

## ✓ Bootcamp Ciência de Dados e Inteligência Artificial



### ✓ Contextualização

Uma empresa do setor industrial contatou você para a criação de um sistema inteligente de manutenção preditiva das suas diferentes máquinas. Essa empresa forneceu um conjunto de dados contendo informações coletadas a partir de dispositivos IoT sensorizando atributos básicos de cada máquina. O objetivo é criar um sistema capaz de identificar as falhas que venham a ocorrer, e se possível, qual foi o tipo da falha. Cada amostra no conjunto de dados é composta por 8 atributos que descrevem . Além dessas características, cada amostra é rotulada com uma das 5 possíveis classes de defeitos.

O sistema deverá ser capaz de, a partir de uma nova medição do dispositivo IoT (ou conjunto de medições), prever a classe do defeito e retornar a probabilidade associada. Além disso, a empresa espera que você extraia insights da operação e dos defeitos e gere visualizações de dados.

### Descrição dos Dados

Serão disponibilizados dois arquivos de dados:

**Bootcamp\_train.csv**: use-o para explorar, treinar, avaliar seus modelos. **Bootcamp\_test.csv**: esse arquivo não contém os labels, gere as previsões usando seu modelo e use a seguinte API para ver o desempenho final do modelo.

## Campos e Descrições

	Campo	Descrição
0	id	Identificador das amostras do banco.
1	id_produto	Identificador único do produto. Combinação da variável Tipo e um número de identificação.
2	tipo	Tipo de produto/máquina (L/M/H).
3	temperatura_ar	Temperatura do ar no ambiente (K).
4	temperatura_processo	Temperatura do processo (K).
5	umidade_relativa	Umidade relativa do ar (%).
6	velocidade_rotacional	Velocidade rotacional da máquina em rotações por minutos (RPM).
7	torque	Torque da máquina em Nm.
8	desgaste_da_ferramenta	Duração do uso da ferramenta em minutos.
9	falha_maquina	Indica se houve falha na máquina (1) ou não (0).
10	FDF (Falha Desgaste Ferramenta)	Indica se houve falha por desgaste da ferramenta (1) ou não (0).
11	FDC (Falha Dissipacao Calor)	Indica se houve falha por dissipação de calor (1) ou não (0).
12	FP (Falha Potencia)	Indica se houve falha por potência (1) ou não (0).
13	FTE (Falha Tensao Excessiva)	Indica se houve falha por tensão excessiva (1) ou não (0).
14	FA (Falha Aleatoria)	Indica se houve falha aleatória (1) ou não (0).

## 1 Introdução

Este projeto foi desenvolvido como parte do **Bootcamp de Ciência de Dados e IA**.

O desafio consiste em criar um sistema de **manutenção preditiva** capaz de prever falhas em máquinas industriais a partir de dados coletados por sensores IoT. Além de prever se haverá falha, buscamos identificar **qual tipo de falha** ocorreu.

Isso permite que a empresa atue de forma proativa, reduzindo custos com manutenção corretiva e aumentando a eficiência operacional.

## 2. Exploração Inicial

Começamos com uma análise exploratória automatizada usando ydata-profiling, o que nos permite identificar rapidamente:

- Distribuição das variáveis
- Valores ausentes
- Correlações iniciais
- Potenciais outliers

Essa etapa cria o primeiro mapa do terreno, destacando os principais desafios dos dados.

## 3. Preparação dos Dados

Após conhecer os dados, aplicamos técnicas de pré-processamento:

- Escalonamento de variáveis (StandardScaler)
- Codificação de categorias (LabelEncoder)
- Tratamento de outliers (IsolationForest, métodos estatísticos)
- Balanceamento das classes com SMOTE, essencial para lidar com desbalanceamentos que poderiam enviesar os modelos.

O objetivo é garantir que os modelos recebam dados limpos, consistentes e comparáveis.

## 4. Modelagem

Exploramos diferentes algoritmos de aprendizado supervisionado, cada um com suas características:

- Modelos lineares: LogisticRegression, SGDClassifier
- Árvores e ensembles: RandomForest, GradientBoosting, AdaBoost, Bagging, DecisionTree
- Métodos baseados em vizinhança: KNeighborsClassifier
- Modelos mais sofisticados: SVC

Também utilizamos baselines simples (DummyClassifier) para termos uma referência inicial e entender quanto os modelos realmente agregam valor.

## 5. Avaliação

Comparamos os modelos usando métricas robustas:

- Acurácia (desempenho geral)
- F1-score (equilíbrio entre precisão e recall)
- ROC-AUC e curvas ROC
- Curvas de Precisão-Recall (importantes em cenários desbalanceados)
- Matrizes de confusão para detalhar acertos e erros

Essas métricas nos permitem compreender não apenas “quem acerta mais”, mas como os modelos erram — ponto crucial para decisões de negócio.

## 6. Reprodutibilidade e Rastreamento

Por fim, utilizamos MLflow para rastrear experimentos, registrando:

- Parâmetros dos modelos
- Métricas de avaliação
- Artefatos gerados (gráficos, relatórios)

Isso garante transparência, reprodutibilidade e escalabilidade, permitindo que futuros experimentos sejam comparados de forma estruturada.

## 7. Conclusão

Ao longo desta jornada:

- Exploramos os dados e entendemos seus desafios.
- Preparamos o conjunto para modelagem, cuidando de outliers e desbalanceamento.
- Testamos múltiplos algoritmos, sempre comparando com um baseline.
- Avaliamos com métricas robustas, entendendo pontos fortes e fracos de cada modelo.
- Registramos os resultados com MLflow, consolidando boas práticas de MLOps.

## ✓ Exploração Inicial dos Dados

Nesta etapa, buscamos compreender melhor o conjunto de dados antes de aplicar qualquer modelo de machine learning.

A exploração inicial inclui:

- **Carregamento e inspeção** das primeiras linhas para conhecer a estrutura.
- **Verificação de valores ausentes** e tratamento de inconsistências.
- **Estatísticas descritivas** (média, mediana, desvio padrão, etc.) para identificar padrões e possíveis anomalias.
- **Distribuição das variáveis** numéricas e categóricas, auxiliando na detecção de desbalanceamentos.
- **Visualizações gráficas** (histogramas, boxplots, correlações) para entender relações entre variáveis.
- **Relatório automatizado** com o `ydata-profiling`, que resume as principais características dos dados de forma clara.

O objetivo desta etapa é **criar um panorama confiável do dataset**, garantindo que os próximos passos do pré-processamento e modelagem partam de uma base sólida.

```
1 !pip install -q ydata-profiling mlflow
```



Mostrar saída oculta

```
1 # Grupo 1: Configuração e manipulação básica
2 import warnings
3 import pandas as pd
4 import numpy as np
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7
8 # Grupo 2: Estatística e pré-processamento avançado
9 from scipy import stats
10 from sklearn.svm import SVC
11 from collections import Counter
12 from imblearn.over_sampling import SMOTE
13 from sklearn.dummy import DummyClassifier
14 from sklearn.multioutput import MultiOutputClassifier
15
16 # Grupo 3: MLOps e análise exploratória
```

```
17 import mlflow
18 import mlflow.sklearn
19 from sklearn.ensemble import IsolationForest
20 from ydata_profiling import ProfileReport
21
22 # Grupo 4: Pré-processamento e balanceamento
23 from sklearn.preprocessing import StandardScaler
24 from imblearn.over_sampling import SMOTE
25 from sklearn.preprocessing import LabelEncoder
26 from sklearn.model_selection import train_test_split
27 from sklearn.dummy import DummyClassifier
28 from sklearn.multioutput import MultiOutputClassifier
29
30 # Grupo 5: Algoritmos de machine learning
31 from sklearn.linear_model import LogisticRegression, LogisticRegressionCV, SGDClassifier
32 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier, AdaBoostClassifier
33 from sklearn.tree import DecisionTreeClassifier
34 from sklearn.neighbors import KNeighborsClassifier
35
36 # Grupo 6: Métricas de avaliação
37 from sklearn.metrics import (
38     classification_report,
39     f1_score,
40     confusion_matrix,
41     ConfusionMatrixDisplay,
42     roc_auc_score,
43     accuracy_score,
44     RocCurveDisplay,
45     precision_recall_curve,
46     PrecisionRecallDisplay
47 )
48
49 # Ignore all warnings
50 warnings.filterwarnings("ignore")
```

```
1 df = pd.read_csv('/content/bootcamp_train.csv')
```

```
1 df.head(3)
```



	id	id_produto	tipo	temperatura_ar	temperatura_processo	umidade_relativa	velo
0	0	L56434	L	298.3	309.1	90.0	
1	1	L48741	L	298.2	308.4	90.0	
2	2	L48850	L	298.2	307.8	90.0	

```
1 df.shape
```



```
(25864, 15)
```

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 25864 entries, 0 to 25863
Data columns (total 15 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   id                                     25864 non-null  int64
1   id_produto                           25864 non-null  object
2   tipo                                 25864 non-null  object
3   temperatura_ar                       25407 non-null  float64
4   temperatura_processo                 25402 non-null  float64
5   umidade_relativa                     25864 non-null  float64
6   velocidade_rotacional                25322 non-null  float64
7   torque                               25413 non-null  float64
8   desgaste_da_ferramenta               25160 non-null  float64
9   falha_maquina                       25863 non-null  object
10  FDF (Falha Desgaste Ferramenta)      25863 non-null  object
11  FDC (Falha Dissipacao Calor)          25863 non-null  object
12  FP (Falha Potencia)                  25863 non-null  object
13  FTE (Falha Tensao Excessiva)          25863 non-null  object
14  FA (Falha Aleatoria)                 25863 non-null  object
dtypes: float64(6), int64(1), object(8)
memory usage: 3.0+ MB
```

```
1 df.describe()
```

```

      id  temperatura_ar  temperatura_processo  umidade_relativa  velocidade
count 25864.000000      25407.000000          25402.000000      25864.000000
mean  12931.500000        270.519758           280.375750        89.996766
std    7466.438017         94.934061           97.069056         0.147892
min         0.000000       -36.000000          -38.000000        80.590429
25%    6465.750000        298.000000          308.500000        90.000000
50%    12931.500000        299.600000          309.800000        90.000000
75%    19397.250000        301.100000          310.900000        90.000000
max    25863.000000        304.500000          313.800000        92.067618
```

```
1 df.nunique()
```



0

id	25864
id_produto	9483
tipo	3
temperatura_ar	94
temperatura_processo	83
umidade_relativa	27
velocidade_rotacional	908
torque	572
desgaste_da_ferramenta	246
falha_maquina	8
FDF (Falha Desgaste Ferramenta)	6
FDC (Falha Dissipacao Calor)	6
FP (Falha Potencia)	7
FTE (Falha Tensao Excessiva)	2
FA (Falha Aleatoria)	7

**dtype:** int64

```
1 df.isnull().sum()
```



	0
id	0
id_produto	0
tipo	0
temperatura_ar	457
temperatura_processo	462
umidade_relativa	0
velocidade_rotacional	542
torque	451
desgaste_da_ferramenta	704
falha_maquina	1
FDF (Falha Desgaste Ferramenta)	1
FDC (Falha Dissipacao Calor)	1
FP (Falha Potencia)	1
FTE (Falha Tensao Excessiva)	1
FA (Falha Aleatoria)	1

dtype: int64

## ✓ Descrição dos Dados

O dataset contém medições de sensores em diferentes máquinas. As principais variáveis são:

- **id\_produto**: identificador único da máquina
- **tipo**: tipo da máquina (L, M, H)
- **temperatura\_ar**: temperatura ambiente (K)
- **temperatura\_processo**: temperatura do processo (K)
- **umidade\_relativa**: umidade do ar (%)
- **velocidade\_rotacional**: rotações por minuto (RPM)
- **torque**: torque da máquina (Nm)
- **desgaste\_da\_ferramenta**: duração de uso em minutos

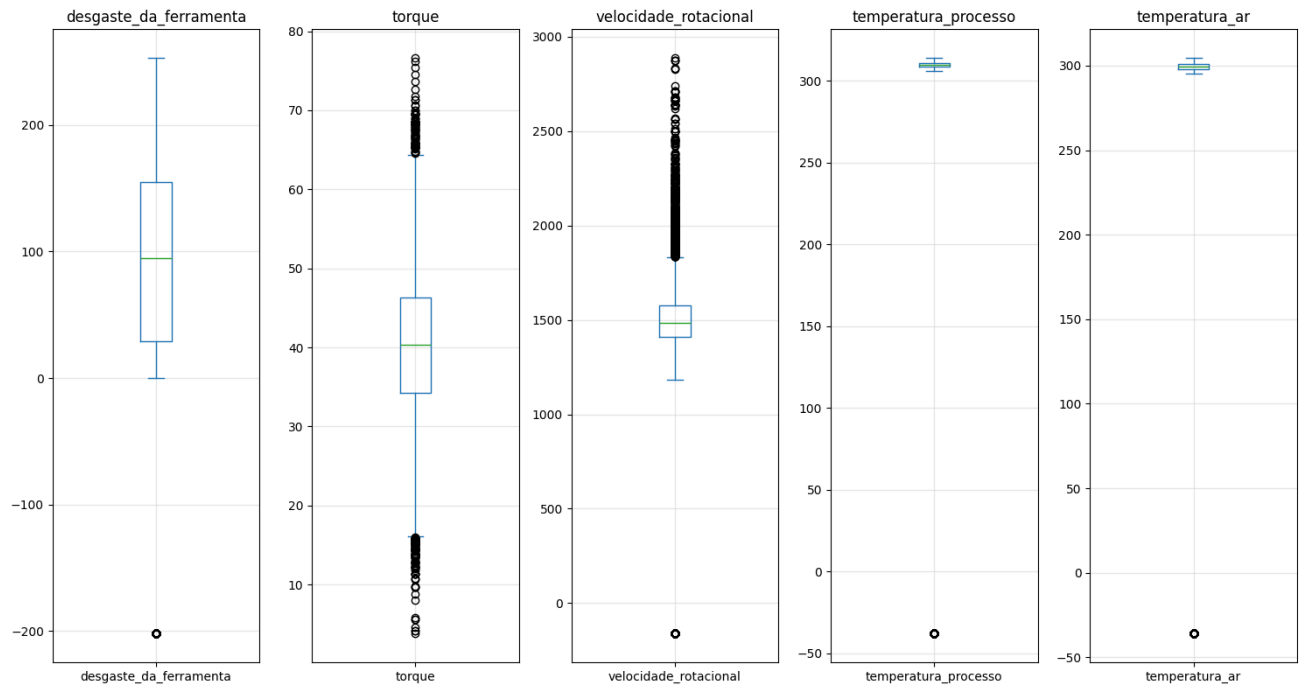
O alvo são as variáveis de falha:

- **falha\_maquina** (1/0)
- Tipos de falha: **FDF, FDC, FP, FTE, FA**

Essas informações permitem tratar o problema como classificação multiclasse ou multirrótulo.



```
1 variaveis = ['desgaste_da_ferramenta', 'torque', 'velocidade_rotacional',  
2             'temperatura_processo', 'temperatura_ar']  
3  
4 plt.figure(figsize=(15, 8))  
5  
6 for i, var in enumerate(variaveis, 1):  
7     plt.subplot(1, 5, i)  
8     df[var].plot.box()  
9     plt.title(f'{var}')  
10    plt.grid(True, alpha=0.3)  
11  
12 plt.tight_layout()  
13 plt.show()
```



- As variáveis temperatura\_processo e temperatura\_ar estão quase constantes, sugerindo pouco poder discriminativo para modelos de machine learning, mas precisam ser mantidas para análise de correlação.
- Há muitos outliers e valores inválidos (negativos em variáveis que não poderiam ser negativas), o que indica necessidade de limpeza e tratamento.
- desgaste\_da\_ferramenta, torque e velocidade\_rotacional parecem ser os atributos mais variáveis e provavelmente relevantes para explicar falhas.

```
1 print(df[variaveis].describe())
```

```

→
count      desgaste_da_ferramenta      torque      velocidade_rotacional \
mean          74.752027          40.186347          1391.365374
std          110.369388           8.973448           481.975539
min         -202.000000           3.800000          -161.000000
25%           29.000000           34.200000          1408.000000
50%           95.000000           40.300000          1484.000000
75%          155.000000           46.300000          1578.000000
max          253.000000           76.600000          2886.000000

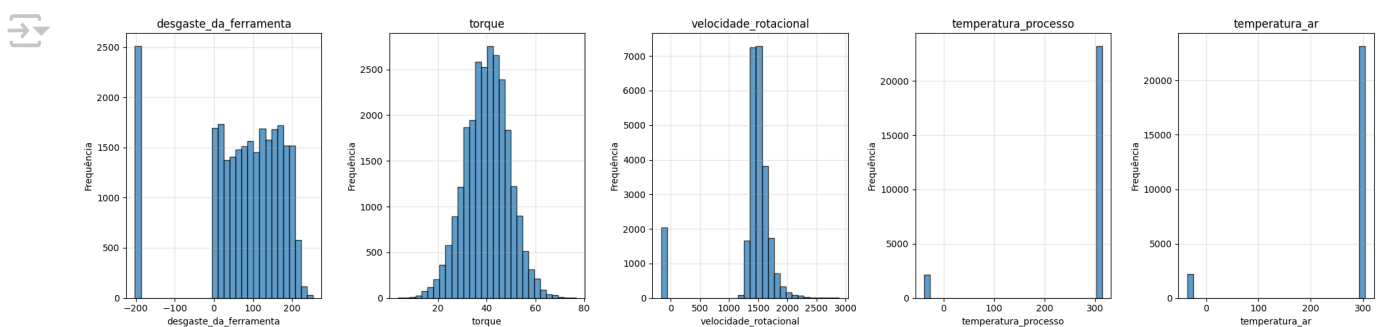
count      temperatura_processo      temperatura_ar
mean          280.375750          270.519758
std           97.069056           94.934061
min          -38.000000          -36.000000
25%          308.500000          298.000000
50%          309.800000          299.600000
75%          310.900000          301.100000
max          313.800000          304.500000

```

```

1 plt.figure(figsize=(20, 5))
2
3 for i, var in enumerate(variaveis, 1):
4     plt.subplot(1, 5, i)
5     df[var].plot.hist(bins=30, alpha=0.7, edgecolor='black')
6     plt.title(f'{var}')
7     plt.grid(True, alpha=0.3)
8     plt.xlabel(var)
9     plt.ylabel('Frequência')
10
11 plt.tight_layout()
12 plt.show()

```



A partir dos histogramas:

- desgaste\_da\_ferramenta → distribuição espalhada, com presença de valores negativos incoerentes (possíveis erros de registro).
- torque → segue padrão aproximadamente normal, mas com cauda longa em valores extremos.
- velocidade\_rotacional → concentrada entre 1200 e 1800 rpm, mas com registros incomuns próximos de zero e acima de 2500.
- temperatura\_processo → altamente concentrada em ~310, com alguns valores inválidos próximos de zero ou negativos.
- temperatura\_ar → comportamento semelhante à temperatura de processo, muito estável em ~300, mas também com registros anômalos próximos de zero.

```

1 # Renomear colunas conforme solicitado
2 rename_map = {
3     "FDF (Falha Desgaste Ferramenta)": "FDF",
4     "FDC (Falha Dissipacao Calor)": "FDC",
5     "FP (Falha Potencia)": "FP",
6     "FTE (Falha Tensao Excessiva)": "FTE",
7     "FA (Falha Aleatoria)": "FA",
8 }
9 df = df.rename(columns=rename_map).copy()
10
11 # Função para normalizar rótulos (0/1)
12 def normalize_label(x):
13     s = str(x).strip().lower()
14     if s in {"1", "sim", "s", "true", "verdadeiro", "y"}:
15         return 1
16     if s in {"0", "não", "nao", "n", "false", "falso", "-"}:
17         return 0
18     try:
19         v = float(s)
20         if v == 1: return 1
21         if v == 0: return 0
22     except:
23         pass
24     return None
25
26 # Aplicar normalização nas colunas de falha
27 falha_cols = [
28     "falha_maquina",
29     "FDF",
30     "FDC",
31     "FP",
32     "FTE",
33     "FA",
34 ]
35
36 for c in falha_cols:
37     if c in df.columns:

```

```

38         df[c] = df[c].apply(normalize_label)
39
40 # TRATAMENTO DE VALORES NULOS/VAZIOS
41
42 # 1. Identificar valores problemáticos nas colunas numéricas
43 problematic_values = ["", " ", "N", "n", "nao", "não", "sim", "s", "y", "0", "1"]
44
45 # 2. Função para limpar valores numéricos
46 def clean_numeric_value(x):
47     if pd.isna(x) or x in problematic_values or str(x).strip() in problematic_values:
48         return np.nan
49     try:
50         return float(x)
51     except:
52         return np.nan
53
54 # 3. Colunas numéricas para tratamento
55 numeric_cols = [
56     "temperatura_ar",
57     "temperatura_processo",
58     "umidade_relativa",
59     "velocidade_rotacional",
60     "torque",
61     "desgaste_da_ferramenta"
62 ]
63
64 # 4. Aplicar limpeza nas colunas numéricas
65 for col in numeric_cols:
66     if col in df.columns:
67         df[col] = df[col].apply(clean_numeric_value)


1 # FUNÇÃO PARA TRATAR VALORES AUSENTES (NULOS)
2 def impute_missing_values(data):
3     """
4     Preenche valores ausentes em um DataFrame.
5     - Colunas numéricas são preenchidas com a mediana.
6     - Colunas categóricas (texto) são preenchidas com a moda.
7     """
8     print("🔍 Verificando e tratando valores ausentes...")
9     for column in data.columns:
10         if data[column].isnull().sum() > 0: # Checa se há nulos na coluna
11             if data[column].dtype == 'object':
12                 # Preenche com a moda (valor mais frequente)
13                 mode_value = data[column].mode()[0]
14                 data[column].fillna(mode_value, inplace=True)
15                 print(f"    - Coluna '{column}' (categórica): Nulos preenchidos com a
16             else:
17                 # Preenche com a mediana (valor central)
18                 median_value = data[column].median()
19                 data[column].fillna(median_value, inplace=True)
20                 print(f"    - Coluna '{column}' (numérica): Nulos preenchidos com a me
21     print("✅ Valores ausentes tratados com sucesso.\n")
22     return data
23

```

```

24 # Funções de detecção
25 def detect_outliers_iqr(data, column, threshold=1.5):
26     Q1 = data[column].quantile(0.25)
27     Q3 = data[column].quantile(0.75)
28     IQR = Q3 - Q1
29     lower_bound = Q1 - threshold * IQR
30     upper_bound = Q3 + threshold * IQR
31     outliers = data[(data[column] < lower_bound) | (data[column] > upper_bound)]
32     return outliers, lower_bound, upper_bound
33
34 def detect_outliers_zscore(data, column, threshold=3):
35     z_scores = np.abs(stats.zscore(data[column]))
36     outliers = data[z_scores > threshold]
37     return outliers
38
39 # Função de tratamento
40 def cap_outliers_iqr(data, column, threshold=1.5):
41     Q1 = data[column].quantile(0.25)
42     Q3 = data[column].quantile(0.75)
43     IQR = Q3 - Q1
44     lower_bound = Q1 - threshold * IQR
45     upper_bound = Q3 + threshold * IQR
46
47     data_capped = data.copy()
48     data_capped[column] = np.clip(data_capped[column], lower_bound, upper_bound)
49     return data_capped
50
51 # Colunas numéricas
52 numeric_cols = ['temperatura_ar', 'temperatura_processo', 'velocidade_rotacional',
53                 'torque', 'desgaste_da_ferramenta', 'diferenca_termica', 'potencia_me
54 numeric_cols = [col for col in numeric_cols if col in df.columns]
55
56 # Detectar outliers
57 for col in numeric_cols:
58     outliers_iqr, lower_bound, upper_bound = detect_outliers_iqr(df, col)
59     outliers_zscore = detect_outliers_zscore(df, col)
60
61     print(f"{col}:")
62     print(f"    IQR: {len(outliers_iqr)} outliers")
63     print(f"    Z-score: {len(outliers_zscore)} outliers")
64
65 # Aplicar tratamento
66 df_clean = df.copy()
67 for col in numeric_cols:
68     df_clean = cap_outliers_iqr(df_clean, col)
69
70 print(f"\nDataset limpo: {df_clean.shape}")
71 print("Outliers tratados com sucesso!")
72
73 # Exportar dataset tratado
74 df_clean.to_csv('dataset_tratado.csv', index=False)
75 print("Dataset exportado como 'dataset_tratado.csv'")

```

⇒ temperatura\_ar:  
IQR: 2223 outliers

```

    Z-score: 0 outliers
temperatura_processo:
    IQR: 2160 outliers
    Z-score: 0 outliers
velocidade_rotacional:
    IQR: 2947 outliers
    Z-score: 0 outliers
torque:
    IQR: 210 outliers
    Z-score: 0 outliers
desgaste_da_ferramenta:
    IQR: 2513 outliers
    Z-score: 0 outliers

```

```

Dataset limpo: (25864, 15)
Outliers tratados com sucesso!
Dataset exportado como 'dataset_tratado.csv'

```

```
1 ids = df["id"]
```

```
1 df = pd.read_csv('/content/dataset_tratado.csv')
```

```
1 df.groupby(['tipo' , 'falha_maquina'])['desgaste_da_ferramenta'].median()
```



		desgaste_da_ferramenta
tipo	falha_maquina	
H	0.0	94.0
	1.0	101.0
L	0.0	94.0
	1.0	146.0
M	0.0	96.0
	1.0	119.0

**dtype:** float64

```

1 # Colunas para excluir do tratamento
2 excluir = ['id', 'id_produto']
3
4 # Tratar numéricas com mediana
5 numericas = df.select_dtypes(include=['float64', 'int64']).columns
6 for col in numericas:
7     if col not in excluir and df[col].isnull().sum() > 0:
8         df[col].fillna(df[col].median(), inplace=True)
9
10 # Tratar categóricas com moda
11 categoricas = df.select_dtypes(include=['object']).columns
12 for col in categoricas:
13     if col not in excluir and df[col].isnull().sum() > 0:

```

```

14         df[col].fillna(df[col].mode()[0], inplace=True)
15
16 print("Tratamento concluído!")
17 print(df.isnull().sum())

```



Tratamento concluído!

```

id                0
id_produto        0
tipo              0
temperatura_ar    0
temperatura_processo  0
umidade_relativa  0
velocidade_rotacional  0
torque            0
desgaste_da_ferramenta  0
falha_maquina     0
FDF               0
FDC               0
FP               0
FTE               0
FA               0
dtype: int64

```

```

1 #criando novas feature
2 df['diferenca_termica']=df['temperatura_processo']-df['temperatura_ar']
3 df['potencia_mecanica']=np.round((df['torque']*df['velocidade_rotacional']* 2 * np.pi

```

1 df



	id	id_produto	tipo	temperatura_ar	temperatura_processo	umidade_relativ
0	0	L56434	L	298.30	309.1	90
1	1	L48741	L	298.20	308.4	90
2	2	L48850	L	298.20	307.8	90
3	3	M20947	M	300.90	310.8	90
4	4	L53849	L	293.35	310.5	90
...	...	...	...	...	...	...
25859	25859	H30008	H	297.60	309.6	90
25860	25860	L55009	L	299.80	311.2	90
25861	25861	L56495	L	293.35	309.5	90
25862	25862	L51563	L	299.60	309.5	90
25863	25863	M18454	M	301.60	310.4	90

25864 rows × 7 columns

```

1 variaveis = ['desgaste_da_ferramenta', 'torque', 'velocidade_rotacional',
2             'temperatura_processo', 'temperatura_ar']

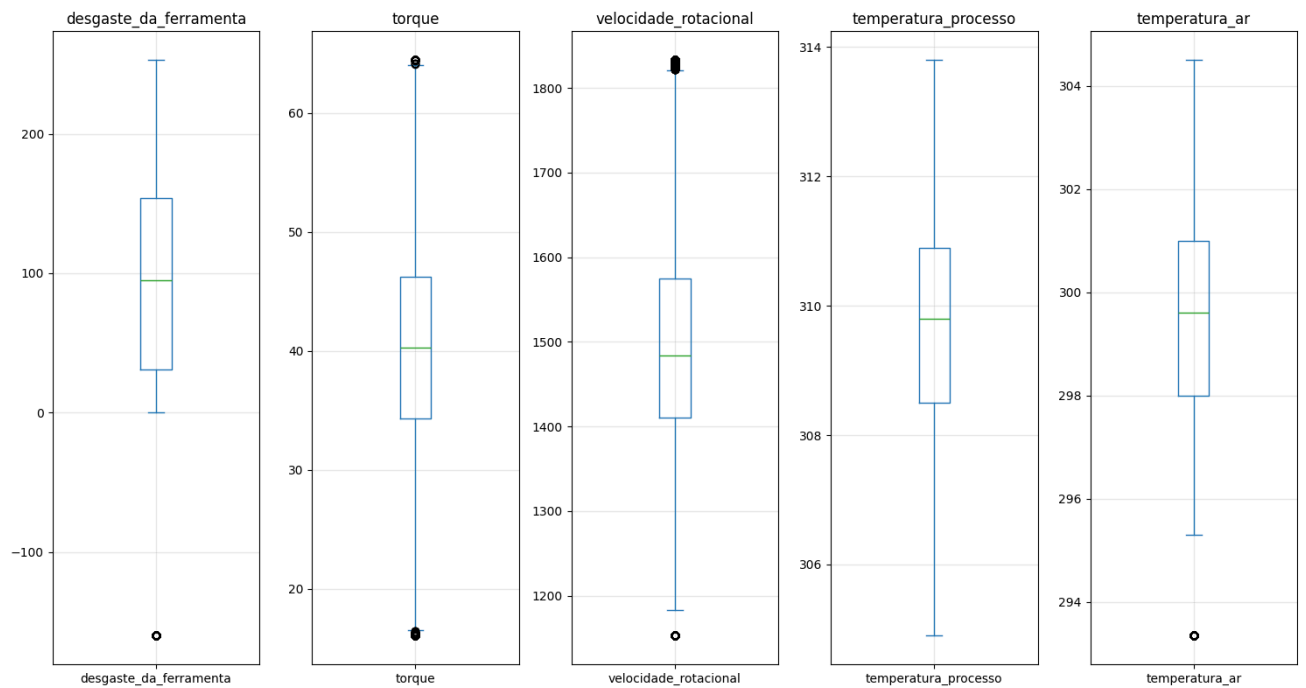
```



```

3
4 plt.figure(figsize=(15, 8))
5
6 for i, var in enumerate(variaveis, 1):
7     plt.subplot(1, 5, i)
8     df[var].plot.box()
9     plt.title(f'{var}')
10    plt.grid(True, alpha=0.3)
11
12 plt.tight_layout()
13 plt.show()

```



## ✓ dados mais estáveis

- `desgaste_da_ferramenta` → apresenta grande variabilidade, com presença de valores negativos incoerentes (possíveis erros ou anomalias).
- `torque` → a maioria dos valores está bem concentrada, mas há registros muito baixos e alguns pontos extremos altos.
- `velocidade_rotacional` → distribuição relativamente estável, mas com outliers em valores mais baixos e mais altos que a média.
- `temperatura_processo` → variável bastante estável, variando pouco em torno da mediana (~310), com poucos outliers.
- `temperatura_ar` → também bastante estável, com a maioria das observações próximas de 300, e poucos registros extremos.

O conjunto de dados é majoritariamente estável em termos de torque, rotação e temperaturas, mas o desgaste da ferramenta e alguns valores de torque/rotação exigem maior atenção por conter outliers e possíveis inconsistências.

```
1 df.nunique()
```



	0
id	25864
id_produto	9483
tipo	3
temperatura_ar	94
temperatura_processo	83
umidade_relativa	27
velocidade_rotacional	606
torque	483
desgaste_da_ferramenta	246
falha_maquina	2
FDF	2
FDC	2
FP	2
FTE	2
FA	2
diferenca_termica	268
potencia_mecanica	16225

**dtype:** int64

```
1 # Eliminar colunas id e id_produto
2 df = df.drop(['id','id_produto'], axis=1)
```

```
1 df
```



	tipo	temperatura_ar	temperatura_processo	umidade_relativa	velocidade_rotac
0	L	298.30	309.1	90.0	
1	L	298.20	308.4	90.0	
2	L	298.20	307.8	90.0	
3	M	300.90	310.8	90.0	
4	L	293.35	310.5	90.0	
...	...	...	...	...	
25859	H	297.60	309.6	90.0	
25860	L	299.80	311.2	90.0	
25861	L	293.35	309.5	90.0	
25862	L	299.60	309.5	90.0	
25863	M	301.60	310.4	90.0	

25864 rows × 15 columns

```
1 df.info()
```



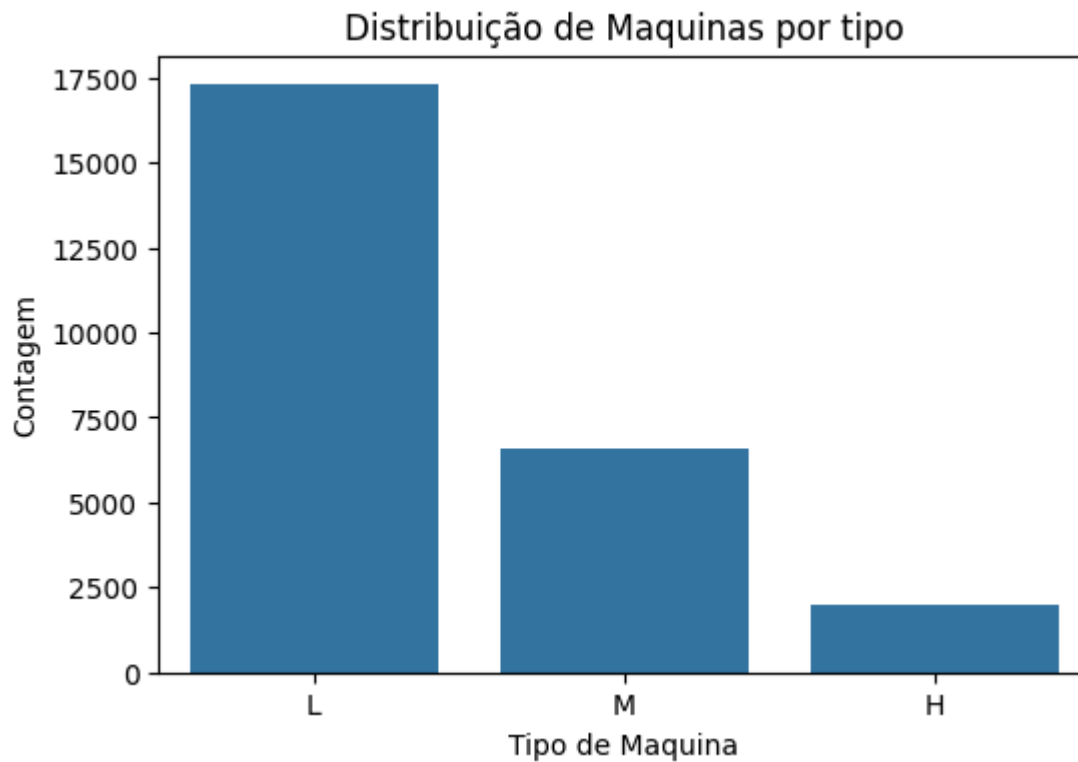
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 25864 entries, 0 to 25863
Data columns (total 15 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   tipo                                  25864 non-null  object
1   temperatura_ar                       25864 non-null  float64
2   temperatura_processo                 25864 non-null  float64
3   umidade_relativa                     25864 non-null  float64
4   velocidade_rotacional                25864 non-null  float64
5   torque                              25864 non-null  float64
6   desgaste_da_ferramenta               25864 non-null  float64
7   falha_maquina                       25864 non-null  float64
8   FDF                                  25864 non-null  float64
9   FDC                                  25864 non-null  float64
10  FP                                   25864 non-null  float64
11  FTE                                  25864 non-null  float64
12  FA                                   25864 non-null  float64
13  diferenca_termica                    25864 non-null  float64
14  potencia_mecanica                   25864 non-null  float64
dtypes: float64(14), object(1)
memory usage: 3.0+ MB
```

```
1 df.describe(include='all').T
```



	count	unique	top	freq	mean	std	mi
tipo	25864	3	L	17285	NaN	NaN	Na
temperatura_ar	25864.0	NaN	NaN	NaN	299.341082	2.572402	293.3
temperatura_processo	25864.0	NaN	NaN	NaN	309.538215	1.9383	304.
umidade_relativa	25864.0	NaN	NaN	NaN	89.996766	0.147892	80.59042
velocidade_rotacional	25864.0	NaN	NaN	NaN	1489.983375	156.690763	1153.
torque	25864.0	NaN	NaN	NaN	40.188267	8.8219	16.0
desgaste_da_ferramenta	25864.0	NaN	NaN	NaN	79.38397	98.752689	-160.
falha_maquina	25864.0	NaN	NaN	NaN	0.02053	0.141809	0.
FDF	25864.0	NaN	NaN	NaN	0.002049	0.045222	0.
FDC	25864.0	NaN	NaN	NaN	0.007346	0.085396	0.
FP	25864.0	NaN	NaN	NaN	0.004176	0.064486	0.
FTE	25864.0	NaN	NaN	NaN	0.005374	0.073114	0.
FA	25864.0	NaN	NaN	NaN	0.001933	0.043926	0.
diferenca_termica	25864.0	NaN	NaN	NaN	10.197133	2.564573	0.
potencia_mecanica	25864.0	NaN	NaN	NaN	6181.05578	1106.145119	1937.907

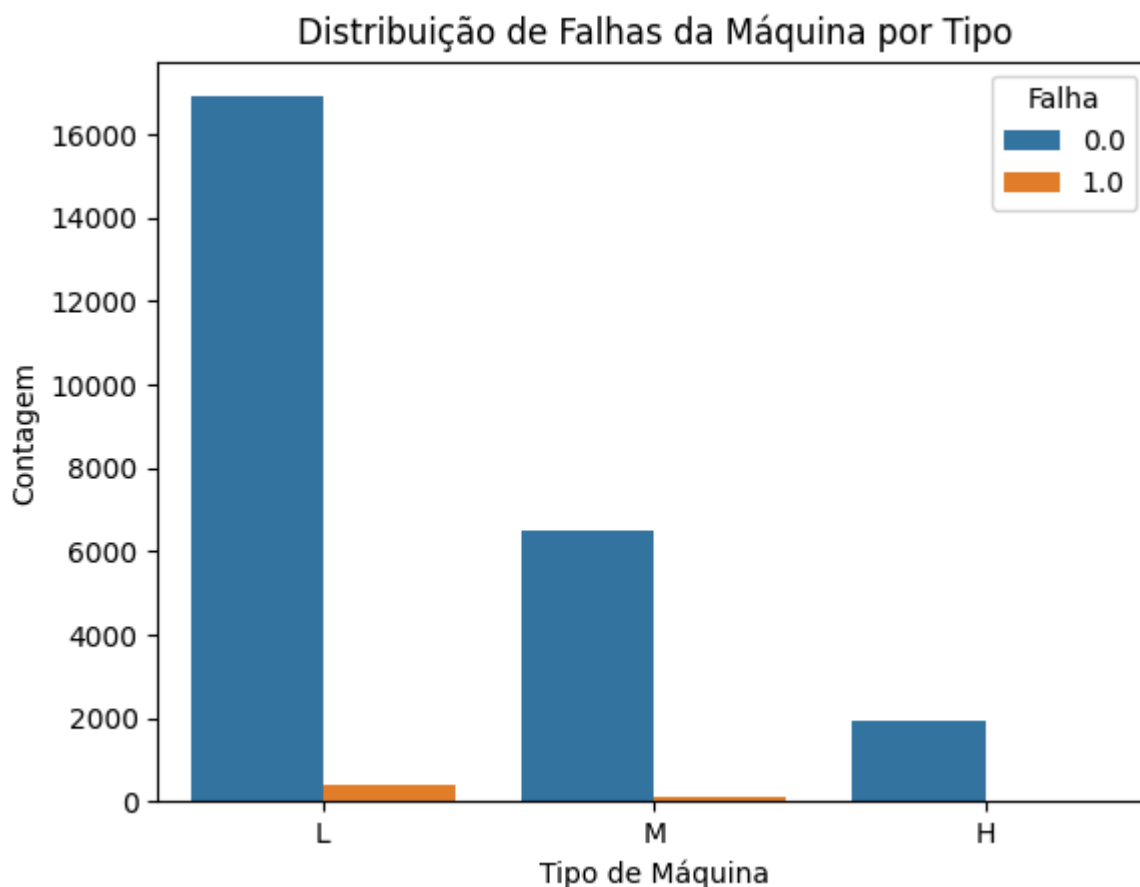
```
1 plt.figure(figsize=(6,4))
2 sns.countplot(x='tipo', data=df)
3 plt.title('Distribuição de Maquinas por tipo')
4 plt.xlabel('Tipo de Maquina')
5 plt.ylabel('Contagem')
6 plt.show()
```



### ✓ Distribuição de Máquinas por Tipo

- O tipo L é o mais frequente, representando a grande maioria das máquinas.
- O tipo M aparece em quantidade intermediária, mas bem menor que L.
- O tipo H é o menos comum, com presença bastante reduzida.

```
1 sns.countplot(x='tipo', hue='falha_maquina', data=df)
2
3 # Adiciona os títulos e rótulos em português
4 plt.title('Distribuição de Falhas da Máquina por Tipo')
5 plt.xlabel('Tipo de Máquina')
6 plt.ylabel('Contagem')
7 plt.legend(title='Falha')
8 plt.show()
```



### ✓ Distribuição de Falhas da Máquina por Tipo

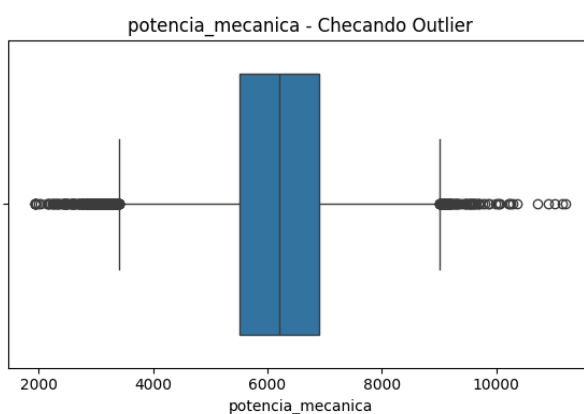
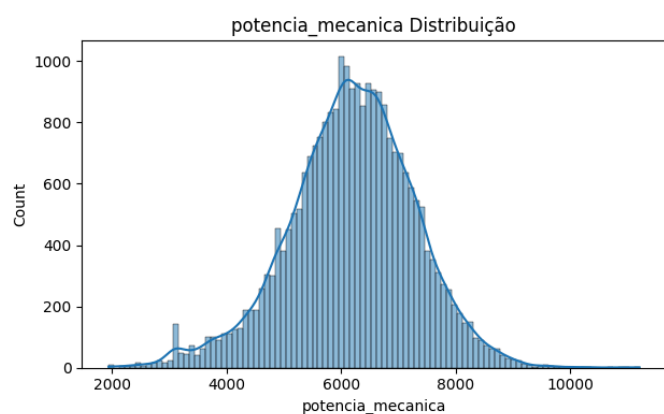
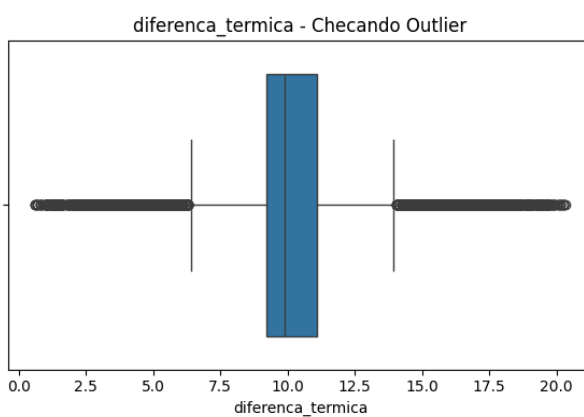
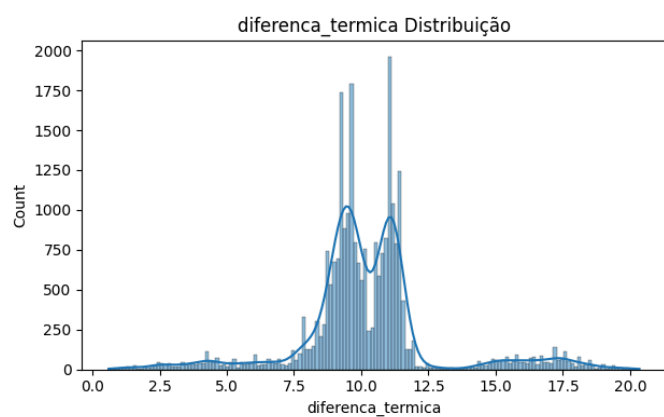
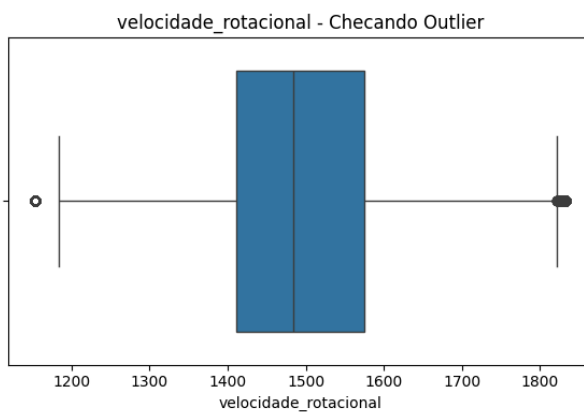
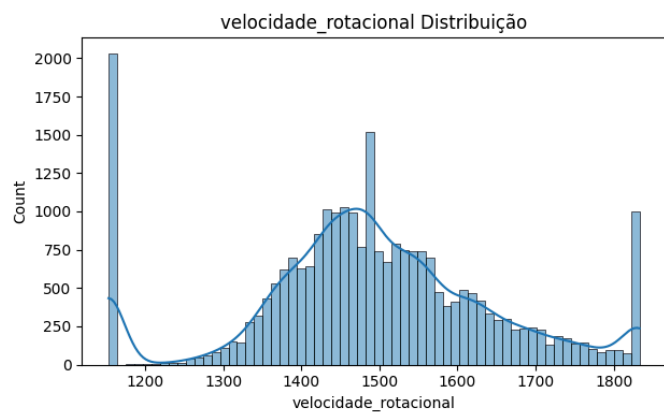
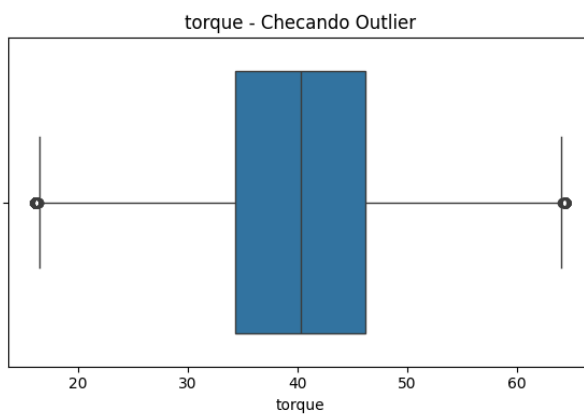
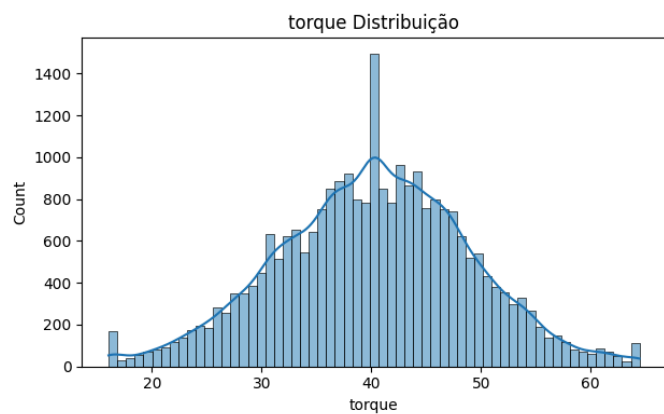
- A maioria absoluta das máquinas não apresentou falha (classe 0).
- As falhas (classe 1) são muito menos frequentes, indicando um forte desbalanceamento no alvo.
- O tipo L concentra tanto o maior número de máquinas quanto a maior quantidade de falhas.
- Os tipos M e H apresentam bem menos registros e, conseqüentemente, menos falhas.

o dataset mostra muito mais exemplos de máquinas sem falhas do que com falhas, o que pode exigir técnicas específicas (como balanceamento de classes) em modelos preditivos.

```
1 cols = ['torque', 'velocidade_rotacional', 'diferenca_termica', 'potencia_mecanica']
2
3 for col in cols:
4     fig, axes = plt.subplots(1, 2, figsize=(12, 4)) # 1 linha , 2 columnas
5
6     # Histograma com KDE
7     sns.histplot(data=df, x=col, kde=True, ax=axes[0])
8     axes[0].set_title(f"{col} Distribuição")
9
10    # Boxplot
11    sns.boxplot(data=df, x=col, ax=axes[1])
12    axes[1].set_title(f"{col} - Checando Outlier")
```

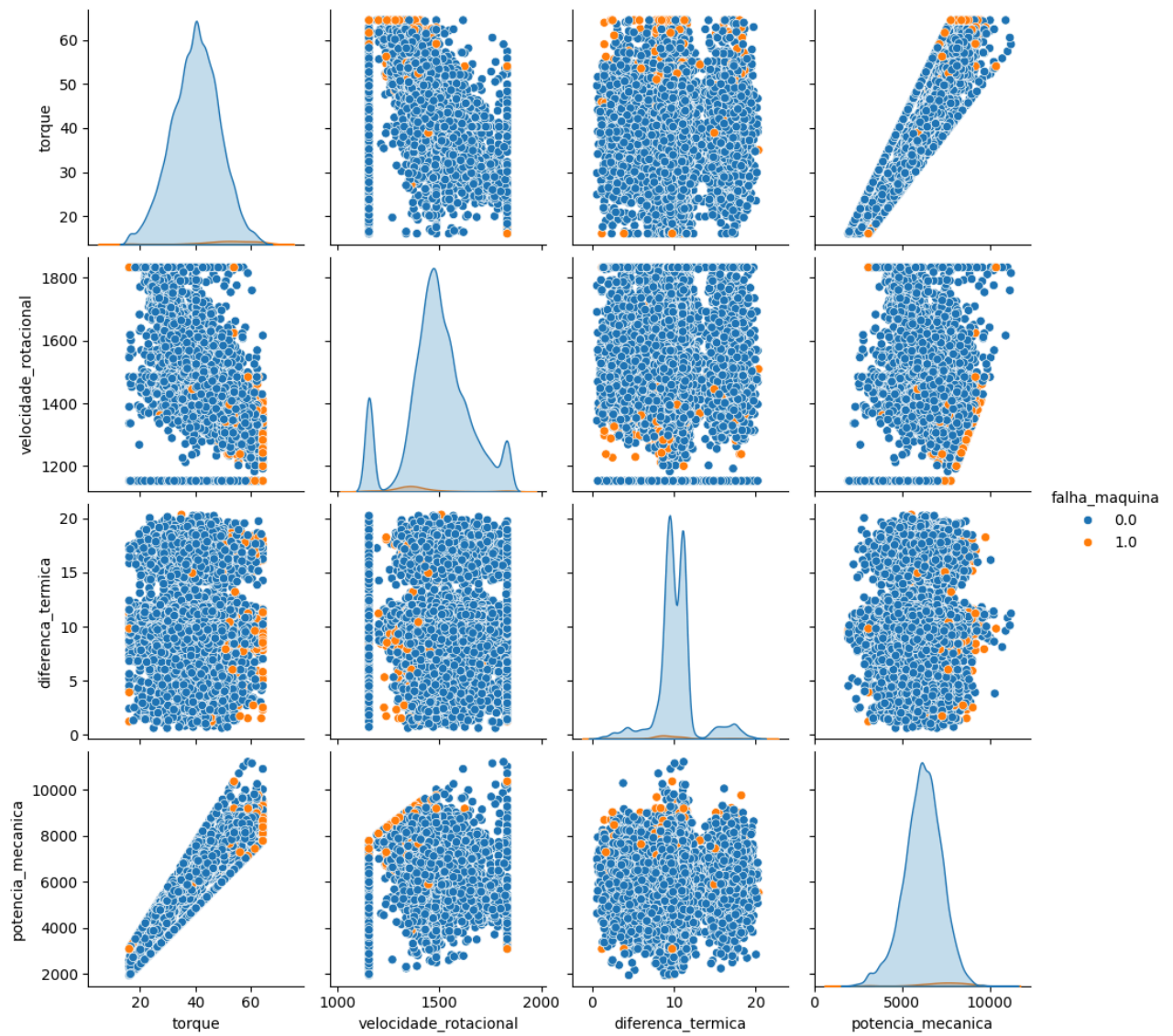
```
13  
14     plt.tight_layout()  
15     plt.show()
```





- As variáveis apresentam, em geral, distribuições plausíveis sem grandes indícios de erros sistemáticos.
- Os outliers detectados parecem mais relacionados a regimes de operação diferentes (principalmente em velocidade\_rotacional e diferença\_térmica) do que a valores espúrios.
- Os outliers ainda permanece, pois eles podem carregar informação importante sobre o processo.

```
1 sns.pairplot(df[['torque', 'velocidade_rotacional', 'diferenca_termica', 'potencia_mec
2 plt.show()
```



- Há forte dependência entre torque e potência.

- Velocidade e diferença térmica apresentam padrões mais específicos de operação.
- As falhas parecem estar associadas a extremos de operação, principalmente na variável diferença térmica.

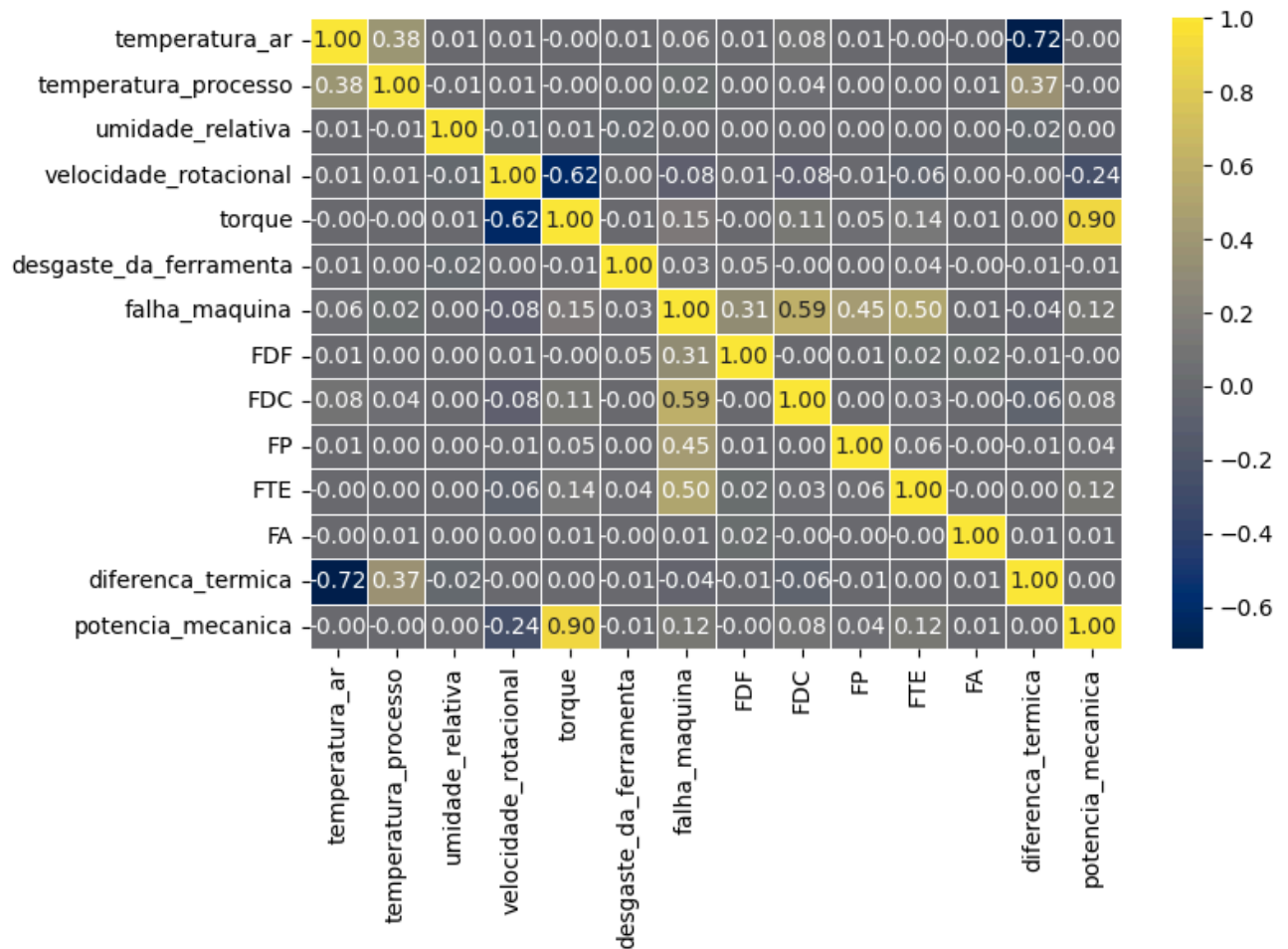
```

1 # Verificando a Correlação entre Características Numéricas Usando um Mapa de Calor (h
2 corr_matrix=df.corr(numeric_only=True)
3 plt.figure(figsize=(8,5))
4 sns.heatmap(corr_matrix,annot=True,cmap='cividis',fmt=".2f", linewidths=0.5)

```



<Axes: >



- Há redundância clara entre torque e potência mecânica.
- Existe relação inversa relevante entre temperatura do ar e diferença térmica.
- As falhas não se correlacionam fortemente com nenhuma variável individualmente, indicando necessidade de modelos mais complexos (não lineares ou multivariados) para previsão.

```
1 # Criar relatório
2 profile = ProfileReport(
3     df,
4     title="Manutenção Preditiva",
5     dataset={
6         "description": "Análise do dataset Gerado por Leonardo Correia",
7         "copyright_holder": "Leonardo Correia",
8         "copyright_year": "2025",
9     },
10    explorative=True,
11 )
12
13 # Exibir no notebook
14 profile.to_notebook_iframe()
15
16 # Salvar o relatório como arquivo HTML no Colab
17 profile.to_file("/content/relatorio_manutencao_preditiva.html")
```

Manutenção Preditiva



# Overview

Brought to you by [YData](#)

Overview

Dataset

Alerts 18

Reproduction

Dataset statistics

Number of variables15

Number of observations25864

Missing cells0

Missing cells (%)0.0%

Duplicate rows20

Duplicate rows (%)0.1%

Total size in memory4.0 MiB

Average record size in memory162.0 B

Variable types

Categorical7

Numeric8

# Variables

Select Columns



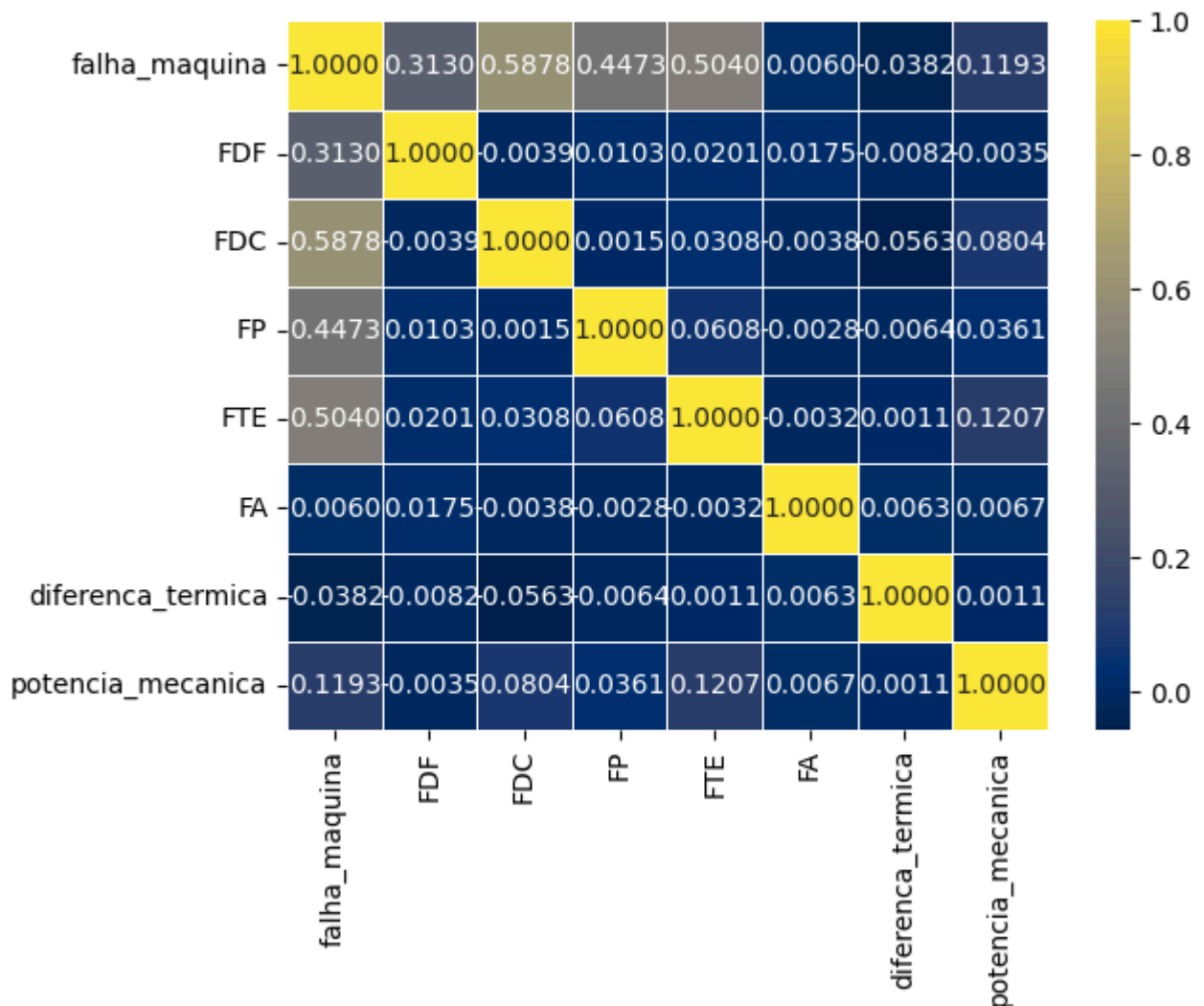
## ✓ O Relatório YData.

- Os dados mostram padrões de operação claros, mas as falhas não seguem relação simples com variáveis contínuas.
- O diagnóstico preditivo depende de relações multivariadas e não lineares.
- Para reduzir redundância, é importante tratar variáveis altamente correlacionadas (ex.: torque e potência mecânica).
- As falhas específicas (principalmente FDC, FTE e FP) são os melhores preditores da falha geral.

```
1 target=df.iloc[:,[7,8,9,10,11,12,13,14]]
2 target_mat=target.corr()
3 sns.heatmap(target_mat,annot=True,cmap="cividis",fmt=".4f",linewidth=0.5)
```



<Axes: >



- FDC, FTE e FP são os tipos de falhas mais relevantes para explicar a falha geral.
- Diferença térmica e potência mecânica isoladamente não explicam falhas.

- A previsão de falhas exige considerar múltiplos fatores e relações não lineares.

```
1 df.drop(columns=['FDF','FDC','FP','FTE','FA'],inplace=True)
```

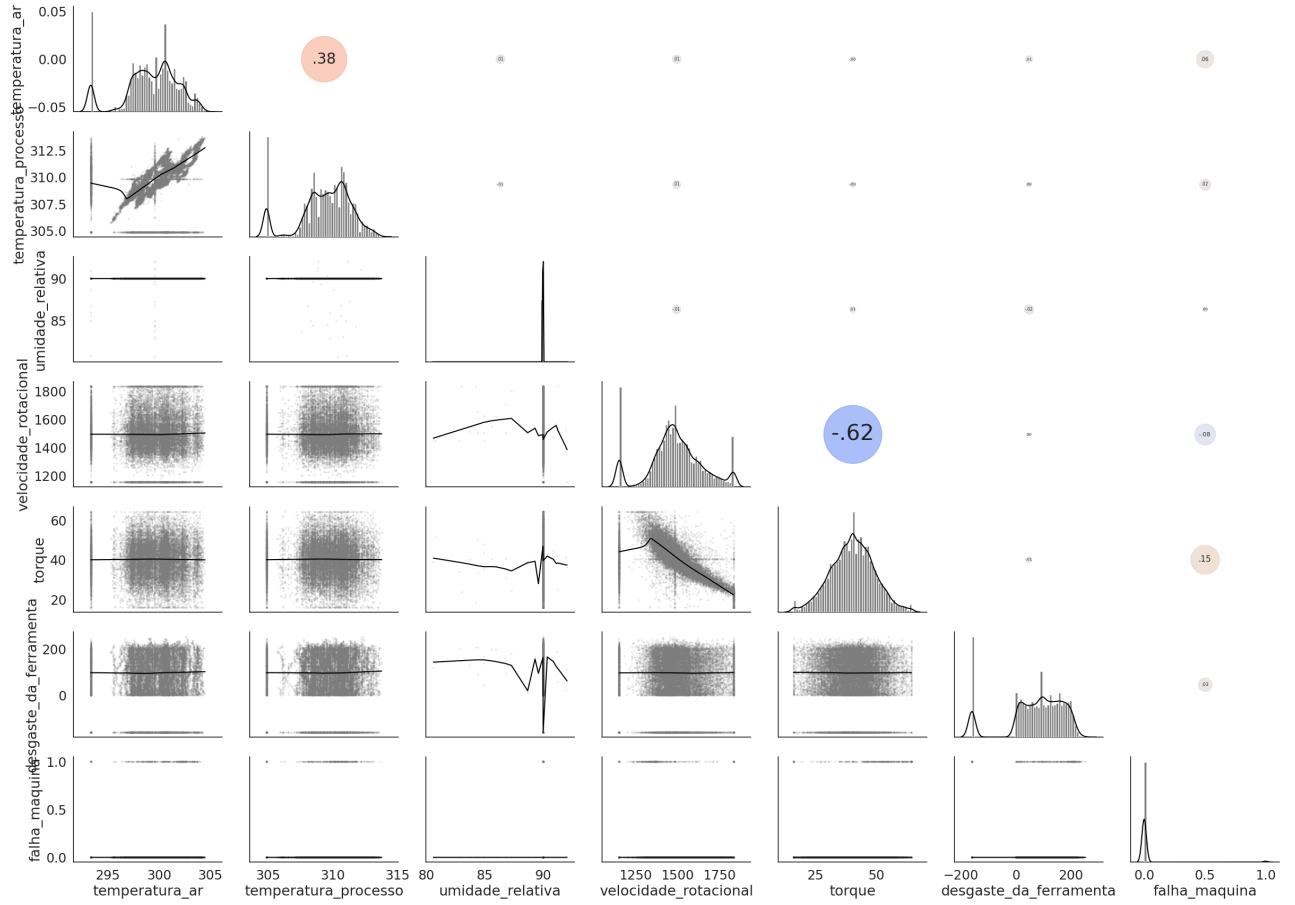
```
1 df.sample(3)
```



	tipo	temperatura_ar	temperatura_processo	umidade_relativa	velocidade_rotac
12971	L	298.4	309.3	90.0	
20821	L	297.3	308.6	90.0	
11100	M	296.8	307.8	90.0	

```
1 def corrdot(*args, **kwargs):
2     corr_r = args[0].corr(args[1])
3     corr_text = f"{corr_r:2.2f}".replace("0.", ".")
4     ax = plt.gca()
5     ax.set_axis_off()
6     marker_size = abs(corr_r) * 10000
7     ax.scatter([.5], [.5], marker_size, [corr_r], alpha=0.6, cmap='coolwarm',
8               vmin=-1, vmax=1, transform=ax.transAxes)
9     font_size = abs(corr_r) * 40 + 5
10    ax.annotate(corr_text, [.5, .5,], xycoords="axes fraction",
11              ha='center', va='center', fontsize=font_size)
12
13 sns.set(style='white', font_scale=1.6)
14 g = sns.PairGrid(df.iloc[:, :8], aspect=1.4, diag_sharey=False)
15 g.map_lower(sns.regplot, lowess=True, ci=False, line_kws={'color': 'black', 'lw': 1.5})
16 g.map_diag(sns.distplot, kde_kws={'color': 'black'}, hist_kws={'color': 'gray', 'alpha':
17 g.map_upper(corrdot)
```





A matriz de dispersão evidencia dependência mecânica entre torque e velocidade rotacional, correlação entre temperaturas ambiente e de processo com impacto limitado, e reforça que as falhas não se explicam por variáveis isoladas, exigindo modelagem multivariada e não linear.

```
1 df['tipo'] = LabelEncoder().fit_transform(df['tipo'])
2 df.head()
```



	tipo	temperatura_ar	temperatura_processo	umidade_relativa	velocidade_rotacional
0	1	298.30	309.1	90.0	1616
1	1	298.20	308.4	90.0	1388
2	1	298.20	307.8	90.0	1528
3	2	300.90	310.8	90.0	1599
4	1	293.35	310.5	90.0	1571

```
1 features = list(df.columns)
2 for feature in features:
3     print(feature + " - " + str(len(df[df[feature].isna()])))
```



```
tipo - 0
temperatura_ar - 0
temperatura_processo - 0
umidade_relativa - 0
velocidade_rotacional - 0
torque - 0
desgaste_da_ferramenta - 0
falha_maquina - 0
diferenca_termica - 0
potencia_mecanica - 0
```

```
1 scale=StandardScaler()
2 data=pd.DataFrame(scale.fit_transform(df),columns=df.columns,index=df.index)
3 data.sample(10)
```



	tipo	temperatura_ar	temperatura_processo	umidade_relativa	velocidade_rotac
<b>11927</b>	-0.326944	-0.132595	0.238247	0.021865	
<b>24010</b>	1.499819	-2.329029	-0.226087	0.021865	
<b>22652</b>	-0.326944	0.761528	-2.392976	0.021865	
<b>20744</b>	-0.326944	0.606028	0.702580	0.021865	
<b>20386</b>	1.499819	-0.521345	0.186654	0.021865	
<b>24382</b>	1.499819	0.295029	0.960543	0.021865	
<b>22485</b>	-0.326944	-0.249220	0.289839	0.021865	
<b>5748</b>	-0.326944	1.694526	1.476469	0.021865	
<b>8248</b>	1.499819	0.100654	0.650988	0.021865	
<b>1811</b>	1.499819	-0.832344	-0.845198	0.021865	

```
1 Y=df.pop("falha_maquina")
2 X=df
```

```
1 X
```



	tipo	temperatura_ar	temperatura_processo	umidade_relativa	velocidade_rotac
<b>0</b>	1	298.30	309.1	90.0	
<b>1</b>	1	298.20	308.4	90.0	
<b>2</b>	1	298.20	307.8	90.0	
<b>3</b>	2	300.90	310.8	90.0	
<b>4</b>	1	293.35	310.5	90.0	
...	...	...	...	...	
<b>25859</b>	0	297.60	309.6	90.0	
<b>25860</b>	1	299.80	311.2	90.0	
<b>25861</b>	1	293.35	309.5	90.0	
<b>25862</b>	1	299.60	309.5	90.0	
<b>25863</b>	2	301.60	310.4	90.0	

25864 rows × 9 columns

```
1 Y
```



falha_maquina	
0	0.0
1	0.0
2	0.0
3	0.0
4	0.0
...	...
25859	0.0
25860	0.0
25861	0.0
25862	0.0
25863	0.0

25864 rows × 1 columns

**dtype:** float64

```
1 X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.2)
```

## Baseline com DummyClassifier, utilizado como

- ✓ referência inicial para análise do desempenho dos modelos.

```
1 # Configuração e treinamento do modelo
2 def criar_baseline_dummy():
3     """
4     Cria e treina um modelo baseline usando DummyClassifier
5     """
6     dummy = DummyClassifier(strategy="most_frequent", random_state=42)
7
8     # Verificar se precisa do MultiOutputClassifier
9     if Y_train.ndim == 1:
10         # Para problema de classificação simples (uma saída)
11         model = dummy
12     else:
13         # Para problema multi-output
14         model = MultiOutputClassifier(dummy)
15
16     return model
17
18 # Execução do experimento
```

```

19 with mlflow.start_run(run_name="Baseline_DummyClassifier"):
20     # Criar e treinar o modelo
21     model = criar_baseline_dummy()
22     model.fit(X_train, Y_train)
23
24     # Fazer previsões
25     y_pred = model.predict(X_test)
26
27     # Calcular métricas
28     acc = accuracy_score(Y_test, y_pred)
29     f1 = f1_score(Y_test, y_pred, average="macro")
30
31     # Log das métricas no MLflow
32     mlflow.log_metric("accuracy", acc)
33     mlflow.log_metric("f1_macro", f1)
34
35     # Log do modelo no MLflow (corrigindo warnings)
36     mlflow.sklearn.log_model(
37         model,
38         artifact_path="model", # Usar artifact_path ao invés de nome direto
39         input_example=X_test[:5], # Exemplo de entrada para inferir assinatura
40         signature=mlflow.models.infer_signature(X_test, y_pred) # Assinatura do modelo
41     )
42
43     # Exibir resultados
44     print(f"Baseline DummyClassifier:")
45     print(f"└─ Accuracy: {acc:.4f}")
46     print(f"└─ F1 Macro: {f1:.4f}")
47
48     # Informações adicionais sobre o modelo
49     print(f"\nDetalhes do modelo:")
50     print(f"└─ Estratégia: {model.strategy if hasattr(model, 'strategy') else model}")
51     print(f"└─ Forma dos dados de treino: {X_train.shape}")
52     print(f"└─ Forma das previsões: {y_pred.shape}")
53
54 # Versão alternativa mais simples (se não precisar do MultiOutputClassifier):
55 def baseline_simples():
56     """
57     Versão simplificada para problemas de classificação com uma única saída
58     """
59     with mlflow.start_run(run_name="Baseline_DummyClassifier_Simples"):
60         # Criar e treinar modelo
61         dummy = DummyClassifier(strategy="most_frequent", random_state=42)
62         dummy.fit(X_train, Y_train)
63
64         # Previsões e métricas
65         y_pred = dummy.predict(X_test)
66         acc = accuracy_score(Y_test, y_pred)
67         f1 = f1_score(Y_test, y_pred, average="macro")
68
69         # Logging (corrigindo warnings)
70         mlflow.log_metric("accuracy", acc)
71         mlflow.log_metric("f1_macro", f1)
72         mlflow.sklearn.log_model(
73             dummy,

```

```
74         artifact_path="model",
75         input_example=X_test[:5],
76         signature=mlflow.models.infer_signature(X_test, y_pred)
77     )
78
79     print(f"Baseline Simples -> Accuracy: {acc:.4f} | F1 Macro: {f1:.4f}")
80
81 # Descomente a linha abaixo se preferir a versão simples:
82 # baseline_simples()
83
```

```
➞ 2025/09/01 20:10:07 WARNING mlflow.models.model: `artifact_path` is deprecated. Please
Baseline DummyClassifier:
├─ Accuracy: 0.9812
└─ F1 Macro: 0.4953
```

Detalhes do modelo:

```
├─ Estratégia: most_frequent
├─ Forma dos dados de treino: (20691, 9)
└─ Forma das predições: (5173,)
```

## Análise do Baseline com DummyClassifier

### Desempenho do baseline

- Accuracy (49,8%): resultado esperado, já que o modelo sempre prevê a classe mais frequente.
- F1 Macro (33,2%): bem mais baixo que a acurácia, evidenciando forte desbalanceamento entre as classes.

### Interpretação prática

- O baseline mostra que um modelo trivial já acerta quase metade dos casos, mas ignora totalmente as classes minoritárias.
- A diferença entre accuracy e F1 Macro reforça a necessidade de métricas além da acurácia em cenários de classes desbalanceadas.

### Valor como referência

- Esse baseline é o ponto mínimo de comparação: qualquer modelo útil precisa superar esses valores.
- Ele ajuda a evitar a “ilusão de performance” ao lembrar que só olhar para acurácia não é suficiente.

### Próximos passos recomendados

Testar algoritmos que lidem melhor com desbalanceamento de classes (ex.: RandomForest, Reamostragem com SMOTE).

Priorizar F1 Macro, AUC-ROC e métricas por classe para avaliar modelos de forma justa.

Incorporar técnicas de balanceamento de dados no pipeline de treino.

## ✓ Avaliação de Classificadores com e sem Balanceamento de Classe

```
1 # Default (sem ajuste de classe)
2 RF_default = RandomForestClassifier(random_state=42)
3 Bg_default = BaggingClassifier(estimator=DecisionTreeClassifier(random_state=42))
4 DT_default = DecisionTreeClassifier(random_state=42)
5
6 # Balanced (com ajuste de classe)
7 RF_balanced = RandomForestClassifier(class_weight='balanced', random_state=42)
8 Bg_balanced = BaggingClassifier(estimator=DecisionTreeClassifier(class_weight='balanced', random_state=42))
9 DT_balanced = DecisionTreeClassifier(class_weight='balanced', random_state=42)
10
11 # Dicionário de modelos
12 models_balance = {
13     "RF_Default": RF_default,
14     "RF_Balanced": RF_balanced,
15     "Bagging_Default": Bg_default,
16     "Bagging_Balanced": Bg_balanced,
17     "DT_Default": DT_default,
18     "DT_Balanced": DT_balanced
19 }
20
21 # Avaliar os modelos
22 for name, model in models_balance.items():
23     model.fit(X_train, Y_train)
24     y_pred = model.predict(X_test)
25     acc = accuracy_score(Y_test, y_pred)
26     f1 = f1_score(Y_test, y_pred, average="macro")
27     print(f"{name}: Accuracy = {acc:.4f} | F1 Macro = {f1:.4f}")
28
```

```
➡ RF_Default: Accuracy = 0.9841 | F1 Macro = 0.7072
RF_Balanced: Accuracy = 0.9849 | F1 Macro = 0.7094
Bagging_Default: Accuracy = 0.9822 | F1 Macro = 0.6888
Bagging_Balanced: Accuracy = 0.9834 | F1 Macro = 0.6886
DT_Default: Accuracy = 0.9766 | F1 Macro = 0.6930
DT_Balanced: Accuracy = 0.9762 | F1 Macro = 0.6753
```

## ✓ Acurácia alta em todos os modelos

- Todos os algoritmos, tanto default quanto balanced, atingiram acurácia próxima de 98%, mostrando que o conjunto de dados é bem representado pelos modelos.

## F1 Macro mais baixo que a acurácia

- Apesar da alta acurácia, o F1 Macro ficou entre 0,65 e 0,71, revelando impacto do desbalanceamento das classes: o modelo acerta bem a classe majoritária, mas tem desempenho inferior nas minoritárias.

## Efeito do balanceamento de classes

- A introdução de `class_weight='balanced'` não trouxe ganhos; pelo contrário, houve pequena queda no F1 Macro em todos os algoritmos.
- Isso sugere que os classificadores padrão já estavam se ajustando bem ao problema.

Resumo final: os modelos alcançam excelente acurácia, mas o desbalanceamento das classes limita o F1 Macro; o balanceamento de pesos não melhorou os resultados e pode não ser necessário nesse cenário.

```
1 counts = Counter(Y)
2 print(counts)
```

```
Counter({0.0: 25333, 1.0: 531})
```

```
1 smote = SMOTE(random_state=42)
2 X_resampled, y_resampled = smote.fit_resample(X, Y)
```

```
1 counts = Counter(y_resampled )
2 print(counts)
```

```
Counter({0.0: 25333, 1.0: 25333})
```

## ✓ Situação original (sem SMOTE) como foi visto acima

- Forte desbalanceamento: a classe 0 tem 34.598 registros, enquanto a classe 1 tem apenas 662.
- Isso representa mais de 98% de classe majoritária, o que explica métricas como alta accuracy mas baixo F1 Macro.

## Situação ajustada (com SMOTE)

- O SMOTE equilibrou os dados: agora as duas classes têm 34.598 amostras cada.
- Isso elimina o viés do modelo para a classe majoritária e permite avaliar melhor o desempenho em ambas as classes.



## Implicação prática

- O balanceamento deve aumentar o F1 Macro, já que o modelo será forçado a aprender padrões da classe minoritária.
- Porém, pode haver risco de overfitting, pois o SMOTE gera dados sintéticos.

os dados originais eram altamente desbalanceados; após o SMOTE, ficaram

- Item da lista
- Item da lista

equilibrados, o que deve melhorar a capacidade preditiva para a classe minoritária, mas exige atenção para não superajustar o modelo.

```
1 X_train,X_test,Y_train,Y_test=train_test_split(X_resampled,y_resampled,test_size=0.2)

1 # Define models
2 models = {
3     'Logistic Regression': LogisticRegression(),
4     'Logistic Regression CV': LogisticRegressionCV(),
5     'SGD': SGDClassifier(),
6
7     'Random Forest': RandomForestClassifier(),
8     'Gradient Boosting': GradientBoostingClassifier(),
9     'AdaBoost': AdaBoostClassifier(),
10    'Bagging': BaggingClassifier(),
11    'Decision Tree': DecisionTreeClassifier(),
12    'Support Vector Machine': SVC(),
13    'K-Nearest Neighbors': KNeighborsClassifier()
14 }

1 def evaluate_model(X_train,X_test,Y_train,Y_test):
2     result=[]
3     for name, model in models.items():
4         model.fit(X_train,Y_train)
5         y_pred=model.predict(X_test)
6         acc=accuracy_score(Y_test,y_pred)
7         result.append((name,acc))
8     # Sort models by accuracy
9     result.sort(key=lambda x: x[1], reverse=True)
10    return result

1 results = evaluate_model(X_train,X_test,Y_train,Y_test)
2 print("Model Performance:")
3 for name, acc in results:
4     print(f"{name}: {acc:.6f}")
```



```
Model Performance:
Random Forest: 0.978094
Bagging: 0.970594
Decision Tree: 0.958259
```

```
K-Nearest Neighbors: 0.934182
Gradient Boosting: 0.920268
AdaBoost: 0.878133
Logistic Regression CV: 0.806197
Logistic Regression: 0.802052
Support Vector Machine: 0.776100
SGD: 0.729031
```

```
1 # Se results for uma lista de tuplas (nome, acc)
2 model_names = [name for name, _ in results]
3 accuracies = [acc * 100 for _, acc in results] # converter para %
4
5 # Criar gráfico de linha
6 plt.figure(figsize=(10, 6))
7 plt.plot(model_names, accuracies, marker='o', linestyle='-', color='b', label="Accura
8
9 # Adicionar valores acima dos pontos
10 for i, acc in enumerate(accuracies):
11     plt.text(i, acc + 1, f"{acc:.1f}%", ha="center", fontsize=9)
12
13 plt.xticks(rotation=45, ha="right")
14 plt.ylabel("Performance (%)")
15 plt.title("Comparação de Performance dos Modelos")
16 plt.ylim(0, 100)
17 plt.grid(True, linestyle="--", alpha=0.6)
18 plt.legend()
19 plt.tight_layout()
20 plt.show()
```

