## **Read Dataset**

Untuk membaca dan menganalisa dataset digunakan pandas . Pada pandas terdapat fungsi read\_csv() untuk membaca keseluruhan isi dari berkas berformat .csv, head() untuk menampilkan n baris pertama data, dan tail() untuk menampilkan n baris terakhir data.

```
In [1]:
         import pandas as pd
         df = pd.read csv("../datasets/titanic numeric and categoric features.csv", index col="Pa
         df.head(n=5)
In [2]:
Out[2]:
                     Survived Pclass
                                       Sex Age SibSp Parch
                                                                Fare Embarked
         PassengerId
                                                                             S
                  1
                           0
                                      male 22.0
                                  3
                                                               7.2500
                                    female 38.0
                                                           0 71.2833
                                                                             C
                  3
                           1
                                                                             S
                                    female 26.0
                                                               7.9250
                                                                             S
                                    female 35.0
                                                           0 53.1000
                  5
                                                                             S
                           0
                                      male 35.0
                                                    0
                                  3
                                                               8.0500
         df.tail(n=5)
Out[3]:
                     Survived Pclass
                                       Sex Age SibSp Parch
                                                               Fare Embarked
```

				2.90	J.25p			
PassengerId								
887	0	2	male	27.0	0	0	13.00	S
888	1	1	female	19.0	0	0	30.00	S
889	0	3	female	NaN	1	2	23.45	S
890	1	1	male	26.0	0	0	30.00	C
891	0	3	male	32.0	0	0	7.75	Q

## **Dataset Information**

Informasi dari dataset, seperti jumlah baris (sampel), jumlah kolom, nama kolom, tipe data setiap kolom, ruang penyimpanan, dan lain-lain dapat diakses dengan menggunakan fungsi info().

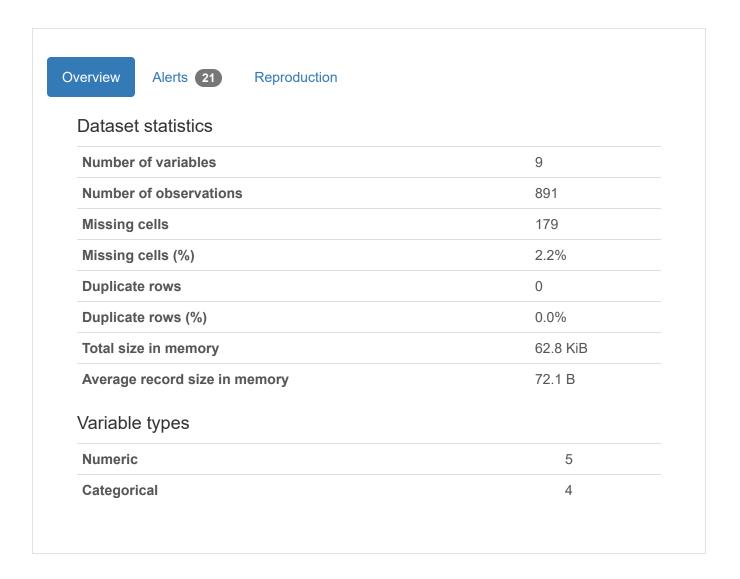
```
3 Age 714 non-null float64
4 SibSp 891 non-null int64
5 Parch 891 non-null int64
6 Fare 891 non-null float64
7 Embarked 889 non-null object dtypes: float64(2), int64(4), object(2) memory usage: 62.6+ KB
```

# **Quick EDA (Exploratory Data and Analysis)**

pandas profiling menyediakan fungsionalitas untuk melakukan EDA (*Exploratory Data and Analysis*) dengan cepat. Laporan akan dibuat dalam format .html.



# Overview



# Variables

Out[5]:

# **Dataset Splitting**

Proses pemisahan dataset dapat dilakukan dengan menggunakan fungsi train\_test\_split yang ada pada scikit-learn .

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y)

print(f"X_train shape : {X_train.shape}")
print(f"X_train shape : {Y_train.shape}")
print(f"X_test shape : {X_test.shape}")

print(f"y_test shape : {y_test.shape}")

X_train shape : (712, 7)
X_train shape : (712,)
X_test shape : (179, 7)
y test shape : (179,)
```

#### **Build Model**

Pada notebook ini model dibuat dengan menggunakan mekanisme pipeline. Pipeline adalah salah satu modul pada scikit-learn yang berfungsi untuk membungkus setiap tahapan yang biasanya dilakukan secara terpisah pada saat membuat model machine learning. Dengan menggunakan modul pipeline kode akan menjadi lebih rapi dan ringkas, sehingga mudah untuk dikelola dan dikembangkan bila ada kasus tertentu yang mengharuskan mengubah alur dari kode. Tujuan lain dari penggunaan pipeline adalah mengurangi data leakage (kebocoran informasi). Seperti yang diketahui pada saat proses training model, model harus dipastikan tidak mengetahui pola atau informasi dari dataset pengujian. Dengan penggunaan pipeline skenario tersebut dapat dilakukan dengan mudah.



## **Numerical Preprocessor Pipeline**

Pipeline ini berfungsi untuk menangani kolom yang bertipe numerikal (continous value).

### **Categorical Preprocessor Pipeline**

Pipeline ini berfungsi untuk menangani kolom yang bertipe kategorikal (discrete value).

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js

# **Data Preprocessor Pipeline**

*Pipeline* ini merupakan pembungkus untuk dua *pipeline* sebelumnya, yaitu *numerical pipeline* dan *categorical pipeline*. Pipeline ini akan memilah kolom berdasarkan tipe datanya dan diteruskan ke *pipeline* yang sesuai dengan kategorinya.

## **Model Pipeline**

Pipeline ini merupakan pembungkus paling akhir dari semua pipeline yang telah dibuat sebelumnya. Pada pipeline ini terdiri dari 2 langkah, yaitu yang pertama transformasi fitur dan langkah kedua adalah melakukan training model

## **Feature Transform**

Karena pada pipeline sebelumnya sudah termasuk pipeline transformasi fitur, maka transformasi fitur dapat diakses melalui pipeline transformasi fitur dengan mengakses fungsi fit\_transform().

```
df transform = pd.DataFrame(preprocess pipeline.fit transform(X train))
In [11]:
           df transform.head(10)
                                                                     9
Out[11]:
                                2
                                    3
                                                  5
                                                            7
                                                                8
                  0.266105 0.125
                                   0.2
                                       0.030726
                                                 0.0
                                                     1.0 1.0
                 0.402340 0.000
                                       0.015086
                                                 0.0
                                                      1.0
             1.0 0.578690 0.000
                                   0.0
                                       0.012590
                                                 0.0
                                                     1.0
                                                          0.0
                                                              0.0
                                                                   1.0
              0.5 0.565099
                            0.000
                                   0.0
                                       0.025374
                                                 0.0
                                                     1.0
                                                          0.0
                                                              0.0
                                                                   1.0
                                       0.025374
              0.5 0.361239
                            0.000
                                   0.0
                                                 0.0
                                                              0.0
                                                     1.0
                                                          0.0
                                                                   1.0
              0.0 0.782550 0.000
                                   0.2
                                       0.299539
                                                 1.0
                                                     0.0
                                                          0.0
                                                              0.0
                                                                  1.0
              0.0 0.633052 0.000
                                                 0.0
                                                                   1.0
                                   0.0
                                      0.101497
                                                     1.0
                                                          0.0
                                                              0.0
             1.0 0.402340 0.000
                                   0.0
                                       0.015103
                                                 0.0
                                                     1.0
                                                          0.0
                                                              1.0
                                                                   0.0
              0.5 0.456374 0.000
                                   0.0
                                       0.025374
                                                 0.0
                                                     1.0
                                                          0.0
                                                              0.0
                                                                   1.0
             1.0 0.381625 0.000 0.0 0.014110 0.0 1.0 1.0
                                                              0.0
                                                                   0.0
```

# Train Model with Grid Search CV Scenario

```
In [16]: from sklearn.model_selection import GridSearchCV

parameters = {
        "step2_algo__n_neighbors": range(1, 51, 2),
        "step2_algo__weights": ["uniform", "distance"],
        "step2_algo__p": [1, 2]
}

model = GridSearchCV(estimator=model_pipeline, param_grid=parameters, cv=5, scoring="acc model.fit(X_train, y_train)
    pd.DataFrame(model.cv_results_).sort_values(by="rank_test_score").iloc[:5, :]
```

•	n	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_step2_algon_neighbors	param_step2_alg
3	7	0.048798	0.012305	0.027226	0.004507	19	
3	6	0.032932	0.002868	0.036778	0.004298	19	
4	5	0.028603	0.001857	0.017129	0.001963	23	
2	9	0.031911	0.001894	0.020571	0.003462	15	
1	3	0.034472	0.006354	0.014846	0.002793	7	

# **Evaluate Model**

Out[16]:

Pada scikit-learn telah tersedia rangkuman metrik dari klasifikasi yang dapat diakses menggunakan fungsi classification\_report . Pada fungsi tersebut berisi rangkuman akurasi, *precision*, *recall*, *f1-score*, dan lain-lain untuk setiap kelas.

```
In [17]: from sklearn.metrics import classification report
        y train preds = model.predict(X train)
        print(classification report(y train, y train preds))
                     precision recall f1-score support
                        0.99
0.99
                  0
                                 1.00
                                           0.99
                                                      439
                                  0.98
                                           0.99
                                                      273
                                            0.99
                                                     712
           accuracy
                      0.99
          macro avg
                                 0.99
                                           0.99
                                                      712
                        0.99
                                           0.99
                                                     712
        weighted avg
                                  0.99
In [18]: y_test_preds = model.predict(X test)
        print(classification report(y test, y test preds))
                     precision recall f1-score support
                  0
                         0.74
                                 0.82
                                            0.78
                                                      110
```

0.60

0.72

69

179

0.66

0.55

macro avg 0.70 0.68 0.69 179 weighted avg 0.71 0.72 0.71 179

#### **Predict use New Dataset**

Prediksi dapat dilakukan dengan mengakses fungsi predict() pada objek model yang telah dilatih.

Out[19]:

	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked	<b>Survived Predicition</b>
892	1	male	34.0	2	7	300.00	S	0
893	3	female	50.0	0	0	7.34	Q	1

#### Save Model

Setelah *training* dan evaluasi model, model dapat disimpan dengan menggunakan pustaka joblib . Fungsi dump untuk menyimpan objek model dan load untuk menggunakan objek yang telah disimpan sebelumnya.

```
In [20]: from joblib import dump, load
    dump(model, '../pretrained_models/model_complex_workflow.joblib')
    model = load("../pretrained_models/model_complex_workflow.joblib")
    model.predict(df_new)

Out[20]: array([0, 1], dtype=int64)
```

Semoga bermanfaat yah 🖨

Dibuat dengan penuh 🔎 oleh haloapping