

NETWORK ANOMALY DETECTION

COLLABORATORS

i	Name	ID
1	Zeyad Zidan	19015709
2	Abdelaziz Mohammed	19015941
3	Omar Khairat	19016063

Problem Statement

The exponential growth of network traffic has led to increased network anomalies, such as cyber-attacks, network failures, and hardware malfunctions. Network anomaly detection is a critical task for maintaining the security and stability of computer networks. This assignment aims to help students understand how K-Means and Normalized Cut algorithms can be used for network anomaly detection.

Understanding Dataset

After downloading the dataset and understanding the format, we split the training and testing data into training data, testing data, training labels, and testing labels.

Now the data is categorical, but to perform K-Means, Normalized Cut, and DB Scan it is required to convert the data into numerical data.

A function to convert each training feature and its corresponding testing feature from categorical to numerical form is provided below.

```

1 import numpy as np
2 from sklearn.preprocessing import LabelEncoder
3
4 def cat_to_num(trcolumn, tscolumn):
5     """
6     Converts 2 categorical columns of the same types into numerical columns
7
8     Args:
9         trcolumn (ndarray): ndarray of values of the first column.
10        tscolumn (ndarray): ndarray of values of the second column.
11
12    Returns:
13        tuple: a tuple of 2 ndarrays
14    """
15    encoder = LabelEncoder()
16    categories = set(np.unique(trcolumn)).union(set(np.unique(tscolumn)))
17    encoder.fit(list(categories))
18    return encoder.transform(trcolumn), encoder.transform(tscolumn)

```

[15] Python

Now the function shall convert the categorical features to numerical ones.

```

1 # Copy the data into another dataframe to convert its categorical values into numerical.
2 num_training = training.copy()
3 num_testing = testing.copy()
4
5 # Convert the categorical features.
6 for i in range(1, 4):
7     values = cat_to_num(num_training.iloc[:, i].values, num_testing.iloc[:, i].values)
8     num_training.isetitem(i, values[0])
9     num_testing.isetitem(i, values[1])
10
11 # Convert the labels.
12 num_trlabels, num_tslabels = cat_to_num(trlabels, tslabels)

```

[16] Python

K-Means

The code for K-Means is provided below in which :

This code implements the K-means algorithm to cluster data into k labels. The algorithm iteratively updates the centroids by assigning each data point to the nearest centroid and then computing the mean of the data points assigned to each centroid.

The function takes three arguments: X (the data matrix), k (the number of clusters), and epsilon (a tolerance fraction for convergence). It returns the final centroids after performing K-means on the data.

The algorithm initializes the centroids by randomly choosing k data points from X. Then it computes the Euclidean distance from each sample to each centroid, assigns each sample to the nearest centroid, and updates the centroids by computing the mean of the assigned samples. It repeats this process until the centroids do not change by more than epsilon.

```

1 def k_means(X, k, epsilon):
2
3     n_samples, n_features = X.shape
4
5     # Randomly choose k data points as the initial centroids
6     centroids = X[np.random.choice(n_samples, k, replace=False)]
7     distances = np.zeros((n_samples, k))
8     labels = np.zeros(n_samples)
9     old_centroids = np.zeros((k, n_features))
10
11     # Continue until the centroids don't change by more than epsilon
12     while np.linalg.norm(centroids - old_centroids) > epsilon:
13         old_centroids = centroids.copy()
14
15         # Calculate the Euclidean distances from each sample to each centroid
16         for i in range(k):
17             distances[:, i] = np.linalg.norm(X - centroids[i], axis=1)
18
19         # Assign each sample to the nearest centroid
20         labels = np.argmin(distances, axis=1)
21
22         # Update the centroids to be the mean of the samples assigned to them
23         for i in range(k):
24             X_i = X[labels == i]
25             if len(X_i) == 0:
26                 centroids[i] = old_centroids[i]
27             else:
28                 centroids[i] = np.mean(X_i, axis=0)
29
30     return centroids

```

```

from sklearn.metrics.cluster import contingency_matrix
from scipy.optimize import linear_sum_assignment

labels = []
train = np.array(num_training_10)
test = np.array(num_testing_10)
results = []
for K in [7,15,23,31,45]:
    centroids = k_means(train, K, 0.001)
    distances = np.linalg.norm(test[:, np.newaxis, :] - centroids, axis=2)
    labels = np.argmin(distances, axis=1)
    contingency = contingency_matrix(num_tslables_10, labels)
    row_ind, col_ind = linear_sum_assignment(-contingency)
    y_pred = np.zeros_like(labels)
    for i, j in zip(row_ind, col_ind):
        y_pred[labels == j] = i
    print("K:",K)
    print("precision:",precision(y_pred.tolist(), num_tslables.tolist()))
    print("recall",recall(y_pred.tolist(), num_tslables.tolist()))
    print("f1_score",f1(y_pred.tolist(), num_tslables.tolist()))
    print("conditional_entropy",conditional_entropy(y_pred.tolist(), num_tslables.tolist()))
    print()

```

Normalized Cut

In the normalized cut, we set the random seed by 42. Using the `train_test_split` function we produce a training dataset representing 0.15% of the dataset. Then we calculate the cosine similarity matrix, degree matrix, and the Laplacian matrix of the training dataset. After

calculating the Laplacian matrix, we get the normalized sorted eigenvectors and perform K-Means on them and finally return the labels of the clustering.

```
39 def ncut(training_data, k):
40     """
41     Splits the data into a training set and testing set with ratio 0.15% for training dataset,
42     then applies the normalized cut algorithm on the reduced training dataset.
43
44     Args:
45         data (pd.DataFrame): pd.DataFrame containing the original dataset.
46         k (int): number of clusters.
47
48     Returns:
49         predicted, true: nparrays of predicted labels after applying the normalized cut algorithm and the true labels.
50     """
51     training = tts(training_data, random_state=42, train_size=0.0015)[0]
52
53     # Get the true labels
54     true = training.iloc[:, 41].values
55
56     # Convert the data into numpy arrays
57     training = np.array(training)
58
59     # Construct the similarity graph
60     S = cosine_similarity(training)
61
62     # Construct the degree matrix
63     degrees = np.sum(S, axis=1)
64     D = np.diag(degrees)
65
66     # Compute Laplacian Matrix
67     L = D - S
68
69     # Compute sorted eigenvectors of the Laplacian Matrix then normalize them
70     values, vectors = np.linalg.eigh(L)
71     eigenvectors = vecsort(vectors, values)
72
73     normalized = norm(eigenvectors[:, :k])
74
75     # Perform K-means clustering on eigenvectors
76     centroids = k_means(normalized, k, 0.01)
77     distances = np.linalg.norm(normalized[:, np.newaxis, :] - centroids, axis=2)
78     predicted = np.argmin(distances, axis=1)
79
80     return predicted, true
81
```

Evaluation of Normalized Cut Results

- Precision = 0.9248672927725603
- Recall = 0.7068238647290678
- F1 Measure = 0.30294307548555777
- Conditional Entropy = 0.3441486877728897

```
1 num_training['labels'] = num_train_labels
2 ncut_clustering = ncut(num_training, 11)
3 ncut_predicted = np.array(ncut_clustering[0])
4 ncut_true = np.array(ncut_clustering[1])

[38] ✓ 1m 16.5s

1 print("Precision:", precision(ncut_predicted.tolist(), ncut_true.tolist()))
2 print("Recall", recall(ncut_predicted.tolist(), ncut_true.tolist()))
3 print("F1 Measure", f1(ncut_predicted.tolist(), ncut_true.tolist()))
4 print("Conditional Entropy", conditional_entropy(ncut_predicted.tolist(), ncut_true.tolist()))

[39] ✓ 0.0s

... Precision: 0.9248672927725603
Recall 0.7068238647290678
cluster: 0 pre: 0.8837735849056094 rec: 0.9545454545454546 f1: 0.8726953467954346
cluster: 1 pre: 0.9421221864951769 rec: 0.13739742086752638 f1: 0.23981993042766525
cluster: 2 pre: 1.0 rec: 0.11969389654149571 f1: 0.21379738968383294
cluster: 3 pre: 0.9552238805970149 rec: 0.08907446068197634 f1: 0.16295353278166774
cluster: 4 pre: 1.0 rec: 0.015309672929714684 f1: 0.03015764222069911
cluster: 5 pre: 0.6820083682008368 rec: 0.11343075852470424 f1: 0.19451073985680192
cluster: 6 pre: 0.9532163742690059 rec: 0.11343075852470424 f1: 0.20273631840796016
cluster: 7 pre: 0.9942528735632183 rec: 0.12038970076548365 f1: 0.21477343265052762
cluster: 8 pre: 0.9491525423728814 rec: 0.03897007654836465 f1: 0.07486631016942782
cluster: 9 pre: 0.9880952380952381 rec: 0.11551844119693806 f1: 0.2068535825451713
cluster: 10 pre: 0.9849986191369606 rec: 0.861664712778429 f1: 0.9192896848024013
F1 Measure 0.30294307548555777
Conditional Entropy 0.3441486877728897
```

DB Scan

The code defines a function called DB Scan that implements the density-based spatial clustering of applications with noise (DBSCAN) algorithm. The function takes three inputs: data, which is a pandas DataFrame containing the data to be clustered, eps, which is the maximum distance between two points for them to be considered as part of the same cluster, and min_samples, which is the minimum number of points required for a cluster to be formed. The function initializes some variables and then loops over each point in the data. For each point, it checks if it has already been assigned to a cluster. If not, it finds all the neighboring points within a distance of eps and determines whether the point is a "core point" (i.e., it has at least min_samples neighbors) or not. If the point is a core point, a new cluster is formed and all its connection points are assigned to the same cluster. If the point is not a core point, it is marked as an outlier. The function also calls another function called expand_cluster to expand the cluster starting from the current core point. The expand_cluster function takes the current point, its neighbors, and some other parameters and recursively adds new points to the cluster until no new points can be added. Finally, the function returns a list of cluster labels for each point in the data, where the label is an integer that indicates the cluster to which the point belongs. If a point is an outlier, its label is -1. We consider the outliers as one of the clusters.

```

vis = []
my_dict = {}
def dbscan(data, eps, min_samples):
    """
    Performs DB Scan clustering algorithm on a given dataset.

    Args:
        data (ndarray): the dataset to be clustered
        eps (float): a tolerance extent.
        min_samples (number): a number indicating the minimum number of samples.

    Returns:
        labels: the labels of the dataset after clustering.
    """
    X = data.values
    global vis
    global my_dict
    vis = [0] * len(X)
    my_dict = {i: [] for i in range(len(X))}
    labels = [0] * len(X)
    cluster_id = 0
    for i in range(len(X)):
        if labels[i] != 0:
            continue
        # Find all neighbors of the current point within eps distance
        neighbors = get_neighbors(X, i, eps)
        # If the point is not a core point, mark it as an outlier
        if len(neighbors) < min_samples:
            labels[i] = -1
            continue
        # Expand the cluster starting from the current core point
        cluster_id += 1
        labels[i] = cluster_id

        expand_cluster(X, labels, i, neighbors, eps, min_samples, cluster_id)

    return labels

```

```

def expand_cluster(X, labels, i, neighbors, eps, min_samples, cluster_id):
    """
    Auxiliary function for the DB Scan to expand the core points clusters.

    Args:
        X (ndarray): the dataset
        labels (ndarray): labels of the clustering
        i (number): cluster's id
        neighbors (ndarray): neighbors of the core point
        eps (float): a tolerance extent
        min_samples (number): the number of minimum samples
        cluster_id (number): the cluster's id
    """
    # Loop over each neighbor of the core point
    for j in neighbors:
        if labels[j] == -1:
            labels[j] = cluster_id
        elif labels[j] == 0:
            labels[j] = cluster_id
            # Find all neighbors of the current point within eps distance
            new_neighbors = get_neighbors(X, j, eps)
            # If the point is a core point, add its neighbors to the list of neighbors
            if len(new_neighbors) >= min_samples:
                neighbors += new_neighbors

```

```
def get_neighbors(X, i, eps):
    """
    This functions gets the neighbour of the ith instance within given epsilon

    Args:
        X (ndarray): the dataset
        i (number): cluster's id
        eps (float): a tolerance extent

    Returns:
        neighbors: ndarray of neighbors of the ith instance
    """
    global vis
    global my_dict
    if vis[i] == 1:
        return my_dict[i]
    neighbors = []
    for j in range(len(X)):
        if i == j:
            continue

        dist = np.linalg.norm(X[i] - X[j])
        if dist <= eps:
            neighbors.append(j)
    vis[i] = 1
    my_dict[i] = neighbors
    return neighbors
```

Evaluation of DBSCAN Results

- Precision = 0.9897917517354023
- Recall = 0.7563745809565455
- F1 Measure = 0.23751813503832142
- Conditional Entropy = 0.08650543363578882

Evaluation

Clusterize

Auxiliary function to annotate the predicted labels with the true labels.

```
1 def clusterize(pred_labels, true_labels):
2     """
3     Annotate the predicted with the true labels
4
5     Args:
6         pred_labels (list): list of predicted labels
7         true_labels (list): list of true labels
8
9     Raises:
10        ValueError: The two lists must be equal.
11
12     Returns:
13        clusters, clusters_set: clusters is a dictionary of the annotation, and clusters_set is a set of the predicted labels.
14     """
15     if len(pred_labels) != len(true_labels):
16         raise ValueError("The two list should be equal")
17     clusters_set = set(pred_labels)
18     # num_clusters = len(set(pred_labels))
19     clusters = {}
20     for cluster in clusters_set:
21         clusters[cluster] = []
22     for i in range(len(pred_labels)):
23         clusters[pred_labels[i]].append(true_labels[i])
24     return clusters, clusters_set
```

Python

Precision

```
1 def precision(pred_labels, true_labels):
2     """
3     Evaluates the precision of a clustering given the true labels.
4
5     Args:
6         pred_labels (list): list of predicted labels
7         true_labels (list): list of true labels
8
9     Returns:
10        res: the precision value of the clustering
11     """
12     clusters, clusters_set = clusterize(pred_labels, true_labels)
13     res = 0
14     for cluster in clusters_set:
15         most_common = max(set(clusters[cluster]), key = clusters[cluster].count)
16         count = clusters[cluster].count(most_common)
17         res += (len(clusters[cluster]) / len(true_labels)) * (count / len(clusters[cluster]))
18     return res
```

Recall

```
1 def recall(pred_labels, true_labels):
2     """
3     Evaluates the recall of a clustering given the true labels.
4
5     Args:
6         pred_labels (list): list of predicted labels
7         true_labels (list): list of true labels
8
9     Returns:
10        res: the recall value of the clustering
11    """
12    clusters, clusters_set = clusterize(pred_labels, true_labels)
13    res = 0
14    r = len(clusters_set)
15    for cluster in clusters_set:
16        most_common = max(set(clusters[cluster]), key = clusters[cluster].count)
17        count = clusters[cluster].count(most_common)
18        count_total = true_labels.count(most_common)
19        res += (len(clusters[cluster]) / len(true_labels)) * (count / count_total)
20    return res
```

F1 Measure

```
1 def f1(pred_labels, true_labels):
2     """
3     Evaluates the F1-Measure of a clustering given the true labels.
4
5     Args:
6         pred_labels (list): list of predicted labels
7         true_labels (list): list of true labels
8
9     Returns:
10        res: the F1-Measure value of the clustering
11    """
12    clusters, clusters_set = clusterize(pred_labels, true_labels)
13    res = 0
14    r = len(clusters_set)
15    for cluster in clusters_set:
16        most_common = max(set(clusters[cluster]), key = clusters[cluster].count)
17        count = clusters[cluster].count(most_common)
18        count_total = true_labels.count(most_common)
19        precision = count / len(clusters[cluster])
20        recall = count / count_total
21        f1 = (2 * precision * recall) / (precision + recall)
22        print(f"cluster: {cluster} pre: {precision} rec: {recall} f1: {f1}")
23        res += (float(f1) / float(r))
24    return res
```

Conditional Entropy

```
1 from math import log2
2 def conditional_entropy(pred_labels, true_labels):
3     """
4     Evaluates the conditional entropy of a clustering given the true labels.
5
6     Args:
7         pred_labels (list): list of predicted labels
8         true_labels (list): list of true labels
9
10    Returns:
11        res: the conditional entropy value of the clustering
12    """
13    clusters, clusters_set = clusterize(pred_labels, true_labels)
14    res = 0
15    true_labels_set = set(true_labels)
16    for cluster in clusters_set:
17        temp = 0
18        for t in true_labels_set:
19            t_count = clusters[cluster].count(t)
20            if t_count != 0:
21                temp += -(t_count / len(clusters[cluster])) * log2(t_count / len(clusters[cluster]))
22        res += (len(clusters[cluster]) / len(true_labels)) * temp
23    return res
```