

Dataset2_Friedman1_output_4

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 = 10
     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	0.281529	0.018827	0.003404	0.679202	0.949131	0.042657	0.724302
1	0.868140	0.656731	0.758583	0.700712	0.638856	0.340029	0.003930
2	0.854184	0.091399	0.782765	0.807115	0.391684	0.409945	0.147166

```

3  0.020821  0.566786  0.217840  0.096170  0.436308  0.916902  0.537897
4  0.337245  0.314863  0.408296  0.605577  0.815081  0.543180  0.154017

```

```

          X7          X8          X9          Y
0  0.692897  0.029772  0.095120  16.652963
1  0.888997  0.659282  0.915307  21.184954
2  0.615591  0.661829  0.474850  14.090380
3  0.345637  0.745413  0.247845   5.235379
4  0.769471  0.870716  0.143051  13.697050

```

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:                Wed, 20 Oct 2021      Prob (F-statistic):         nan
Time:                19:59:25      Log-Likelihood:        330.50
No. Observations:         10      AIC:                -641.0
Df Residuals:             0      BIC:                -638.0
Df Model:                 9
Covariance Type:         nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	1.11e-16	inf	0	nan	nan	nan
x1	1.2577	inf	0	nan	nan	nan
x2	1.2081	inf	0	nan	nan	nan
x3	-1.6026	inf	-0	nan	nan	nan
x4	0.7478	inf	0	nan	nan	nan
x5	-0.6235	inf	-0	nan	nan	nan
x6	-1.1416	inf	-0	nan	nan	nan
x7	-0.8985	inf	-0	nan	nan	nan
x8	-0.5751	inf	-0	nan	nan	nan
x9	-0.1083	inf	-0	nan	nan	nan
x10	-0.7975	inf	-0	nan	nan	nan

```

=====
Omnibus:                1.702      Durbin-Watson:          2.650
Prob(Omnibus):           0.427      Jarque-Bera (JB):        1.186
Skew:                   -0.692      Prob(JB):                0.553
Kurtosis:                2.035      Cond. No.                33.5
=====

```

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 1.110223e-16

x1 1.257726e+00

x2 1.208062e+00

x3 -1.602593e+00

x4 7.477876e-01

x5 -6.235258e-01

x6 -1.141630e+00

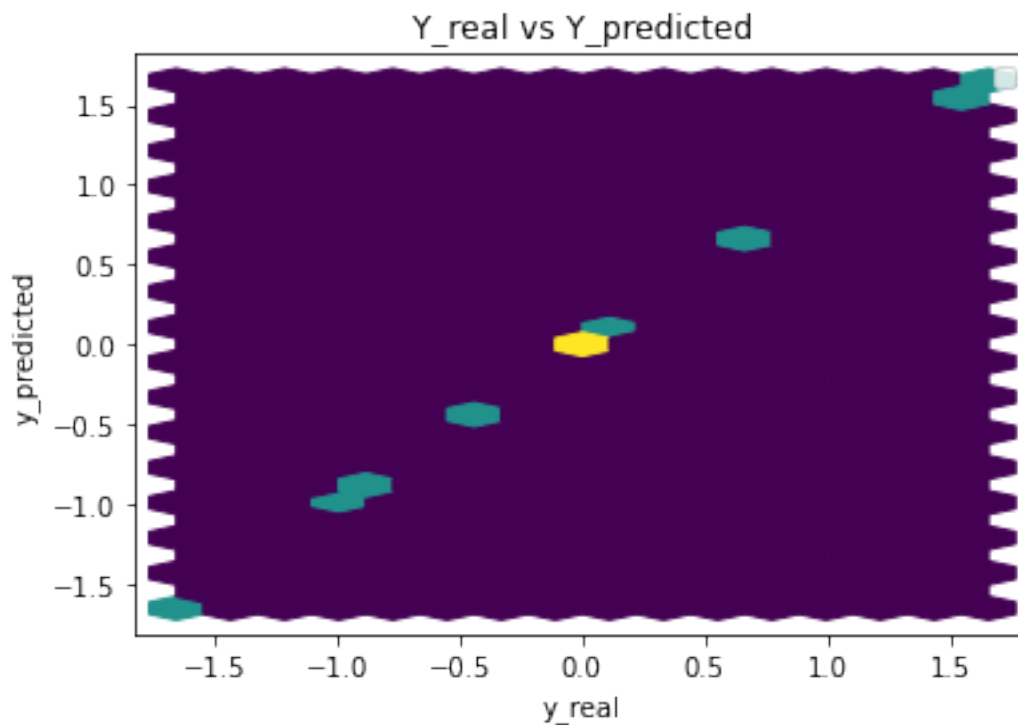
x7 -8.984535e-01

x8 -5.751330e-01

x9 -1.083431e-01

x10 -7.974747e-01

dtype: float64



Performance Metrics

Mean Squared Error: 1.148822026651847e-30

Mean Absolute Error: 8.23646706393788e-16

Manhattan distance: 8.23646706393788e-15

Euclidean distance: 3.3894277196185305e-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

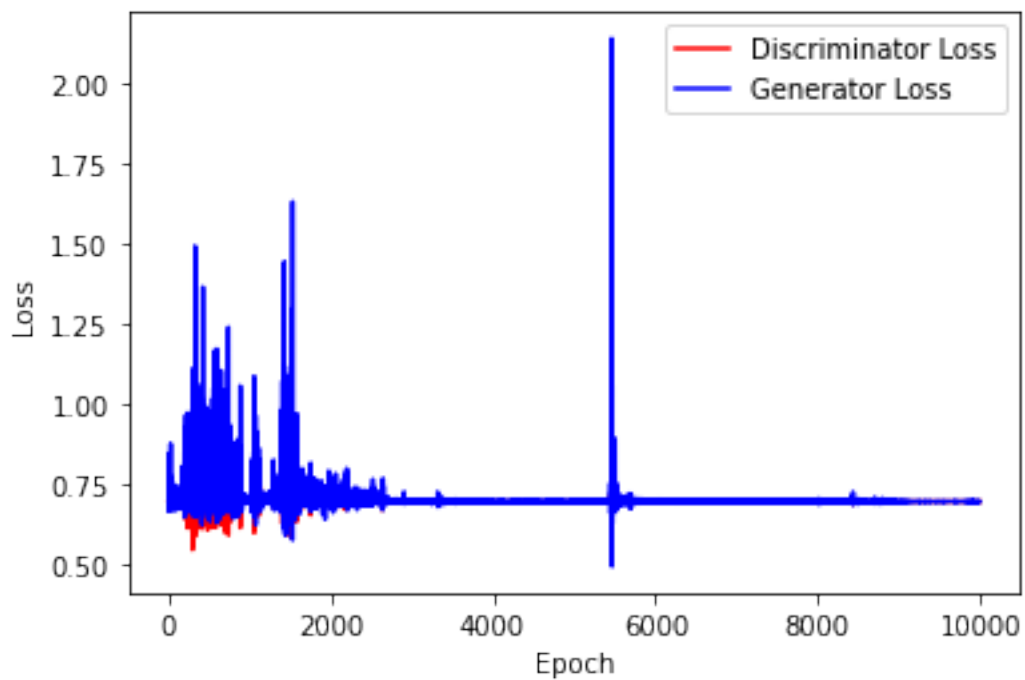
```
[9]: generator = network.Generator(n_features+2)
     discriminator = network.Discriminator(n_features+2)

     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))
```

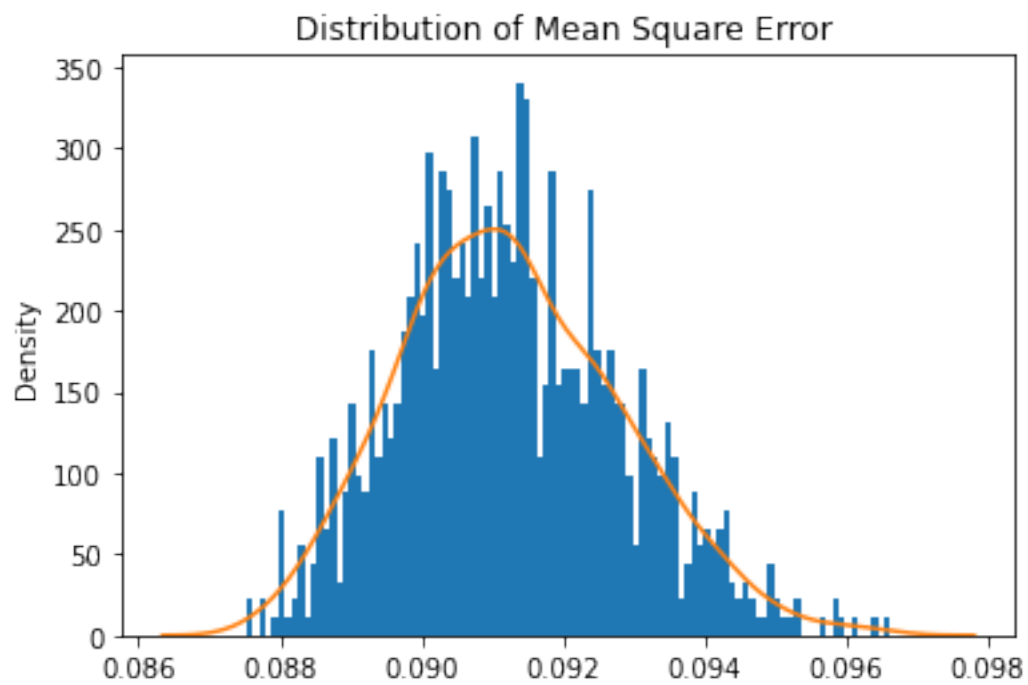
```
[10]: print(generator)
      print(discriminator)
```

```
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()
)
```

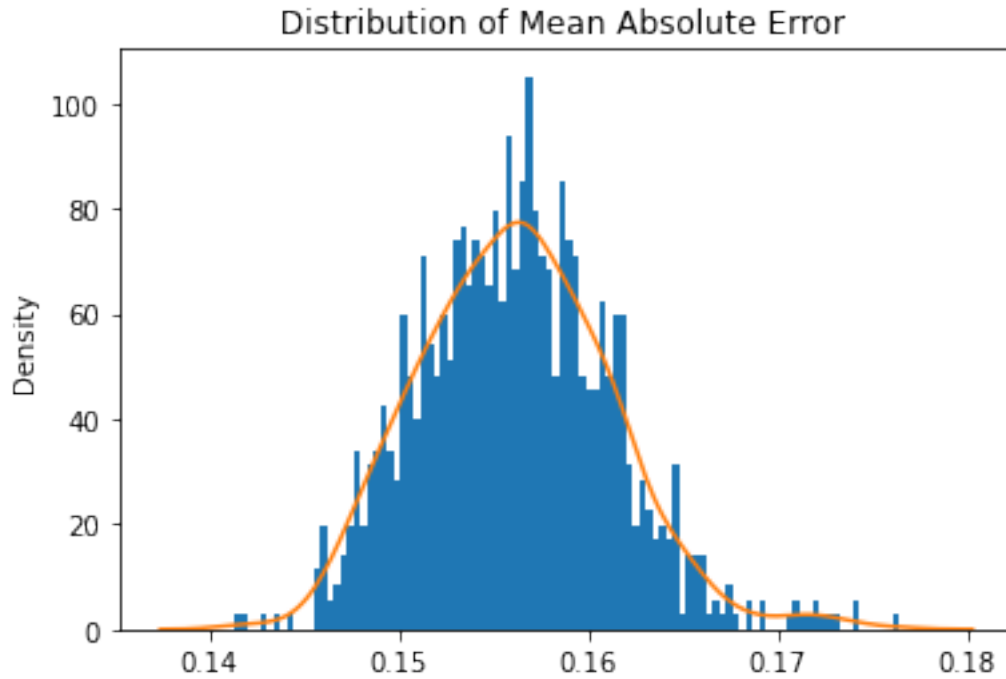
```
[11]: train_test.
     ↪training_GAN(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,
     ↪n_epochs,criterion,device)
```



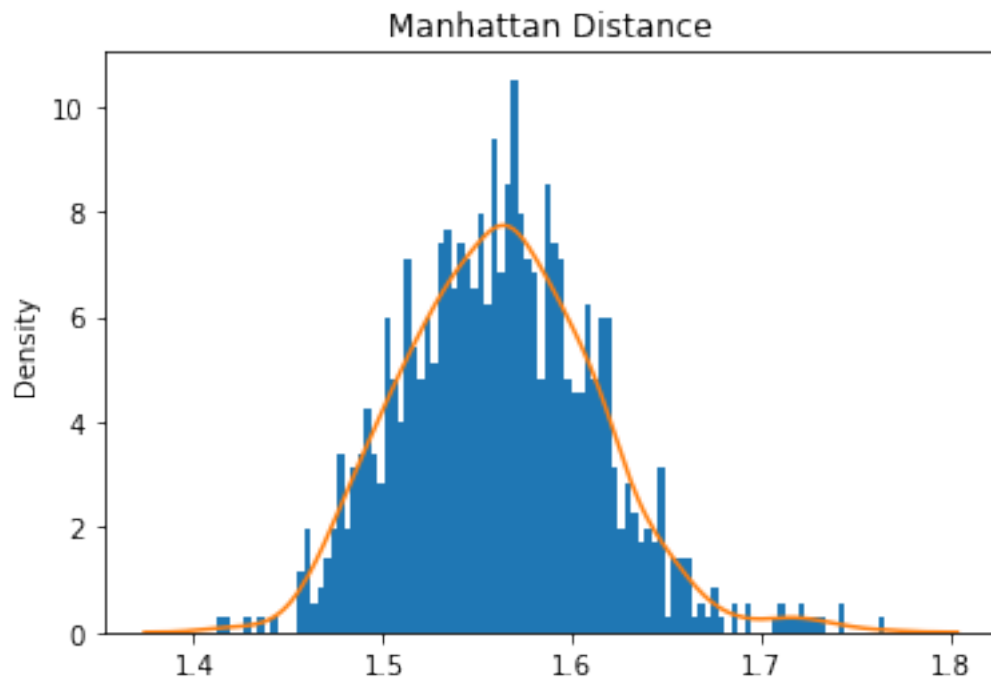
```
[12]: train_test.test_generator(generator,real_dataset,device)
```



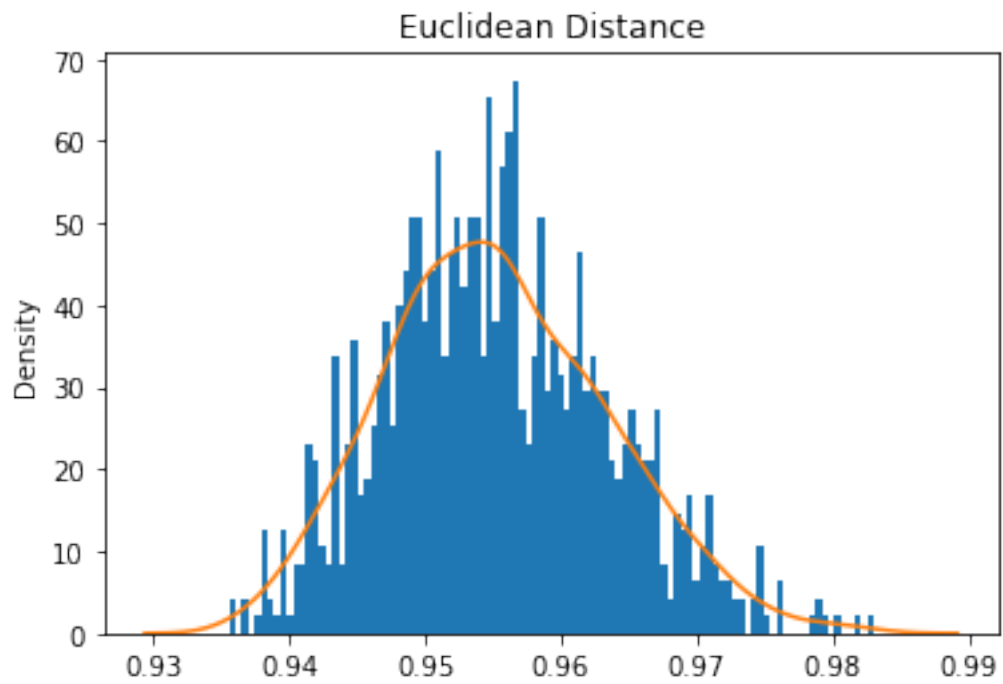
Mean Square Error: 0.09125961415312712



Mean Absolute Error: 0.15615501274988056

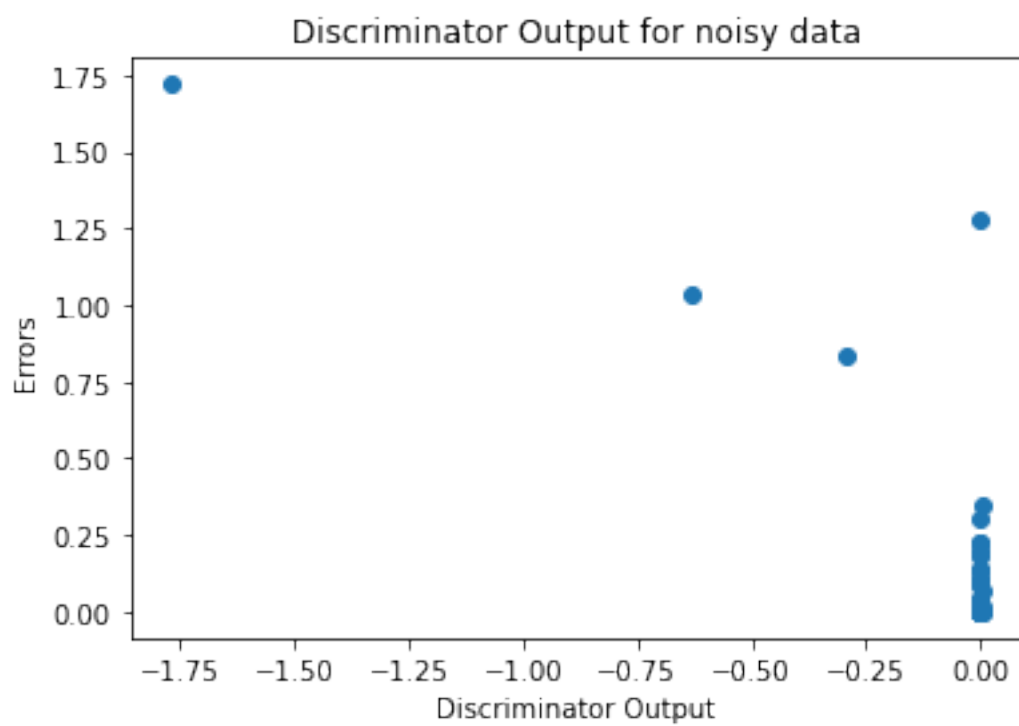
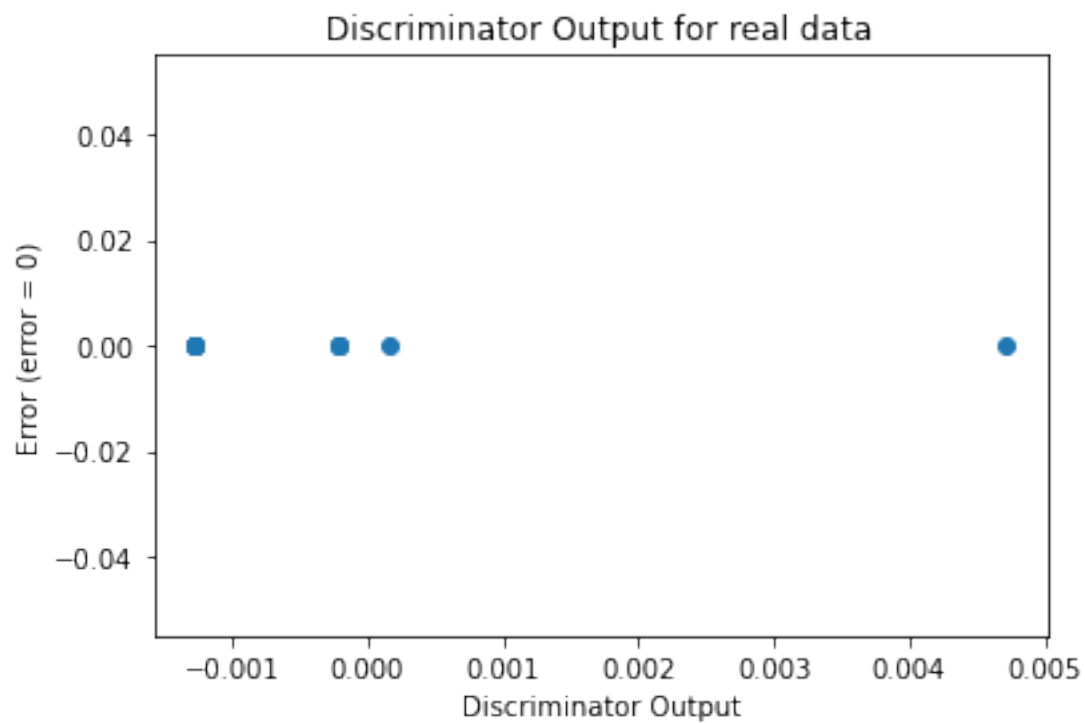


Mean Manhattan Distance: 1.5615501274988055



Mean Euclidean Distance: 0.9552637362202471

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

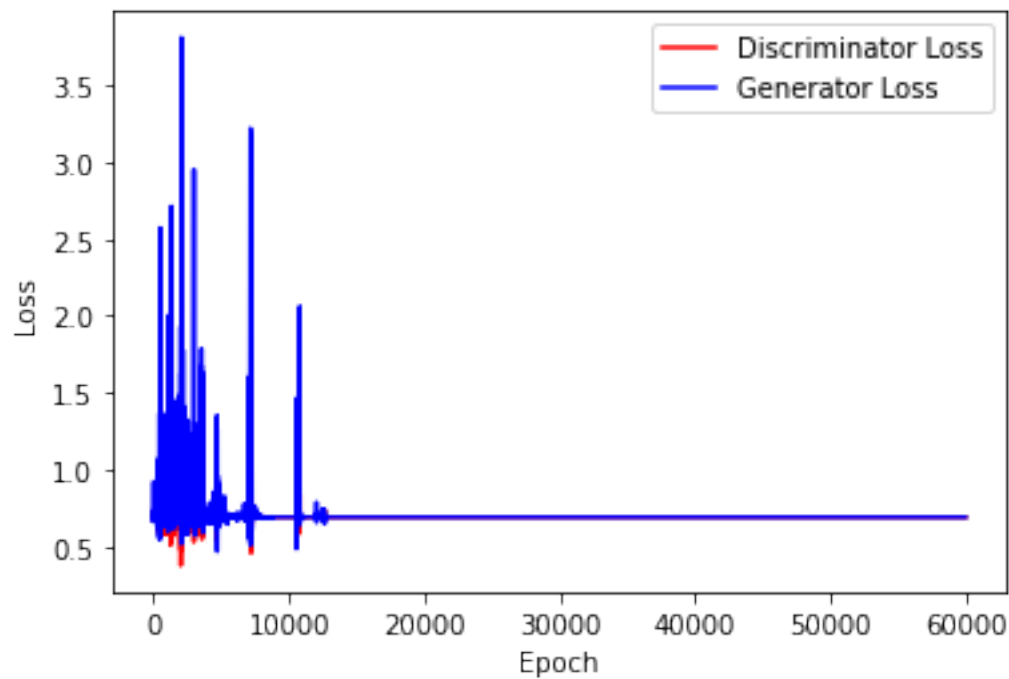



Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

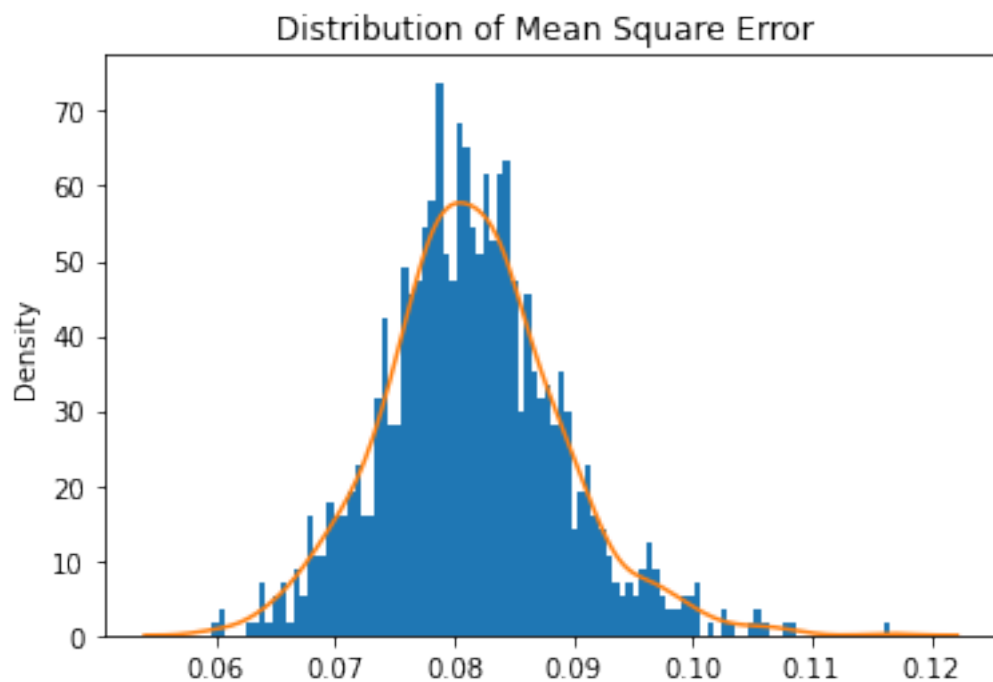
```
[14]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator.parameters(), lr=0.01, betas=(0.5,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

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

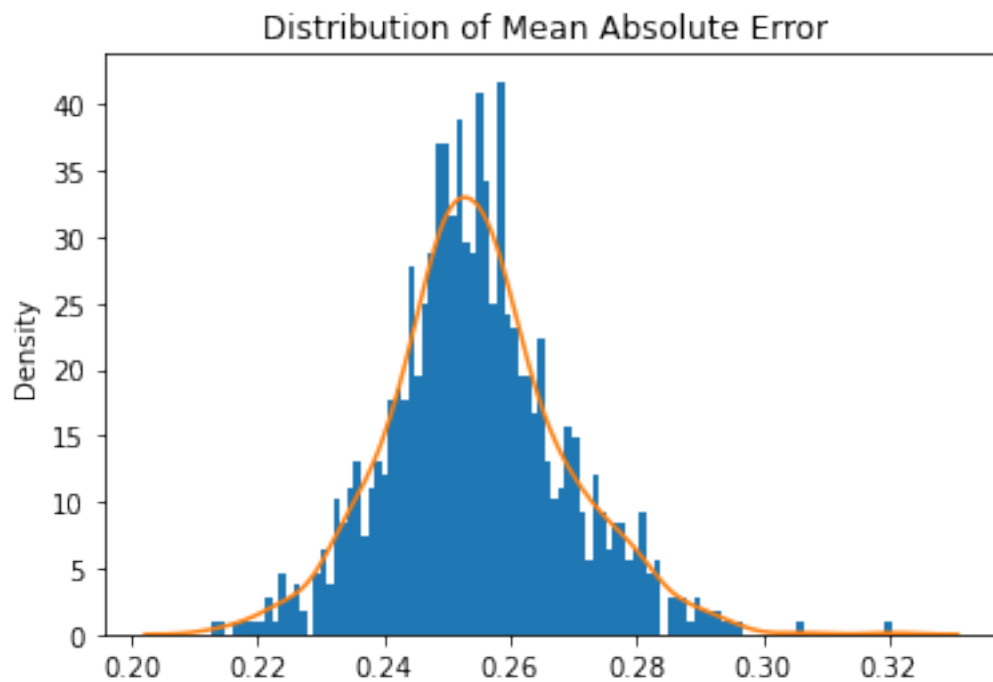
Number of epochs needed 30000



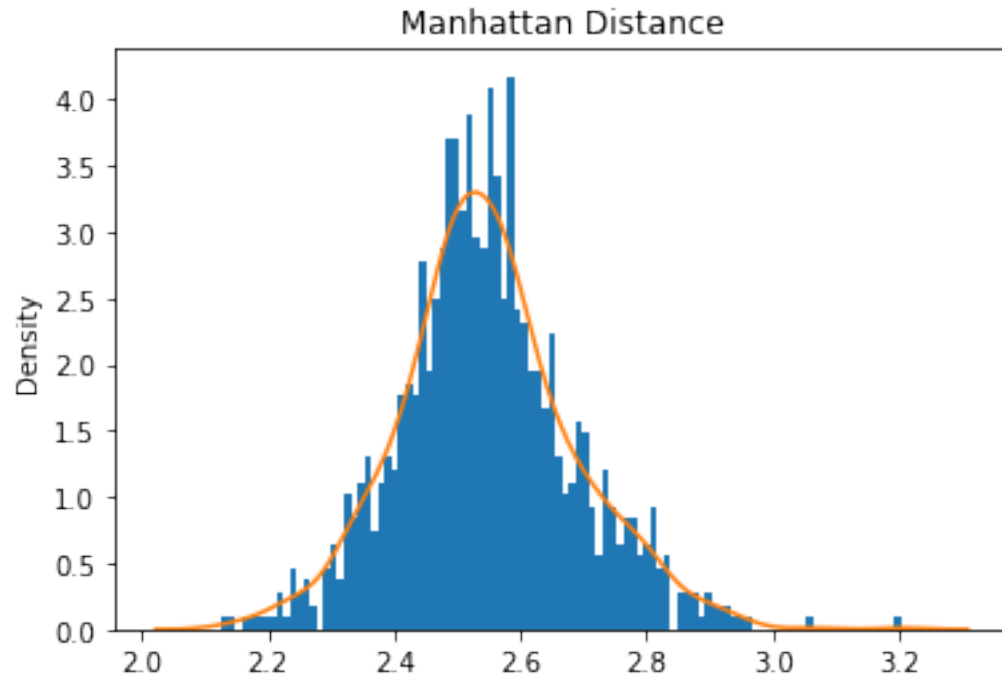
```
[16]: train_test.test_generator(generator,real_dataset,device)
```



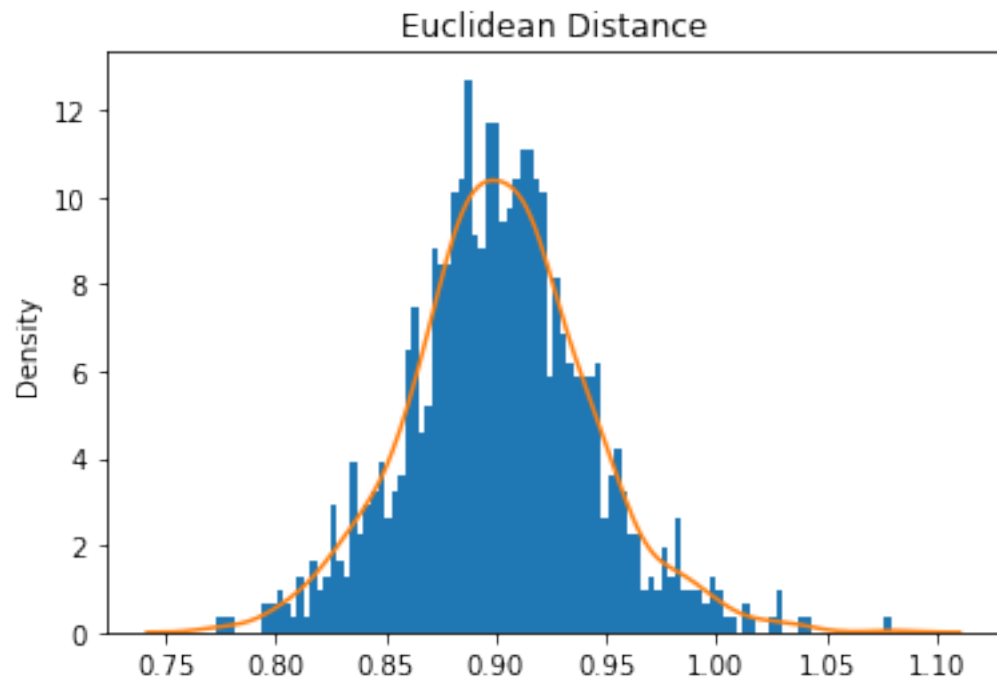
Mean Square Error: 0.0814982731036312



Mean Absolute Error: 0.254736946503818



Mean Manhattan Distance: 2.54736946503818



Mean Euclidean Distance: 0.9018278394981034

2 ABC GAN Model

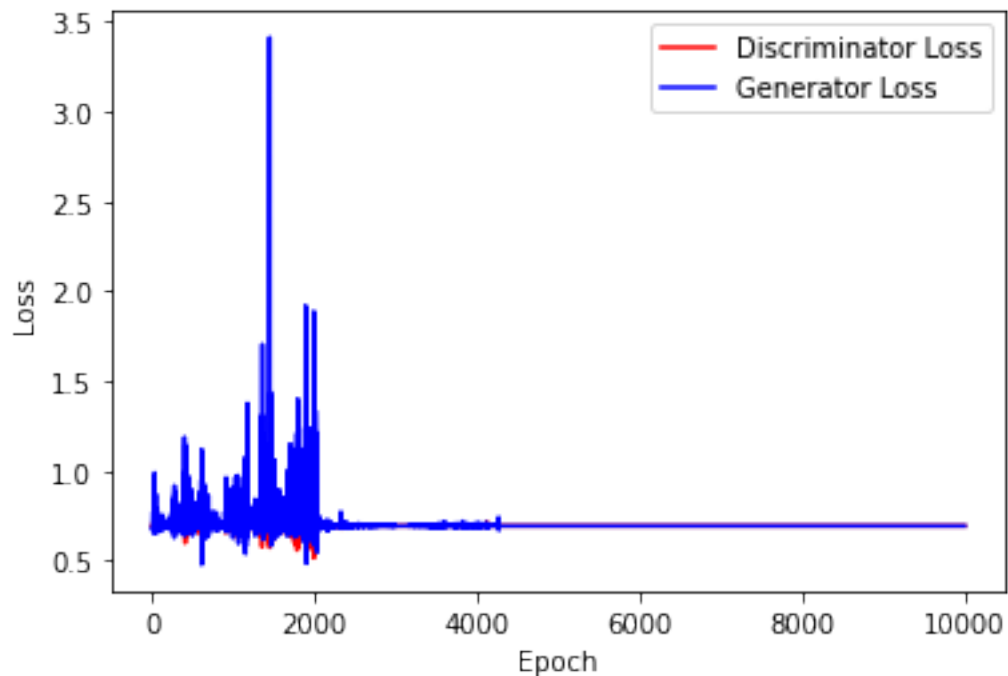
2.0.1 Training the network

Training ABC-GAN for `n_epochs` number of epochs

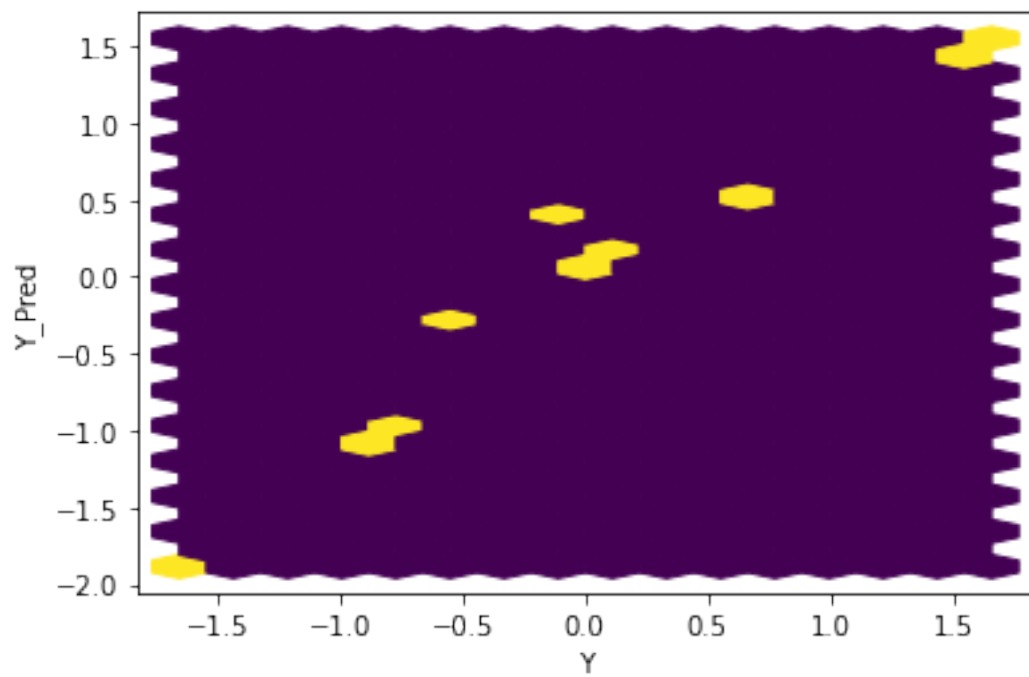
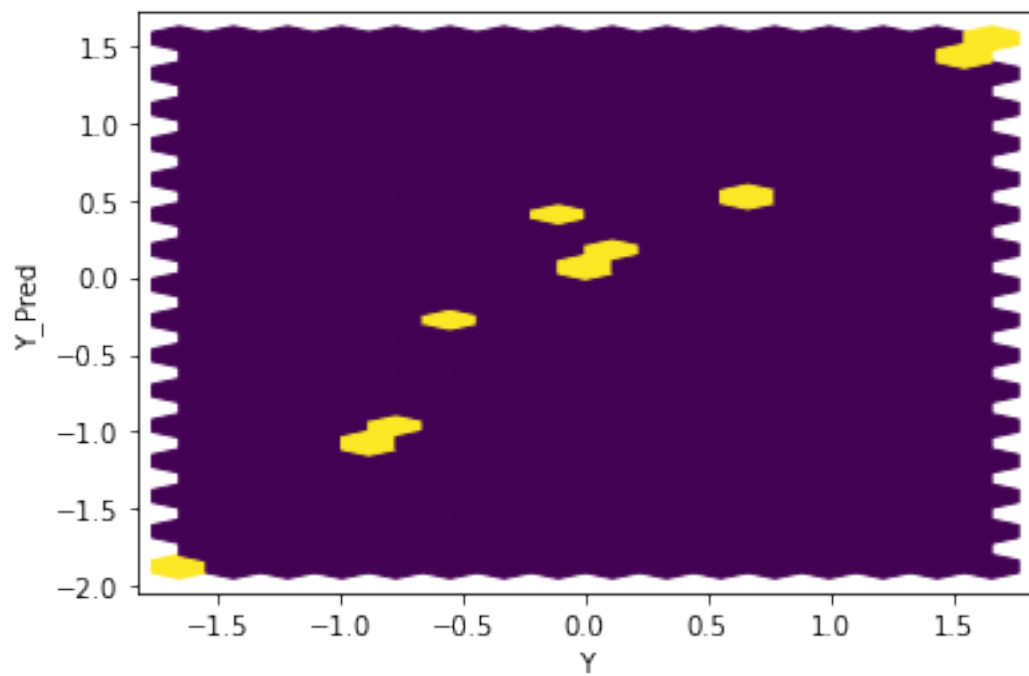
```
[17]: gen = network.Generator(n_features+2)
      disc = network.Discriminator(n_features+2)

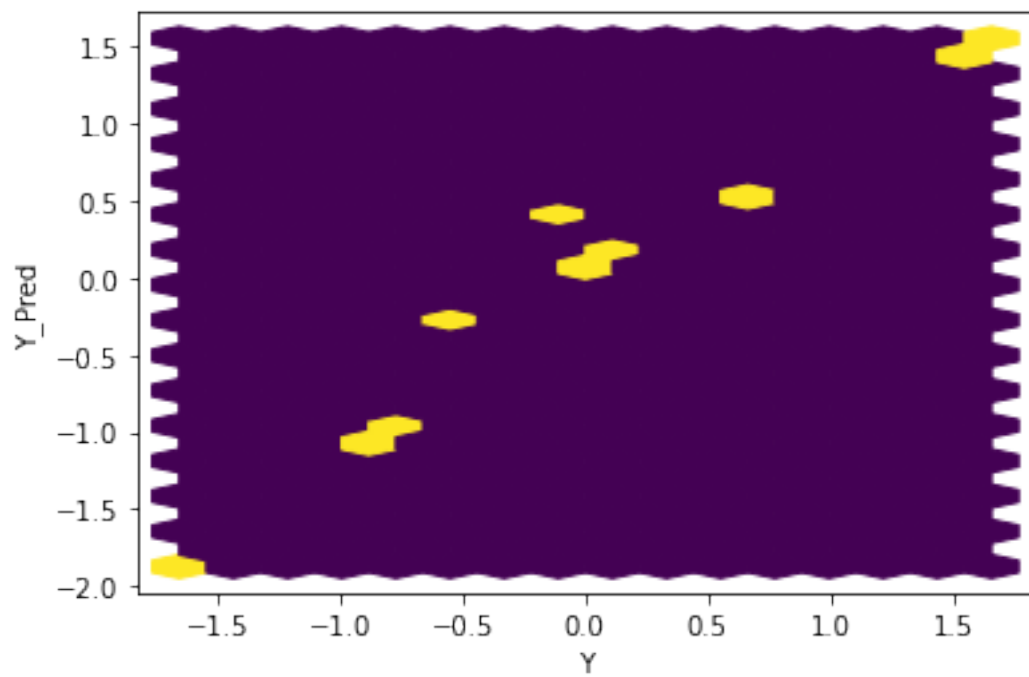
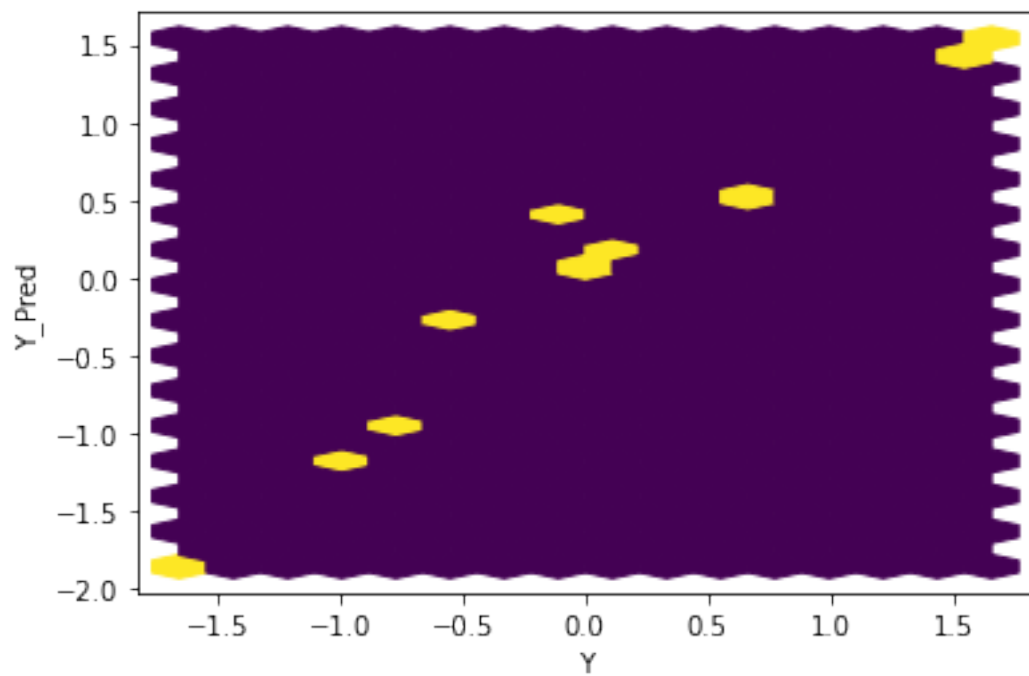
      criterion = torch.nn.BCEWithLogitsLoss()
      gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
      disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))

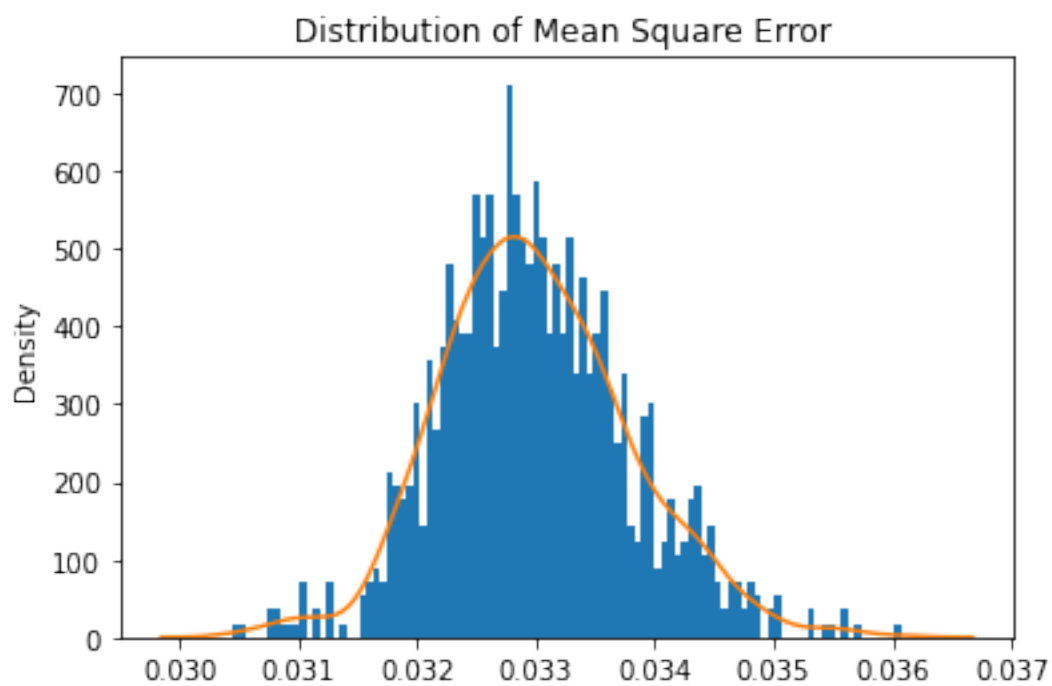
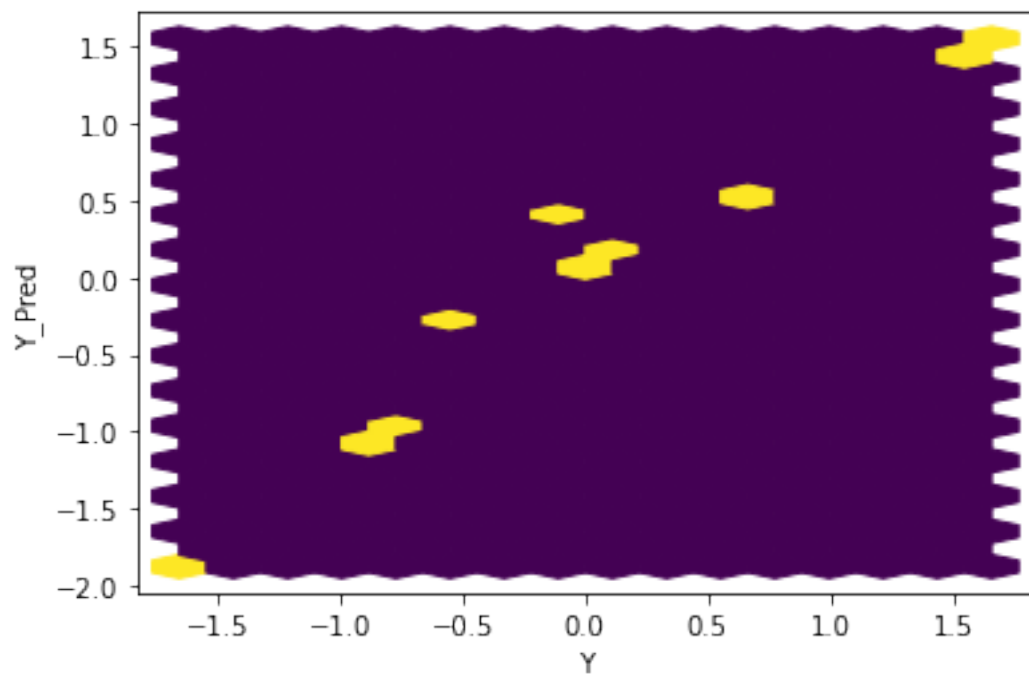
[18]: ABC_train_test.training_GAN(disc, gen, disc_opt, gen_opt, real_dataset,
      ↪ batch_size, n_epochs, criterion, coeff, mean, variance, device)
```



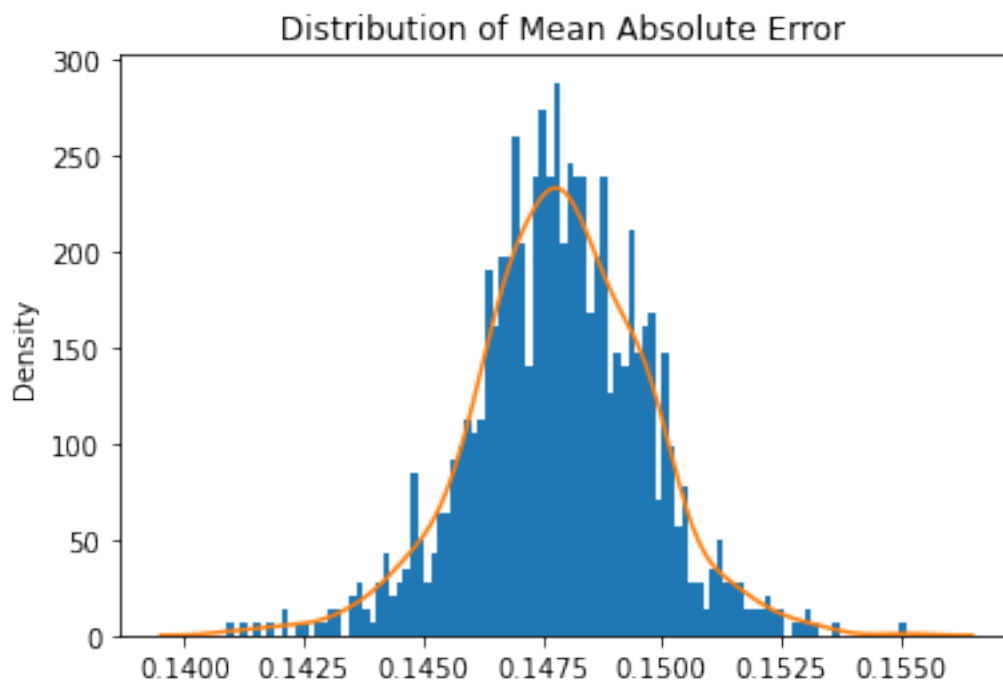
```
[19]: ABC_train_test.test_generator(gen, real_dataset, coeff, mean, variance, device)
```





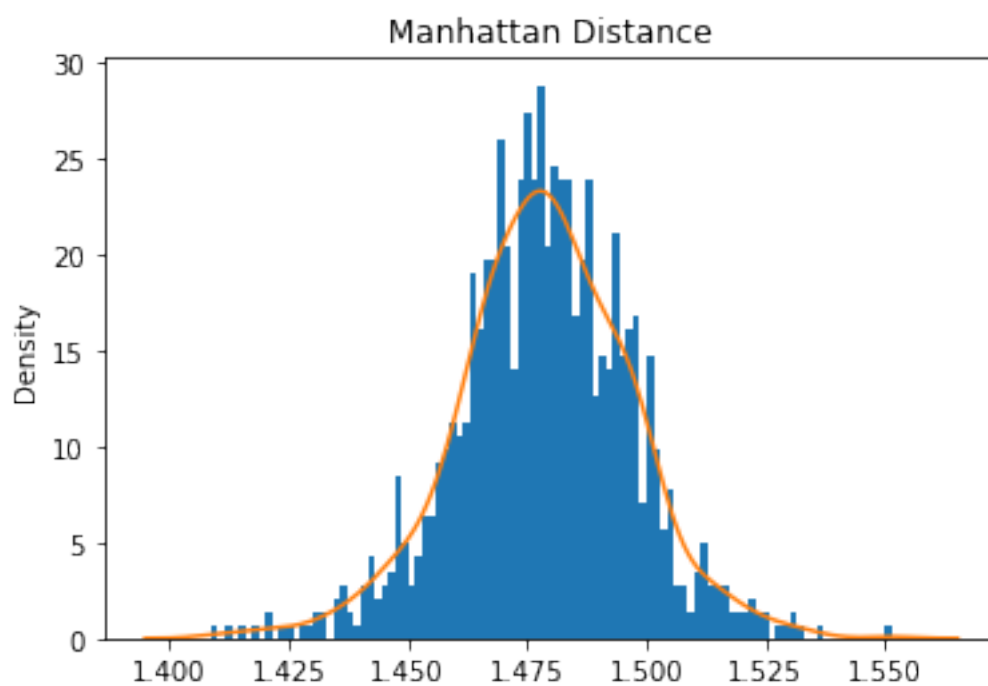


Mean Square Error: 0.032991318890674144

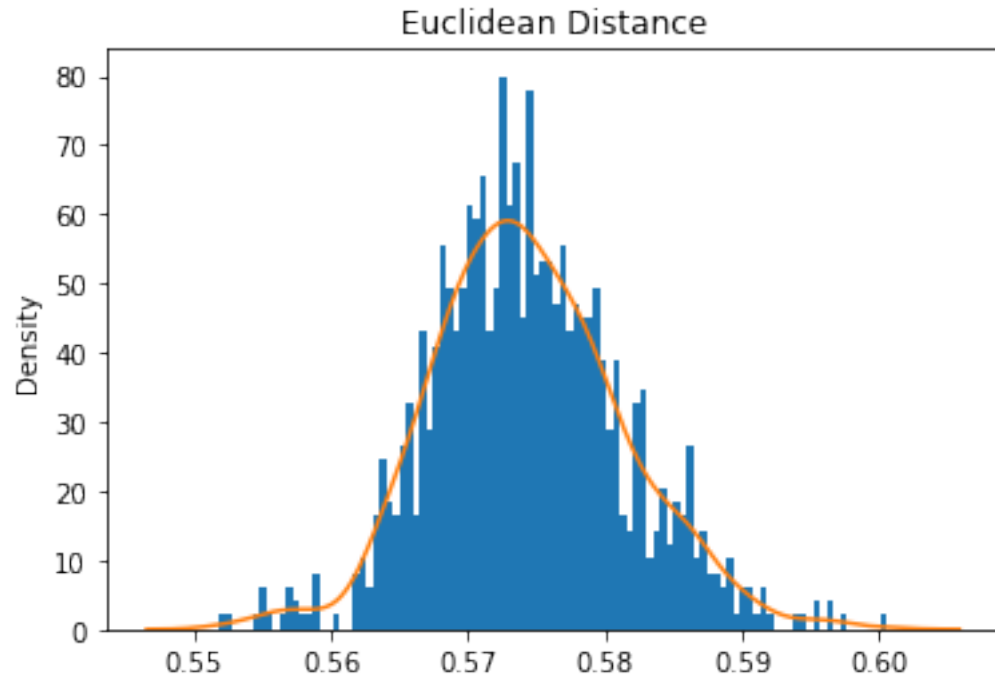


Mean Absolute Error: 0.14785735341757536

Mean Manhattan Distance: 1.4785735341757535

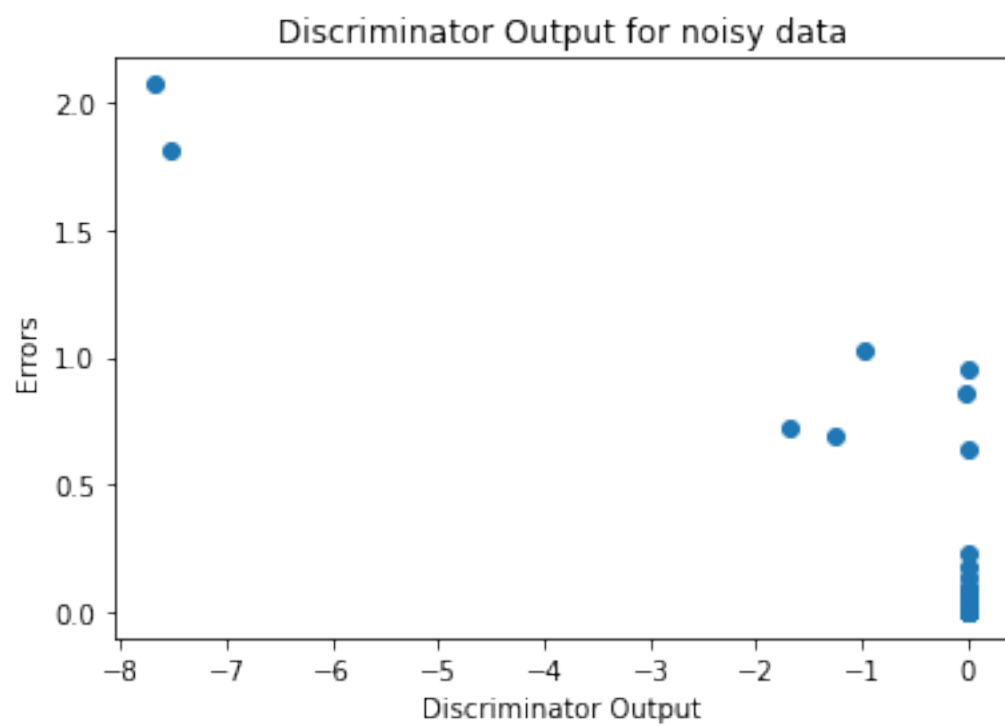
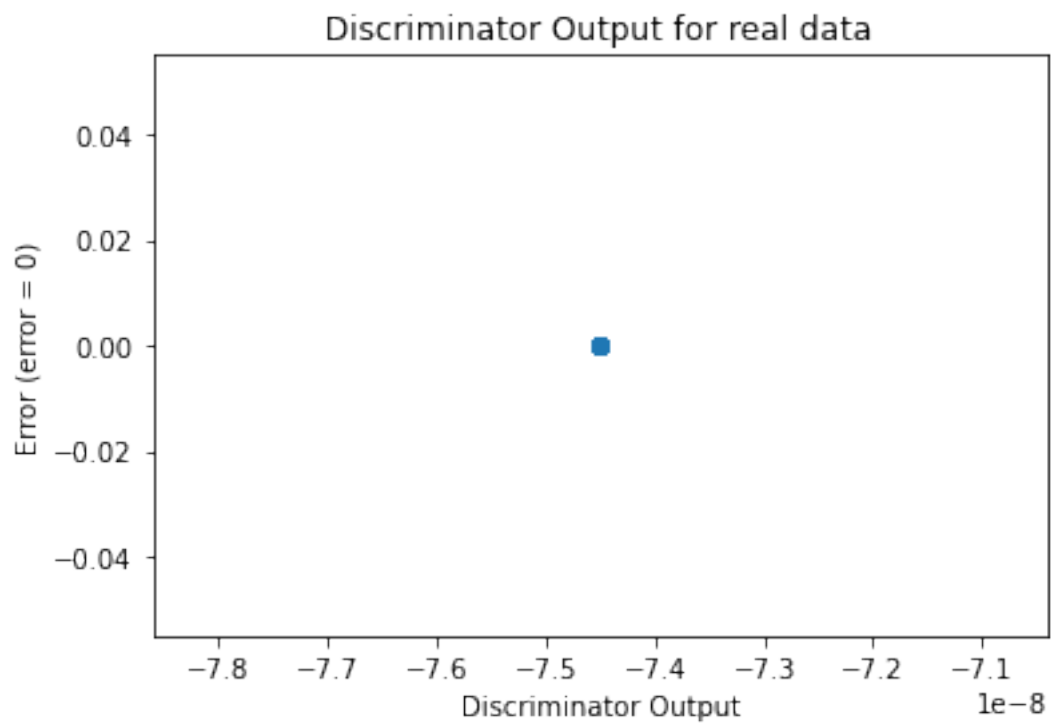


Mean Euclidean Distance: 0.5743385569829554



Sanity Checks

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



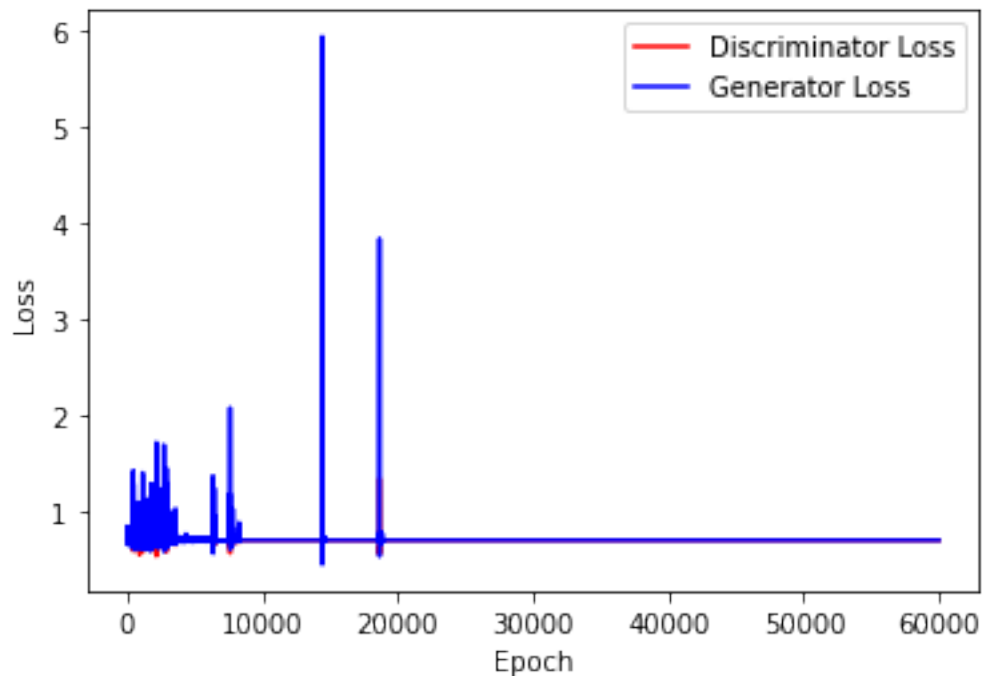
Training GAN until mse of y_{pred} is > 0.1 or $n_{\text{epochs}} < 30000$

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

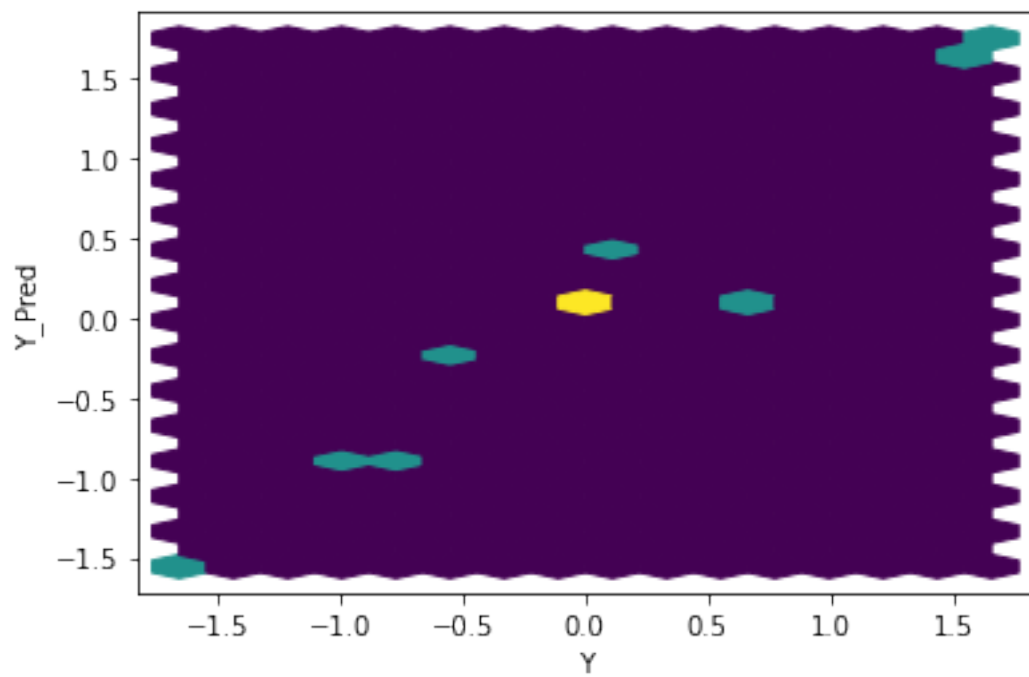
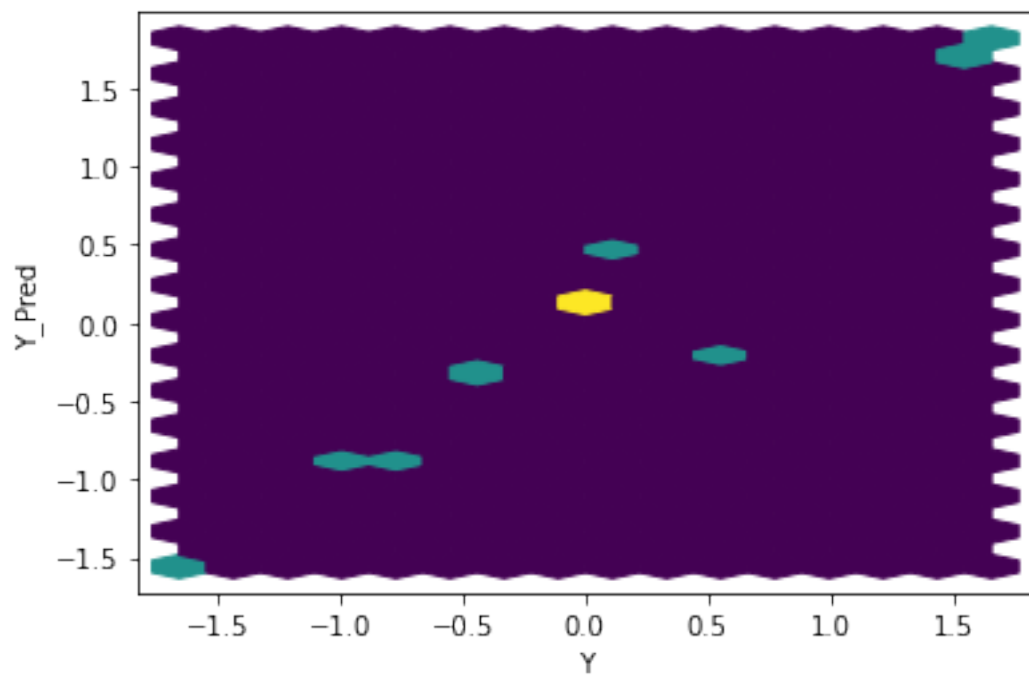
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))
```

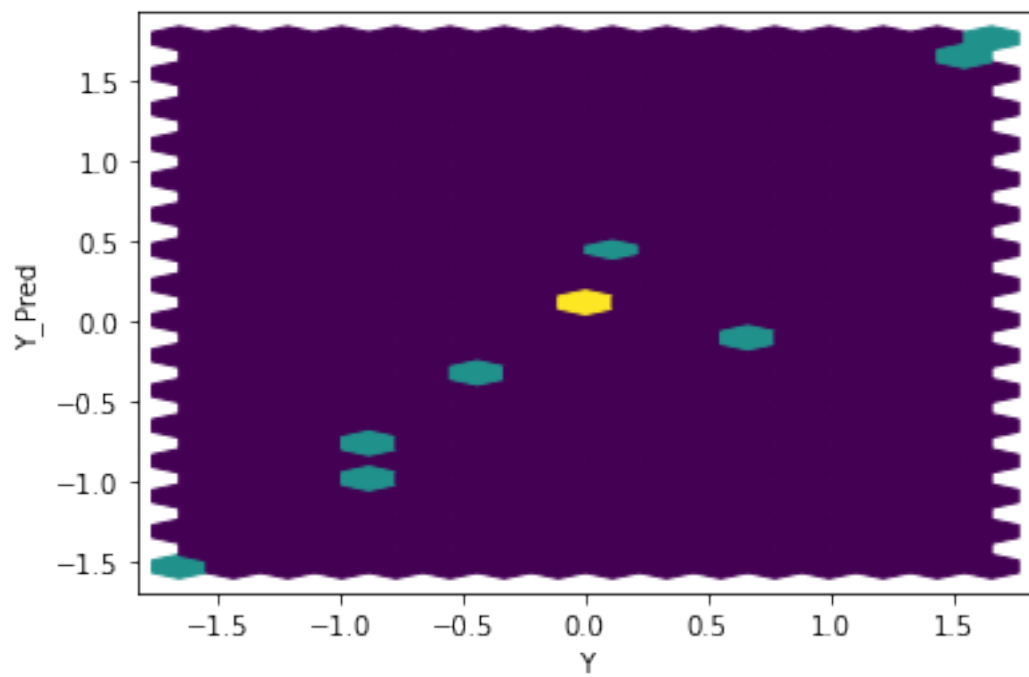
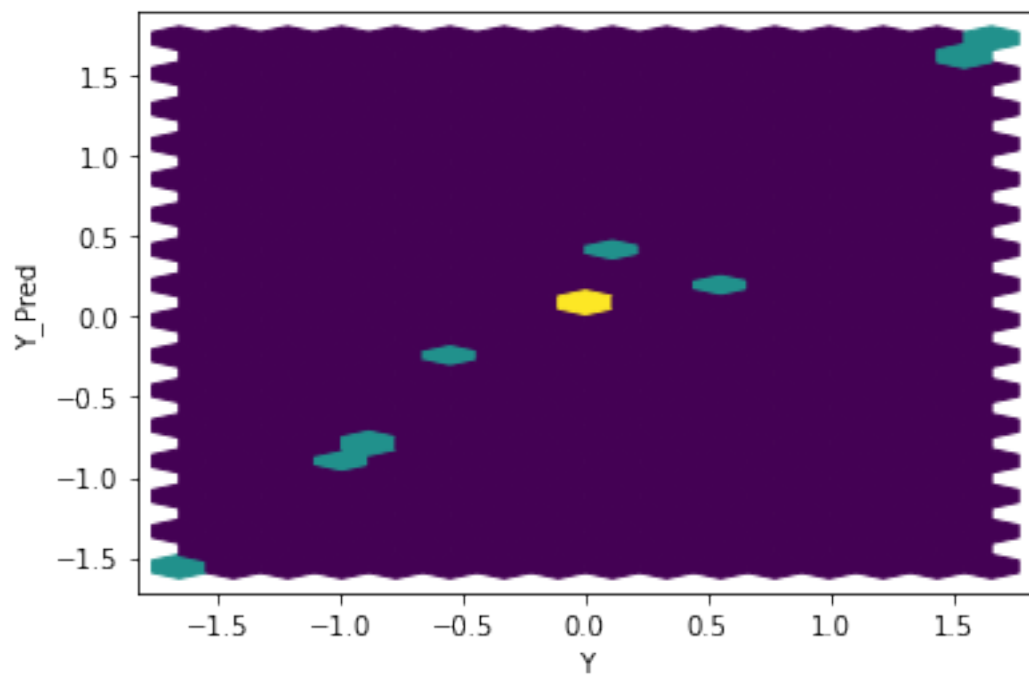
```
[22]: ABC_train_test.
      ↪ training_GAN_2(disc,gen,disc_opt,gen_opt,real_dataset,batch_size,
      ↪ error,criterion,coeff,mean,variance,device)
```

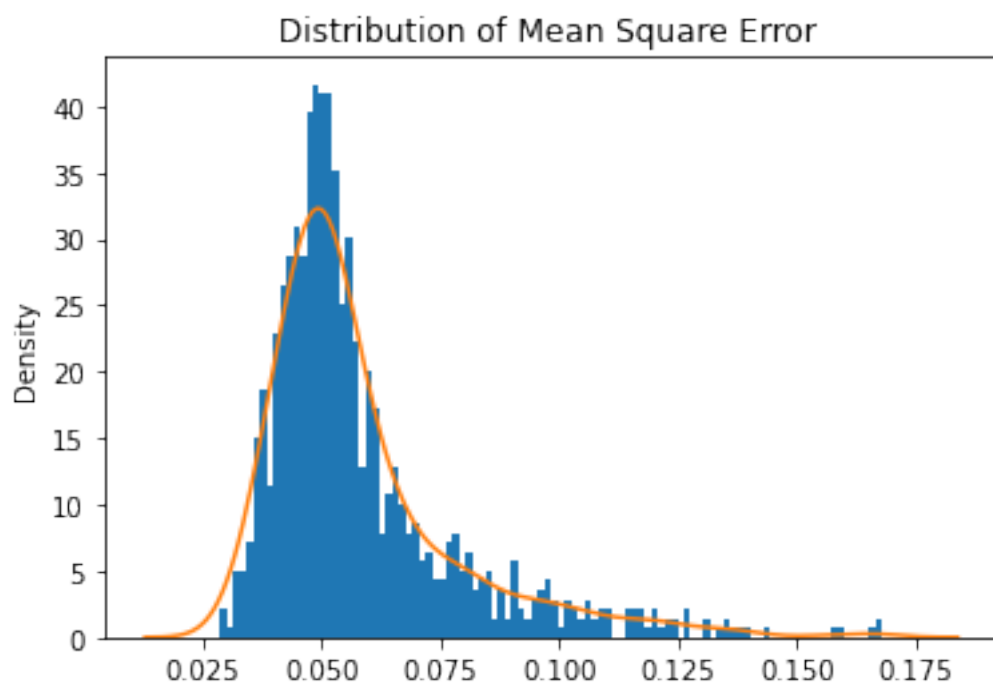
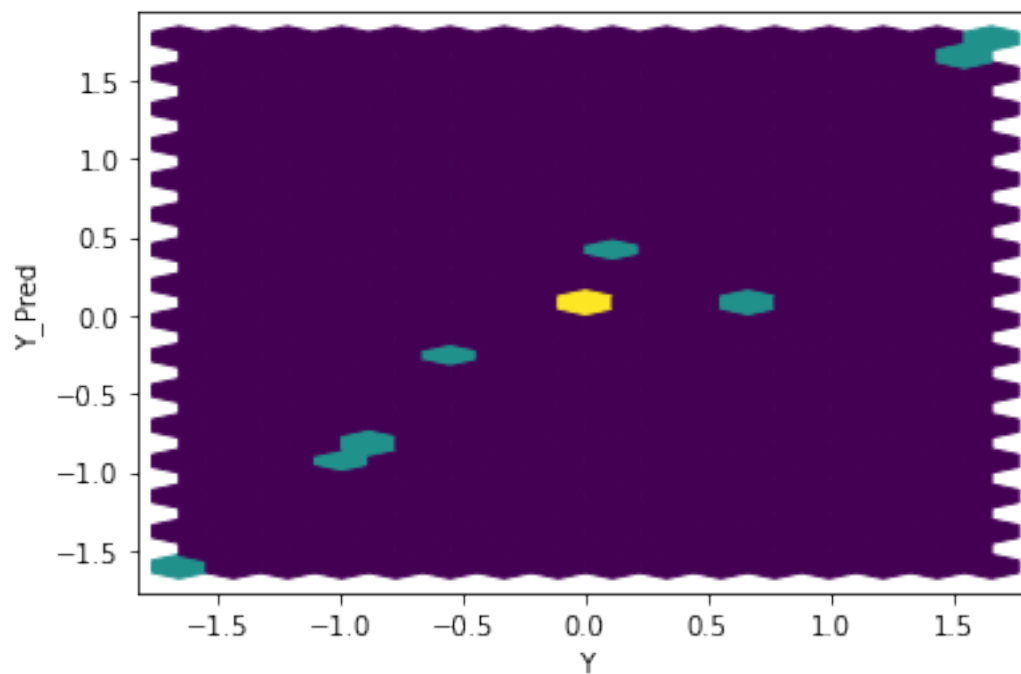
Number of epochs 30000



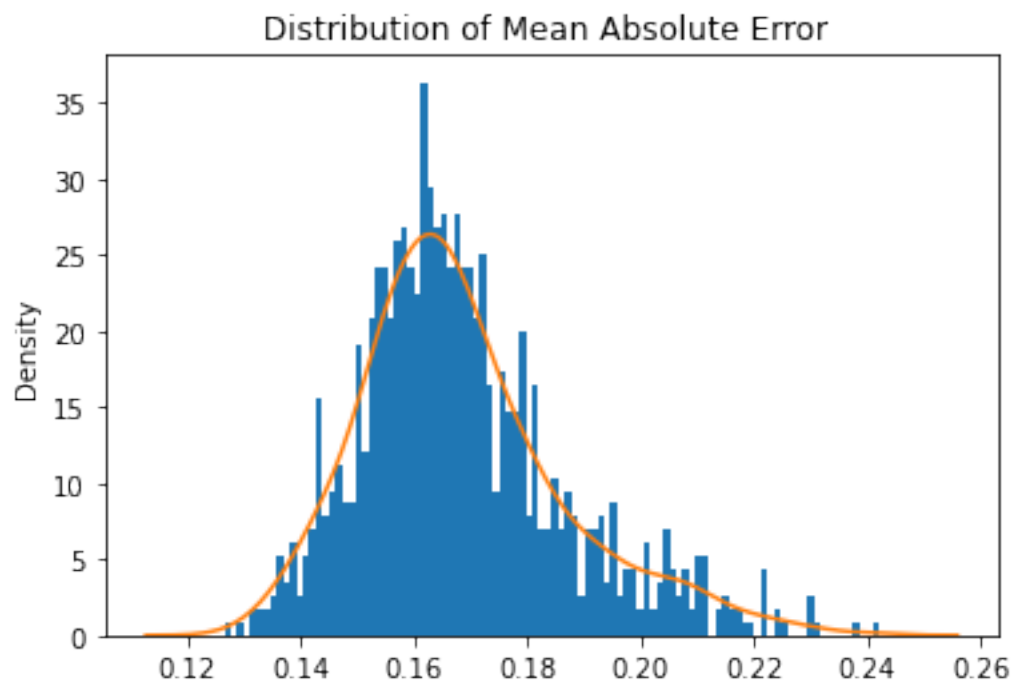
```
[23]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```



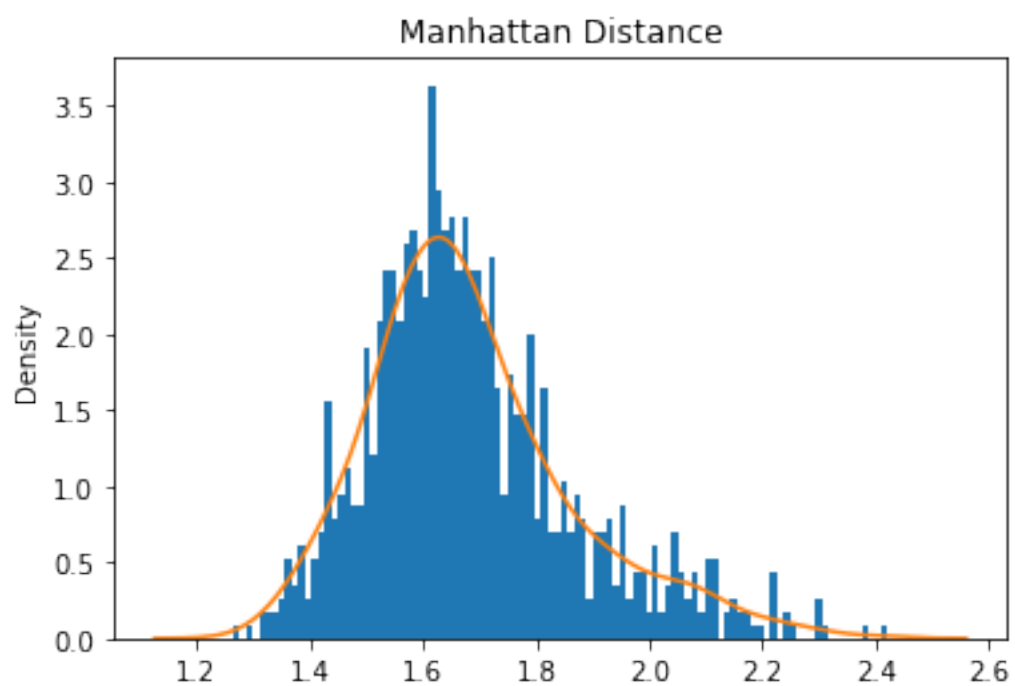




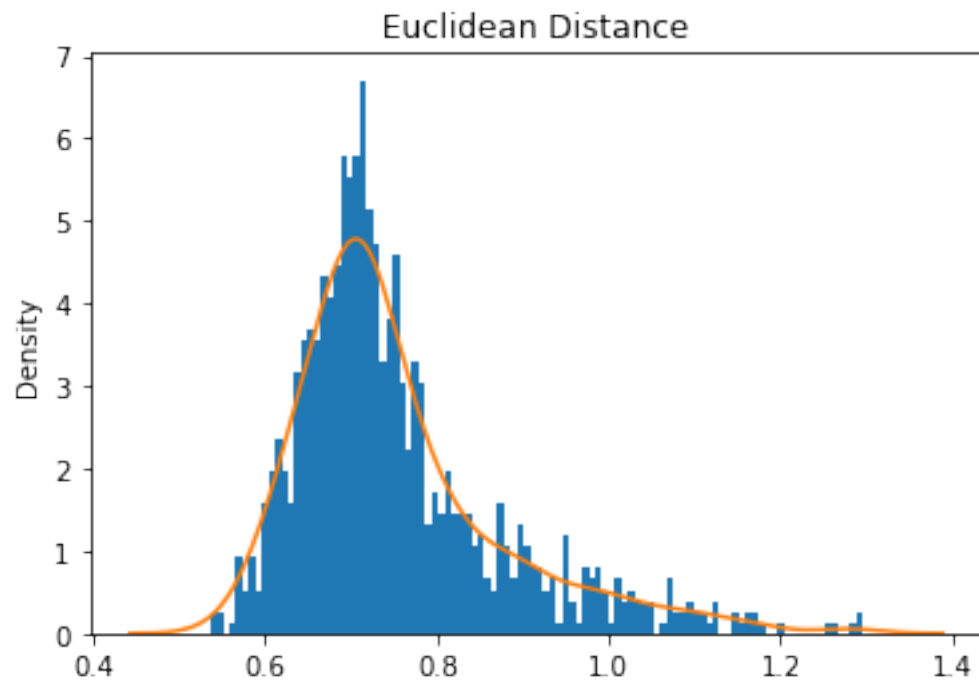
Mean Square Error: 0.058318151012330434



Mean Absolute Error: 0.16839529609903695
Mean Manhattan Distance: 1.6839529609903694



Mean Euclidean Distance: 0.7536669691987697



[]: