· Lab - 2.

Demonstrate the steps to build a machine learning model that predicts the median housing price using the california housing price dataset.

1. Perform describe & info steps

data = pd.read_csv ("Housing.csv")
data. describe()
data. info()

2. Plot the histograms of each feature (Indicate what clow histogram indicate on median-income & house-median-age)
data. hist (bins = 50, figsize = (20,15))
plt. show ()

3. Demonstrate the process of creating a test set (write the difference b/w random & stratified test set)

Random split = Splits data randomly without considering feature distribution. This can lead to skewed test sets if the dataset is imbalanced (e.g: if median-income has uneven representation).

Stratified split: Ensures the test set reflects the distribution of a key feature (median-income). This is crucial when a feature strongly correlates with the target (home price), reducing sampling bias

4. What does geographical feature graph indicate
w.r.t housing price & location
The graph shows California's outline, with
higher house values (red/yellow).
concentrated near the coast.
Inland areas (green/blue) have lower values.
This indicates a strong location-bound.
influence on housing-price, with proximity
to the coast correlating with higher values

5. Plot graph to show feature correlation with
housing price, which feature correlate to movement.
And lyse what graph indicates

Maximum correlation: median-income
The scatter plot of median-income vs median-
house values shows a clear positive trend -
higher incomes correspond to higher house value.

6. hint the features that could be combined to
improve correlation.

Features to combine

① total-rooms & households → room-per-household
② total-bedrooms & total-rooms → bedrooms-per-room
③ population & household → population-per-household

Observation: room-per-household improves correlation
slightly over total-rooms. The plot shows
positive trend, but with more scatter,
indicating less predictive close than median-income

7. hilt the featurs that needs to cleaned.

total bedrooms has missing values (207 missing)

median = data ["total bedrooms"]. median ()
data ["total bedrooms"]. fillna (median.inplace = true)

8. Is there any catagorical data that needs to be convented to numerical? If so explain the cuthred used to convent. E c

Yes, there is a catagorical data ocean proximity use one-hot encoding (since it nominal, not ordinal).
Convents catagory Nearbay, Ocean into binary columns

encoder = Onehot tworlk (sparse = false)
ocean encoded = encoder: fit transform (data [[ocean proximity]])

9. Discuss the importance of feature scaling

* Many ML algorithm are sensitive to feature like total room & median income
* Scaling ensure all features contribut equally improving convergence & accuracy

Methods: Standardization, Normalization

10. Design a pipeline manipulating (missing values feature scaling & encoding).

1. Numeric pipeline:
+ fills missing values (total-bedrooms) with median

+ Adds derived features (rooms-per-household, bedrooms-per-room, population-per-household).

+ Scales standardizes numeric features to zero mean & unit variance.

2. categorical pipeline:

One Hot encoder: convert ocean-proximity into binary columns.