

Dataset2_Friedman1_output_5

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 = 0
     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	0.320567	0.407436	0.494248	0.887519	0.687441	0.374252	0.701353
1	0.870573	0.422217	0.096807	0.806529	0.116306	0.620213	0.162936
2	0.121992	0.959921	0.026261	0.960916	0.240006	0.965835	0.699085

```

3  0.472821  0.948381  0.896898  0.968297  0.123775  0.322727  0.624299
4  0.145426  0.315487  0.690114  0.103206  0.009039  0.500428  0.583773

```

```

          X7          X8          X9          Y
0  0.642739  0.971524  0.210137  16.443805
1  0.333939  0.039234  0.321735  21.249520
2  0.579075  0.164716  0.278566  18.907658
3  0.731854  0.738494  0.225107  23.201949
4  0.337531  0.525335  0.651208   3.343592

```

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:                20:05:04      Log-Likelihood:          325.96
No. Observations:                10      AIC:                   -631.9
Df Residuals:                    0      BIC:                   -628.9
Df Model:                        9
Covariance Type:                nonrobust
=====
              coef      std err          t      P>|t|      [0.025      0.975]
-----
const          2.22e-16         inf         0         nan         nan         nan
x1              0.3402         inf         0         nan         nan         nan
x2              0.8040         inf         0         nan         nan         nan
x3              0.0220         inf         0         nan         nan         nan
x4              0.5338         inf         0         nan         nan         nan
x5              0.3440         inf         0         nan         nan         nan
x6             -0.1498         inf        -0         nan         nan         nan
x7             -0.3580         inf        -0         nan         nan         nan
x8             -0.2925         inf        -0         nan         nan         nan
x9              0.0163         inf         0         nan         nan         nan
x10            -0.1335         inf        -0         nan         nan         nan
=====
Omnibus:                3.588      Durbin-Watson:           2.073
Prob(Omnibus):          0.166      Jarque-Bera (JB):        1.116
Skew:                  -0.245      Prob(JB):                0.572
Kurtosis:              1.439      Cond. No.                25.0
=====

```

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 2.220446e-16

x1 3.401771e-01

x2 8.040467e-01

x3 2.195216e-02

x4 5.337661e-01

x5 3.440083e-01

x6 -1.498168e-01

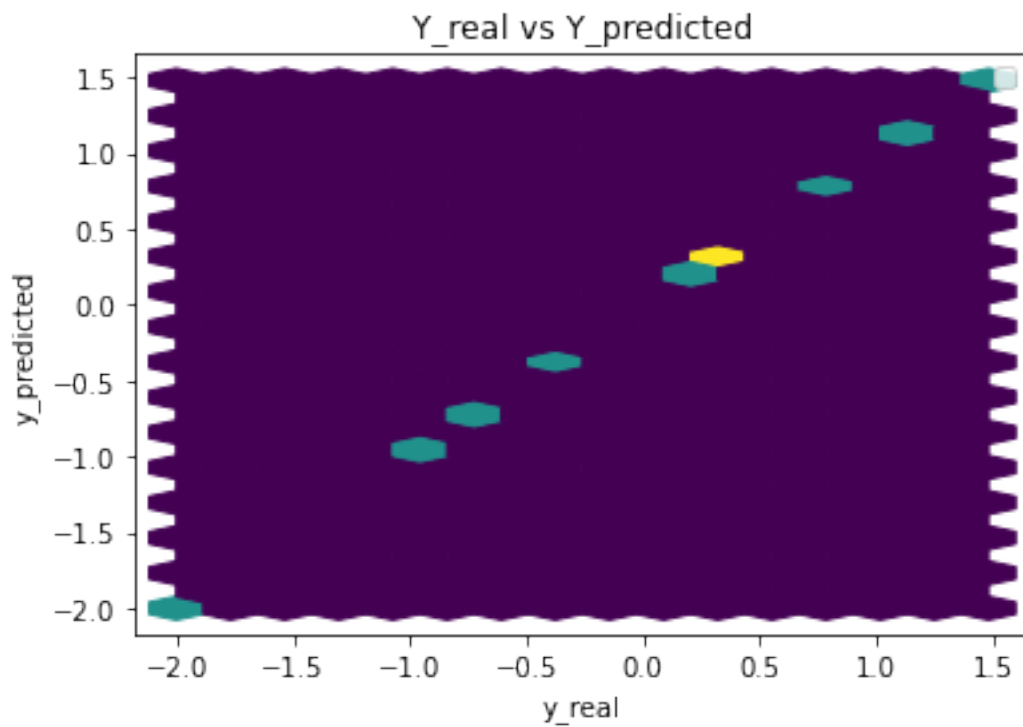
x7 -3.579580e-01

x8 -2.925108e-01

x9 1.631096e-02

x10 -1.335153e-01

dtype: float64



Performance Metrics

Mean Squared Error: 2.8482192761553952e-30

Mean Absolute Error: 1.4932499681208355e-15

Manhattan distance: 1.4932499681208355e-14

Euclidean distance: 5.336871064730153e-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

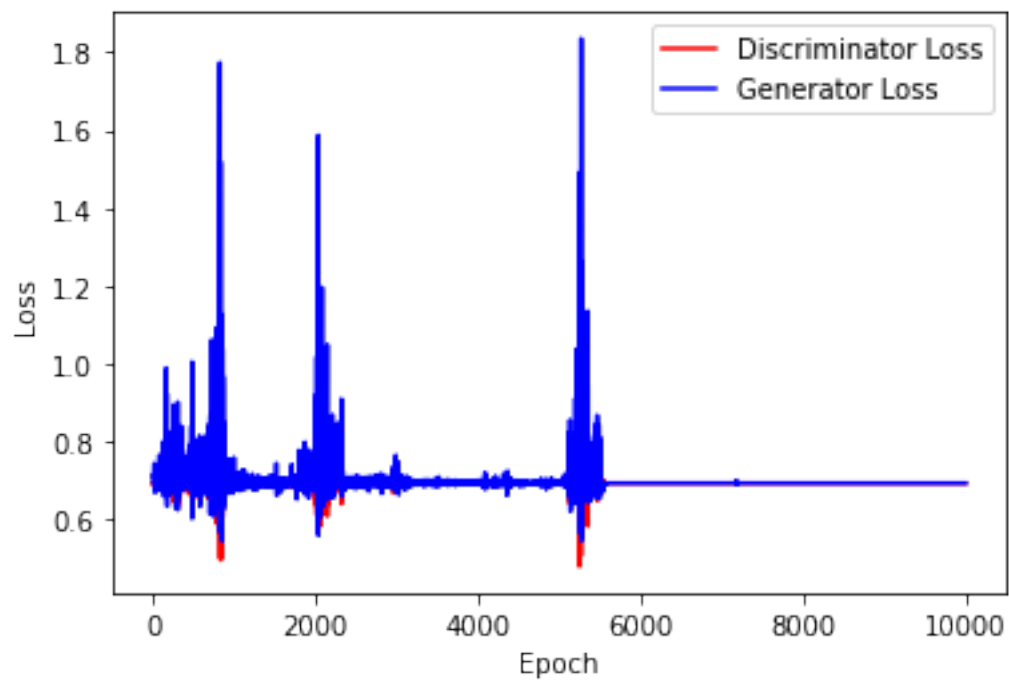
```
[9]: generator = network.Generator(n_features+2)
     discriminator = network.Discriminator(n_features+2)

     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))
```

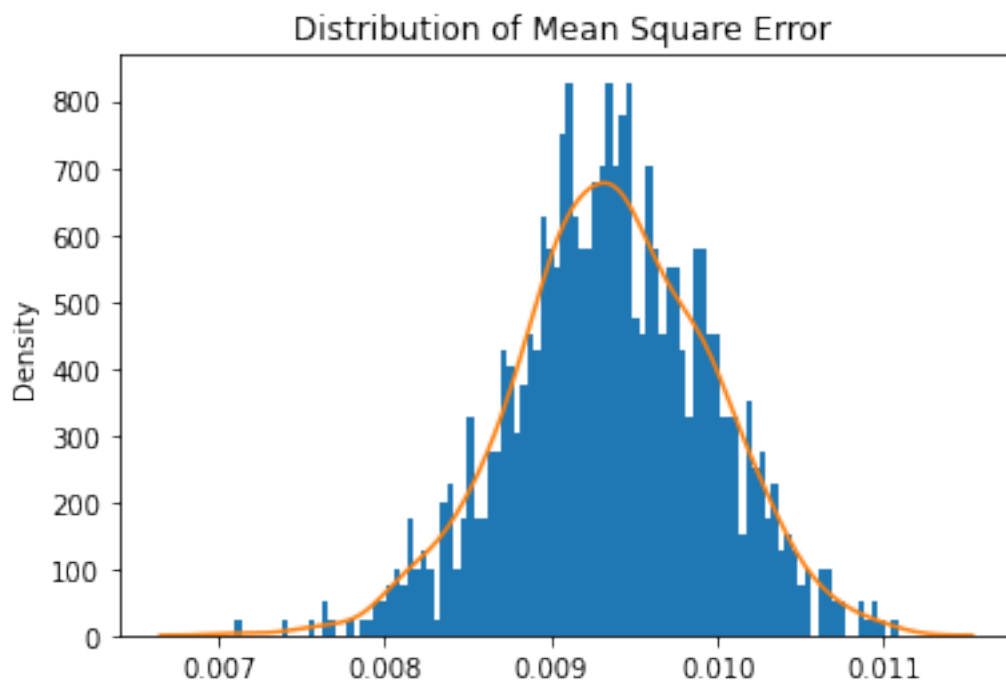
```
[10]: print(generator)
      print(discriminator)
```

```
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()
)
```

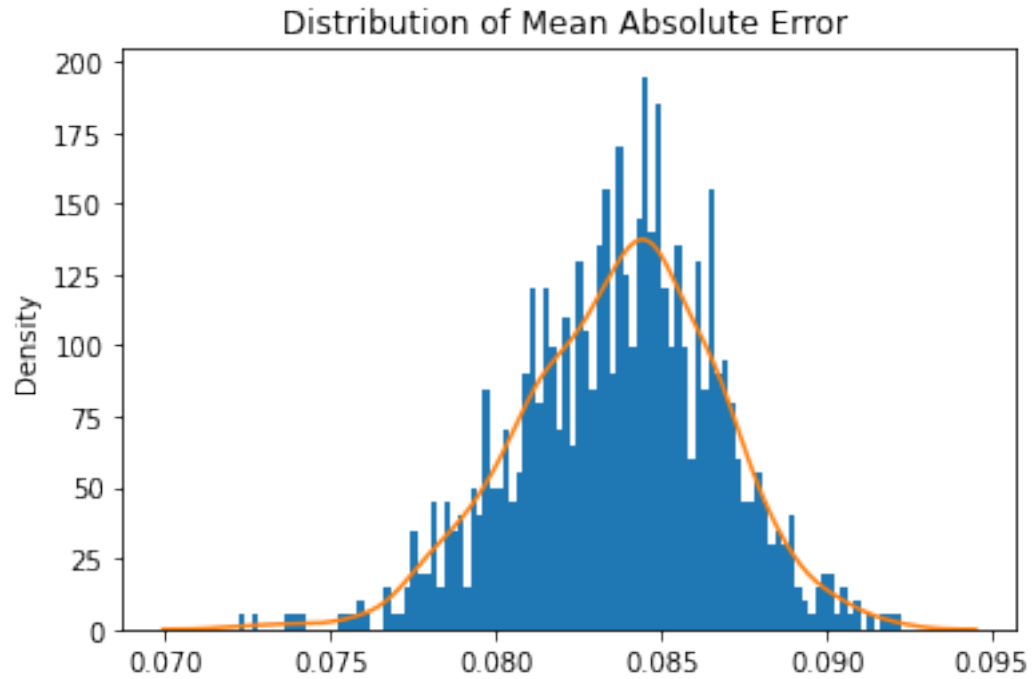
```
[11]: train_test.
     ↪training_GAN(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,
     ↪n_epochs,criterion,device)
```



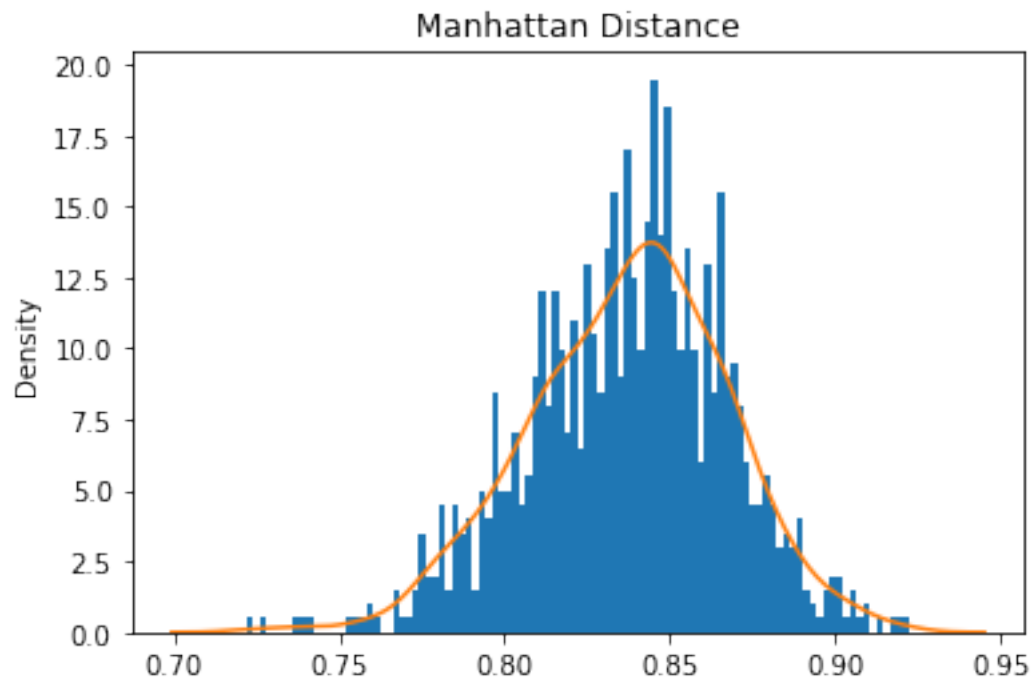
```
[12]: train_test.test_generator(generator,real_dataset,device)
```



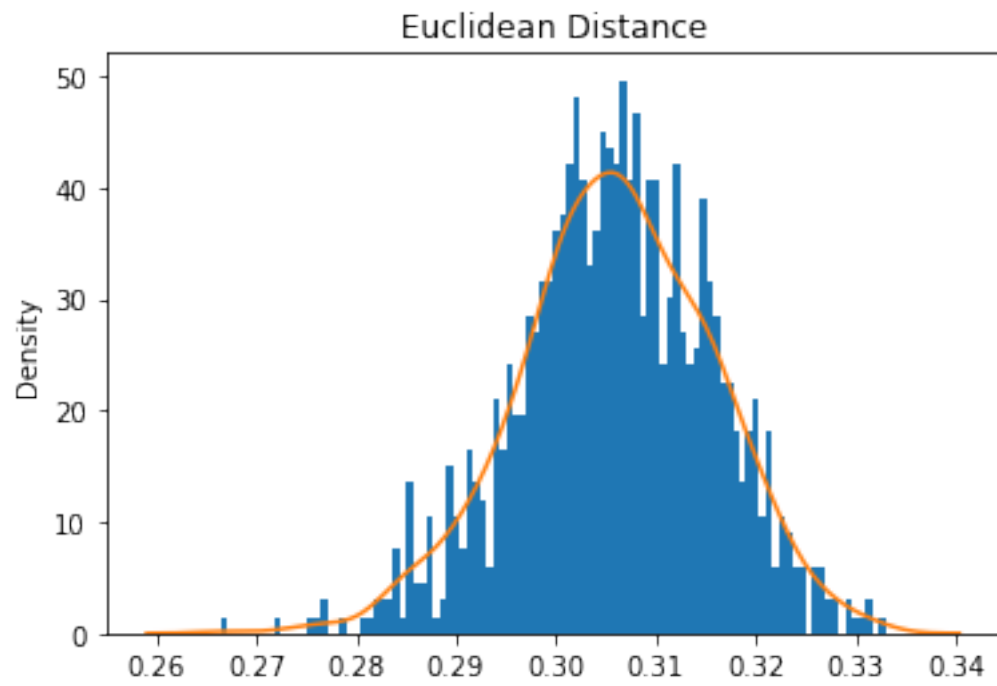
Mean Square Error: 0.009363350902948004



Mean Absolute Error: 0.08371758110150695

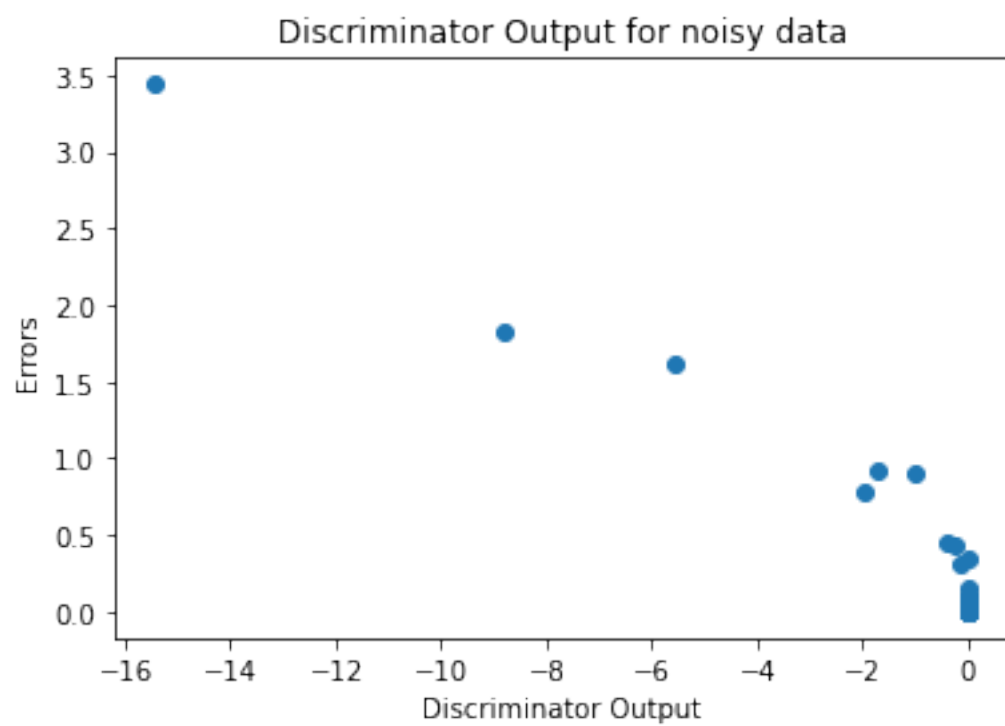
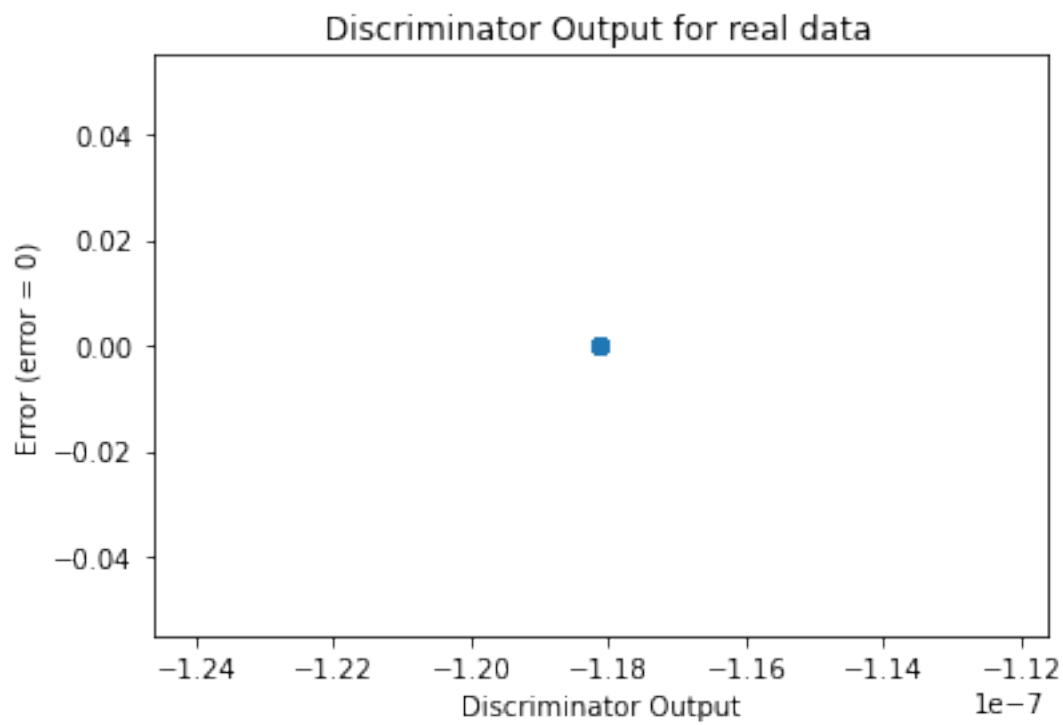


Mean Manhattan Distance: 0.8371758110150694



Mean Euclidean Distance: 0.30583998132636514

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

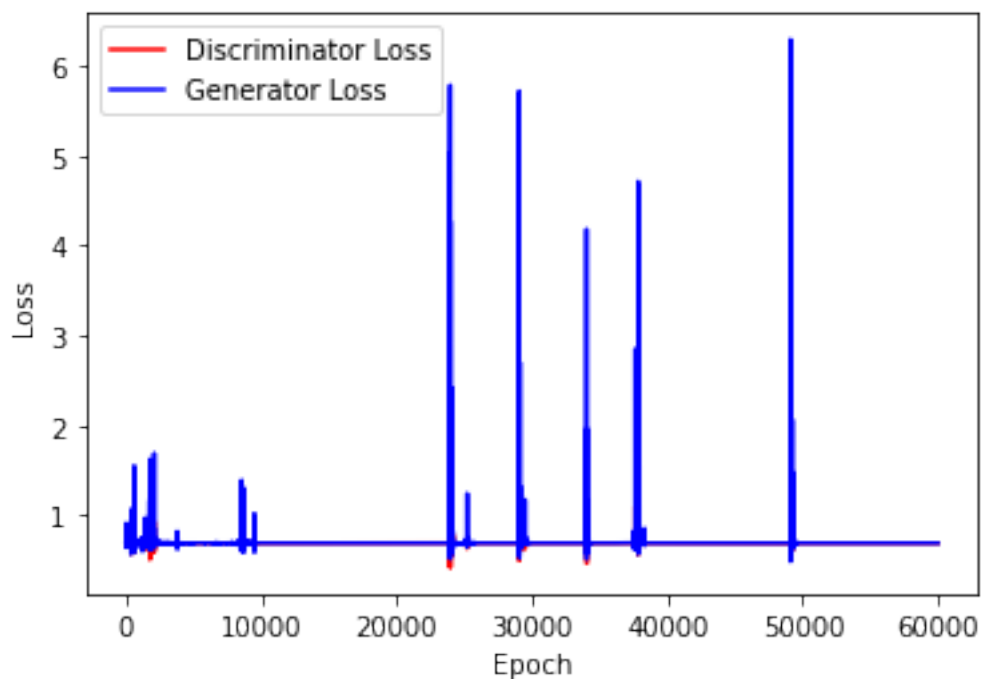



Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

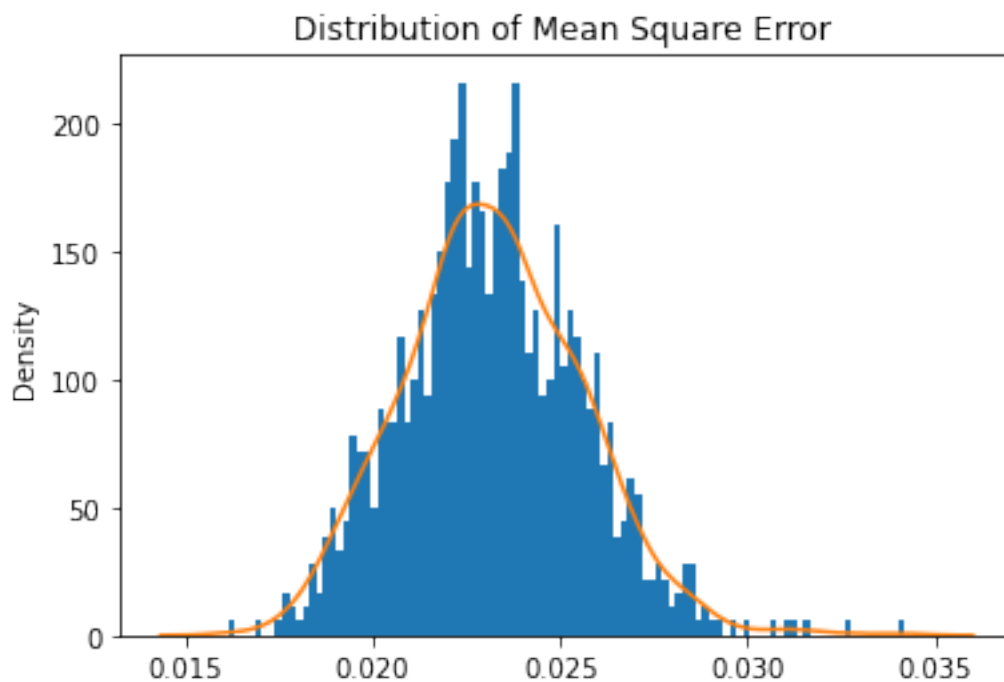
```
[14]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator.parameters(), lr=0.01, betas=(0.5,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

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

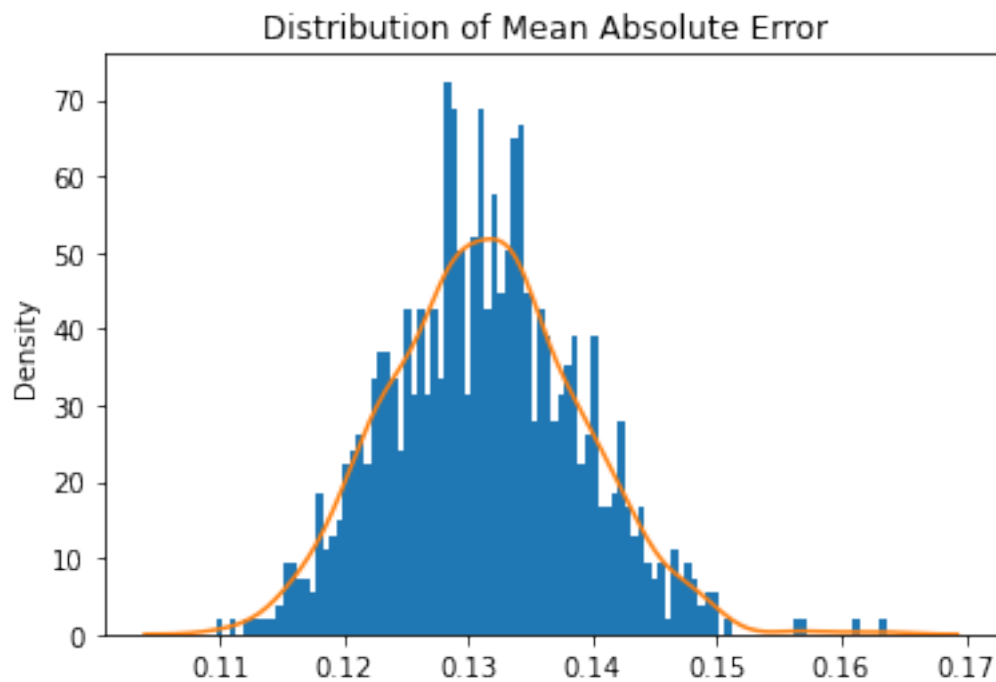
Number of epochs needed 30000



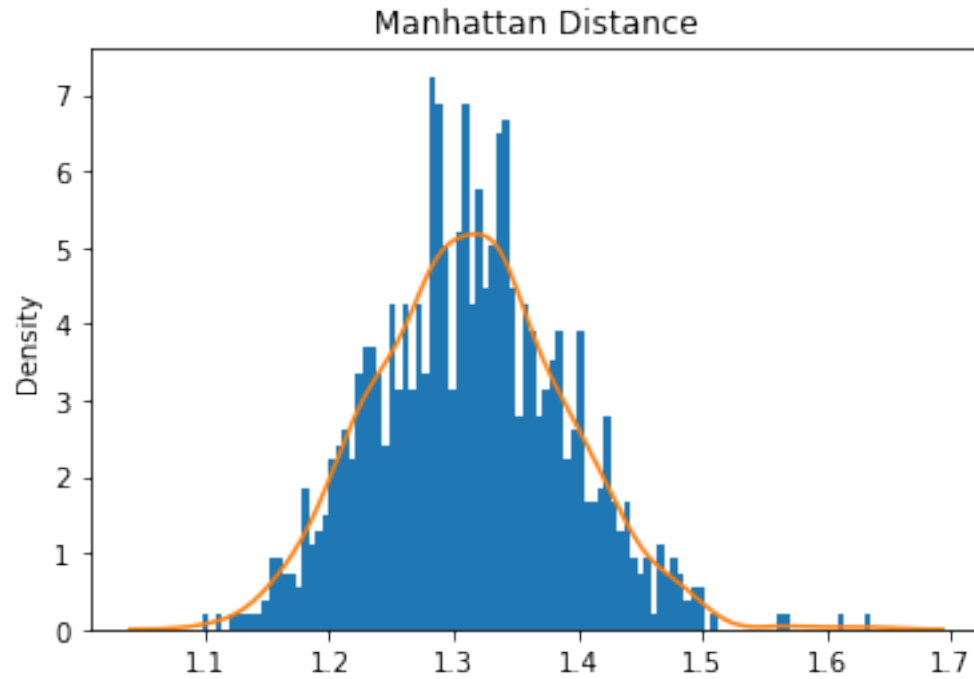
```
[16]: train_test.test_generator(generator,real_dataset,device)
```



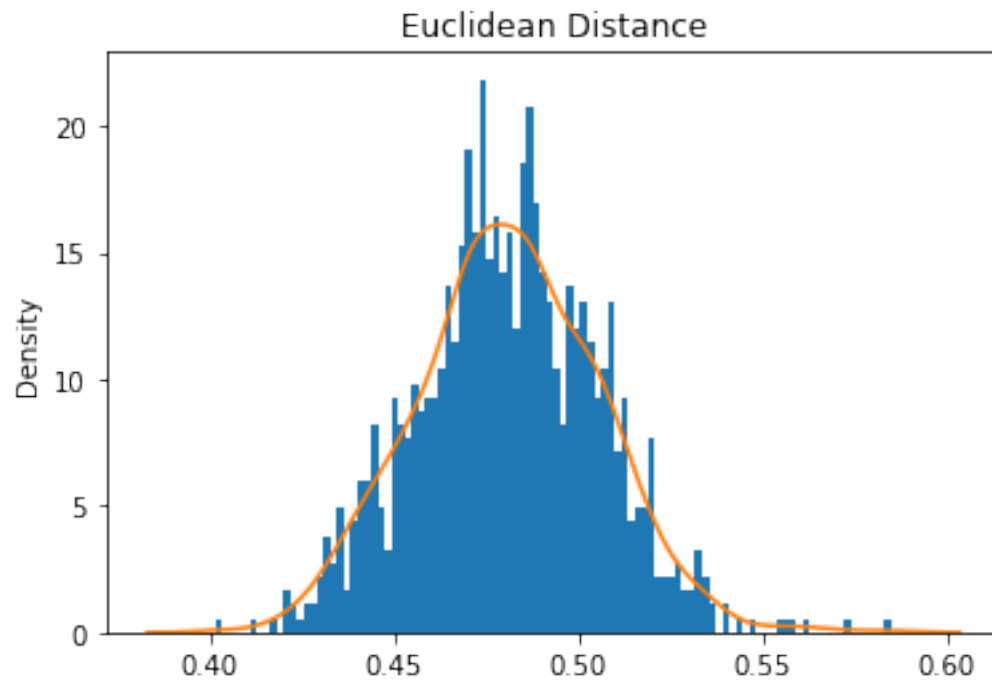
Mean Square Error: 0.023195370538341045



Mean Absolute Error: 0.1313441815651953



Mean Manhattan Distance: 1.3134418156519532



Mean Euclidean Distance: 0.48097909851401455

2 ABC GAN Model

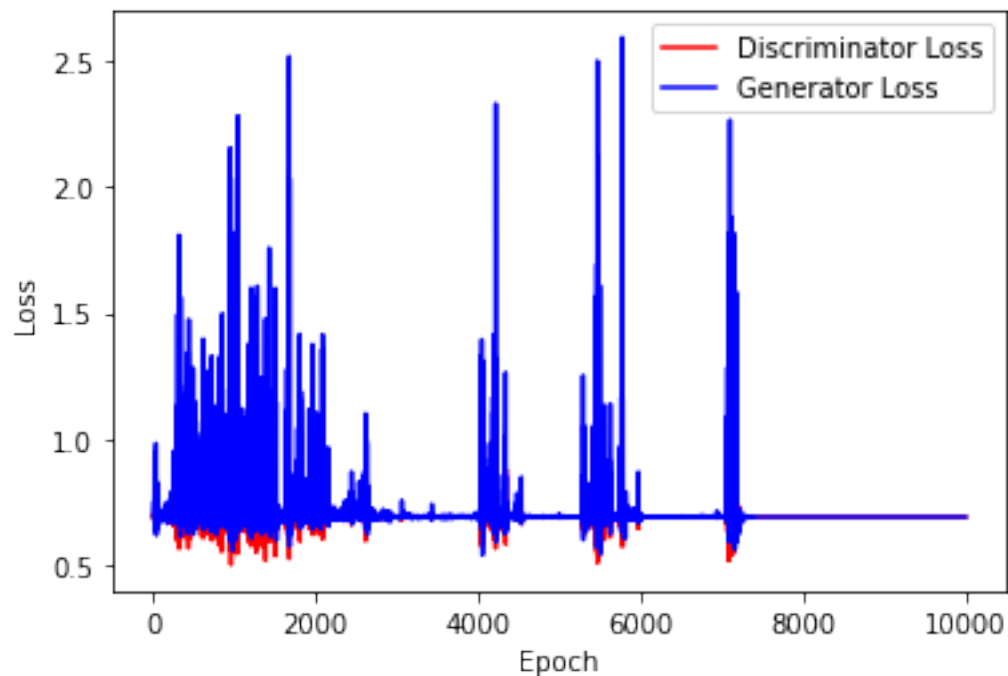
2.0.1 Training the network

Training ABC-GAN for `n_epochs` number of epochs

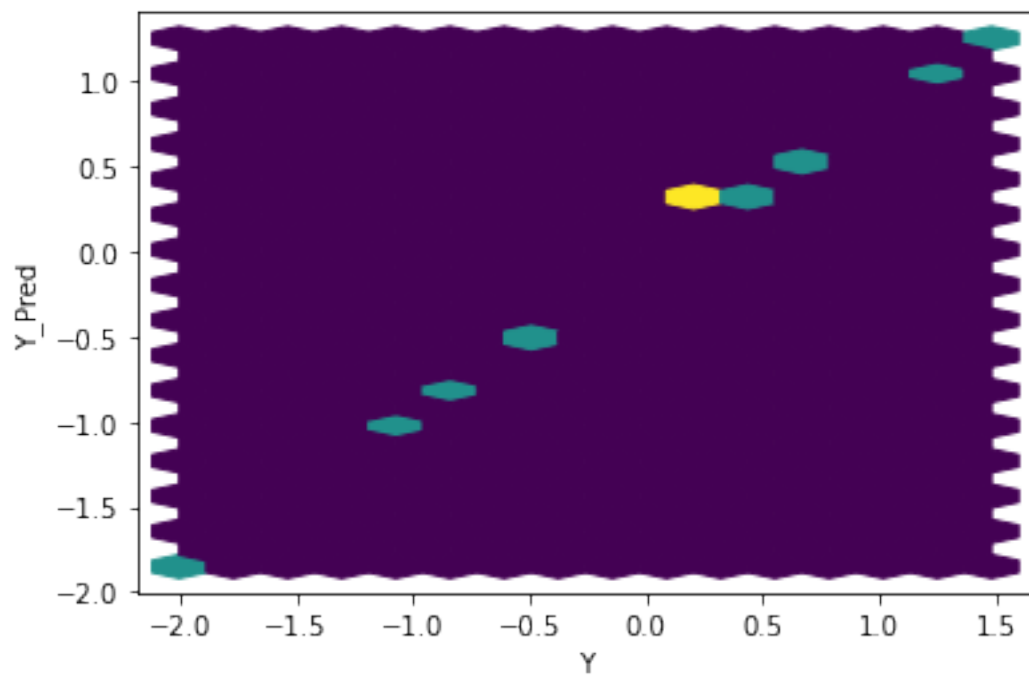
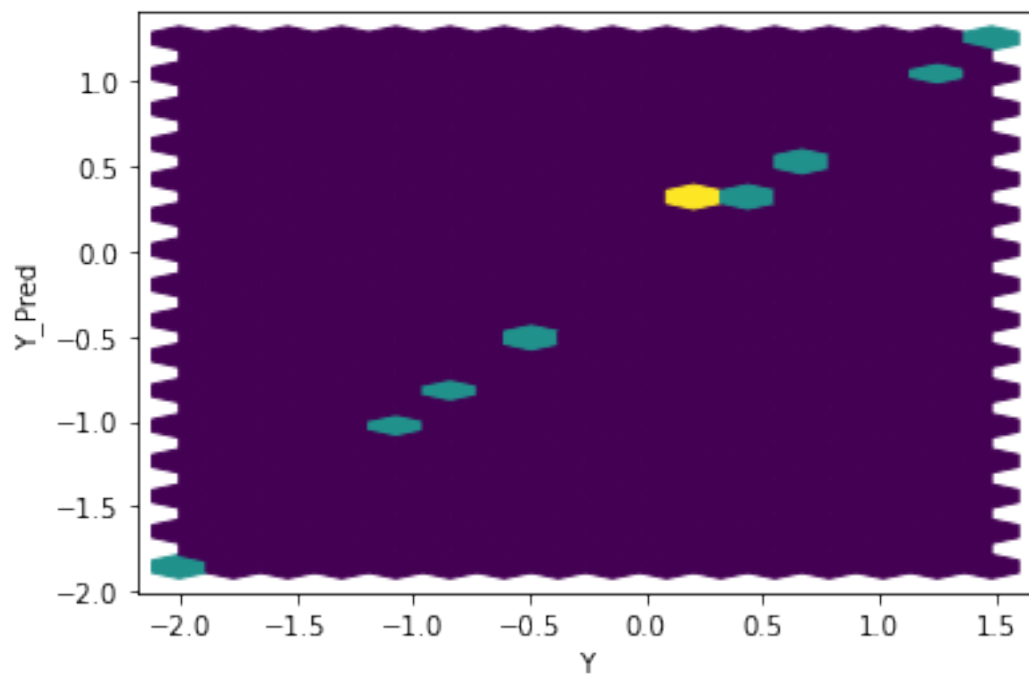
```
[17]: gen = network.Generator(n_features+2)
      disc = network.Discriminator(n_features+2)

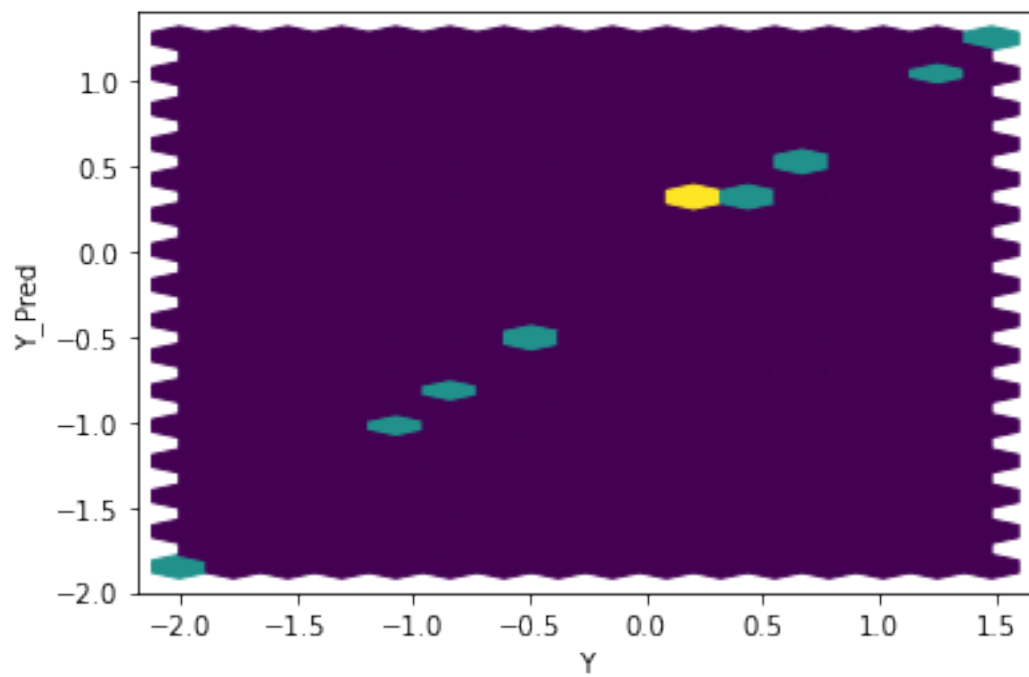
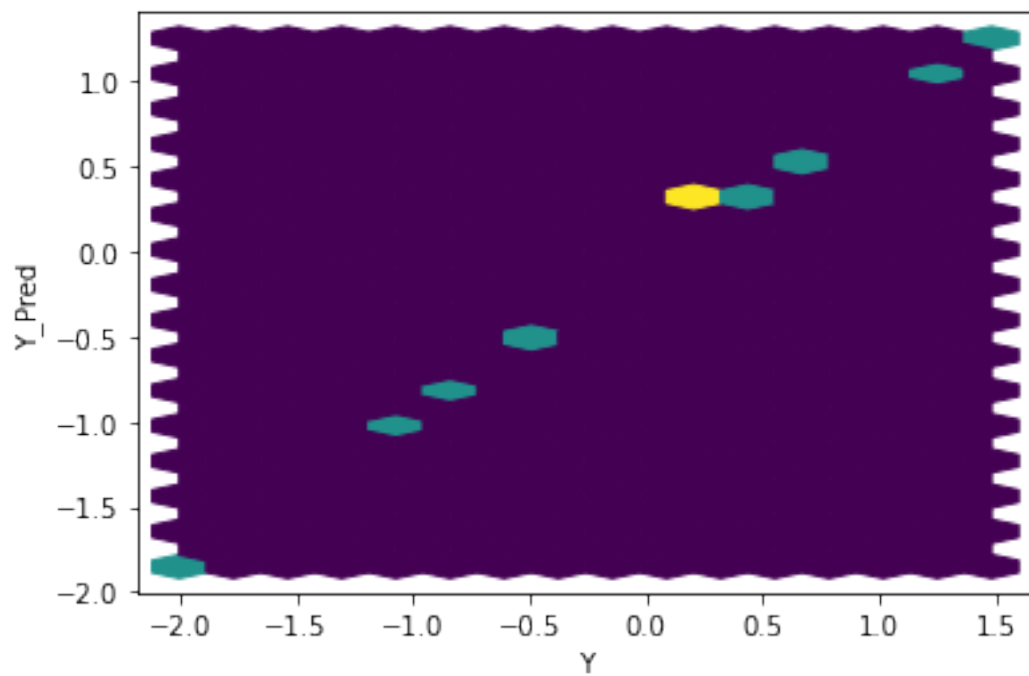
      criterion = torch.nn.BCEWithLogitsLoss()
      gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
      disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))

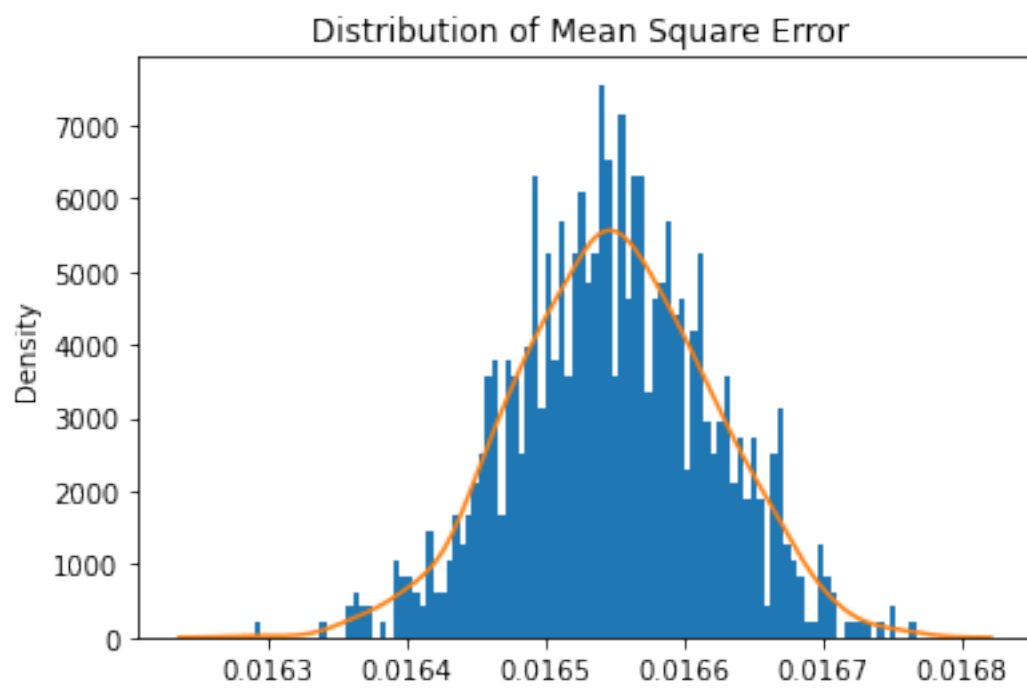
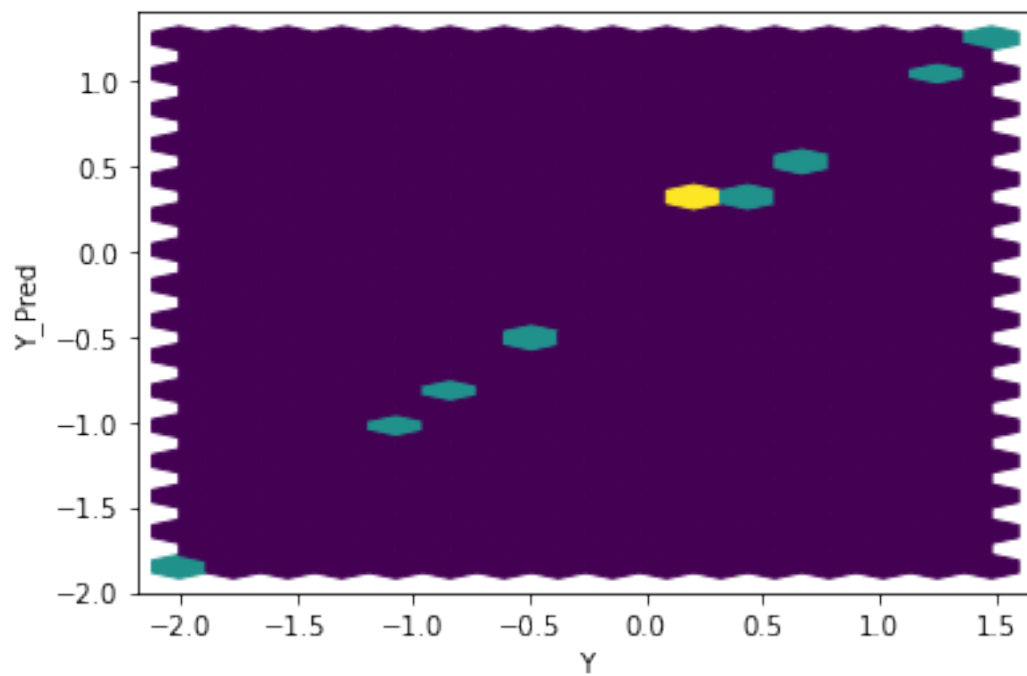
[18]: ABC_train_test.training_GAN(disc, gen,disc_opt,gen_opt,real_dataset,
      ↪batch_size, n_epochs,criterion,coeff,mean,variance,device)
```



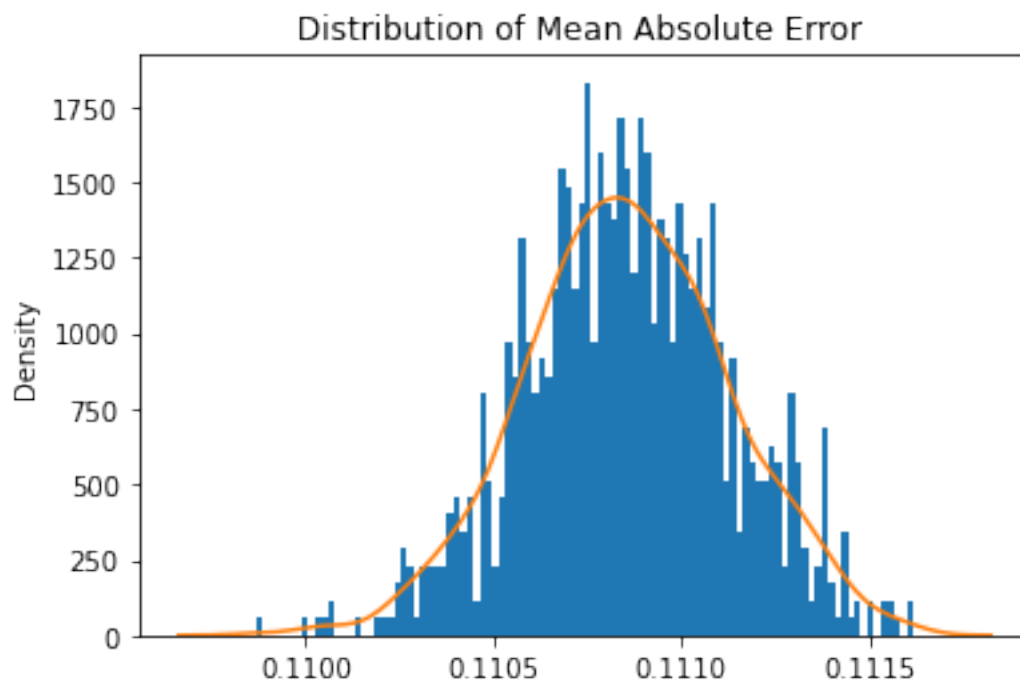
```
[19]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```





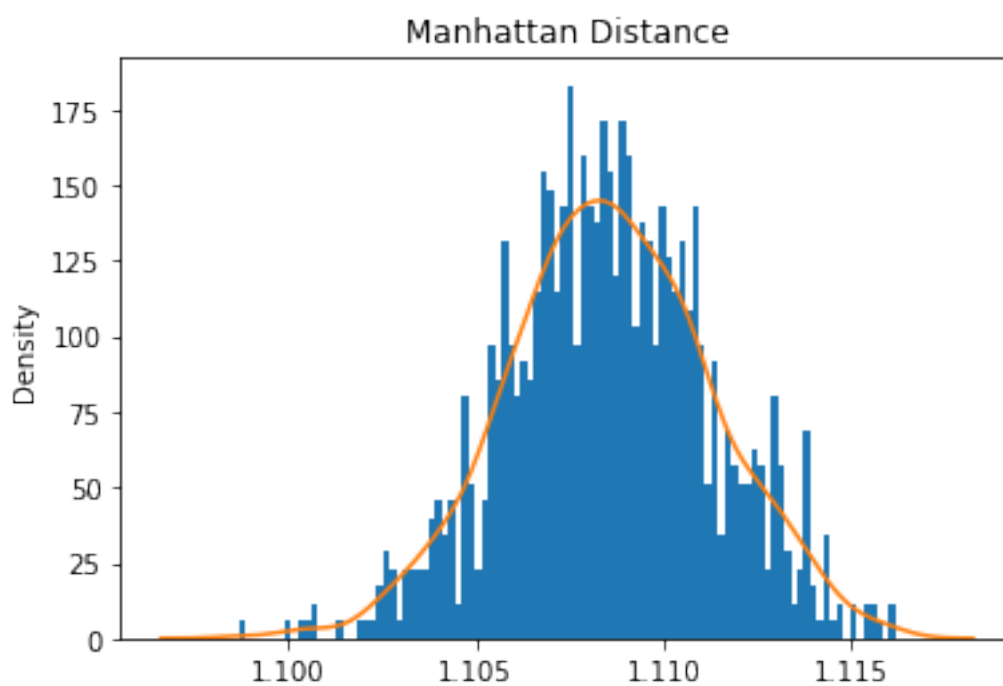


Mean Square Error: 0.016549121208843404

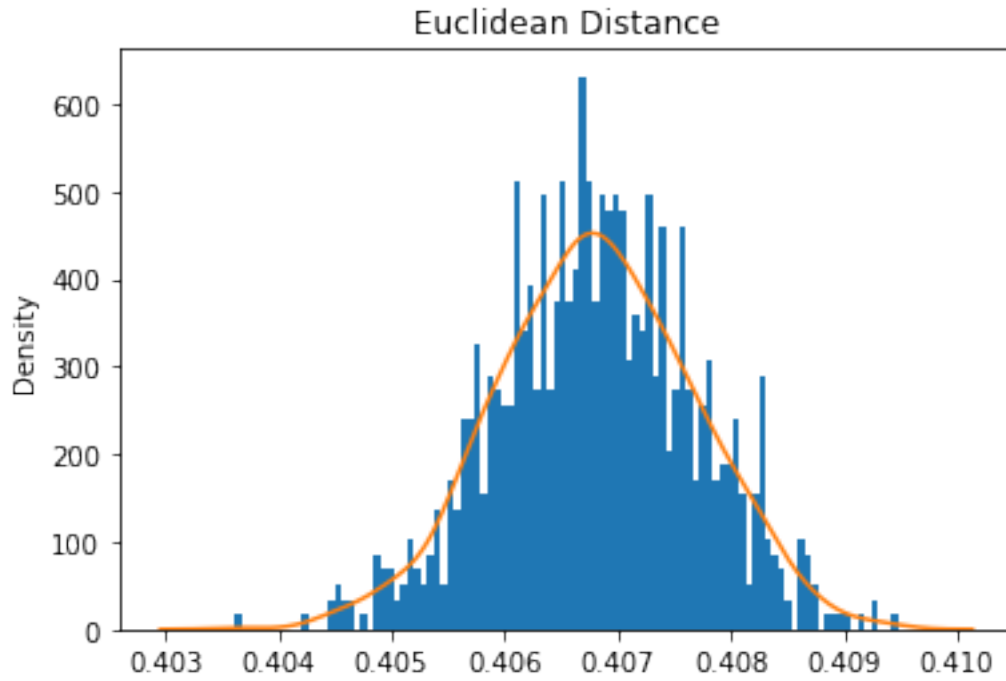


Mean Absolute Error: 0.1108549439251423

Mean Manhattan Distance: 1.1085494392514228

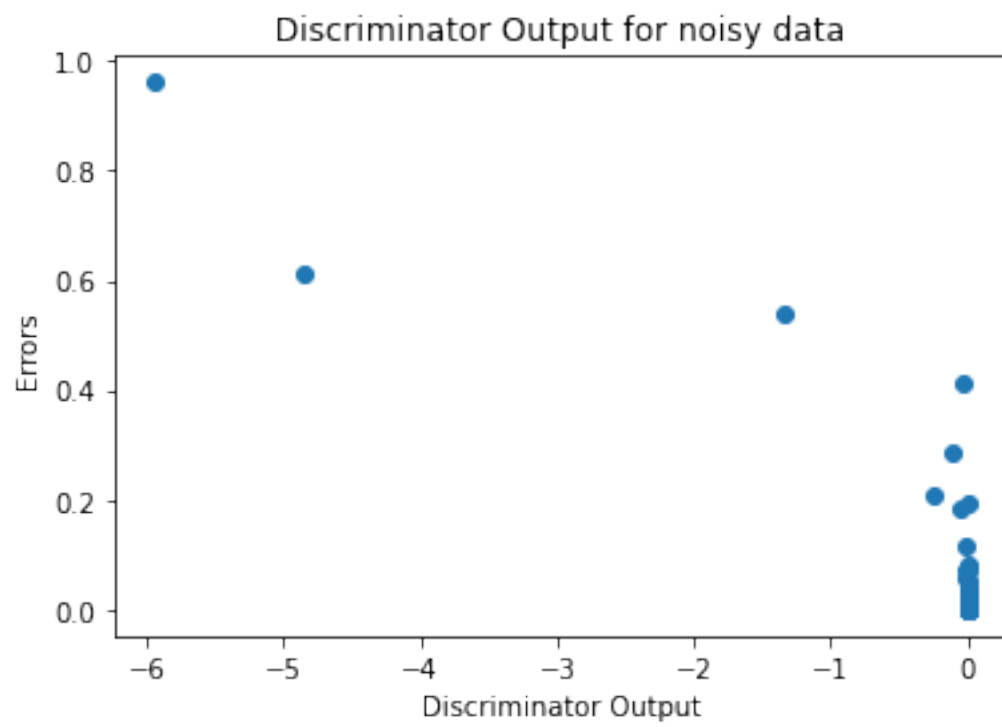
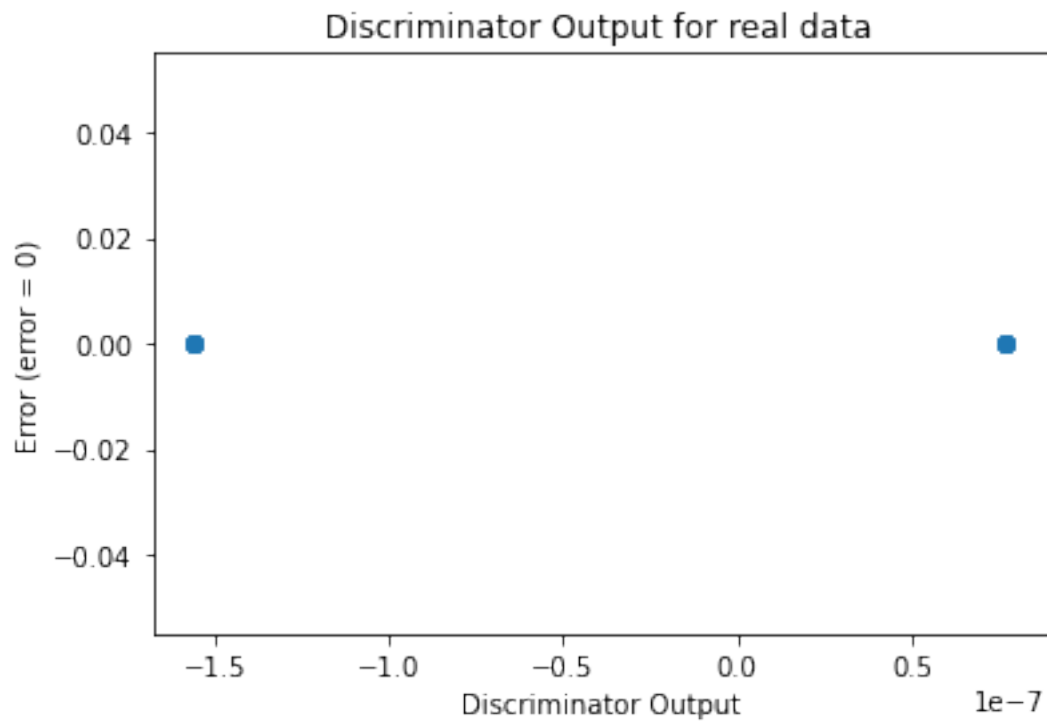


Mean Euclidean Distance: 0.40680517100364105



Sanity Checks

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



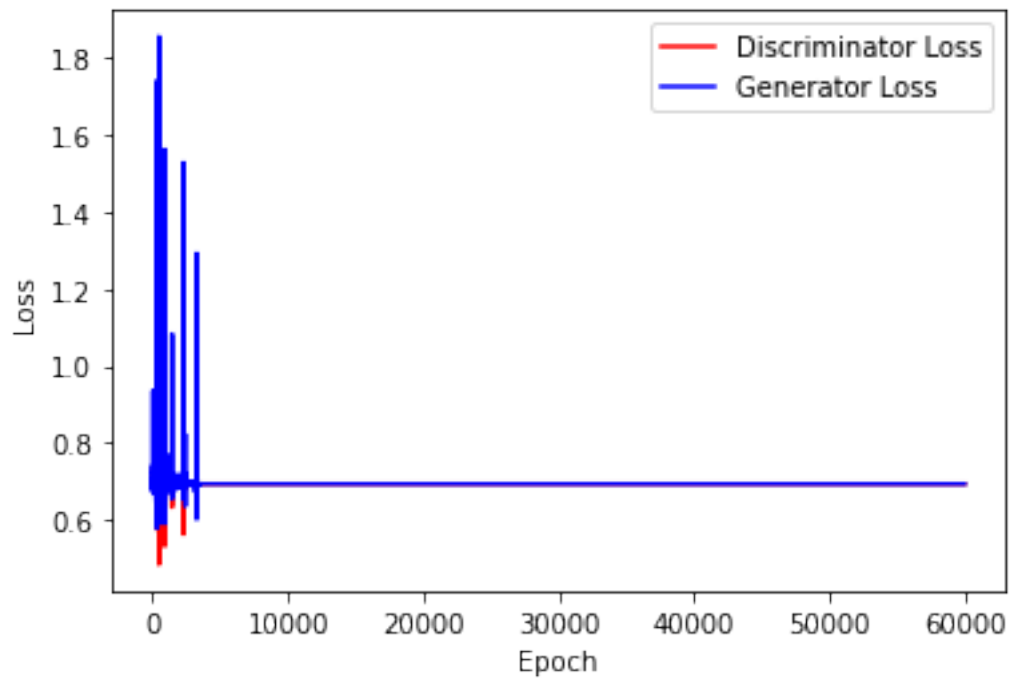
Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

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

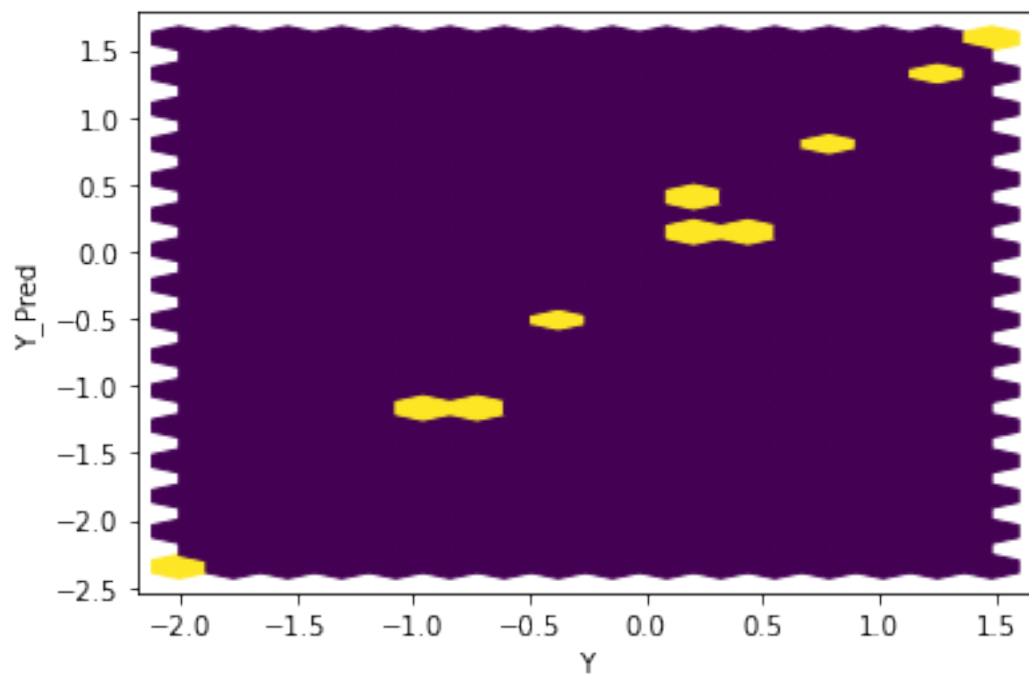
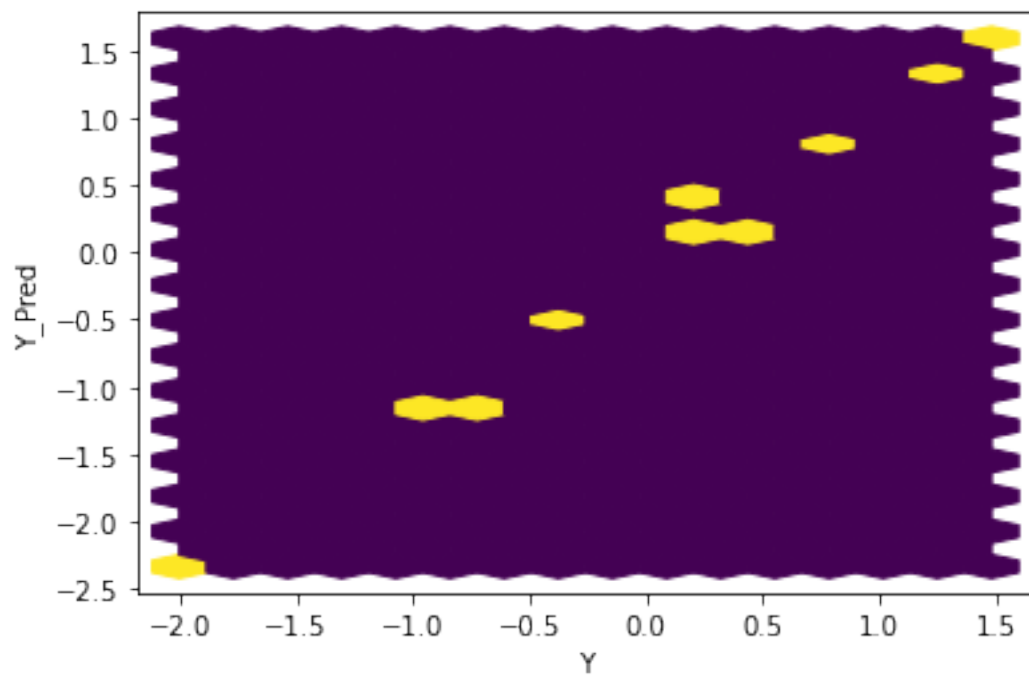
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))
```

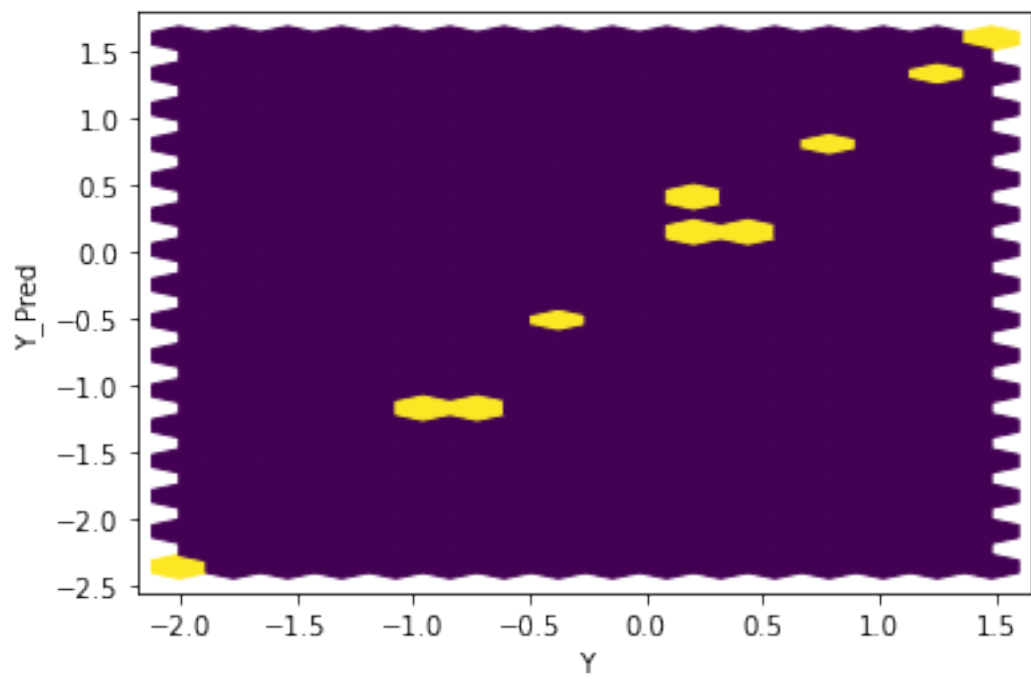
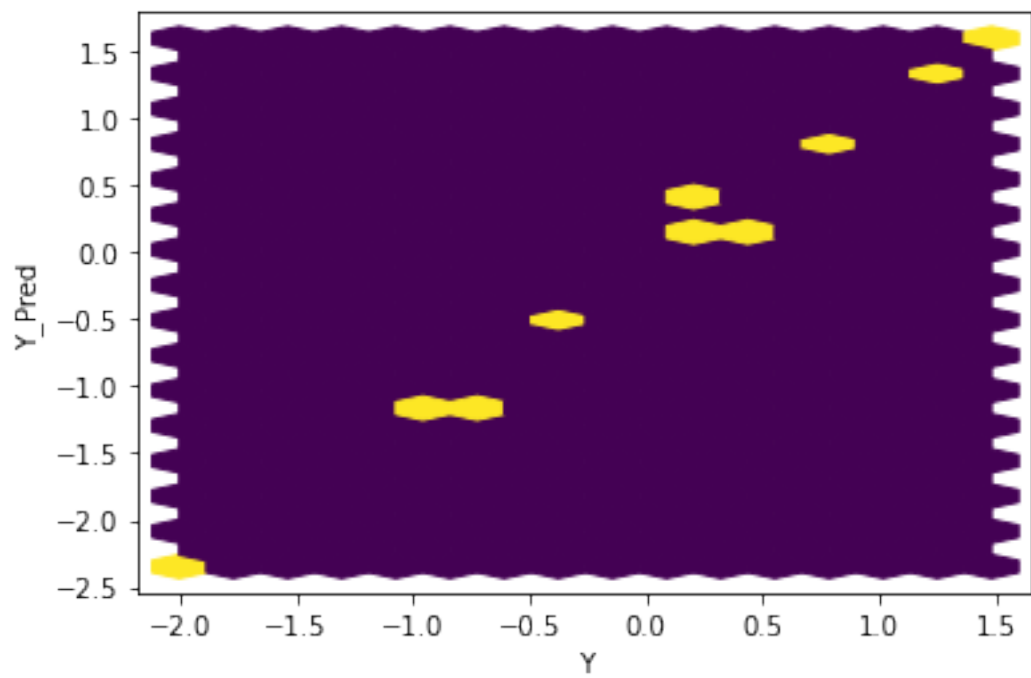
```
[22]: ABC_train_test.
      ↪ training_GAN_2(disc,gen,disc_opt,gen_opt,real_dataset,batch_size,
      ↪ error,criterion,coeff,mean,variance,device)
```

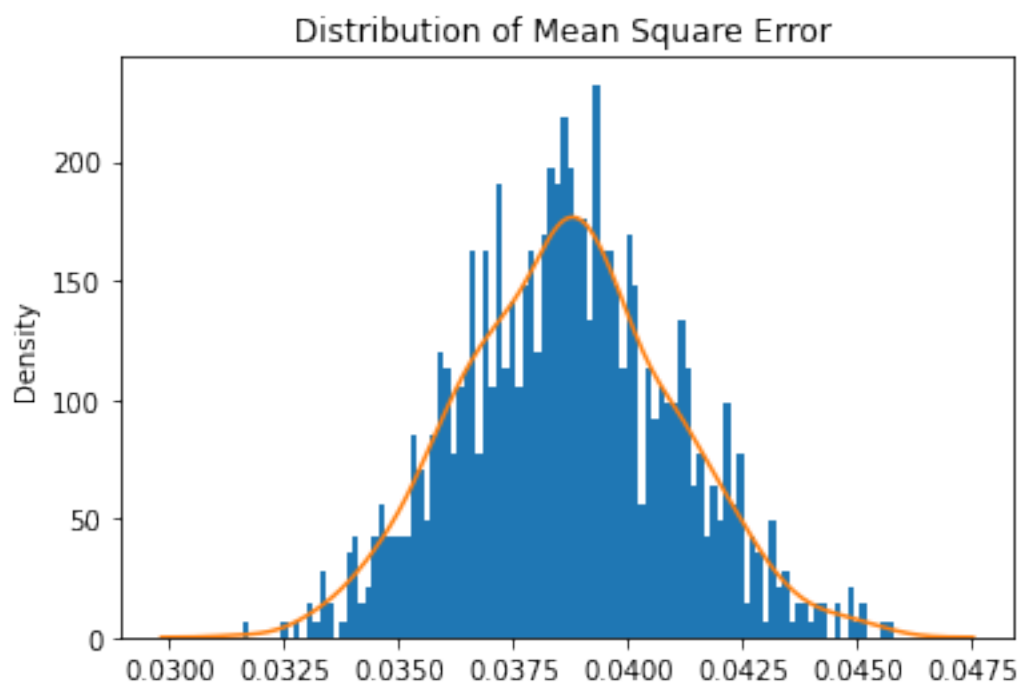
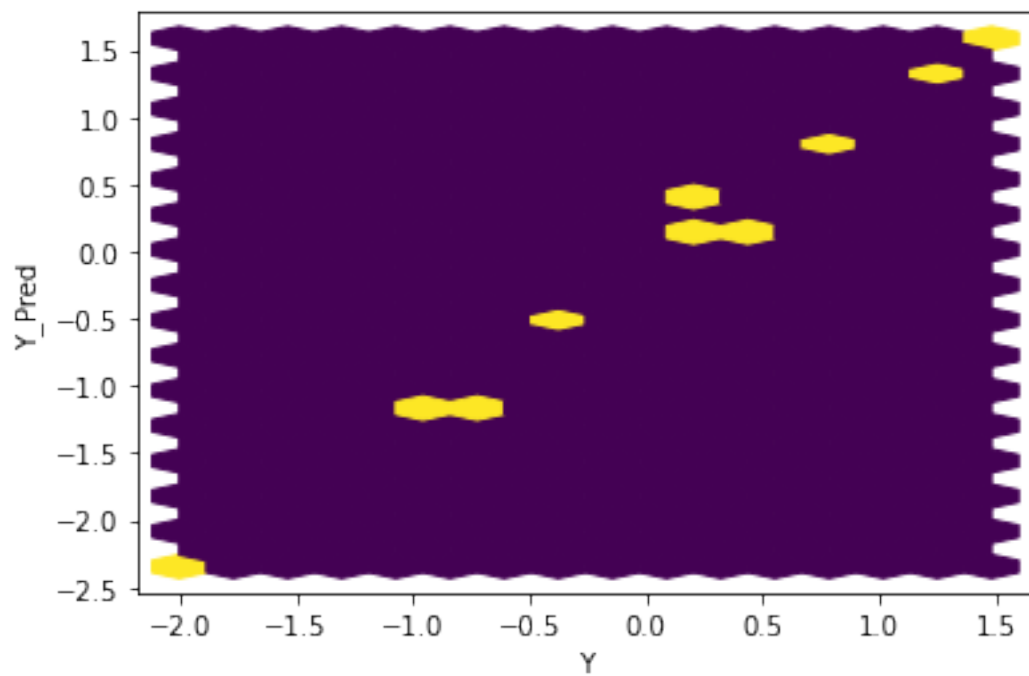
Number of epochs 30000



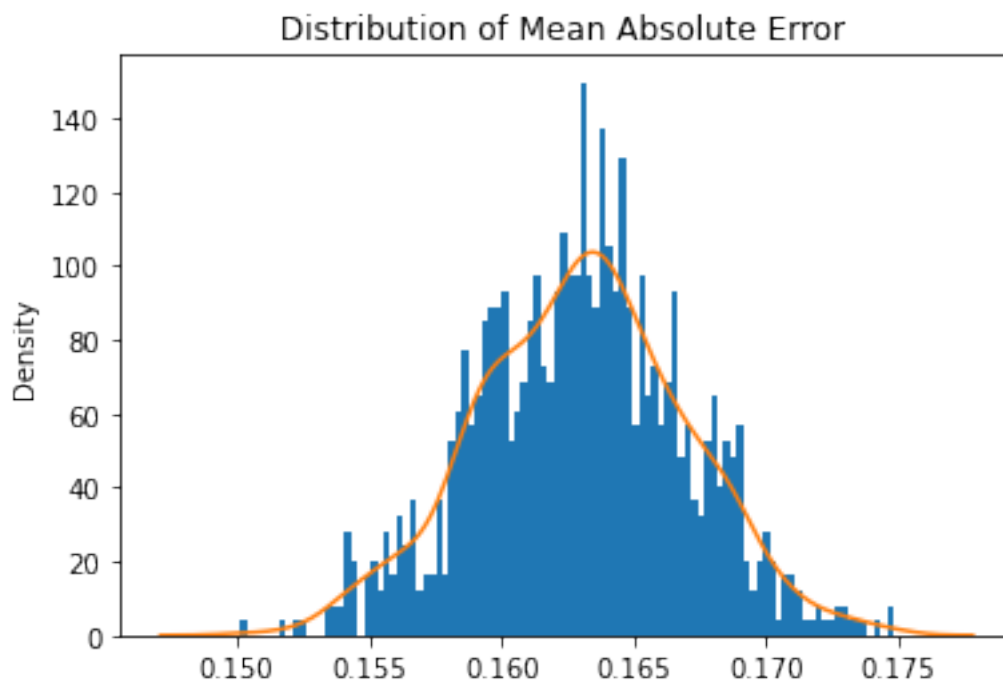
```
[23]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```





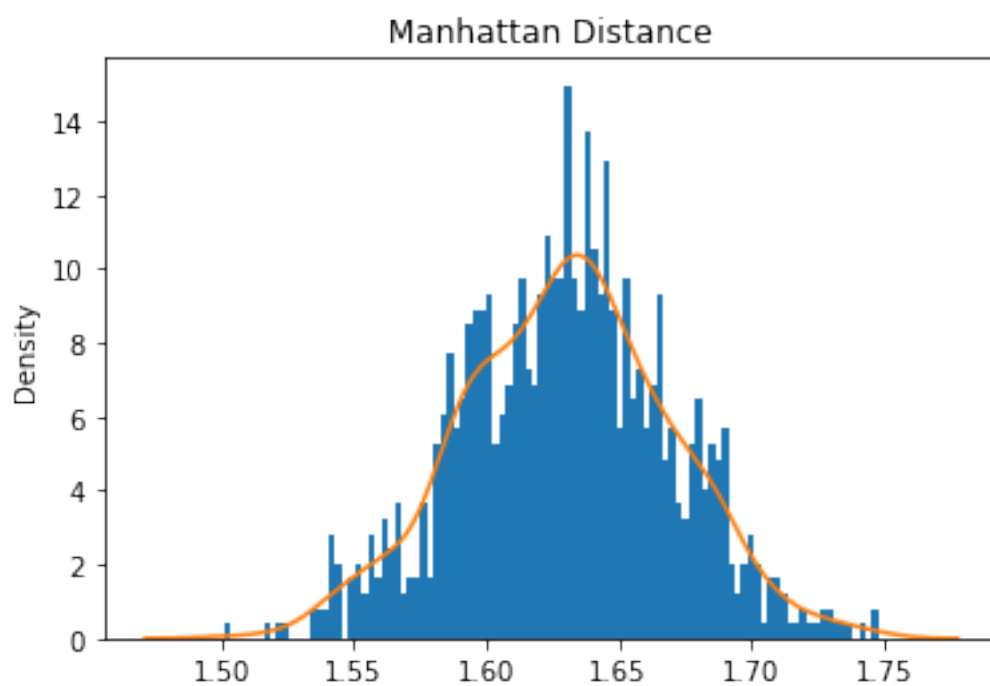


Mean Square Error: 0.038680920402443836

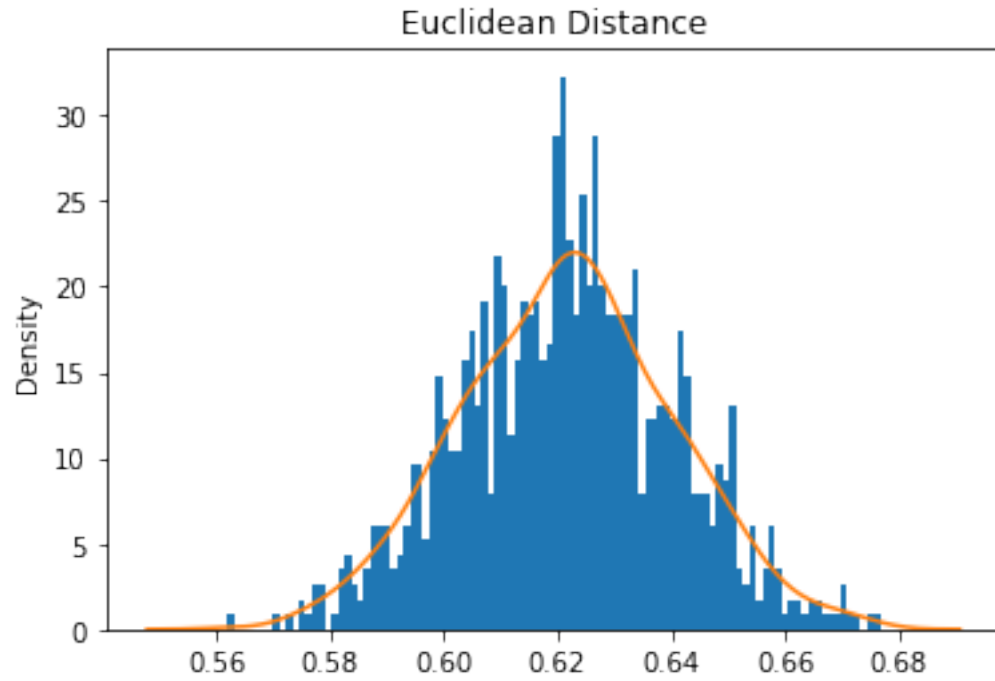


Mean Absolute Error: 0.1630047065794468

Mean Manhattan Distance: 1.630047065794468



Mean Euclidean Distance: 0.6216597382299895



[]: