

Dataset2_Friedman1_output_10

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 = 100
     n_features = 10
     mean = 0
     variance = 0.1

```

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	0.713345	0.694135	0.109061	0.660299	0.216982	0.121298	0.753592
1	0.662808	0.419202	0.694175	0.982231	0.498710	0.675789	0.773940
2	0.933667	0.489171	0.826761	0.748329	0.009052	0.882996	0.971777

```

3  0.182255  0.936999  0.506063  0.947099  0.812895  0.053196  0.103375
4  0.701398  0.017712  0.503944  0.767280  0.638988  0.415207  0.261121

```

```

          X7          X8          X9          Y
0  0.954970  0.382936  0.393201  20.932255
1  0.920067  0.610729  0.849193  20.808128
2  0.234279  0.536423  0.557053  19.562277
3  0.025612  0.075345  0.596363  18.748718
4  0.233957  0.078030  0.006233  11.256395

```

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:                0.793
Model:                          OLS    Adj. R-squared:           0.770
Method:                        Least Squares    F-statistic:           34.08
Date:                          Wed, 20 Oct 2021    Prob (F-statistic):      3.03e-26
Time:                          20:34:06    Log-Likelihood:          -63.158
No. Observations:              100    AIC:                     148.3
Df Residuals:                   89    BIC:                     177.0
Df Model:                       10
Covariance Type:                nonrobust
=====
              coef      std err          t      P>|t|      [0.025      0.975]
-----
const      -3.261e-16      0.048    -6.76e-15      1.000      -0.096      0.096
x1           0.3578      0.050       7.102      0.000       0.258      0.458
x2           0.4031      0.051       7.961      0.000       0.303      0.504
x3          -0.0044      0.050      -0.090      0.929      -0.103      0.094
x4           0.5466      0.051      10.811      0.000       0.446      0.647
x5           0.2959      0.051       5.779      0.000       0.194      0.398
x6          -0.0061      0.051      -0.118      0.906      -0.108      0.096
x7           0.0158      0.052       0.304      0.762      -0.087      0.119
x8           0.0878      0.050       1.766      0.081      -0.011      0.186
x9           0.0746      0.050       1.503      0.136      -0.024      0.173
x10          0.0626      0.049       1.270      0.207      -0.035      0.160
=====
Omnibus:                 35.519    Durbin-Watson:           1.959
Prob(Omnibus):            0.000    Jarque-Bera (JB):        102.801
Skew:                    -1.210    Prob(JB):                 4.75e-23
Kurtosis:                 7.338    Cond. No.                  1.61
=====

```

Notes:

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

Parameters: const -3.261280e-16

x1 3.577852e-01

x2 4.031217e-01

x3 -4.446373e-03

x4 5.465560e-01

x5 2.958617e-01

x6 -6.067945e-03

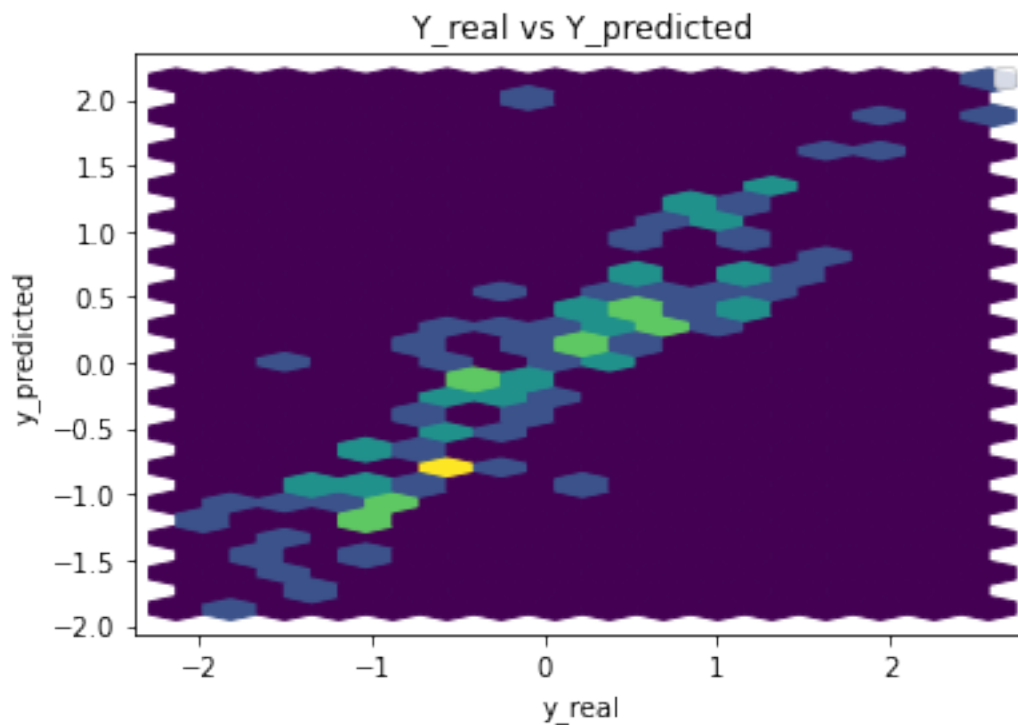
x7 1.580033e-02

x8 8.776097e-02

x9 7.457342e-02

x10 6.255050e-02

dtype: float64



Performance Metrics

Mean Squared Error: 0.20706486582455558

Mean Absolute Error: 0.32780247948845537

Manhattan distance: 32.780247948845535

Euclidean distance: 4.55043806489612

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

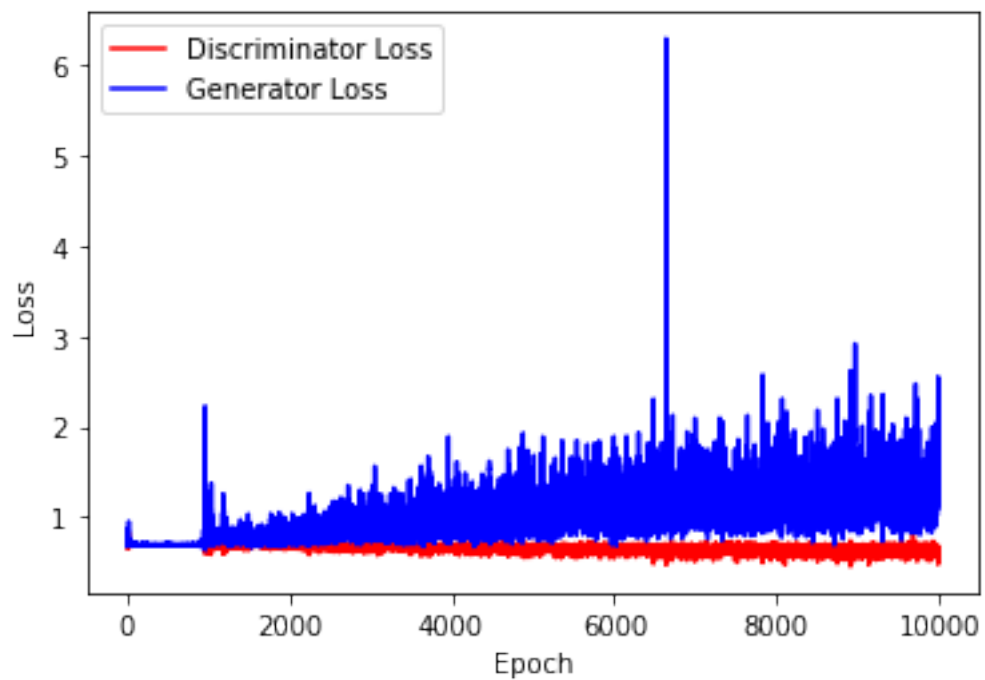
```
[9]: generator = network.Generator(n_features+2)
     discriminator = network.Discriminator(n_features+2)

     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))
```

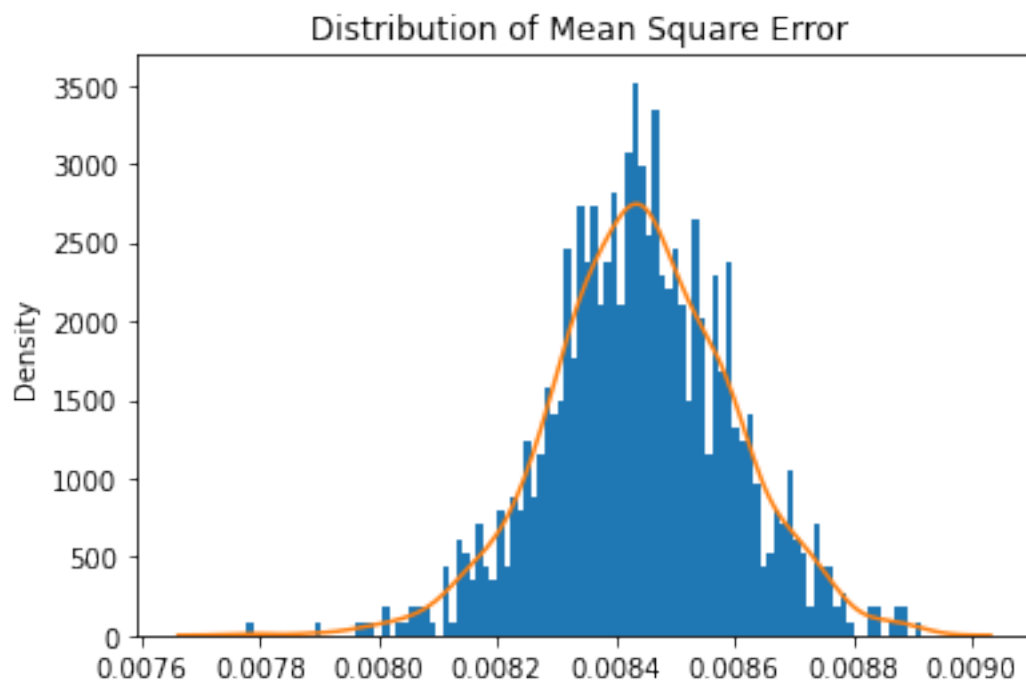
```
[10]: print(generator)
      print(discriminator)
```

```
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()
)
```

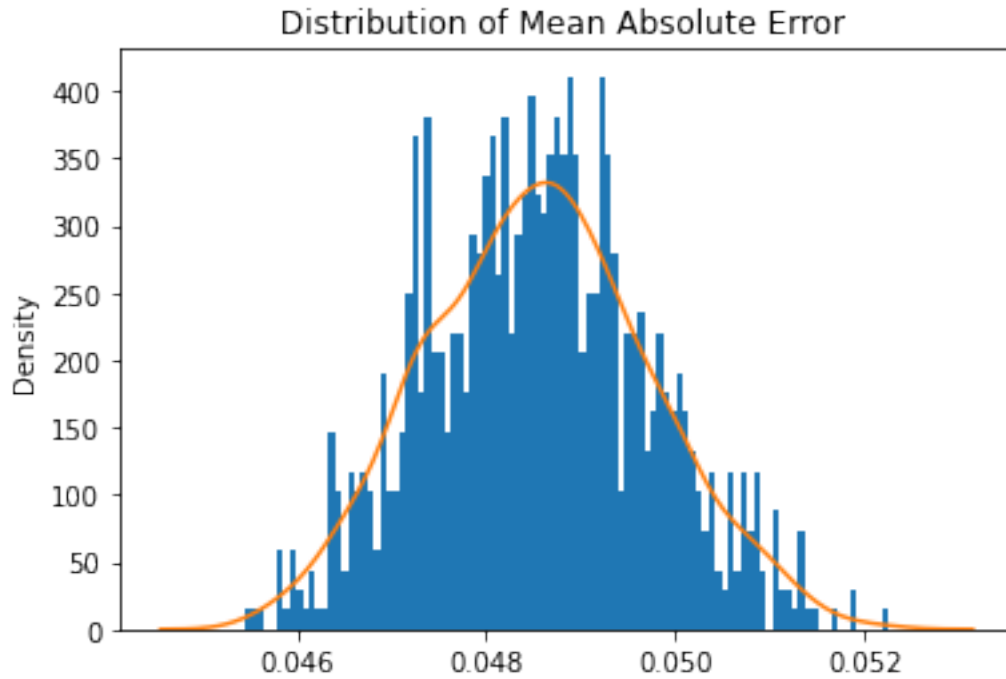
```
[11]: train_test.
     ↪training_GAN(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,
     ↪n_epochs,criterion,device)
```



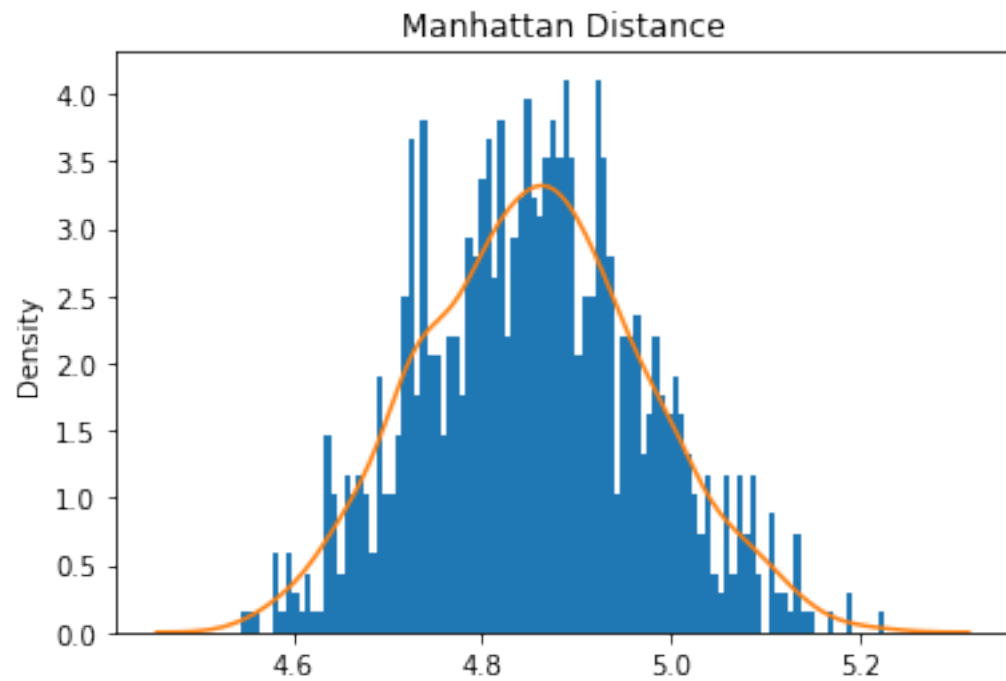
```
[12]: train_test.test_generator(generator,real_dataset,device)
```



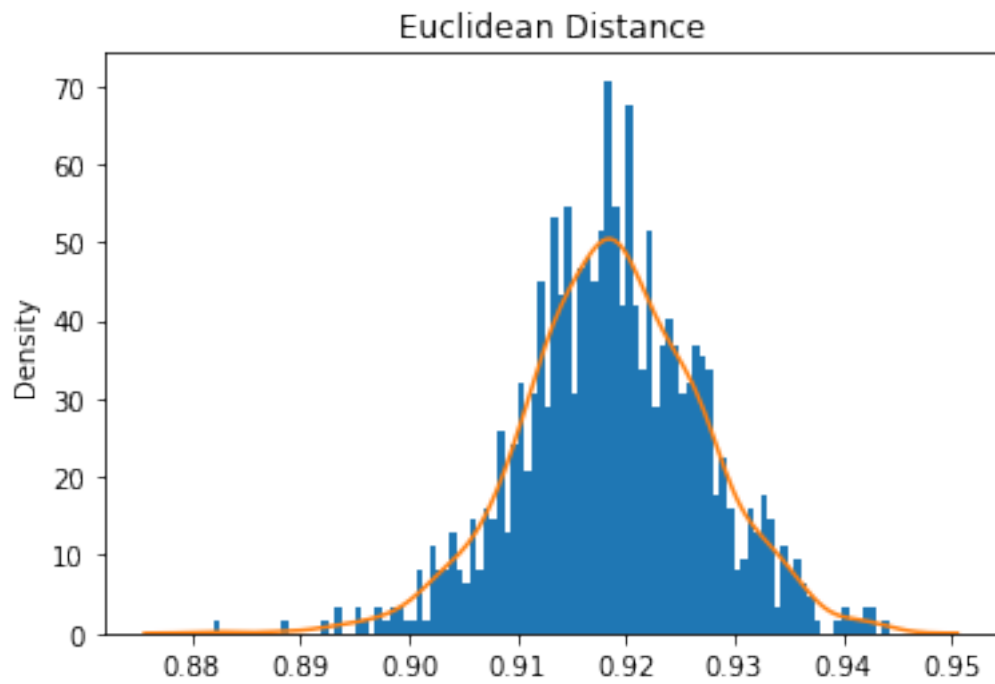
Mean Square Error: 0.008442594869358501



Mean Absolute Error: 0.048555512084439395

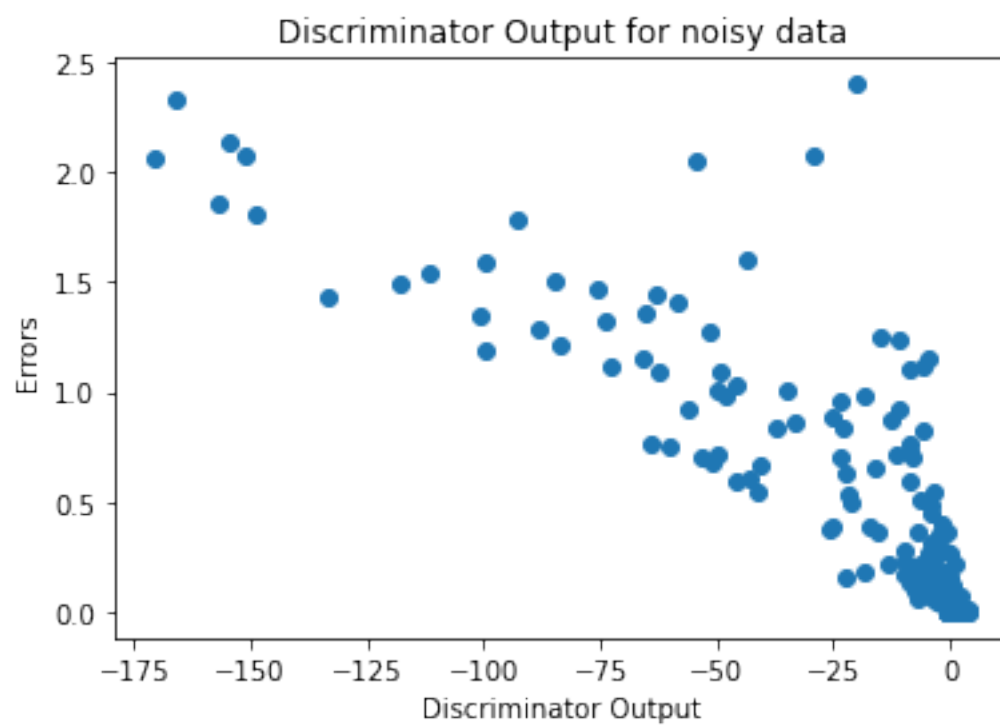
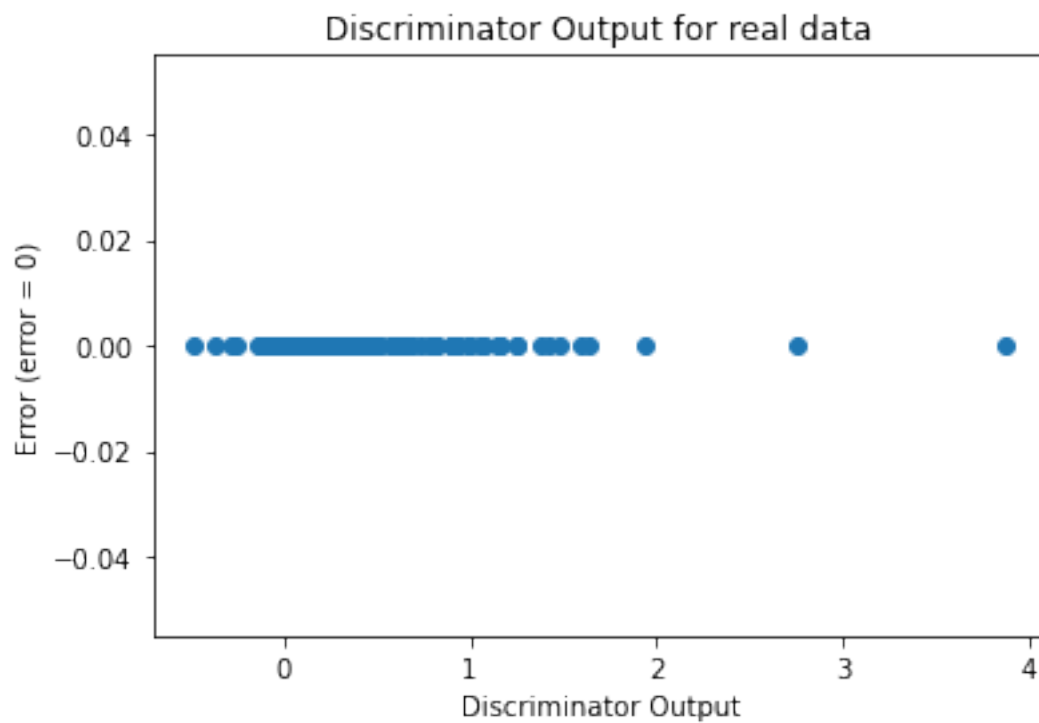


Mean Manhattan Distance: 4.85555120844394



Mean Euclidean Distance: 0.9187977216758765

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

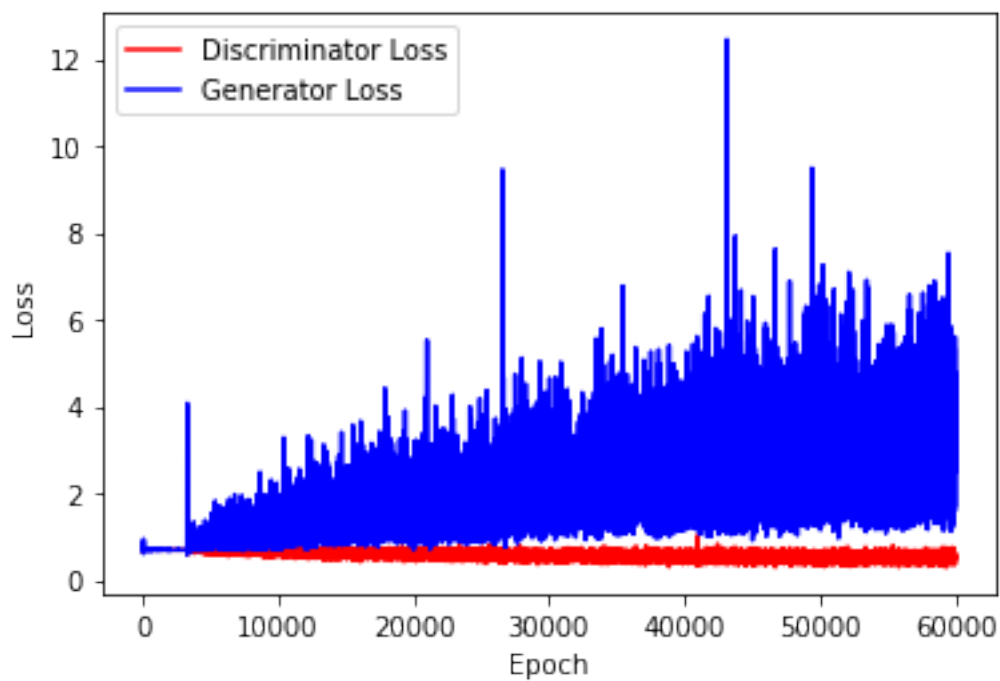



Training GAN until mse of y_pred is > 0.1 or $n_epochs < 30000$

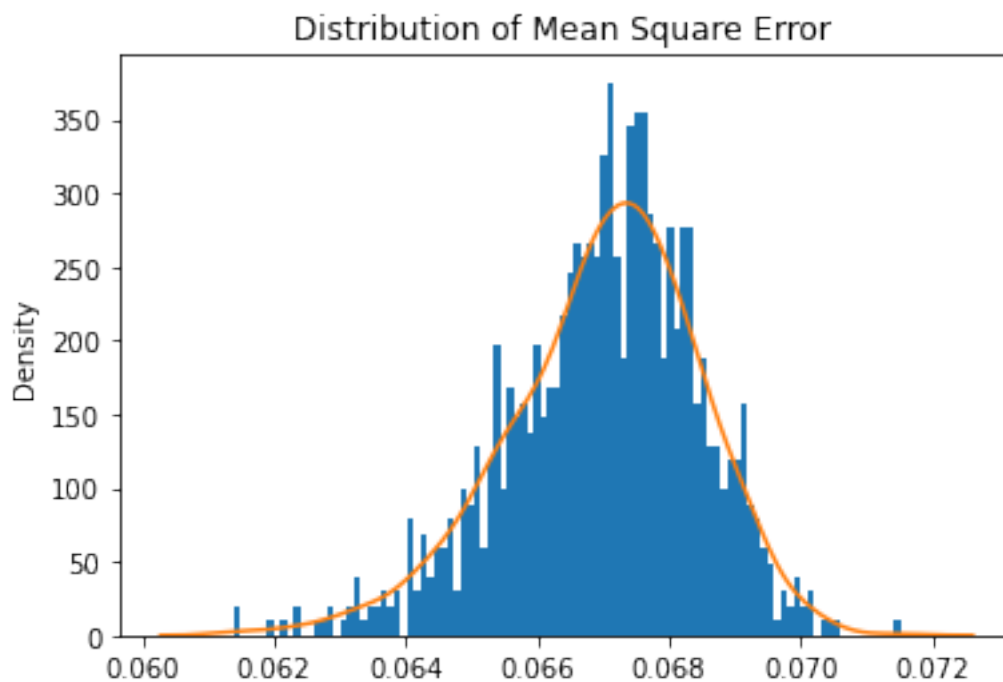
```
[14]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator.parameters(), lr=0.01, betas=(0.5,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

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

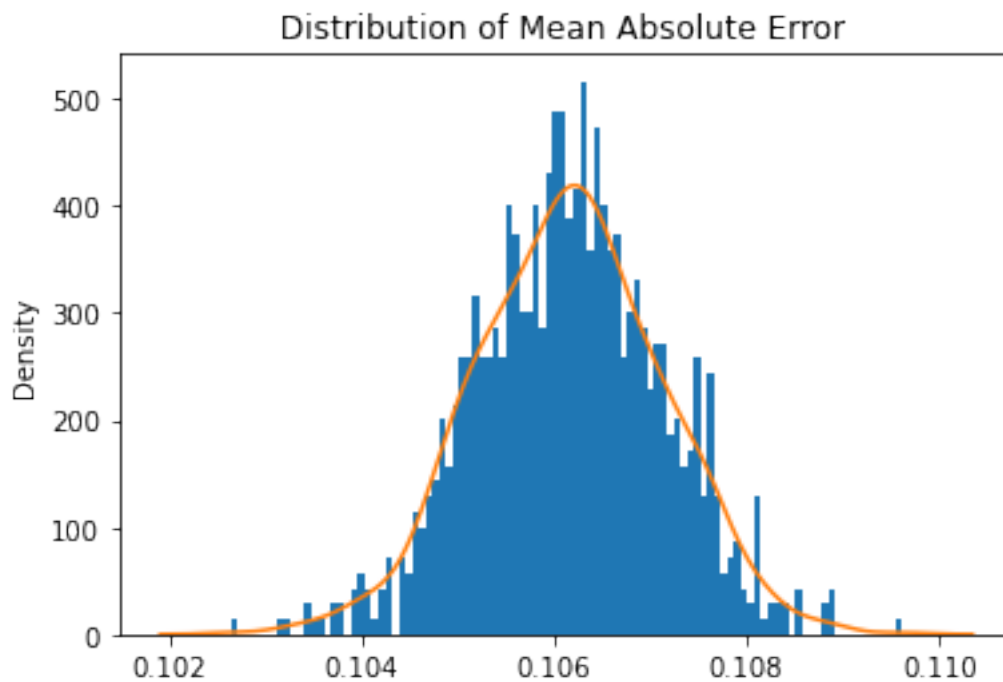
Number of epochs needed 30000



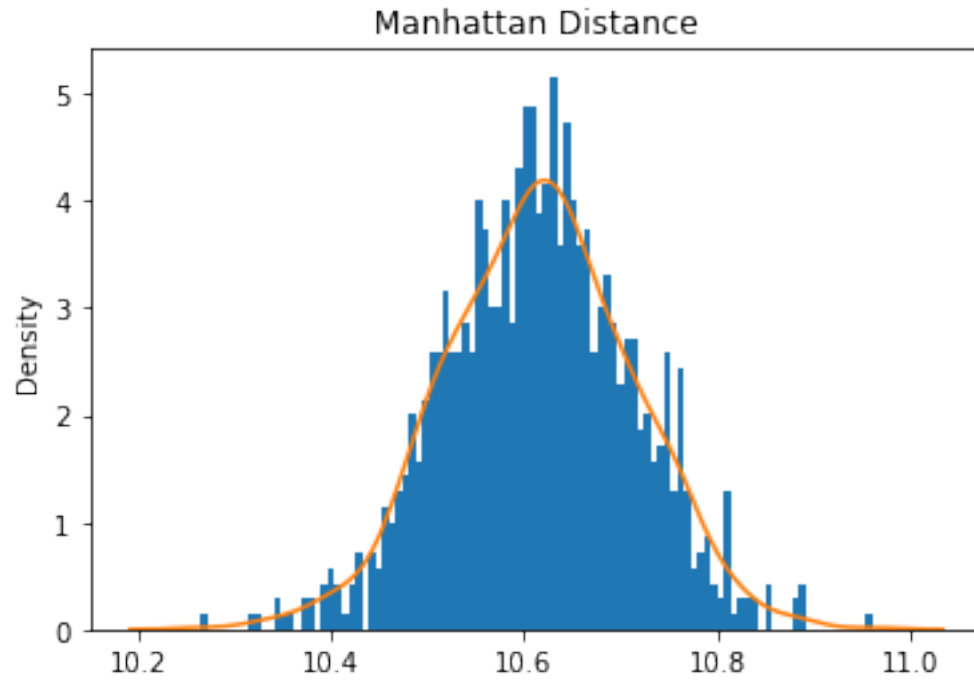
```
[16]: train_test.test_generator(generator,real_dataset,device)
```



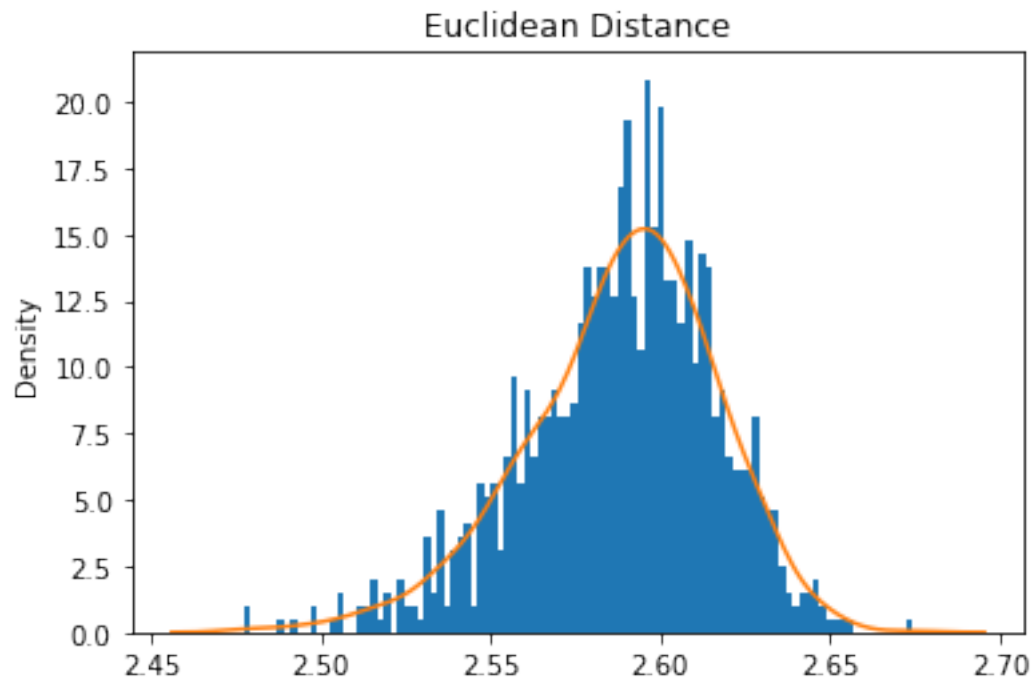
Mean Square Error: 0.06700097326587692



Mean Absolute Error: 0.10616071399521082



Mean Manhattan Distance: 10.616071399521083



Mean Euclidean Distance: 2.588300347312939

2 ABC GAN Model

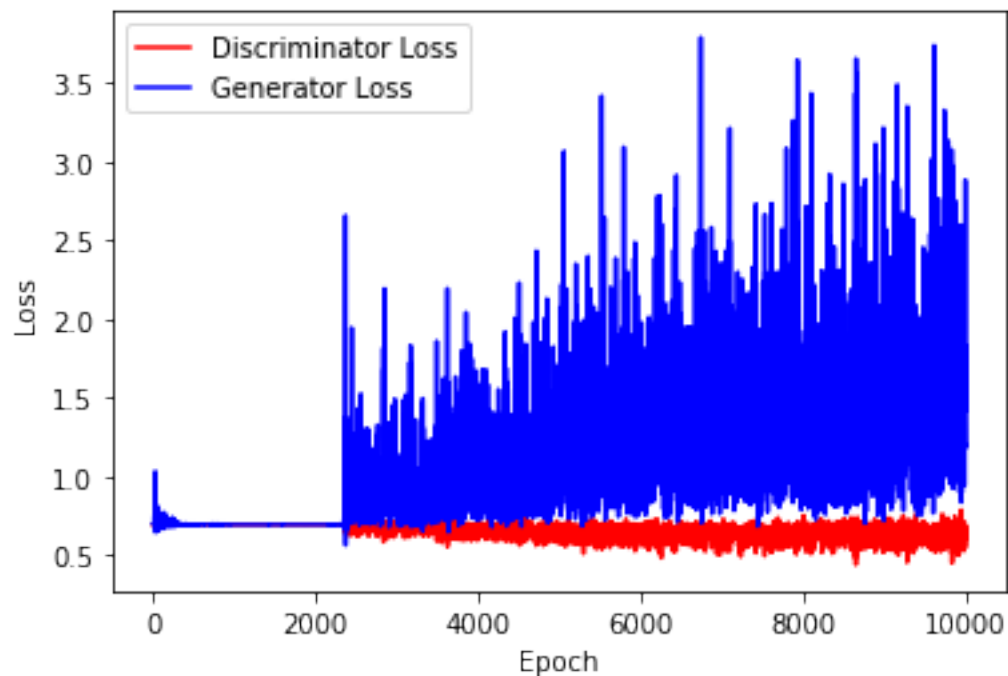
2.0.1 Training the network

Training ABC-GAN for `n_epochs` number of epochs

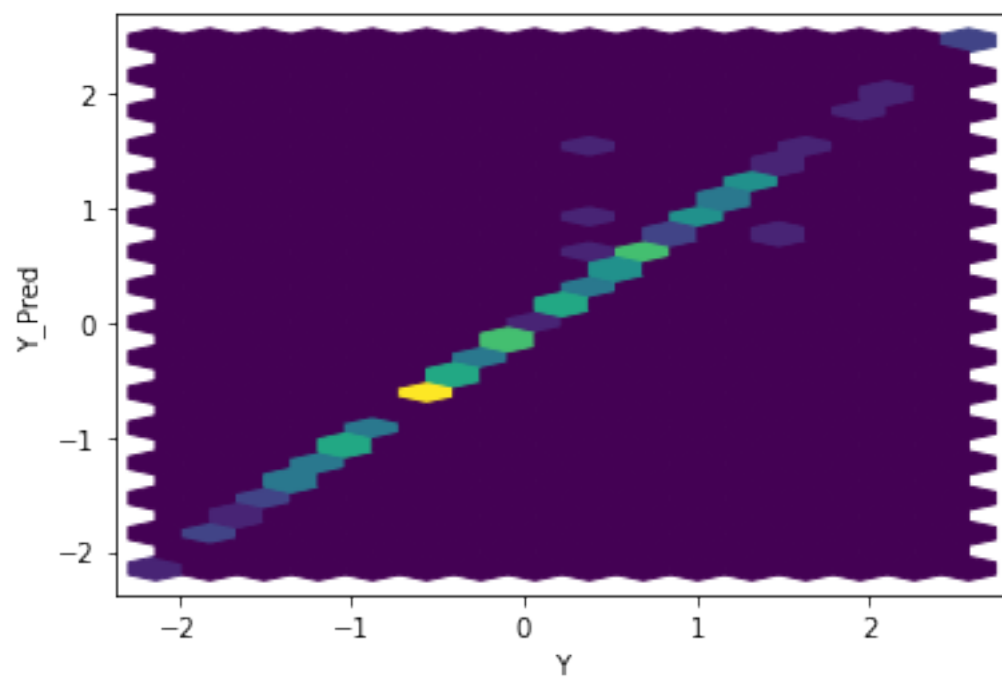
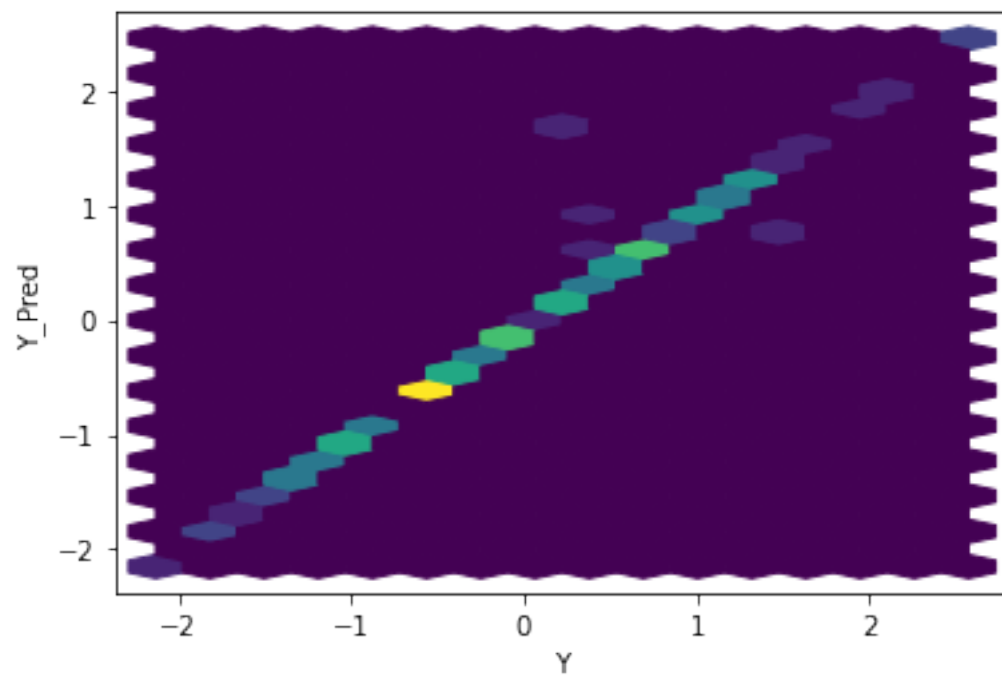
```
[17]: gen = network.Generator(n_features+2)
      disc = network.Discriminator(n_features+2)

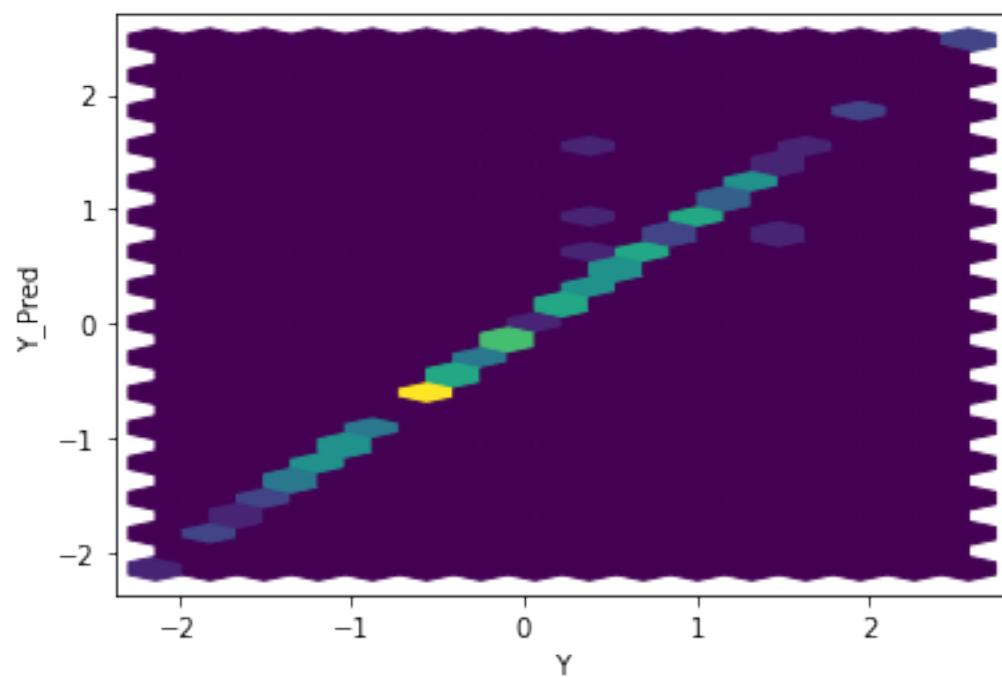
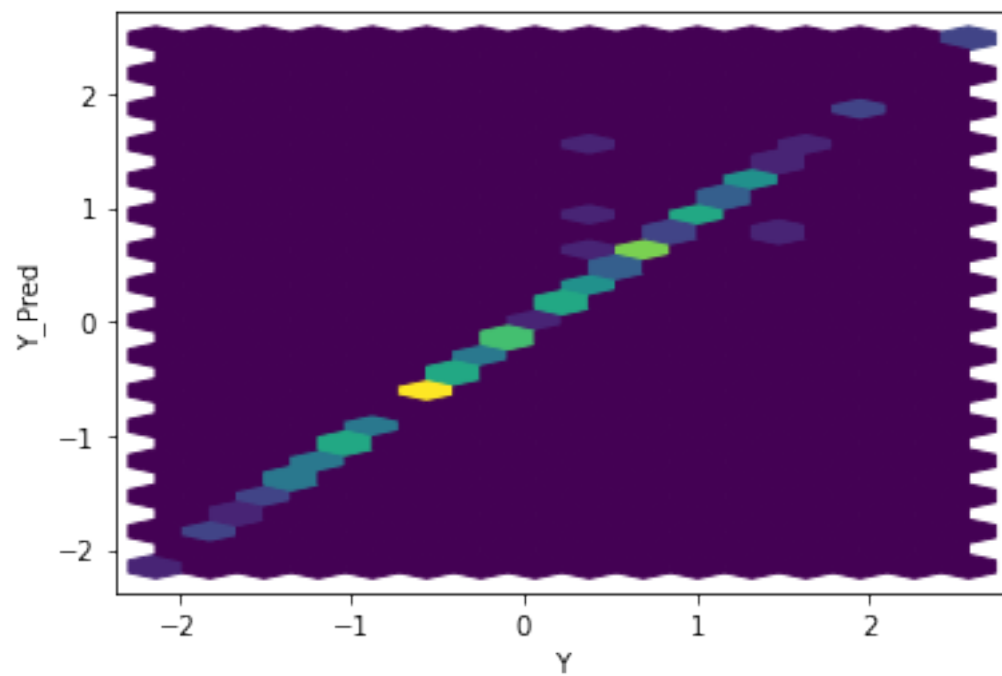
      criterion = torch.nn.BCEWithLogitsLoss()
      gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
      disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))

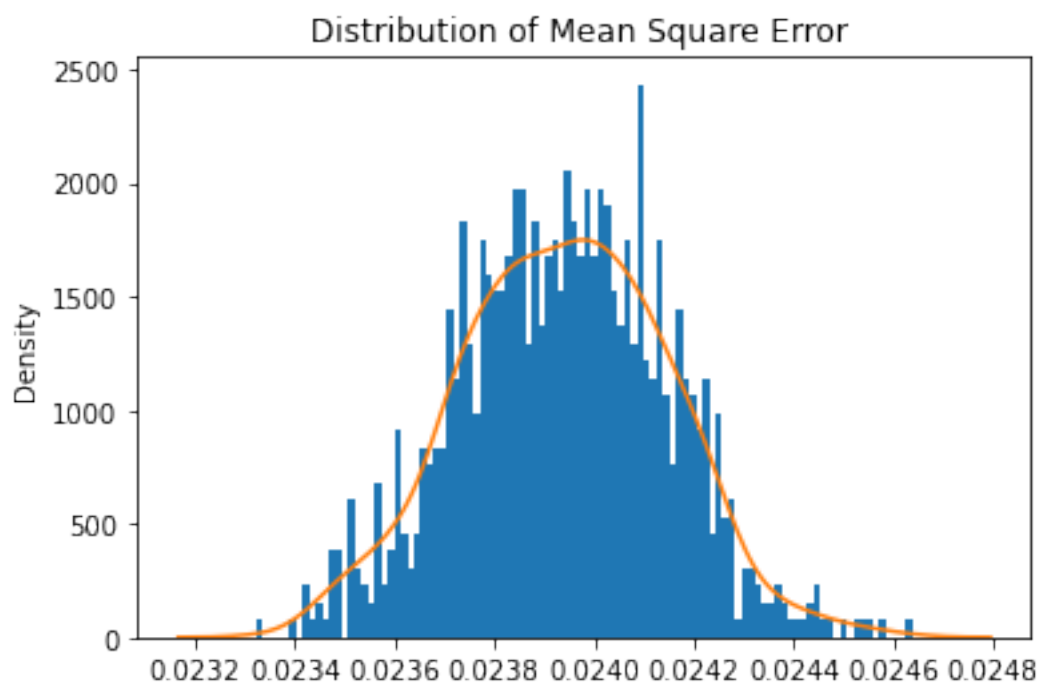
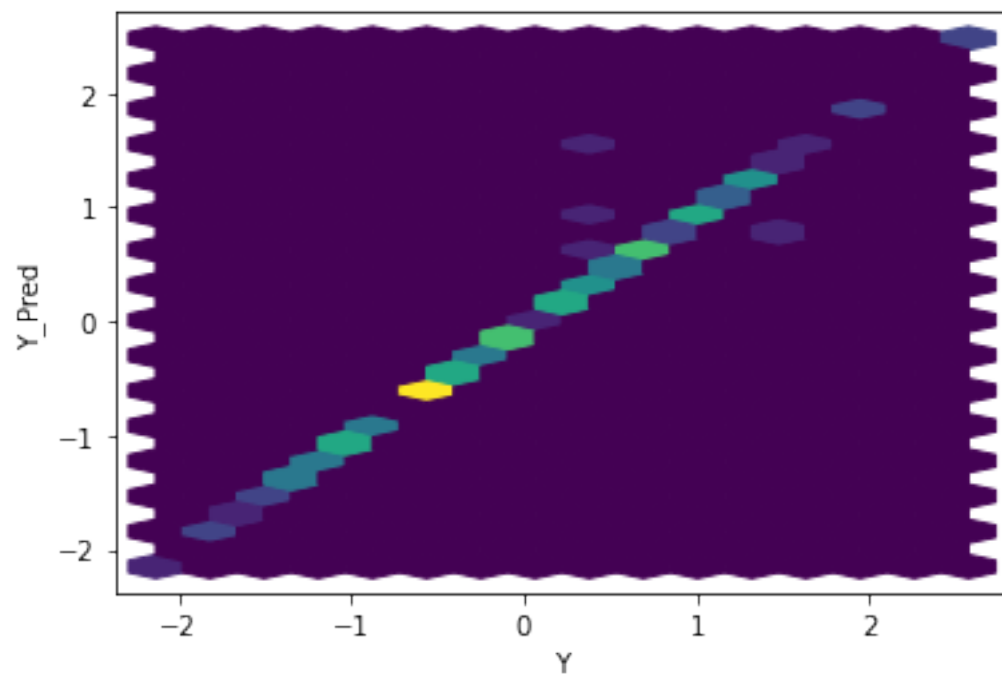
[18]: ABC_train_test.training_GAN(disc, gen,disc_opt,gen_opt,real_dataset,
      ↪batch_size, n_epochs,criterion,coeff,mean,variance,device)
```



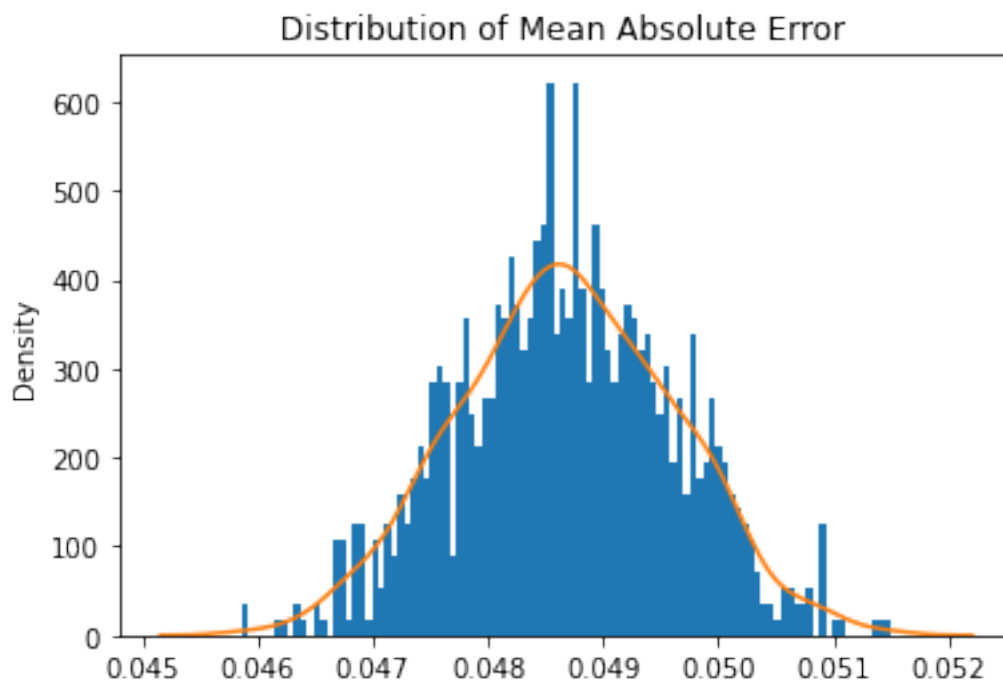
```
[19]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```



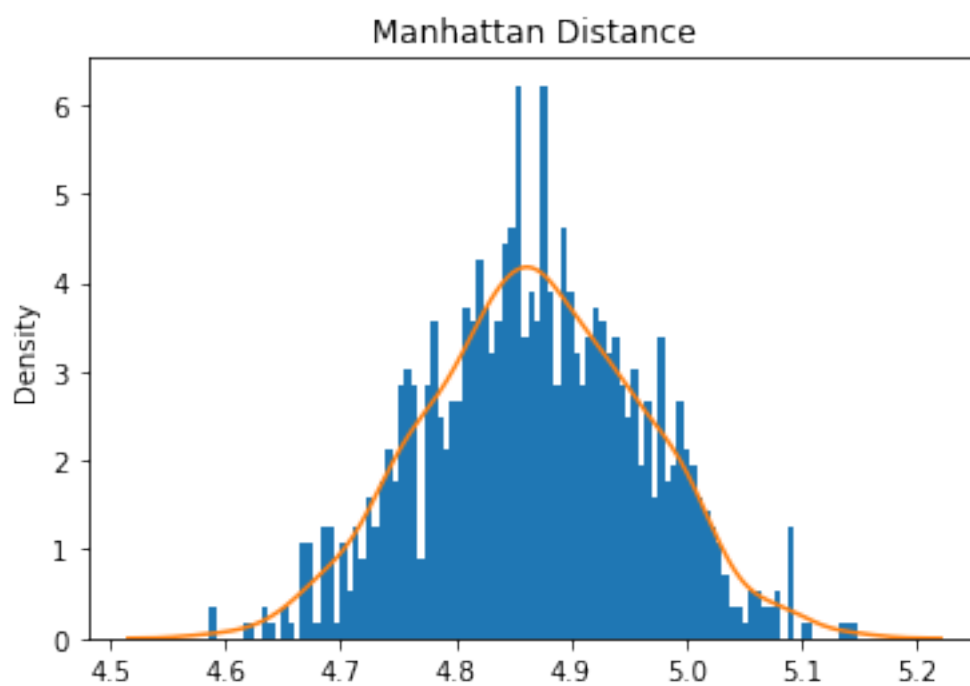




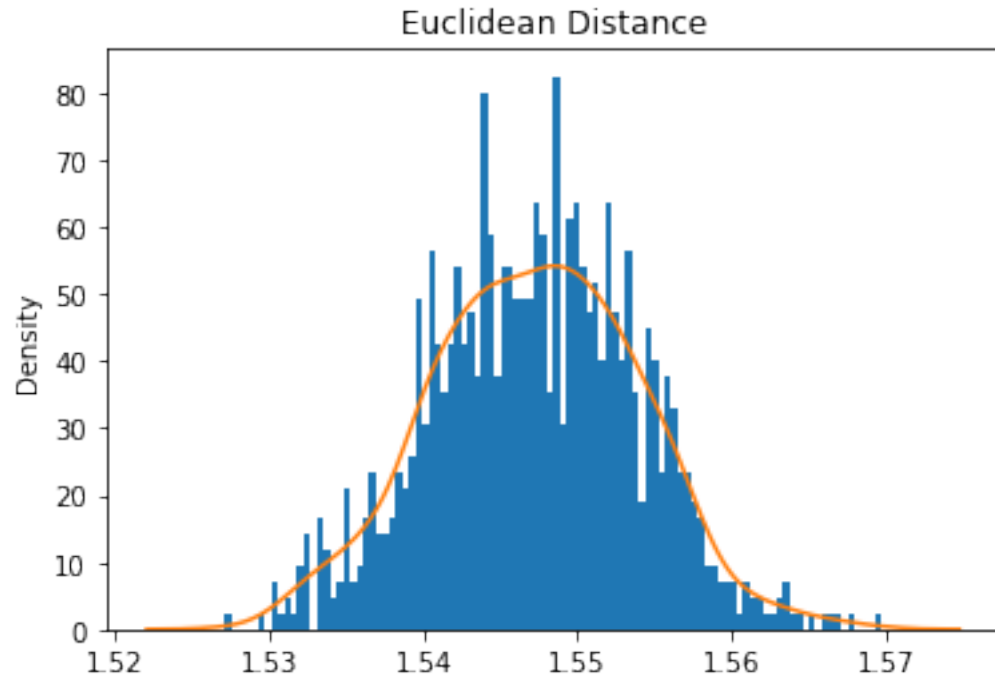
Mean Square Error: 0.0239390650589317



Mean Absolute Error: 0.04867806139353663
Mean Manhattan Distance: 4.867806139353663

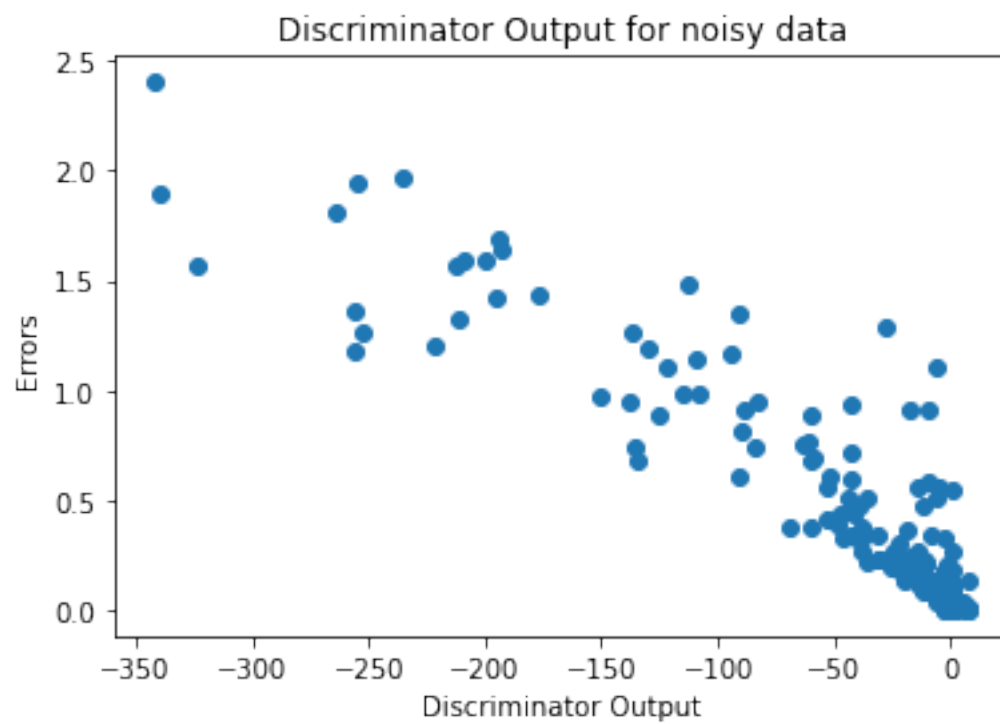
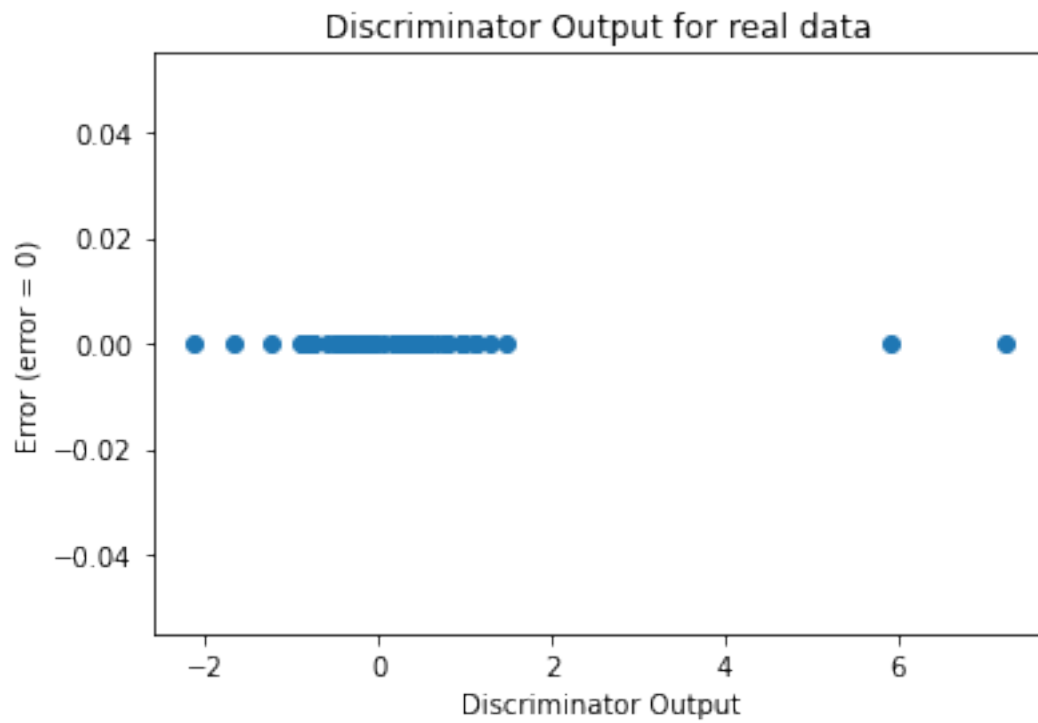


Mean Euclidean Distance: 1.5472107609014114



Sanity Checks

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



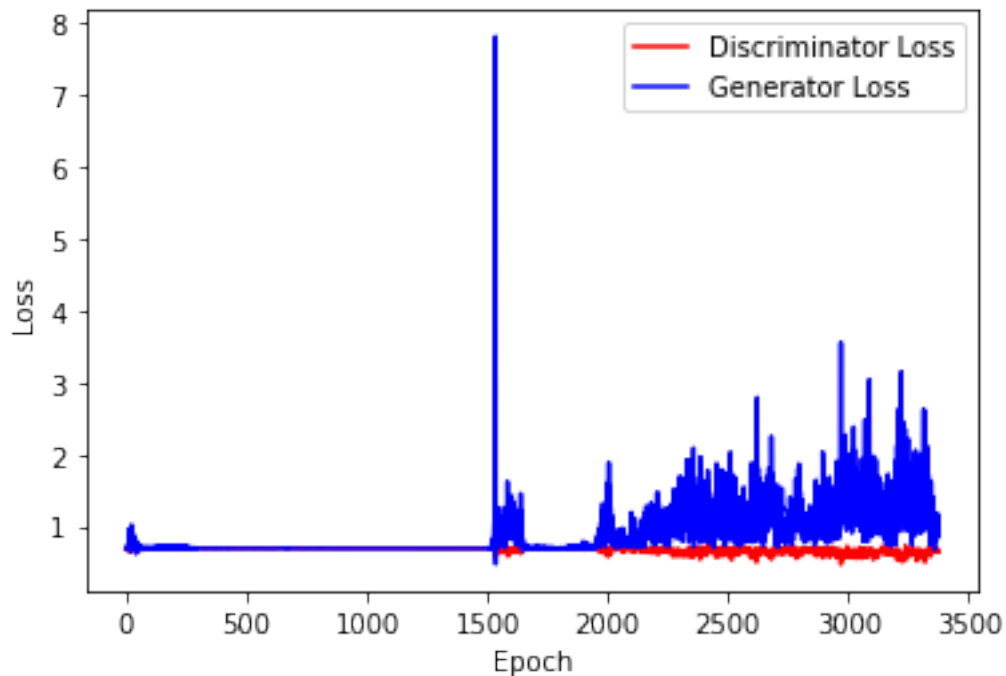
Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

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

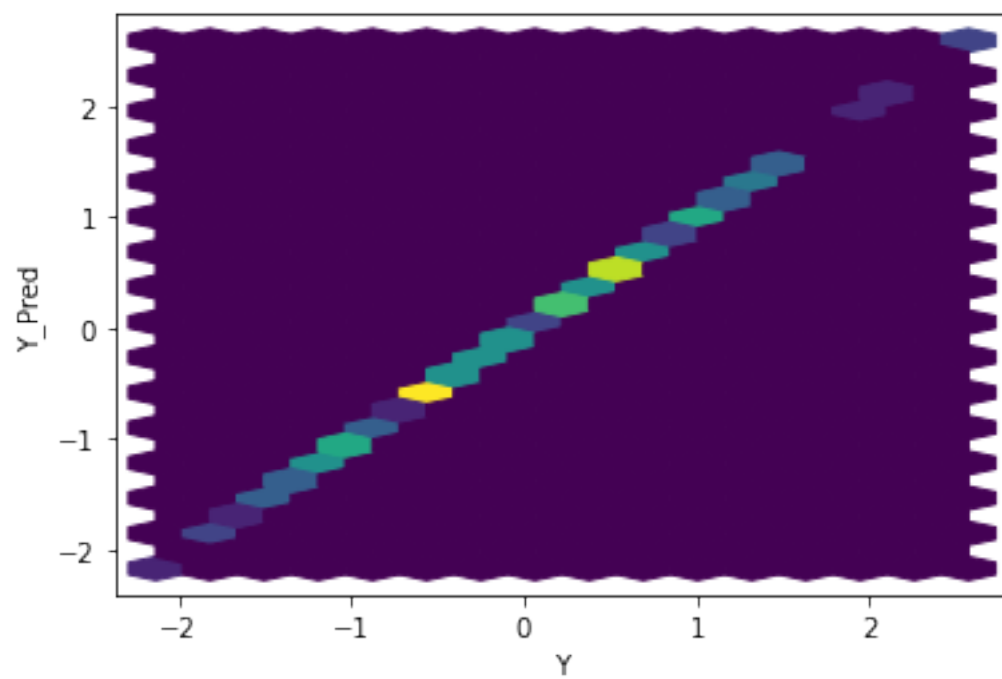
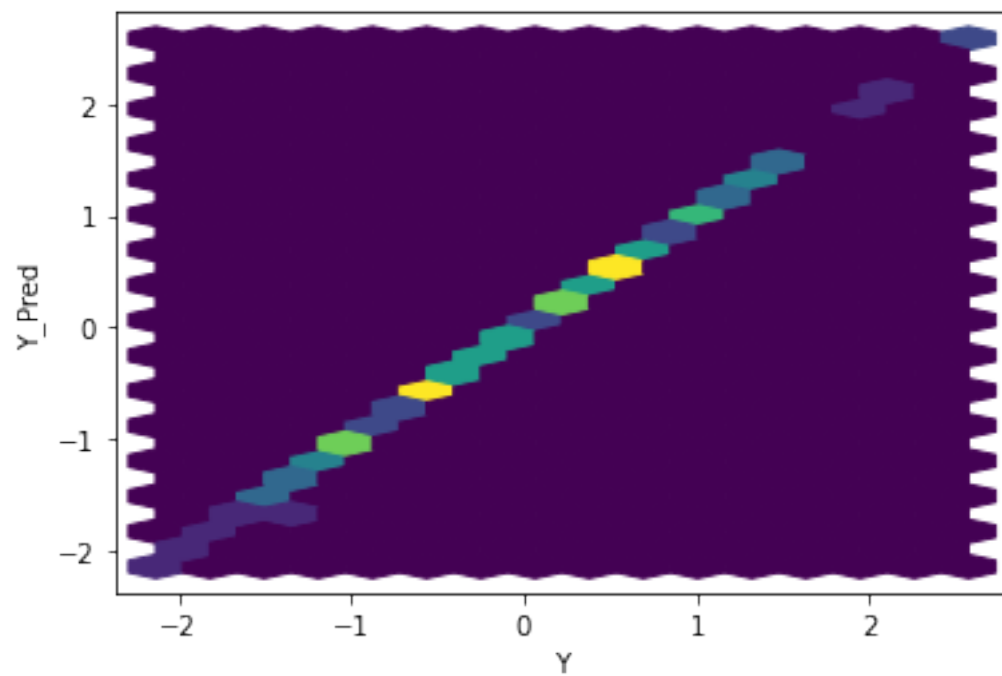
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))
```

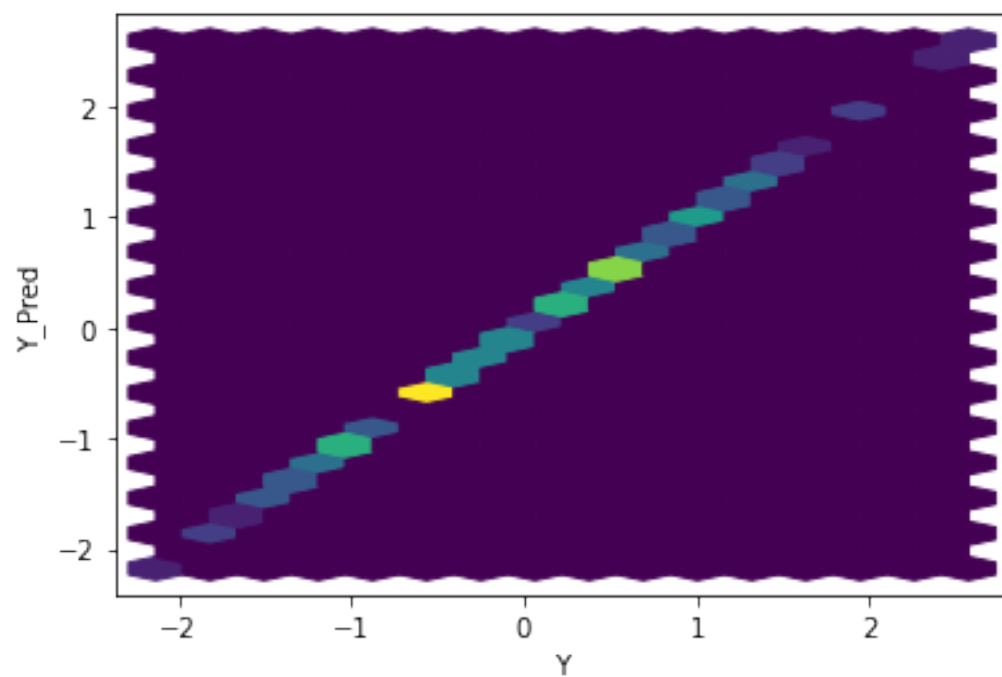
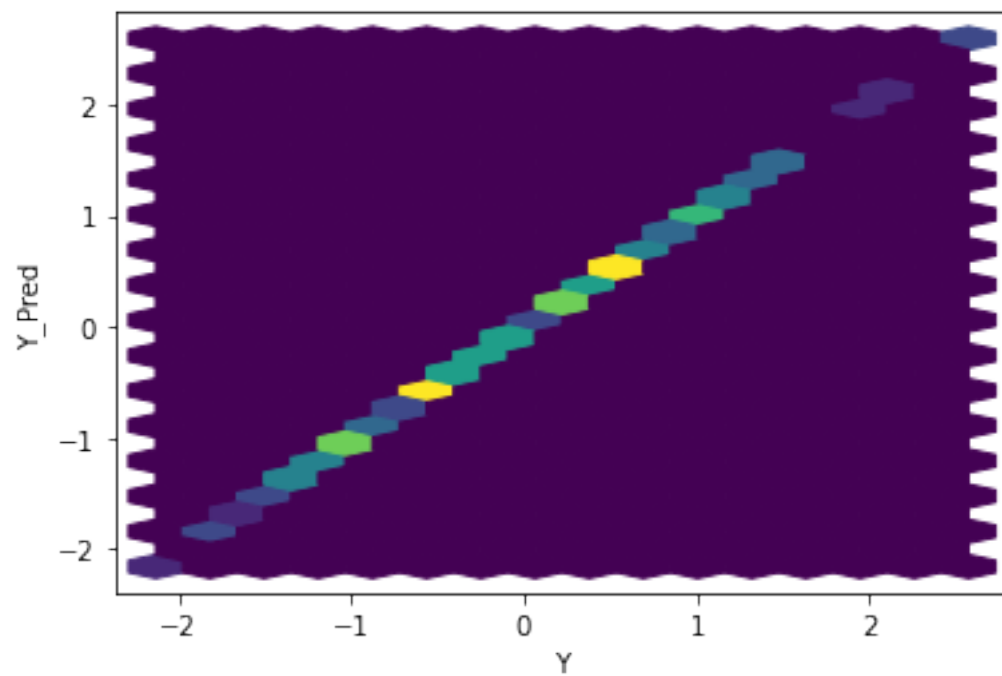
```
[22]: ABC_train_test.
      ↪ training_GAN_2(disc,gen,disc_opt,gen_opt,real_dataset,batch_size,
      ↪ error,criterion,coeff,mean,variance,device)
```

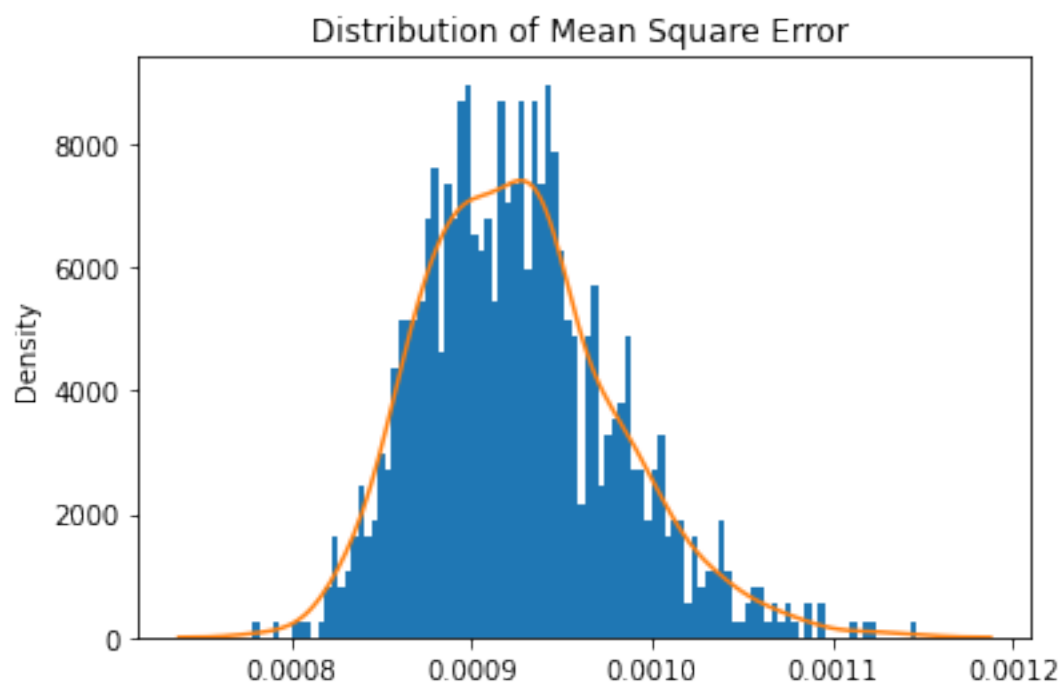
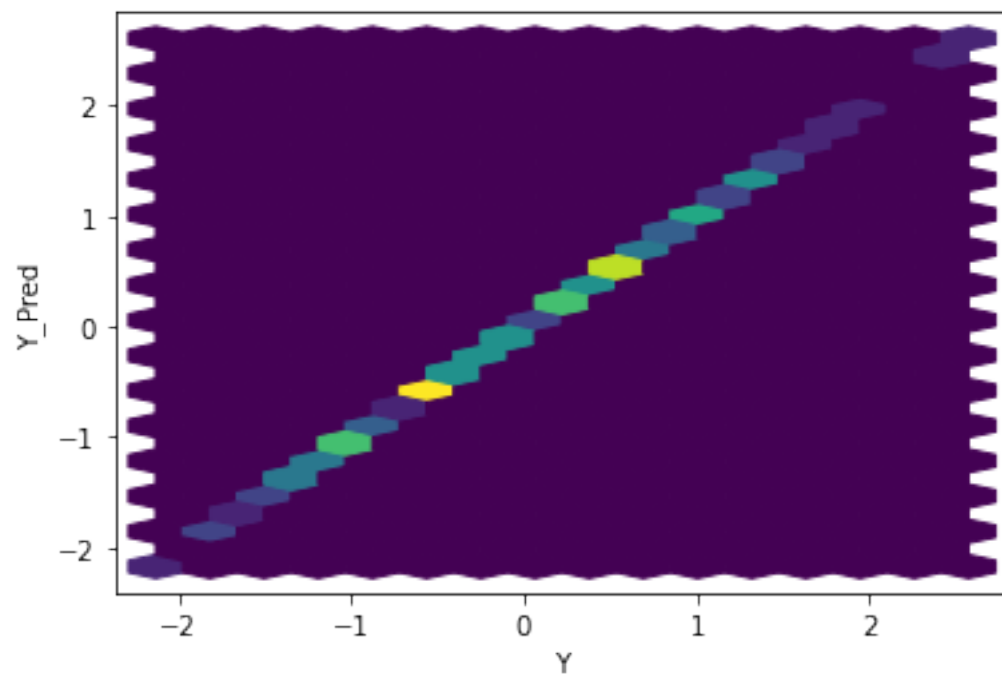
Number of epochs 1688



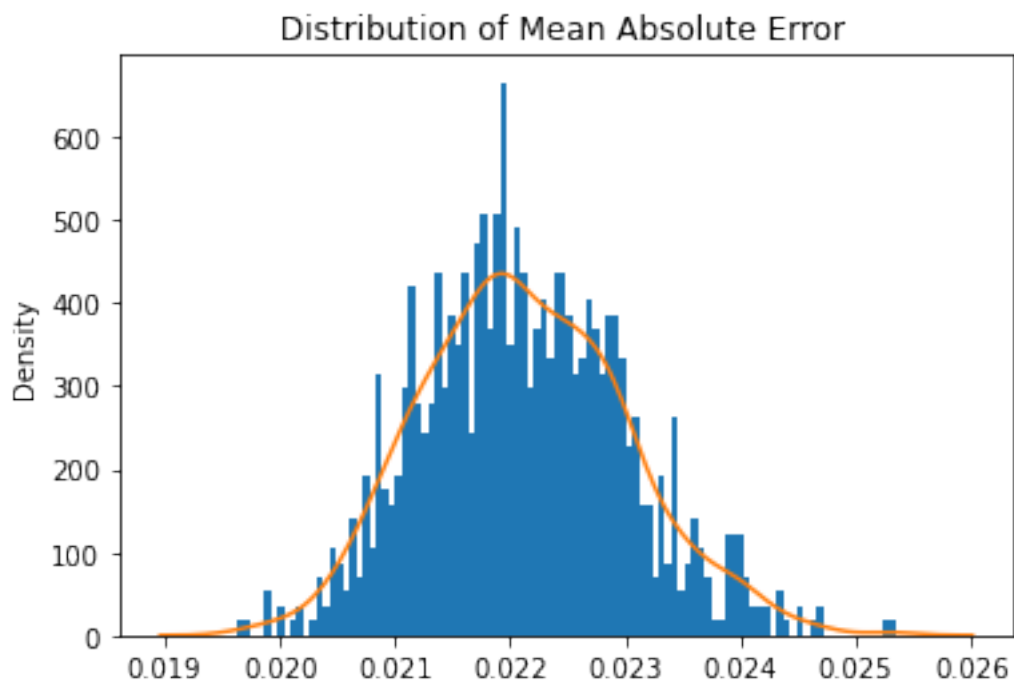
```
[23]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```





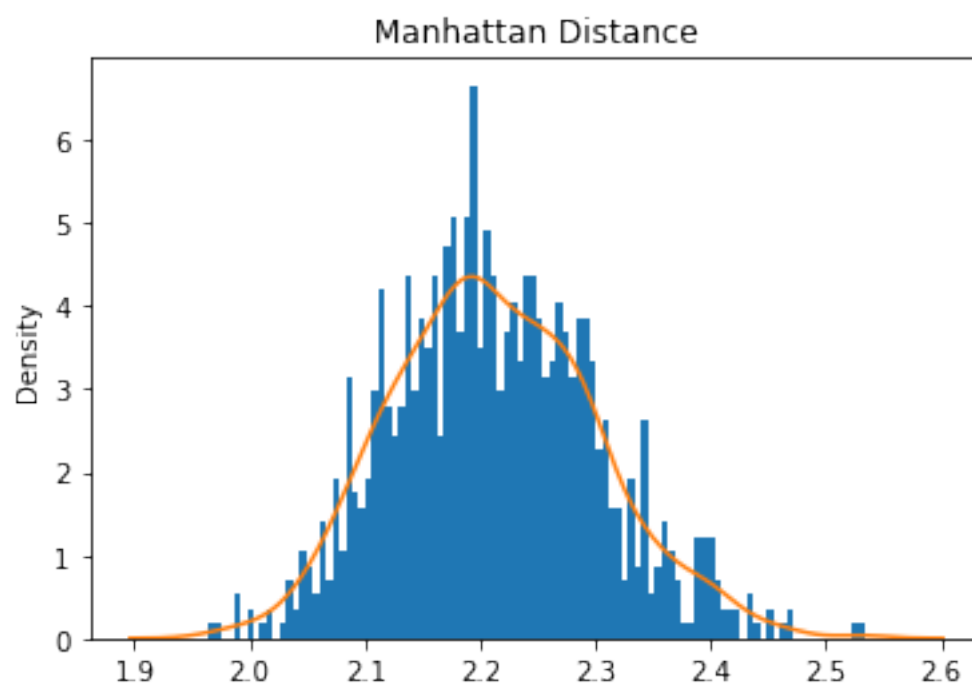


Mean Square Error: 0.0009258706826262589

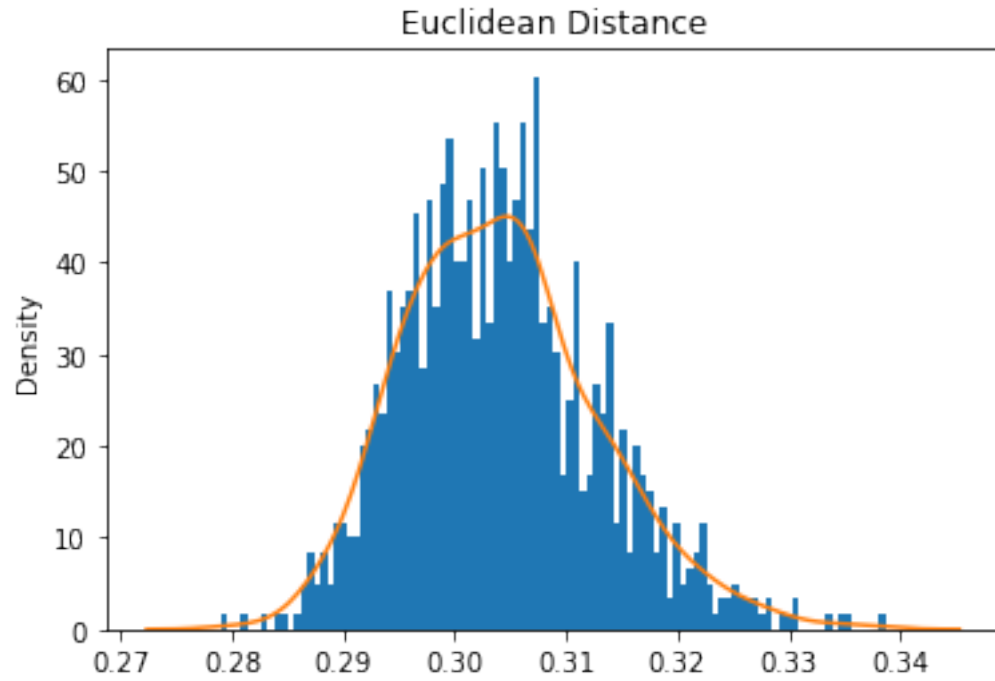


Mean Absolute Error: 0.022123692467100917

Mean Manhattan Distance: 2.2123692467100917



Mean Euclidean Distance: 0.30415406292519875



[]: