

# Dataset1-Regression\_output\_8

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.  $Y = 1 + \beta_1 x_1 + \beta_2 x_2 + \beta_2 x_3 + \dots + \beta_n x_n + N(0, \sigma)$  where  $\sigma = 0.1$
    2.  $\beta_i \sim N(0, \sigma^*)$  when  $\mu = 0$  else  $\beta_i \sim N(\beta_i^*, \sigma^*)$  where  $\beta_i^*$ s are coefficients obtained from statistical model
    3.  $\sigma^*$  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  $\beta^*$  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 = 1
     variance = 0.01

```

### 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	1.025662	-1.899881	0.722342	0.373928	0.025516	-0.286861	-0.038484
1	0.140417	-1.629063	-0.086283	1.051090	-1.230093	1.225208	0.874823
2	0.745632	-1.013213	-1.393919	0.481407	0.131400	1.014415	0.651640
3	-0.733528	0.191755	0.154356	-1.515971	-1.096288	0.878339	0.038619
4	0.732404	-0.110956	1.100714	-1.508329	-0.110869	-1.705385	2.132598
	X8	X9	X10	Y			

```

0  0.642875  1.279171  0.316527  111.964030
1 -0.742069 -2.048548  0.353430  -49.482990
2  0.976793 -0.413256 -0.888886  129.092016
3 -0.453928  0.052814 -1.198723 -229.953197
4  1.104004  0.737947 -0.837665   6.645053

```

## 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:              6.456e+07
Date:                Tue, 19 Oct 2021    Prob (F-statistic):      1.03e-300
Time:                  23:24:32    Log-Likelihood:          647.96
No. Observations:          100      AIC:                      -1274.
Df Residuals:              89      BIC:                      -1245.
Df Model:                  10
Covariance Type:            nonrobust
=====
               coef      std err          t      P>|t|      [0.025      0.975]
-----
const          9.714e-17  3.94e-05  2.47e-12    1.000   -7.82e-05   7.82e-05
x1              0.4249   4.11e-05  1.03e+04    0.000    0.425    0.425
x2              0.2455   4.04e-05  6073.384    0.000    0.245    0.246
x3              0.0028   4.23e-05   66.385    0.000    0.003    0.003
x4              0.2320   4.04e-05  5743.177    0.000    0.232    0.232
x5              0.3735   4.28e-05  8725.767    0.000    0.373    0.374
x6              0.4082   3.97e-05  1.03e+04    0.000    0.408    0.408
x7              0.4174   4.09e-05  1.02e+04    0.000    0.417    0.417
x8              0.1183    4.1e-05  2883.937    0.000    0.118    0.118
x9              0.2999   4.08e-05  7357.503    0.000    0.300    0.300
x10             0.3485   4.11e-05  8484.676    0.000    0.348    0.349
=====
Omnibus:                  4.870    Durbin-Watson:              2.256
Prob(Omnibus):            0.088    Jarque-Bera (JB):          4.598
Skew:                    0.525    Prob(JB):                  0.100
Kurtosis:                3.038    Cond. No.                  1.68
=====

```

Notes:

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

```

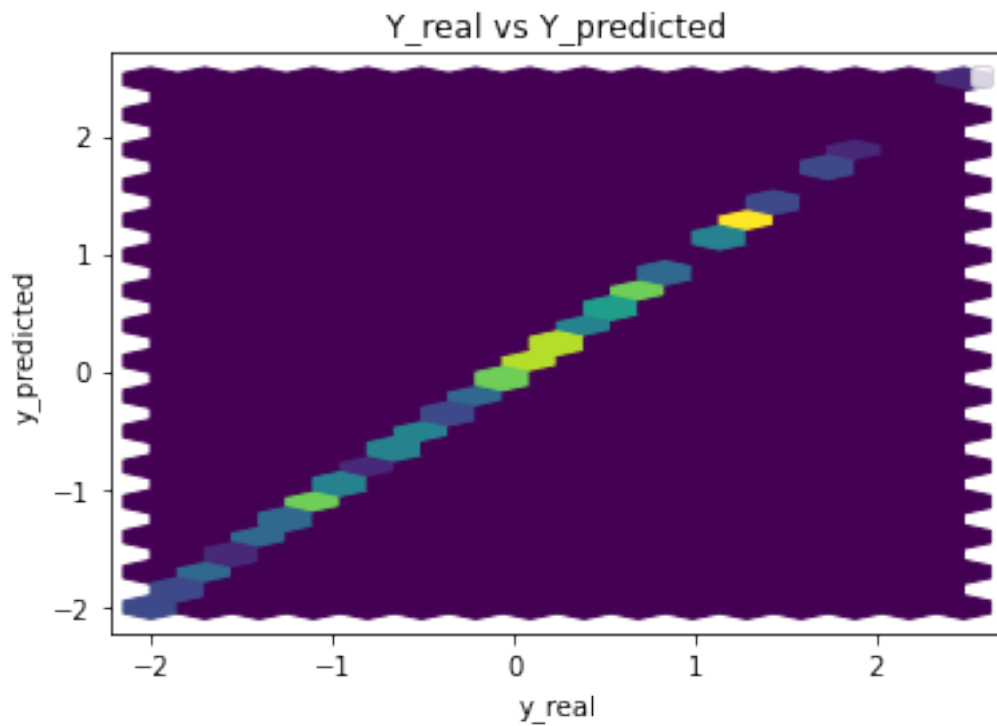
Parameters:  const      9.714451e-17
x1          4.249335e-01

```

```

x2      2.454511e-01
x3      2.808859e-03
x4      2.320214e-01
x5      3.735171e-01
x6      4.081686e-01
x7      4.173838e-01
x8      1.183379e-01
x9      2.998583e-01
x10     3.485329e-01
dtype: float64

```



#### Performance Metrics

```

Mean Squared Error: 1.3784920857399033e-07
Mean Absolute Error: 0.0002915312169664944
Manhattan distance: 0.02915312169664944
Euclidean distance: 0.003712804985102104

```

### 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

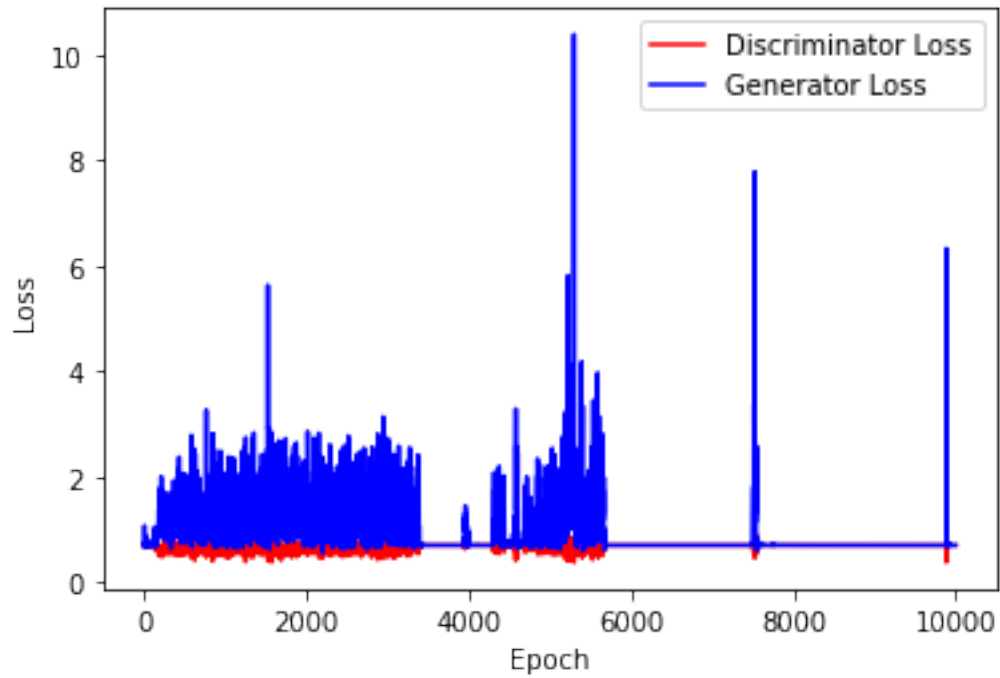
```
[9]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)

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))
```

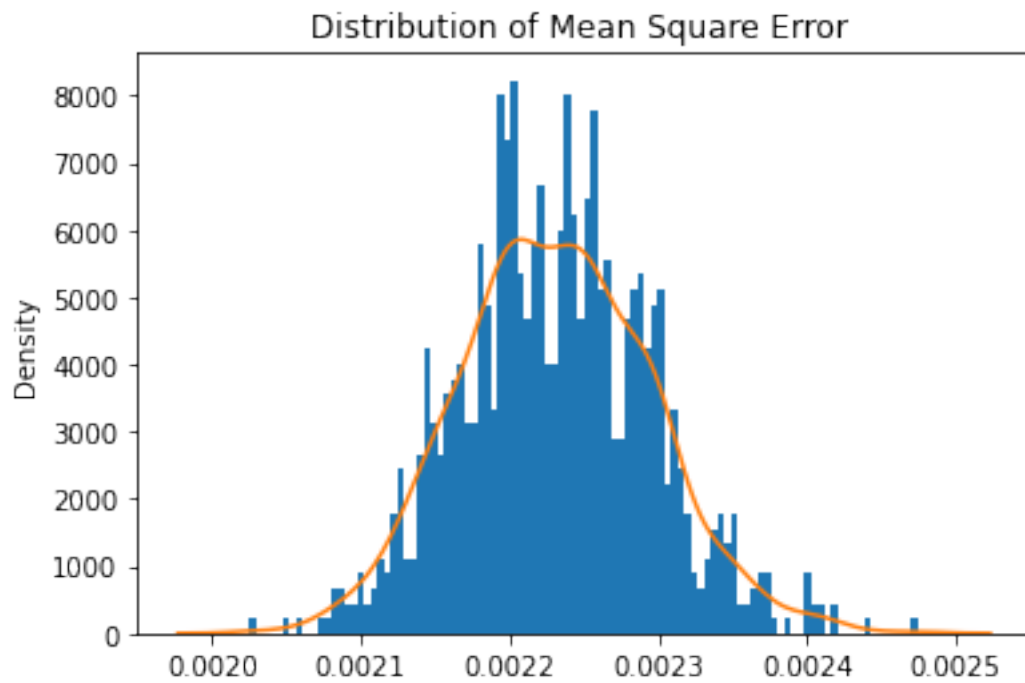
```
[10]: print(generator)
print(discriminator)
```

```
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()
)
```

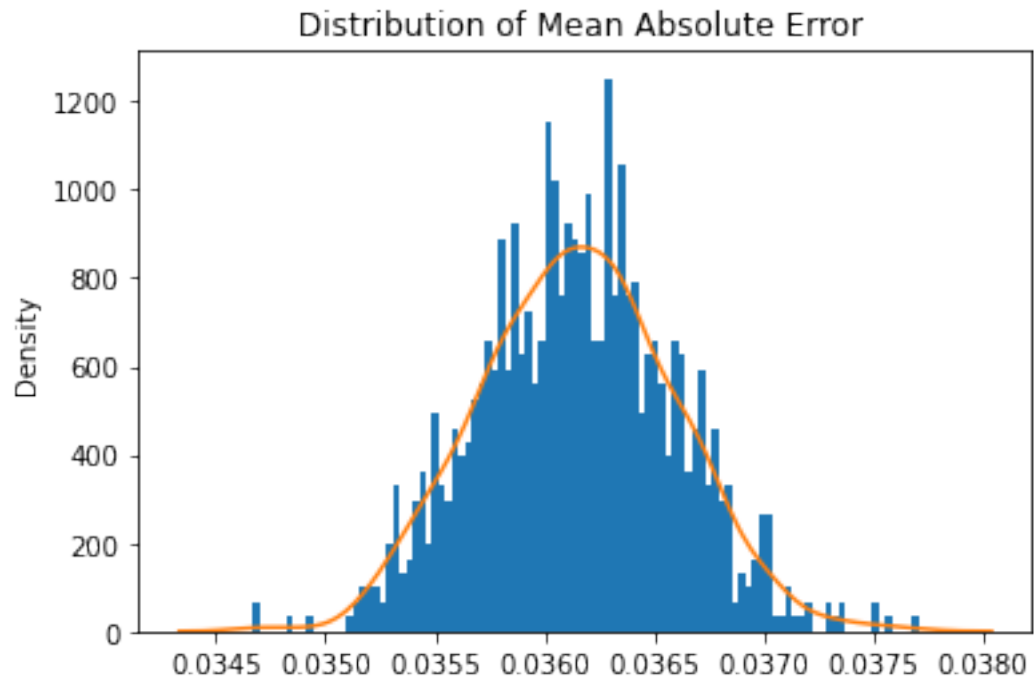
```
[11]: train_test.
→training_GAN(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,
→n_epochs,criterion,device)
```



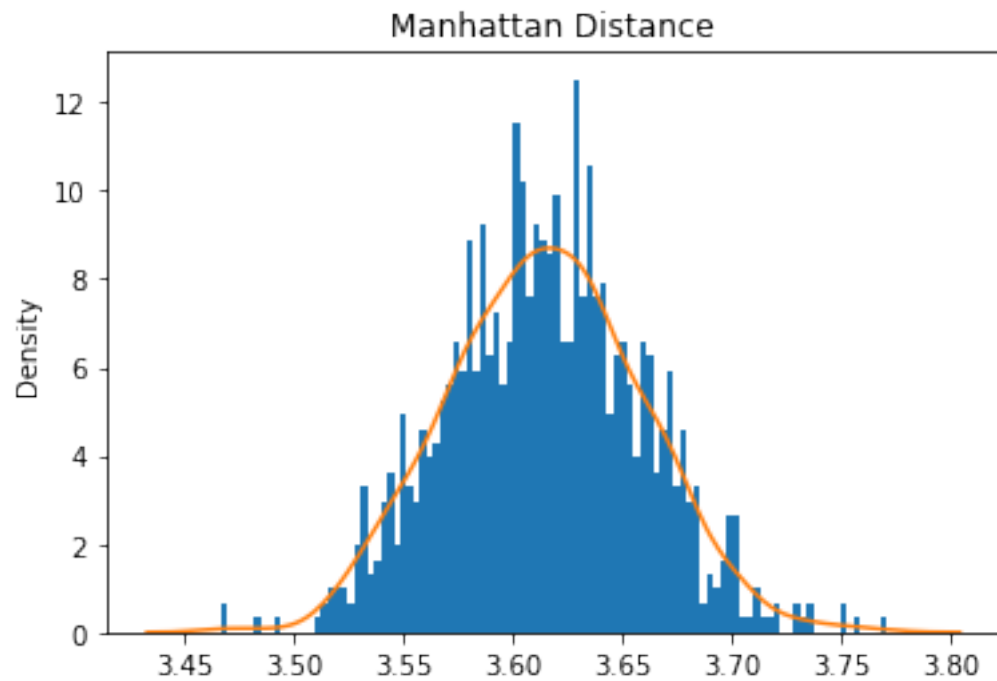
```
[12]: train_test.test_generator(generator,real_dataset,device)
```



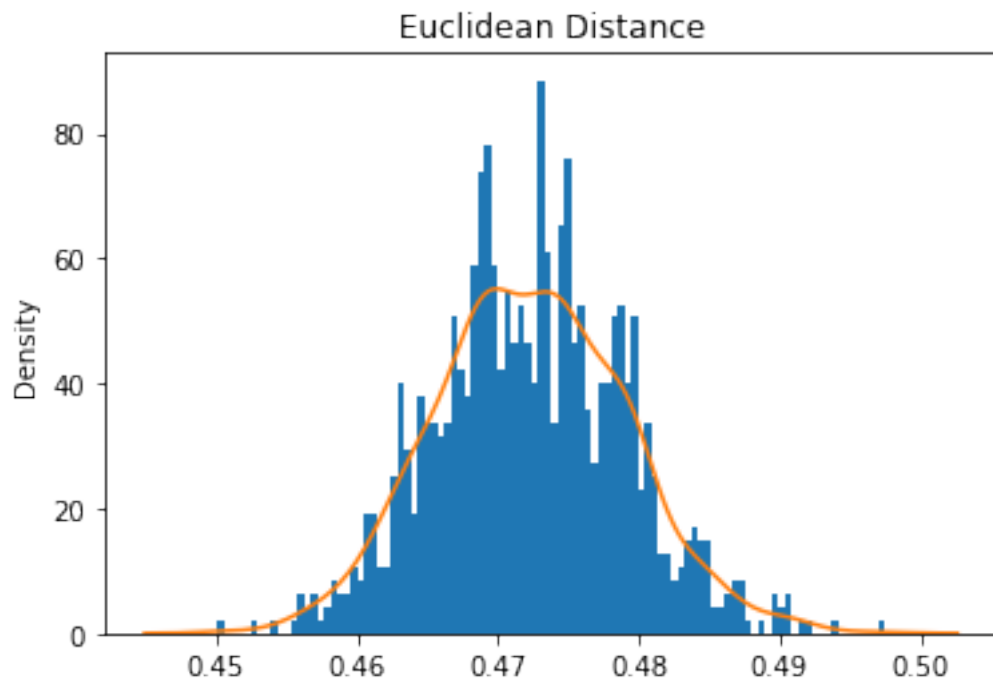
Mean Square Error: 0.002231230401942124



Mean Absolute Error: 0.036147032448649405



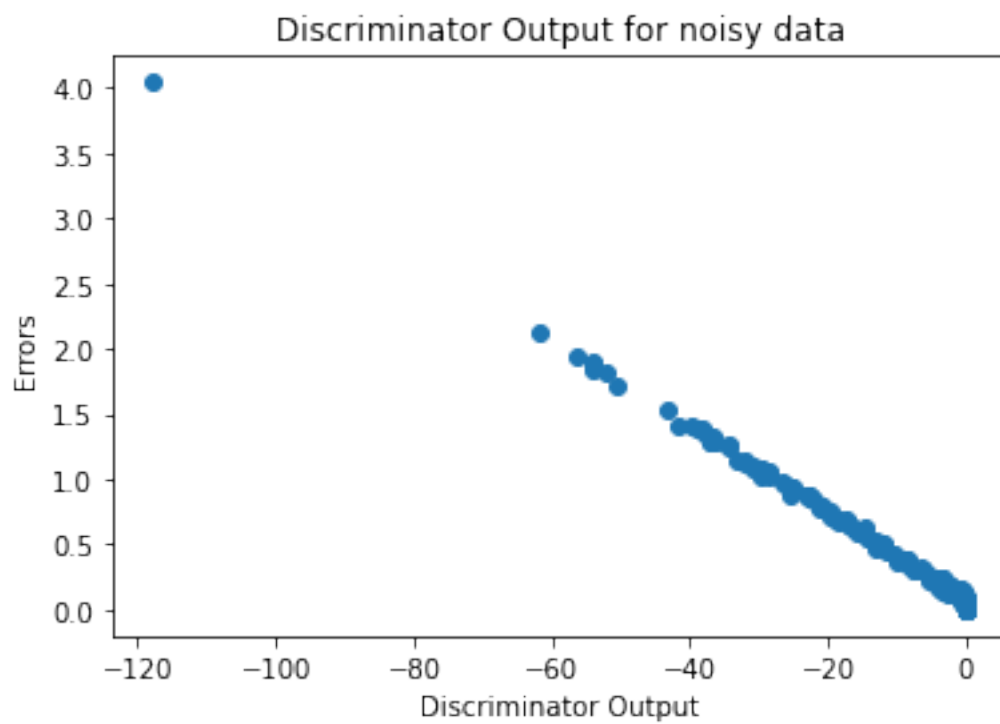
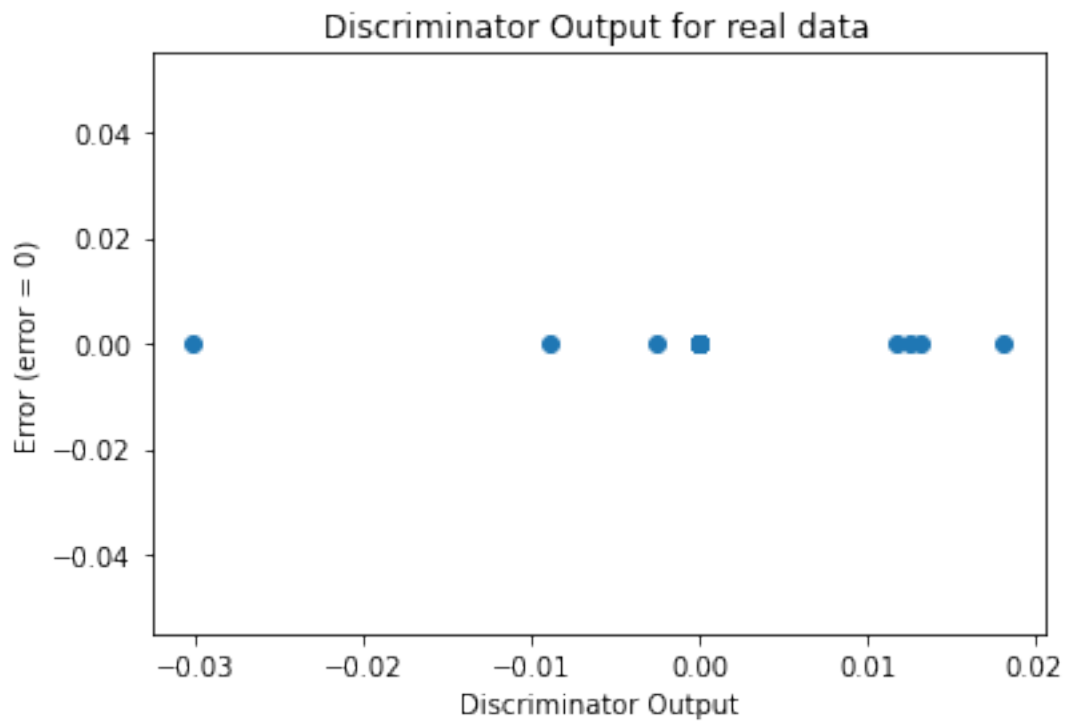
Mean Manhattan Distance: 3.6147032448649408



Mean Euclidean Distance: 0.47231046113665526

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



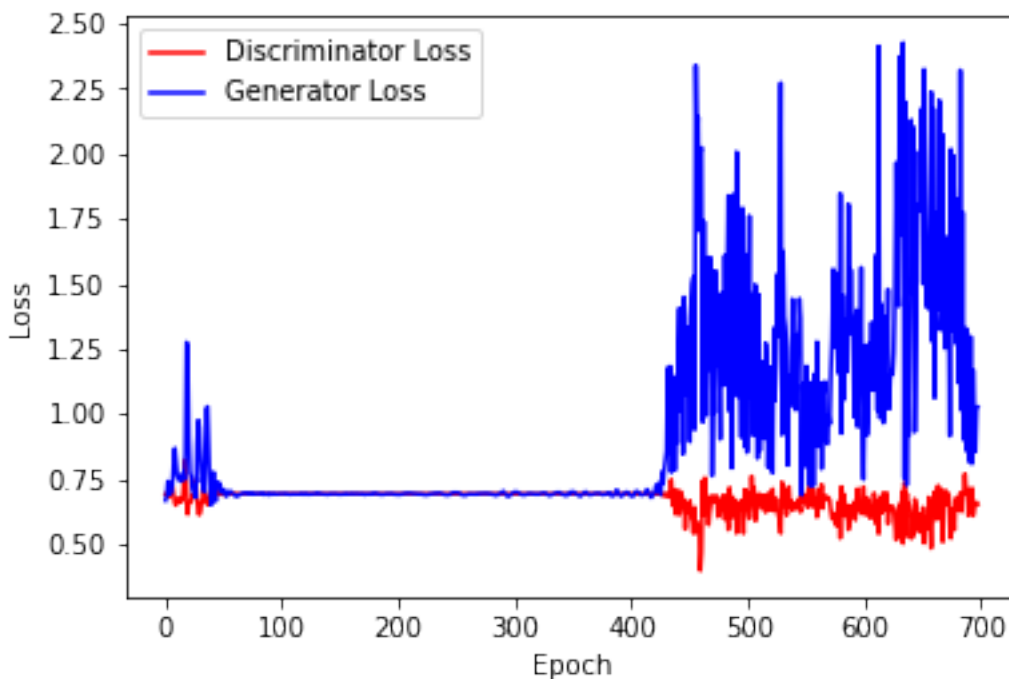


Training GAN until mse of y\_pred is  $> 0.1$  or n\_epochs  $< 30000$

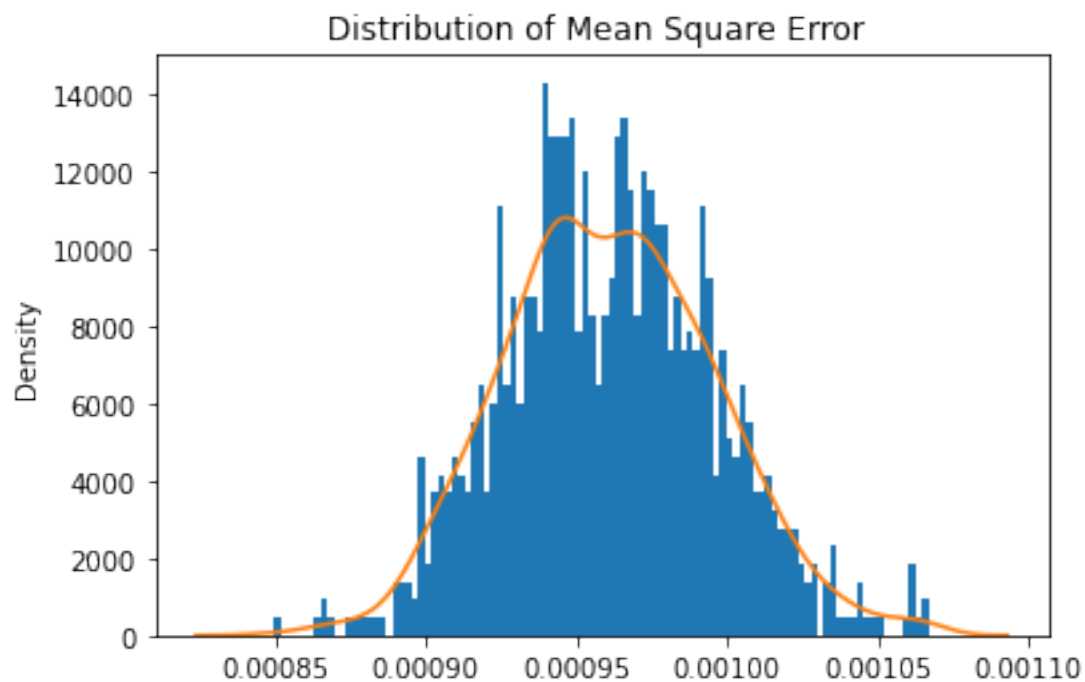
```
[14]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator.parameters(), lr=0.01, betas=(0.5,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

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

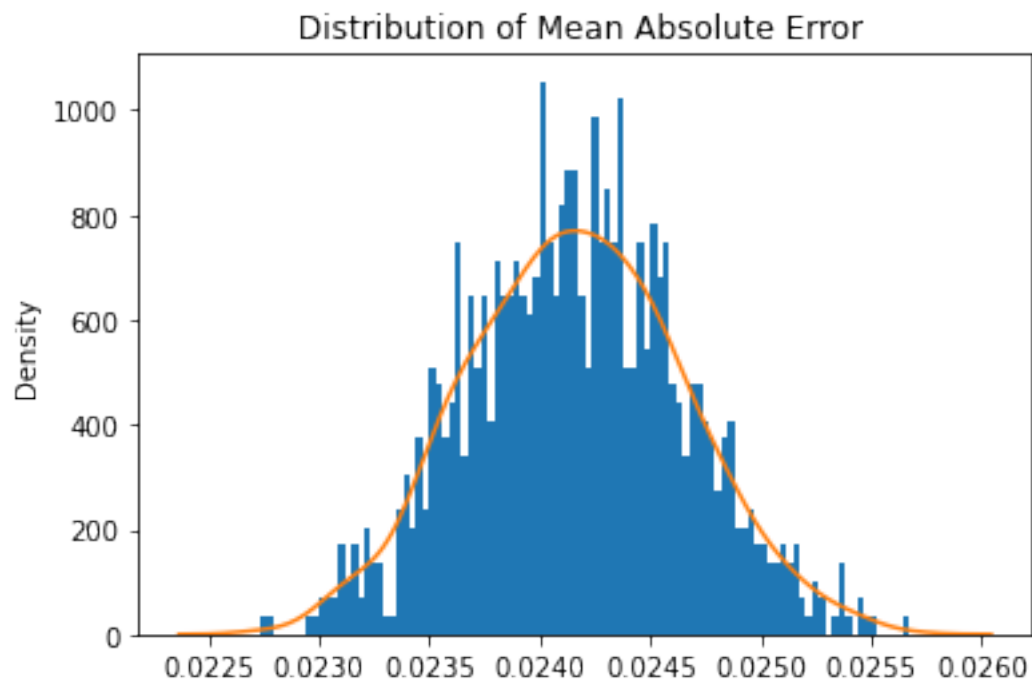
Number of epochs needed 349



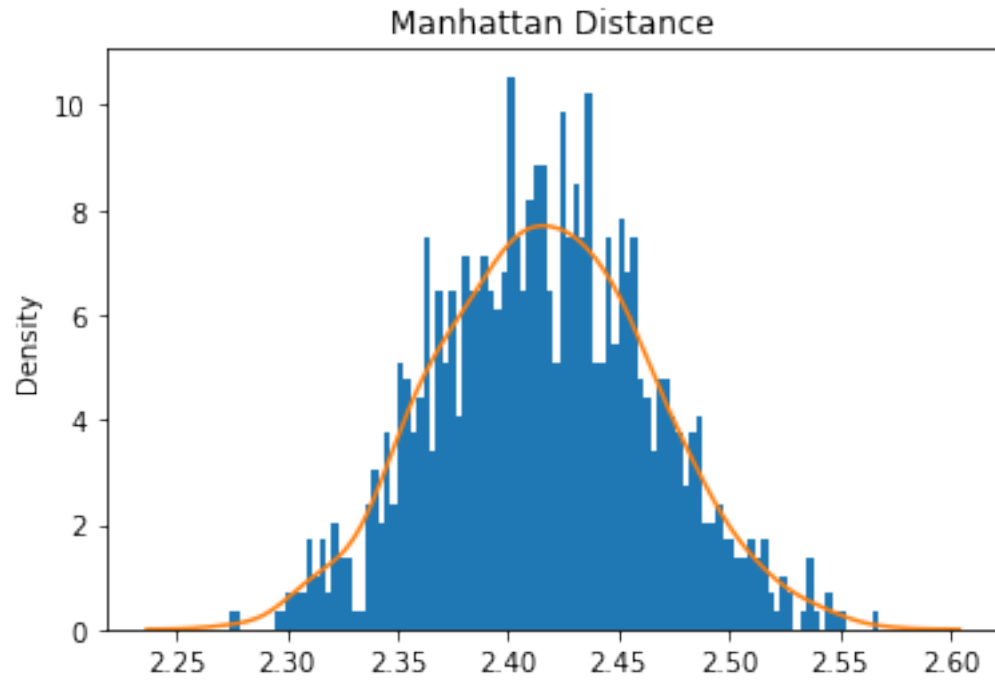
```
[16]: train_test.test_generator(generator,real_dataset,device)
```



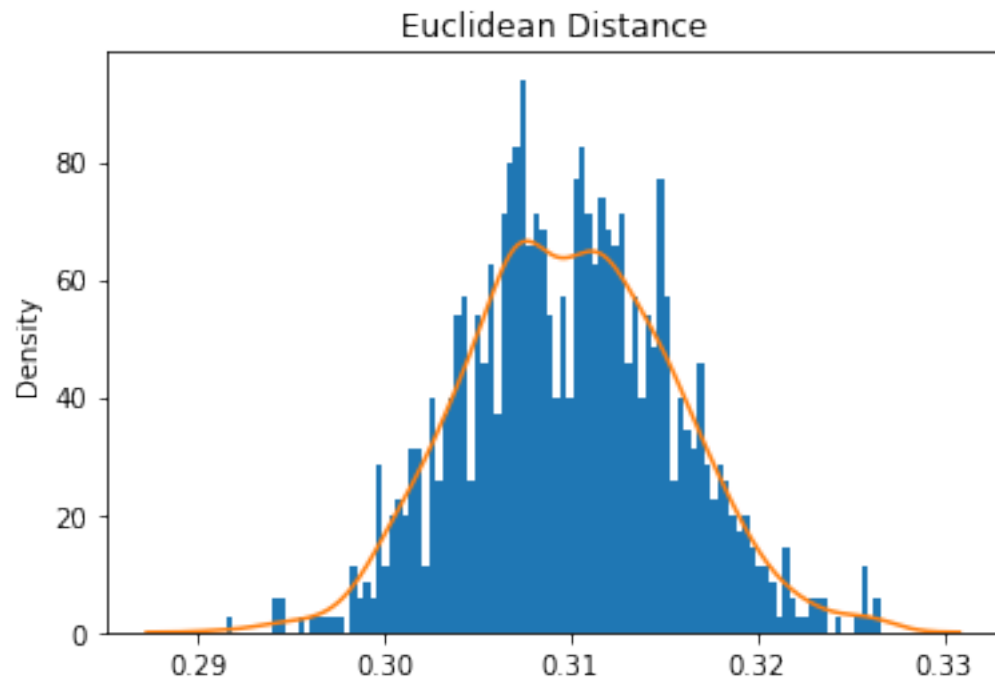
Mean Square Error: 0.0009608574757496455



Mean Absolute Error: 0.0241625481233187



Mean Manhattan Distance: 2.41625481233187



Mean Euclidean Distance: 0.3099265489498143

## 2 ABC GAN Model

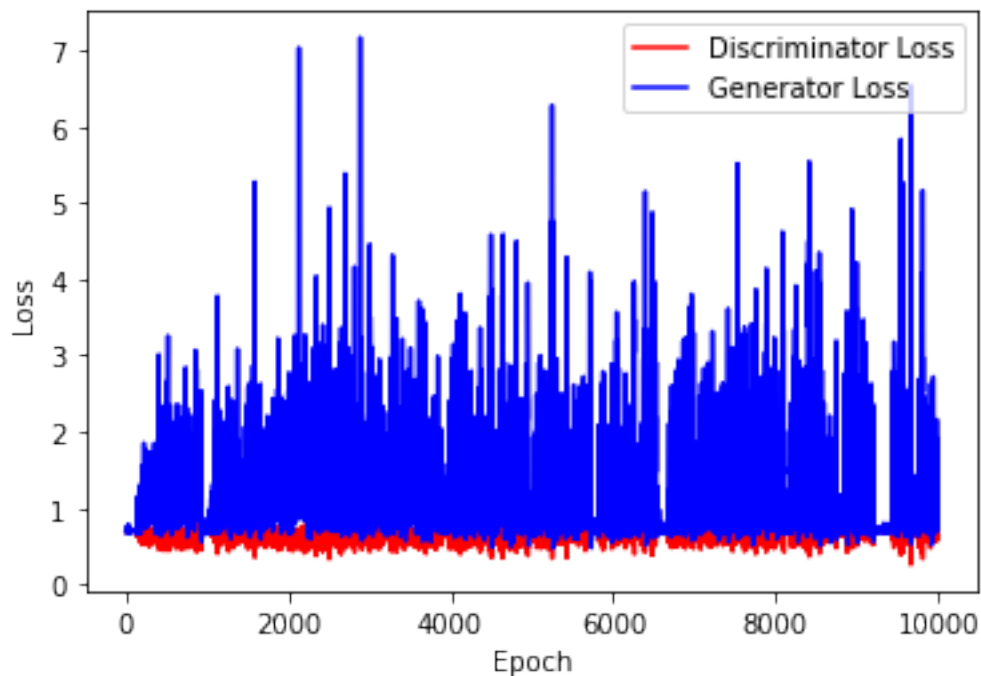
### 2.0.1 Training the network

Training ABC-GAN for `n_epochs` number of epochs

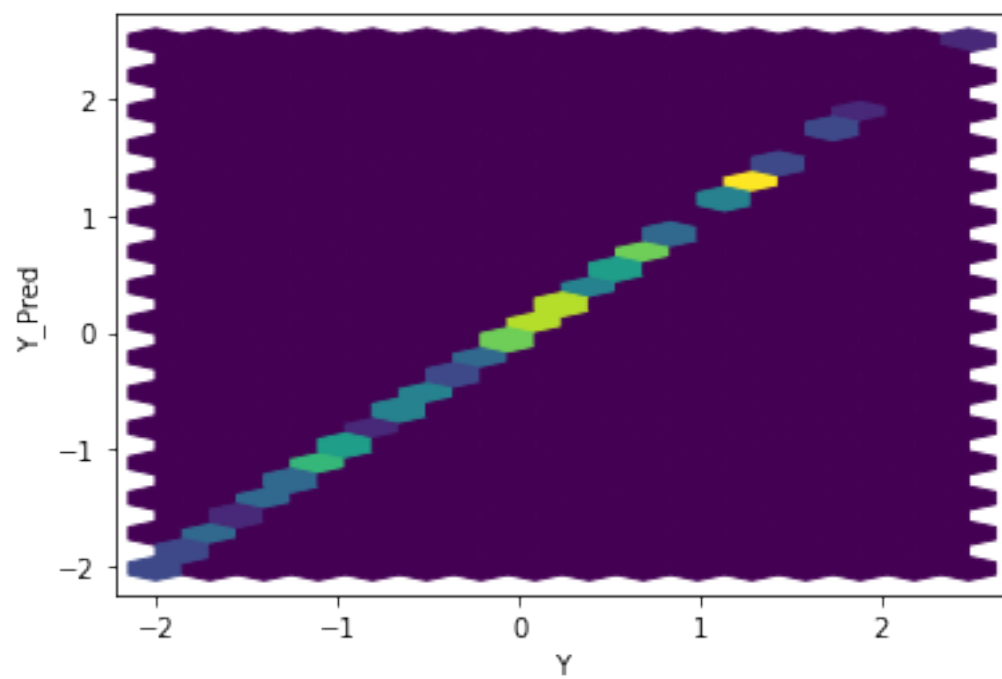
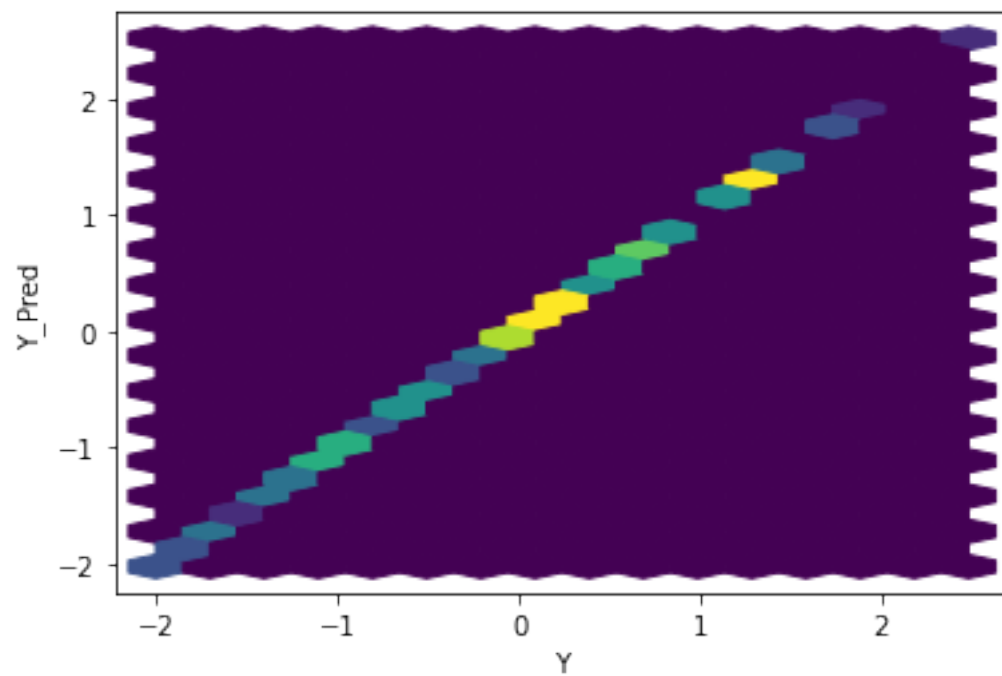
```
[17]: gen = network.Generator(n_features+2)
      disc = network.Discriminator(n_features+2)

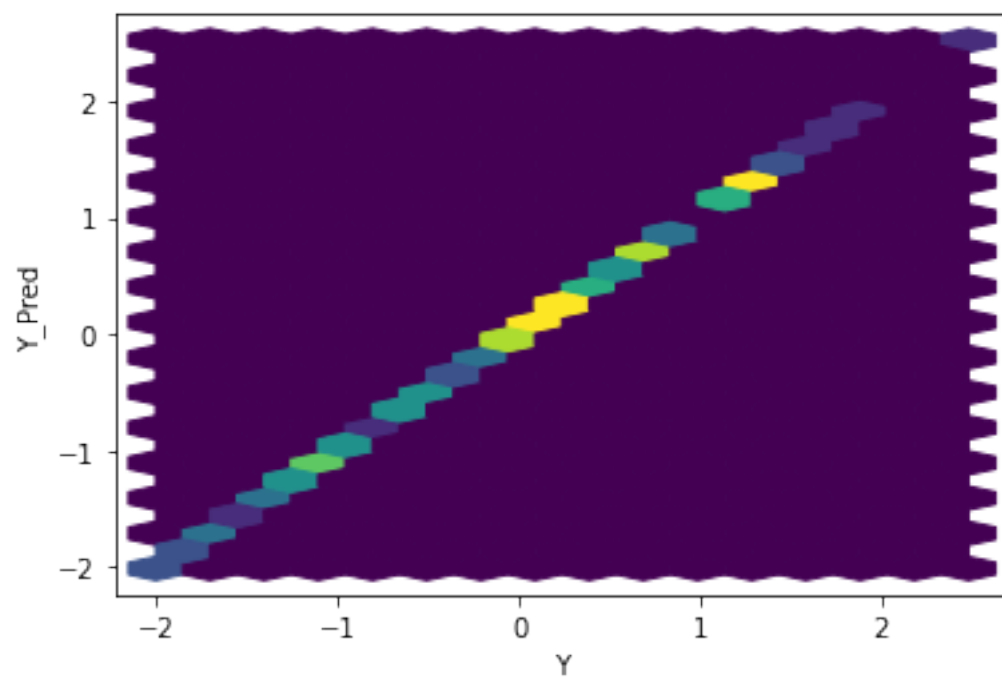
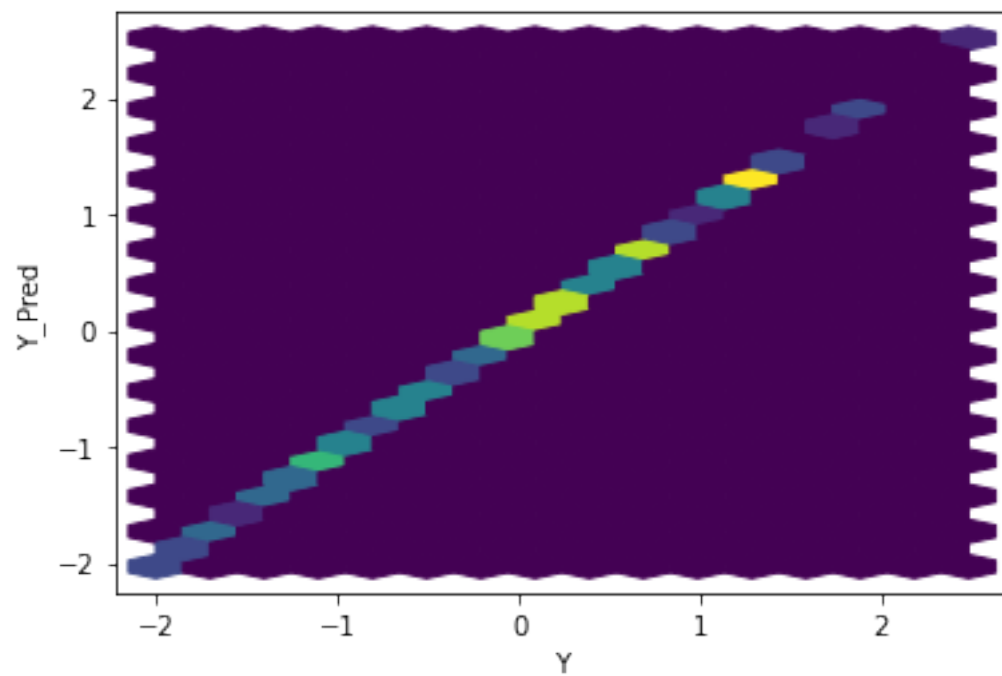
      criterion = torch.nn.BCEWithLogitsLoss()
      gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
      disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))

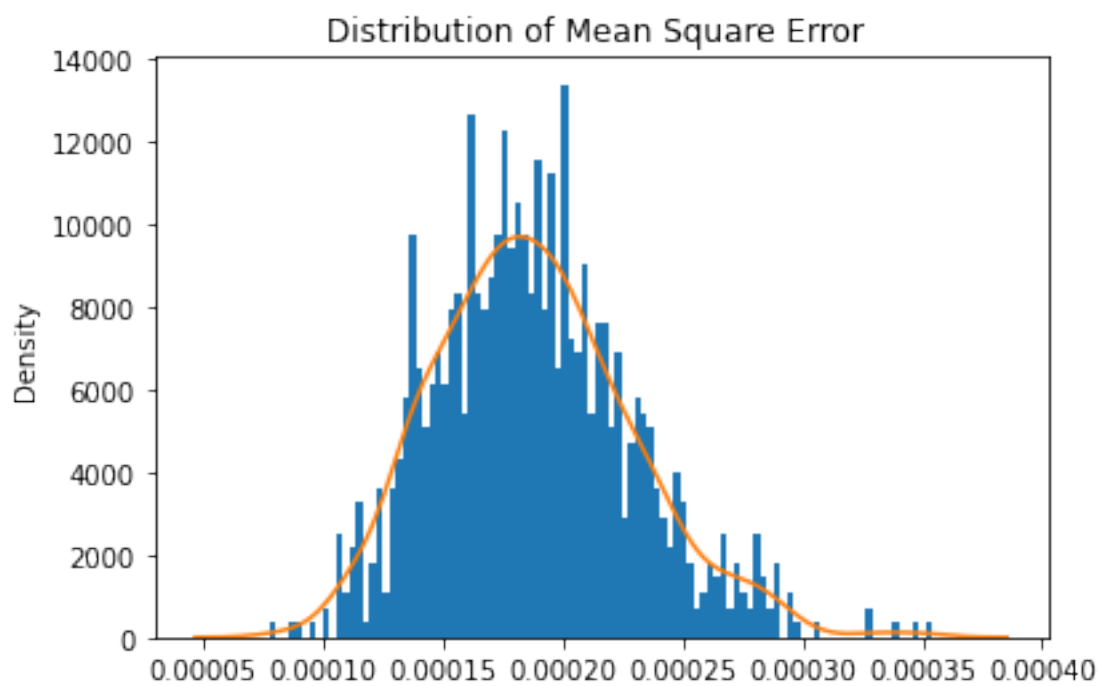
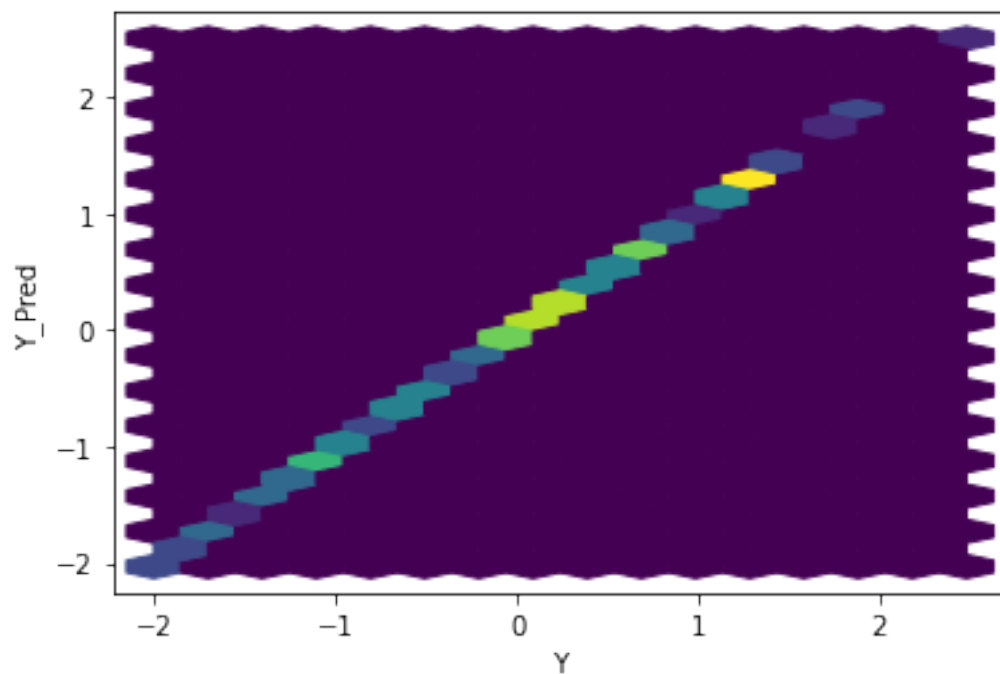
[18]: ABC_train_test.training_GAN(disc, gen,disc_opt,gen_opt,real_dataset,
      ↪batch_size, n_epochs,criterion,coeff,mean,variance,device)
```



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

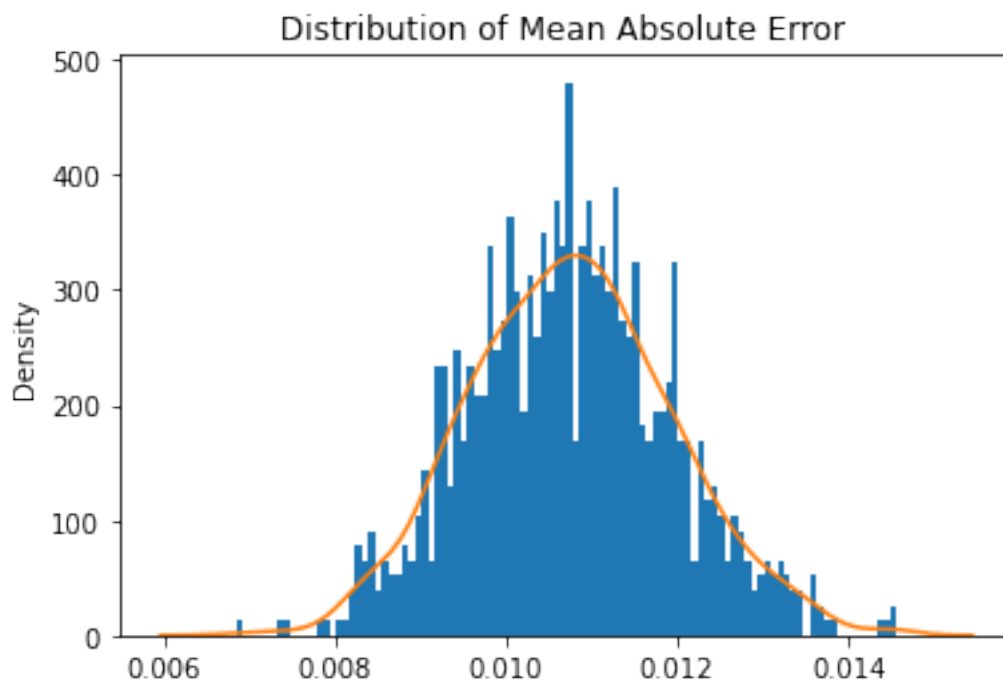






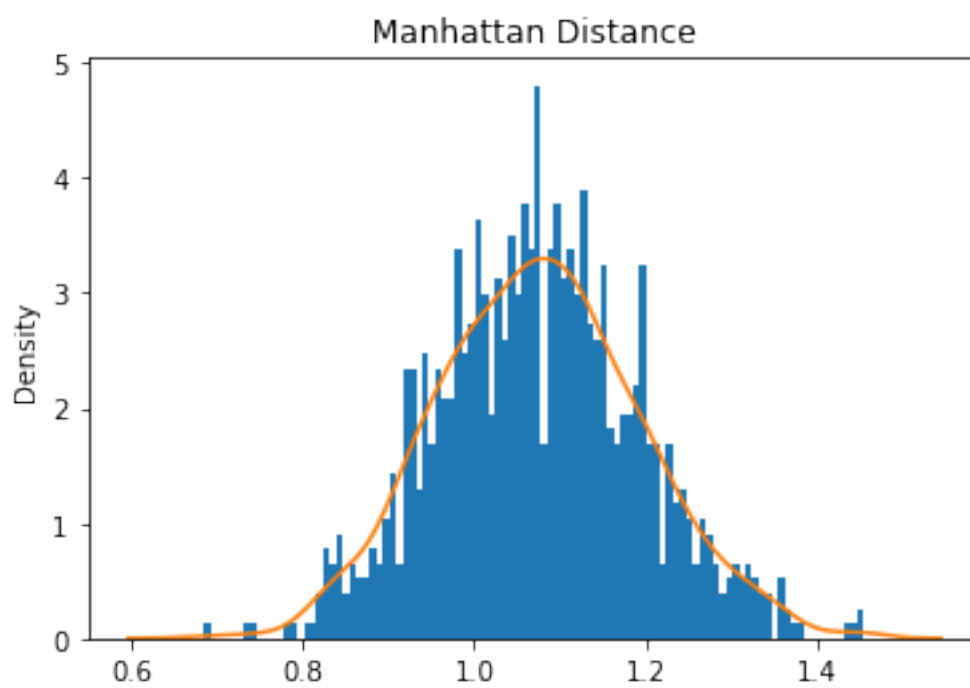
Mean Square Error: 0.00018764875007634972



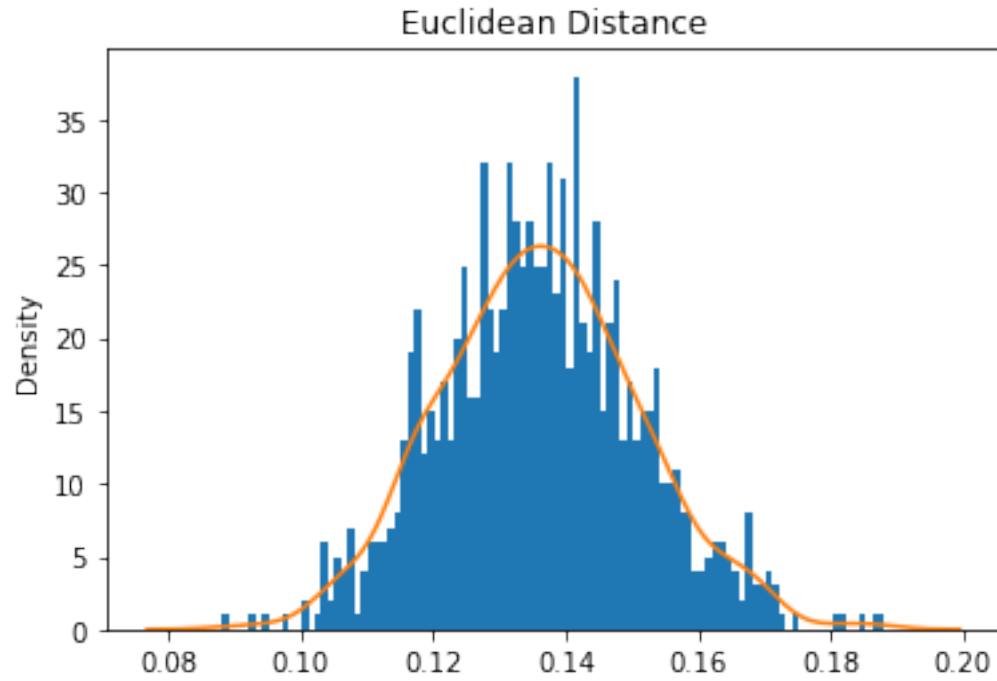


Mean Absolute Error: 0.010752491418197752

Mean Manhattan Distance: 1.0752491418197752

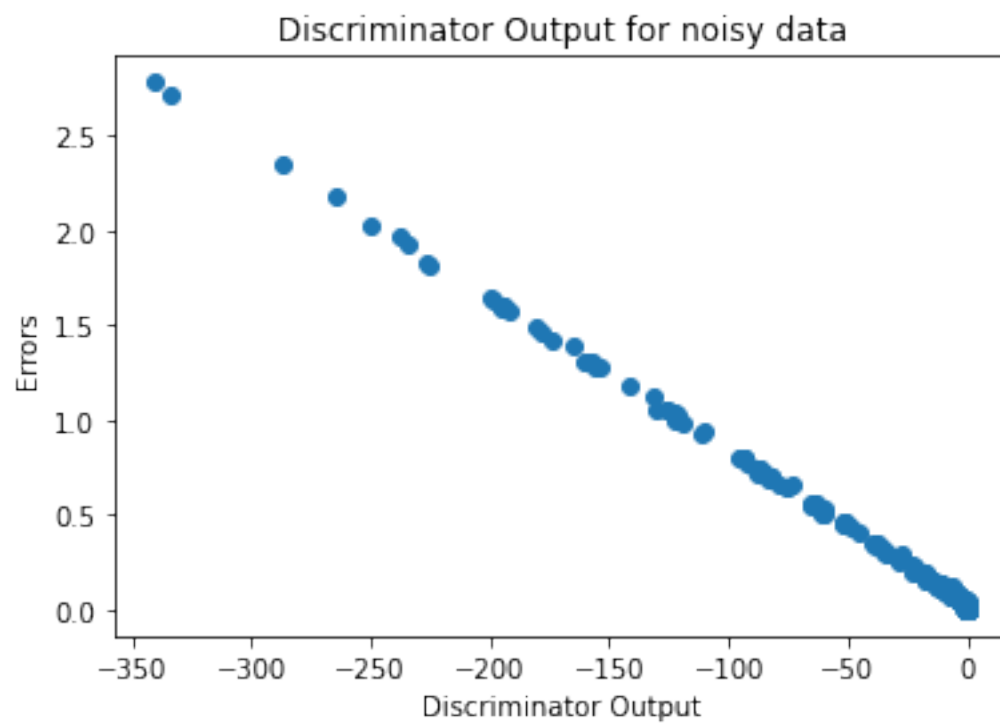
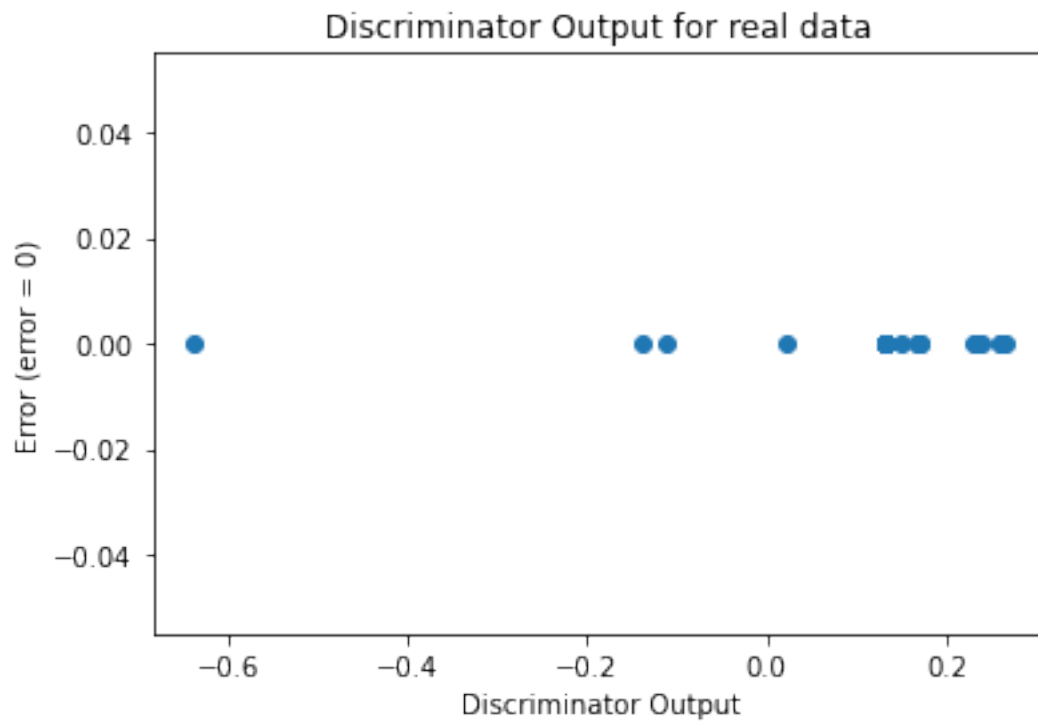


Mean Euclidean Distance: 0.13615867677333449



### Sanity Checks

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



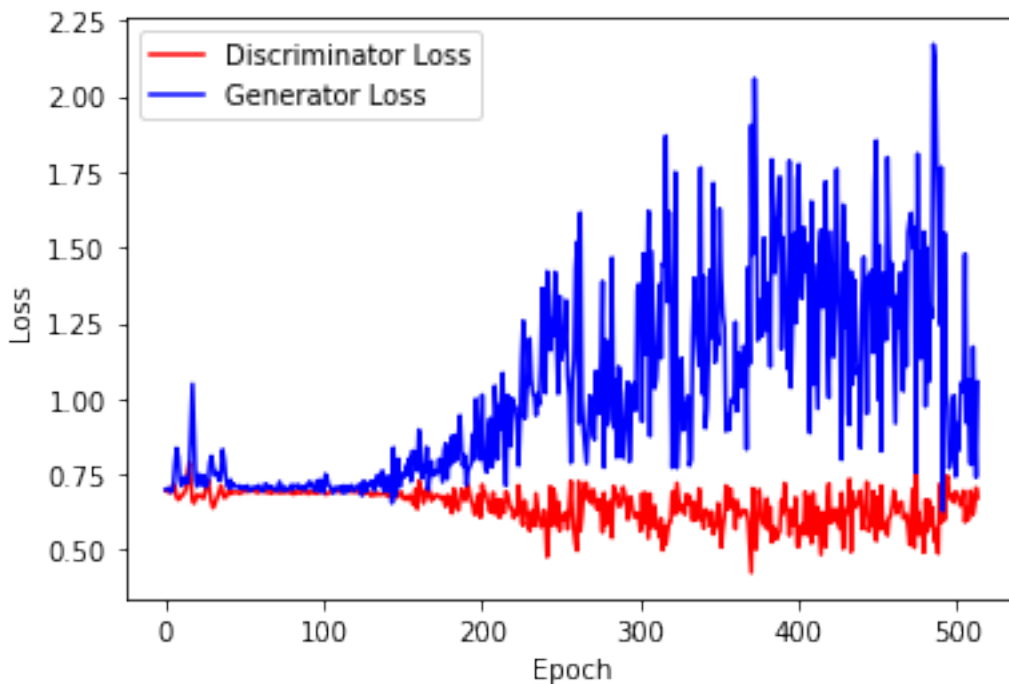
Training GAN until mse of y\_pred is  $> 0.1$  or n\_epochs  $< 30000$

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

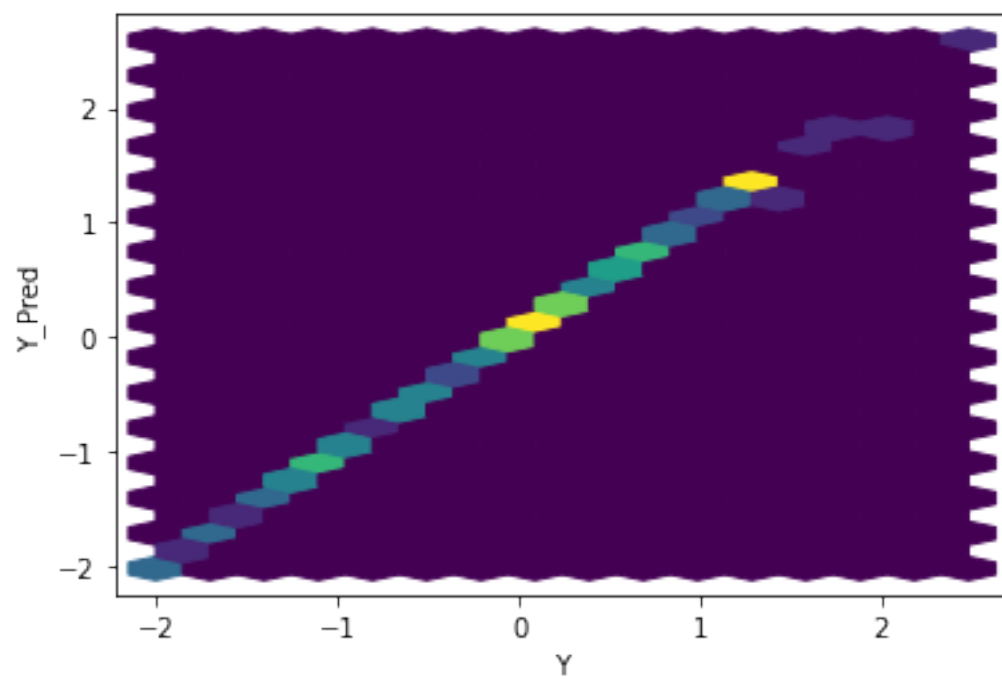
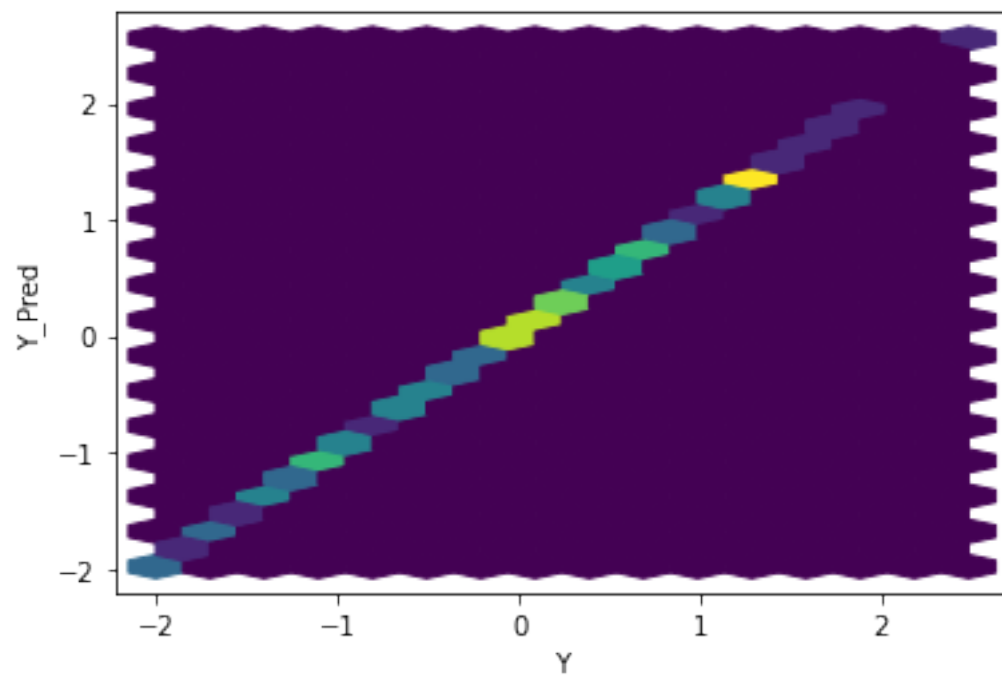
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))
```

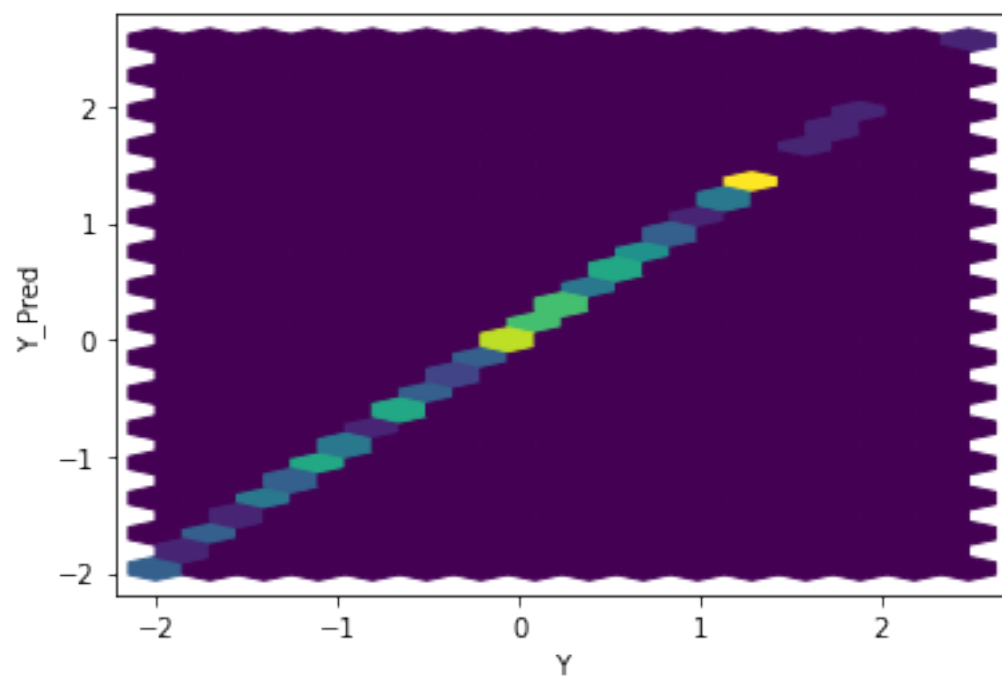
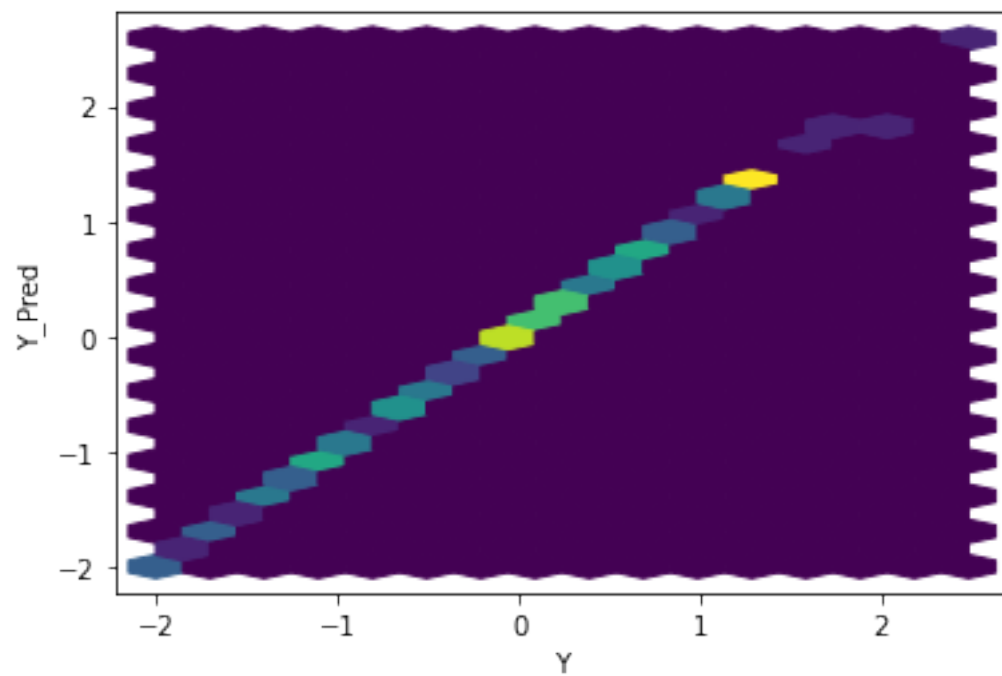
```
[22]: ABC_train_test.
      ↪ training_GAN_2(disc,gen,disc_opt,gen_opt,real_dataset,batch_size,
      ↪ error,criterion,coeff,mean,variance,device)
```

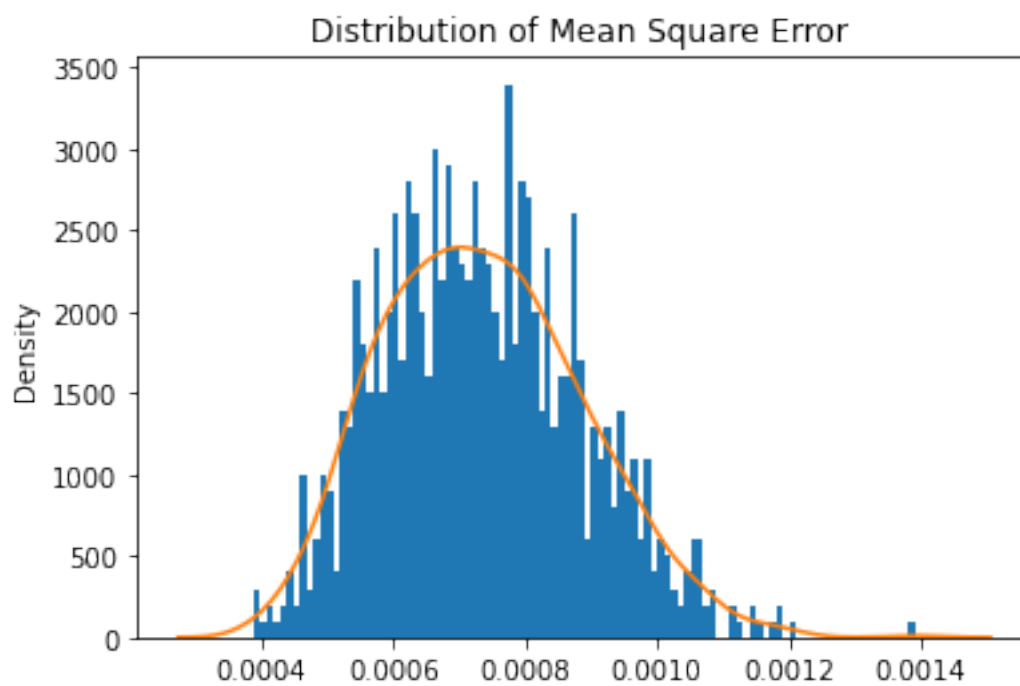
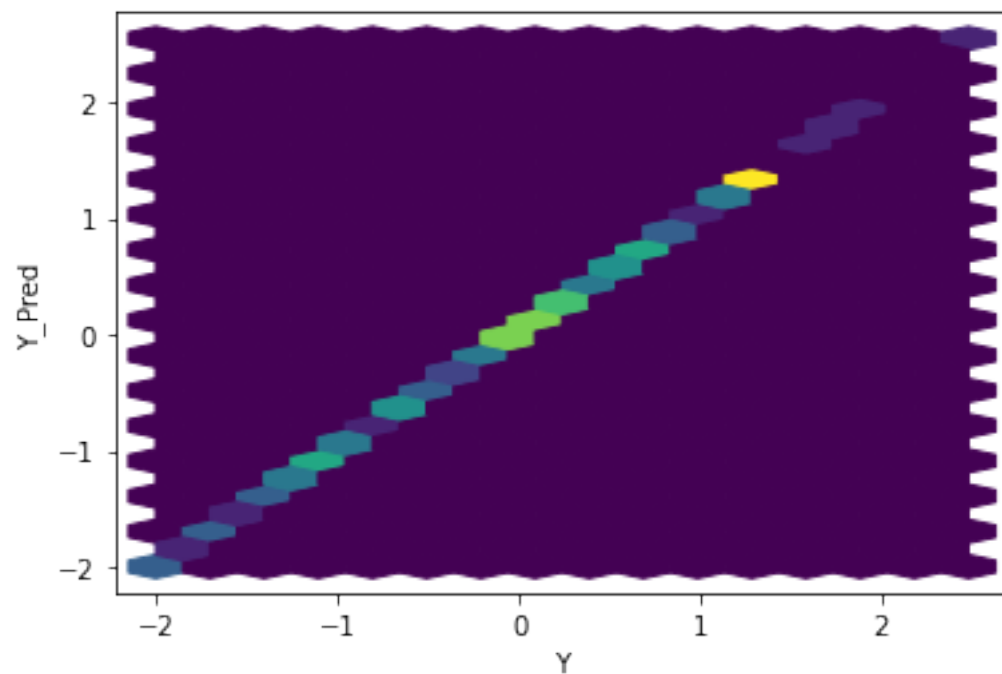
Number of epochs 257



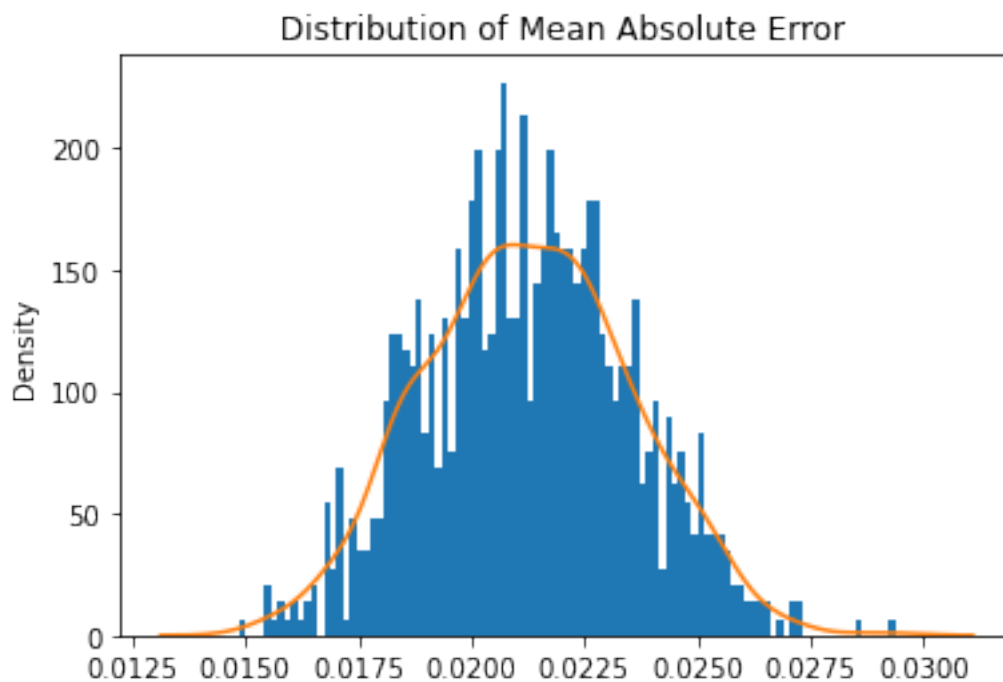
```
[23]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```





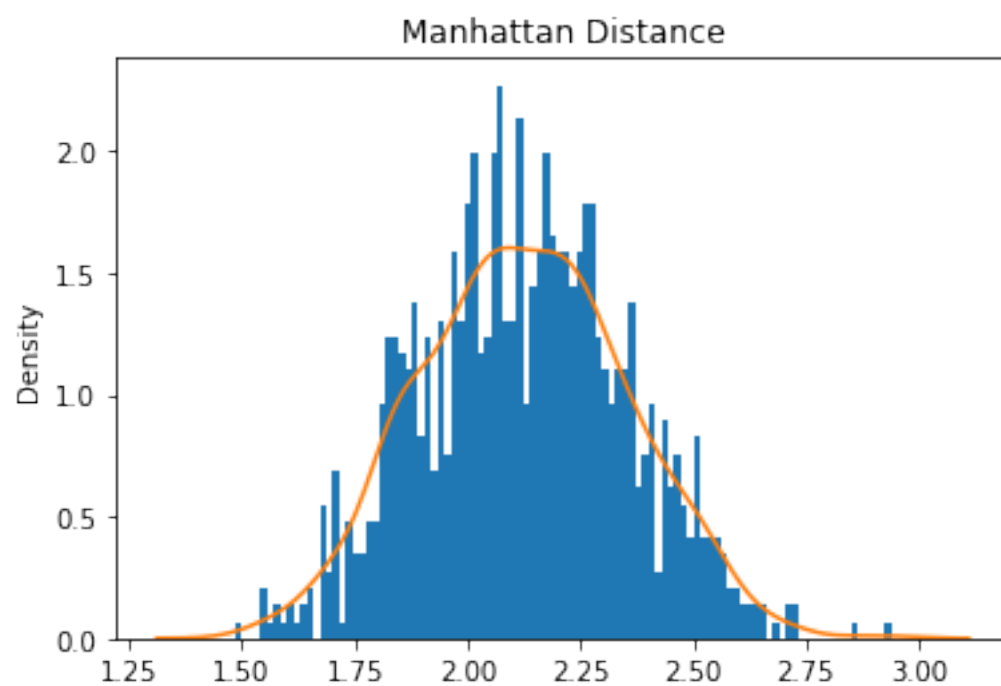


Mean Square Error: 0.0007353198819635694



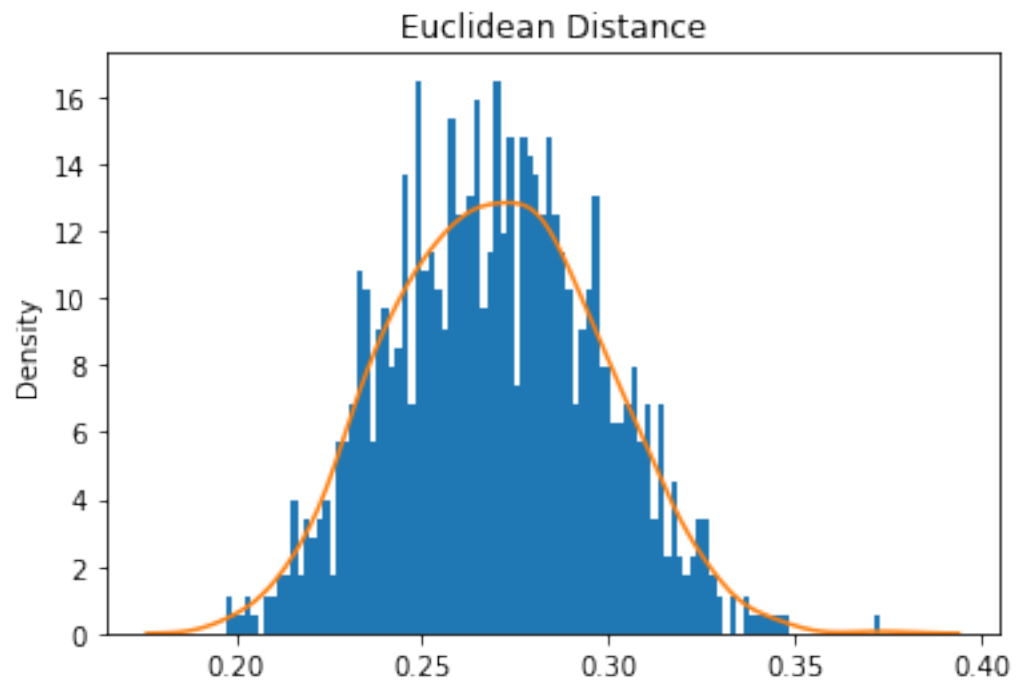
Mean Absolute Error: 0.02125035571185872

Mean Manhattan Distance: 2.125035571185872





Mean Euclidean Distance: 0.2697281570391835



[ ]: