

Dataset1-Regression_output_10

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 = 0
     variance = 0.1

```

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	-0.772485	-1.420635	0.762644	-0.364055	-1.600409	1.005007	0.463959	
1	1.364214	0.383486	0.382380	0.613427	0.800188	-1.425757	-0.716276	
2	0.264398	0.395630	0.400624	-1.484984	-2.252866	1.228557	1.768419	
3	1.532085	0.021507	0.156135	1.159594	-0.652339	2.340618	0.354511	
4	-0.399384	-1.342096	0.302469	-0.252437	0.280556	-0.823839	0.382775	
	X8	X9	X10	Y				

```

0  0.303395 -0.289320  1.625926 -111.256934
1  0.254746  0.122893 -0.375292  206.610811
2 -1.192514  0.271815 -0.078237 -162.165728
3  0.376661  0.018858 -2.335494   56.044218
4 -0.153399 -1.431796  0.958126 -178.029006

```

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:                2.445e+07
Date:                Tue, 19 Oct 2021  Prob (F-statistic):        6.02e-282
Time:                        23:27:41  Log-Likelihood:             599.40
No. Observations:                100      AIC:                    -1177.
Df Residuals:                    89      BIC:                    -1148.
Df Model:                        10
Covariance Type:                nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	-1.279e-17	6.4e-05	-2e-13	1.000	-0.000	0.000
x1	0.2696	6.8e-05	3964.672	0.000	0.269	0.270
x2	0.5433	6.77e-05	8024.837	0.000	0.543	0.543
x3	0.5160	6.71e-05	7692.194	0.000	0.516	0.516
x4	0.3477	6.9e-05	5041.039	0.000	0.348	0.348
x5	0.3554	6.79e-05	5234.487	0.000	0.355	0.356
x6	0.0486	6.79e-05	715.716	0.000	0.048	0.049
x7	0.0177	6.54e-05	270.718	0.000	0.018	0.018
x8	0.3229	6.67e-05	4839.480	0.000	0.323	0.323
x9	0.4774	6.76e-05	7058.155	0.000	0.477	0.478
x10	0.2676	6.78e-05	3949.716	0.000	0.268	0.268

```

=====
Omnibus:                        1.154  Durbin-Watson:                1.810
Prob(Omnibus):                  0.562  Jarque-Bera (JB):                1.057
Skew:                          -0.064  Prob(JB):                        0.590
Kurtosis:                      2.513  Cond. No.:                      1.73
=====

```

Notes:

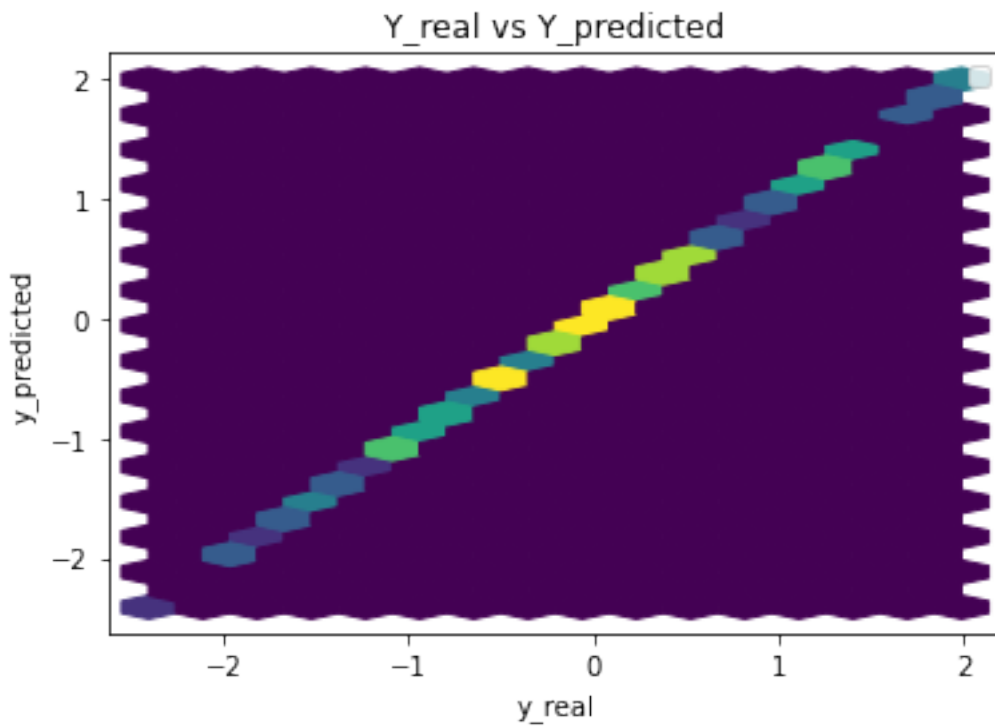
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```

Parameters:  const  -1.279359e-17
x1          2.696219e-01

```

```
x2      5.432785e-01
x3      5.159818e-01
x4      3.477148e-01
x5      3.554402e-01
x6      4.856773e-02
x7      1.769769e-02
x8      3.228900e-01
x9      4.773986e-01
x10     2.676349e-01
dtype: float64
```



Performance Metrics

```
Mean Squared Error: 3.6405507389605817e-07
Mean Absolute Error: 0.0004860682224953361
Manhattan distance: 0.04860682224953361
Euclidean distance: 0.0060336976548055355
```

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

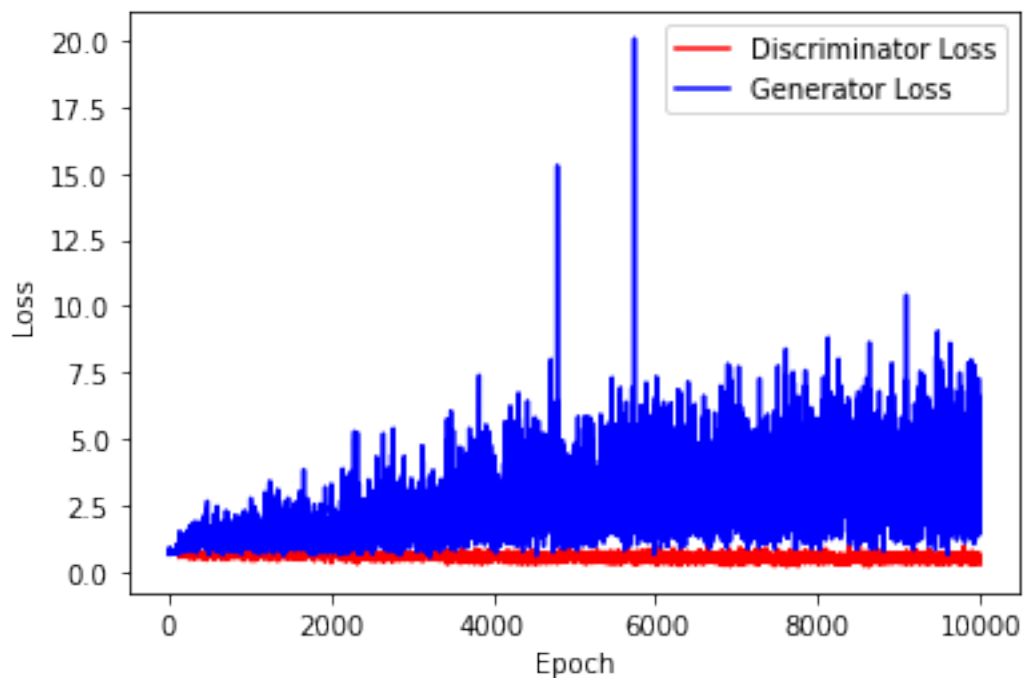
```
[9]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)

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))
```

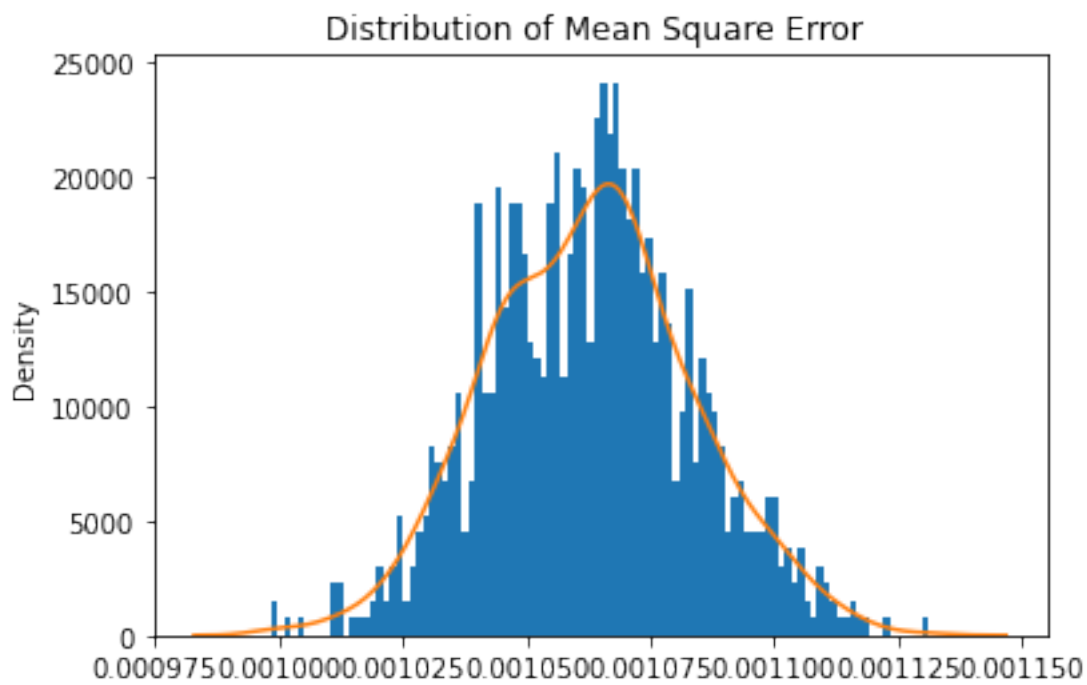
```
[10]: print(generator)
print(discriminator)
```

```
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()
)
```

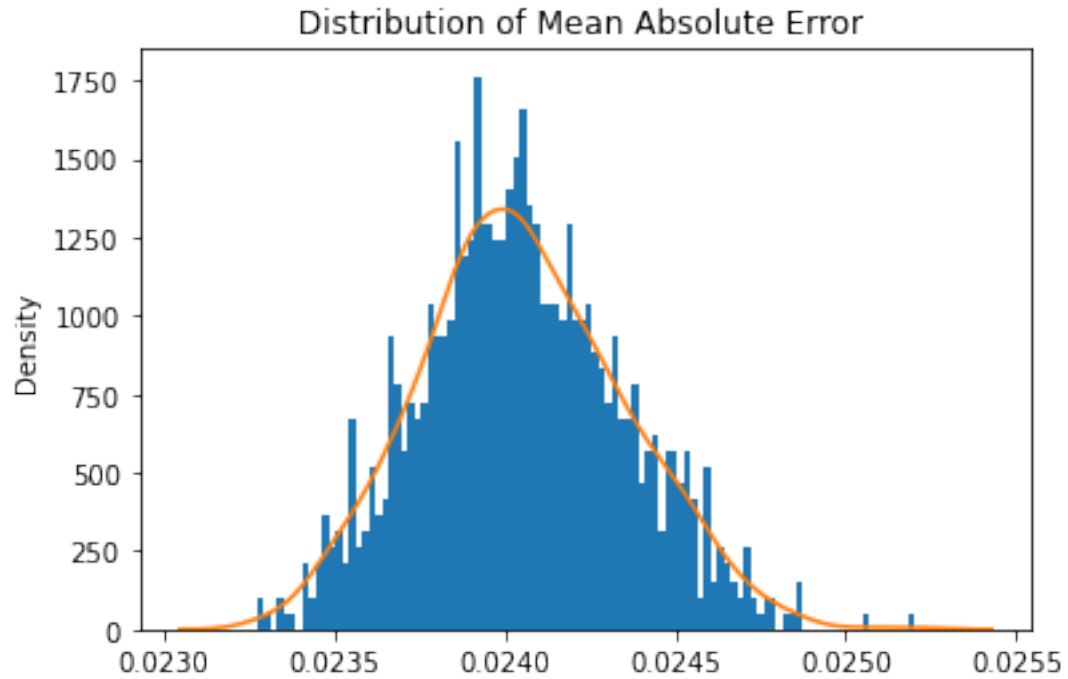
```
[11]: train_test.
→training_GAN(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,
→n_epochs,criterion,device)
```



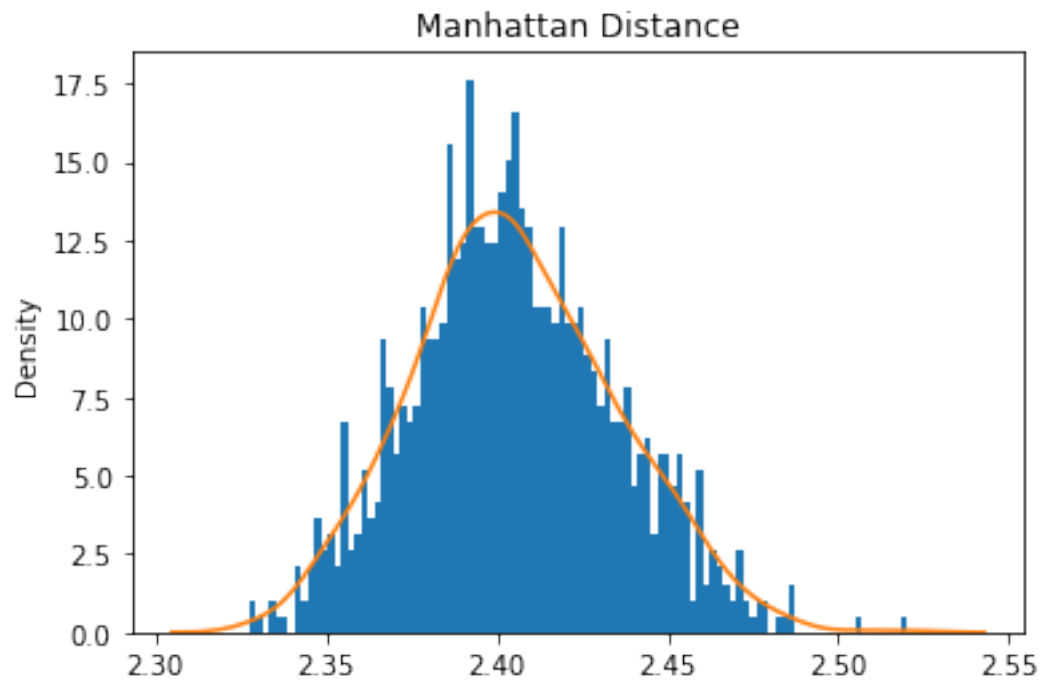
```
[12]: train_test.test_generator(generator,real_dataset,device)
```



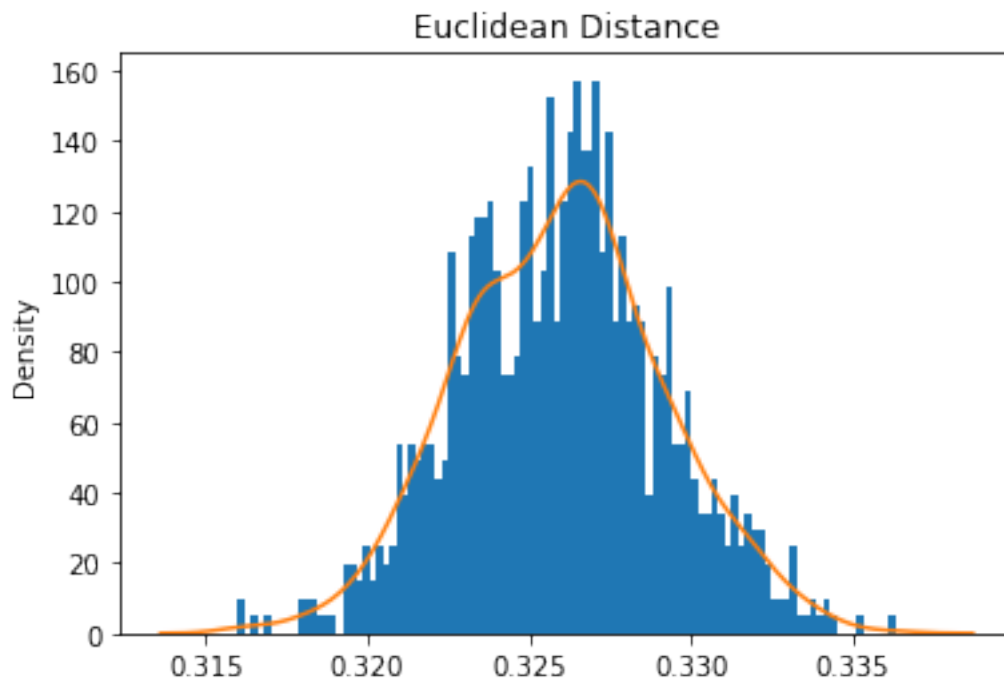
Mean Square Error: 0.001062809431385779



Mean Absolute Error: 0.02405153436873108

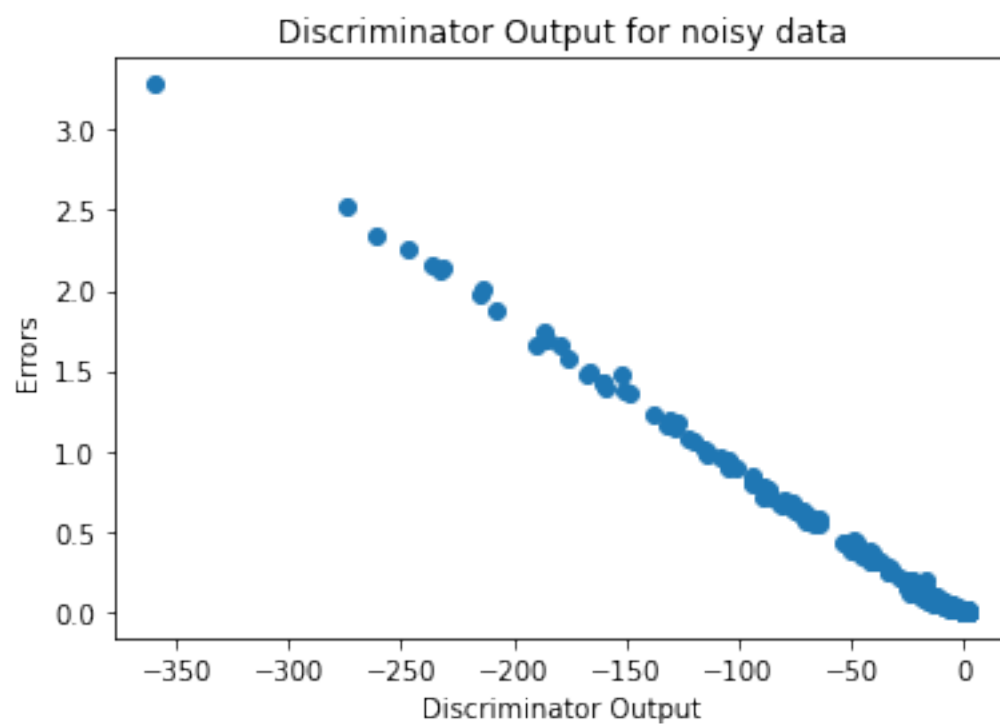
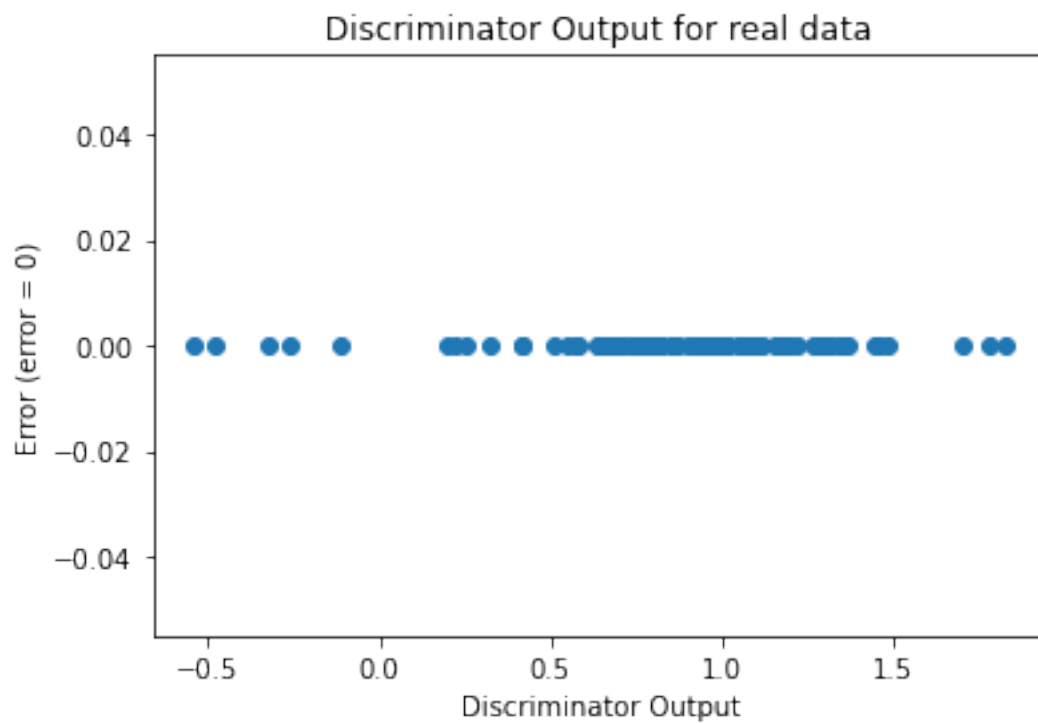


Mean Manhattan Distance: 2.4051534368731082



Mean Euclidean Distance: 0.3259921643066699

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

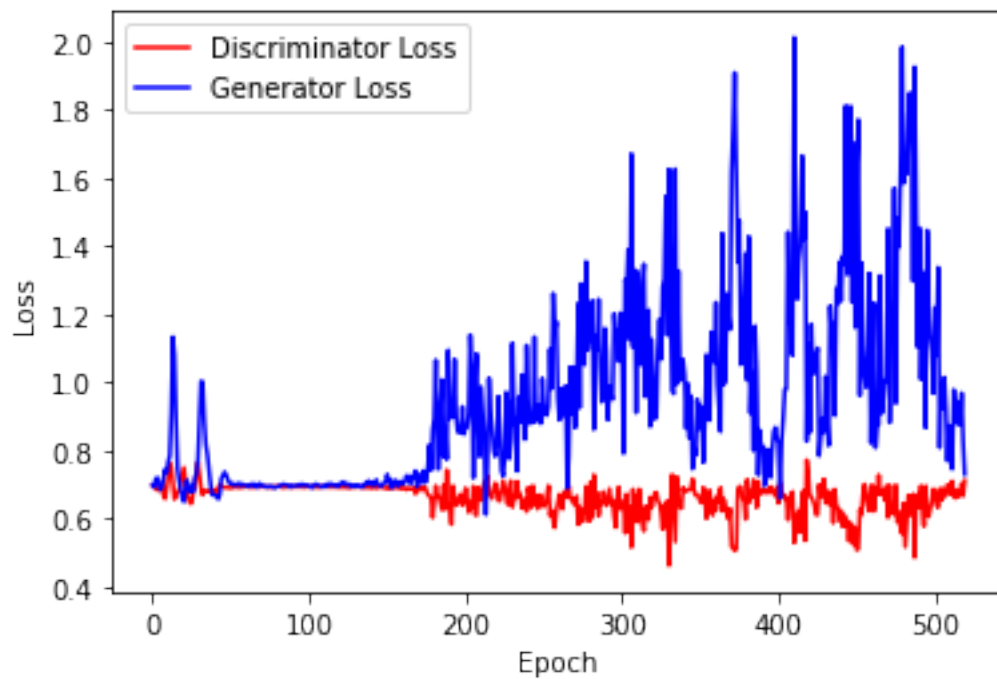



Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

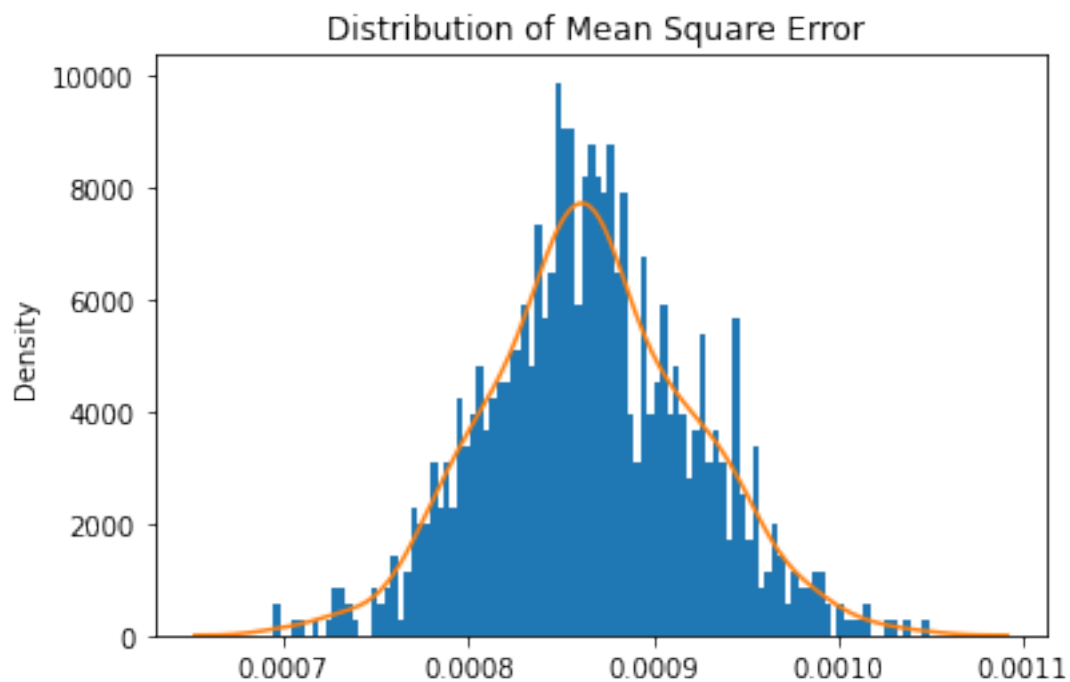
```
[14]: generator = network.Generator(n_features+2)
discriminator = network.Discriminator(n_features+2)
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(generator.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(discriminator.parameters(), lr=0.01, betas=(0.5,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

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

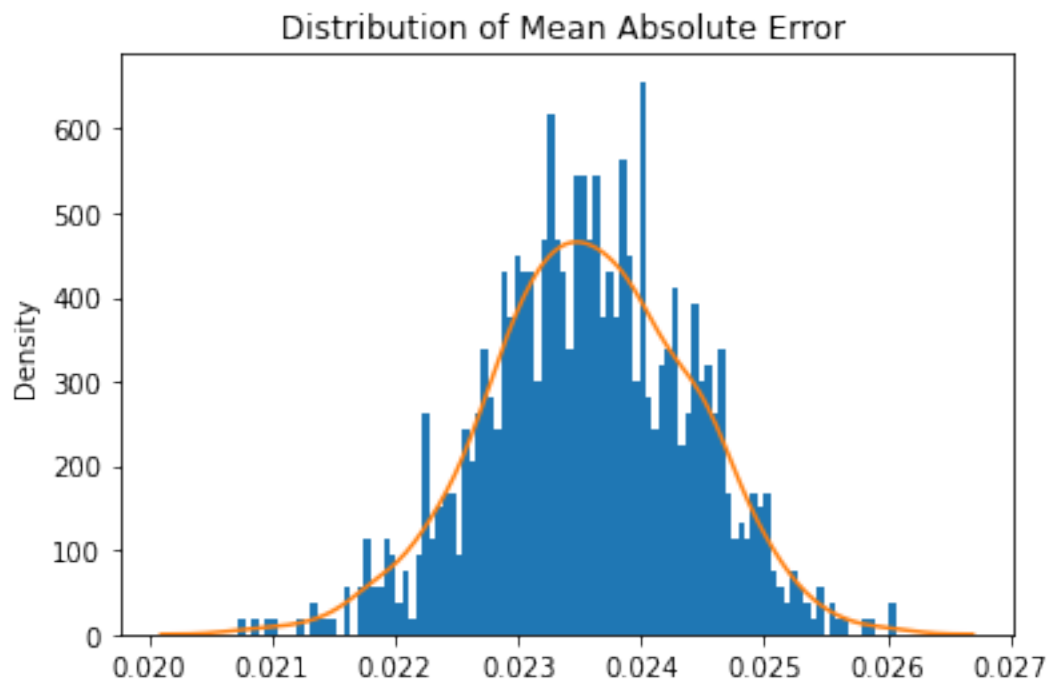
Number of epochs needed 260



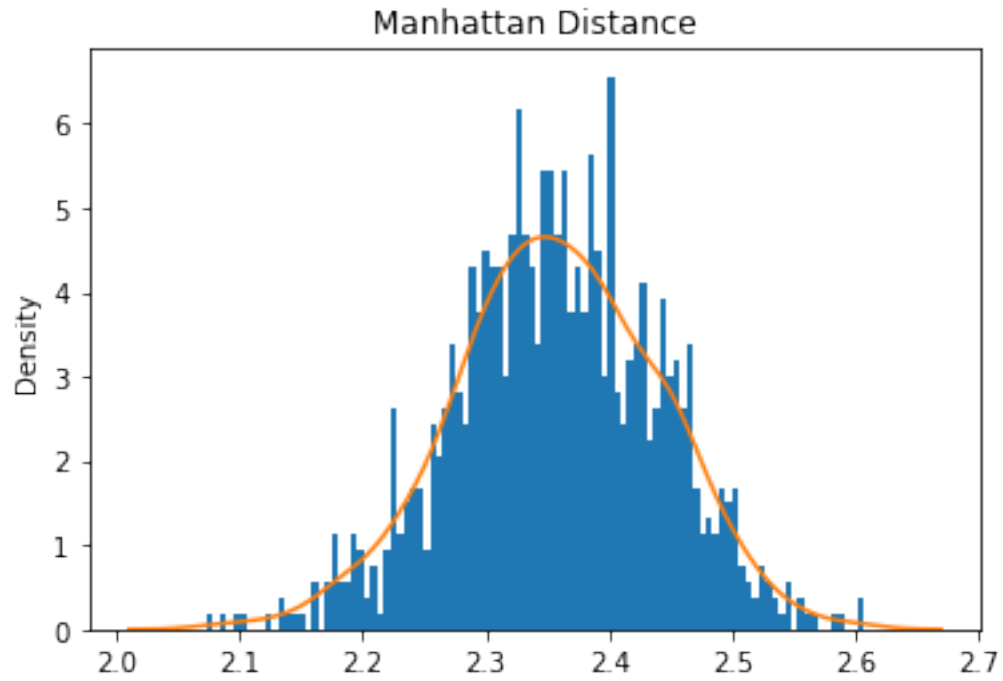
```
[16]: train_test.test_generator(generator,real_dataset,device)
```



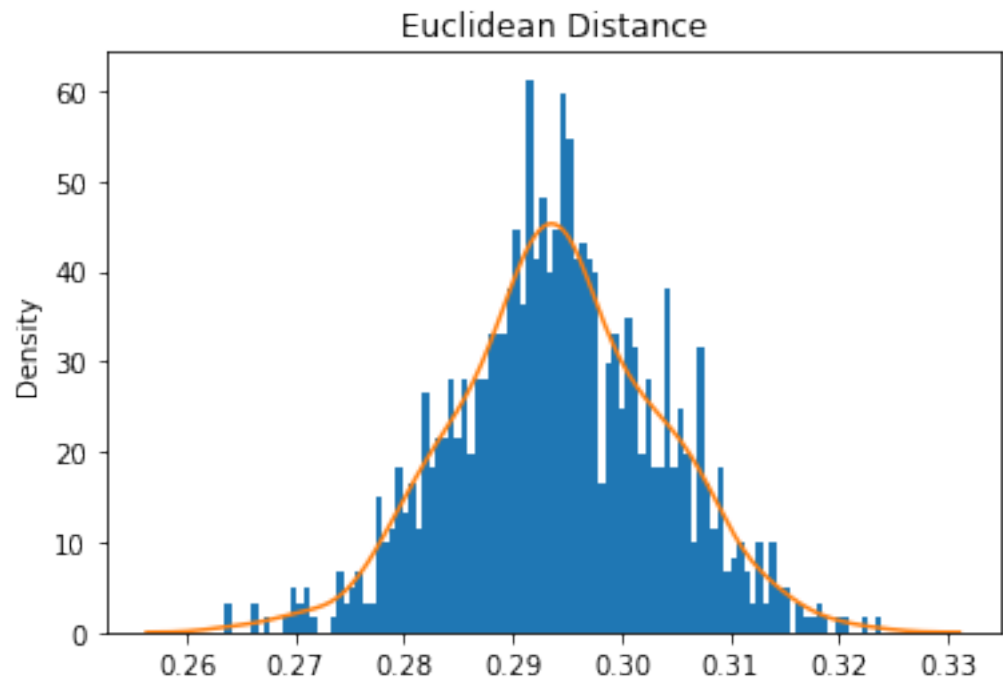
Mean Square Error: 0.000866021471304425



Mean Absolute Error: 0.023569348094947637



Mean Manhattan Distance: 2.3569348094947635



Mean Euclidean Distance: 0.29412822639055236

2 ABC GAN Model

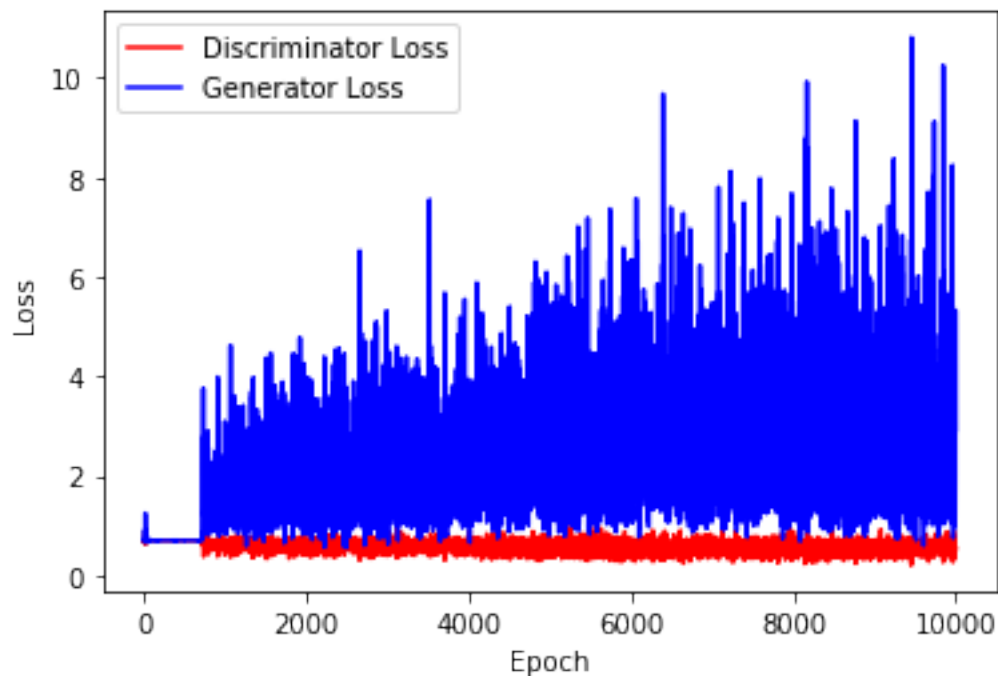
2.0.1 Training the network

Training ABC-GAN for `n_epochs` number of epochs

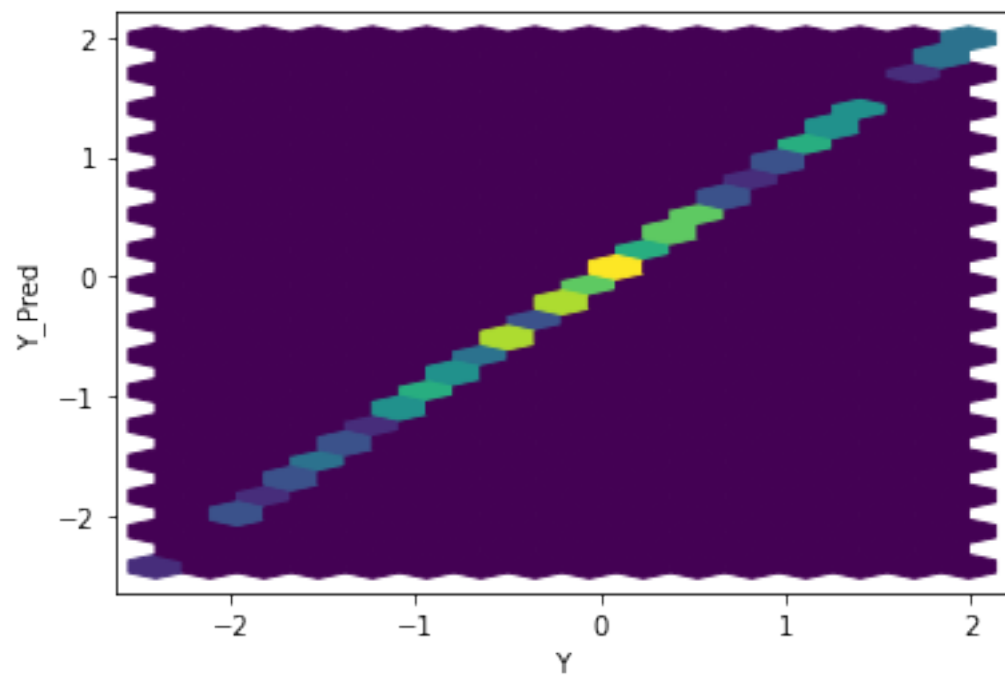
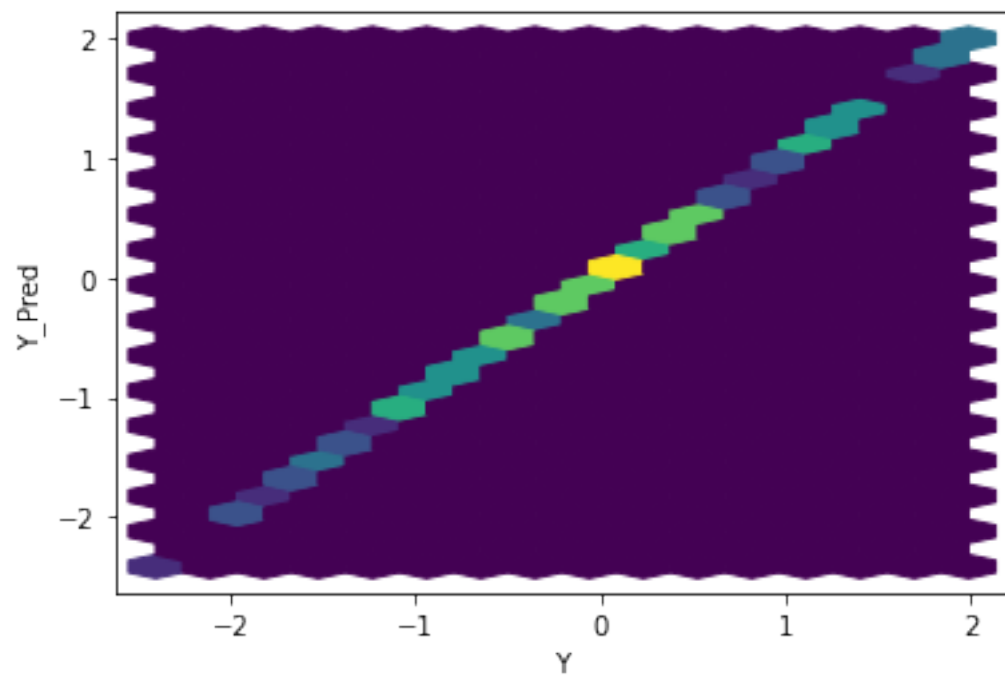
```
[17]: gen = network.Generator(n_features+2)
      disc = network.Discriminator(n_features+2)

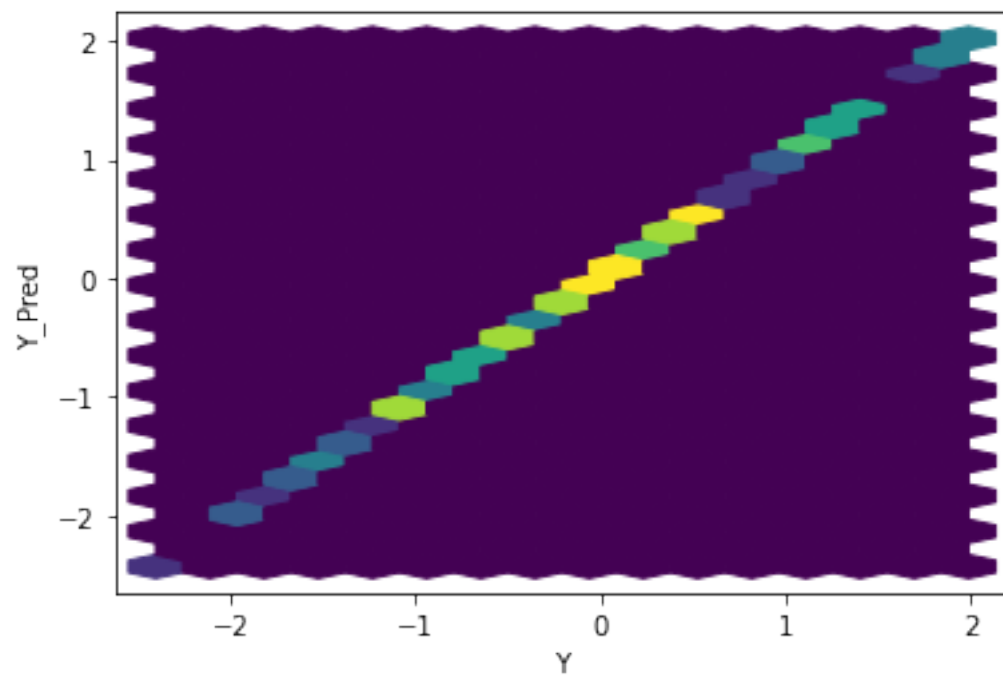
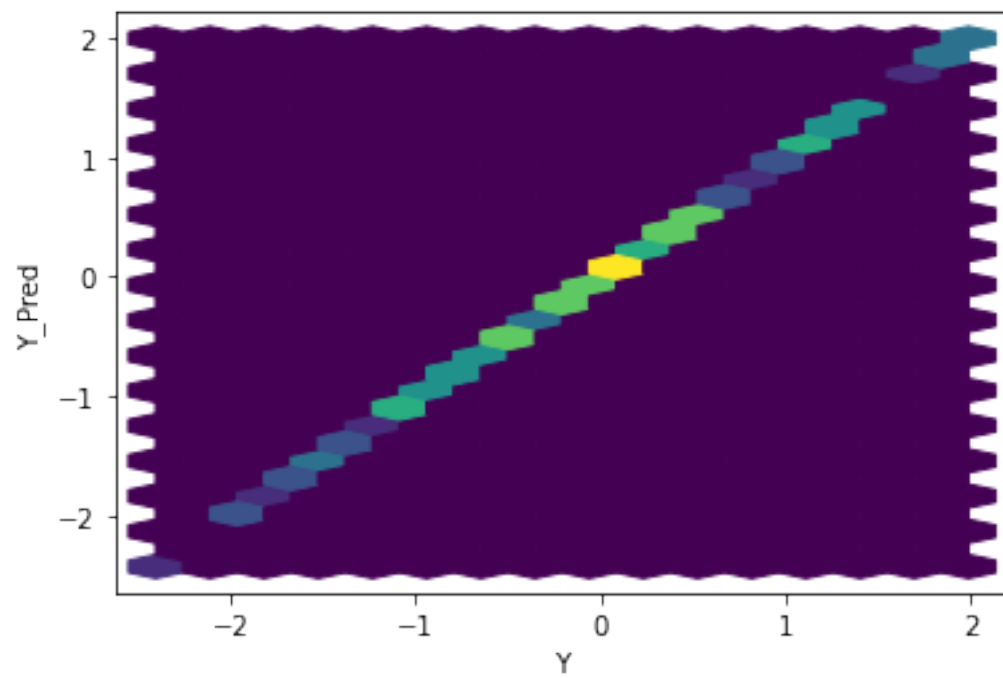
      criterion = torch.nn.BCEWithLogitsLoss()
      gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
      disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))

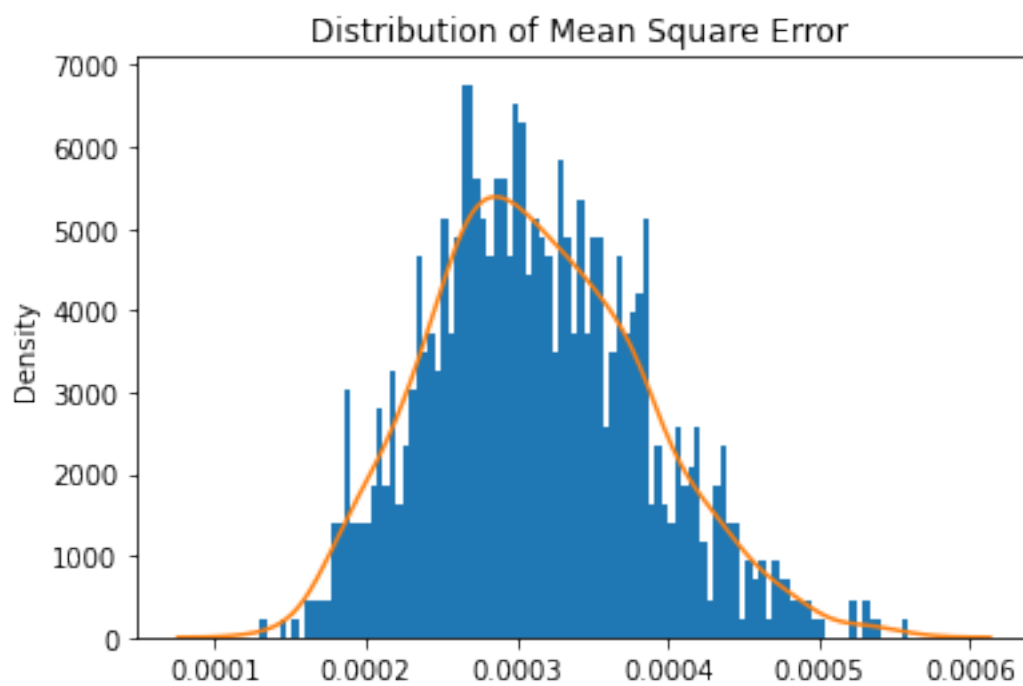
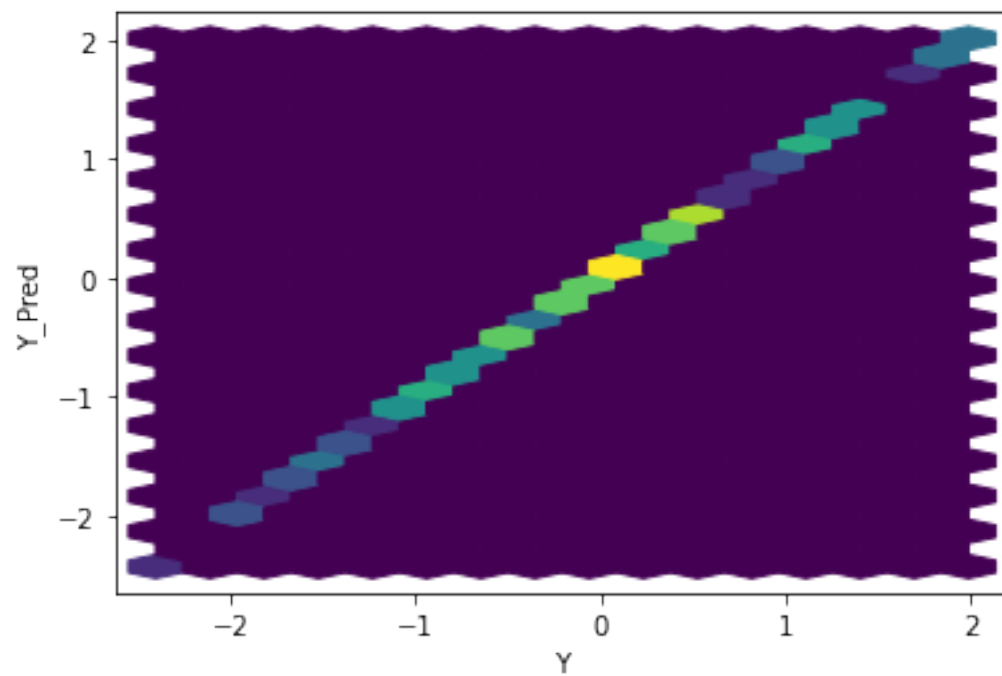
[18]: ABC_train_test.training_GAN(disc, gen,disc_opt,gen_opt,real_dataset,
      ↪batch_size, n_epochs,criterion,coeff,mean,variance,device)
```



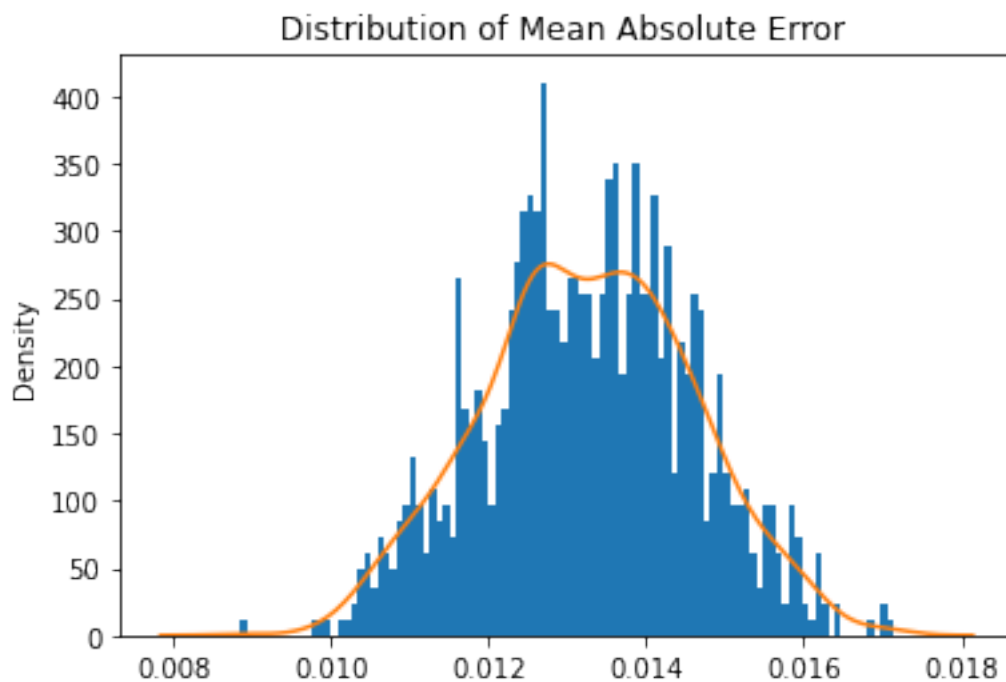
```
[19]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```



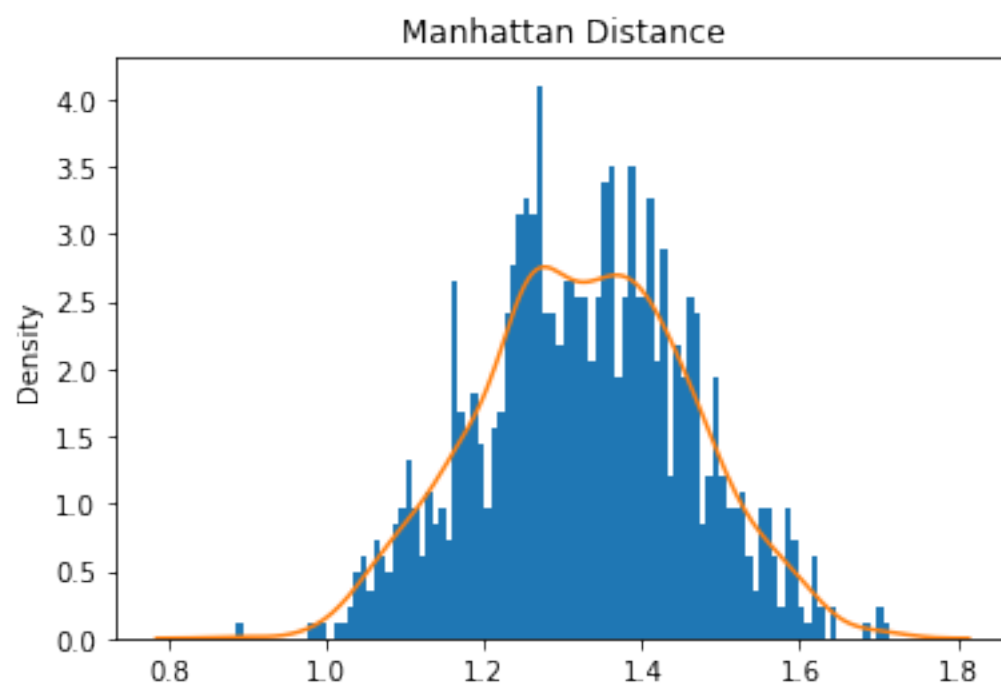




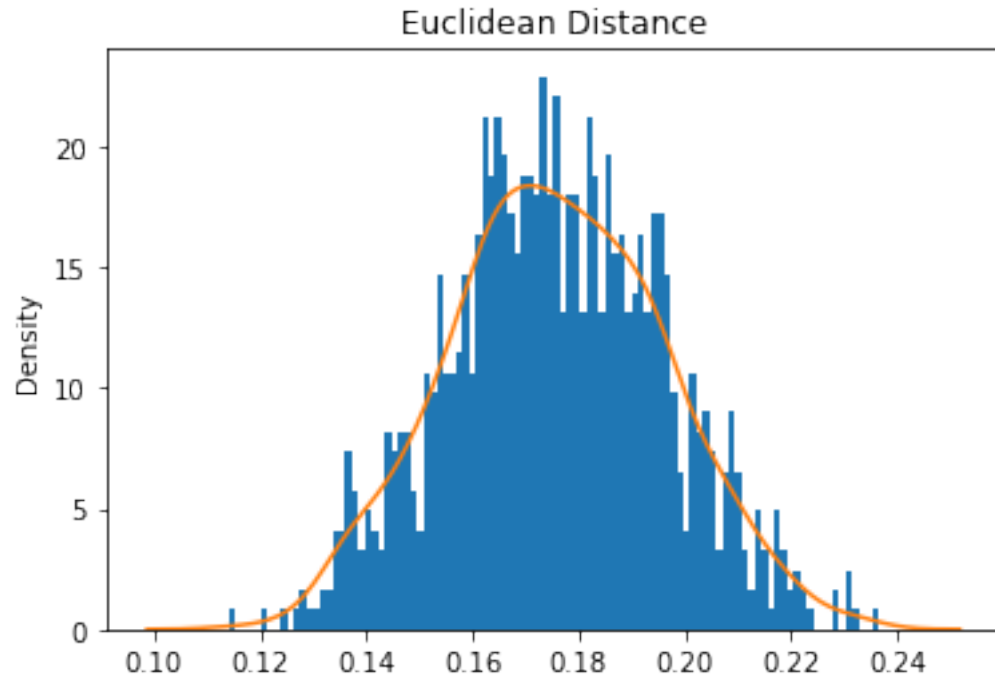
Mean Square Error: 0.0003131597013423791



Mean Absolute Error: 0.013259097947441041
Mean Manhattan Distance: 1.325909794744104

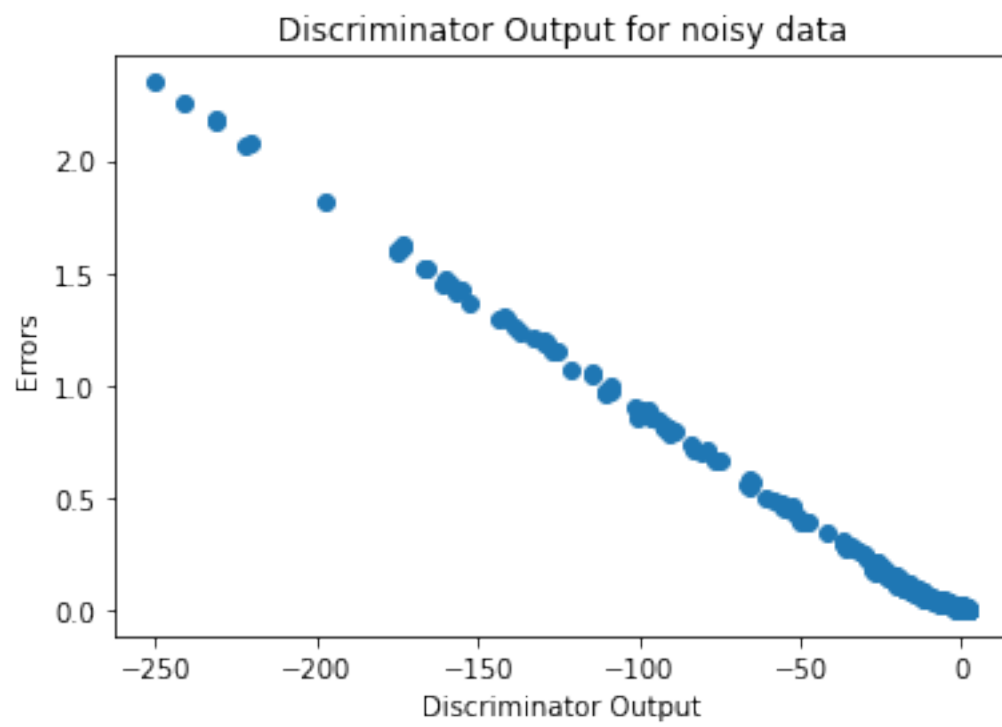
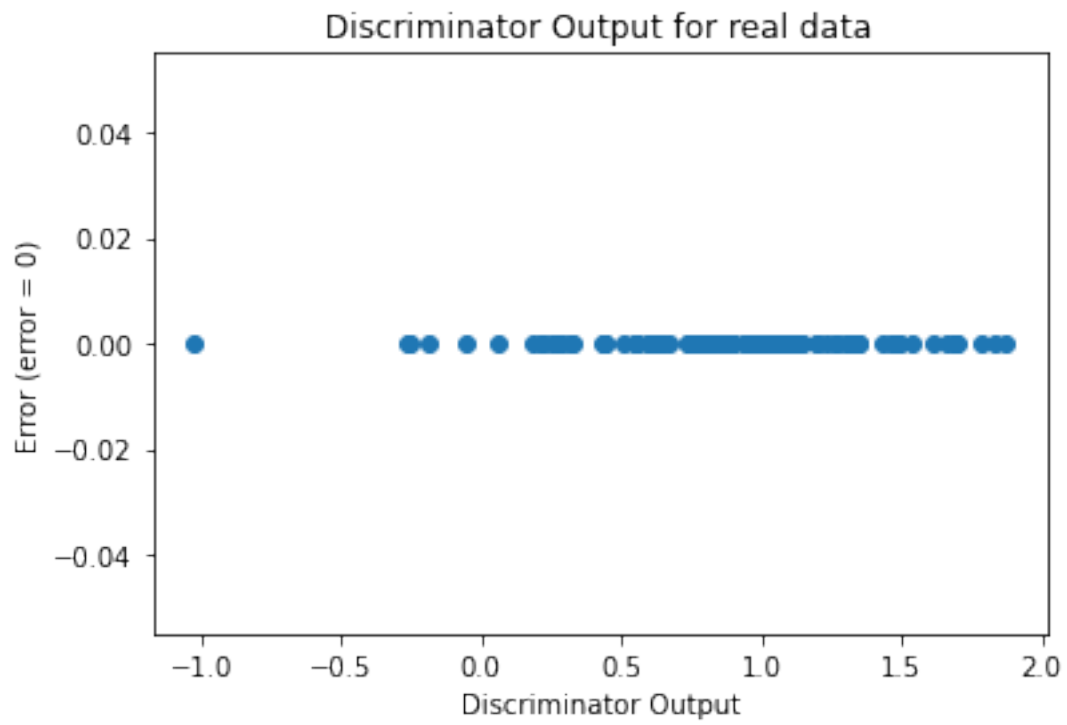


Mean Euclidean Distance: 0.17578680509304517



Sanity Checks

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



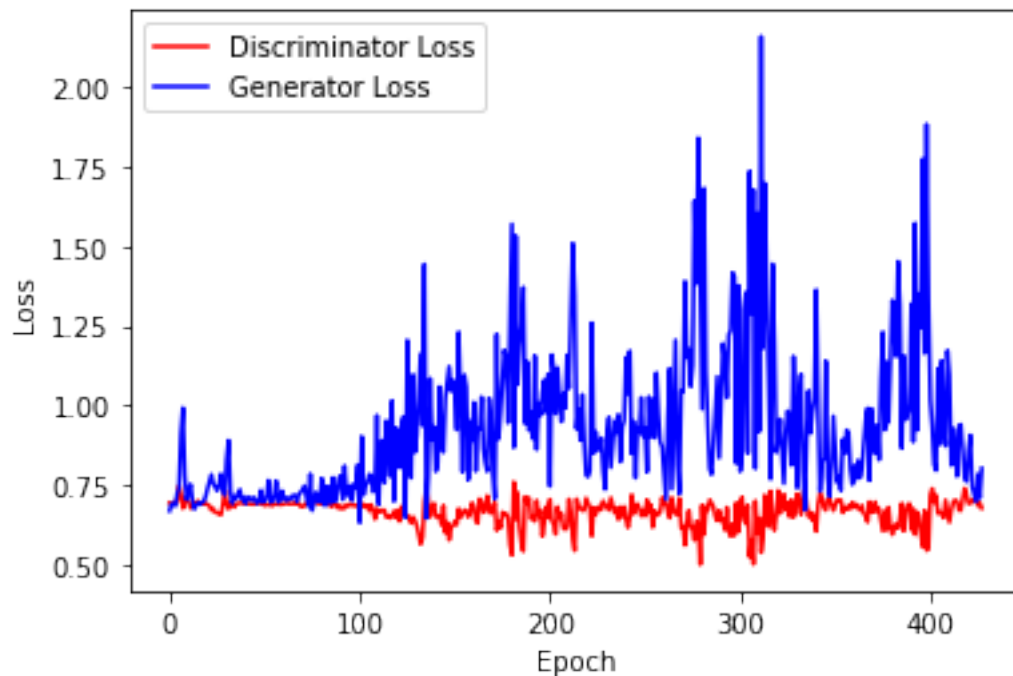
Training GAN until mse of y_pred is > 0.1 or n_epochs < 30000

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

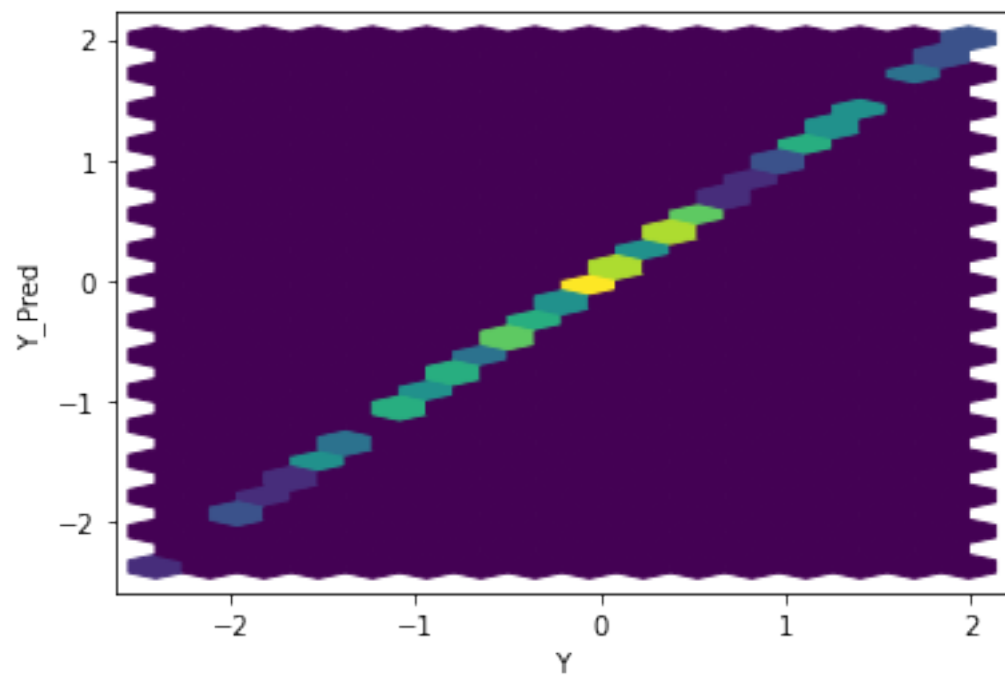
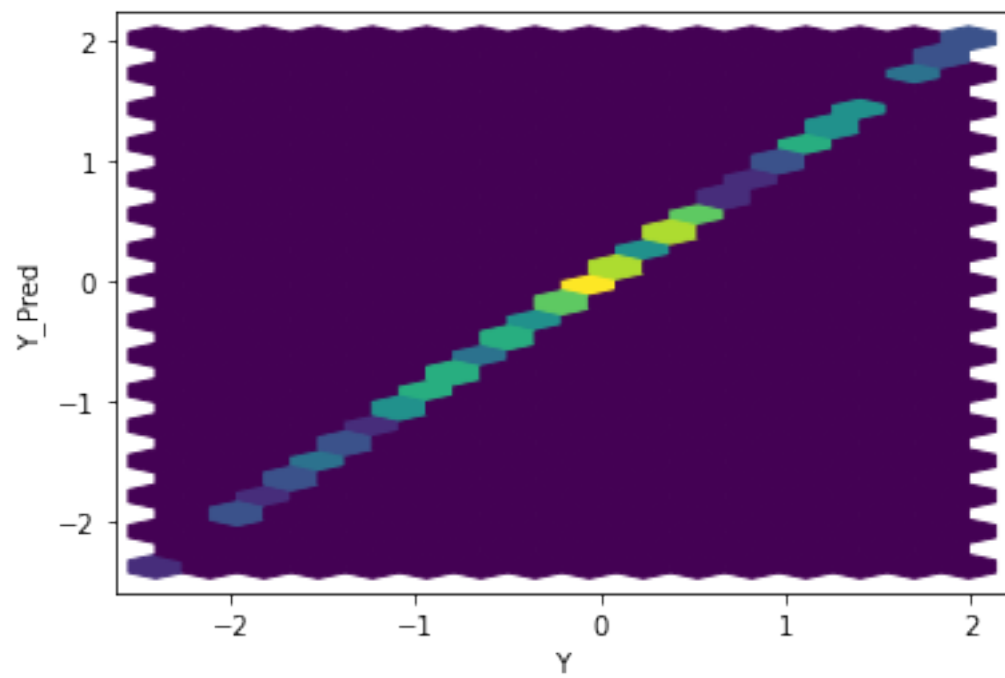
criterion = torch.nn.BCEWithLogitsLoss()
gen_opt = torch.optim.Adam(gen.parameters(), lr=0.01, betas=(0.5, 0.999))
disc_opt = torch.optim.Adam(disc.parameters(), lr=0.01, betas=(0.5, 0.999))
```

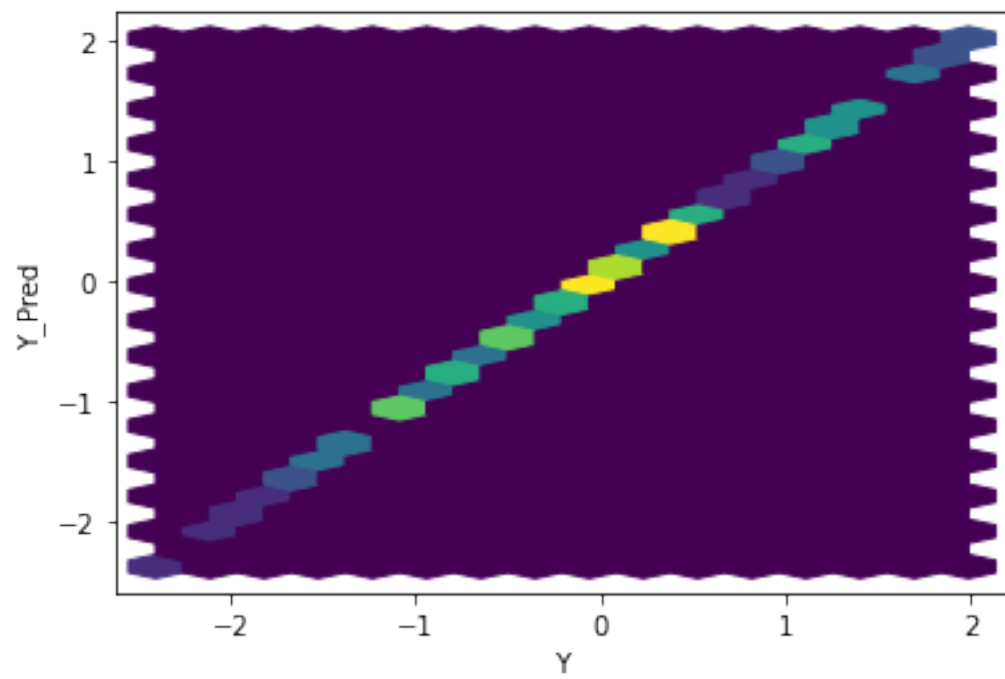
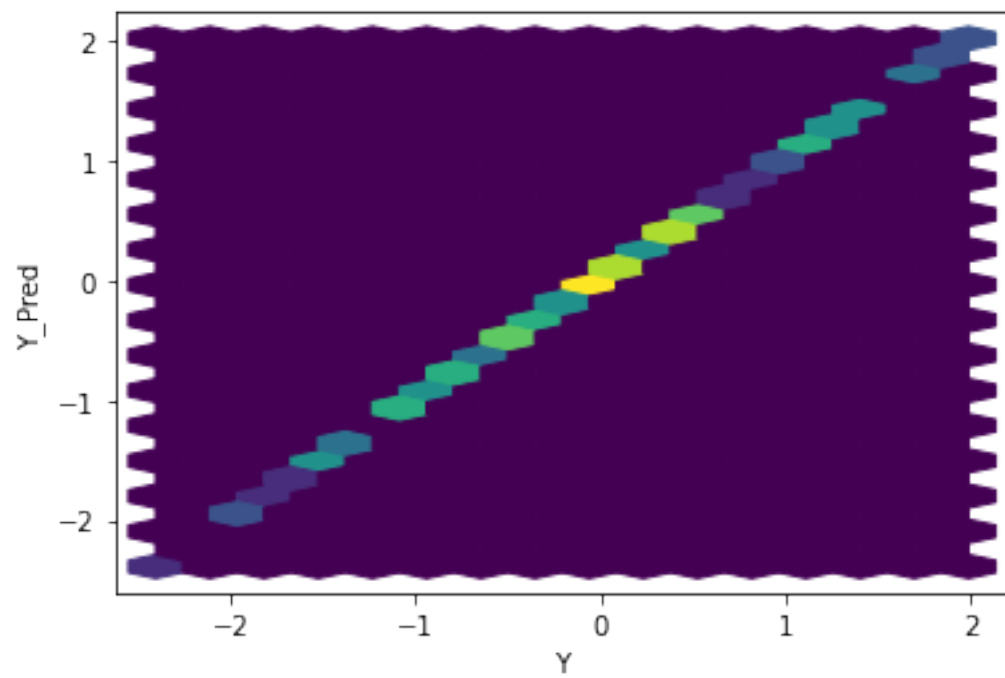
```
[22]: ABC_train_test.
      ↪ training_GAN_2(disc,gen,disc_opt,gen_opt,real_dataset,batch_size,
      ↪ error,criterion,coeff,mean,variance,device)
```

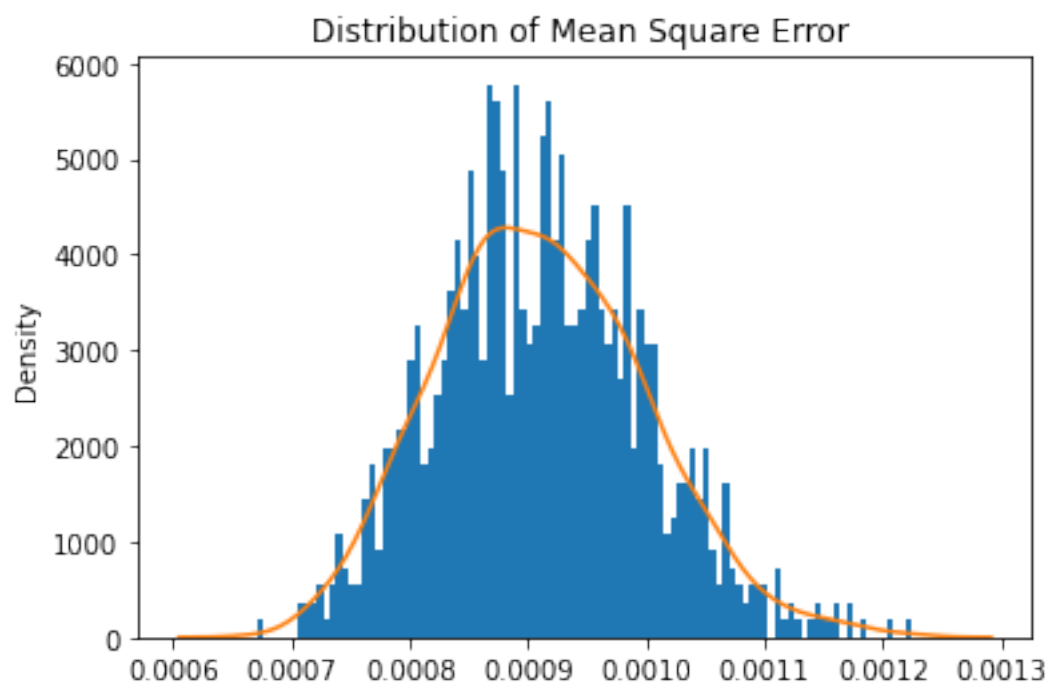
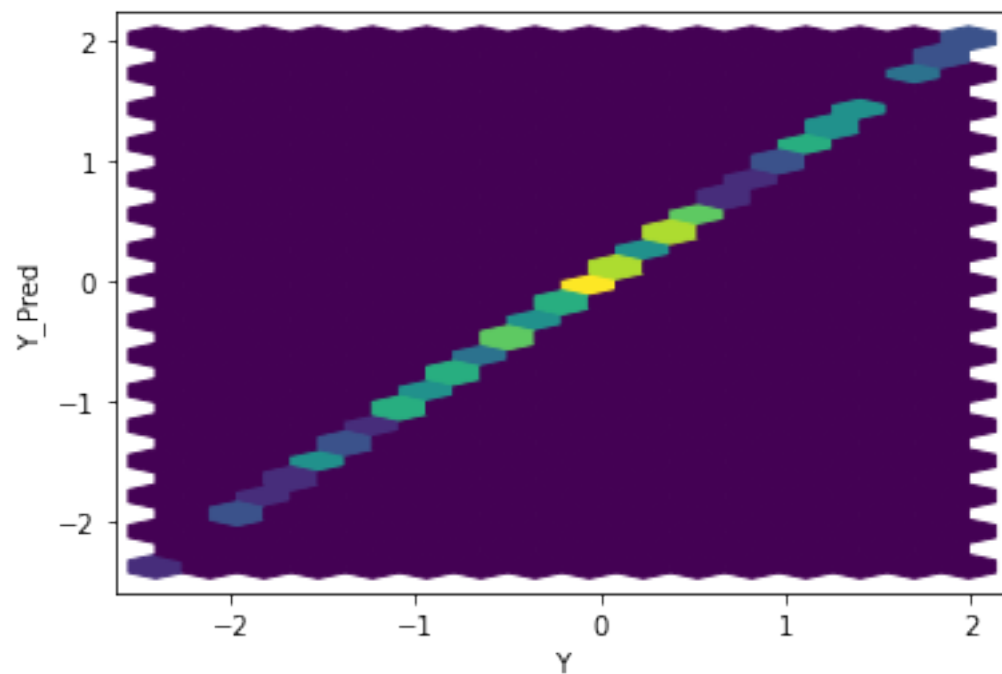
Number of epochs 214



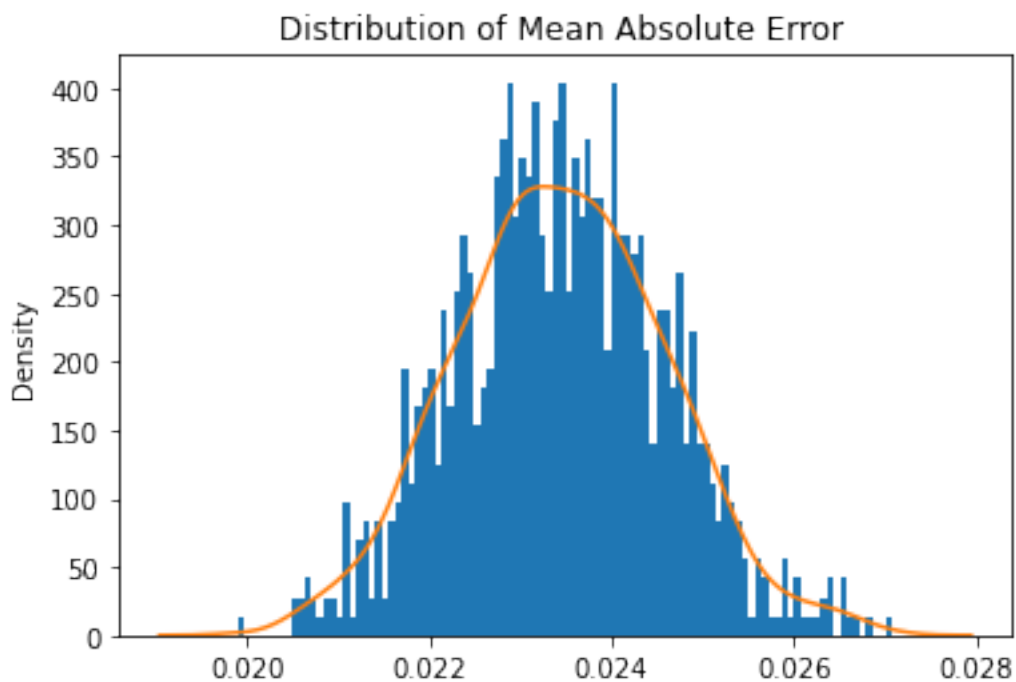
```
[23]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```



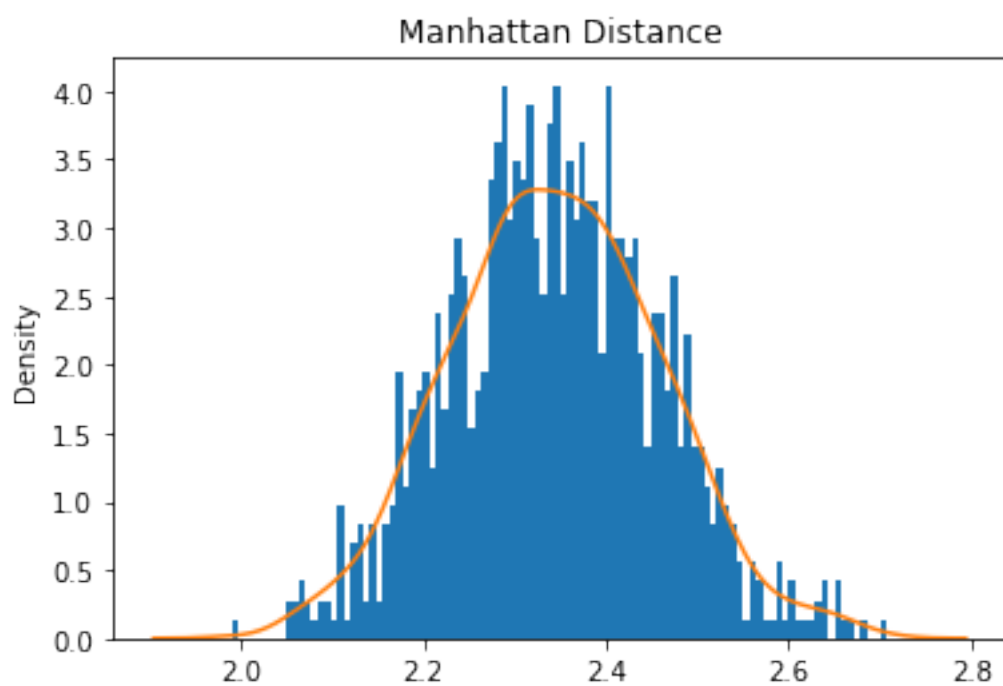




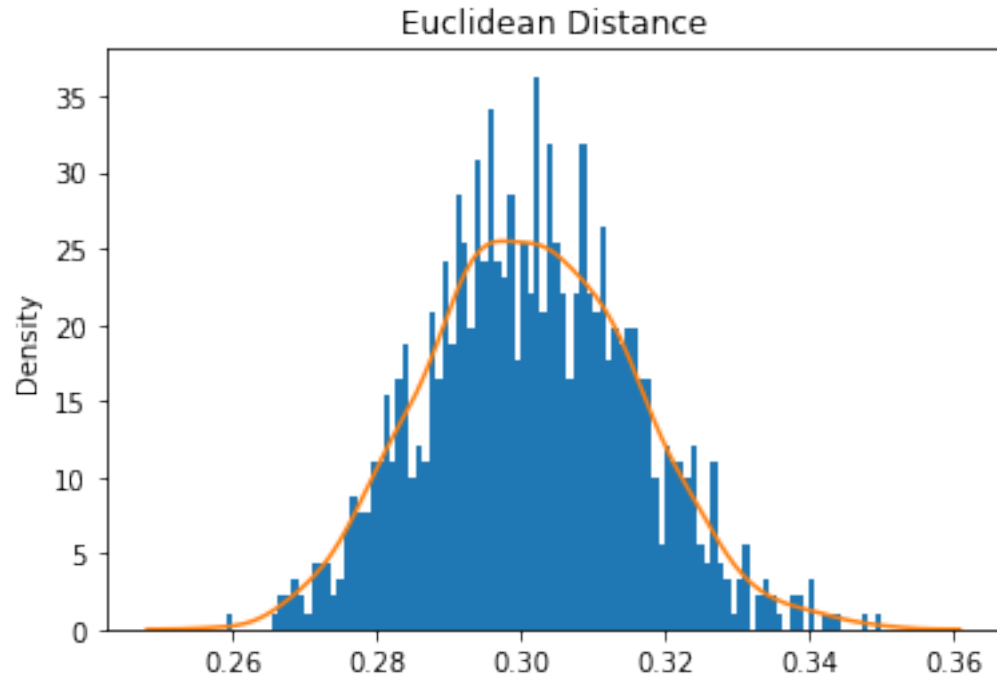
Mean Square Error: 0.0009105367547558498



Mean Absolute Error: 0.023437042755261064
Mean Manhattan Distance: 2.343704275526106



Mean Euclidean Distance: 0.3014004486313472



[]: