# Machine Learning

① Data. pre-processing : load a dataset, handle
missing values, encode categorical data,
and normalize / standardize features.


```
Import pandas as pd
Import numpy as np.
from sklearn. Impute Import SimpleImputer
from sklearn. preprocessing Import One Hot
 Encoder, Standard Scaler.
from sklearn. Compose Import ColumnTrans fo
 -rmer
from sklearn. pipeline. Import pipeline.
data = pd. Dataframe ({ 'Age': [25, np.nan,
 35, 40, 29], 'Salary': [50000, 6000, np.nan
 80000, 52000], 'Department': ['Sales',
 'Engineering', 'HP', np.nan, 'Sales'],
 'purchased': ['yes', 'No', 'yes', 'No', 'yes']
 })
paint ("Original dataset: \n")
paint (data)
X = data.drop ('purchased', axis = 1)
y = data ['purchased']
numerical _cols = X. select-dtypes (Include =
 ['int64', 'float64']). columns. tolist ()
categorical - cols = X. select-dtypes (Include =
     = ['object']). columns. tolist()
numerical _pipeline = pipeline (steps = [
    ('imputer', SimpleImputer (strategy = 'mean')),
   ('scaler', StandardScaler ())
 ])
```

Computation + One-hat encoding)
categorical -pipeline = pipeline(steps =[
   ('Imputer', Simple·Imputer (strategy
   = "most _ frequent")),
   ('Encoder', onehat encoder.(handle- unknown
      = "Ignore"))
])

preprocessor = Column Transformer (trans
-formers = [
   ("num", numerical -pipeline, numerical
   - cols),
   ('cat', categorical -pipeline, categorical
   - cols)])

x - processed = preprocessor . fit - trans
- farm (x)

print ("Inprocessed features (after
   handling missing values, encoding,
   and Scaling):\n")
print(x_processed . toarray() if
   has attr (x - processed, "toarray")
   else x-processed .

names  encoded - feature - names.
preprocessor •named _ trans formers
   - ['cat'] ['Encoder'].get - feature
-name - out ( categorical - cols)
all . features _ names = numerical - cols
+ Encoded - features -nams . tolist()
Print ("In processed features Names :")
print ( all- feature - names)

**OUTPUT :**

**Original Dataset.**

|   | Age | Salary | Department | purchased |
|---|-----|--------|------------|-----------|
| 0 | 25.0 | 50000 | Sales | Yes. |
| 1 | Nan. | 60000.0 | Enginering | No. |
| 3 | 35.0 | Nan | HR. | Yes. |

**Proceeed feature Names.**

['Age', 'Salary', 'Department_Enginering', 'department_HR', 'Department_Sales'].

**proceeed features (after handling missing values, encoding & scaling):**

```
[[-1.41421356   -1.8321596    1.   0.  0.]
 [   0.            0.          0.   1.  0.]
 [ 0.70710678   -0.50709255   0.   0.  1.]
 [ 1.41421356    1.5212 7766   1.   0.  0.]
 [-0.70710678   -0.16903085   1.   0.  0.]]
```