# ECE469 - Introduction to ML
## *Midterm - Part 1*

Chase A. Lotito - *SIUC Undergraduate*

October 21, 2024

**Question 1.** *Regression in machine learning.*

**Solution.**

*(A)*



Figure 1: housing.info()

*(B)*



Figure 2: housing.describe()
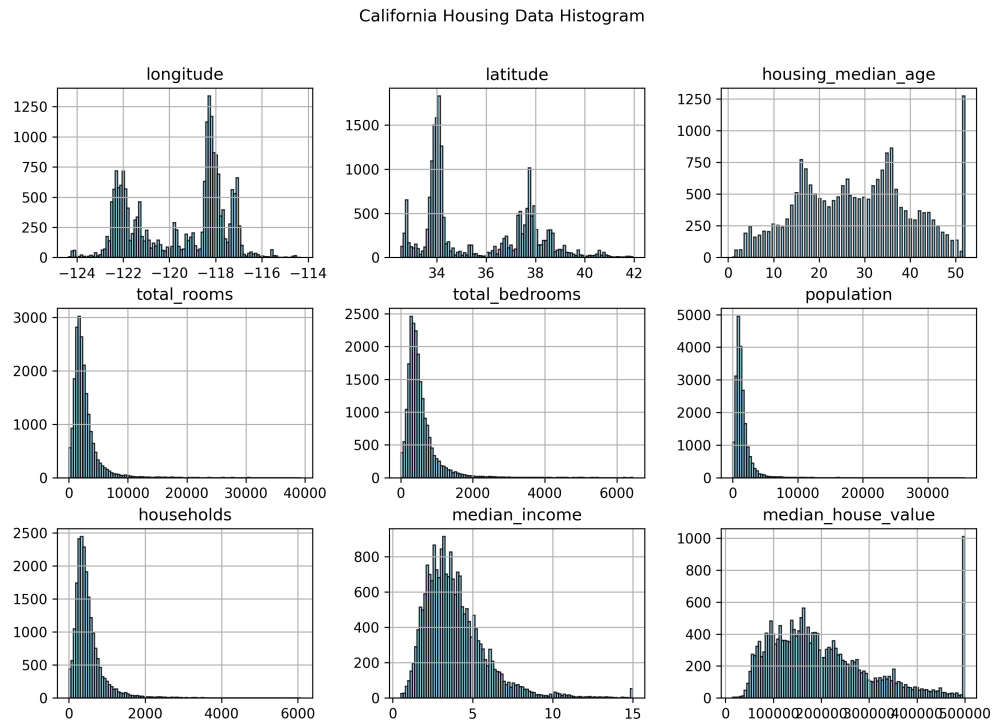
*(C)*

California Housing Data Histogram



Figure 3: Housing Dataset Histogram

*(D)*

Using *SimpleImputer* from SciKit-Learn to replace missing data with the median value, and using *StandardScaler* from SciKit-Learn to standardize the dataset, the dataset looks as follows:

```
PREPROCESSED HOUSING DATA
-------------------------
      people_per_house  bedrooms_ratio  rooms_per_house  ...  households  median_income  median_house_value
0            -0.049597       -1.029988         0.628559  ...   -0.977033       2.344766            2.129631
1            -0.092512       -0.888897         0.327041  ...    1.669961       2.332238            1.314156
2            -0.025843       -1.291686         1.155620  ...   -0.843637       1.782699            1.258693
3            -0.050329       -0.449613         0.156966  ...   -0.733781       0.932968            1.165100
4            -0.085616       -0.639087         0.344711  ...   -0.629157      -0.012881            1.172900
...                ...             ...              ...  ...         ...            ...                 ...
20635        -0.049110        0.165994        -0.155023  ...   -0.443449      -1.216128           -1.115804
20636         0.005021        0.021671         0.276881  ...   -1.008420      -0.691593           -1.124470
20637        -0.071735        0.021134        -0.090318  ...   -0.174042      -1.142593           -0.992746
20638        -0.091225        0.093467        -0.040211  ...   -0.393753      -1.054583           -1.058608
20639        -0.043682        0.113275        -0.070443  ...    0.079672      -0.780129           -1.017878

[20640 rows x 12 columns]
```

Figure 4: Preprocessed dataset

Below is the python code used for Question 1:

Listing 1: File: q1.py

```python
"""
Southern Illinois University Carbonale
Department of Electrical Engineering
--------------------------------------
```

```python
ECE469: Intro to Machine Learning
Midterm Exam: Question 1
Chase Lotito

10/14/2024

"California Housing Prices"
"""

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

from pandas.core.common import random_state
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import PolynomialFeatures
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error

# import housing data from repository
url = "https://github.com/ageron/data/raw/main/housing/housing.csv"
housing = pd.read_csv(url)        # Store housing data in DataFrame
temp = housing

# important variables
plots_url = "./plots/q1/"

# (a) Use info() method to identify the attributes of this data-set
print('HOUSING INFO\n----------')
housing.info()


# (b) Use describe() method to identify and peek at a summary of
#       the numerical attributes.
housing_desc = housing.describe()
print('\nHOUSING DESCRIPTION\n------------------')
print(housing_desc)


# (c) Use hist() method on the whole dataset and plot a histogram
#       for each numerical attribute. Notice that many histograms
#       are skewed right.
housing.hist(bins=100, color='skyblue', alpha=0.8, edgecolor='black', figsize
    =(12,8))
plt.suptitle('California Housing Data Histogram')
#plt.savefig(plots_url + 'housing_histogram.png', dpi=300)
plt.close()

# (d) Clean and normalize/standardize the data-set to make it appropriate
#       for training a regression model. Creating training and test sets.
#       Create a copy of the data with only the numerical attributes by
#       excluding the text attribute ocean proximity from the data-set.

# remove ocean_proximity from the dataset
housing = housing.drop('ocean_proximity', axis=1)

# clean data via imputer; replacing missing values with median or mean
```

```python
63   # this satisfies part (i) since total_bedrooms now is complete
64   simple_imputer = SimpleImputer(strategy='median')
65   housing_imputed = simple_imputer.fit_transform(housing)
66   housing = pd.DataFrame(housing_imputed, columns=housing.columns)
67
68   #  STICKING PART H HERE SO I CAN ADD THESE BEFORE SPLITTING DATASET
69   #  (h) Add three new attributes; (i) rooms per house = total rooms/households,
70   #      (ii) bedrooms ratio = total bedrooms/total rooms, and
71   #      (iii) people per house = population/households.
72
73   # assign housing DataFrame attributes to arrays
74   total_rooms = housing['total_rooms'].values
75   households = housing['households'].values
76   total_bedrooms = housing['total_bedrooms'].values
77   population = housing['population'].values
78
79   # calculate new attributes
80   rooms_per_house = total_rooms / households
81   bedrooms_ratio = total_bedrooms / total_rooms
82   people_per_house = population / households
83
84   # assign new attributes to housing DataFrame
85   # using .insert(0, ...) to stack each in front
86   housing.insert(0, 'rooms_per_house', rooms_per_house)
87   housing.insert(0, 'bedrooms_ratio', bedrooms_ratio)
88   housing.insert(0, 'people_per_house', people_per_house)
89
90   # standardize housing dataset
91   standard_scaler = StandardScaler()
92   housing_scaled = standard_scaler.fit_transform(housing)
93   housing = pd.DataFrame(housing_scaled, columns=housing.columns)
94   print('\nPREPROCESSED HOUSING DATA\n-------------------------')
95   print(housing)
96
97   # extract input features and output target
98   x = housing[[
99       'people_per_house',
100      'bedrooms_ratio',
101      'rooms_per_house',
102      'longitude',
103      'latitude',
104      'housing_median_age',
105      'total_rooms',
106      'total_bedrooms',
107      'population',
108      'households',
109      'median_income'
110      ]].values
111  y = housing['median_house_value'].values
112
113  # split into training and testing datasets
114  test_ratio = 0.2
115  x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=test_ratio,
         random_state=42)
116
117  # (e) Because this data-set includes geographical information (latitude
118  #      and longitude), you are asked to create a scatterplot of all the
119  #      districts to visualize the geographical data in 2D space.
120
```

4

```python
121  # get arrays containing lattitude and longitudes
122  lat = housing[['latitude']].values
123  long = housing[['longitude']].values
124
125  # plot them as a scatterplot
126  plt.figure(figsize=(10,8))
127  plt.scatter(lat, long, marker='o', s=0.75, c='#32a852', alpha=0.95, label='2D Map'
          )
128  plt.xlabel('Latitude')
129  plt.ylabel('Longitude')
130  plt.title(r'Normalized 2D Space Visualization of $housing.csv$')
131  plt.legend()
132  plt.gca().set_facecolor((0,0.1,0.8, 0.1))
133  #plt.savefig(plots_url + '2d_housing_scatter.png', dpi=300)
134  plt.close()
135
136
137  # (f) Compute the standard correlation coefficient between every pair of
138  #      attributes using the corr() method.
139  corr_matrix = housing.corr()
140  corr_to_median_house_val = corr_matrix['median_house_value'].sort_values(ascending
          =False)
141  print('\nCORRELATION of _____ TO median_house_value\n
          ------------------------------------------')
142  print(corr_to_median_house_val)
143
144
145  # (g)
146  pd.plotting.scatter_matrix(housing, alpha=0.2, figsize=(20,16))
147  plt.suptitle('Scatter Matrix for $housing.csv$')
148  plt.tight_layout()
149  # plt.savefig(plots_url + 'housing_scatter_matrix.png', dpi=80)
150  plt.close()
151
152
153  # (i)
154
155  # initalize linear model
156  linear_regressor = LinearRegression()
157
158  # train linear model
159  linear_regressor.fit(x_train, y_train)
160
161  # test linear model
162  y_train_predicted = linear_regressor.predict(x_train)
163  y_test_predicted = linear_regressor.predict(x_test)
164
165
166  # now test for different dataset sizes
167  train_sizes = np.linspace(0.1, 1.0, 10)  # Dataset sizes from 10% to 100% of the
          training set
168  train_errors = []
169  test_errors = []
170
171  for train_size in train_sizes:
172      # Split data into training and test sets
173      X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.2,
              random_state=42)
174
```

```python
175        # Use only a portion of the training set defined by `train_size`
176        X_train_subset = X_train[:int(train_size * len(X_train))]
177        y_train_subset = y_train[:int(train_size * len(y_train))]
178
179        # Train the model
180        lin_reg = LinearRegression()
181        lin_reg.fit(X_train_subset, y_train_subset)
182
183        # Make predictions
184        y_train_pred = lin_reg.predict(X_train_subset)
185        y_test_pred = lin_reg.predict(X_test)
186
187        # Calculate the errors
188        train_mse = mean_squared_error(y_train_subset, y_train_pred)
189        test_mse = mean_squared_error(y_test, y_test_pred)
190
191        # Store the errors
192        train_errors.append(train_mse)
193        test_errors.append(test_mse)
194
195  # PLOTTING RESULTS
196
197  plt.figure(figsize=(16,9))
198
199  # plot predicted against actual
200  ax1 = plt.subplot(1, 2, 1)
201  ax1.scatter(y_test, y_test_predicted, c='orange', marker='x', label='Learned Model
          ', alpha=0.6)
202  ax1.plot([min(y_test), max(y_test)], [min(y_test_predicted), max(y_test_predicted)
          ], c='blue', label='Ideal Model')
203  plt.xlabel('$\mathbf{y}_{test}$ Actual')
204  plt.ylabel('$\mathbf{y}_{test}$ Predicted')
205  plt.title('Learned Model Spread')
206  plt.legend()
207
208  ## Plot the Learned Weights Coefficients of the model
209  #weights = linear_regressor.coef_
210  #feature_names = ['people_per_house', 'bedrooms_ratio', 'rooms_per_house','
          longitude',
211  #                 'latitude', 'housing_median_age', 'total_rooms', 'total_bedrooms
          ',
212  #                 'population', 'households', 'median_income']
213  #
214  #ax2 = plt.subplot(1, 3, 2)
215  #ax2.barh(feature_names, weights)
216  #plt.xlabel("Model Weights $w_i$")
217  #plt.title("Linear Regression Feature Weights")
218
219  # Plot the learning curve
220  ax3 = plt.subplot(1, 2, 2)
221  ax3.plot(train_sizes, train_errors, label='Training Error', color='blue')
222  ax3.plot(train_sizes, test_errors, label='Test Error', color='green')
223  plt.title('Learning Curve')
224  plt.xlabel('Training Set Size (Ratio)')
225  plt.ylabel('Mean Squared Error')
226  plt.subplots_adjust(wspace=0.2)
227
228  plt.suptitle('Linear Model Results for $housing.csv$')
229  plt.legend()
```

```
230  plt.savefig(plots_url + "linear_model_results.png", dpi=300)
231  plt.close()
```

**Question 2.** *Classification in machine learning.*
**Solution.**
Below is the python code used for Question 2:

Listing 2: File: q2.py

```
 1  """
 2  Southern Illinois University Carbonale
 3  Department of Electrical Engineering
 4  -------------------------------------
 5  ECE469: Intro to Machine Learning
 6  Midterm Exam: Question 2
 7  Chase Lotito
 8
 9  10/15/2024
10
11  "MNIST"
12  """
13
14  import pandas as pd
15  import numpy as np
16  import matplotlib.pyplot as plt
17
18  from sklearn.datasets import fetch_openml
19  from sklearn.model_selection import train_test_split
20  from sklearn.preprocessing import StandardScaler
21  from sklearn.linear_model import LogisticRegression
22  from sklearn.neighbors import KNeighborsClassifier
23  from sklearn.metrics import root_mean_squared_error, accuracy_score, f1_score
24  from scipy.ndimage import shift
25
26  from sklearn.decomposition import PCA
27
28  # fetch the MNIST dataset
29  mnist = fetch_openml('mnist_784', as_frame=False)
30
31  # extract input features and output target
32  X = mnist['data']
33  y = mnist['target']
34
35  # the target values in y are strings, so we must
36  # first convert them to integers
37  y = y.astype(int)
38
39
40
41  # (b) Write a function that can shift an MNIST image
42  #     in any direction. Do this in all directions for
43  #     the training set, and append them to it.
44  def quad_direction_enricher(X: np.array, y: np.array, size: int, px: int):
45      """
46      To enrich an image dataset with 4 sets of
47      the original set shifted in all directions px.
48      (up, down, left, right)
49      """
50      temp = []
51      y_temp = []
```

7

```python
52        for k in range(0, 4, 1):
53            for i in range(0, size, 1):
54                img = X[i].reshape(28,28)
55
56                if (k == 0):
57                    img = shift(img, (px, 0))   # shift up
58                elif (k == 1):
59                    img = shift(img, (-px, 0))  # shift down
60                elif (k == 2):
61                    img = shift(img, (0, px))   # shift right
62                elif (k == 3):
63                    img = shift(img, (0, -px))  # shift left
64                else:
65                    print('ERROR: image shift bounds error')
66
67                temp.append(img)
68                y_temp.append(y[i])
69
70        #enriched_X = np.array([img.flatten() for img in temp])
71        enriched_X = np.array([img.flatten() for img in temp])
72        enriched_Y = np.array(y_temp)
73        return enriched_X, enriched_Y
74
75 # enrich X
76 X, y = quad_direction_enricher(X, y, len(X), px=1)
77
78 print(f'len(X) = {len(X)}, len(y) = {len(y)}')
79
80 subset_size = 30000
81 # plot the examples of elements in the enriched dataset
82 #ax1 = plt.subplot(2, 4, 1)
83 #ax1.imshow(X[1].reshape(28,28))
84 #plt.title('Original NMIST')
85 #ax2 = plt.subplot(2, 4, 2)
86 #ax2.imshow(X[2].reshape(28,28))
87 #ax3 = plt.subplot(2, 4, 3)
88 #ax3.imshow(X[3].reshape(28,28))
89 #ax4 = plt.subplot(2, 4, 4)
90 #ax4.imshow(X[4].reshape(28,28))
91 #ax5 = plt.subplot(2, 4, 5)
92 #ax5.imshow(X[subset_size - 1].reshape(28,28))
93 #plt.title('Shifted NMIST', loc='center')
94 #ax6 = plt.subplot(2, 4, 6)
95 #ax6.imshow(X[2*subset_size - 2].reshape(28,28))
96 #ax7 = plt.subplot(2, 4, 7)
97 #ax7.imshow(X[3*subset_size - 3].reshape(28,28))
98 #ax8 = plt.subplot(2, 4, 8)
99 #ax8.imshow(X[4*subset_size - 4].reshape(28,28))
100 #plt.suptitle('Example Elements of Enriched MNIST')
101 #plt.show()
102 #plt.close()
103
104 # split dataset into training and testing sets
105 test_ratio = 0.2
106 X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=test_ratio,
        random_state=42)
107
108
109 # scale down size of dataset
```

```python
110  test_size = int(subset_size*test_ratio)
111  train_size = int(subset_size*(1-test_ratio))
112  X_train = X_train[:train_size]
113  X_test = X_test[:test_size]
114  y_train = y_train[:train_size]
115  y_test = y_test[:test_size]
116
117  # next we want to standardize our data, but not
118  # targets, as to preserve y\in(0,9)
119  standard = StandardScaler()
120  X_train = np.array(X_train).reshape(len(X_train), -1)
121  X_test = np.array(X_test).reshape(len(X_test), -1)
122  X_train = standard.fit_transform(X_train)
123  X_test = standard.transform(X_test)
124
125  # train logisitic regression classifier
126
127  pca = PCA(n_components=0.95)
128  X_train_pca = pca.fit_transform(X_train)
129  X_test_pca = pca.transform(X_test)
130  X_train = X_train_pca
131  X_test = X_test_pca
132
133  tolerance = 1e-3
134  classify = LogisticRegression(solver='lbfgs', penalty='l2', C=2, max_iter=1000,
         random_state=0)
135  classify.fit(X_train, y_train)
136
137  # predict using model
138  y_test_pred = classify.predict(X_test)
139  y_train_pred = classify.predict(X_train)
140
141  # calc MSE
142  test_rmse = root_mean_squared_error(y_test, y_test_pred)
143  train_rmse = root_mean_squared_error(y_train, y_train_pred)
144
145  # determine model accuracy
146  test_accuracy = accuracy_score(y_test, y_test_pred) * 100
147  train_accuracy = accuracy_score(y_train, y_train_pred) * 100
148
149  # determine model f1 score
150  test_f1 = f1_score(y_test, y_test_pred, average='macro')
151  train_f1 = f1_score(y_train, y_train_pred, average='macro')
152
153  print('MNIST CLASSIFICATION REPORT')
154  print('##########################')
155  print(f'Dataset Size : {subset_size}')
156  print(f'Enriched Size : {len(X)}')
157  print(f'Probabilistic Model Tolerance     : {tolerance*100}%')
158  print(f'Probabilistic Test      RMSE      : {test_rmse:.2f}')
159  print(f'Probabilistic Test  Accuracy      : {test_accuracy:.2f}%')
160  print(f'Probabilistic Test  F1 Score      : {test_f1:.2f}')
161  print(f'Probabilistic Train     RMSE      : {train_rmse:.2f}')
162  print(f'Probabilistic Train Accuracy      : {train_accuracy:.2f}%')
163  print(f'Probabilistic Train F1 Score      : {train_f1:.2f}')
164
165  # (c) KNN-based algorithms belong to the class of non-probabilistic classifiers.
166  #     You are asked to design a KNN-based classifier to classify handwritten
167  #     digits (0-9) in MNIST data-set.
```

```python
# initialize the knn classifier
knn = KNeighborsClassifier(n_neighbors=5, metric='minkowski')

# train knn classifier
knn.fit(X_train, y_train)

# predict using knn classifier
y_test_pred_knn = knn.predict(X_test)
y_train_pred_knn = knn.predict(X_train)

# calc MSE
test_rmse_knn = root_mean_squared_error(y_test, y_test_pred_knn)
train_rmse_knn = root_mean_squared_error(y_train, y_train_pred_knn)

# determine model accuracy
test_accuracy_knn = accuracy_score(y_test, y_test_pred_knn) * 100
train_accuracy_knn = accuracy_score(y_train, y_train_pred_knn) * 100

# determine model f1 score
test_f1_knn = f1_score(y_test, y_test_pred_knn, average='macro')
train_f1_knn = f1_score(y_train, y_train_pred_knn, average='macro')

# add KNN model results to classification report
print(f'Non-probabilistic Test     RMSE  : {test_rmse_knn:.2f}')
print(f'Non-probabilistic Test  Accuracy : {test_accuracy_knn:.2f}%')
print(f'Non-probabilistic Test  F1 Score : {test_f1_knn:.2f}')
print(f'Non-probabilistic Train     RMSE  : {train_rmse_knn:.2f}')
print(f'Non-probabilistic Train Accuracy : {train_accuracy_knn:.2f}%')
print(f'Non-probabilistic Train F1 Score : {train_f1_knn:.2f}')
```