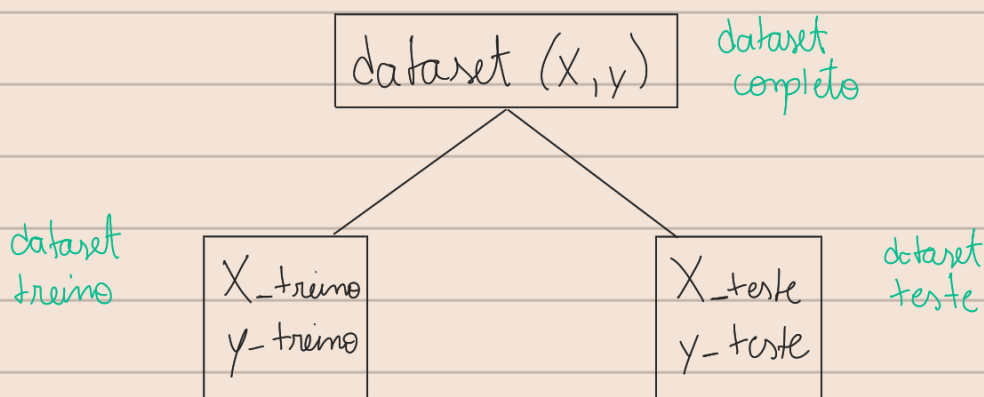


Etapas de um projeto de Machine Learning

Etapa 1: análise exploratória preliminar

- foco é conhecer a natureza das features individualmente
 - estatística descritiva
 - visualizações
- descobrir "anomalias" nos dados: outliers, erros, artefatos (naturais e etc)
"buracos" e etc
- Não investigar interrelações

Etapa 2: repartição treino e teste



Etapa 3: análise exploratória posterior

- Foco: analisar interrelações
 - correlações
 - visualizações conjuntas

As Etapas 1, 2, 3 correspondem às atividades "Business Understanding" e "Data understanding". Ao término das etapas 1-3, temos decisões sobre: objetivos, limpeza, (filtroagem) e transformações.

Etapa 4: modelagem

- Filtroagem: eliminar
 - não entra nas "Pipelines" do sklearn
- Transformação e modelagem
 - Pipeline do sklearn
- Escolha do melhor modelo

A Etapa 4 corresponde às atividades "Data preparation" e "modeling" do CRISP-DM.

Etapa 5: certificação

- Medir desempenho do modelo escolhido
- Decidir se pode entrar em **produção** / **"Evaluation"** do CRISP-DM

≡

Corresponde ao

Etapa 6: produção

- treino final e avanço de modo treinado
- para responsabilidade para equipes **MLOPS / DEVOPS** / **"Deployment"** do CRISP-DM

"Evaluation"

≡

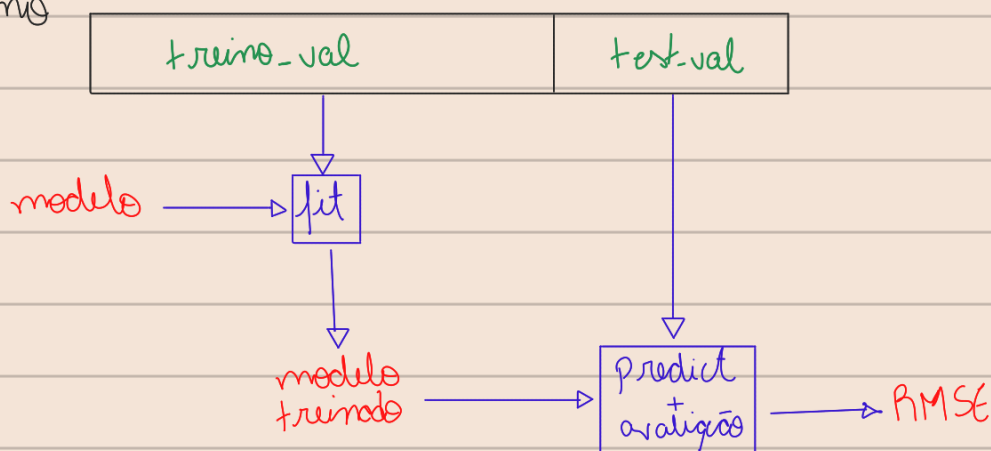
Corresponde ao

CRISP-DM	Etapas	Nível
Business Understanding	Etapa 1	pré Nível 2
Data Understanding	Etapa 2	
	Etapa 3	
Data prep / Modeling	Etapa 4	nível 2: escolha mod
Evaluation	Etapa 5	nível 1: certificação
Deployment	Etapa 6	nível 0: deploy

Escolha do modelo

Estratégia: separação treino_val / teste_val

treino



* em um pipeline no fit, os transformers aprendem e transformam caso não seja o último, caso Pipe tenha um modelo os transformers vão aplicar fit-transform. Caso tenha um modelo podemos dar fit-predict, caso não apenas fit-transform

* Toda busca por modelos tem que ter
Dummy
RAW Reg Lin } benchmark de "basta"

Estratégia: validação cruzada

treino	1	2	3	4	5	$cv = 5$ cross validation
--------	---	---	---	---	---	---------------------------------

round	treino_val	teste_val	} 5 RMSE
1	2 3 4 5	1	
2	3 4 5 1	2	
3	4 5 1 2	3	
4	5 1 2 3	4	
5	1 2 3 4	5	