

Dataset1-Regression_output_6

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. $Y = 1 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n + N(0, \sigma)$ 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 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

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	-1.164009	0.075330	0.207145	0.793251	-0.210853	0.568334	0.689743
1	1.561790	-0.523979	0.085241	0.772850	1.608976	1.384342	-1.769365
2	-0.136666	-2.211915	-0.465028	-1.305688	0.101391	0.703763	1.240547
3	-0.377286	0.083002	-0.587552	-1.391446	0.231825	-0.837071	-0.201160
4	1.216609	-1.674807	3.317122	1.345199	0.245888	0.045475	-0.722295
	X8	X9	X10	Y			

```

0  1.328015  0.184057 -0.351956  238.192437
1 -0.098416 -0.574409  0.695244   97.977653
2  1.934852 -0.304928  0.281480   53.889129
3  0.341796 -1.236662  1.788613 -237.557294
4  0.134729 -0.811826  0.129812  274.776791

```

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:             5.096e+07
Date:                Tue, 19 Oct 2021  Prob (F-statistic):       3.84e-296
Time:                  23:21:21      Log-Likelihood:          636.13
No. Observations:                100      AIC:                  -1250.
Df Residuals:                    89      BIC:                  -1222.
Df Model:                        10
Covariance Type:                nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	-4.163e-17	4.43e-05	-9.4e-13	1.000	-8.8e-05	8.8e-05
x1	0.1715	4.53e-05	3782.368	0.000	0.171	0.172
x2	0.2008	4.59e-05	4376.364	0.000	0.201	0.201
x3	0.4351	4.51e-05	9646.652	0.000	0.435	0.435
x4	0.3444	4.64e-05	7425.271	0.000	0.344	0.344
x5	0.2687	4.73e-05	5683.277	0.000	0.269	0.269
x6	0.2571	4.61e-05	5581.485	0.000	0.257	0.257
x7	0.3526	4.66e-05	7562.614	0.000	0.352	0.353
x8	0.4313	4.75e-05	9084.265	0.000	0.431	0.431
x9	0.3219	4.67e-05	6893.523	0.000	0.322	0.322
x10	0.1204	4.8e-05	2506.362	0.000	0.120	0.121

```

=====
Omnibus:                        1.411      Durbin-Watson:           1.980
Prob(Omnibus):                  0.494      Jarque-Bera (JB):        0.884
Skew:                          0.162      Prob(JB):                0.643
Kurtosis:                      3.328      Cond. No.:               1.64
=====

```

Notes:

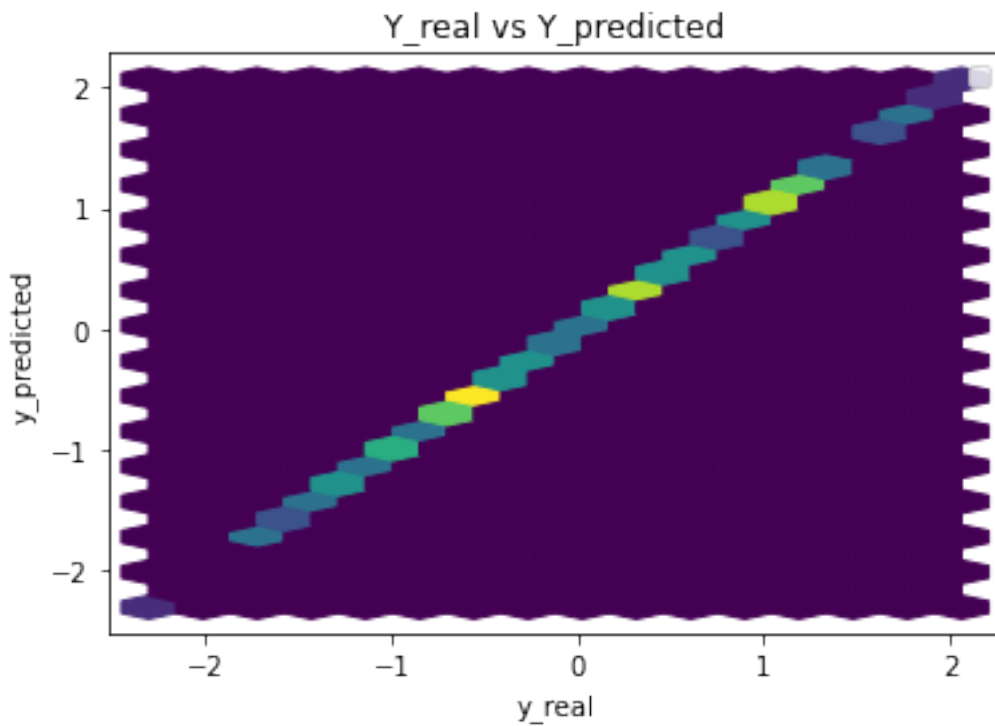
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```

Parameters:  const  -4.163336e-17
            x1      1.714537e-01

```

```
x2      2.008155e-01
x3      4.350730e-01
x4      3.444022e-01
x5      2.686729e-01
x6      2.571092e-01
x7      3.525552e-01
x8      4.312585e-01
x9      3.219233e-01
x10     1.204105e-01
dtype: float64
```



Performance Metrics

```
Mean Squared Error: 1.7465267099958894e-07
Mean Absolute Error: 0.00032931217013065046
Manhattan distance: 0.03293121701306505
Euclidean distance: 0.004179146695194952
```

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

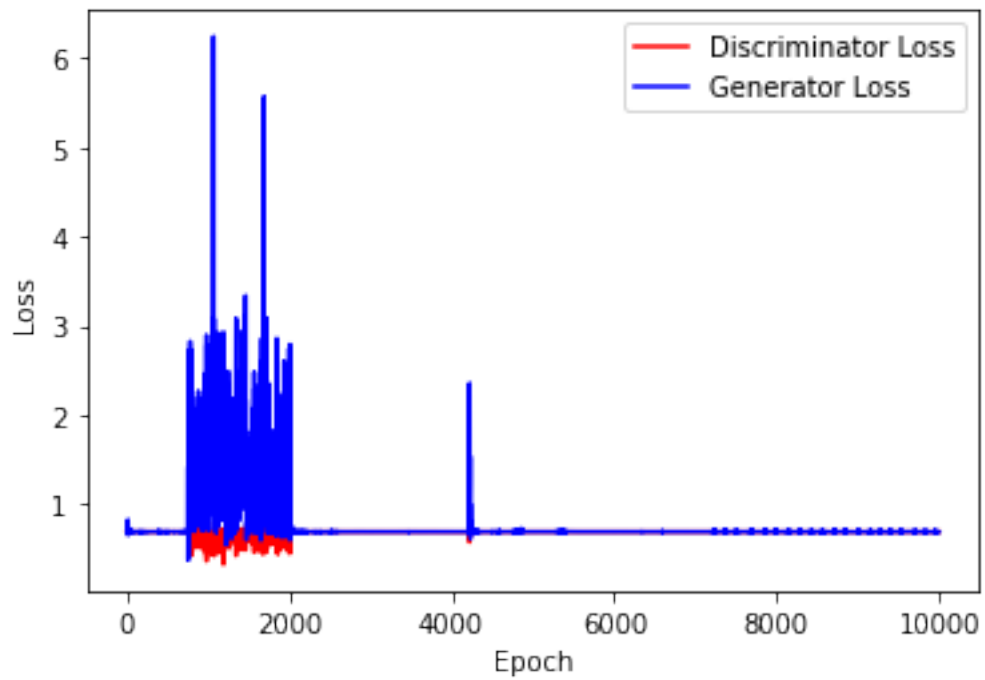
```
[9]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)

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))
```

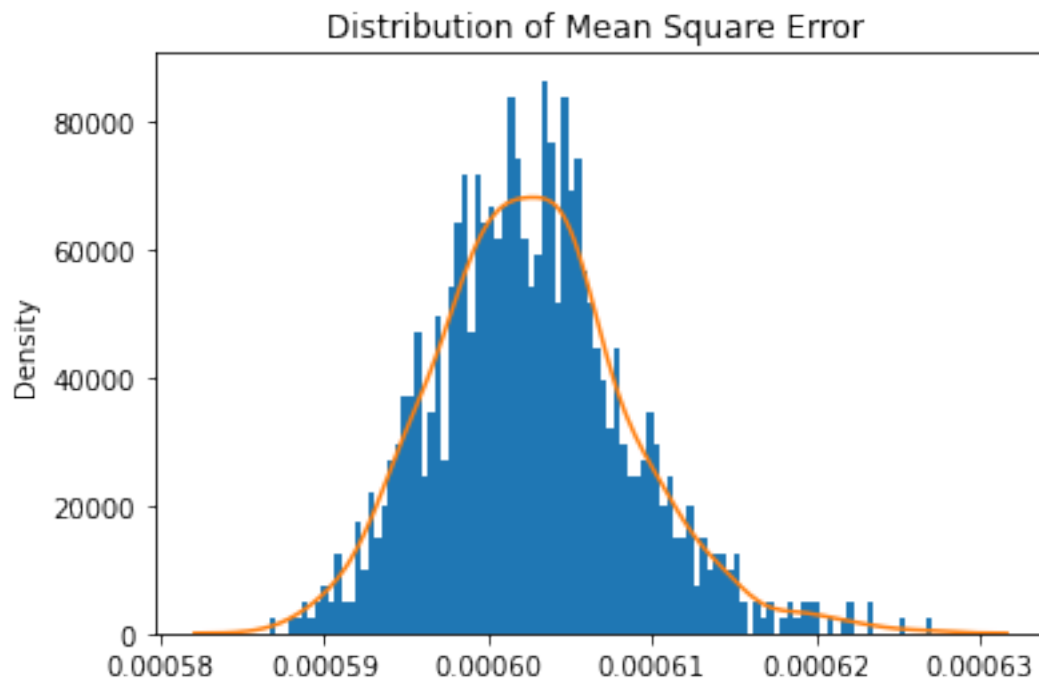
```
[10]: print(generator)
print(discriminator)
```

```
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()
)
```

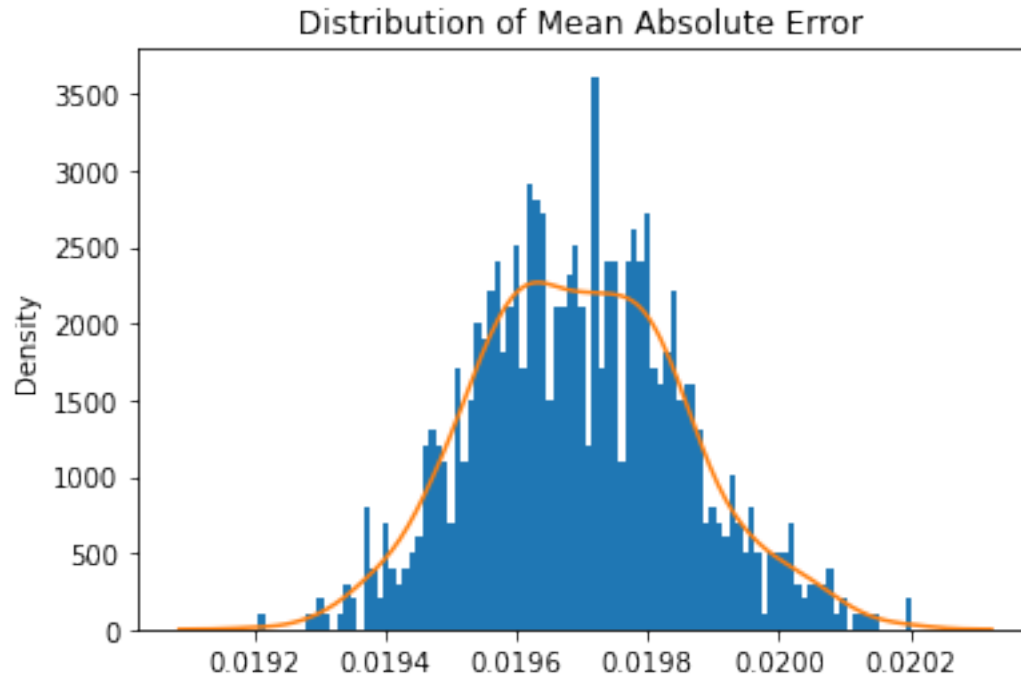
```
[11]: train_test.
→training_GAN(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,
→n_epochs,criterion,device)
```



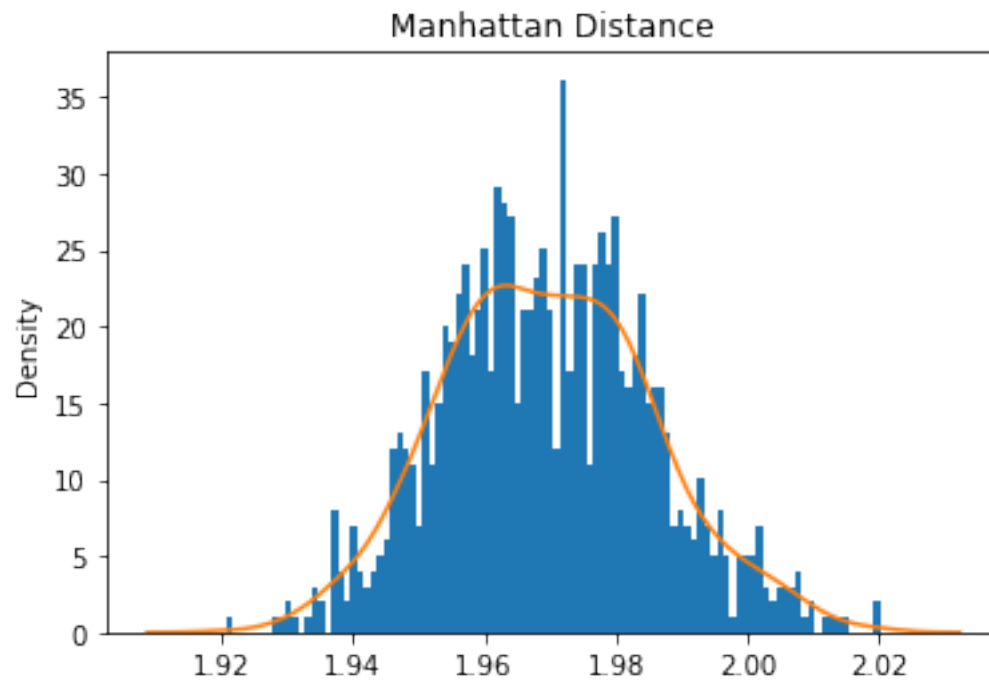
```
[12]: train_test.test_generator(generator,real_dataset,device)
```



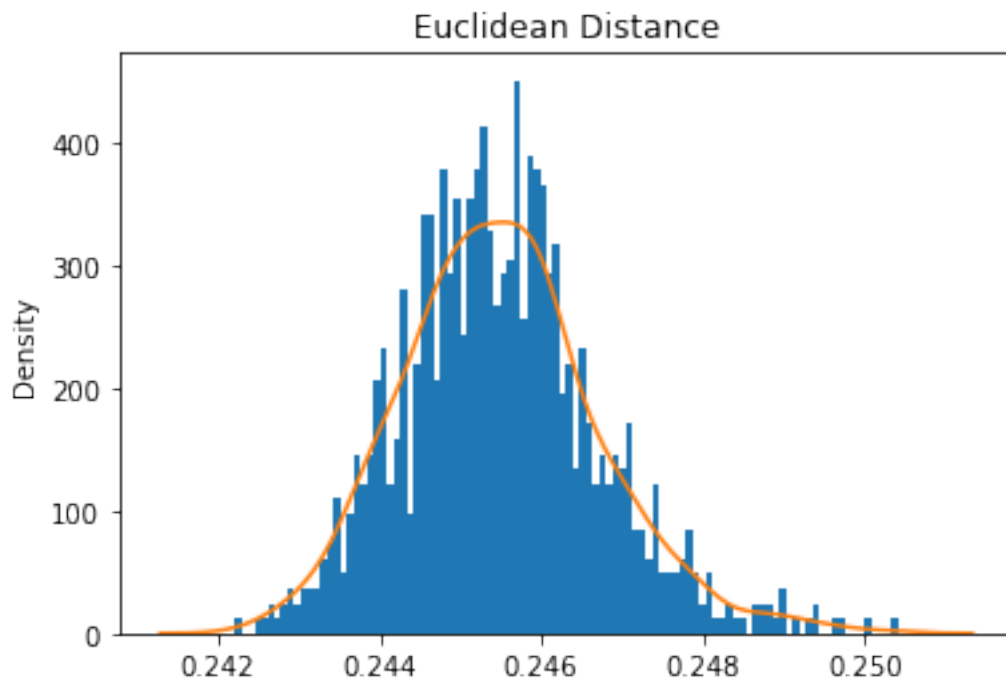
Mean Square Error: 0.0006026870181693523



Mean Absolute Error: 0.019696912256106733

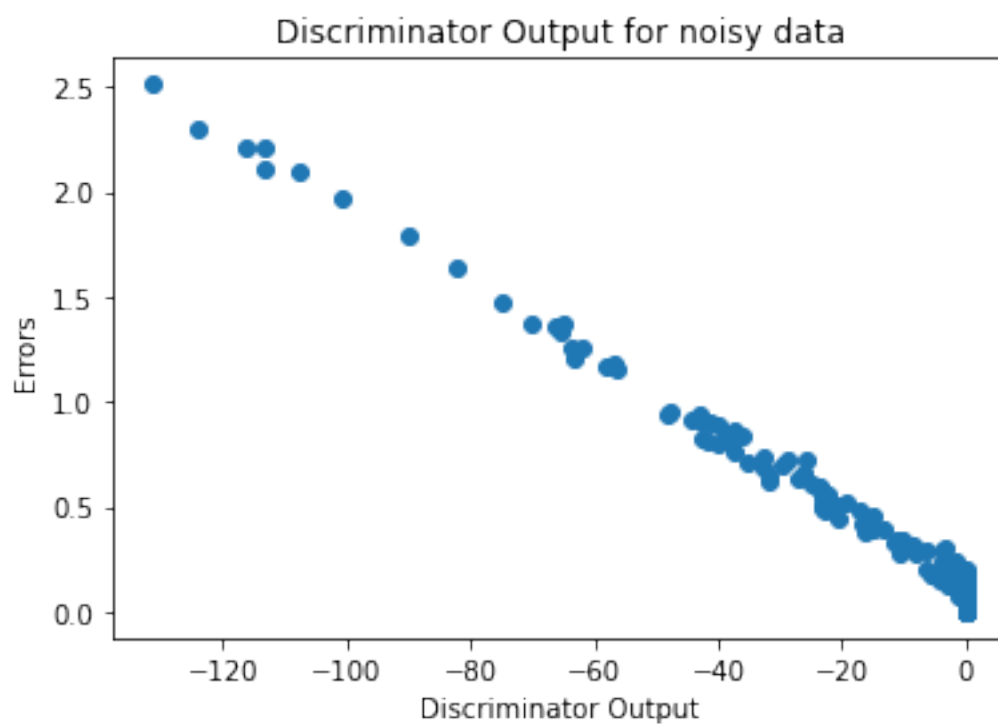
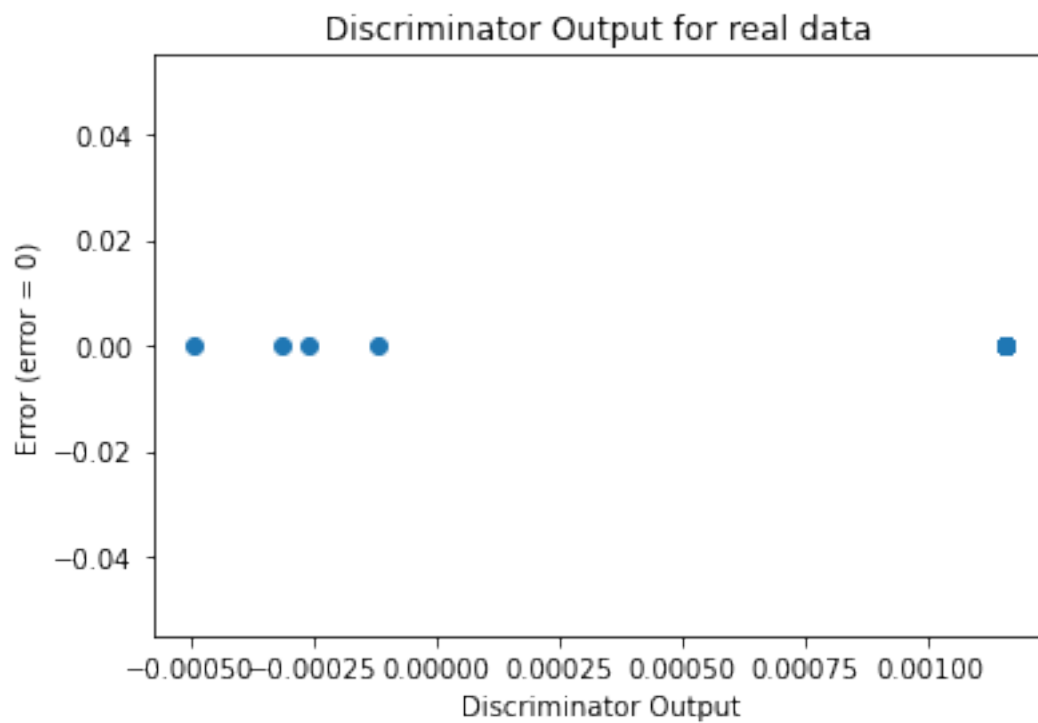


Mean Manhattan Distance: 1.9696912256106733



Mean Euclidean Distance: 0.24549382893681101

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

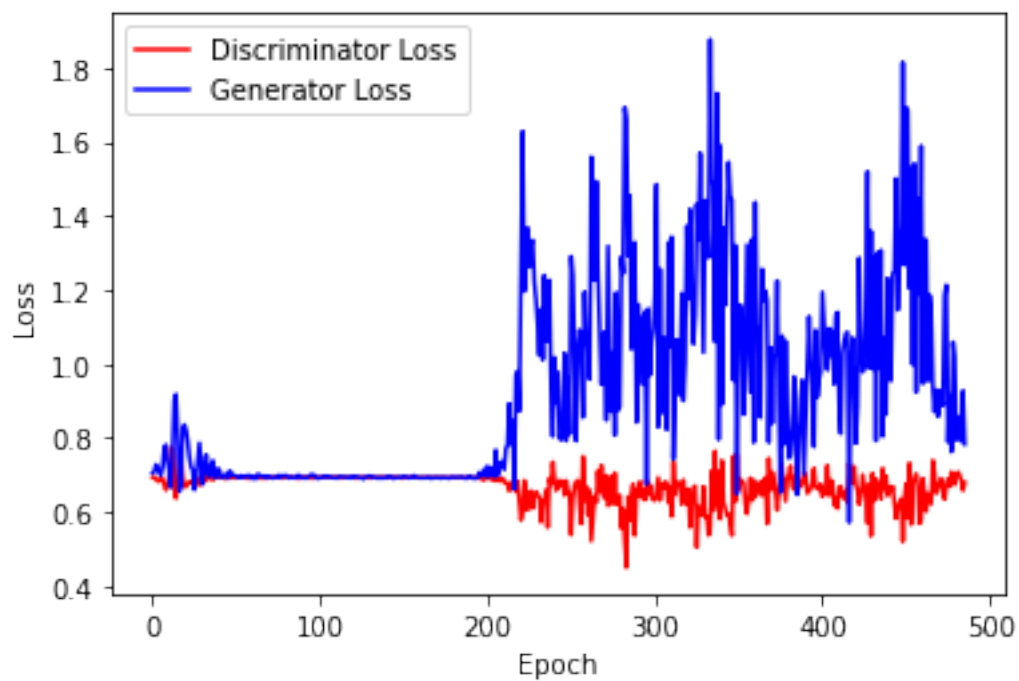



Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

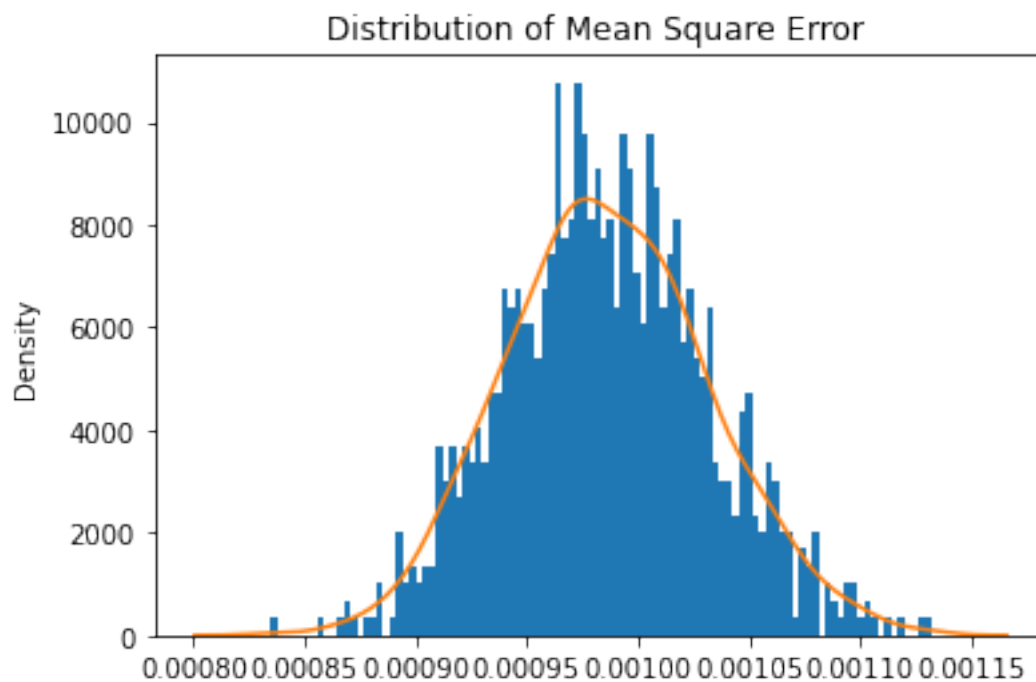
```
[14]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator.parameters(), lr=0.01, betas=(0.5,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

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

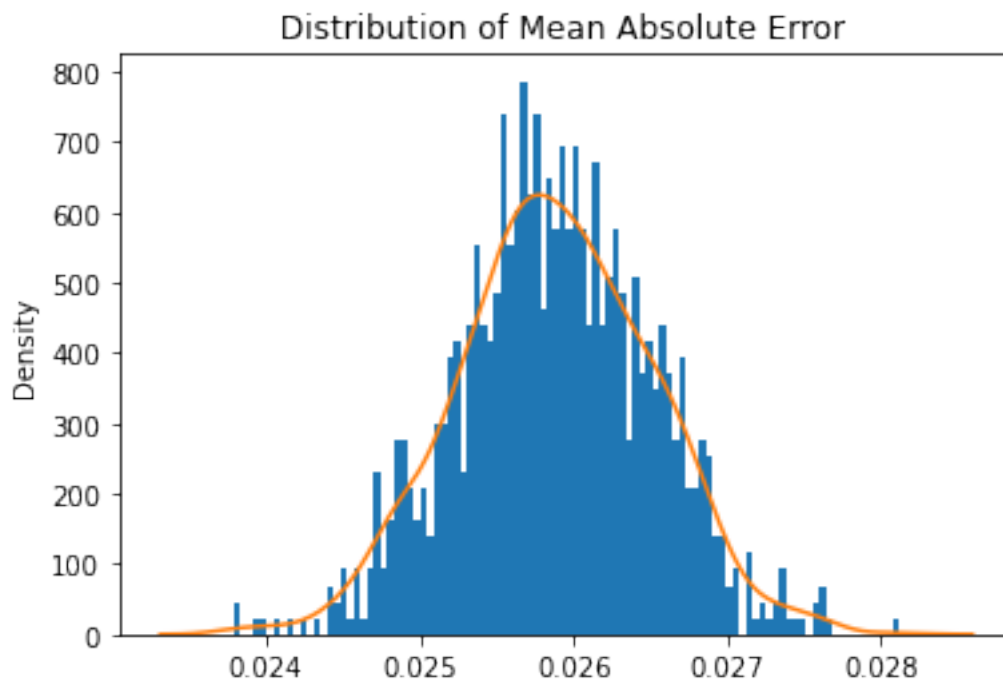
Number of epochs needed 243



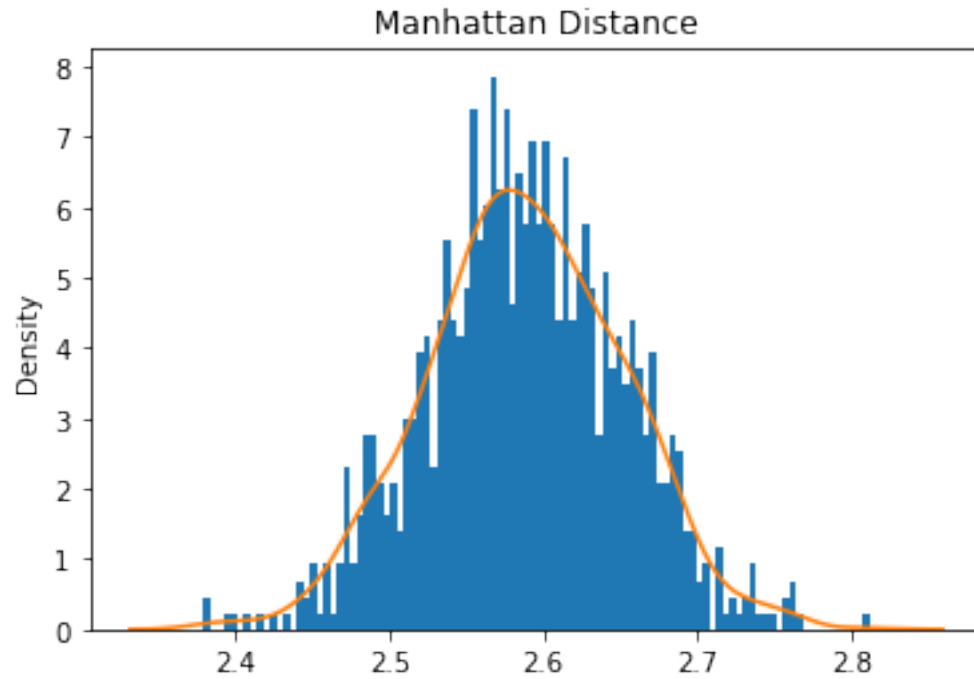
```
[16]: train_test.test_generator(generator,real_dataset,device)
```



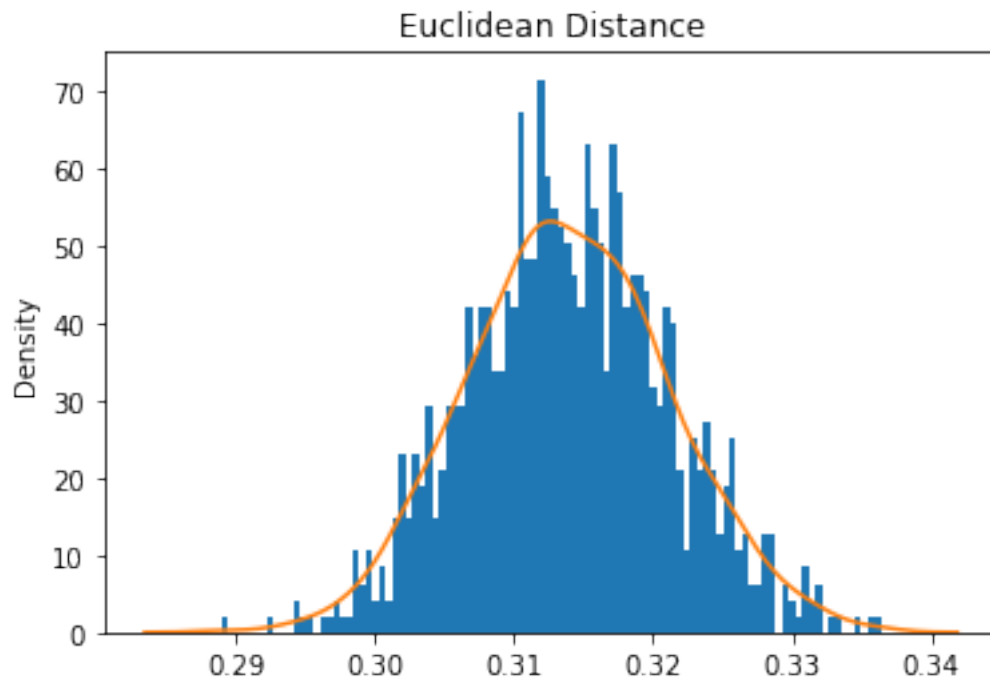
Mean Square Error: 0.0009861241193612274



Mean Absolute Error: 0.025877772430852056



Mean Manhattan Distance: 2.5877772430852057



Mean Euclidean Distance: 0.3139422989068535

2 ABC GAN Model

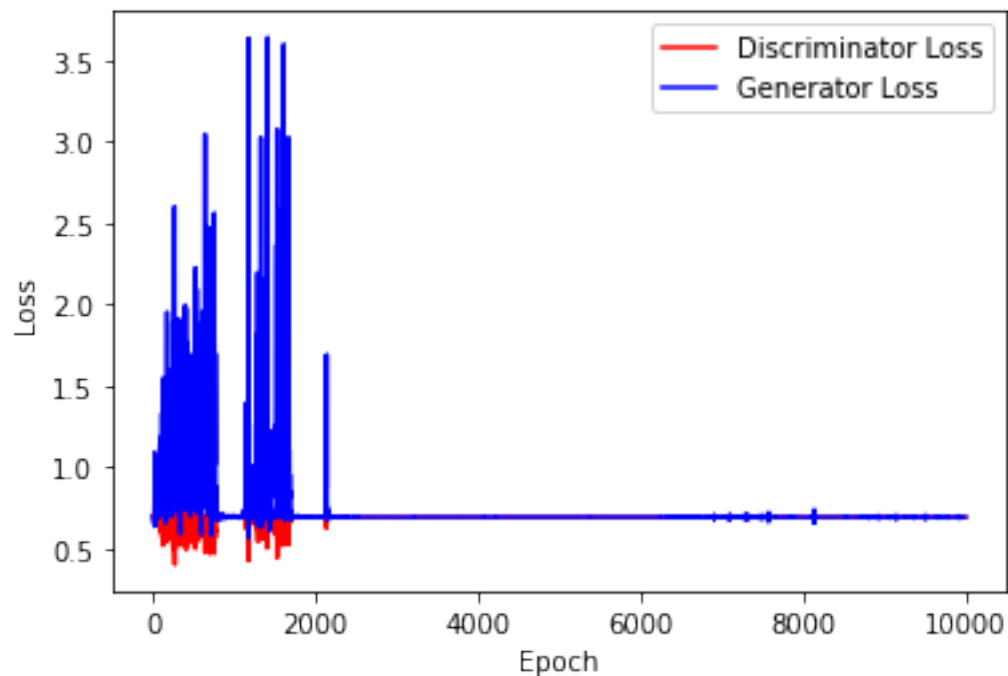
2.0.1 Training the network

Training ABC-GAN for `n_epochs` number of epochs

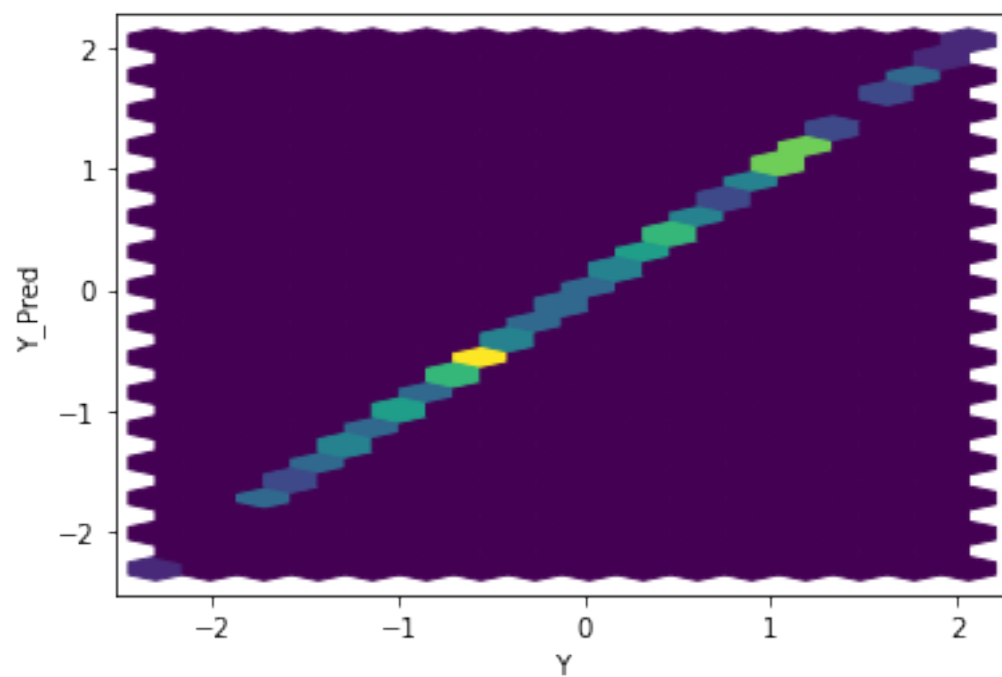
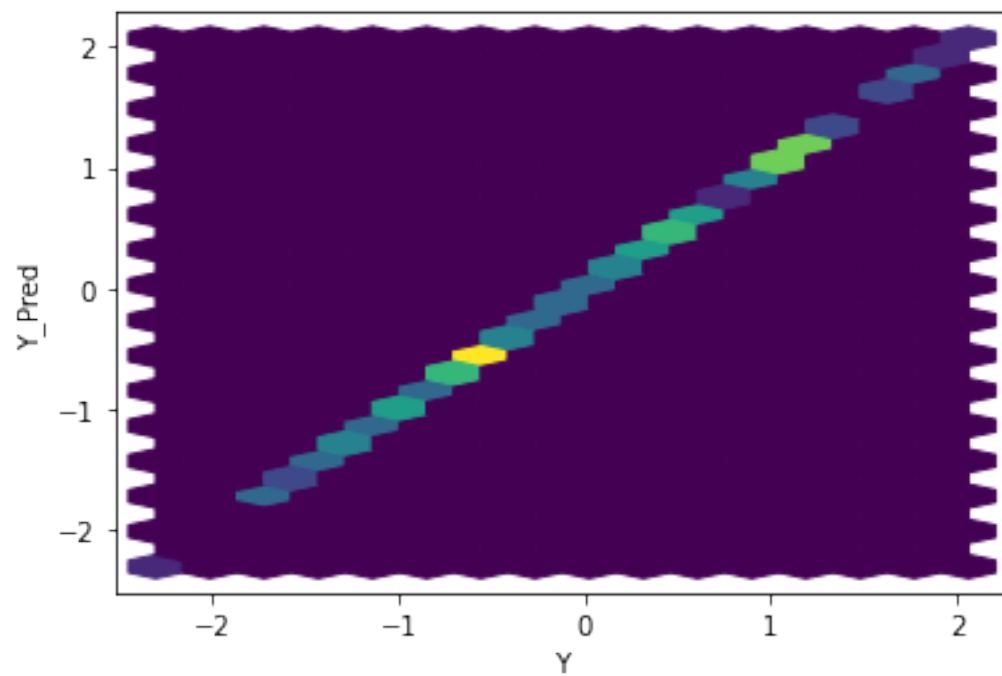
```
[17]: gen = network.Generator(n_features+2)
      disc = network.Discriminator(n_features+2)

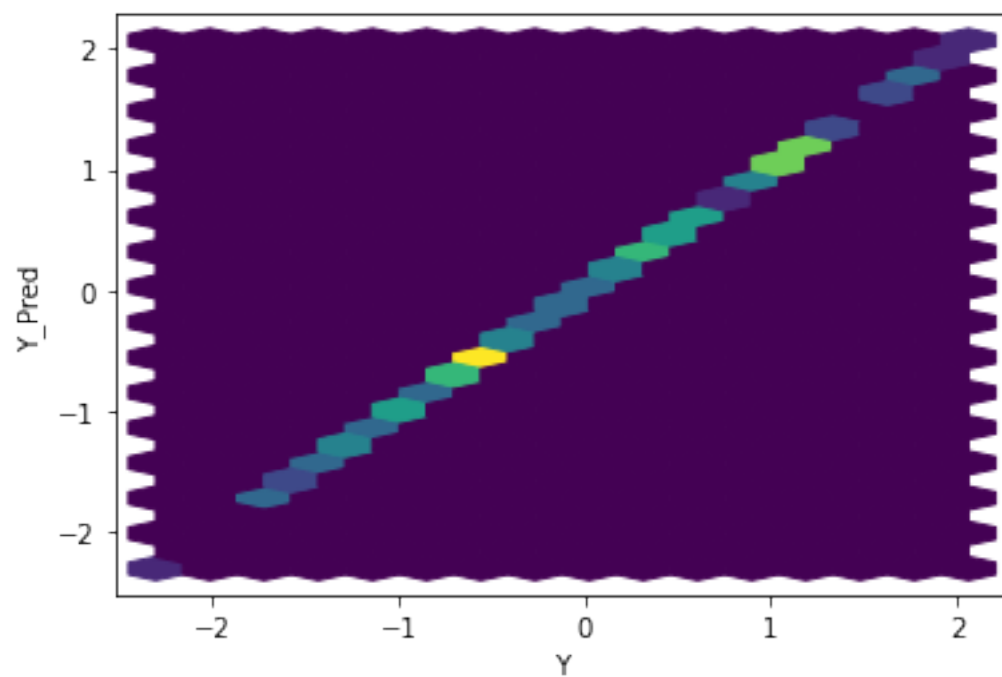
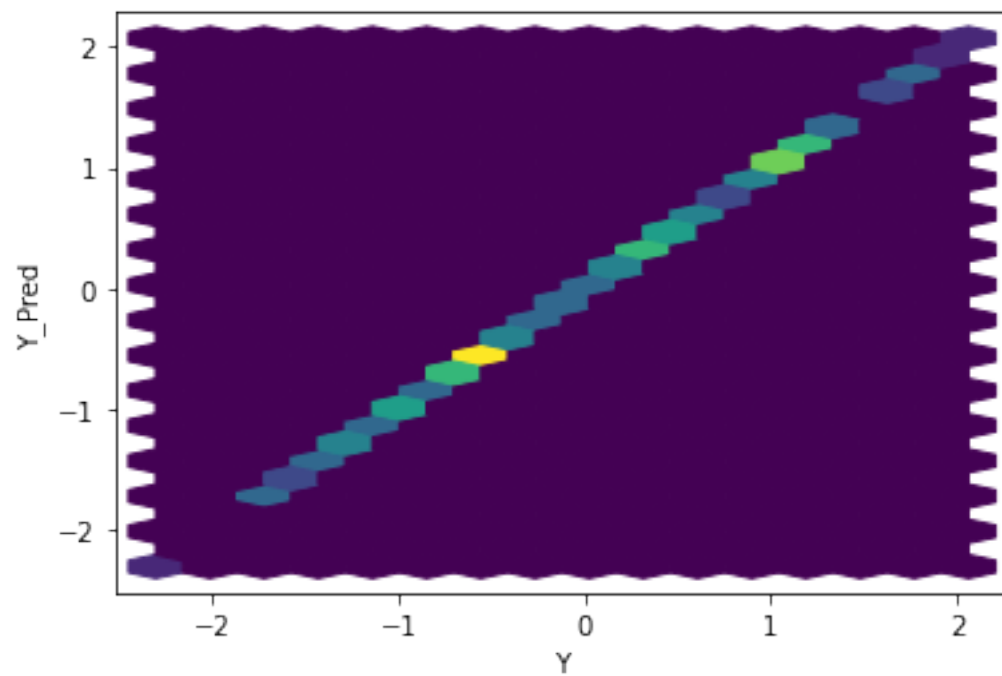
      criterion = torch.nn.BCEWithLogitsLoss()
      gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
      disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))

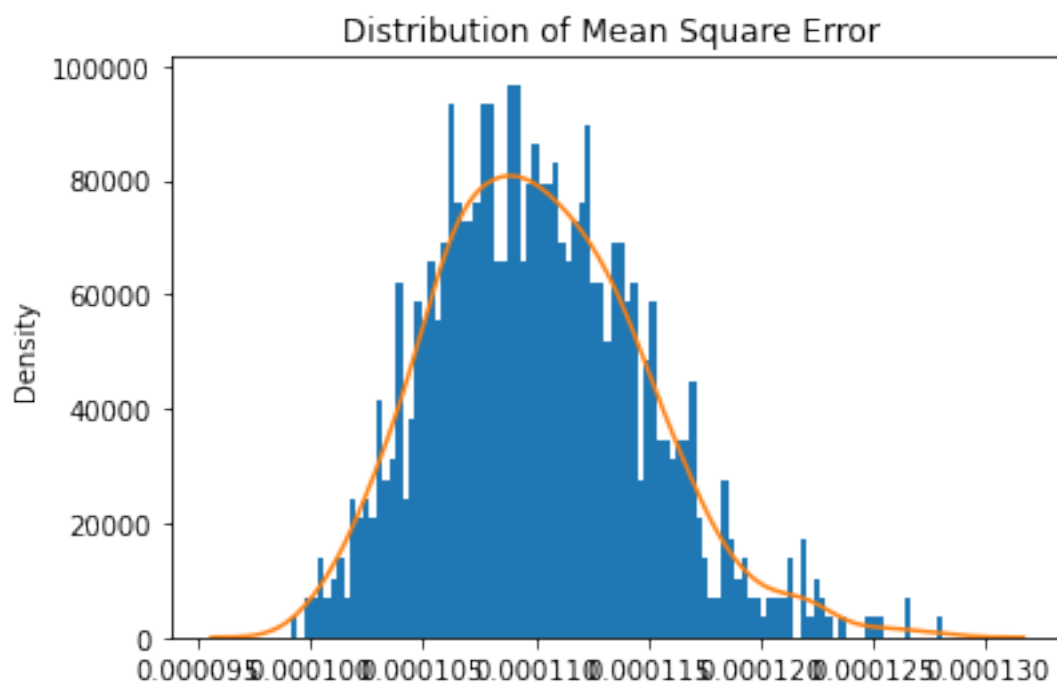
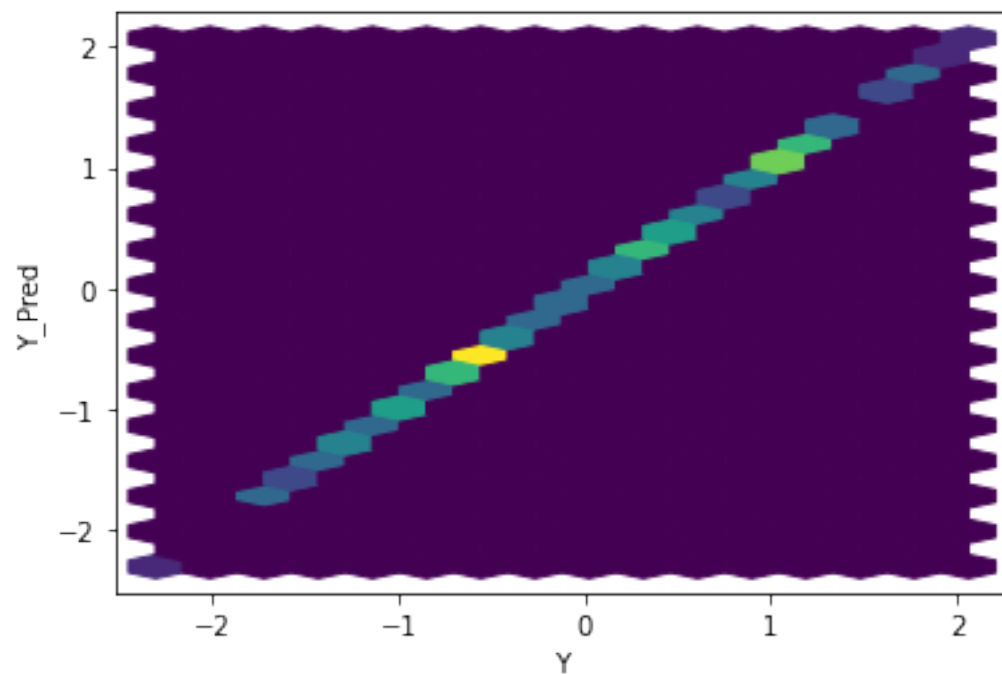
[18]: ABC_train_test.training_GAN(disc, gen,disc_opt,gen_opt,real_dataset,
      ↪batch_size, n_epochs,criterion,coeff,mean,variance,device)
```



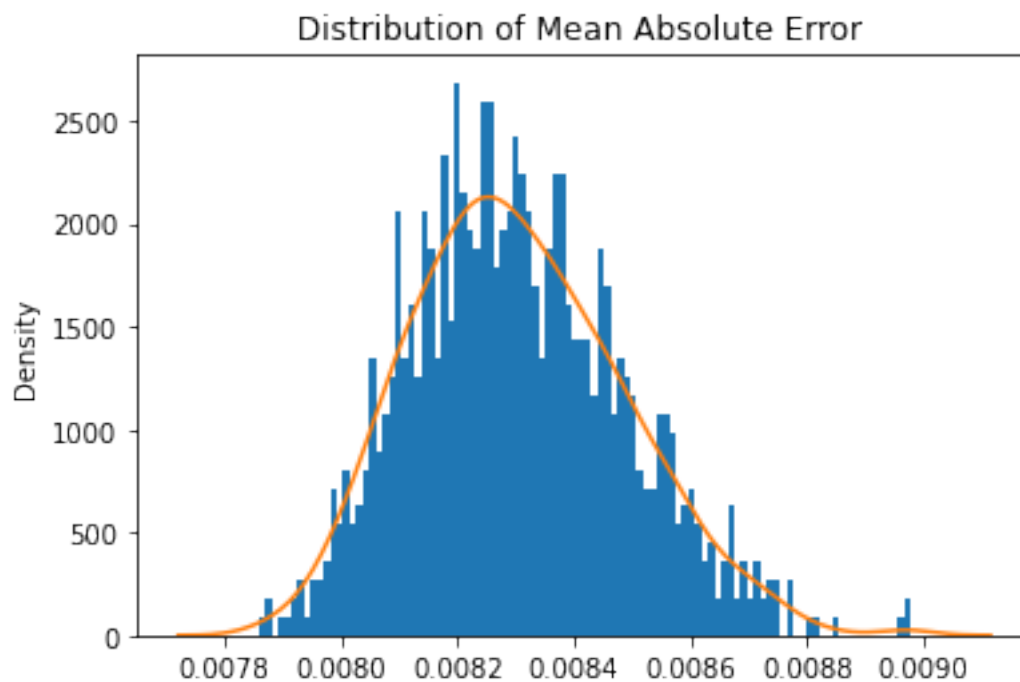
```
[19]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```





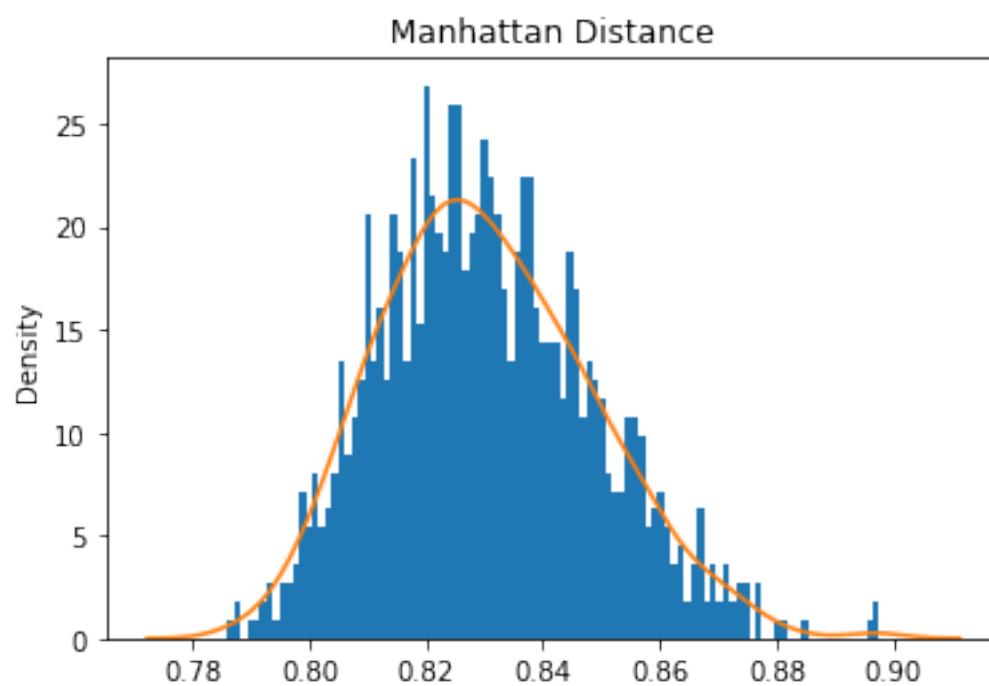


Mean Square Error: 0.00011015109824227252

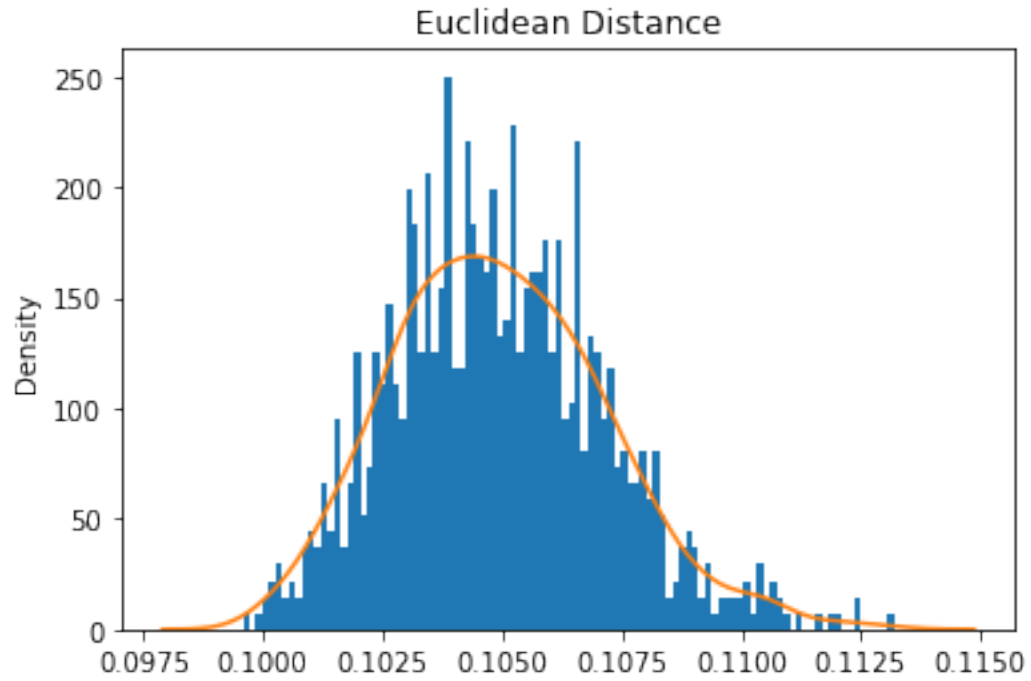


Mean Absolute Error: 0.008303945847675204

Mean Manhattan Distance: 0.8303945847675205

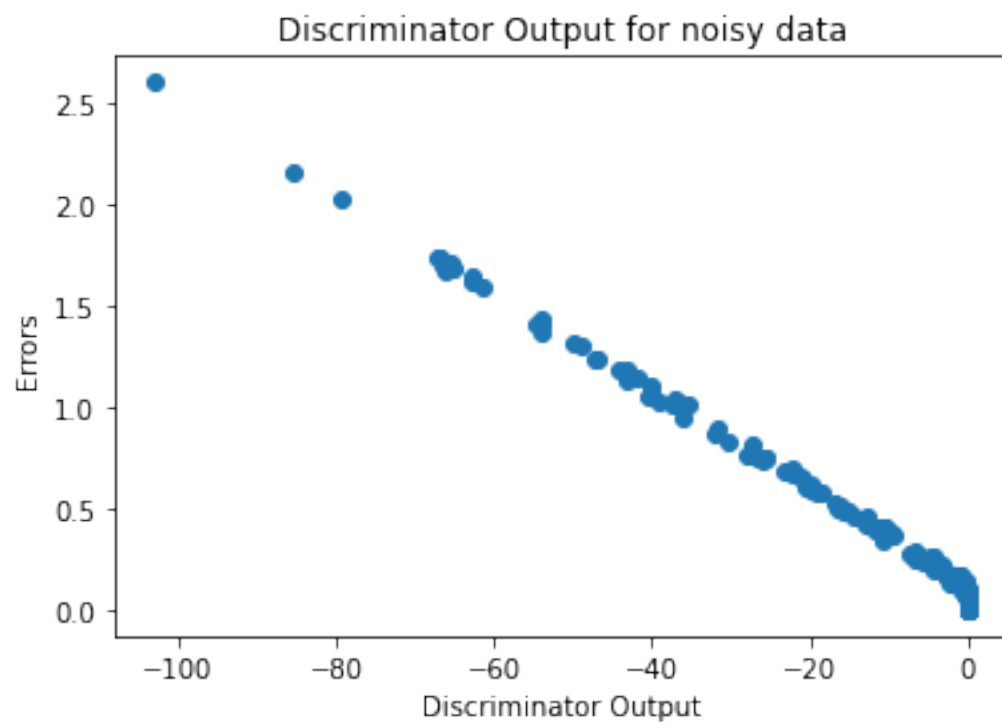
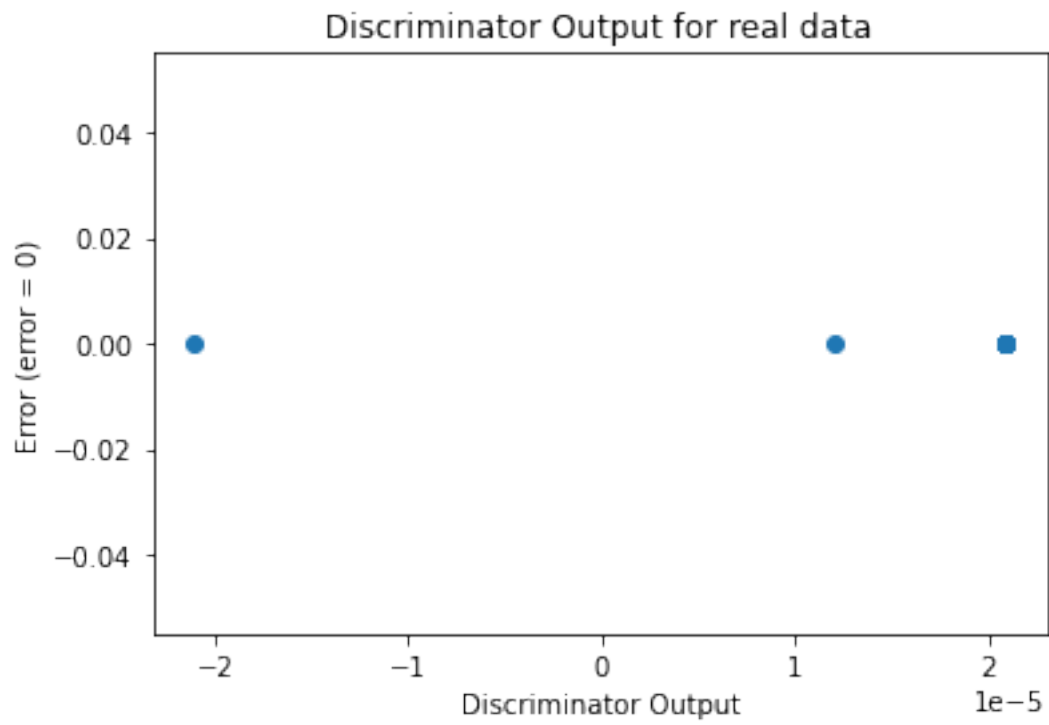


Mean Euclidean Distance: 0.10492915194815627



Sanity Checks

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



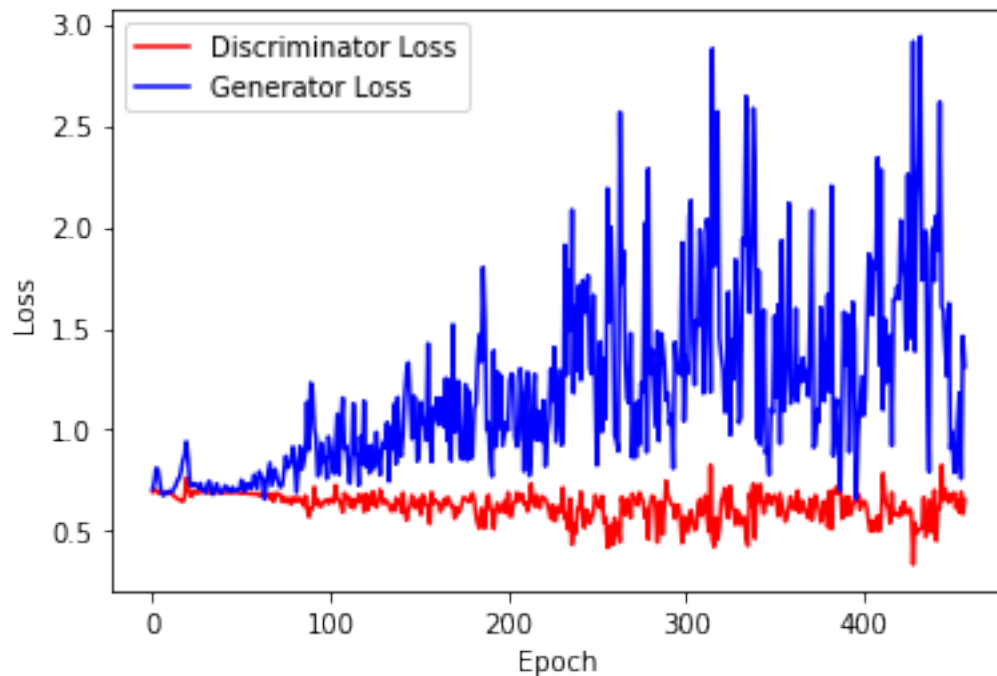
Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

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

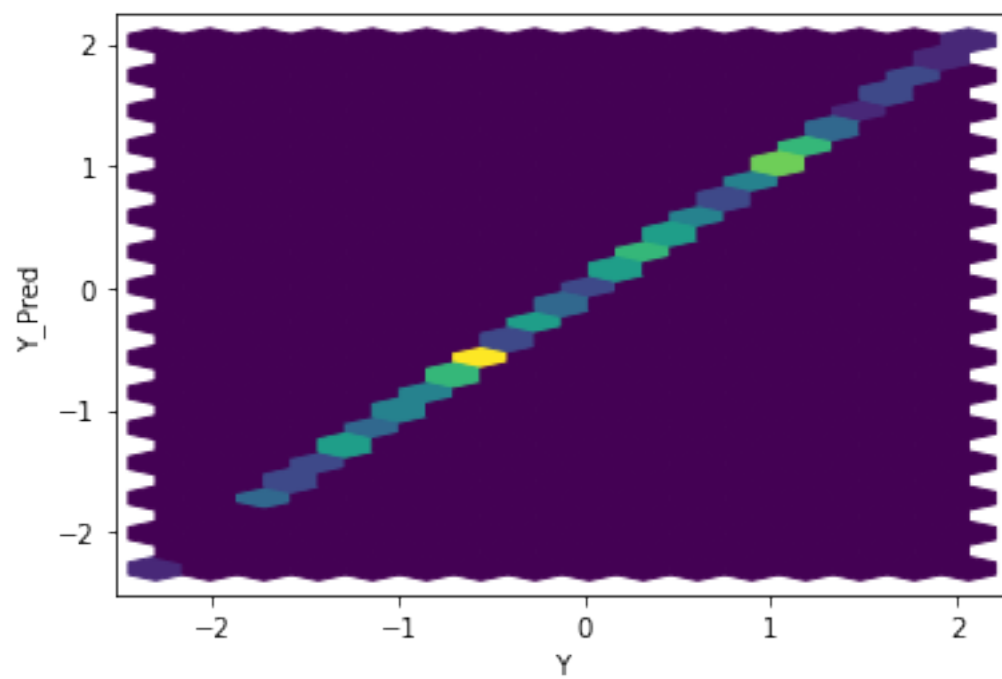
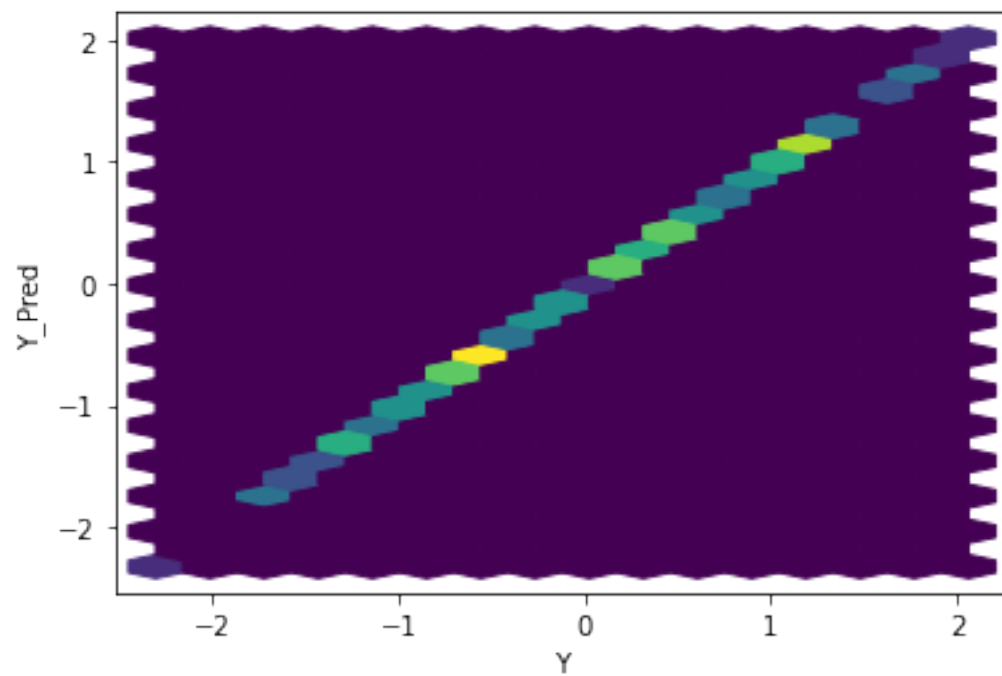
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))
```

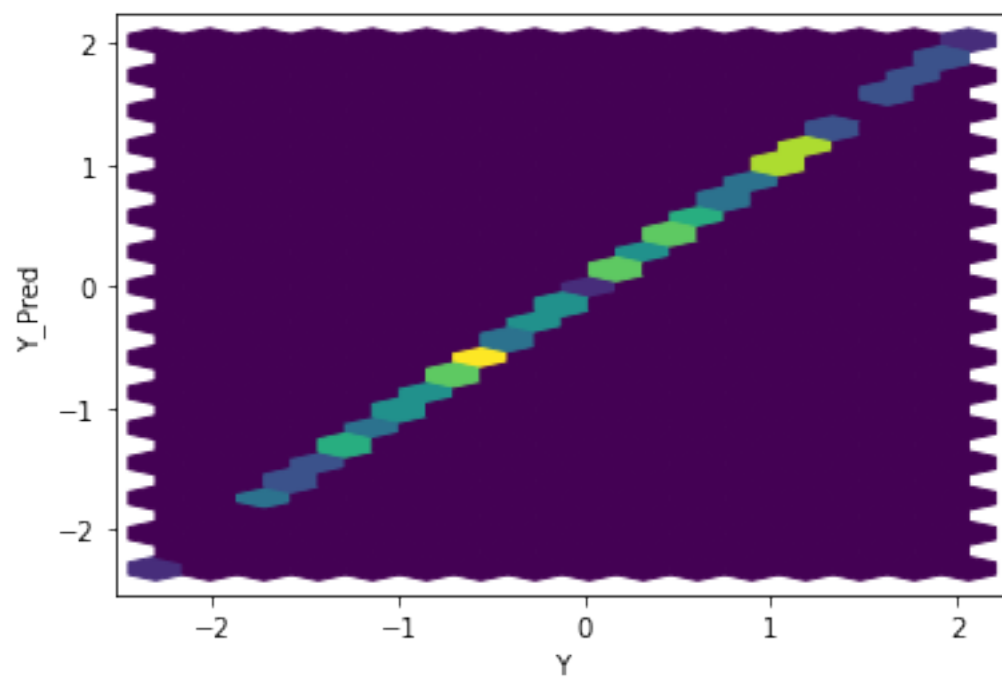
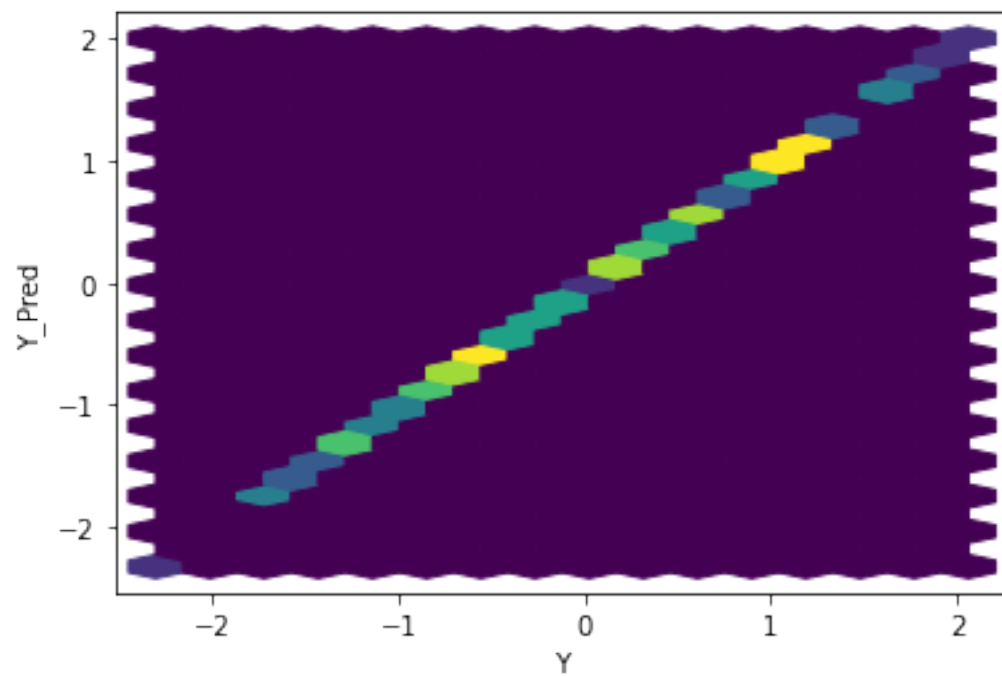
```
[22]: ABC_train_test.
      ↪ training_GAN_2(disc,gen,disc_opt,gen_opt,real_dataset,batch_size,
      ↪ error,criterion,coeff,mean,variance,device)
```

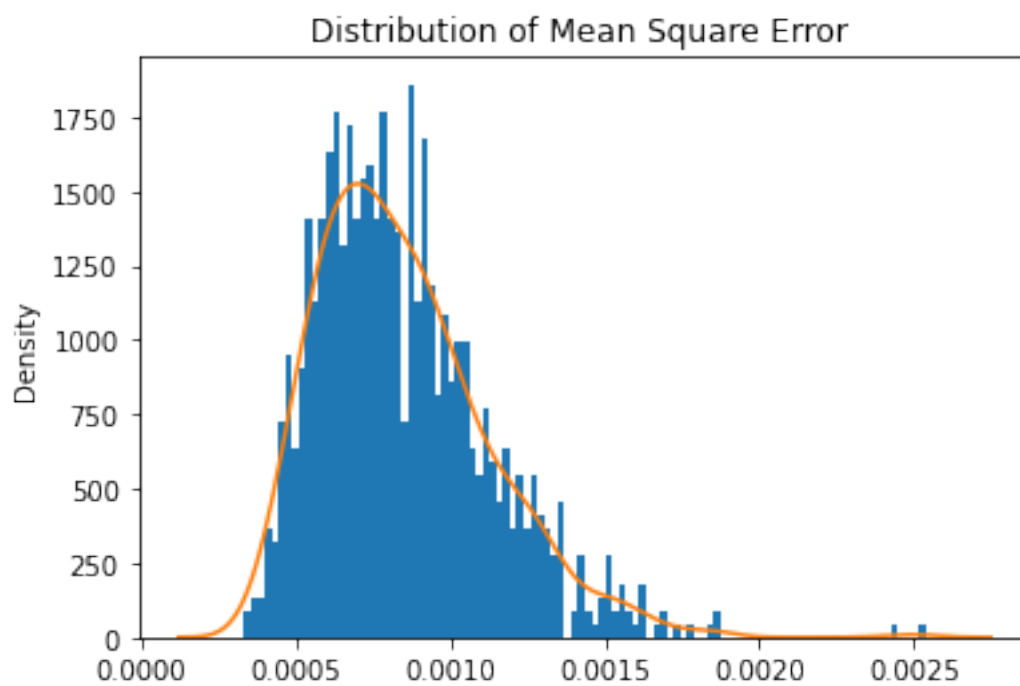
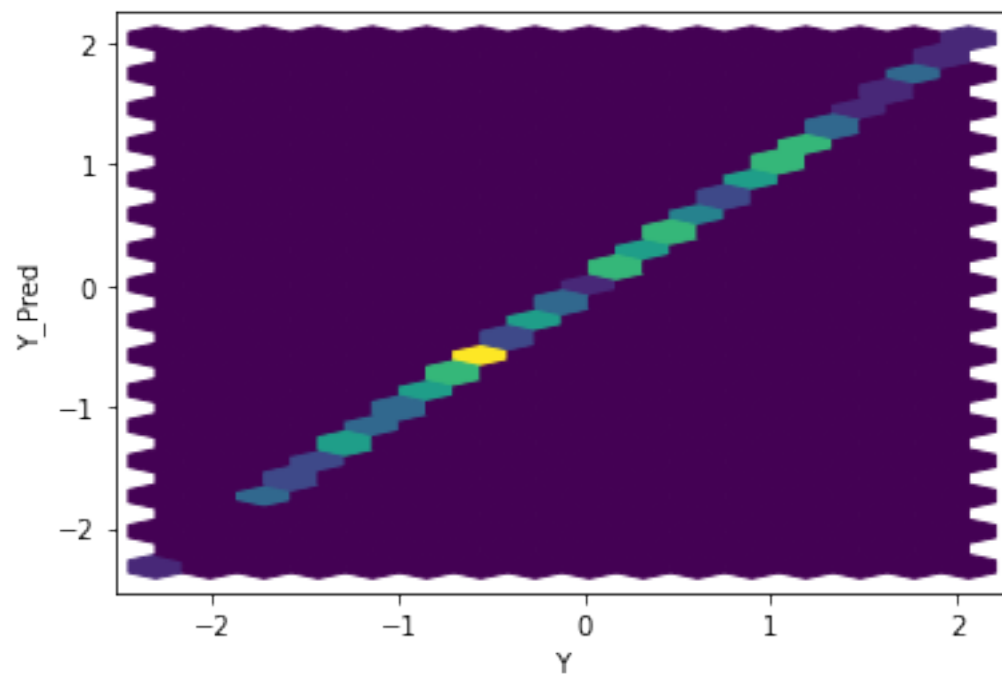
Number of epochs 229



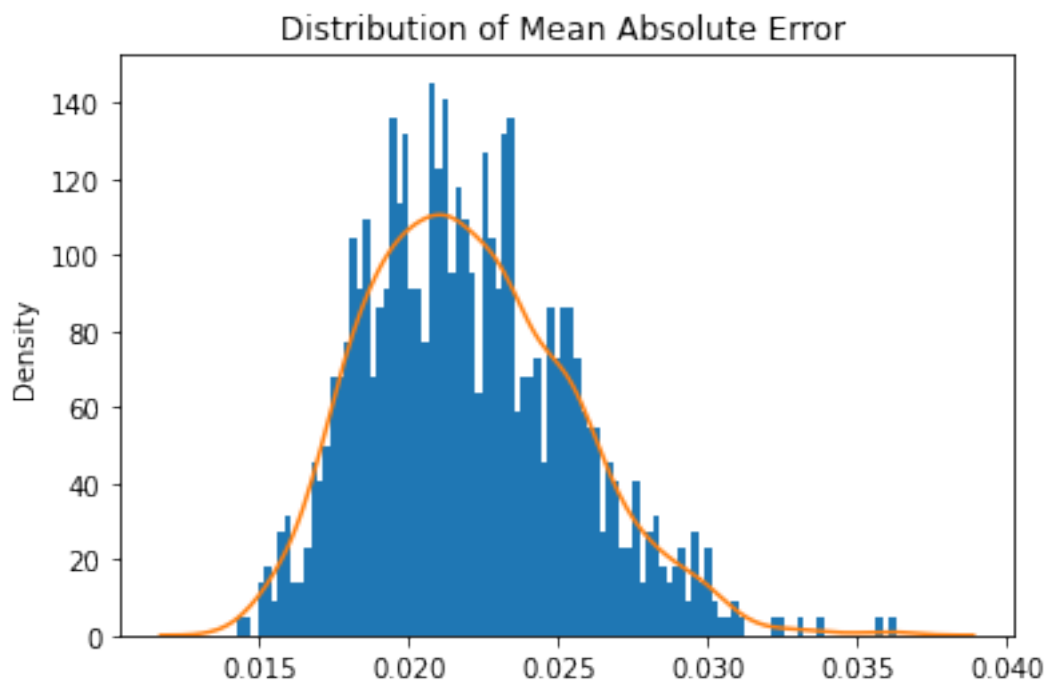
```
[23]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```





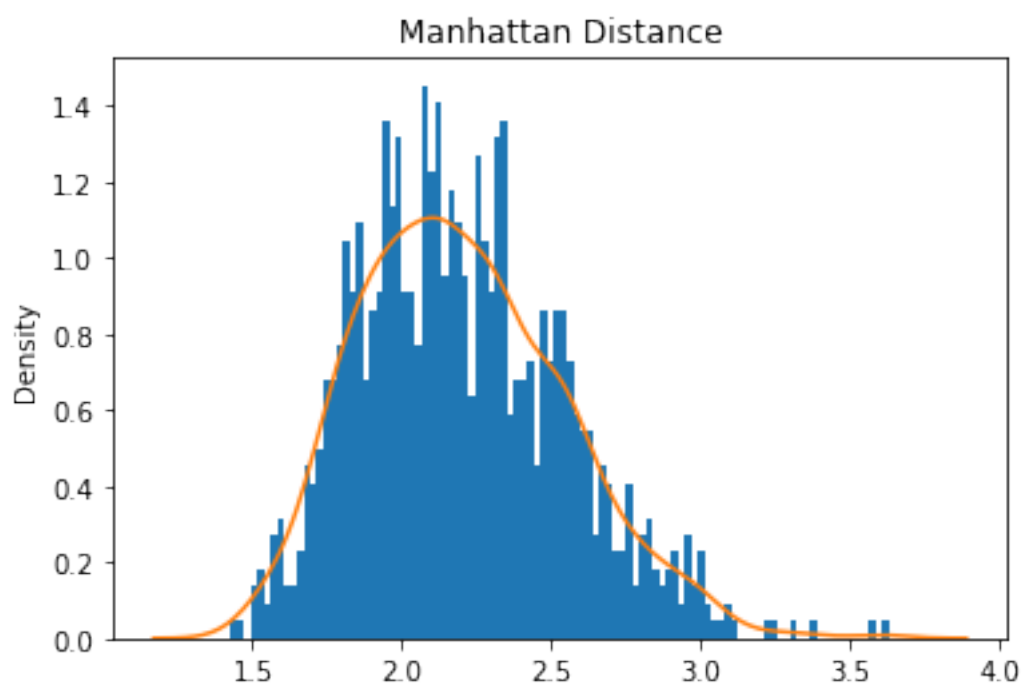


Mean Square Error: 0.0008472464905586359

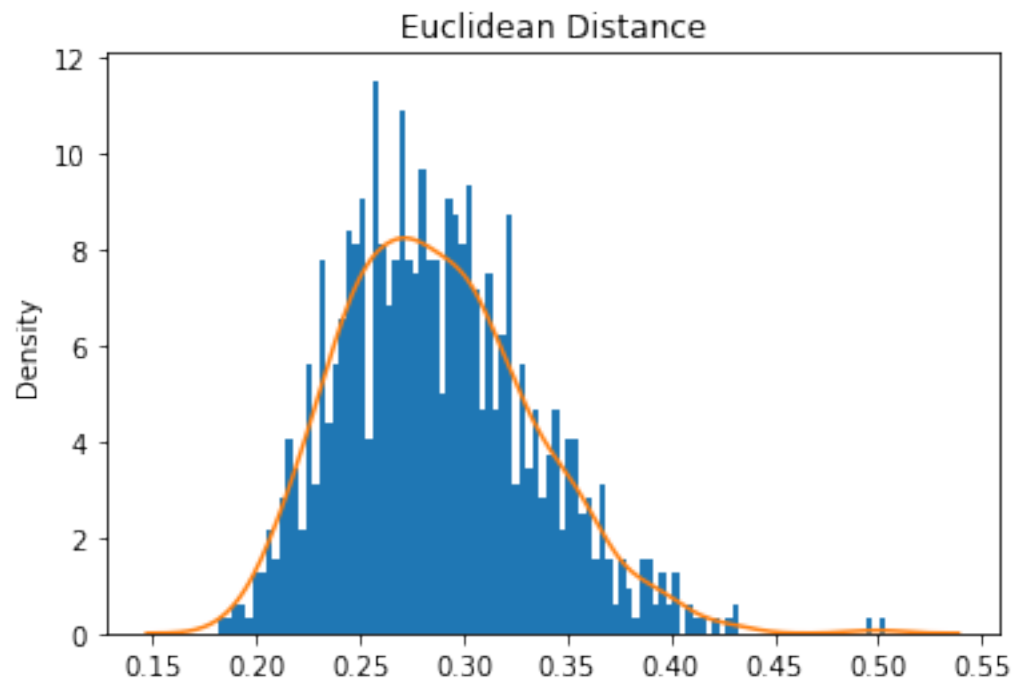


Mean Absolute Error: 0.022022543726041913

Mean Manhattan Distance: 2.2022543726041914



Mean Euclidean Distance: 0.28732539150062336



[]: