





Análise de Dados em Python para Não Programadores

Patrícia Aguiar Moreira

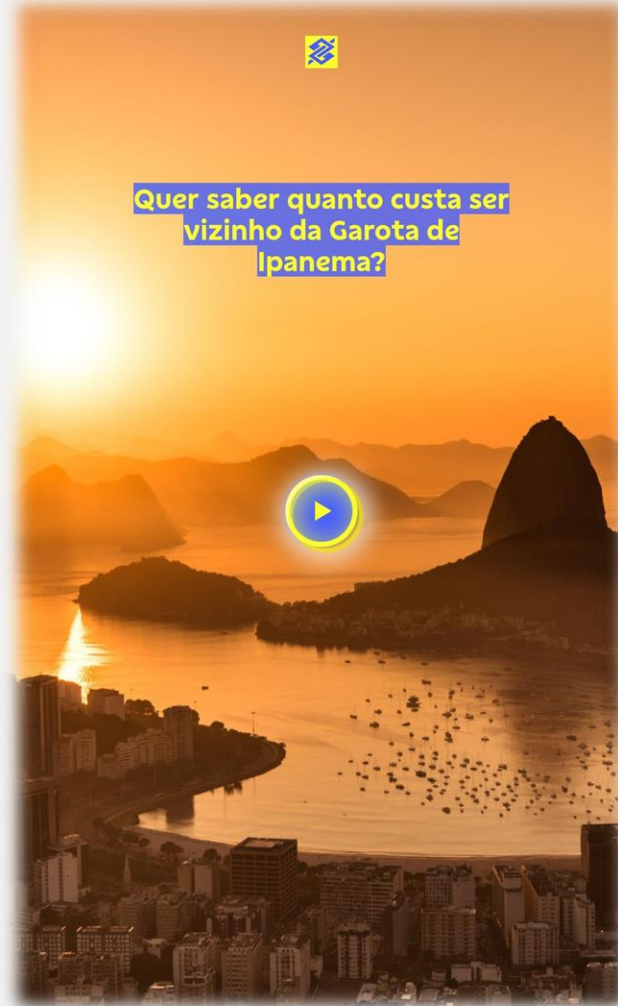
Estudo de Caso: Valoração de Imóveis



Quanto custa o imóvel dos seus sonhos?



Já se imaginou morando
cercado pela natureza e quer
saber quanto custa?



Quer saber quanto custa ser
vizinho da Garota de
Ipanema?

O que vamos ver?

- Acesso a base de dados
- Insights a partir de visualizações
- Desenvolvimento de um modelo preditivo
- Simulação de uso do modelo

Onde vamos realizar as análises?

Linguagem de Programação e Ambiente



[https://colab.research.google.com/github/
mpatricia/workshop/blob/main/Notebook.
ipynb](https://colab.research.google.com/github/mpatricia/workshop/blob/main/Notebook.ipynb)



ANALYTICSLABB

<https://labblite.bb.com.br/>

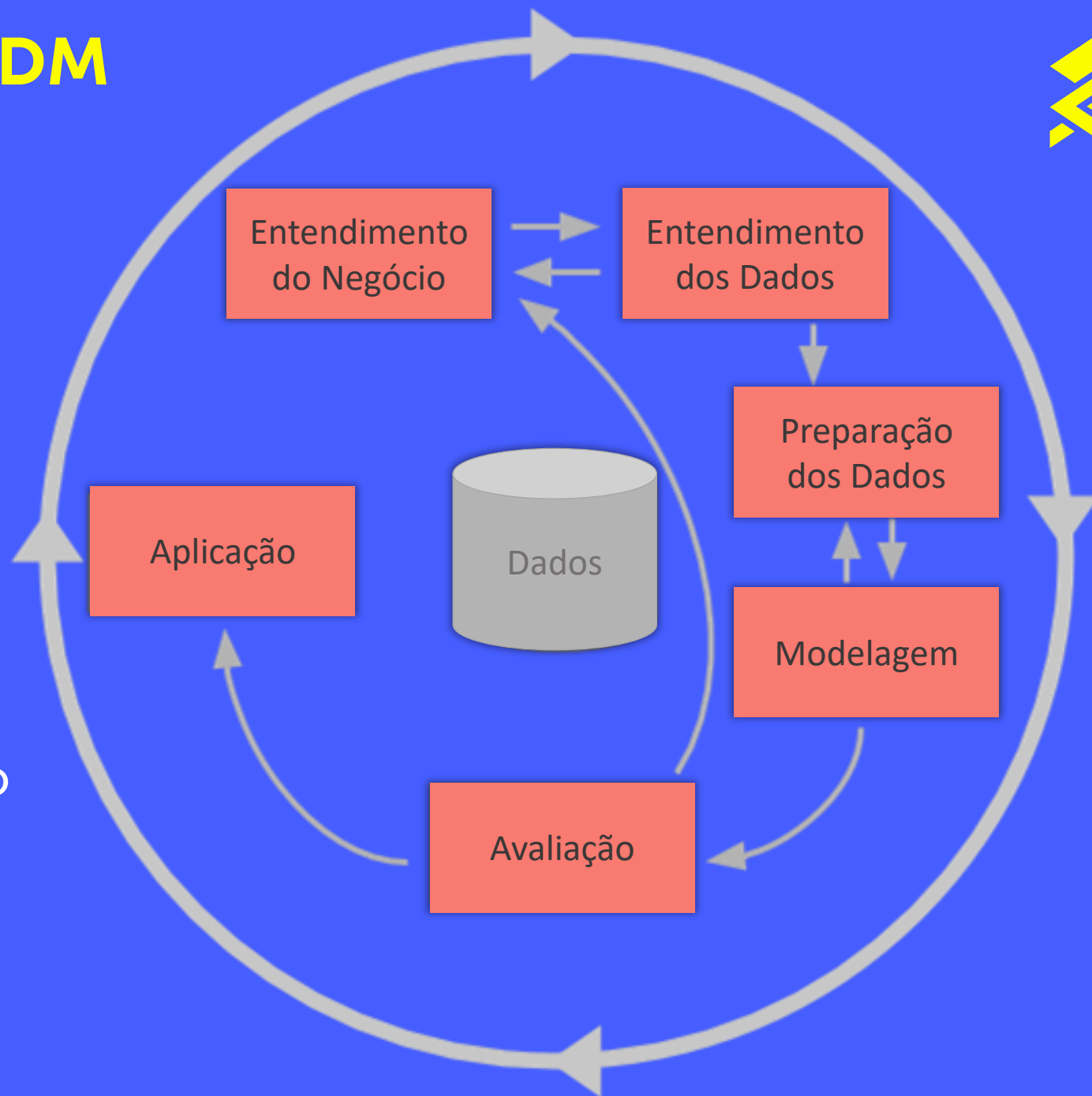


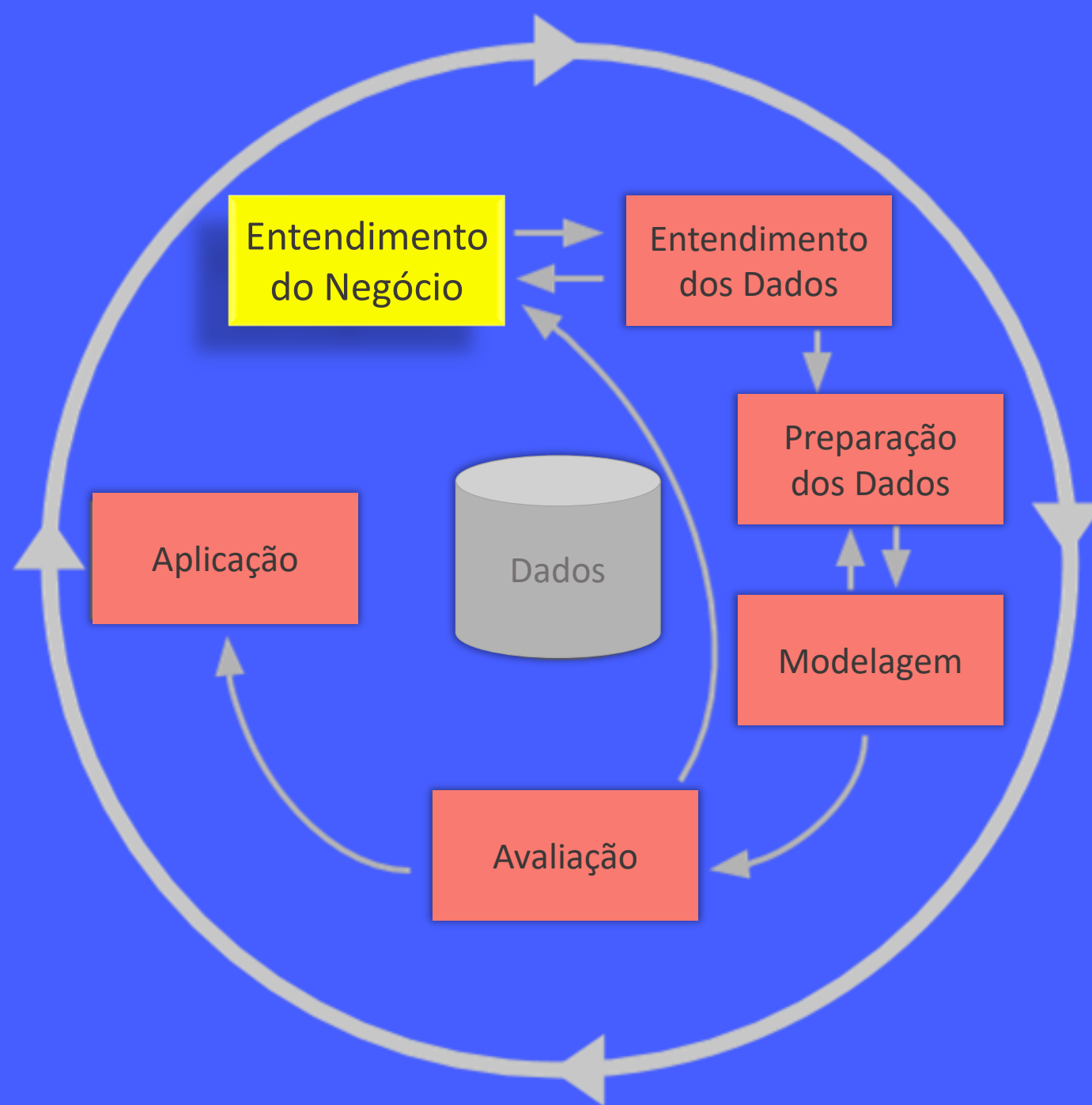
Metodologia CRISP-DM

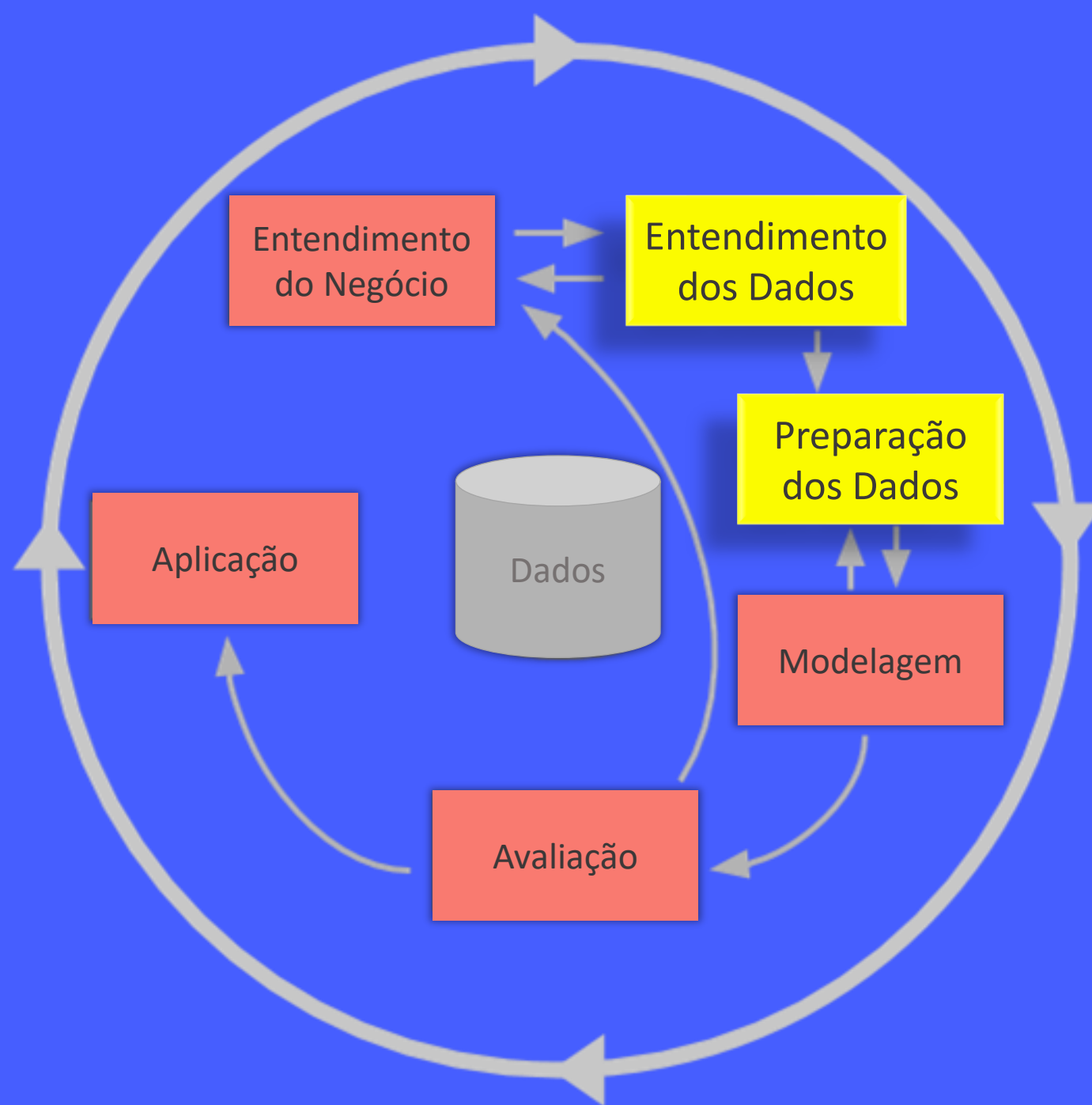


CRoss
Industry
Standard
Process for
Data **M**ining

É um modelo de processo
para análise de dados







Entendimento e Preparação dos Dados



Fontes de Origem → Dados Estruturados

Atributos/Características
nas colunas

ID	Tem vista para orla?	Índice de Qualidade	Número de Quartos	Preço
1	S	3	2	500.000,00
2	N	7	1	600.000,00
3	N	3	4	2.000.000,00
4	S	11	6	5.000.000,00
5	N	5	2	700.000,00

Observações/Exemplos
nas linhas

Identificador
da observação

Inputs
Variáveis Explicativas

Output
Variável Resposta/Target

Entendimento e Preparação dos Dados



Qual o tratamento adequado?

Excluir atributos
não relevantes

Os algoritmos não interpretam
textos (variáveis qualitativas),
é preciso codificá-los em números.

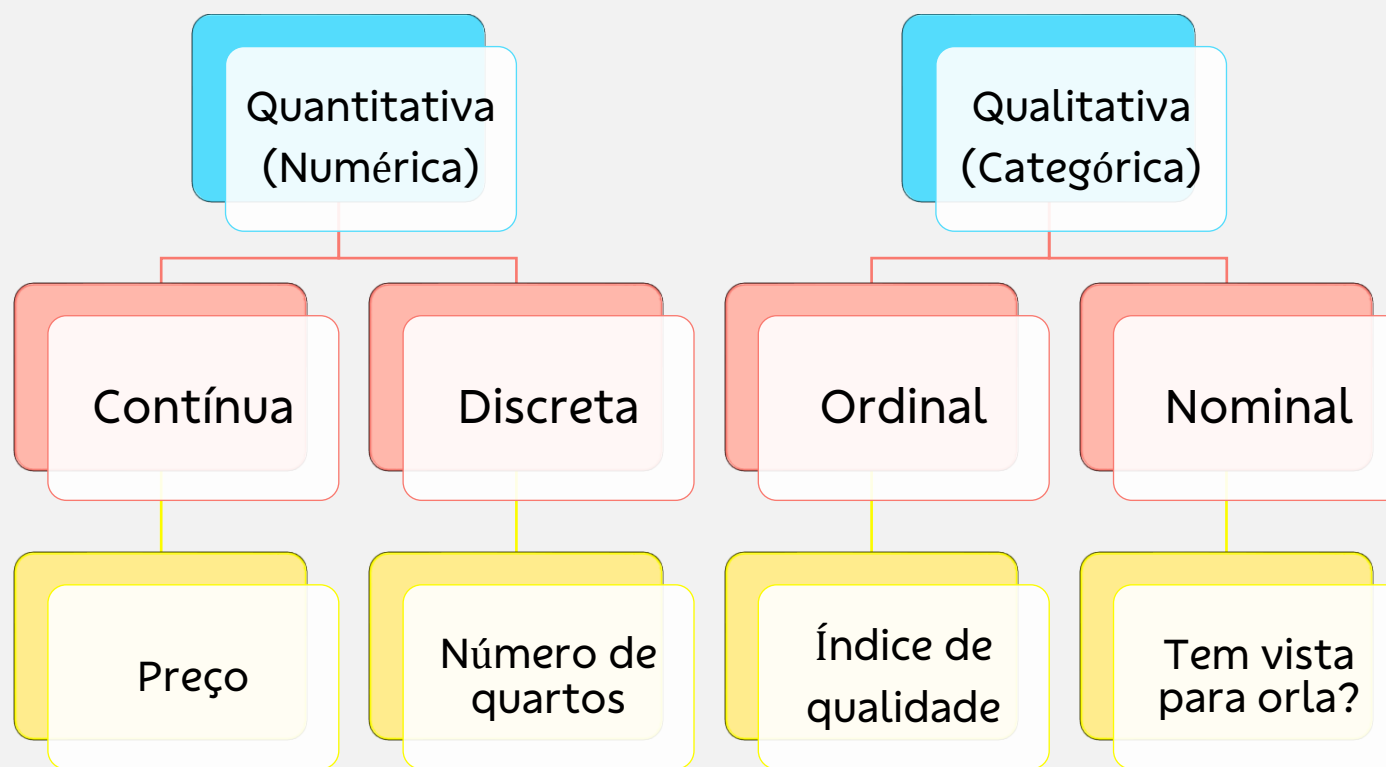
ID	Tem vista para orla?	Índice de Qualidade	Número de Quartos	Preço
1	S	3	2	500.000,00
2	N	7	1	600.000,00
3	N	3	4	2.000.000,00
4		11	6	5.000.000,00
5	N	5	2	700.000,00

Valores ausentes são informação,
se foram observados nos dados
treino, é provável observar nos
novos dados. Como tratá-los?

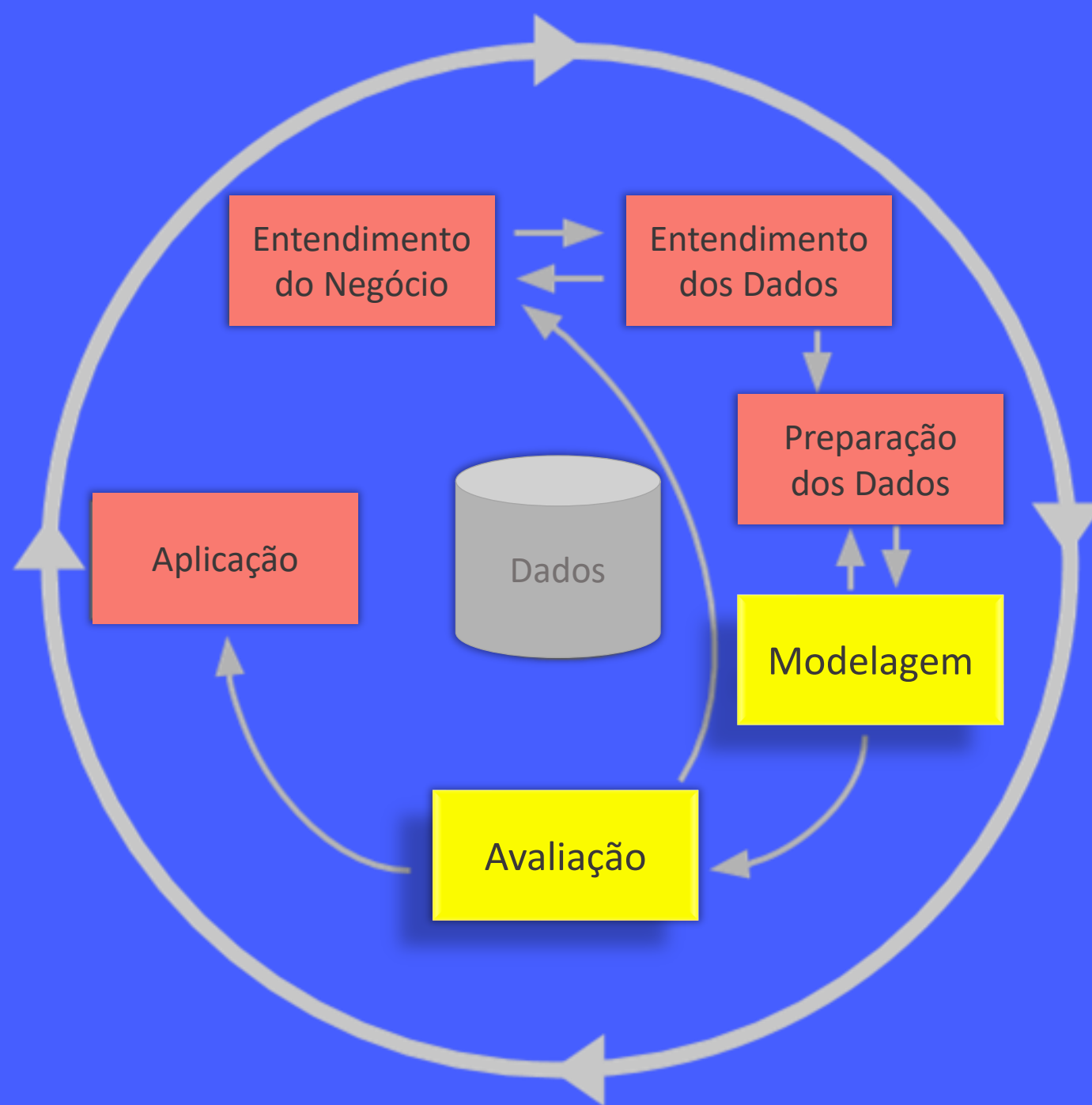
Entendimento e Preparação dos Dados



Qual o tipo das variáveis?

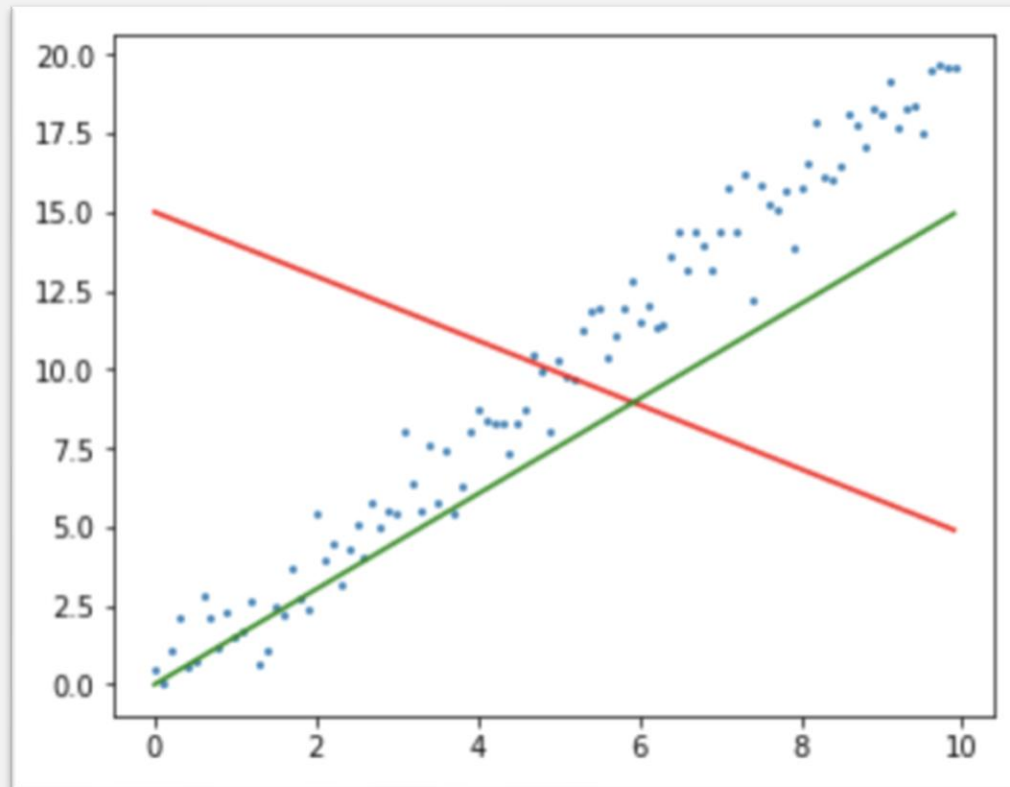






Modelagem

Como prever y a partir de x?



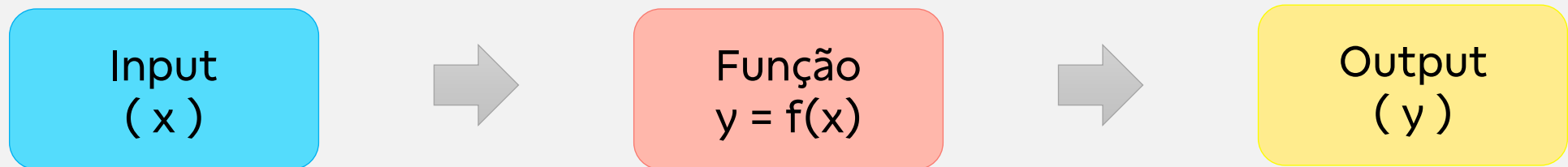
Qual reta melhor descreve a relação entre x e y?

$$y = 15 - 1x$$

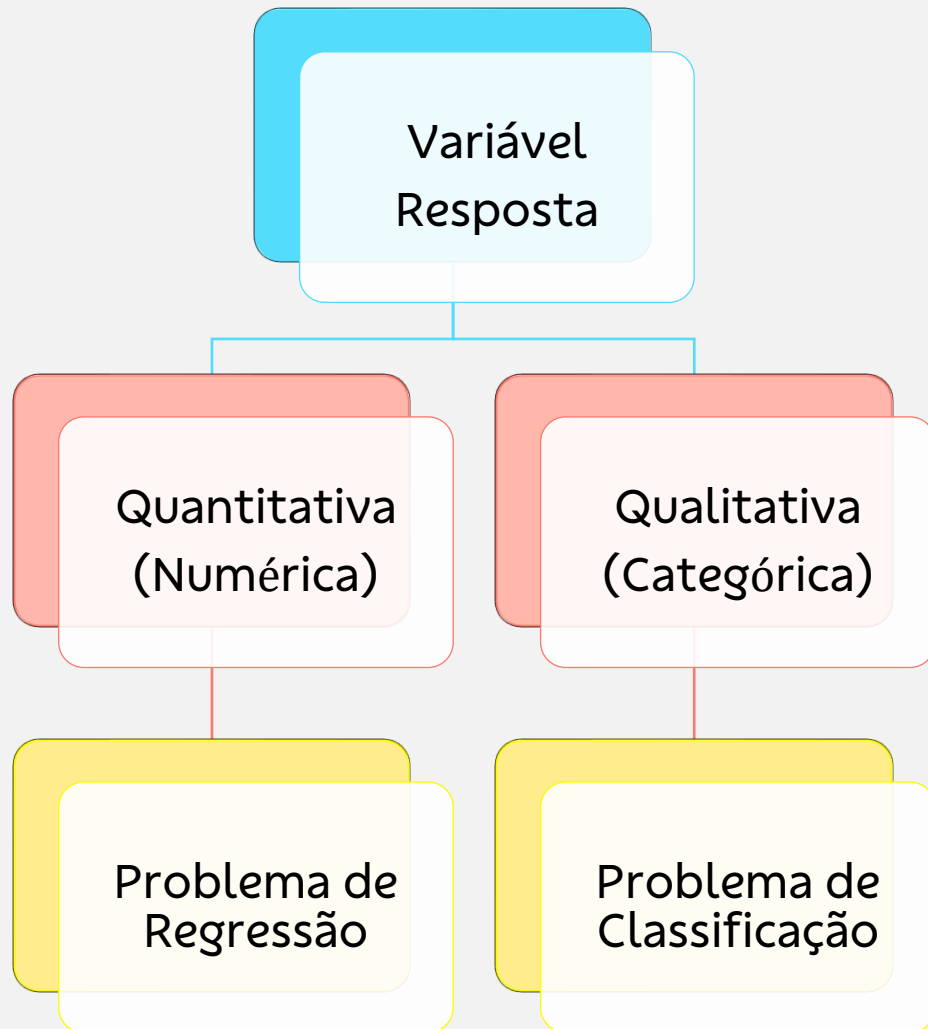
$$y = 1,5x$$



Temos um fenômeno em estudo onde cada valor \mathbf{y} pode ser encontrado por uma função desconhecida $\mathbf{y = f(x)}$:

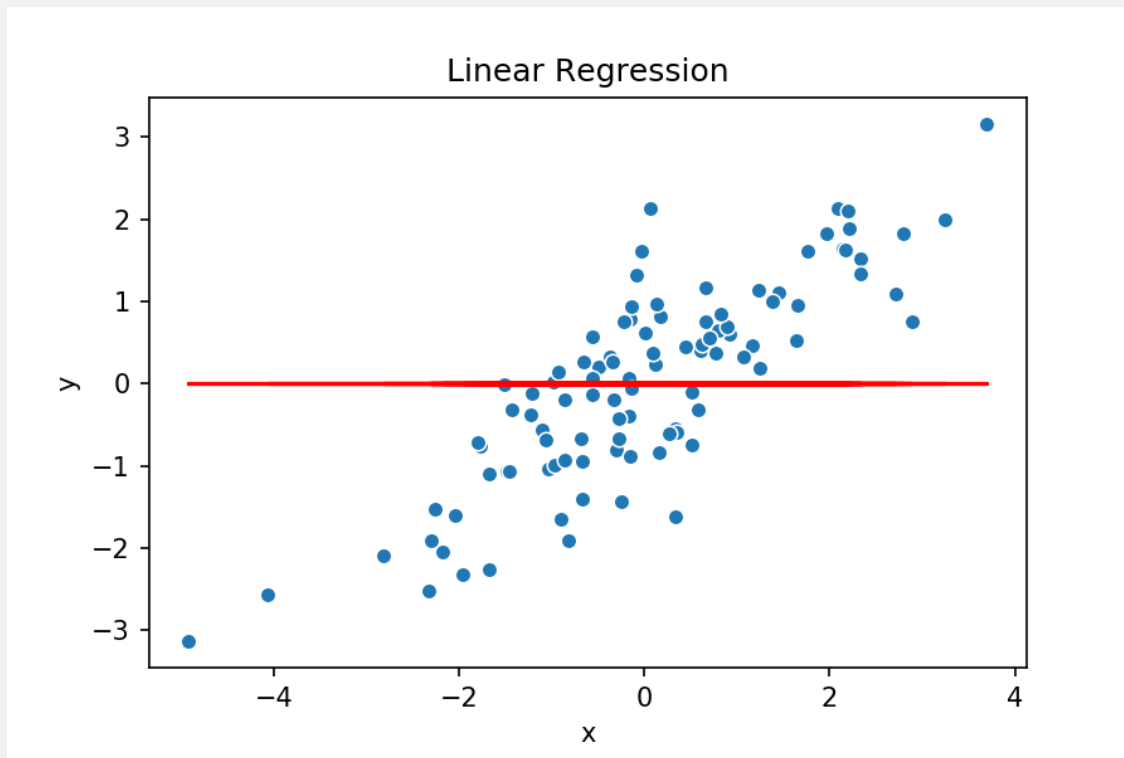


O objetivo do treinamento é estimar uma função h (hipótese) que se aproxime da função verdadeira f .



Estimadores para Regressão:

- Regressão Linear
- KNN
- Árvore de Decisão
- Random Forest
- Gradient Boosting



Regressor que faz previsões a partir de uma **combinação linear** dos preditores:

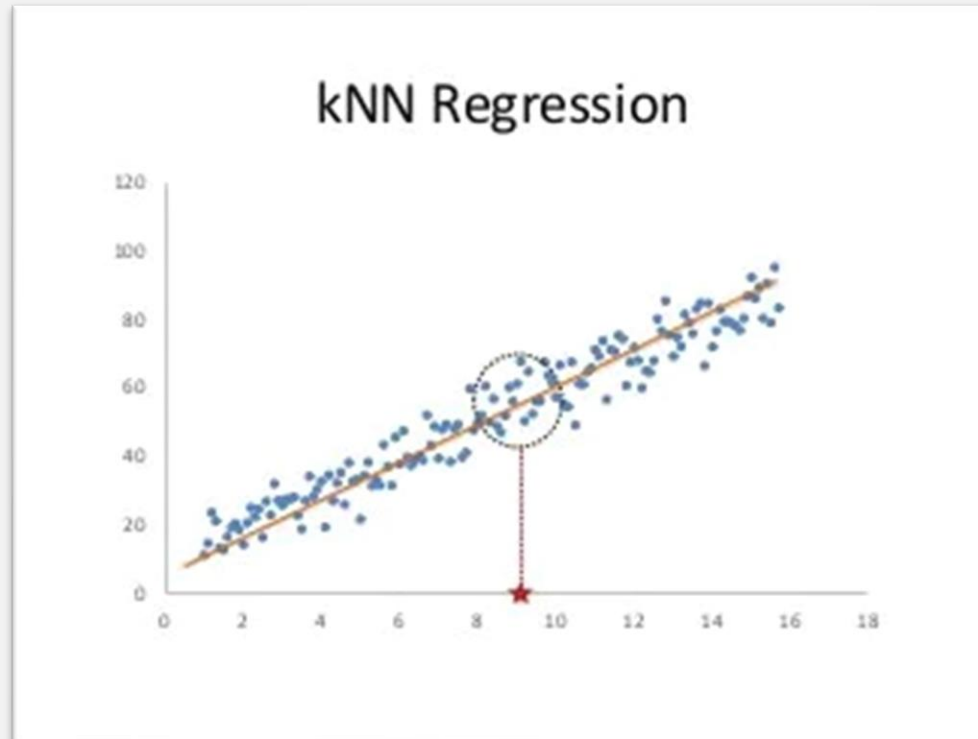
$$y = \alpha + \beta X$$

Modelagem

KNN - K Nearest Neighbors

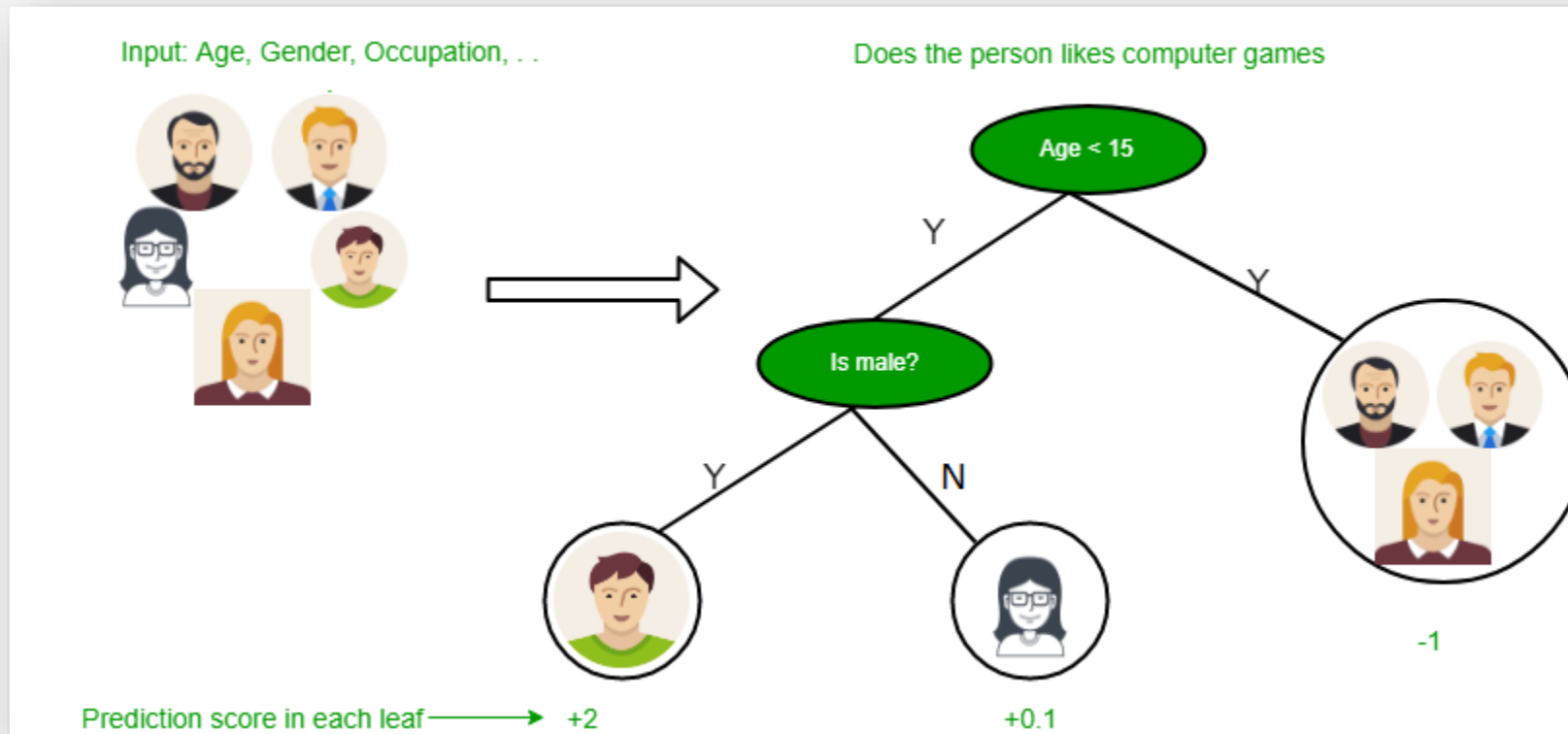


Regressor que faz previsões baseadas nos K vizinhos mais próximos.



Modelagem

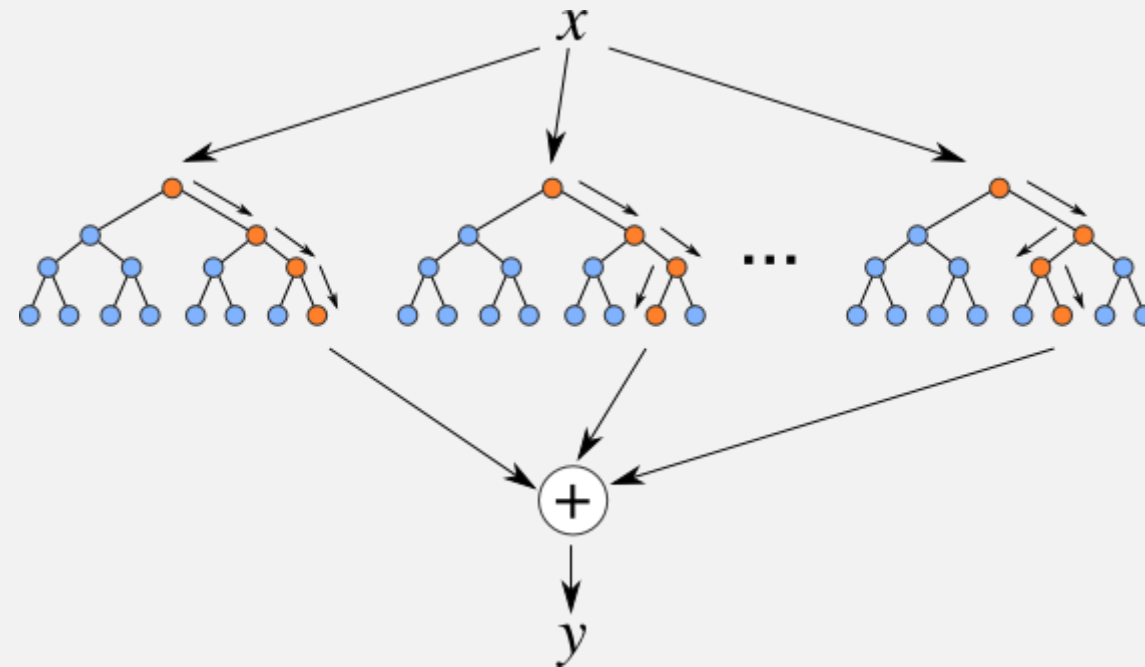
Árvore de Decisão



Regressor que, a cada passo, divide as observações em grupos de acordo com a característica que melhor separa os valores da variável target.

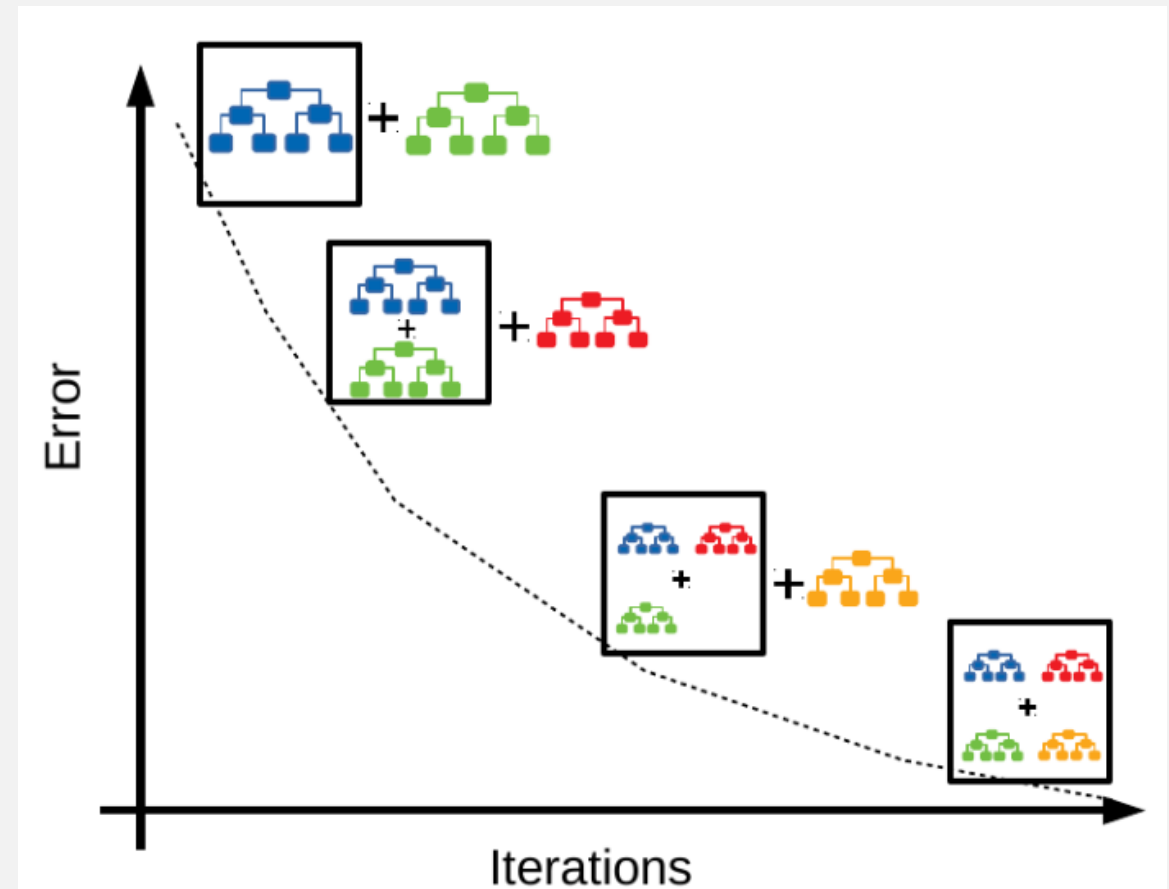


Regressor que consulta vários regressores “fracos” treinados em amostras de variáveis e observações, e faz a previsão com base em uma média das previsões dos regressores “fracos”.



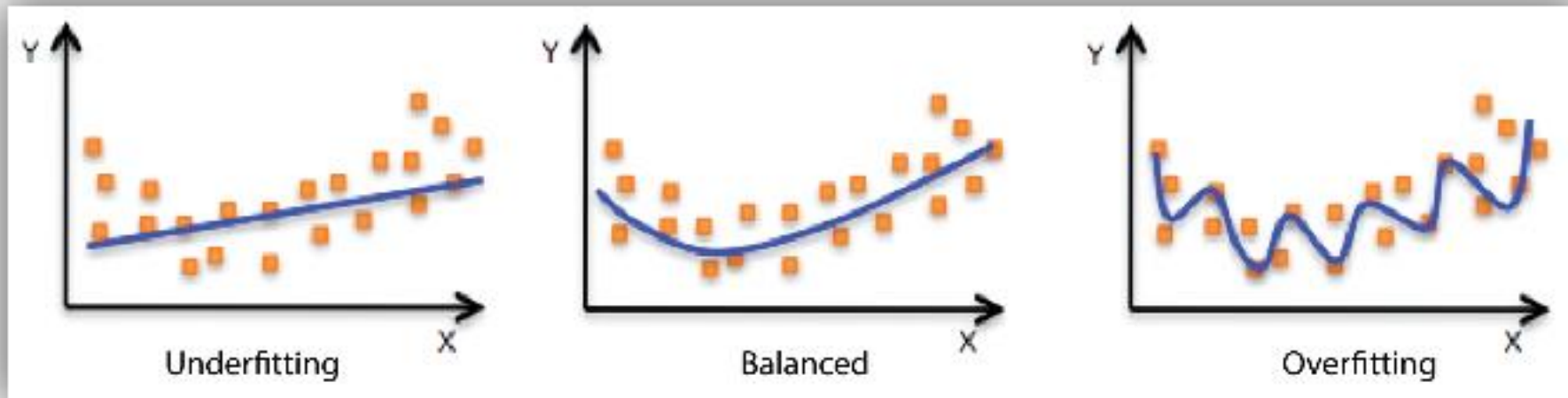


Regressor que faz uso de vários regressores “fracos” em sequência, sendo que cada regressor adicionado é treinado para corrigir o erro dos anteriores.



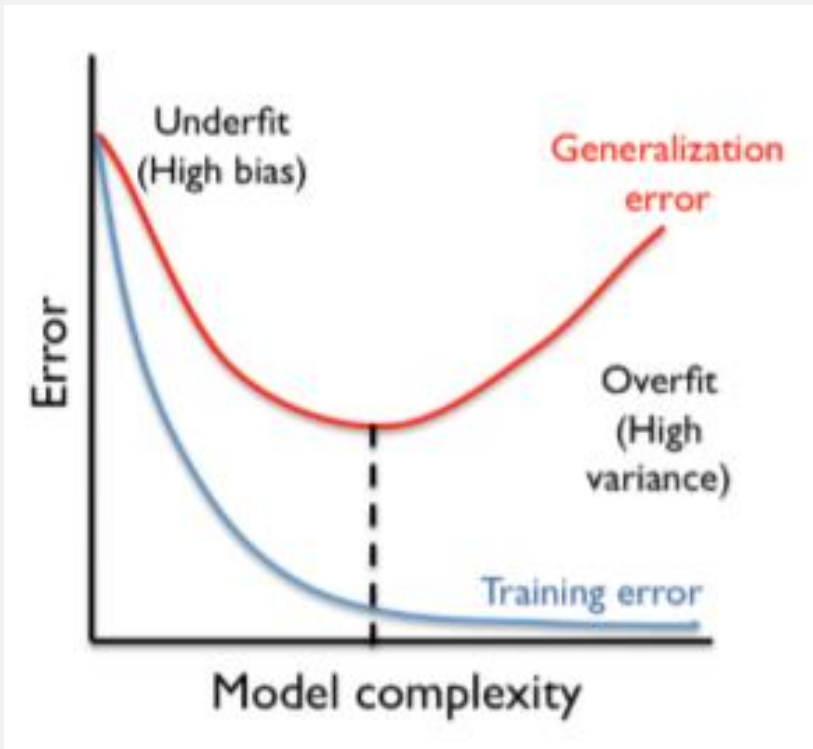


O estimador vai treinar um modelo que minimiza o erro das previsões nos dados de treino, mas será que é capaz de generalizar para novos casos?



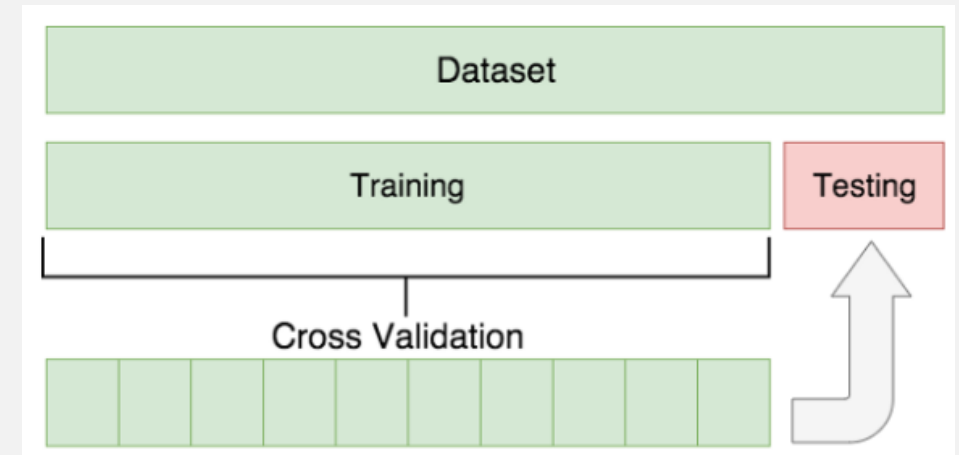


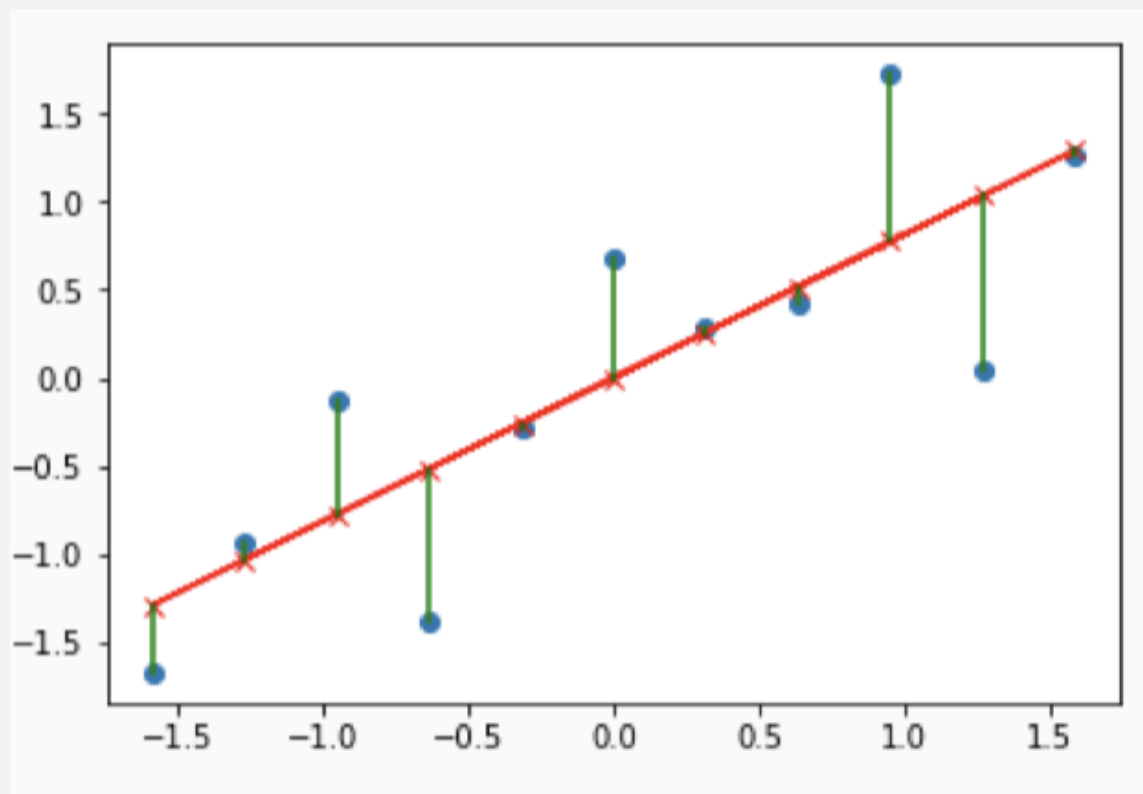
Underfitting/Overfitting



Solução

Treinar algoritmo de aprendizado em parte dos dados e reservar a outra parte para mensurar a métrica de performance.





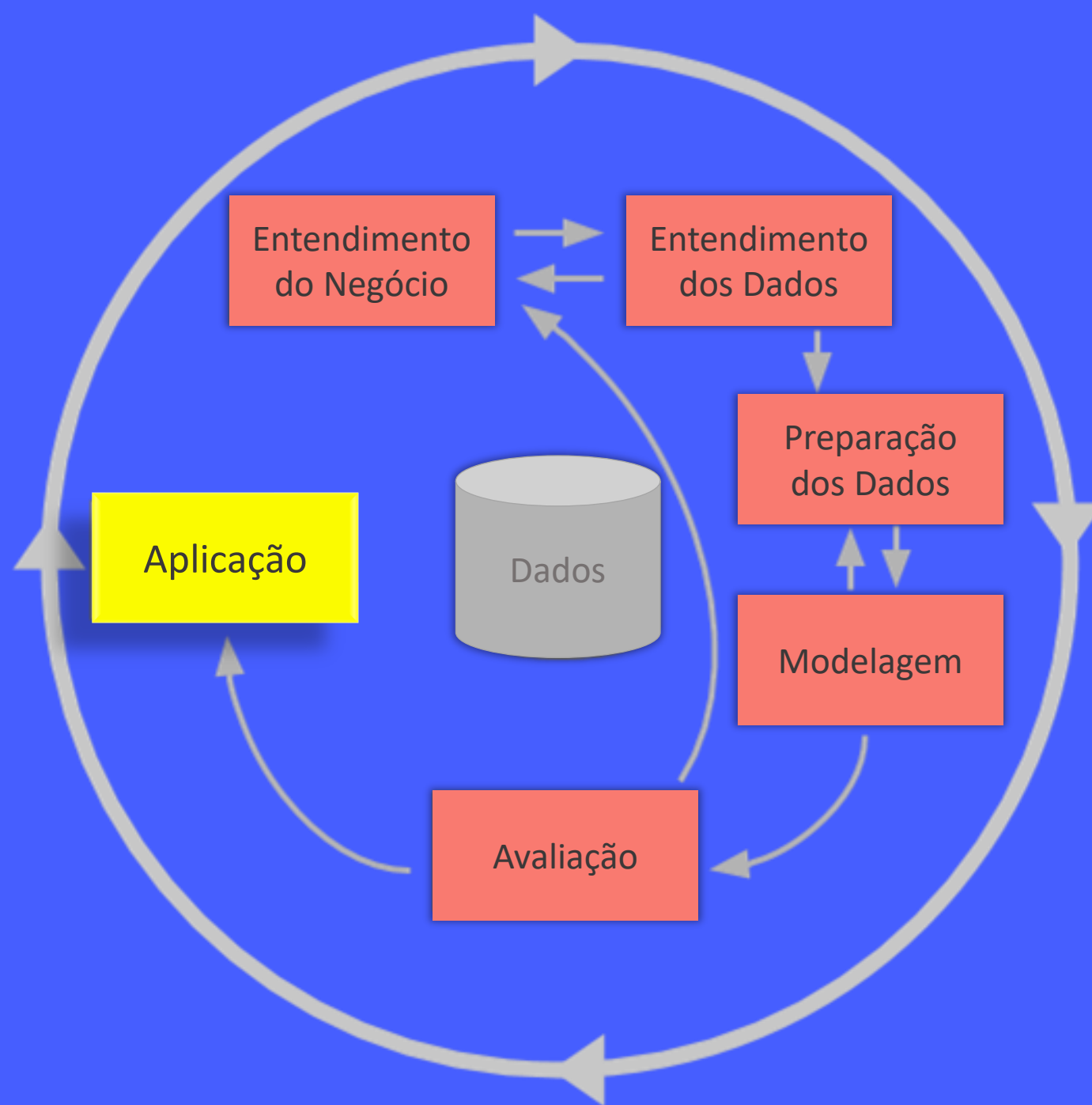
Erro absoluto mediano:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|$$

R²: quanto da variância da variável resposta é explicado pelas variáveis explicativas

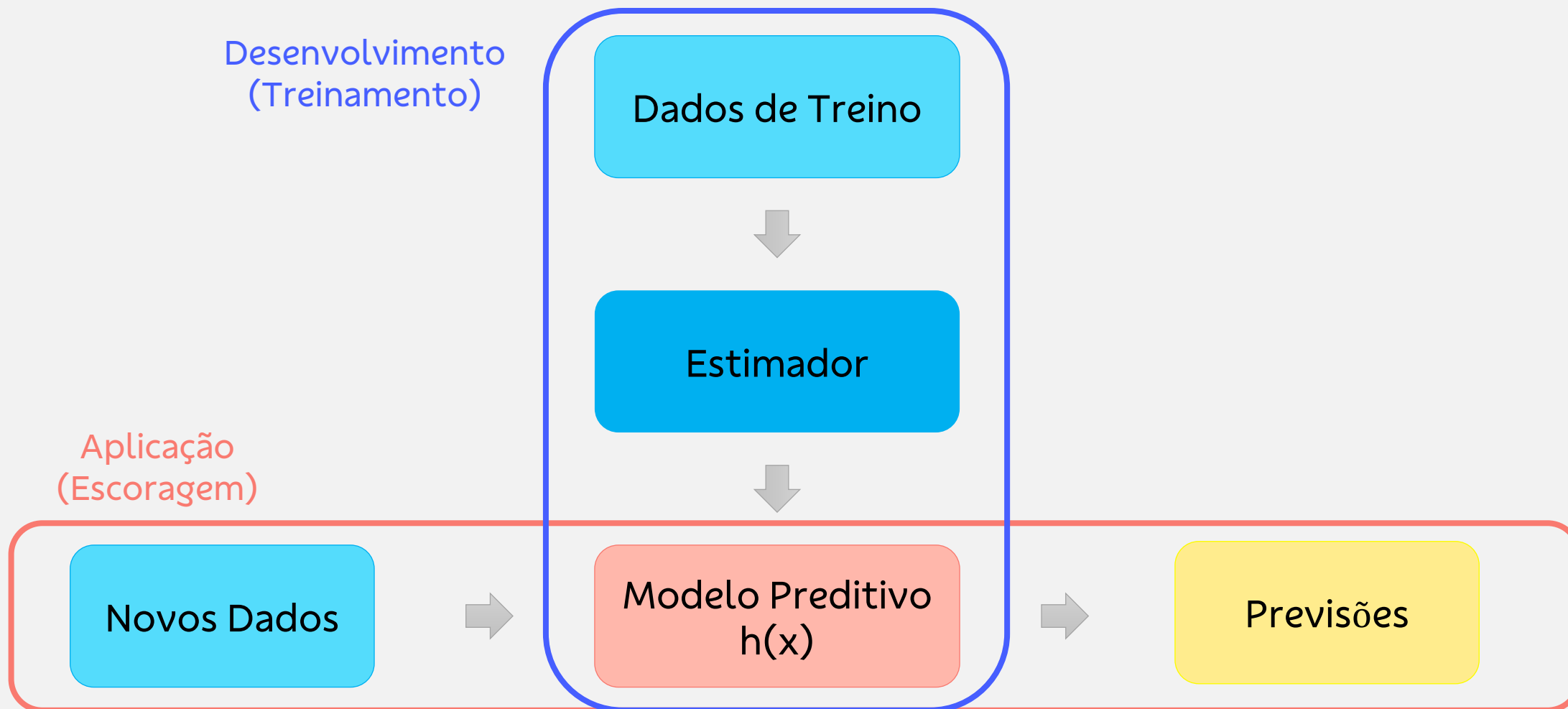
$$R^2 = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$





Aplicação

Uso do modelo para tomada de decisão



Expectativa x Realidade







Comunidade Data Hackers

<https://datahackers.com.br/>

- Fórum (Slack)
- Podcasts (Google, Spotify)
- Blog (Medium)
- Newsletter semanal

Kaggle

<https://www.kaggle.com/>

- Bases de dados reais e fictícias
- Competições com premiações
- Compartilhamento de soluções entre usuários
- Cursos online gratuitos

Coursera

- Algoritmos de Machine Learning:
<https://www.coursera.org/learn/machine-learning>



Pra tudo
que você
imaginar