



Projeto

Machine Learning

Risco de Crédito

Framework



1 Definir o Problema

HOME CREDIT

A Home Credit disponibilizou **246mil** registros de aplicações de empréstimos. Cada aplicação com **122** características, entre elas:

- Valor do empréstimo
- Valor das parcelas
- Patrimônio do aplicante
- Em qual horário o empréstimo foi solicitado
- Há quanto tempo o aplicante trabalha
- + 6 Datasets de informações de aplicações anteriores
- + outras 118 características

O **objetivo** é usar esses dados para construir um Modelo de Machine Learning para prever se um aplicante terá dificuldades em pagar um empréstimo. Cascadeando:

- O primeiro foco do Modelo é distinguir entre “tipos” de aplicantes. No nosso caso: bons pagadores de maus pagadores. Em termos de métricas, buscamos um **ROC-AUC** alto.
- Como esse é um projeto de aprendizado proposto no Degree de Data Science da Ada/Let’s Code, o segundo foco está em experimentar e testar diversos aspectos de **Machine Learning**.

Desafios e Soluções

Desafio

O tamanho do armazenamento dos datasets ultrapassa 1GB. Limitando, por exemplo, o armazenamento no Github.

Solução

Usar um template de pastas e desenvolver uma automação usando 'make'

Desafio

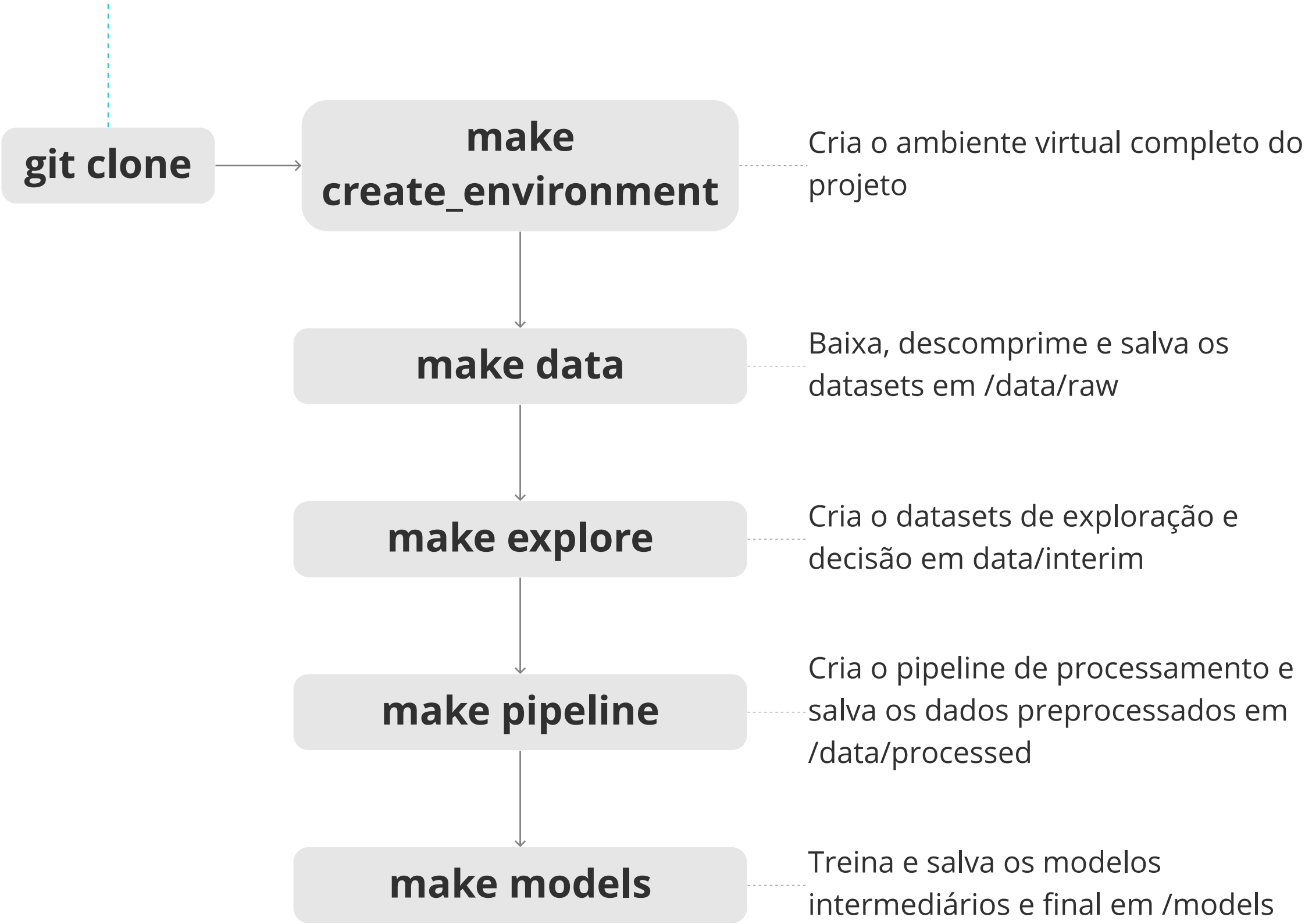
Múltiplos datasets com múltiplos tipos de relações e tipos de características

Solução

Fazer um "Pre-overview" de todos os datasets para guiar a Análise Exploratória dos Dados (EDA)

Um dos focos do projeto foi **reprodutibilidade** e **automação**.

Para reproduzir todos os passos do projeto:



Exploração e Decisões

Valores Nulos

Decisão
Rejeitar colunas
> 30% de NaNs

Tipo de Dado

Categórico

Decisão
Transformar com
OrdinalEncoder
da biblioteca
feature-engine

Numérico

Ruído
Outliers (Por
IQR e Z-score)

Decisão
Rejeitar registros
com Z-score > +-3

Distribuição
Testar normalidade,
calcular Skewness
& Kurtosis e
correlação

Decisão
Agrupar colunas
relacionadas

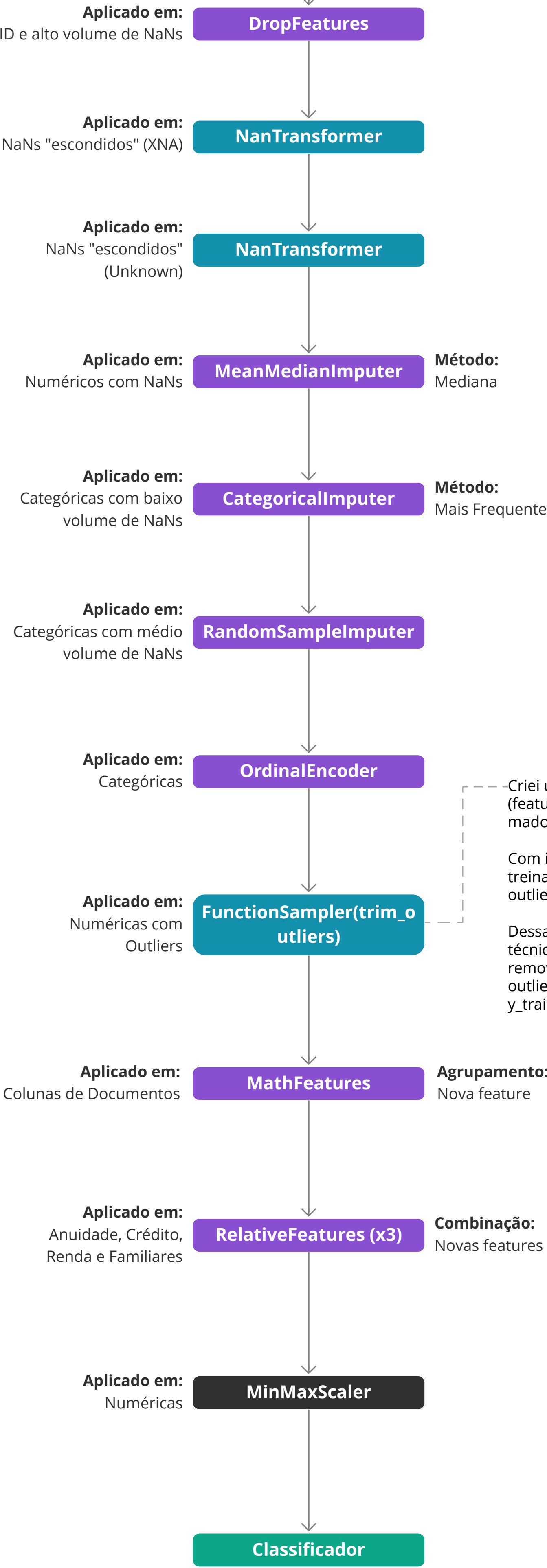
Decisão
Combinar colunas
para criar novas
features

4

Preparar os Dados

-  scikit-learn
-  Feature-engine
-  Customizado

Pipeline



5

Selecionar o Modelo

Modelos Testados

Métricas calculadas com validação cruzada

	Tempo de Treino	ROC-AUC
DummyClassifier	91s	50.0%
RegressionClassifier	105s	73.9%
KNeighborsClassifier	105s	58.1%
DecisionTreeClassifier	114s	53.7%
LinearSVC	149s	74.0%
RandomForestClassifier	185s	71.0%
AdaBoostClassifier	133s	74.4%
GradientBoostingClassifier	217s	75.1%
XGBClassifier	208s	74.7%
CatBoostClassifier	204s	75.6%
LGBMClassifier	92s	75.5%
Chosen Models for Hyperparameter tuning and combination		

6

Ajustar o Modelo

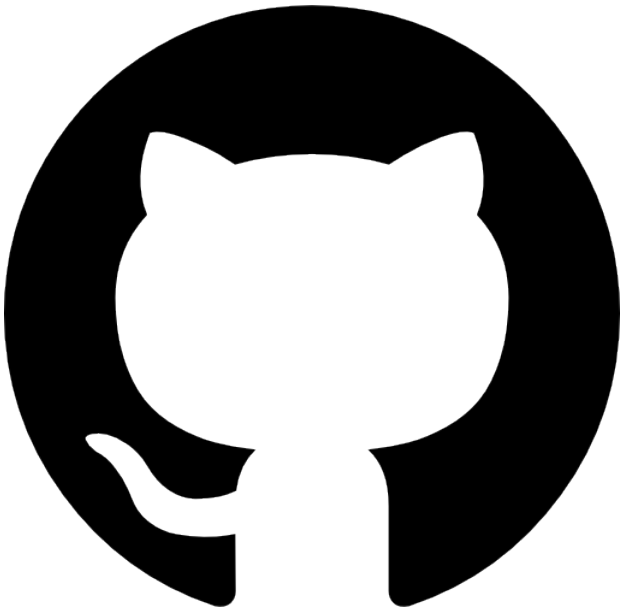
Modelo Final



Modelo Final

ROC-AUC do set de validação

75.58%



<https://github.com/ewerthonk>