

Pandas library -Data Preproceeing

Aim:

To write the Python program to understand and perform data preprocessing on the given dataset.

Algorithm:

1. Load the dataset and inspect its structure using .info() and .head().
2. Identify missing values in categorical and numerical columns.
3. Handle missing values using mode (for categorical), median (for age), and mean (for salary).
4. Encode categorical variables using one-hot encoding.
5. Concatenate encoded columns with the original dataset.
6. Prepare the cleaned dataset for further analysis or modeling.

Program:

```
[5]: import numpy as np
import pandas as pd
df=pd.read_csv(r"C:\Users\siddesh\Downloads\pre_process_datasample.csv")
df
```

```
[5]:
```

	Country	Age	Salary	Purchased
0	France	44.0	72000.0	No
1	Spain	27.0	48000.0	Yes
2	Germany	30.0	54000.0	No
3	Spain	38.0	61000.0	No
4	Germany	40.0	NaN	Yes
5	France	35.0	58000.0	Yes
6	Spain	NaN	52000.0	No
7	France	48.0	79000.0	Yes
8	Germany	50.0	83000.0	No
9	France	37.0	67000.0	Yes

```
[6]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype  
---  --
 0   Country     10 non-null    object  
 1   Age         9 non-null     float64 
 2   Salary      9 non-null     float64 
 3   Purchased   10 non-null    object  
dtypes: float64(2), object(2)
memory usage: 452.0+ bytes
```

```
[7]: df.Country.mode()
```

```
[7]: 0    France
Name: Country, dtype: object
```

```
•[12]: df['Country'] = df['Country'].fillna(df['Country'].mode()[0])
df['Age'] = df['Age'].fillna(df['Age'].median())
df['Salary'] = df['Salary'].fillna(round(df['Salary'].mean()))
df
```

```
[12]:
```

	Country	Age	Salary	Purchased
0	France	44.0	72000.0	No
1	Spain	27.0	48000.0	Yes
2	Germany	30.0	54000.0	No
3	Spain	38.0	61000.0	No
4	Germany	40.0	63778.0	Yes
5	France	35.0	58000.0	Yes
6	Spain	38.0	52000.0	No
7	France	48.0	79000.0	Yes
8	Germany	50.0	83000.0	No
9	France	37.0	67000.0	Yes

```
[9]: pd.get_dummies(df.Country)
```

```
[9]:
```

	France	Germany	Spain
0	True	False	False
1	False	False	True
2	False	True	False
3	False	False	True
4	False	True	False
5	True	False	False
6	False	False	True
7	True	False	False
8	False	True	False
9	True	False	False

```
[10]: updated_dataset=pd.concat([pd.get_dummies(df.Country),df.iloc[:,[1,2,3]]],axis=1)
updated_dataset
```

```
[10]:
```

	France	Germany	Spain	Age	Salary	Purchased
0	True	False	False	44.0	72000.0	No
1	False	False	True	27.0	48000.0	Yes
2	False	True	False	30.0	54000.0	No
3	False	False	True	38.0	61000.0	No
4	False	True	False	40.0	63778.0	Yes
5	True	False	False	35.0	58000.0	Yes
6	False	False	True	38.0	52000.0	No
7	True	False	False	48.0	79000.0	Yes
8	False	True	False	50.0	83000.0	No
9	True	False	False	37.0	67000.0	Yes

```
[11]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Country     10 non-null     object
1   Age         10 non-null     float64
2   Salary      10 non-null     float64
3   Purchased   10 non-null     object
dtypes: float64(2), object(2)
memory usage: 452.0+ bytes
```

Result:

Thus, the Python program is executed successfully for preprocessing the given dataset.