

# Dataset2\_Friedman1\_output\_11

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  $\beta_i^*$ 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 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  $\beta^*$  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 = 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	0.587534	0.715051	0.006270	0.959201	0.277320	0.660776	0.807930
1	0.773957	0.711747	0.747548	0.084442	0.017888	0.837288	0.742134
2	0.189832	0.032933	0.123804	0.921659	0.136819	0.364092	0.935144

```

3  0.902807  0.801179  0.518884  0.573619  0.311466  0.172338  0.410989
4  0.512190  0.175756  0.612142  0.973839  0.130293  0.019953  0.941551

```

```

          X7          X8          X9          Y
0  0.525464  0.209956  0.810544  25.597487
1  0.917705  0.380164  0.830480  11.876152
2  0.735924  0.305125  0.428422  12.866592
3  0.855542  0.122956  0.841124  14.912152
4  0.475752  0.738644  0.100452  13.500895

```

## 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:          0.716
Model:                OLS      Adj. R-squared:      0.684
Method:              Least Squares      F-statistic:      22.43
Date:                Wed, 20 Oct 2021      Prob (F-statistic):      2.62e-20
Time:                20:39:54      Log-Likelihood:      -78.965
No. Observations:      100      AIC:              179.9
Df Residuals:          89      BIC:              208.6
Df Model:              10
Covariance Type:      nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	5.169e-16	0.056	9.15e-15	1.000	-0.112	0.112
x1	0.3995	0.060	6.605	0.000	0.279	0.520
x2	0.3981	0.058	6.813	0.000	0.282	0.514
x3	-0.1006	0.059	-1.691	0.094	-0.219	0.018
x4	0.6723	0.060	11.203	0.000	0.553	0.791
x5	0.2779	0.059	4.725	0.000	0.161	0.395
x6	0.0096	0.058	0.166	0.868	-0.105	0.124
x7	-0.0621	0.058	-1.069	0.288	-0.178	0.053
x8	0.1024	0.057	1.784	0.078	-0.012	0.217
x9	0.0195	0.059	0.329	0.743	-0.098	0.137
x10	0.0561	0.059	0.952	0.344	-0.061	0.173

```

=====
Omnibus:              1.459      Durbin-Watson:          1.824
Prob(Omnibus):        0.482      Jarque-Bera (JB):        1.460
Skew:                 -0.281      Prob(JB):                0.482
Kurtosis:             2.813      Cond. No.                1.59
=====

```

Notes:

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

Parameters: const 5.169476e-16

x1 3.995024e-01

x2 3.980634e-01

x3 -1.006107e-01

x4 6.722560e-01

x5 2.778574e-01

x6 9.579475e-03

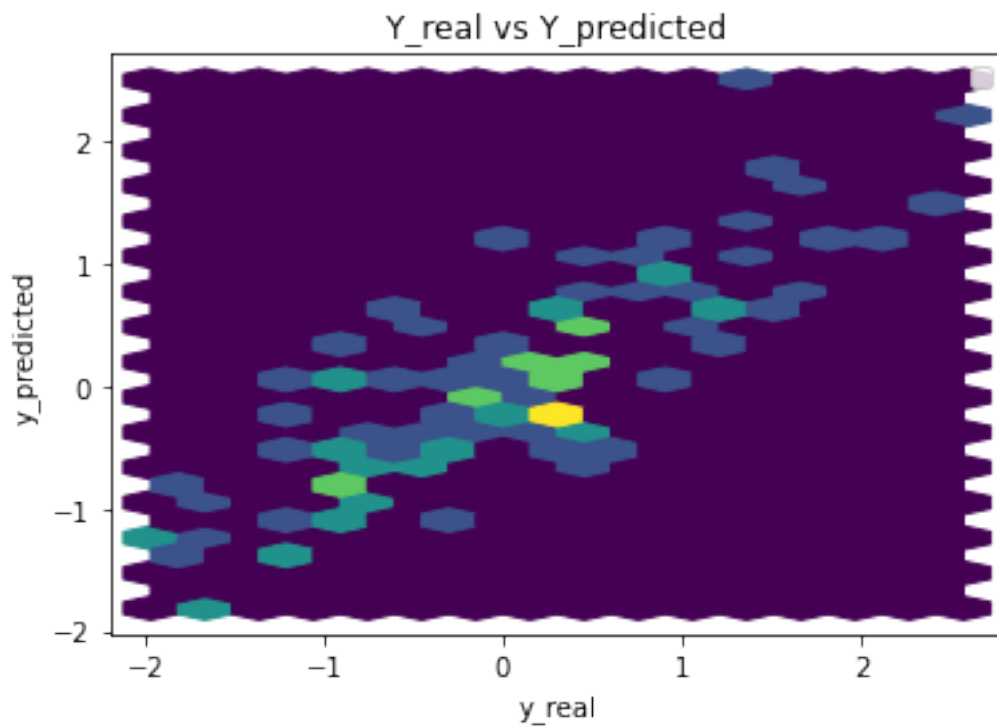
x7 -6.213483e-02

x8 1.024370e-01

x9 1.946343e-02

x10 5.606967e-02

dtype: float64



Performance Metrics

Mean Squared Error: 0.2840579184435031

Mean Absolute Error: 0.4092796782381576

Manhattan distance: 40.92796782381575

Euclidean distance: 5.329708420199955

## 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

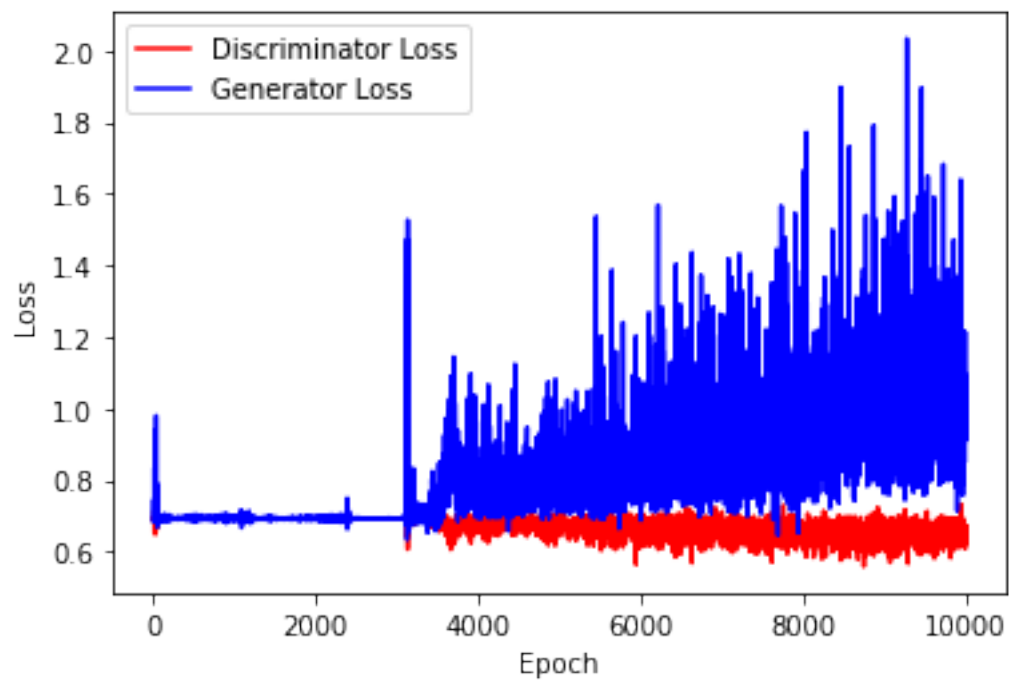
```
[9]: generator = network.Generator(n_features+2)
     discriminator = network.Discriminator(n_features+2)

     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))
```

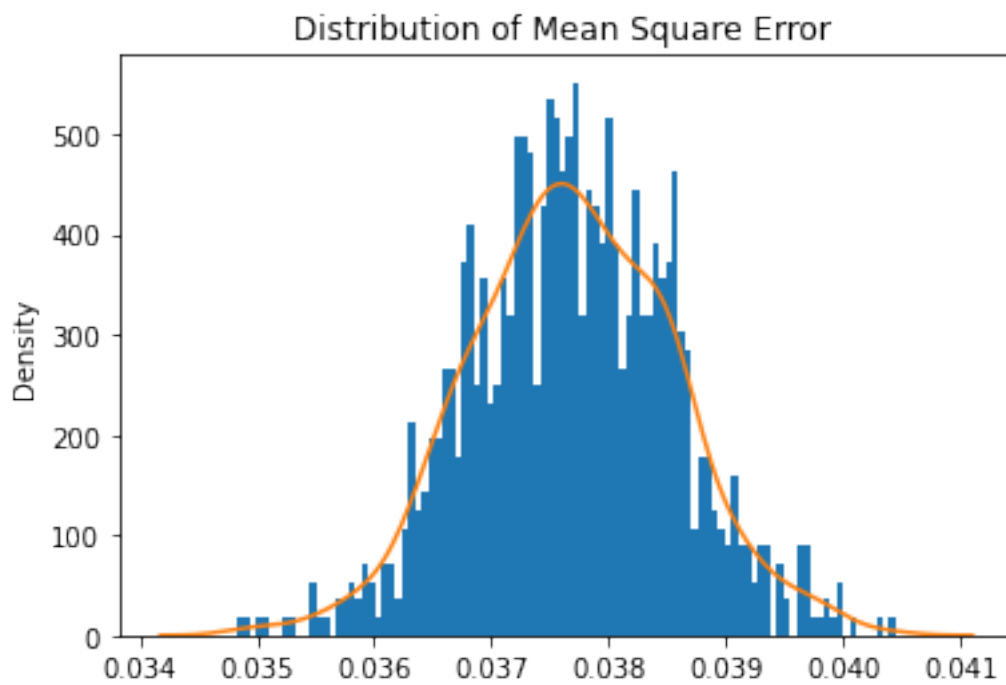
```
[10]: print(generator)
      print(discriminator)
```

```
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()
)
```

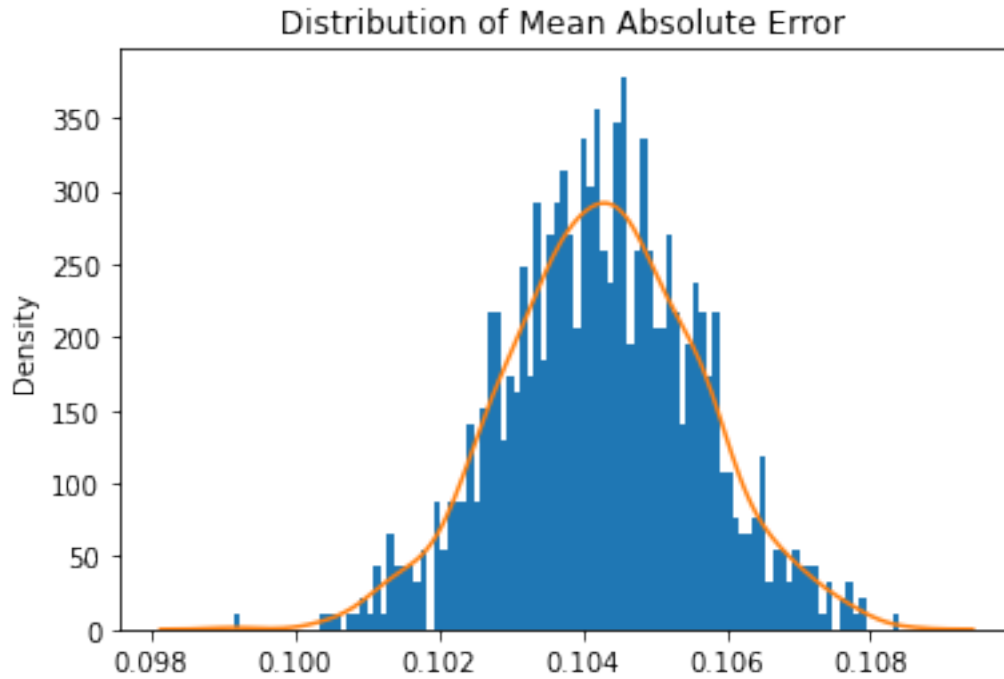
```
[11]: train_test.
     ↪training_GAN(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,
     ↪n_epochs,criterion,device)
```



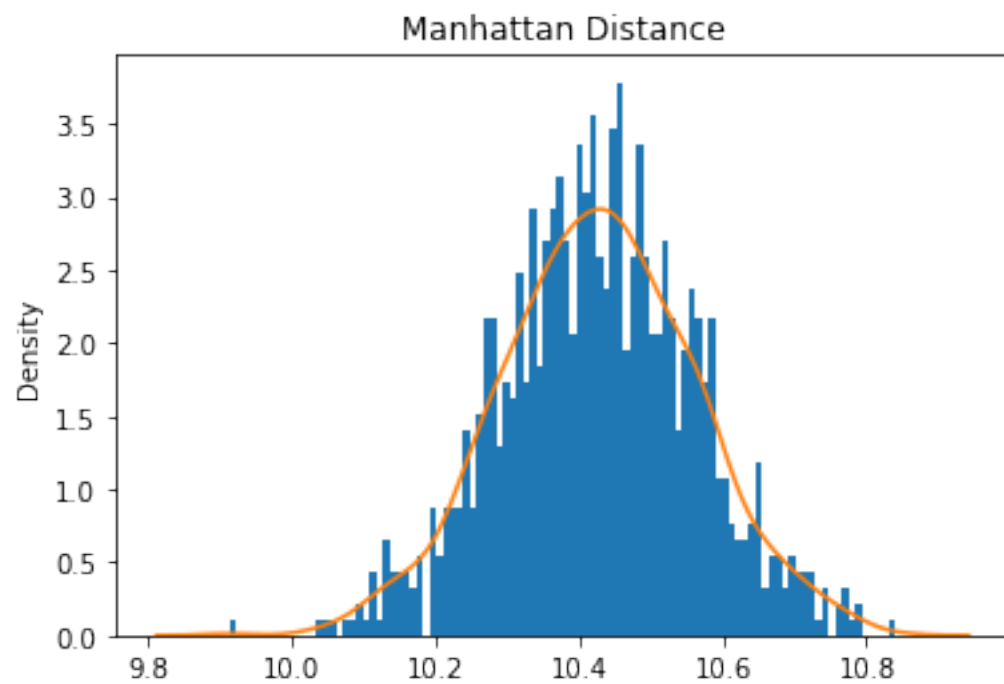
```
[12]: train_test.test_generator(generator,real_dataset,device)
```



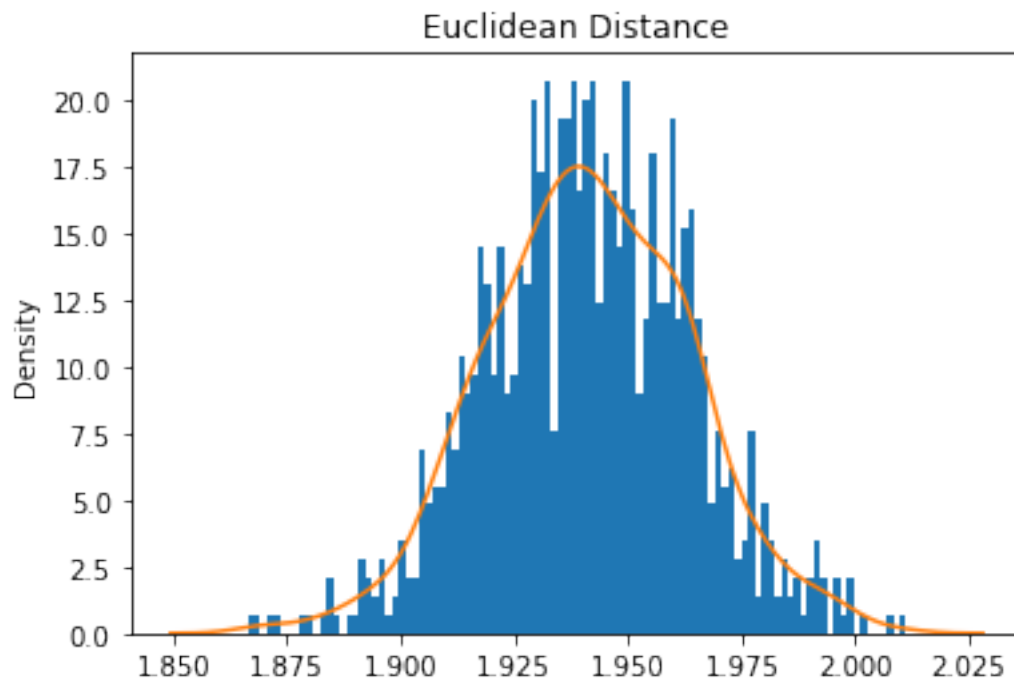
Mean Square Error: 0.03768611089559864



Mean Absolute Error: 0.10426962545337155



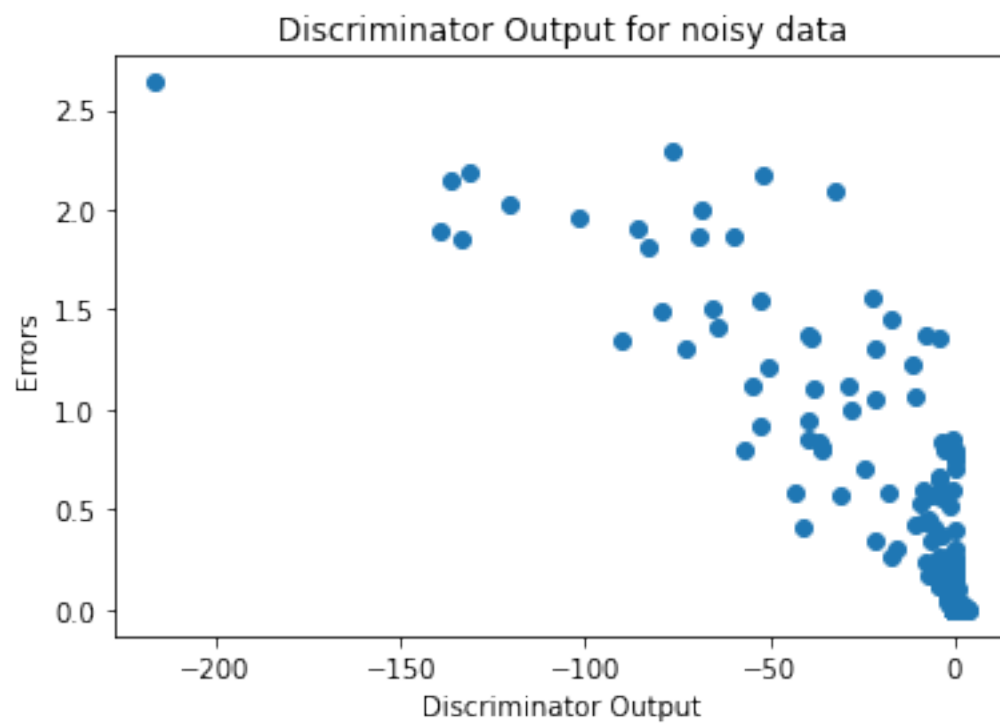
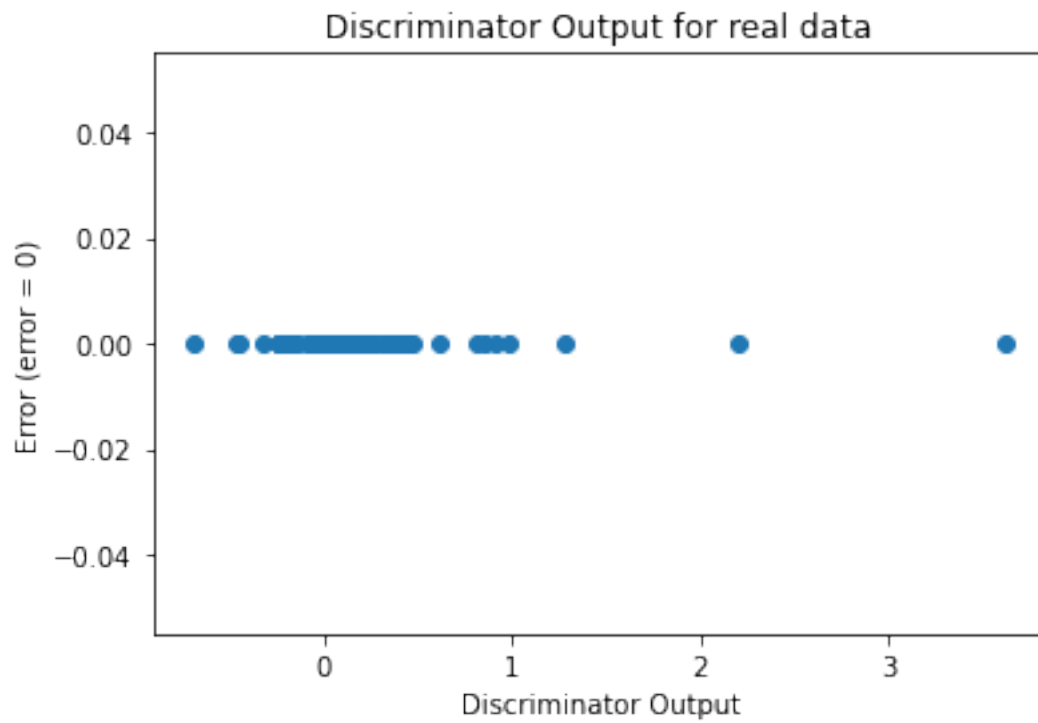
Mean Manhattan Distance: 10.426962545337155



Mean Euclidean Distance: 1.94116232170135

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



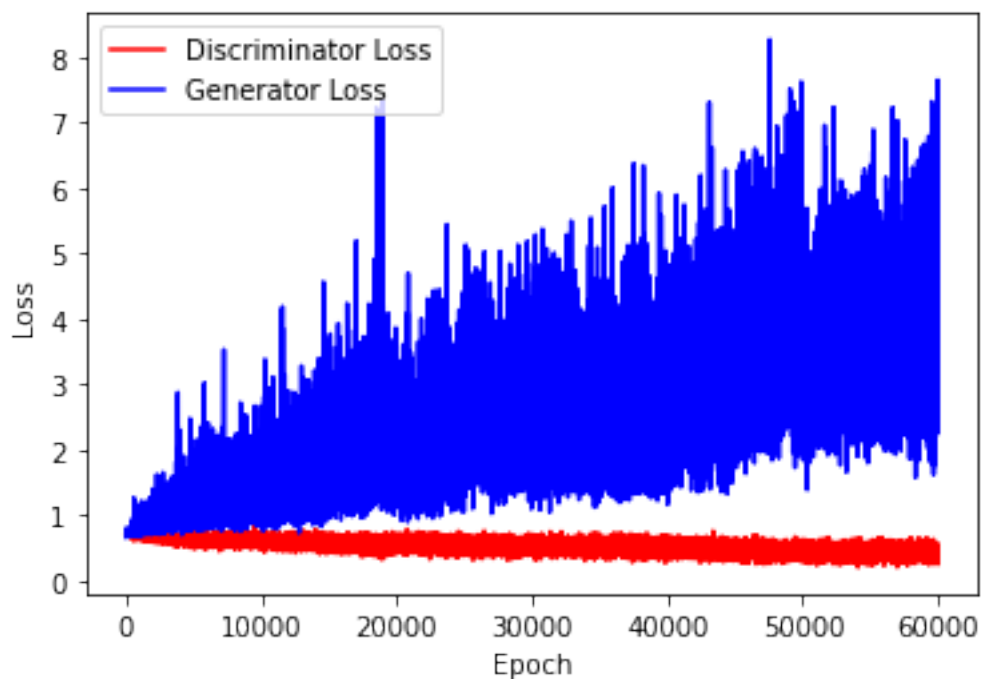


Training GAN until mse of  $y\_pred$  is  $> 0.1$  or  $n\_epochs < 30000$

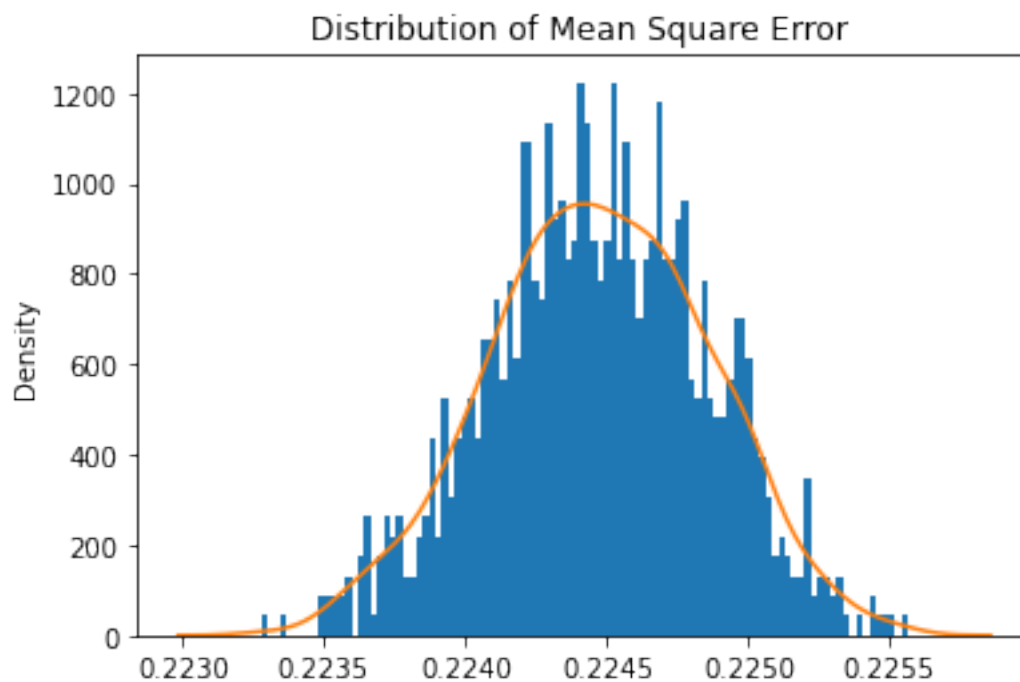
```
[14]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator.parameters(), lr=0.01, betas=(0.5,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

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

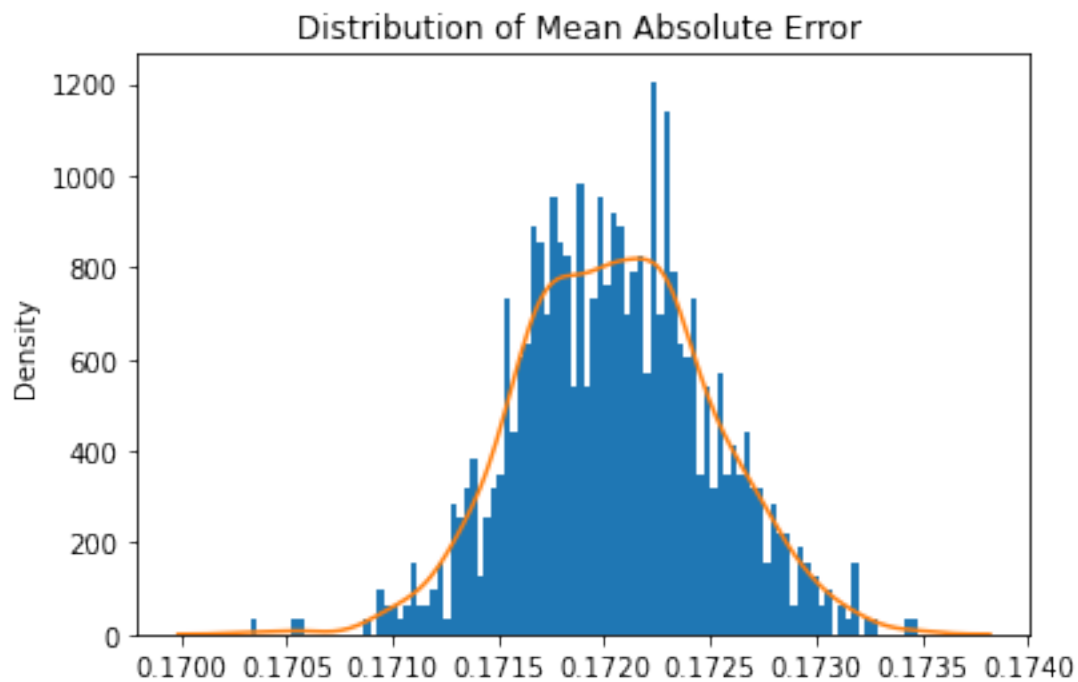
Number of epochs needed 30000



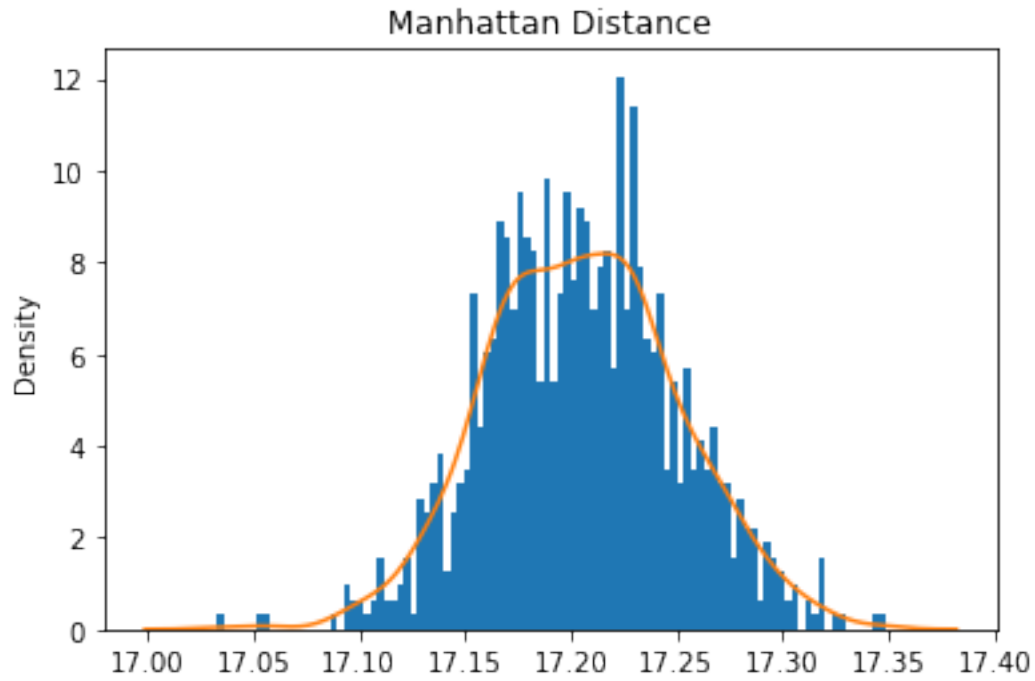
```
[16]: train_test.test_generator(generator,real_dataset,device)
```



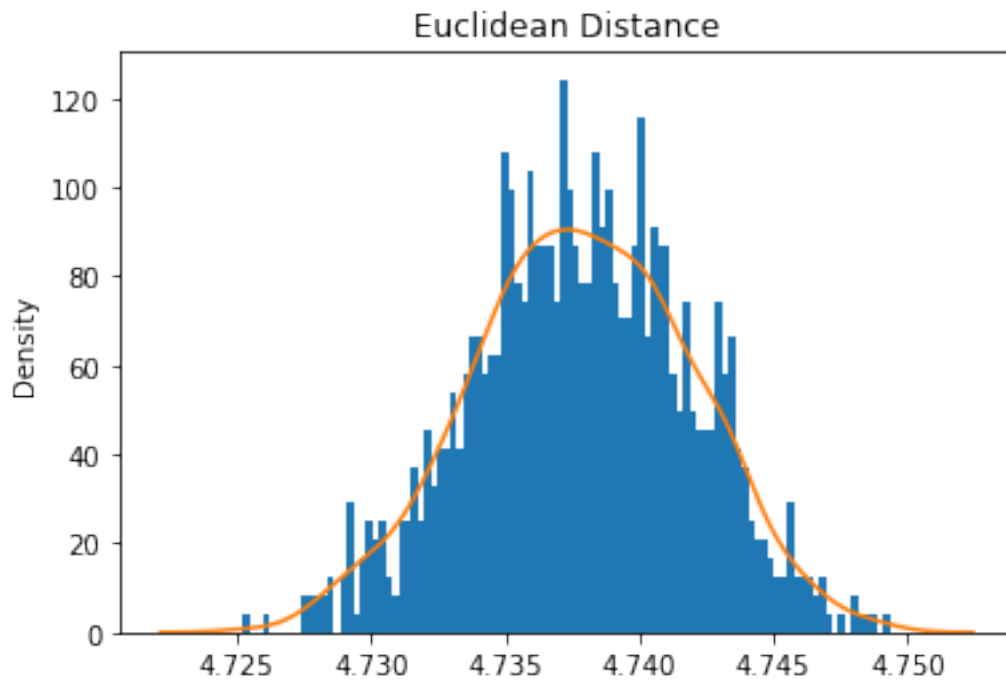
Mean Square Error: 0.22447027233679934



Mean Absolute Error: 0.17205319996088744



Mean Manhattan Distance: 17.205319996088743



Mean Euclidean Distance: 4.737827625426793

## 2 ABC GAN Model

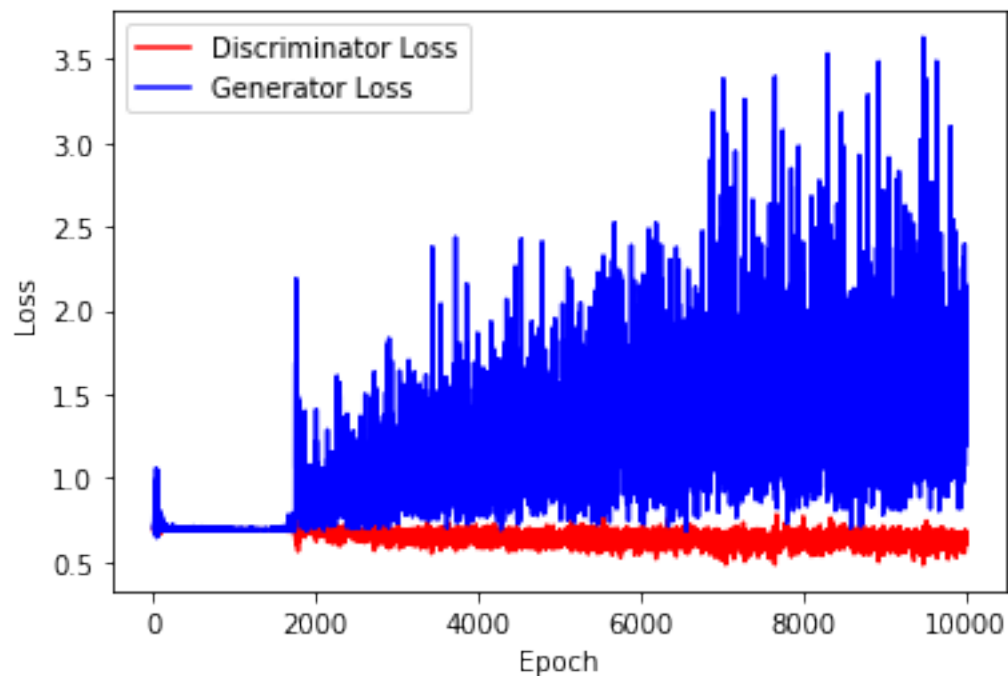
### 2.0.1 Training the network

Training ABC-GAN for `n_epochs` number of epochs

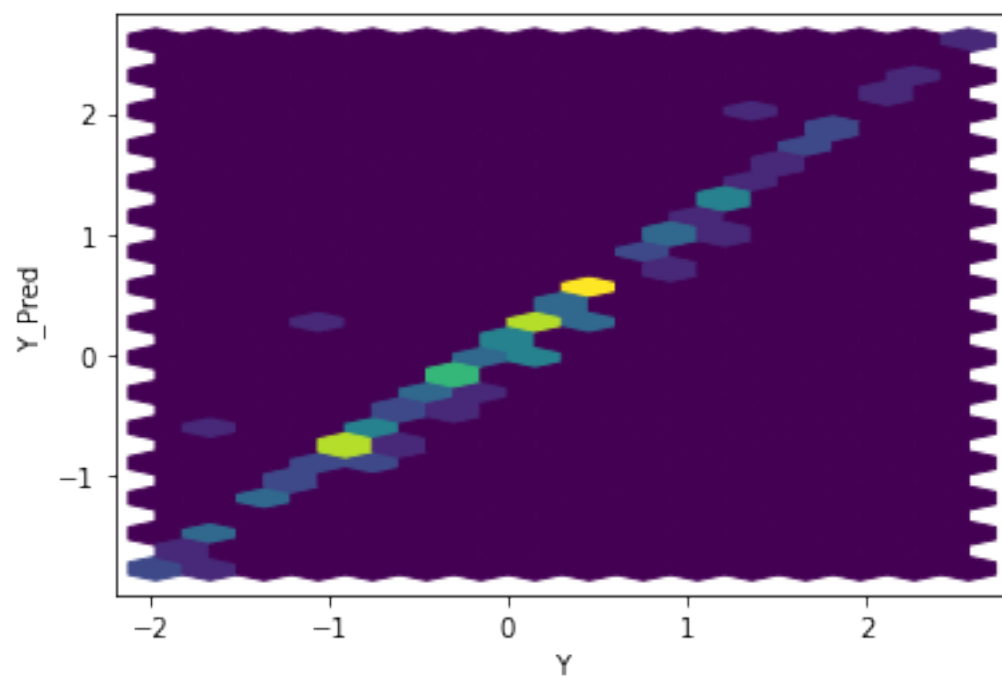
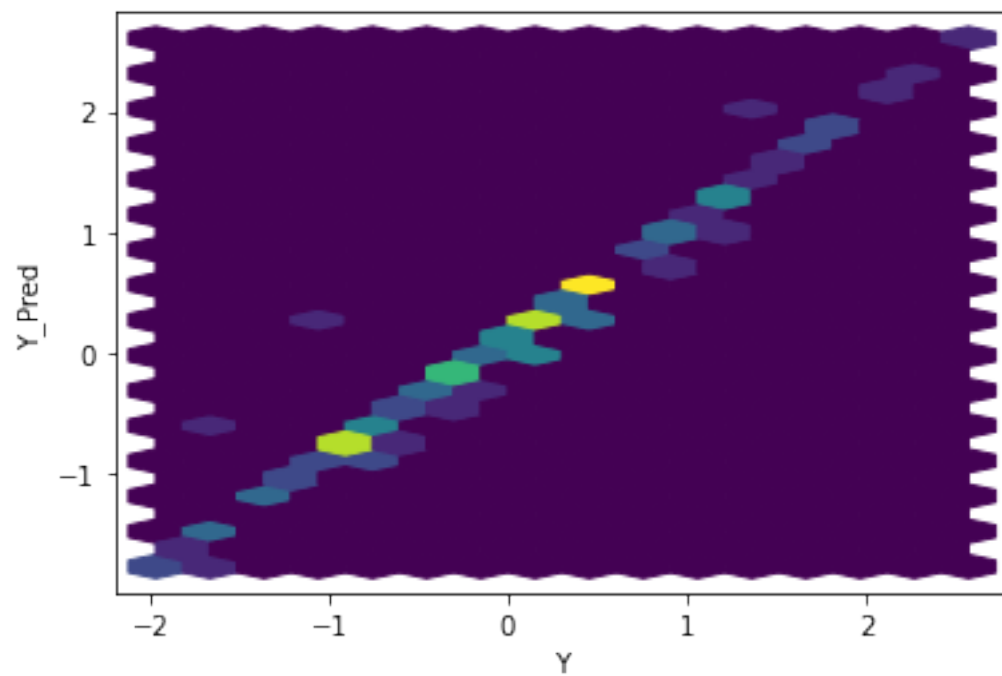
```
[17]: gen = network.Generator(n_features+2)
      disc = network.Discriminator(n_features+2)

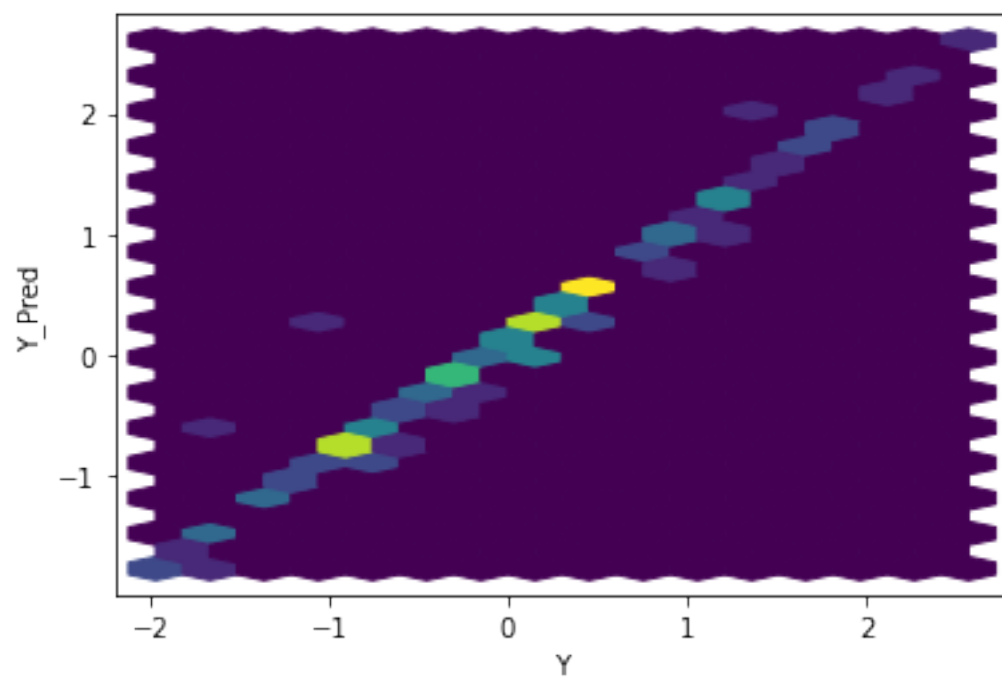
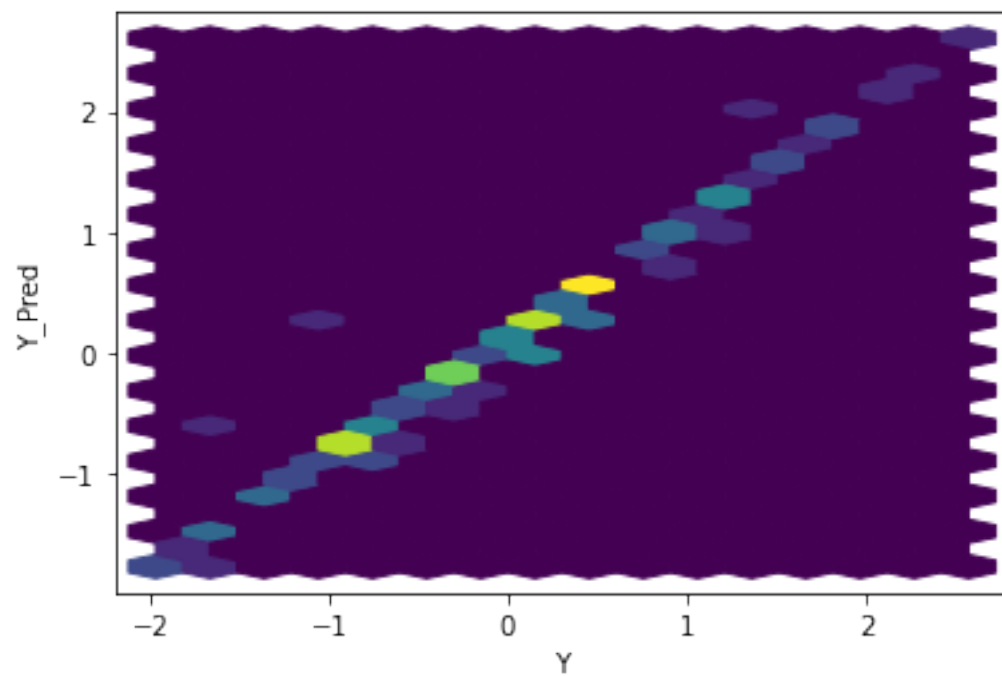
      criterion = torch.nn.BCEWithLogitsLoss()
      gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
      disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))

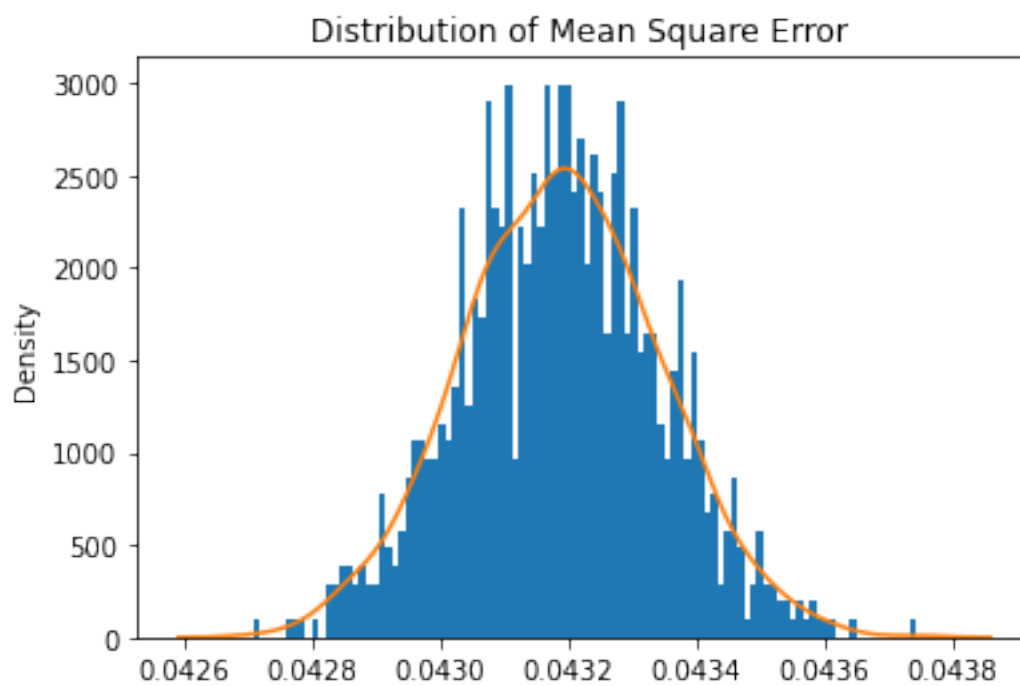
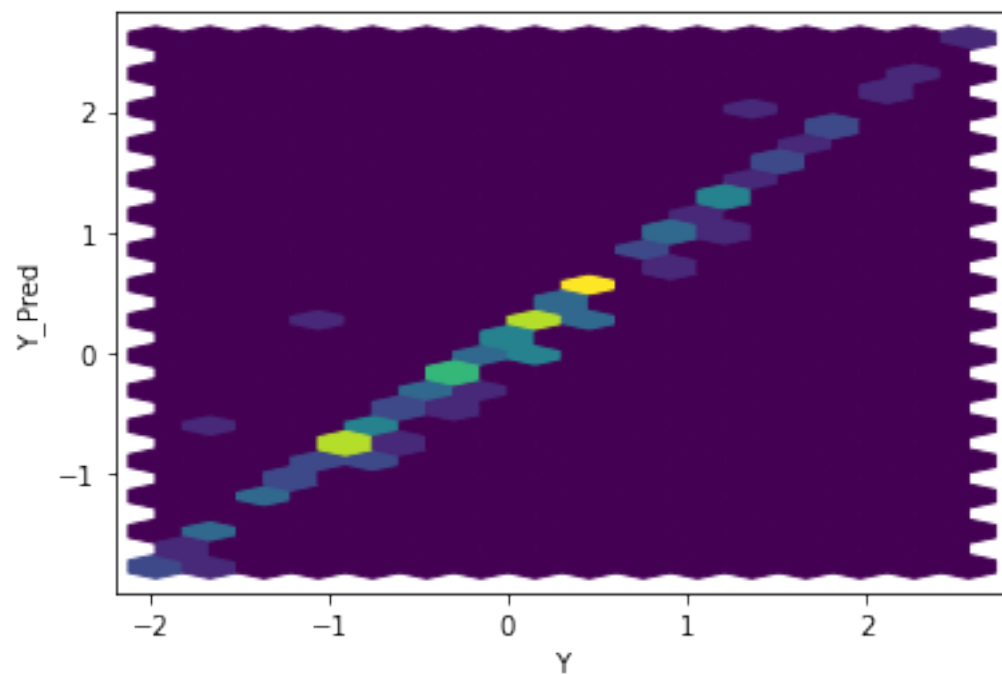
[18]: ABC_train_test.training_GAN(disc, gen,disc_opt,gen_opt,real_dataset,
      ↪batch_size, n_epochs,criterion,coeff,mean,variance,device)
```



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

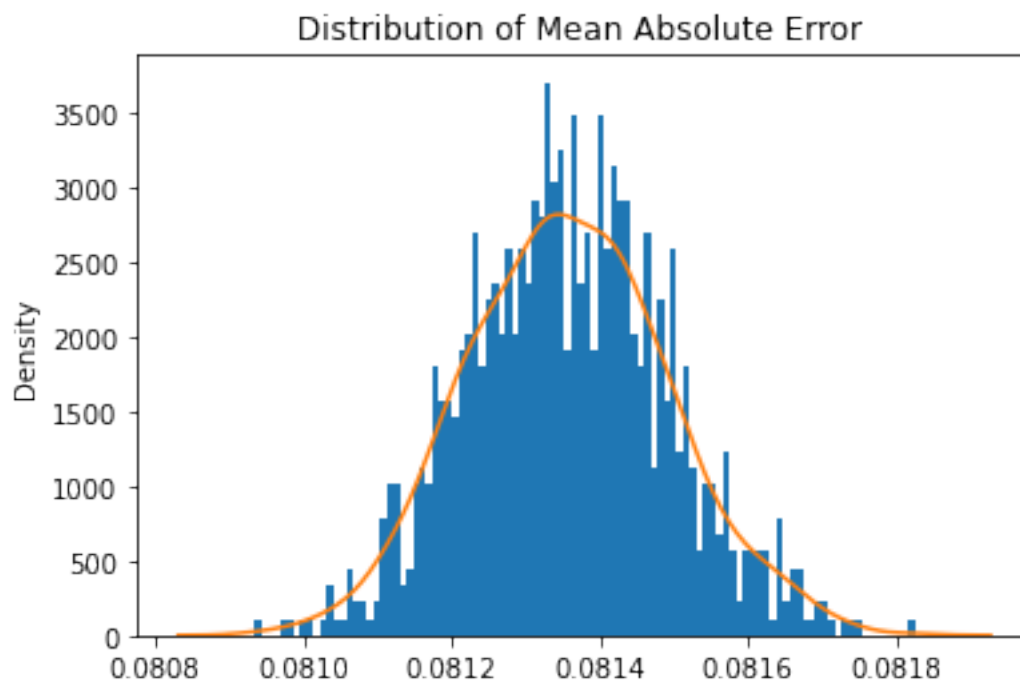




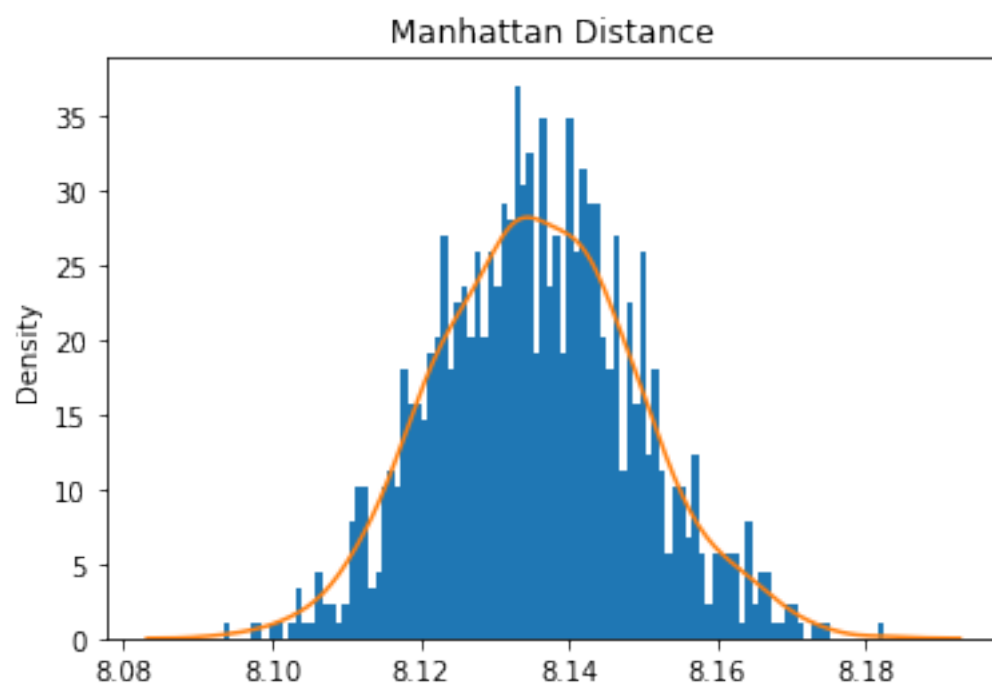


Mean Square Error: 0.04318766302579169

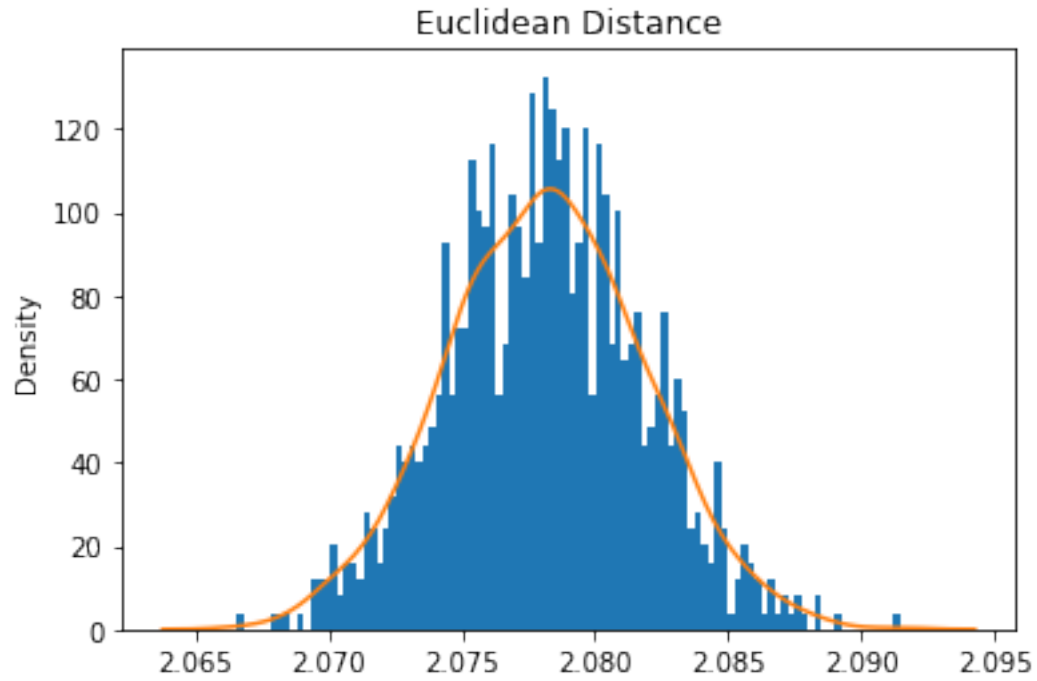




Mean Absolute Error: 0.08135732149664313  
Mean Manhattan Distance: 8.135732149664312

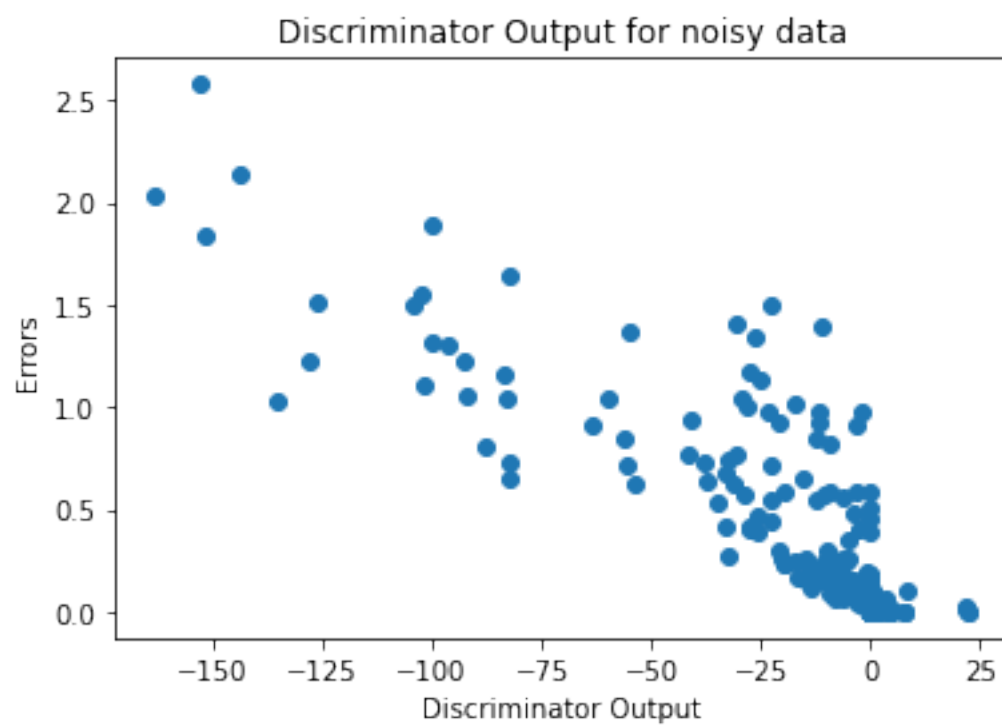
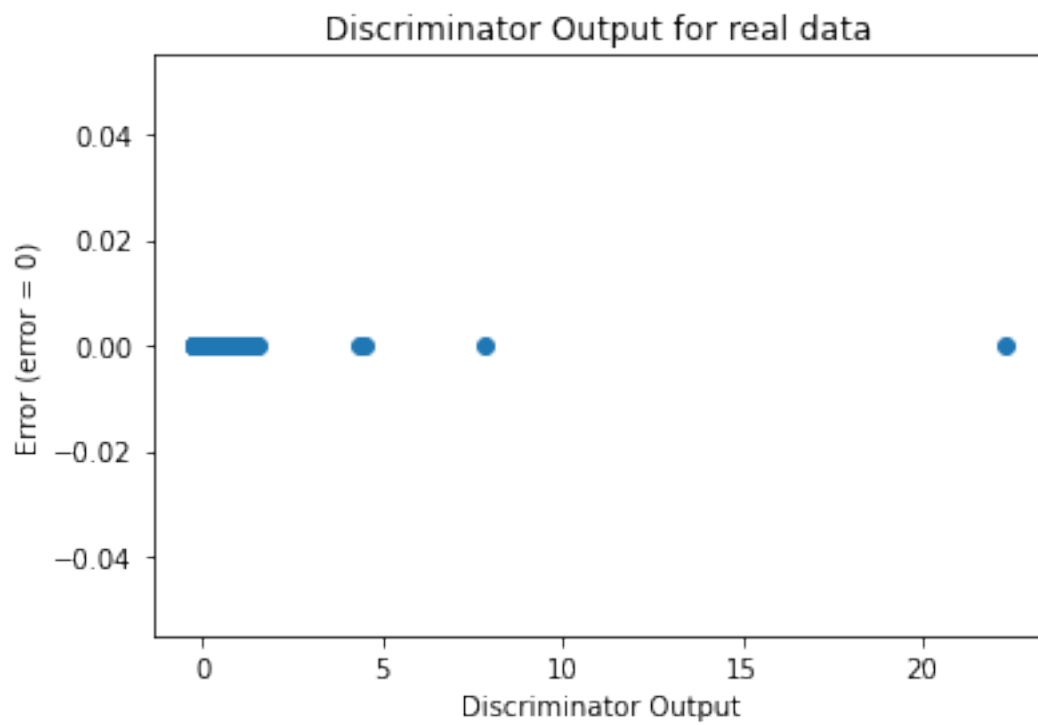


Mean Euclidean Distance: 2.078160866216521



### Sanity Checks

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



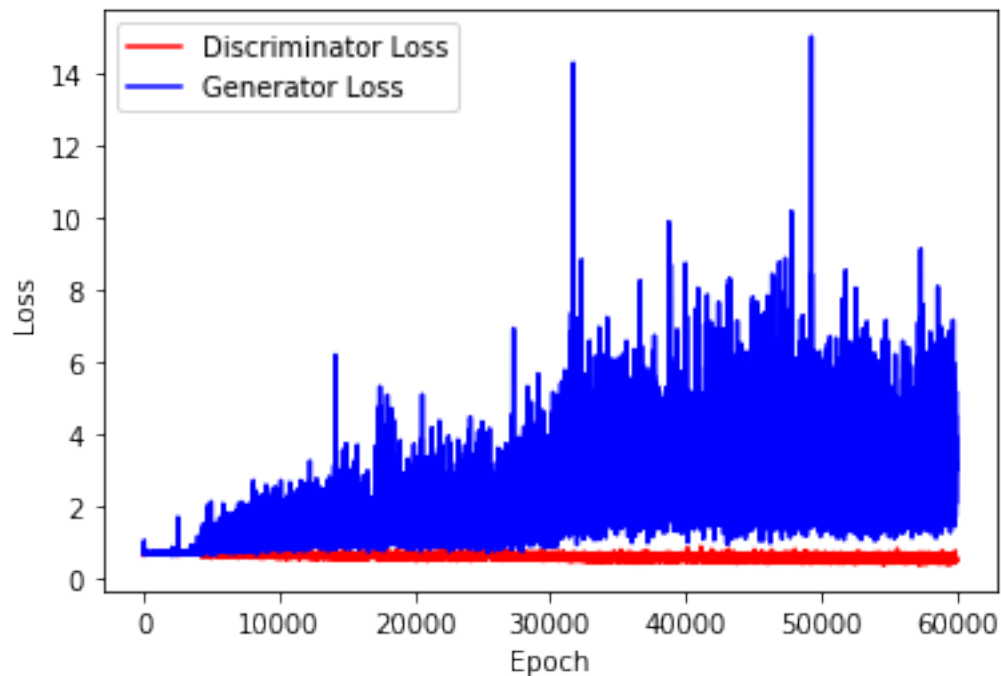
Training GAN until mse of y\_pred is  $> 0.1$  or n\_epochs  $< 30000$

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

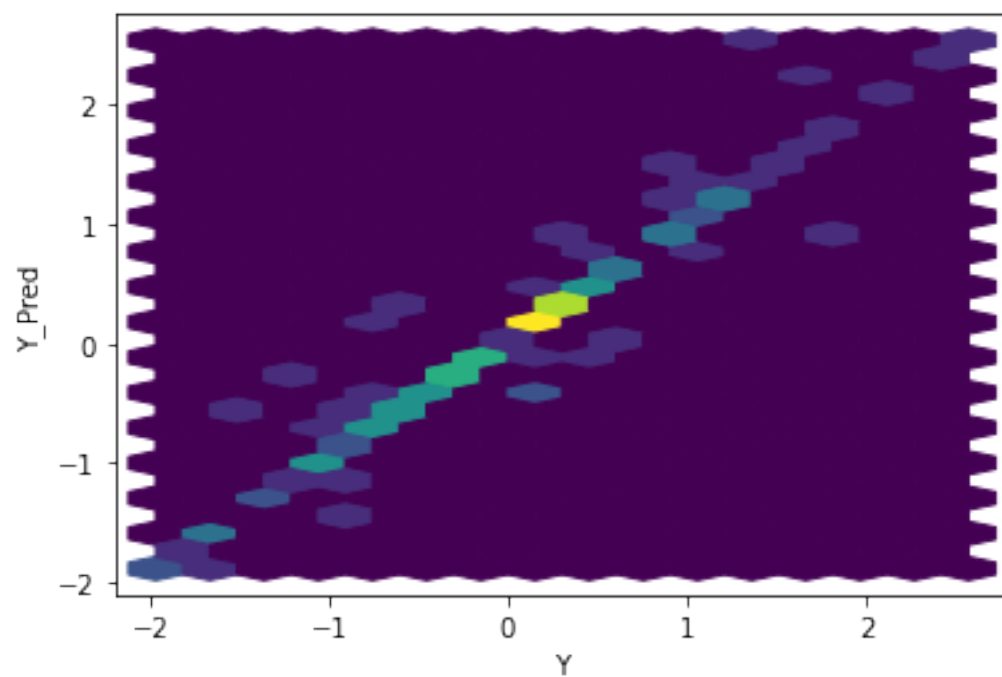
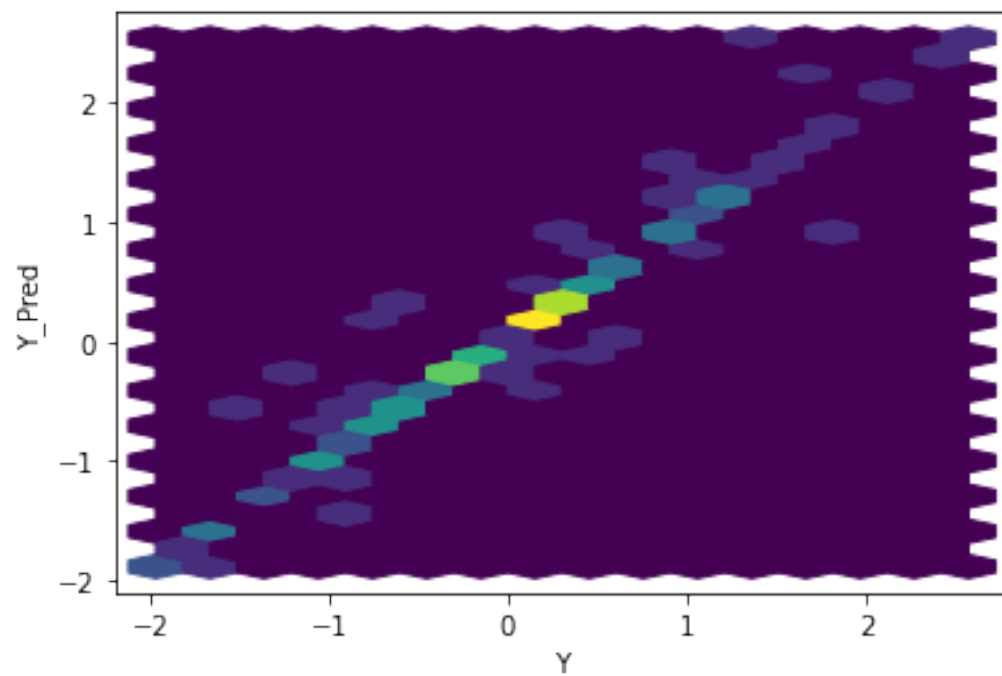
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))
```

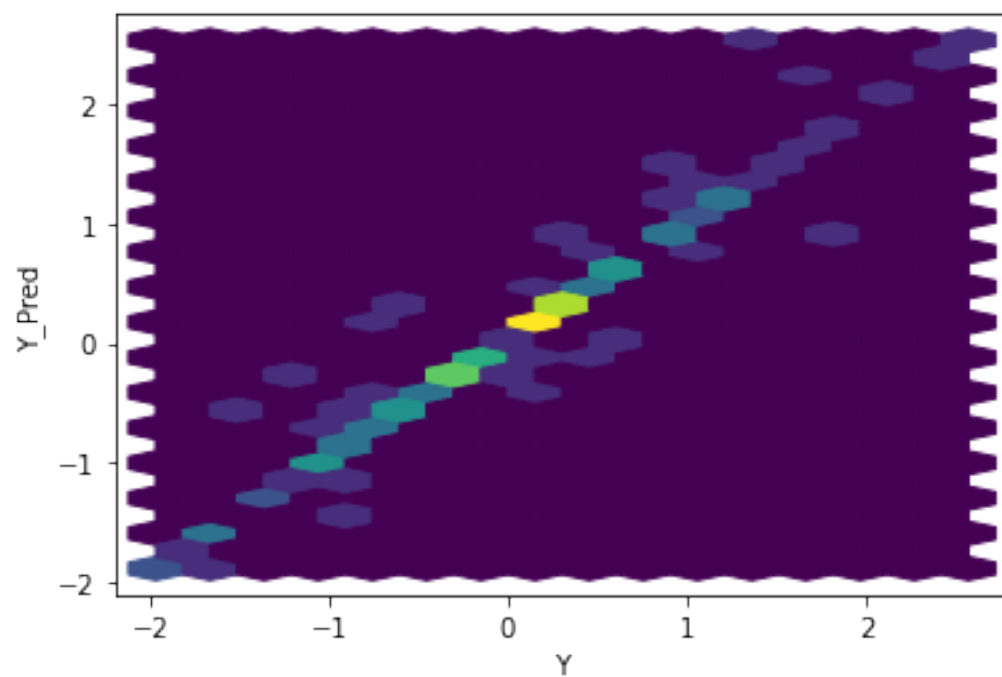
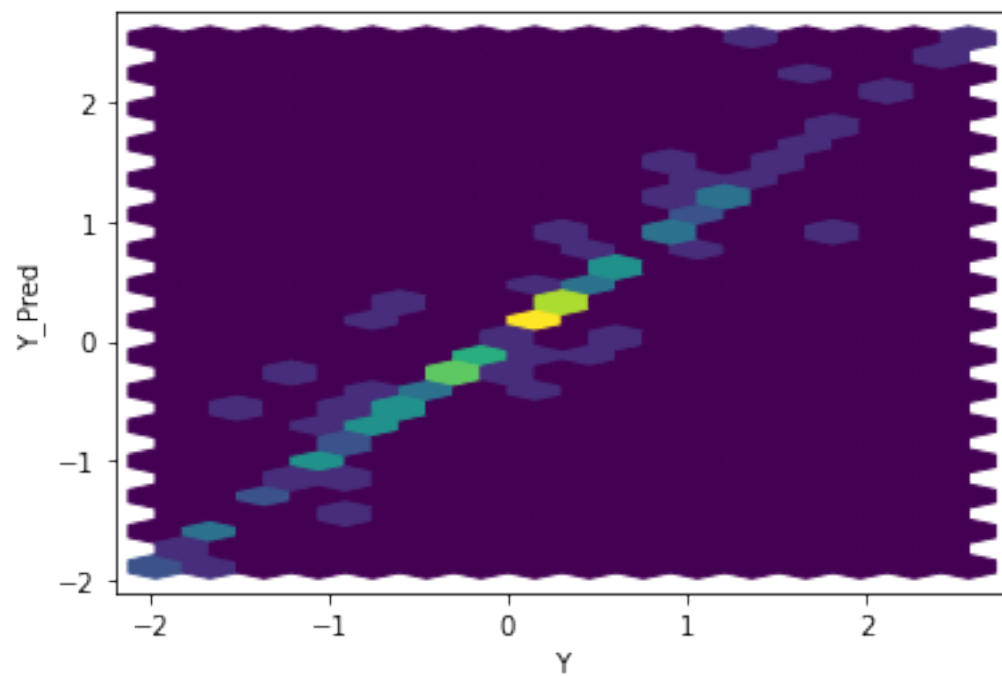
```
[22]: ABC_train_test.
      ↪ training_GAN_2(disc,gen,disc_opt,gen_opt,real_dataset,batch_size,
      ↪ error,criterion,coeff,mean,variance,device)
```

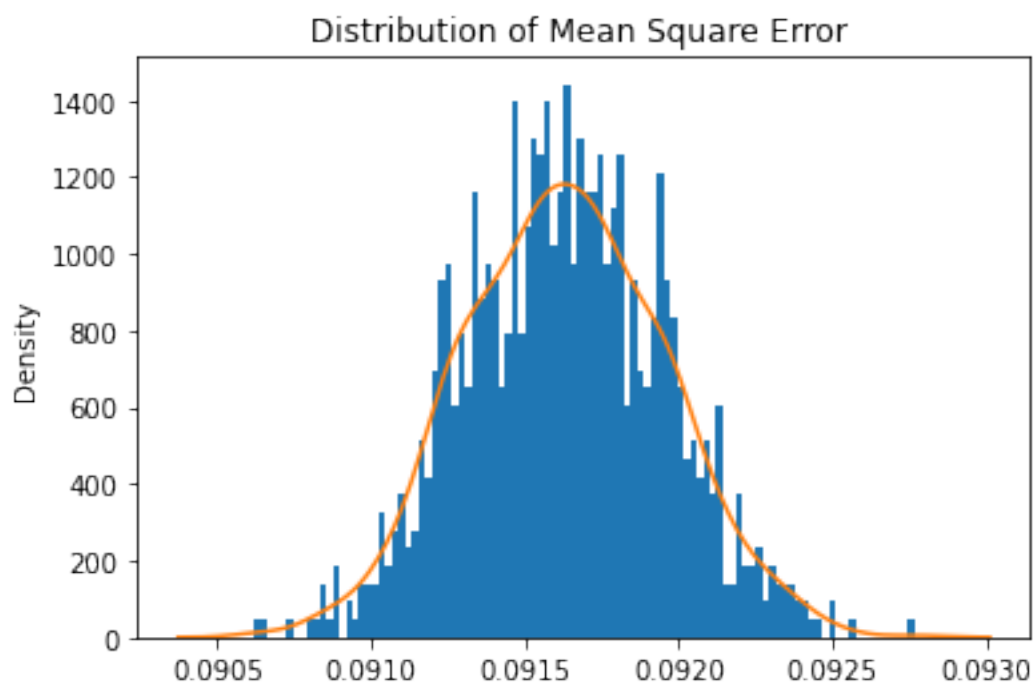
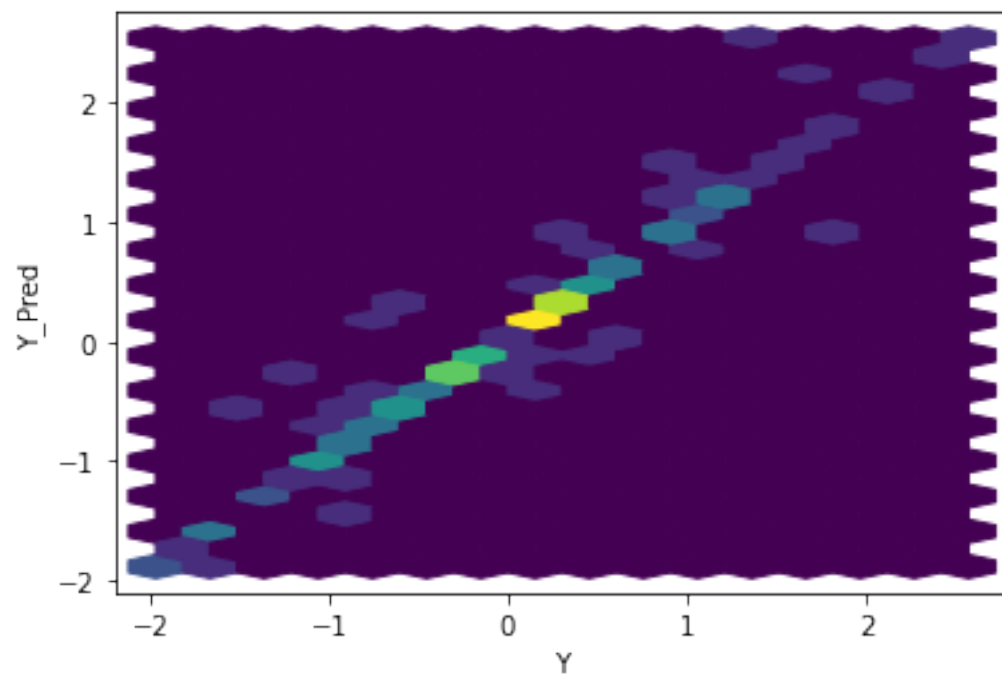
Number of epochs 30000



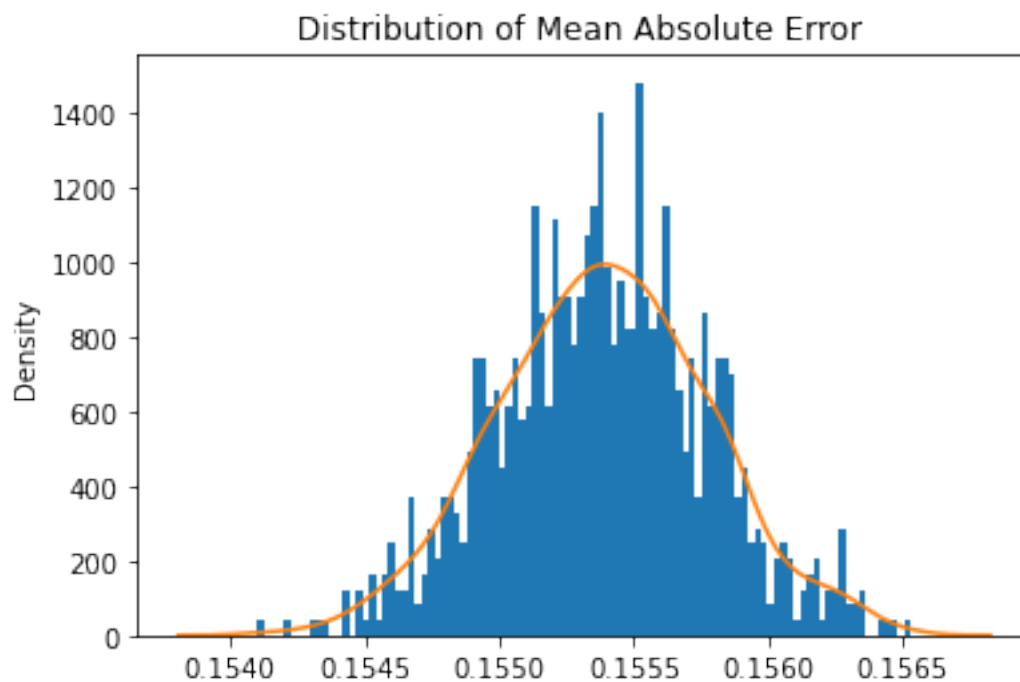
```
[23]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```



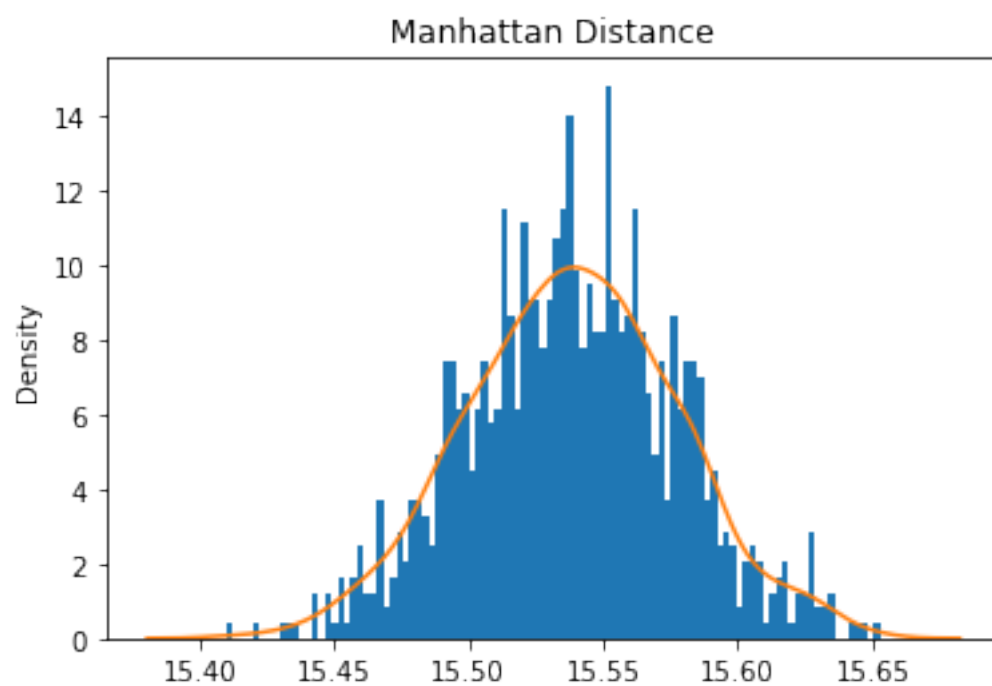




Mean Square Error: 0.09162953329869841

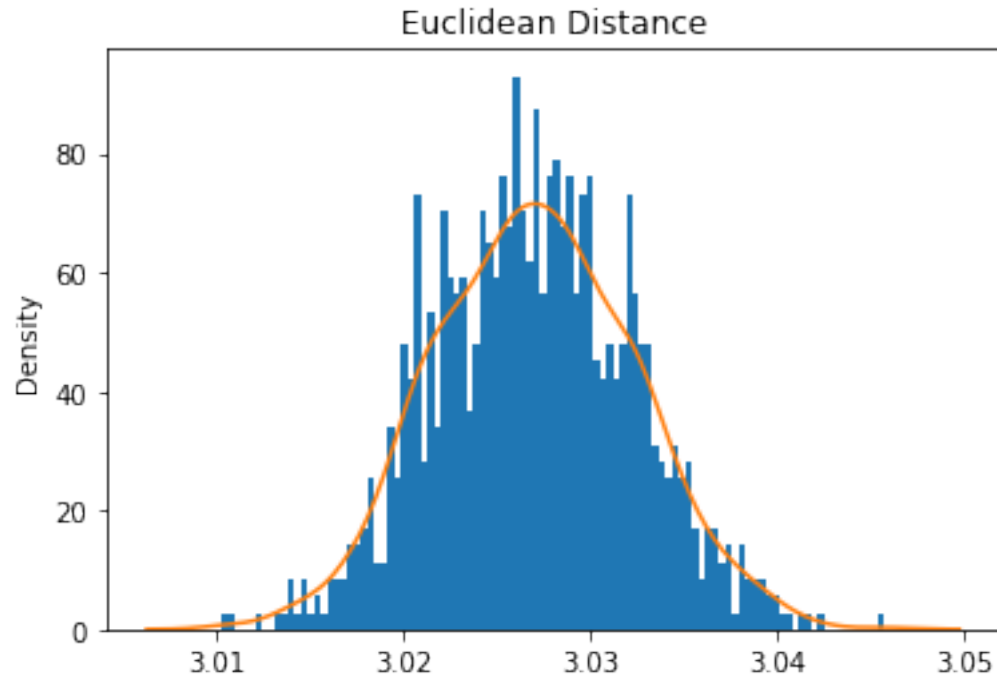


Mean Absolute Error: 0.15538604737251996  
Mean Manhattan Distance: 15.538604737251998





Mean Euclidean Distance: 3.0270323489319373



[ ]: