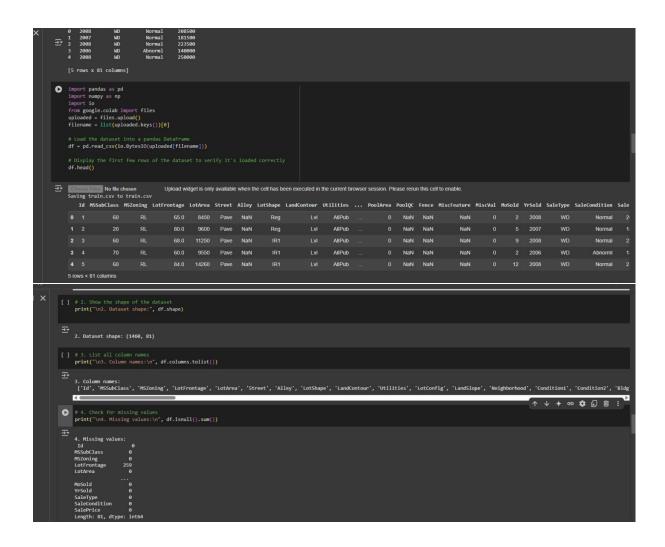
Name = PAYAL S BAJANTRI

PRN = 202401120064

**BATCH = CS8-68** 

## Data set =

https://drive.google.com/file/d/1c26XS7KvN\_IOnZwtU
JrUZQRzZTCKzisz/view?usp=drivesdk



```
# 5. Display data types of each column
print("\n5. Data types:\n", df.dtypes)
             5. Data types:
Id
MSSubClass
MSZoning
LotFrontage
LotArea
              MoSold int64
YrSold int64
SaleType object
SaleCondition object
SalePrice int64
Length: 81, dtype: object
0
            6. Statistical summary:

Count 1460.000000 1460.000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.00000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.00000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.0000000 1201.000000000 1201.00000000 1201.00000000 1201.00000
                                                                                                                                                                                      YearBuilt YearRemoddd
1460.000000 1460.000000
1971.267888 1984.867573
30.202904 20.645407
1872.000000 1950.000000
1954.000000 1967.000000
1973.000000 1964.000000
2010.000000 2004.0000000
2010.0000000 2010.000000
                                       OverallCond
1460.000000
5.575342
1.112799
1.000000
5.000000
6.000000
9.000000
                                                                                                                                                                                                                                        BsmtFinSF1
1460.000000
443.639726
456.098091
0.000000
0.000000
383.500000
712.250000
5644.000000
                                                                                                                                                                                           MasVnrArea
1452.000000
103.685262
181.066207
0.000000
0.000000
0.000000
              count
mean
std
min
25%
50%
75%
     # 7. Count of unique values in each column
print("\n7. Unique values per column:\n", df.nunique())
      7. Unique values per column:
Id 1460
      MSSubClass
MSZoning
LotFrontage
                                                                                       110
       MoSold
         YrSold
       SaleType
SaleCondition
       SalePrice 663
Length: 81, dtype: int64
          numeric_cols = df.select_dtypes(include='number').columns
         print("\n8. Correlation matrix (numerical features only):\n", df[numeric_cols].corr())
   YearBuilt
YearRemodAdd
MasVnrArea
BsmtFinSF1
                                                                       -0.012713
-0.021998
                                                                                                                                  0.027850
0.040581
                                                                                                                                                                                              0.123349 0.014228
0.088866 0.013788
                                                                                                                                                                                                                                                                                                        0.572323
0.550684
                                                                        -0.050298
-0.005024
-0.005968
                                                                                                                                0.022936
-0.069836
-0.065649
                                                                                                                                                                                              0.193458
0.233633
                                                                                                                                                                                                                                          0.104160
0.214103
0.111170
                                                                                                                                                                                                                                                                                                       0.411876
0.239666
-0.059119
         BsmtFinSF2
                                                                                                                                                                                               0.049900
                                                                                                                                                                                             0.132644 -0.002618
0.392075 0.260833
0.457181 0.299475
0.080177 0.050986
       BsmtUnfSF
TotalBsmtSF
                                                                         -0.007940
-0.015415
                                                                                                                                -0.140759
-0.238518
                                                                                                                                                                                                                                                                                                        0.308159
0.537808
         1stFlrSF
                                                                          0.010496
0.005590
                                                                                                                                -0.251758
0.307886
                                                                                                                                                                                                                                                                                                         0.476224
```

```
[ ] # 11. Create a new column 'Price_per_SqFt'
if 'GrLivArea' in df.columns and 'salePrice' in df.columns:
    df['Price_per_SqFt'] = df['salePrice'] / df['GrLivArea']
    print('\n1. Added 'Price_per_SqFt' column.")
else:
    print('\n1. Columns 'GrLivArea' or 'SalePrice' missing.")

11. Added 'Price_per_SqFt' column.

[ ] # 12. Find the property with the highest 'SalePrice'
if 'SalePrice' in df.columns:
    print('\n12. Property with highest SalePrice'\n', df.loc[df['SalePrice'].idomax()])
else:
    print('\n12. 'SalePrice' column not found.")

22. Property with highest SalePrice:
    1d 602
MSSubClass 60
MSZoning RL
Lotfrontage 104.8
Lotfrontage 104.8
Lotfrontage 21535

YrSold 2007
SaleCype MD
SaleCype MD
SaleCype MD
SaleCype MD
SaleCype TSALOSAMON
SaleCype TSALOSAMON
Normal
SalePrice 7550000
Price_per_Sqft 174.0304001
Name: 601, Length: 82, dtype: object
```

```
# 13. Find the property with the lowest 'SalePrice'
if 'SalePrice' in df.columns:
                         print("\n13. Property with lowest SalePrice:\n", df.loc[df['SalePrice'].idxmin()])
else:
                      13. Property with lowest SalePrice:
Id 496
MSSubClass 30
MSZoning C (all)
LotFrontage 60.0
LotArea 7879
                                                                                                       ...
2009
WD
Abnorml
34900
                       Yrsold 2009
SaleType WD
SaleCondition Abnorm1
SalePrice 34900
Price_per_SqFt 48.472222
Name: 495, Length: 82, dtype: object
   [ ] # 14. Average 'SalePrice' for houses built after 2000
    if 'YearBuilt' in df.columns and 'SalePrice' in df.columns:
        print("\n14. Average SalePrice for houses built after 2000:\n", df[df['YearBuilt'] > 2000]['SalePrice'].mean())
                                        print("\n14. Required columns not found.")
    14. Average SalePrice for houses built after 2000: 244527.4587912088
               # 15. Number of houses with 'GarageArea' greater than 500 if 'GarageArea' in df.column: print('Nn5. Houses with GarageArea > 500:\n", df[df['GarageArea'] > 500].shape[0]) else:
[ ] # 16. Median 'LotArea'
if 'LotArea' in df.columns:
    print("\n16. Median LotArea:", df['LotArea'].median())

# 17. Standard deviation of 'SalePrice'
if 'SalePrice' in df.columns:
    print('N17. SalePrice standard deviation:", df['SalePrice'].std())
else:
    mint for the standard deviation of the standard 
 [] # 18. Normalize 'SalePrice' column
if 'SalePrice' in df.Columns:
df['SalePrice | Normalized'] = (df['SalePrice'] - df['SalePrice'].mean()) / df['SalePrice'].std()
print("\n18. Normalized 'SalePrice' column added.")
 18. Normalized 'SalePrice' column added
```