

Dataset2_Friedman1_output_6

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. $Y = 1 + \beta_1 x_1 + \beta_2 x_2 + \beta_2 x_3 + \dots + \beta_n x_n + N(0, \sigma)$ where $\sigma = 0.1$
 2. $\beta_i \sim N(0, \sigma^*)$ when $\mu = 0$ else $\beta_i \sim N(\beta_i^*, \sigma^*)$ where 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 = 1
     variance = 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	0.187587	0.134677	0.847462	0.612660	0.137046	0.687641	0.988331
1	0.922862	0.423734	0.177819	0.782755	0.756410	0.574923	0.597466
2	0.430331	0.799804	0.430384	0.164755	0.617313	0.576486	0.021633

```

3  0.090216  0.198686  0.332492  0.878699  0.947600  0.980393  0.012332
4  0.108835  0.081441  0.891636  0.446162  0.704788  0.837229  0.918084

```

```

      X7      X8      X9      Y
0  0.850751  0.186227  0.412563  10.262900
1  0.121094  0.459999  0.557252  23.303144
2  0.397491  0.635276  0.546216  13.492939
3  0.034171  0.898897  0.723021  14.598230
4  0.812956  0.232613  0.137706  11.287338

```

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:          0.852
Model:                  OLS    Adj. R-squared:      0.836
Method:                 Least Squares    F-statistic:      51.41
Date:                  Wed, 20 Oct 2021    Prob (F-statistic):  1.14e-32
Time:                  20:10:42    Log-Likelihood:     -46.219
No. Observations:      100    AIC:              114.4
Df Residuals:          89    BIC:              143.1
Df Model:               10
Covariance Type:       nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	6.939e-16	0.041	1.7e-14	1.000	-0.081	0.081
x1	0.4281	0.047	9.166	0.000	0.335	0.521
x2	0.3924	0.044	8.979	0.000	0.306	0.479
x3	0.0482	0.044	1.101	0.274	-0.039	0.135
x4	0.5386	0.044	12.193	0.000	0.451	0.626
x5	0.2971	0.042	7.062	0.000	0.213	0.381
x6	-0.0189	0.043	-0.444	0.658	-0.104	0.066
x7	-0.0667	0.044	-1.531	0.129	-0.153	0.020
x8	-0.0553	0.045	-1.238	0.219	-0.144	0.033
x9	0.0623	0.044	1.427	0.157	-0.024	0.149
x10	0.0081	0.043	0.188	0.851	-0.077	0.093

```

=====
Omnibus:              3.000    Durbin-Watson:          2.054
Prob(Omnibus):        0.223    Jarque-Bera (JB):        2.348
Skew:                 -0.307    Prob(JB):                0.309
Kurtosis:             3.432    Cond. No.:               1.85
=====

```

Notes:

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

Parameters: const 6.938894e-16

x1 4.280757e-01

x2 3.924343e-01

x3 4.823562e-02

x4 5.386116e-01

x5 2.970731e-01

x6 -1.894973e-02

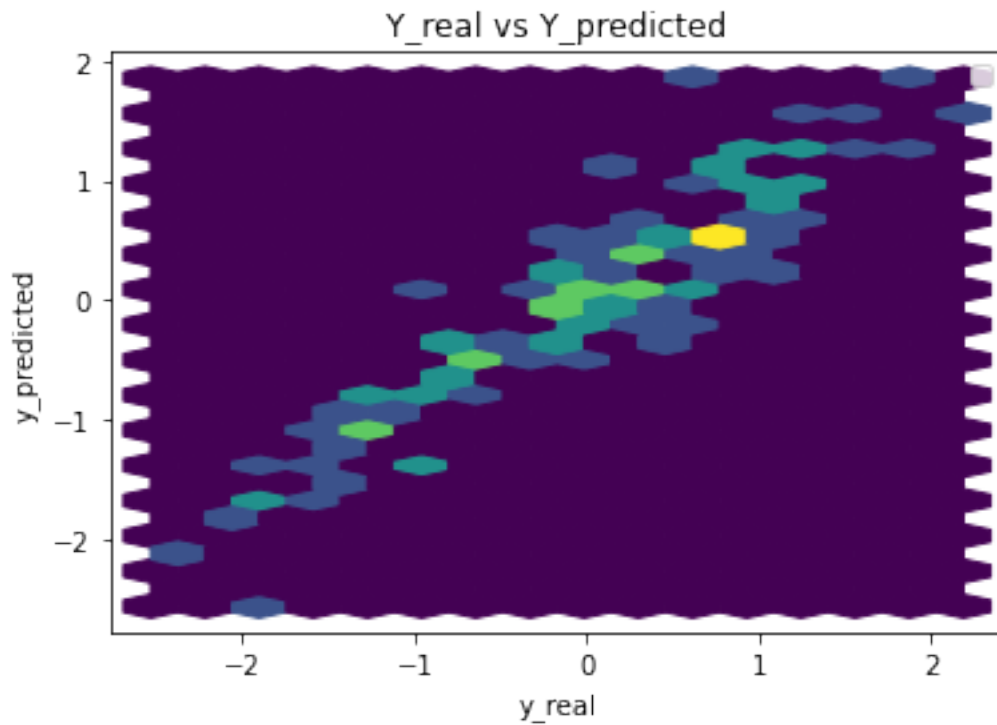
x7 -6.672992e-02

x8 -5.534752e-02

x9 6.226599e-02

x10 8.068316e-03

dtype: float64



Performance Metrics

Mean Squared Error: 0.14756329678882776

Mean Absolute Error: 0.3094021888439953

Manhattan distance: 30.940218884399517

Euclidean distance: 3.841396839547143

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

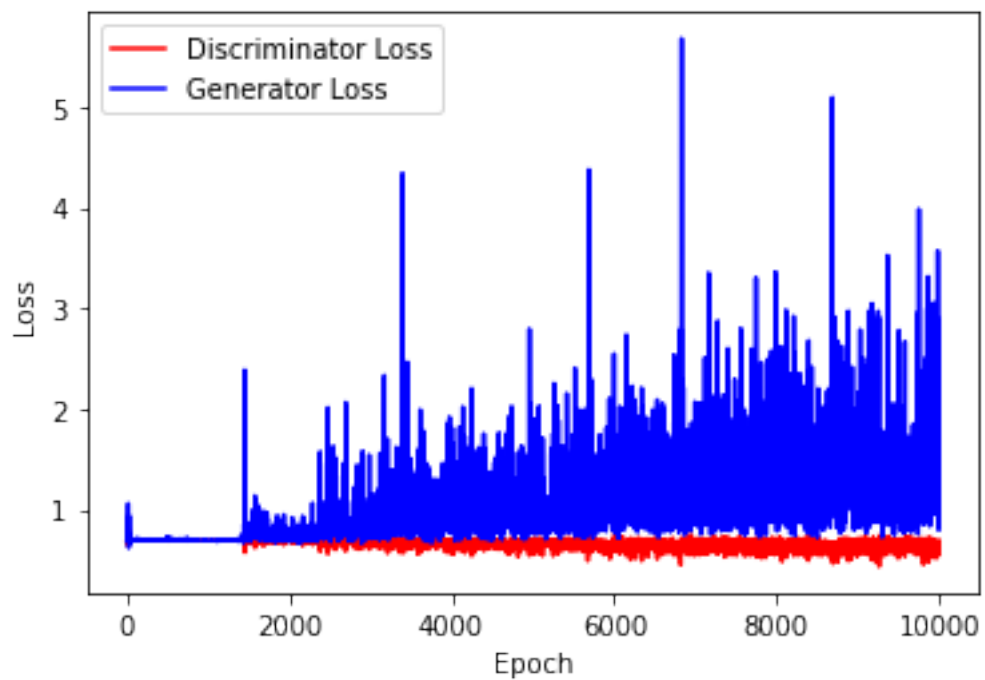
```
[9]: generator = network.Generator(n_features+2)
     discriminator = network.Discriminator(n_features+2)

     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))
```

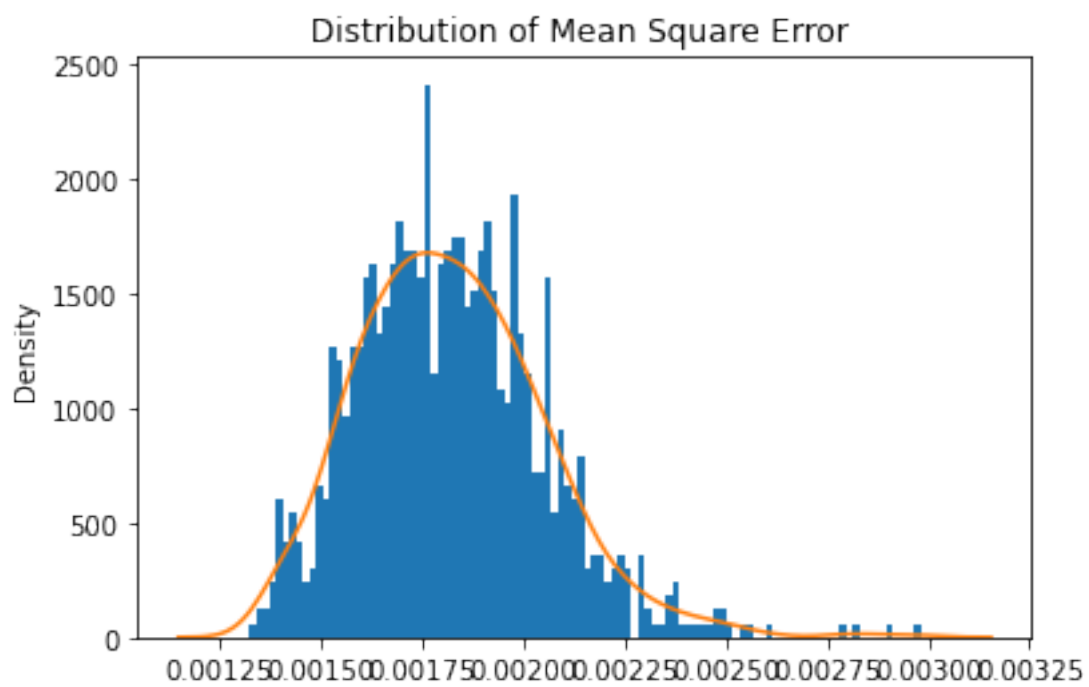
```
[10]: print(generator)
      print(discriminator)
```

```
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()
)
```

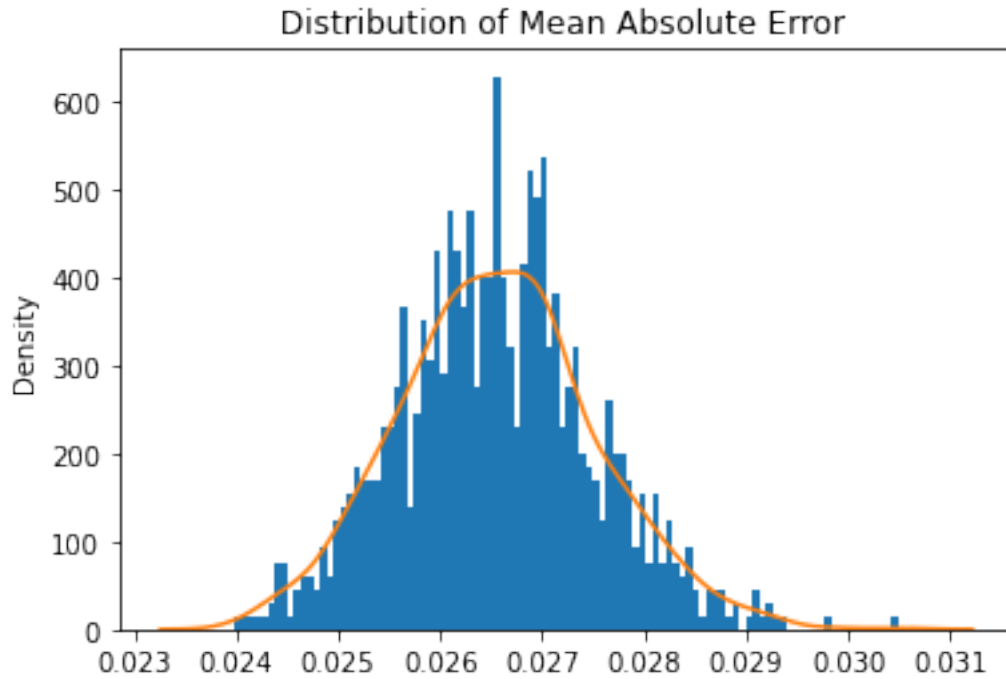
```
[11]: train_test.
     ↪training_GAN(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,
     ↪n_epochs,criterion,device)
```



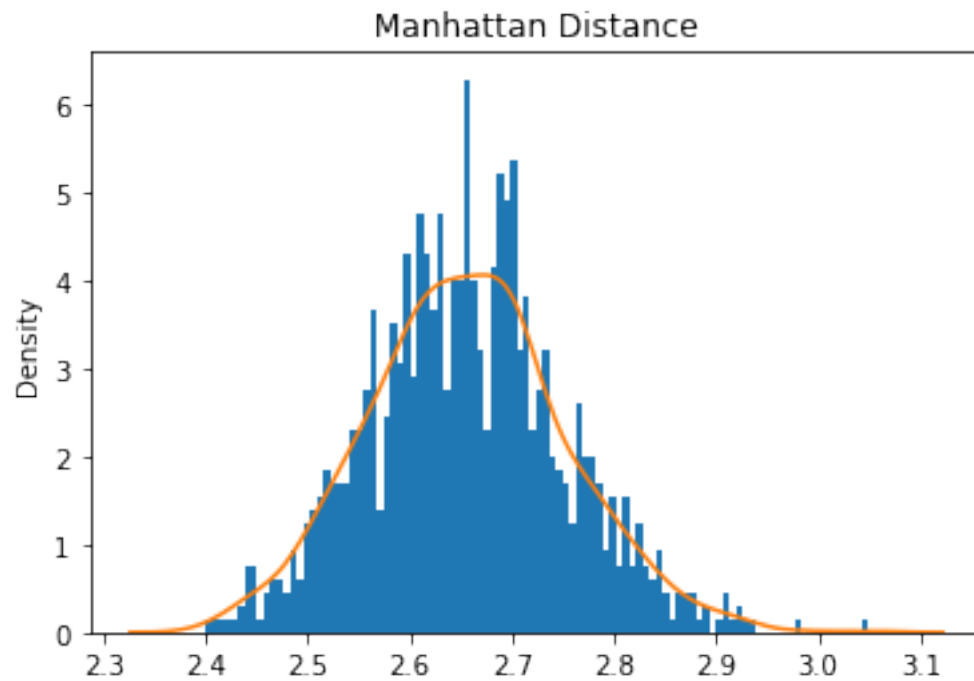
```
[12]: train_test.test_generator(generator,real_dataset,device)
```



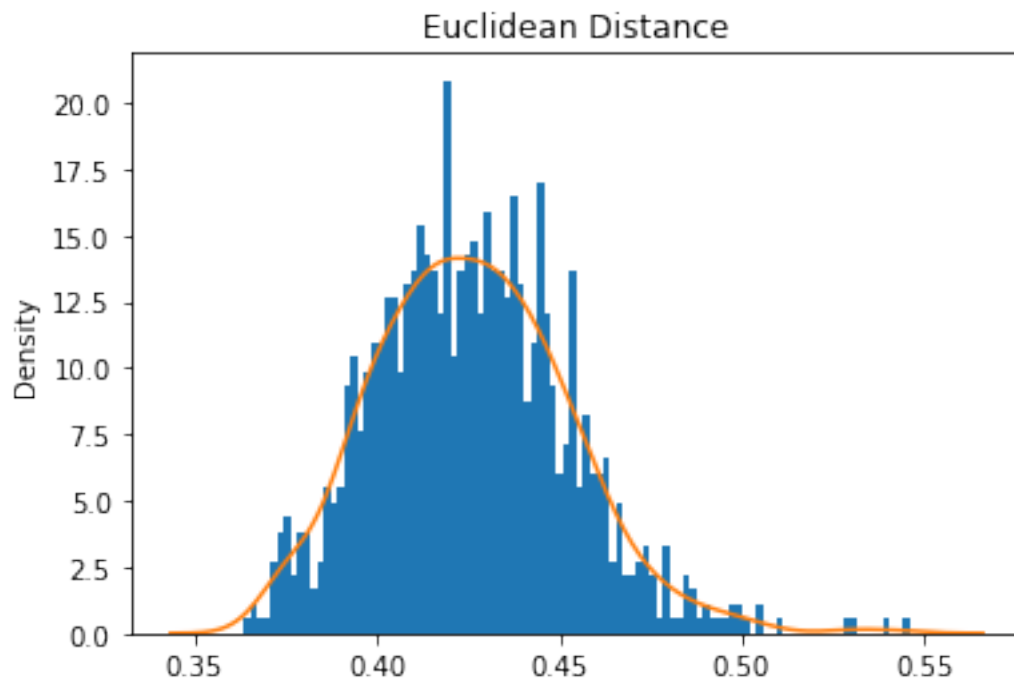
Mean Square Error: 0.001821917559121434



Mean Absolute Error: 0.026559316127616913

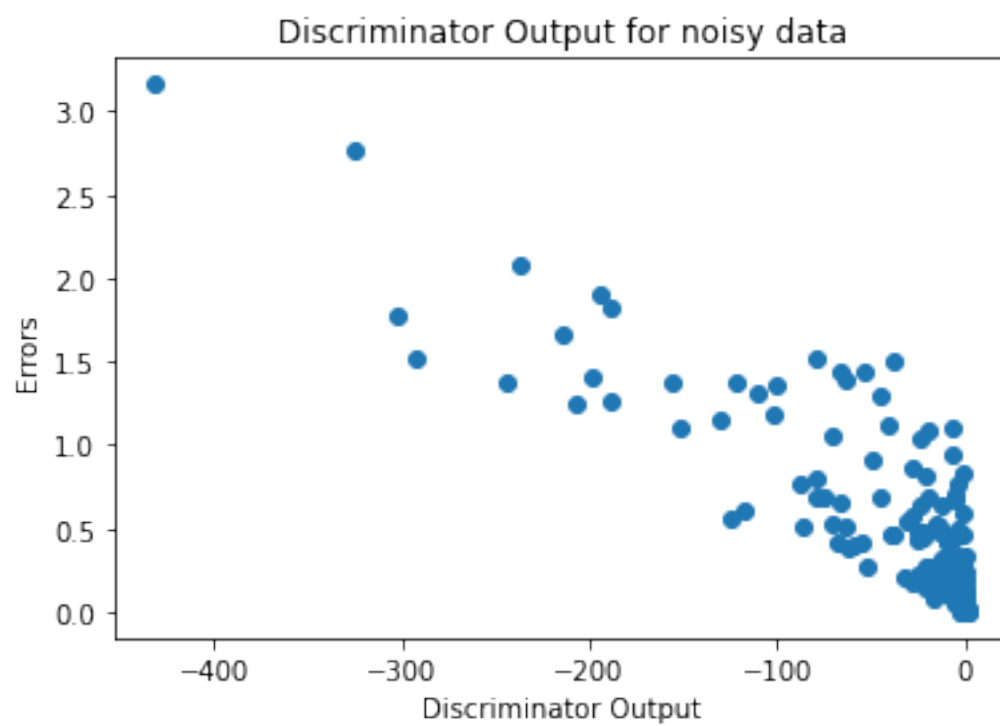
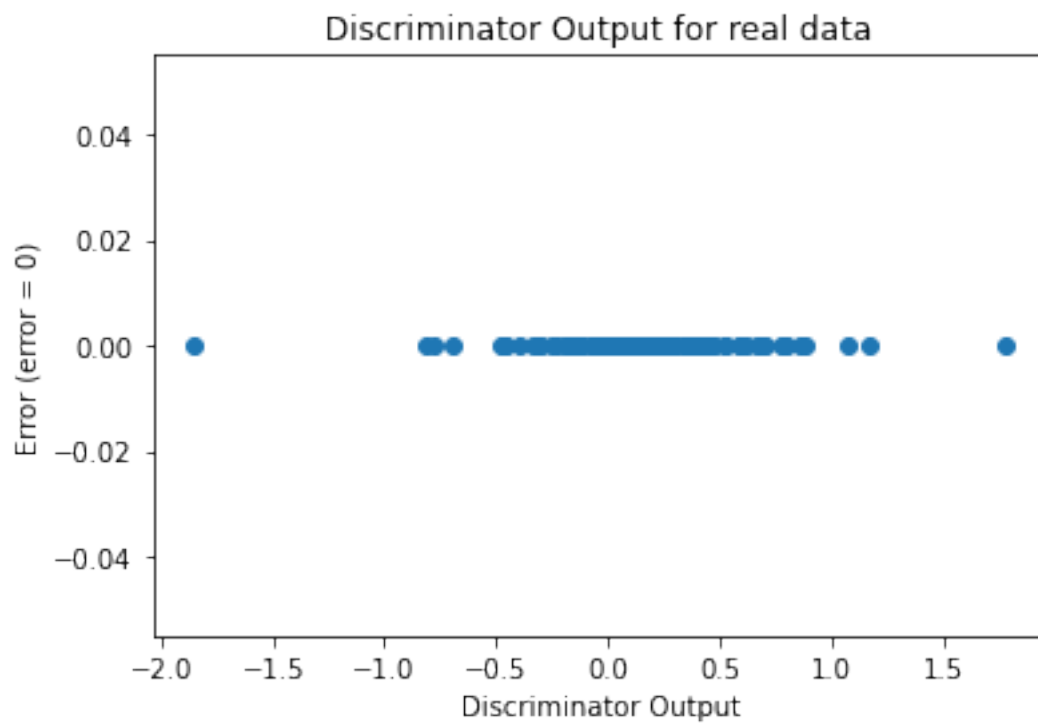


Mean Manhattan Distance: 2.6559316127616914



Mean Euclidean Distance: 0.42599847390394147

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

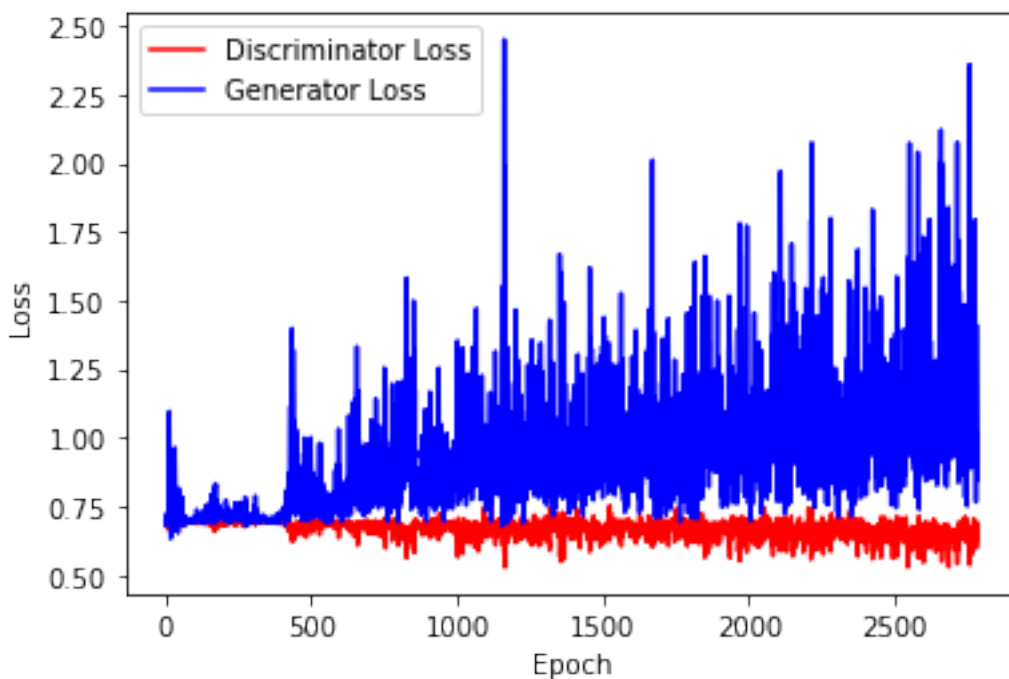



Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

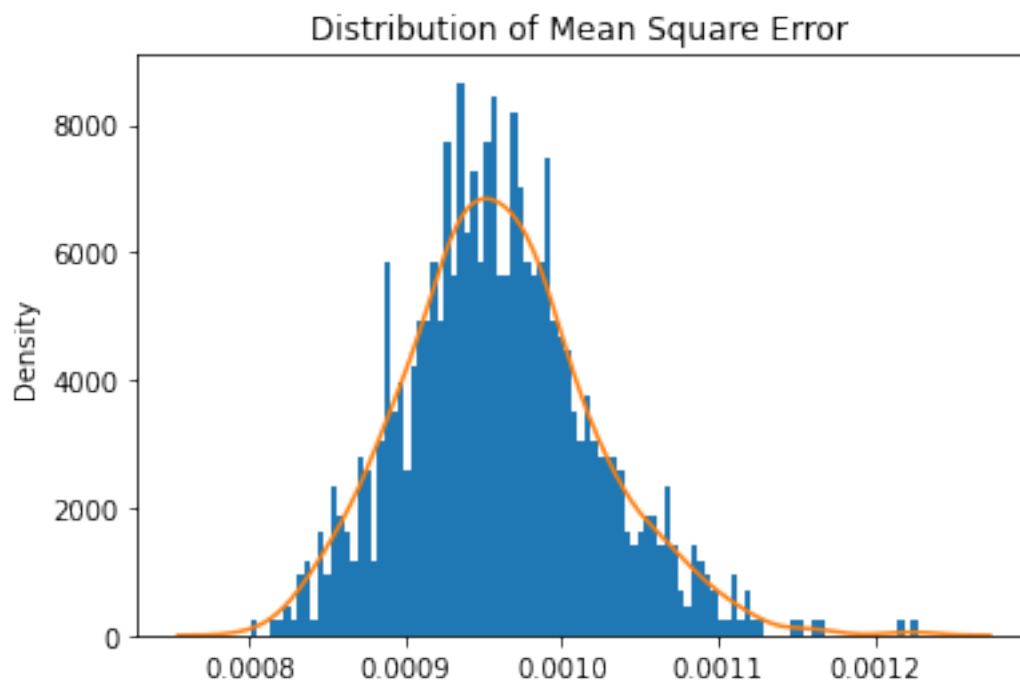
```
[14]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator.parameters(), lr=0.01, betas=(0.5,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

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

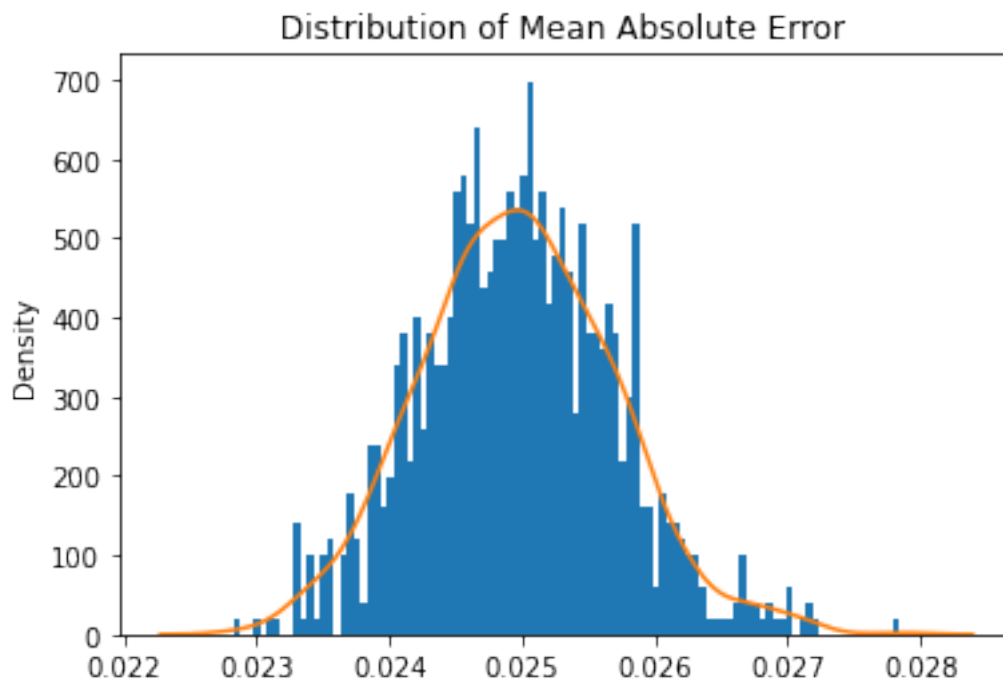
Number of epochs needed 1393



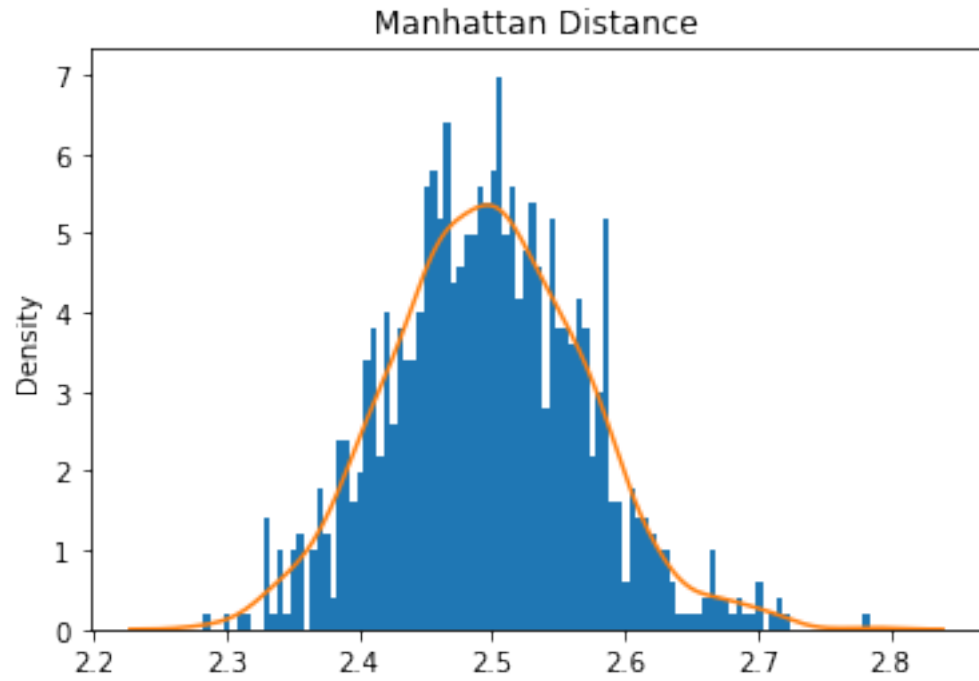
```
[16]: train_test.test_generator(generator,real_dataset,device)
```



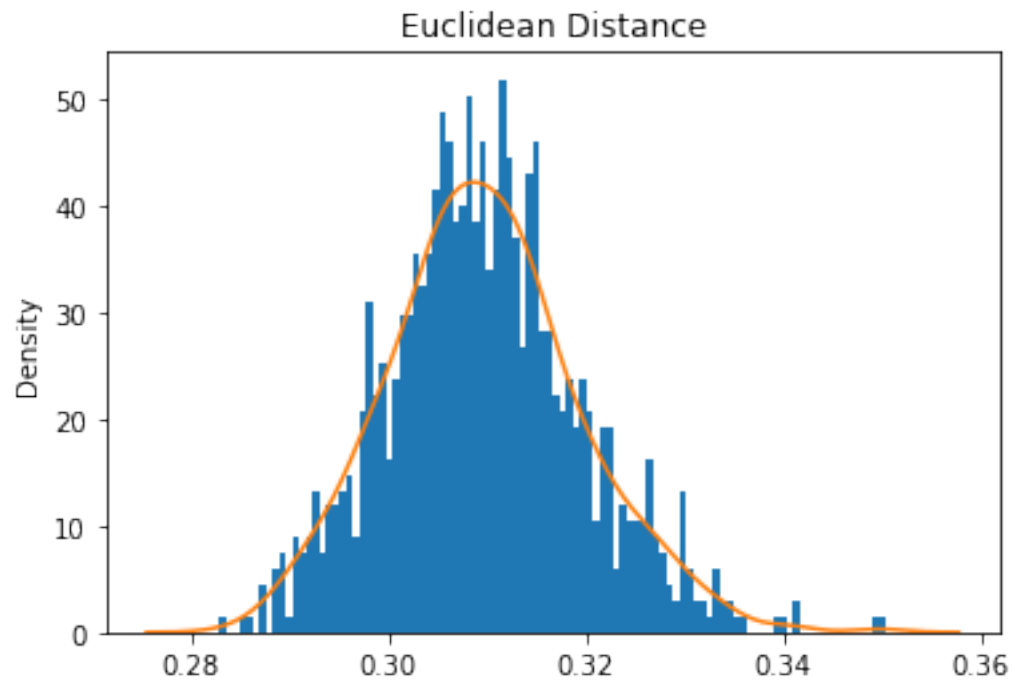
Mean Square Error: 0.0009602660747440065



Mean Absolute Error: 0.024967915923874824



Mean Manhattan Distance: 2.4967915923874826



Mean Euclidean Distance: 0.30972660059549284

2 ABC GAN Model

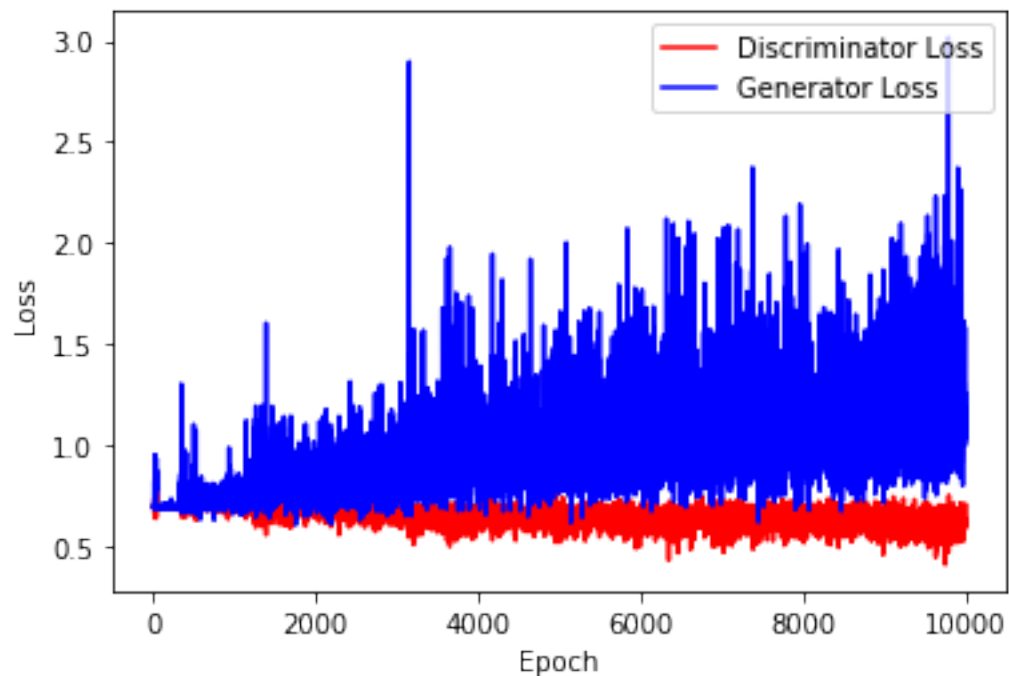
2.0.1 Training the network

Training ABC-GAN for `n_epochs` number of epochs

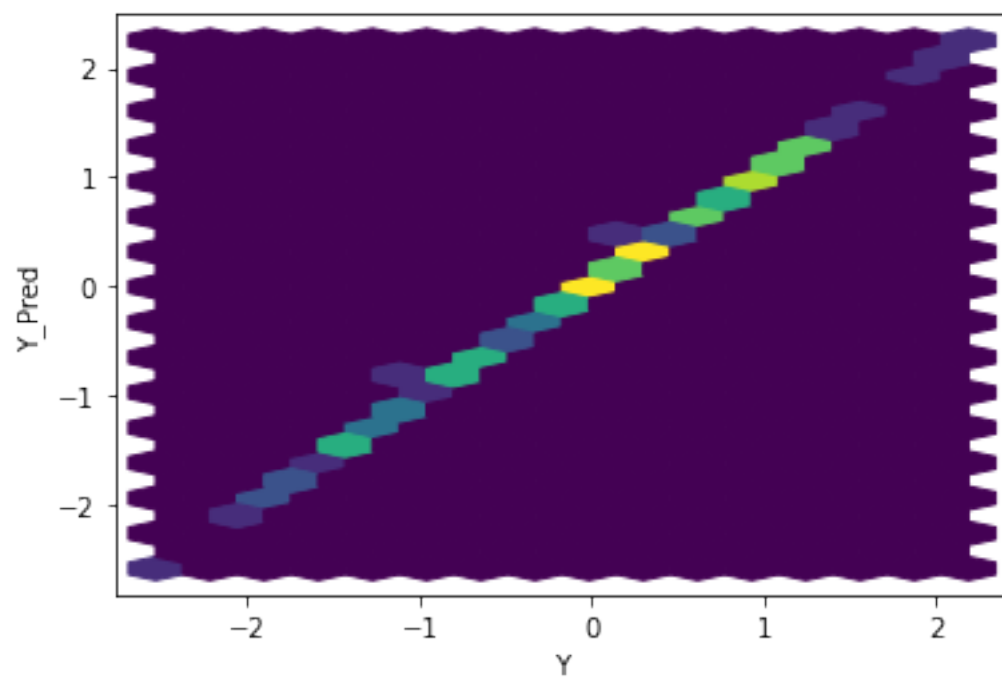
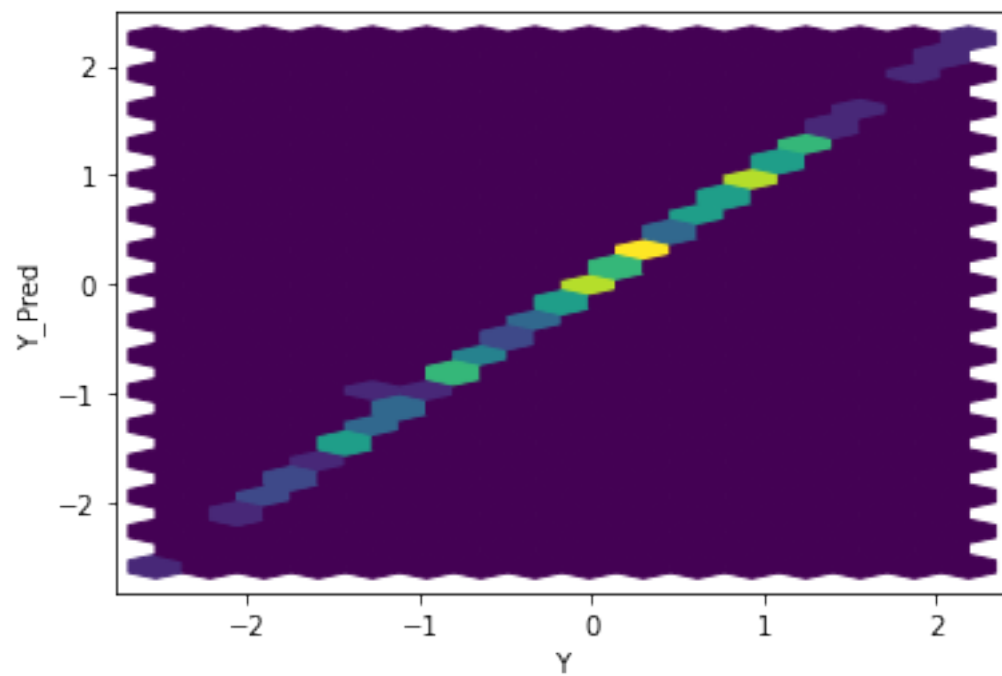
```
[17]: gen = network.Generator(n_features+2)
      disc = network.Discriminator(n_features+2)

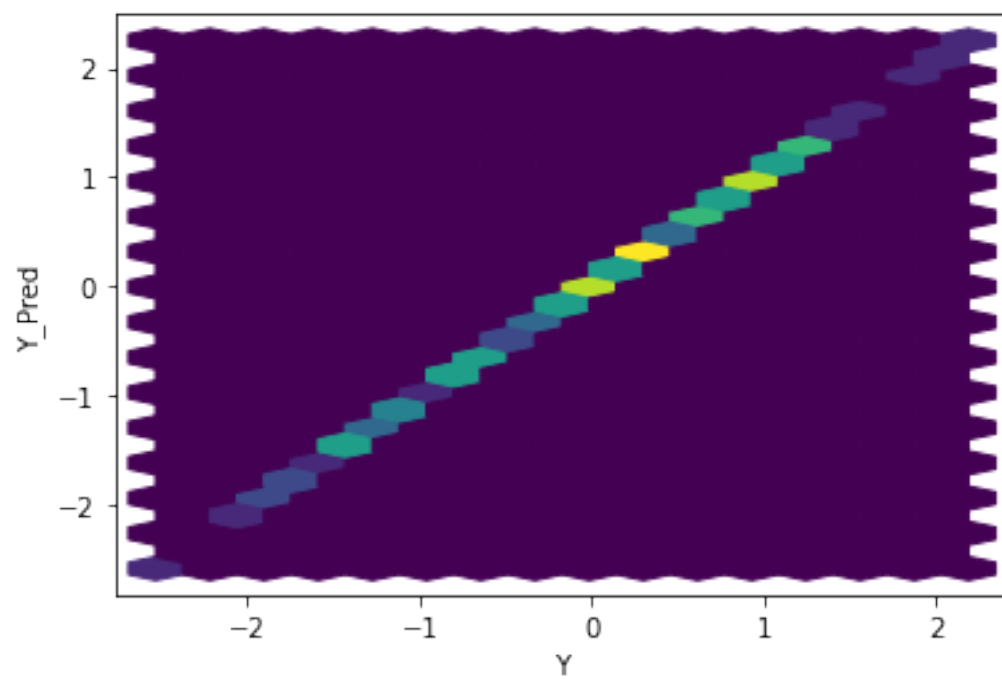
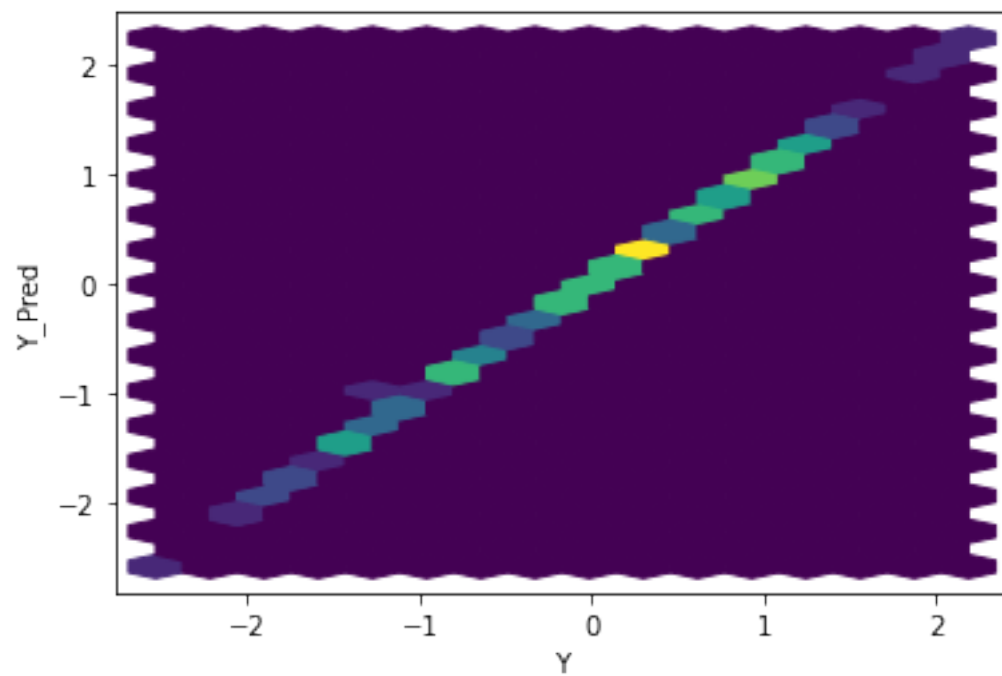
      criterion = torch.nn.BCEWithLogitsLoss()
      gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
      disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))

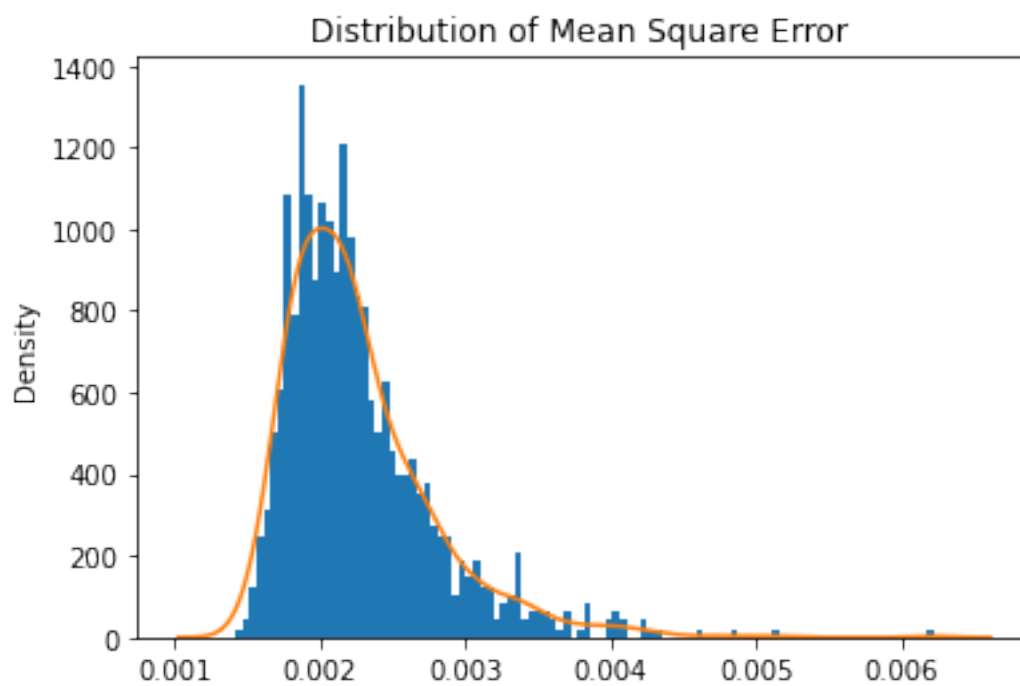
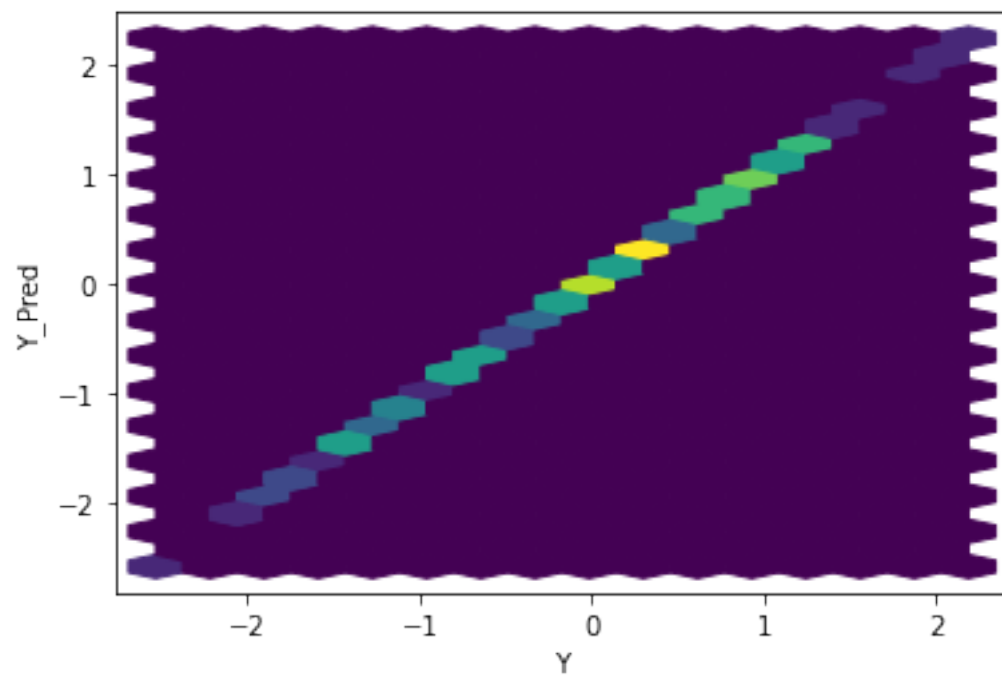
[18]: ABC_train_test.training_GAN(disc, gen, disc_opt, gen_opt, real_dataset,
      ↪ batch_size, n_epochs, criterion, coeff, mean, variance, device)
```



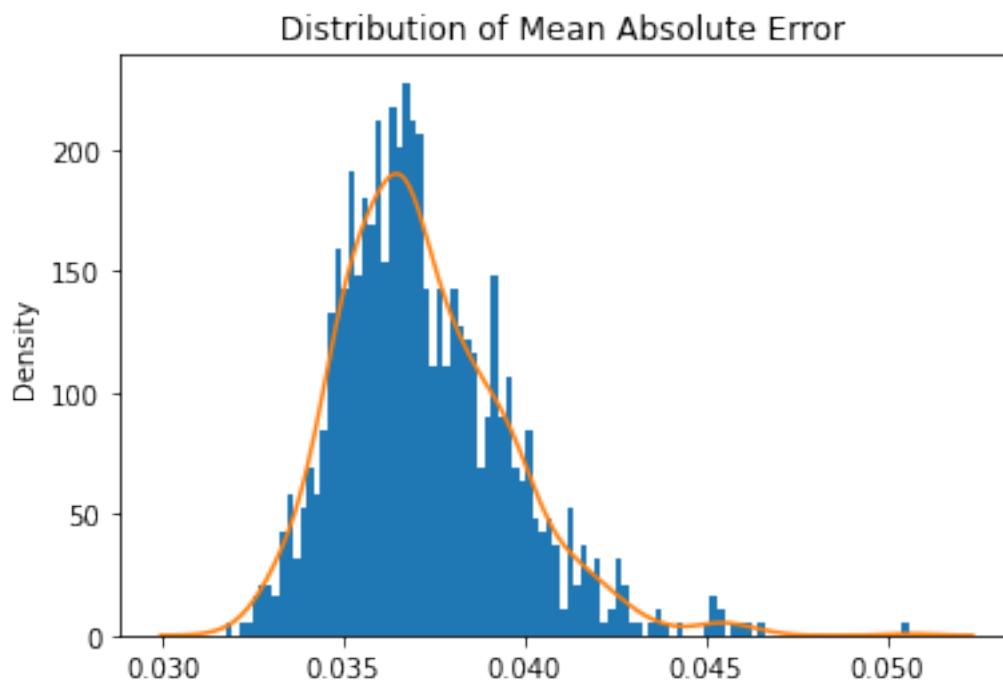
```
[19]: ABC_train_test.test_generator(gen, real_dataset, coeff, mean, variance, device)
```



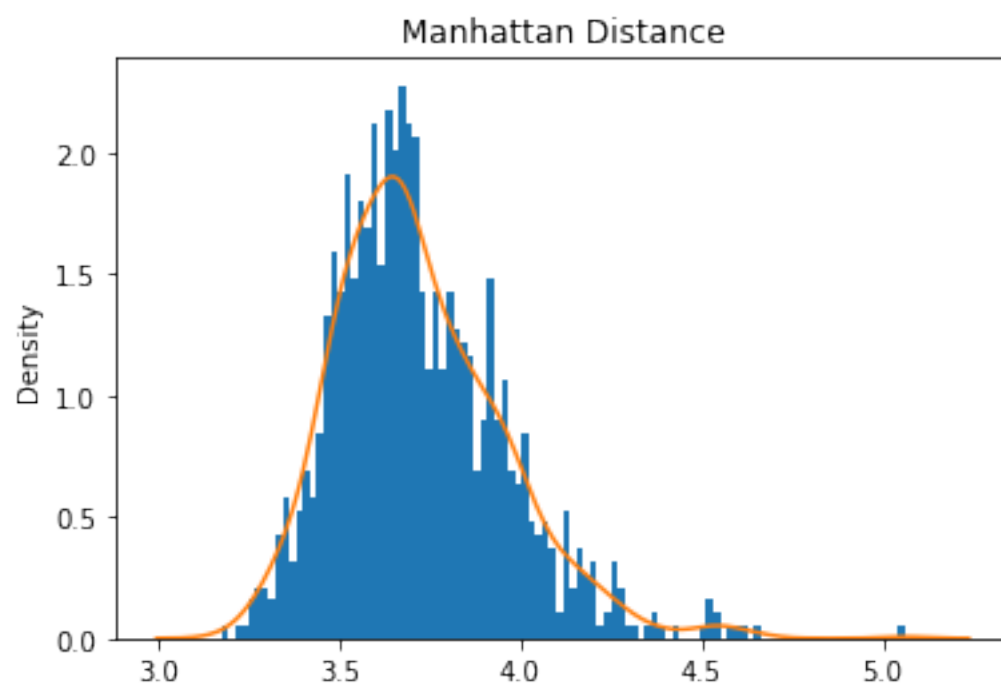




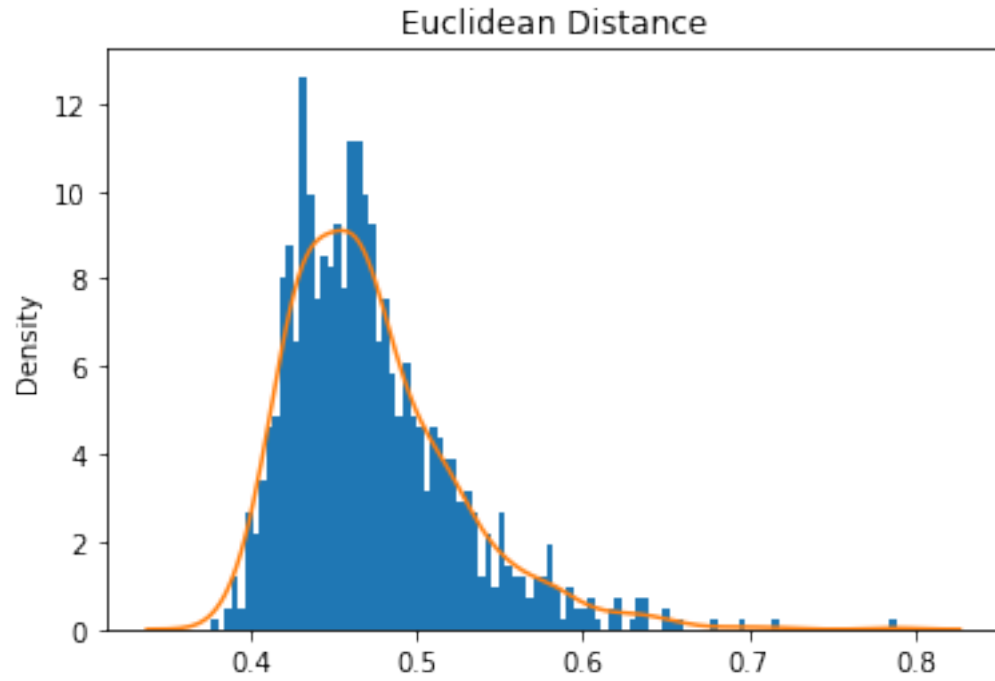
Mean Square Error: 0.002264222917581302



Mean Absolute Error: 0.037210186062548306
Mean Manhattan Distance: 3.721018606254831

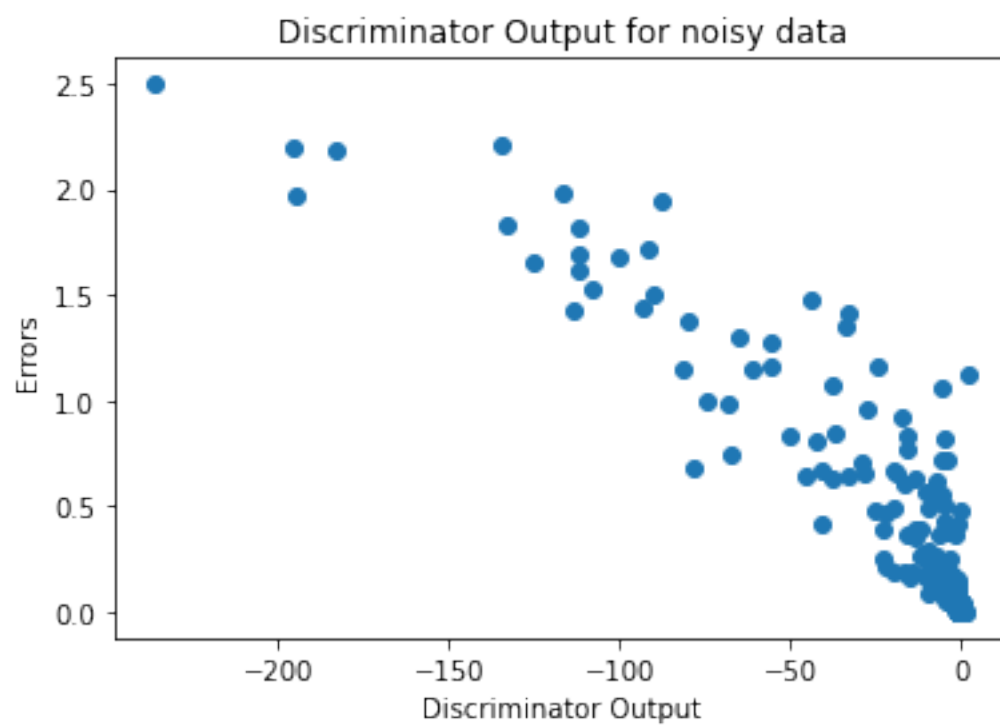
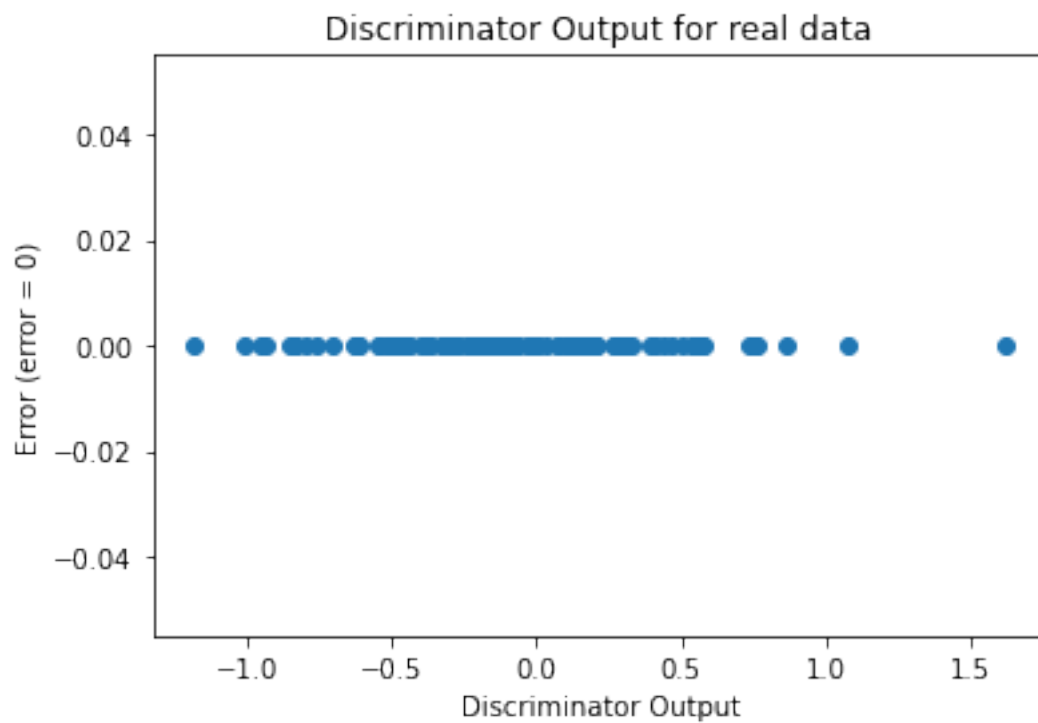


Mean Euclidean Distance: 0.47308031109808474



Sanity Checks

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



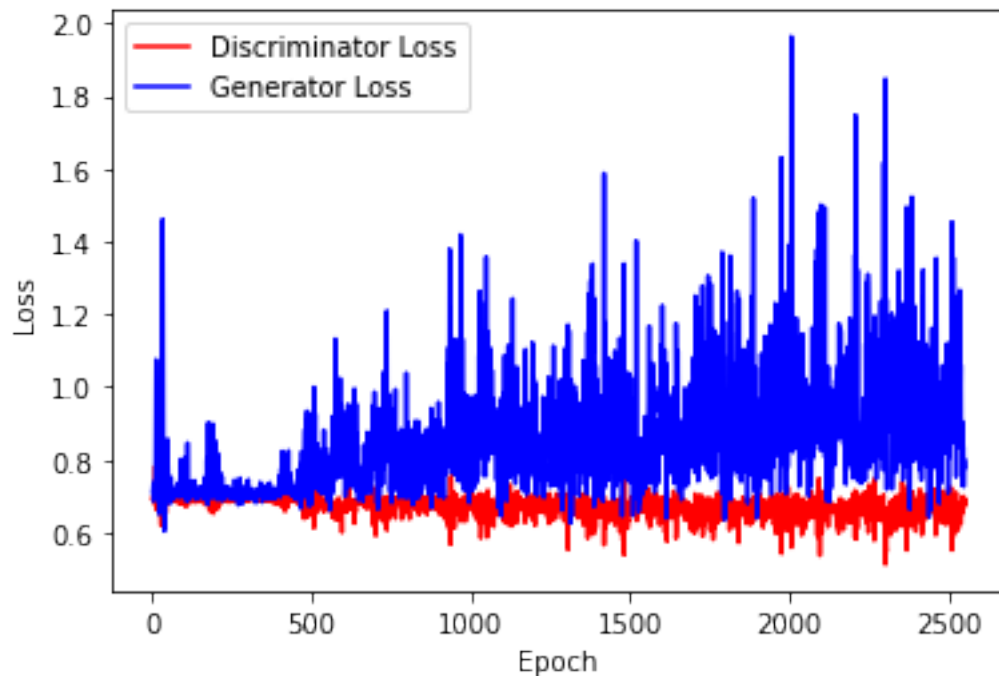
Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

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

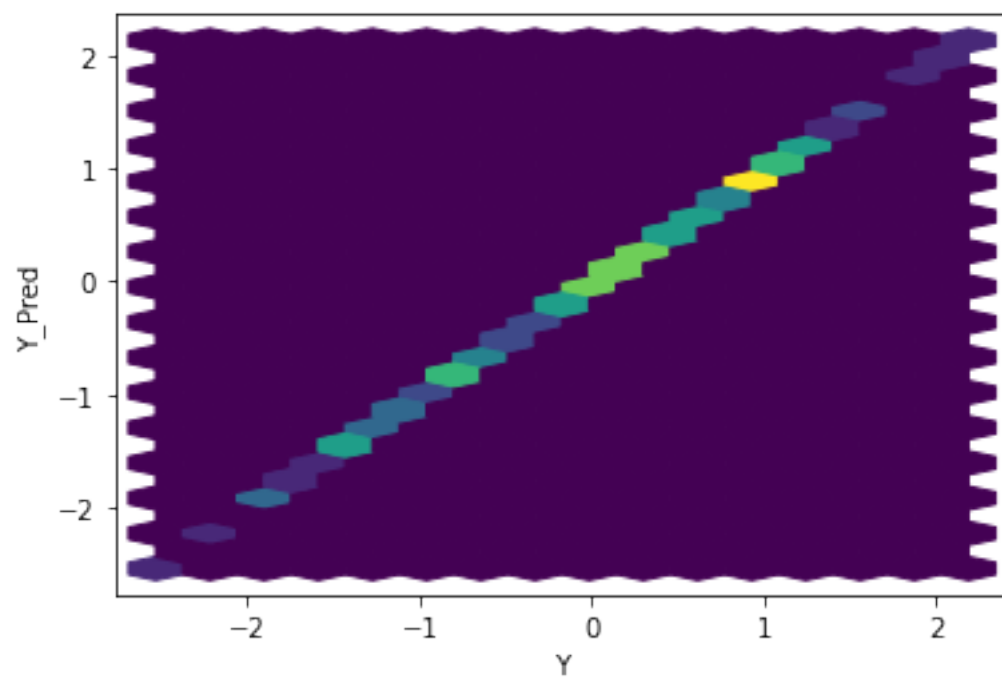
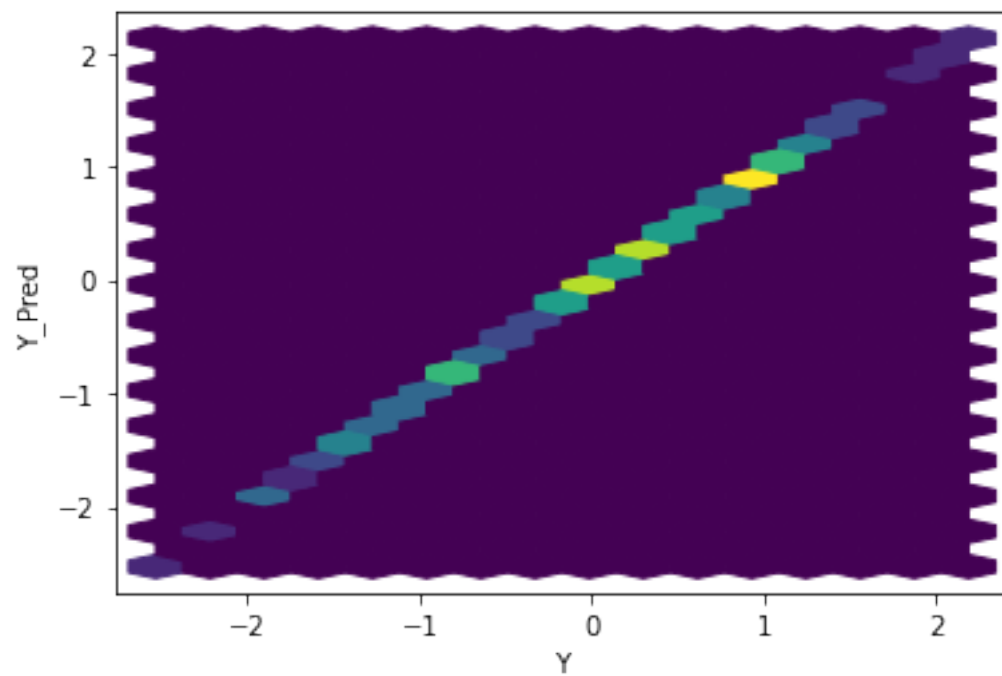
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))
```

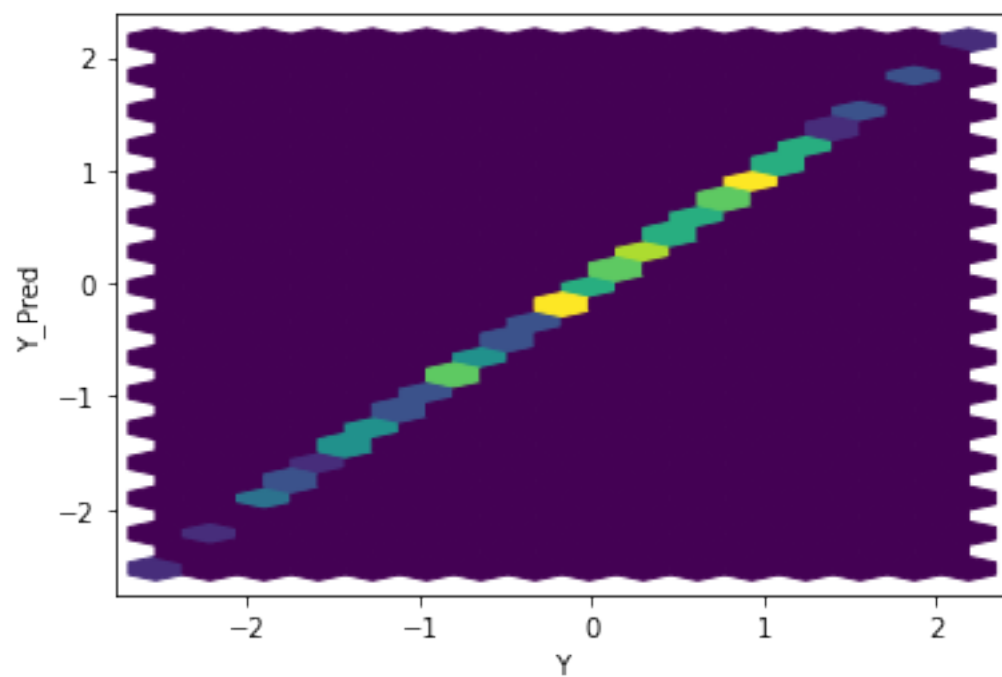
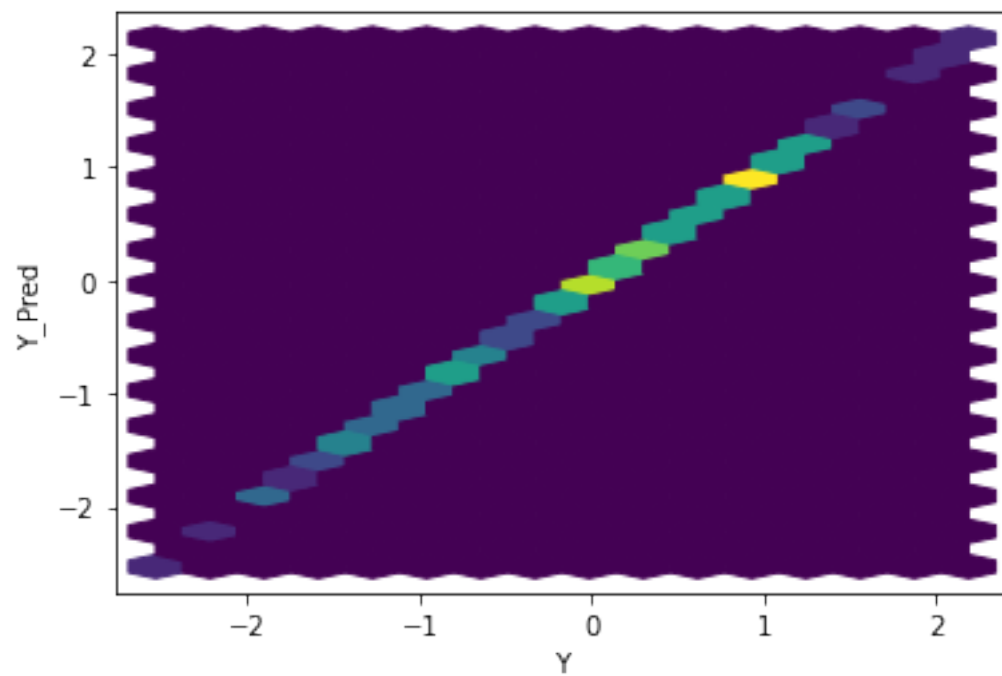
```
[22]: ABC_train_test.
      ↪ training_GAN_2(disc,gen,disc_opt,gen_opt,real_dataset,batch_size,
      ↪ error,criterion,coeff,mean,variance,device)
```

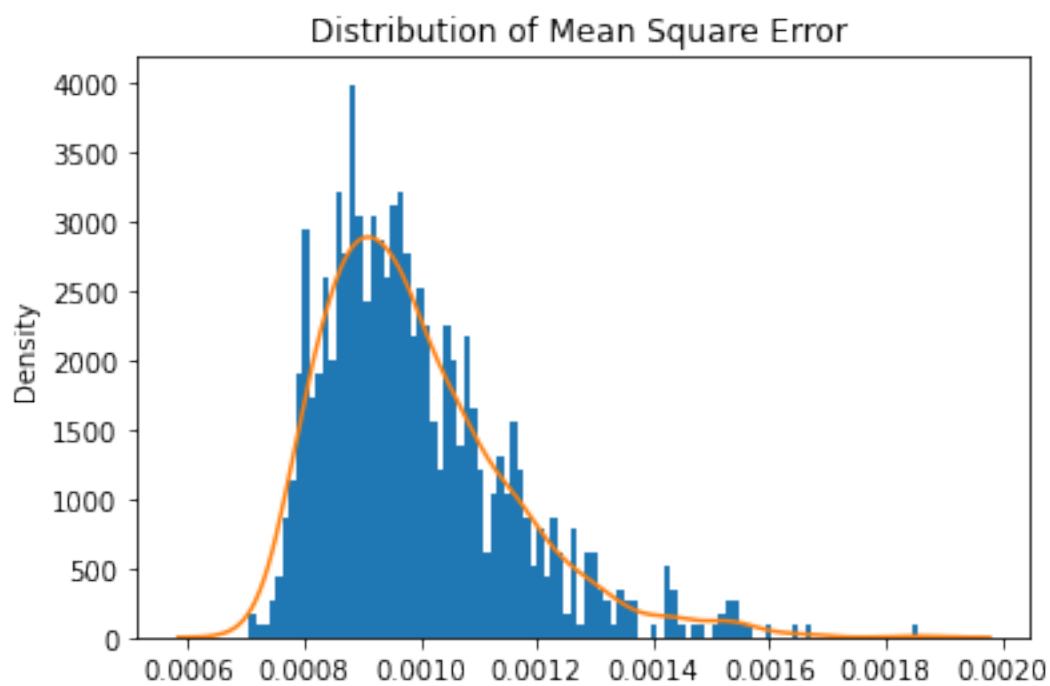
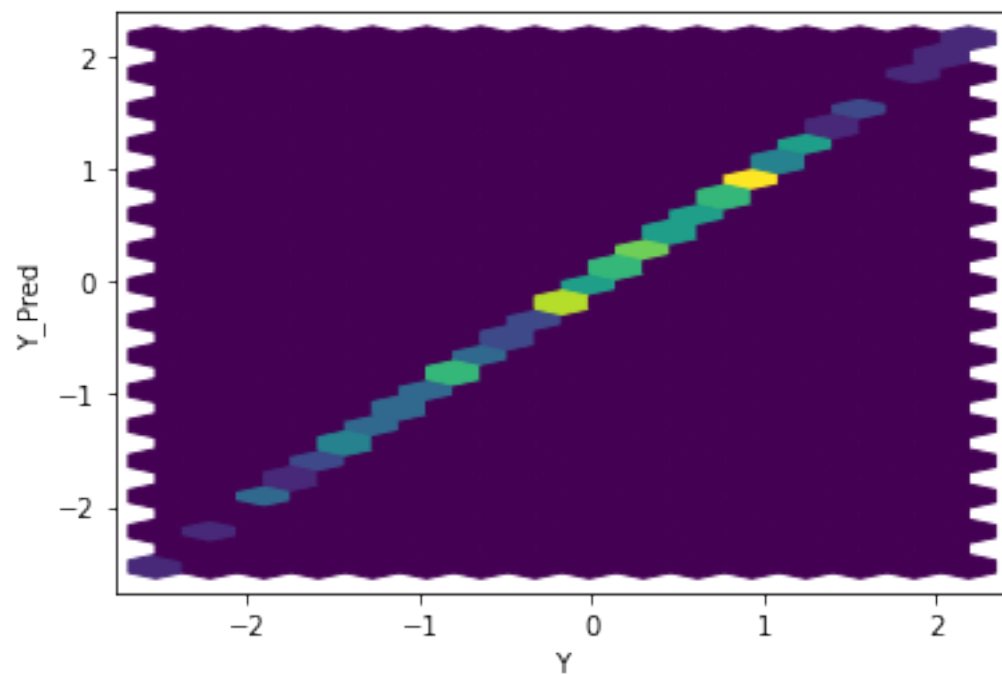
Number of epochs 1276



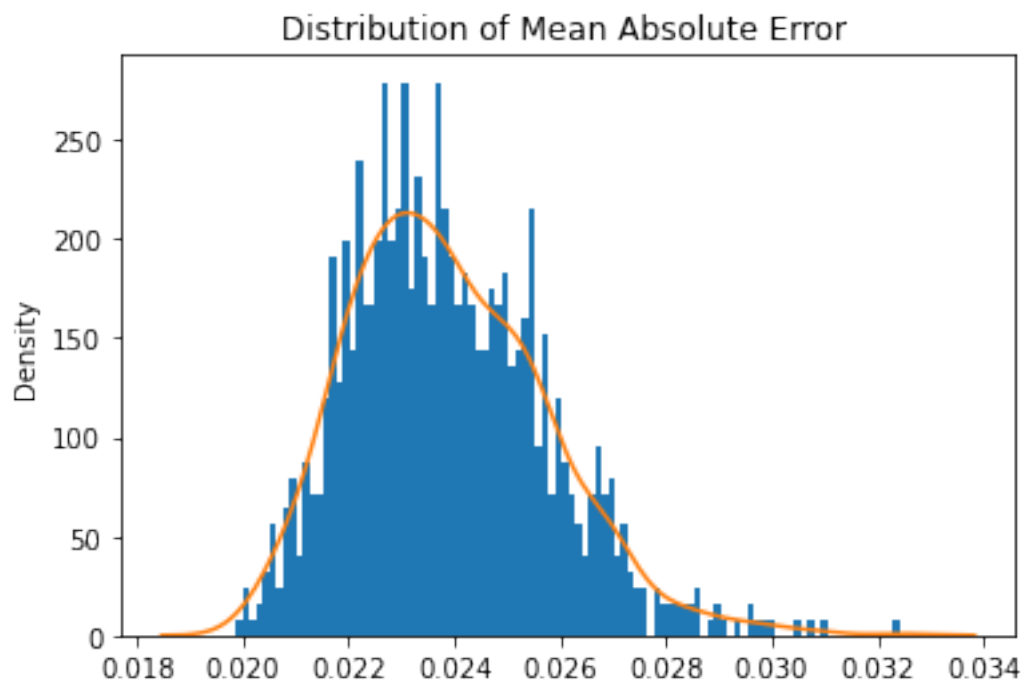
```
[23]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```





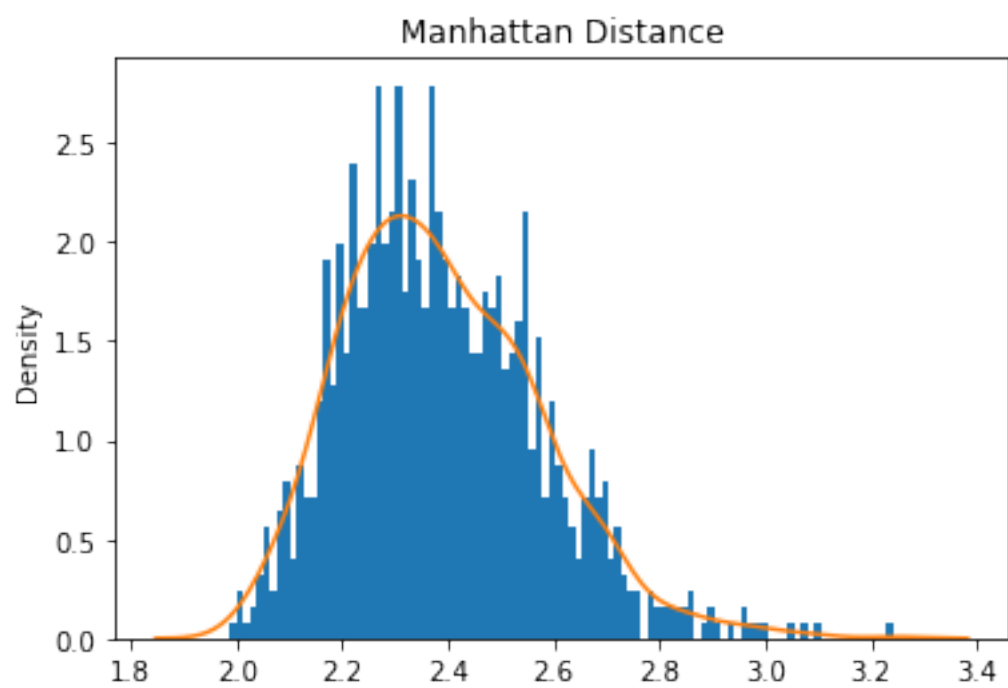


Mean Square Error: 0.0009921230884490796

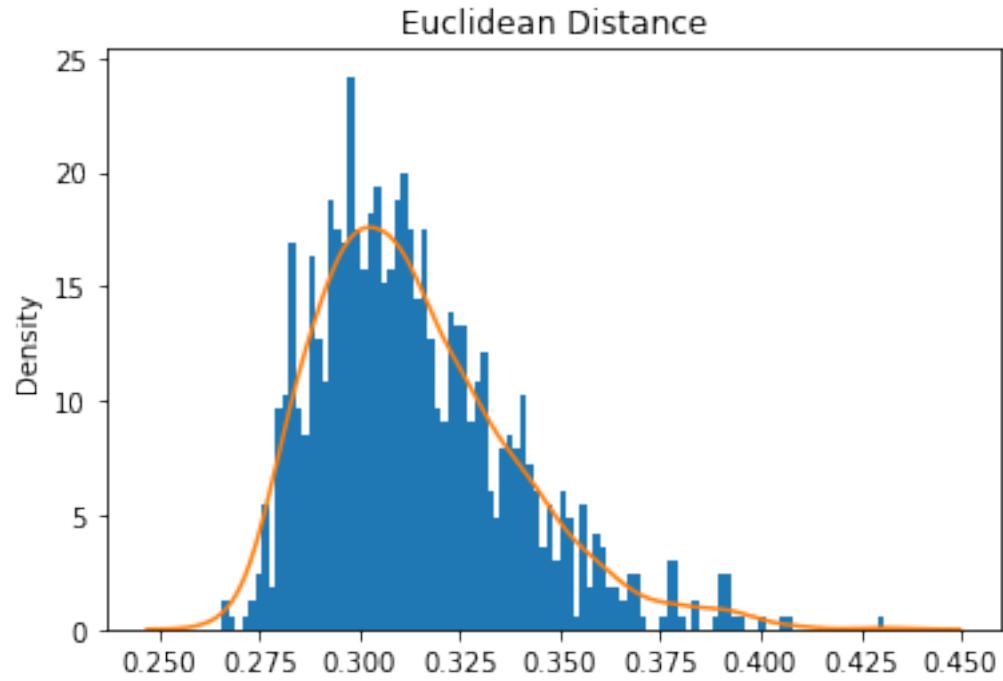


Mean Absolute Error: 0.023885777876460925

Mean Manhattan Distance: 2.3885777876460925



Mean Euclidean Distance: 0.31399841692887076



[]: