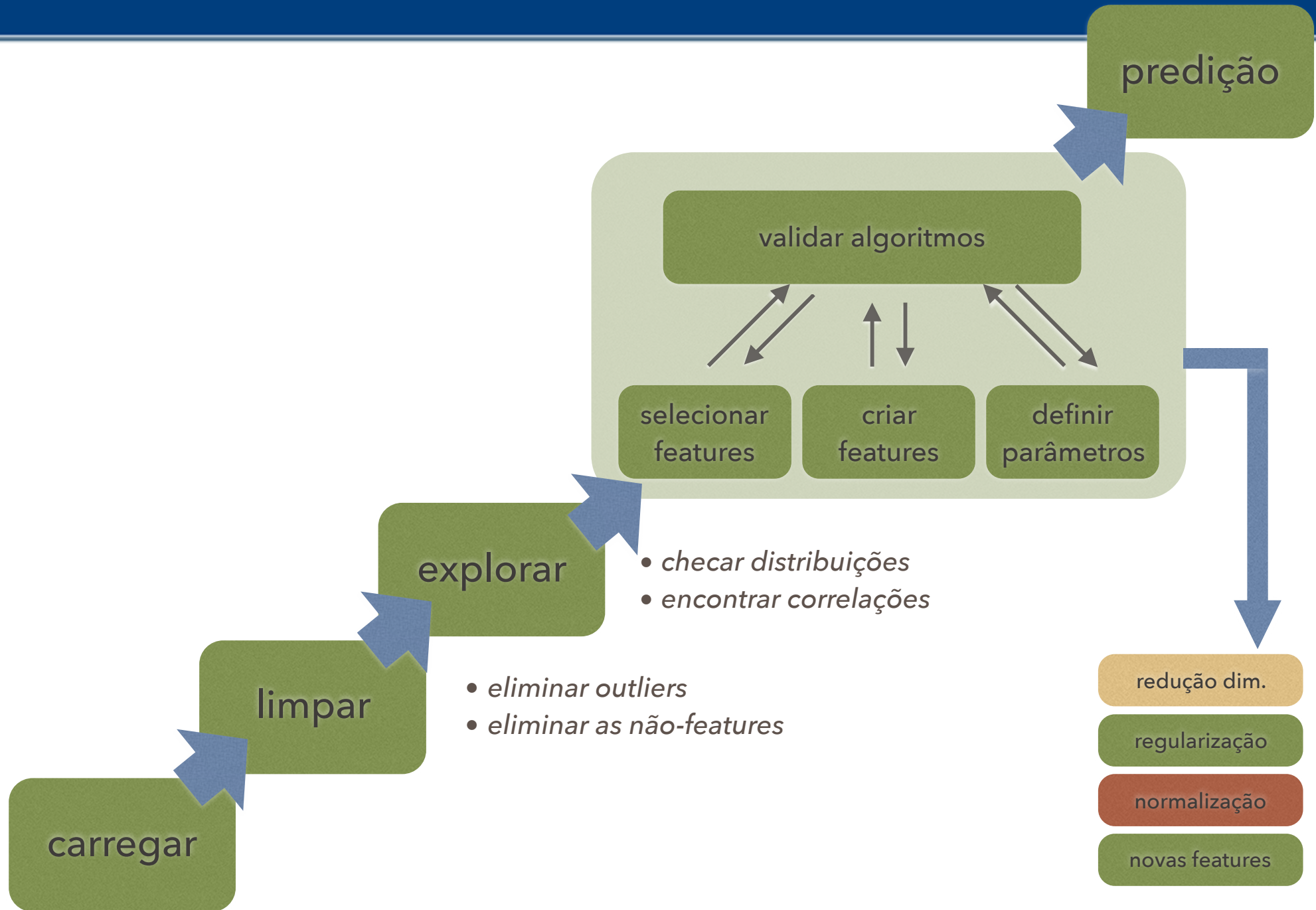


Introdução a Ciencia de Dados

Hitoshi Nagano, Ph.D.

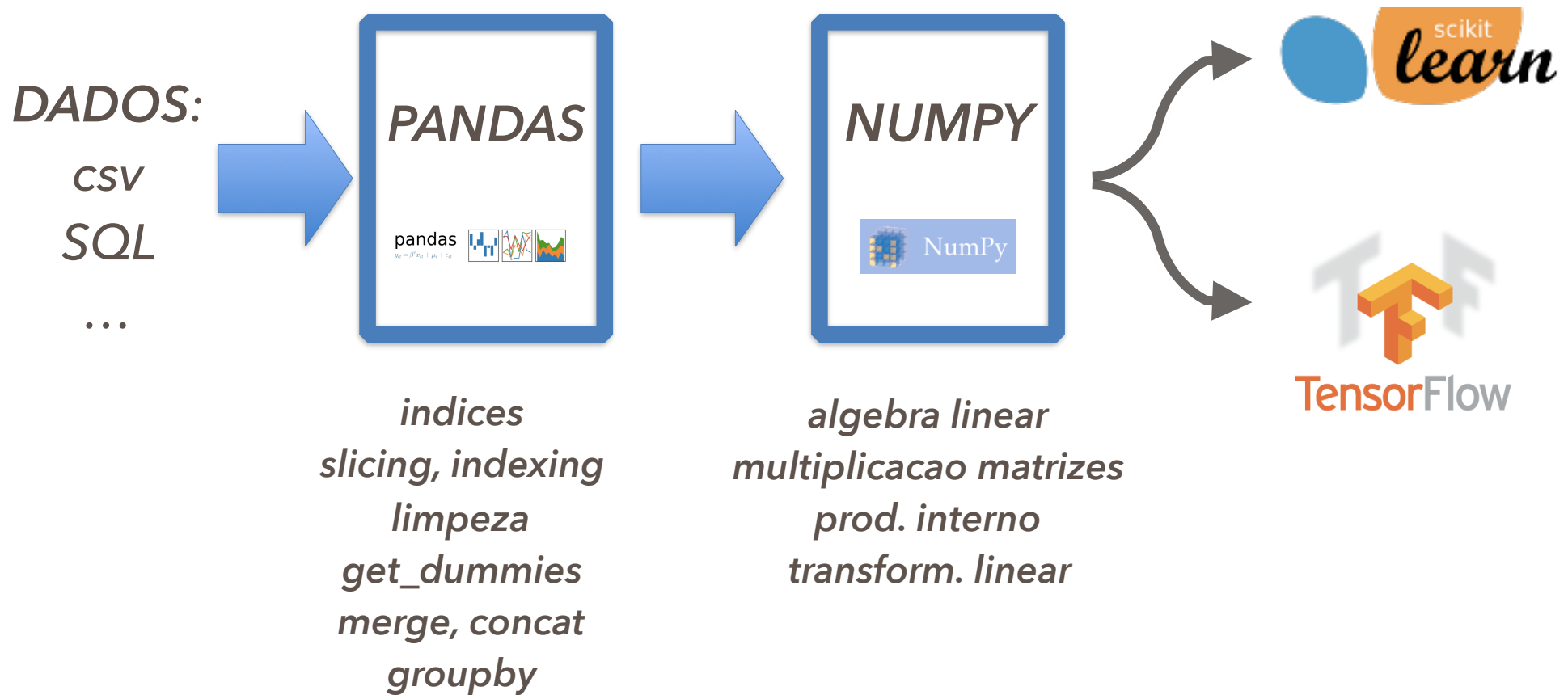


Workflow de projetos em ciencia de dados

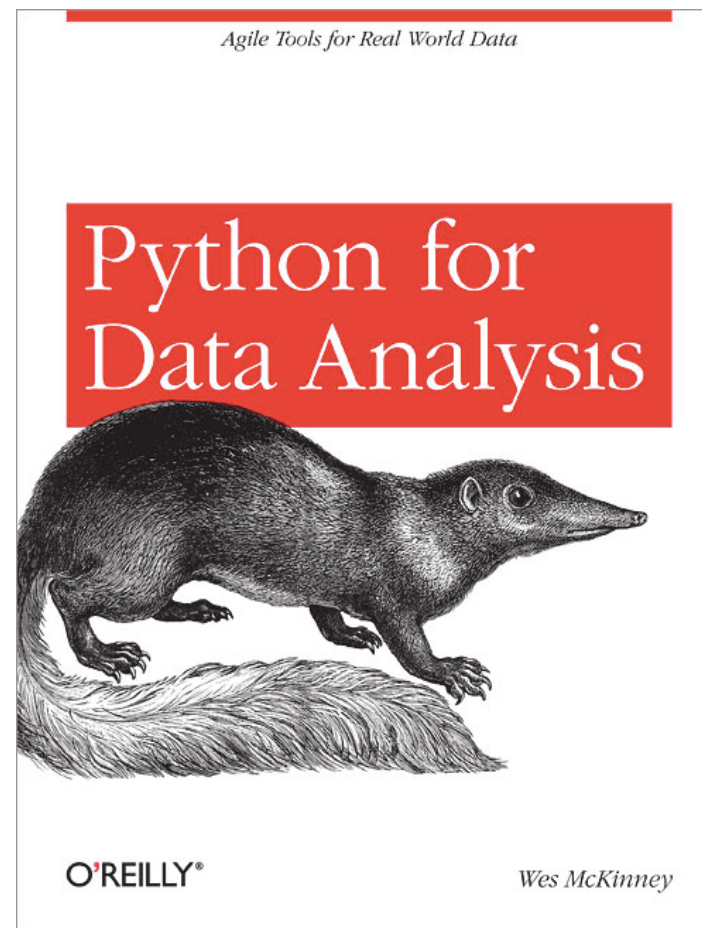


- carregar dados
- checar distribuições
- existem correlações?
- remover outliers
- imputação
- limpeza

PANDAS, NUMPY E SCIKIT-LEARN



- <http://shop.oreilly.com/product/0636920023784.do>



- métodos para limpeza, imputação
- exploração dos dados
- manipulação de variáveis para feature engineering

⇒ qualidade dos dados (lembrar GIGO)

⇒ melhora da performance

VARIÁVEIS E DISTRIBUIÇÕES

- numérica:
 - discreta
 - contínua
- categórica
 - ordinais (ordenáveis)
 - nominais (não ordenáveis)

- ➡ discreta
 - bernoulli
 - binomial
 - poisson
- ➡ contínua
 - normal
 - exponencial

- ➡ paramétricas
- ➡ não-paramétricas

... univariada:

- ➔ média
- ➔ variancia
desvio padrão
- ➔ mediana, quartil
- ➔ moda

... multivariada

- ➔ covariância
- ➔ correlação

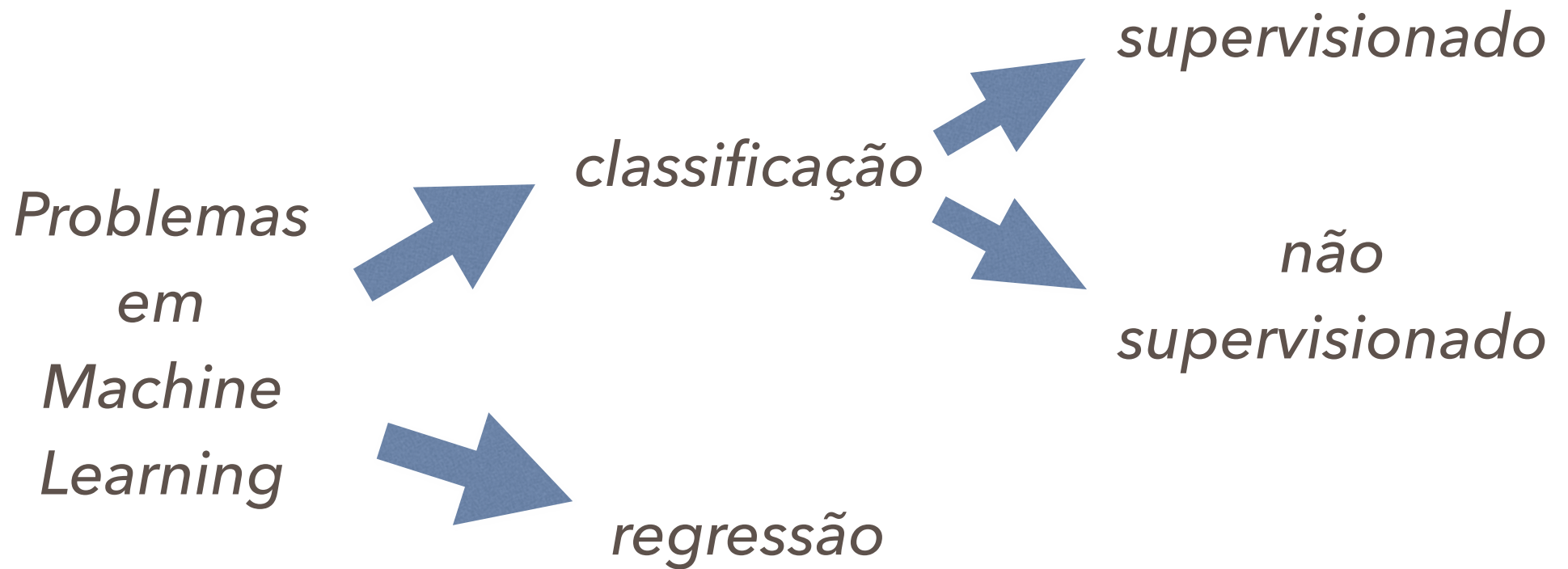
OUTLIERS E IMPUTAÇÃO

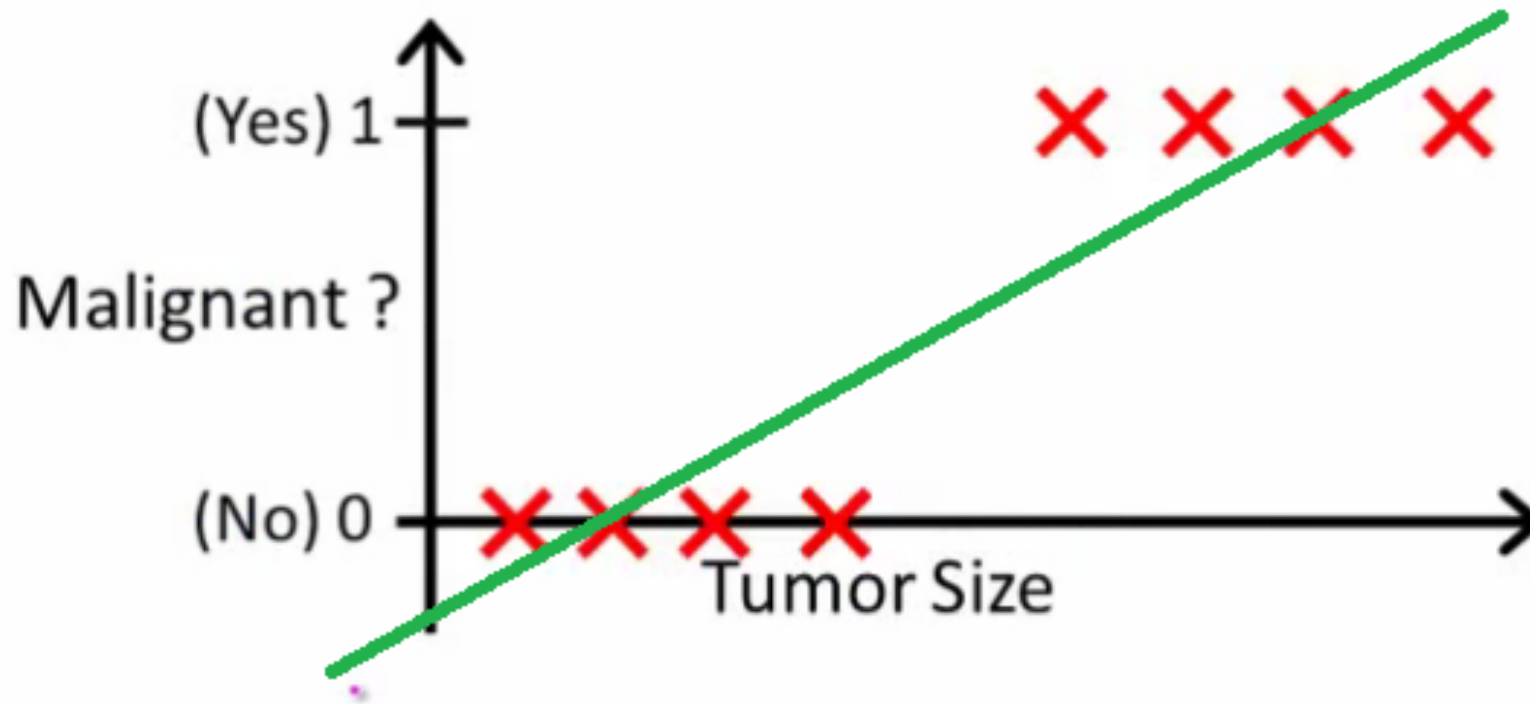
- eventos estranhos,
combinação de eventos analisar entender
- sensores quebrados ignorar
- entrada manual, digitação errada ignorar re-inserir
- erros no processamento ignorar corrigir

- treinar
- verificar as observações com maior erro residual
- retirar essas observações ≈10%
- rodar de novo

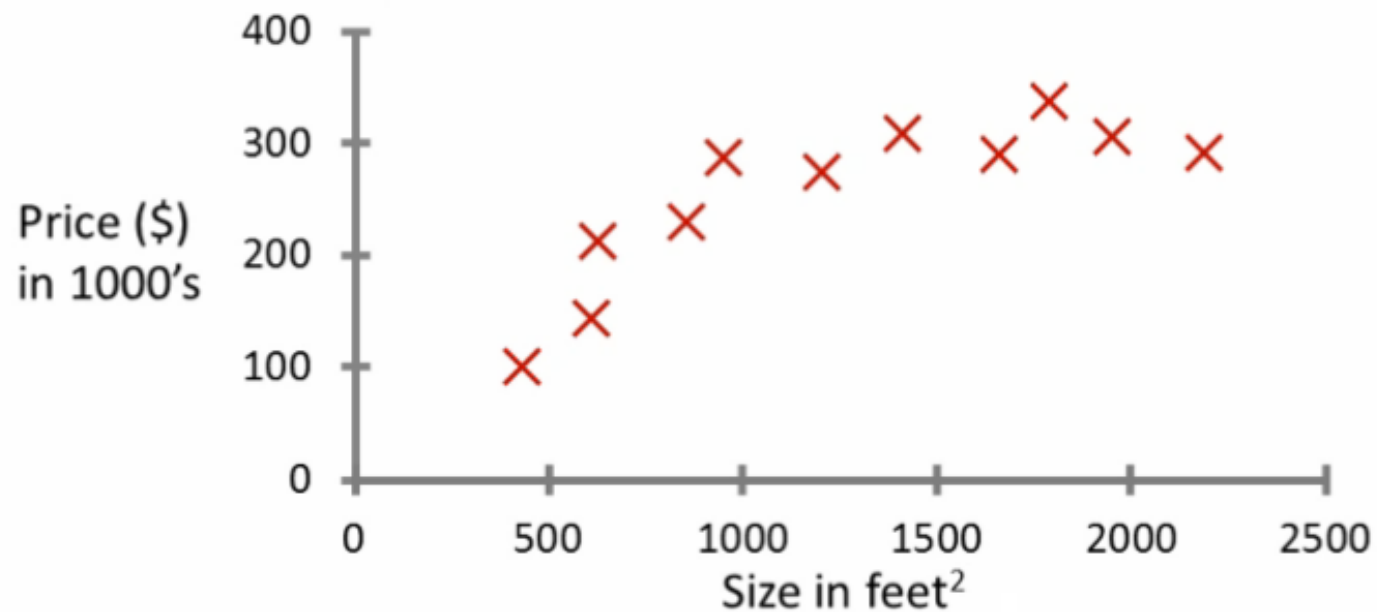
- valor não disponível (NaN, NULL)
- valor claramente errado, equiv. a não disponível
- substituição por um valor unico, pela média,
ou outro critério

TIPOS DE PROBLEMAS

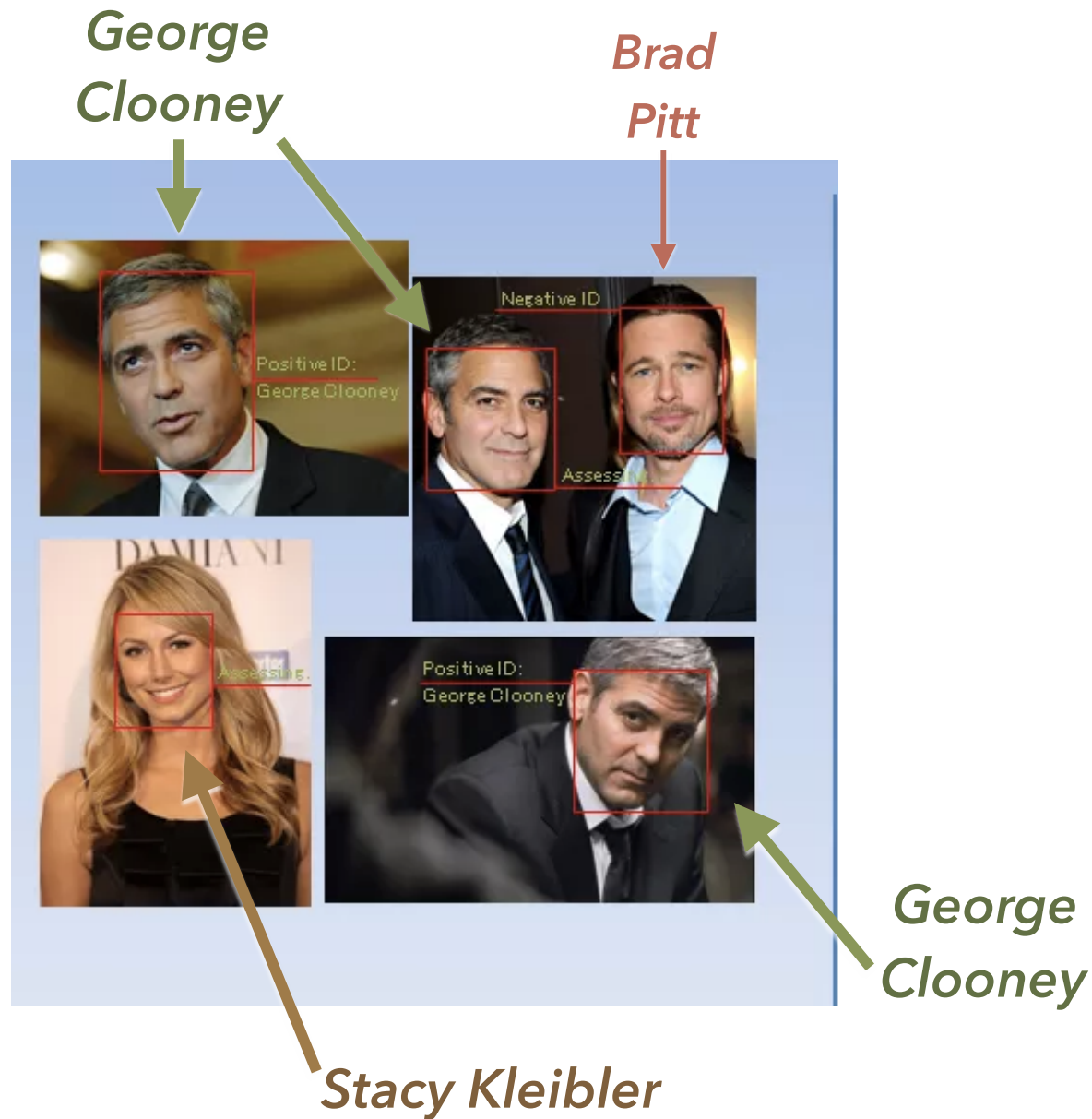




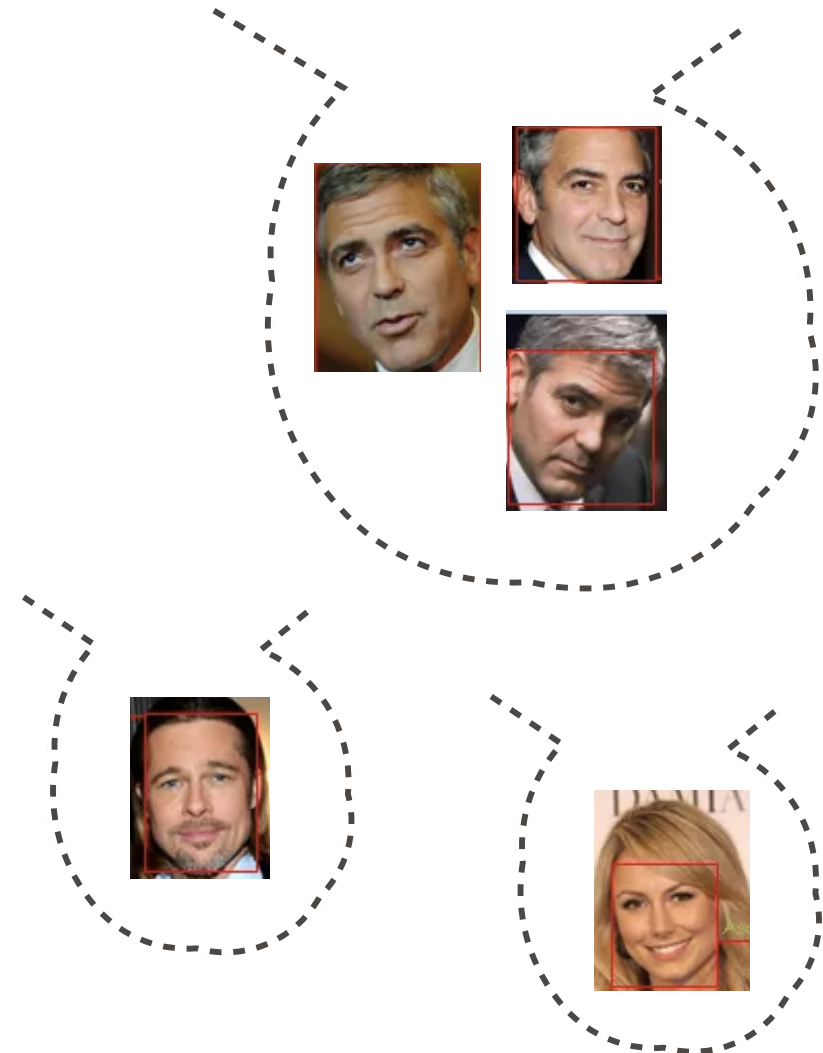
Housing price prediction.



SUPERVISIONADA OU NÃO-SUPERVISIONADA?



SUPERVISIONADA OU NÃO-SUPERVISIONADA?





<https://www.continuum.io/>

IMPORTANTE: escolher versão para Python3

- características
- tipos
- funções e classes
- principais estruturas
- macetes:
 - list comprehension e dict comprehension
 - generators
 - lambda
 - sintaxes compactas e úteis
 - medindo tempo de execução