

# MACHINE LEARNING E DATA SCIENCE COM R DE A À Z

INSTRUTOR: JONES GRANATYR [MESTRADO E DOUTORADO  
EM IA]

MINHAS ANOTAÇÕES DE AULA # SÃO ANOTAÇÕES QUE  
FAZEM AO LONGO DO CURSO

PROF. CÍCERO QUARTO - CICERO@ENGCOMP.UEMA.BR

DEPARTAMENTO DE ENGENHARIA DA COMPUTAÇÃO  
UNIVERSIDADE ESTADUAL DO MARANHÃO (UEMA)

IA (INTELIGÊNCIA ARTIFICIAL) → PARA O INSTRUTOR,  
IA É O CAMPO DE TI  
QUE + ESTÁ CRESCEN-  
DO ATUALMENTE

POR QUE DE A À Z?

↑  
VAI DE CONTEÚDOS  
+ BÁSICOS INTRODUTÓRIOS  
DE MACHINE LEARNING  
↳ ATÉ CONCEITOS + AVANÇADOS  
USANDO A LINGUAGEM R

## 1. CONTEÚDO DO CURSO



~~CONTEÚDO~~

INTRODUÇÃO # CONCEITOS BÁSICOS SOBRE  
APRENDIZAGEM DE MÁQUINA

### PARTE 1 - CLASSIFICAÇÃO

PRÉ-PROCESSAMENTO, NAIVE BAYES, ÁRVORES DE  
DECISÃO, REGRAS, INSTÂNCIAS, REGRESSÃO LOGÍSTICA,  
MÁQUINAS DE VETORES DE SUPORTE, REDES NEURAIS  
ARTIFICIAIS, AVALIAÇÃO DE ALGORITMOS, COMBINAÇÃO  
E REJEIÇÃO DE CLASSIFICADORES, BASE DE DADOS  
BOM OU MAL PAGADOR, BASE DE DADOS PARA  
PREVISÃO DA CLASSE DO SALÁRIO.

### PARTE 2 - REGRESSÃO

REGRESSÃO LINEAR, REGRESSÃO POLINOMIAL, ÁRVO-  
RES DE DECISÃO, RANDOM FOREST, VETORES DE SU-  
PORTE, REDES NEURAIS ARTIFICIAIS, BASE DE DADOS  
PLANO DE SAÚDE, BASE DE DADOS PREVISÃO DO  
PREÇO DE CASAS.

### PARTE 3 - REGRAS DE ASSOCIAÇÃO

ALGORITMO APRIORI, ALGORITMO ECLAT, BASE DE DADOS  
DO MEREADO

## PARTE 4 - AGRUPAMENTO

K-MEANS, AGRUPAMENTO HIERÁRQUICO, DBSCAN, BASE DE DADOS DE GASTOS NO CARTÃO DE CRÉDITO.

## PARTE 5 - TÓPICOS COMPLEMENTARES

**Observação:** O instrutor ressalta que é recomendável fazer primeiro a parte 1 e só depois fazer a parte 2, mesmo que o aluno queira pulá-la diretamente para a parte 2. Ou seja, o instrutor destaca que é melhor entender classificação para depois passar para regressão.

### Pré-requisitos

- Conhecimento sobre lógica de programação, principalmente estruturas condicionais e de repetição;
- Conhecimentos básicos em R são desejáveis, embora seja possível acompanhar o curso sem saber essa linguagem em profundidade;
- Conhecimentos básicos sobre instalação de softwares básicos;
- Nível: todos os níveis

## APLICAÇÕES DE MACHINE LEARNING

### RECONHECIMENTO FACIAL

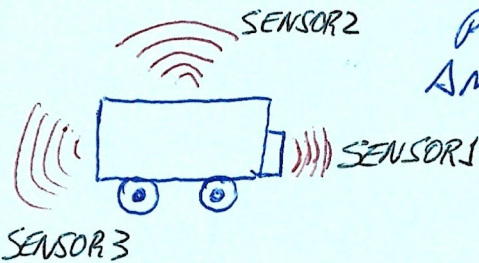


### KINECT (MICROSOFT) # REPRODUÇÃO DE MOVIMENTOS POR CAPTURA DE IMAGENS

↳ O COMPUTADOR VAI (VIA CÂMERA) APRENDER COMO VOCÊ JOGA E AI ELE PODE EUFAR UMA JOGADA SUA, OU SEJA, MAREAR UM GOL DE CERVEJA AO COBRAR UM ESCANTEIO.

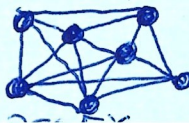
### HUMANOIDES DO GOOGLE → DETECTAM PERCEPÇÕES AO SEU REDOR PARA AGIR SOBRE O AMBIENTE

### CARROS AUTÔNOMOS → USAM SENSORES PARA DETECTAR PERCEPÇÕES E AGIR SOBRE O AMBIENTE



### ROBÔ MÉDICO → NOVOS MEDICAMENTOS, PREVER DOENÇAS NO FUTURO

### COLABORAÇÃO → REDE DE COMUNICAÇÃO COLABORATIVA PARA SISTEMAS DE RECOMENDAÇÃO. EX, NETFLIX RECOMENDA FILMES AFINS PARA ATENDER TENDÊNCIAS DE GASTOS



RECOMENDAR MÚSICAS BASEADO NO PERFIL DO CLIENTE (USUÁRIO) QUE ELE TEM MAIS TENDÊNCIA DE GOSTAR;

■ HARMONY → SITE SÓ DE ENCONTROS. VOCÊ FAZ O CADASTRO DE SUAS CARACTERÍSTICAS PESSOAIS, FÍSICAS, SOCIAIS, GEOGRÁFICAS, DENTRE OUTRAS PARA ENCONTRAR PARCEIROS OU PARCEIRAS AFINES COM O SEU PERFIL E GOSTOS, OBJETIVANDO QUE O RELACIONAMENTO SEJA CERTO

■ AMAZON.COM → QUEM COMPROU O PRODUTO A, TAMBÉM COMPROU O PRODUTO B E/OU C E AÍ SUGERE A RECOMENDAÇÃO PARA VOCÊ

■ NASA → PARA EXPLORAR UM PLANETA DESCONHECIDO, ENVIA UM ROBÔ E ESTE VAI REALIZAR VÁRIAS AÇÕES (DETECTAR OBSTÁCULOS, TIRAR FOTOS DO AMBIENTE,)

■ FACEBOOK ADS → ANÚNCIOS DO FACEBOOK, PERSONALIZANDO PARA AS PESSOAS QUE QUEREM OU PRETENDEM COMPRAR UM DETERMINADO PRODUTO

■ TEXT-TO-SPEAK → ESCRITO → SOM (REPRODUÇÃO DO SAÍ A TRADUÇÃO EM SOM) (TEXT-TO-SPEAK EM SOM) OU AO CONTRÁRIO → STEPHEN HAWKING



■ CIA → CENTRAL INTELLIGENCE AGENCY  
↳ USAM ALGORITMOS DE APRENDIZAGEM DE MÁQUINA PARA DETECTAR QUEM TEM OU NÃO CHANCE DE SER UM TERRORISTA, CRIMINOSO

■ REALIDADE VIRTUAL (RV) → USO DE ÓCULOS PARA REPRODUZIR A REALIDADE.

■ GOOGLE → USA ALGORITMOS DE APRENDIZAGEM DE MÁQUINA PARA AUTOMATIZAR BUSCAS

ALGUMAS NOTÍCIAS      SOFTWARES PARA INFECTAR O COMPUTADOR      ← MALICIOUS SOFTWARE (SOFTWARES MALICIOSOS)

• COMO A IA PODE AJUDAR PREVENIR MALWARE  
VIRUS, WORMS, CAVALOS DE TROIA, SPYWARE, ETC...

• SITE PORNOGRÁFICO VAI USAR IA PARA IDENTIFICAR ATRIZES E ATORES

• IA PODE AJUDAR A REDUZIR ROUBO DE ELETRICIDADE NO BRASIL

• IA TENTA DESCOBRIR QUEM TRAIU ANNE FRANK

• GOOGLE USA IA PARA ENCONTRAR VÍRUS NA PLAY STORE

• TESLA E AMD ESTÃO CRIANDO CHIP DE IA PARA CARROS AUTÔNOMOS (TENDÊNCIA EM ALTA)

• MICROSOFT PIX EDITA IMAGENS ATRAVÉS DE IA;

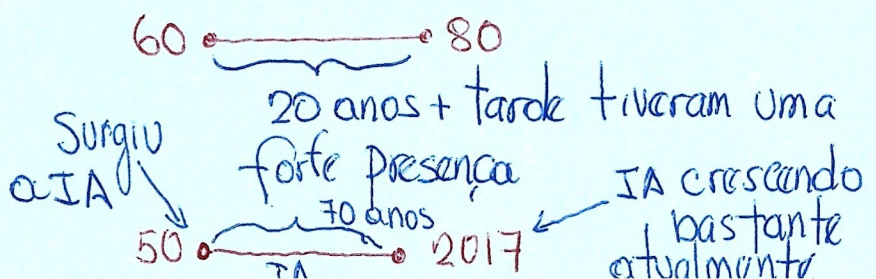
• IA DETECTA MAL DE ALZHEIMER UMA DÉCADA ANTES DE SINTOMAS APARECEREM

- IA ESTÁ ESCRIBENDO O FIM DE GAME OF THRONES
- IA CONSEGUE REERIR FASES DE SUPER MÁRIO;
- ROBÔ FAZ EM SEGUNDOS O QUE DEMORAVA 360 MIL HORAS PARA UM ADVOGADO;
- IA, UMA DAS TENDÊNCIAS QUE REVOLUCIONAM O SERVIÇO DE ATENDIMENTO AO CLIENTE
- CANTORES ESTÃO COMENDO UM ALBUM INTEIRO USANDO IA

PARA FINALIZAR O TÓPICO DE AVULS SOBRE APLICAÇÕES DA IA, O INSTRUTOR RESGATOU UMA MATÉRIA DA REVISTA VEJA INTITULADA "DE MÃOS DADAS COM A INTELIGÊNCIA ARTIFICIAL"

Por que aprender Machine Learning?

1. O instrutor faz um resgate da linha do tempo dos computadores desde quando surgiram na década de 60, com expansão na década de 80.



Redes Neurais + Deep Learning => Forte tendência de uso

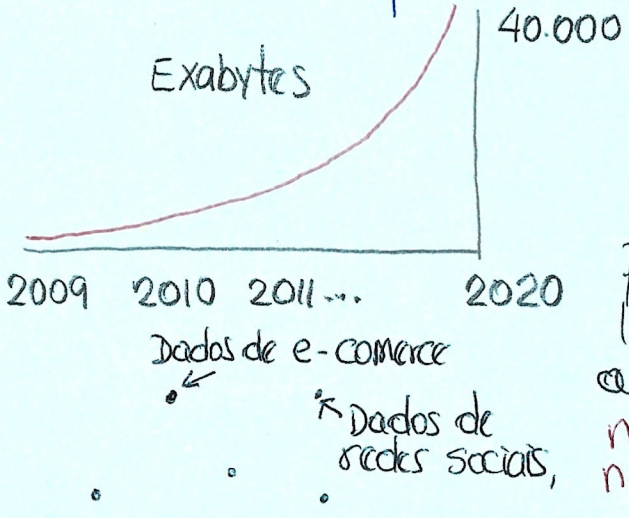
1956 -> Um HD de computador precisava ser transportado de avião, pois era muito grande e pesado

1980 -> HD de 10 Megabyte, entanto <, \$3,495

2017 -> Seagate Ultra de 256GB, entretanto + caro

Storage Limits # Neste item, o instrutor resgatou a linha do tempo da Capacidade de armazenamento de dados baseado em HD (Hard Disk), Flash memory e Bacterial DNA.

\* Exabytes > Capacidade de armazenamento



Moral da História: temos cada vez + dados e menos capacidade de profissionais para manipular esses dados... e aí entra o Machine Learning para ajudar a manipular esses dados.

Empresas Brasileiras que precisam de profissionais de Análise de dados e de Machine Learning.

- Itaú e/ou qualquer outro banco estão desenvolvendo aplicações que usam Machine Learning [BRÁCI];

- IBM #Contém Laboratório de estudos nas áreas:

- ▶ Gestão Inteligente de Recursos Naturais (petróleo, mineração, agricultura);
- ▶ Computação Cognitiva com ênfase em análise de dados sociais;
- ▶ Big-Data;
- ▶ Visual analytics e Comprehension;
- ▶ Ciência e Tecnologia para Aplicações Industriais voltada para ~~Aplicações~~ recursos naturais, ciências da vida e internet das Coisas.

### Lab da IBM no Brasil

atua de forma integrada c/a rede global da IBM Research, formada por 12 Laboratórios no mundo onde trabalham mais de 3 mil cientistas.

IBM.biz/brazilresearchlab

Deep Learning  
↓

- NVIDIA # Trabalha muito c/ RNAs e DL  
• GPU (Unidade de processamento Gráfico)

Disponibilizam vagas no LinkedIn nas áreas de IA → Machine Learning

### Algumas Notícias

- Empresas recorrem à IA para melhorar negócios;
- As 100 startups que lideram a revolução da IA;
- Sales force vai investir U\$50 milhões em IA;
- Brasil está atrasado em IA;
- A IA no centro da guerra dos smartphones;
- Intel trabalha com Facebook no projeto de chip para IA

- Brainwave é o novo projeto para IA da MS, sendo que a MS (Microsoft) abriu no Brasil um laboratório somente de IA;
- Ainda estamos na Internet discada de onde a IA chegará;
- China quer se tornar líder em IA até 2025;
- Mastercard adquire empresa que desenvolve software de IA;
- US\$300 a US\$500 mil dólares?? # Um profissional de IA, nos EUA, pode ganhar isso/ano

## TERMINOLOGIA DA ÁREA DE IA

### INTELIGÊNCIA ARTIFICIAL (IA)

▶▶ **SISTEMAS ESPECIALISTAS** → CONSIDERADOS UM DOS PRIMEIROS ESFORÇOS DA ÁREA DE IA QUE CONSISTE DA CONSTRUÇÃO DE SISTEMAS DE BASE DE CONHECIMENTO BASEADO NO CONHECIMENTO DE UM ESPECIALISTA HUMANO → PEGA TODO O CONHECIMENTO QUE O ESPECIALISTA TEM E JOGA PARA O SISTEMA COMPUTACIONAL, PARA ESTE SER CAPAZ DE DAR AS RESPOSTAS A PERGUNTAS FEITAS

- ▶▶ **VISÃO COMPUTACIONAL** → SIMULAR A VISÃO HUMANA PARA DETECÇÃO DE OBJETOS, RECONHECIMENTO FACIAL, TEM APLICAÇÕES NAS ÁREAS DE SEGURANÇA, ROBOTICA, ETC...
- ▶▶ **MACHINE LEARNING** → APRENDIZAGEM DE MÁQUINA, CONSISTINDO DE VÁRIOS ALGORITMOS QUE VÃO FAZER OS COMPUTADORES A APRENDEREM ATRAVÉS DE UMA BASE DE DADOS
- ▶▶ **PROCESSAMENTO DE LINGUAGEM NATURAL** ENVOLVE BASICAMENTE O COMPUTADOR ENTENDER A SEMÂNTICA DE UMA FRASE, TANTO ESCRITA COMO FALADA (INTERPRETAÇÃO DE TEXTO)
- ▶▶ **ALGORITMOS GENÉTICOS** →
- ▶▶ **SISTEMAS MULTIAGENDES**
- ▶▶ **MINERAÇÃO DE DADOS** # O INSTRUTOR RESSALTA QUE É UMA ≠ SUÁIL ENTRE MACHINE LEARNING E MINERAÇÃO DE DADOS.
- ▶▶ **RESOLUÇÃO DE PROBLEMAS POR MEIO DE BUSCA** EXEMPLOS, COMO ROTAS (GOOGLE MAPS) PARA ENCONTRAR MELHOR CAMINHOS
- ▶▶ **LÓGICA FUZZY (NEBULOSA)** → TEM MUITA APLICAÇÃO NA ÁREA INDUSTRIAL
- ▶▶ **RACIOCÍNIO BASEADO EM CASOS**
- ▶▶ **REDES NEURAIS**

▶▶ **ROBÓTICA** → PODE SER CONSIDERADA UMA ÁREA INDEPENDENTE, ENTRETANTO PARA FAZERMOS APLICAÇÕES MAIS INTERESSANTE VAI TER QUE INTEGRAR IA NESSES ROBÔES

• INTELIGÊNCIA ARTIFICIAL (IA)

• ÁREA DA CIÊNCIA DA COMPUTAÇÃO RESPONSÁVEL PELO DESENVOLVIMENTO DE SISTEMAS QUE SIMULEM A CAPACIDADE HUMANA DE RESOLVER PROBLEMAS

TEMOS NO BRASIL O EVENTO ENIAC (ENCONTRO NACIONAL DE IA E COMPUTACIONAL.

• TERMO + GERAL

• INTELIGÊNCIA COMPUTACIONAL

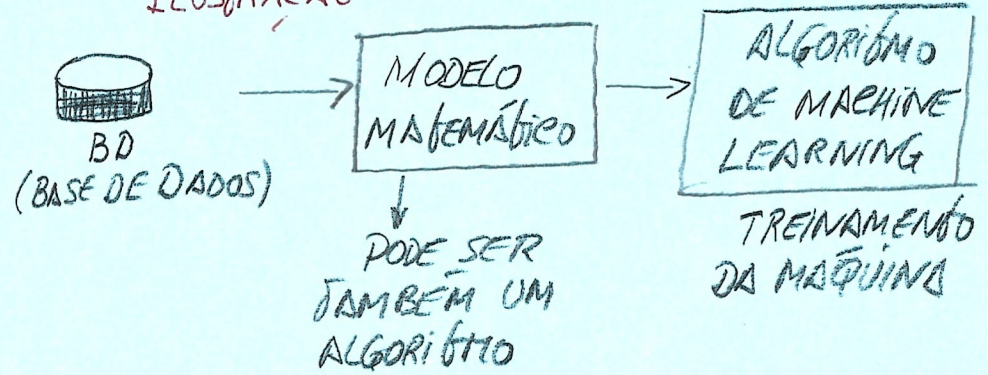
- REDES NEURAIS, COMPUTAÇÃO EVOLUCIONÁRIA, LÓGICA NEBULOSA (FUZZY)
- BOM CANDIDATO PARA SUBSTITUIR O TERMO "IA"

NOTA DO INSTRUCTOR } O INSTRUCTOR FRISA QUE MUITOS CIENTISTAS DA COMPUTAÇÃO NÃO GOSTAM E/OU EVITAM USAR O TERMO IA, POIS ESTES ACREDITAM QUE ESTE TERMO NÃO REPRESENTA OU NÃO TRAZ NADA DE COMPUTAÇÃO EM "INTELIGÊNCIA ARTIFICIAL". O INSTRUCTOR PONTUA QUE ∃ UMA CONFERÊNCIA DO IEEE QUE É ESPECIALIZADA EM "INTELIGÊNCIA COMPUTACIONAL" E ESTA SÓ ACEITA TRABALHOS QUE ABORDAM SOBRE RNA, COMPUTAÇÃO EVOLUCIONÁRIA, LÓGICA

NEBULOSA (FUZZY).

- MACHINE LEARNING (APRENDIZAGEM DE MÁQUINA)
- MÉTODOS MATEMÁTICOS PARA TREINAR ALGORITMOS.

ILUSTRAÇÃO



• DATA MINING (MINERAÇÃO DE DADOS)

• EXTRAIR CONHECIMENTO DE BASE DE DADOS, USANDO MÉTODOS DE APRENDIZAGEM DE MÁQUINA

• REDES NEURAIS

• É UM TIPO DE APRENDIZAGEM DE MÁQUINA É UM ALGORITMO QUE PODE SER UTILIZADO PARA FAZER O COMPUTADOR APRENDER

• DEEP LEARNING (APRENDIZADO PROFUNDO)

- MUITO MAIS DADOS E PROCESSADORES MAIS POTENTES
- REDE NEURAL COM MUITAS CAMADAS.

- Big Data # demanda de algoritmos de Machine Learning para extrair conhecimentos
- Imenso volume de dados dos dados
- Ciência de Dados
  - Exploração e análise de dados
  - Ciência de Computação + Estatística
  - Machine Learning

### MÉTODOS PREDITIVOS

### Aprendizagem de Máquina (MACHINE LEARNING)

Métodos Preditivos

Métodos Derivativos

Classificação

Regressão

Associação

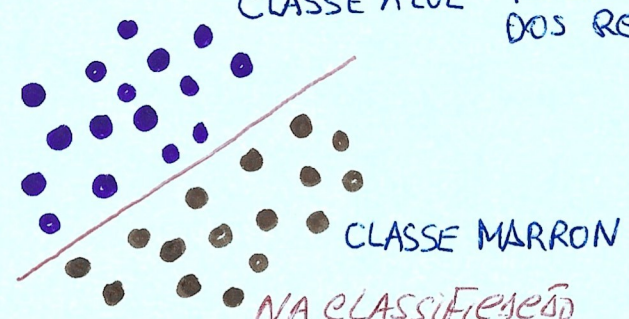
AGRUPAMENTO

DETECÇÃO DE DESVIOS

PADRÕES SEQUENCIAIS

SUMARIZAÇÃO

**CLASSIFICAÇÃO** # DEFINE CLASSES PARA ESSA UM DOS REGISTROS



NA CLASSIFICAÇÃO TEMOS RÓTULOS.

EXEMPLOS:

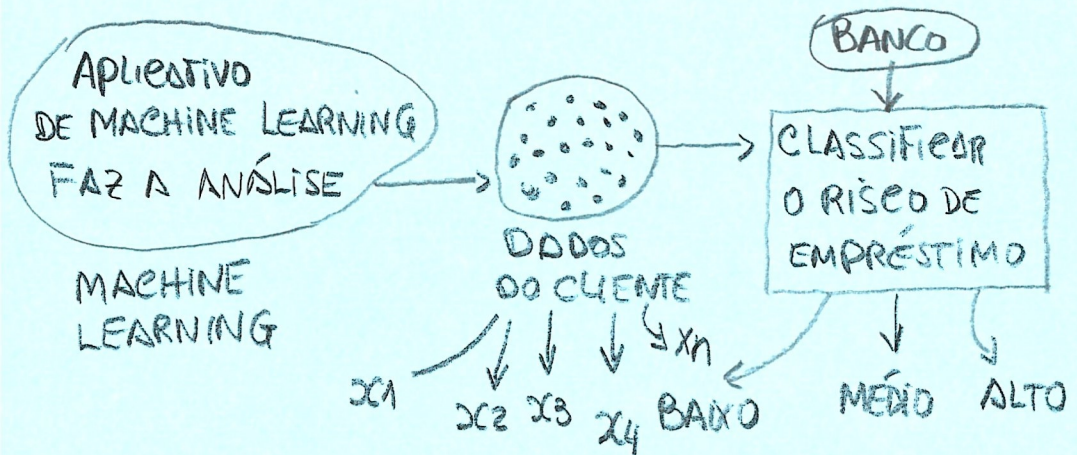
- **MARKETING DIRETO** → VOCÊ TEM UM NOVO PRODUTO E NÃO SABE PARA ONDE MANDAR?

O PROBLEMA É CLASSIFICAR SE A PESSOA VAI COMPRAR OU NÃO O PRODUTO

- **INSATISFAÇÃO DE CLIENTES** → PODEMOS CLASSIFICAR SE O CLIENTE ESTÁ SATISFEITO OU INSATISFEITO COM ALGUM PRODUTO, COM O ATENDIMENTO, PRESTAÇÃO DE SERVIÇOS

- **RISCO DE CRÉDITO** → BASTANTE UTILIZADO PELOS BANCOS PARA SABER SE O CLIENTE APRESENTA RISCO BAIXO, MÉDIO OU ALTO PARA LIBERAR UM EMPRÉSTIMO.





- **FILTROS DE SPAM** → UTILIZADOS EM TODOS OS SISTEMAS DE E-MAIL, PEGANDO UMA MENSAGEM E CLASSIFICANDO ESTA COMO SPAM, PROMOÇÕES, SOCIAIS, ... OU SEJA EM PÁSTAS CORRESPONDENTES.
- **SEPARAÇÃO DE NOTÍCIAS** → DADO VÁRIOS TEXTOS OU DOCUMENTOS PODEMOS SEPARAR AS NOTÍCIAS EM ESPORTE, POLÍTICA, SAÚDE, TECNOLOGIAS, ECONOMIA, ...
- **RECONHECIMENTO DE VOZ**: UMA VOZ PERTENCE A UMA PESSOA A, B, C, ...
- **RECONHECIMENTO DE FACE**: UMA DETERMINADA FACE DE UM SISTEMA DE SEGURANÇA PERTENCE À PESSOA A, B, C, ...
- **PREVISÃO DE DOENÇAS** → BASEADO NOS DADOS DO PACIENTE FAZER UMA PREVISÃO

SE O PACIENTE VAI DESENVOLVER OU NÃO VAI DESENVOLVER UMA DETERMINADA DOENÇA.

**REGRESSÃO** → FAZER PREVISÃO DE VALORES NUMÉRICOS. POR EXEMPLO, QUANTO UMA EMPRESA GANHARIA DE LUCRO, BASEADO NOS DADOS HISTÓRICOS

QUANTITATIVA ← VARIÁVEL NUMÉRICAS

NA REGRESSÃO TEMOS NÚMEROS

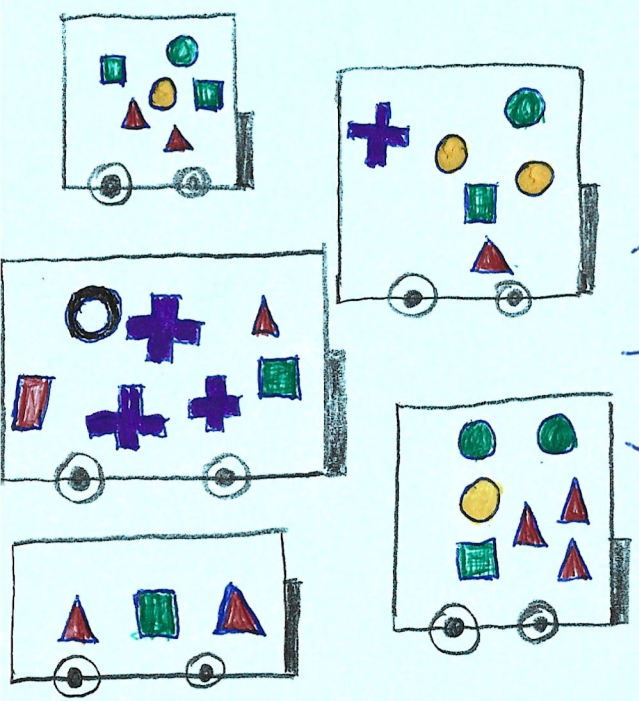
⇒ FUDO DEPOIS DA SEFA SÃO NÚMEROS.

- ALGUNS EXEMPLOS:
- GASTOS PROPAGANDA ⇒ VALOR DE VENDAS
  - TEMPERATURA, UMIDADE E PRESSÃO DO AR ⇒ VELOCIDADE DO VENTO
  - FATORES EXTERNOS → VALOR DO DÓLAR (\$)
  - RESULTADOS DO EXAME → PROBABILIDADE DE UM PACIENTE SOBREVIVER
  - RISCO DE INVESTIMENTO
  - GASTOS NO CARTÃO DE CRÉDITO, HISTÓRICO → LIMITE DO CARTÃO. OU SEJA, O CLIENTE NÃO PRECISA PEDIR PARA AUMENTAR SEU LIMITE DE CRÉDITO, POIS BASEADO EM SEU HISTÓRICO DE COMPRAS E PAGAMENTO DAS FATURAS, ISSO É PREVISÃO PELO MÉTODO DE REGRESSÃO.
  - Valores anteriores → Valores de produtos

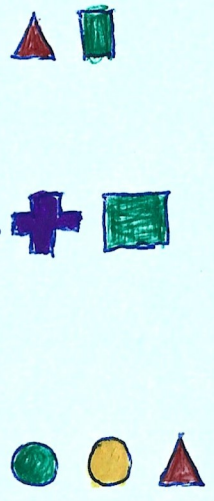
# MÉTODOS DESCRITIVOS

- Associação
- Agrupamento
- Detecção de desvios
- Padrões sequenciais
- Sumarização

## ASSOCIAÇÃO



OBJETIVO É ENCONTRAR ALGUM TIPO DE ASSOCIAÇÃO



CARRINHOS DE COMPRAS DE SUPERMERCADO

REGRAS DE ASSOCIAÇÃO (RA)  
 $RA_s = \{RA_1, RA_2, RA_3, RA_4, RA_5, \dots, RA_n\}$

RA1: QUEM COMRA UM PRODUTO  $\blacktriangle$   $\rightarrow$  COMRA TAMBEM  $\blacksquare$

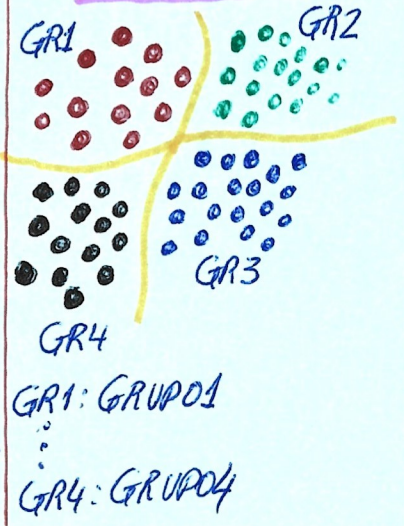
RA2: QUEM COMRA UM PRODUTO  $\bullet$  E  $\circ$   $\rightarrow$  COMRA TAMBEM  $\blacktriangle$

## APLICAÇÕES DAS REGRAS DE ASSOCIAÇÃO:

- PRATELEIRAS DE SUPERMERCADO
  - PROMOÇÕES COM ITENS QUE SÃO VENDIDOS EM CONJUNTO
- PLANEJAR CATALÓGOS DAS LOJAS E FOLHETOS DE PROMOÇÕES
- CONTROLE DE EVASÃO EM UNIVERSIDADES # PODEMOS ANALISAR OS DADOS DAS PESSOAS QUE DESISTEM DE UM DETERMINADO CURSO. POR EXEMPLO, SE A NOTA FOR X, Y: RENDA E Z: CURSO, ENTÃO ELA TEM + CHANCE DE DESISTIR OU UMA DETERMINADA PORCENTAGEM

## AGRUPAMENTO

$\rightarrow$  OBJETIVO É ANALISAR OS DADOS E ENCONTRAR GRUPOS. EXEMPLOS DESTA ÁREA:



$\rightarrow$  SEGMENTAÇÃO DE MERCADO  
 $\rightarrow$  PEGA UMA BASE DE DADOS DE CLIENTES E SEGMENTAR PARA MANDAR PROPAGANDAS SÓMENTE PARA OS CLIENTES ESPECÍFICOS, OU SEJA, MANDAR UMA PROPAGANDA DE UM NOVO VIDEO GAME PARA UM CLIENTE QUE ADORA VIDEO GAMES.

GR1: GRUPO1  
 $\vdots$   
 GR4: GRUPO4

→ → ENCONTRAR GRUPOS DE CLIENTES QUE IRÃO COMPRAR UM PRODUTO (MALA DIRETA)

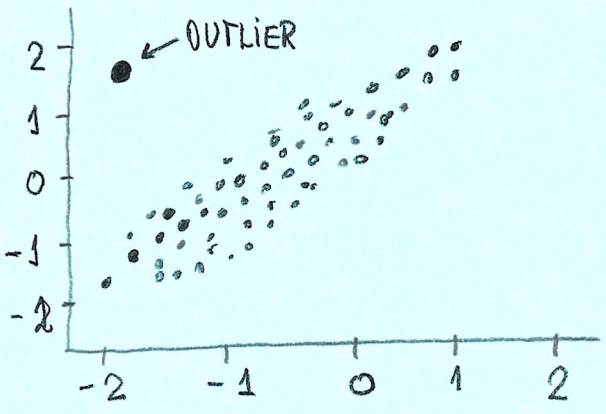
→ → AGRUPAMENTO DE DOCUMENTOS/NOTÍCIAS # PASSAR COMO PARÂMETROS VÁRIOS DOCUMENTOS E O ALGORITMO VAI AGRUPAR POR CONTEÚDO AQUELES CONTEÚDOS QUE SÃO PARECIDOS E AGRUPAR EM PASTAS NOTÍCIAS ESPECÍFICAS, COMO: NOTÍCIAS DE ESPORTE NA PASTA DE ESPORTE, ECONOMIA NA PASTA ECONOMIA, E ASSIM SUCESSIVAMENTE.

→ → AGRUPAMENTO DE PRODUTOS SIMILARES A PARTIR DE UMA BASE DE DADOS, RODA-SE ~~UMA~~ UM ALGORITMO DE MACHINE LEARNING PARA FAZER A SEPARAÇÃO DOS PRODUTOS

→ → PERFIS DE CLIENTES (NETFLIX) # MELHORAR AS RECOMENDAÇÕES DE FILMES FUTUROS EM FUNÇÃO DO PERFIL DO CLIENTE BASEADO NO HISTÓRICO DE FILMES ASSISTIDOS POR ESTE CLIENTE

→ → ANÁLISE DE REDES SOCIAIS → ENCONTRAR GRUPOS DE USUÁRIOS MAIS INFLUENTES DA REDE E ATRAVÉS DESTES ENVIAR PROMOÇÕES PARA ESTES USUÁRIOS INFLUENTES PODEREM DIVULGAR TAIS PRODUTOS

# DETECÇÃO DE DESVIOS (OUTLIERS)



## APLICAÇÕES

- FRAUDE EM CARTÃO DE CREDITO, OU SEJA, SAIU DO PADRÃO DE COMPRA QUE VOCÊ TEM;
- INTRUSÃO EM REDES # DETECÇÃO DE INVASÃO EM UMA REDE
- USO DE ENERGIA ELÉTRICA, ÁGUA OU TELEFONE # SE O CONSUMO DESTES SERVIÇOS ESTÁ AUMENTANDO MUITO PODE SER UM PROBLEMA, POIS SAIU DO PADRÃO. NO CASO DO CONSUMO ELEVADO DE ÁGUA, DIFERENTE DOS MESES ANTERIORES, PODE SER UM INDICATIVO DE UM VASAMENTO. NO CASO DE ENERGIA ELÉTRICA, PODE SER UM INDICATIVO DE ROUBO DE ENERGIA. EM AMBOS OS CASOS, A EMPRESA DE FORNECIMENTO PODE AÇIONAR UM TÉCNICO PARA INVESTIGAR O DESVIO DE PADRÃO
- Desempenho de atletas (doping) # De mãos de todo o histórico de um atleta e em uma determinada corrida, um atleta que tem um padrão de fazer uma corrida em 15 segundos e ele faz em 10 segundos,



ou seja, saiu muito do padrão dele. Nesse caso, seria interessante chamar esse atleta para fazer um exame de doping, de forma a detectar um produto ilegal que o fez aumentar seu desempenho além do padrão

### • Monitorar máquinas em um data center

Se num data Center alguma máquina não estiver funcionando corretamente ou fora do padrão, o sistema vai gerar um outlier, provocando um desempenho abaixo do padrão.

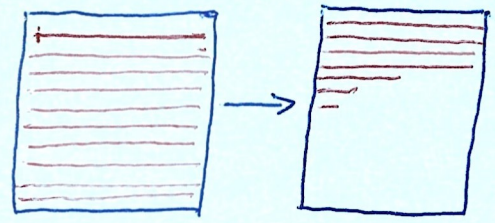
### DESCOBERTA DE PADRÕES SEQUENCIAIS

Técnicas interessantes que muitas lojas podem utilizar, como por exemplo você compra o 1º livro do Harry Potter, e depois de dois meses um e-mail é enviado para você da empresa de compra ofertando + livros do Harry Potter

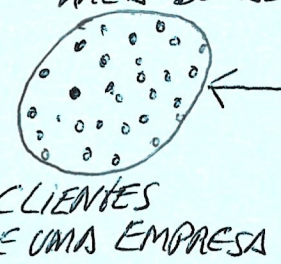
- LIVRARIAS, LOJAS DE EQUIPAMENTOS DE ATLETISMO, COMPUTADORES;
- MARKETING DIRECIONADO PARA PESSOAS QUE TEM MAIORES CHANCES DE ADQUIRIR UM NOVO PRODUTO;

- PREVENÇÃO DE DOENÇAS → SE VOCÊ SABE QUE DETERMINADOS SISTEMAS VÃO CAUSAR UMA DOENÇA, ENTÃO ISSO É UMA SEQUÊNCIA. POR EXEMPLO, SINTOMAS A, B, C, ... N VÃO CAUSAR DOENÇA X
- NAVEGAÇÃO EM SITES → VOCÊ PODE Mapear a sequência de menus acessados em seu site

### SUMARIZAÇÃO



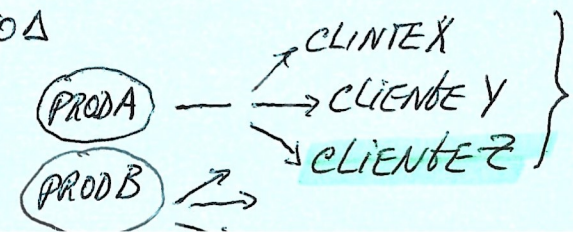
• SÃO OUVINHAS DO PROGRAMA HOMENS NA FAIXA DE 25 A 35 ANOS, COM NIVEL SUPERIOR E QUE TRABALHAM NA ÁREA DE ADMINISTRAÇÃO



APLICAR UM ALGORITMO DE SUMARIZAÇÃO

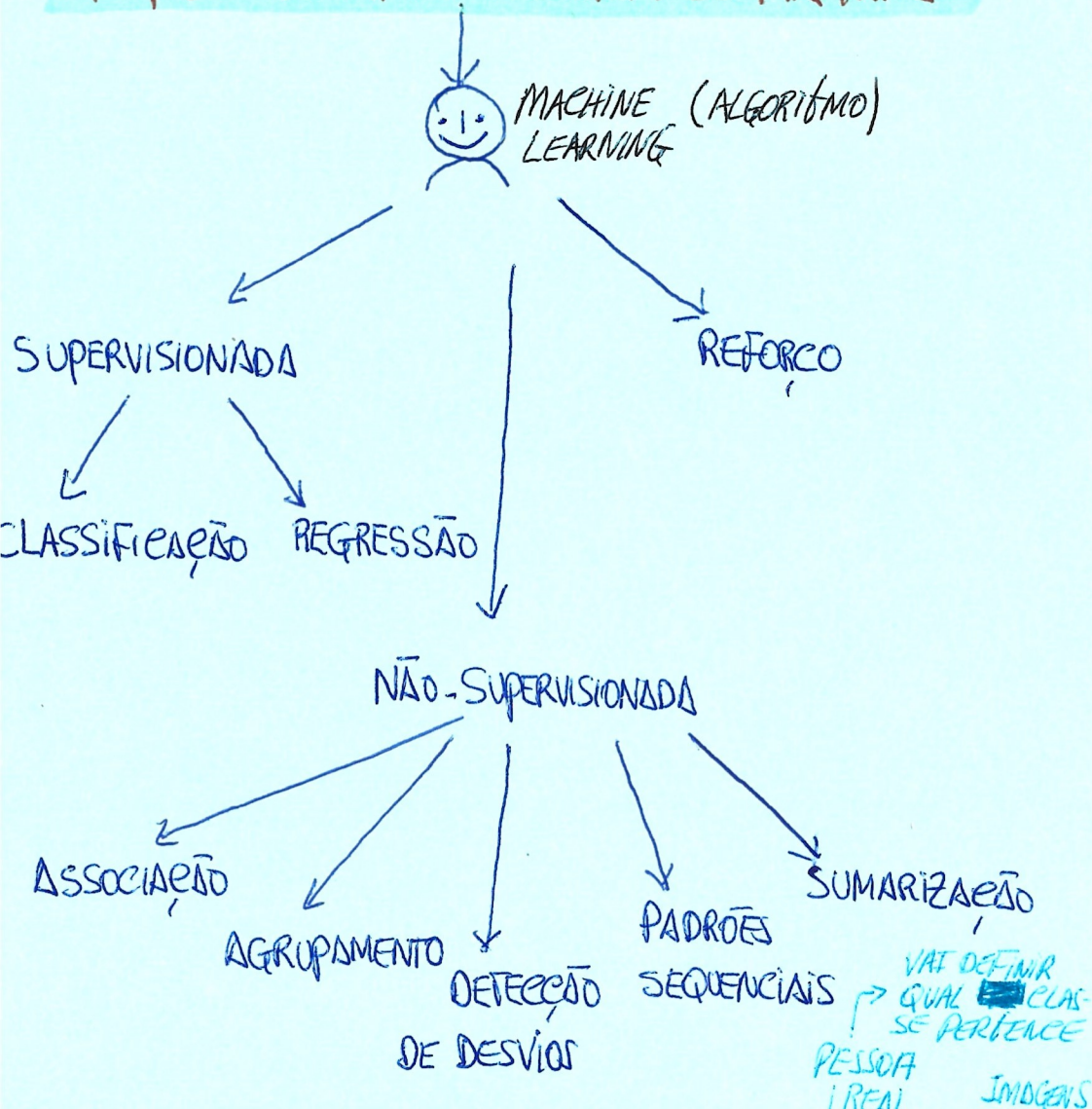
VAI PEGAR TODOS OS CLIENTES COM ESSE DETERMINADO PERFIL.

• SEGMENTAÇÃO DE MERCADO → VOCÊ PODE FAZER UMA DETERMINADA VENDA DE UM PRODUTO EM FUNÇÃO DE DETERMINADAS CARACTERÍSTICAS DA PESSOA

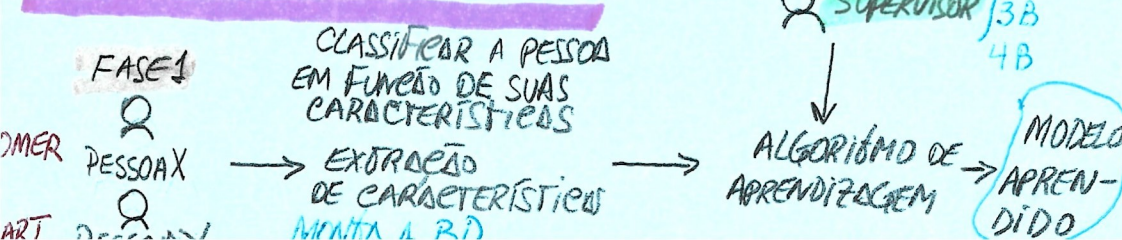


O ALGORITMO DE SUMARIZAÇÃO ACHOU QUE VAI PARA O CLIENTE Z

# TIPOS DE APRENDIZAGEM DE MÁQUINA



## APRENDIZAGEM SUPERVISIONADA



## FASE 2

PESSOA Y → EXTRAÇÃO DE CARACTERÍSTICAS → MODELO APRENDIDO → BART SIMPSON

ENTÃO, PRECISAMOS FAZER A EXTRAÇÃO DAS CARACTERÍSTICAS DAS IMAGENS E DEPOIS CLASSIFICAR O QUE TEM EM UMA IMAGEM E NÃO TEM NA OUTRA (CABECA REDONDA, CABECA ESPETADA, TEM BIGODE, CALEÇA, SHORT, COR DA ROUPA, ...)

OU SEJA, EXTRAIR TODAS AS CARACTERÍSTICAS DOS PERSONAGENS

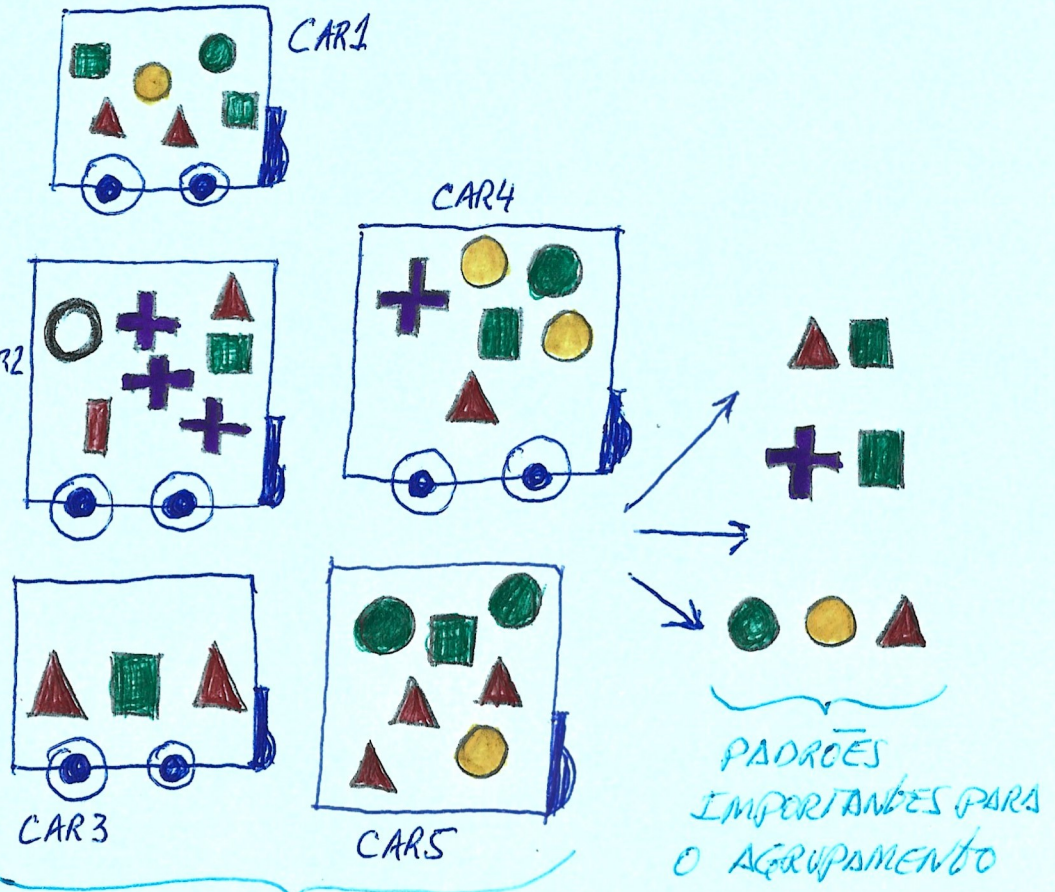
MODELO APRENDIDO → CONHECIMENTO QUE O ALGORITMO GEROU

FASE 2 → NÃO SEI A QUAL CLASSE A IMAGEM É, OU SEJA, NÃO SEI SE É UMA IMAGEM DO HOMER OU DO BART. VAMOS EXTRAIR AS CARACTERÍSTICAS E DAR PARA O MODELO APRENDIDO ESSAS CARACTERÍSTICAS PARA QUE O ALGORITMO (MODELO APRENDIDO) DÊ A RESPOSTA DE QUEM É A IMAGEM LIDA.

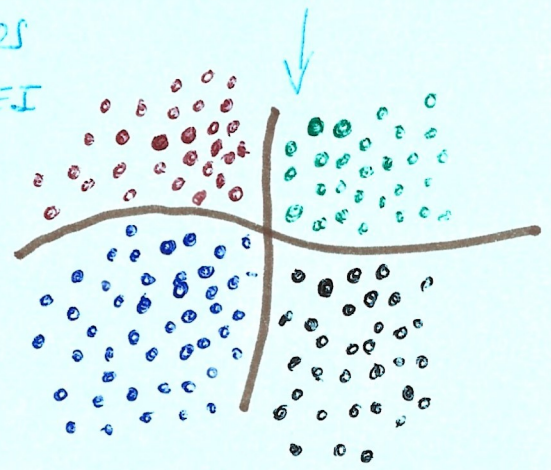
## APRENDIZAGEM NÃO-SUPERVISIONADA

- ANALISAR AUTOMATICAMENTE OS DADOS (ASSOCIAÇÃO, AGRUPAMENTO);
- NECESSITA ANÁLISE PARA DETERMINAR O SIGNIFICADO

- DO DOS PADRÕES ENCONTRADOS



DESCREVENDO OS DADOS DA BD QUE EU NÃO SEI O QUE TEM NELA



APRENDIZAGEM POR REFORÇO

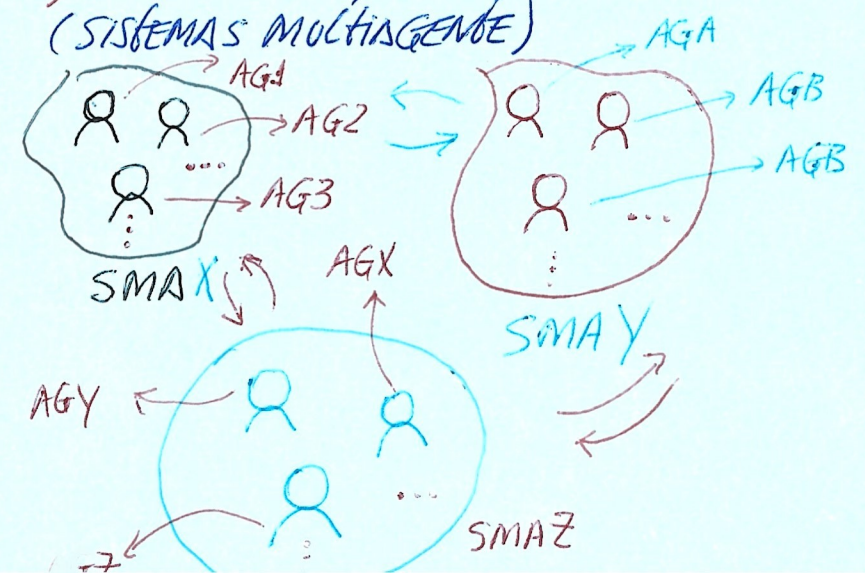
- APRENDER COM AS INTERAÇÕES COM O AMBIENTE (CAUSA E EFEITO)



APRENDER COM A SUA PRÓPRIA EXPERIÊNCIA

- APRENDER COM SUA PRÓPRIA EXPERIÊNCIA
- ROBÔ COLETANDO LIXO APRENDENDO A ANDAR EM UM AMBIENTE.

NOTA DO INSTRUCTOR } A APRENDIZAGEM POR REFORÇO TEM MUITA APLICAÇÃO NA ÁREA DE SMA (SISTEMAS MÚLTIAGENTES)

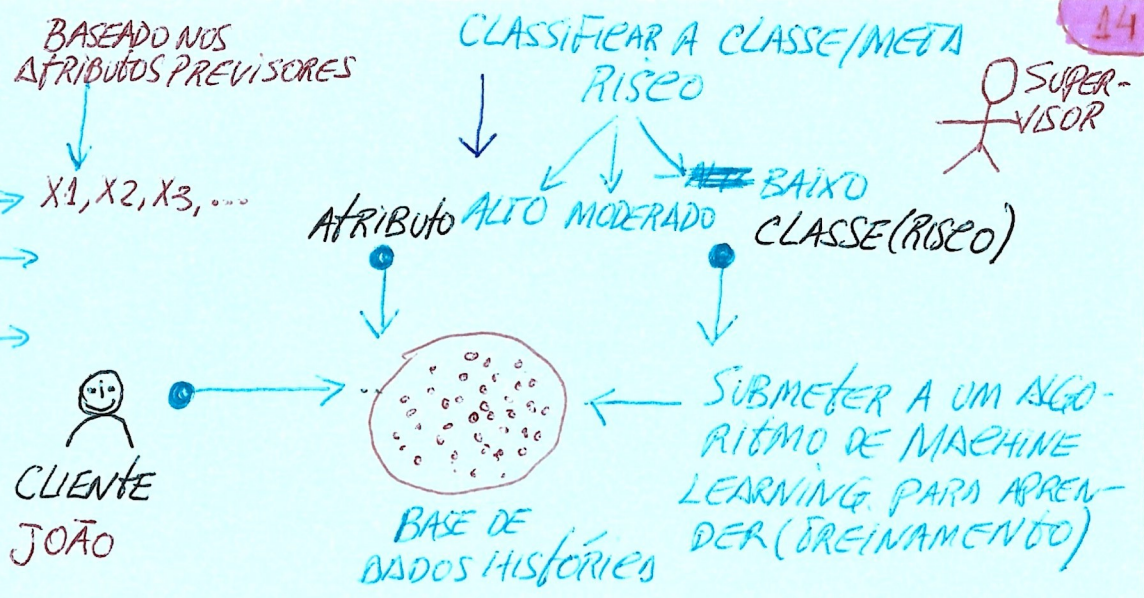
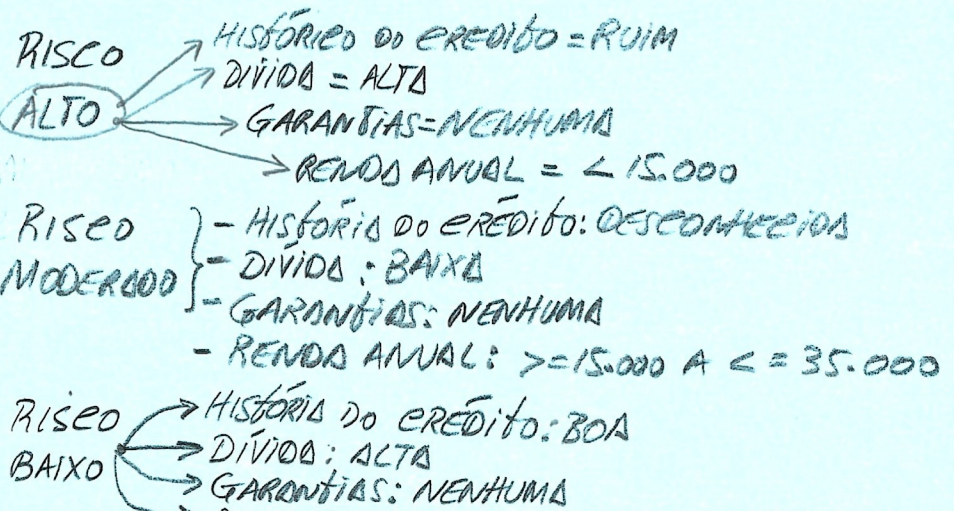


O MÉTODO PREDITIVO CLASSIFICAÇÃO # PARA ESSE TÓPICO DE APRENDIZAGEM DE MÁQUINA, UTILIZAREMOS UMA BASE DE DADOS BANCÁRIO COM OS ATRIBUTOS ABAIXO.

BASE DE DADOS: RISCO DE CRÉDITO

HISTÓRIA DO CRÉDITO	DÍVIDA	GARANTIAS	RENDA ANUAL
RUIM	ALTA	NENHUMA	< 15.000
DESCONHECIDA	BAIXA	ADEQUADA	>= 15.000 A <= 35.000
BOA		BENS, FAIS COMO: CARROS, CASAS, APTS, TERRENOS, ...	> 35.000

RISCO (circled) → ATRIBUTOS PREVISORES  
 RISCO → META/CLASSE  
 HISTÓRIA DO CRÉDITO, DÍVIDA, GARANTIAS, RENDA ANUAL → TREINAMENTO



- ALGORITMO DE NAVE BAYES → GERA UMA TABELA DE PROBABILIDADES
- ÁRVORE DE DECISÃO → GERA UM ARVORE
- REDES NEURAIIS → GERAM PESOS
- ...
- N

AGORA, A PARTIR DE UMA NOVA BASE DE DADOS, FAZEMOS AS CLASSIFICAÇÕES PARA OS NOVOS CLIENTES.

TESTE

HISTÓRIA DO CRÉDITO	DÍVIDA	GARANTIAS	RENDA ANUAL	RISCO
RUIM	ALTA	ADEQUADA	< 15.000	?
DESCONHECIDA	ALTA	ADEQUADA	< 15.000	?
DESCONHECIDA	BAIXA	NENHUMA	> 35.000	?
BOA	ALTA	ADEQUADA	>= 15.000 A <= 35.000	?

# CLASSIFICAÇÃO (VENDA DE LIVROS)

$x_1$	$x_2$	$x_3$	$Y$
SEXO	PAIS	IDADE	COMPRAR
M	FRANÇA	25	SIM
M	INGLATERRA	21	SIM
F	FRANÇA	23	SIM
F	INGLATERRA	34	SIM
F	FRANÇA	30	NÃO
M	ALEMANHA	21	NÃO
M	ALEMANHA	20	NÃO
F	ALEMANHA	18	NÃO
F	FRANÇA	34	NÃO
...	...	...	...



← ATRIBUTOS PREVISORES

$x_1$	$x_2$	$x_3$	META
SEXO	PAIS	IDADE	COMPRAR
M	FRANÇA	38	?
F	INGLATERRA	25	?
M	ALEMANHA	55	?
F	FRANÇA	20	?
...	...	...	...

TREINAMENTO

TESTE

# CLASSIFICAÇÃO (PREVER O ESPORTE)

COR DOS OLHOS	CASADO	SEXO	CABELO	ESPORTE
CASTANHO	SIM	M	LONGO	FUTEBOL
AZUL	SIM	M	CURTO	FUTEBOL
CASTANHO	SIM	M	LONGO	FUTEBOL
CASTANHO	NÃO	F	LONGO	AERÓBICA
CASTANHO	NÃO	F	LONGO	AERÓBICA
AZUL	NÃO	M	CURTO	FUTEBOL
...	...	...	...	...

TREINAMENTO

COR DOS OLHOS	CASADO	SEXO	CABELO	ESPORTE
CASTANHO	SIM	M	CURTO	?
CASTANHO	NÃO	M	LONGO	?
AZUL	NÃO	F	LONGO	?
AZUL	SIM	M	LONGO	?
...	...	...	...	...

TESTE

NOTA DO INSTRUCTOR } ESSA TÉCNICA É MUITA USADA EM TESTE VOCACIONAL. POR EXEMPLO, UMA UNIVERSIDADE QUANDO VAI FAZER DIVULGAÇÃO DE VESTIBULAR APLICA O SUPRA-CITADO TESTE E APARTIR DAÍ, CRIA UMA BASE DE DADOS COM TODAS AS PESSOAS QUE VÃO FAZER O TESTE E COM BASE NAS CARACTERÍSTICAS DAS PESSOAS PREVER QUAL CURSO A PESSOA VAI FAZER.





# CLASSIFICAÇÃO (JOGAR TÊNIS)

BASE DE DADOS HISTÓRICAS  
[CARACTERÍSTICAS CLIMÁTICAS]

- ATRIBUTOS PREVISORES →
- TEMPO
  - TEMPERATURAS
  - UMIDADE
  - VENTO

META/CLASSE → JOGAR TÊNIS

TEMP	TEMPERATURA	UMIDADE	VENTO	JOGAR
ENSOLARADO	QUENTE	ALTA	FRACO	NÃO
ENSOLARADO	QUENTE	ALTA	FORTE	NÃO
NUBLADO	QUENTE	ALTA	FRACO	SIM
CHUVOSO	MODERADA	ALTA	FRACO	SIM
CHUVOSO	AGRADÁVEL	NORMAL	FORTE	NÃO
NUBLADO	AGRADÁVEL	NORMAL	FORTE	SIM
ENSOLARADO	MODERADA	ALTA	FRACO	NÃO
CHUVOSO	MODERADA	NORMAL	FRACO	SIM
⋮	⋮	⋮	⋮	⋮

TREINAMENTO

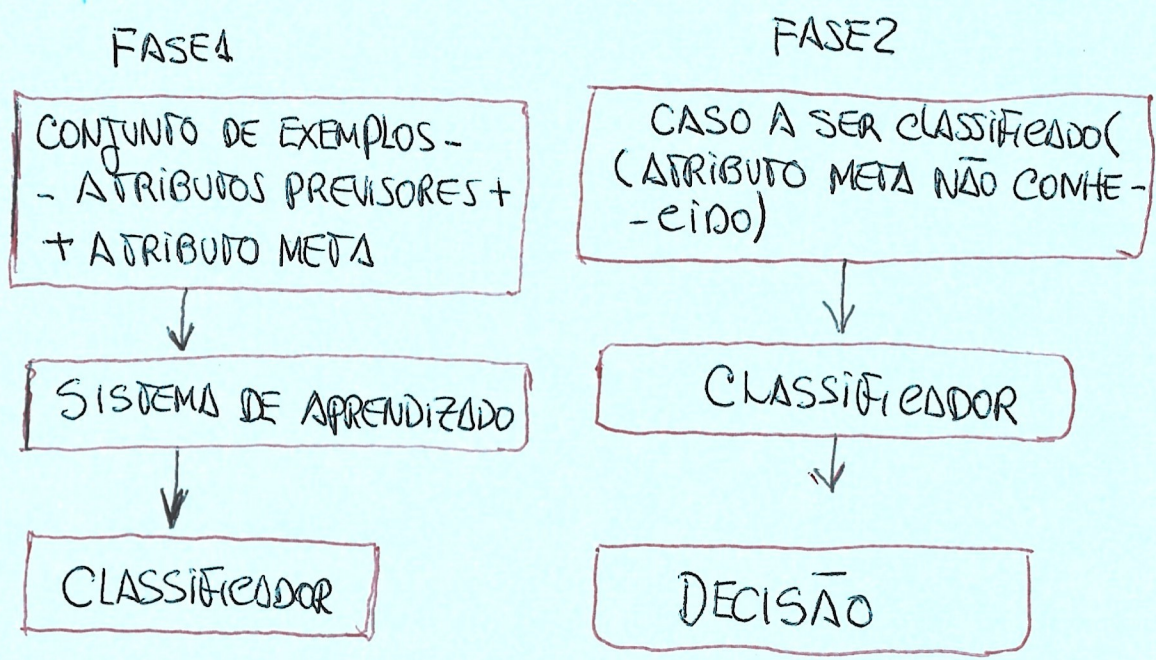
Tempo	TEMPERATURAS	UMIDADE	VENTO	JOGAR TÊNIS
ENSOLARADO	MODERADA	NORMAL	FORTE	?
CHUVOSO	AGRADÁVEL	NORMAL	FRACO	?
NUBLADO	QUENTE	NORMAL	FRACO	?
NUBLADO	AGRADÁVEL	ALTA	FORTE	?

TESTE

## ATENÇÃO → CLASSIFICAÇÃO

- CADA REGISTRO PERTENCE A UMA CLASSE E POSSUI UM CONJUNTO DE ATRIBUTOS PREVISORES;
- OBJETIVA-SE DESECOBRIR UM RELACIONAMENTO ENTRE OS ATRIBUTOS PREVISORES E O ATRIBUTO META;
  - .. OS ALGORITMOS DE MACHINE LEARNING VÃO TENTAR ENCONTRAR ESSES RELACIONAMENTOS
- O VALOR DO ATRIBUTO META É CONHECIDO (A PRENDIZAGEM SUPERVISIONADA)

## Representação da classificação (Método Indutivo)



NOTA: IR PARA A PÁGINA 12 E VER ILUSTRAÇÃO DA REPRESENTAÇÃO DA CLASSIFICAÇÃO APLICADA NA RESOLUÇÃO DE UMA TAREFA DO MUNDO REAL.

# INSTALAÇÃO DO R E DO RSTUDIO

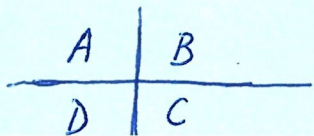
• NA BARRA DO BROWSER, DIGITAR DOWNLOAD R E CLICAR EM THE R PROJECT FOR STATISTICAL COMPUTING (O PROJETO R PARA COMPUTAÇÃO ESTADÍSTICA) → [HTTPS://WWW.R-PROJECT.ORG](https://www.r-project.org). DEPOIS, CLICAR EM ~~DOWNLOAD~~ (CRAN), DEPOIS ESCOLHER O ESPELHO [HTTPS://CLOUD.R-PROJECT.ORG/](https://cloud.r-project.org/) E DEPOIS, ESCOLHER O SISTEMA OPERACIONAL USADO. NO MEU CASO, SO WINDOWS, DEPOIS CLICAR EM BASE (EM SUBDIRETÓRIOS).

(RX64) → CARREGAR ESSE ÍCONE, VIA MENU DE PROGRAMAS (JANELA DO WINDOWS)

> 1+1  
[1] 2

Nota do Instrutor } NÃO VAMOS TRABALHAR COM O R, E SIM COM UMA IDE RSTUDIO → FERRAMENTA + POPULAR PARA TRABALHAR COM O R

DOWNLOAD RSTUDIO # BAIXAR O ARQUIVO DE INSTALAÇÃO E EXECUTÁ-LO ATÉ À SUA CONCLUSÃO.



ONDE:  
A, B, C, D SÃO AS INTERFACES DA IDE RSTUDIO

# PRÉ-PROCESSAMENTO NAS BASES DE DADOS

↳ SÃO PROCESSOS QUE DEVEM SER EXECUTADOS ANTES DE FAZERMOS APLICAÇÃO DOS ALGORITMOS DE APRENDIZAGEM DE MÁQUINA.

## CONTEÚDO

- TIPOS DE VARIÁVEIS
- BASE DE DADOS
  - DADOS DE CRÉDITO
  - CENSO
- CARREGAMENTO DAS BASES DE DADOS
- VALORES INCONSISTENTES (POR EXEMPLO, IDADE < 0 (NEGATIVA) → NÃO PODE ACONTECER ISSO. PORTANTO, PRECISAMOS CORRIGIR ESSE PROBLEMA.
- VALORES FALTANTES (DADOS AUSENTES, NÃO DISPONÍVEIS) E ISSO PREJUDICA O APRENDIZADO DOS ALGORITMOS
- ESCALONAMENTO DE ATRIBUTOS → PODE TER ATRIBUTOS NUMÉRICOS FORA DE ESCALA EM RELAÇÃO AOS DEMAIS, PARA CIMA OU PARA BAIXO
- TRANSFORMAÇÃO DE VARIÁVEIS CATEGÓRICAS → VARIÁVEIS DE TEXTO E NESSE CASO VAI PRECISAR TRANSFORMAR PARA NÚMEROS, POIS ALGUNS OU A MAIORIA DOS ALGORITMOS DE APRENDIZAGEM DE MÁQUINA NÃO LEM TEXTOS
- INTRODUÇÃO À AVALIAÇÃO DE ALGORITMOS
- BASE DE TREINAMENTO E BASE DE TESTE

# TIPOS DE VARIÁVEIS

# existem algoritmos de maquinas que não reconhecem variáveis numéricas, como é o caso de regras de associação.

## NUMÉRICAS

Contínua  
Números reais  
temperatura, altura, peso, salário

Discreta  
conjunto de valores finitos (inteiros)  
CONTAGEM DE ALGUMA COISA  
1 BAIXO  
2 MEDIO BAIXO  
3 MÉDIO  
4 MEDIO ALTO  
5. ALTO

## CATEGÓRICAS {string}

NOMINAL  
DADOS NÃO MENSURÁVEIS  
SEM ORDENAÇÃO:  
COR DOS OLHOS, GÊNERO, ID, NOME

ORDINAL  
CATEGORIZADO SOB UMA ORDENAÇÃO.  
TAMANHO, P, M E G

### BASE DE DADOS DE CRÉDITO

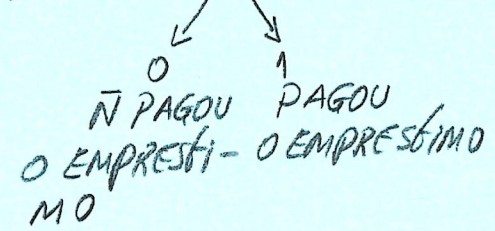
↓  
A BASE DE DADOS A SER USADA COMO ESTUDO DE CASO PARA O DESENVOLVIMENTO DO CURSO.

# VARIÁVEIS

CLIENTID: ID DO CLIENTE  
INCOME: RENDA  
LOAN: EMPRÉSTIMO (DÍVIDA)  
DEFAULT: CLASSE

} ATRIBUTOS PREVISORES

} ATRIBUTO CLASSE



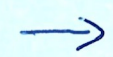
0: NÃO PAGOU A DÍVIDA  
1: PAGOU A DÍVIDA

## NO IDE DO R STUDIO

- FILE
  - NEW FILE
  - R SCRIPT
  - SALVAR O ARQUIVO COM O NOME: PRE\_PROCESSAMENTO\_CREDIT\_DATA
- ESTE ARQUIVO PASSA TER EXTENSÃO .R

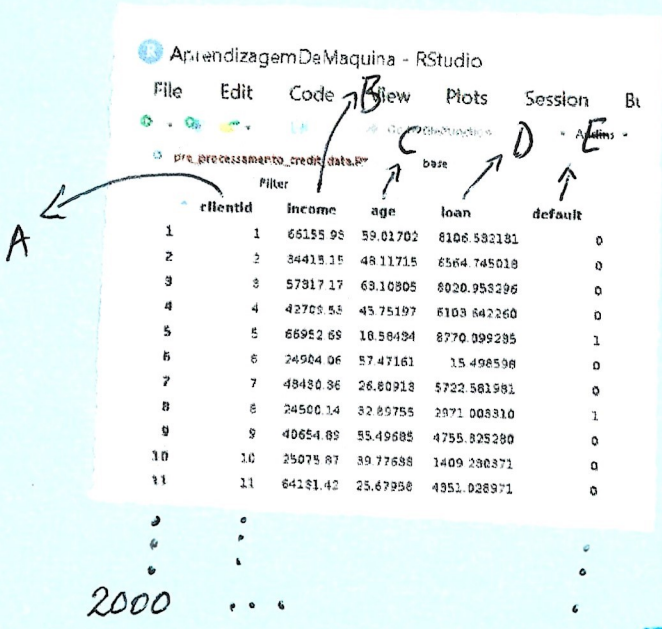
NO MENU "MORE", ESCOLHER A OPÇÃO "SET AS WORKING DIRECTORY"  
> SETWD("D:/MACHINE LEARNING E DATA SCIENCE COM R")

# CARREGANDO A BASE DE DADOS .CSV QUE ESTÁ NA PASTA SUPRACITADA PARA DENTRO DO R  
# NA ÁREA DE SCRIPT DO ARQUIVO "PRE\_PROCESSAMENTO\_CREDIT\_DATA.R", DIGITAR OS COMANDOS QUE SEGUEM:



BASE ← READ.esv('CREDIT.DATS.esv') # SELECIONAR TUDO E EXECUTAR

2000 OBS. OF 5 VARIABLES



FAREMOS AGORA A VINCULAÇÃO DOS ATRIBUTOS DA BASE COM QUELES TIPOS DE VARIÁVEIS DESCRITOS NA PÁGINA 18

A: ID DO CLIENTE

↓  
CHAVE PRIMÁRIA

~~VARIÁVEL NOMINAL~~  
VARIÁVEL CATEGÓRICA

NOTA DO INSTRUCTOR } PARA OS ALGORITMOS DE MACHINE LEARNING NOMINAL ESSE TIPO DE CAMPO NÃO É IMPORTANTE PARA FAZER UM TIPO DE PREVISÃO. PORTANTO, VAMOS APAGAR ESSE CAMPO DA BASE DE DADOS BASE.

BASE \$ CLIENTID ← NULL

↓  
BASE 2000 OBS. OF 4 VARIABLES

B: INCOME (RENDIA) → VARIÁVEL NUMÉRICAS CONTÍNUAS

C: AGE (IDADE) → VARIÁVEL NUMÉRICAS CONTÍNUAS

D: LOAN (DÍVIDA/EMPRESSTIMO) → VARIÁVEL NUMÉRICAS CONTÍNUAS

E: DEFAULT (CLASSE/META) → VARIÁVEL NUMÉRICAS DISCRETAS.

SUMMARY (BASE)

	INCOME	AGE	LOAN	DEFAULT
MIN.	: 2004	: -52.42	: 1.378	: 0.0000
1st QU.	: 32796	: 28.99	: 1939.709	: 0.0000
MEDIAN	: 45789	: 41.32	: 3974.719	: 0.0000
MEAN	: 45332	: 40.81	: 4444.370	: 0.1415
3rd QU.	: 57791	: 52.59	: 6432.411	: 0.0000
MAX	: 69996	: 63.97	: 13766.051	: 1.0000

NA'S: 3

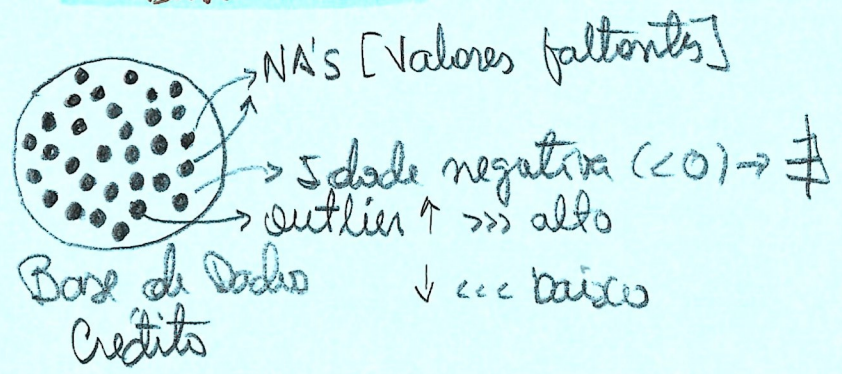
Tabela de resumo de estatística da BD base

1st QU: 1ª fatia da base = 25% da base  
 Median (mediana) = 50% da base  
 3rd QU: 75% da base

PRÓXIMA AULA SERÁ:

↳ Tratamento de valores inconsistentes - base crédito.

# Tratamento de Valores Inconsistentes - Base Créditos



A partir do Comando `summary(base)` podemos perceber que tem-se no Campo "age" um valor negativo para a idade, que é inconsistente

`age = -52.42` → uma inconsistência e isso pode influenciar negativamente no desempenho e aprendizagem do algoritmo.

Nos registros 16, 22, 27 apresentam `age < 0`. Desta forma, precisamos achar uma forma de corrigir essa inconsistência.

Nota do Instrutor } Para ver se tem + registros com "age" < 0, siga os comandos.

No script criado, digite:  
`base[base$age < 0, ]`

	Income	age	loan	default
16	50501.73	-28.21	3977.	0
22	32197.62	-52.42	4244.	0
27	63287.04	-36.49	9595.	0
NA	NA	NA	NA	NA
NA.1	NA	NA	NA	NA
NA.2	NA	NA	NA	NA

`base[base$age < 0, 1:2] #NA (Not Available)`  
NA → 29, 31 e 32

`base[base$age < 0 & !is.NA(base$age), ]`

	Income	age	loan	default
16	50501.73	-28.21	3977.	0
22	32197.62	-52.42	4244.	0
27	63287.04	-36.49	9595	0

`IDADE_INVALIDOS <- base[base$AGE < 0 & !is.NA(BASE$AGE), ]`

`IDADE_INVALIDOS`

NOTA DO INSTRUCTOR } A BASE DE DADOS NÃO PODE FICAR COM ESSES TRÊS VALORES COM ~~NA~~ IDADES INVÁLIDAS

## MANEIRAS DE TRATAR DADOS INVÁLIDOS (AGE < 0)

1ª OPÇÃO: APAGAR A COLUNA INTEIRA # PEGAR A COLUNA "AGE" E APAGAR A COLUNA INTEIRA, OU SEJA EXCLUIR A VARIÁVEL "AGE". ENTRETANTO, ESSA NÃO É UMA OPÇÃO OU SOLUÇÃO RECOMENDÁVEL, POIS NÃO É POR CAUSA DE 3 REGISTROS, VAMOS PERDER 1997 REGISTROS BONS, SE FOSSE O INVERSO, AÍ SIM, SERIA VIÁVEL.

2ª OPÇÃO: APAGAR SOMENTE OS REGISTROS COM PROBLEMAS

```
BASE <- BASE[BASE$AGE > 0, ] # A BASE VAI FICAR COM 1997 REGISTROS
```

# ESSA OPÇÃO É MENOS PIOR QUE A 1ª OPÇÃO, ENTANTO, VAMOS PERDER DADOS.

3ª OPÇÃO: PREENCHER OS DADOS MANUALMENTE

↳ PEGAR OS 3 REGISTROS, ENTRAR EM CONTATO COM ESSAS TRÊS PESSOAS E PERGUNTAR AS IDADES DOS MESMOS # É UMA OPÇÃO + CORRETA, ENTRETANTO É + BRABAL E MANUAL.

4ª OPÇÃO: CALCULAR A MÉDIA DA IDADE E PREENCHER OS VALORES INVÁLIDOS COM O VALOR DA MÉDIA.

```
MEAN(BASE$AGE)
```

[1] NA # O R retorna NA porque existem NA's na base. Portanto, devemos informar ao R que não queremos NA's.

```
mean(BASE$age, na.rm=TRUE)
```

[1] 40.80756 # Entretanto esse valor não é verdadeiro, pois ele considera o cálculo da média usando os valores negativos em sua age < 0. Temos que dizer ao R para não incluir os valores de age < 0.

```
MEAN(BASE$AGE[BASE$AGE > 0], NA.rm=TRUE)
```

[1] 40.92777 → ESSA É A MÉDIA REAL QUE VAMOS CONSIDERAR NOS CÁLCULOS PARA RESOLVER O PROBLEMA DE DADOS INVÁLIDOS.

```
BASE$AGE <- IFELSE(BASE$AGE < 0, 40.92, BASE$AGE)
```

	INCOME	AGE	LOAN	DEFAULT
16	50501.73	40.9200	3977....	0
22	32197.62	40.9200	4244....	0
27	63287.04	40.9200	9595....	0

```
MEAN(BASE$AGE, NA.rm=TRUE)
```

[1] 40.92769 # MÉDIA ATUALIZADA

NOTA DO INSTRUCTOR } COLOCANDO O CURSOR SOBRE UM COMANDO E APERTEAR "F1", MOSTRA O HELP SOBRE ESSE COMANDO.

# TRATAMENTO DE VALORES FALTANTES - BASE CRÉDITO

↳ NO ATRIBUTO PREVISOR "AGE", TEMOS TRÊS VALORES NA'S, OS REGISTROS 29, 31 E 32.

BASE[IS.NA(BASE\$AGE),] # VAI MOSTRAR OS REGISTROS COM NA'S, CONFORME TRAZIDO ABAIXO.

	INCOME	AGE	LOAN	DEFAULT
29	59417.81	NA	2082.626	0
31	48528.85	NA	6155.785	0
32	23526.30	NA	2862.040	0

```
BASE$AGE <- ifelse(is.na(BASE$AGE), mean(BASE$AGE, na.rm = TRUE), BASE$AGE)
```

# ESCALONAMENTO DE ATRIBUTOS - BASE CRÉDITO

CAMPO 15 } REGISTRO INCOME → 63061 - 27268 → 35793  
 14 }  
 CAMPO 61 } REGISTRO AGE → 61 - 39 = 22  
 39 }

PRECISAMOS COLOCAR OS CAMPOS "INCOME" E "AGE" NUMA MESMA ESCALA DIMENSIONAL

↳ ALGORITMOS BASEADOS EM DISTÂNCIA DEVE SER LEVADO EM CONSIDERAÇÃO ESSA QUESTÃO DE DISTÂNCIA (ESCALONAMENTO)

↳ DEIXAR OS ATRIBUTOS NUM MESMO PADRÃO.

## TÉCNICAS DE ESCALONAMENTO

### • NORMALIZAÇÃO (NORMALIZATION)

$$X = \frac{X - \text{MÍNIMO}(X)}{\text{MÁXIMO}(X) - \text{MÍNIMO}(X)}$$

PRINCIPALMENTE PARA TRATAR OUTLIER

### • PADRONIZAÇÃO (STANDARDISATION)

$$X = \frac{X - \text{MÉDIA}(X)}{\text{DESVIO PADRÃO}(X)}$$

← ESSA TÉCNICA É A MAIS RECOMENDADA PARA ESCALONAMENTO.

ONDE:

X É O VALOR DO RESPECTIVO ATRIBUTO QUE QUEREMOS NORMALIZAR. SUPONHAMOS QUE QUEIRAMOS NORMALIZAR UMA DETERMINADA IDADE. POR EXEMPLO: VALOR QUE QUEREMOS NORMALIZAR

☺ → AGE = 40 }  $\frac{40 - 20}{60 - 20} = 0,5$  VALOR MÍNIMO (AGE)

PACIENTE

## # FAZENDO O ESCALONAMENTO SOMENTE DOS ATRIBUTOS PREVISORES

`BASE[, 1:3] ← SCALE(BASE[, 1:3])`

## BASE DE DADOS DO CENSO

↳ ∃ na WEB A BASE UCI MACHINE LEARNING REPOSITÓRIO (REPOSITORY)

ARCHIVE.ICS.UCI.EDU/ML/INDEX.PHP → REPOSITÓRIO DE MACHINE LEARNING PARA DIVERSAS TAREFAS DE MINERAÇÃO DE DADOS.

A BASE QUE USA-REMOS É A "ADULT"

←

ADULT DATA SET → Em DOWNLOAD, clicar em Data Folder, depois em adult.data

CRIANDO NOVO SCRIPT PARA ESSA NOVA BD

↳ SALVANDO O SCRIPT COM O NOME DE "PRE PROCESSAMENTO\_CENSO"

1. IMPORTANDO A BASE DE DADOS PARA O R;

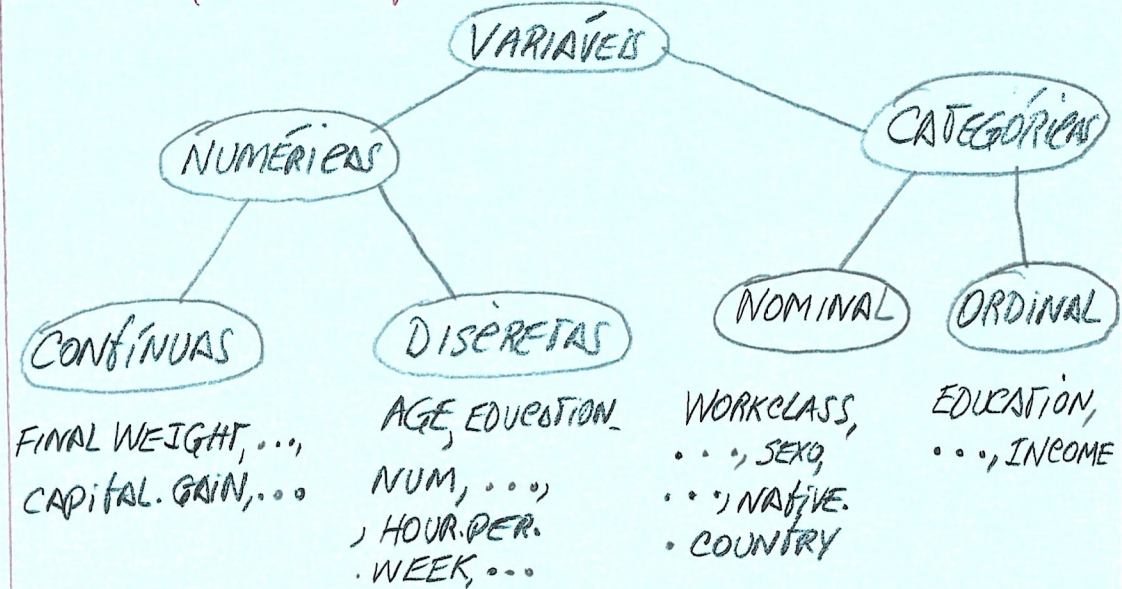
`BASE_CENSO ← READ.ESV('CENSUS.ESV')`

↳ 30162 OBS. OF 16 VARIABLES

2. VAMOS APAGAR A COLUNA "X" (ID), POIS ESTE NÃO SERVIRÁ PARA AS PREVISÕES.

`BASE_CENSO$X ← NULL`

3. VINCULAÇÃO DE TIPOS DE VARIÁVEIS # PASSO 3



NOTA DO CRIADOR } NESTA BASE, NÃO TEMOS VALORES NA'S E VALORES INCONSISTENTES. OU SEJA, A BASE ESTÁ MAIS ORGANIZADA QUE A BASE "CREDIT.DATA".

## TRANSFORMAÇÃO DE ATRIBUTOS CATEGÓRICOS - BASE CENSO

↳ PRECISAMOS TRANSFORMAR ATRIBUTOS CATEGÓRICOS PARA NUMÉRICOS, POIS ALGUNS ALGORITMOS DE MACHINE LEARNING NÃO RECONHECEM ESSES TIPOS DE VARIÁVEIS NO MODELO.



# ESCALONAMENTO DE ATRIBUTOS - CENSO

PORTANTO, NA BASE DE DADOS CENSUS TEMOS VÁRIOS ATRIBUTOS NOMINAIS QUE PRECISAM SER TRANSFORMADOS PARA ATRIBUTOS NUMÉRICOS DO TIPO DISCRETO. POR EXEMPLO, ATRIBUTOS GÊNERO, M X F, PRECISAMOS COLOCAR UM VALOR NUMÉRICO PARA A MÁQUINA ENTENDER. M → 0 E F → 1 OU VICE-VERSA. PARA FAZ, VAMOS FAZER UM COMANDO DO R, CONFORME TRAZIDO ABAIXO:

```
TABLE(BASE_CENSO$SEX) <-
```

FEMALE	MALE
9782	20380

```
UNIQUE(BASE_CENSO$SEX) <-
```

"MALE" "FEMALE"

## VAMOS FAZER A MODIFICAÇÃO DO TIPO DE ATRIBUTO

```
base_censo$sex <- factor(nare_censo$sex, levels =  
= c('Female', 'Male'), labels = c(0, 1))
```

↑ ESPAÇO } OLHAR COMO ESTÁ A SAÍDA NO COMANDO UNIQUE(BASE\_CENSO\$SEX)

```
BASE_CENSO[15, NA(BASE_CENSO$SEX)]  
DATA FRAME WITH 0 COLUMNS AND 30462 ROWS
```

FINAL.WEIGHT      EDUCATION

77516      10  
83811      40

# MUITO GRANDE DE ESCALAS. ISSO É PROBLEMA, POIS

HA ALGORITMOS DE MACHINE LEARNING (KNN) QUE LEVA EM CONSIDERAÇÃO AS DISTÂNCIAS.

TEMOS QUE COLOCAR ESSES ATRIBUTOS NA MESMA ESCALA.

```
BASE_CENSO <- SCALE(BASE_CENSO)
```

```
> BASE_CENSO <- SCALE(BASE_CENSO)
```

ERROR IN COLMEANS(X, NA.rm = TRUE): 'X' MUST BE NUMERIC

```
BASE_CENSO[, 1] <- SCALE(BASE_CENSO[, 1])
```

- WORKCLASS → 2
- EDUCATION → 4
- MARITAL STATUS → 6
- OCCUPATION → 7
- RELATIONSHIP → 8
- RACE → 9
- SEX → 10
- NATIVE.COUNTRY → 14
- INCOME → 15

TODOS ESSES SÃO DO TIPO ORIGINAL FACTOR, E PRECISAM, DESTA FORMA, SEREM CONVERTIDOS PARA NUMERIC ANTES DO ESCALONAMENTO.



- ③ FINAL.WEIGHT
- ⑤ EDUCATION.NUM
- ⑪ CAPITAL.GAIN
- ⑫ CAPITAL.LOOS
- ⑬ HOUR.PER.WEEK

ATRIBUTOS NUMÉRICOS QUE NÃO PRECISAM SER TRANSFORMADOS E PODEM SER ESCALONADOS LOGO DIRETAMENTE, CONFORME ABAIXO

- BASE\_CENSO[, 3] ← SCALE(BASE\_CENSO[, 3]) ↓
- BASE\_CENSO[, 5] ← SCALE(BASE\_CENSO[, 5]) ↓
- BASE\_CENSO[, 11:13] ← SCALE(BASE\_CENSO[, 11:13]) ↓



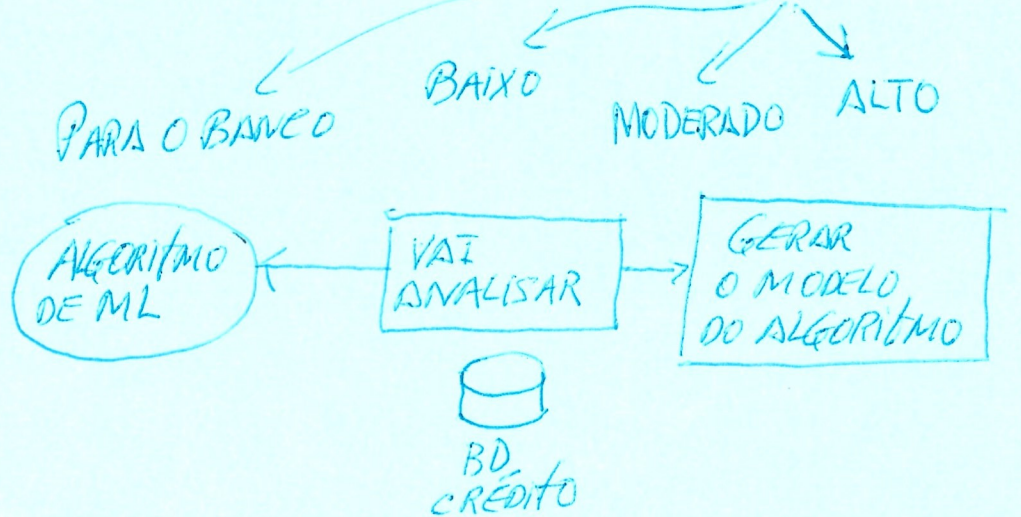
COM ISSO, CONCLUIMOS O ESCALONAMENTO DOS ATRIBUTOS NUMÉRICOS DA BASE DE DADOS CENSUS

### INTRODUÇÃO A AVALIAÇÃO DE ALGORITMOS

→ VAMOS APRENDER SE O DESEMPENHO DO ALGORITMO É BOM OU RUIM. VAMOS DIVIDIR AS BASES DE DADOS EM UMA PORÇÃO SÓTENTE PARA TREINAMENTO E UMA PORÇÃO SÓMENTE PARA TESTE

# Para esse tópico, vamos usar a Base de crédito, já trazida e disectada na página 14, cujas

Os atributos preditores são: História do crédito, dívida, garantias e renda anual, de forma a prever a meta (classe = Risco).

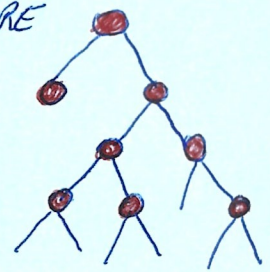


POR EXEMPLO, O ALGORITMO NAIVE BAYES VAI GERAR UMA TABELA DE PROBABILIDADES, CONFORME ILUSTRAÇÃO ABAIXO:

%	%	%	} PROBABILIDADES DAS CLASSE P(A), P(M) e P(B)
%	%	%	
%	%	%	

↑ RISCO ALTO      ↓ RISCO MODERADO      ↓ RISCO BAIXO

JÁ O ALGORITMO DE DECISÃO, VAI CONSTRUIR UMA ARVORE

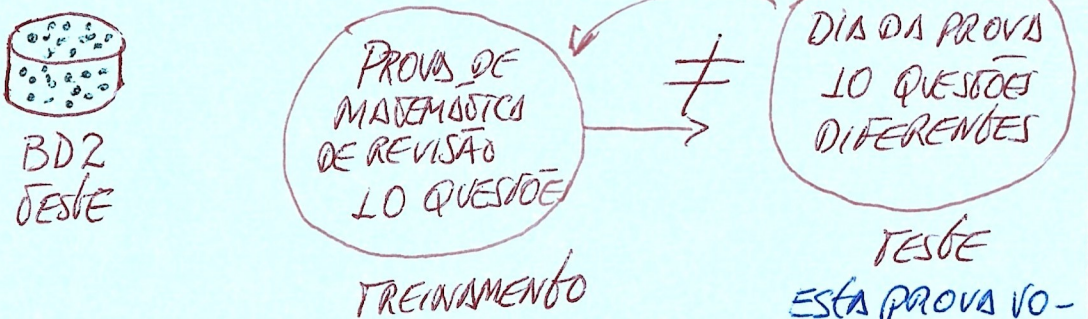
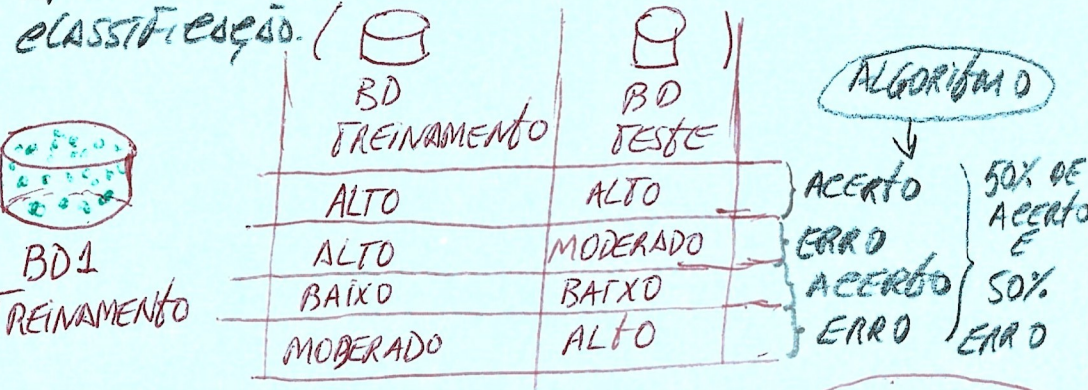


DEPOIS QUE EU QUISER CLASSIFICAR UM NOVO REGISTRO (CLIENTE), EU NÃO SEI A QUE CLASSE ESSES NOVOS REGISTROS PERTENCEM

TABELA:

HISTÓRIA DO CRÉDITO	DÍVIDA	GARANTIAS	RENDA ANUAL	RISCO
X1	X2	X3	X4	
RUIM	ALTA	ADEQUADA	<15.000	?
DESPONHÍVEL	ALTA	ADEQUADA	>35.000	?
BOA	BAIXA	NENHUMA	>=15.000 &lt;=35.000	?

SUBMETEREMOS ESSES NOVOS REGISTROS DO MODELO APRENDIDO DO ALGORITMO E AI ESTE FARÁ A NOVA CLASSIFICAÇÃO.



$$\text{EFICIÊNCIA DO ALGORITMO} = \frac{\text{Nº DE ACERTOS}}{\text{Nº TOTAL DE REGISTROS}}$$

BASE DE TREINAMENTO E TESTE

NESTE TÓPICO FAREMOS A DIVISÃO DAS BASES DE DADOS DE TREINAMENTO E TESTE

PRÉ-PROCESSAMENTO\_CREDIT.DAT.R

```
IMPLEMENTAÇÃO1 ← READ.ESV('CREDIT.DAT.ESV')
IMPLEMENTAÇÃO1 $CLIENTID = NULL
SUMMARY(IMPLEMENTAÇÃO1)
IDADE_INVALIDA ← IMPLEMENTAÇÃO1[IMPLEMENTAÇÃO1 $
$AGE < 0 & !IS.NA(IMPLEMENTAÇÃO1 $
$AGE), ]
IMPLEMENTAÇÃO1 $AGE ← NULL
IMPLEMENTAÇÃO1 ← IMPLEMENTAÇÃO1 $AGE > 0, ]
MEAN(IMPLEMENTAÇÃO1 $AGE NA.RM = TRUE)
MEAN(IMPLEMENTAÇÃO1 $AGE [IMPLEMENTAÇÃO1 $AGE > 0,
NA.RM = TRUE])
```

```
IMPLEMENTAÇÃO1 $AGE ← IFELSE(IMPLEMENTAÇÃO1 $AGE < 0,
40.92, IMPLEMENTAÇÃO1 $AGE)
IMPLEMENTAÇÃO1 [IS.NA(IMPLEMENTAÇÃO1 $AGE), ]
IMPLEMENTAÇÃO1 $AGE ← IFELSE(IS.NA(IMPLEMENTAÇÃO1 $
$AGE), MEAN(IMPLEMENTAÇÃO1 $AGE,
NA.RM = TRUE), IMPLEMENTAÇÃO1 $AGE)
IMPLEMENTAÇÃO1 [, 1:3] = SCALE(BASE[, 1:3])
```

ONDE: EXECUTAR SOMENTE OS COMANDOS →

VAMOS AGORA FAZER UMA INSTALAÇÃO ADICIONAL DE UM PACOTE DO R PARA QUE CONSIGAMOS FAZER A DIVISÃO DA BASE DE DADOS EM BASE DE TREINAMENTO E BASE DE TESTE.

```
INSTALL_PACKAGES("caTools") # INSTALANDO O PACOTE
library(caTools) # USANDO O PACOTE
```

```
set.seed(1)
```

```
divisao <- sample.split(BASE$default, splitRatio = 0.75)
```

NOTA DO INSTRUCTOR: NESSA BASE "DIVISÃO" TEM VÁRIOS "TRUE" E VÁRIOS "FALSE", QUE SIGNIFICAM OS REGISTROS QUE VÃO ENTRAR NA BASE (TRUE) E OS QUE NÃO VÃO ENTRAR NA BASE (FALSE).

```
BASE_TREINAMENTO <- subset(BASE, divisao == TRUE)
```

```
BASE_TESTE <- subset(BASE, divisao == FALSE)
```

OBSERVAÇÃO: QUANDO FORTOS TRABALHAR COM A BASE CENSO, CONSIDERAR 0.85 PARA A VARIÁVEL DIVISÃO E O RESTANTE PARA TESTE

MUDAR PARA INCOME

```
set.seed(1)
```

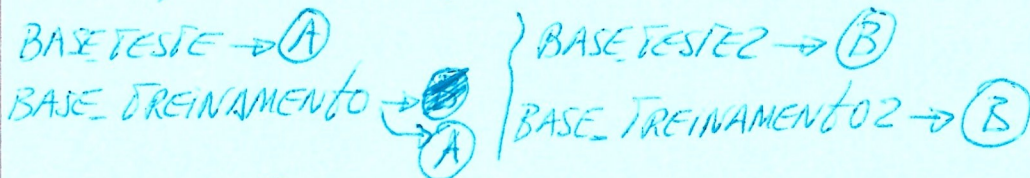
```
divisao <- sample.split(BASE$default, splitRatio = 0.85)
```

```
BASE_TREINAMENTO <- subset(BASE, divisao == TRUE)
```

```
BASE_TESTE <- subset(BASE, divisao == FALSE)
```



- BASE → REFERE-SE À BASE DE DADOS CREDIT. DATA
- BASE\_CENSO → REFERE-SE À BASE DE DADOS CENSUS-DATA



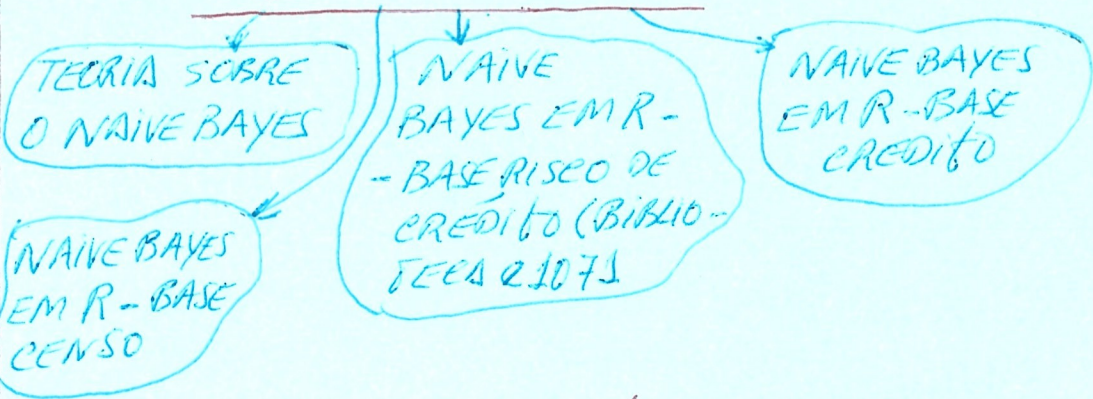
DIVISÃO → A

DIVISÃO2 → B

USADO EM TAREFA COMO CLASSIFICAÇÃO E MINERAÇÃO DE TEXTOS.

### MÓDULO APRENDIZAGEM BAYESIANA

#### CONTEÚDO

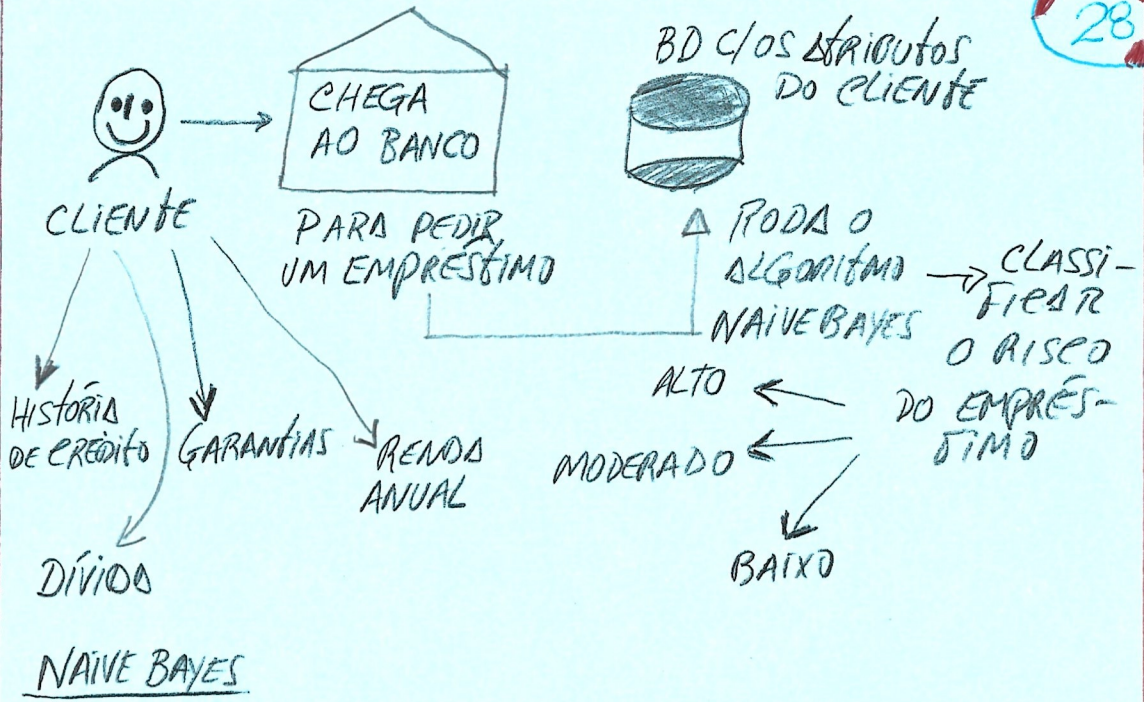
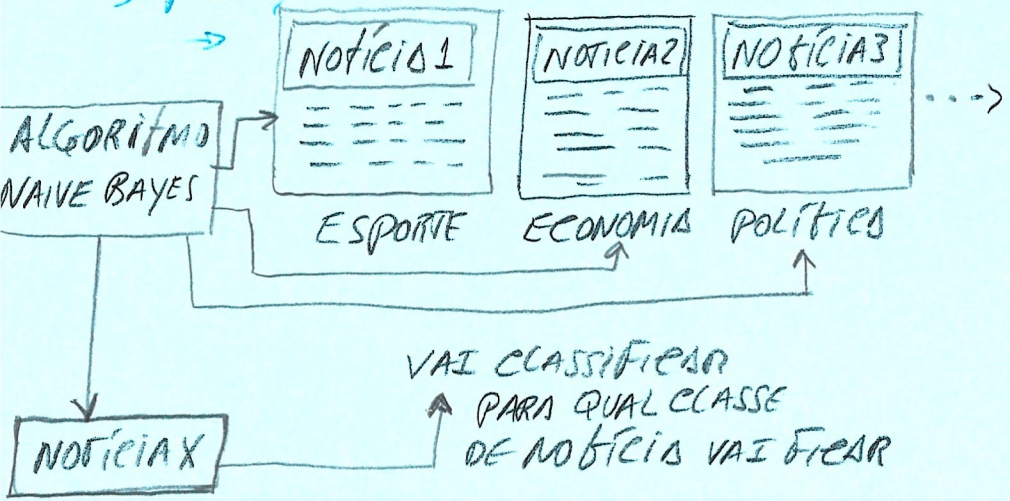


### NAIVE BAYES - INTRODUÇÃO # BASEADO EM PROBABILIDADES

- ABERDAGEM PROBABILÍSTICA (TEOREMA DE BAYES)
- EXEMPLOS
  - FILTROS DE SPAM;

• MINERAÇÃO DE EMOCÕES (😊; 😞; 😡, ...)  
 DADA UMA FRÁSE, MINERAR AS EMOCÕES QUE SE ENCONTRAM NO TEXTO ESCRITO

• SEPARAÇÃO DE DOCUMENTOS →



— BASE ORIGINAL ATRIBUÍDOS PREVISORES → CLASSES/META

HISTÓRIA DO CRÉDITO	DÍVIDA	GARANTIAS	RENDAMENTO ANUAL	RISCO
RUIM	ALTA	NENHUMA	<15.000	ALTO
DESCONHECIDA	ALTA	NENHUMA	>=15.000 &lt;=35.000	ALTO
DESCONHECIDA	BAIXA	NENHUMA	>=15.000 &lt;=35.000	MODERADO
DESCONHECIDA	BAIXA	NENHUMA	>35.000	BAIXO
DESCONHECIDA	BAIXA	ADEQUADA	>35.000	BAIXO
RUIM	BAIXA	NENHUMA	<15.000	ALTO

RISCO DE CRÉDITO	HISTÓRIA DO CRÉDITO			DÍVIDA		GARANTIAS		RENDAMENTO ANUAL		
	BOA 5	DESCONH. 5	RUIM 4	ALTA 7	BAIXA 7	NENHUMA 11	ADEQUADA 3	<15000 3	>=15000 &lt;=35000 4	>35000 7
ALTO 6/14	1/6	2/6	3/6	4/6	2/6	6/6	0	3/6	2/6	1/6
MODERADO 3/14	4/3	4/3	4/3	4/3	2/3	2/3	4/3	0	2/3	4/3
BAIXO 5/14	3/5	2/5	0	2/5	3/5	3/5	2/5	0	0	5/5

6 : SIGNIFICA 6 REGISTROS DA CLASSE ALTA DE UM TOTAL DE 14 REGISTROS E ASSIM SUCESSIVAMENTE PARA OS OUTROS VALORES DA TABELA

# NAIVE BAYES - APRENDIZAGEM

RISCO DE CRED.	HIST. DE CRED.		DÍVIDA		GARANT.		REVENHA ANUAL			HIST/CRÉDITO	RISCO	
	BOA	DESCONH.	RUIM	ALTA	BAIXA	NENH.	ADÉQ.	A	B	C		
ALTO 6/14	5	5	4	7	7	11	3	3	4	7	DESCONH.	ALTO
MODER.											DESCONH.	MOD.
BAIXO 5/14											DESCONH.	ALTO
											DESCONH.	BAIXO
											RUIM	ALTO
											RUIM	MOD.
											BOA	BAIXO
											BOA	ALTO
											BOA	MOD.
											BOA	BAIXO
											RUIM	ALTO

A: < 15.000  
 B: >= 15.000  
 C: > 35.000

TABELA DE PROBABILIDADES


↓  
 MACHINE LEARNING → DE NAIVE BAYES

ALGORITMO


# NAIVE BAYES - CLASSIFICAÇÃO

↳ ESSE ALGORITMO TRABALHA ANALISANDO UMA BASE HISTÓRICA DE CRÉDITO PARA A CONSTRUÇÃO DE UMA TABELA DE PROBABILIDADES, COM O OBJETIVO DE CLASSIFICAR NOVOS CLIENTES SOCIALIZANDO UMA LINHA DE CRÉDITO.

POR EXEMPLO, SUPONHAMOS QUE FENHAMOS UM NOVO REGISTRO (CLIENTE) COM OS ATRIBUTOS ABAIXO, E VAMOS SABER QUAL SERIA A CLASSE DELE.

CLIENTE 

- HISTÓRIA = BOA
- DÍVIDA = ALTA
- GARANTIAS = NENHUMA
- RENDA = > 35.000

RISCO DE CRED.	BOA	ALTA	NENHUMA	> 35.000	RISCO
ALTO 6/14	5	7	11	7	BAIXO 
MODERAD. 3/14	1/6	4/6	6/6	1/6	
BAIXO 5/14	1/3	1/3	2/3	1/3	

$P(ALTO) = 6/14 * 1/6 * 4/6 * 6/6 * 1/6 = 0,0079 = 0,79\%$   
 $P(MODERADO) = 3/14 * 1/3 * 1/3 * 2/3 * 1/3 = 0,0052$   
 $P(BAIXO) = 5/14 * 3/5 * 2/5 * 3/5 * 5/5 = 0,0514$   
 $\sum P(A, M, B) = 0,0079 + 0,0052 + 0,0514 = 0,0645 \underline{100\%}$

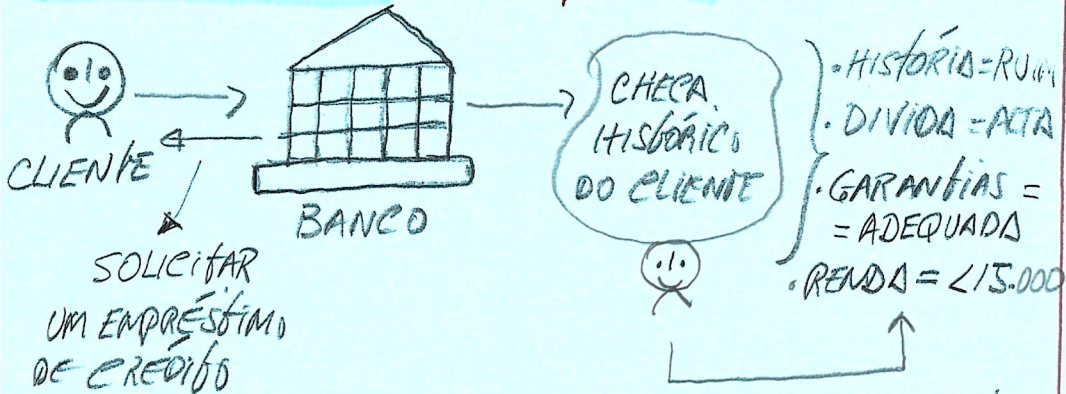
$P(ALTO) = 0,0079 / 0,0645 * 100 = 12,24\%$   
 $P(MOD.) = 0,0052 / 0,0645 * 100 = 8,07\%$   
 $P(BAIXO) = 0,0514 / 0,0645 * 100 = 79,69\%$

$P(B) = 79,69\%$

RESULTADO FINAL → CRÉDITO CONCE-  
 DIDO, POIS O RISCO  
 É BAIXO, PODENDO SER CONCEPIDO  
 100% OU +

PODE SER USADO EM CIMA DO VALOR DO EMPREST.

# NAIVE BAYES - CORREÇÃO LAPLACIANA



RODAR O ALGORITMO DE MACHINE LEARNING PARA ANALISAR O RISCO DE CONCESSÃO DE CRÉDITO.

FAZENDO O MESMO PROCEDIMENTO FEITO NA PÁGINA 29, ACHAREMOS O SEGUINTE RESULTADO FINAL. ENTRETANTO, ESSE EXEMPLO TEM UMA PARTICULARIDADE, CONFORME DESCRITO ABAIXO.

$$\begin{aligned}
 P(\text{ALTO}) &= 6/14 * 3/6 * 4/6 * 0 * 3/6 = 0 \\
 P(\text{MODERADO}) &= 3/14 * 1/3 * 1/3 * 1/3 * 0 = 0 \\
 P(\text{BAIXO}) &= 5/14 * 0 * 2/5 * 2/5 * 0 = 0
 \end{aligned}$$

OU SEJA, AS FRES P(N) = 0  
 NESSE CASO, NÃO VAMOS FAZER A PREVISÃO DE RISCO DO FINANCIAMENTO

PARA RESOLVER ESSE PROBLEMA, RECORREMOS A UM ARBITRÍCIO DO ALGORITMO CHAMADO DE CORREÇÃO LAPLACIANA.

# CORREÇÃO LAPLACIANA → ADICIONARÁ UM REGISTRO

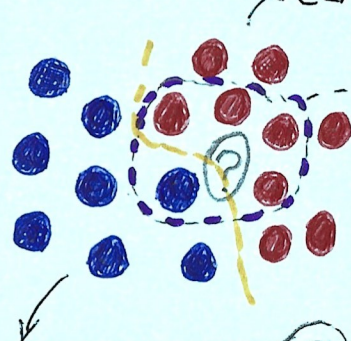
NO REGISTRO "RUIM", VALOR = 0, VAI FICAR 1/6. EM BOA E DESCONHECIDA, VAI FICAR 3/6 E 2/6, RESPECTIVAMENTE. NA COLUNA "RISCO DE CRÉDITO", VALOR "BAIXO" VAI FICAR 6/15. A QUANTIDADE DE "RUIM" = 4, VAI MUDAR PARA 5.

$$P(\text{BAIXO}) = 6/15 * 1/6 * 3/6 * 2/6 * ?$$

NOTA DO INSTRUTOR } ESSA AULA É SOMENTE ILUSTRATIVA DE COMO FUNCIONA ESSA CORREÇÃO. PARA TER OS VALORES EXATOS PARA SEREM USADOS NA FÓRMULA, O IDEAL É ADICIONAR OS NOVOS REGISTROS NA BASE DE DADOS E DEPOIS CONSTRUIR A TABELA DE PROBABILIDADES.

# NAIVE BAYES - + CONCEITOS

CLASSE 1: VERMELHA



$$P(\text{VERMELHA}) = 9/17$$

$$P(\text{AZUL}) = 8/17$$

OS ALGORITMOS NAIVE BAYES TEM UM PARÂMETRO CHAMADO DE RADIUS = RAIO QUE SÃO ALGUNS PONTOS QUE ELE VAI PEGAR QUE ESTÃO DENTRO DESSA ÁREA CHAMADA "RADIUS" PEGANDO ALGUNS VALORES, OU SEJA 4 VERMELHOS E 1 AZUL

CLASSE 2: AZUL

$$P'(\text{VERMELHO}) = 4/9$$

$$P'(\text{AZUL}) = 1/8$$

PORTANTO, A BOLINHA VAI SER VERMELHA  
 VOCÊ TEM UM PONTO QUE NÃO SABEMOS A QUE CLASSE PERTENCE

PEP' SÃO CHAMADAS DE PROBABILIDADES APRIORI

$$P''(\text{VERMELHO}) = 9/17 * 4/9 = 0,24$$

$$P''(\text{AZUL}) = 8/17 * 1/8 = 0,06$$

# ... NAIVE BAYES - + CONCEITOS ... CONT...

P" SÃO CHAMADAS DE PROBABILIDADES POSTERORI

## VANTAGENS X DESVANTAGENS

### VANTAGENS

- RÁPIDO;
- SIMPLICIDADE DE INTERPRETAÇÃO;
- TRABALHA COM ALTAS DIMENSÕES;
- BOAS PREVISÕES EM BASES PEQUENAS

### DESVANTAGENS

- COMBINAÇÃO DE CARACTERÍSTICAS (ATRIBUTOS INDEPENDENTES) - CADA PAR DE CARACTERÍSTICAS SÃO INDEPENDENTES - NEM SEMPRE É VERDADE.

BOA \* ALTA \* < 15 → INDEPENDENTES  
 RENDA | PAR  
 DIVIDA } INDEPENDENTE

## IMPLEMENTAÇÃO NAIVE BAYES EM R DE CREDITO

```
setwd("D:/MACHINE LEARNING E DATA SCIENCE COM R") # SET AS WORKING DIRECTORY
```

NO CONSOLE DO RSTUDIO, VAMOS CRIAR NOVO SCRIPT DE COMANDOS DE IMPLEMENTAÇÃO →

SALVANDO-O COM O NOME DE "NAIVE\_BAYES\_RISEO - credito.R" # A EXTENSÃO .R, O R SALVA AUTOMATICAMENTE.

```
BASE2 <- read.csv('RISEO_CREDITO.csv')
BASE2
```

lista_id	divida	garantias	renda	risco	
1	ruim	alta	nenhuma	0_15	alto
2	desconhecido	alta	nenhuma	15_35	alto
3	desconhecido	baixa	nenhuma	15_35	moderado
4	desconhecido	baixa	nenhuma	acima_35	alto
5	desconhecido	baixa	nenhuma	acima_35	baixo
6	desconhecido	baixa	adequada	acima_35	baixo
7	ruim	baixa	nenhuma	0_15	alto
8	ruim	baixa	adequada	acima_35	moderado
9	boa	baixa	nenhuma	acima_35	baixo
10	boa	alta	adequada	acima_35	baixo
11	boa	alta	nenhuma	0_15	alto
12	boa	alta	nenhuma	15_35	moderado
13	boa	alta	nenhuma	acima_35	baixo
14	ruim	alta	nenhuma	15_35	alto

```
install.packages("e1071") # INSTALANDO PACOTE.
library(e1071) # CARREGANDO O PACOTE
```

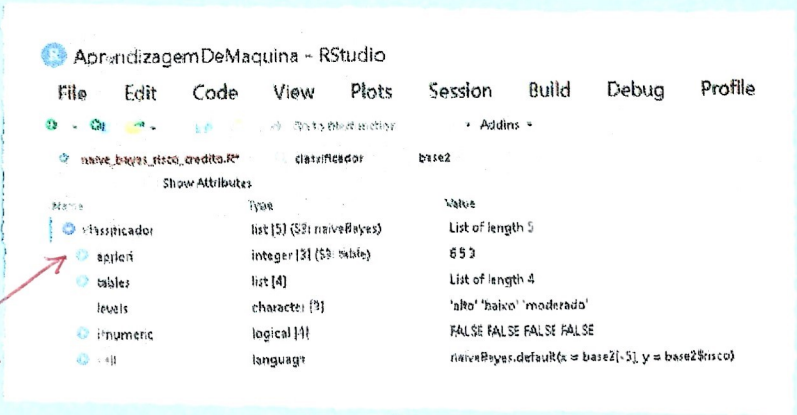
# CRIANDO O CLASSIFICADOR

```
CLASSIFICADOR <- naiveBayes(X = BASE2[-5], Y = BASE2$risco)
```

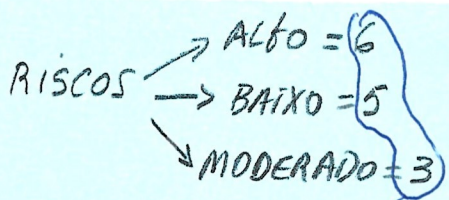


DEPOIS QUE CRIAR O CLASSIFICADOR NAIVE BAYES NO SCRIPT, RODA-SE O COMANDO ACIMA E GERA-SE A TABELA (QUADRO) A SEGUIR

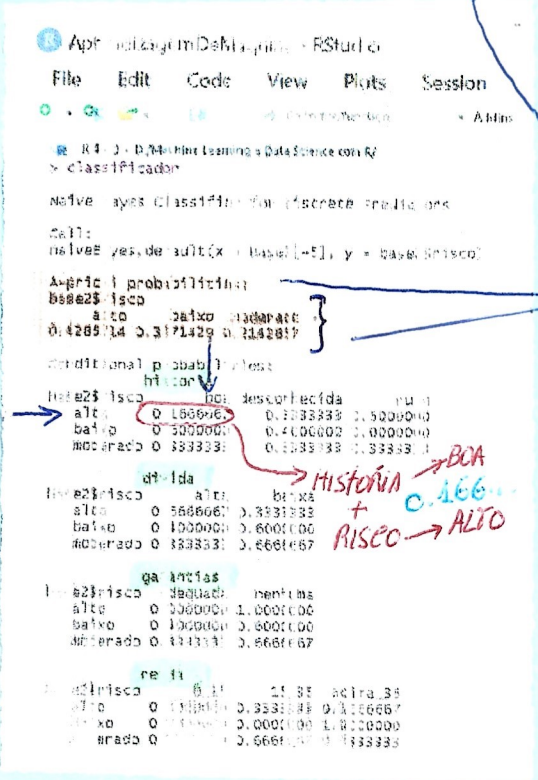




EXPANDIR  
NO-SE,  
TEREMOS:



O TREINAMENTO DO ALGORITMO É JUSTAMENTE A CRIAÇÃO DA TABELA DE P(n).



SE JOGAR ESSES VALORES NA TABELA DE PROBABILIDADES CONSTRUÍDA PELO ALGORITMO (VER PÁGINA 29) VAMOS ACHAR ESSES MESMOS VALORES.

VAMOS SUPOR QUE QUERAMOS CLASSIFICAR UM NOVO REGISTRO, CONFORME TRAZIDO ABAIXO:

# HISTÓRIA: BOA, DÍVIDA: ALTA, GARANTIAS: NENHUMA, RENDA: >35

# HISTÓRIA: RUIM, DÍVIDA: ALTA, GARANTIAS: ADEQUADA, RENDA: <15

QUAL SERÁ A CLASSE PARA ESSES DOIS NOVOS REGISTROS

- ALTO
- BAIXO
- MODERADO

RISCO?

PARA ISSO, VAMOS PRECISAR CRIAR UM NOVO DATAFRAME COM ESSES NOVOS REGISTROS, CONFORME TRAZIDO ABAIXO.

```

historia <- c("boa")
divida <- c("alta")
garantias <- c("nenhuma")
renda <- c("acima_35")
df <- data.frame(historia,
                 divida, garantias,
                 renda)
    
```

```

previsao <- predict(classificador, newdata =
                    = df)
    
```

PARA VER O RESULTADO, BASTA DIGITAR NO PROMPT DO RSTUDIO O COMANDO `print(previsao)` E SAÍ O RESULTADO ABAIXO:

[1] BAIXO. # PORTANTO, O CLIENTE 1 TEM RISCO BAIXO PARA O FINANCIAMENTO DE CRÉDITO.

... QUERENDO VER AS PROBABILIDADES, TEM-SE:  
 PREVISAO <- PREDICT(CLASSIFICADOR, NEWDATA =  
 =DF, 'RAW')  
 PRINT(PREVISAO) # SAIDA MOSTRADA ABAIXO

	ALTO	BAIXO	MODERADO
[1,]	0.1227496	0.79417349	0.08183306
[2,]	0.8993525	0.000719482	0.09992805

### NAIVE BAYES EM R - BASE CREDITO

↗ CRIA-SE UM NOVO ARQUIVO DE SCRIPT COM O NOME:  
 NAIVE\_BAYES\_CREDITO.DAT COM AS LINHAS DE  
 CÓDIGO ABAIXO:

```
# LEITURA DA BASE DE DADOS
BASE <- READ.ESV('CREDITO.DAT.ESV')

# APAGA A COLUNA CLIENTID
BASE$CLIENTID = NULL

# TRATA VALORES INCONSISTENTES
BASE$AGE <- IFELSE(BASE$AGE < 0, 40.92,
, BASE$AGE)

# FAZENDO O ESCALONAMENTO
BASE[, 1:3] <- SCALE(BASE[, 1:3])
```



# DIVISÃO ENTRE TREINAMENTO E TESTE  
 LIBRARY(CATools)

```
SET.SEED(1)
DIVISAO <- SAMPLE.SPLIT(BASE$INCOME,
, SPLITRATIO = 0.75)
```

```
BASE.TREINAMENTO <- SUBSET(BASE, DIVISAO ==
== TRUE)
```

```
BASE.TESTE <- SUBSET(BASE, DIVISAO == FALSE)
```

~~VAMOS CRIAR UM NOVO SCRIPT~~

```
CLASSIFICADOR2 <- NAIVE.BAYES(X = BASE.TREINAMENTO[,
[-4], Y = BASE.TREINAMENTO$DEFAULT)
```

```
PRINT(CLASSIFICADOR2) # MOSTRA QUE 86% É À CLASSE
0 E 14% É À CLASSE 1. SENDO
QUE "0" NÃO PAGA O FINANCIAMENTO E "1" PO-
GUE O FINANCIAMENTO BANCÁRIO
```

NOTA DO INSTRUCTOR } ESSA BASE DE TREINAMENTO TEM SOMENTE  
 DADOS (VARIÁVEIS NUMÉRICAS)

↓  
 OU SEJA, A TABELA DE PROBABILIDADES TAMBÉM É GERADA

# VAMOS AGORA FAZER AS NOVAS PREVISÕES USANDO  
 A BASE DE TESTE JÁ CRIADA (VER PÁGINA 37)



... CONTINUAÇÃO

```
PREVISOES <- predict(classificador, newdata =
  BASE_TESTE[-4])
```

```
VIEW(PREVISOES)
```

NAME	TYPE	VALUE
PREVISOES	FACTOR	FACTOR WITH 2 LEVELS: "0", "1"

DEFAULT → OS VALORES DAS BASES DE TREINAMENTO  
 1 E DE DESDE ESTÃO COMO VARIÁVEIS (VALO-  
 0 RES NUMÉRICOS) E NA REALIDADE ELAS SÃO  
 0 NA REALIDADE VARIÁVEIS CATEGÓRICAS.  
 1  
 ⋮

# ENCODE DA CLASSE

```
BASE$DEFAULT <- factor(BASE$DEFAULT, levels =
  = c(0, 1))
```

NOTA DO INSTRUCTOR } VAMOS REFAZER TUDO NOVAMENTE PARA ADE-  
 QUAR A BASE ORIGINAL JÁ TRANSFORMANDO  
 AS VARIÁVEIS, INCLUSIVE A CLASSE META QUE  
 É A DEFAULT.

# LEITURA DA BASE DE DADOS

```
BASE <- read.esv('credit_data.esv')
```

# APAGANDO A COLUNA CLIENTID

```
BASE$CLIENTID <- NULL
```

# TRATANDO VALORES INCONSISTENTES

```
BASE <- ifelse(BASE$AGE < 0, 40.92, BASE$AGE)
```

# TRATANDO VALORES FALTANTES (NAs)

```
BASE$AGE <- ifelse(is.na(BASE$AGE), mean(
  (BASE$AGE, na.rm = TRUE), BASE$
  $AGE))
```

# FAZENDO O ESCALONAMENTO

```
BASE[, 1:3] <- scale(BASE[, 1:3])
```

# ENCODE DA CLASSE

```
BASE$DEFAULT <- factor(BASE$DEFAULT, levels =
  = c(0, 1))
```

# DIVISÃO DA BASE EM TREINAMENTO E TESTE

```
library(caretools)
```

```
set.seed(1)
```

```
DIVISAO <- sample.split(BASE$INCOME, splitratio =
  = 0.75)
```

# CRIANDO AS BASES TREINAMENTO E TESTE

```
BASE_TREINAMENTO <- subset(BASE, DIVISAO ==
  == TRUE)
```

```
BASE_TESTE <- subset(BASE, DIVISAO == FALSE)
```

```
library(e1071)
```

```
CLASSIFICADOR <- naiveBayes(X = BASE_TREINAMEN-
```



... CONT

SENSITIVITY: 0.9418

SPECIFICITY: 0.8868

POS PRED VALUE: 0.9859

NEG PRED VALUE: 0.6438

PREVALENCE: 0.8940

DETECTION RATE: 0.8420

DETECTION PREVALENCE: 0.8540

BALANCED ACCURACY: 0.9143

'POSITIVE' CLASS: 0

VAMOS SALVAR OS RESULTADOS EM UM ARQUIVO

↳ RESULTADOS\_CREDIT\_DATA.DXF\*

1 RESULTADOS CREDIT DATA

2 -----

3

4 0.936 - NAIVE BAYES - INCONSISTENTES + FALTANTES + ESCALONAMENTO

- NAIVE BAYES - INCONSISTENTES + FALTANTES

- NAIVE BAYES - SEM PRÉ-PROCESSAMENTO

# LEITURA DA BASE DE DADOS

BASE <- READ.ESV('CREDIT\_DATA.esv')

# APAGA-SE A COLUNA CLIENTID

~~BASE~~ BASE\$CLIENTID <- NULL

# ENCODE DA CLASSE

BASE\$DEFAULT <- FACTOR(BASE\$DEFAULT, LEVELS =  
= C(0, 1))

# DIVISÃO ENTRE TREINAMENTO E TESTE

LIBRARY(CATTOOLS)

SET.SEED(1)

DIVISAO <- SAMPLE.SPLIT(BASE\$INCOME, SPLITPRO-  
P = 0.75)

BASE\_TREINAMENTO <- SUBSET(BASE, DIVISAO == TRUE)  
BASE\_TESTE <- SUBSET(BASE, DIVISAO == FALSE)

LIBRARY(e1071)

CLASSIFICADOR <- NAIVE.BAYES(X = BASE\_TREINAMENTO[  
[ -4 ], Y = BASE\_TREINAMENTO\$DEFAULT)

# PRINT(CLASSIFICADOR)

PREDISOES <- PREDICT(CLASSIFICADOR, NEWDATA =  
= BASE\_TESTE[ -4 ])

MATRIZ\_CONFUSAO <- TABLE(BASE\_TESTE[ , 4 ], PRE-  
VISOES)

# PRINT(MATRIZ\_CONFUSAO)

# INSTALL.PACKAGES("earet")

LIBRARY(earet)

CONFUSION MATRIX (MATRIZ\_CONFUSAO)



... CONT

SEM PRÉ-PROCESSAMENTO → ACCURACY = 0.932

### NAIVE BAYES EM R - BASE CENSO

★ VAMOS CRIAR NOVO SCRIPT E SALVAR O MESMO COM O NOME "NAIVE\_BAYES\_CENSO.R"...

DEPOIS DE RODARMOS TODO O TEMPLAYTE, VAMOS AGORA IMPLEMENTAR O ALGORITMO NAIVE BAYES PARA A BASE CENSO.

# 1º Vamos começar a biblioteca necessária  
library(e1071)

# Criando o nome classificador Naive Bayes  
classificador ← naiveBayes(x = base\_treina-  
mentos[-15], y = base\_treina-  
mentos\$income)

print(classificador) # mostra a tabela de pro-  
babilidades

# Criando as previsões  
previsoes ← predict(classificador, newdata =  
base\_teste[-15])

previsoes # digite no prompt de comando.

# Vamos agora comparar os valores da base de teste que já sabemos com os valores previstos de forma a ver o quão o algoritmo acertou e errou.

# Vamos construir a matriz de confusão  
matriz\_confusao ← table(base\_teste[-15],  
previsoes)  
print(matriz\_confusao)

	0	1
0	3166	232
1	576	550

Erros do algoritmo } (1,0)  
                          (0,1)

Acertos do algoritmo } (0,0)  
                          (1,1)

- Acertou 3166 "0", ou seja, classificados como "0";
- Teve 232 erros para a classe "0";
- Teve 576 erros para a classe "1", ou seja errou + que acertou (550) para esta classe

# Corrigendo a biblioteca "caret"

library(caret) # permite ver o percentual de acerto.  
confusionMatrix(matriz\_confusao)

Confusion Matrix and Statistics

	Previsoes	
	0	1
0	3166	232
1	576	550



... cont

Accuracy: 0.8214  
 95% CI: (0.8099, 0.8325)  
 No Information Rate: 0.8271  
 P-Value [ACC > NSR]: 0.8512  
 Kappa: 0.468

McNemar's Test P-Value: < 2e-16

Sensitivity: 0.9461  
 Specificity: 0.7033  
 Pos Pred Value: 0.9317  
 Neg Pred Value: 0.4985  
 Prevalence: 0.8271  
 Detection Rate: 0.6978  
 Detection Prevalence: 0.7511  
 Balanced Accuracy: 0.7747  
 'Positive' class: 0

# Podemos, também, fazer os mesmos procedimentos para avaliar os resultados

### # Resultados Censos

- Naive Bayes - Categóricas + escalonamentos
- Naive Bayes - Categóricas
- Naive Bayes - Escalonamentos
- Naive Bayes - sem pré-proc.amentos

Nota do instrutor } O algoritmo Naive Bayes não utiliza cálculos de distâncias, por isso que não é necessário fazer o escalonamento caso não queira.

### Referências Complementares

- O instrutor recomenda seu curso Mineração de Emoções em Textos com Python e NLTK para quem deseja usar o algoritmo Naive Bayes para tarefas de classificação;
- O livro Thoughtful Machine Learning de Matthew Kirk: o Capítulo 4 apresenta mais sobre a teoria do Naive Bayes e mostra como contribuir um classificador de spam;
- Livro Data Algorithms de Mahmoud Perron: o Capítulo 14 apresenta explicações sobre a teoria do Naive Bayes.



# Aprendizagem por Árvores de Decisão

Outro paradigma de aprendizado de máquina, o qual tem o conteúdo programático abaixo:

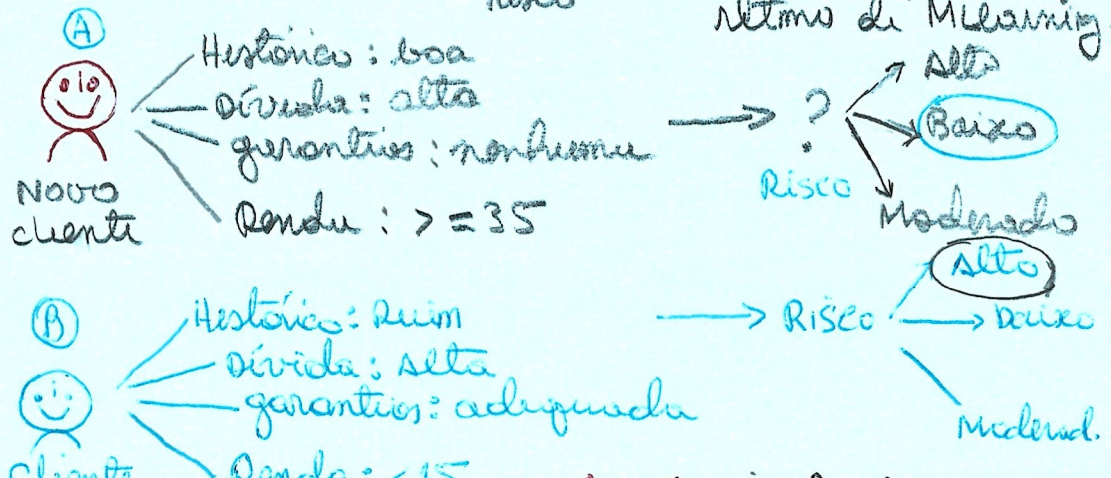
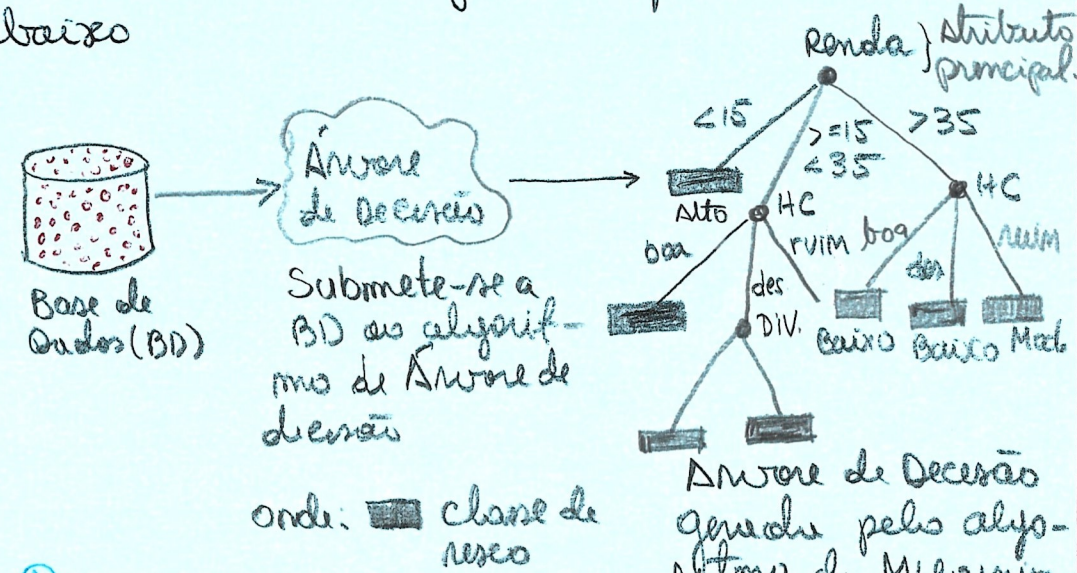
- Teoria sobre árvores de decisão;
- Árvores de decisão em R - base risco de crédito;
- Árvores de decisão em R - base censo;
- Podas em árvores de decisão (base censo);
- Random Forests (Florestas Aleatórias) # Essa técnica é usada no Kinect da MS para detecção de movimentos

## 1 INTRODUÇÃO

\* Moderação  
Classe meta

Base Original	História do crédito	Dívida	Garantias	Renda Anual	Risco
Desconhecida	Ruim	alta	nenhuma	< 15 000	alto
Desconhecida	Desconhecida	alta	nenhuma	$\geq 15 000 \wedge \leq 35 000$	alto
Desconhecida	Desconhecida	baixa	nenhuma	$\geq 15 000 \wedge \leq M^*$ $\leq 35 000$	alto
Desconhecida	Desconhecida	baixa	nenhuma	> 35 000	alto
	...	...	...	...	...

A ideia é fazer a previsão do risco em "alto", "moderado" ou "baixo" baseado nos atributos preditores (histórico do crédito, dívida, garantias e renda anual) usando algoritmo de árvore de decisão, veja o esquema estrutural abaixo



Nota: É importante definir qual vai ser o nó raiz, atributo principal previsor.



CLASSE META RISCO

A	A	M	A	B	B	A	M	B	B	A	M	B	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---

ONDE:

A: ALTO; M: MODERADO; B: BAIXO

ALTO = 6/14

MODERADO = 3/14

BAIXO = 5/14

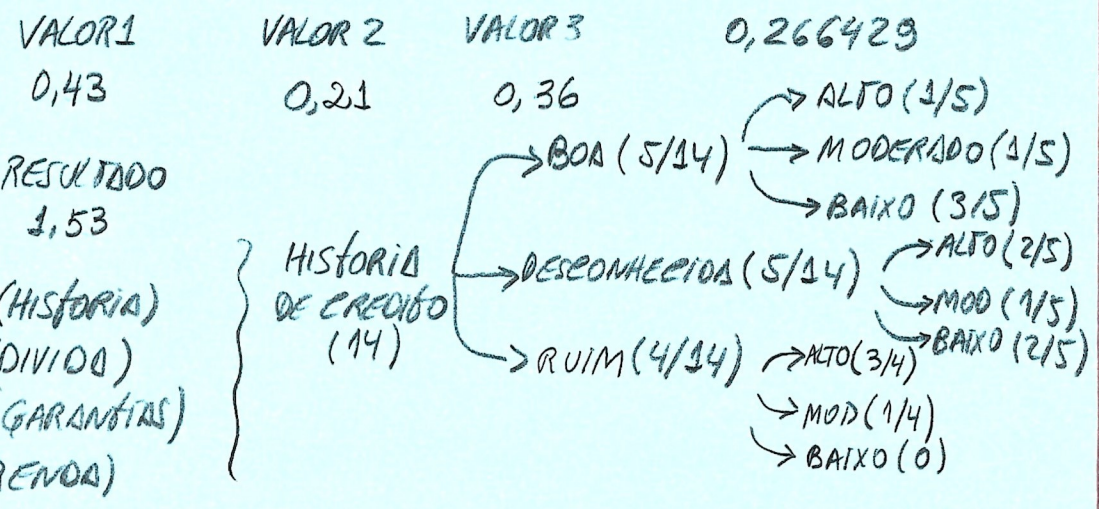
$$ENTROPY(S) = \sum_{i=1}^C -P_i \log_2 P_i$$

$$E(S) = -6/14 * \log_2(6/14; 2) - 3/14 * \log_2(3/14; 2) - 5/14 * \log_2(5/14; 2) = 1,53$$

NO EXCEL, FAÇA-SE:

ENTROPIA

GANHO DE INFORMAÇÃO



$$E(S) = -1/5 * \log_2(1/5; 2) - 4/5 * \log_2(4/5; 2) - 3/5 * \log_2(3/5; 2) = 1,37 \text{ \# HISTORIA DE CREDITO BOA}$$

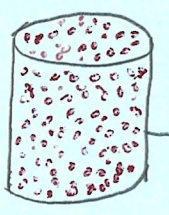
$$E(S) = -2/5 * \log_2(2/5; 2) - 1/5 * \log_2(1/5; 2) - 2/5 * \log_2(2/5; 2) = 1,52 \text{ \# HISTORIA DE CREDITO DESONHECIDA}$$

$$E(S) = -3/4 * \log_2(3/4; 2) - 1/4 * \log_2(1/4; 2) - 0 * \log_2(0; 2) = 0,81$$

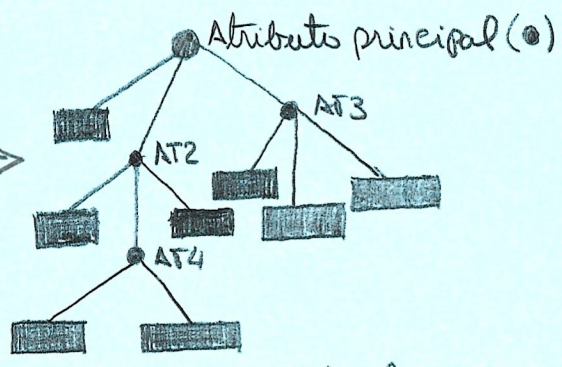
Naive Bayes

%	%	%	%	%
%	%	%	%	%
%	%	%	%	%
%	%	%	%	%

Tabela de Probabilidades



Base Treinamento

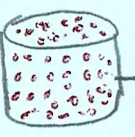


Árvore de decisão

Registros % Acerto

Naive Bayes

Árvore de decisão



Base Teste

ÁRVORE DE DECISÃO - APRENDIZAGEM ↓

$$Entropy(S) = \sum_{i=1}^C -P_i \log_2 P_i$$

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

... CONT

$$GAIN(S, A) = ENTROPY(S) - \sum_{V \in VALUES(A)} \frac{|S_V|}{|S|} ENTROPY(S_V)$$

$$GANHO(HISTORIA) = 1,53 - (5/14 * 1,37) - (5/14 * 1,52) - (4/14 * 0,81) = \boxed{0,26}$$

DÍVIDA/RISCO

$$E(S) = -4/7 * LOG(4/7; 2) - 1/7 * LOG(1/7; 2) - 2/7 * LOG(2/7; 2) = 1,38$$

$$GANHO(DÍVIDA) = 1,53 - (7/14 * 1,38) - (7/14 * 1,56) = \boxed{0,06}$$

GARANTIAS/RISCO

$$E(S) = -6/11 * LOG(6/11; 2) - 2/11 * LOG(2/11; 2) - 3/11 * LOG(3/11; 2) = 1,44$$

$$GANHO(GARANTIAS) = 1,53 - (11/14 * 1,44) - (3/14 * 0,92) = \boxed{0,20}$$

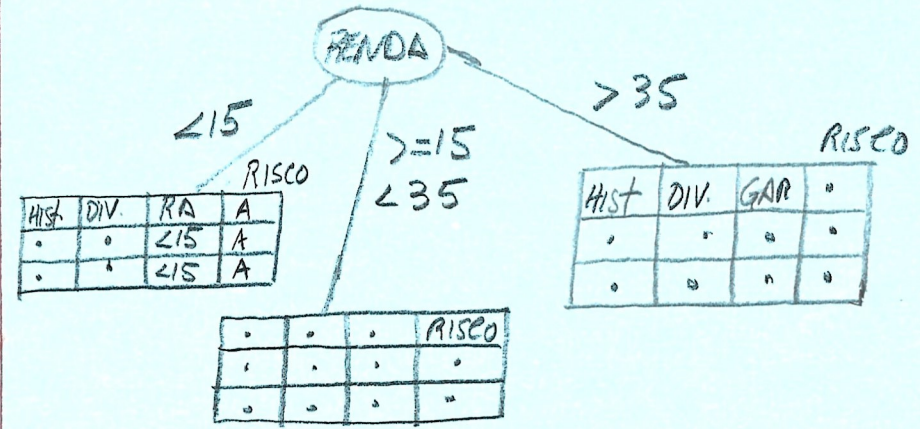
RENDA/RISCO

$$E(S) = -3/3 * LOG(3/3; 2) - 0 * LOG(0; 2) - 0 * LOG(0; 2) = 0,0$$

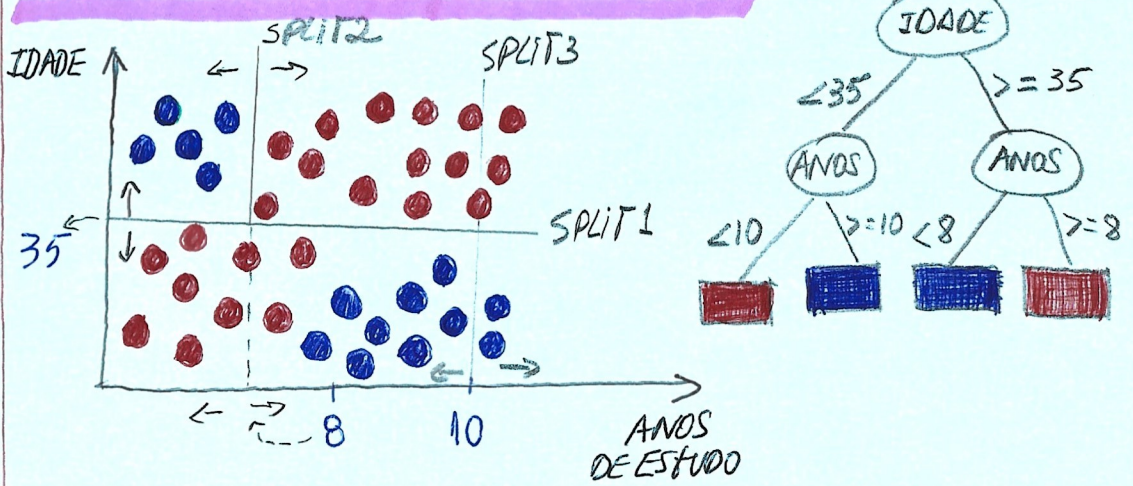
$$GANHO(RENDA) = 1,53 - (3/14 * 0,00) - (4/14 * 1,0) - (4/14 * 1,00) - (7/14 * 1,15) = \boxed{0,66}$$

RESUMO GANHO(S)

- HISTORIA DE CRÉDITO = 0,26
  - DÍVIDA = 0,06
  - GARANTIAS = 0,20
  - **RENDA = 0,66**
- Portanto, o nó RAÍZ DA ARVORE SERÁ "RENDA"



ÁRVORES DE DECISÃO - + CONCEITOS



A SEGUIR, VEREMOS O CONCEITO DE PODA EM ÁRVORES DE DECISÃO



# PODA EM ÁRVORES DE DECISÃO

## BIAS (VIÉS)

- ERROS POR CLASSIFICAÇÃO ERRADA

## VARIÂNCIA

- ERROS POR SENSIBILIDADE PEQUENAS A MUDANÇAS NA BASE DE TREINAMENTO;

- PODE LEVAR A OVERFITTING

ACONTECE QUANDO O ALGORITMO SE ADAPTA FA DETALHES AOS DADOS DE TREINAMENTO, OU SEJA, QUANDO DARMOS UMA NOVA BASE DE DADOS, O ALGORITMO COMETE MUITOS ERROS NO DESEMPENHO.

O ALGORITMO TEM UM ÓTIMO DESEMPENHO DE TREINAMENTO.

➤ É SEMELHANTE A UM ALUNO QUE DECORA UM CONTEÚDO EM VEZ DE APRENDER.

## VANTAGENS

- FÁCIL INTERPRETAÇÃO;
- NÃO PRECISA NORMALIZAÇÃO OU PADRONIZAÇÃO;
- RÁPIDO PARA CLASSIFICAR NOVOS REGISTROS

## DESvantagens

- GERAÇÃO DE ÁRVORES MUITO COMPLEXAS
- PEQUENAS MUDANÇAS NOS DADOS PODE MUDAR A ÁRVORE (PODA PODE AJUDAR);
- PROBLEMA NP-COMPLETO PARA CONSTRUIR A ÁRVORE

- ERAM MUITO POPULARES EM MEADOS DOS ANOS 90

- Upgrades como random forest (florestas randômicas) melhoram o desempenho (usado no Kinect da Microsoft) # Kinect é um aplicativo para capturar os movimentos do usuário e na sua reprodução.

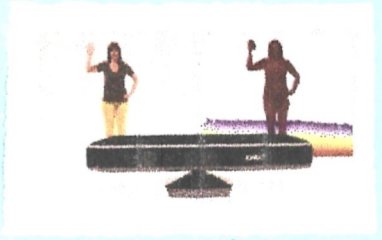


Imagem do Kinect da MS-Soft (google imagens)

- CART - Classification And Regression Trees

## Árvores de decisão em R - Base Rtree de

James neste tópico, em Crédito  
plomenta uma árvore de decisão no R usando como Base de dados a base Rtree de Crédito

1º passo: Criar novo arquivo de script, e depois salvar como "arvore-decisao-rtree-creditos.R".

> base<-read.esv('rtree-creditos.esv') # Na área de trabalho do RStudio, podemos observar que base tem 14 obs de 5 variáveis.

Nota do instalador } Para implementarmos árvores de decisão no R, precisamos instalar o pacote RPART

```

> install.packages("rpart")
> library(rpart)

# Quando o classificador
> classificador <- rpart(formula = ruseo ~.,
                          data = base)

> print(classificador)
n = 14
node), split, n, loss, yval, (yprob)
* denotes terminal node

```

1) root 14 8 alto (0.42855714 0.3571429 0.2142857)\*

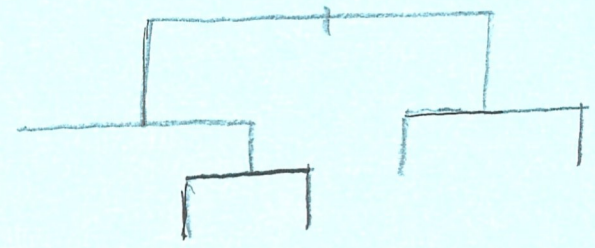
> plot(classificador) # após a execução deste comando, tem-se em esc:

Erro em plot.rpart(classificador): fit is not a tree, just a root. # Isso ocorre porque a nome base de dados é muito pequena. Para resolver esse problema, vamos passar um parâmetro adicional dentro da modelagem do classificador, conforme é mostrado a seguir

```

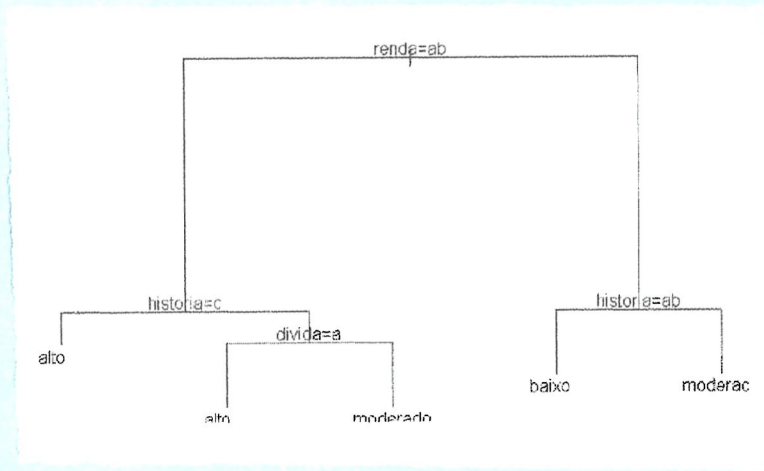
> classificador <- rpart(formula = ruseo ~.,
                        data = base, control = rpart.control(
                          minbucket = 1)) # Toda uma linha é o
> print(classificador) # aprendizado de máquina.
> plot(classificador)

```



# Para esboçar os ramos das nós, procede da forma abaixo

```
> text(classificador)
```



Nota do instalador } Para desajustarmos a árvore de como foi -  
 mais + interessante, vamos instalar outro  
 pacote chamado Rpart.plot

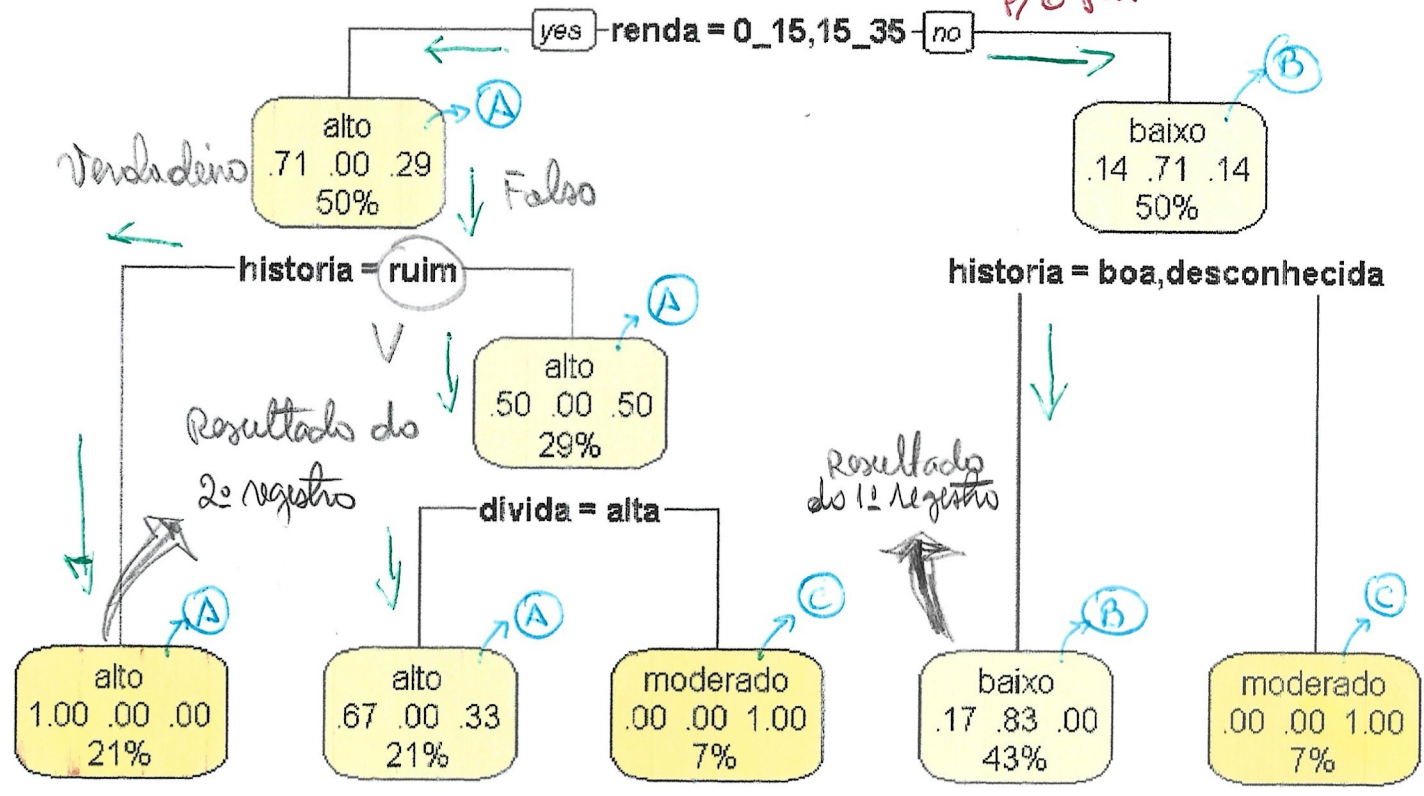
```

> install.packages("rpart.plot")
> library(rpart.plot)

```

> report.plot (classificador) ↵

- alto # A → novo crédito
- baixo # B → novo crédito
- moderado # C → vende



```
> renda2 = c("acima_35", "0_15")
```

```
> df <- data.frame(historia, divida, garantias, renda)
```

# Vamos agora classificar esses dois registros:

```
> prevoes <- predict(classificador, newdata = df)
```

```
> print(prevoes)
```

	alto	baixo	moderado
1	0.1666667	0.8333333	0
2	1.0000000	0.0000000	0

# Criando com data frame novo para testar o classificador

```
> historia <- c("boa", "ruim")
> divida <- c("alta", "alta")
> garantias <- c("nenhuma", "adequada")
```

ÁRVORES DE DECISÃO EM R - BASE CRÉDITO

# Criar novo arquivo de script com o nome "oware\_decisao\_credit\_dadu.RData"

```
# Leitura da base de dados
base <- read.esv("credit_data.esv")
```

```
# apaga a coluna clientid
base$clientid <- NULL
```



```
# Trata os valores inconsistentes
> base$age <- ifelse (base$age < 0, 40.92, base$age)
```

```
# Trata os valores faltantes
> base$age <- ifelse (is.na (base$age), mean (
  (base$age, na.rm = TRUE),
  , base$age)
```

```
# Faz o escalonamento
> base[, 1:3] <- scale (base[, 1:3])
```

```
# Divide base em base treinamento e teste
library (caTools)
```

```
set.seed (1)
> divisao <- sample.split (base$income,
  , splitRatio = 0.75)
```

```
> base_treino <- subset (base, divisao ==
  == TRUE)
```

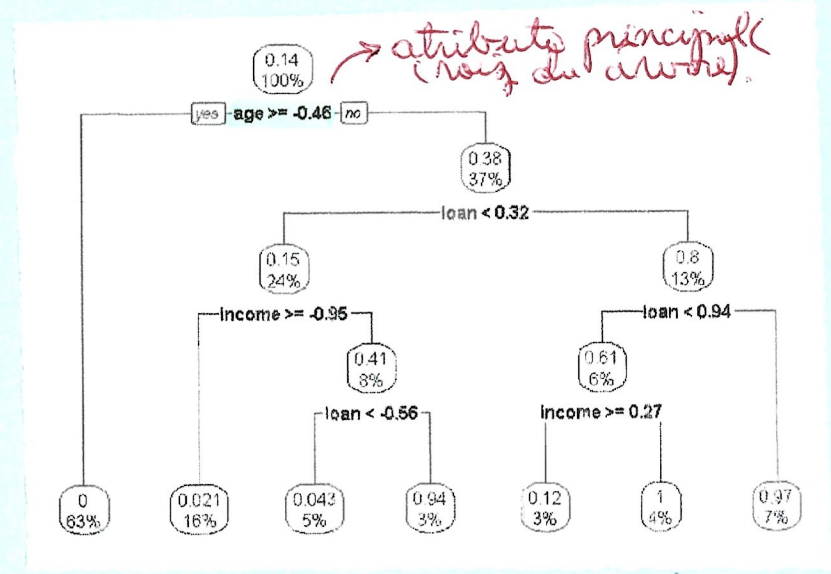
```
> base_teste <- subset (base, divisao == FALSE)
```

```
library (rpart)
> classificador <- rpart (formula = default ~.,
  , data = base_treino)
```

```
> print (classificador)
```



```
> library (rpart.plot)
> rpart.plot (classificador)
```



atributo principal (raiz do arvore)

Árvore de decisão: O algoritmo concluiu que o atributo "age" é o atributo principal

Nota do instrutor } Para que a idel seja principal não aparece < 0, basta refazer o pré-processamento sem isolar o escalonamento.

# Vamos agora submeter a base teste para a árvore de decisão construída. Esse teste vai servir para avaliar o desempenho da árvore de decisão construída na classificação dos registros.

```
> previsoes <- predict (classificador, newdata =
  = base_teste [-4])
```

```
# Encode de classe
> base$default <- factor (base$default,
  , levels = c (0, 1))
```



> library(caret) ↵  
 > confusionMatrix(matrix\_confusao) ↵

Base teste					Previsões		
	income	age	loan	default		0	1
4	-0.16302429	0.364045494	0.54484373	0	4	1.00000000	0.00000000
6	-1.42507390	1.240421432	-1.45427744	0	6	1.00000000	0.00000000
7	0.21829826	-1.055401088	0.41371763	0	7	0.00000000	1.00000000
15	1.23760692	-3.130236093	-0.05177375	0	15	1.00000000	0.00000000
18	-0.13895210	1.503231576	-0.38307884	0	18	1.00000000	0.00000000
20	0.94921304	-1.120717011	0.17112519	0	20	0.97916667	0.02083333
21	-1.55872582	-0.691701288	-1.08236711	0	21	0.95652174	0.04347826
29	0.90323913	3.000000000	-0.77550929	0	29	1.00000000	0.00000000
35	0.85530462	-3.321137171	-0.89181618	0	35	1.00000000	0.00000000
37	-1.35269971	1.302256115	-0.76370309	0	37	1.00000000	0.00000000
41	1.46984348	-0.606027373	0.850386060	0	41	0.88035238	0.11964762
43	1.20157166	3.612136204	-0.02874973	0	43	1.00000000	0.00000000
46	-0.25170358	-1.720334695	-1.07018594	0	46	0.97916667	0.02083333
52	0.85509890	-0.981509232	-0.36119440	0	52	0.97916667	0.02083333
51	-0.53543536	0.957212918	-0.76008667	0	61	1.00000000	0.00000000
60	-1.01900904	-1.410034319	0.49334986	1	68	0.00000000	1.00000000
70	0.84318658	0.651861900	0.05332618	0	70	1.00000000	0.00000000
					72	0.00000000	1.00000000
					76	1.00000000	0.00000000

Valores das probabilidades das classes.

# Para retornar se é da classe 0 ou 1, veja o código abaixo:

> previsoes <- predict(classificador, newdata = base\_teste[-4], type = 'class')

# Vamos agora gerar a matriz de confusão

> matrix\_confusao <- table(base\_teste[-4], previsoes) ↵

> print(matrix\_confusao) ↵

	Previsões	
	0	1
0	420	7
1	11	62

Considerações:

- acertou 420 "0" e 62 "1";
- Errou 11 "1" e 7 "0";
- Se  $\Sigma$  os erros, a máquina errou somente 18 registros de 427.

```
> library(caret)
> confusionMatrix(matrix_confusao)
Confusion Matrix and Statistics

previsoes
  0      1
0 420    7
1  11   62

Accuracy : 0.964
95% CI : (0.9437, 0.9785)
No Information Rate : 0.882
P-value [Acc > NIR] : 1.334e-14

Kappa : 0.8523

McNemar's Test P-value : 0.4795

sensitivity : 0.9745
specificity : 0.8986
pos Pred Value : 0.9836
neg Pred Value : 0.8493
prevalence : 0.8620
detection Rate : 0.8400
detection Prevalence : 0.8540
balanced Accuracy : 0.9365

'positive' class : 0
```

Percentual de acerto da árvore de decisão.

Desempenho da máquina (classificador) em novos dados (registros) (base teste).

Comparando o desempenho do classificador baseado em Naive Bayes e Árvore de decisão.

Resultados Credit Data

- 0.936 - Naive Bayes - Inconsistentes + faltam + Especial
- 0.936 - Naive Bayes - Inconsistentes + faltantes
- 0.932 - Naive Bayes - sempre - processamento
- 0.964 - Árvore - Inconsistentes + faltantes + realocam.
- 0.964 - Árvore - Inconsistentes + faltantes
- 0.964 - Árvore - sem pré-processamento

Próximo tópico, estudaremos Árvore de decisão em R - base censo.



# ÁRVORES DE DECISÃO EM R - BASE CENSO

1º PASSO } PRE-PROCESSAMENTO DO TEMPLATE BASE CENSO, CONFORME DESCRITO ABAIXO:

■ CRIAR NOVO SCRIPT E SALVÁ-LO COM O NOME "ARVORE.DECISAO.CENSUS.RDATA" A EXTENSÃO O R COLOCA AUTOMATICAMENTE.

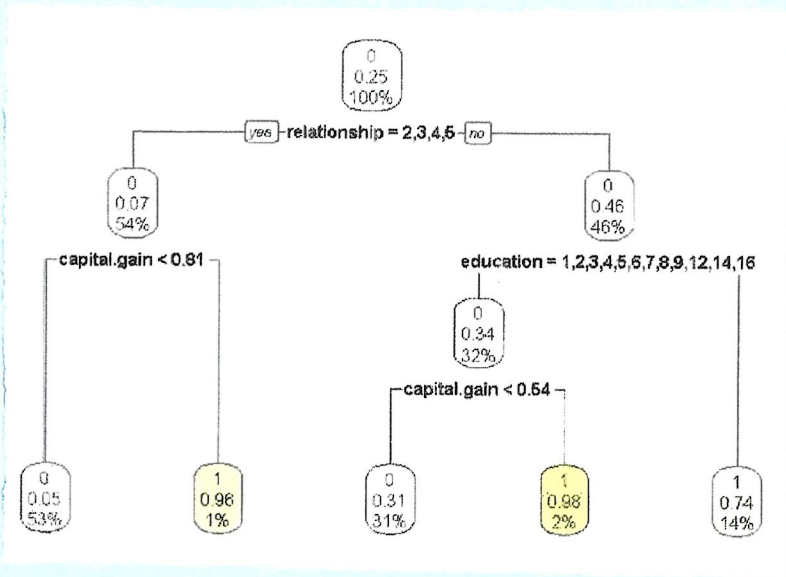
NOTA DO INSTRUCTOR } A BASE TESTE É PARA FAZER A AVALIAÇÃO DO ALGORITMO.

■ LIBRARY(RPART)

> CLASSIFICADOR <- RPART(FORMULA = INCOME N., DATA = BASE TREINAMENTO) <-

# VISUALIZANDO A ÁRVORE DE DECISÃO

> RPART.PLOT(CLASSIFICADOR) <- # MOSTRA ÁRVORE.



## # FAZENDO AS PREVISÕES

> PREVISOES <- PREDICT(CLASSIFICADOR, NEWDATA = BASE\_TESTE[-15], TYPE = 'CLASS')

NOTA DO INSTRUCTOR } TYPE = 'CLASS' É PARA MOSTRAR A CLASSE META EM '0' E '1' E NÃO AS PROBABILIDADES DAS CLASSES EM CADA REGISTRO.

> PREVISOES # MOSTRA TODAS AS PREVISÕES PARA OS REGISTROS DO BASE\_TESTE.

## # GERANDO A MATRIZ DE CONFUSÃO

> MATRIZ\_CONFUSAO <- TABLE(BASE\_TESTE[-15], PREVISOES) <-

> PRINT(MATRIZ\_CONFUSAO)

		PREVISOES	
		0	1
0	3216	182	
1	564	562	

> LIBRARY(CARET) # VER O PERCENTUAL DE ACERTO

> CONFUSIONMATRIX(MATRIZ\_CONFUSAO) <-

ACCURACY: 0.8351

95% CI: (0.824, 0.8458)

NO INFORMATION RATE: 0.8355

KAPPA: 0.5025



# COMPARATIVO ENTRE O DESEMPENHO DOS ALGORITMOS

## RESULTADOS CENSUS

- 0.8214 - NAIVE BAYES - CATEGÓRICAS + ESCALONAMENTO
- 0.8216 - NAIVE BAYES - CATEGÓRICAS
- 0.8214 - NAIVE BAYES - SEM PRÉ-PROCESSAMENTO
- #-----#-----
- 0.8351 - ÁRVORE DE DECISÃO - CATEGÓRICAS + ESCALONAM.
- 0.8351 - ÁRVORE DE DECISÃO - CATEGÓRICAS
- 0.8351 - ÁRVORE DE DECISÃO - ESCALONAMENTO
- 0.8351 - ÁRVORE DE DECISÃO - SEM PRÉ-PROCESSAMENTO.

## ÁRVORES DE DECISÃO EM R - PODA BASE CENSO

RODA-SE TODO O SCRIPT "ARVORE\_DECISAO - CENSUS.RDATA"

# LOGO ABAIXO DO CLASSIFICADOR, VAMOS COLAR UM CÓDIGO ADICIONAL, CONFORME MOSTRADO ABAIXO:

```
> PODA <- CLASSIFICADOR$CP TABLE[WHICH.MIN(
  (CLASSIFICADOR$CP TABLE[, "XERROR"],
  ), "CP"] # CRIA PODA = 0.01 COMO MENOR ERRO
```

```
> PRINT (CLASSIFICADOR$CP TABLE) <
```

	CP	NSPLIT	REL ERRO	XERROR	XSTD
1	0.13232529	0	1.0000000	1.0000000	0.0108...
2	0.06424318	2	0.7353494	0.7353494	0.00970...
...	0.01000000	4	0.6342839	0.6342839	0.0091...

## > PRUNE (CLASSIFICADOR, PODA)

```
> prune(classificador, poda)
n= 25638
node), split, n, loss, yval, (yprob)
# denotes terminal node
1) root 25638 6382 0 (0.75107263 0.24892737)
2) relationship=2,3,4,5 13852 955 0 (0.93105689 0.06894311)
4) capital.gain< 0.8076171 13597 710 0 (0.94778260 0.05221740) *
5) capital.gain<= 0.8076171 255 10 1 (0.03921569 0.96078431) *
3) relationship=1,6 11786 5427 0 (0.52953844 0.46046156)
6) education=1,2,3,4,5,6,7,8,9,12,14,16 8261 2820 0 (0.65863697 0.34136303)
12) capital.gain< 0.5405489 7835 2402 0 (0.69342693 0.30657307) *
13) capital.gain<= 0.5405489 426 8 1 (0.01877934 0.98122066) *
7) education=10,11,13,15 3525 918 1 (0.26042553 0.73957447) *
```

## ÁRVORES DE DECISÃO - RANDOM FOREST



CLASSIFICADOR 3 → DIZ QUE ORISEO É MODER.  
 CLASSIFICADOR 1 → DIZ QUE ORISEO É ALTO  
 CLASSIFICADOR 2 → DIZ QUE ORISEO É BAIXO  
 CLASSIFICADOR N  
 ... COMBINAR AS DECISÕES

+ DE 01 ÁRVORE, DEZENAS DE ÁRVORES

# Este algoritmo utiliza a decisão de várias árvores.

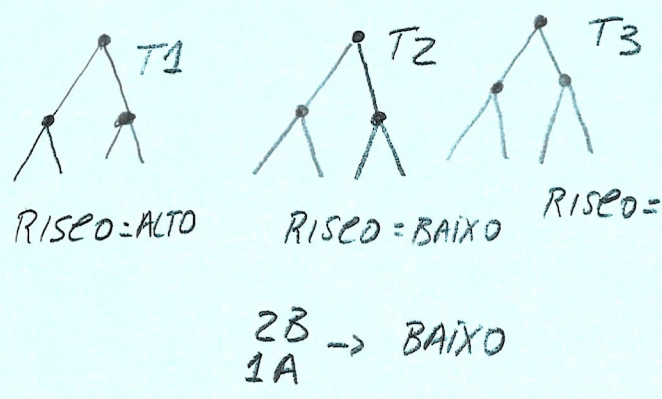
## RANDOM FOREST

- ENSEMBLE LEARNING (APRENDIZAGEM EM CONJUNTO)
- "A IDÉIA É CONSULTAR DIVERSOS PROFISSIONAIS PARA TOMAR UMA DECISÃO"
- VÁRIOS ALGORITMOS JUNTOS PARA CONSTRUIR UM ALGORITMO MAIS "FORTE";
- USA A MÉDIA (REGRESSÃO) OU VOTOS DA MAIORIA (CLASSIFICAÇÃO) PARA DAR A RESPOSTA FINAL



Nota do instrutor } A QUANTIDADE DE ÁRVORES DE DECISÃO É UMA PARAMETRIZAÇÃO DO ALGORITMO, OU SEJA, VOCÊ PASSA PARA O ALGORITMO A QUANTIDADE DE ÁRVORES QUE VOCÊ QUER.

SUPosição:



Portanto, → A RESPOSTA FINAL DO ALGORITMO VAI SER BAIXO.

RANDOM FOREST

HISTÓRICO DE CREDITO	DÍVIDA	GARANTIAS	RENDA ANUAL	RISCO
RUIM	ALTA	NENHUMA	<15.000	ALTO
DESCONHECIDA	ALTA	NENHUMA	>=15.000	ALTO
DESCONHECIDA	Baixa	Nenhuma	>=15000 a med. <= 35000	
DESCONHECIDA	Baixa	Nenhuma	>35000	Baixo
⋮	⋮	⋮	⋮	⋮

Nota do instrutor } A floresta é dita "Random" porque escolhe de forma aleatória  $K$  atributos

... para comparação da métrica de pureza / impureza (impureza de gini/entropia)

Suposição:

$K=3$  #NÚMERO DE ATRIBUTOS

ÁRVORES = 3



• RENDA  
• HISTÓRIA DE CREDITO  
• DÍVIDAS

A BASE DE DADOS SÓ TEM 4 ATRIBUTOS, ENTRETANTO,  $K=3$ . NESSE CASO, HAVERÁ

A ESCOLHA ALEATÓRIA DOS ATRIBUTOS PARA A CONSTRUÇÃO DAS ÁRVORES DE DECISÃO.

Artigo: REAL-TIME HUMAN POSE RECOGNITION IN PARTS FROM SINGLE DEPTH IMAGES

RANDOM FOREST EM R - BASE CREDITO

- ▶ PASSO 1: ABRIR O ARQUIVO "TEMPLATE\_CREDITO.DAT";
- ▶ PASSO 2: COPIAR TODO O CÓDIGO DESSE SCRIPT;
- ▶ PASSO 3: CRIAR NOVO SCRIPT;
- ▶ PASSO 4: SALVAR O SCRIPT COM O NOME DE "RANDOM-FOREST\_CREDITO.DAT";
- ▶ PASSO 5: COLAR O CONTEÚDO, PASSO 2, NO SCRIPT NOVO;
- ▶ PASSO 6: EXECUTAR TODO O CÓDIGO;
- ▶ PASSO 7: INSTALAR O PACOTE "RANDOMFOREST";

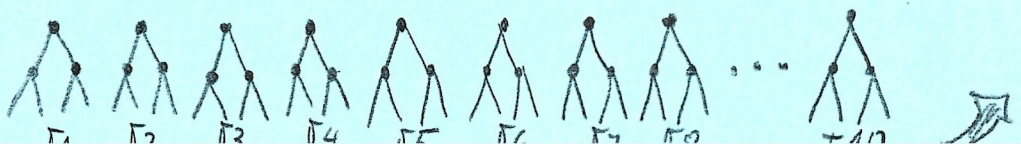


```
# INSTALANDO O PACOTE "RANDOMFOREST"
> INSTALL.PACKAGES("RANDOMFOREST")
> # CARREGANDO PACOTE NECESSÁRIO
> LIBRARY(RANDOMFOREST)
> CRIANDO O CLASSIFICADOR
> CLASSIFICADOR <- RANDOMFOREST(X = BASE_
  - TREINAMENTO[-4], y = BASE_TREINAMEN
  TO$DEFAULT, N.TREE = 10)
```

NOTA DO INSTRUCTOR } ANTES DE EXECUTAR O CLASSIFICADOR, VOLTAR NO SCRIPT E INCLUIR A LINHA DE CÓDIGO, ~~DEPOIS~~ DEPOIS DA LINHA DE COMANDO ABAIXO:

```
# ESCALONAMENTO
> BASE[, 1:13] <- SCALE(BASE[, 1:13])
# ENCODE PARA A CLASSE
BASE$DEFAULT <- FACTOR(BASE$DEFAULT,
  , LEVELS = C(0, 1))
```

NOTA DO INSTRUCTOR } O TREINAMENTO JÁ ESTÁ PRONTO A PARTIR DA CRIAÇÃO DO CLASSIFICADOR, COM UMA FLORESTA DE 10 ÁRVORES.



```
# CRIANDO AS PREVISÕES
> PREVISOES <- PREDICT(CLASSIFICADOR, NEWDATA =
  = BASE_TESTE[-4])
> PREVISOES
```

NA BASE TESTE ESSE VALOR DEU 0, OU SEJA, O ALGORITMO ACERTOU. E ISSO QUE VAMOS COMPARAR. ESSAS PREVISÕES FORAM FEITAS PELO CLASSIFICADOR

previsões

4	6	7	15	18	20	21	29	35	37	41	43	46	52	61	68	70	72	76
99	104	109	111	121	125	135	139	150	162	164	165	172	173	176	180	183	185	187
1	1	0	0	0	0	1	0	0	0	1	0	0	0	0	1	0	1	0
215	218	219	225	230	243	250	251	252	260	264	268	280	283	293	300	306	308	313
331	334	340	341	353	355	357	359	362	366	373	378	384	387	388	390	393	403	404
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
429	430	431	436	440	448	451	458	466	472	477	481	483	484	489	492	499	506	508
0	1	0	1	0	0	0	0	0	0	1	1	0	1	0	0	0	0	0
530	542	547	551	554	555	559	560	565	570	575	579	582	589	593	595	601	602	605
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
641	643	645	648	650	652	660	662	667	671	673	681	682	690	693	697	699	701	702
0	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0	0	0
752	759	761	763	769	772	774	780	781	785	793	799	801	803	808	814	817	818	824
0	1	0	0	0	0	0	0	1	1	1	0	0	0	0	0	0	0	0
851	854	872	873	880	882	883	884	886	887	901	902	912	920	922	923	925	926	932
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
957	968	969	971	974	980	983	989	996	1004	1009	1011	1014	1016	1017	1022	1024	1026	1034
0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0
1040	1075	1086	1089	1096	1100	1109	1115	1120	1122	1124	1125	1126	1131	1139	1143	1145	1149	1151

A IDEIA VAI ATÉ OS 500 REGISTROS DA BASE TESTE AGORA É COMPARAR COM OS VALORES DA BASE TESTE QUE JÁ SABEMOS OS VALORES DA CLASSE META PARA VER QUAL FOI O DESEMPENHO DO ALGORITMO DE MACHINE LEARNING, OU SEJA, QUAL FOI SEU PERCENTUAL DE ACERTO E DE ERRO.



### #CRIANDO A MATRIZ DE CONFUSÃO

```
> MATRIZ_CONFUSAO <- TABLE(BASE_TESTE[4],
                             PREVISOES)
> PRINT(MATRIZ_CONFUSAO)
```

	PREVISOES	
	0	1
0	423	4
1	7	66

#O ALGORITMO SÓ ERROU 11 REGISTROS DOS 500 REGISTROS DO DADOS DA BASE\_TESTE.

→ ERROS

↙ ACERTOS

### #MOSTRANDO AS ESTATÍSTICAS DE DESEMPENHO

```
LIBRARY(CARET) #CARREGA O PACOTE NECESSÁRIO
CONFUSIONMATRIX(MATRIZ_CONFUSAO) ↙
```

### CONFUSION MATRIX AND STATISTICS

	PREVISOES	
	0	1
0	423	4
1	7	66

Accuracy: 0,978  
95% CI: (0,961, 0,989)

NO INFORMATION RATE : 0.86  
P-VALUE [ACC > NIR] : < 2e-16  
KAPPA : 0.9102  
MCNEMAR'S TEST P-VALUE: 0.5465

SENSITIVITY : 0.9837  
SPECIFICITY : 0.9429  
POS PRED VALUE : 0.9906  
NEG PRED VALUE : 0.9041  
PREVALENCE : 0.8600  
DETECTION RATE : 0.8460  
DETECTION PREVALENCE : 0.8540  
BALANCED ACCURACY : 0.9633  
'POSITIVE' CLASS : 0

NOTA DO INSTRUCTOR } REFAZER O TESTE PARA n=15 ÁRVORES # A Accuracy: 0.982, ou seja aumentou.

# DEPOIS DISTO, FAZ-SE O COMPARATIVO DE DESEMPENHO DE ALGORITMO PARA AS TÉCNICAS NAIVE BAYES, ÁRVORE, RANDOM FOREST 30, BASEADO NOS TIPOS DE PREPARAÇÃO DE DADOS EM:

- INCONSISTENTES + FALTANTES + ESCALONAMENTO;
- INCONSISTENTES + FALTANTES; SEM PRÉ-PROCESSAMENTO DE DADOS PRÉVIOS

## RESULTADOS CREDIT DATA

- 0.936 - NAIVE BAYES - INCONSISTENTES + FALTANTES + ESCALONAMENTO
- 0.936 - NAIVE BAYES - INCONSISTENTES + FALTANTES
- 0.932 - NAIVE BAYES - SEM PRÉ-PROCESSAMENTO
- \* 0.964 - ÁRVORE - INCONSISTENTES + FALTANTES + ESCALONAMENTO
- \* 0.964 - ÁRVORE - INCONSISTENTES + FALTANTES
- \* 0.964 - ÁRVORE - SEM PROCESSAMENTO
- RANDOM FOREST 30 - INCONSISTENTES + FALTANTES + ESCALONAMENTO (0.984)
- RANDOM FOREST 30 - INCONSISTENTES + FALTANTES (0.984)
- RANDOM FOREST 30 - SEM PRÉ-PROCESSAMENTO (0.988)

Nota do instrutor } PARA O ALGORITMO "RANDOM FOREST", NECESSARIAMENTE, PRECISA-SE CORRIGIR VALORES FALTANTES (NAs), POIS SENÃO DAR ERRO NA HORA DA EXECUÇÃO DOS CÓDIGOS.

## IMPLEMENTAÇÃO DO RANDOM FOREST EM R - BASE CENSO

- 1º PASSO: PEGAR O ARQUIVO TEMPLATE\_CENSUS.R;
- 2º PASSO: COPIAR TODO O CÓDIGO E COLAR NUM NOVO SCRIPT, SALVANDO ESSE COM O NOME "RANDOM\_FOREST\_CENSUS". APÓS ISSO, RODAR O CÓDIGO, CUJA A SAÍDA É MOSTRADA ABAIXO.

```
BASE      30162 OBS. OF 15 VARIABLES
BASE_TESTE 4524 OBS. OF 15 VARIABLES
BASE_TREINAMENTO 25638 OBS. OF 15 VARIABLES
VALUES
DIVISAO   LOG1 [1:30162] TRUE TRUE TRUE ...
```

```
> LIBRARY(RANDOMFOREST)
> SET.SEED(1)
> CLASSIFICADOR <- RANDOMFOREST(X = BASE_TREINAMENTO[-15], Y = BASE_TREINAMENTO$INCOME,
  N.TREE = 10)
# CRIANDO A VARIÁVEL PREVISÕES
PREVISÕES <- PREDICT(CLASSIFICADOR, NEWDATA = BASE_TESTE[-15])
# MOSTRANDO PREVISÕES
> PREVISÕES
```

BASE ORIGINAL INCOME	PREVISÕES (MACHINE LEARNING) INCOME
0	# ACERTO
0	# ACERTO
0	# ACERTO
1	# ACERTO
1	# ERRO

# CRIANDO A MATRIZ DE CONFUSÃO

```

> MATRIZ_CONFUSAO <- TABLE(BASE_TESTE[, 15],
                             PREVISOES)
> PRINT(MATRIZ_CONFUSAO)
  
```

	PREVISOES	
	0	1
0	3135	263
1	429	697

ERROS → (0, 1) and (1, 0)  
 ACERTOS → (0, 0) and (1, 1)

# VENDO ESTATÍSTICAS DE DESEMPENHO DO MACHINE LEARNING

```

> LIBRARY(CARET) <-
> CONFUSIONMATRIX(MATRIZ_CONFUSAO) <- # SAÍDA ABAIXO
CONFUSIONM MATRIX AND STATISTICS
  
```



```

... PREVISOES
      0  1
0  3135 263
1  429 697

ACCURACY : 0.847
95% CI : (0.8362, 0.8574)
NO INFORMATION RATE : 0.7878
P-VALUE [ACC > NSR] : < 2.2e-16
KAPPA : 0.5697
MCNEMAR'S TEST P-VALUE : 3.55e-10
  
```

```

SENSITIVITY : 0.8796
SPECIFICITY : 0.7260

POS PRED VALUE : 0.9226
NEG PRED VALUE : 0.6190
PREVALENCE : 0.7878
DETECTION RATE : 0.6930
DETECTION PREVALENCE : 0.7511
BALANCED ACCURACY : 0.8028
'POSITIVE' CLASS : 0
  
```

## REFERÊNCIAS COMPLEMENTARES

- DATA MINING WITH DECISION TREES: THEORY AND APPLICATIONS (MACHINE PERCEPTION AND ARTIFICIAL INTELLIGENCE) DE ODED Z. MAIMON: LIVRO DE FÁCIL COMPREENSÃO E FOCADO EM ÁRVORES DE DECISÃO.
- PROGRAMS FOR MACHINE LEARNING (MORGAN KAUFMANN SERIES IN MACHINE LEARNING) DE J. ROSS QUINLAN: UM DOS LIVROS MAIS CLÁSSICOS SOBRE O ASSUNTO
- DECISION TREES AND RANDOM FOREST: A VISUAL INTRODUCTION FOR BEGINNERS: A SIMPLE GUIDE TO MACHINE LEARNING WITH DECISION TREES DE CHRIS SMITH E MARK KONING: UM DOS LIVROS MAIS FÁCEIS E DIDÁTICOS SOBRE O ASSUNTO.

# CONTINUA NA PÁGINA 55 COM  
REGRAS DE ASSOCIAÇÃO