

Dataset1-Regression_output_10

November 2, 2021

1 Dataset 1 - Regression

1.1 Experiment Details

The aim of the experiment is to verify if the: 1. ABC_GAN model corrects model misspecification
2. ABC_GAN model performs better and converges faster than a simple C-GAN model

In the experiment we predict the distribution that represents the real data and simulate realistic fake data points using statistical model, C-GAN and ABC-GAN model with 3 priors. We analyze and compare their performance using metrics like mean squared error, mean absolute error, manhattan distance and euclidean distance between y_{real} and y_{pred}

The models are as follows:

1. The statistical model assumes the distribution $Y = \beta X + \mu$ where $\mu \sim N(0, 1)$
2. The Conditional GAN consists of
 1. Generator with 2 hidden layers with 100 nodes each and ReLu activation.
 2. Discriminator with 2 hidden layers with 25 and 50 nodes and ReLu activation. We use Adam's optimiser and BCE Logit Loss to train the model. The input to the Generator of the GAN is (x,e) where x are the features and $e \sim N(0, 1)$. The discriminator output is linear.
3. The ABC GAN Model consists of
 1. ABC generator is defined as follows:
 1. $Y = 1 + \beta_1 x_1 + \beta_2 x_2 + \beta_2 x_3 + \dots + \beta_n x_n + N(0, \sigma)$ where $\sigma = 0.1$
 2. $\beta_i \sim N(0, \sigma^*)$ when $\mu = 0$ else $\beta_i \sim N(\beta_i^*, \sigma^*)$ where β_i^* s are coefficients obtained from statistical model
 3. σ^* takes the values 0.01, 0.1 and 1
 2. C-GAN network is as defined above. However the input to the Generator of the GAN is (x, y_{abc}) where y_{abc} is the output of the ABC Generator.

1.2 Import Libraries

```
[1]: import warnings
warnings.filterwarnings('ignore')
```

```
[2]: import train_test
import ABC_train_test
import regressionDataset
import network
```

```

import statsModel
import performanceMetrics
import dataset
import sanityChecks
import torch
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import norm
from torch.utils.data import Dataset, DataLoader
from torch import nn

```

1.3 Parameters

General Parameters

1. Number of Samples
2. Number of features

ABC-Generator parameters are as mentioned below: 1. mean : 1 ($\beta \sim N(\beta^*, \sigma)$ where β^* are coefficients of statistical model) or 1 ($\beta \sim N(0, \sigma)$) 2. std : $\sigma = 1, 0.1, 0.01$ (standard deviation)

```

[3]: n_features = 10
     n_samples= 100

     #ABC Generator Parameters
     mean = 1
     variance = 0.001

```

```

[4]: # Parameters
     n_samples = 100
     n_features = 10
     mean = 0
     variance = 0.1

```

1.4 Dataset

Generate a random regression problem

$Y = 1 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n + N(0, \sigma)$ where $\sigma = 0.1$

```

[5]: X,Y = regressionDataset.regression_data(n_samples,n_features)

```

	X1	X2	X3	X4	X5	X6	X7 \
0	0.057305	-0.950282	1.746588	0.009800	1.699977	1.140455	1.087520
1	-1.473913	1.017400	0.318073	1.629009	-0.183935	0.092896	-1.707042
2	-1.382404	0.970019	0.091930	0.255133	-0.227017	-0.274624	1.714881
3	-1.361411	2.017341	0.222785	-0.460051	-0.257047	1.244095	-1.366925
4	0.988020	0.538109	0.739545	0.155162	-0.299528	-0.836996	1.658489
	X8	X9	X10	Y			

```

0 -0.191977  1.451915  2.651454  640.568290
1 -0.857250  0.228631  0.218040  -39.111641
2  0.792337 -0.297779  0.863654   93.563855
3  0.285689 -0.895574 -0.049159 -198.743485
4  0.556129  0.812952 -0.529187  258.250845

```

1.5 Stats Model

```
[6]: [coeff,y_pred] = statsModel.statsModel(X,Y)
```

No handles with labels found to put in legend.

```

                                OLS Regression Results
=====
Dep. Variable:                  Y      R-squared:                1.000
Model:                        OLS    Adj. R-squared:            1.000
Method:                    Least Squares  F-statistic:            4.438e+07
Date:                Tue, 02 Nov 2021    Prob (F-statistic):      1.81e-293
Time:                  18:28:55    Log-Likelihood:         629.21
No. Observations:          100    AIC:                   -1236.
Df Residuals:              89    BIC:                   -1208.
Df Model:                  10
Covariance Type:            nonrobust
=====
               coef      std err          t      P>|t|      [0.025      0.975]
-----
const      2.776e-17   4.75e-05   5.85e-13   1.000   -9.43e-05   9.43e-05
x1          0.3902    5.06e-05  7704.224   0.000    0.390    0.390
x2          0.1218    4.95e-05  2462.375   0.000    0.122    0.122
x3          0.5013    4.85e-05  1.03e+04   0.000    0.501    0.501
x4          0.3713    4.88e-05  7611.220   0.000    0.371    0.371
x5          0.3469    4.85e-05  7148.629   0.000    0.347    0.347
x6          0.1511    4.92e-05  3073.005   0.000    0.151    0.151
x7          0.3807    4.94e-05  7699.774   0.000    0.381    0.381
x8          0.0785    4.91e-05  1599.985   0.000    0.078    0.079
x9          0.2470    4.84e-05  5107.072   0.000    0.247    0.247
x10         0.3401    4.9e-05   6946.198   0.000    0.340    0.340
=====
Omnibus:                 1.596    Durbin-Watson:           1.888
Prob(Omnibus):           0.450    Jarque-Bera (JB):        1.534
Skew:                   -0.296    Prob(JB):                0.464
Kurtosis:               2.865    Cond. No.                 1.52
=====

```

Notes:

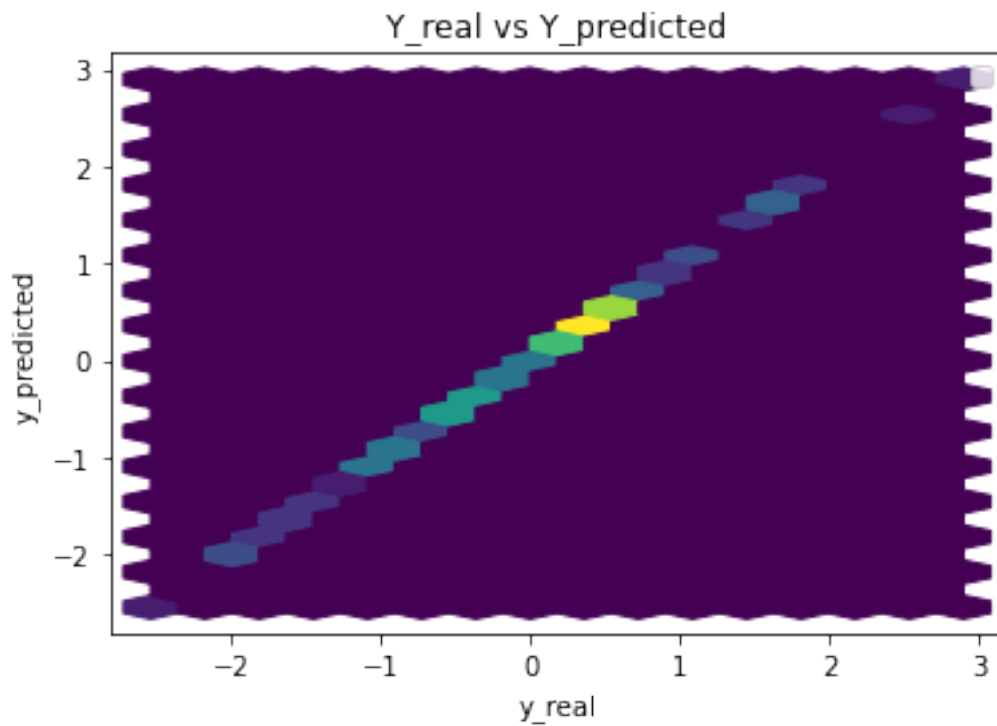
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```

Parameters:  const      2.775558e-17
x1          3.901733e-01

```

```
x2      1.217682e-01
x3      5.013407e-01
x4      3.712575e-01
x5      3.468962e-01
x6      1.511231e-01
x7      3.807135e-01
x8      7.848161e-02
x9      2.470308e-01
x10     3.400863e-01
dtype: float64
```



Performance Metrics

```
Mean Squared Error: 2.0055520171492347e-07
Mean Absolute Error: 0.0003560384852008077
Manhattan distance: 0.03560384852008077
Euclidean distance: 0.004478338996937631
```

1.6 Common Training Parameters (GAN & ABC_GAN)

```
[7]: n_epochs = 5000
     error = 0.001
     batch_size = n_samples
```

1.7 GAN Model

```
[8]: real_dataset = dataset.CustomDataset(X,Y)
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
```

Training GAN for n_epochs number of epochs

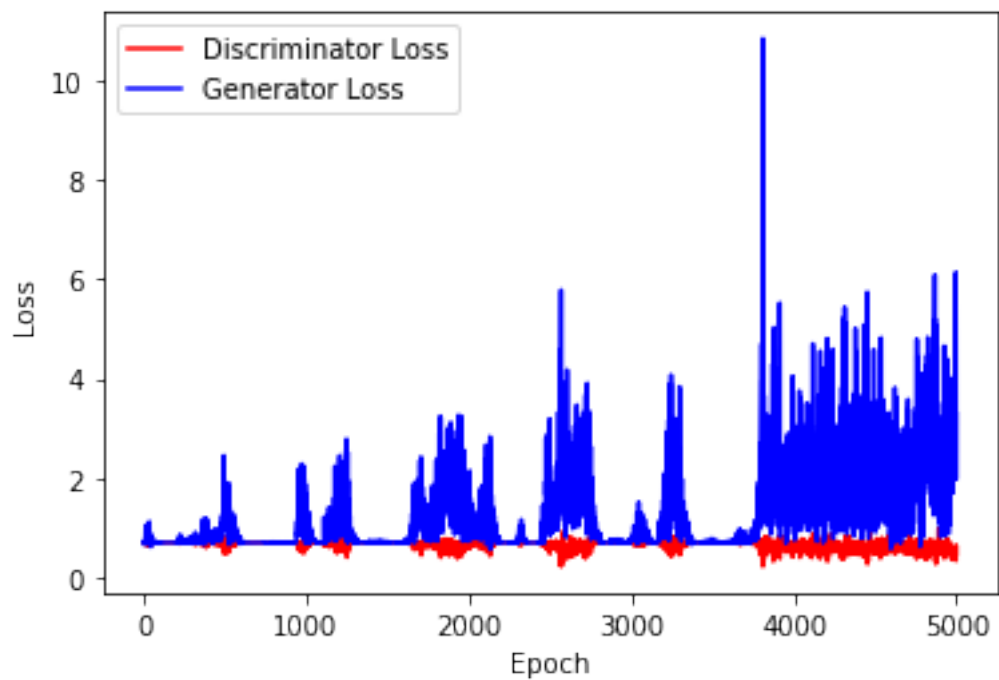
```
[9]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)

criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator.parameters(), lr=0.01, betas=(0.5, 0.
→999))
```

```
[10]: print(generator)
print(discriminator)
```

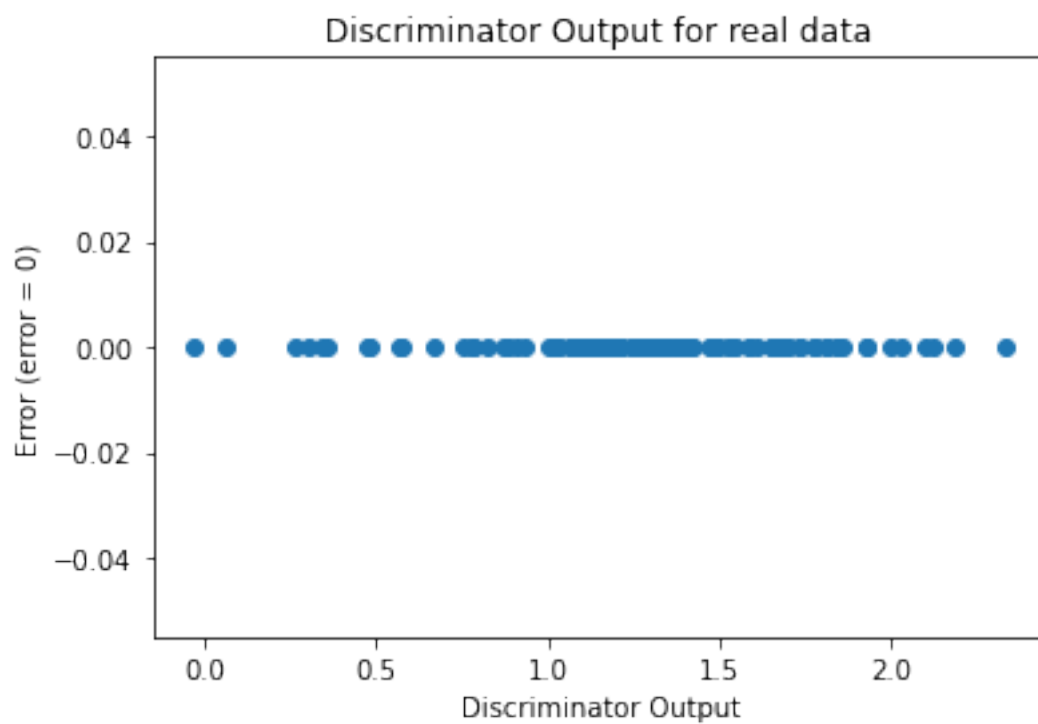
```
Generator(
  (hidden1): Linear(in_features=12, out_features=100, bias=True)
  (hidden2): Linear(in_features=100, out_features=100, bias=True)
  (output): Linear(in_features=100, out_features=1, bias=True)
  (relu): ReLU()
)
Discriminator(
  (hidden1): Linear(in_features=12, out_features=25, bias=True)
  (hidden2): Linear(in_features=25, out_features=50, bias=True)
  (output): Linear(in_features=50, out_features=1, bias=True)
  (relu): ReLU()
)
```

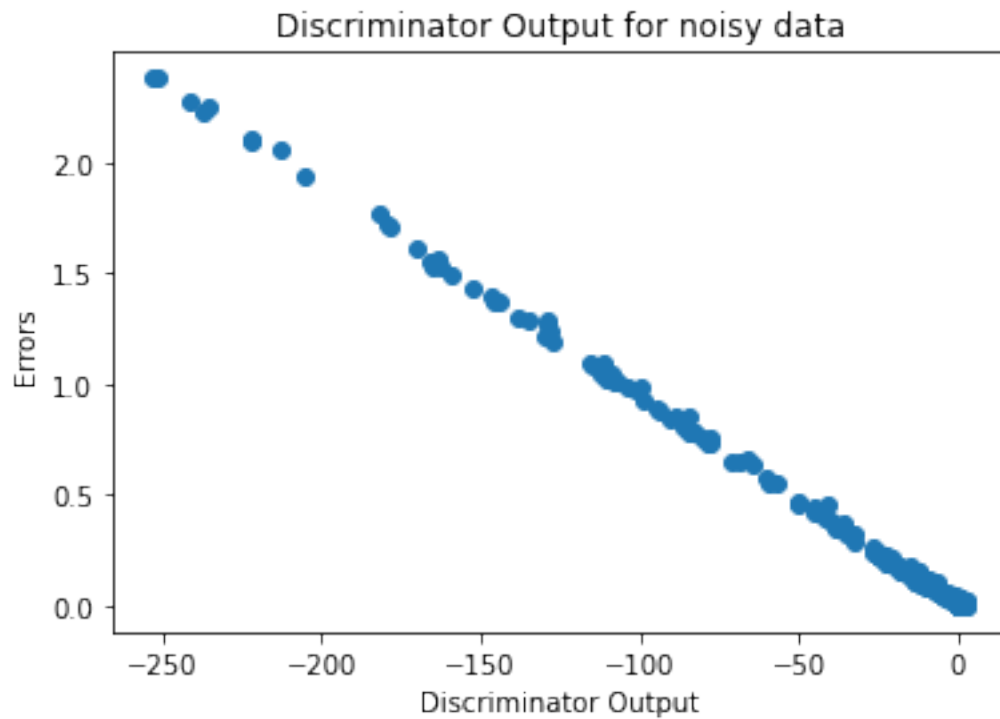
```
[11]: train_test.
→training_GAN(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,
→n_epochs,criterion,device)
```



```
[12]: GAN1_metrics = train_test.test_generator(generator,real_dataset,device)
```

```
[13]: sanityChecks.discProbVsError(real_dataset,discriminator,device)
```



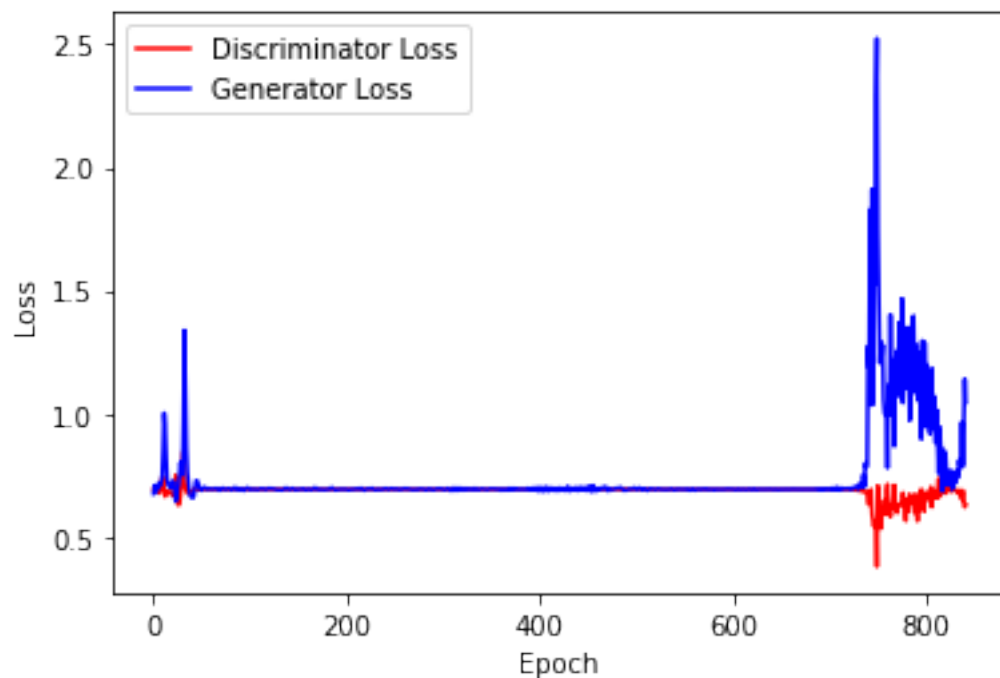


Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

```
[14]: generator2 = network.Generator(n_features+2)
discriminator2 = network.Discriminator(n_features+2)
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator2.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator2.parameters(), lr=0.01, betas=(0.5, 0.
↪999))
```

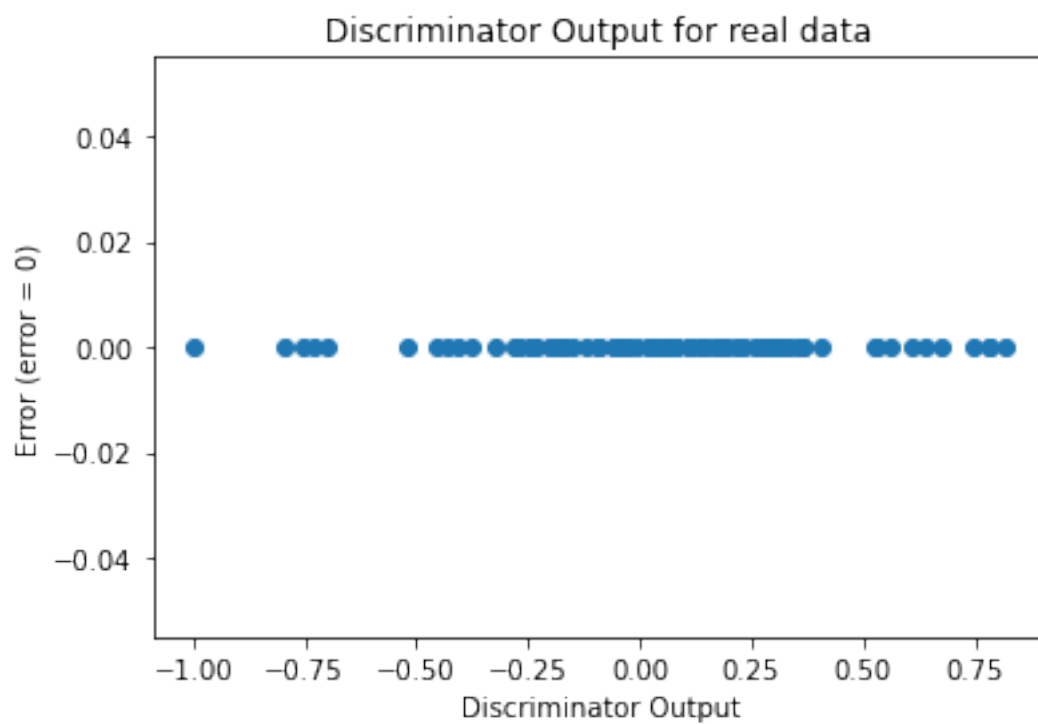
```
[15]: train_test.
↪training_GAN_2(discriminator2,generator2,disc_opt,gen_opt,real_dataset,batch_size,error,cri
```

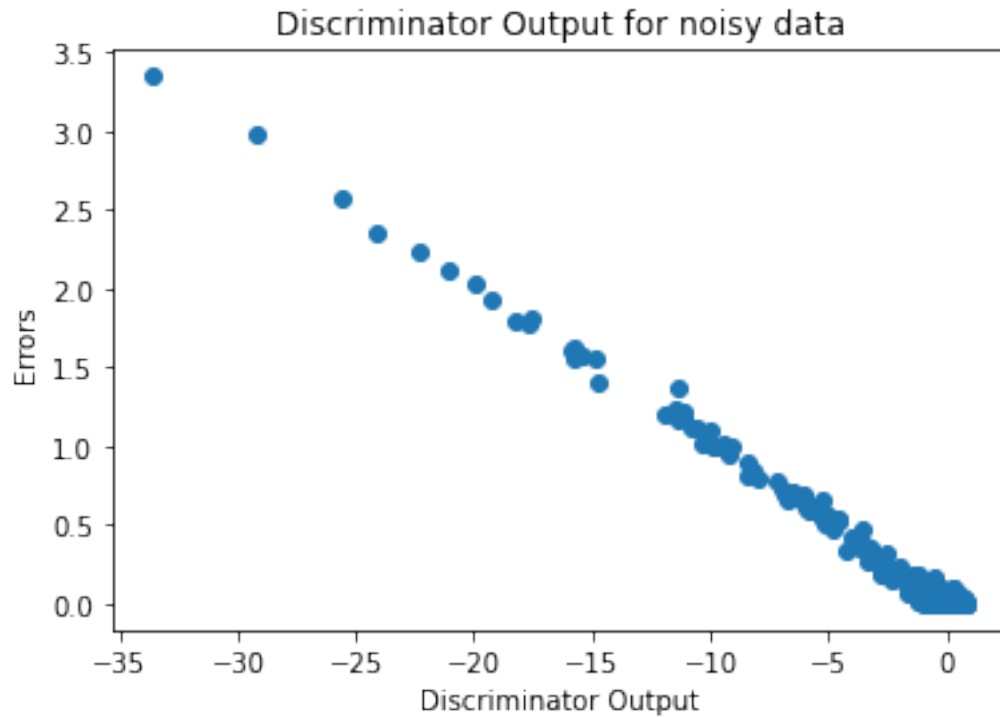
Number of epochs needed 842



```
[16]: GAN2_metrics=train_test.test_generator_2(generator2,real_dataset,device)
```

```
[17]: sanityChecks.discProbVsError(real_dataset,discriminator2,device)
```





2 ABC GAN Model

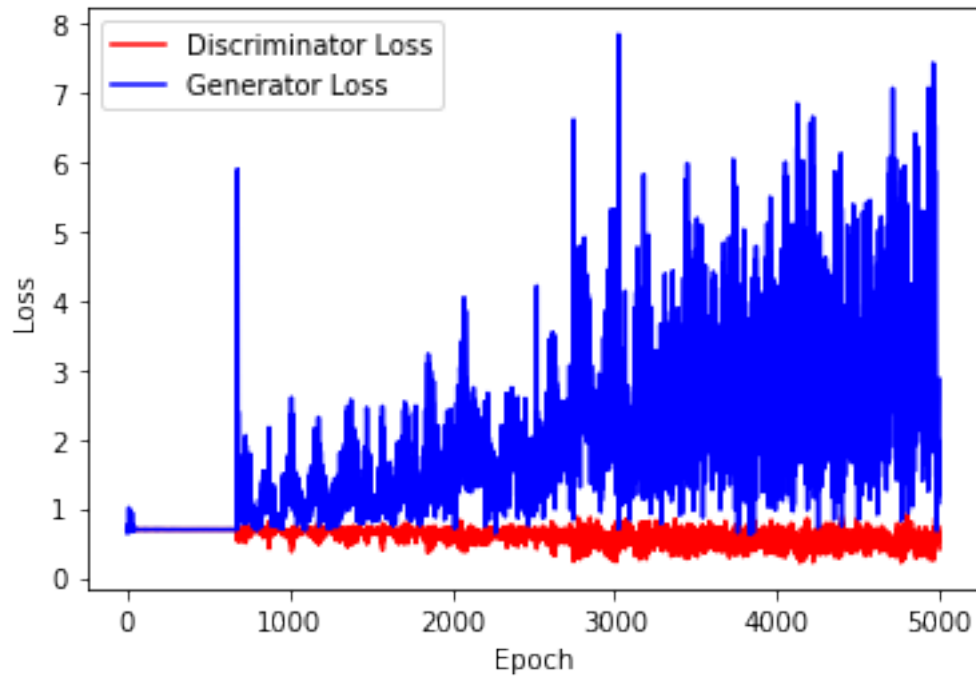
2.0.1 Training the network

Training ABC-GAN for `n_epochs` number of epochs

```
[18]: gen = network.Generator(n_features+2)
      disc = network.Discriminator(n_features+2)

      criterion = torch.nn.BCEWithLogitsLoss()
      gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
      disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))

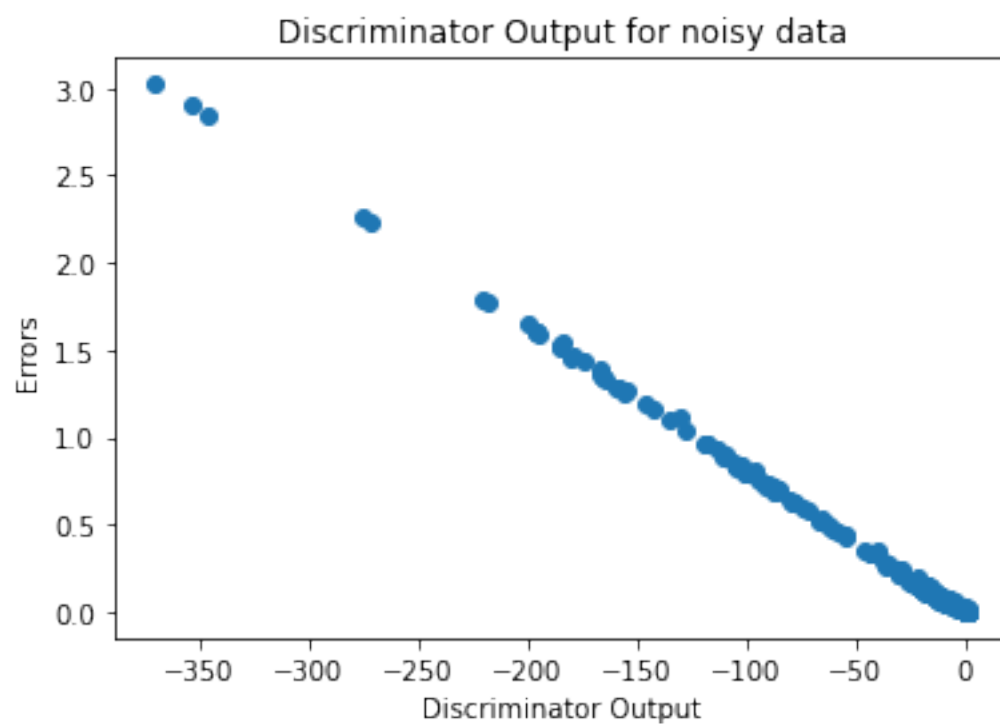
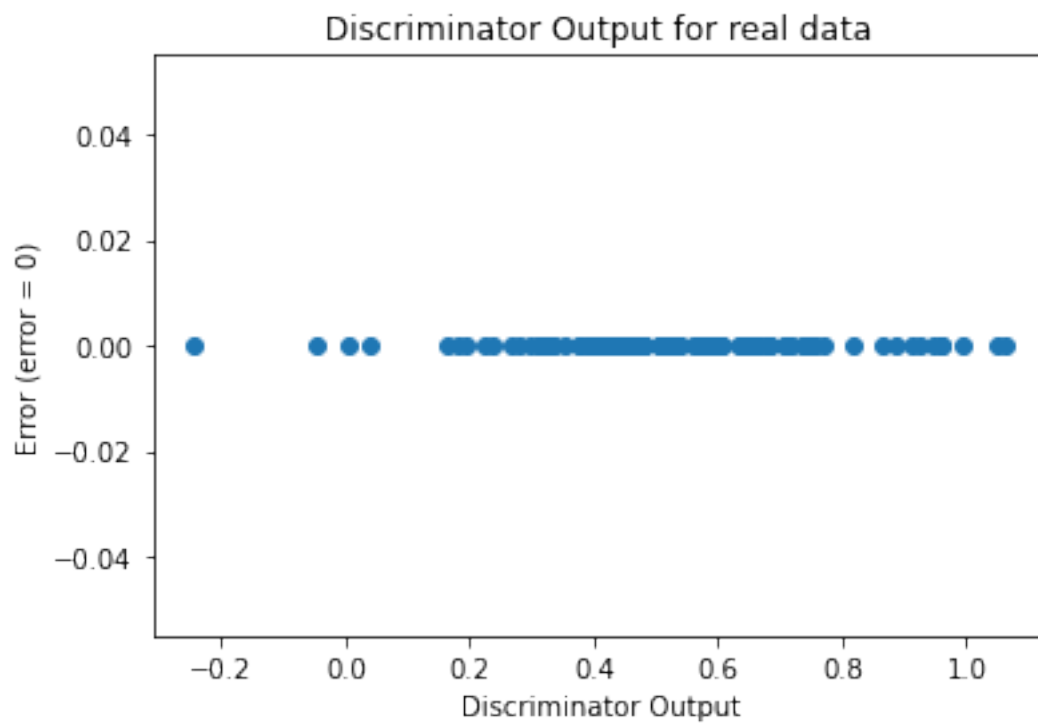
[19]: ABC_train_test.training_GAN(disc, gen,disc_opt,gen_opt,real_dataset,
      ↪batch_size, n_epochs,criterion,coeff,mean,variance,device)
```



```
[20]: ABC_GAN1_metrics=ABC_train_test.  
      ↪ test_generator(gen,real_dataset,coeff,mean,variance,device)
```

Sanity Checks

```
[21]: sanityChecks.discProbVsError(real_dataset,disc,device)
```



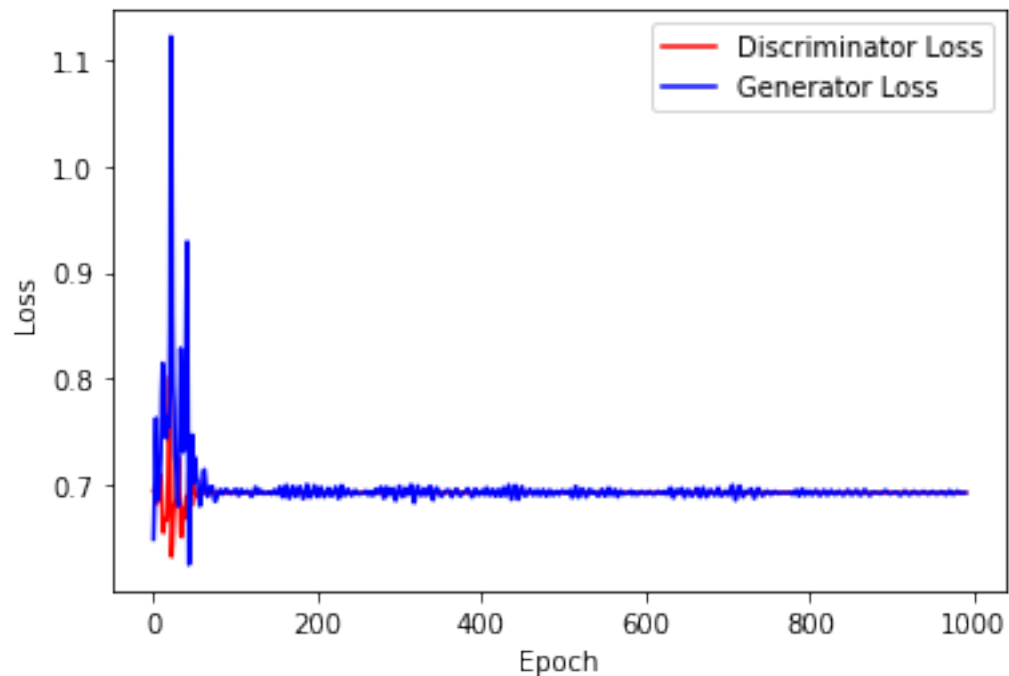
Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

```
[22]: gen2 = network.Generator(n_features+2)
disc2 = network.Discriminator(n_features+2)

criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(gen2.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(disc2.parameters(), lr=0.01, betas=(0.5, 0.999))
```

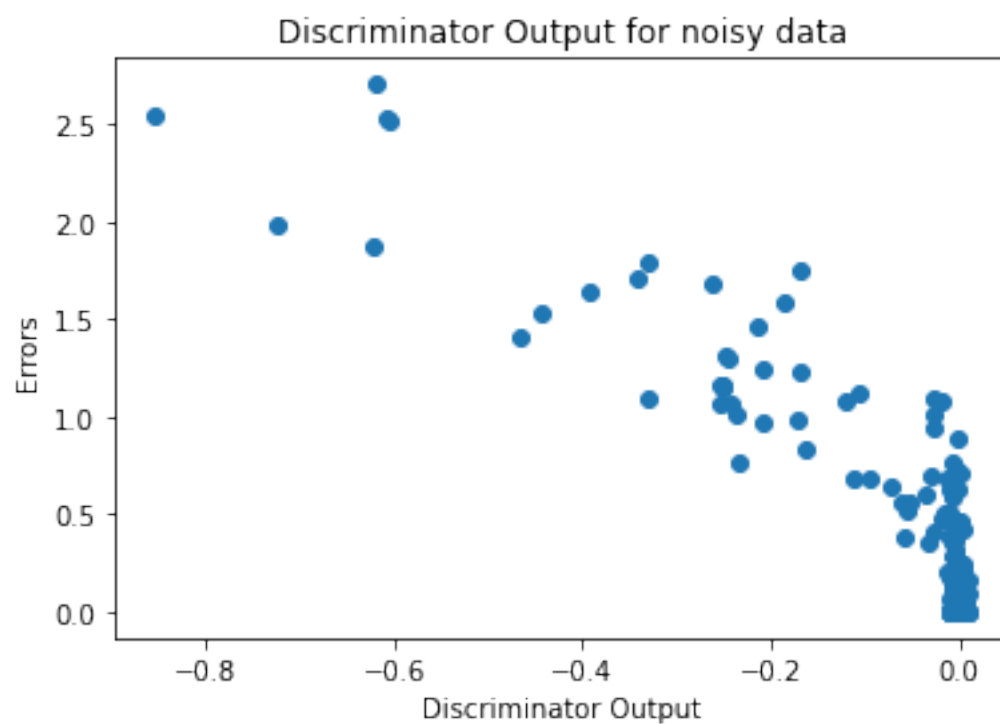
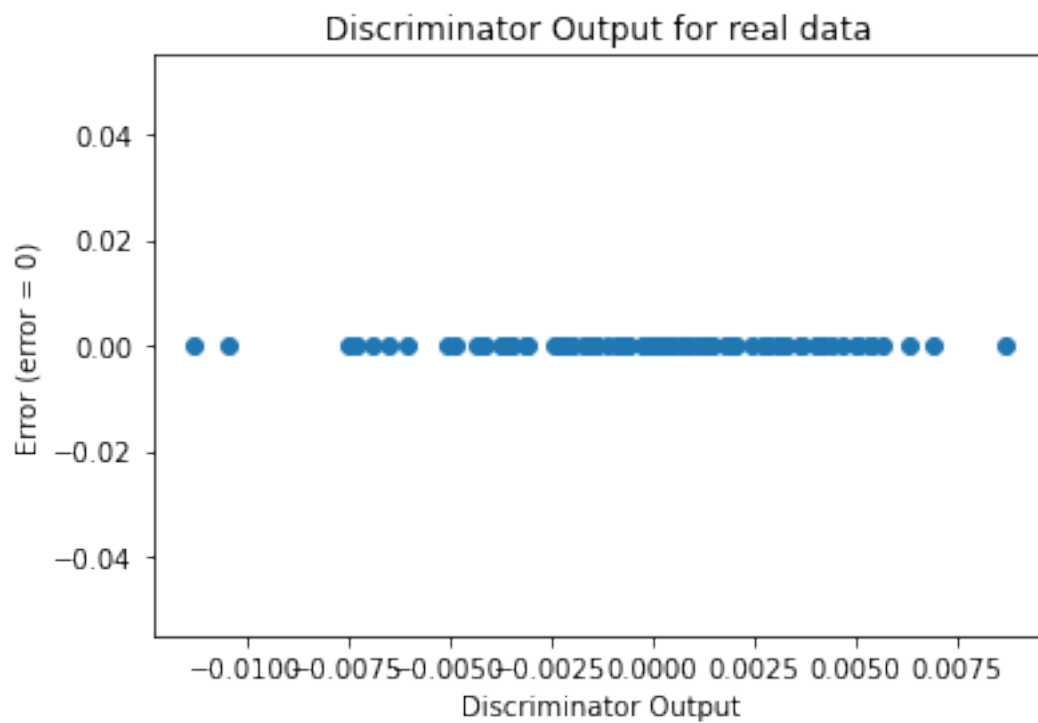
```
[23]: ABC_train_test.
      ↪ training_GAN_2(disc2,gen2,disc_opt,gen_opt,real_dataset,batch_size,
      ↪ error,criterion,coeff,mean,variance,device)
```

Number of epochs 991



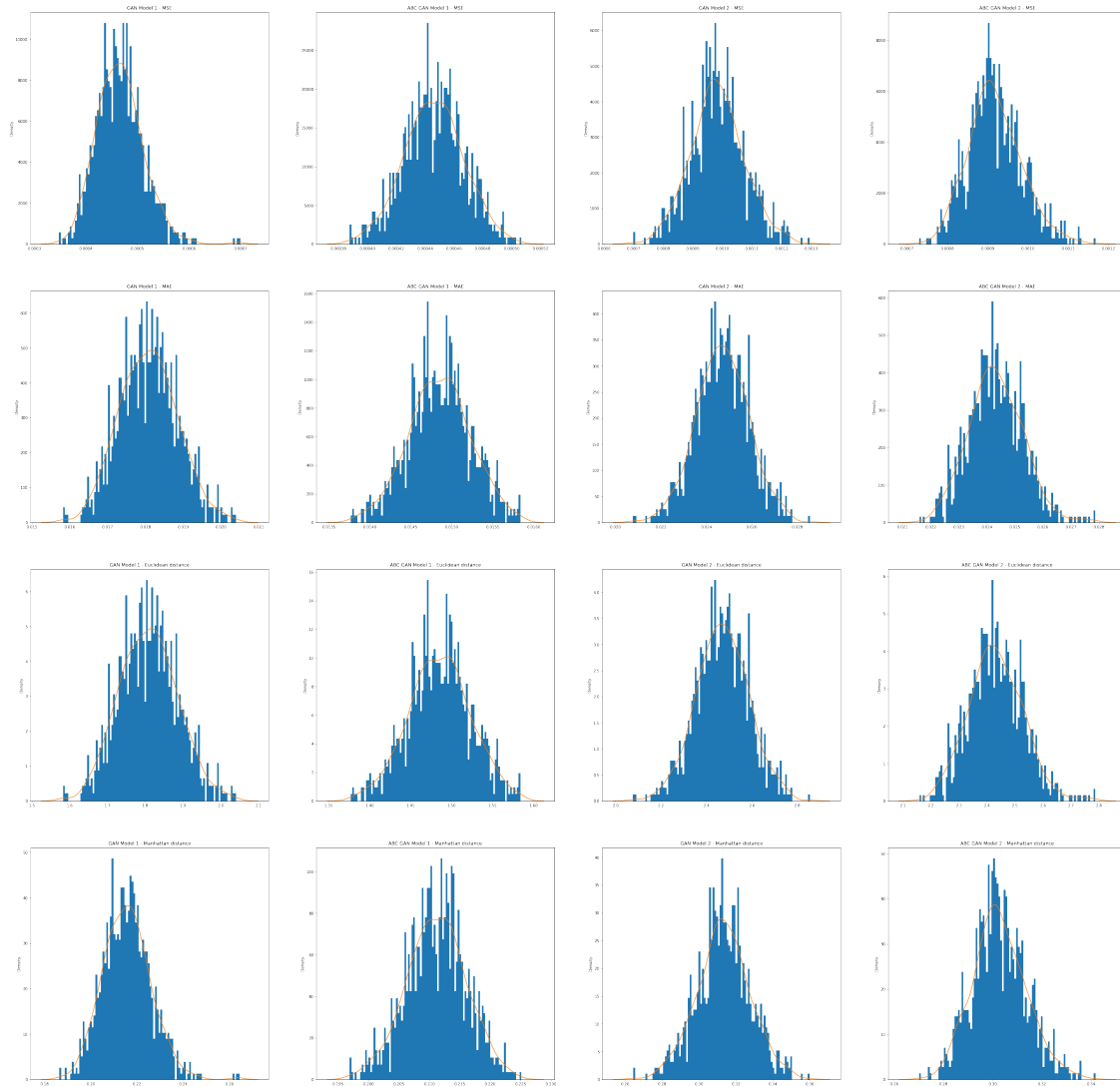
```
[24]: ABC_GAN2_metrics=ABC_train_test.
      ↪ test_generator_2(gen2,real_dataset,coeff,mean,variance,device)
```

```
[25]: sanityChecks.discProbVsError(real_dataset,disc2,device)
```



3 Model Analysis

```
[26]: performanceMetrics.  
      ↪ modelAnalysis(GAN1_metrics,ABC_GAN1_metrics,GAN2_metrics,ABC_GAN2_metrics)
```



```
[ ]:
```