# MSAI349: Machine Learning
# Homework 2

Sayantani Bhattacharya, Harrison Bounds, Andrew Kwolek, Sharwin Patil

November 3, 2024

# Contents

# 1 Distance Metrics, K-Nearest Neighbors, K-Means-Clustering

## 1.1 Distance Metrics

We implemented 4 distance metrics that are capable of comparing any 2 $n$ dimensional vectors to each other:

1. Euclidean Distance

2. Cosine Similarity

3. Hamming Distance

4. Pearson Distance

Each of these metrics were implemented without the use of python libraries and were verified against the output of those metrics offered by libraries such as `sklearn`, `numpy`, and `scipy`. These tests were written using `pytest` and can be run with the command: `pytest -s` *(-s flag will show stdout for debugging)*.

### 1.1.1 Euclidean Distance

Euclidean distance scales similarly in low vs high dimensions and is a good option for low dimension comparisons. Additionally, euclidean distance maps "physical" proximity which may be useful for data that has geometric features such as image data.

$$d(x,y)_{\text{euclidean}} = \sqrt{\sum (x_i - y_i)^2} \tag{1}$$

```
euclidean_distance = np.sqrt(np.sum(np.subtract(a, b)**2))
```

### 1.1.2 Cosine Similarity

Cosine Similarity will measure how similar two vectors are by finding the "angle" between them. Despite this concept being impossible to visualize in higher dimensions, cosine similarity excels in high dimensional data, making it a good metric for image data, natural language processing, and other high-dimensional data. Note that the value is bounded between $[-1, 1]$ so the magnitude of the data isn't encapsulated in the return value.

$$d(x,y)_{\text{cosim}} = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2} \sqrt{\sum y_i^2}} \tag{2}$$

```
cosim = np.dot(a, b) / np.sqrt(np.sum(a**2)) * np.sqrt(np.sum(b**2))
```

### 1.1.3 Hamming Distance

Hamming distance counts the number of differences between the vectors (Both vectors must be the same dimension and size). Hamming distance is useful for categorical data, error detection, and comparing strings (edit distance for spell-checking).

$$d(x,y)_{\text{hamming}} = \sum \delta(x_i, y_i), \; \delta = 1 \text{ if } x_i \neq y_i \text{ else } \delta = 0 \tag{3}$$

```
hamming_distance = np.array([ai != bi for ai, bi in zip(a, b)]).sum()
```

### 1.1.4 Pearson Distance

Pearson distance measures the linear relationship between the vectors.

$$d(x,y)_{\text{pearson}} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} \tag{4}$$

```
pearson_distance = np.sum((a - a_bar)*(b - b_bar)) /  np.sqrt(np.sum((a - a_bar)**2)) *
↪  np.sqrt(np.sum((b - b_bar)**2))
```

## 1.2 K-Nearest Neighbors Classifier

When implementing K-Nearest Neighbors, our only hyper-parameter was $k$. When attempting to use less observations during the distance comparison, there was a drop in accuracy compared to using all the of the given observations. Dimensionality reduction improved the algorithm's efficiency but had no impact (positive or negative) on performance so dimensionality reduction was always done before using KNN.

We got satisfactory results with $k = 5$ using euclidean distance as our distance metric. Dimensionality reduction didn't impact our results at all, however made execution much faster ($\approx 15\%$ faster). The results observed for the mnist testing and validation datasets upon KNN classification:

```
Testing Data Metrics (euclidean):
    Accuracy: 0.9800000000000001
    Precision: 0.9025205627705628
    Recall: 0.8905366767735188
    F1_score: 0.8900603439819568
Validation Data Metrics (euclidean):
    Accuracy: 0.8150000000000001
    Precision: 0.08002930605013658
    Recall: 0.07614514712340799
    F1 Score: 0.07582969280837747
```
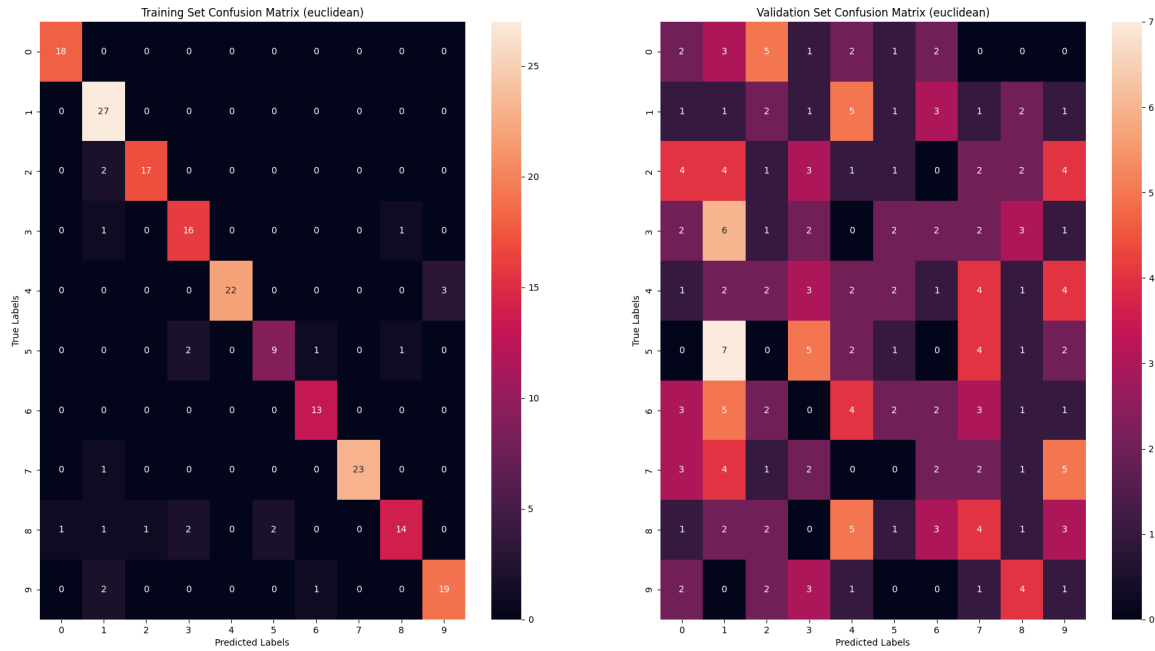
Figure 1: Confusion Matrices for KNN Classifier using the Euclidean distance metric
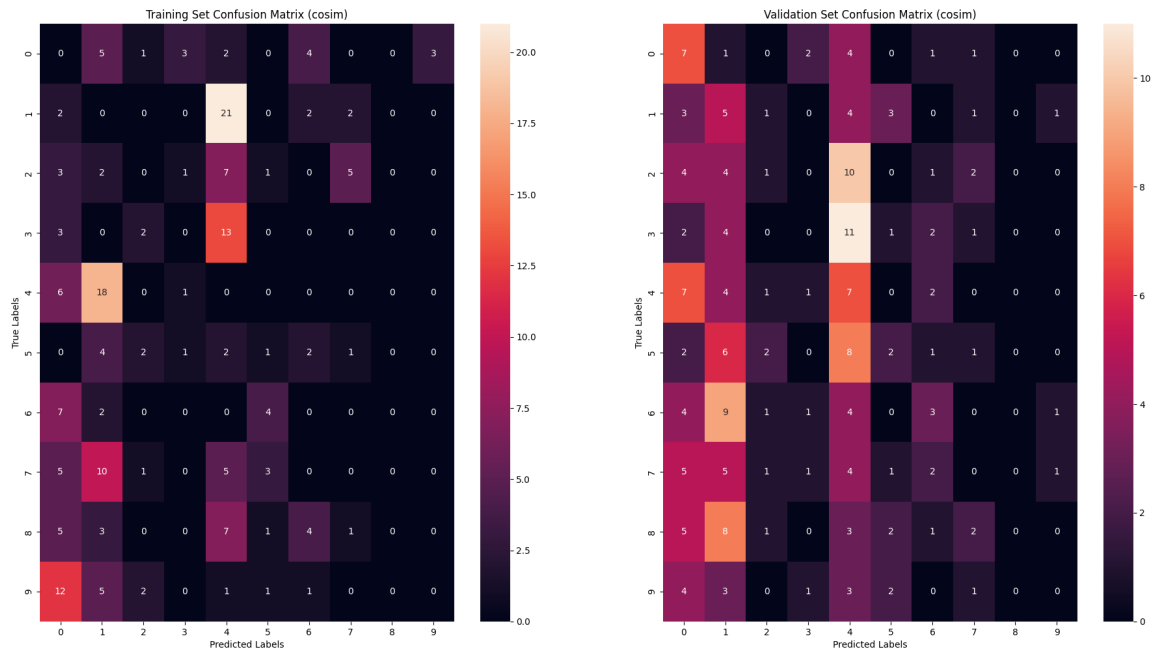


Figure 2: Confusion Matrices for KNN Classifier using the Cosine Similarity distance metric

## 1.3   K-Means Classifier

In order to determine how well our clusters aligned with the labels in `mnist_test.csv`, we created a metric that calculates the total accuracy of the clusters based on how many points are correctly labeled. We do this by looping through the K clusters and assigning each cluster a label based on the mode of the true labels in the cluster. Then we calculate how many points are correctly classified based on the assigned cluster label. Ideally, each cluster would only contain points of the same label when run on the testing set. We tried this for different K values and found different accuracies.

```
Accuracy(K=10, Euclidean): 0.625
Accuracy(K=36, Euclidean): 0.775
```

When K was increased, the overall accuracy was improved, however the computational cost also increased with higher K values resulting in trade-offs between speed and accuracy. Better performance at higher K values could be attributed to the points in mnist forming clusters that aren't fully representative of the 10 labels, but rather there are multiple clusters embedded in each label, resulting in better accuracy with a higher K. Euclidean distance and cosine similarity performed relatively the same, so we chose to use Euclidean. Some unique design choices we made include using variance thresholding for dimensionality reduction and a random sampling technique for initializing centroids which calculates probabilities based on distances from the currently initialized centroids. This prevents centroids from being initialized very close to one another which breaks the algorithm.

# 2   Movie Reccomender

## 2.1   Collaborative Filter

The collaborative filter is made up of multiple steps to ensure a good recommendation is given to the target user. First, a user item matrix is built where the rows are user ids and the columns are movie ids. Each value in this matrix is the user's rating for that particular movie. This is made with the training set for each user, combined into one single matrix for all training data. Second, we take each row as a vector, and calculate the cosine similarity with a function that is used in part 1. Doing this gets us a value that tells us how similar each user is to one another. We can get the K top most similar users and look at their movies that the target user hasn't seen. To do this, we look through these top K users and assign each rating a weight based on the actual rating, and the similarity with the target user. This ensures that the most similar users have higher weighted ratings, which will give more accurate recommendations. Finally, this weighted ratings list returns the top M movies that the target user has not rated, but will enjoy based on its most similar K users.

K and M are the hyper parameters for this collaborative filtering, where K is the number of similar users closest to the target user, and M is the number of movies to recommend to the target user. Tuning K to a higher number would allow us to choose from a

wider variety of movies, increasing the performance of our model. M is the number of movies that is recommended to the target user. Increasing M can have a negative affect on the model because each new recommendation must pass the evaluation standards. Here we are calculating the precision, recall, and f1 score. With less variety, increasing M might make your movies fall below a certain rating threshold set (another hyper-parameter), decreasing the performance of the model. Yet, if there are many good movies to recommend, this is a good way to decrease your false negatives and achieve a higher recall and f1 score.

Overall, testing the model turned out to be tricky. With a truth label between 1-5, there was no clear-cut way to accumulate true positives, false positives, and true negatives (each of the values needed to evaluate the model). So, a rating threshold was set and if the recommended movie exceeded this rating threshold, the value was counted as a true positive. Another roadblock was the differences in the training, testing, and validation sets. If a movie recommended in the training set did not appear in the validation or testing sets, it would impact the model negatively, influencing the evaluation metrics even if the recommended movies were accurate.

```
Testing Data Metrics (cosine similarity):
    Precision: 0.09090909090909091,
    Recall: 0.001937984496124031,
    F1 Score: 0.0018975332068311196
Validation Data Metrics (cosine similarity):
    Precision: 0.5,
    Recall: 0.008988764044943821,
    F1 Score: 0.008830022075055188
```

Our metrics perform poorly because the movies recommended in the training set are not found in the test set, therefore there is no way to represent that value. This negatively impacts both the training and validation metrics.

## 2.2   Collaborative Filter with demographic information

The approach of including demographic information enhances the recommended system by making it more personalized. The metric used is the normalized version of the similarity, to ensure the metric is comparable with respect to the rating metric and the range of age variation in the dataset.

We added weight values to each metric (rating value and age), with the weights being the hyper parameters for the filter. Tuning them would determine the dominance of each metric over deciding the final movie to recommend. To check, if we set the age weight to be 0, we get the same result as above. Upon tuning we got a stabilized result at 50-50 weight to both the metrics, and improved the performance by a good margin.

```
Testing Data Metrics (cosine similarity):
    Precision: 0.3125,
    Recall: 0.007246376811594203,
```

```
    F1 Score: 0.00708215297450425
Validation Data Metrics (cosine similarity):
    Precision: 0.36363636363636365
    Recall: 0.007029876977152899,
    F1 Score: 0.00689655172413793
```

# 3    Code Appendix

## 3.1    Distance Metrics

```python
import numpy as np


# returns Euclidean distance between vectors and b
def euclidean(a, b):
    """ Returns the euclidean distance between vectors a and b

    Args:
        a (np.ndarray): A vector of any dimension
        b (np.ndarray): A vector of any dimension

    Returns:
        float: The euclidean distance between the two vectors

    Raises:
        ValueError: If the given vectors are different dimensions
    """
    return np.sqrt(np.sum(np.subtract(a, b)**2))


# returns Cosine Similarity between vectors a and b
def cosim(a, b):
    """ Returns the cosine similarity between vectors a and b

    Args:
        a (np.ndarray): A vector of any dimension
        b (np.ndarray): A vector of any dimension

    Returns:
        float: The cosine similarity between the two vectors

    Raises:
        ValueError: If the given vectors are different dimensions
    """

    if not in_same_dimension(a, b):
        print(
            f"Given vectors have different shapes: " +
            f"{np.shape(a)} != {np.shape(b)}"
        )
        raise ValueError(
            "Cosine Similarity requires 2 identically-shaped vectors"
        )
    numerator = np.dot(a, b)
    denominator = np.sqrt(np.sum(a**2)) * np.sqrt(np.sum(b**2))

    # If denominator is 0, cosim is undefined not 0
    if denominator == 0:
        return np.nan

    return numerator / denominator


# returns Pearson Correlation between vectors a and b
def pearson(a: np.ndarray, b: np.ndarray):
    """ Returns the pearson correlation between vectors a and b
```

```python
58          Args:
59              a (np.ndarray): A vector of any dimension
60              b (np.ndarray): A vector of any dimension
61
62          Returns:
63              float: The pearson correlation between the two vectors
64
65          Raises:
66              ValueError: If the given vectors are different dimensions
67          """
68          a_bar = np.sum(a)/len(a)
69          b_bar = np.sum(b)/len(b)
70
71          numerator = np.sum((a - a_bar)*(b - b_bar))
72          denominator = np.sqrt(np.sum((a - a_bar)**2)) * \
73              np.sqrt(np.sum((b - b_bar)**2))
74
75          if numerator == 0 or denominator == 0:
76              return 0
77
78          return numerator / denominator
79
80
81      def hamming(a: np.ndarray, b: np.ndarray) -> int:
82          """ Returns the Hamming distance between vectors a and b
83
84          Args:
85              a (np.ndarray): A vector of any dimension
86              b (np.ndarray): A vector of any dimension
87
88          Returns:
89              int: The Hamming distance between the two vectors
90
91          Raises:
92              ValueError: If the given vectors are different dimensions
93          """
94          # Ensure that the two vectors occupy the same dimension
95          if not in_same_dimension(a, b):
96              print(
97                  f"Given vectors have different shapes: " +
98                  f"{np.shape(a)} != {np.shape(b)}"
99              )
100             raise ValueError("Hamming requires 2 identically-shaped vectors")
101         # Create a vector
102         comparison_vector = np.array([ai != bi for ai, bi in zip(a, b)])
103         return comparison_vector.sum()
104
105
106     def in_same_dimension(a: np.ndarray, b: np.ndarray) -> bool:
107         """ Determines if the two given vectors are in the
108             same dimension or not
109
110         Args:
111             a (np.ndarray): The first vector of any dimension
112             b (np.ndarray): The second vector of any dimension
113
114         Returns:
115             bool: Whether the 2 vectors are the same dimension or not
116         """
117         return np.shape(a) == np.shape(b)
```

## 3.2   K-Nearest-Neighbors Implementation

```python
from distance_metrics import euclidean, cosim
from starter import read_data, reduce_data, reduce_query
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

# returns a list of labels for the query dataset based upon observations in the train
 ↪  dataset.
# labels should be ignored in the training set
# metric is a string specifying either "euclidean" or "cosim".
# All hyper-parameters should be hard-coded in the algorithm.


def knn(train: list, query: list, metric: str, k: int = 5, debug: bool = False) -> list:
    """
    Returns a list of labels for the query dataset based upon observations in the train
     ↪  dataset.

    Note that the length of the labels returned is the same as the length of the query
     ↪  dataset
    since each query is assigned a label.

    Args:
        train (list): The training dataset or examples that KNN will utilize to
         ↪  calculate distance and assign labels
        query (list): The dataset of queries that KNN will assign labels to
         ↪  [query_label, [list(pixels)]]
        metric (str): The distance metric to use for KNN. Either 'euclidean' or 'cosim'
        k (int, optional): The number of neighbors to consider. Defaults to 5.
        debug (bool, optional): Whether to print debug information. Defaults to False.

    Raises:
        ValueError: If the distance metric is not 'euclidean' or 'cosim'
        ValueError: If the query data is not the same size as the data in the training
         ↪  set.

    Returns:
        list: The labels assigned to each query in the query dataset
    """
    f_d = None
    if metric == 'euclidean':
        f_d = euclidean
    elif metric == 'cosim':
        f_d = cosim
    else:
        raise ValueError('Invalid distance metric given')
    if debug:
        print(
            f'K-Nearest Neighbors using {metric} distance metric and k={k}, ' +
            f'{len(train)} training examples and {len(query)} queries:'
        )
    labels = []
    for q_label, q in query:
        # Sort neighbors using distance metric and take the first k entries
        nearest_neighbors = sorted(
            [t for t in train], key=lambda x: f_d(x[1], q)
        )[:k]
        labels_for_neighbor = [int(t[0]) for t in nearest_neighbors]
        # Find the most common label among the k closest neighbors
```

```python
54                 most_common_label = np.argmax(np.bincount(labels_for_neighbor))
55                 if debug:
56                     print(
57                         f'Query {q}\n' +
58                         f'Nearest neighbors: {nearest_neighbors}\n' +
59                         f'Labels for neighbors: {labels_for_neighbor}\n' +
60                         f'Most common label: {most_common_label}\n' +
61                         f'Expected label: {q_label}'
62                     )
63                 labels.append(most_common_label)
64         return labels
65
66
67     def evaluate_knn_accuracy(labels: list, query: list) -> tuple:
68         """
69         Calculates and prints metrics, i.e. Accuracy, Precision, Recall and F1 Score for a
        ↪   trained KNN model.
70
71         Args:
72             labels (list): The labels assigned to each query in the query dataset by the KNN
            ↪   model, whose accuracy is being measured.
73             query (list): The dataset of queries that KNN will assign labels. [query_label,
            ↪   [list(pixels)]]
74
75         Returns:
76             tuple: A tuple containing the accuracy, precision, recall, and F1 score of the
            ↪   KNN model
77
78         """
79         # For 100% accuracy, diagonal elements of confusion matrix need to be non-zero and
        ↪   rest all needs to be 0.
80         expected_result = [row[0] for row in query]
81         confusion_matrix = generate_confusion_matrix(labels, expected_result)
82         num_labels = confusion_matrix.shape[0]
83         metrics = []   # (accuracy, precision, recall, f1_score)
84         for i in range(num_labels):
85             # True Positives: Diagonal elements (Correctly classified)
86             tp = confusion_matrix[i][i]
87             # False Negatives: Everything in this row except TP since it is not classified
            ↪   as i
88             fn = np.sum(confusion_matrix[i, :]) - tp
89             # False Positives: Everything in this column except TP since it is not
            ↪   classified as i
90             fp = np.sum(confusion_matrix[:, i]) - tp
91             # True Negatives: Everything except TP, FN, FP
92             tn = np.sum(confusion_matrix) - (tp + fn + fp)
93             accuracy = (tp + tn) / np.sum(confusion_matrix)
94             precision = tp / (tp + fp) if tp + fp > 0 else 0
95             recall = tp / (tp + fn) if tp + fn > 0 else 0
96             f1_score = (2 * precision * recall) / \
97                 (precision + recall) if precision + recall > 0 else 0
98             metrics.append(
99                 [accuracy, precision, recall, f1_score]
100            )
101        accuracy, precision, recall, f1_score = np.mean(metrics, axis=0)
102        print(f"Accuracy: {accuracy}")
103        print(f"Precision: {precision}")
104        print(f"Recall: {recall}")
105        print(f"F1_score: {f1_score}")
106        return (accuracy, precision, recall, f1_score)
107
```

```python
108
109  def generate_confision_matrix(labels: list, expected_result: list):
110      """
111      Returns the confusion matrix with the input of label and expected result.
112
113      Args:
114          labels (list): The labels assigned to each query in the query dataset by the KNN
                 ↪   model, whose accuracy is being measured.
115          expected_result (list): The correct label values of the query dataset.
116
117      Returns: Confusion matrix (a 2D array): is a square (n*n) matrix, with n = number of
             ↪   label options, as the union of knn and actual label of the querry set.
118              In this case, if the input data is sufficiently large: CM -> 10*10
                     ↪
119      """
120      n = len(set(expected_result))
121      confusion_matrix = np.zeros((n, n), dtype=int)
122      for expected_label, predicted_label in zip(expected_result, labels):
123          confusion_matrix[int(expected_label)][int(predicted_label)] += 1
124      return confusion_matrix
125
126
127  def display_confusion_matrices(train_matrix, validation_matrix, metric):
128      fig, axs = plt.subplots(1, 2, figsize=(10, 10))
129      axes = axs.flatten()
130      sns.heatmap(train_matrix, annot=True, fmt='d', ax=axes[0])
131      axes[0].set_xlabel("Predicted Labels")
132      axes[0].set_ylabel("True Labels")
133      axes[0].set_title(f"Training Set Confusion Matrix ({metric})")
134      sns.heatmap(validation_matrix, annot=True, fmt='d', ax=axes[1])
135      axes[1].set_xlabel("Predicted Labels")
136      axes[1].set_ylabel("True Labels")
137      axes[1].set_title(f"Validation Set Confusion Matrix ({metric})")
138      plt.tight_layout()
139      plt.show()
140
141
142  def run_knn():
143      # Parse the MNIST dataset
144      mnist_training_data = read_data("mnist_train.csv")
145      mnist_testing_data = read_data("mnist_test.csv")
146      mnist_validation_data = read_data("mnist_valid.csv")
147
148      print(
149          f'Training Data Size: {len(mnist_training_data)}\n' +
150          f'Testing Data Size: {len(mnist_testing_data)}\n' +
151          f'Validation Data Size: {len(mnist_validation_data)}'
152      )
153
154      # Before using training data, we may need to run dimensionality reduction on it
155      # to reduce the number of features. We should try reduce() that we wrote
156      # but we should try other methods that the assignment reccomends as well:
157      # grayscale to binary, dimension scaling, etc.
158      reduced_training_data, train_features = reduce_data(mnist_training_data)
159      reduced_testing_data, test_features = reduce_data(mnist_testing_data)
160      reduced_validation_data, valid_features = reduce_data(
161          mnist_validation_data)
162
163      test_query = reduce_query(mnist_testing_data, train_features)
164      valid_query = reduce_query(mnist_validation_data, train_features)
165
```

```python
166        # Run training data through KNN and receive the labels for each query
167        # We might have to modify KNN so the query is [label, list(pixels)] instead of just
      ↪    list(pixels)
168        # so that we can compare the assigned label to the actual label
169        # Not actually sure if this is how we do this
170        predicted_labels = knn(
171            train=reduced_training_data,
172            query=test_query,
173            metric='euclidean',
174            k=5
175        )
176
177        training_matrix = generate_confision_matrix(
178            predicted_labels, [q[0] for q in mnist_testing_data]
179        )
180
181        (accuracy, precision, recall, f1_score) = evaluate_knn_accuracy(
182            labels=predicted_labels,
183            query=test_query
184        )
185
186        print(
187            f'Test Data Metrics (euclidean):\n' +
188            f'Accuracy: {accuracy}\n' +
189            f'Precision: {precision}\n' +
190            f'Recall: {recall}\n' +
191            f'F1 Score: {f1_score}'
192        )
193
194        validation_matrix = generate_confision_matrix(
195            predicted_labels, [q[0] for q in mnist_validation_data]
196        )
197
198        (accuracy, precision, recall, f1_score) = evaluate_knn_accuracy(
199            labels=predicted_labels,
200            query=mnist_validation_data
201        )
202
203        print(
204            f'Validation Data Metrics (euclidean):\n' +
205            f'Accuracy: {accuracy}\n' +
206            f'Precision: {precision}\n' +
207            f'Recall: {recall}\n' +
208            f'F1 Score: {f1_score}'
209        )
210
211        display_confusion_matrices(training_matrix, validation_matrix, 'euclidean')
212
213        predicted_labels = knn(
214            train=reduced_training_data,
215            query=test_query,
216            metric='cosim',
217            k=5
218        )
219
220        training_matrix = generate_confision_matrix(
221            predicted_labels, [q[0] for q in mnist_testing_data]
222        )
223
224        (accuracy, precision, recall, f1_score) = evaluate_knn_accuracy(
225            labels=predicted_labels,
```

```
226            query=reduced_testing_data
227        )
228
229        print(
230            f'Test Data Metrics (cosim):\n' +
231            f'Accuracy: {accuracy}\n' +
232            f'Precision: {precision}\n' +
233            f'Recall: {recall}\n' +
234            f'F1 Score: {f1_score}'
235        )
236
237        validation_matrix = generate_confision_matrix(
238            predicted_labels, [q[0] for q in mnist_validation_data]
239        )
240
241        (accuracy, precision, recall, f1_score) = evaluate_knn_accuracy(
242            labels=predicted_labels,
243            query=reduced_validation_data
244        )
245
246        print(
247            f'Validation Data Metrics (cosim):\n' +
248            f'Accuracy: {accuracy}\n' +
249            f'Precision: {precision}\n' +
250            f'Recall: {recall}\n' +
251            f'F1 Score: {f1_score}'
252        )
253
254        display_confusion_matrices(training_matrix, validation_matrix, 'cosim')
255
256
257    if __name__ == "__main__":
258        run_knn()
```

## 3.3   K-Means Clustering Implementation

```
1   import numpy as np
2   from distance_metrics import euclidean, cosim
3   from starter import reduce_data, reduce_query, read_data
4
5
6   np.random.seed(30)
7
8
9   def initialize_centroids(k, data):
10       centroids = []
11
12       ind = np.random.randint(0, len(data))
13       centroids.append(data[ind][1])
14
15       for i in range(1, k):
16           distances = []
17           for x in data:
18               to_centroid = []
19               for c in centroids[:i]:
20                   to_centroid.append(euclidean(x[1], c))
21               distances.append(np.min(to_centroid))
22
23           distances = np.array(distances)
24
```

```python
25          probabilities = distances**2
26          probabilities /= np.sum(probabilities)
27
28          next_centroid = np.random.choice(len(data), p=probabilities)
29          centroids.append(data[next_centroid][1])
30
31      return np.array(centroids)
32
33
34  # returns a list of labels for the query dataset based upon observations in the train
    ↪   dataset.
35  # labels should be ignored in the training set
36  # metric is a string specifying either "euclidean" or "cosim".
37  # All hyper-parameters should be hard-coded in the algorithm.
38  def kmeans(train, query, metric, k=10, threshold=0.01):
39      max_iters = 100
40      labels = []
41      train_reduced, removed_features = reduce_data(train, threshold=threshold)
42      centroids = initialize_centroids(k, train_reduced)
43
44      for it in range(max_iters):
45          distances = []
46          for c in centroids:
47              centroid_dist = []
48              for x in train_reduced:
49                  if metric == "euclidean":
50                      centroid_dist.append(euclidean(c, x[1]))
51                  elif metric == "cosim":
52                      centroid_dist.append(cosim(c, x[1]))
53              distances.append(np.array(centroid_dist))
54          distances = np.array(distances)
55
56          # Assign clusters based on minimum distance to centroids
57          if metric == "euclidean":
58              clusters = np.argmin(distances, axis=0)
59          elif metric == "cosim":
60              clusters = np.argmax(distances, axis=0)
61
62          new_centroids = []
63          for i in range(k):
64              cluster_group = []
65              for j in range(len(train_reduced)):
66                  if clusters[j] == i:
67                      cluster_group.append(train_reduced[j][1])
68
69              cluster_group = np.array(cluster_group)
70
71              if len(cluster_group) > 0:
72                  centroid = cluster_group.mean(axis=0)
73              else:
74                  centroid = np.zeros(len(train_reduced[0][1]))
75
76              new_centroids.append(centroid)
77
78          new_centroids = np.array(new_centroids)
79
80          if np.all(new_centroids == centroids):
81              print(f"Exited at {it}")
82              break
83          else:
84              centroids = new_centroids
```

```
85
86         query_reduced = reduce_query(query, removed_features)
87
88         query_distances = []
89         for c in centroids:
90             centroid_dist = []
91             for q in query_reduced:
92                 if metric == "euclidean":
93                     centroid_dist.append(euclidean(q[1], c))
94                 elif metric == "cosim":
95                     centroid_dist.append(cosim(q[1], c))
96             query_distances.append(np.array(centroid_dist))
97         query_distances = np.array(query_distances)
98
99         if metric == "euclidean":
100            classes = np.argmin(query_distances, axis=0)
101        elif metric == "cosim":
102            classes = np.argmax(query_distances, axis=0)
103
104        for c in classes:
105            labels.append(int(c))
106
107        return labels
108
109
110 def calculate_clustering_accuracy(labels, test_data, k=10):
111     label_mapping = {}
112     correct = 0
113     true_labels = []
114     for x in test_data:
115         true_labels.append(int(x[0]))
116
117     for c in range(k):
118         indices = []
119         for i, x in enumerate(labels):
120             if x == c:
121                 indices.append(i)
122         cluster_labels = []
123         for x in indices:
124             cluster_labels.append(true_labels[x])
125         if len(cluster_labels) > 0:
126             vals, count = np.unique(
127                 np.array(cluster_labels), return_counts=True)
128             common = vals[np.argmax(count)]
129             label_mapping[c] = [common, cluster_labels]
130
131     for key in label_mapping.keys():
132         id = label_mapping[key][0]
133         values = label_mapping[key][1]
134         for val in values:
135             if int(id) == val:
136                 correct += 1
137
138     acc = correct / len(true_labels)
139     print(f"K-Means Clustering Accuracy: {acc}")
140     return acc
141
142
143 def main():
144     mnist_training_data = read_data("mnist_train.csv")
145     mnist_testing_data = read_data("mnist_test.csv")
```

```
146        mnist_validation_data = read_data("mnist_valid.csv")
147
148        labels = kmeans(mnist_training_data,
149                        mnist_testing_data, "euclidean", k=36, threshold=3.0)
150        calculate_clustering_accuracy(labels, mnist_testing_data, k=36)
151
152
153  if __name__ == "__main__":
154      main()
```

## 3.4   Collaborative Filtering Implementation

```
1   import numpy as np
2   import pandas as pd
3   from distance_metrics import cosim
4
5   # Hyperparameters
6
7   K = 10
8   M = 7
9   # Weight values for different metrics of the dataset.
10  rating_weight = 3
11  age_weight = 3
12
13
14  # rating_threshold: is to discretize the value for calculating the precision etc
    ↪   metrics.
15  rating_threshold = 4
16  target_users = [405, 655, 13]
17  user_ids = []
18  similar_users = []
19  similarity_dict = {}
20  ratings_vector_1 = []
21  ratings_vector_2 = []
22
23  total_matrix = pd.read_csv("movielens.txt", delimiter='\t')
24  user_age = dict(zip(total_matrix['user_id'], total_matrix['age']))
25
26  user_a_train = pd.read_csv("train_a.txt", delimiter='\t')
27  user_age[user_a_train['user_id'].values[0]] = user_a_train['age'].values[0]
28  user_b_train = pd.read_csv("train_b.txt", delimiter='\t')
29  user_age[user_b_train['user_id'].values[0]] = user_b_train['age'].values[0]
30  user_c_train = pd.read_csv("train_c.txt", delimiter='\t')
31  user_age[user_c_train['user_id'].values[0]] = user_c_train['age'].values[0]
32
33  user_a_valid = pd.read_csv("valid_a.txt", delimiter='\t')
34  user_b_valid = pd.read_csv("valid_b.txt", delimiter='\t')
35  user_c_valid = pd.read_csv("valid_c.txt", delimiter='\t')
36
37  user_a_test = pd.read_csv("test_a.txt", delimiter='\t')
38  user_b_test = pd.read_csv("test_b.txt", delimiter='\t')
39  user_c_test = pd.read_csv("test_c.txt", delimiter='\t')
40
41  combined_data = pd.concat([user_a_train, user_b_train, user_c_train])
42  combined_data_valid = pd.concat([user_a_valid, user_b_valid, user_c_valid])
43  combined_data_test = pd.concat([user_a_test, user_b_test, user_c_test])
44
45
46  user_total_matrix = total_matrix.pivot_table(index='user_id', columns='movie_id',
    ↪   values='rating').fillna(0)
```

```
47
48
49   user_item_matrix = combined_data.pivot_table(index='user_id', columns='movie_id',
     ↪   values='rating').fillna(0)
50   user_valid_matrix = combined_data_valid.pivot_table(index='user_id', columns='movie_id',
     ↪   values='rating').fillna(0)
51   user_test_matrix = combined_data_test.pivot_table(index='user_id', columns='movie_id',
     ↪   values='rating').fillna(0)
52
53
54   # print("user item matrix: ", user_item_matrix)
55   # print("user total matrix: ", user_total_matrix)
56
57   print("user ids: ", user_ids)
58
59   def normalised_metric_similarity(user_id_first:int, user_id_second:int):
60       max_age_difference = max(user_age.values()) - min(user_age.values())
61       age_similarity = 1 - abs(user_age.get(user_id_first) - user_age.get(user_id_second))
         ↪   / max_age_difference
62       return age_similarity
63
64   #Calculate similarity between users
65   for id1 in target_users:
66       print("user id: ", id1)
67       for id2 in user_total_matrix.index:
68           if id1 == id2:
69               continue
70           for movie_id in user_total_matrix.columns:
71               if movie_id in user_item_matrix.columns and movie_id in
                 ↪   user_total_matrix.columns:
72                   rating1 = user_item_matrix.loc[id1, movie_id]
73                   rating2 = user_total_matrix.loc[id2, movie_id]
74                   if rating1 > 0 and rating2 > 0:
75                       ratings_vector_1.append(rating1)
76                       ratings_vector_2.append(rating2)
77
78           # Calculate the similarity
79           rating_sim = cosim(np.array(ratings_vector_1), np.array(ratings_vector_2))
80           age_sim = normalised_metric_similarity(id1,id2)
81           user_sim = rating_sim * rating_weight + age_sim * age_weight
82           similar_users.append((id2, user_sim))
83
84       similar_users = sorted(similar_users, key=lambda x: x[1], reverse=True)
85       # Keep only top K similar users
86       similarity_dict[id1] = similar_users[:K]
87
88   # Recommend movies for each user
89   true_positive = 0
90   false_positive = 0
91   false_negative = 0
92
93   for user_id in target_users:
94       weighted_ratings = {}
95
96       top_k_users = similarity_dict[user_id][:K]
97       #print("Top K Users: ", top_k_users)
98       #target_user_ratings = user_item_matrix.loc[user_id]
99
100      for sim_user_id in top_k_users:
101          sim_user_ratings = user_total_matrix.loc[sim_user_id[0]]
102          similarity_score = sim_user_id[1]
```

```
103
104            for movie_id in user_total_matrix.columns:
105                rating = sim_user_ratings[movie_id]
106                if movie_id in user_item_matrix.columns:
107                    if user_item_matrix.loc[user_id][movie_id] == 0 and rating > 0:
108                        if movie_id not in weighted_ratings:
109                            weighted_ratings[movie_id] = 0.0
110                        weighted_ratings[movie_id] += rating * similarity_score
111
112        best_movies = sorted(weighted_ratings, key=weighted_ratings.get, reverse=True)[:M]
113
114        print(f"The best movies for user {user_id} based off of {K} similar users:
       ↪  {best_movies}")
115
116        # Evaluate using the validation set
117        for best_movie_id in best_movies:
118            if best_movie_id in user_test_matrix.columns:
119                actual_rating = user_test_matrix.loc[user_id, best_movie_id]
120                if actual_rating >= rating_threshold:
121                    true_positive += 1
122                else:
123                    false_positive += 1
124
125                user_ratings = user_test_matrix.loc[user_id]
126                for test_movie_id, user_rating in user_ratings.items():
127                    if test_movie_id not in best_movies and user_rating >= rating_threshold:
128                        false_negative += 1
129
130  # Evaluation Metrics
131  precision = true_positive / (true_positive + false_positive)
132  recall = true_positive / (true_positive + false_negative)
133  f1_score = (precision * recall) / (precision + recall)
134
135  print(f"Precision: {precision}\nRecall: {recall}, \nF1 Score: {f1_score}")
```

### 3.4.1 Dimensionality Reduction Implementation

```
1   import numpy as np
2   from copy import deepcopy
3
4
5   np.random.seed(30)
6
7
8   def reduce_data(data_set, threshold=0.01):
9       """ Returns the reduced dataset using variance thresholding
10
11      Args:
12          data_set (ndarray(int, ndarray)): processed data
13          threshold (float): variance threshold, default 0.01
14
15      Returns:
16          tuple of (ndarray(int, ndarray), list): Reduced dataset and removed features
17      """
18      data_cp = deepcopy(data_set)
19      features = np.array([feature[1] for feature in data_cp])
20      variances = np.var(features, axis=0)
21      removed_features = [index for index, variance in enumerate(
22          variances) if variance < threshold]
23
```

```
24        for entry in data_cp:
25            entry[1] = np.delete(entry[1], removed_features)
26
27        return data_cp, removed_features
28
29
30  def reduce_query(data_set, removed_features):
31      """ Returns the reduced query point
32
33      Args:
34          (int, ndarray): image
35          removed_features (list): list of removed features
36
37      Returns:
38          (int, ndarray): Reduced image
39      """
40      query_cp = deepcopy(data_set)
41      for entry in query_cp:
42          entry[1] = np.delete(entry[1], removed_features)
43
44      return query_cp
45
46
47  def read_data(file_name: str) -> list:
48
49      data_set = []
50      with open(file_name, 'rt') as f:
51          for line in f:
52              line = line.replace('\n', '')
53              tokens = line.split(',')
54              label = tokens[0]
55              attribs = []
56              for i in range(784):
57                  attribs.append(tokens[i+1])
58              data_set.append([label, np.array(attribs, dtype=float)])
59      return (data_set)
60
61
62  def show(file_name, mode):
63
64      data_set = read_data(file_name)
65      for obs in range(len(data_set)):
66          for idx in range(784):
67              if mode == 'pixels':
68                  if data_set[obs][1][idx] == '0':
69                      print(' ', end='')
70                  else:
71                      print('*', end='')
72              else:
73                  print('%4s ' % data_set[obs][1][idx], end='')
74              if (idx % 28) == 27:
75                  print(' ')
76          print('LABEL: %s' % data_set[obs][0], end='')
77          print(' ')
```