional

Um paradigma de programação que não permite mudança de estados e dados mutáveis.

Computação como sequência de funções matemáticas.

λ -calculus: formalização matemática de computação.

Turing completa.

Propriedades interessantes para processamento em paralelo e distribuído.

Dada a expressão:

$$x + f(x)$$

Podemos substituir $f(x)$ pelo valor avaliado sem alterar o resultado final:

$$x + y$$

Funções que permitem uma transparência referencial são chamadas de funções puras.

$f(x)$ sempre retornará o mesmo valor ao passar o mesmo valor de x .

Contra Exemplo

```
1 int f(int * x) {  
2     *x = *x + 1;  
3     return *x;  
4 }
```

```
1 System.currentTimeMillis()
```

Em linguagens funcionais puras uma variável não pode ter seu valor alterado.

Uma variável x passada para uma função $f()$ não pode ter seu valor alterado por ela.

Como as variáveis são imutáveis, não existe conflito de acesso.

Se dois processos tentam acessar o conteúdo de uma mesma variável, tenho certeza que ambos receberão o mesmo valor.

A pureza das funções e expressões matemáticas, permite que elas sejam avaliadas em qualquer ordem:

$$f(x) = \sin(x) + \cos(x) * \text{sqrt}(x)$$

Não importa quem vai ser calculado primeiro, *sin*, *cos* ou *sqrt*, o resultado final será o mesmo!

$$f(x) = \sin(x) + \cos(x) * \text{sqrt}(x)$$

Cada *thread* pode ficar responsável por computar uma parte da expressão.

Por muito tempo o uso dessas linguagens foi acadêmico.
Atualmente muitas empresas já fazem uso desse paradigma.

Erlang - criado pela Ericsson para programar sistemas tolerante a falhas.

OCaml - uso em sistemas financeiros, programas livres de bugs.

Scala - elementos funcionais no Java.

Closure - dialeto do Lisp que roda no Java Virtual Machine.

Haskell - linguagem acadêmica que encontrou seu lugar na indústria (anti-spam do Facebook).

Cenário que recebemos um fluxo contínuo de dados e queremos processá-los em tempo real.

Como tratar?

Algoritmos de aprendizado de máquina que aprendem incrementalmente.

O algoritmo deve se adaptar durante o processo de aprendizagem (mudança de tendências).

Apenas uma quantidade limitada de dados pode residir em memória.

Maior parte dos algoritmos de Aprendizado de Máquina assumem uma entrada de dados bem definida.

Geralmente um vetor multidimensional.

Muitas bases de dados de interesse não possuem esse formato:

- Textos
- Imagens
- Som
- Vídeos

Necessidade de transformar tipos desestruturados em uma estrutura vetorial.

Cuidado para não perder informação no meio do caminho!

Pense em exemplo de tarefas que demanda a combinação de dados de imagem, som e texto.

Conclusão

Trabalhar com grandes quantidades de dados traz diversos problemas desafiadores.

O simples volume de dados não implica em maior facilidade na extração de informação.

Algumas das formas de tratamento dos problemas serão apresentados nas próximas aulas.