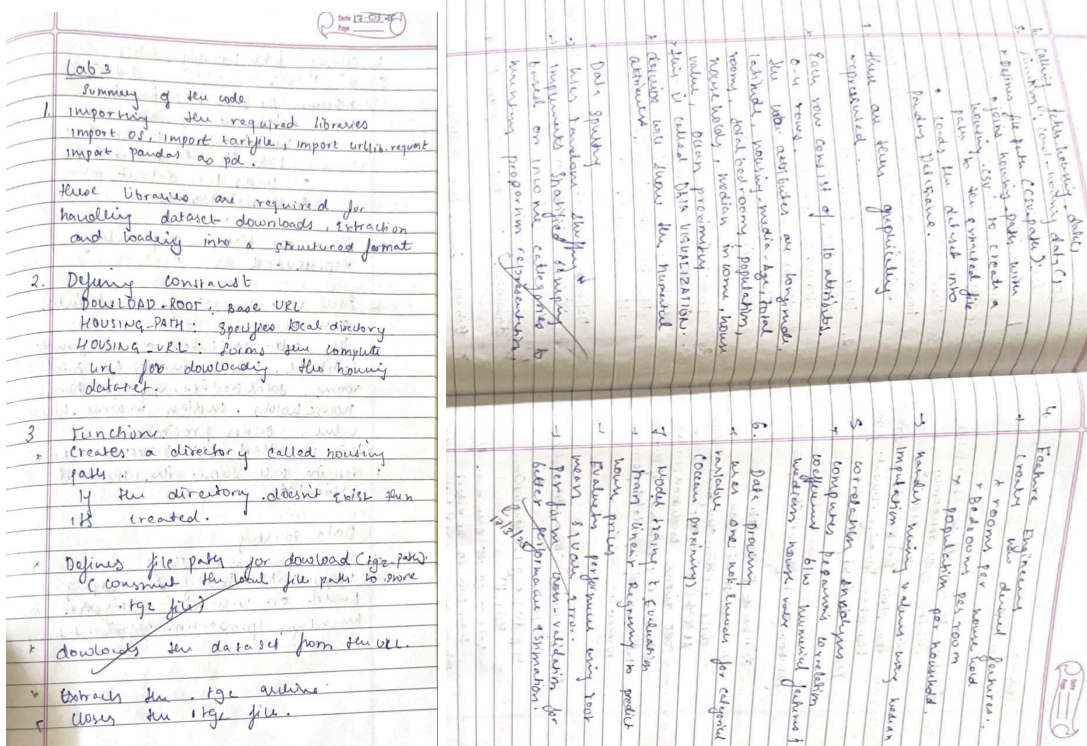# Demonstrate various data pre-processing techniques for a given dataset

**Code:**

```python
import pandas as pd
def load_housing_data():
  csv_path = r'/content/housing.csv'
  return pd.read_csv(csv_path)


housing = load_housing_data()
```

```python
import numpy as np
def split_train_test(data, test_ratio):
  shuffled_indices = np.random.permutation(len(data))
  test_set_size = int(len(data) * test_ratio)
  test_indices = shuffled_indices[:test_set_size]
  train_indices = shuffled_indices[test_set_size:]
  return data.iloc[train_indices], data.iloc[test_indices]
train_set, test_set = split_train_test(housing, 0.2)
```

```python
from sklearn.model_selection import StratifiedShuffleSplit
```

```python
split = StratifiedShuffleSplit(n_splits=1, test_size=0.2, random_state=42)
for train_index, test_index in split.split(housing,
housing["income_cat"]):
  strat_train_set = housing.loc[train_index]
  strat_test_set = housing.loc[test_index]

for set_ in (strat_train_set, strat_test_set):
  set_.drop("income_cat", axis=1, inplace=True)
X = imputer.transform(housing_num)
housing_tr = pd.DataFrame(X, columns=housing_num.columns,
index=housing_num.index)
housing_cat = housing[["ocean_proximity"]]
housing_cat.head(10)
from sklearn.linear_model import LinearRegression
lin_reg = LinearRegression()
lin_reg.fit(housing_prepared, housing_labels)
some_data = housing.iloc[:5]
some_labels = housing_labels.iloc[:5]
some_data_prepared = full_pipeline.transform(some_data)
print("Predictions:", lin_reg.predict(some_data_prepared))
print("Labels:", list(some_labels))

from sklearn.metrics import mean_squared_error
housing_predictions = lin_reg.predict(housing_prepared)
lin_mse = mean_squared_error(housing_labels, housing_predictions)
lin_rmse = np.sqrt(lin_mse)
print(f"lin_rmse: {lin_rmse}")
```

Output:
```
[[0. 1. 0. 0. 0.]
 [0. 0. 0. 0. 1.]
 [0. 1. 0. 0. 0.]
 ...
 [1. 0. 0. 0. 0.]
 [1. 0. 0. 0. 0.]
 [0. 1. 0. 0. 0.]]
[[1.]
 [4.]
 [1.]
 [4.]
```

```
[0.]
[3.]
[0.]
[0.]
[0.]
[0.]]
```

```
[array(['<1H OCEAN', 'INLAND', 'ISLAND', 'NEAR BAY', 'NEAR OCEAN'],
       dtype=object)]
```

**ocean_proximity**

| | |
|---|---|
| **12655** | INLAND |
| **15502** | NEAR OCEAN |
| **2908** | INLAND |
| **14053** | NEAR OCEAN |
| **20496** | <1H OCEAN |
| **1481** | NEAR BAY |
| **18125** | <1H OCEAN |
| **5830** | <1H OCEAN |
| **17989** | <1H OCEAN |
| **4861** | <1H OCEAN |