

Dataset1-Regression_output_3

October 19, 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. $Y = 1 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n + N(0, \sigma)$ 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 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 = 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	0.156774	1.787233	-0.438206	0.362037	-0.055433	-0.781328	-0.479756
1	-0.227270	-0.240804	-1.251164	0.804242	-1.938773	-1.092935	0.173623
2	-0.741364	-0.601604	-3.169273	-0.795221	1.154517	-0.200235	0.666830
3	-0.491081	-1.462484	1.618935	1.183166	-1.241891	0.210289	0.433036
4	-0.199979	-1.095289	-0.414224	-0.594502	0.157827	0.030406	1.032034
	X8	X9	X10	Y			

```

0  0.277893  0.053214  0.137308  103.695097
1 -0.256248  0.989091 -0.492101 -221.905297
2  2.089661  2.084155  0.192048  154.350884
3 -0.757491  0.363914  0.053216 -105.638324
4  0.635609  1.485456 -1.106509 -66.341732

```

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:                Tue, 19 Oct 2021      Prob (F-statistic):          nan
Time:                  22:49:23      Log-Likelihood:            329.30
No. Observations:                10      AIC:                      -638.6
Df Residuals:                    0      BIC:                      -635.6
Df Model:                        9
Covariance Type:                nonrobust
=====
               coef      std err          t      P>|t|      [0.025      0.975]
-----
const          7.772e-16         inf         0         nan         nan         nan
x1              0.1092         inf         0         nan         nan         nan
x2              0.4560         inf         0         nan         nan         nan
x3              0.0539         inf         0         nan         nan         nan
x4              0.7028         inf         0         nan         nan         nan
x5              0.6005         inf         0         nan         nan         nan
x6              0.4352         inf         0         nan         nan         nan
x7              0.2314         inf         0         nan         nan         nan
x8              0.3563         inf         0         nan         nan         nan
x9              0.1468         inf         0         nan         nan         nan
x10             0.2073         inf         0         nan         nan         nan
=====
Omnibus:                  1.235      Durbin-Watson:              1.432
Prob(Omnibus):            0.539      Jarque-Bera (JB):            0.927
Skew:                    0.566      Prob(JB):                    0.629
Kurtosis:                2.029      Cond. No.                    14.8
=====

```

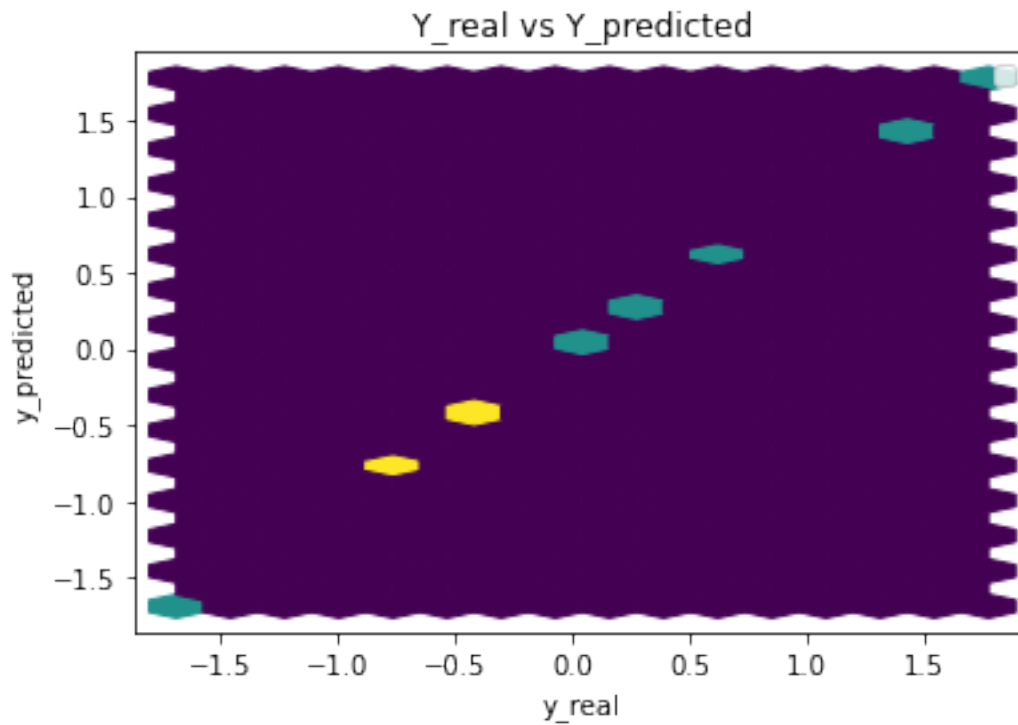
Notes:

[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 7.771561e-16

```
x1      1.091747e-01
x2      4.559758e-01
x3      5.393262e-02
x4      7.027847e-01
x5      6.005446e-01
x6      4.351942e-01
x7      2.314477e-01
x8      3.562659e-01
x9      1.468212e-01
x10     2.072599e-01
dtype: float64
```



Performance Metrics

```
Mean Squared Error: 1.4610537892359059e-30
Mean Absolute Error: 1.041527974976475e-15
Manhattan distance: 1.041527974976475e-14
Euclidean distance: 3.822373332415223e-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

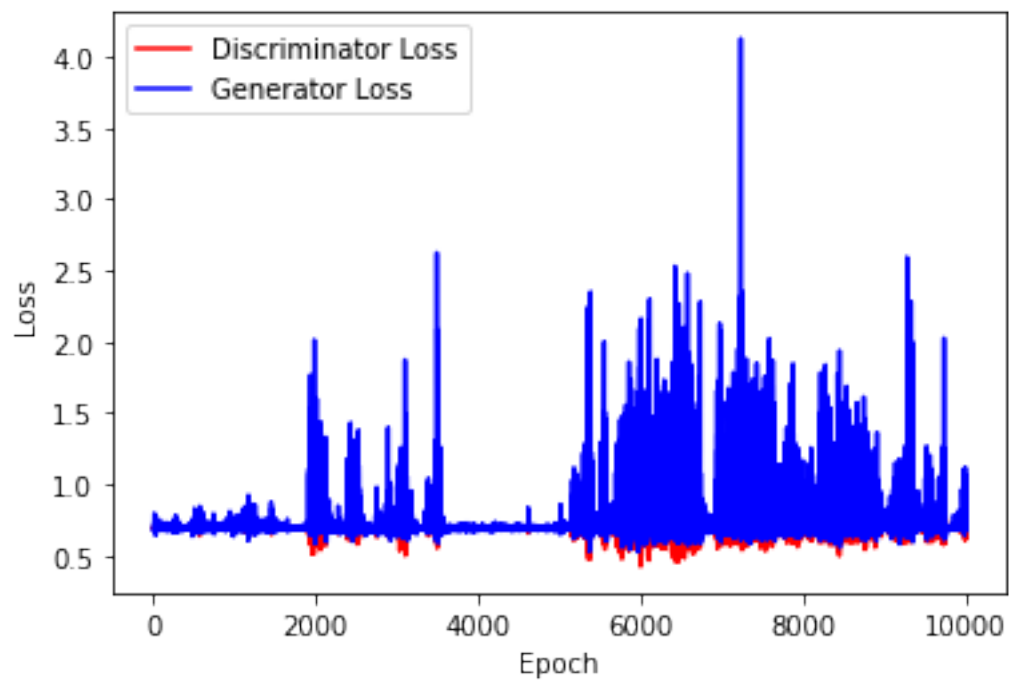
```
[9]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)

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))
```

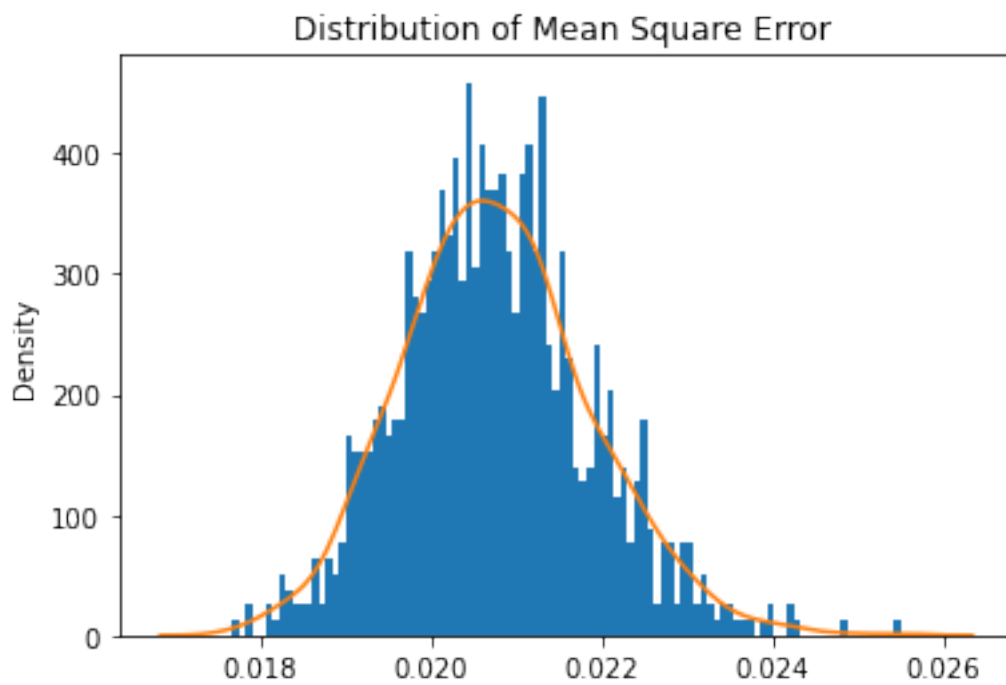
```
[10]: print(generator)
print(discriminator)
```

```
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()
)
```

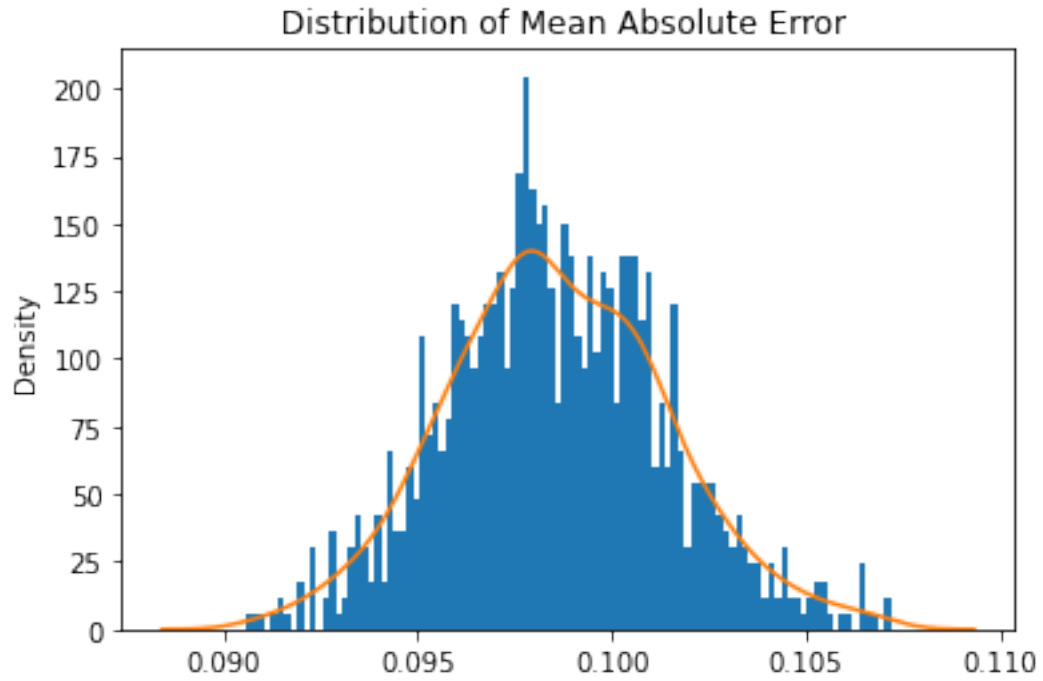
```
[11]: train_test.
→training_GAN(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,
→n_epochs,criterion,device)
```



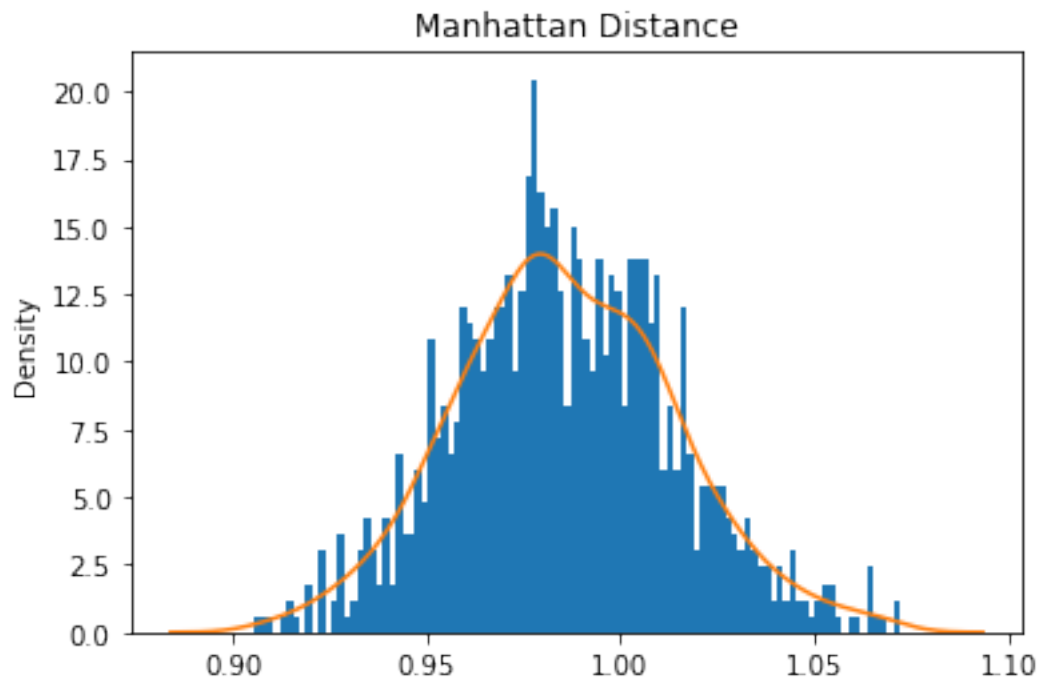
```
[12]: train_test.test_generator(generator,real_dataset,device)
```



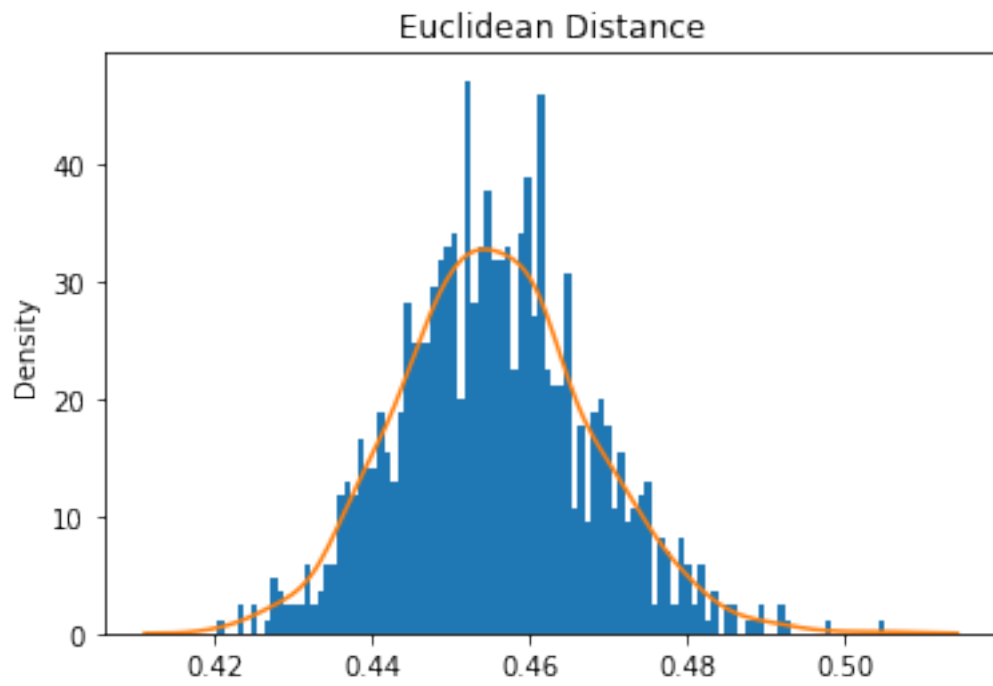
Mean Square Error: 0.020761243660592153



Mean Absolute Error: 0.09857589429616928

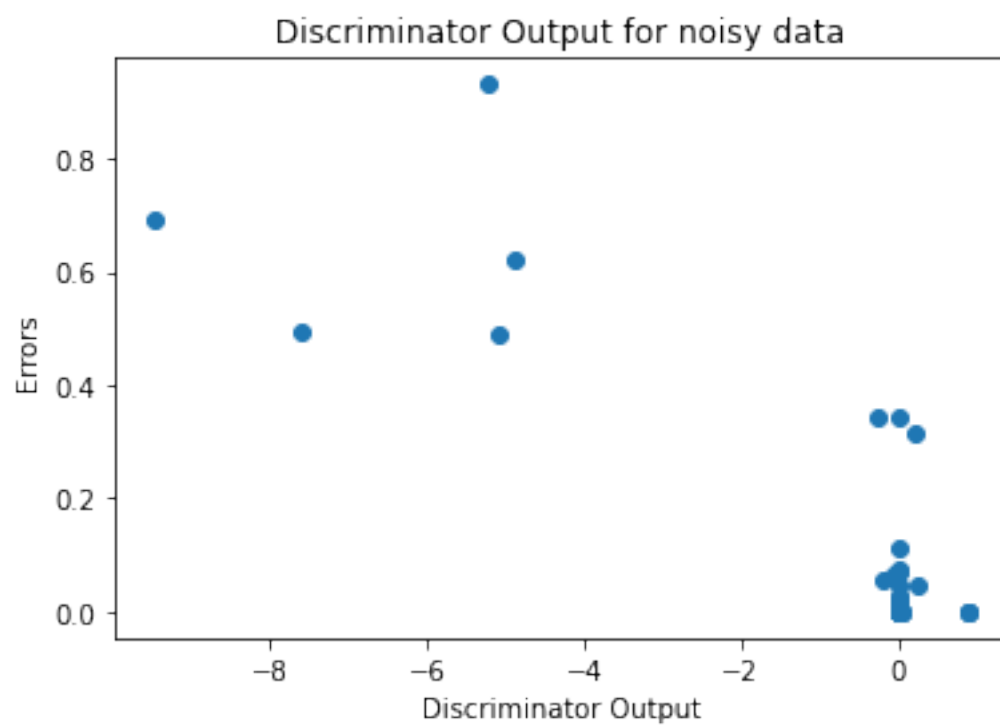
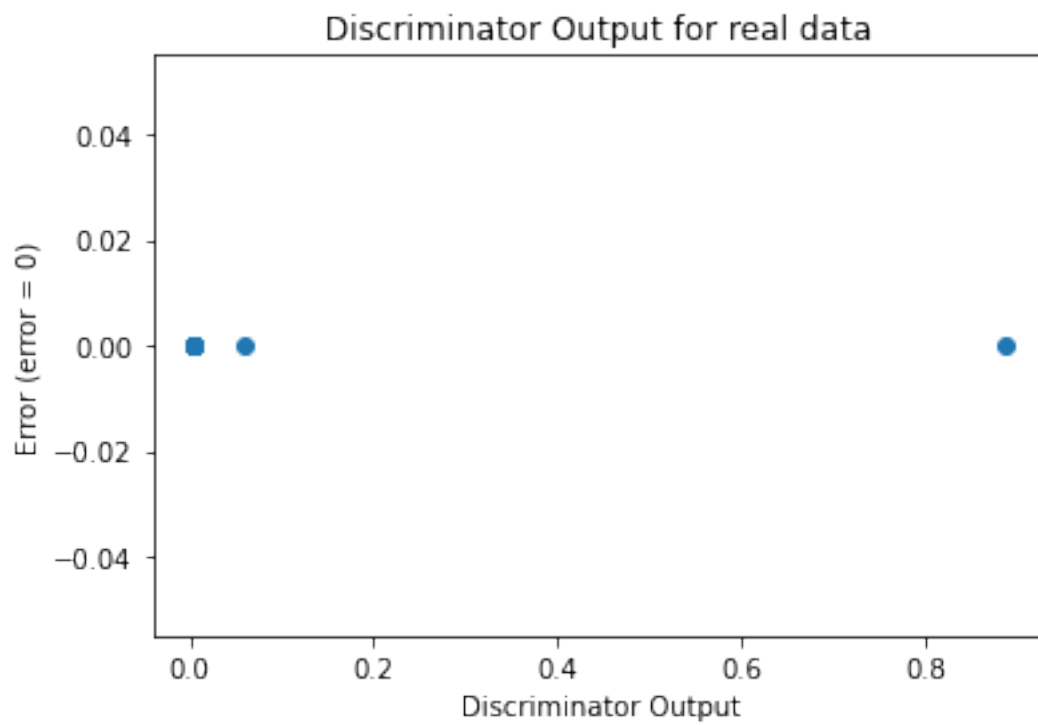


Mean Manhattan Distance: 0.9857589429616929



Mean Euclidean Distance: 0.45548278129039615

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

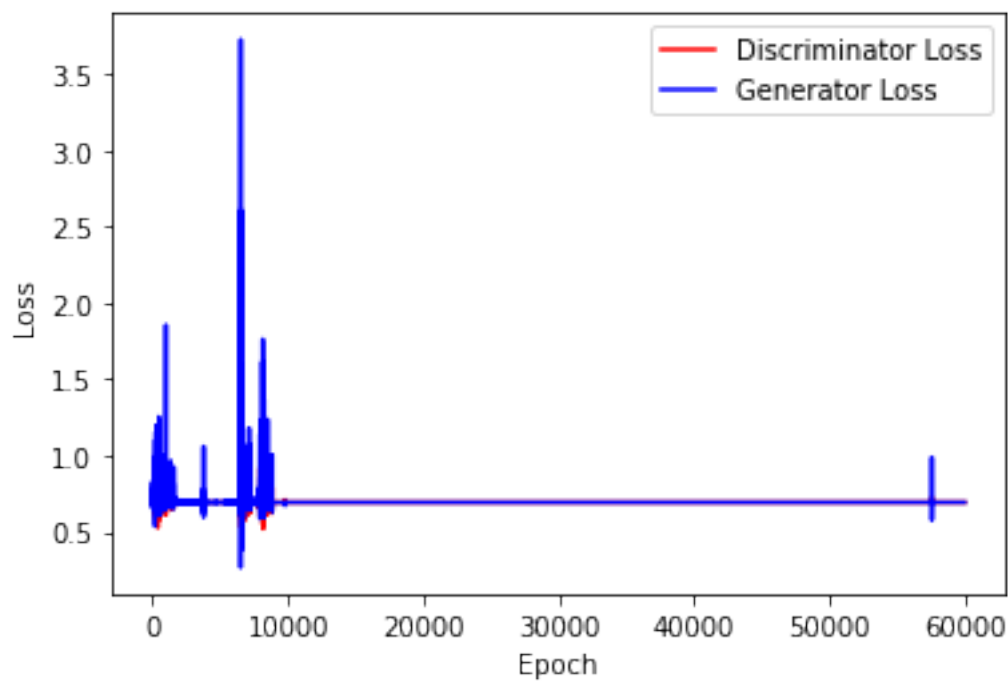



Training GAN until mse of y_{pred} is > 0.1 or $n_{\text{epochs}} < 30000$

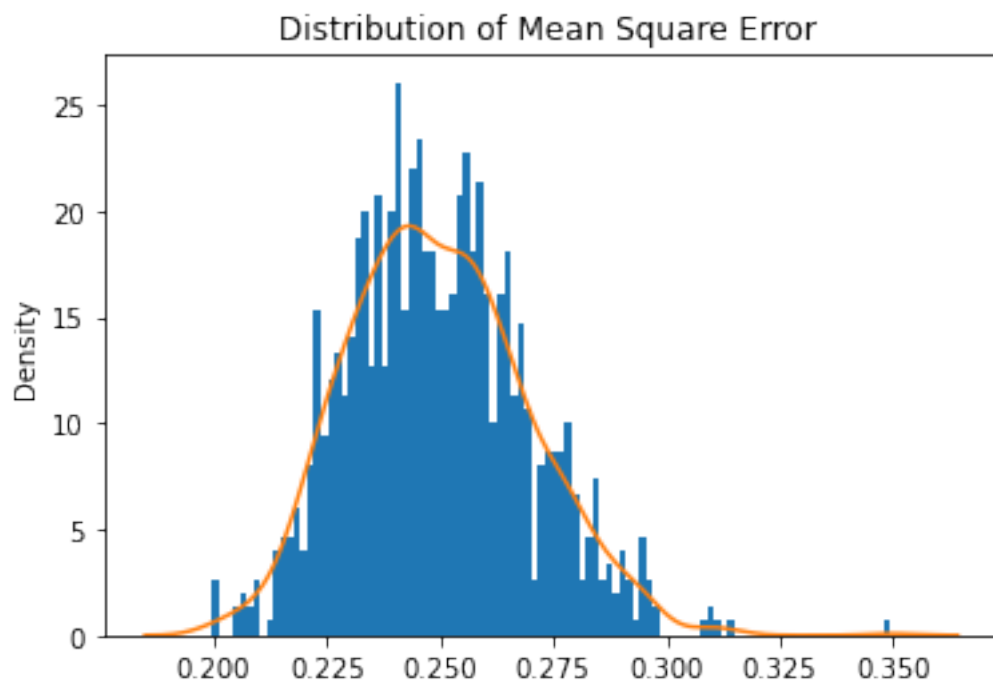
```
[14]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator.parameters(), lr=0.01, betas=(0.5,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

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

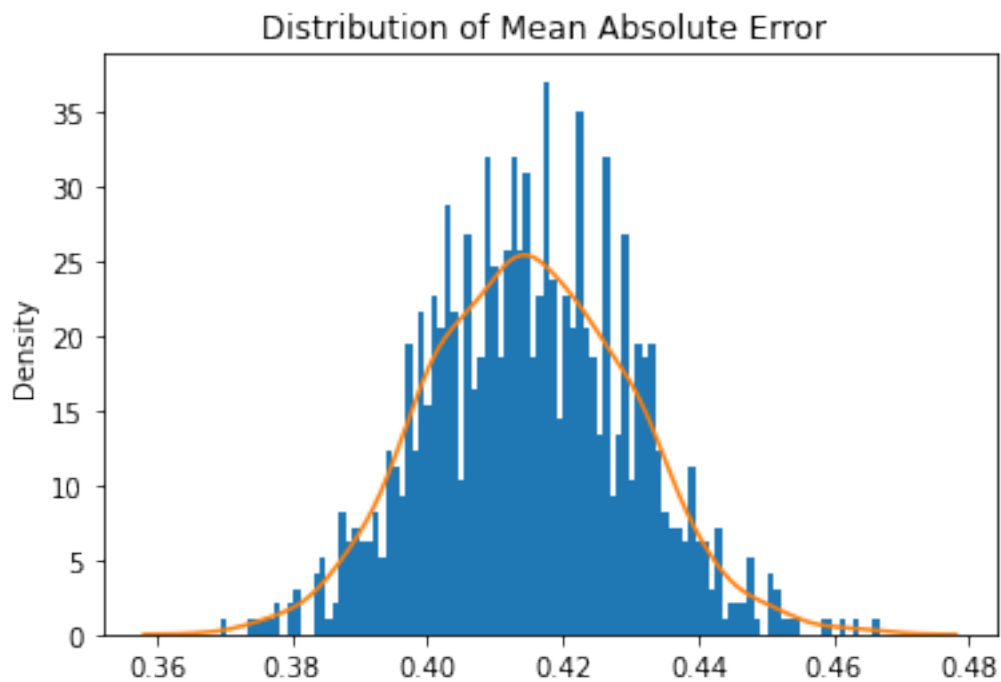
Number of epochs needed 30000



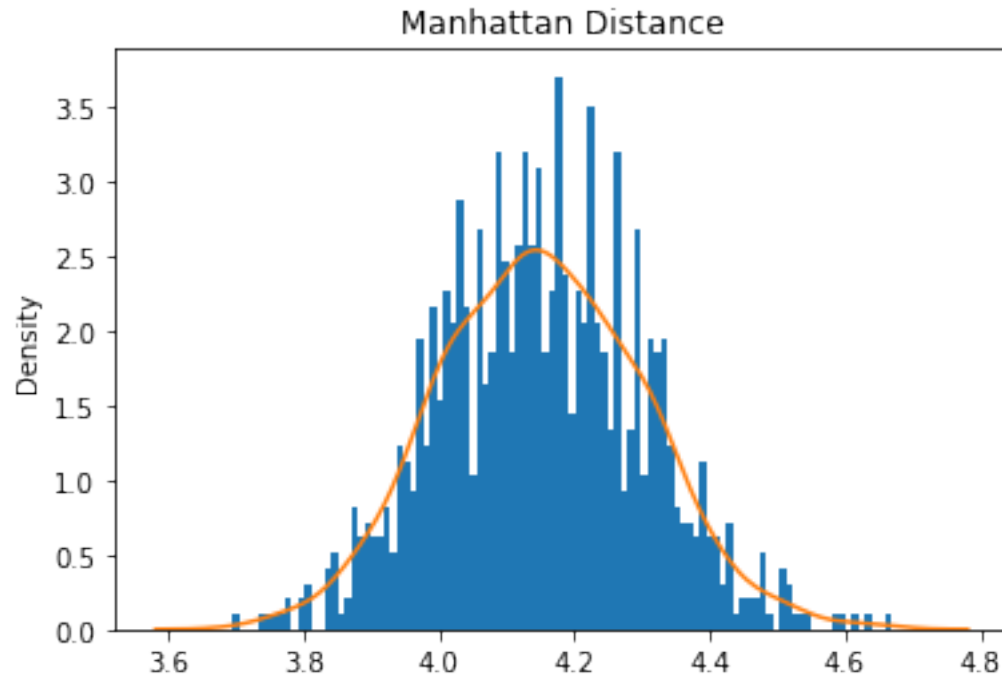
```
[16]: train_test.test_generator(generator,real_dataset,device)
```



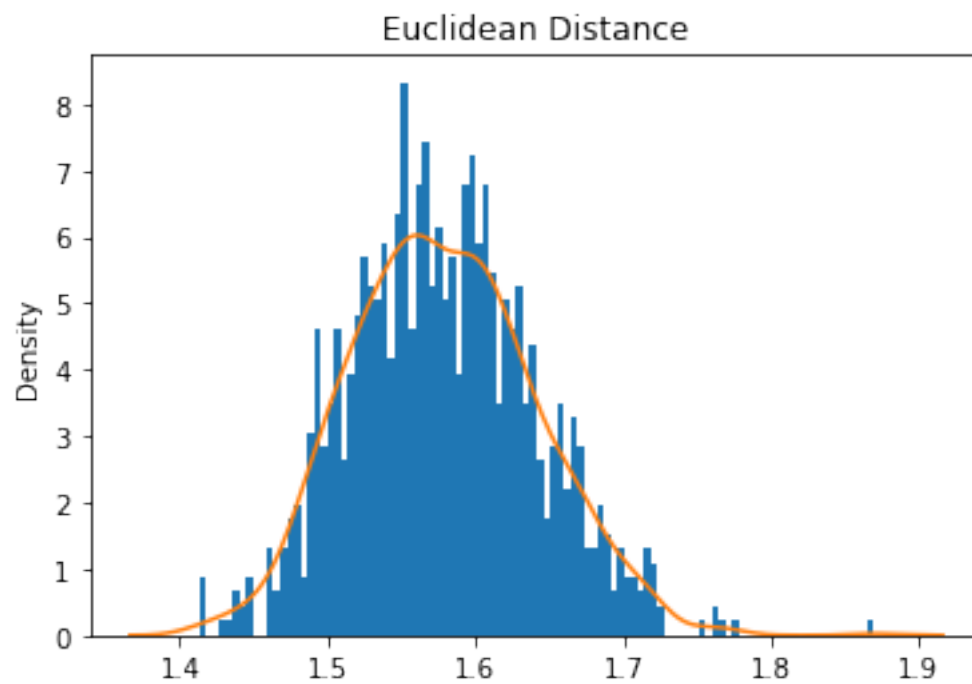
Mean Square Error: 0.24957797262035514



Mean Absolute Error: 0.4150769299000502



Mean Manhattan Distance: 4.150769299000502



Mean Euclidean Distance: 1.5785815749761782

2 ABC GAN Model

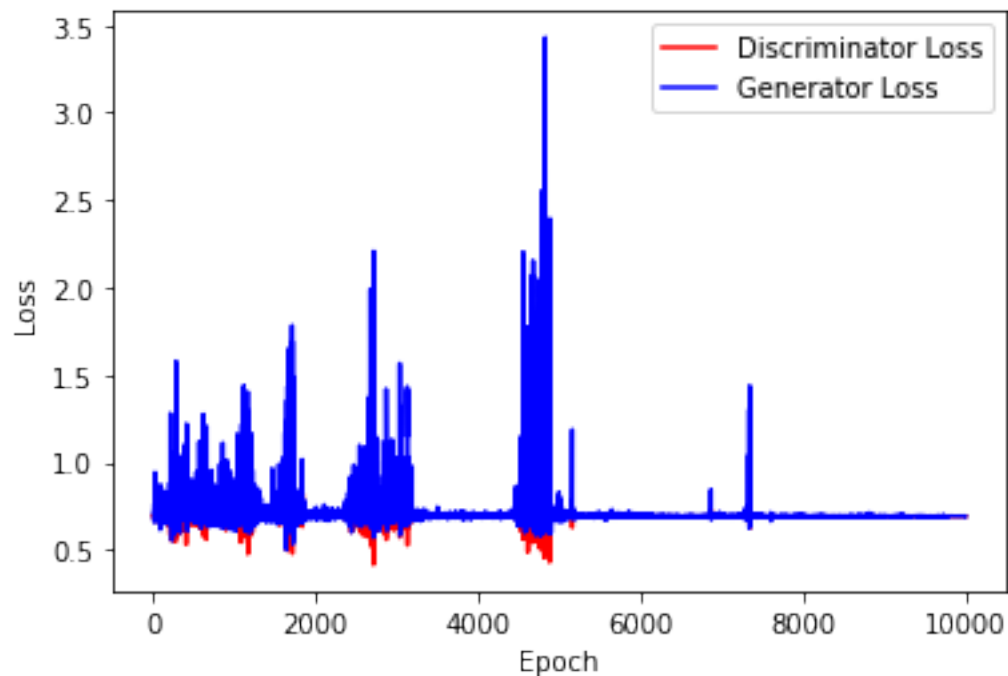
2.0.1 Training the network

Training ABC-GAN for `n_epochs` number of epochs

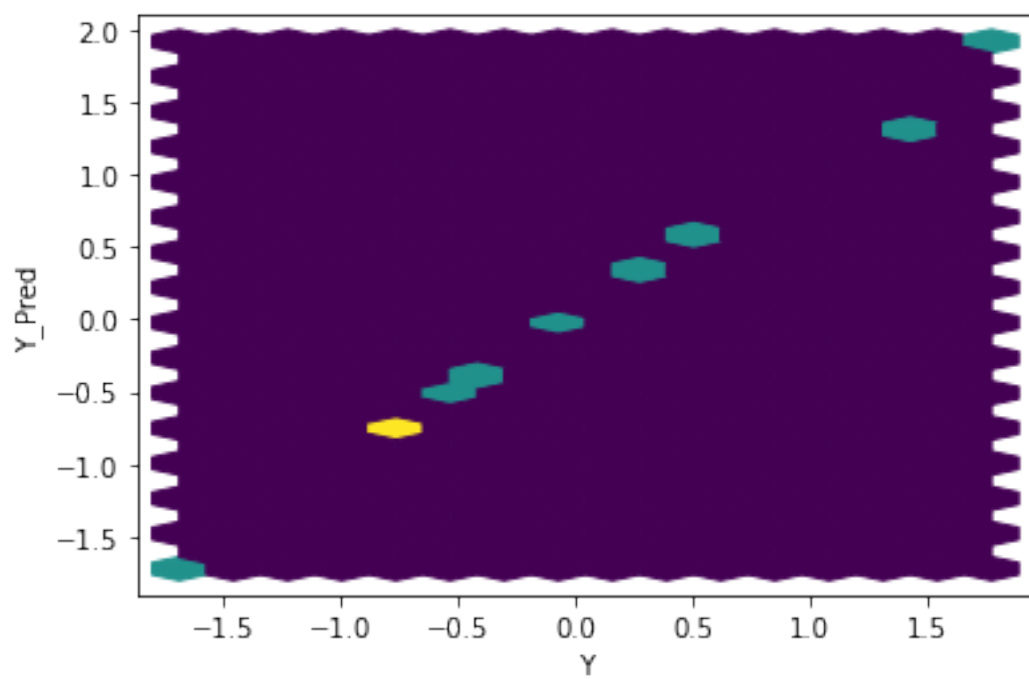
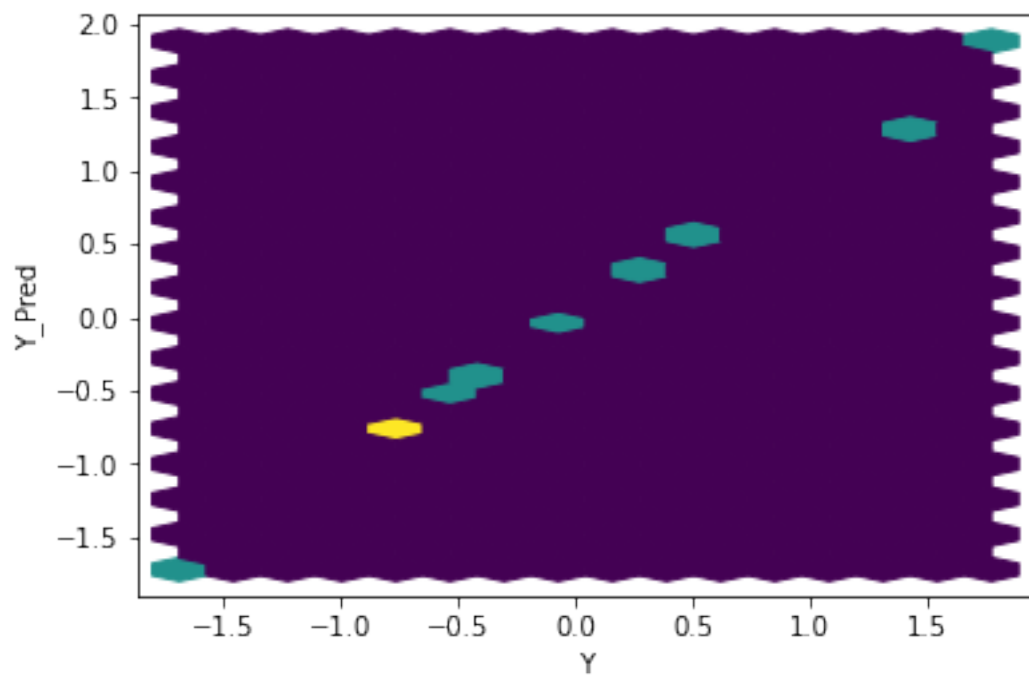
```
[17]: gen = network.Generator(n_features+2)
      disc = network.Discriminator(n_features+2)

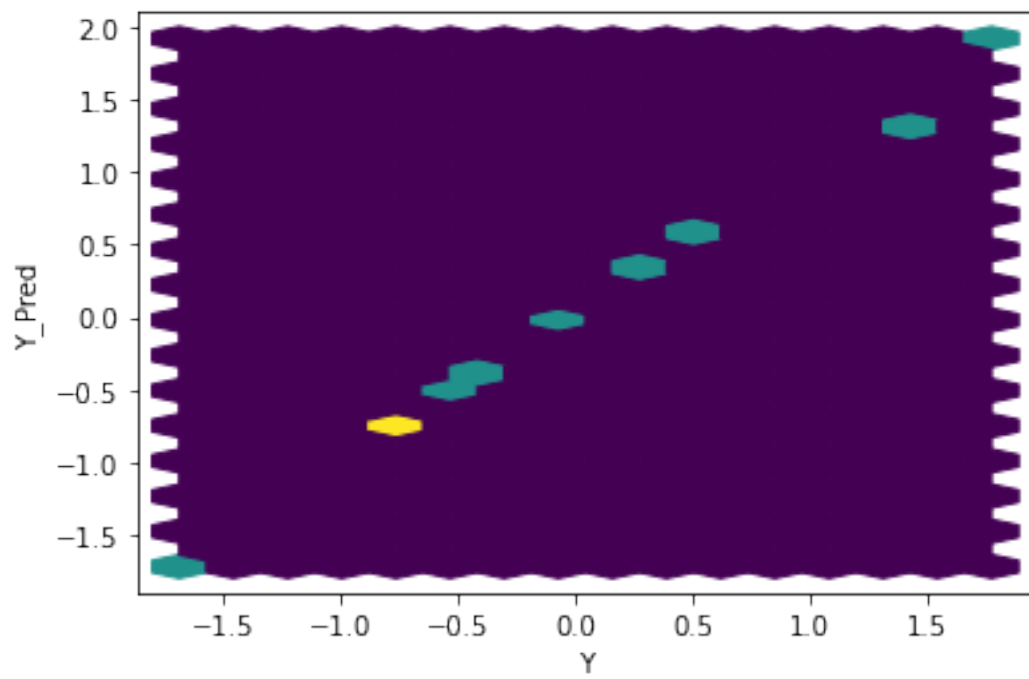
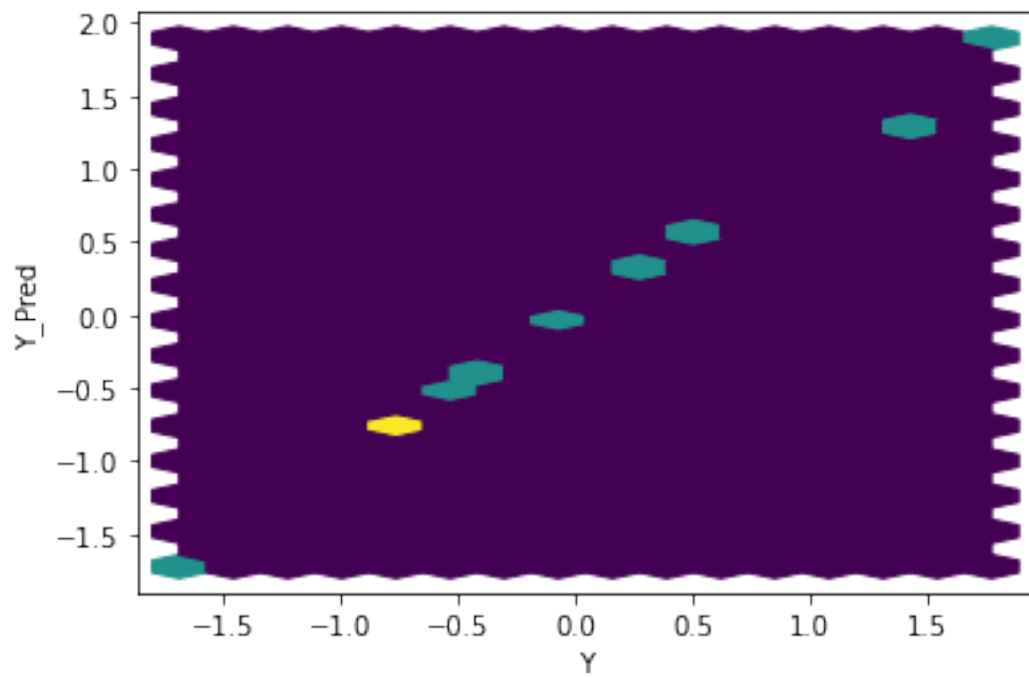
      criterion = torch.nn.BCEWithLogitsLoss()
      gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
      disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))

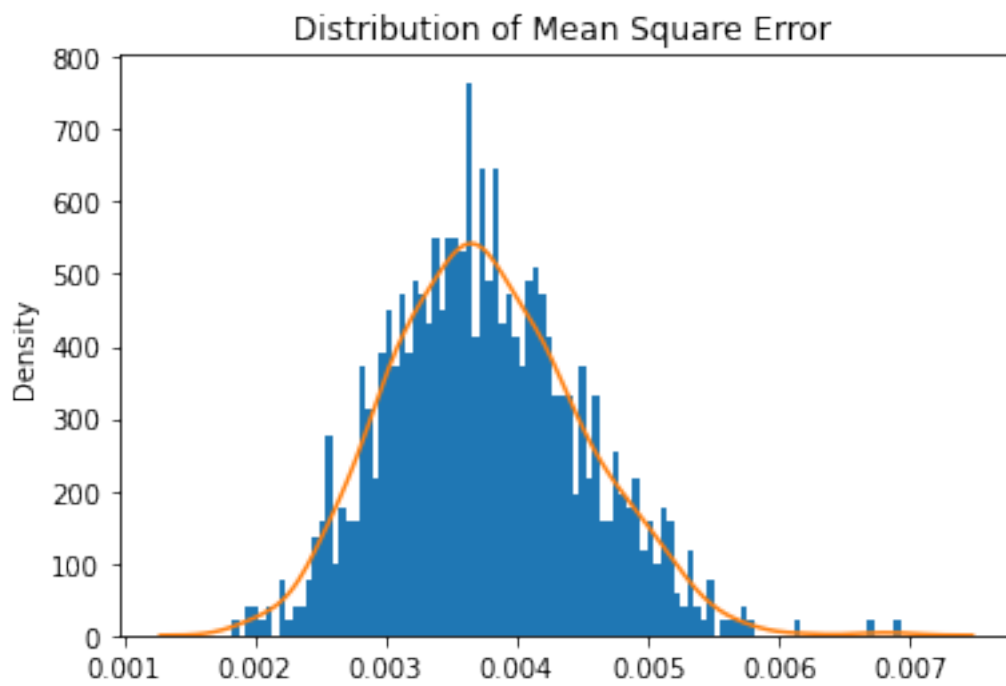
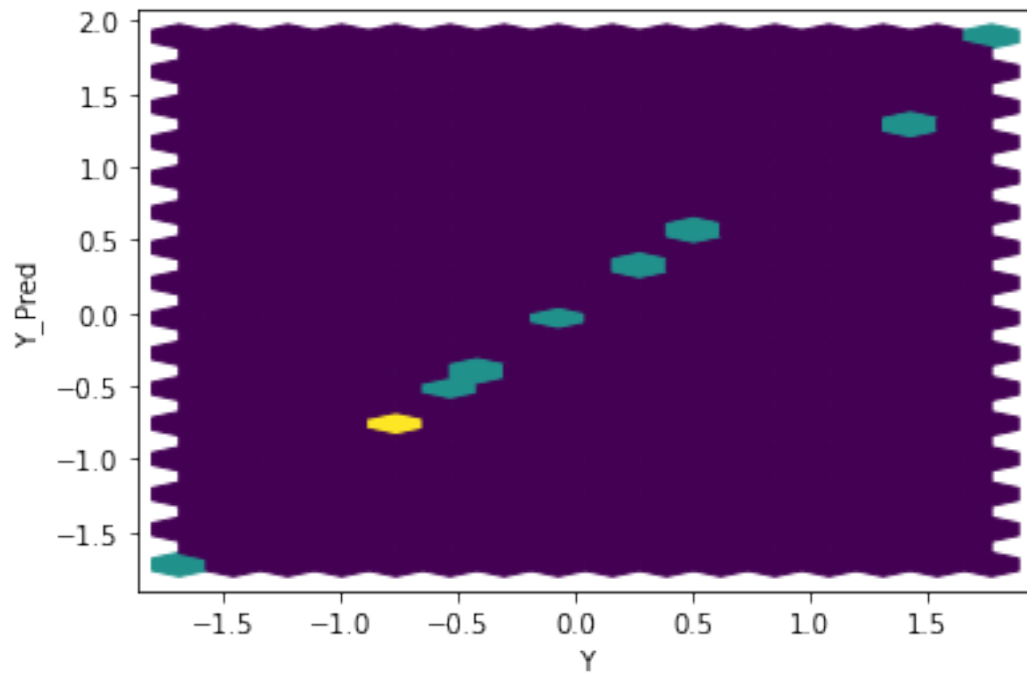
[18]: ABC_train_test.training_GAN(disc, gen,disc_opt,gen_opt,real_dataset,
      ↪batch_size, n_epochs,criterion,coeff,mean,variance,device)
```



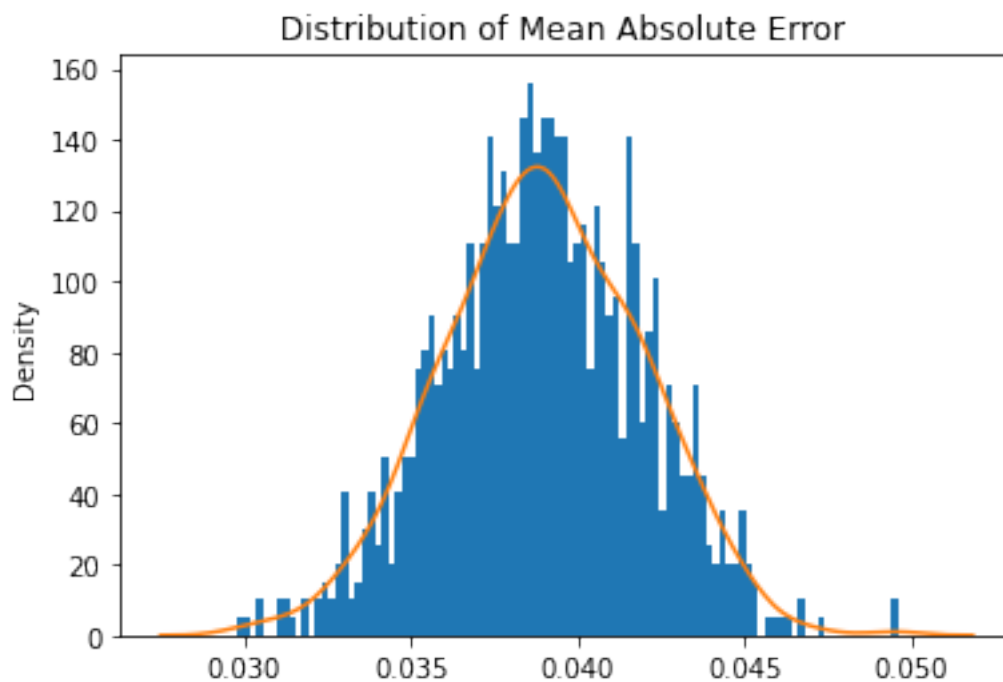
```
[19]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```





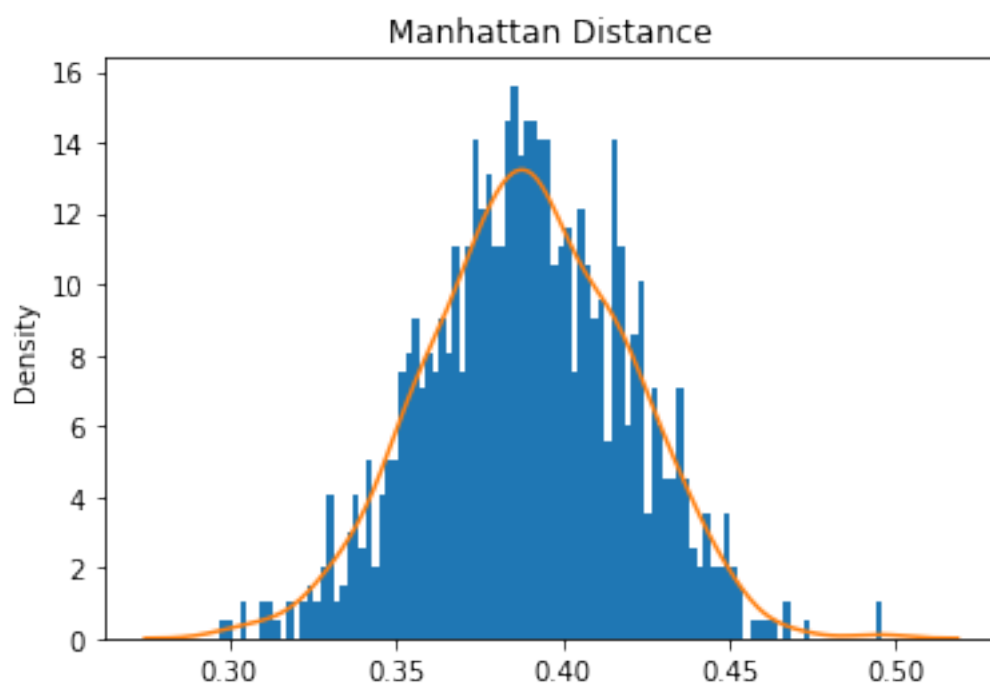


Mean Square Error: 0.0037592004502225336

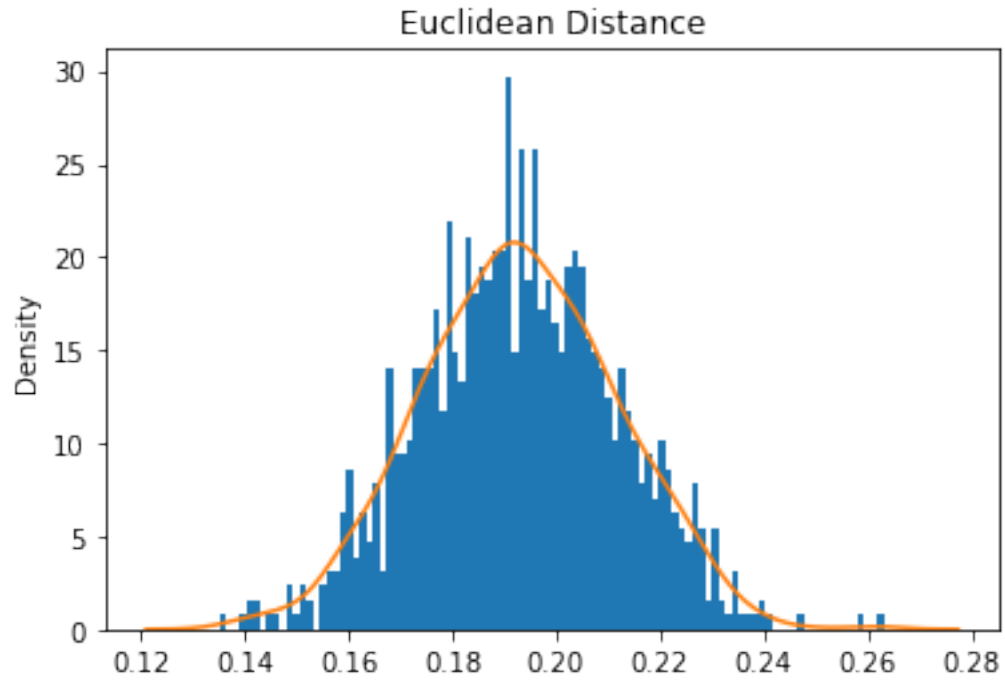


Mean Absolute Error: 0.03894110944420099

Mean Manhattan Distance: 0.38941109444200994

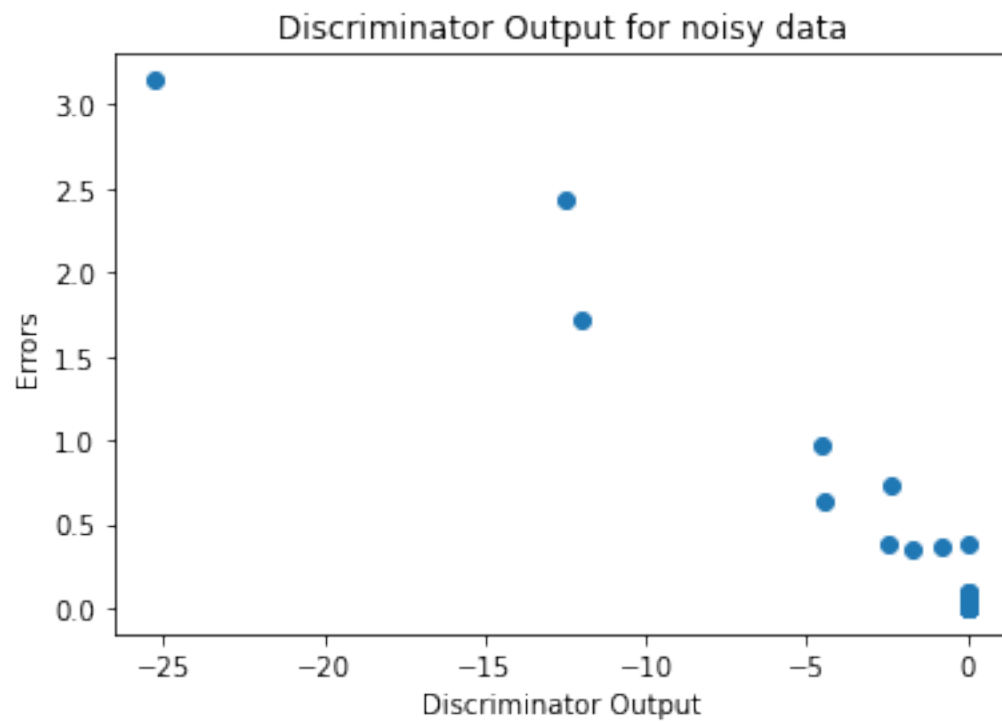
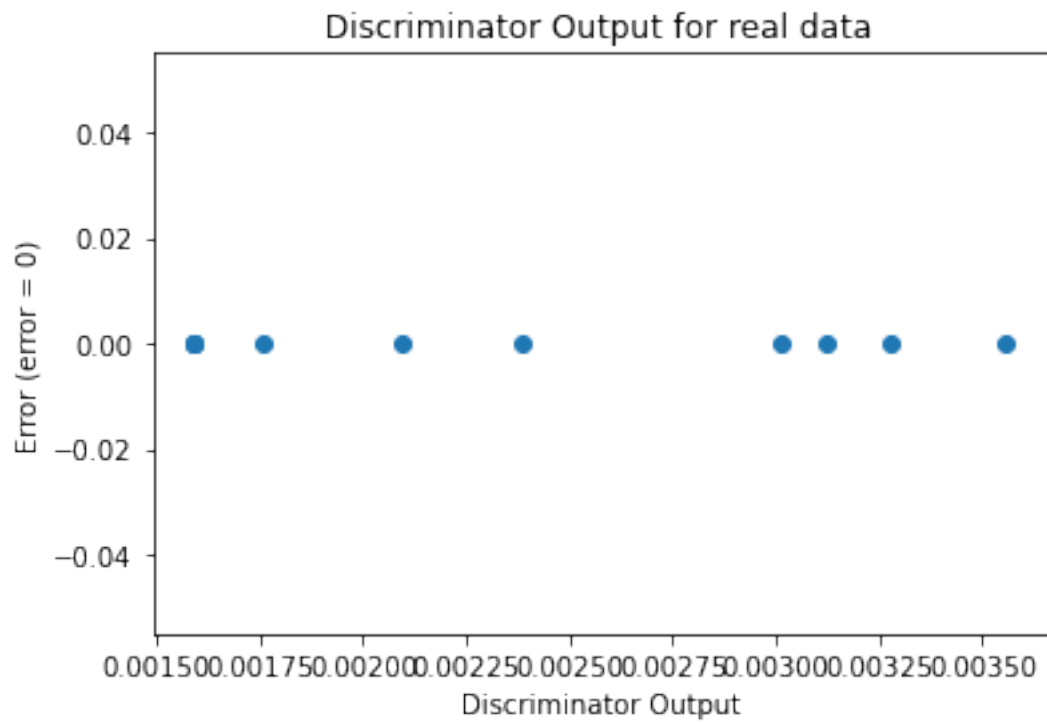


Mean Euclidean Distance: 0.1929655901071764



Sanity Checks

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



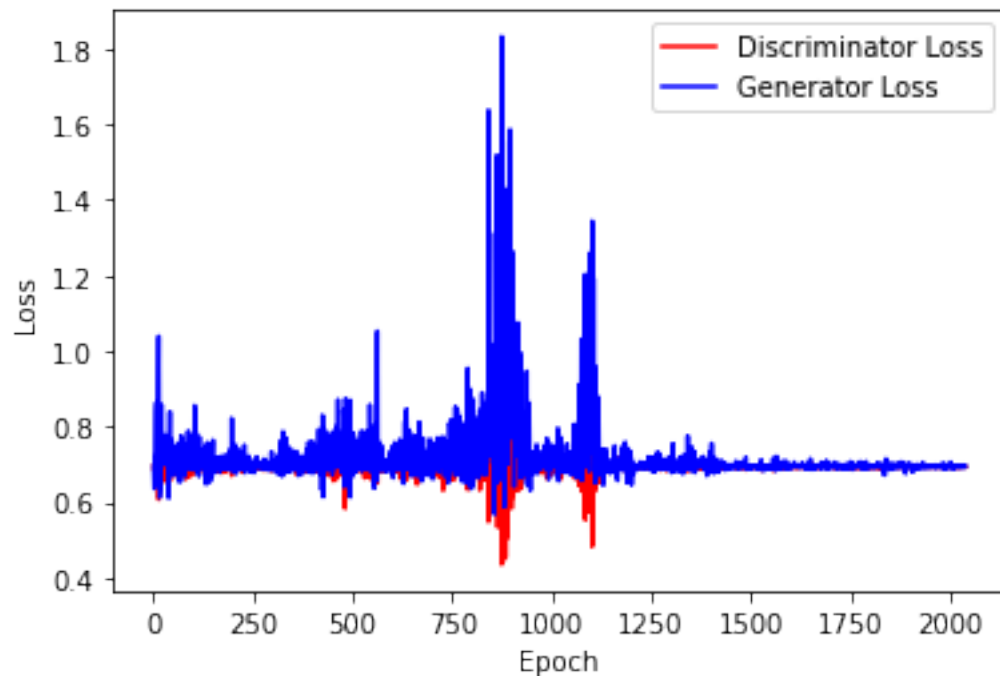
Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

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

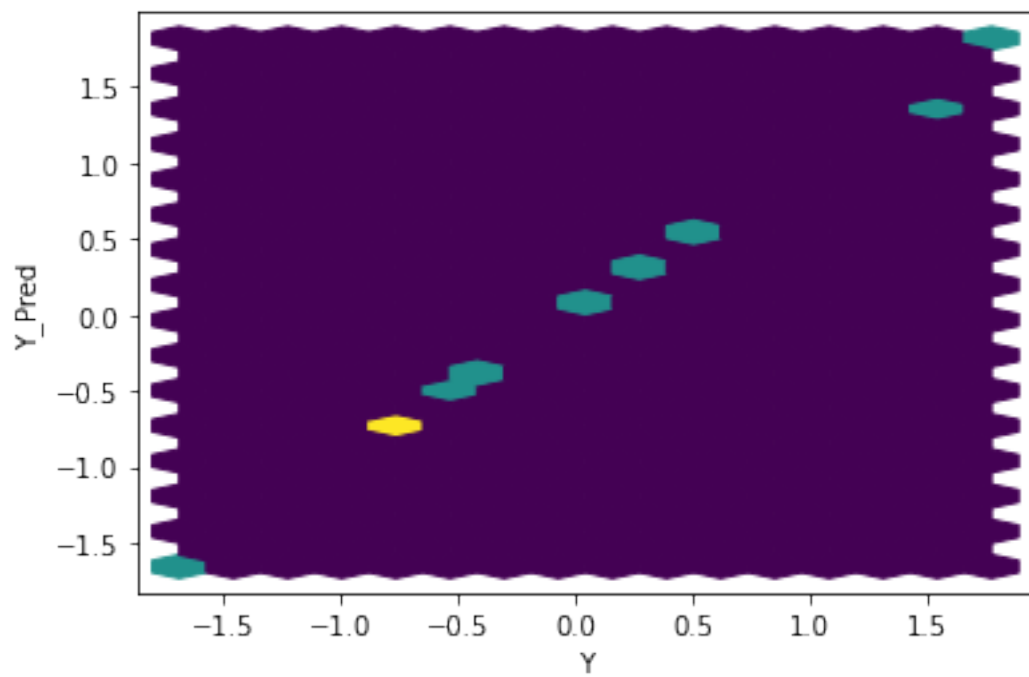
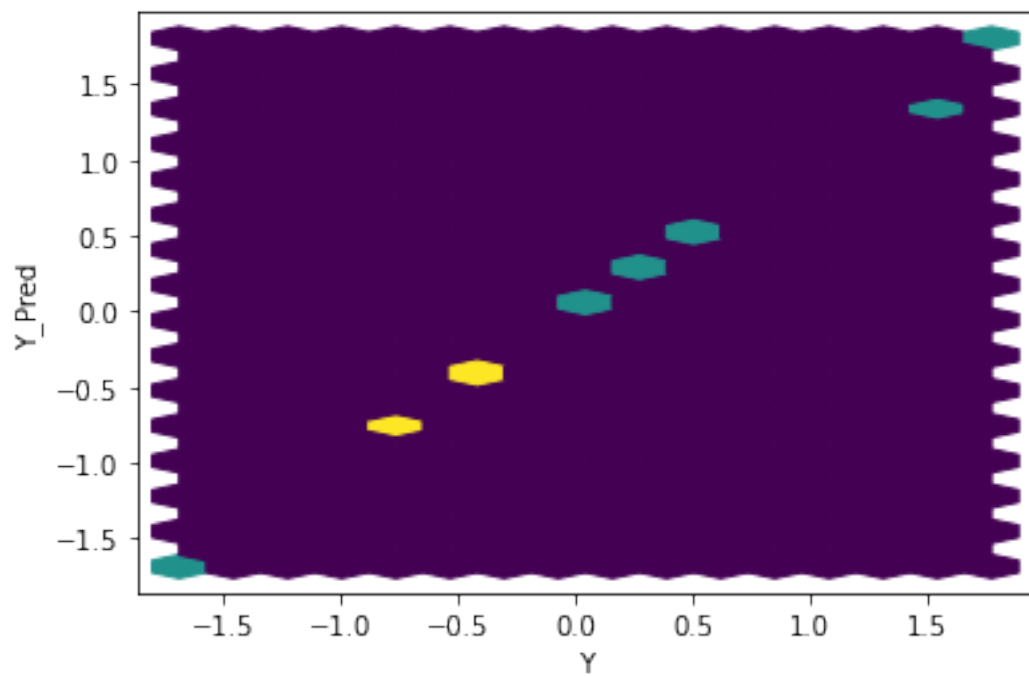
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))
```

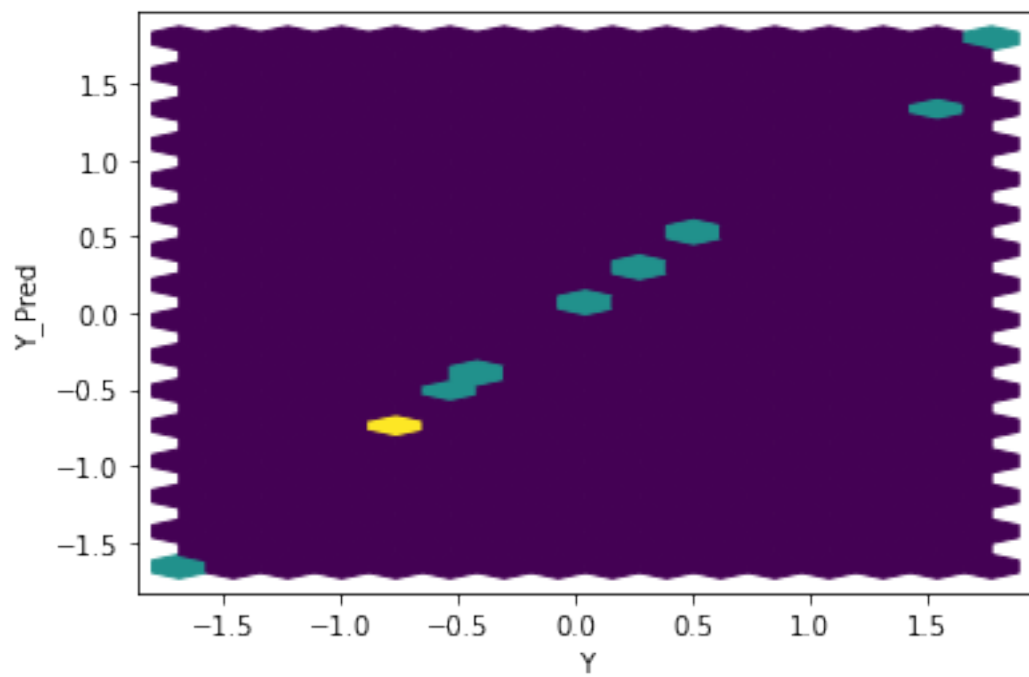
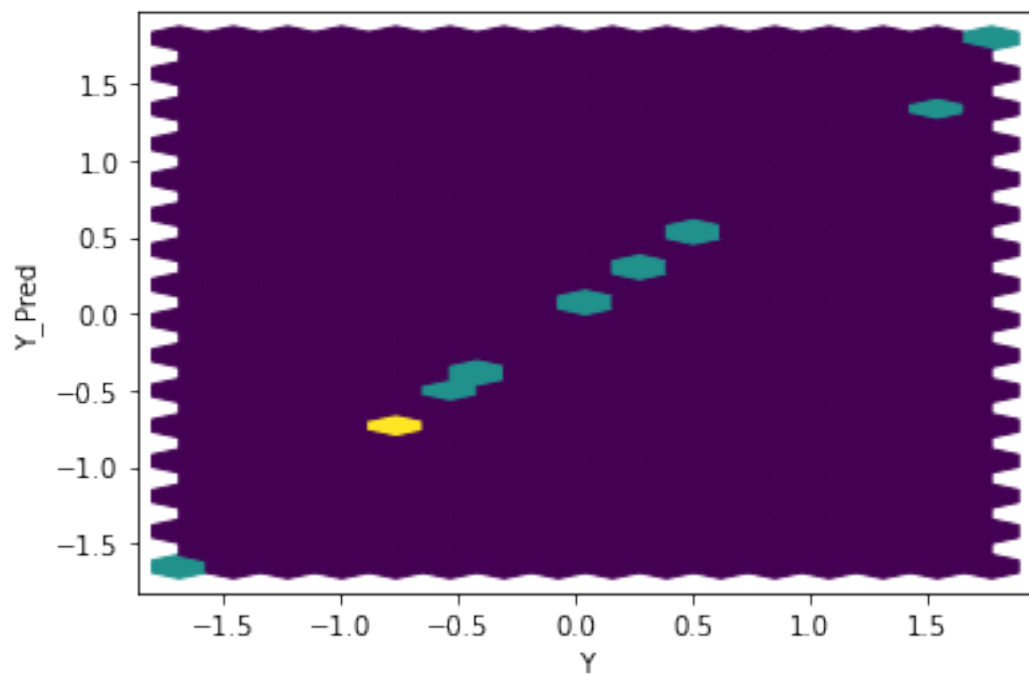
```
[22]: ABC_train_test.
      ↪ training_GAN_2(disc,gen,disc_opt,gen_opt,real_dataset,batch_size,
      ↪ error,criterion,coeff,mean,variance,device)
```

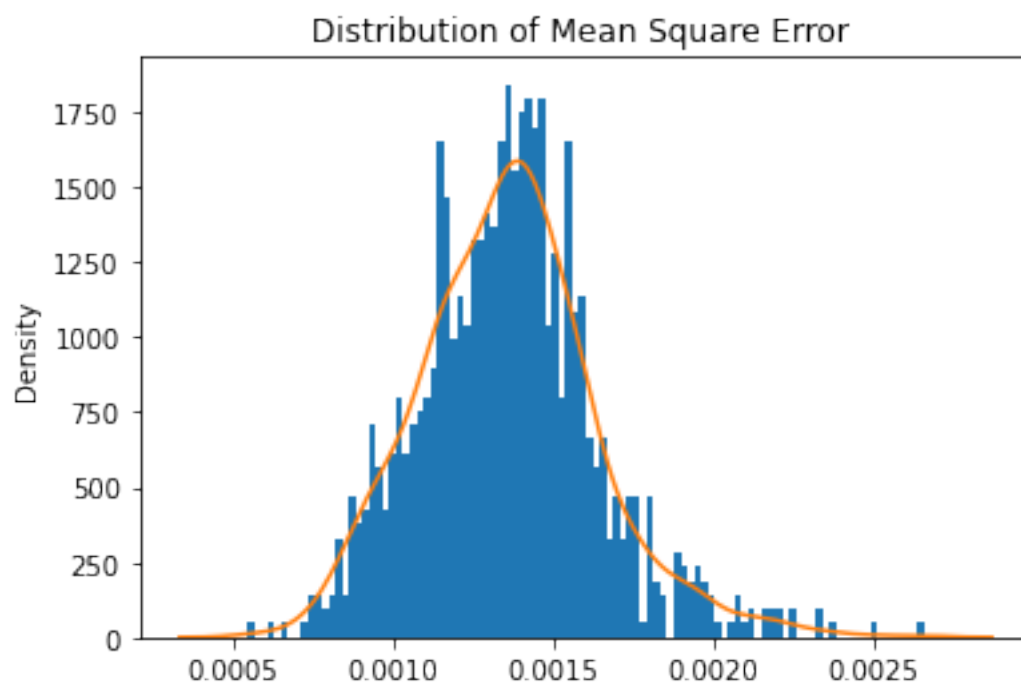
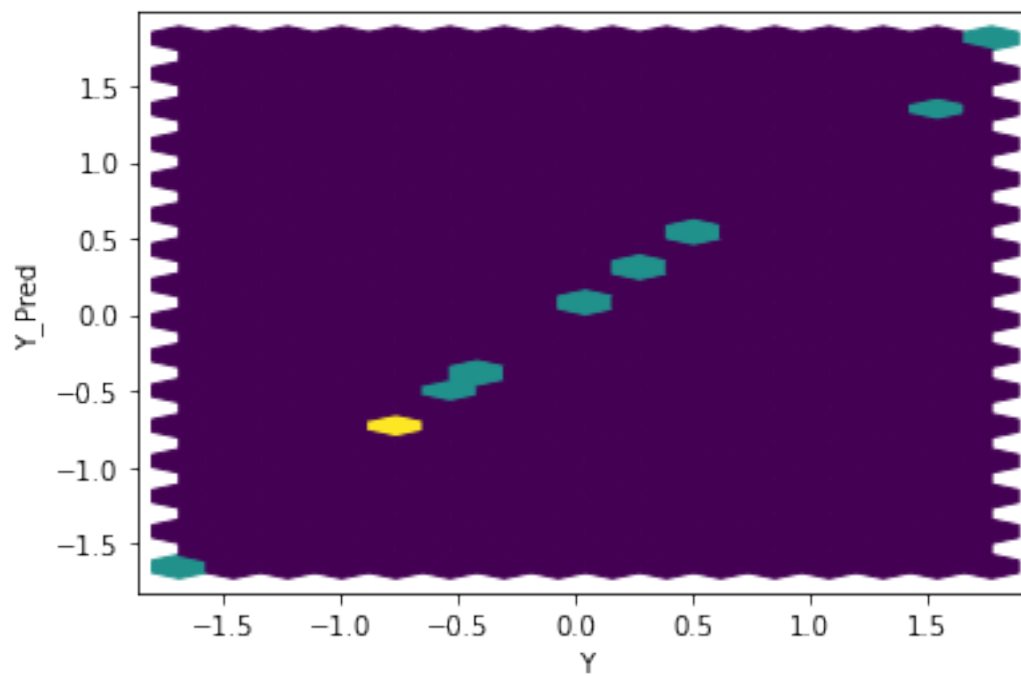
Number of epochs 1020



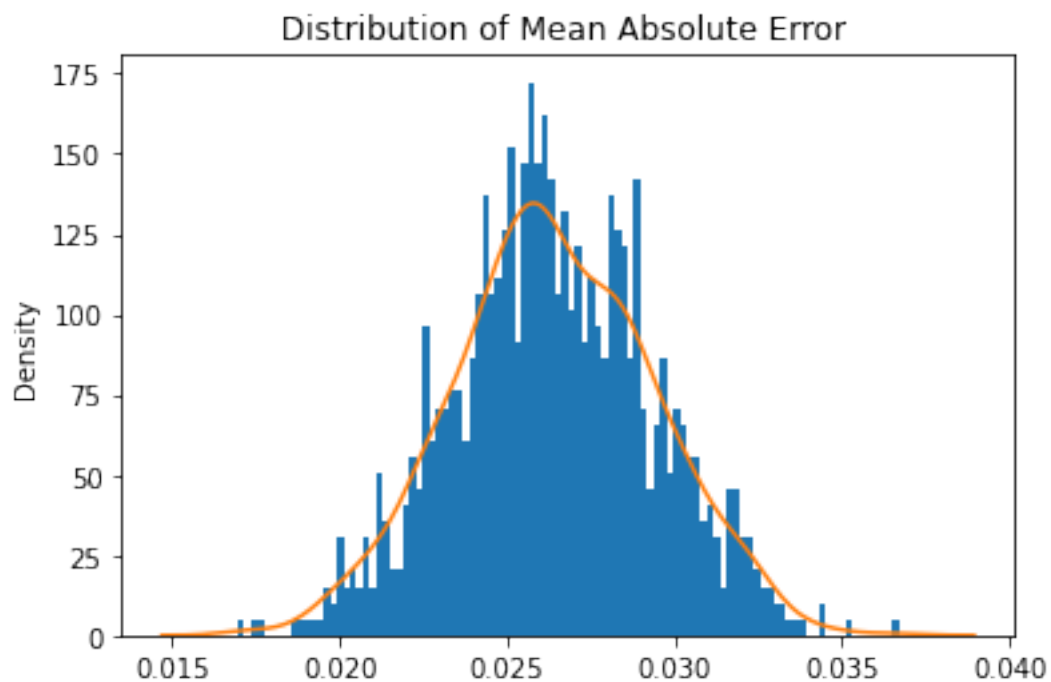
```
[23]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```





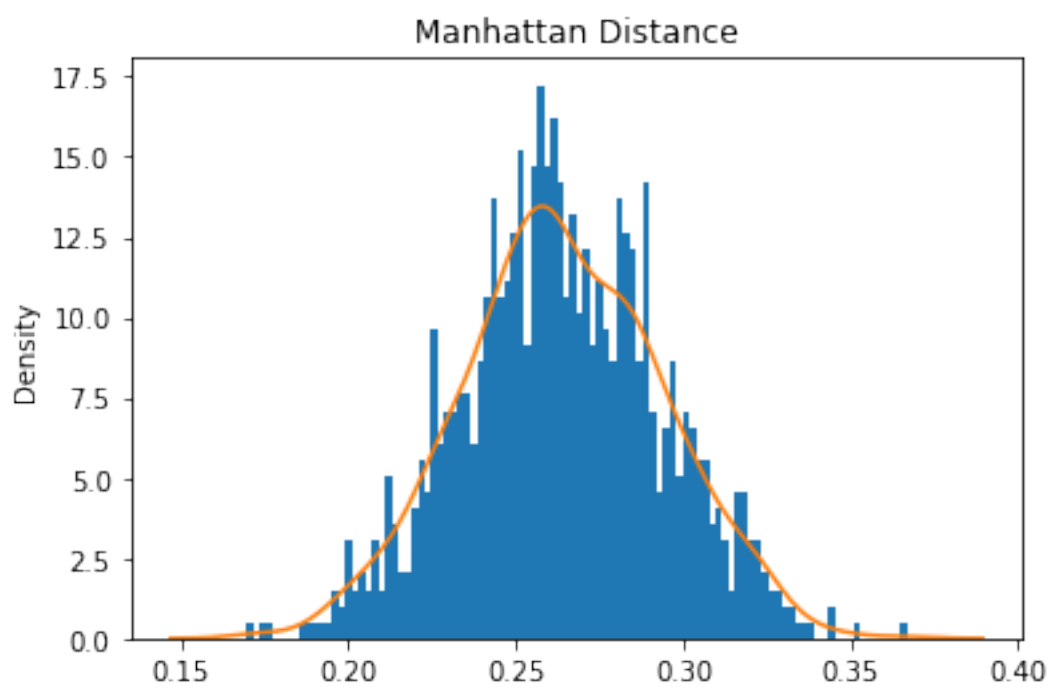


Mean Square Error: 0.0013581491915936777

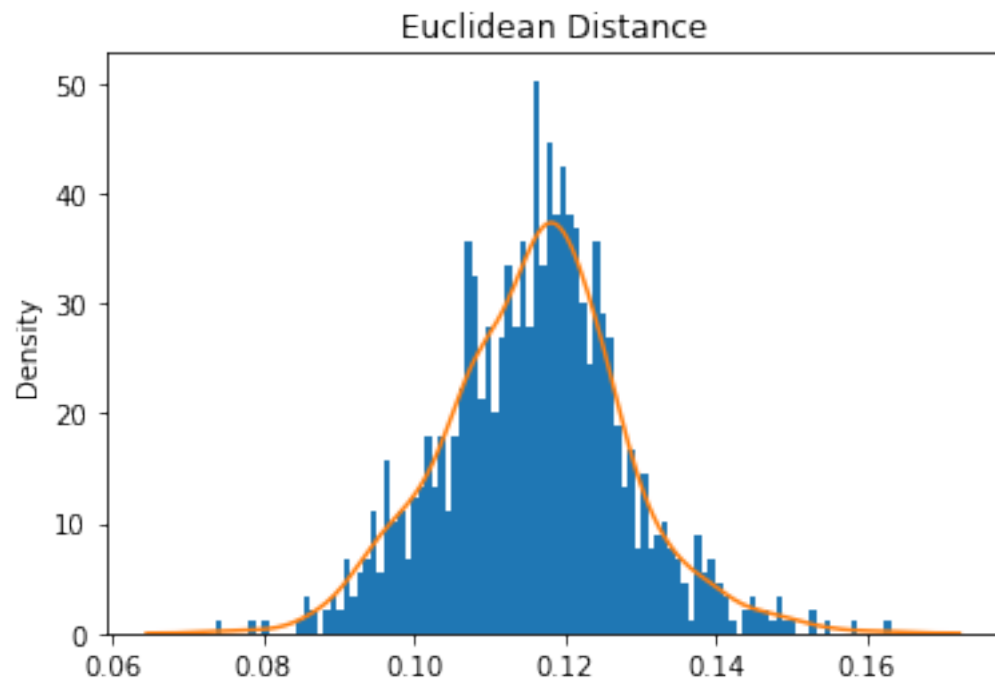


Mean Absolute Error: 0.026388249799609183

Mean Manhattan Distance: 0.26388249799609187



Mean Euclidean Distance: 0.11592284000209037



[]: