

Dataset2_Friedman1_output_1

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 = 1
     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	0.245489	0.325043	0.294875	0.836413	0.769018	0.069491	0.895824
1	0.752714	0.277079	0.370823	0.318056	0.935359	0.167154	0.356545
2	0.872698	0.703273	0.051094	0.068922	0.862233	0.209455	0.334793

```

3  0.617839  0.087360  0.600906  0.103527  0.659043  0.981351  0.950721
4  0.428738  0.372519  0.913631  0.098390  0.297539  0.261434  0.528085

```

```

          X7          X8          X9          Y
0  0.017168  0.465338  0.879917  15.430355
1  0.767923  0.967421  0.231402  14.240882
2  0.531339  0.487945  0.649879  18.508117
3  0.602827  0.625784  0.197146   6.212618
4  0.365814  0.666167  0.394064  10.720961

```

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:                19:50:28      Log-Likelihood:        329.92
No. Observations:         10      AIC:                -639.8
Df Residuals:             0      BIC:                -636.8
Df Model:                 9
Covariance Type:         nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	-5.343e-16	inf	-0	nan	nan	nan
x1	0.3674	inf	0	nan	nan	nan
x2	0.1046	inf	0	nan	nan	nan
x3	-0.5237	inf	-0	nan	nan	nan
x4	0.0243	inf	0	nan	nan	nan
x5	-0.2069	inf	-0	nan	nan	nan
x6	-0.8082	inf	-0	nan	nan	nan
x7	0.0638	inf	0	nan	nan	nan
x8	0.0823	inf	0	nan	nan	nan
x9	-0.2179	inf	-0	nan	nan	nan
x10	-0.0232	inf	-0	nan	nan	nan

```

=====
Omnibus:                0.898      Durbin-Watson:          2.669
Prob(Omnibus):          0.638      Jarque-Bera (JB):        0.744
Skew:                   -0.455      Prob(JB):                0.689
Kurtosis:               2.023      Cond. No.                 25.6
=====

```

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 -5.342948e-16

x1 3.673549e-01

x2 1.045949e-01

x3 -5.236686e-01

x4 2.434223e-02

x5 -2.068663e-01

x6 -8.081926e-01

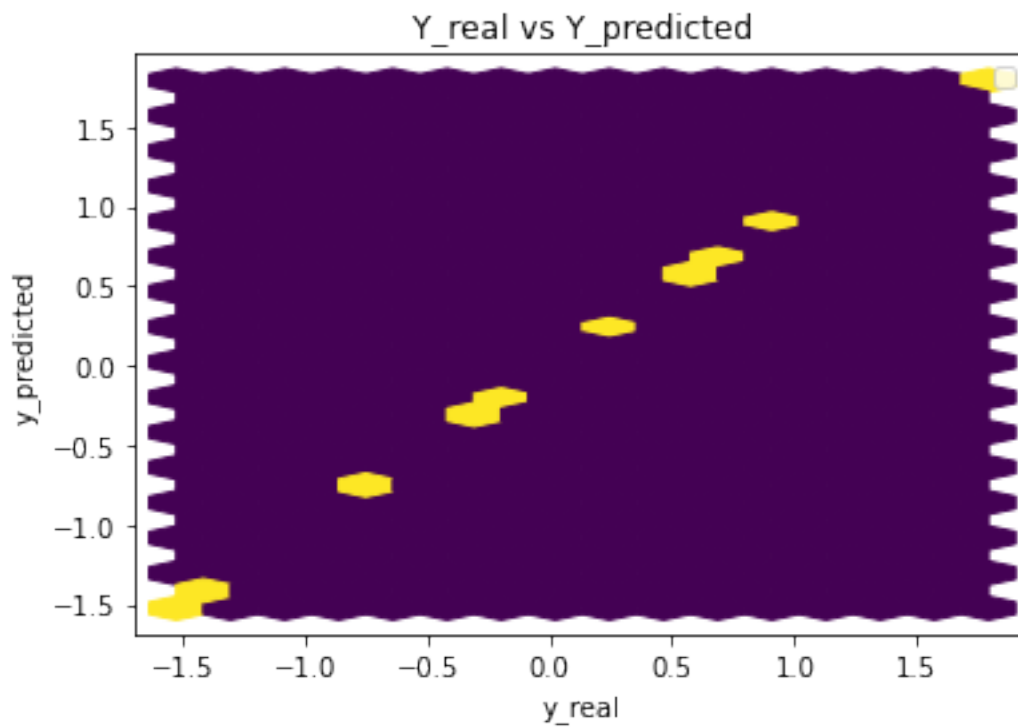
x7 6.376433e-02

x8 8.229014e-02

x9 -2.179198e-01

x10 -2.316199e-02

dtype: float64



Performance Metrics

Mean Squared Error: 1.2909123231238765e-30

Mean Absolute Error: 9.68669588985449e-16

Manhattan distance: 9.686695889854491e-15

Euclidean distance: 3.592926833549323e-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

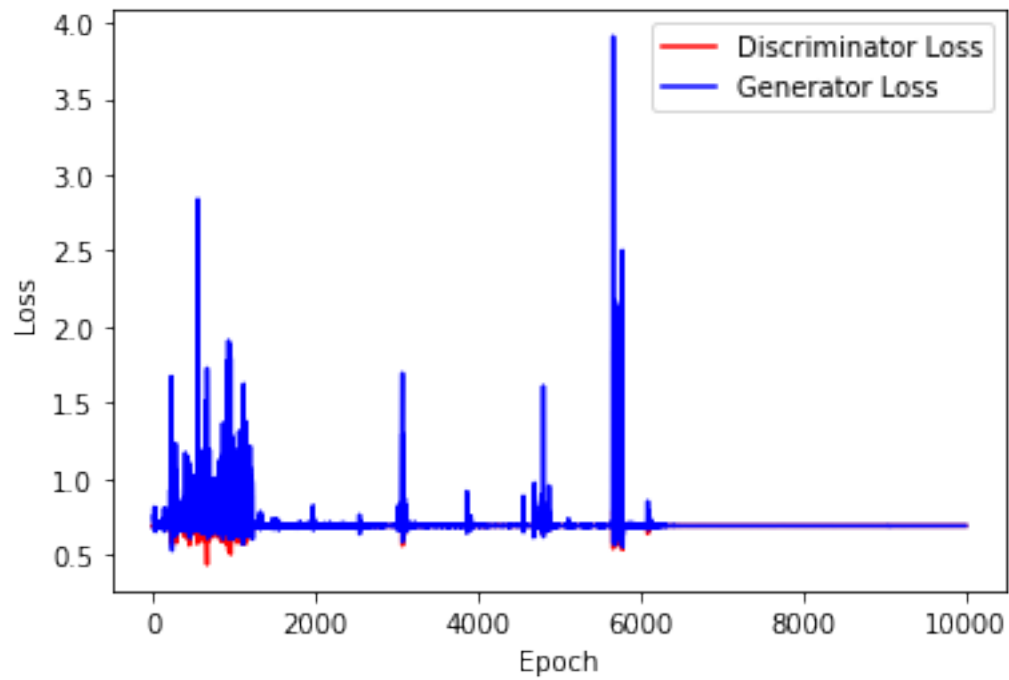
```
[9]: generator = network.Generator(n_features+2)
     discriminator = network.Discriminator(n_features+2)

     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))
```

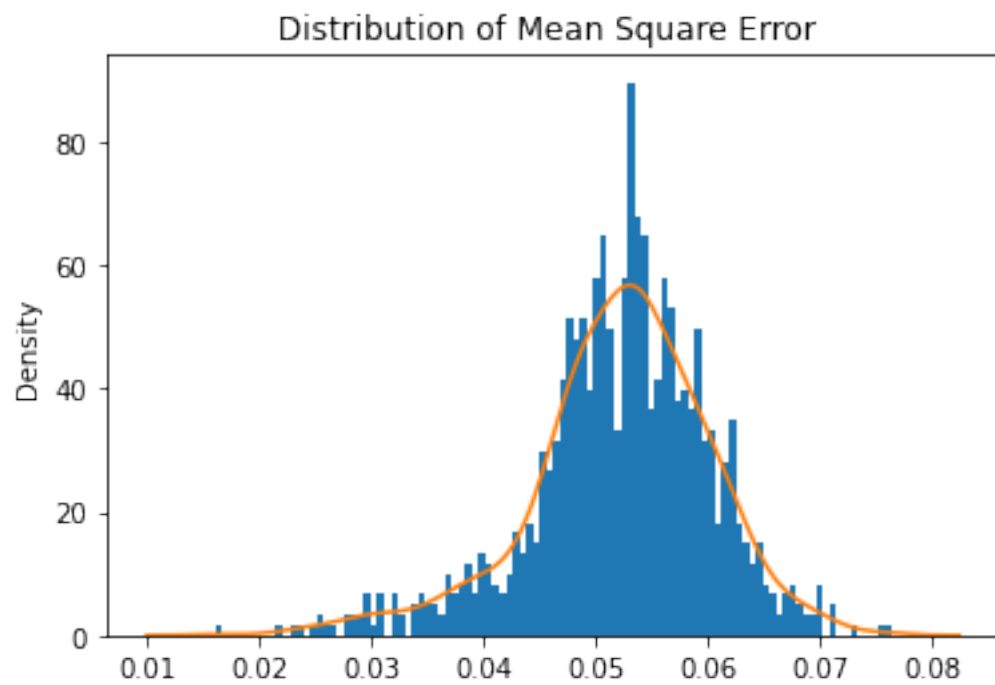
```
[10]: print(generator)
      print(discriminator)
```

```
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()
)
```

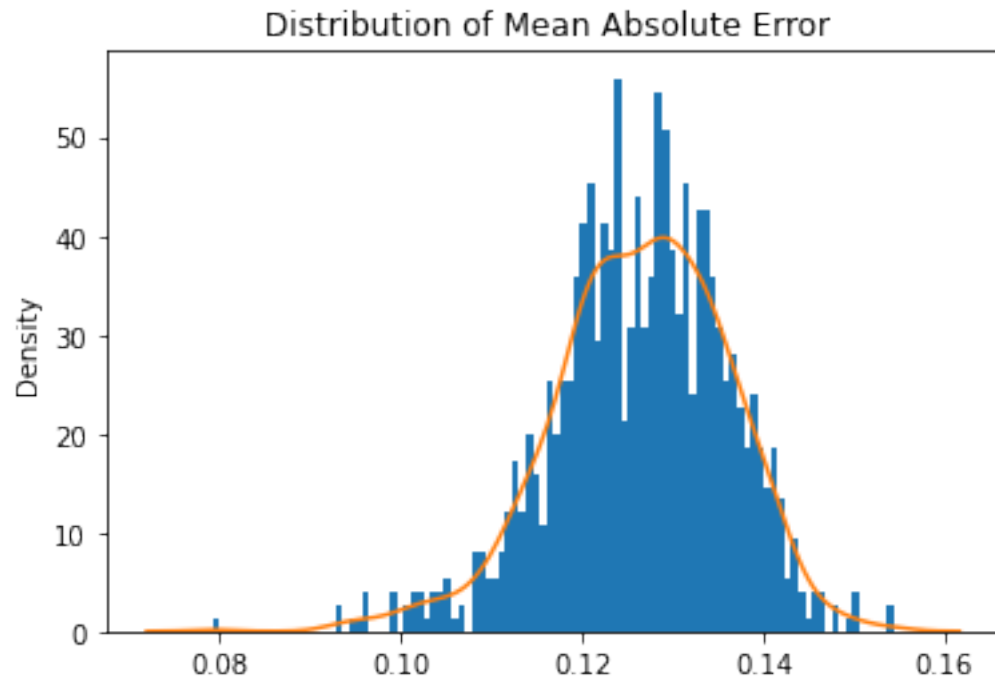
```
[11]: train_test.
     ↪training_GAN(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,
     ↪n_epochs,criterion,device)
```



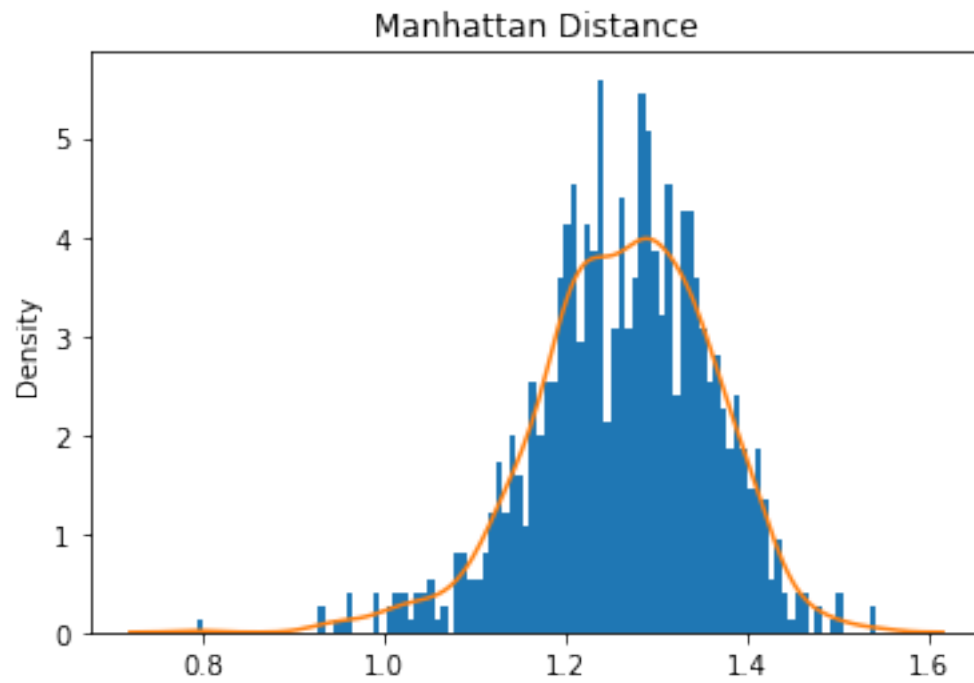
```
[12]: train_test.test_generator(generator,real_dataset,device)
```



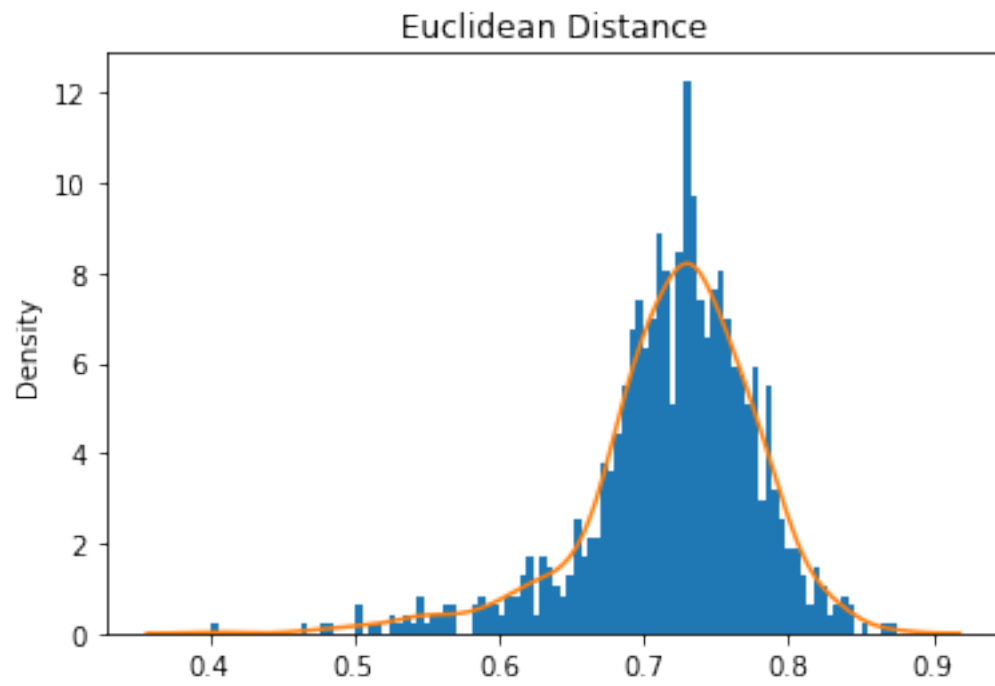
Mean Square Error: 0.052229789696043666



Mean Absolute Error: 0.12641233491003515

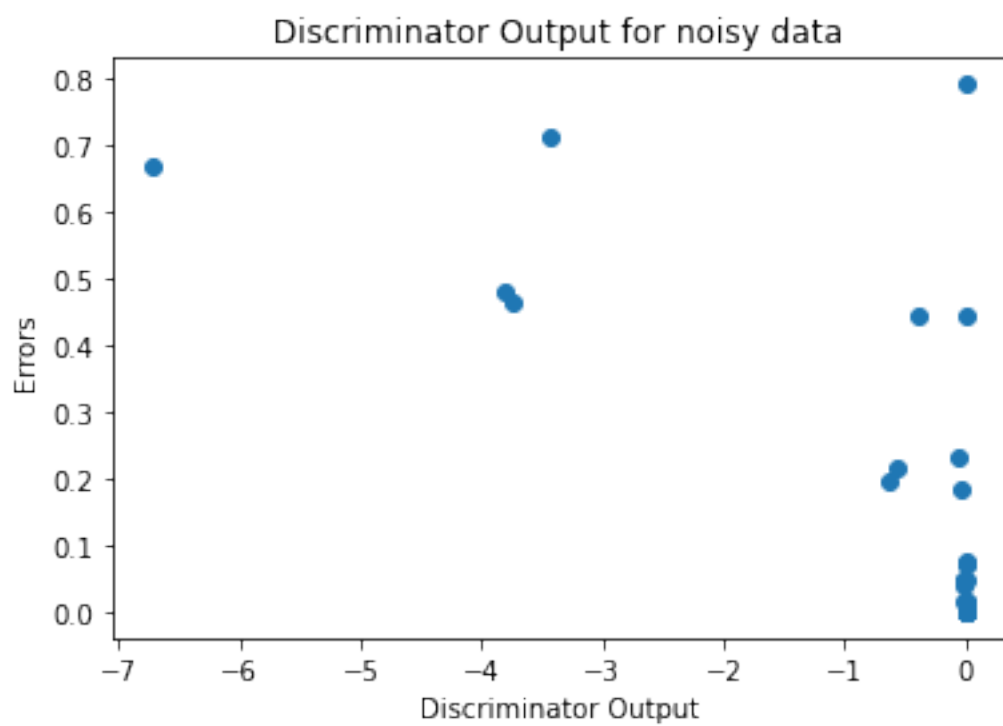
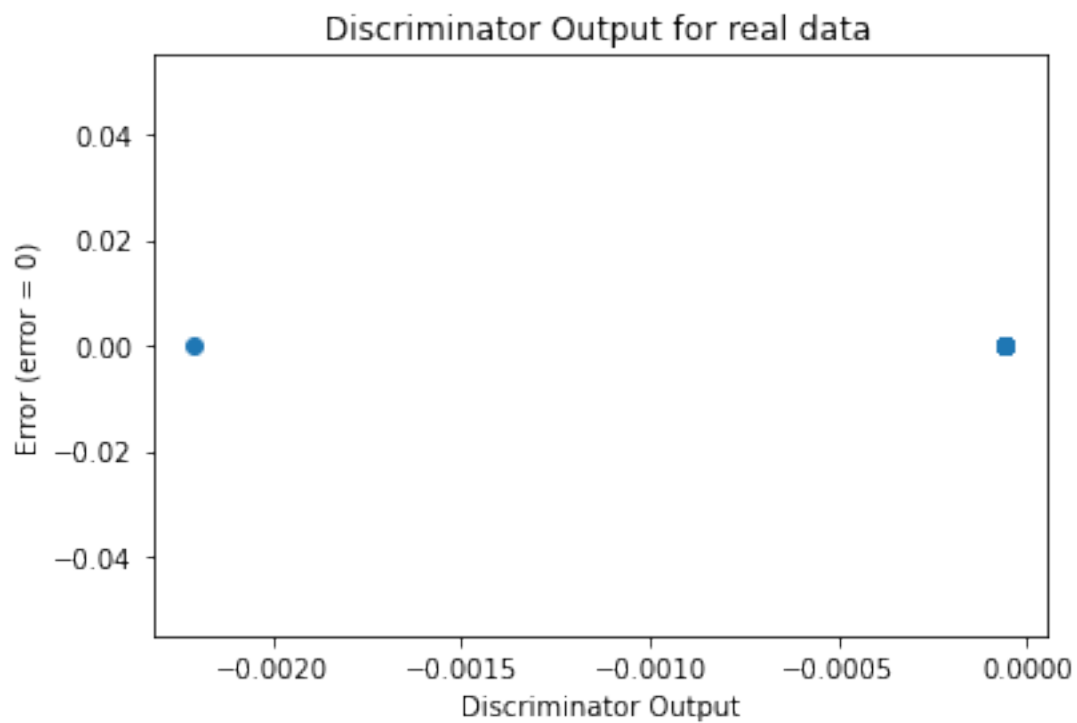


Mean Manhattan Distance: 1.2641233491003514



Mean Euclidean Distance: 0.720307347557266

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

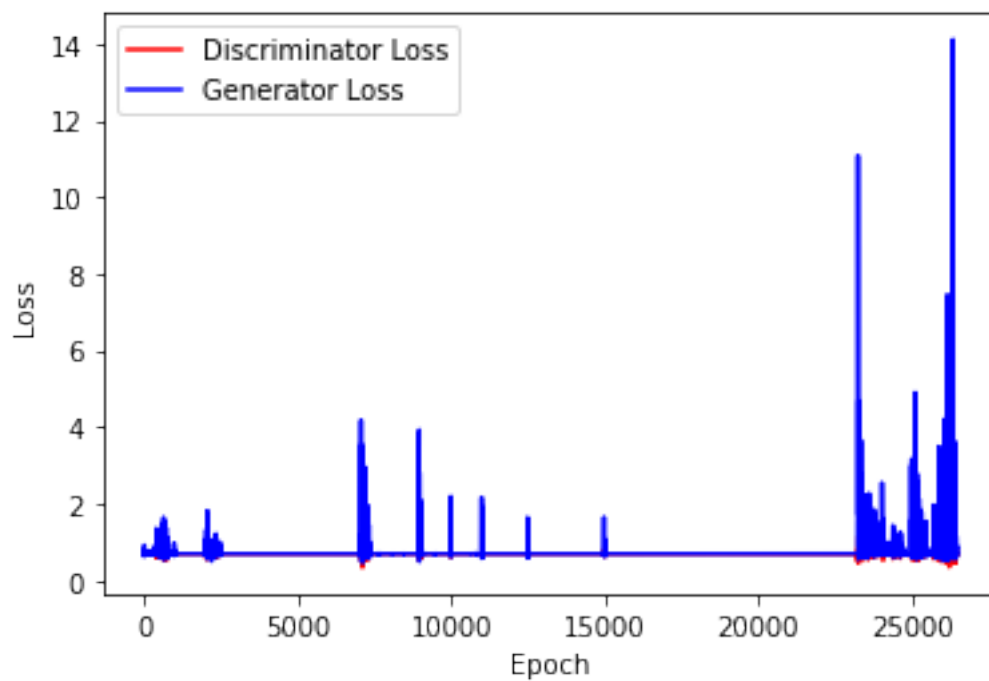



Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

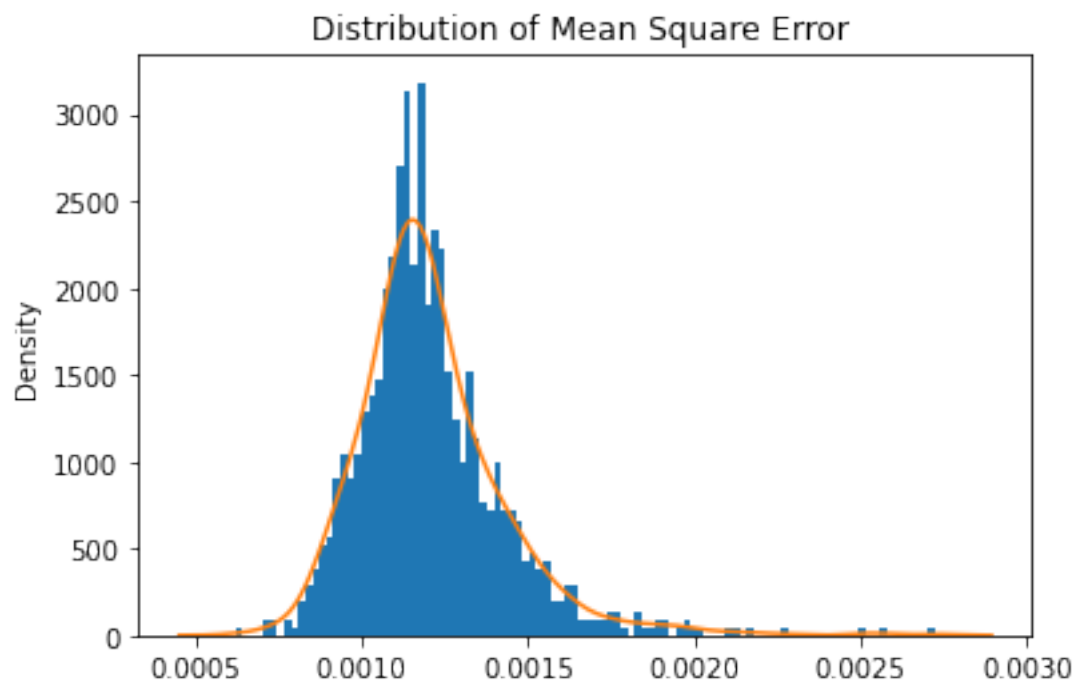
```
[14]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator.parameters(), lr=0.01, betas=(0.5,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

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

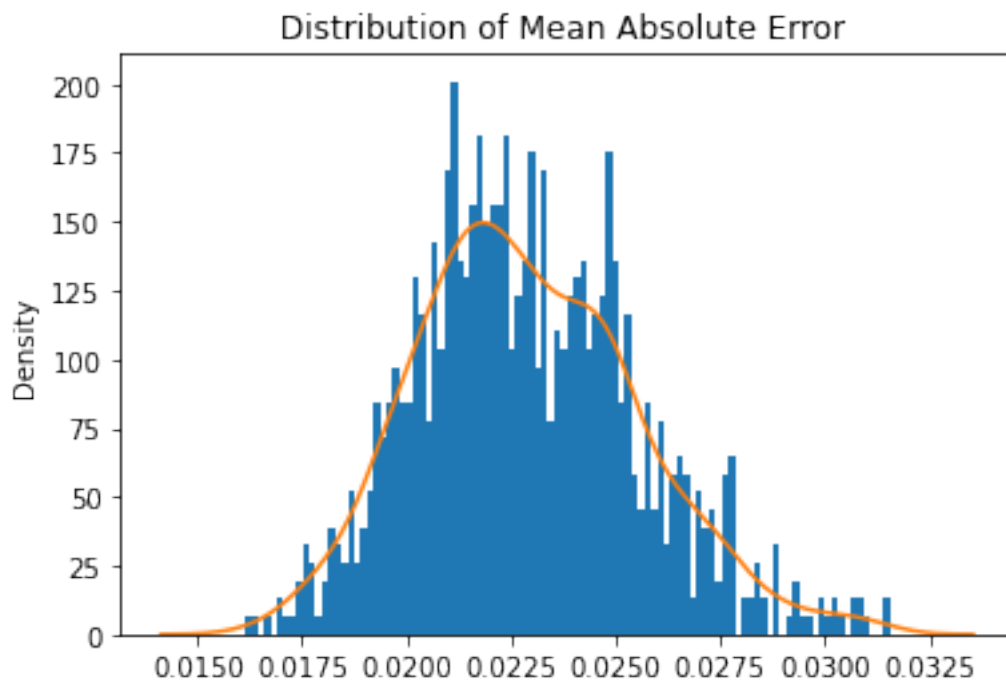
Number of epochs needed 13241



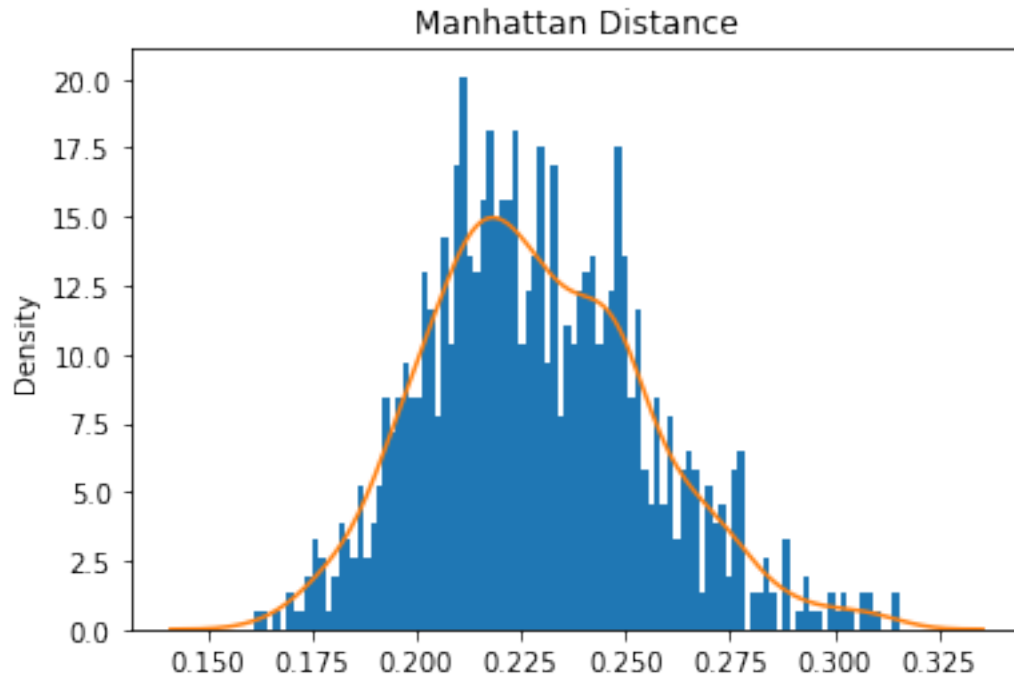
```
[16]: train_test.test_generator(generator,real_dataset,device)
```



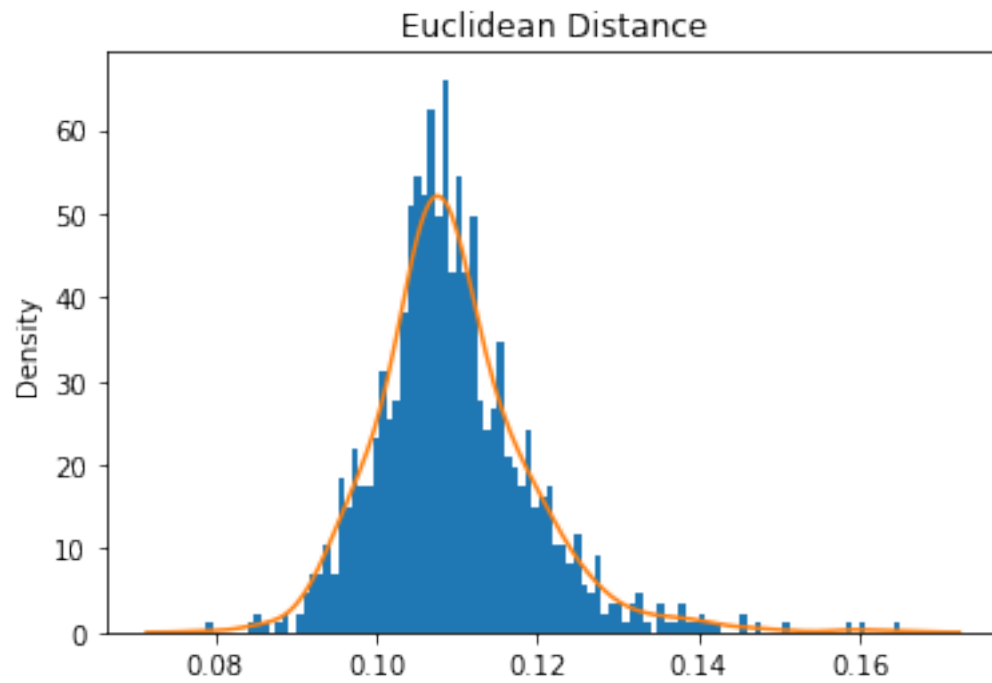
Mean Square Error: 0.0012096561543101353



Mean Absolute Error: 0.022875970375537874



Mean Manhattan Distance: 0.2287597037553787



Mean Euclidean Distance: 0.10955465578817648

2 ABC GAN Model

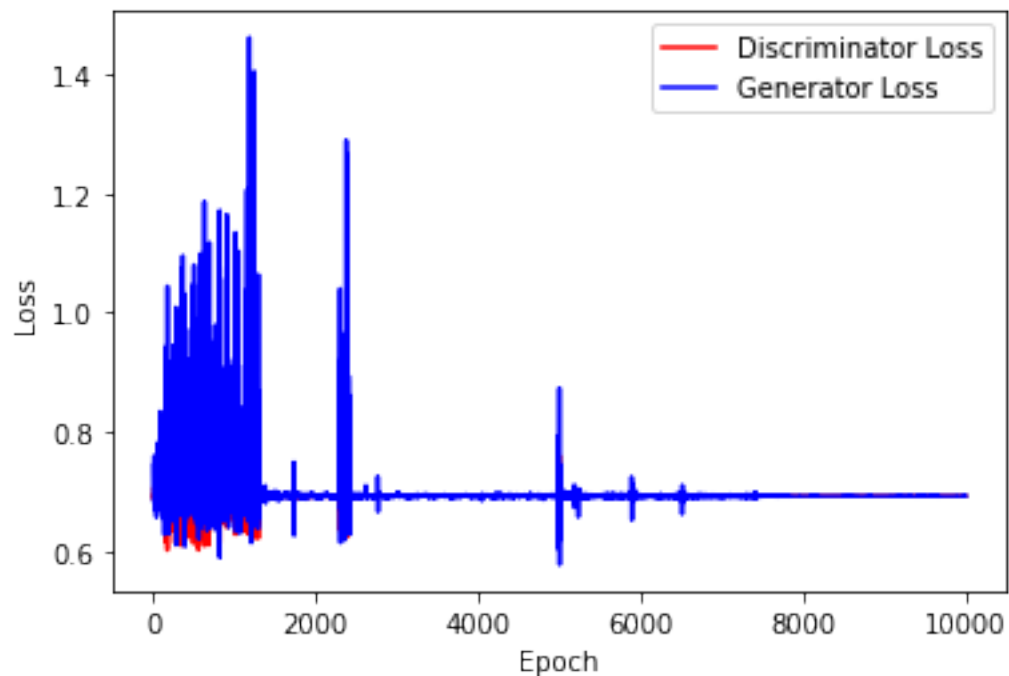
2.0.1 Training the network

Training ABC-GAN for `n_epochs` number of epochs

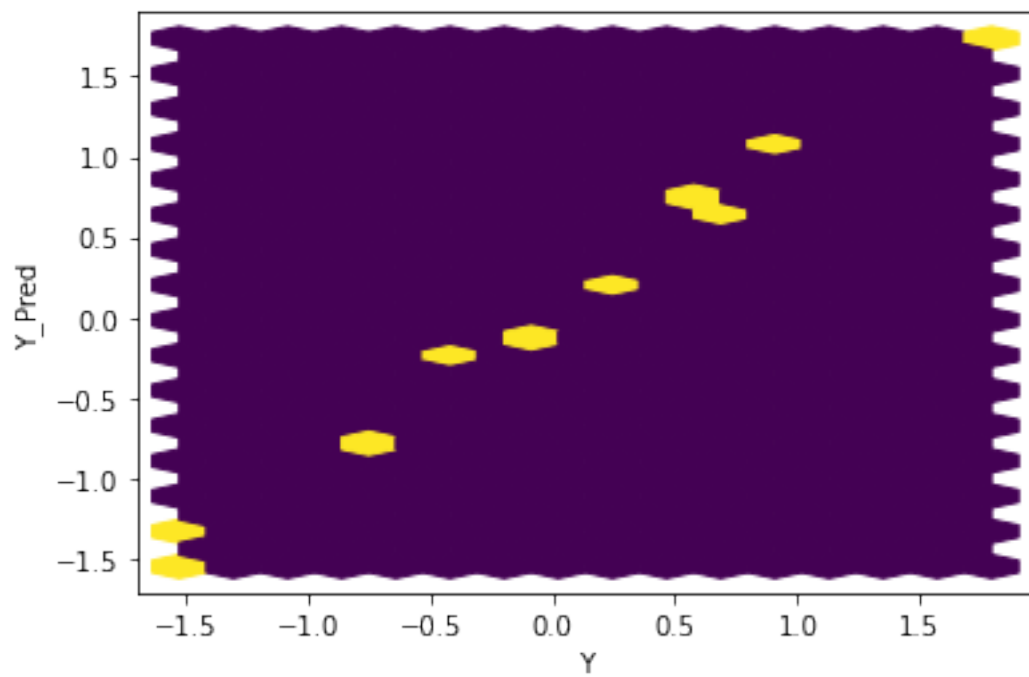
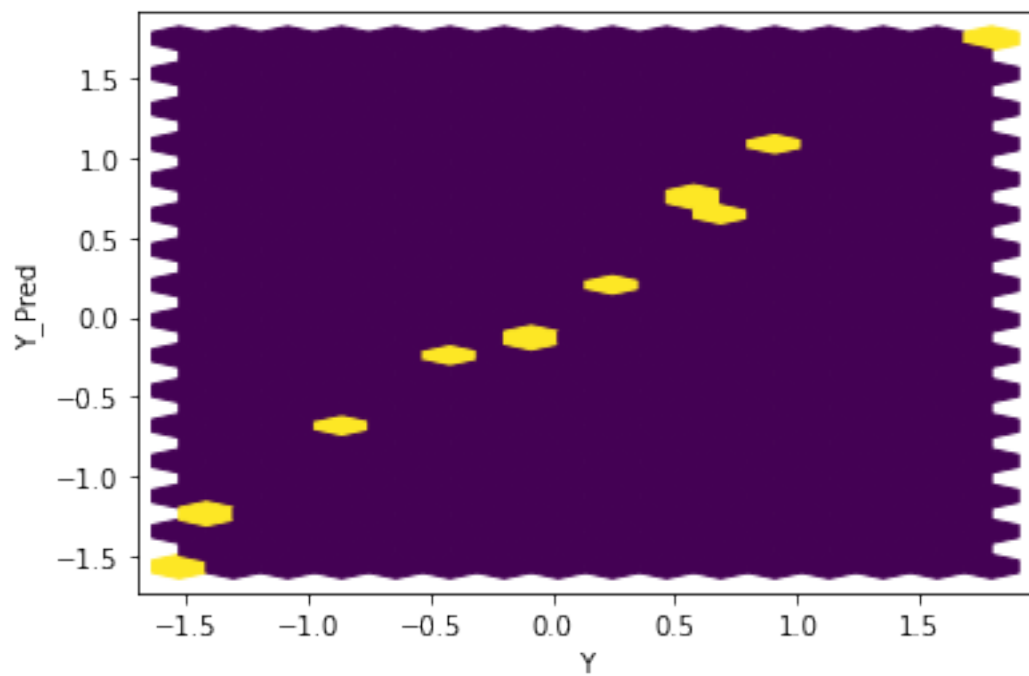
```
[17]: gen = network.Generator(n_features+2)
      disc = network.Discriminator(n_features+2)

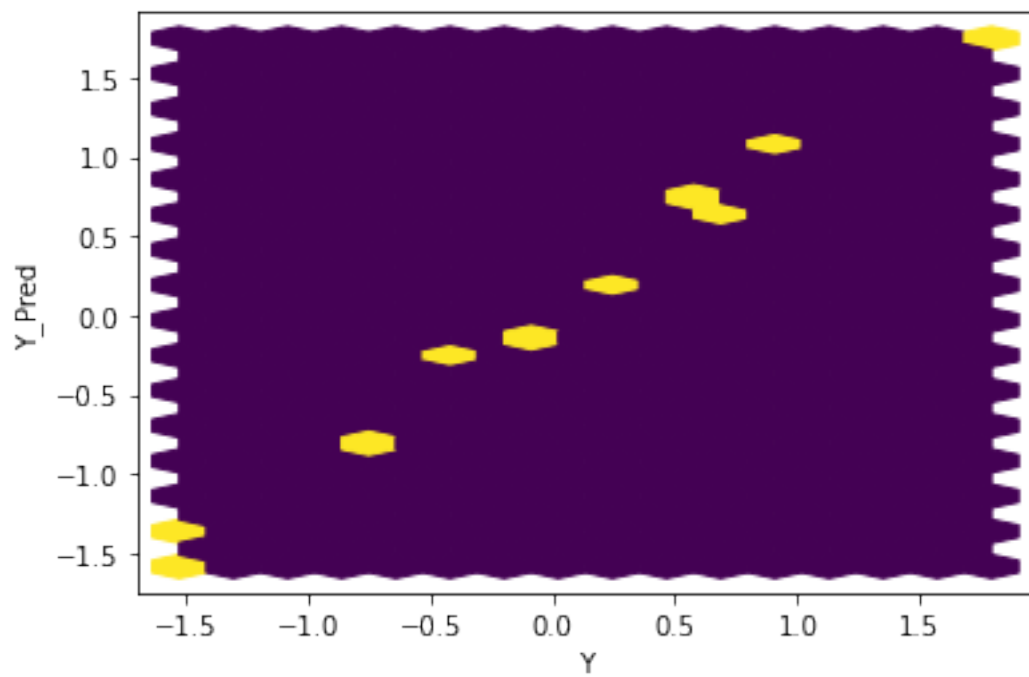
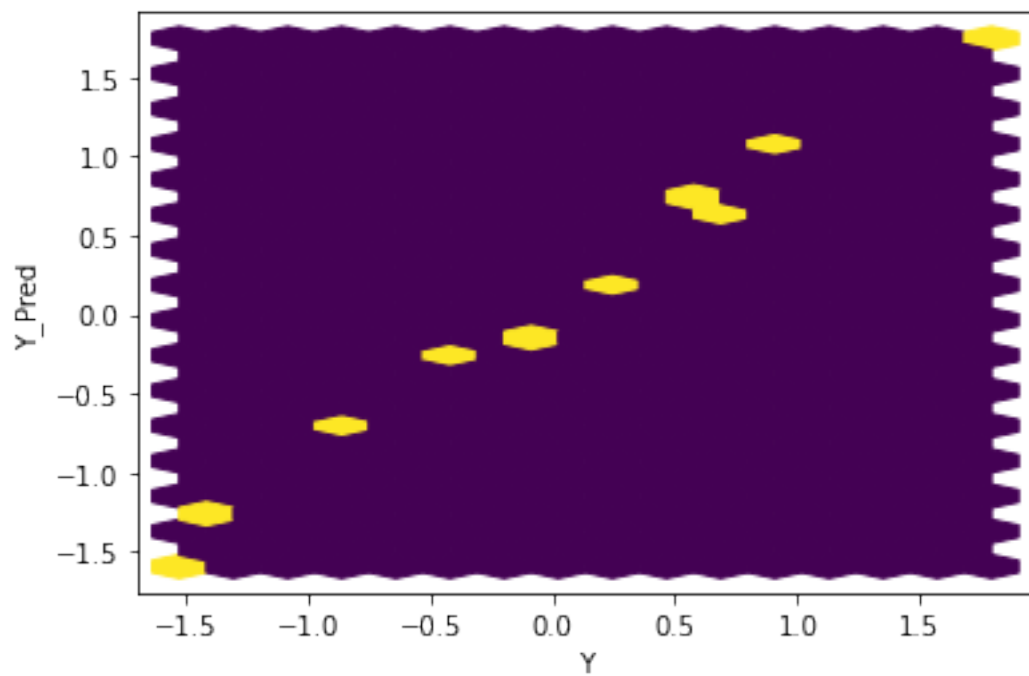
      criterion = torch.nn.BCEWithLogitsLoss()
      gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
      disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))

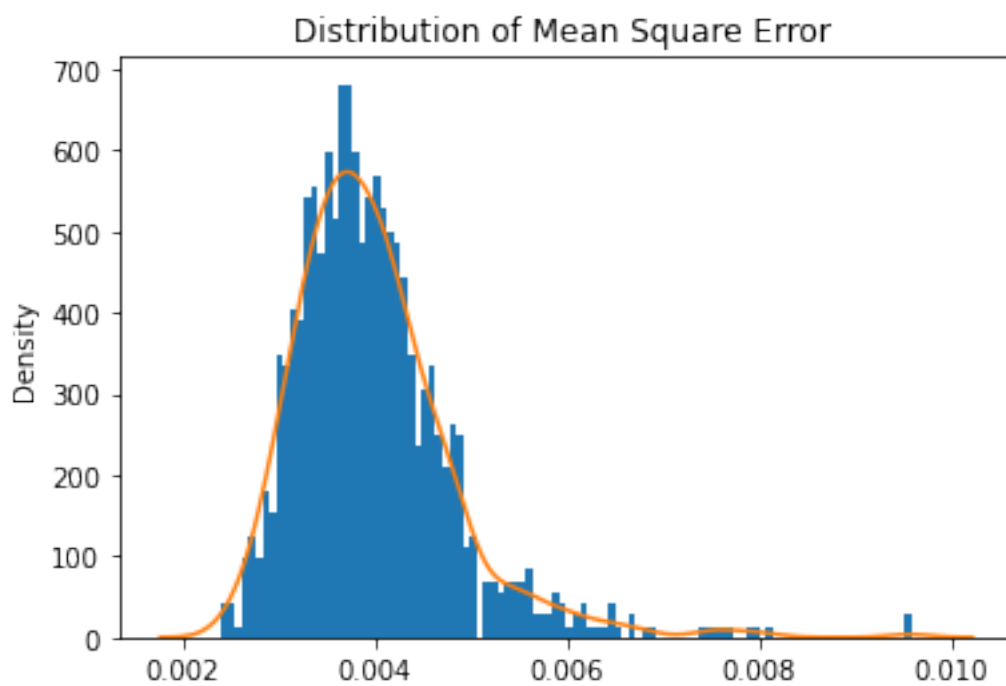
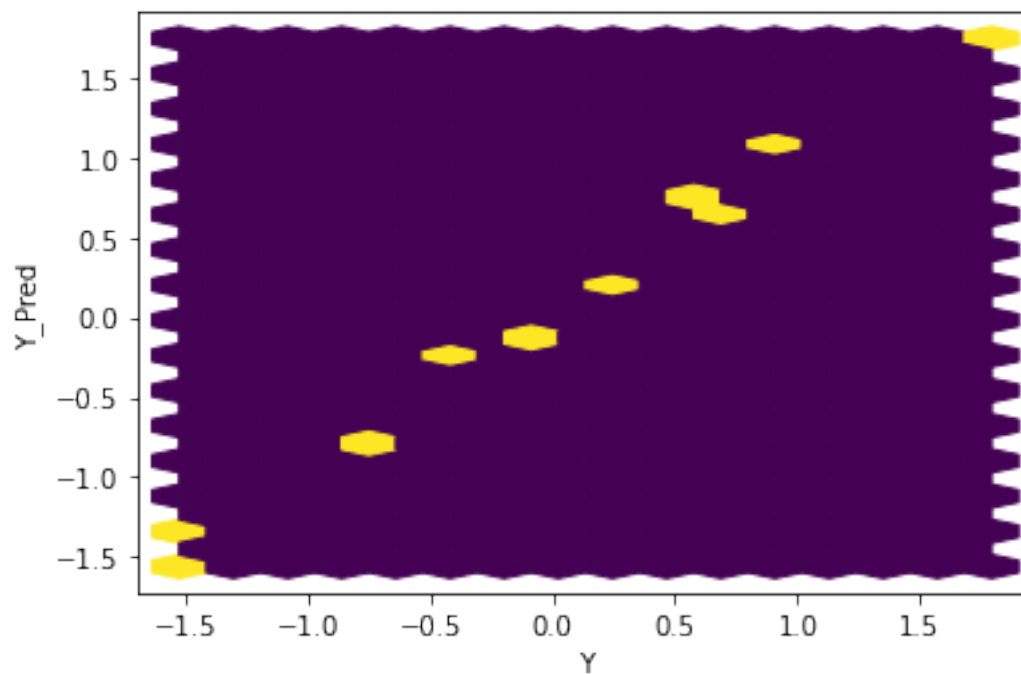
[18]: ABC_train_test.training_GAN(disc, gen, disc_opt, gen_opt, real_dataset,
      ↪ batch_size, n_epochs, criterion, coeff, mean, variance, device)
```



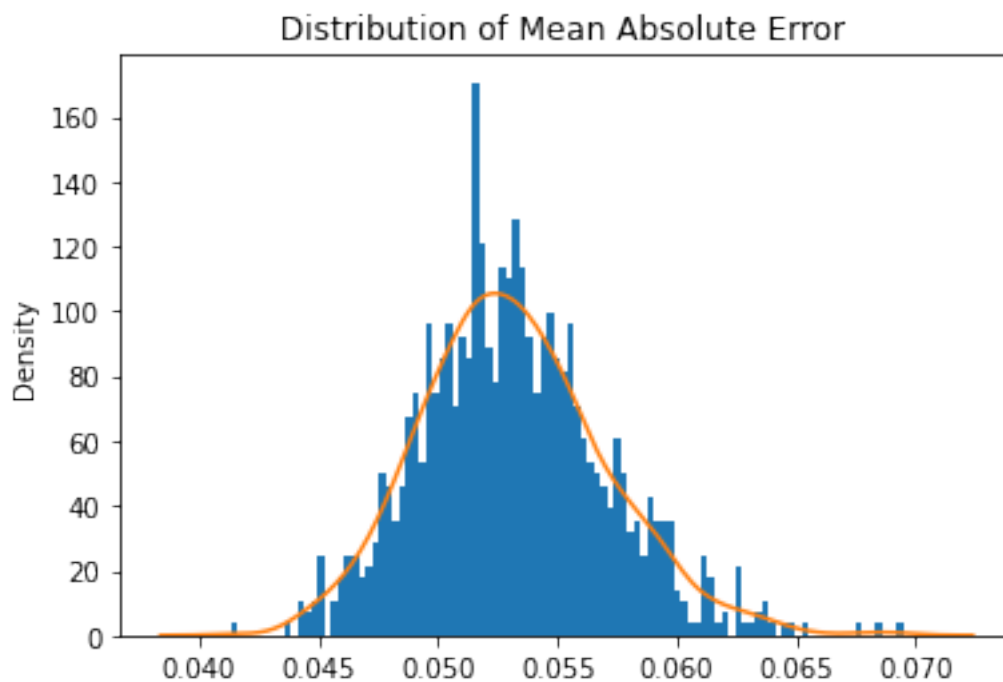
```
[19]: ABC_train_test.test_generator(gen, real_dataset, coeff, mean, variance, device)
```





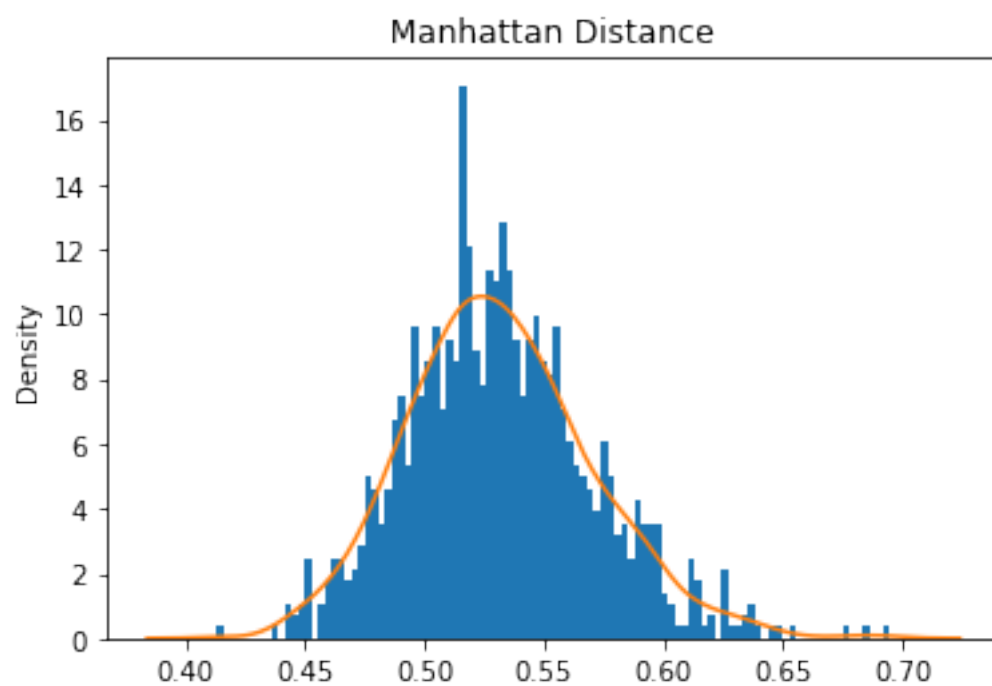


Mean Square Error: 0.003976970537413215

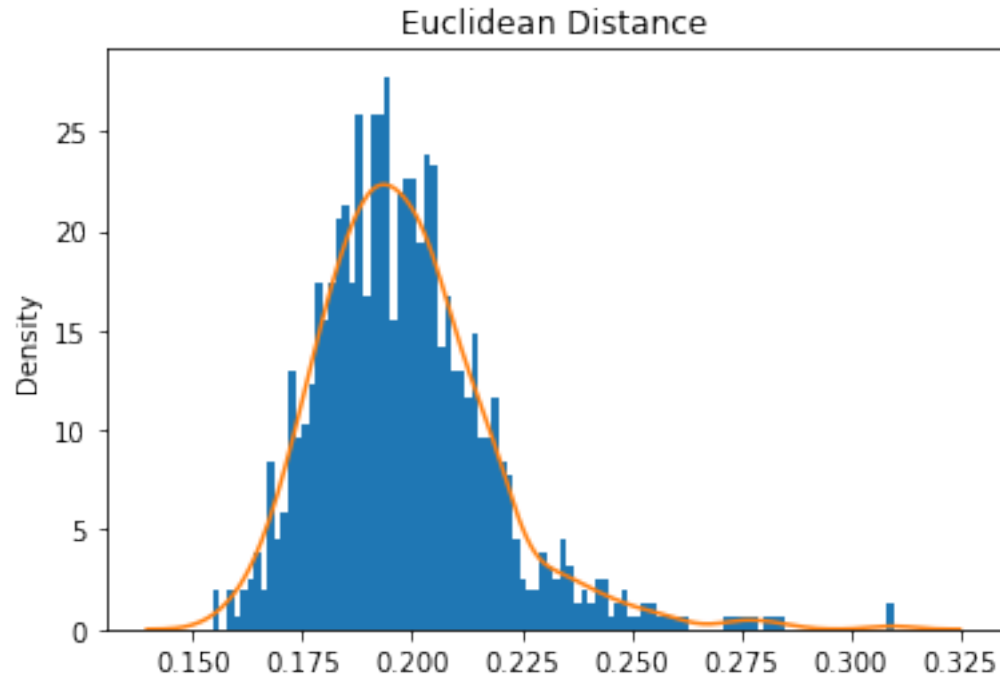


Mean Absolute Error: 0.05314368405193091

Mean Manhattan Distance: 0.5314368405193091

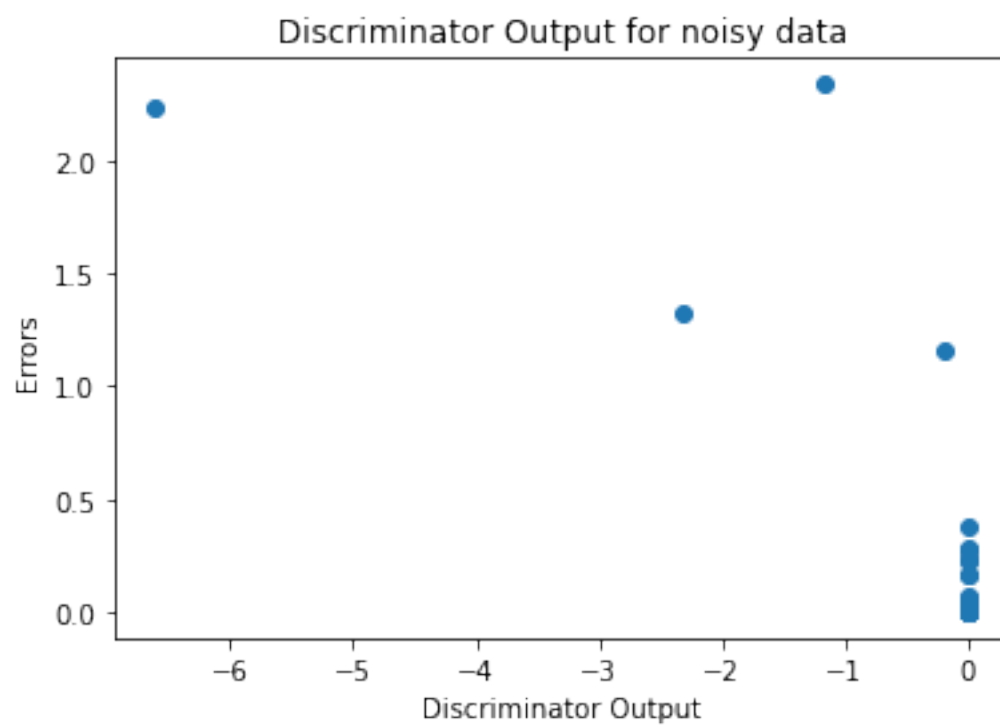
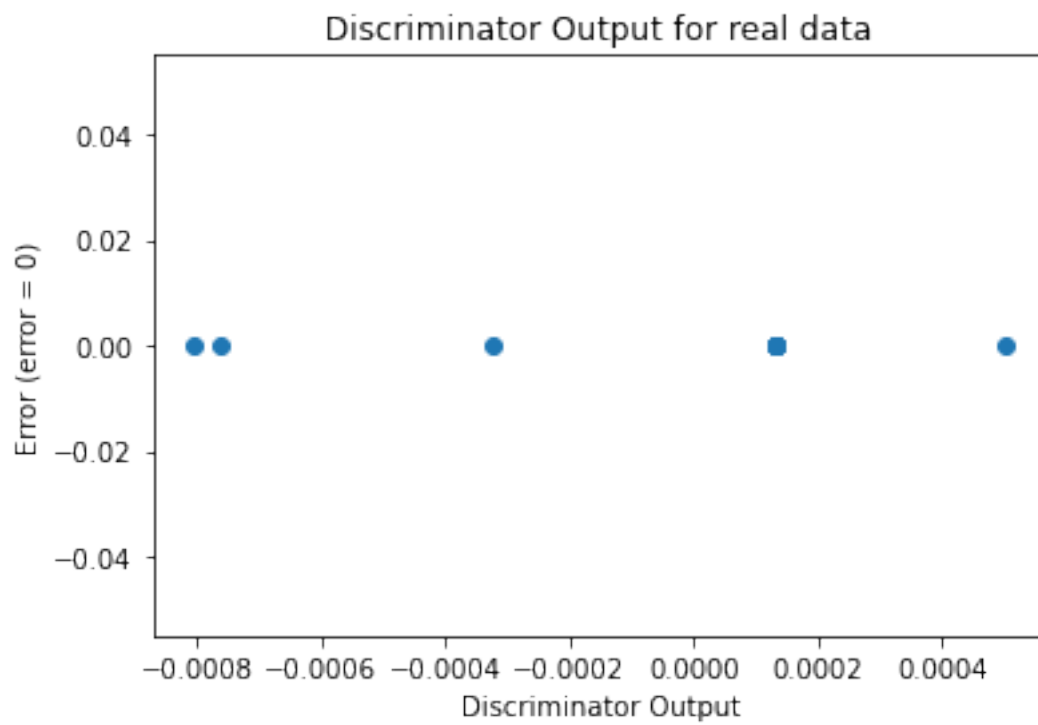


Mean Euclidean Distance: 0.19842150602520428



Sanity Checks

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



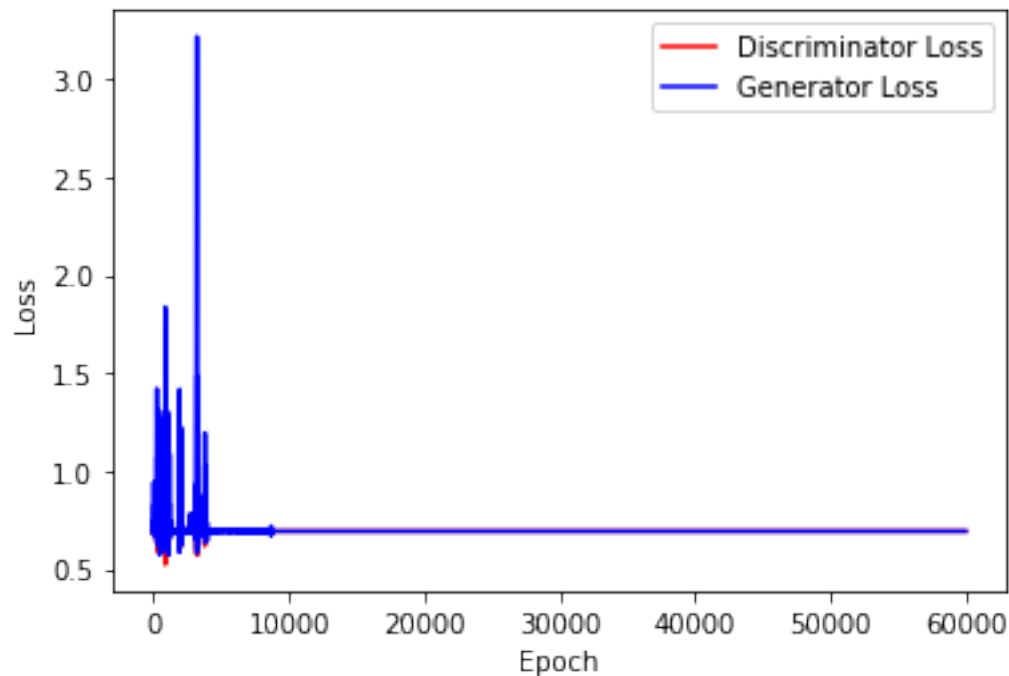
Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

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

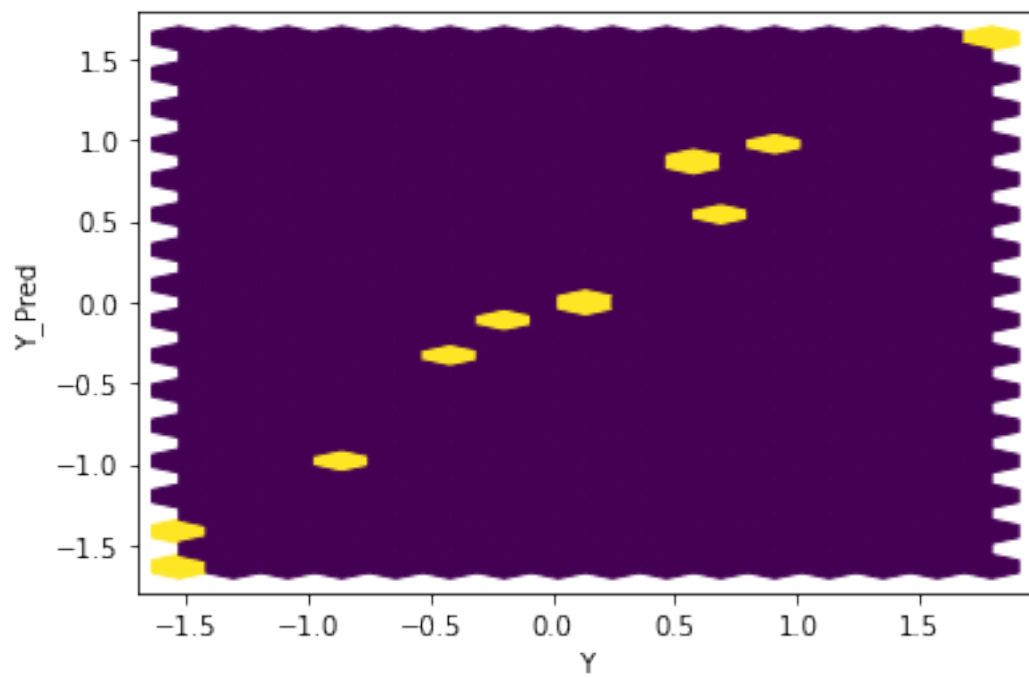
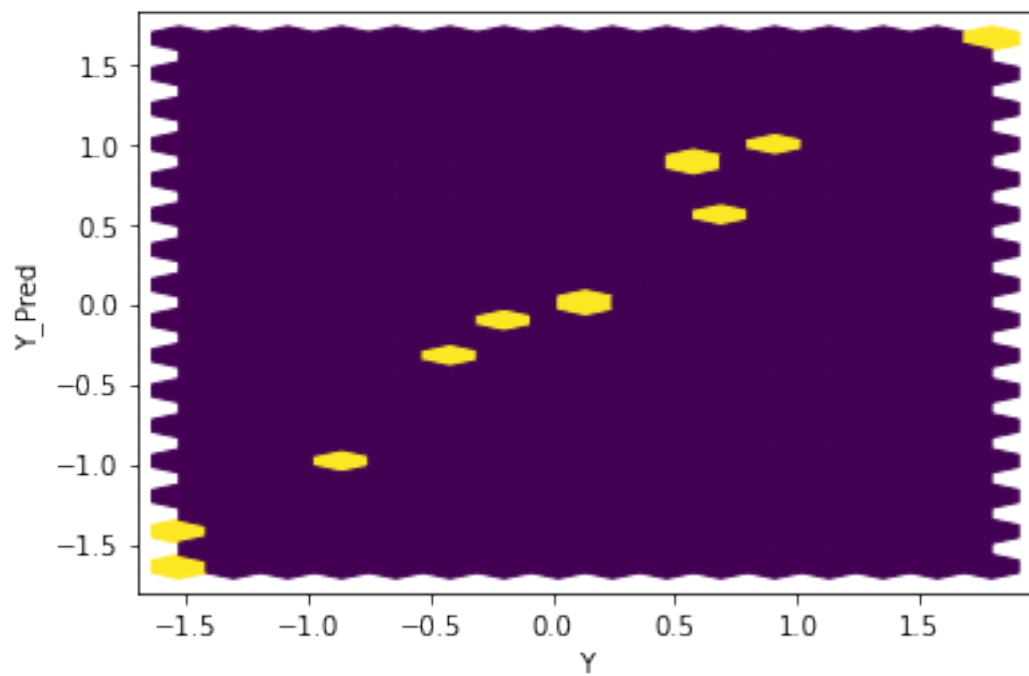
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))
```

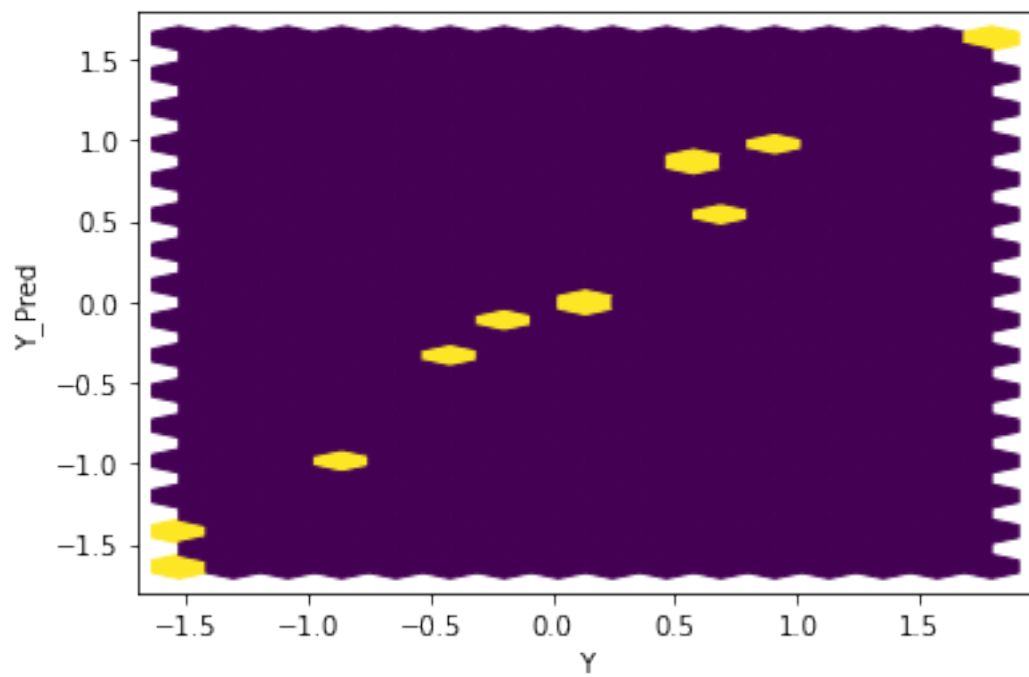
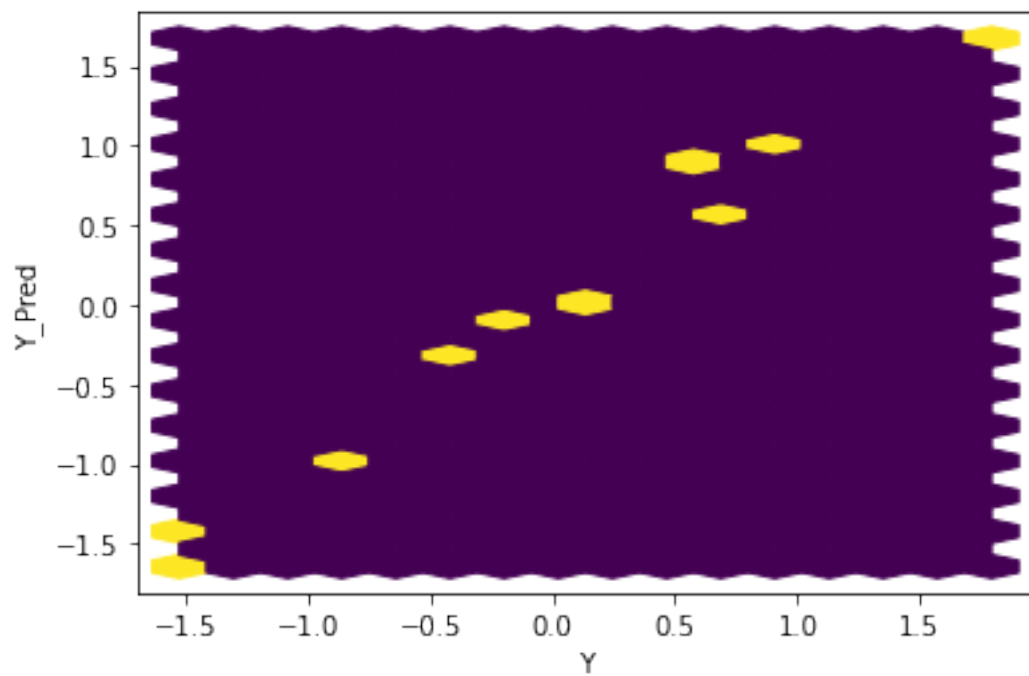
```
[22]: ABC_train_test.
      ↪ training_GAN_2(disc,gen,disc_opt,gen_opt,real_dataset,batch_size,
      ↪ error,criterion,coeff,mean,variance,device)
```

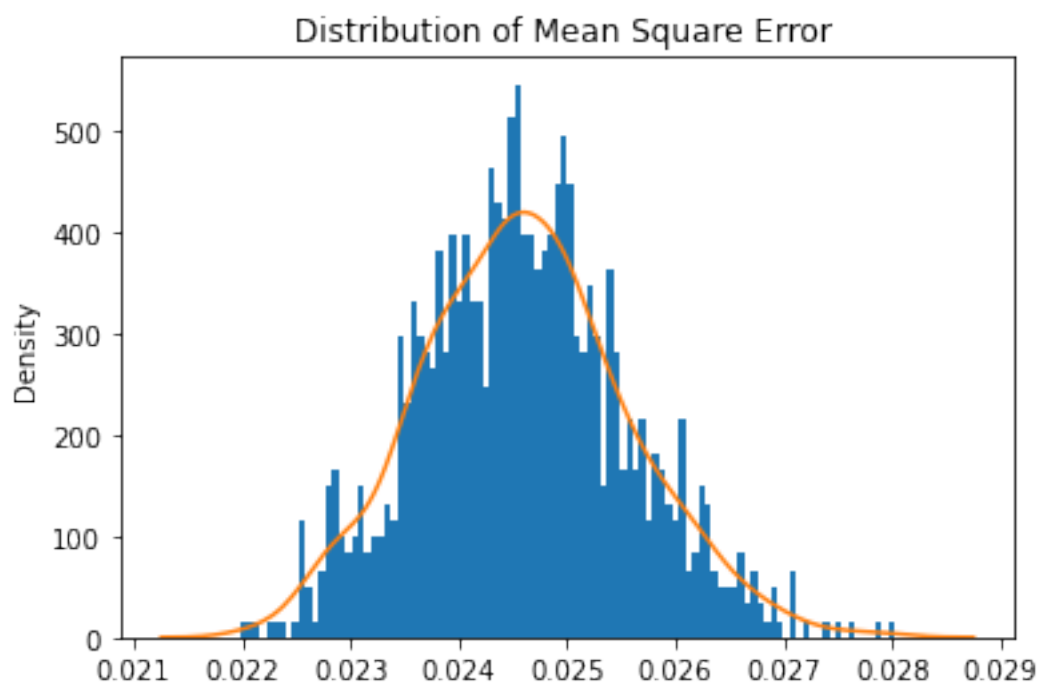
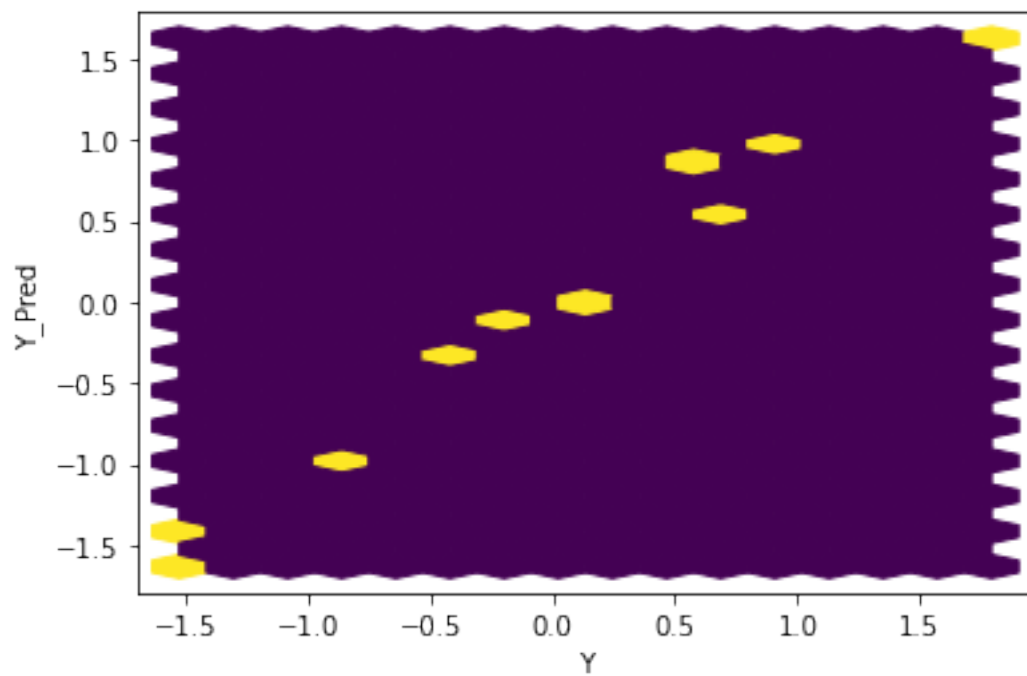
Number of epochs 30000



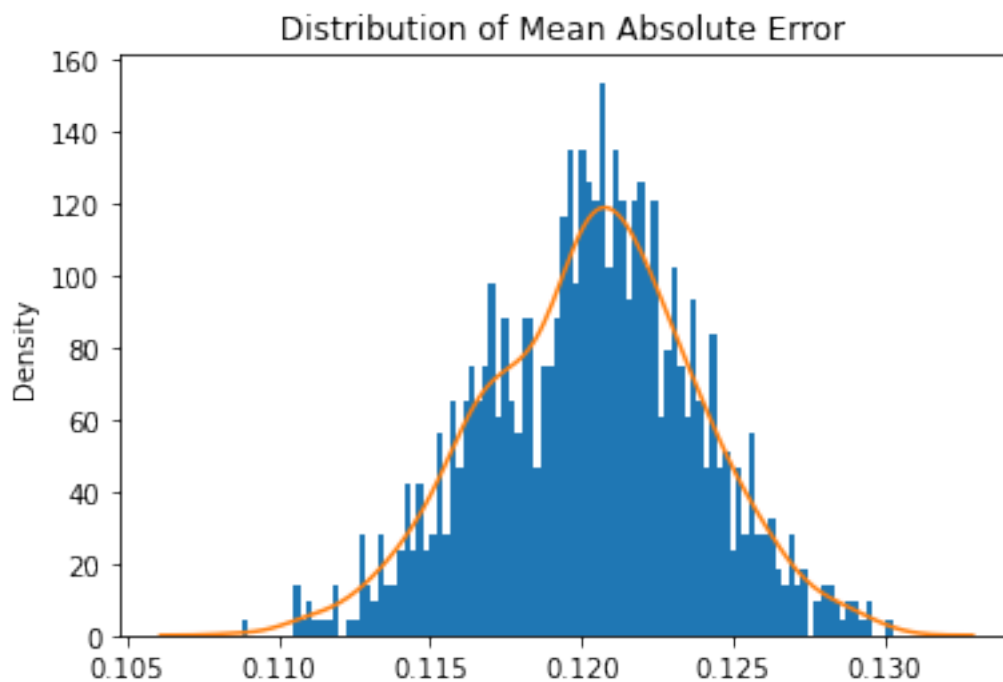
```
[23]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```





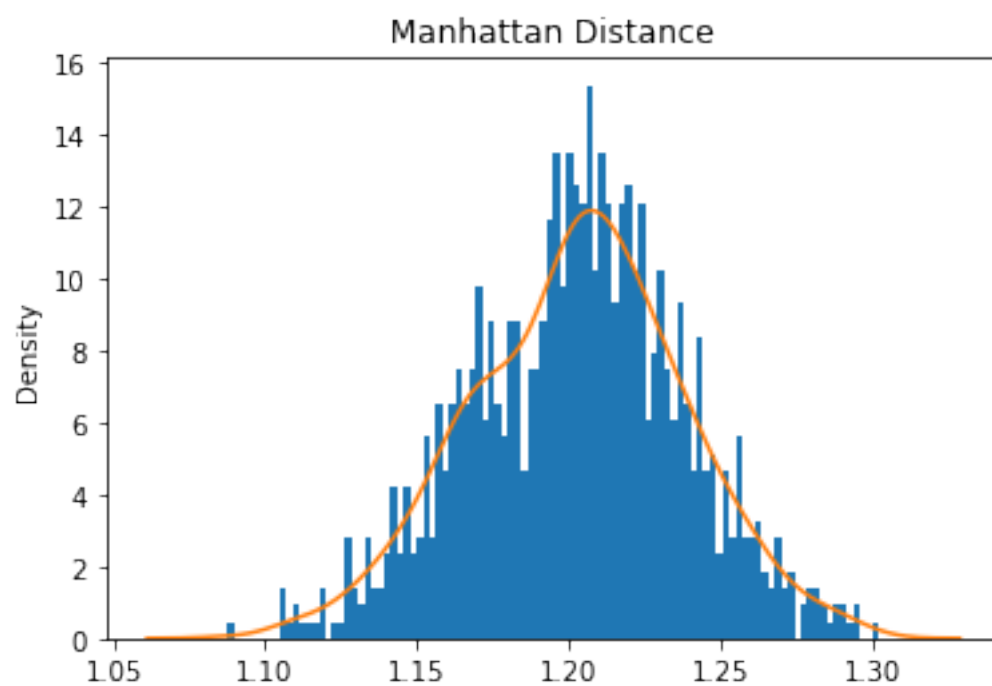


Mean Square Error: 0.024600701288361164

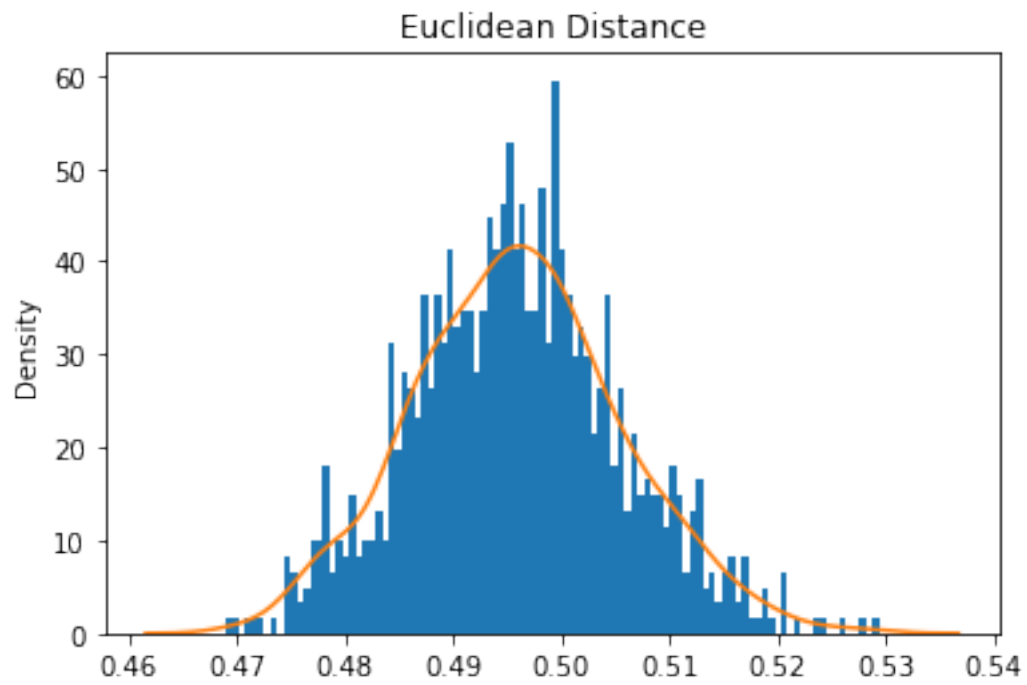


Mean Absolute Error: 0.1202955771021545

Mean Manhattan Distance: 1.2029557710215448



Mean Euclidean Distance: 0.49589572897342116



[]: