

CENG 409

Introduction to Machine Learning

Fall '2023-2024

Homework 2

Due date: December 6, 2023, 23:55

Objectives

This assignment aims to fulfill the following objectives:

- To familiarize you with the cross-validation technique, which is crucial for model tuning and evaluation, and to enable you to gain hands-on experience with it.
- To provide hands-on experience with the K Nearest Neighbor method and how to tune its hyperparameters with cross-validation.
- To familiarize you with Kmeans and Kmedoid algorithms at the implementation level for clustering problems and provide hands-on experience in pinpointing the number of clusters with the elbow method.
- To familiarize you with the Hierarchical Agglomerative Clustering algorithm and Silhouette method for clustering and evaluating clustering results, respectively.
- To familiarize you with the t-SNE, PCA, and UMAP dimensionality reduction methods and provide hands-on experience to visualize data instances in 1-D, 2-D, and 3-D spaces.

The assignment involves both coding and writing a report. You can prepare your report in Latex or Word (the report format is up to you). There should be a single report submitted for all parts.

Part 1

When a problem to be solved has abundant data to form training, validation, and test datasets, a machine model could be trained on the training dataset, its hyperparameters (if there are any) could be determined with the validation dataset, and its expected generalization performance could be measured on the test dataset. Separation of a dataset into three reduces the data amount for training, especially if the dataset contains a relatively small number of instances (the fewer instances there are in a training dataset, the worse generalization performance we expect since the training dataset fails to represent the actual data generation process's data distribution). One way to alleviate this problem is to omit the validation dataset and perform model training and testing on a training dataset and a test dataset (instead of dividing a whole dataset into three, we divide it into two parts, namely, the training dataset and testing dataset). The test dataset is used to assess the trained model's expected generalization performance.

One crucial point regarding the test dataset is that it should contain enough instances to yield reliable results. With the training dataset, we need to carry out both the parameter search and hyperparameter tuning procedures. We could extract a validation dataset, which will be arbitrary and reduce the training sample amount. We could pick instances randomly to form a validation dataset, but by chance, we may have selected some outlier or noisy data instances, so the results attained on this single validation dataset may not be reliable for parameter tuning. Instead, we can form multiple validation datasets and average results achieved on them for parameter tuning. This is the fundamental approach considered by the cross-validation technique. The training dataset is divided into K parts. For each hyperparameter configuration, $K-1$ parts are used for training the model, and the remaining partition is used to assess the performance of the model with that hyperparameter configuration in an iterative manner (first 1st partition is set as the test dataset, 2nd, 3rd, 4th... K^{th} partitions are used for training. Second, for the same hyperparameter value, 2nd partition is considered to be the test dataset, and the rest (1st, 3rd, 4th... K^{th}) is regarded as the training dataset, this procedure is repeated for each partition and each hyperparameter configuration). In total, K many test datasets and training datasets are used for each hyperparameter configuration, and the average performance score over these test datasets provides an idea about the generalization capability of that hyperparameter configuration. This average score can be used to pick the best hyperparameter configuration. After determining the best-performing hyperparameter values, the model can be trained with all K partitions, and the trained model's generalization performance could be measured on the original test dataset (the test dataset we set aside just before the beginning of cross-validation). The cross-validation procedure is summarized in the following figure.

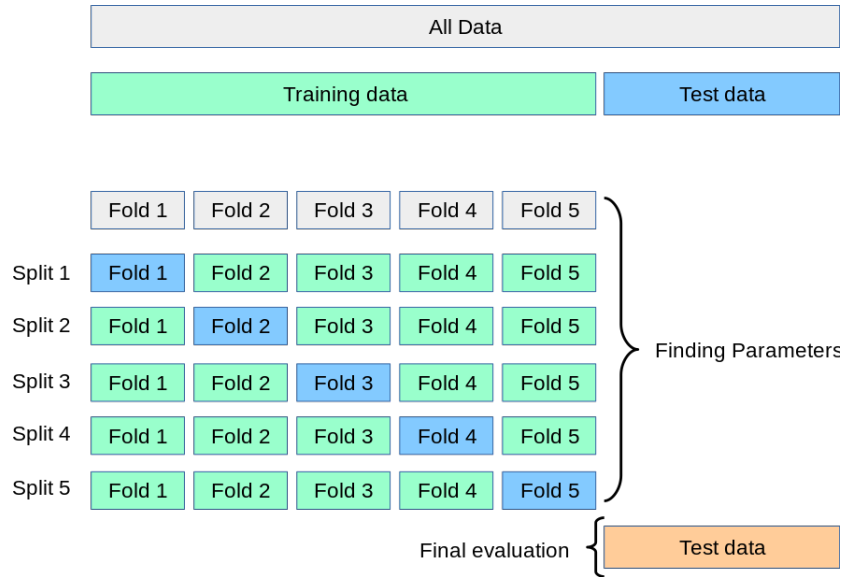


Figure 1: Cross-validation technique applied on a dataset for hyperparameter tuning. A dataset is divided into training (referred to as TRAIN1) and test (referred to as TEST1) parts (the training part is used for both training the model and tuning hyperparameters, whereas the test part is solely used to assess the generalization performance of the trained model. The training part is further divided into K parts, and for each hyperparameter configuration, one part is designated as a test dataset, and the remaining ones are regarded as a training dataset in an iterative manner. Depending on the average performance score attained on these partitions, the best hyperparameter configuration could be determined. After identification of the best-performing hyperparameter values, the model is trained on the training part (TRAIN1), and its final performance (generalization) score is calculated on the test part (TEST1). Generally, K values of 5, 10, or 30 are used. The image has been taken from [1]

Since $K-1$ parts are used during training, the reduction in the number of training samples is smaller compared to designating a separate validation dataset. In addition, because many validation datasets are used to assess the performance of a hyperparameter configuration, a more reliable (unbiased) result is obtained. Of course, all these come with an extra computation cost (for each hyperparameter configuration, the model must be trained K times). If the model being trained involves some randomness, the cross-validation procedure should be repeated many times to eliminate random results (even if the model is deterministic, it is still better to repeat cross-validation by shuffling the original dataset to obtain more reliable results further). In other words, the results attained during cross-validation should be statistically significant. Typically, K values of 5, 10, or 30 are considered.

Despite all its merits, there are mainly two pitfalls that should be considered with cross-validation. The first one is that while partitioning the training part into K parts, there is a risk of introducing a class imbalance problem, especially if we pick partition instances randomly (for example, for a classification task, for one class, no data instance may be present in a partition which is considered to be a test dataset). The partitions should be formed depending on the original class distribution of the training part to eliminate this problem. Such a strategy is coined as stratification [2] (Chapter 19). The second pitfall emerges when the cross-validation technique is used on a whole dataset for both model hyperparameter tuning and model performance assessment (without sparing a test dataset out of the entire dataset). Especially if the problem to be solved features a small dataset and sparing a portion of it as a test dataset may not be possible (Sample size may be too small), direct application of the cross-validation technique could be considered on the whole dataset. With such an approach, the same dataset is used for both hyperparameter tuning and model evaluation, which leads to optimistic evaluation results [3, 4, 5]. Since the hyperparameters are tuned to perform well on the dataset, the model can be considered to overfit the dataset. Since the evaluation is done on the same dataset, the results may be misleading [4] (It may not provide reliable results about the model's generalization capability, it is likely to be optimistic). This serious problem can be alleviated via the nested cross-validation technique [4, 5]. The K -nearest neighbor (KNN) method is a lazy method (also, it is a non-parametric density estimation method [2] (Chapter 18)) which can be used for both classification and regression problems, it is "lazy" in the sense that it does not require training (it does not have parameters to tune like MLPs, non-parametric). For a given instance whose label is not known, it finds the K nearest neighbors of that instance from a training data set depending on a distance/similarity metric/measure and labels the given data sample. Labeling can be done via majority voting or weighted majority voting (some other strategies could also be utilized; for instance, when labels of the neighbors are considered to be an input, a neural network could be trained for the final decision). In the first strategy, every nearest neighbor has equal importance (their decisions count equally) on the final label decision, and the data sample is labeled with the label that has been voted the most by the neighbors. In the second strategy (a weighted score is calculated for each K nearest neighbor depending on their distance/similarity to the given data sample), closer neighbors have more contribution (since they are closer to the given data sample or they are more similar to it) to the final decision. For classification problems, with majority voting, the class that attains the most votes is considered the predicted label of a given point. With weighted majority voting, the class label with the highest cumulative weight is returned as the resulting label (different weight combination strategies could also be considered). For a regression task, with majority voting, values of the K nearest neighbors are averaged. With weighted voting, a weighted average of the values can be calculated as the label of a given instance (other combination schemes could be considered, too).

KNN requires tuning two main hyperparameters: K and distance/similarity metric. The optimal hyperparameter configuration varies depending on data distribution and the characteristics of the features of data samples (e.g., number of features, discrete/continuous-valued, etc.). For $K=1$, the method is called nearest neighbor (NN), which is a commonly used form of KNN for classification and regression tasks. A given data sample is labeled with the label of the closest data point. As the distance measure, several metrics (and more) could be considered: Minkowski, Minkowski ($p=2$ leads to Euclidian distance), cosine, and Mahalanobis. The cosine similarity is based on the angle between two vectors. The smaller the angle between two vectors, the more similar these vectors are to each other. Mahalanobis distance [6] takes

into account the possible feature correlations (which affect the data distribution) of data samples by incorporating the covariance matrix in calculating the distance between two data points (one feature may lead data to spread out/elongate in its direction, the covariance matrix captures such distributions) (If the covariance matrix is the identity matrix, Mahalanobis distance becomes equivalent to the Euclidian distance). Here are the formulas for these distance/similarity measures:

$$\text{Minkowski}(\vec{x}, \vec{y}, p) = \left(\sum_{i=1}^d |x_i - y_i|^p \right)^{1/p}, \text{ where } d \text{ is the dimension of } x \text{ and } y \text{ vectors.}$$

$$\text{Cosine}(\vec{x}, \vec{y}) = 1 - \frac{\vec{x} \cdot \vec{y}}{\|\vec{x}\| * \|\vec{y}\|}, \text{ where } \cdot \text{ denotes the dot product operation.}$$

$$\text{Mahalanobis}(\vec{x}, \vec{y}, S^{-1}) = \sqrt{(\vec{x} - \vec{y})^T S^{-1} (\vec{x} - \vec{y})}, \text{ where } S^{-1} \text{ is the inverse of the covariance matrix of a dataset.}$$

Part 1 Specifications

- You are expected to implement the following class (whose file name should be **Distance.py**) for the distance/similarity metrics:

```

1 class Distance:
2     @staticmethod
3     def calculateCosineDistance(x, y):
4         pass
5     @staticmethod
6     def calculateMinkowskiDistance(x, y, p=2):
7         pass
8     @staticmethod
9     def calculateMahalanobisDistance(x,y, S_minus_1):
10        pass

```

`S_minus_1` parameter of `calculateMahalanobisDistance` specifies the inverse of the covariance matrix of a dataset. This matrix can be calculated with **numpy.cov** [7] and **numpy.linalg.inv** [8] for a given dataset. All the datasets in this assignment are 2D numpy arrays. Each data sample (instance) is represented with a row vector.

- For the KNN method, you are expected to implement the following class (whose file name should be **Knn.py**):

```

1 class KNN:
2     def __init__(self, dataset, data_label, similarity_function,
3 similarity_function_parameters=None, K=1):
4         """
5         :param dataset: dataset on which KNN is executed, 2D numpy array
6         :param data_label: class labels for each data sample, 1D numpy
7 array
8         :param similarity_function: similarity/distance function, Python
9 function
10        :param similarity_function_parameters: auxiliary parameter or
11 parameter array for distance metrics
12        :param K: how many neighbors to consider, integer
13        """
14        self.K = K
15        self.dataset = dataset
16        self.dataset_label = data_label
17        self.similarity_function = similarity_function
18        self.similarity_function_parameters =
19 similarity_function_parameters

```

```

15
16     def predict(self, instance):
17         pass

```

The `similarity_function` is assigned to one of the `Distance` class functions, `similarity_function_parameters` variable holds additional parameters for the similarity/distance function (e.g., it holds p for the Minkowski distance or S^{-1} for the Mahalanobis distance). The `predict` function takes a data instance (as numpy array) and returns the predicted label for that instance via **majority voting** among the K nearest neighbors.

- You are provided with a classification problem dataset and a partial implementation that loads the dataset (**Knnexperiment.py**). On this dataset, you are expected to perform 10-fold cross-validation (with stratification) for hyperparameter tuning (via grid search). You are free to determine which hyperparameter configurations to test (there should be at least five configurations). Please repeat this cross-validation procedure five times (by shuffling the original dataset) and compute confidence intervals of **accuracy** performance scores for each hyperparameter configuration. In your report, please specify these hyperparameter values and attained confidence intervals for each hyperparameter configuration. In addition, please add some comments on how you have picked the best-performing hyperparameter values. For stratified cross-validation, you may refer to the following Scikitlearn documentation [9]. For data shuffling with numpy, you may refer to [10].
- You may prefer to consider other performance metrics (it is up to you) (e.g., precision, recall, f1 score, etc.).
- Your whole experimentation code for this part should be implemented in **Knnexperiment.py**.
- For this part, you are expected to make use of your own implementation for KNN. Library implementations are not allowed.

Note: Although the KNN method does not feature any randomness (multiple runs of a particular hyperparameter setting will yield the same/identical results), in this part, by repeating the cross-validation procedure multiple times, we aim to reduce the dependency on a single fixed partitioning, whose results may be biased and misleading. By shuffling the dataset, we introduce different partitionings to the method, and by averaging the attained performance scores, we obtain more statistically reliable results. In other words, we want our results not to depend on a single partitioning of the dataset.

Part 2

The unsupervised learning paradigm aims to extract knowledge from data samples lacking label information. One crucial problem in this paradigm is to cluster data samples depending on their similarities/dissimilarities/closeness and extract insightful knowledge from these groupings. To this end, this and the following parts are devoted to the clustering problem and clustering analysis. The Kmeans algorithm is a special case of the expectation-maximization procedure where the likelihood of data samples is maximized with respect to a model parameter by considering hidden/latent variables (the factors that we don't observe directly) [2] (Chapter 7), [11] (Chapter 9). It consists of two main steps: the expectation (E) step and the maximization (M) step. The procedure starts with an initial guess of the parameters. Later, the procedure maximizes the likelihood of data samples by iterating through the E and M steps. In the E step, hidden variables are estimated with the current parameter values, and these estimated hidden variable values are used to improve the current parameter values. For Kmeans, the information that specifies a data point belongs to which cluster (cluster label) is a hidden variable, and the mean vectors of clusters are parameters. Kmeans begins its operations by initializing the cluster means randomly. Then, it estimates the cluster labels of data samples by finding the closest cluster mean among the current

cluster mean vectors and assigning them to the nearest cluster (it labels the data samples with the label of the closest cluster mean) (this constitutes the E-step for Kmeans, we estimate labels of data samples by using the current cluster means). Next, the current cluster mean vectors are updated with the current cluster labels (an average (mean) vector of data points in a cluster is calculated. This constitutes the M step(the current cluster means are updated with the current cluster labels). In other words, in the E step, the current mean values are fixed, and the cluster labels are updated. The cluster labels are fixed in the M step, and the cluster means are updated.

For Kmeans, these operations can be formulated as the following loss function [2] (Chapter 7):

$$L = \sum_s^n \sum_k^K l_k^s \|x_s - \mu_k\|^2 \quad (1)$$

where n is the number of data samples (x_s is the s^{th} data sample), K is the number of clusters (μ_k is the mean vector of the k^{th} cluster), l_k^s specifies cluster label for a data sample. If an instance s belongs to a cluster k , l_k^s is 1 for that cluster (for the other clusters, it is 0). Minimizing this loss function corresponds to finding compact clusters whose data points are close to their cluster centers and not scattered around. Furthermore, we can interpret the loss function as the sum of inner cluster variances [12] (Chapter 5) ($Variance = \mathbb{E}[(x - \mu)^2]$). Hence, Kmeans aims to minimize variance within each cluster (it aims to form clusters whose inner variances are minimal).

As the main distance/similarity measure Euclidean distance is considered for Kmeans (the method requires a centroid notation to work (which is available in Euclidean space), datasets may contain data features for which averaging/calculation mean is meaningless, e.g., marital status: single, married (mean/average cannot be defined))). This introduces bias to Kmeans and leads Kmeans to perform well on the spherical-shaped datasets (if clusters are of arbitrary shape, it is highly likely to perform poorly). Such a restriction leads to the Kmedoid method [13] (4), which considers a data sample as the cluster representative rather than the mean of cluster data points (for Kmeans, clusters are summarized/identified as a mean vector of data points, for Kmedoids a single data point is considered). In the M step, instead of calculating a centroid (mean of cluster data samples), the data point closest to all cluster data points is picked to be the new cluster center. Similarly, data points are labeled in the E step depending on their closeness to these cluster representatives. With this method, other distance/similarity metrics could be considered for forming clusters. The loss function for the method is a generalization of that of Kmeans:

$$L = \sum_s^n \sum_k^K l_k^s d(x_s, c_k) \quad (2)$$

where $d(x, y)$ is an arbitrary distance/similarity metric (function, possibly not Euclidean distance), c_k is the representative data point for the cluster k . Minimizing this loss function corresponds to forming clusters whose data points are close to the cluster representative data point (more similar data points are grouped).

The expectation maximization procedure may get stuck in local optima while maximizing the likelihood of data samples. The similar handicap is valid for Kmeans (since it is a specific instance of the EM procedure). So Kmeans is not guaranteed to yield an optimal clustering result for a given dataset (neither is Kmedoid).

Kmeans (also Kmedoid) strongly depends on the initially picked cluster centers while generating a final clustering result. Several strategies could be considered to this end. We could run Kmeans (by selecting the initial cluster centers randomly) many times (e.g., 50 times) and pick the one that attains the smallest loss error (Equation 1) after running. Or we can choose instances that are distant from each other depending on a fixed distance threshold value. A better approach is from the Kmeans++ algorithm [14]: picking initial clusters with a probability score depending on the distance between data samples. Kmeans++ differs from Kmeans solely in determining initial cluster representatives. As the first step, Kmeans++ picks an instance randomly as the first cluster center. The rest of the cluster centers (the

remaining $K-1$ centers) are sampled with a probability score depending on the distances between data samples and already picked cluster centers (D^2 weighting). After selecting all cluster centers, standard E and M steps of Kmeans are applied. With this strategy, Kmeans++ has been shown to complete faster and yield better clustering results compared to Kmeans [14].

Along with how to pick initial cluster centers, the most critical hyperparameter of Kmeans is the number of clusters to form. If we have knowledge about the data generation process, we might know the exact number of clusters beforehand (e.g., in a handwritten digit dataset, ten different clusters are expected). If we cannot access such information, we can project the dataset into a lower dimension (e.g., 2D, 3D) and visually determine the number of clusters from the projected data. Another possibility is to apply information-theoretic measures [15]. As a last resort (but not the least), we can use the elbow method, which requires plotting the clustering loss (Equation 1) for possible several K values and suggests picking the K value where an elbow shape occurs (after this point we do not further see a sharp decrease in the loss score).

Both K-means and Kmeans++ algorithms utilize the Euclidean distance to measure similarity/dissimilarity between data instances, which restricts their applicability to data instances with many features (high dimensional) and mixed feature values (discrete and continuous). From the Euclidean distance perspective, data instances look alike in a high dimensional space, which is a severe problem likely to lead to inaccurate clustering results. For example, let's consider the following 4-D data instances: $[3, 3, 3, 0]$, $[0, 1, 0, \sqrt{6}]$, $[0, 0, 0, \sqrt{5}]$. These instances are equally distant to $[1, 1, 1, 0]$. As can be seen, this data instance is a scaled form of the first data instance. We expect it to be more similar to the first data instance (they share the same direction if we regard them as vectors), but due to the Euclidean distance measure, the other data instances are not different from the first one. Another important consideration is that when some features of a data instance have discrete values, we cannot utilize the Euclidean distance to quantify similarity/dissimilarity between data instances. For example, one feature may represent occupation or marital status, and we cannot calculate a squared distance between different occupations or marital statuses. On the other hand, since Kmedoids is capable of working with an arbitrary similarity/distance measure, it can easily be employed with these types of data instances. Furthermore, since other distance metrics could be utilized, clustering results may be better than those obtained via Euclidean distance. However, such an appealing merit comes with a computation overhead. Kmeans calculates the cluster centers by averaging data instances in clusters. In contrast, Kmedoids has to measure the similarity between each possible data instance pair in clusters to find cluster centers.

As humans, we are very accustomed to 1-D, 2-D, and 3-D spaces and understand relationships between objects and their behavior in these spaces. When data instances have a dimensionality higher than 3, we cannot fully observe them, and our common sense knowledge about 1-D, 2-D, and 3-D spaces may no longer apply to those data instances. For instance, most of the volume of a sphere with radius 1 in high dimensions lies in a thin shell near the surface of the sphere [11] (Chapter 1), which is against our common sense with spheres in 3-D. Data instances for different problem types can have different dimensionality/number of features, and every feature of data instances may not help solve a problem. The number of features plays an essential role in the computation and storage demands of a problem. It directly or indirectly determines the complexity of the problem, thus the number of parameters in an ML algorithm. For instance, if data instances of the problem/task lie in a high dimensional space, say 500, the function to be learned by an MLP to solve the problem should be more complex compared to the function when they had lower dimensionality, say 3, (it is more complex in the sense that there are more dimensions thus the search space for solution functions is larger. In this regard, more weights and hidden layers may be needed to learn the complex function). A similar argument could be made for other methods, such as decision trees, SVMs, etc as well. The more features, the more branching there are in a decision tree (the bigger the tree is).

To reduce the complexity and storage demands of a problem, we can utilize dimensionality reduction methods [2] (Chapter 6). After reducing the dimensionality of data instances, we can use less complex ML algorithms to solve the problem. However, dimensionality reduction could be considered a data compression operation since it tries to represent a higher dimensional data point with a lower dimensional

data point, and information loss could occur. Thus, helpful information for a problem solution could be lost, which can render the problem unsolvable (algorithms can perform very poorly due to the fact that data instances do not contain enough information after the reduction operation). In general, the success of applying dimensionality reduction operation is task/problem, data (dataset), and algorithm dependent.

We have previously mentioned that one way to determine the number of clusters in a clustering problem is to utilize dimensionality reduction methods. Simply, we can map/transform original data instances as 1-D, 2-D, or 3-D data instances and visualize them. We can pick a K value depending on the number of clusters formed in the visualization. The visualization quality depends on data instances (their distribution) and the dimensionality reduction method utilized. Various dimensionality reduction algorithms handle the dimensionality reduction problem by approaching it with different assumptions such as Principal Component Analysis (PCA), Multidimensional Scaling, Linear Discriminant Analysis (LDA), Isomap, Locally Linear Embedding (LLE) [2] (Chapter 6), Self-Organizing Maps (SOM) [12] (Chapter 9), T-distributed Stochastic Neighbor Embedding (t-SNE) [16, 17], Uniform Manifold Approximation and Projection (UMAP) [18, 19], Autoencoder [20, 21]. Since there are many within the scope of this assignment, we focus on the three widely used methods: PCA, t-SNE, and UMAP. PCA is the simplest among them and aims to linearly project data instances onto particular vectors such that the projected data instances have the highest variance (scatter) [2] (Chapter 6), which turns out to find eigenvectors of the covariance matrix of data instances. T-SNE, on the other end, aims to model the distribution of data instances in a high dimension and maps these data instances into a low dimension (1-D, 2-D, or 3-D) by preserving the modeled data distribution [16]. Like t-SNE, UMAP aims to form a topological representation of data instances in a high dimension and map these data instances to a low dimensional space by preserving the modeled representation [18].

In this part, you are expected to implement Kmeans and Kmedoids algorithms (by randomly picking initial cluster centers for both methods) and employ them on two datasets by utilizing the elbow technique and dimensionality reduction methods (namely, PCA, t-SNE, and UMAP).

```

Initialize  $\mathbf{m}_i, i = 1, \dots, k$ , for example, to  $k$  random  $\mathbf{x}^t$ 
Repeat
  For all  $\mathbf{x}^t \in \mathcal{X}$ 
     $b_i^t \leftarrow \begin{cases} 1 & \text{if } \|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$ 
  For all  $\mathbf{m}_i, i = 1, \dots, k$ 
     $\mathbf{m}_i \leftarrow \sum_t b_i^t \mathbf{x}^t / \sum_t b_i^t$ 
Until  $\mathbf{m}_i$  converge

```

Figure 2: Pseudo code for Kmeans [2] (Chapter 7). Equivalently, b_i^t s could be checked for convergence. When none of b_i^t s change, the algorithm terminates.

- 1a. Take one center c_1 , chosen uniformly at random from \mathcal{X} .
- 1b. Take a new center c_i , choosing $x \in \mathcal{X}$ with probability $\frac{D(x)^2}{\sum_{x \in \mathcal{X}} D(x)^2}$.
- 1c. Repeat Step 1b. until we have taken k centers altogether.
- 2-4. Proceed as with the standard **k-means** algorithm.

Figure 3: Pseudo code for Kmeans++ [14]. $D(x)$ denotes the shortest distance from data point x to the closest cluster center, which has already been chosen.

1. Initialize exemplars: $\{z^{(1)}, \dots, z^{(k)}\} \subseteq \{x^{(1)}, \dots, x^{(n)}\}$ (exemplars are k points from the original dataset)
2. Repeat until there is no further change in cost:
 - (a) for each j : $C^j = \{i : x^{(i)}\text{'s closest exemplar is } z^{(j)}\}$
 - (b) for each j : set $z^{(j)}$ to be the point in C^j that minimizes $\sum_{i \in C^j} d(x^{(i)}, z^{(j)})$

Figure 4: Pseudo code for Kmedoids [27]

Part 2 Specifications

- For Kmeans, you are expected to implement the following class (whose file name should be **Kmeans.py**):

```

1 class KMeans:
2     def __init__(self, dataset, K=2):
3         """
4         :param dataset: 2D numpy array, the whole dataset to be clustered
5         :param K: integer, the number of clusters to form
6         """
7         self.K = K
8         self.dataset = dataset
9         # each cluster is represented with an integer index
10        # self.clusters stores the data points of each cluster in a
11        dictionary
12        self.clusters = {i: [] for i in range(K)}
13        # self.cluster_centers stores the cluster mean vectors for each
14        cluster in a dictionary
15        self.cluster_centers = {i: None for i in range(K)}
16        # you are free to add further variables and functions to the class
17
18        def calculateLoss(self):
19            """Loss function implementation of Equation 1"""
20
21        def run(self):
22            """Kmeans algorithm implementation"""
23            return self.cluster_centers, self.clusters, self.calculateLoss()

```

The constructor takes two parameters: `dataset` and `K`. The dataset to be fed is a 2D numpy array where each row vector represents a single data point (the total number of instances equals the number of rows of this array). In the constructor, two variables (of type Python dictionary) are defined to store cluster points and centers: `self.clusters` and `self.cluster_centers`. Each cluster is represented with an integer index (the whole index range is 0 to `K-1`). Mainly, you are expected to complete the functions: `calculateLoss` and `run`. The `calculateLoss` function calculates the loss function of Kmeans (Equation 1). The `run` function implements Kmeans, and after running, it returns the cluster centers, cluster data points, and the loss value for the final clustering. You can add other variables and auxiliary functions to the class implementation if you need.

- For Kmedoids, you are expected to implement the following class (whose file name should be **Kmedoids.py**):

```

1 class KMemoids:
2     def __init__(self, dataset, K=2, distance_metric="cosine"):
3         """
4         :param dataset: 2D numpy array, the whole dataset to be clustered
5         :param K: integer, the number of clusters to form

```

```

6         """
7         self.K = K
8         self.dataset = dataset
9         self.distance_metric = distance_metric
10        # each cluster is represented with an integer index
11        # self.clusters stores the data points of each cluster in a
dictionary
12        # In this dictionary, you can keep either the data instance
themselves or their corresponding indices in the dataset (self.dataset)
.
13        self.clusters = {i: [] for i in range(K)}
14        # self.cluster_medoids stores the cluster medoid for each cluster
in a dictionary
15        # # In this dictionary, you can keep either the data instance
themselves or their corresponding indices in the dataset (self.dataset)
.
16        self.cluster_medoids = {i: None for i in range(K)}
17        # you are free to add further variables and functions to the class
18
19    def calculateLoss(self):
20        """Loss function implementation of Equation 2"""
21
22    def run(self):
23        """Kmedoids algorithm implementation"""
24        return self.cluster_medoids, self.clusters, self.calculateLoss()

```

The constructor takes three parameters: `dataset`, `K` and `distance_metric`. The dataset to be fed is a 2D numpy array where each row vector represents a single data point (the total number of instances equals the number of rows of this array). Since Kmedoids can work with different distance measures, the distance measure to be used is determined via the `distance_metric` parameter. For this assignment, you are expected to implement the Kmedoids algorithm with cosine distance. In the constructor, two variables (of type Python dictionary) are defined to store cluster points and cluster centers (medoids): `self.clusters` and `self.cluster_medoids`. Each cluster is represented with an integer index (the whole index range is 0 to K-1). Mainly, you are expected to complete the functions: `calculateLoss` and `run`. The `calculateLoss` function calculates the loss function of Kmedoids (Equation 2). The `run` function implements Kmedoids, and after running, it returns the cluster medoids, cluster data points, and the loss value for the final clustering. If you need to add other variables and auxiliary functions, you can add them to the class implementation. Since the Kmedoids algorithm works solely with data instances (centers/medoids are also data instances), in the variables `self.clusters` and `self.cluster_medoids`, you can store the data instances themselves or their corresponding indices in the dataset (`self.dataset`). Keeping indices may be helpful if we want to create a distance table/dictionary that provides the distance value between two data instances (after calculating a distance value between two data instances, their distance value can be stored in this table/dictionary for future distance queries to reduce computation time, which results in increased memory usage compared to Kmeans).

- You are given two datasets and two Python files, namely: **Kmeansexperiment.py** and **Kmedoidsexperiment.py** (which already load the datasets). In these files, you are expected to implement the elbow method for the methods separately by employing the two datasets. To get a loss value for a particular K value (for both methods), you can run the algorithms ten times (with the same K values) and pick the lowest loss value as the result. To further alleviate randomness, this procedure can be repeated ten times from scratch to obtain an average loss value for each K value (similarly, a confidence interval can be calculated with these scores if need be). Depending on these average loss values, the most suitable K values could be picked for both datasets with both methods. For instance, let's say we have two datasets, A and B; first, we would like to consider Kmeans, and we

would like to determine a suitable K value for A . By starting from 2 to (say) 10 we run Kmeans on A . First, we run for $K=2$ 10 times on A . Since each time different initial cluster centers are picked, it is highly likely to get different loss values for each run. After completing ten runs, we choose the lowest loss score for $K=2$, say α_1 . We repeat the same procedure for the second time ($K=2$, ten runs, picking the smallest loss value) to get α_2 . We repeat this until we get α_{10} . The average loss value calculated from $\alpha_1, \alpha_2, \alpha_3 \dots \alpha_{10}$ is for the loss value of $K=2$ on the dataset A . These steps need to be repeated for the rest of the K values ($K=3, K=4, \dots$). With these average loss values (on dataset A), a K versus loss graph could be drawn (it provides insight into the performance of Kmeans on A), and a suitable K value for A can be determined. Similarly, the whole process can be applied to dataset B , and Kmedoids could be employed similarly for clustering. Confidence intervals calculated via α_i s could also be shown in K versus loss graphs to represent experiment results better and provide more detailed information.

- In your report, please provide the K versus Loss graphs and comment on the best number of clusters for each method and dataset (4 graphs, two datasets, and two algorithms).
- For this part, you are expected to make use of your own implementation for Kmeans and Kmedoids. Library implementations are not allowed.
- For this part, you can consider ten as the highest K value to test during experiments.
- In the report, please provide all plots of the elbow method (K versus Loss) with confidence intervals (you can provide confidence intervals as a table and depict them on K versus loss graphs) and comment about the results attained for both methods and datasets separately (e.g., what is the most suitable cluster number for each dataset (for each method individually)?).
- You are given a Python file named **part2_dimensionalityreduction.py**, which already loads the datasets of this part. In this file, you are expected to employ the three dimensionality reduction methods to visualize the datasets in 2-D. You can try different hyperparameter configurations for t-SNE and UMAP to enhance visualization results (please ensure that the hyperparameter configurations you have tested should be explicit in the file). In the report, please provide the "best" dimensionality reduction results (as matplotlib scatter plots) of each method for both datasets (the definition of "best" is up to you) (if you would like, you can add 3-D visualizations too). Also, please briefly discuss whether the number of clusters you have identified via the elbow method matches the number of clusters in visualizations for both datasets separately. For t-SNE and UMAP implementation hyperparameters, you can refer to [17] and [19], respectively.
- In your report, please provide a worst-case running time analysis for both Kmeans and Kmedoids with respect to the number of data points (N), data sample vector dimension (d), cluster number (K), and the number of iterations (I).

Note: In order to better grasp the relation between EM and Kmeans, the reader is encouraged to refer to Gaussian mixture models (GMM), which can be thought be the generalization of Kmeans where clusters are represented with Gaussian distributions (with separate mean vectors and covariance matrices) and the labeling of data points to clusters are performed probabilistically [2] (Chapter 7), [11] (Chapter 9), [22]. Kmeans is a specific instance of GMM where the covariance matrix is the identity matrix (this is the reason why Kmeans is biased towards spherical clusters) and labeling is done deterministically (by measuring distance to the cluster centers) (for the closest cluster the probability score is 1 and for the other clusters, it is 0).

Note: Kmedoids algorithm has an extra hyperparameter compared to Kmeans: distance measure. Within the scope of this assignment, we have considered it as fixed (as cosine). But in general we need to consider other possible metrics and test them during experiments. To this end, we can repeat the experiment procedure described above for each possible distance measure from scratch and report the results of the distance measure that attains the smallest average loss value at the elbow location.

Part 3

Kmeans and Kmeans++ methods introduced in the previous part are intrinsically random and tend to form spherical clusters. Even if Kmedoids can be utilized to alleviate this tendency, it still strongly depends on the initial configuration of clusters, which requires multiple runs to obtain satisfactory results. Instead of determining initial cluster centers (in other words, enforcing initial cluster centers for a dataset), we may look forward to capturing natural groupings (due to data distribution) within a dataset itself without assuming anything about where clusters may be formed (in a sense, by introducing initial cluster centers, we dictate how clusters should be formed). This is the basic approach that the hierarchical clustering approach considers. Hierarchical clustering (HC) aims to form clusters in a bottom-up or top-down approach without the notion of a cluster center. In this part, we consider the bottom-up approach (Hierarchical Agglomerative Clustering (HAC)). With HAC, initially, every data instance is considered a cluster. By starting from these clusters iteratively, the two closest/the most similar clusters are merged until a single cluster is formed. How to merge two clusters (linkage criterion) becomes the main problem since a similarity/distance measure between clusters needs to be defined, which is a harder problem than instance-instance distance/similarity measurement. To this end, in this part, three linkage criteria are considered [2] (Chapter 7), [27]: single, complete, and average. Let $C_i = \{x_m^i | m = 1 \dots M\}$ and $C_j = \{x_l^j | l = 1 \dots L\}$ denote two clusters with M and L elements (instances), respectively. x_m^i denotes the m^{th} data instance of cluster i , similarly, x_l^j is l^{th} instance of cluster j . With these cluster definitions, the linkage criteria are defined as follows [2] (Chapter 7), [27]:

$$\text{Single Linkage}(C_i, C_j) = \min_{m,l} D(x_m^i, x_l^j)$$

$$\text{Complete Linkage}(C_i, C_j) = \max_{m,l} D(x_m^i, x_l^j)$$

$$\text{Average Linkage}(C_i, C_j) = \frac{1}{N * M} \sum_m^M \sum_l^L D(x_m^i, x_l^j)$$

where D is instance-instance similarity/distance function/metric.

A tree-like diagram called a dendrogram could be constructed by keeping track of distance values at which two clusters are merged. By picking a suitable distance value on this structure, any demanded number of clusters could be obtained (it could be likened to a cut operation on the dendrogram; each branch on the cut represents a cluster).

To assess the quality of the clustering results of HAC, we can utilize the silhouette method [23, 24] (the silhouette method could also be utilized for other clustering methods). A silhouette score provides an idea of how close a data point is to its cluster instances and how well it is separated from other clusters. The silhouette score for m^{th} data instance of cluster i (x_m^i) is calculated as follows:

$$\alpha(x_m^i) = \frac{1}{|C_i| - 1} \sum_{n, n \neq i} D(x_m^i, x_n^i)$$

$$\beta(x_m^i) = \min_{j \neq i} \frac{1}{|C_j|} \sum_l D(x_m^i, x_l^j)$$

$$s(x_m^i) = \begin{cases} \frac{\beta(x_m^i) - \alpha(x_m^i)}{\max(\alpha(x_m^i), \beta(x_m^i))} & |C_i| > 1 \\ 0 & |C_i| = 1 \end{cases}$$

where $|C_i|$ denotes the number of elements of cluster i . $\alpha()$ finds the average distance within a cluster with respect to a given data point. $\beta()$ gets the average distance of the given point to the closest cluster instances. The value range of the silhouette score is $[-1, 1]$. A value of 1 indicates that the sample point is close to its cluster and far away from other clusters. On the other hand, -1 indicates that the data

point may be in the wrong cluster. The average silhouette value of all data points gives an idea about the compactness of overall clustering.

$$\text{Average silhouette value} = \frac{\sum_c \sum_l s(x_l^c)}{\sum_c |C_c|}$$

Furthermore, when silhouette values of individual data instances of each cluster are sorted and displayed on a single plot [23, 24], we could elicit a very insightful idea about how well the clustering results are, and we could determine the optimal number of clusters. When silhouette values of data points of all clusters are above the average silhouette value or close to it, and clusters are of similar size (they contain a similar number of data points), such clustering can be considered desirable [24]. Readers are encouraged to refer to [26] (Chapter 2) for more detail regarding the silhouette method application.

Part 3 Specifications

- For this part, you are given a dataset (as a 2D numpy array) for HAC clustering and a partial implementation that loads it, namely **part3.py**.
- You are expected to perform clustering experiments via scikitlearn's **AgglomerativeClustering** [25] on the given dataset and report the results.
- You are expected to try different HAC hyperparameter values (namely for linkage criterion and instance-instance distance/similarity measure). You should consider the following hyperparameter values: ['single', 'complete'] and ['euclidean', 'cosine'], respectively (in total, four hyperparameter configurations).
- For each hyperparameter configuration, please add the resulting dendrogram plots in your report.
- For each hyperparameter configuration, please perform silhouette analysis for the K values of 2, 3, 4, and 5 to get the best K value. Among these four best configurations (('single', 'euclidean', K_1), ('single', 'cosine', K_2), ('complete', 'euclidean', K_3), ('complete', 'cosine', K_4)), report the one that attains the highest average silhouette score.
- In your report, please add all silhouette value plots and comment on them.
- You can make use of code segments provided in the scikitlearn documentation.
- The whole experimentation code of this part should be in the **part3.py** file, and this file should be submitted.
- Similar to Part 2, you are given a Python file named **part3_dimensionalityreduction.py**, which already loads the dataset of this part. In this file, you are expected to employ the three dimensionality reduction methods for visualization of the dataset of this part in 2-D. You can try different hyperparameter configurations for t-SNE and UMAP to enhance visualization results (please ensure that the hyperparameter configurations you have tested are explicit in the file). In the report, please provide the "best" dimensionality reduction results (as matplotlib scatter plots) of each method (the definition of "best" is up to you) (if you would like, you can add 3-D visualizations too). Also, please briefly discuss whether the number of clusters you have identified via the silhouette method matches the number of clusters in visualizations for the dataset. For t-SNE and UMAP implementation hyperparameters, you can refer to [17] and [19], respectively.
- In your report, please provide a worst-case run time analysis for HAC with respect to the number of data points (N) and dimension (D) of data instances and comment on which clustering method (Kmeans or HAC) you would prefer to use with a dataset consisting of 1 million data points each of which has a dimension of 120000 (e.g., 200x200 RGB image).

Regulations

1. You are expected to write your code in Python using scikitlearn, numpy, copy, pickle, umap-learn, and Matplotlib libraries.
2. Falsifying results or changing the composition of training, validation, and test data is strictly forbidden, and you will receive 0 if this is the case. Your programs will be examined to see if the reported results are obtained from them and if it is working correctly.
3. **Commenting:** Since all implementations are going to be inspected manually, comments are of great importance for the evaluation procedure. Please add extensive explanatory comments in all of your implementations.
4. **Late Submission:** You have a total of 3 late days for all homework without receiving a penalty. Penalization will be in effect as soon as you have depleted your quota. The late submission penalty will be calculated using $5d^2$; that is, 1 day late submission will cost you 5 points, 2 days will cost you 20 points, and 3 days will cost you 45 points. No late submission is accepted after reaching a total of 3 late days.
5. **Cheating:** Using any piece of code that is not your own is strictly forbidden and constitutes cheating. This includes friends, previous homework, or the internet. However, example code snippets shared on Scikitlearn's website can be used. **We have a zero-tolerance policy for cheating.** People involved in cheating will be punished according to university regulations.
6. **Discussion:** You must follow ODTUClass for discussions and possible updates/corrections/clarifications on a daily basis. **For the previous assignment, we received many questions that could have been asked on the discussion forum. Please ask your questions on ODTUClass unless you really think that your question is private and does not concern anyone.**
7. **Evaluation:** Your assignment is going to be graded manually.

Submission

Submissions will be done via the ODTUClass system. You are expected to upload a single pdf file named **report.pdf** for all parts. For Part 1, you are expected to upload **Distance.py**, **Knn.py**, and **Knnexperiment.py** files. For Part 2 you are expected to upload **Kmeans.py** and **Kmedoids.py**, **Kmeansexperiment.py**, **Kmedoidsexperiment.py** and **part2_dimensionalityreduction.py** files. For Part 3, you are expected to upload **part3.py** and **part3_dimensionalityreduction.py**.

References

- [1] https://scikit-learn.org/stable/modules/cross_validation.html
- [2] Introduction to Machine Learning, 2nd edition, Ethem Alpaydm
- [3] Cawley, G.C.; Talbot, N.L.C. On over-fitting in model selection and subsequent selection bias in performance evaluation. J. Mach. Learn. Res 2010,11, 2079-2107.
- [4] https://scikit-learn.org/stable/auto_examples/model_selection/plot_nested_cross_validation_iris.html
- [5] <https://machinelearningmastery.com/nested-cross-validation-for-machine-learning-with-python/>
- [6] https://en.wikipedia.org/wiki/Mahalanobis_distance

- [7] <https://numpy.org/doc/stable/reference/generated/numpy.cov.html>
- [8] <https://numpy.org/doc/stable/reference/generated/numpy.linalg.inv.html>
- [9] https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.StratifiedKFold.html
- [10] <https://numpy.org/doc/stable/reference/random/generated/numpy.random.shuffle.html>
- [11] Pattern Recognition and Machine Learning, Christopher M. Bishop.
- [12] Neural Networks and Learning Machines, 3rd edition, Simon Haykin.
- [13] Kaufman, Leonard; Rousseeuw, Peter J. (1990-03-08), "Partitioning Around Medoids (Program PAM)", Wiley Series in Probability and Statistics, Hoboken, NJ, USA: John Wiley and Sons, Inc., pp. 68–125.
- [14] Arthur, D.; Vassilvitskii, S. (2007). "k-means++: the advantages of careful seeding" (PDF). Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms. Society for Industrial and Applied Mathematics Philadelphia, PA, USA. pp. 1027–1035.
- [15] https://en.wikipedia.org/wiki/Determining_the_number_of_clusters_in_a_data_set
- [16] Van der Maaten, L. and Hinton, G., 2008. Visualizing data using t-SNE. Journal of machine learning research, 9(11).
- [17] <https://scikit-learn.org/stable/modules/generated/sklearn.manifold.TSNE.html>
- [18] McInnes, L., Healy, J. and Melville, J., 2018. Umap: Uniform manifold approximation and projection for dimension reduction. arXiv preprint arXiv:1802.03426.
- [19] <https://umap-learn.readthedocs.io/en/latest/>
- [20] Kramer, M.A., 1991. Nonlinear principal component analysis using autoassociative neural networks. AIChE journal, 37(2), pp.233-243.
- [21] <https://en.wikipedia.org/wiki/Autoencoder>
- [22] <https://scikit-learn.org/stable/modules/mixture.html>
- [23] [https://en.wikipedia.org/wiki/Silhouette_\(clustering\)](https://en.wikipedia.org/wiki/Silhouette_(clustering))
- [24] https://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html
- [25] <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.AgglomerativeClustering.html>
- [26] Kaufman, Leonard and Rousseeuw, Peter. (1990). Finding Groups in Data: An Introduction To Cluster Analysis. 10.2307/2532178.
- [27] Lecture notes
- [28] Announcements Page
- [29] Discussion Forum