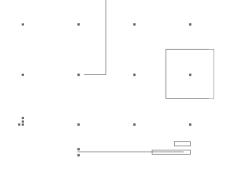
# THE SCIENCE SCIENCE

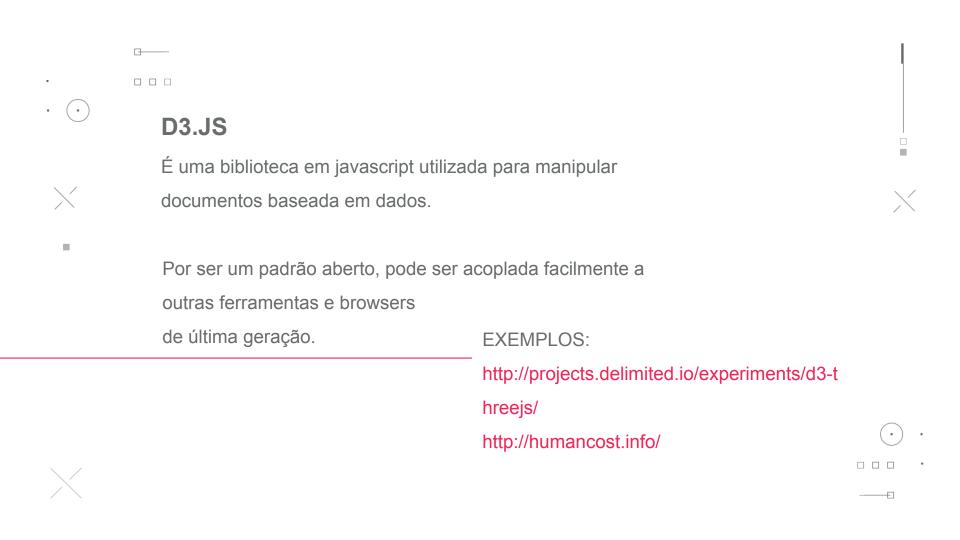
MACHINE LEARNING E DATA MINING



# BIG DATA SCIENCE PARTE 4

DIÓGENES JUSTO





#### **ScatterPlot**

```
install.packages("scatterD3")

library(scatterD3)

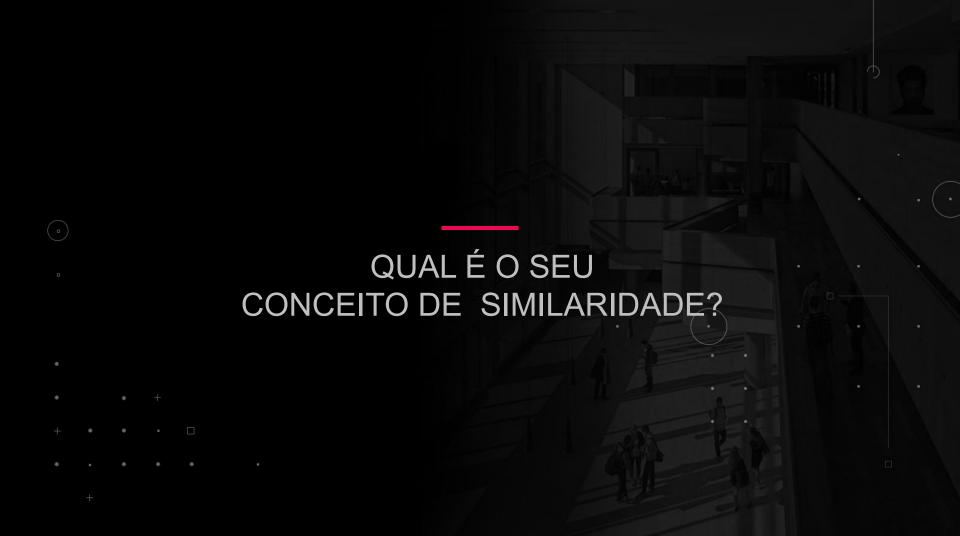
Utiliza a D3.js para criar gráficos de dispersão mais ricos e dinâmicos.

d <- mtcars

scatterD3(x = d$wt, y = d$mpg, lab = rownames(d),

col_var=d$cyl, symbol_var = d$am, symbol_lab =

"Manual transmission")
```



•

•

DISTÂNCIA ENTRE DOIS PONTOS

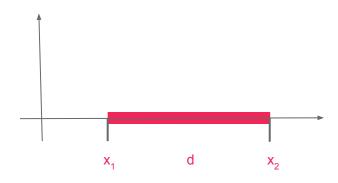
Para o caso unidimensional:

$$d = x_2 - x_1$$



• • [

+



•

•

DISTÂNCIA ENTRE DOIS PONTOS

E para o caso bidimensional?

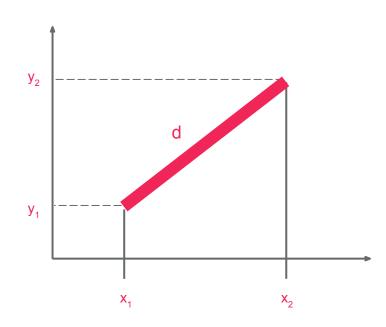
$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

•

•

+

• •



+

## CLUSTERIZAÇÃO

- X X
  - UMA COLEÇÃO DE OBJETOS COM:
    - Alta similaridade (relacionados) no grupo
    - Baixa similaridade (não relacionados)
  - APRENDIZAGEM NÃO SUPERVISIONADA (sem classes pré-definidas):
    - Contrasta com classificação (supervisionada)
  - APLICAÇÕES COMUNS:
    - Análise de distribuições
    - Etapa intermediária (precede aprendizagem supervisionada, por exemplo)
    - Identificação de outliers

## : = × × ×

#### • CONCEITO DE SIMILARIDADE:

Dois objetos são similares se eles te um grau
 alto de proximidade, isto é, pequena distância
 na classe que está sendo analisada

## SIMILARIDADE

#### • EXEMPLO:

- Alunos estudiosos ou não: distância da nota
- Pessoas parecidas: distância entre pontos da face
- Detecção de movimento: distância entre similaridade de quadros (frames) aumentou



## SIMILARIDADE

- FUNÇÃO SIMILARIDADE (EM GERAL):
  - Número real

 $\times$ 

- Entre 0 (completamente diferentes) e 1 (idênticos)
- DISSIMILARIDADE (DISTÂNCIA)
  - Entre 0 (idênticos) e 1 (completamente diferentes)

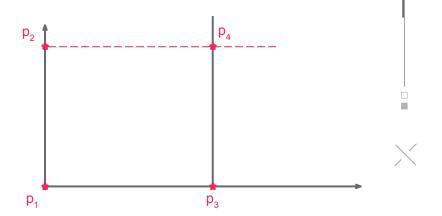
 $\cdot$   $\cdot$ 

## DISTÂNCIA NUMÉRICA



Ponto	Х	Υ
p1	0	0
p2	0	1
рЗ	1	0
p4	1	1

Matriz de dados



	p1	p2	р3	p4
p1	0			
p2	1	0		
р3	1	1,41	0	
p4	1,41	1	1	0

Matriz de dissimilaridade (distância Euclidiana)











.

## DISTÂNCIA (OU MÉTRICA) DE MINKOWSKI\*:

$$d(i,j) = \sqrt{p} |x_{i1} - x_{i1}|^{p} + |x_{i2} - x_{j2}|^{p} + ... + |x_{i1} - x_{j1}|^{p}$$

$$x_{j1}|^{p}$$











#### PROPRIEDADES:

$$d(i,j) > 0$$
, se  $i <> j$  e  $d(i,i) = 0$   
 $d(i,j) = d(j,i)$   
 $d(i,j) <= d(i,k) + d(k,j)$ 

(Positividade)

(Simetria)

Casos específicos de p (chama-se norma L-p)

Norma L-1: distância de Manhattan

Norma L-2: distância Euclidiana

Norma L-max ou L-inf: distância "Suprema"













## OUTRAS DISTÂNCIAS

- DISTÂNCIA PARA ATRIBUTOS BINÁRIOS:
  - Coeficiente de Jaccard (coerência)
- Distância para ATRIBUTOS CATEGÓRICOS
- Distância entre VARIÁVEIS ORDINAIS
- SIMILARIDADE do cosseno (utilizado em text mining)



## MÉTODO K-MEANS\*

Cada cluster é representado por um centro (centróide) Arbitra-se K clusters, inicie:

- Aleatoriamente escolha K centróides
- Faça até critério de convergência
  - Associe cada ponto ao centróide mais próximo
  - Para cada cluster, calcule a nova posição do centróide como ponto médio do cluster

















Melhor escolha dos K centróides iniciais

- K-Means++, K-Means inteligente

Diferentes tipos de clusters

- K-Medois, K-Medians, K-Modes

Aplicando tranformações

- K-Means ponderado, Kernel K-Means





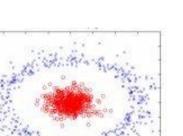








# OUTRAS DISTÂNCIAS -





 K-Medoids/K-Median: menos sensível a outliers.

 Kernel K-Means: clusters em regiões não lineares





Veja uma animação que mostra o funcionamento do K-Means:

http://shabal.in/visuals/kmeans/1.html



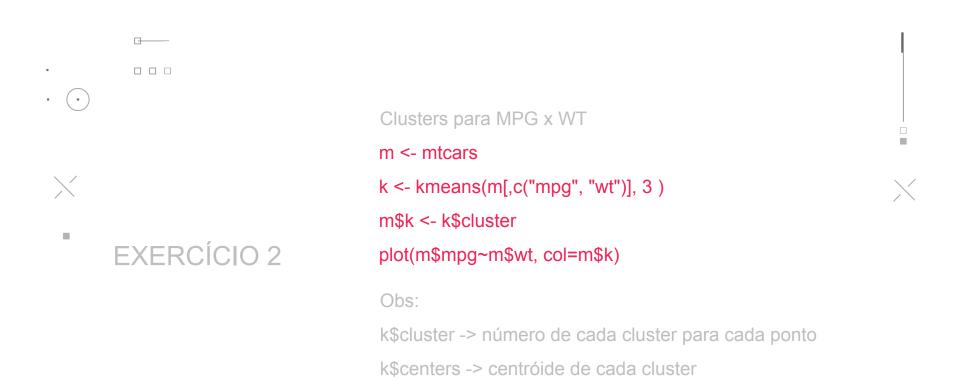
















A clusterização como etapa prévia de processamento, gerando uma nova característica (variável) categórica.

```
summary(lm(mpg~wt), data=m)
summary(lm(mpg~wt + k, data=m)
```











## ÁRVORES DE DECISÃO COM TITANIC

```
library(party)
df <- read.csv("train.csv")</pre>
set.seed(33)
df<- df[sample(nrow(df)) ,]</pre>
train \leftarrow df[1:600,]
test <- df[601:891,]
mDT <- ctree(Survived ~ Sex,data=train)
pred <- predict(mDT, newdata=test)</pre>
table(pred, test$Survived)
```

## ÁRVORES DE DECISÃO COM TITANIC

```
mDT <- ctree(Survived ~ Sex+Age+Pclass+Fare, data= train)
```

pred <- predict(mDT, newdata=test)</pre>

table(pred, test\$Survived)

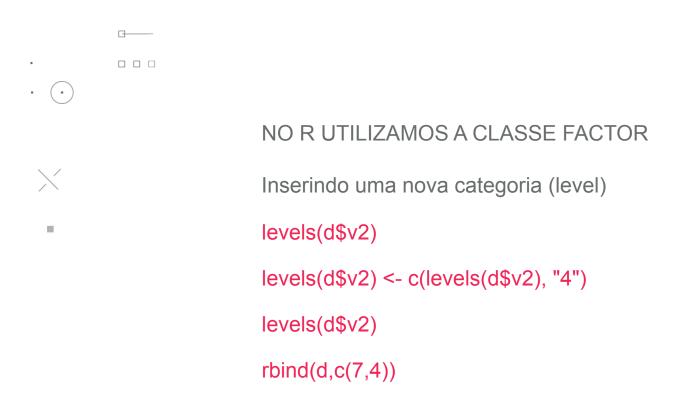
## ÁRVORES DE DECISÃO COM TITANIC

```
mDT <- ctree(Survived ~ Sex+Age+Pclass+Fare, data=
train)
pred <- predict(mDT, newdata=test)
table(pred, test$Survived)</pre>
```

```
#Regressor x Classificador: as.factor
endógena
df$Survived F <- as.factor(df$Survived)</pre>
```



```
____
NO R UTILIZAMOS A CLASSE FACTOR
               Função factor() factor(c(1,2,3,3,3,1))
                    Como inserir dados em uma variável factor?
               v1<-1:6
                v2 < -factor(c(1,2,3,3,3,1))
                  d<-data.frame(cbind(v1, v2))
                str(d)
                  d$v2 <- as.factor(d$v2)
                 rbind(d,c(7,4))
```









AVALIAÇÃO DE MODELOS

Como avaliar se um modelo prevê corretamente os dados?

 PREVISTO

 PVC
 NORMAL

 PVC
 0.78
 0.22

 NORMAL
 0.16
 0.84

MATRIZ DE CONFUSÃO DOS RESULTADOS

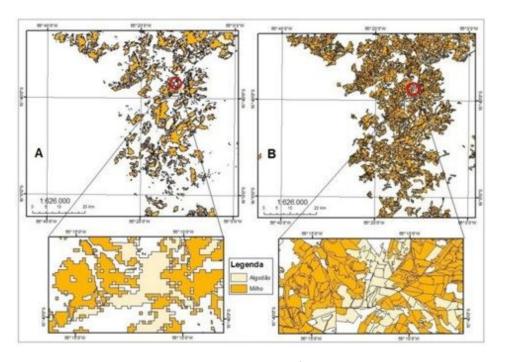
AVALIAÇÃO DE MODELOS

Obs.1: na fonte não é citado se o real é coluna ou linha. Inseri arbitrariamente para utilizar o exemplo.

Obs.2: reparem que os dados foram apresentados como percentual (probabilidade) em linha. Isto é, 78% dos PVCs reais foram previsto corretamente. 22% foram estimados erroneamente como normais.

Fonte: http://slideplayer.com.br/slide/1266306/

Classificação de arritmias cardíacas



Fonte: MAPEAMENTO DE CULTURAS AGRÍCOLAS POR IMAGENS DE SATÉLITE - COSTA, SOUTO e Zeilhofer

http://www.sinageo.org.br/2012/trabalhos/10/10-379-577.html

**VERDADE DE CAMPO** ALGODÃO MILHO OUTROS TOTAL ALGODÃO 819 6.388 5.569 CLASSIFICADOR MILHO 3.552 3.552 761 OUTROS 5.674 6.435 5.569 5.119 5.674 16.375 TOTAL

•

.

+

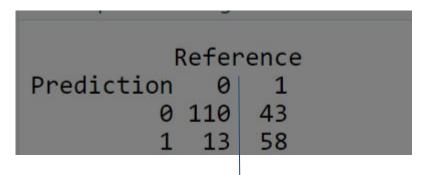
.

### Arv Decisão

```
> # Matriz de confusão
> table(prev, teste$Sur

prev 0 1
0 162 64
1 5 60
>
```

## Regr Logist



.

MODELOS DE CLASSIFICAÇÃO – AVALIAÇÃO

Confusion Matrix para Titanic:

- 891 passageiros (em Train)
- Acerto 78% => ? Passsageiros?
- ~ 326

x1 + x2 + x3 = ?

•

SURVIVED 326 x1

NOT SURVIVED x2 x3

**PREVISTO** 

(Estimado no modelo)

+

.

•

•

MODELOS DE CLASSIFICAÇÃO – AVALIAÇÃO

Como avaliar se um modelo prevê corretamente os dados?

- Matriz de erro ou Confusion Table

CONDIÇÃO
POSITIVA

TESTE
POSITIVO

TESTE
POSITIVO

TESTE
NEGATIVO

TESTE
NEGATIVO

TESTE
Negative—Erro1)

CONDIÇÃO
NEGATIVA

FP (False Positive Erro 2)

TN (True
Negative)

+

.

#### MODELOS DE CLASSIFICAÇÃO – MÉTRICAS

#### Métricas comparativas:



ACC (Accuracy) : ACC = (TP + TN)/Tot

- % acertos do total

Precision: P = TP / (TP + FP)

- % acertos das minhas "afirmativas"

Sensitivity (Recall): R = TP / (TP + FN)

 % acertos frente a condição positiva (testa presença da característica em elementos que se assume que tenham)

#### MODELOS DE CLASSIFICAÇÃO – MÉTRICAS

Métricas comparativas:

CONDIÇÃO
POSITIVA

TESTE
POSITIVO

TESTE
POSITIVO

TESTE
NEGATIVO

TESTE
NEGATIVO

TESTE
NEGATIVO

TO NEGATIVA

FP (False Positive)

TN (True
Negative)

Specificity = TN / (FP + TN)

- % de True Negative

(testa ausência da característica em elementos que se assume que não a tenham)

$$F1 = 2*TP / (2*TP + FP + FN)$$

- Métrica geral para comparar modelos

#### **CURVA ROC**

http://www.scielo.br/scielo.php?script =sci\_arttext&pid=S0034-8910200600 0600021

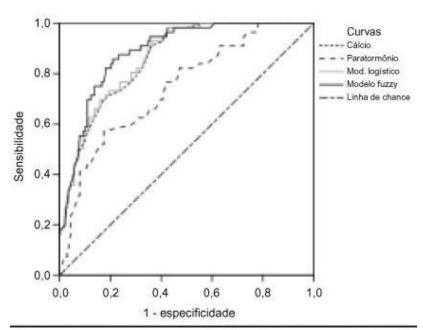


Figura 3 - Curvas ROC: Desempenho de cálcio, paratormônio, modelo logístico e modelo fuzzy na indicação de cintilografia das paratiróides.



### REGRESSÃO LOGÍSTICA

- É um Modelo de Regressão Generalizado (ou MLG, modelo linear generalizado), assim como regressão de Poisson, entre outros
- Utiliza, assim como demais MLG, um algoritmo que aplica o método da máxima verossimilhança (proposto por Nelder e Wedderburn)

# REGRESSÃO LOGÍSTICA

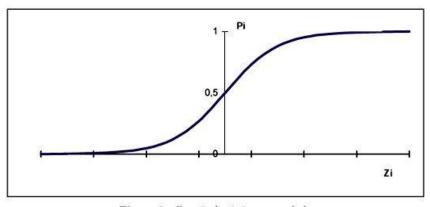
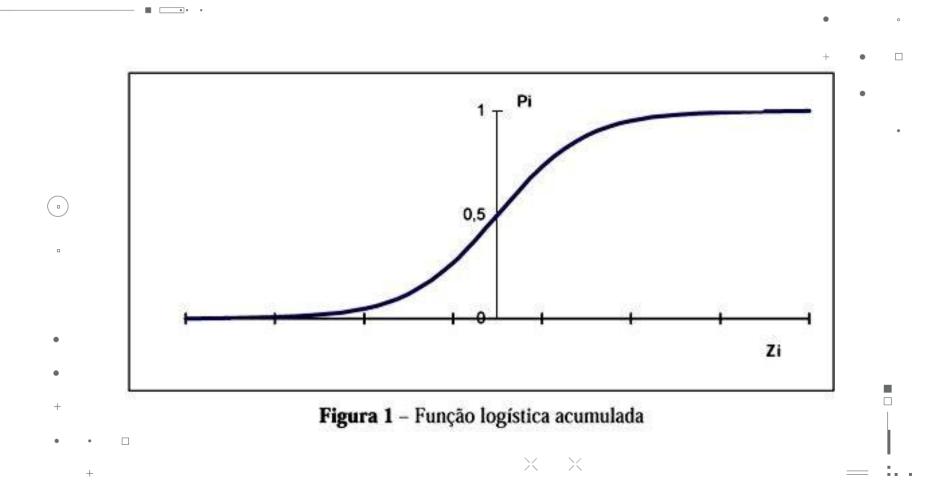


Figura 1 – Função logística acumulada

## REGRESSÃO LOGÍSTICA

- Trabalha bem como classificador
- Tem o seguinte formato:  $y = \frac{1 + e^{-f(x)}}{1 + e^{-f(x)}}$
- Como modelos uma função linear (f(X)=ax+b) basta substituir na fórmula acima e temos sua formulação clássica.



#### EXERCÍCIO

#### TITANIC COM REGRESSÃO LOGÍSTICA

```
install.packages("e1071")
library(e1071)
mRL <- glm(Survived_F~Sex+Age+Pclass+Fare, data=
train, family = binomial())</pre>
```

pred <- predict(mRL, newdata=test)
table(pred, test\$Survived)</pre>

table(ifelse(pred<0,0,1), test\$Survived)

3300030



.

Uma evolução da técnica de árvores de decisão, que procura resolver o problema de overfit. Criado por Brieman & Cutler é uma técnica mista (ensemble learning) que compreende entre outras técnicas prune (poda) e bagging (bootstrap aggregating).

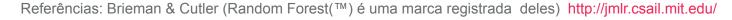
#### Em geral:

- Corrige overfit de árvores de decisão
- Trabalha bem com pequenas amostras
- Tem boa precisão com test-sets e real-sets diferentes do train-set











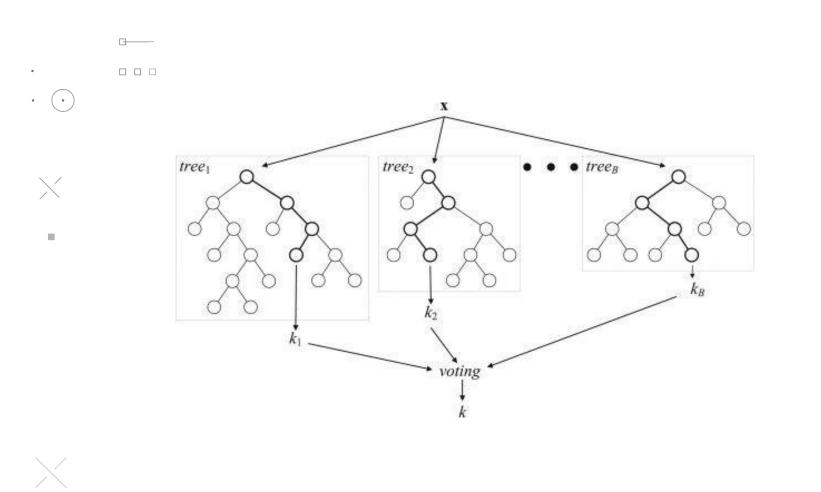
#### CONCEITO DE FUNCIONAMENTO

elaboração de um grande número de árvores de decisão (uma floresta), escolhidas aleatoriamente. É estabelecido de um critério de decisão para escolha do melhor resultado.















# TÉCNICAS UTILIZADAS

- Validação do erro com 1/3 da amostra (out of bag);
- Avaliação da relevância das variáveis para escolha;
- Avaliação da importância em relação a hierarquia;
- Execução de interações aleatórias;
- Gera os "votos" nos ramos pela avaliação de proximidade (similar ao K-Nearest, K-NN);

+

# TÉCNICAS UTILIZADAS

- Usa o cálculo de auto-funções (matricial) para acelerar em escala (grande volume de dados);
- Analisa outliers;
- Faz balanceamento de erros e trata problemas com fitting

#### EXERCÍCIO

#### TITANIC COM RANDOM FOREST

install.packages("randomForest")

library(randomForest)

mRF <- randomForest(Survived\_F~Sex+Age+Pclass+Fare, data=train, importance=TRUE, proximity=TRUE)

pred <- predict(mRF, newdata=test)
table(pred, test\$Survived)</pre>

45697056





# SVM SUPPORT VECTOR MACHINES

Como separar dois conjuntos de dados por uma linha? E da forma ótima?

Conceito de margem: max (distância entre conjuntos)

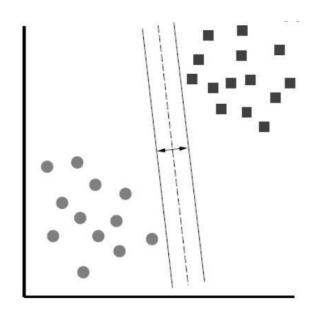
http://www.svms.org/tutorials/Berwick2003.pdf

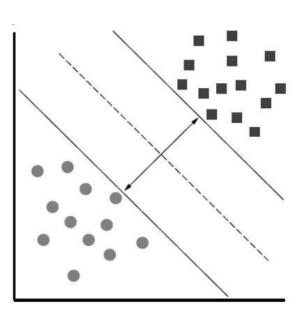






.





45697056

9500.6950

. . .



É uma "nova geração" de algoritmos, que apresenta evoluções em relação ao uso de redes neurais e outros (classificadores):

- Utilizado para classificação de SPAM's de emails (junto com o classificador de Naive Bayes)
- "Estado da arte" para utilização em Text Mining







# · • •

#### PONTOS FORTES

Não necessita percorrer todos os pontos, somente a "fronteira" da linha (ou hiperplano, plano n-dimensional) de decisão: tem melhor performance sobre grandes volumes.







•

•

•

#### ENCONTRAR A, B E C TAIS QUE:

Para vermelho, ax + by > c, e;

Para verde, ax + by < c.

http://www.svms.org/tutorials/Berwick2003.pdf

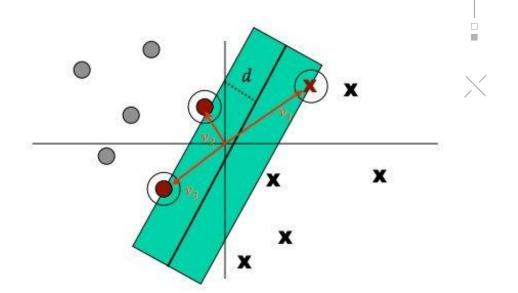
•

+

• •

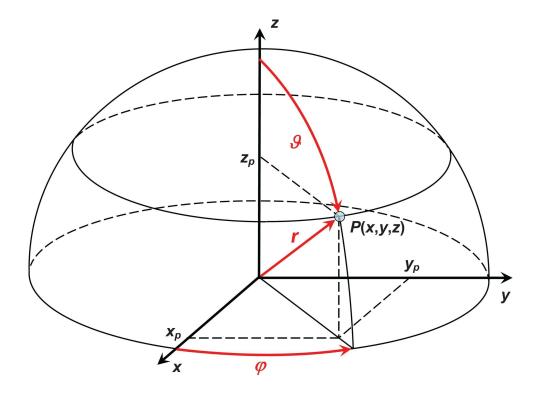












Coordenadas polares! (Cálculo 2)

http://www.svms.org/tutorials/Berwick2003.pdf

#### OUTRAS TRANSFORMAÇÕES

$$K(x,y) = (x \cdot y + 1)^{p}$$

$$K(x,y) = \exp \left\{ -\left[ \begin{array}{c|c} \| & \| \\ & 20 \end{array} \right] \right\}$$

$$K(x,y) = tanh(kx \cdot y + \delta)$$

- • [
  - +

•

•

f<sup>st</sup> is polynomial (includes x . X as special casa)

2<sup>nd</sup> is radial basis function (gaussians)

3<sup>rd</sup> is sigmoid (neural net activation function)

# Obrigado



profDiogenes.Justo@fiap.com.br



Copyright © 2018 | Diógenes Justo

Todos os direitos reservados. Reprodução ou divulgação total ou parcial deste documento, é expressamente proibido sem consentimento formal, por escrito, do professor/autor.







Técnicas de solução: multiplicador (ou operador) de Lagrange (lagrangiano). Problema matemático de otimização: encontrar máximos e mínimos de funções sujeitas a restrições.

#### LIMITE X SEM LIMITES











## ESCALANDO SOLUÇÕES

- Em ordem de crescimento:
- Hadoop / SPARK
- BD SQL e NoSQL em nuvem (Scale-Up Scale-out)
- Bancos de dados SQL e NoSQL (Scale-Up)
- CSV/Desktop Processing



Notícias, ações e muito mais





últimas mercados onde investir minhas finanças imóveis franquias negócios carreir

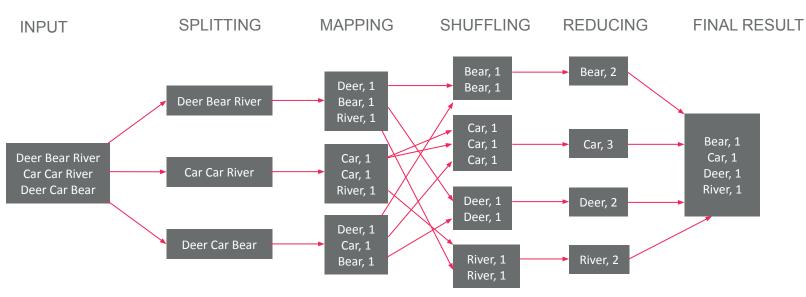
Home > Bloomberg > Mercados > Notícia Bloomberg

# O supercomputador de US\$ 10 para hedge funds que está arrasando em Wall Street

Prodígios da matemática e da computação estão levando a ciência dos dados para novos níveis em Wall Street

http://www.infomoney.com. br/bloomberg/mercados/noticia/4054301/supercomputador-para-hedge-funds-que-esta-arrasando-wall-street

#### THE OVERALL MapReduce WORD COUNT PROCESS



. . .

http://research.google.com/archive/mapreduce-osdi04-slides/index.html

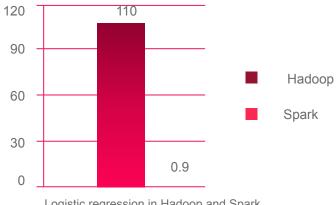
Hadoop e Spark

#### **SPEED**

. . .

Run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk.

Spark has an advanced DAG execution engine that supports cyclic data flow and in-memory computing.



Hadoop e Spark

#### EASE OF USE

Write applications quickly in Java, Scala, Python, R.

Spark offer over 80 high-level operators that make it easy to build parallel apps. And you can use it *interactively* from the Scala, Python and R shells.

```
text_file - spark.textFile("hdfs://...")
```

```
text_file. flatMap(lamba line: line.split())
.map(lambda word: (word, 1))
.reduceByKey (lambda a, b: a+b)
```

Word count in Spark's Python API

997/59



# SPARK 1.6 (COM SPARKR)

A partir da versão 1.4, lançada em jun/2015, incorporado o SparkR

https://spark.apache.org/docs/1.6.0/sparkr.html







# Looking forward

We have many other features planned for SparkR in upcoming releases: these include support for high level machine learning algorithms and making SparkR DataFrames a stable component of Spark.

The SparkR package represents the work of many contributors from various organizations including AMPLab, Databricks, Alteryx and Intel. We'd like to thank our contributors and users who tried out early versions of SparkR and provided feedback. If you are interested in SparkR, do check out our talks at the upcoming Spark Summit 2015.

### **# CLUSTER SPARK - CUSTO CLOUD**

3.xlarge			aws.amazon.com/pt/ec2/pric			
		PERÍODO	DE 1 ANO			
Opção de pagamento	Inicial	Mensalmente*	Efetiva por hora**	Economia em relação ao On- Demand	On-Demand por hora	
Sem taxas iniciais	\$0	\$138.7	\$0.19	29%		
Adiantado parcial	\$842	\$51.1	\$0.1662	38%	\$0.266 por hora	
Adiantado integral	\$1428	\$0	\$0.1631	39%		
		PERÍODO	DE 3 ANOS			
Opção de pagamento	Inicial	Mensalmente*	Efetiva por hora**	Economia em relação ao On- Demand	On-Demand por hora	
Adiantado parcial	\$1345	\$43.8	\$0.1112	58%	\$0.266 por hora	
Adiantado integral	\$2746	\$0	\$0.1045	61%		

### zM3



Esta família inclui os tipos de instância M3, que oferecem recursos equilibrados de computação, memória e rede, e são uma boa opção para diversos aplicativos.

#### Recursos:

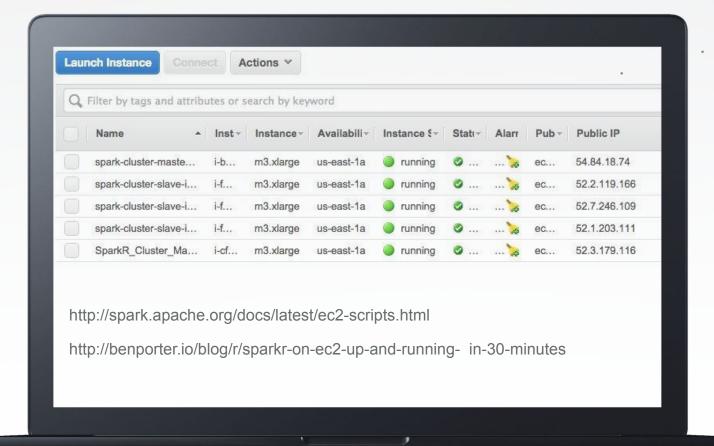
- Processadores Intel Xeon E5-2670 v2 (Ivy Bridge) de alta frequência\*
- Armazenamento em instância baseado em SSD para alto desempenho de E/S
- Equilíbrio entre recursos de computação, memória e rede

Modelo	vCPU	Mem (GiB)	Armazenamento em SSD (GB)
m3.medium	1	3,75	1 x 4
m3.large	2	7,5	1 x 32
m3.xlarge	4	15	2 x 40
m3.2xlarge	8	30	2 x 80

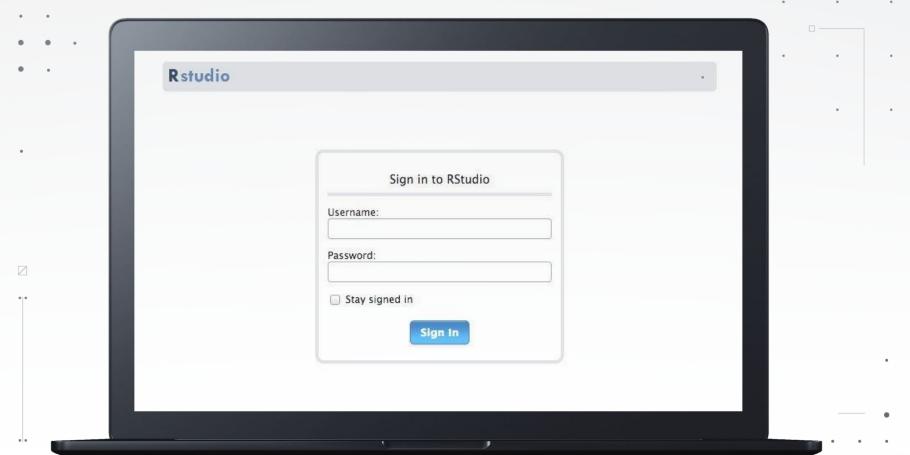
#### Casos de uso

Panaga do dados do pagueno o mádio parto, tarafas do processamento do dados que evigam memário

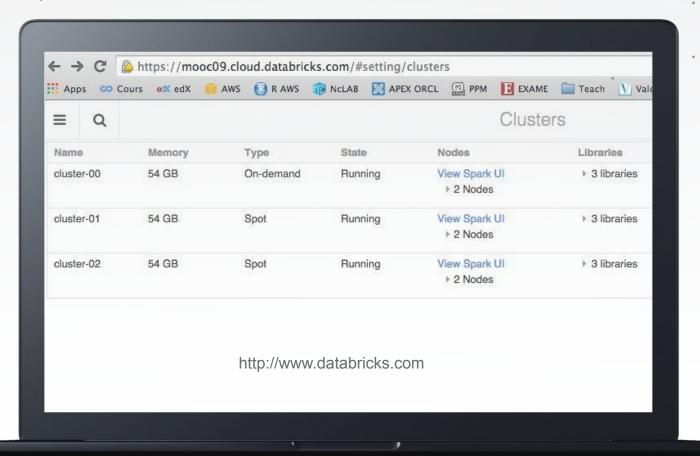
### CLUSTER SparkR+RStudio no EC2



\* CLUSTER SparkR+RStudio no EC2



### **CLUSTER Spark NA DATABRICKS**



### Otimizadas para memória

### R3

As instâncias R3 são otimizadas para aplicativos que usam muita memória e têm o custo mais baixo por GiB de RAM dentre os tipos de instância do Amazon EC2.

#### Recursos:

- Processadores Intel Xeon E5-2670 v2 (Ivy Bridge) de alta frequência
- · Custo mais baixo por GiB de RAM
- · Armazenamento em SSD
- · Suporte a redes aperfeiçoadas

Modelo	vCPU	Mem (GiB)	Armazenamento em SSD (GB)
r3.large	2	15.25	1 x 32
r3.xlarge	4	30.5	1 x 80
r3.2xlarge	8	61	1 x 160
r3.4xlarge	16	122	1 x 320
r3.8xlarge	32	244	2 x 320

Scale-up pode ser suficiente

# **Experiment Setup:** http://nerds.airbnb.com/redshift-performance-cost/

Redshift — 16 node cluster (\$13.60 per hour — \$0.85 per hour per node) Node Type: dw.hs1.xlarge CPU: 4.4 EC2 Compute Units (2 virtual cores) per node Memory: 15 GiB per node I/O Performance: Moderate Platform: 64-bit Storage: 3 HDD with 2 TB of storage per node

Hive/EMR — 44 node cluster (\$57 per hour) Node Type: Cluster Compute Quadruple Extra Large Memory: 23.00 GB CPU: 33.5 (2xIntel Xeon X5570) Storage: 1690 GB (2x840 GB) Platform: 64-bit I/O Performance: Very High

#### Test 1:

Runtime Hive: 28 minutes Redshift: less than 6 minutes

Simple range query against a giant table with 3 billion rows, we saw 5x performance improvement over Hive!

Test 2: For the following slightly more complex query that has two joins with millions of rows.

Há alternativas par

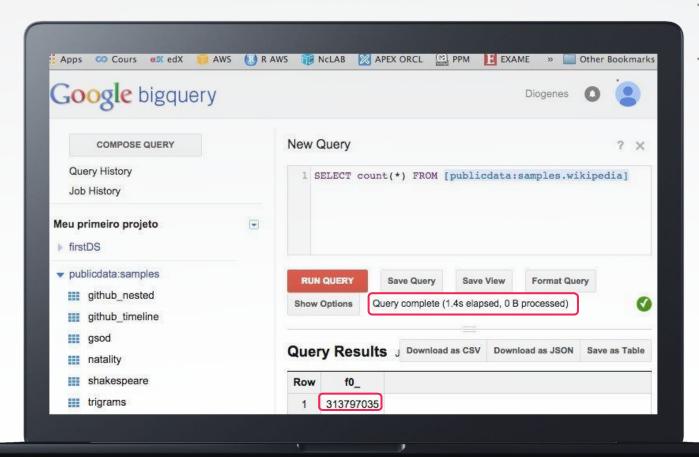
Runtime Hive: 182 seconds Redshift: 8 seconds

Há alternativas para escalar soluções em NoSQL fora Hadoop's flavor

Redshift is 20x faster than the Hive version of the query!

Results As shown above the performance gain is pretty significant, and the cost saving is even more impressive \$13.60/hour versus \$57/hour. This is hard to compare due to the different pricing models, but check out pricing

### EU PRECISO DE UM CLUSTER?





## CONSIDERAÇÕES PARA USO DO R EM SCALE-UP

Verificar os pacotes adequados:

- Monitorar uso de memória (pryr, mem\_used())
- Avaliar tempo de processamento
- Read.csv: usar ff package
- Uso do data.table

# • EXERCÍCIO

Trabalhando com volumes de dados, monitore memória

one <- 1:10e5

print(object.size(one), units = "GB")

#Ou acesse Environment, opção grid

install.packages("pryr")

library(pryr)

mem\_used()

#Remover objetos rm(one)









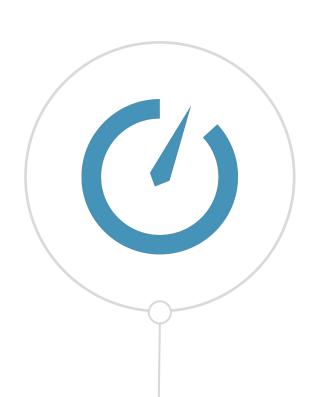


# **EXERCÍCIO**

Monitore tempo de execução

#Opção 1: com execução Sys.time() # Comandos Sys.time()

#Opçã 2: sem execução system.time(#Comandos)









# EXERCÍCIO

Carga de dados com ff package (não utiliza memória principal, deposita o data.frame em disco)

install.packages("ff")
library(ff)

ldf <- read.csv.ffdf(file="large.csv",</pre>

nrows=1000000,sep="#", header=TRUE, first. rows=5000, next.rows=50000, VERBOSE=TRUE)

#next.rows divide em leituras parciais #Verbose mostra para bloco











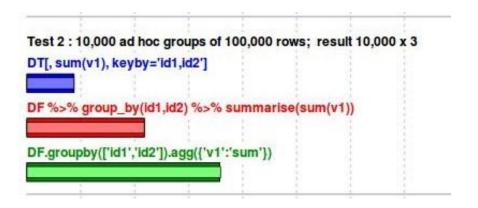
# **EXERCÍCIO**

Uso de data.table (ao invés de data.frame)

install.packages("data.table")
library(data.table)

dt <- as.data.table(df)</pre>

https://github.com/Rdatatable/data.tab le/wiki/Benchmarks-%3A-Grouping









## RECOMENDAÇÕES PARA ESTUDOS

- Processamento
  - Considere seriamente MapR: Hadoop (PRD) e Spark (Test) massivo?
  - Cursos MMDS (Stanford) e Intro Data Science (Bill Howe)
  - Estudar algos de redução de dimensão (PCA, SVD)

- Embasamento acadêmico ou desenvolvimento
  - algos?
- Processo Data
  - Cursos MMDS, Machine Learning (Andrew Ng), Science Data Science (John

Minning (Illinois at UC) Hopkins)