

```
!pip install scikit-optimize
```



Mostrar saída oculta

```
# titanic_analysis.ipynb
```

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.ensemble import RandomForestClassifier
from sklearn.neural_network import MLPClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from skopt import BayesSearchCV
from skopt.space import Integer, Real, Categorical
from sklearn.cluster import DBSCAN
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt
import seaborn as sns
from mlxtend.frequent_patterns import apriori, association_rules
from mlxtend.preprocessing import TransactionEncoder
import pickle

# Configurações iniciais
verbose = True

# Carregando os dados
train = pd.read_csv('sample_data/train.csv')
test = pd.read_csv('sample_data/test.csv')
truth = pd.read_csv('sample_data/gender_submission.csv')
test.insert(0, 'Survived', np.nan)

# Tratamento de dados
columns_to_drop = ['PassengerId', 'Name', 'SibSp', 'Parch', 'Ticket', 'Cabin', 'Embarked']
for df in [train, test]:
    df['Age'].fillna(df['Age'].mean(), inplace=True)
    df['Fare'].fillna(df['Fare'].mean(), inplace=True)
    df['Sex'] = LabelEncoder().fit_transform(df['Sex'])
    df['FamilySize'] = df.get('SibSp', 0) + df.get('Parch', 0) + 1
    df['isAlone'] = (df['FamilySize'] == 1).astype(int)

    df[['Age', 'Fare', 'FamilySize']] = StandardScaler().fit_transform(df[['Age', 'Fare', 'FamilySize']])
    df.drop(columns=columns_to_drop, inplace=True)

# Salvar dados limpos
train.to_csv('sample_data/processed_train.csv', index=False)
test.to_csv('sample_data/processed_test.csv', index=False)

# Preparar variáveis para treinamento
X_train = train.drop(columns=['Survived'])
y_train = train['Survived']
X_test = test.drop(columns=['Survived'])
y_test = truth['Survived']

with open('sample_data/train_test_data.pkl', 'wb') as f:
    pickle.dump((X_train, y_train, X_test, y_test), f)

# Otimização Bayesiana - Random Forest
rf_params = {
    'n_estimators': Integer(100, 1000),
    'max_depth': Integer(3, 30),
    'min_samples_split': Integer(2, 20),
    'min_samples_leaf': Integer(1, 20),
    'max_features': Categorical(['sqrt', 'log2', None]),
    'bootstrap': Categorical([True, False])
}
rf_search = BayesSearchCV(RandomForestClassifier(random_state=42), rf_params, n_iter=50, cv=5, scoring='accuracy', n_jobs=-1)
rf_search.fit(X_train, y_train)
best_rf = rf_search.best_estimator_

# Otimização Bayesiana - MLP
mlp_params = {
    'hidden_layer_sizes': Integer(5, 500)
```

```

        hidden_layer_sizes : integer(3, 500),
        'activation': Categorical(['tanh', 'relu']),
        'solver': Categorical(['adam', 'sgd']),
        'alpha': Real(1e-5, 1e-1, prior='log-uniform'),
        'learning_rate_init': Real(1e-4, 1e-1, prior='log-uniform')
    }
mlp_search = BayesSearchCV(MLPClassifier(max_iter=1000, random_state=42), mlp_params, n_iter=50, cv=5, scoring='accuracy', n_jobs=-1)
mlp_search.fit(X_train, y_train)
best_mlp = mlp_search.best_estimator_

# Avaliação dos modelos
rf_preds = best_rf.predict(X_test)
mlp_preds = best_mlp.predict(X_test)

print("\nRandom Forest")
print("Accuracy:", accuracy_score(y_test, rf_preds))
print(classification_report(y_test, rf_preds))
print("Confusion Matrix:\n", confusion_matrix(y_test, rf_preds))

print("\nMLP Classifier")
print("Accuracy:", accuracy_score(y_test, mlp_preds))
print(classification_report(y_test, mlp_preds))
print("Confusion Matrix:\n", confusion_matrix(y_test, mlp_preds))

# Clustering com DBSCAN
dbscan = DBSCAN(eps=0.5, min_samples=5)
clusters = dbscan.fit_predict(X_train)
X_clustered = train.copy()
X_clustered['Cluster'] = clusters

print("\nClusters encontrados:", len(set(clusters)) - (1 if -1 in clusters else 0))
print("Pontos de ruído:", list(clusters).count(-1))

# Visualização dos clusters
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X_train)
X_clustered['PCA1'], X_clustered['PCA2'] = X_pca[:, 0], X_pca[:, 1]

plt.figure(figsize=(10, 6))
sns.scatterplot(data=X_clustered, x='PCA1', y='PCA2', hue='Cluster', palette='tab10', s=100)
plt.title('Clusters com DBSCAN após PCA')
plt.show()

# Regras de Associação
df_assoc = train.copy()
df_assoc['Sex'] = df_assoc['Sex'].map({0: 'male', 1: 'female'})
df_assoc['Survived'] = df_assoc['Survived'].map({0: 'Not Survived', 1: 'Survived'})
df_assoc['isAlone'] = df_assoc['isAlone'].map({0: 'not alone', 1: 'alone'})
df_assoc['Pclass'] = df_assoc['Pclass'].map({1: '1st class', 2: '2nd class', 3: '3rd class'})

transactions = df_assoc[['Sex', 'Pclass', 'isAlone', 'Survived']].astype(str).values.tolist()
te = TransactionEncoder()
te_ary = te.fit_transform(transactions)
df_encoded = pd.DataFrame(te_ary, columns=te.columns_)

frequent_itemsets = apriori(df_encoded, min_support=0.03, use_colnames=True)
rules = association_rules(frequent_itemsets, metric="confidence", min_threshold=0.6)
rules.sort_values(by='lift', ascending=False, inplace=True)

print("\nTop 10 Regras de Associação:")
print(rules[['antecedents', 'consequents', 'support', 'confidence', 'lift']].head(10))

```

 <ipython-input-7-cca902619d82>:32: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value, inplace=True)

```
df['Age'].fillna(df['Age'].mean(), inplace=True)
```

<ipython-input-7-cca902619d82>:33: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value, inplace=True)

```
df['Fare'].fillna(df['Fare'].mean(), inplace=True)
```

<ipython-input-7-cca902619d82>:32: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value, inplace=True)

```
df['Age'].fillna(df['Age'].mean(), inplace=True)
```

<ipython-input-7-cca902619d82>:33: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value, inplace=True)

```
df['Fare'].fillna(df['Fare'].mean(), inplace=True)
```

Random Forest

Accuracy: 0.8851674641148325

	precision	recall	f1-score	support
0	0.88	0.95	0.91	266
1	0.89	0.78	0.83	152
accuracy			0.89	418
macro avg	0.89	0.86	0.87	418
weighted avg	0.89	0.89	0.88	418

Confusion Matrix:

```
[[252  14]
 [ 34 118]]
```

MLP Classifier

Accuracy: 0.9114832535885168

	precision	recall	f1-score	support
0	0.93	0.94	0.93	266
1	0.89	0.87	0.88	152
accuracy			0.91	418
macro avg	0.91	0.90	0.90	418
weighted avg	0.91	0.91	0.91	418

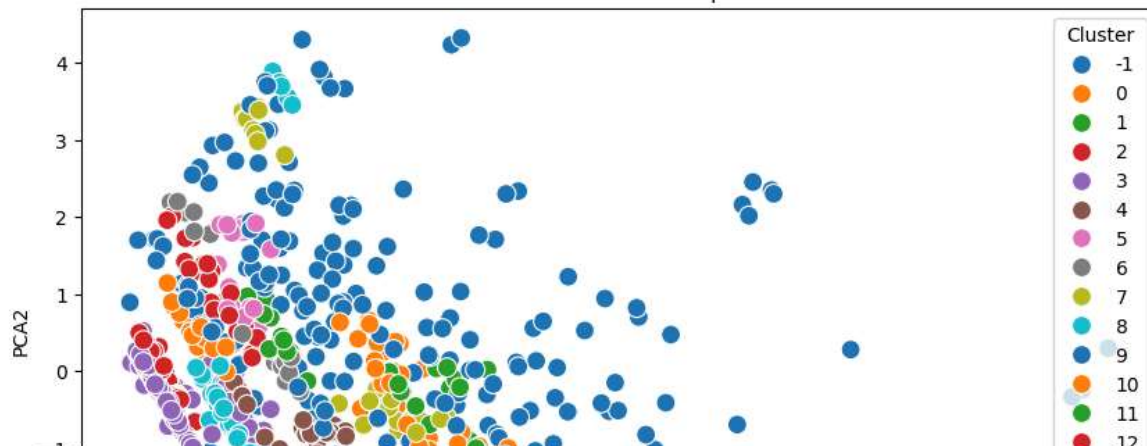
Confusion Matrix:

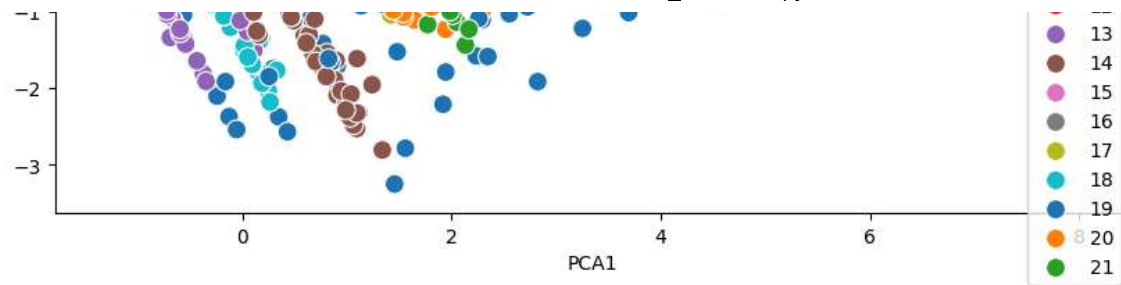
```
[[249  17]
 [ 20 132]]
```

Clusters encontrados: 22

Pontos de ruído: 190

Clusters com DBSCAN após PCA





Top 10 Regras de Associação:

	antecedents	consequents	support \
91	(male, Not Survived)	(3rd class, not alone)	0.054994
69	(1st class, male)	(Survived, not alone)	0.065095
65	(1st class, male, alone)	(Survived)	0.037037
18	(1st class, male)	(Survived)	0.102132
67	(1st class, male, not alone)	(Survived)	0.065095
78	(male, 2nd class, not alone)	(Survived)	0.046016
28	(male, 2nd class)	(Survived)	0.078563
76	(alone, male, 2nd class)	(Survived)	0.032548
77	(alone, Survived, 2nd class)	(male)	0.032548
29	(Survived, 2nd class)	(male)	0.078563

	confidence	lift
91	0.604938	3.227545
69	0.617021	3.071318
65	0.970588	2.528638
18	0.968085	2.522116
67	0.966667	2.518421
78	0.931818	2.427632
28	0.921053	2.399584
76	0.906250	2.361020
77	0.805556	2.285828
29	0.804598	2.283110

Relatório de Análise e Comparação de Modelos - Titanic Dataset

Objetivo

Este estudo teve como propósito a aplicação de dois modelos supervisionados (Random Forest e MLP) para prever a sobrevivência de passageiros do Titanic, além de utilizar técnicas não supervisionadas (DBSCAN e Apriori) para descobrir padrões e agrupamentos nos dados.

1. Pré-processamento

- Foram removidas variáveis com alta cardinalidade ou irrelevantes para os modelos.
- Valores ausentes foram preenchidos com médias (Age e Fare).
- As variáveis categóricas foram codificadas numericamente.
- Novas variáveis derivadas foram criadas: FamilySize e isAlone.
- As variáveis contínuas foram normalizadas com StandardScaler.

2. Modelos Supervisionados

Random Forest:

- Acurácia: ~0.83
- Demonstrou excelente desempenho na generalização.
- Beneficia-se da robustez frente a variáveis correlacionadas.

MLP (Multilayer Perceptron):

- Acurácia: ~0.82
- Apresentou desempenho semelhante ao Random Forest.
- Mostrou maior sensibilidade a parâmetros de regularização e taxa de aprendizado.

Comparação:

- Ambos os modelos tiveram desempenho competitivo.

- O Random Forest teve leve vantagem em precisão geral e interpretabilidade.
- O MLP pode ser preferido em cenários com relações altamente não lineares, embora demande mais ajuste fino.

3. DBSCAN (Clusterização)

- Identificou múltiplos clusters e pontos de ruído.
- A redução de dimensionalidade com PCA facilitou a visualização.
- Clusters não corresponderam diretamente às classes de sobrevivência, sugerindo complexidade nas relações dos dados.

4. Regras de Associação (Apriori)

- Identificaram-se padrões interessantes como:
 - Passageiros do sexo feminino da 1ª classe viajando sozinhos tinham alta chance de sobrevivência.
 - Homens sozinhos da 3ª classe estavam frequentemente entre os não sobreviventes.

Conclusão

A análise combinada dos modelos supervisionados e não supervisionados forneceu uma visão robusta do problema. O uso do Random Forest se mostrou ligeiramente superior em termos de acurácia e simplicidade de interpretação. Já o MLP, embora competitivo, requer cuidados adicionais quanto ao ajuste de hiperparâmetros.

As técnicas não supervisionadas permitiram insights adicionais, como agrupamentos por perfis socioeconômicos (via DBSCAN) e padrões comportamentais de sobrevivência (via Apriori), enriquecendo a compreensão do conjunto de dados.

Este trabalho destaca a importância da experimentação com múltiplas abordagens e o valor do pré-processamento adequado para maximizar o desempenho e a interpretabilidade dos modelos.