

CS 354 - Assignment 7

Arnav Jain - 220002018

Question 1

Code:

```
q1.cpp x
Question 1 > q1.cpp > SOM > SOM(int, float, float)
1  #include <iostream>
2  #include <vector>
3  #include <random>
4  #include <cmath>
5  #include <fstream>
6
7  void save_dataset(const std::vector<float>& dataset, const std::string& filename) {
8      std::ofstream out(filename);
9      if(out) {
10         for(const auto& val : dataset) {
11             out << val << "\n";
12         }
13     }
14 }
15
16 class SOM {
17 private:
18     std::vector<float> weights;
19     float learning_rate;
20     float radius;
21     int num_neurons;
22
23     float neighborhood(float dist, float sigma) {
24         return exp(-(dist*dist)/(2*sigma*sigma));
25     }
26
27 public:
28     SOM(int n_neurons, float lr, float init_radius) :
29         num_neurons(n_neurons),
30         learning_rate(lr),
31         radius(init_radius)
32     {}
33
34     std::random_device rd;
35     std::mt19937 gen(rd());
36     std::uniform_real_distribution<float> dist(0.0f, 1.0f);
37
38     weights.resize(n_neurons);
39     for(auto& w : weights) {
40         w = dist(gen);
41     }
42 }
```

```

43     int find_bmu(float input) {
44         int bmu_idx = 0;
45         float min_dist = INFINITY;
46
47         for(int i=0; i<num_neurons; ++i) {
48             float dist = fabs(input - weights[i]);
49             if(dist < min_dist) {
50                 min_dist = dist;
51                 bmu_idx = i;
52             }
53         }
54         return bmu_idx;
55     }
56
57     void train(float input, int iteration, int max_iter) {
58         int bmu = find_bmu(input);
59
60         float current_lr = learning_rate * (1.0 - (float)iteration/max_iter);
61         float current_radius = radius * (1.0 - (float)iteration/max_iter);
62
63         for(int i=0; i<num_neurons; ++i) {
64             float distance_to_bmu = fabs(i - bmu);
65             float influence = neighborhood(distance_to_bmu, current_radius);
66             weights[i] += current_lr * influence * (input - weights[i]);
67         }
68     }
69
70     void save_weights(std::ofstream& out, int iteration) const {
71         out << "iteration " << iteration << ": ";
72         for(size_t i = 0; i < weights.size(); ++i) {
73             out << weights[i] << " ";
74         }
75         out << "\n";
76     }
77
78     const std::vector<float>& get_weights() const { return weights; }
79 };

```

```

81 int main() {
82     const int NUM_NEURONS = 15;
83     const int EPOCHS = pow(10 , 6);
84     const float INIT_LR = 0.01;
85     const float INIT_RADIUS = NUM_NEURONS/2.0f;
86
87     SOM som(NUM_NEURONS, INIT_LR, INIT_RADIUS);
88
89     std::vector<float> dataset(100);
90     std::random_device rd;
91     std::mt19937 gen(rd());
92     std::uniform_real_distribution<float> dist(0.0f, 10.0f);
93     for(auto& val : dataset) {
94         val = dist(gen);
95     }
96
97     save_dataset(dataset, "dataset.txt");
98
99     std::ofstream weight_file("weights_history.txt");
100
101     for(int epoch=0; epoch<EPOCHS; ++epoch) {
102         int random_idx = std::rand() % dataset.size();
103         float input = dataset[random_idx];
104
105         som.train(input, epoch, EPOCHS);
106
107         // Save weights for every iteration
108         if(weight_file) {
109             som.save_weights(weight_file, epoch);
110         }
111     }
112
113     std::cout << "Training complete. Check weights_history.txt and dataset.txt\n";
114     return 0;
115 }

```

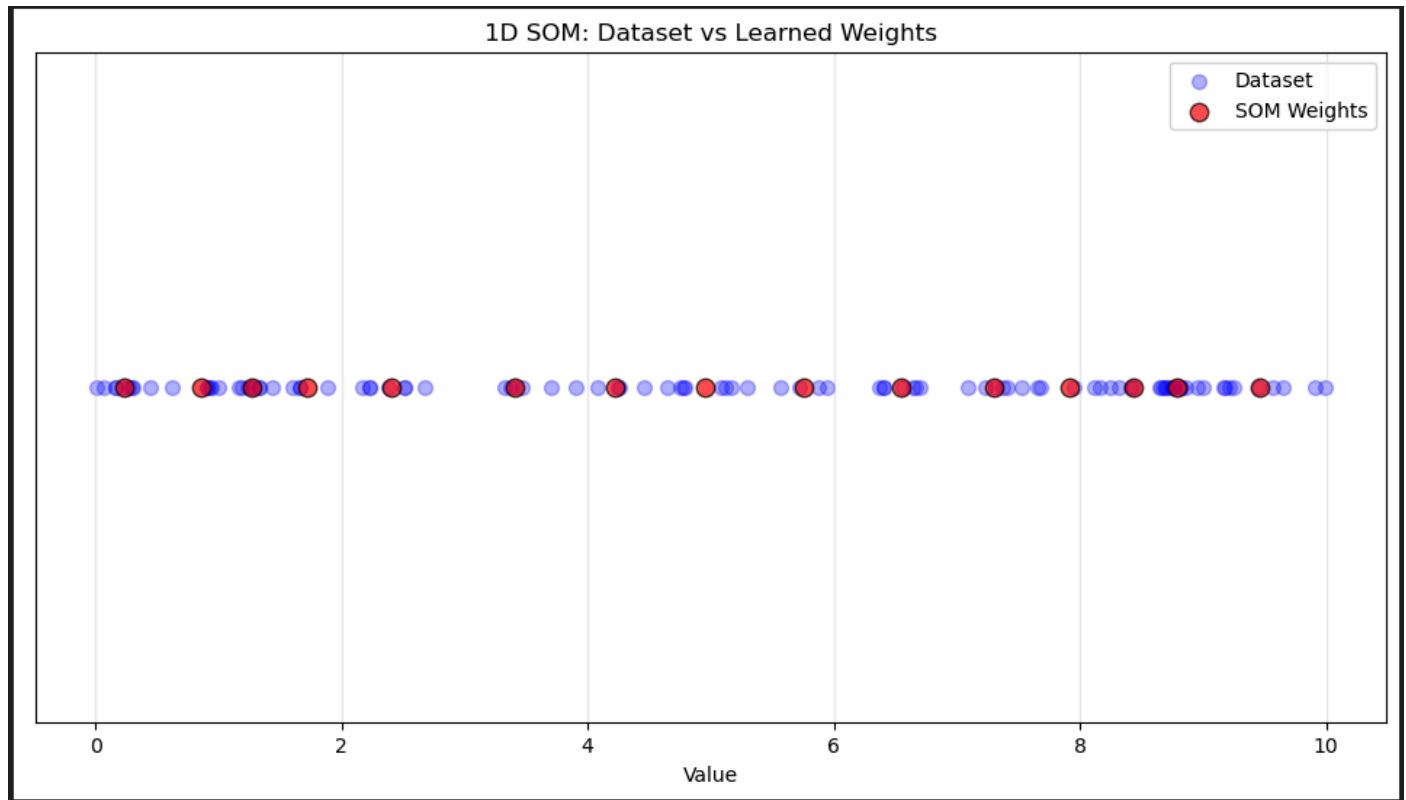
Dataset Generated:

dataset.txt	
Question 1	> dataset.txt
1	3.82685
2	1.88795
3	3.3214
4	1.24071
5	7.68586
6	8.85502
7	6.64583
8	7.53061
9	3.41195
10	4.75947
11	0.0119643
12	0.221545
13	1.26753
14	5.72585
15	0.169204
16	5.30185
17	9.21849
18	3.46402
19	2.51592
20	5.0806
21	1.16353
22	8.64934
23	8.15961
24	8.3182
25	0.629595
26	0.936581
27	4.45675
28	8.70864
29	3.36936
30	0.2659
31	2.22883
32	0.999832

Weights of SOM:

weights_history.txt	
Question 1	> weights_history.txt
999985	iteration 999984: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999986	iteration 999985: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999987	iteration 999986: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999988	iteration 999987: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999989	iteration 999988: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999990	iteration 999989: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999991	iteration 999990: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999992	iteration 999991: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999993	iteration 999992: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999994	iteration 999993: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999995	iteration 999994: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999996	iteration 999995: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999997	iteration 999996: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999998	iteration 999997: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
999999	iteration 999998: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608
1000000	iteration 999999: 0.236256 0.862613 1.27465 1.72522 2.40995 3.40857 4.218 4.95643 5.755 6.54902 7.30039 7.9191 8.4324 8.79113 9.45608

Visualisation of SOM:



Question 2

Code:

```
q2.cpp x
Question 2 > q2.cpp > main()
1  #include <iostream>
2  #include <fstream>
3  #include <vector>
4  #include <string>
5  #include <sstream>
6  #include <cmath>
7  #include <random>
8  #include <algorithm>
9  #include <numeric>
10
11  using Matrix = std::vector<std::vector<double>>>;
12  using namespace std;
13
14  Matrix load_csv(const std::string& filename) {
15      Matrix data;
16      std::ifstream file(filename);
17      std::string line;
18
19      // Skip header
20      if(file.good()) std::getline(file, line);
21
22      while(std::getline(file, line)) {
23          std::vector<double> row;
24          std::stringstream ss(line);
25          std::string value;
26
27          while(std::getline(ss, value, ',')) {
28              try {
29                  row.push_back(std::stod(value));
30              } catch(const std::exception& e) {
31                  std::cerr << "Error parsing value: " << value << "\n";
32              }
33          }
34          if(!row.empty()) data.push_back(row);
35      }
36      return data;
37  }
```

```

39 void zscore_normalize(Matrix& data) {
40     const size_t num_features = data[0].size();
41     const size_t num_samples = data.size();
42
43     for(size_t col=0; col<num_features; col++) {
44         // Calculate mean
45         double sum = std::accumulate(data.begin(), data.end(), 0.0,
46             [col](double s, const std::vector<double>& row) { return s + row[col]; });
47         double mean = sum / num_samples;
48
49         // Calculate standard deviation
50         double sq_diff = std::accumulate(data.begin(), data.end(), 0.0,
51             [col, mean](double s, const std::vector<double>& row) {
52                 return s + pow(row[col] - mean, 2);
53             });
54         double std_dev = sqrt(sq_diff / num_samples);
55
56         // Normalize
57         for(auto& row : data) {
58             row[col] = (row[col] - mean) / std_dev;
59         }
60     }
61 }
62
63 class KohonenSOM {
64     Matrix weights;
65     int input_size;
66     int output_size;
67     double initial_radius;
68
69 public:
70     KohonenSOM(int n, int m) : input_size(n), output_size(m),
71         initial_radius(std::max(m, n)/2.0) {
72         std::random_device rd;
73         std::mt19937 gen(rd());
74         std::uniform_real_distribution<double> dist(-1.0, 1.0);
75
76         weights.resize(input_size, std::vector<double>(output_size));
77         for(auto& row : weights) {
78             std::generate(row.begin(), row.end(), [&]() { return dist(gen); });
79         }
80     }

```

```

82     int find_bmu(const std::vector<double>& input) {
83         int bmu = 0;
84         double min_dist = std::numeric_limits<double>::max();
85
86         for(int j=0; j<output_size; j++) {
87             double dist = 0.0;
88             for(int i=0; i<input_size; i++) {
89                 dist += pow(input[i] - weights[i][j], 2);
90             }
91             if(dist < min_dist) {
92                 min_dist = dist;
93                 bmu = j;
94             }
95         }
96         return bmu;
97     }
98
99     void train(const Matrix& data, int epochs) {
100         const double initial_learning_rate = 0.5;
101         const double time_constant = epochs / log(initial_radius);
102
103         for(int epoch=0; epoch<epochs; epoch++) {
104             double learning_rate = initial_learning_rate * exp(-epoch/static_cast<double>(epochs));
105             double radius = initial_radius * exp(-epoch/time_constant);
106
107             for(const auto& sample : data) {
108                 int bmu = find_bmu(sample);
109
110                 // Update neighborhood
111                 for(int j=0; j<output_size; j++) {
112                     double distance_to_bmu = abs(j - bmu);
113                     if(distance_to_bmu <= radius) {
114                         double influence = exp(-pow(distance_to_bmu, 2)/(2*pow(radius, 2)));
115
116                         for(int i=0; i<input_size; i++) {
117                             weights[i][j] += learning_rate * influence *
118                                 (sample[i] - weights[i][j]);
119                         }
120                     }
121                 }
122             }
123         }
124     }

```



```

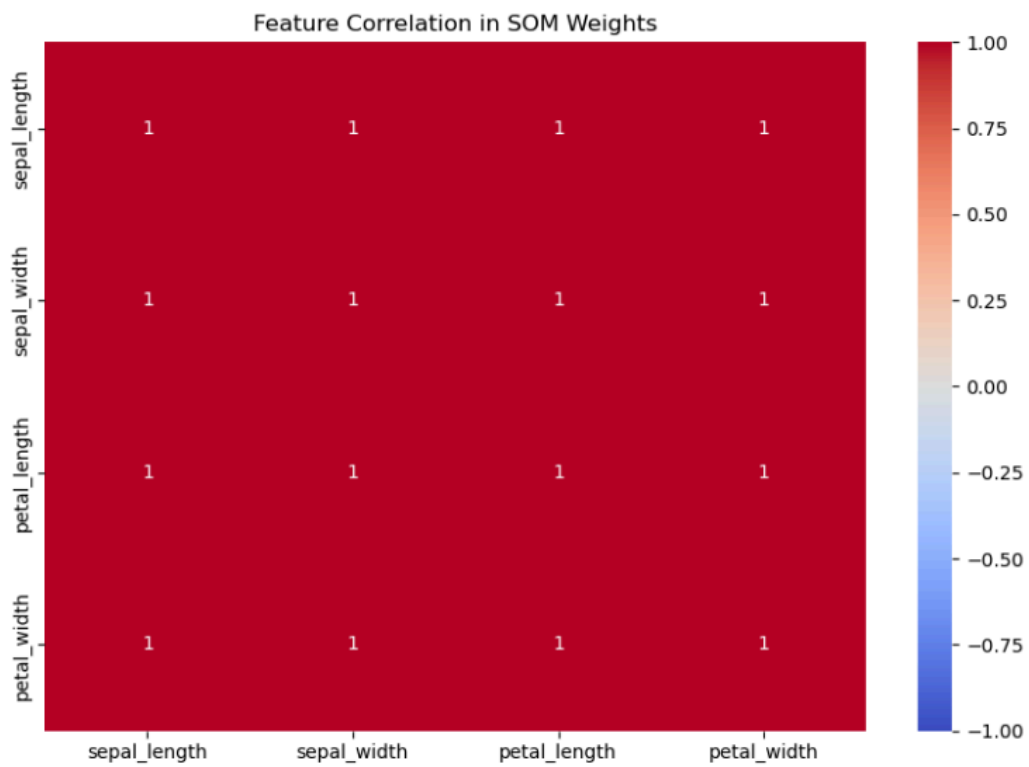
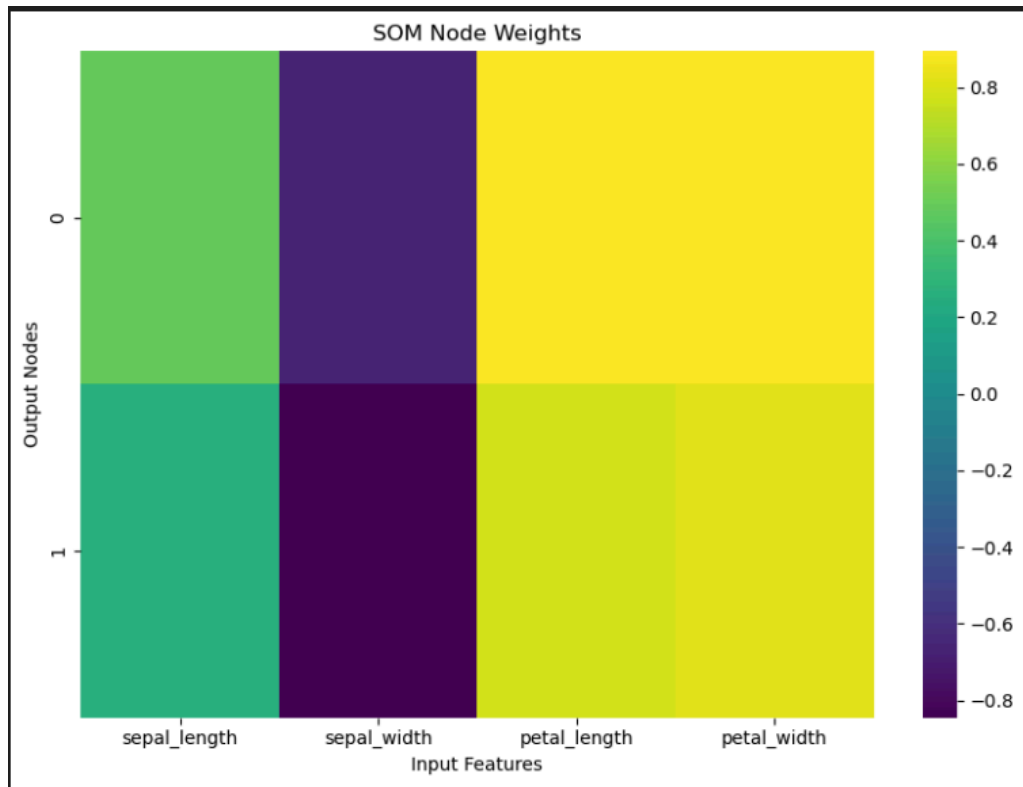
126     void save_weights(const std::string& filename) {
127         std::ofstream outfile(filename);
128         for(const auto& row : weights) {
129             for(size_t i=0; i<row.size(); i++) {
130                 outfile << row[i] << (i < row.size()-1 ? "," : "");
131             }
132             outfile << "\n";
133         }
134     }
135 };
136
137 int main() {
138     Matrix dataset = load_csv("iris_train.csv");
139     zscore_normalize(dataset);
140
141     const int num_features = dataset[0].size();
142     cout<<"Sample Size: "<<dataset.size()<<endl;
143     cout<<"Feature Size: "<<num_features<<endl;
144     const int output_nodes = 2;
145
146     KohonenSOM som(num_features, output_nodes);
147     som.train(dataset, 50);
148     som.save_weights("som_weights.csv");
149
150     return 0;
151 }

```

SOM Weights:

som_weights.csv	
Question 2 >	som_weights.csv > data
1	0.49111,0.255439
2	-0.665502,-0.844412
3	0.882156,0.772489
4	0.893402,0.819389
5	

Heatmaps:



Question 3

```
from sklearn.datasets import load_iris
import pandas as pd

data = load_iris()
iris_df = pd.DataFrame(data.data, columns=data.feature_names)
iris_df['target'] = data.target
```

```
iris_df.head()
```

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0

```
!pip install minisom
```

```
Collecting minisom
  Downloading minisom-2.3.5.tar.gz (12 kB)
  Preparing metadata (setup.py) ... done
Building wheels for collected packages: minisom
  Building wheel for minisom (setup.py) ... done
  Created wheel for minisom: filename=MiniSom-2.3.5-py3-none-any.whl si
  Stored in directory: /root/.cache/pip/wheels/19/db/95/5e53bc2b88a3282
Successfully built minisom
Installing collected packages: minisom
Successfully installed minisom-2.3.5
```

```
iris = iris_df.iloc[:, :-1]
iris.head()
```

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2

```
from sklearn.preprocessing import MinMaxScaler
from minisom import MiniSom

# Normalization
scaler = MinMaxScaler()
normalized_data = scaler.fit_transform(iris_df.iloc[:, :-1])

# SOM initialization
som = MiniSom(x=10, y=10, input_len=4, sigma=1.0, learning_rate=0.5)
som.random_weights_init(normalized_data)
som.train_random(normalized_data, 100)
```

```
import numpy as np

def calculate_bmu_distances(som, data):
    return [np.linalg.norm(sample - som.get_weights()[som.winner(sample)])
            for sample in data]

distances = calculate_bmu_distances(som, normalized_data)
iris_df['bmu_distance'] = distances
```

```
from scipy.stats import zscore

iris_df['z_score'] = zscore(distances)
anomalies = iris_df[iris_df['z_score'].abs() > 2]
```

```
print(len(anomalies))
anomalies
```

7

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target	bmu_distance	z_score
41	4.5	2.3	1.3	0.3	0	0.262963	2.631764
109	7.2	3.6	6.1	2.5	2	0.240869	2.302556
114	5.8	2.8	5.1	2.4	2	0.228376	2.116404
117	7.7	3.8	6.7	2.2	2	0.417237	4.930500
118	7.7	2.6	6.9	2.3	2	0.329537	3.623743
122	7.7	2.8	6.7	2.0	2	0.254917	2.511880
131	7.9	3.8	6.4	2.0	2	0.424832	5.043677

```

import pandas as pd
import numpy as np

COLUMNS = [
    'duration', 'protocol_type', 'service', 'flag', 'src_bytes', 'dst_bytes',
    'land', 'wrong_fragment', 'urgent', 'hot', 'num_failed_logins',
    'logged_in', 'num_compromised', 'root_shell', 'su_attempted', 'num_root',
    'num_file_creations', 'num_shells', 'num_access_files', 'num_outbound_cmds',
    'is_host_login', 'is_guest_login', 'count', 'srv_count', 'error_rate',
    'srv_error_rate', 'rerror_rate', 'srv_rerror_rate', 'same_srv_rate',
    'diff_srv_rate', 'srv_diff_host_rate', 'dst_host_count', 'dst_host_srv_count',
    'dst_host_same_srv_rate', 'dst_host_diff_srv_rate', 'dst_host_same_src_port_rate',
    'dst_host_srv_diff_host_rate', 'dst_host_error_rate', 'dst_host_srv_error_rate',
    'dst_host_rerror_rate', 'dst_host_srv_rerror_rate',
]

def load_and_preprocess_kdd(file_path):
    df = pd.read_csv(
        file_path,
        header=None,
        names=COLUMNS,
        dtype={
            'protocol_type': 'category',
            'service': 'category',
            'flag': 'category',
        }
    )

    print("Data loaded successfully. First 5 rows:")
    print(df.head())

    categorical_cols = df.select_dtypes(include=['category']).columns

    print()
    print("Categorical column analysis:")
    for col in categorical_cols:
        unique_vals = df[col].unique()
        print(f"{col}: {len(unique_vals)} unique values")
        print(f"Sample values: {list(unique_vals[:3])}")
        print()

```

```

# Vectorization strategy
for col in categorical_cols:

    # One-hot encode
    if len(df[col].unique()) < 10:
        df = pd.concat([
            df.drop(col, axis=1),
            pd.get_dummies(df[col], prefix=col)
        ], axis=1)
    else:
        # Label encode
        df[col] = df[col].astype('category').cat.codes

num_cols = df.select_dtypes(include=np.number).columns
df[num_cols] = df[num_cols].astype('float32')

return df

```

```
features = load_and_preprocess_kdd('/content/kddcup.testdata.unlabeled_10_percent')
```

Data loaded successfully. First 5 rows:

	duration	protocol_type	service	flag	src_bytes	dst_bytes	land	\
0	0	udp	private	SF	105	146	0	
1	0	udp	private	SF	105	146	0	
2	0	udp	private	SF	105	146	0	
3	0	udp	private	SF	105	146	0	
4	0	udp	private	SF	105	146	0	

	wrong_fragment	urgent	hot	...	dst_host_count	dst_host_srv_count	\
0	0	0	0	...	255	254	
1	0	0	0	...	255	254	
2	0	0	0	...	255	254	
3	0	0	0	...	255	254	
4	0	0	0	...	255	254	

	dst_host_same_srv_rate	dst_host_diff_srv_rate	\
0	1.0	0.01	
1	1.0	0.01	
2	1.0	0.01	
3	1.0	0.01	
4	1.0	0.01	

```

dst_host_same_src_port_rate  dst_host_srv_diff_host_rate  \
0                               0.00                               0.0
1                               0.00                               0.0
...

```

```

flag: 11 unique values
Sample values: ['SF', 'RSTR', 'S1']

```

Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output [settings](#)...

```

if features is not None:
    print()
    print("Processed features shape:", features.shape)
    print("Sample preprocessed data:")
    print(features.~iloc[:3, :5])

```

```

Processed features shape: (311029, 43)
Sample preprocessed data:
  duration  service  flag  src_bytes  dst_bytes
0      0.0     12.0   5.0     105.0     146.0
1      0.0     12.0   5.0     105.0     146.0
2      0.0     12.0   5.0     105.0     146.0

```

```

from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()
som_input = scaler.fit_transform(features)

```

```

# SOM initialization
som = MiniSom(x=50, y=50, input_len=43, sigma=1.0, learning_rate=0.5)
som.random_weights_init(som_input)
som.train_random(som_input, 1000)

```



```

import numpy as np

def calculate_bmu_distances(som, data):
    return [np.linalg.norm(sample - som.get_weights()[som.winner(sample)])
            for sample in data]

distances = calculate_bmu_distances(som, som_input)
features['bmu_distance'] = distances

from scipy.stats import zscore

features['z_score'] = zscore(distances)
anomalies = features[features['z_score'].abs() > 2]

print(len(anomalies))
print(anomalies)

```

```

... 9817
      duration  service  flag  src_bytes  dst_bytes  land  wrong_fragment  \
6      0.0      4.0    5.0    29.0      0.0    0.0      0.0
38     20.0      8.0    5.0   232.0     765.0  0.0      0.0
44      0.0      9.0    5.0   615.0      0.0    0.0      0.0
71      1.0     13.0    5.0  1018.0     333.0  0.0      0.0
79      0.0      9.0    5.0   884.0      0.0    0.0      0.0
...      ...      ...    ...      ...      ...    ...      ...
310138  0.0      4.0    5.0    46.0     134.0  0.0      0.0
310139  0.0      4.0    5.0    44.0      44.0  0.0      0.0
310140  0.0      4.0    5.0    44.0      44.0  0.0      0.0
310186  0.0      4.0    5.0    45.0     127.0  0.0      0.0
310907  0.0     12.0    5.0   105.0     147.0  0.0      0.0

      urgent  hot  num_failed_logins  ...  dst_host_srv_diff_host_rate  \
6      0.0  0.0      0.0  ...      0.00
38      0.0  4.0      0.0  ...      0.00
44      0.0  0.0      0.0  ...      0.00
71      0.0  0.0      0.0  ...      0.02
79      0.0  0.0      0.0  ...      0.05
...      ...  ...      ...  ...      ...
310138  0.0  0.0      0.0  ...      0.01
310139  0.0  0.0      0.0  ...      0.00
310140  0.0  0.0      0.0  ...      0.00
310186  0.0  0.0      0.0  ...      0.00
...
310186      True    0.287537  2.102637
310907      True    0.459558  3.531831

```

[9817 rows x 45 columns]

Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output [settings](#)...

Github for code:

<https://github.com/arnavjain2710/Computational-Intelligence-Lab-CS354N/tree/main/LAB%207>