

# Data Science do ZERO

Pipelines

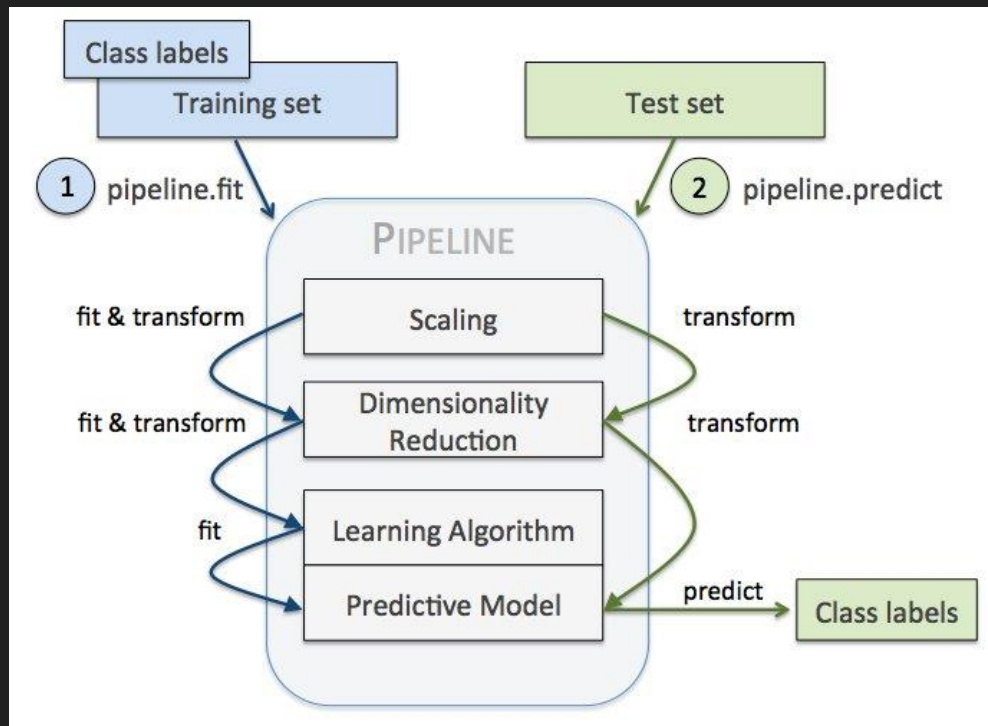
# Pipelines

- Pipeline é um recurso útil para trabalhar com automação de processos em tarefas de Machine Learning.
- Esse recurso é usado para automatizar passos de forma sequencial.
- Permite a padronização e evita erros de manipulação de conjuntos de dados.



# Pipelines

- Com pipelines garantimos que conjuntos de dados de treino e teste são processados de forma correta.
- Com a automatização dos fluxos evitamos erros de manipulação de dados e ajuda na reprodução de código.
- A aplicação sequencial das etapas dos dados de treino e teste da segurança e robustez nos processos.



# Pipelines

- É possível criar diversos Pipelines com características distintas.
- Pipelines usando Scikit-learn permite encapsular etapa de pré-processamento de dados aplicação de um algoritmo de Machine Learning.
- É possível encapsular mais de um pré-processador.

```
1 pip_1 = Pipeline([
2     ('scaler', StandardScaler()),
3     ('clf', svm.SVC())
4 ])
5
6 pip_2 = Pipeline([
7     ('min_max_scaler', MinMaxScaler()),
8     ('clf', svm.SVC())
9 ])
10
11 pip_3 = Pipeline([
12     ('scaler', StandardScaler()),
13     ('pca', PCA(n_components=2)),
14     ('clf', svm.SVC(kernel='rbf'))
15 ])
16
17 pip_4 = Pipeline([
18     ('scaler', StandardScaler()),
19     ('clf', svm.SVC(kernel='poly'))
20 ])
21
22 pip_5 = Pipeline([
23     ('scaler', StandardScaler()),
24     ('clf', svm.SVC(kernel='linear'))
25 ])
```

Hands on!