

# Deep Autoencoding Gaussian Mixture Model for Unsupervised Anomaly Detection

```
import numpy as np
import pandas as pd
import torch
from data_loader import *
from main import *
from tqdm import tqdm
```

## KDD Cup 1999 Data (10% subset)

This is the data set used for The Third International Knowledge Discovery and Data Mining Tools Competition, which was held in conjunction with KDD-99 The Fifth International Conference on Knowledge Discovery and Data Mining. The competition task was to build a network intrusion detector, a predictive model capable of distinguishing between "bad" connections, called intrusions or attacks, and "good" normal connections. This database contains a standard set of data to be audited, which includes a wide variety of intrusions simulated in a military network environment.

```
data = pd.read_csv("kddcup.data_10_percent_corrected", header=None, names=['duration', 'protocol_type', 'service', 'flag',
```

## Pre-processing

Following the paper, since the "normal" only comprises of approximately 20% of the entries, the "normal" data were considered as anomalies instead.

```
data.loc[data["type"] != "normal.", 'type'] = 0
data.loc[data["type"] == "normal.", 'type'] = 1
```

Next, the categorical variables are converted to a one hot encoding representation. My implementation is a bit different from the original paper in this aspect. Since I am only using the 10% subset to generate the columns, I get 118 features instead of 120 as reported in the paper.

```
one_hot_protocol = pd.get_dummies(data["protocol_type"])
one_hot_service = pd.get_dummies(data["service"])
one_hot_flag = pd.get_dummies(data["flag"])

data = data.drop("protocol_type",axis=1)
data = data.drop("service",axis=1)
data = data.drop("flag",axis=1)

data = pd.concat([one_hot_protocol, one_hot_service, one_hot_flag, data],axis=1)
data.head()
```

[illegible]

2	0	1	0	0	0	0	0	0	0	0	...
3	0	1	0	0	0	0	0	0	0	0	...
4	0	1	0	0	0	0	0	0	0	0	...

5 row s × 119 columns

```
proportions = data["type"].value_counts()
print(proportions)
print("Anomaly Percentage",proportions[1] / proportions.sum())
```

```
0    396743
1     97278
Name: type, dtype: int64
Anomaly Percentage 0.19691065764410826
```

```
#proportions_alfa = data["type"].value_counts(normalize=True)
#print(proportions_alfa)
```

Normalize all the numeric variables.

```
cols_to_norm = ["duration", "src_bytes", "dst_bytes", "wrong_fragment", "urgent",
                "hot", "num_failed_logins", "num_compromised", "num_root",
                "num_file_creations", "num_shells", "num_access_files", "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_serror_rate", "dst_host_srv_serror_rate", "dst_host_rerror_rate", "dst_host_srv_rerror_rate" ]

#data.loc[:, cols_to_norm] = (data[cols_to_norm] - data[cols_to_norm].mean()) / data[cols_to_norm].std()
min_cols = data.loc[data["type"]==0 , cols_to_norm].min()
max_cols = data.loc[data["type"]==0 , cols_to_norm].max()

data.loc[:, cols_to_norm] = (data[cols_to_norm] - min_cols) / (max_cols - min_cols)
```

I saved the preprocessed data into a numpy file format and load it using the pytorch data loader.

```
np.savez_compressed("kdd_cup",kdd=data.as_matrix())
```

```
C:\Users\cncluser\Anaconda3\lib\site-packages\ipykernel_launcher.py:1: FutureWarning: Method .as_matrix will be removed in
    """Entry point for launching an IPython kernel.
```

I initially implemented this to be ran in the command line and use argparse to get the hyperparameters. To make it runnable in a jupyter notebook, I had to create a dummy class for the hyperparameters.

```
class hyperparams():
    def __init__(self, config):
        self.__dict__.update(**config)
    defaults = {
```

```

'lr' : 1e-4,
'num_epochs' : 200,
'batch_size' : 1024,
'gmm_k' : 4,
'lambda_energy' : 0.1,
'lambda_cov_diag' : 0.005,
'pretrained_model' : None,
'mode' : 'train',
'use_tensorboard' : False,
'data_path' : 'kdd_cup.npz',

'log_path' : './dagmm/logs',
'model_save_path' : './dagmm/models',
'sample_path' : './dagmm/samples',
'test_sample_path' : './dagmm/test_samples',
'result_path' : './dagmm/results',

'log_step' : 194//4,
'sample_step' : 194,
'model_save_step' : 194,
}

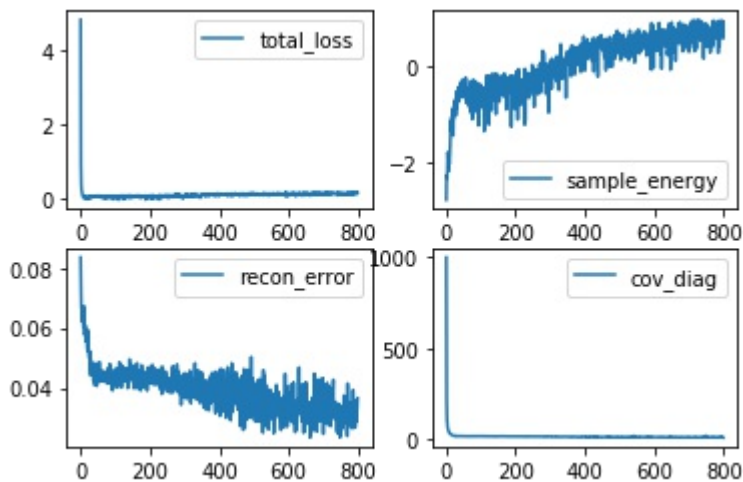
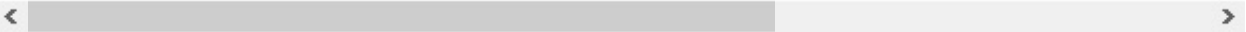
```

```

solver = main(hyperparams(defaults))
accuracy, precision, recall, f_score = solver.test()

```

Elapsed 1:30:00.990700/0:00:00.417624 -- 0:00:00.417624 , Epoch [200/200], Iter [192/194], lr 0.0001, total\_loss: 0.1735,



```

phi tensor([0.0022, 0.0053, 0.0024, 0.9900]) mu tensor([[ -0.2008,  1.0478, -0.0134],
[ -0.2553,  1.5953,  0.0554],
[ -0.1808,  1.0545, -0.0075],
[  2.0896,  0.3388,  0.4585]]) cov tensor([[[15.9189, -1.8339,  3.5659],
[ -1.8339,  2.4321, -0.3004],
[  3.5659, -0.3004,  0.8291]],

[[12.9747, -1.1675,  2.8148],
[ -1.1675,  3.9427, -0.1141],
[  2.8148, -0.1141,  0.6946]],

[[15.8620, -1.8119,  3.5485],
[ -1.8119,  2.4394, -0.2971],
[  3.5485, -0.2971,  0.8260]],

[[14.5045, -1.8173,  3.2677],
[ -1.8173,  0.4342, -0.4000],

```

```
[ 3.2677, -0.4000, 0.7402]]])
```

100%|██| 194/194 [00:10<00:00, 18.23it/s]

=====TEST MODE=====

N: 198371

phi :

```
tensor([[6.6445e-04, 2.1110e-03, 8.6473e-04, 9.9635e-01],  
        grad_fn=<DivBackward0>)
```

mu :

```
tensor([[-0.2979, 2.3099, 0.0913],  
        [-0.3379, 2.8321, 0.1679],  
        [-0.2949, 2.1895, 0.0803],  
        [ 2.0863, 0.3287, 0.4579]], grad_fn=<DivBackward0>)
```

cov :

```
tensor([[[[10.4720, -0.1645, 2.2967],  
          [-0.1645, 5.2059, 0.2242],  
          [ 2.2967, 0.2242, 0.5920]],  
  
         [[ 7.5691, 0.5594, 1.5819],  
          [ 0.5594, 5.3244, 0.2878],  
          [ 1.5819, 0.2878, 0.4418]],  
  
         [[10.9752, -0.3034, 2.4145],  
          [-0.3034, 5.0558, 0.1932],  
          [ 2.4145, 0.1932, 0.6160]],  
  
         [[14.5482, -1.7952, 3.2824],  
          [-1.7952, 0.3579, -0.3949],  
          [ 3.2824, -0.3949, 0.7454]]]), grad_fn=<DivBackward0>)
```

Threshold : 6.946145534515381

Accuracy : 0.9746, Precision : 0.9681, Recall : 0.9542, F-score : 0.9611

I copy pasted the testing code here in the notebook so we could play around the results.

**Incrementally compute for the GMM parameters across all training data for a better estimate**

```

solver.data_loader.dataset.mode="train"
solver.dagmm.eval()
N = 0
mu_sum = 0
cov_sum = 0
gamma_sum = 0

for it, (input_data, labels) in enumerate(solver.data_loader):
    input_data = solver.to_var(input_data)
    enc, dec, z, gamma = solver.dagmm(input_data)
    phi, mu, cov = solver.dagmm.compute_gmm_params(z, gamma)

    batch_gamma_sum = torch.sum(gamma, dim=0)

    gamma_sum += batch_gamma_sum
    mu_sum += mu * batch_gamma_sum.unsqueeze(-1) # keep sums of the numerator only
    cov_sum += cov * batch_gamma_sum.unsqueeze(-1).unsqueeze(-1) # keep sums of the numerator only

    N += input_data.size(0)

train_phi = gamma_sum / N
train_mu = mu_sum / gamma_sum.unsqueeze(-1)
train_cov = cov_sum / gamma_sum.unsqueeze(-1).unsqueeze(-1)

print("N:",N)

```

```
print("phi :\n",train_phi)
print("mu :\n",train_mu)
print("cov :\n",train_cov)
```

```
N: 198371
phi :
  tensor([[6.6445e-04, 2.1110e-03, 8.6473e-04, 9.9636e-01],
         grad_fn=<DivBackward0>)]
mu :
  tensor([[[-0.2979,  2.3099,  0.0913],
          [-0.3379,  2.8321,  0.1679],
          [-0.2949,  2.1895,  0.0803],
          [ 2.0863,  0.3287,  0.4579]], grad_fn=<DivBackward0>)]
cov :
  tensor([[[[10.4748, -0.1583,  2.2963],
          [-0.1583,  5.2358,  0.2256],
          [ 2.2963,  0.2256,  0.5921]],

          [[ 7.5787,  0.5723,  1.5813],
          [ 0.5723,  5.3620,  0.2891],
          [ 1.5813,  0.2891,  0.4420]],

          [[10.9772, -0.2983,  2.4141],
          [-0.2983,  5.0827,  0.1944],
          [ 2.4141,  0.1944,  0.6160]],

          [[14.5466, -1.7950,  3.2821],
          [-1.7950,  0.3579, -0.3949],
          [ 3.2821, -0.3949,  0.7453]]], grad_fn=<DivBackward0>)]
```

```
train_energy = []
train_labels = []
train_z = []
for it, (input_data, labels) in enumerate(solver.data_loader):
    input_data = solver.to_var(input_data)
    enc, dec, z, gamma = solver.dagmm(input_data)
    sample_energy, cov_diag = solver.dagmm.compute_energy(z, phi=train_phi, mu=train_mu, cov=train_cov, size_average=False)

    train_energy.append(sample_energy.data.cpu().numpy())
    train_z.append(z.data.cpu().numpy())
    train_labels.append(labels.numpy())

train_energy = np.concatenate(train_energy,axis=0)
train_z = np.concatenate(train_z,axis=0)
train_labels = np.concatenate(train_labels,axis=0)
```

## Compute the energy of every sample in the test data

```
solver.data_loader.dataset.mode="test"
test_energy = []
test_labels = []
test_z = []
for it, (input_data, labels) in enumerate(solver.data_loader):
    input_data = solver.to_var(input_data)
    enc, dec, z, gamma = solver.dagmm(input_data)
    sample_energy, cov_diag = solver.dagmm.compute_energy(z, size_average=False)
    test_energy.append(sample_energy.data.cpu().numpy())
    test_z.append(z.data.cpu().numpy())
    test_labels.append(labels.numpy())

test_energy = np.concatenate(test_energy,axis=0)
test_z = np.concatenate(test_z,axis=0)
```

```
test_labels = np.concatenate(test_labels,axis=0)
```

```
combined_energy = np.concatenate([train_energy, test_energy], axis=0)
combined_z = np.concatenate([train_z, test_z], axis=0)
combined_labels = np.concatenate([train_labels, test_labels], axis=0)
```

**Compute for the threshold energy. Following the paper I just get the highest 20% and treat it as an anomaly. That corresponds to setting the threshold at the 80th percentile.**

```
thresh = np.percentile(combined_energy, 100 - 20)
print("Threshold :", thresh)
```

```
Threshold : 5.002402305603027
```

```
pred = (test_energy>thresh).astype(int)
gt = test_labels.astype(int)
```

```
from sklearn.metrics import precision_recall_fscore_support as prf, accuracy_score
```

```
accuracy = accuracy_score(gt,pred)
precision, recall, f_score, support = prf(gt, pred, average='binary')
```

```
print("Accuracy : {:.4f}, Precision : {:.4f}, Recall : {:.4f}, F-score : {:.4f}".format(accuracy,precision, recall, f_
```



```
Accuracy : 0.9746, Precision : 0.9687, Recall : 0.9537, F-score : 0.9611
```

## Visualizing the z space

It's a little different from the paper's figure but I assume that's because of the small changes in my implementation.

```
from mpl_toolkits.mplot3d import Axes3D
import matplotlib.pyplot as plt
%matplotlib notebook
fig = plt.figure()
ax = fig.add_subplot(111, projection='3d')
ax.scatter(test_z[:,1],test_z[:,0], test_z[:,2], c=test_labels.astype(int))
ax.set_xlabel('Encoded')
ax.set_ylabel('Euclidean')
ax.set_zlabel('Cosine')
plt.show()
```

<IPython.core.display.Javascript object>

