

Dataset2_Friedman1_output_2

October 20, 2021

1 Dataset 2 - Friedman 1

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. $Y = 1 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n + N(0, \sigma)$ 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 friedman1Dataset
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 = 10
     n_features = 10
     mean = 1
     variance = 0.01

```

1.4 Dataset

Friedman 1 Dataset

- $y(X) = 10 * \sin(\pi * X_0 * X_1) + 20 * (X_2 - 0.5) * 2 + 10 * X_3 + 5 * X_4 + noise * N(0, 1)$.
- Only 5 features used to calculate y
- Noise is Gaussian
- 1000 datapoints and 10 features used in the following experiment

```

[5]: X, Y = friedman1Dataset.friedman1_data(n_samples, n_features)

```

	X0	X1	X2	X3	X4	X5	X6 \
0	0.930455	0.130249	0.832503	0.140883	0.663382	0.983597	0.308009
1	0.690604	0.996593	0.582398	0.910236	0.175759	0.685408	0.215353
2	0.196575	0.840676	0.101300	0.401513	0.950036	0.280379	0.902188

```

3  0.091559  0.361221  0.218400  0.623043  0.101839  0.255072  0.127227
4  0.935940  0.368371  0.163757  0.962080  0.814816  0.294416  0.642267

```

```

          X7          X8          X9          Y
0  0.278949  0.223103  0.040534  10.763963
1  0.027935  0.322671  0.251794  18.339439
2  0.631450  0.460941  0.812878  16.926175
3  0.213472  0.036957  0.511442   9.327138
4  0.360630  0.003228  0.208410  24.786829

```

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:              nan
Method:                    Least Squares      F-statistic:              nan
Date:                Wed, 20 Oct 2021      Prob (F-statistic):          nan
Time:                19:54:53      Log-Likelihood:            328.40
No. Observations:                10      AIC:                      -636.8
Df Residuals:                    0      BIC:                      -633.8
Df Model:                        9
Covariance Type:                nonrobust
=====
              coef      std err          t      P>|t|      [0.025      0.975]
-----
const          -1.11e-16         inf         -0         nan         nan         nan
x1              1.0591         inf          0         nan         nan         nan
x2              0.8453         inf          0         nan         nan         nan
x3             -0.9697         inf         -0         nan         nan         nan
x4             -0.0362         inf         -0         nan         nan         nan
x5              0.1527         inf          0         nan         nan         nan
x6              0.0686         inf          0         nan         nan         nan
x7             -0.1303         inf         -0         nan         nan         nan
x8              0.0223         inf          0         nan         nan         nan
x9             -0.2495         inf         -0         nan         nan         nan
x10             0.0299         inf          0         nan         nan         nan
=====
Omnibus:                1.342      Durbin-Watson:                2.060
Prob(Omnibus):          0.511      Jarque-Bera (JB):              0.720
Skew:                   0.622      Prob(JB):                      0.698
Kurtosis:               2.578      Cond. No.                      23.9
=====

```

Notes:

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

[2] The input rank is higher than the number of observations.

Parameters: const -1.110223e-16

x1 1.059110e+00

x2 8.453255e-01

x3 -9.697476e-01

x4 -3.621278e-02

x5 1.526613e-01

x6 6.857384e-02

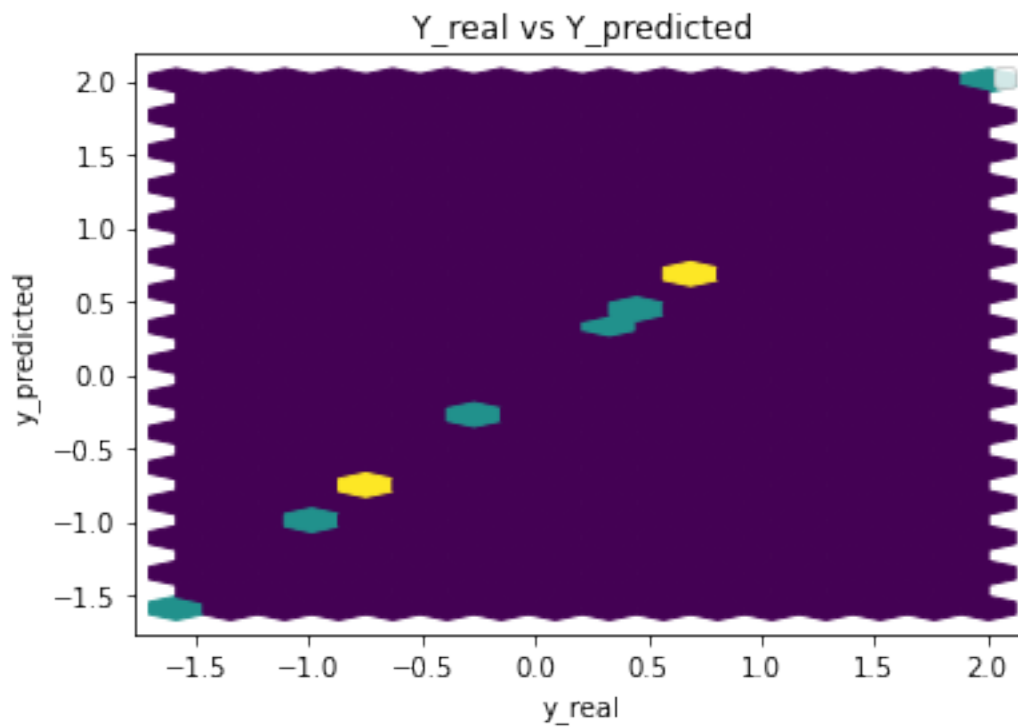
x7 -1.303480e-01

x8 2.228488e-02

x9 -2.494594e-01

x10 2.989860e-02

dtype: float64



Performance Metrics

Mean Squared Error: 1.7490525382947122e-30

Mean Absolute Error: 1.1435297153639113e-15

Manhattan distance: 1.1435297153639112e-14

Euclidean distance: 4.182167546015716e-15

1.6 Common Training Parameters (GAN & ABC_GAN)

```
[7]: n_epochs = 5000
     error = 0.001
     batch_size = n_samples//2
```

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

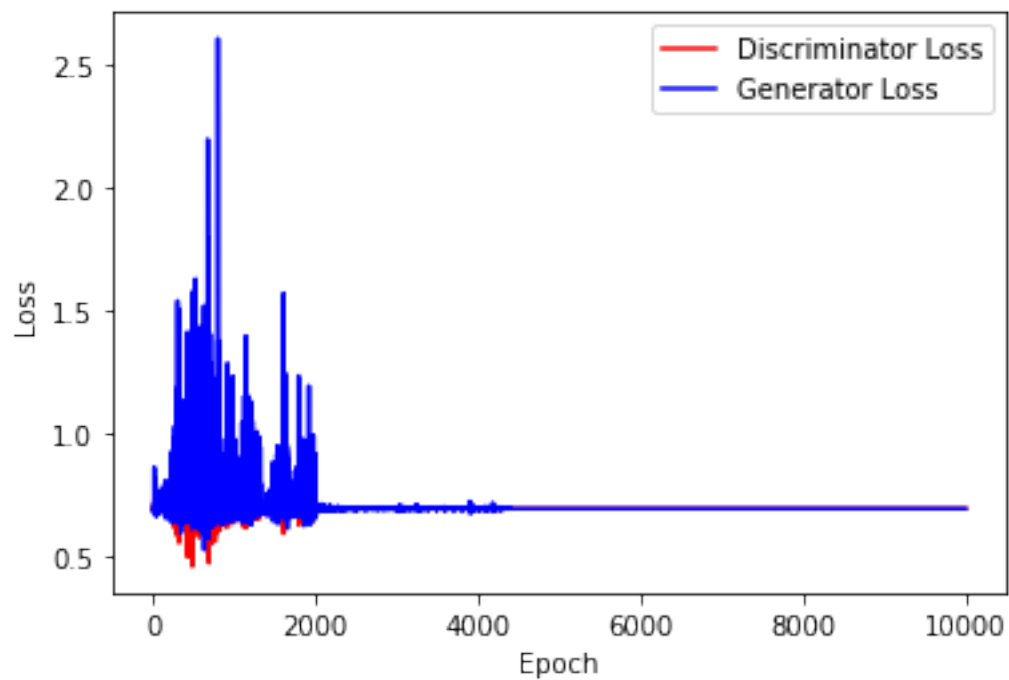
```
[9]: generator = network.Generator(n_features+2)
     discriminator = network.Discriminator(n_features+2)

     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))
```

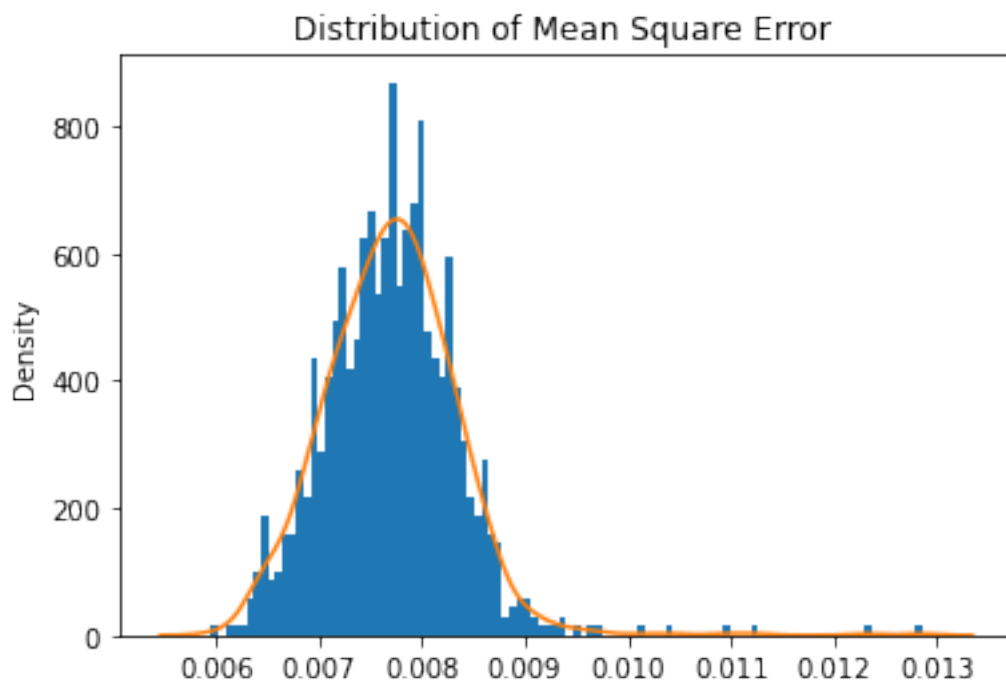
```
[10]: print(generator)
      print(discriminator)
```

```
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()
)
```

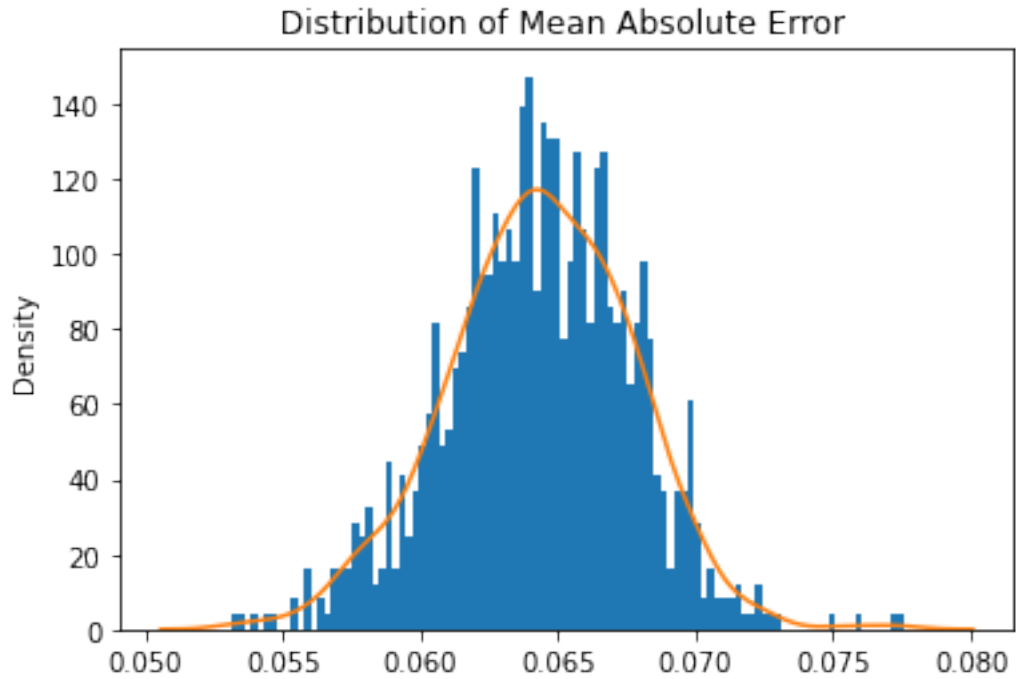
```
[11]: train_test.
     ↪training_GAN(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,
     ↪n_epochs,criterion,device)
```



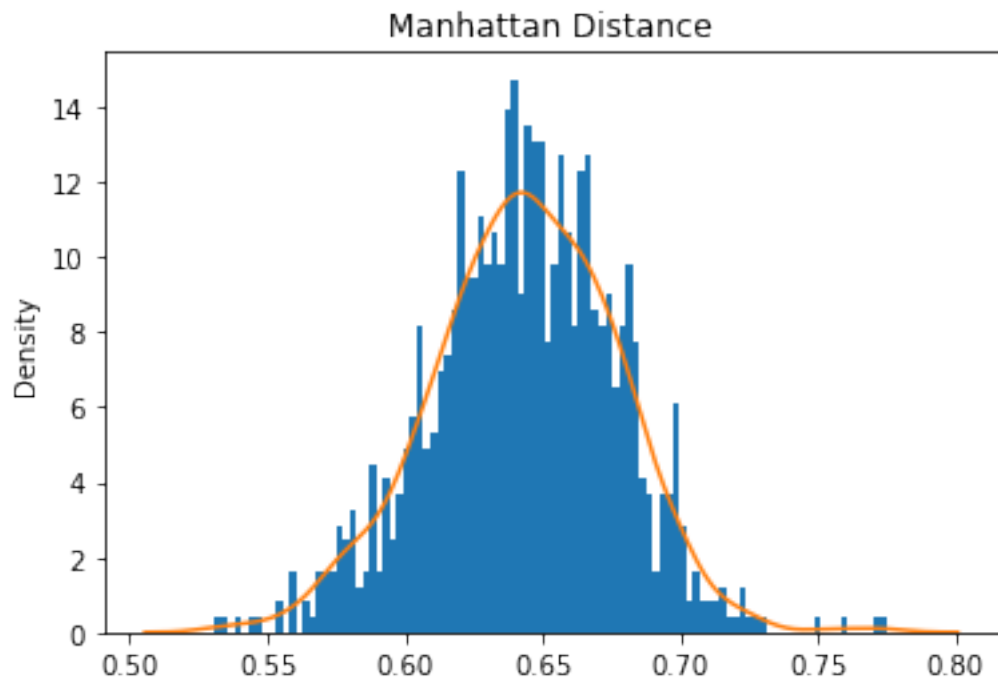
```
[12]: train_test.test_generator(generator,real_dataset,device)
```



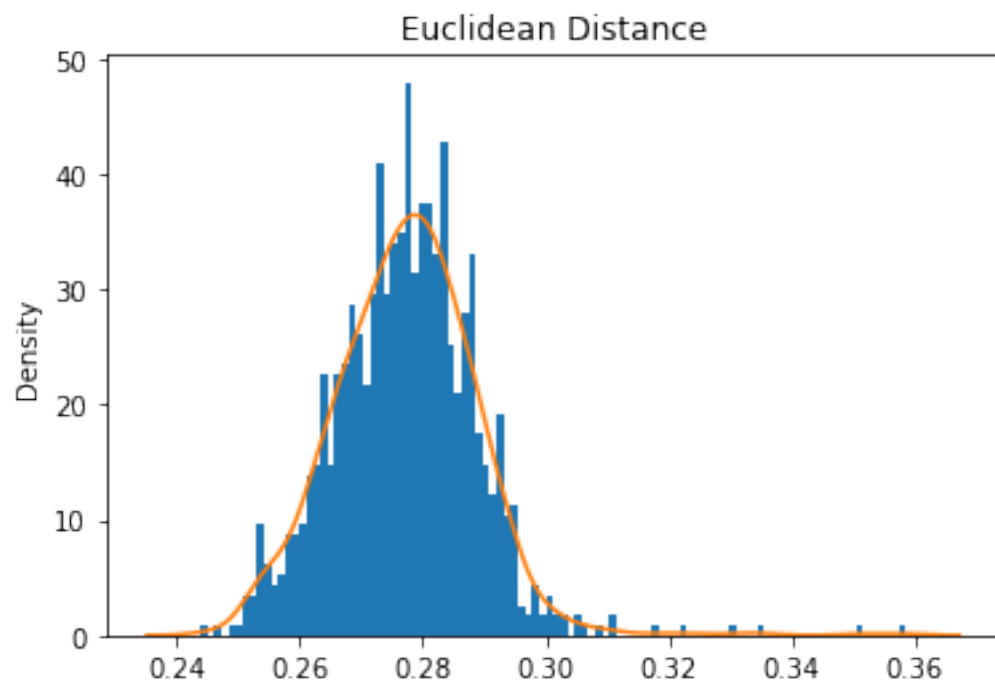
Mean Square Error: 0.007699408372407021



Mean Absolute Error: 0.06432816494703293

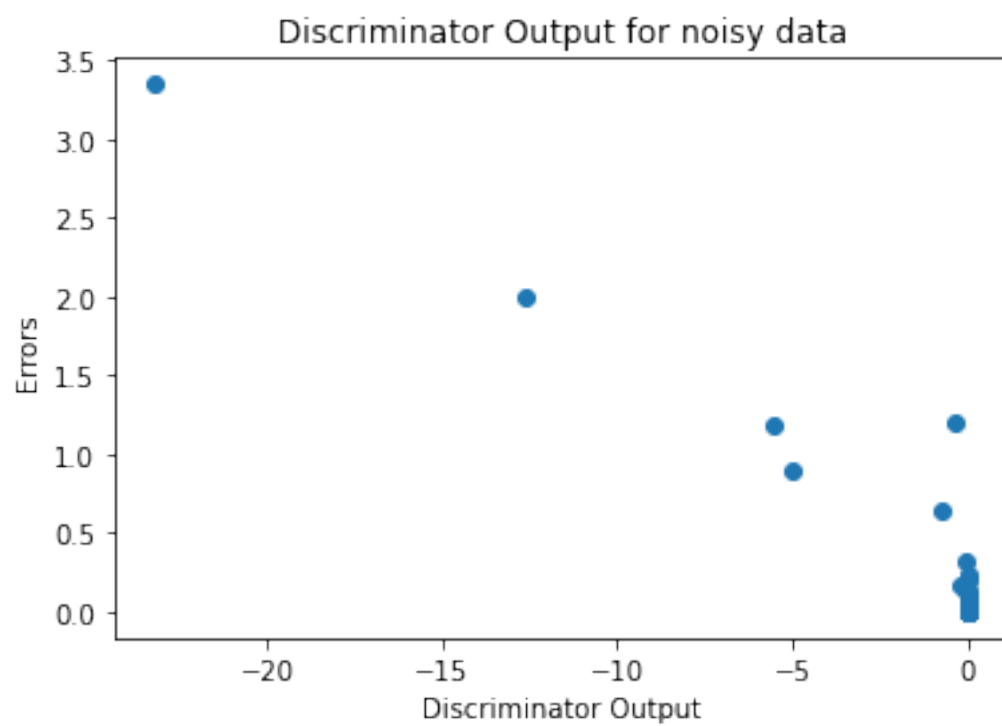
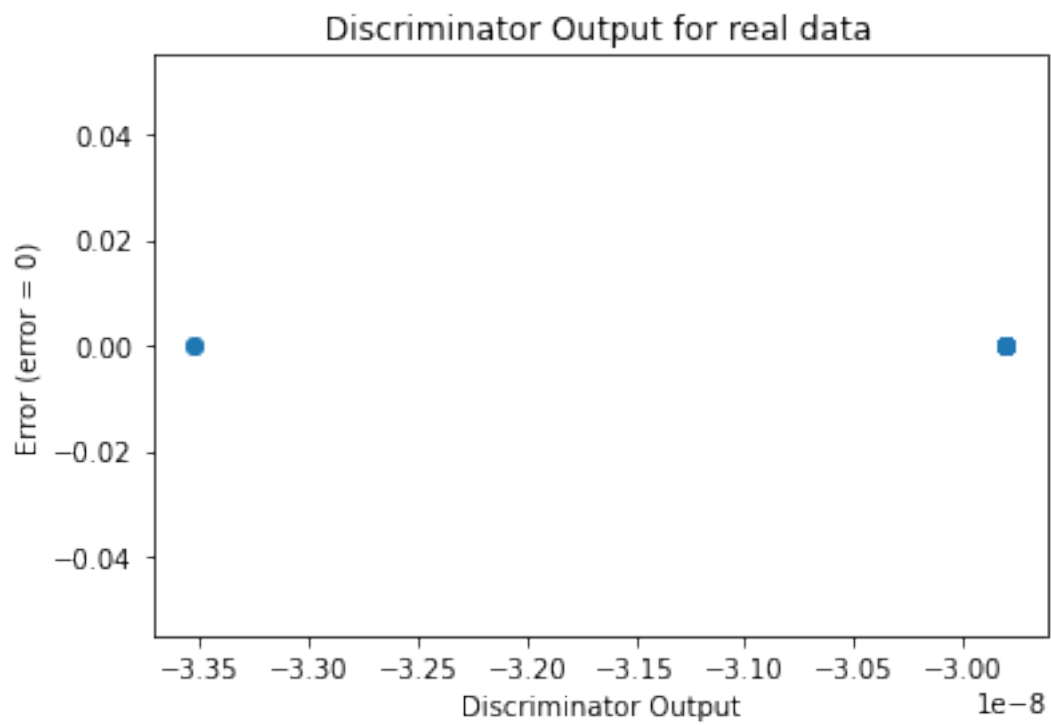


Mean Manhattan Distance: 0.6432816494703293



Mean Euclidean Distance: 0.27723686048482166

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

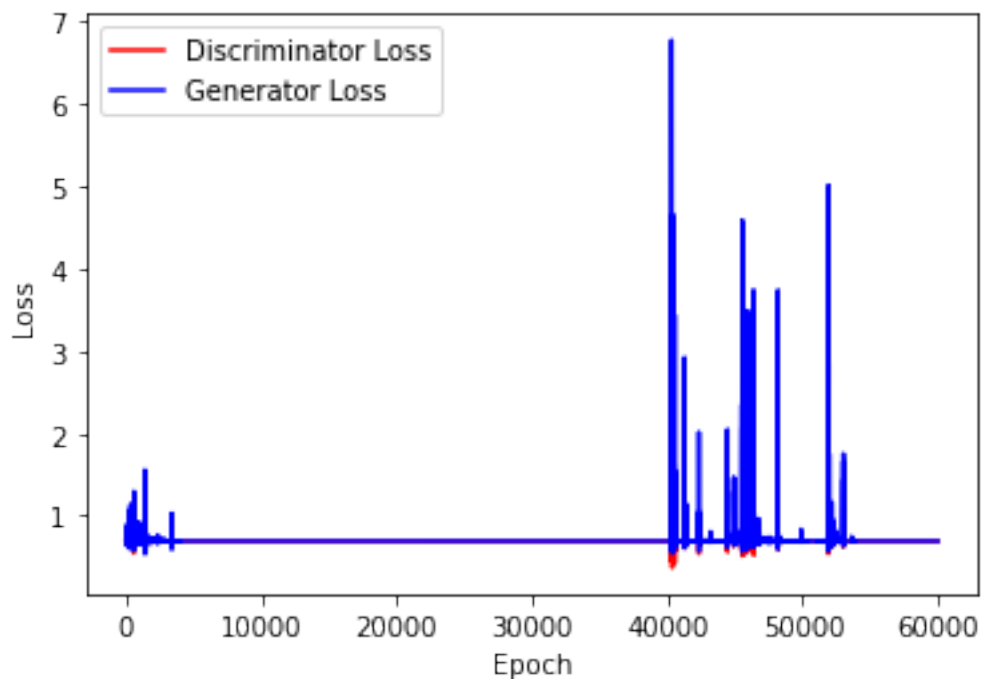



Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

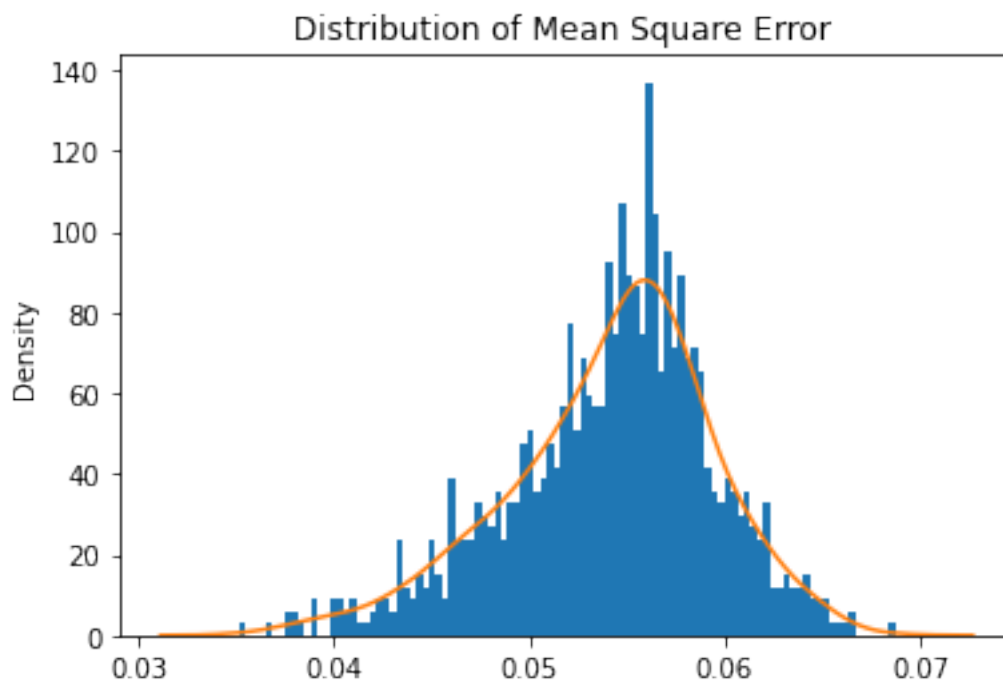
```
[14]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator.parameters(), lr=0.01, betas=(0.5,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

```
[15]: train_test.
↪training_GAN_2(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,error,crite
```

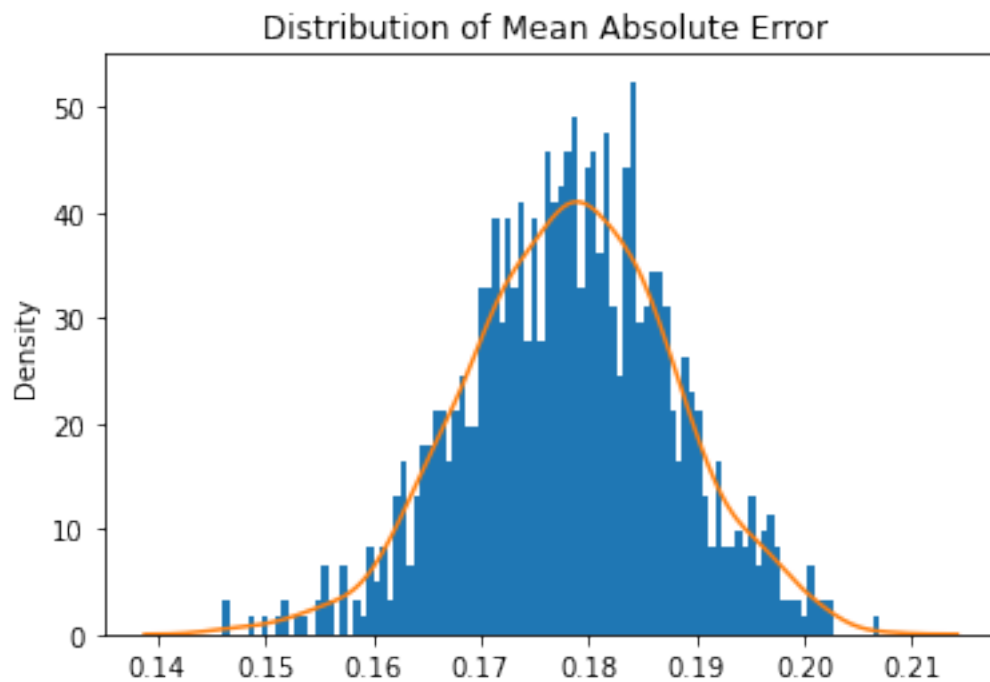
Number of epochs needed 30000



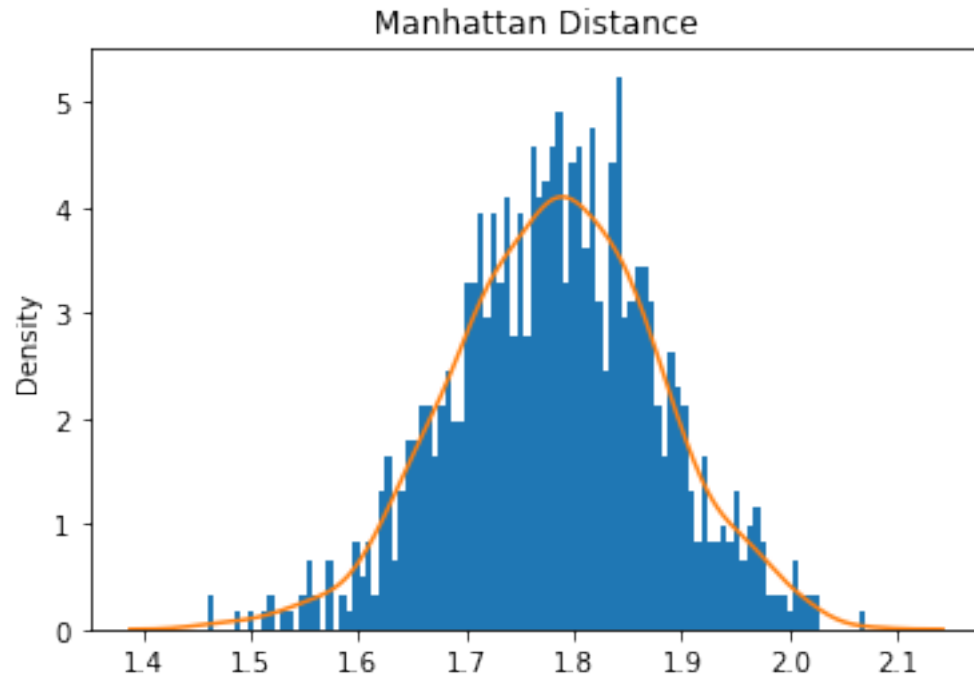
```
[16]: train_test.test_generator(generator,real_dataset,device)
```



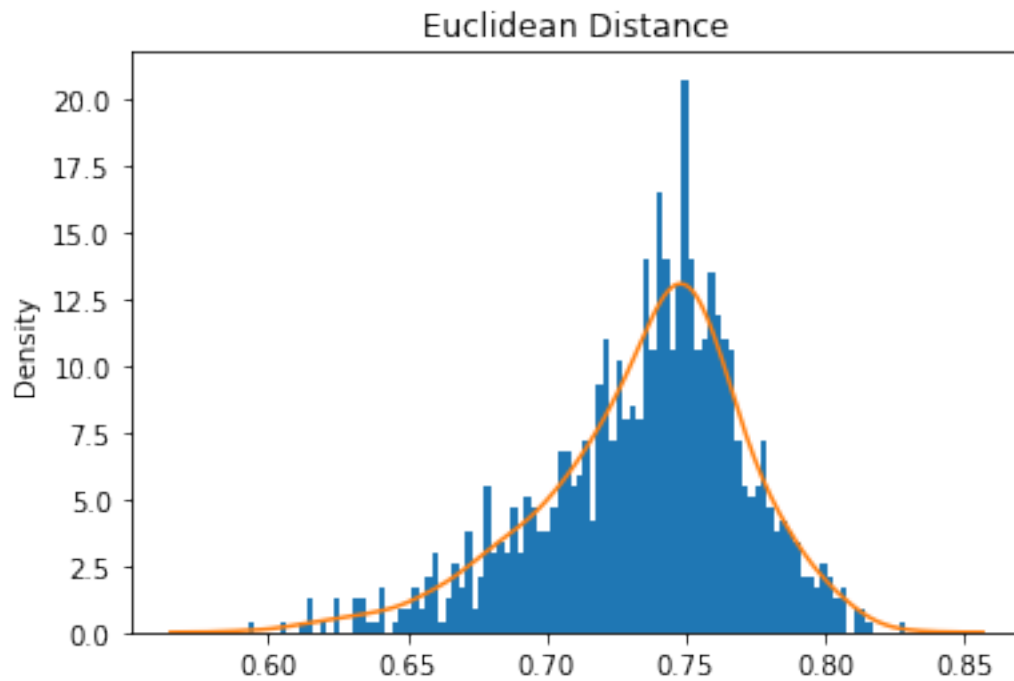
Mean Square Error: 0.05415171025173744



Mean Absolute Error: 0.17829674466103315



Mean Manhattan Distance: 1.7829674466103316



Mean Euclidean Distance: 0.7349637314456919

2 ABC GAN Model

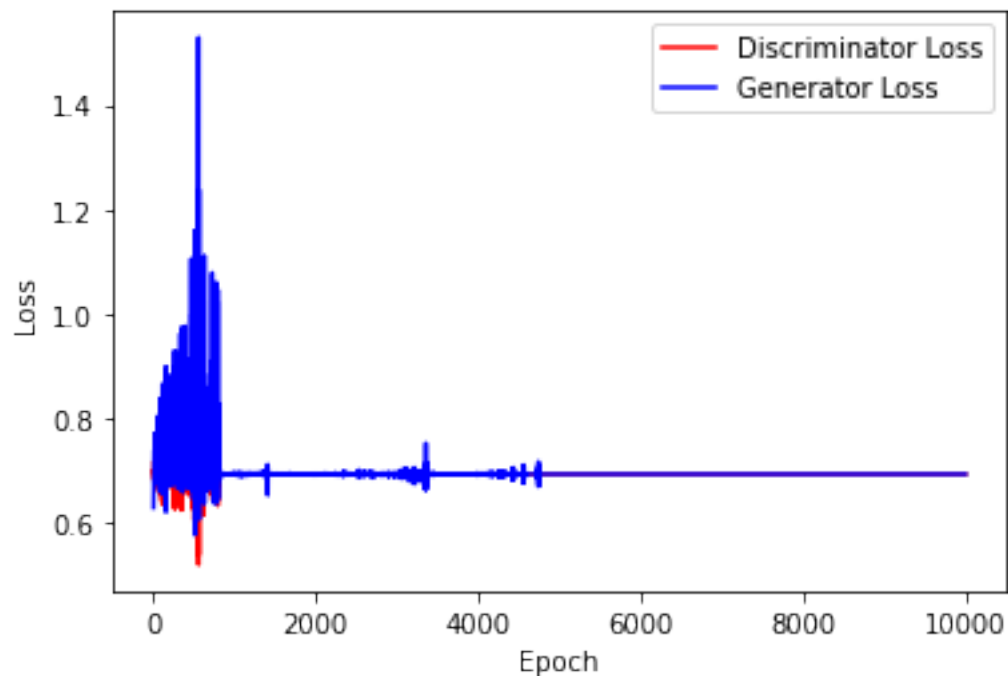
2.0.1 Training the network

Training ABC-GAN for `n_epochs` number of epochs

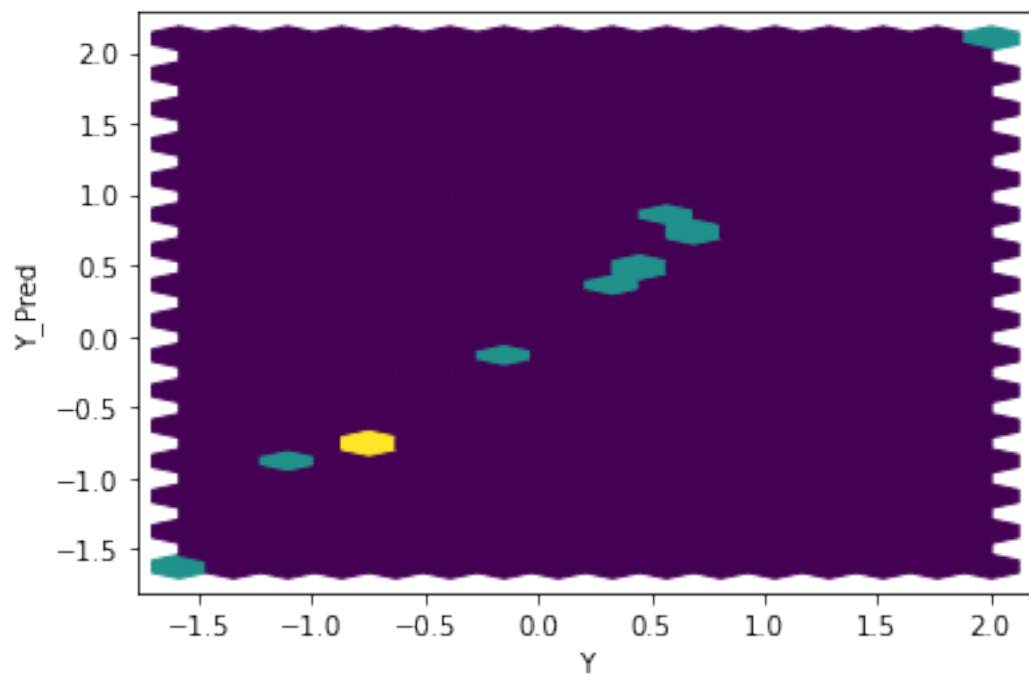
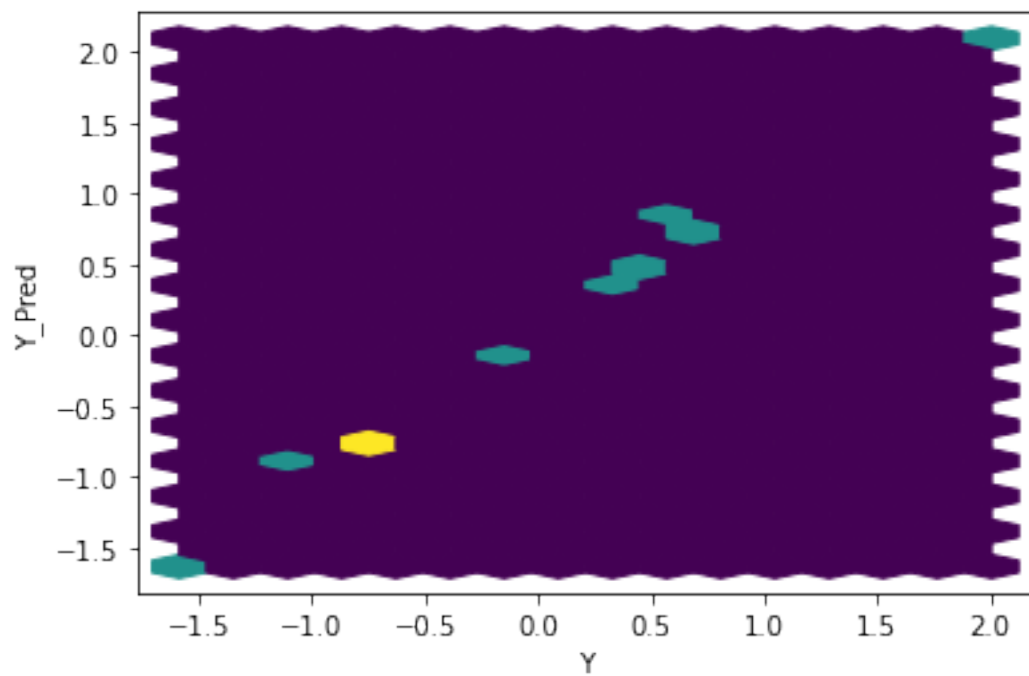
```
[17]: gen = network.Generator(n_features+2)
      disc = network.Discriminator(n_features+2)

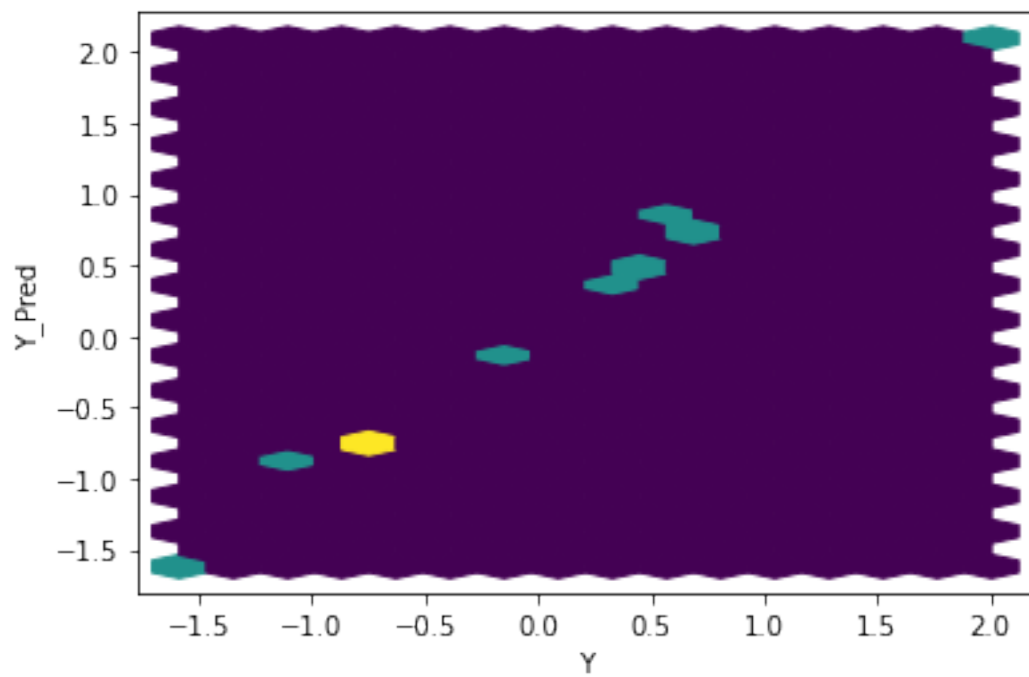
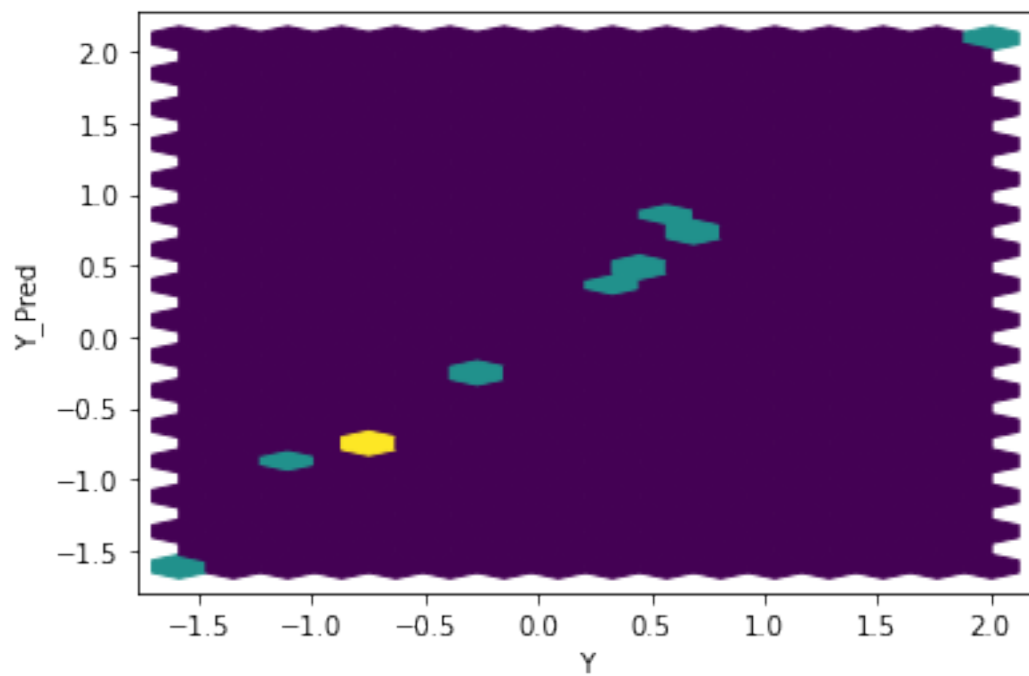
      criterion = torch.nn.BCEWithLogitsLoss()
      gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
      disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))

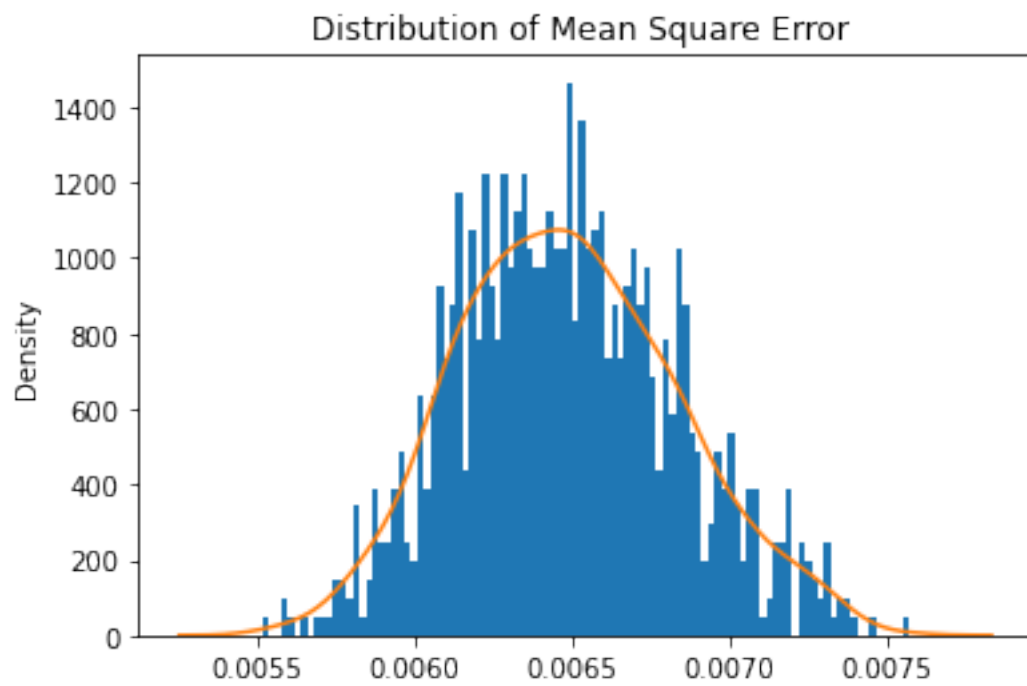
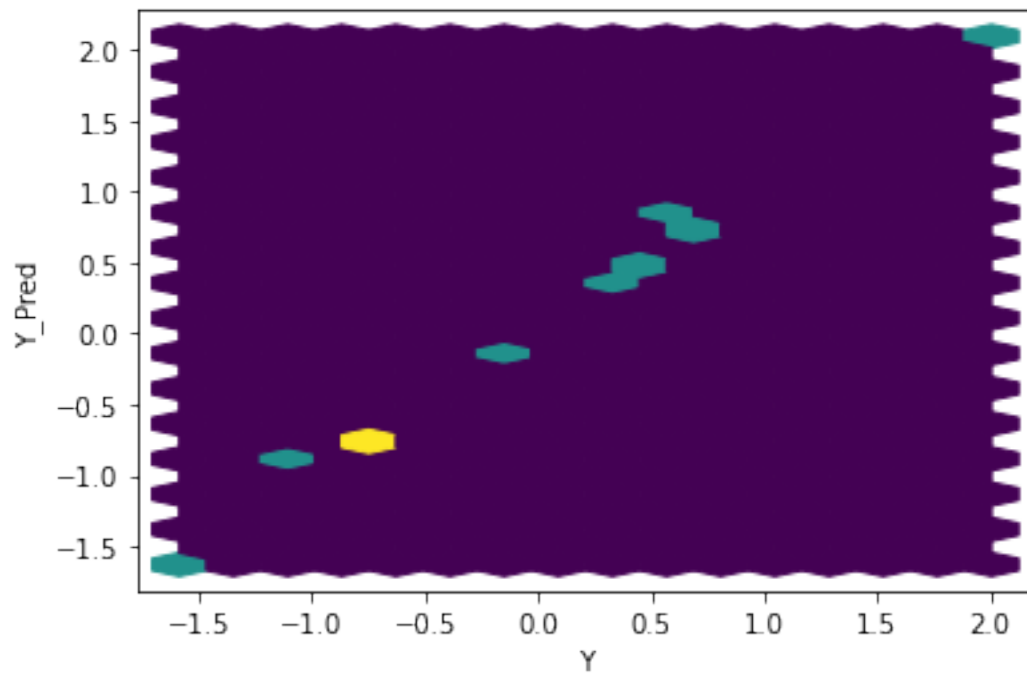
[18]: ABC_train_test.training_GAN(disc, gen,disc_opt,gen_opt,real_dataset,
      ↪batch_size, n_epochs,criterion,coeff,mean,variance,device)
```



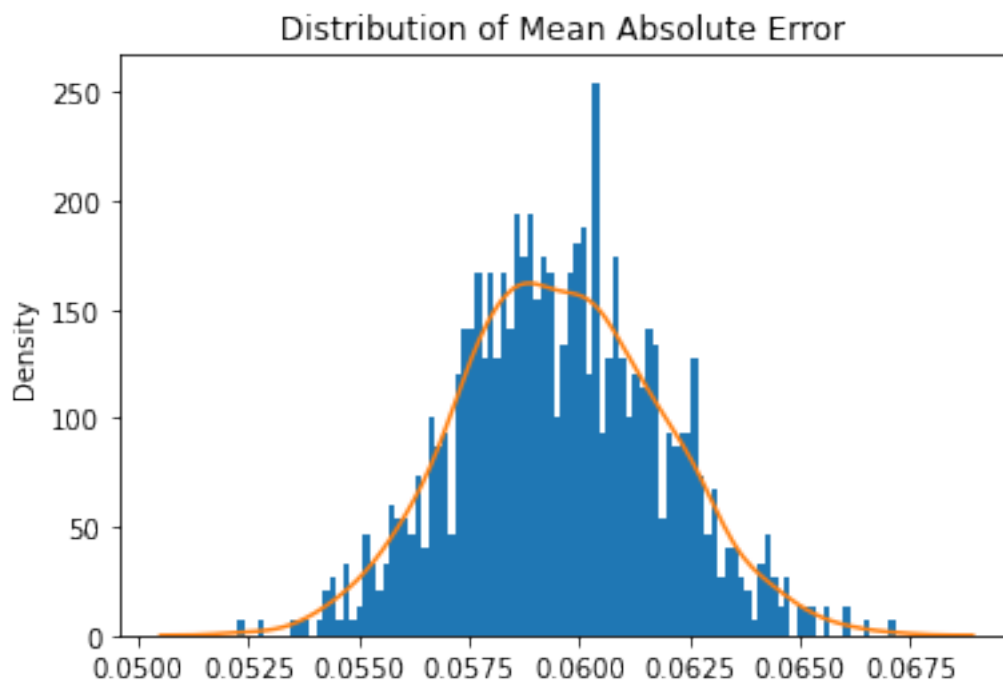
```
[19]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```





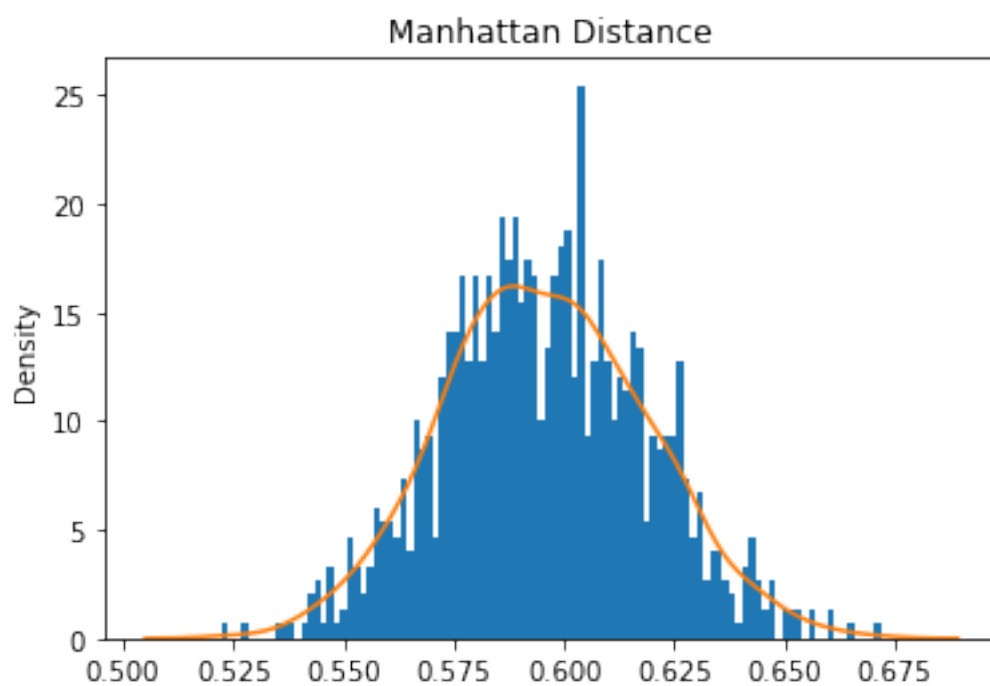


Mean Square Error: 0.00648691094163164

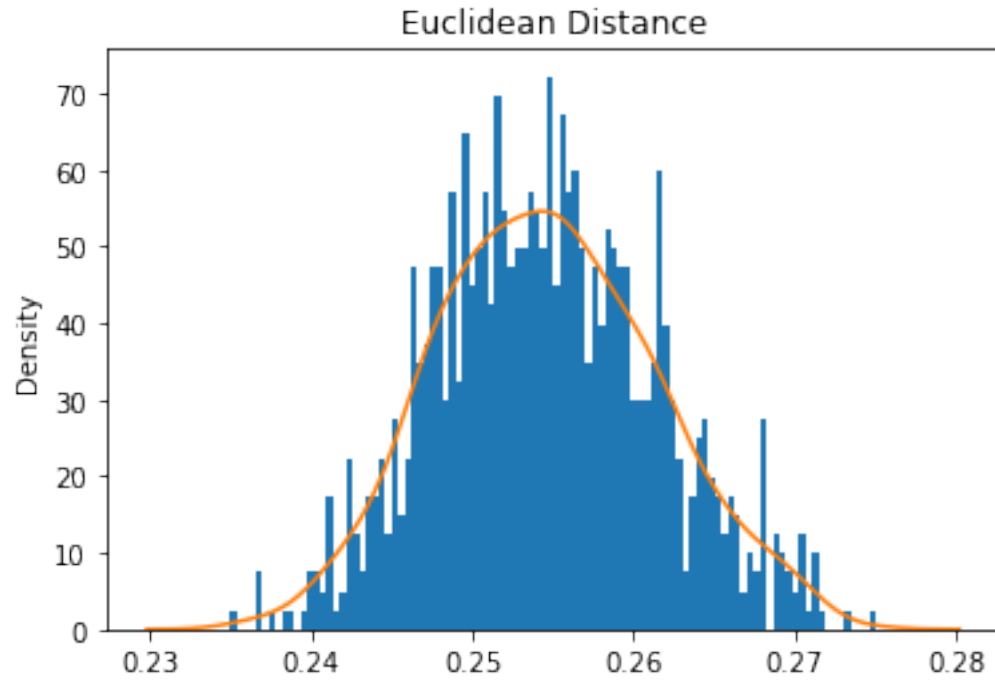


Mean Absolute Error: 0.05956322122812271

Mean Manhattan Distance: 0.5956322122812271

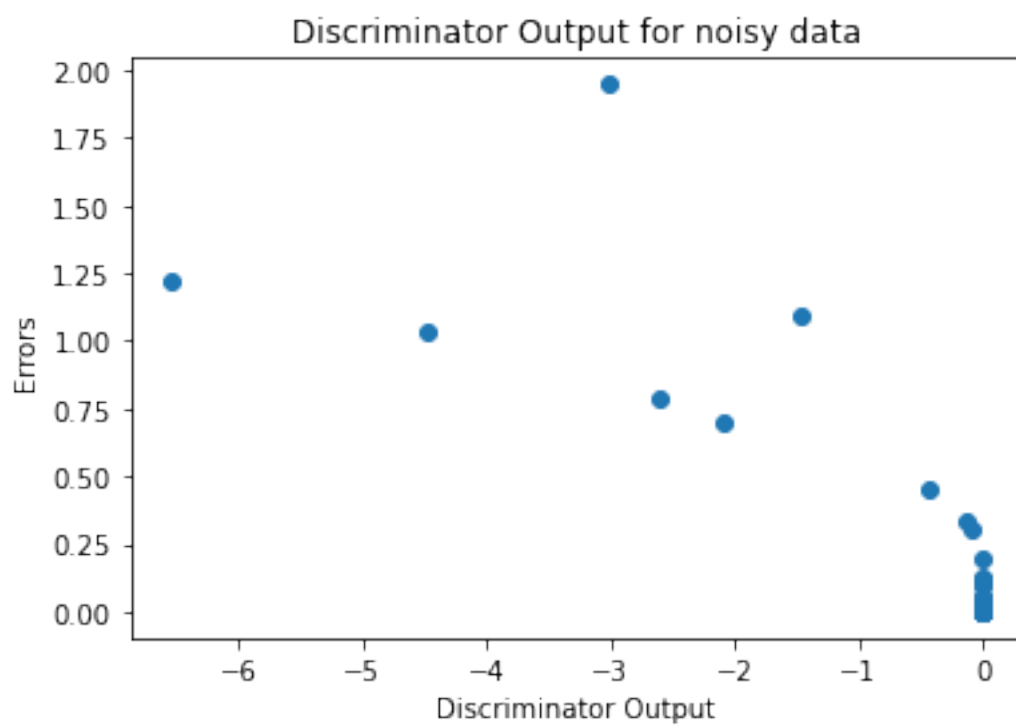
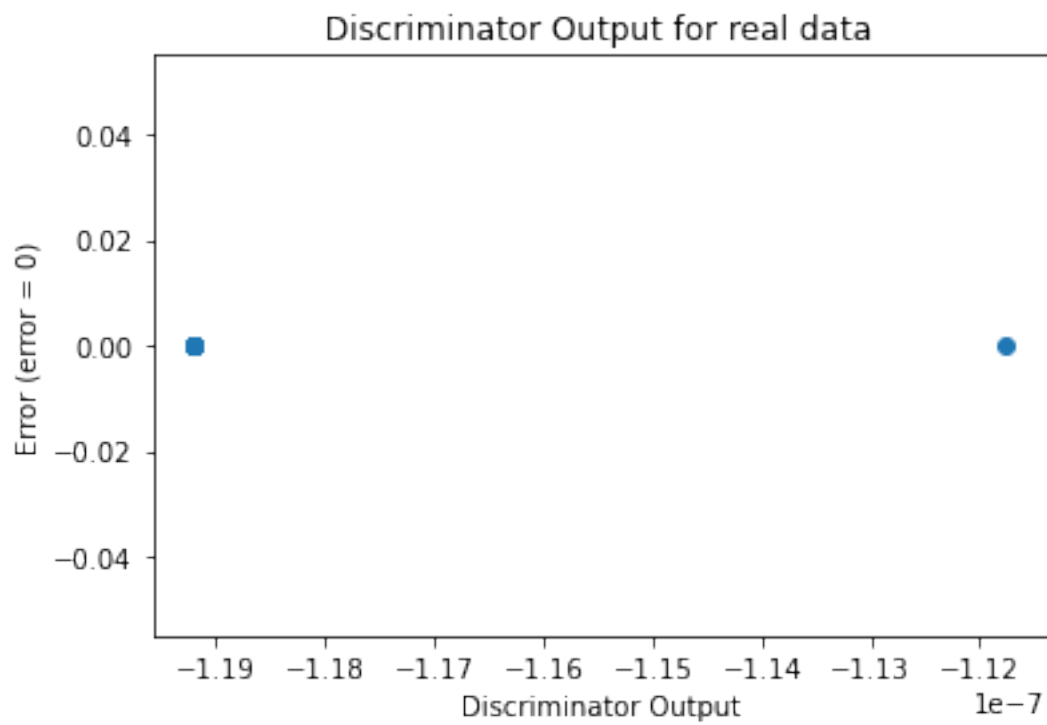


Mean Euclidean Distance: 0.25460236425964233



Sanity Checks

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



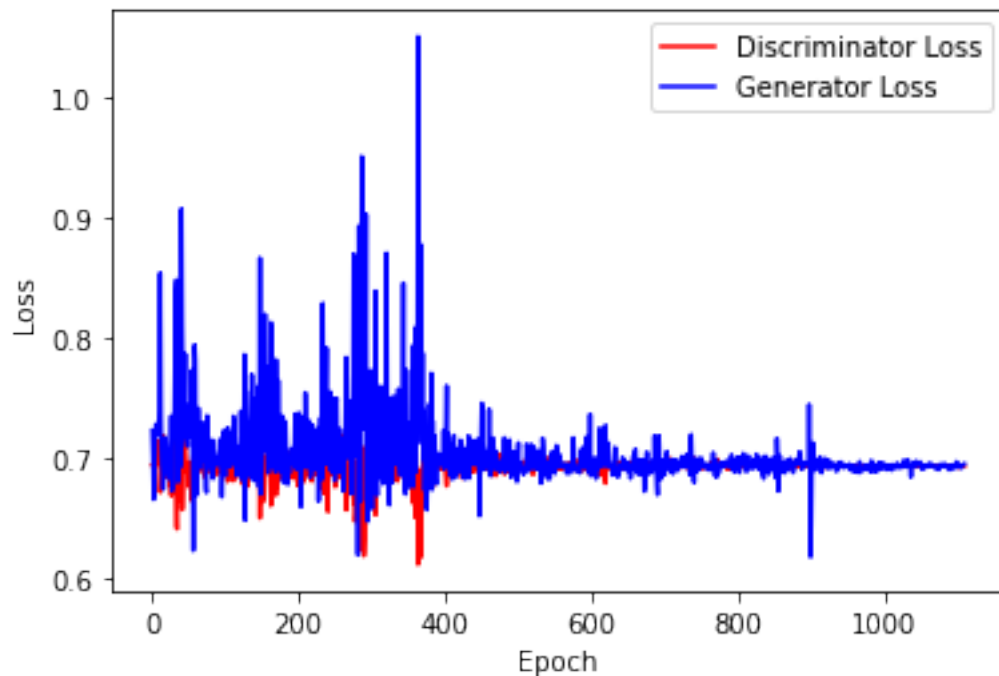
Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

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

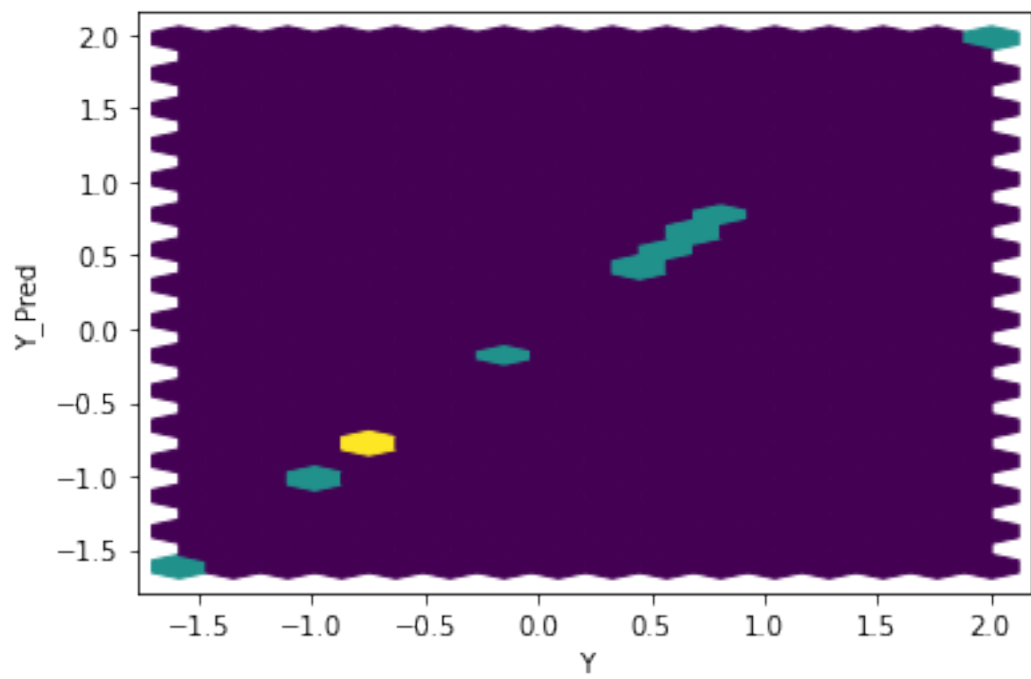
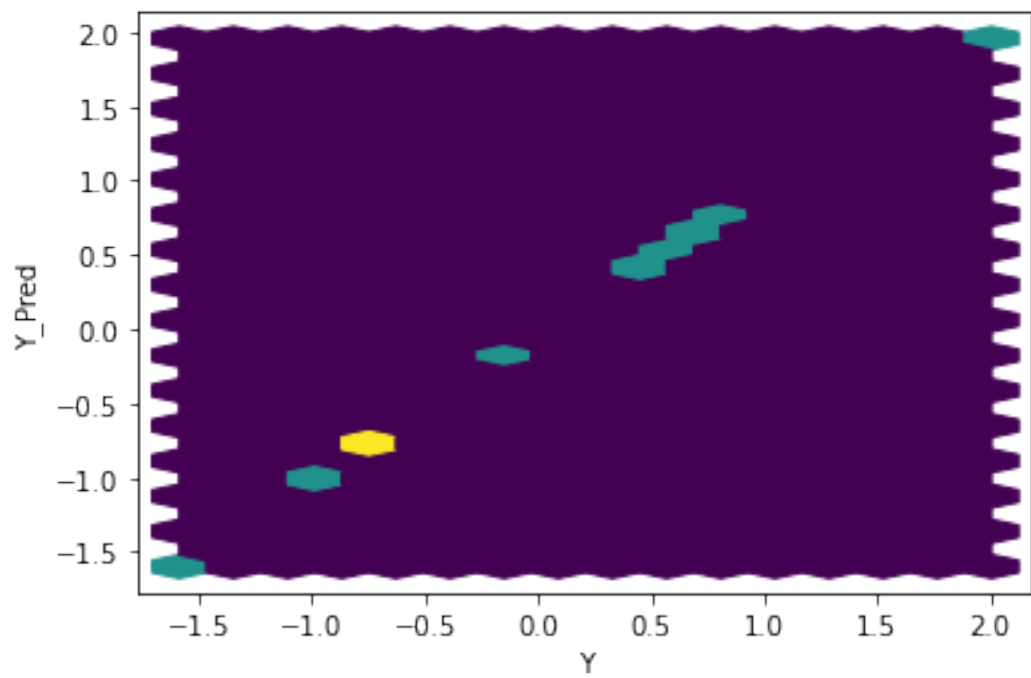
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))
```

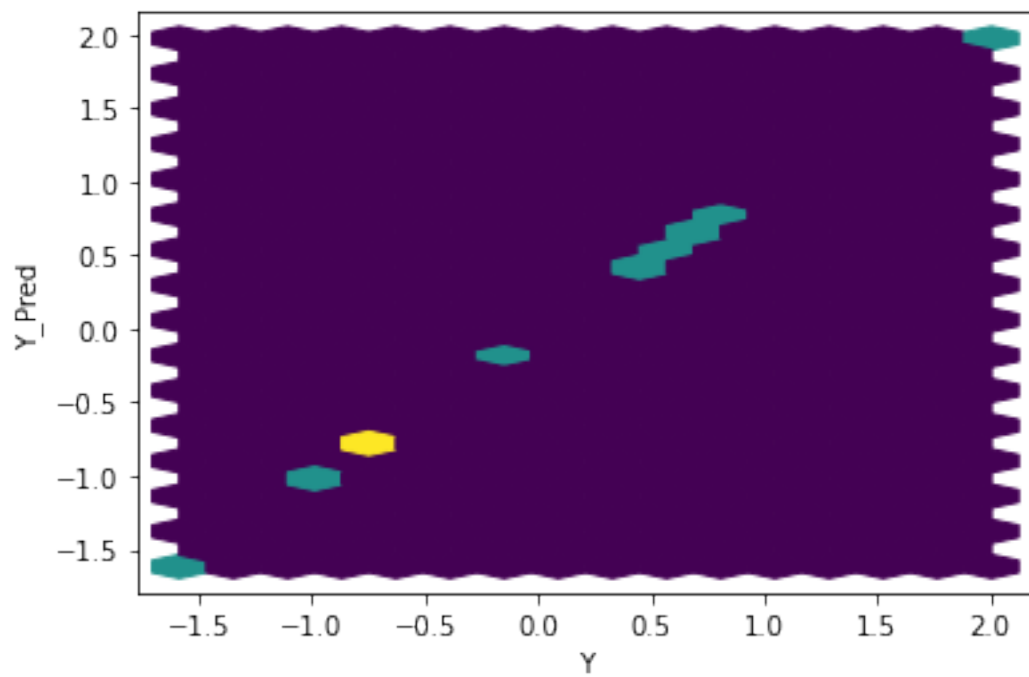
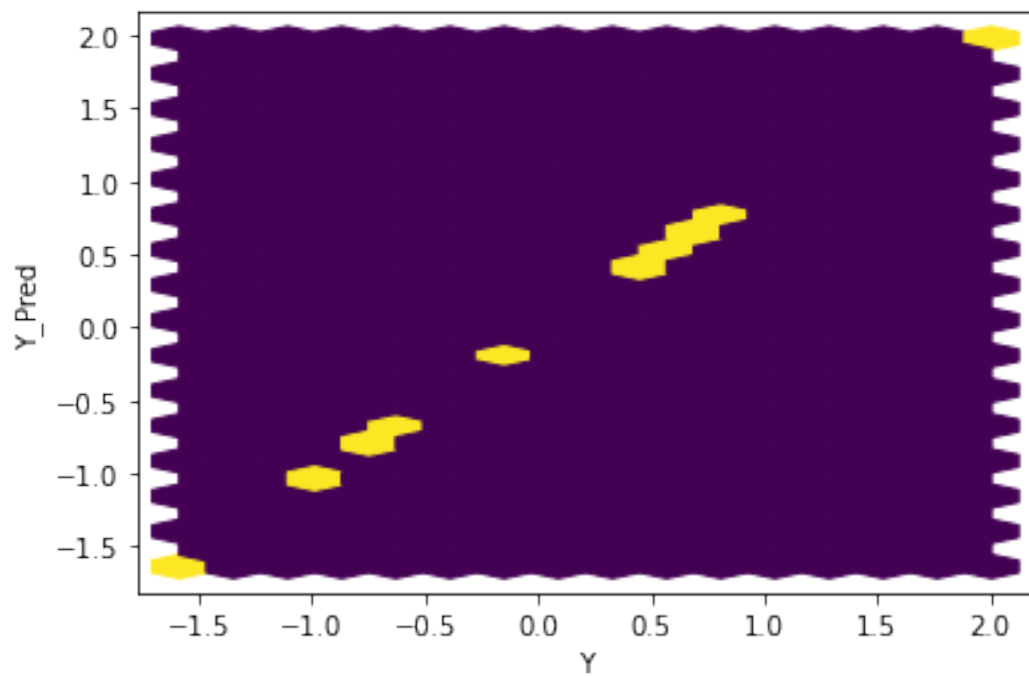
```
[22]: ABC_train_test.
      ↪ training_GAN_2(disc,gen,disc_opt,gen_opt,real_dataset,batch_size,
      ↪ error,criterion,coeff,mean,variance,device)
```

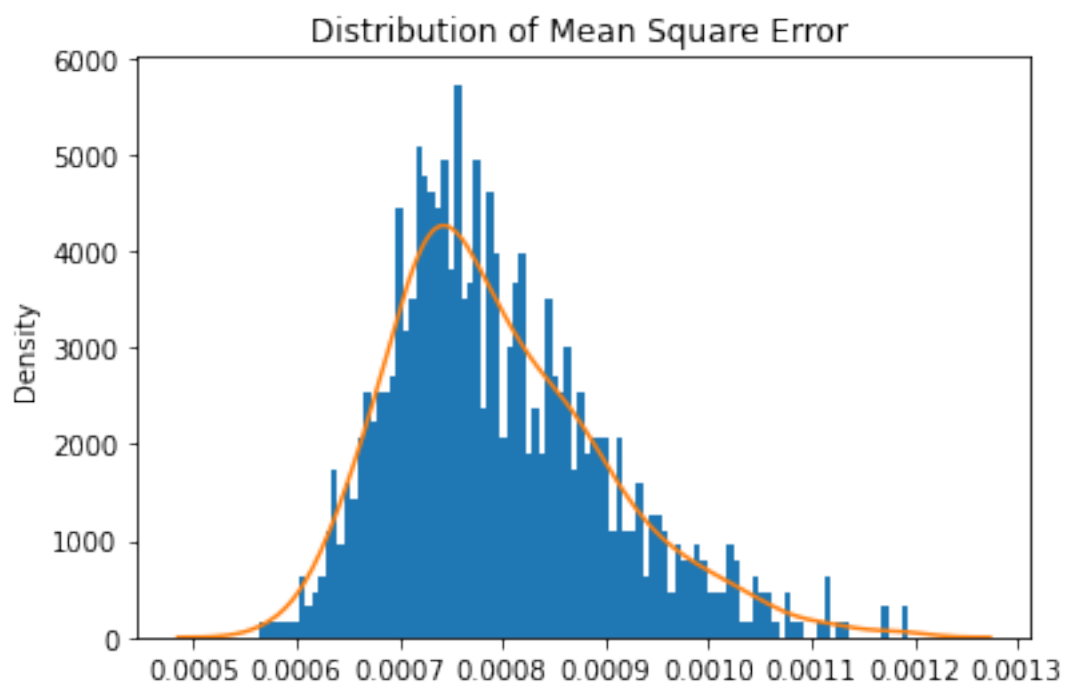
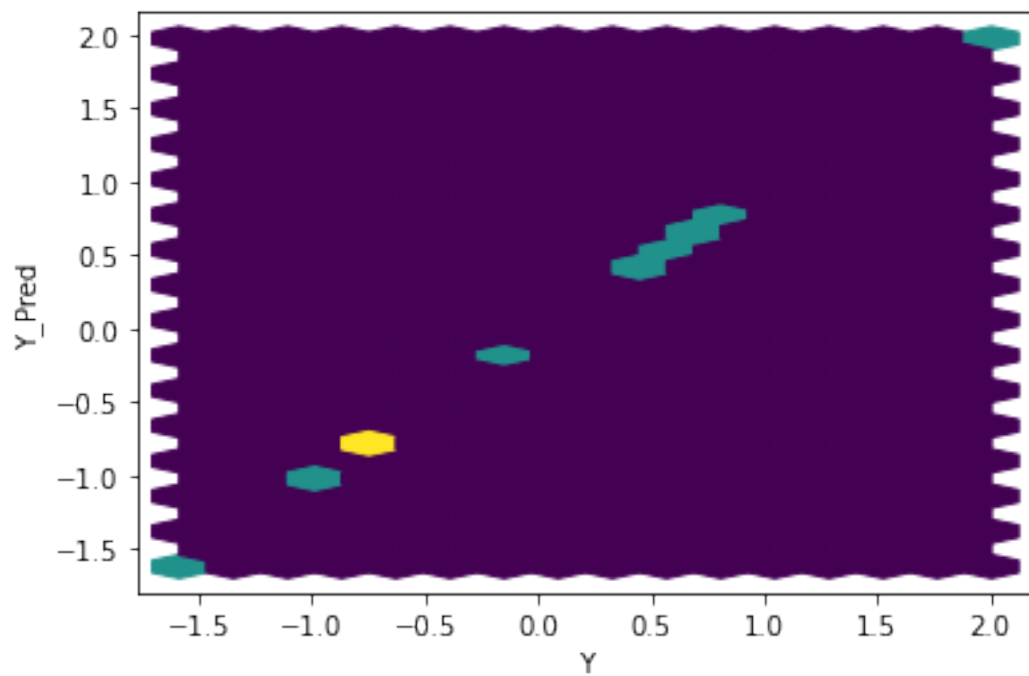
Number of epochs 555



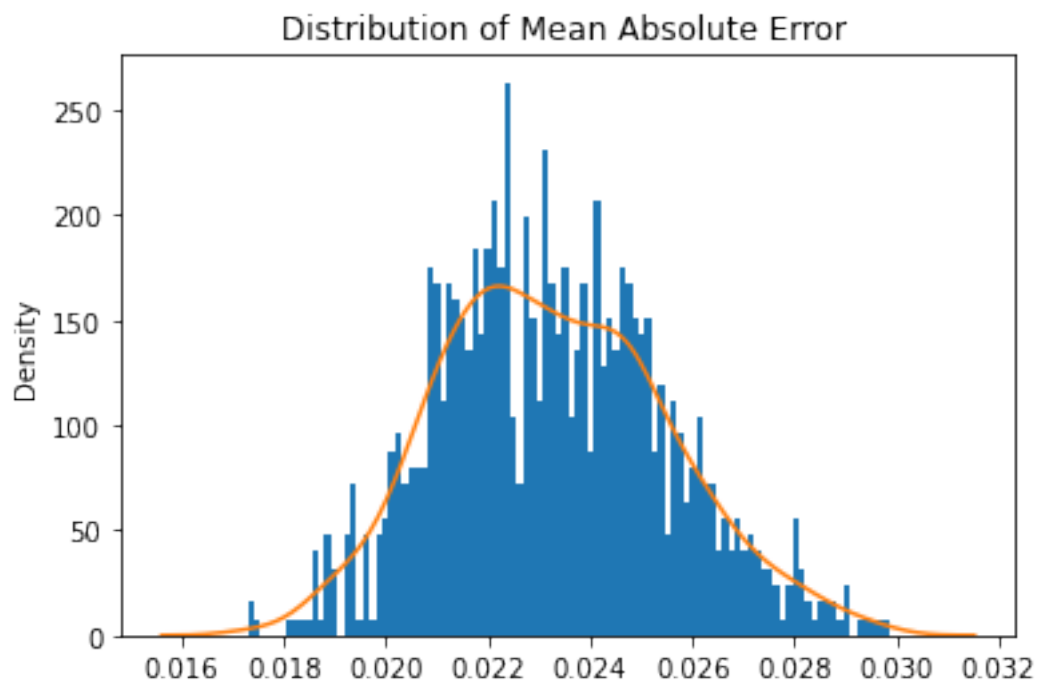
```
[23]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```





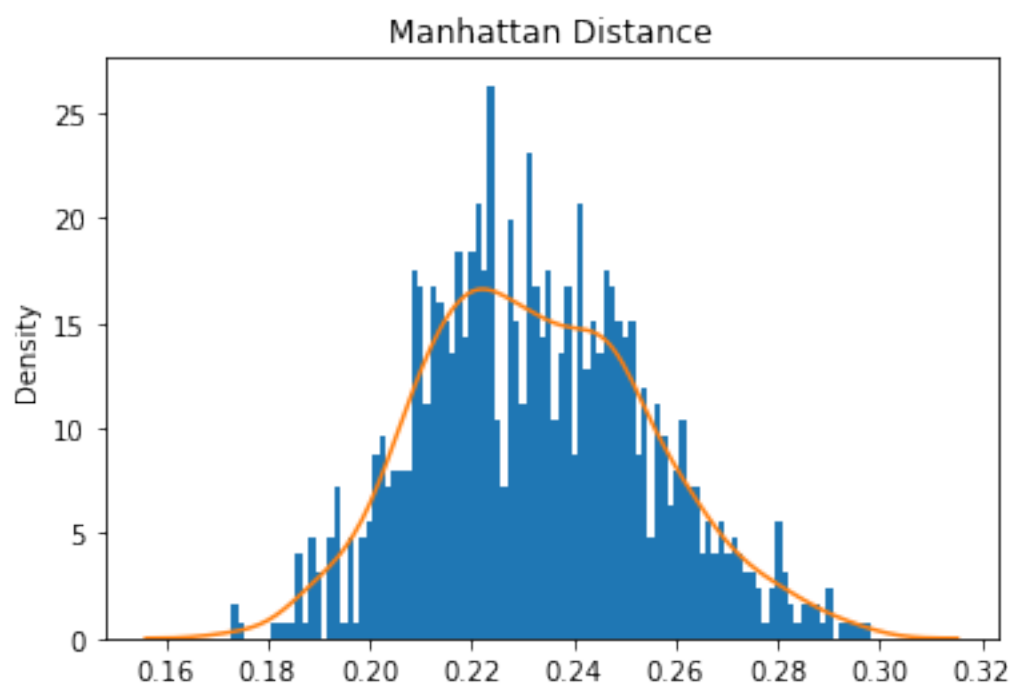


Mean Square Error: 0.000796056572088712

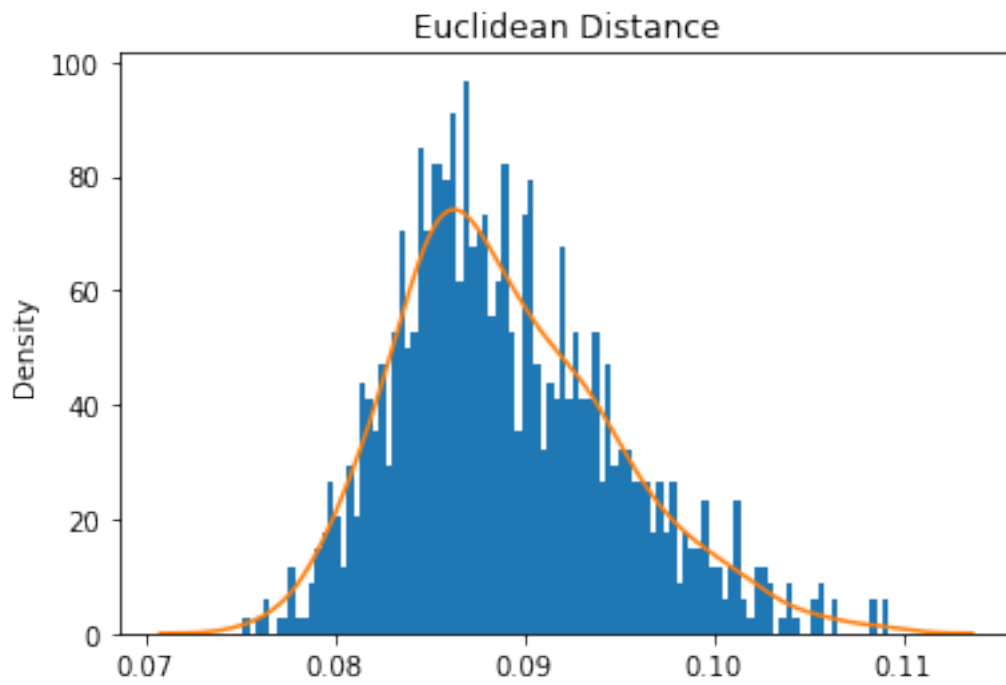


Mean Absolute Error: 0.02327692009806633

Mean Manhattan Distance: 0.2327692009806633



Mean Euclidean Distance: 0.08903167282273954



[]: