**Question 1:** Select ALL correct choices.

[**1.1**] Ans: **B**

[**1.2**] Ans: **A**, **B**

[**1.3**] Ans: **A**, **C**

[**1.4**] Ans: **A**, **B**, **D**

[**1.5**] Ans: **A**, **B**, **C**

**Question 2:** A linear ML model can be written as:

$$f(\mathbf{x}, \mathbf{w}) = \sum_{i=0}^{n} w_i x_i = \mathbf{w}^T \mathbf{x}$$

The loss function can be written as:

$$J(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^{n} \left[ f(\mathbf{x}^{(i)}, \mathbf{w}) - y^{(i)} \right]^2$$

[**2.1**] Show analytically that the optimal weight vector that minimizes the cost function $J(\mathbf{w})$ is:

$$\mathbf{w}^* = \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

**Solution.**

Since our model is linear, we can write the cost function as:

$$J(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^{n} \left[ \mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)} \right]^2 \tag{1}$$

The product $\mathbf{w}^T \mathbf{x}^{(i)} = \mathbf{X}\mathbf{w}, \forall i \in \{1, 2, \cdots, n\}$ since the LHS implies the matrix mulplitcation of the RHS, as $\mathbf{X}$ is the matrix of all input entries $\mathbf{x}^{(i)}$. So, Eq. 1 can be written as:

$$J(\mathbf{w}) = \frac{1}{m} \left[ \mathbf{X}\mathbf{w} - \mathbf{y} \right]^2 \tag{2}$$

The inside of the brackets is just a vector, and the square of a vector is the norm of a vector, so we can reduce Eq. 2:

$$J(\mathbf{w}) = \frac{1}{m} \| \mathbf{X}\mathbf{w} - \mathbf{y} \| \tag{3}$$

$$= \frac{1}{m} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \tag{4}$$

$$= \frac{1}{m} (\mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{y}^T \mathbf{y}) \tag{5}$$

Now to optimize the cost with respect to weights, we can take the gradient of $J(\mathbf{w})$ w.r.t. the weights $\mathbf{w}^{(i)}$.

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \frac{1}{m} \nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y}) \tag{6}$$

$$= \frac{1}{m} \nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w}) \tag{7}$$

$$= \frac{1}{m} [\nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}) - \nabla_{\mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{y}) - \nabla_{\mathbf{w}} (\mathbf{y}^T \mathbf{X} \mathbf{w}))] \tag{8}$$

The gradients $\nabla_{\mathbf{w}}$ are simply derivatives of each matrix function w.r.t. $\mathbf{w}$, which can be computed using equations (69) and (81) from *The Matrix Cookbook* [2].

$$= \frac{1}{m} \left[ \frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}) - \frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^T \mathbf{X}^T \mathbf{y}) - \frac{\partial}{\partial \mathbf{w}} (\mathbf{y}^T \mathbf{X} \mathbf{w})) \right] \tag{9}$$

$$= \frac{1}{m} ((\mathbf{X}^T \mathbf{X} + (\mathbf{X}^T \mathbf{X})^T) \mathbf{w} - \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}) \tag{10}$$

$$= \frac{1}{m} ((\mathbf{X}^T \mathbf{X} + \mathbf{X}^T (\mathbf{X}^T)^T) \mathbf{w} - 2 \mathbf{X}^T \mathbf{y}) \tag{11}$$

$$= \frac{1}{m} ((\mathbf{X}^T \mathbf{X} + \mathbf{X}^T \mathbf{X}) \mathbf{w} - 2 \mathbf{X}^T \mathbf{y}) \tag{12}$$

$$= \frac{1}{m} (2 \mathbf{X}^T \mathbf{X} \mathbf{w} - 2 \mathbf{X}^T \mathbf{y}) \tag{13}$$

$$= \frac{2}{m} (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}) \tag{14}$$

We find the minimum when $\nabla_{\mathbf{w}} J(\mathbf{w}) = 0$,

$$\frac{2}{m} (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}) = 0 \tag{15}$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} = 0 \tag{16}$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y} \tag{17}$$

$$(\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X} \mathbf{w}) = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y}) \tag{18}$$

$$\mathbf{I} \mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \tag{19}$$

$$\boxed{\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{w}^*} \tag{20}$$

Without advanced methods, multiplying two $(n \times n)$ matrices is of $O(n^3)$ complexity [1]. For very large data sets, computing $\mathbf{w}$ would be enormously computationally expensive.

[**2.2**] Develop a pseudo-codes for implementing the batch gradient descent, stochastic gradient descent, and mini-batch gradient descent algorithms to train the above linear model.

**Solution.**

Batch Gradient Descent

```
1  a = learning rate
2  N = size of dataset
3
4  # Psuedocode for Batch Gradient Descent
5  REPEAT UNTIL CONVERGENCE
6  {
7  for all j:
8      temp_j = w_j - (a / N) * sum( i=1..N, gradient_w(J(
           x_i, w_j) )
9
10 w_j = temp_j
11 }
```

Stochastic Gradient Descent

```
1  a = learning rate
2
3  # Psuedocode for Stochastic Gradient Descent
4  REPEAT UNTIL CONVERGENCE
5  {
6  for all i:
7      w_j = w_j - a * gradient_w(J(x_i, w_j))
8  }
```

Mini-Batch Gradient Descent

```
1  a = learning rate
2  b = batch size
3  m = rows in training set
4
5  # Psuedocode for Mini-Batch Gradient Descent
6  REPEAT UNTIL CONVERGENCE
7  {
8  # i indexes will be multiples of batch size
9  for i = (1*b + 1, 2*b + 1, ... , (m-1)*b + 1):
10     w_j = w_j - (a / b) * sum(k=i..i+(b-1), gradient_w(J(
           x_k, w_j)))
11 }
```

[**2.3**] Discuss the performance versus computational complexity of each of the above algorithms.

**Solution.**

Batch-Gradient Descent (BDG) is the most compuationally complex since the algorithm has to sum over the entire dataset *every single update* of $w_j$. Stochastic Gradient Descent (SDG) is much faster than BDG since each update of $w_j$ only looks over a single data row in the $X$ matrix, instead of all of them. Mini-Batch Gradient Descent (MBDG) is inbetween BDG and SDG in terms of computational

complexity, since MBDG updates $w_j$ after summing the gradients of a specified batch size $b \in (1, \cdots, N)$.

A ranking of each algorithm from least compuationally complex to least is as follows:

(1) Stochastic

(2) Mini-Batch

(3) Batch

**Question 3:** Housing data preprocessing.

**Solution.**

```python
1  # Chase Lotito - SIUC Fall 2024 - ECE469: Intro to
       Machine Learning
2  # HW1 - Question 3
3
4  import numpy as np
5  import pandas as pd
6  from sklearn.preprocessing import OrdinalEncoder    # For
       encoding categorical features
7  from sklearn.impute import SimpleImputer            # For
       adding missing values
8  from sklearn.preprocessing import StandardScaler    # For
       standardizing data
9
10 # (A) Get housing data
11 RAW_DATA = 'https://github.com/ageron/data/raw/main/
       housing/housing.csv'
12 housing = pd.read_csv(RAW_DATA)
13
14 # (B) Choose input features and output features (saved
       into numpy.ndarray type)
15 X = housing[
16         ['longitude',
17          'latitude',
18          'housing_median_age',
19          'total_rooms',
20          'total_bedrooms',
21          'population',
22          'households',
23          'median_income',
24          'ocean_proximity']
25     ].values
26 Y = housing[['median_house_value']].values
27
28 # (C) Ocean Proximity is a categorical feature. Drop it
       or transform into numerical values (encode).
29
```

```
30  # Isolate the ocean_proximity data in input data X
31  ocean_proximity = X[:,8].reshape(-1,1)    # reshape(-1,1)
        to make 2D array for Ordinal
32  # Initalize the ordinal encoder
33  ordinal_encoder = OrdinalEncoder()
34  # Encode the ocean_proximity strings into numerical data
35  encoded_ocean = ordinal_encoder.fit_transform(
        ocean_proximity)
36  # Put the encoded version of ocean_proximity into input
        data X
37  X[:,8] = encoded_ocean.flatten()          # flatten to add
         1D version of array back into X
38
39  # (D) Clean the dasta by either dropping or replacing
        missing values
40
41  # Initialized SimpleImputer, will use the median to add
        missing entries
42  simple_imputer = SimpleImputer(strategy='median')
43
44  # Change X np ndarray into a Pandas Dataframe to use
        SimpleImputer
45  dX = pd.DataFrame(X)
46  dY = pd.DataFrame(Y)
47
48  # Perform SimpleImputer transformation, for both inputs X
         and outputs Y
49  imputed_data = simple_imputer.fit_transform(dX)
50  X = imputed_data
51  imputed_data = simple_imputer.fit_transform(dY)
52  Y = imputed_data
53
54  # (E) Carry out feature scaling either via normalization
        or standardization.
55  std_scaler = StandardScaler()
56  scaled_data = std_scaler.fit_transform(X)
57  X = scaled_data
58  scaled_data = std_scaler.fit_transform(Y)
59  Y = scaled_data
60
61  # (F) Create a training dataset and testing dataset
62  def shuffle_and_split_data(data, test_ratio):
63      shuffled_indices = np.random.permutation(len(data))
64      test_set_size = int(len(data) * test_ratio)
65      test_indices = shuffled_indices[:test_set_size]
66      train_indices = shuffled_indices[test_set_size:]
67      return data.iloc[train_indices], data.iloc[
            test_indices]
```

```
68
69  # first, recombine X and Y into augmented matrix
70  housing_data = np.hstack((X,Y))
71
72  # split into testing and training set (both outputted as
        pd.DataFrames)
73  housing_training, housing_testing =
        shuffle_and_split_data(pd.DataFrame(housing_data),
        0.2)
74  print('TRAINING:')
75  print(housing_training)
76  print('TESTING:')
77  print(housing_testing)
```

# References

[1] Andy He and Evan Williams. Computational complexity of matrix multiplication. `https://www.cs.cornell.edu/courses/cs6810/2023fa/Matrix.pdf`, Fall 2023. Accessed: 2024-09-10.

[2] Kaare Brandt Petersen and Michael Syskind Pederson. The matrix cookbook. Distributed by University of Waterloo, November 2012.