

## Previsão de Inadimplência com Machine Learning

```
9 # Previsão de inadimplência de clientes com Machine Learning e Power BI
10 #
11 #
12 # Definindo a pasta de trabalho
13 setwd("E:/Cursos/PowerBI_Data_Science/capitulo15")
14 getwd()
15
16 # Definição do problema
17 # Leio o manual em pdf no capítulo 15 do curso como a definição do problema
18
19 # Instalando os pacotes para o projeto
20 # Obs: os pacotes precisam ser instalados apenas uma vez
21
22 install.packages("Amelia")
23 install.packages("caret")
24 install.packages("ggplot2")
25 install.packages("dplyr")
26 install.packages("reshape")
27 install.packages("randomForest")
28 install.packages("e1071")
29
30 # Carregando os pacotes
31
32 library(Amelia)
33 library(ggplot2)
34 library(caret)
35 library(reshape)
36 library(randomForest)
37 library(dplyr)
38 library(e1071)
39
40 # Carregando o dataset
41 # Fonte: https://archive.ics.uci.edu/ml/datasets/default+of+credit+card+clients
42 dados_clientes <- read.csv("dados/dataset.csv")
43
44 # Visualizando os dados e sua estrutura
45
46 view(dados_clientes)
47 dim(dados_clientes)
48 str(dados_clientes)
49 summary(dados_clientes)
50
51 ##### Análise exploratória, limpeza e transformação #####
52
53 # Removendo a primeira coluna ID
54 dados_clientes$ID <- NULL
55 dim(dados_clientes)
56 view(dados_clientes)
57
58 # Renomeando a coluna de classe
```

### Visão Geral

Este script foi desenvolvido como parte do curso "Microsoft Power BI Para Data Science, Versão 2.0" da Data Science Academy.

O objetivo é prever a inadimplência de clientes utilizando técnicas de Machine Learning integradas ao Power BI.

### Dependências

O script requer os seguintes pacotes R:

```
```R
```

```
install.packages(c("Amelia", "caret", "ggplot2", "dplyr", "reshape", "randomForest",  
"e1071", "devtools"))
```

```
devtools::install_github('cran/DMwR')
```

```
```
```

### Fluxo de Trabalho

## 1. Pré-processamento dos Dados

-Carregamento: Dados de crédito de clientes (fonte: UCI Machine Learning Repository)

-Limpeza:

- Remoção de coluna ID
- Renomeação de colunas para melhor legibilidade
- Tratamento de valores ausentes

-Transformação:

- Conversão de variáveis categóricas (gênero, escolaridade, estado civil)
- Discretização de idade em faixas etárias
- Conversão de variáveis de pagamento para fatores

## 2. Análise Exploratória

- Visualização da distribuição de inadimplentes vs. não-inadimplentes
- Análise de proporções entre classes
- Geração de gráficos exploratórios com ggplot2

## 3. Modelagem Preditiva

-Divisão dos Dados: 75% treino / 25% teste (amostragem estratificada)

-Versões do Modelo:

- 1.Modelo Inicial: Random Forest com todas as variáveis
- 2.Modelo Balanceado: Aplicação de SMOTE para tratar desbalanceamento de classes
- 3.Modelo Otimizado: Utilizando apenas variáveis mais importantes

-Métricas de Avaliação:

- Matriz de confusão
- Precision, Recall e F1-Score

#### 4. Implementação Final

- Seleção do melhor modelo (versão 3)
- Serialização do modelo para uso em produção
- Exemplo de aplicação com novos dados de clientes

#### Estrutura de Arquivos

...

/projeto

/dados

dataset.csv      Dados brutos dos clientes

/modelo

modelo\_v3.rds      Modelo serializado

script.R      Este script de análise

...

#### Como Executar

1. Definir diretório de trabalho (linha 8)
2. Instalar pacotes necessários (apenas na primeira execução)
3. Executar o script sequencialmente

#### Resultados Principais

- Modelo final alcançou F1-Score de [valor a ser preenchido após execução]

- Variáveis mais importantes identificadas:
  1. PAY\_0 (status de pagamento mais recente)
  2. PAY\_2
  3. PAY\_3
  4. PAY\_AMT1 (valor do pagamento mais recente)

### Integração com Power BI

O modelo serializado (modelo\_v3.rds) pode ser carregado no Power BI para:

- Criação de dashboards interativos
- Previsões em tempo real
- Análise de risco de crédito

### Observações

- Para reproduzir exatamente os mesmos resultados, manter o `set.seed(12345)`
- O dataset original contém informações sensíveis - tratar com confidencialidade
- Versão dos dados utilizada: 17/04/2021

### Referências

- Documentação original do curso: Capítulo 15
- Dataset: Yeh, I. C., & Lien, C. H. (2009). The comparisons of data mining techniques for the predictive accuracy of probability of default of credit card clients. Expert Systems with Applications.