

# **INTRODUÇÃO À MANIPULAÇÃO E TRATAMENTO DE DADOS COM R**

IDRIS DA SILVA SANTOS

DEPARTAMENTO DE ESTATÍSTICA E CIÊNCIAS ATUARIAIS – CCET – UFS

26 DE NOVEMBRO DE 2025

# ROTEIRO DA AULA

- 1 Introdução e Objetivos
- 2 O Dialetos Tidyverse
- 3 Ambiente e Pacotes
- 4 Importação e Diagnóstico
- 5 Limpeza de Dados (Data Wrangling)
- 6 Análise Exploratória
- 7 Modelagem: Árvore de Decisão
- 8 Avaliação de Desempenho
- 9 Conclusão

# INTRODUÇÃO E OBJETIVOS

# CRONOGRAMA DO WORKSHOP

Este curso foi desenhado para cobrir o fluxo completo de trabalho científico e analítico:

## Dia 1: Hoje (Manipulação de Dados)

**Foco:** Tratamento de dados e preparação.

Aprenderemos a transformar dados brutos ("sujos") em uma base sólida para análise usando **R**.

## Dia 2: Amanhã (Visualização Interativa)

**Foco:** Business Intelligence (BI).

Criação de dashboards dinâmicos com dados usando **PowerBI**.

## Dia 3: Encerramento (Escrita Técnica)

**Foco:** Comunicação e Publicação.

Como estruturar e formatar relatórios técnicos profissionais e acadêmicos utilizando **LaTeX**.

# CONTEXTUALIZAÇÃO

## ■ O Problema Real:

- ▶ Dados do mundo real raramente vêm prontos para análise.
- ▶ Eles são "sujos", desorganizados e inconsistentes.

## ■ O Fluxo de Ciência de Dados:

Importar → **Limpar/Arrumar** → Explorar → Modelar → Comunicar

## ■ Ferramenta: Linguagem R + Ecossistema *Tidyverse*.

# OBJETIVOS DE APRENDIZAGEM

Ao final desta oficina, você será capaz de:

1. **Importar dados** diagnosticando erros comuns (separadores, encoding).
2. **Limpar dados** usando o dialeto *Tidyverse* (`dplyr`, `tidyr`, `stringr`).
3. **Tratar inconsistências** como tipos errados, duplicatas e valores ausentes (NA).
4. **Visualizar** distribuições básicas com `ggplot2`.
5. **Treinar um modelo** simples de Árvore de Decisão.

# O DIALETO TIDYVERSE

# O ECOSISTEMA TIDYVERSE

## O que é?

- Uma coleção de pacotes R desenhados para ciência de dados.
- Compartilham uma filosofia comum de design, gramática e estrutura de dados.
- **Filosofia:** O código deve ser legível por humanos, não apenas por computadores.

## O Operador Pipe (%>%)

- Lê-se como "então" ou "e depois".
- Conecta as ações em uma sequência lógica.
- Evita o aninhamento confuso de funções:  $f(g(h(x)))$  vira  $x \rightarrow h \rightarrow g \rightarrow f$ .

# A GRAMÁTICA DE MANIPULAÇÃO: OS VERBOS

O pacote dplyr propõe que a manipulação de dados seja feita através de **\*\*verbos\*\*** intuitivos. Imagine que você está escrevendo uma frase:

- **filter( )**: *Filtrar linhas baseadas em condições (seleciona casos).*
- **select( )**: *Seleciona colunas pelo nome (seleciona variáveis).*
- **mutate( )**: *Modifica ou cria novas colunas (transformação).*
- **arrange( )**: *Arranja (ordena) as linhas.*
- **summarise( )**: *Sumariza os dados em uma estatística (média, total).*

## Exemplo de "Frase" em R:

```
1 dados %>%
2   filter(diagnostico == "M") %>% # Pegue os malignos, E DEPOIS...
3   select(area_mean, texture_mean) %>% # Seleccione essas colunas, E DEPOIS...
4   mutate(razao = area_mean / texture_mean) # Crie esta nova variavel.
```

# **AMBIENTE E PACOTES**

# PREPARANDO O AMBIENTE

Vamos instalar e carregar as bibliotecas essenciais.

- tidyverse: Manipulação e visualização.
- janitor: Limpeza de nomes de colunas.
- vroom: Leitura rápida de dados.
- tidymodels: Framework de modelagem.

```
1 pacotes <- c("janitor", "tidyverse", "tidymodels", "vroom", "corrr", "GGally", "rpart
  .plot")
2 instalados <- installed.packages()
3 for(pacote in pacotes) {
4   if (!(pacote %in% instalados)) install.packages(pacote)
5 }
6 library(tidyverse)
7 library(tidymodels)
8 library(janitor)
9 library(vroom)
10 library(corrr)
11 library(GGally)
12 library(rpart.plot)
```

# IMPORTAÇÃO E DIAGNÓSTICO

# IMPORTAÇÃO INICIAL (O CAOS)

Tentativa ingênua de ler o arquivo CSV. O que pode dar errado?

```
1 # Leitura inicial
2 dados <- vroom::vroom("./wbcד_dirty.csv")
3
4 # Visualizar a estrutura bruta
5 dados %>% pillar::glimpse()
```

## Problemas observados:

- Colunas numéricas lidas como texto (chr).
- Mistura de separadores decimais (. e ,).
- Aviso de *parsing error*.

# DIAGNOSTICANDO ERROS DE ESTRUTURA

O R nos avisou sobre problemas. Vamos investigar onde eles estão.

```
1 # Investigando falhas de leitura  
2 dados %>% vroom::problems()
```

## O Diagnóstico:

- O erro aponta para a linha 3.
- Esperava-se 32 colunas, mas o parser encontrou mais (ou caracteres inesperados que quebraram a estrutura).

# INTERVENÇÃO CIRÚRGICA (STRINGS)

Às vezes, o erro está na formação do arquivo (aspas não fechadas, por exemplo). Vamos tratar o arquivo como texto puro para corrigir.

```
1 # Lendo como texto puro
2 linhas <- readr::read_lines("./wbcד_dirty.csv")
3
4 # Identificando tokens problemáticos (Regex Avançado)
5 tokens <- stringr::str_split(linhas[3], ',')
6 tokens %>%
7   purrr::map(grep, pattern = '(?:(?!").(,|$))+', perl = TRUE) %>%
8   purrr::map(~tokens[[1]][.])
9
10 # Corrigindo aspas desbalanceadas na linha 3
11 linhas[3] <- linhas[3] %>%
12   stringr::str_replace('(\d+\.\d+),(\d+\.\d+)', '\1,\2')
13
14 # Salvando a versão corrigida
15 linhas %>% readr::write_lines("./wbcד_01.csv")
```

# LIMPEZA DE DADOS (DATA WRANGLING)

# REIMPORTAÇÃO E NOMES DE COLUNAS

Agora importamos o arquivo corrigido, especificando que o decimal é ponto (padrão US). Em seguida, padronizamos os nomes das colunas.

```
1 # Importando com locale correto
2 dados <- vroom::vroom(
3   "./wbcד_01.csv",
4   locale = locale(decimal_mark = ".") # Importante!
5 )
6
7 dados %>% pillar::glimpse()
8
9 # Verificando nomes originais (ruins)
10 dados %>% names()
11
12 # Limpando nomes automaticamente (snake_case)
13 dados <- janitor::clean_names(dados)
14 dados %>% names()
```

# TRATANDO STRINGS "N/A"

Temos valores "N/A"(texto) que deveriam ser NA (valor nulo do R). Isso impede cálculos numéricos.

```
1 # Substituindo "N/A" por NA real em todas as colunas de texto
2 dados <- dados %>%
3   dplyr::mutate(across(where(is_character), ~na_if(., "N/A")))
4
5 # Verificando a mudança
6 dados %>% pillar::glimpse()
```

# CONVERTENDO TEXTO PARA NÚMEROS

Muitas colunas ainda são texto (chr) por causa de sujeiras (vírgulas no lugar de pontos, caracteres estranhos). Vamos forçar a conversão (parse).

```
1 # Aplicando parse_number em colunas de texto (exceto ID e Diagnosis)
2 dados <- dados %>%
3   dplyr::mutate(
4     across(
5       where(is.character) & !c(id, diagnosis),
6       ~parse_number(
7         str_replace_all(.x, ",", ".") # troca "," por "."
8       )
9     )
10   )
11
12 dados %>% pillar::glimpse()
```

# REMOVENDO DUPLICATAS

Linhas repetidas enviesam a análise estatística.

```
1 # Verificando duplicatas
2 dados %>% janitor::get_dupes()
3
4 # Removendo duplicatas (mantendo apenas linhas distintas)
5 dados <- dados %>% dplyr::distinct()
6
7 # Confirmando limpeza
8 dados %>% janitor::get_dupes()
```

# PADRORIZANDO CATEGORIAS (DIAGNOSIS)

A coluna alvo possui inconsistências: "M", "m", "B", "b", "Benign". Precisamos uniformizar.

```
1 # Verificando valores unicos
2 dados$diagnosis %>% str_unique()
3
4 # Padronizando com case_when
5 dados <- dados %>%
6   dplyr::mutate(
7     diagnosis = case_when(
8       diagnosis %in% c("M", "m") ~ "M",
9       diagnosis %in% c("B", "b", "Benign") ~ "B"
10      ),
11     diagnosis =forcats::as_factor(diagnosis) # Convertendo para Fator
12   )
13
14 dados$diagnosis %>% unique()
```

# IMPUTAÇÃO DE DADOS FALTANTES

Remover linhas com NA pode reduzir demais o dataset. Vamos imputar pela mediana de cada grupo (Benigno/Maligno).

```
1 # Imputacao por mediana por grupo
2 dados <- dados %>%
3   dplyr::group_by(diagnosis) %>%
4   dplyr::mutate(across(
5     where(is.numeric),
6     ~replace_na(.x, median(.x, na.rm = TRUE))))
7   ) %>%
8   ungroup()
9
10 dados %>% pillar::glimpse()
```

# ANÁLISE EXPLORATÓRIA

## SUMARIZANDO OS DADOS

Agora que tratamos os problemas, podemos usar a função `summary` (sumário) para exibir algumas medidas descritivas. Aqui selecionamos as colunas que terminam com `_mean`.

```
1 dados %>%
2   dplyr::select(ends_with("_mean")) %>%
3   summary()
```

# O DESAFIO DA VISUALIZAÇÃO EM MASSA

Façamos o histograma do raio médio, `radius_mean`.

```
1 dados %>%
2   ggplot(aes(x = smoothness_mean)) +
3   geom_histogram(bins = nclass.Sturges(dados$radius_mean)) +
4   labs(
5     x = "Raio médio",
6     y = "Frequência"
7   ) +
8   theme_classic()
```

# O DESAFIO DA VISUALIZAÇÃO EM MASSA

**Cenário:** Temos 30 colunas numéricas de medidas (radius\_mean, texture\_mean, etc).

**O Problema:** O ggplot2 espera *uma* coluna específica para mapear no eixo X. Se quisermos ver a distribuição de todas as variáveis, teríamos que copiar e colar o código do gráfico 30 vezes:

```
1 # Inviavel fazer isso 30 vezes:  
2 ggplot(dados, aes(x = radius_mean)) + geom_histogram()  
3 ggplot(dados, aes(x = texture_mean)) + geom_histogram()  
4 # ...  
5
```

**A Solução:** Transformar a estrutura dos dados de **Wide** (Largo) para **Long** (Longo), transformando os *nomes das colunas* em uma variável categórica.

# TRANSFORMAÇÃO: WIDE → LONG

Usamos a função `pivot_longer` do pacote `tidyverse`. Ela "empilha" as colunas selecionadas.

- **cols:** Quais colunas vamos empilhar?
- **names\_to:** Nome da nova coluna que guardará os "títulos".
- **values\_to:** Nome da nova coluna que guardará os números.

```
1 # Criando um dataset longo
2 dados_longos <- dados %>%
3   tidyr::pivot_longer(
4     cols = matches("_*(mean|se|worst)*"),
5     names_to = "variavel",
6     values_to = "valor"
7   )
8
9 # Resultado: O dataset fica mais estreito, mas muito mais comprido!
```

# VISUALIZAÇÃO EM MASSA: PRIMEIRA TENTATIVA

Vamos usar o `facet_wrap` para criar um gráfico por variável. Por padrão, o `ggplot` obriga todos os gráficos a terem os **mesmos eixos**.

```
1 dados_longos %>%
2   dplyr::filter(variavel %>% grepl("_mean", x = .)) %>%
3   ggplot() +
4   geom_histogram(
5     aes(x = valor),
6     bins = nclass.Sturges(dados$radius_mean), # regra de Sturges
7     color = "black", fill = "white"
8   ) +
9   facet_wrap(~variavel) + # 0 padrao e scales = "fixed"
10  theme_classic()
```

## O Problema:

- `texture_mean` vai até 2500.
- `concavity_mean` vai até 0.08.
- Resultado: As variáveis pequenas ficam "esmagadas" em uma única barra à esquerda.

# AJUSTANDO AS ESCALAS (SCALES FREE)

Para corrigir, usamos o argumento scales = "free". Isso libera cada "faceta" para ter seus próprios limites de eixo X e Y.

```
1 dados_longos %>%
2   dplyr::filter(grepl("_mean", variavel)) %>%
3   ggplot() +
4   geom_histogram(
5     aes(x = valor),
6     bins = nclass.Sturges(dados$radius_mean),
7     color="black", fill="white"
8   ) +
9   # A correcao magica:
10  facet_wrap(~variavel, scales = "free") +
11  theme_classic()
```

**Resultado:** Agora conseguimos ver a forma da distribuição de cada variável, independentemente da sua magnitude.

# ANÁLISE DE DISTRIBUIÇÃO E OUTLIERS

Antes de separar por diagnósticos, vamos olhar a distribuição global de cada medida usando **Boxplots**.

**Objetivo:** Identificar a variabilidade total e a presença de valores extremos (*outliers*).

```
1 dados_longos %>%
 2   dplyr::filter(grepl("_mean", variavel)) %>%
 3   ggplot() +
 4   geom_boxplot(aes(y = valor)) +
 5   facet_wrap(~variavel, scales = "free") +
 6   theme_classic()
```

*Nota: Vemos muitos outliers (pontos pretos), mas não sabemos se são casos graves ou apenas ruído.*

# IDENTIFICANDO PADRÕES (SEPARAÇÃO POR GRUPO)

**Motivação:** Para construir um classificador, precisamos de variáveis que se comportem de maneira **diferente** entre Benignos e Malignos.

Vamos separar os boxplots por cor (diagnosis). Se as caixas se separarem verticalmente, a variável pode ser uma boa preditora.

```
1 dados_longos %>%
  2   dplyr::filter(grepl("_mean", variavel)) %>%
  3   ggplot() +
  4   # Adicionamos cor baseados no diagnóstico
  5   geom_boxplot(aes(y = valor, color = diagnosis)) +
  6   facet_wrap(~variavel, scales = "free") +
  7   labs(color = "Diagnóstico") +
  8   theme_classic()
```

**Conclusão Visual:** Variáveis como radius\_mean e fractal\_dim\_mean separam bem os grupos.

# VISUALIZAÇÃO MULTIVARIADA (PAIRS PLOT)

Às vezes, queremos ver a relação entre várias variáveis ao mesmo tempo (dispersão) e suas distribuições. O pacote GGally automatiza isso.

*Nota: Selecionamos apenas algumas colunas, pois fazer isso com todas as 30 travaria o plot.*

```
1 dados %>%
2   # Selecionando variaveis de interesse
3   dplyr::select(diagnosis, radius_mean, texture_mean, fractal_dim_mean) %>%
4
5   # ggpairs cria uma matriz de graficos
6   GGally::ggpairs(
7     aes(color = diagnosis, alpha = 0.5) # Cores por grupo
8   ) +
9   theme_classic()
```

# CORRELAÇÃO NÃO-PARAMÉTRICA (SPEARMAN)

O gráfico anterior sugere relações não-lineares. Vamos quantificar a associação entre o **Raio** e a **Textura** usando correlação de posto (Spearman), que é mais robusta a outliers que a de Pearson.

Usaremos o pacote `corrr`, que devolve um *tibble* (tabela) em vez de uma matriz, facilitando a leitura.

```
1 # Calculando a correlacao de Spearman
2 dados %>%
3   dplyr::select(radius_mean, texture_mean) %>%
4   corrr::correlate(method = "spearman", quiet = TRUE)
```

**Por que Spearman?** Porque nossas variáveis podem não ser linearmente associadas.

# MODELAGEM: ÁRVORE DE DECISÃO

# DIVISÃO TREINO E TESTE

Preparando os dados para o aprendizado de máquina.

```
1 set.seed(20251125) # Reprodutibilidade
2
3 # Divisão 80/20
4 split <- rsample::initial_split(dados, prop = 0.8)
5 treino <- rsample::training(split)
6 teste <- rsample::testing(split)
```

# TREINANDO O MODELO

Usando a engine rpart (Recursive Partitioning).

```
1 # Especificacao do modelo
2 tree_spec <- parsnip::decision_tree(
3   mode = "classification",
4   engine = "rpart",
5   tree_depth = 4 # Profundidade maxima
6 )
7
8 # Ajuste do modelo
9 tree_fit <- tree_spec %>%
10   fit(diagnosis ~ ., data = treino)
```

# VISUALIZANDO A ÁRVORE

Interpretando as regras criadas pelo modelo.

```
1 # Plot da arvore  
2 rpart.rpart::rpart.plot(tree_fit$fit, type = 4, roundint = FALSE)
```

# AVALIAÇÃO DE DESEMPENHO

# ENTENDENDO A MATRIZ DE CONFUSÃO

A base de toda avaliação em classificação é comparar o **Preditivo vs Real**.

|          |             | Realidade   |   |
|----------|-------------|---|---|
|          |             | Maligno (M)   | Benigno (B)   |
| Predição | Maligno (M) | <b>Verdadeiro Positivo</b><br>(O modelo acertou o câncer) | Falso Positivo<br>(Alarme falso)                                |
|          | Benigno (B) | Falso Negativo<br>(Erro grave: câncer não detectado)      | <b>Verdadeiro Negativo</b><br>(O modelo acertou que é saudável) |

- **Verdadeiro Positivo/Negativo (Verde):** O modelo acertou.
- **Falso Negativo (Erro Tipo II):** O erro mais perigoso na medicina (deixar um paciente doente ir para casa).
- **Falso Positivo (Erro Tipo I):** Gera estresse, mas pode ser verificado com exames adicionais.

# MÉTRICAS RESUMO: ACURÁCIA E KAPPA

Olhar a matriz inteira pode ser difícil. Usamos métricas para resumir a qualidade:

## 1. Acurácia (Accuracy)

$$\frac{\text{Total de Acertos}}{\text{Total de Casos}}$$

- **Pró:** Intuitiva.
- **Contra:** Pode enganar em dados desbalanceados (ex: se 99% dos pacientes são saudáveis, chutar "saudável" sempre dá 99% de acurácia).

## 2. Kappa de Cohen

$$\frac{\text{Acurácia Observada} - \text{Esperada ao Acaso}}{1 - \text{Esperada ao Acaso}}$$

- Mede a concordância **além do acaso**.
- 1 = Perfeito; 0 = Equivalente a jogar uma moeda.
- Útil para dados médicos e crucial para psicométricos.

# CALCULANDO NO R (PACOTE YARDSTICK)

O pacote yardstick (parte do tidymodels) padroniza esses cálculos.

```
1 # Gerar Predicoes no conjunto de Teste
2 predicoes <- tree_fit %>%
3   predict(teste) %>%
4   dplyr::pull(.pred_class)
5
6 # Gerar Matriz de Confusao
7 teste %>%
8   dplyr::mutate(predicoes = predicoes) %>%
9   conf_mat(truth = diagnosis, estimate = predicoes)
10
11 # Calcular Metricas (Acuracia e Kappa)
12 metricas <- metric_set(accuracy, kap) # acuracia e \kappa
13 teste %>%
14   dplyr::mutate(predicoes = predicoes) %>%
15   metricas(truth = diagnosis, estimate = predicoes)
```

# CONCLUSÃO

# RECAPITULANDO

1. **Importamos** dados brutos e identificamos erros de estrutura.
2. **Corrigimos** problemas de codificação e estrutura do CSV.
3. **Limpamos** nomes, tipos de dados, strings inconsistentes e valores nulos.
4. **Exploramos** os dados visualmente para entender as relações.
5. **Modelamos** o problema usando uma Árvore de Decisão simples e interpretável.

**Próximos Passos:** Boosting, testar modelos mais robustos (Florestas Aleatórias) e refinar o tratamento de dados.

OBRIGADO!

## Dúvidas?

Link para o material: <https://github.com/nueidris/semac2025>