# Dataset1-Regression_output_2

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

```python
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 $\beta^*$ 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 = 10
     n_features = 10
     mean = 1
     variance = 0.01
```

## 1.4 Dataset

Generate a random regression problem

$Y = 1 + \beta_1 x_1 + \beta_2 x_2 + \beta_2 x_3 + ... + \beta_n x_n + N(0, \sigma)$ where $\sigma = 0.1$

```
[5]: X,Y = regressionDataset.regression_data(n_samples,n_features)
```

```
          X1         X2         X3         X4         X5         X6         X7  \
0   2.127794  -0.937079  -0.847480   1.019000  -0.667996  -0.416958   0.550997
1  -0.519529  -0.889712  -0.911586   1.549650   1.174325   1.297342   1.037248
2   0.965425   0.132183   0.027848   1.709736   0.516358   0.640920   0.152326
3  -0.680702   0.516085  -0.318750  -0.030849  -2.409276   0.005109  -1.189061
4   0.644779   1.361650  -2.030017  -0.691918   0.173916   0.476185   0.007653

          X8         X9        X10              Y
```

```
0 -0.226009  0.830226  0.958656   182.808266
1 -0.358518  1.356455 -0.992639   205.664595
2 -0.492850  1.596811  0.060344   296.660455
3  1.885670  0.342553 -1.013560  -212.178790
4  0.261403  0.949675 -1.292560    24.066260
```

## 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:                Tue, 19 Oct 2021   Prob (F-statistic):                nan
Time:                        22:48:17   Log-Likelihood:                 328.27
No. Observations:                  10   AIC:                            -636.5
Df Residuals:                       0   BIC:                            -633.5
Df Model:                           9
Covariance Type:            nonrobust
==============================================================================
                 coef    std err          t      P>|t|      [0.025      0.975]
------------------------------------------------------------------------------
const         2.776e-17        inf          0        nan         nan         nan
x1               0.4549        inf          0        nan         nan         nan
x2              -0.0387        inf         -0        nan         nan         nan
x3              -0.0317        inf         -0        nan         nan         nan
x4               0.5405        inf          0        nan         nan         nan
x5               0.7191        inf          0        nan         nan         nan
x6               0.0664        inf          0        nan         nan         nan
x7              -0.1587        inf         -0        nan         nan         nan
x8               0.1699        inf          0        nan         nan         nan
x9               0.1187        inf          0        nan         nan         nan
x10              0.1859        inf          0        nan         nan         nan
==============================================================================
Omnibus:                        2.538   Durbin-Watson:                   1.695
Prob(Omnibus):                  0.281   Jarque-Bera (JB):                0.975
Skew:                           0.244   Prob(JB):                        0.614
Kurtosis:                       1.550   Cond. No.                         34.7
==============================================================================

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    2.775558e-17
```
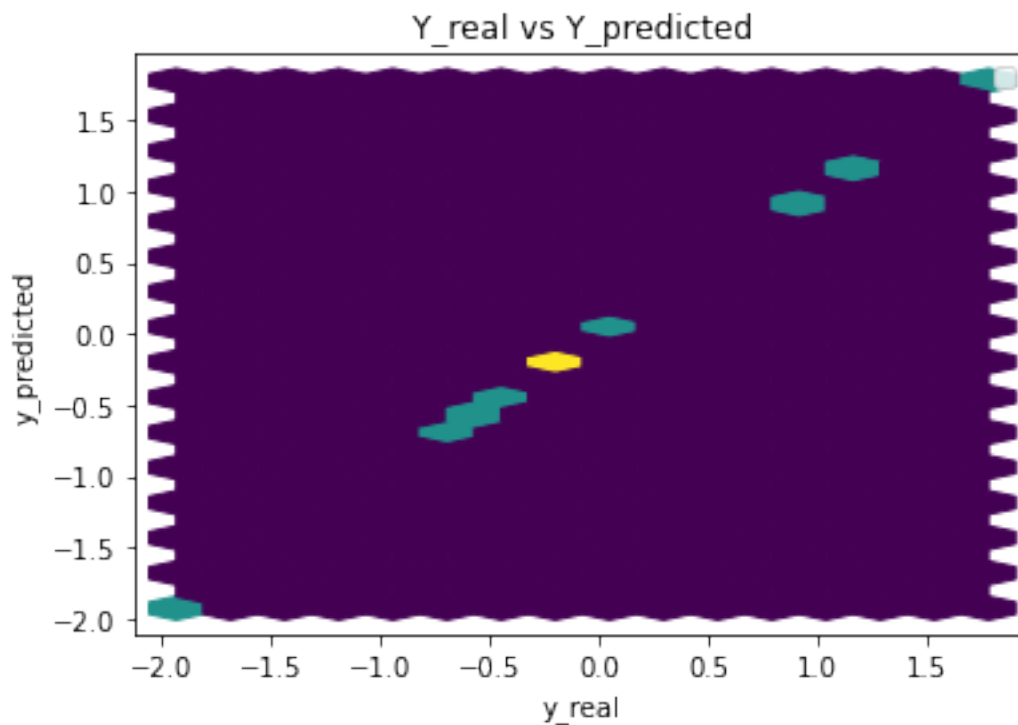
```
x1         4.548814e-01
x2        -3.869486e-02
x3        -3.174064e-02
x4         5.404605e-01
x5         7.191115e-01
x6         6.644852e-02
x7        -1.586722e-01
x8         1.698507e-01
x9         1.186533e-01
x10        1.859416e-01
dtype: float64
```



```
Performance Metrics
Mean Squared Error: 1.797374120599388e-30
Mean Absolute Error: 1.2281842209915794e-15
Manhattan distance: 1.2281842209915794e-14
Euclidean distance: 4.2395449291160815e-15
```

## 1.6 Common Training Parameters (GAN & ABC_GAN)

```python
[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**

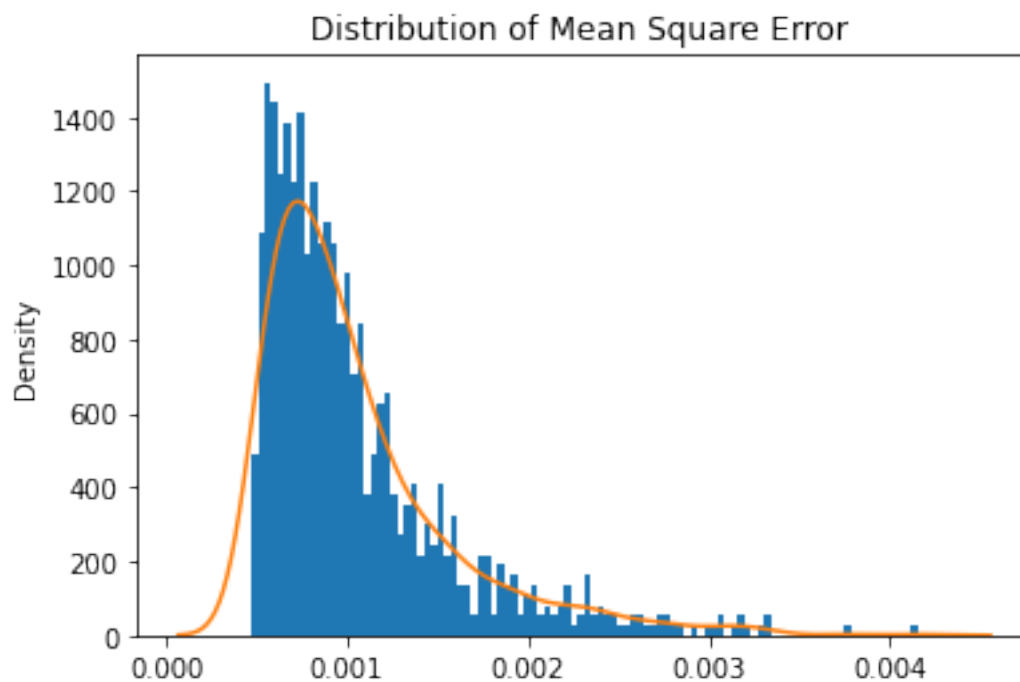```
[9]: generator = network.Generator(n_features+2)
     discriminator = network.Discriminator(n_features+2)

     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)
```
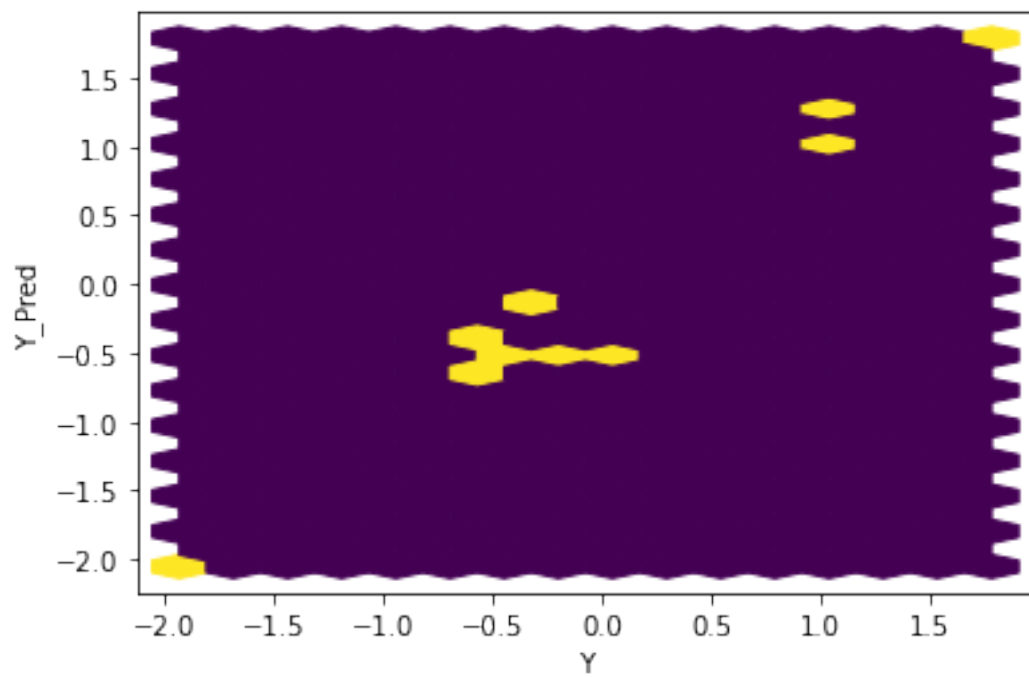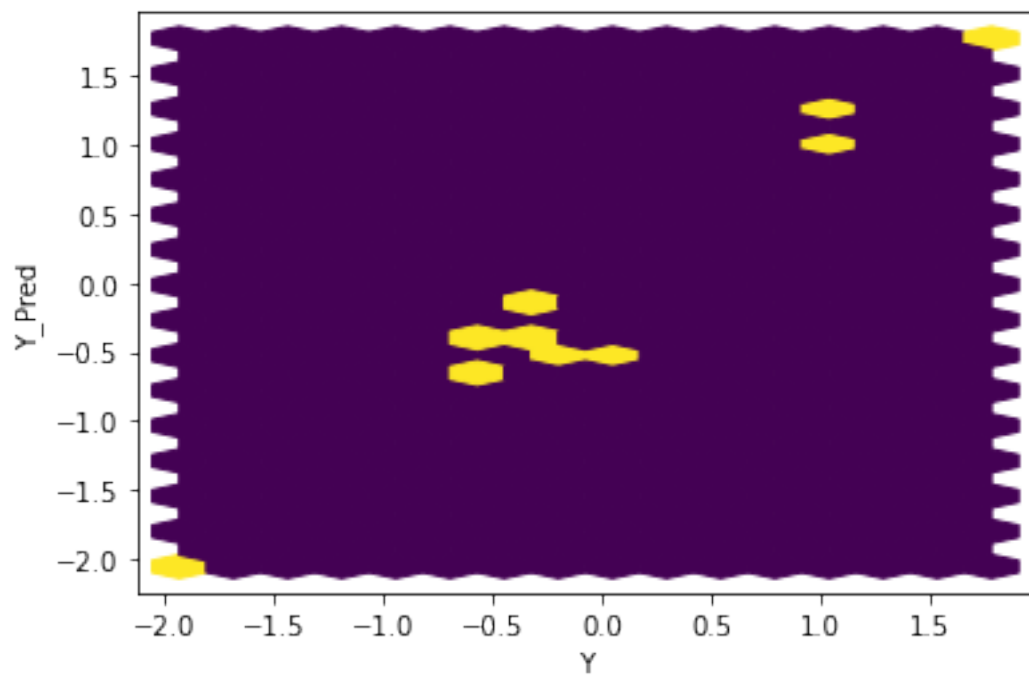


Distribution of Mean Square Error

Mean Square Error: 0.06246369048383713

## Distribution of Mean Absolute Error



Mean Absolute Error: 0.22428416901808232

## Manhattan Distance

Mean Manhattan Distance: 2.242841690180823

## Euclidean Distance



Mean Euclidean Distance: 0.7894223323455603

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

Discriminator Output for real data



Discriminator Output for noisy data

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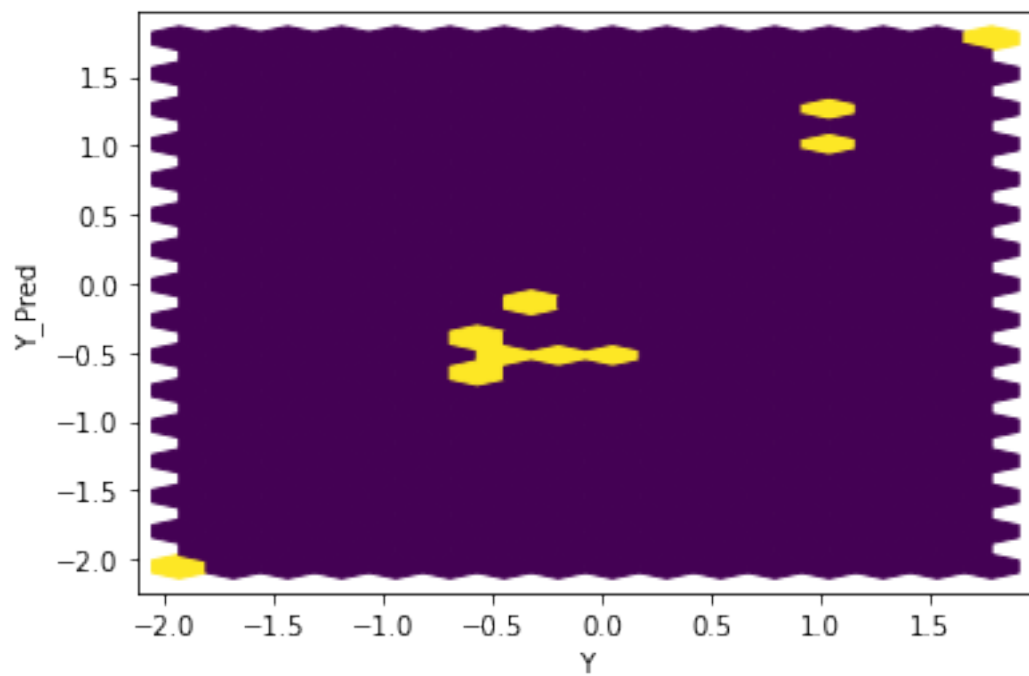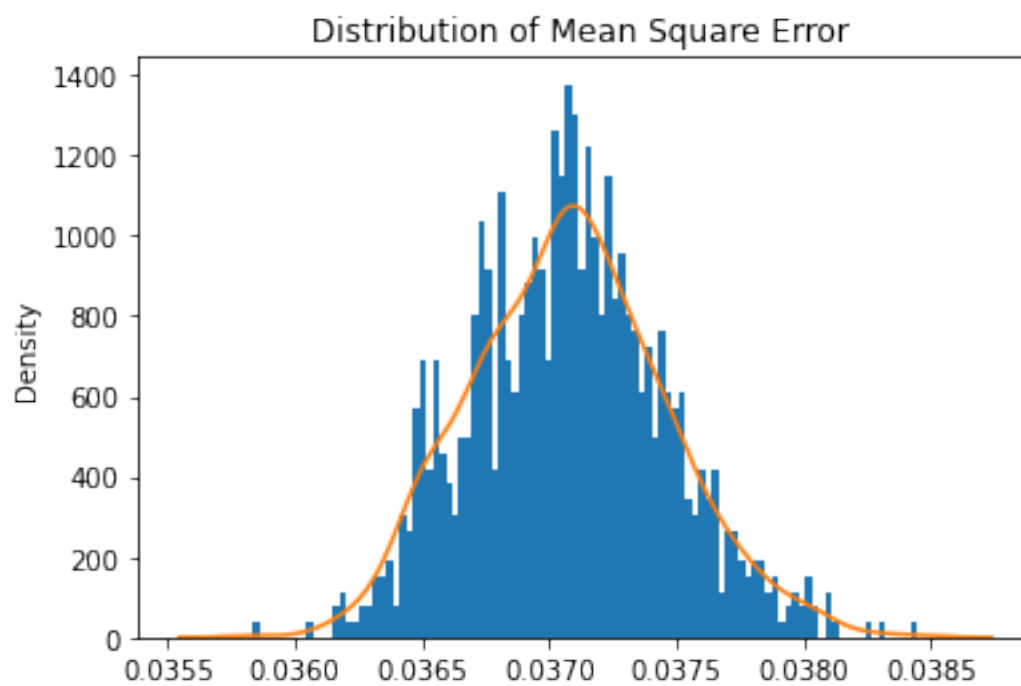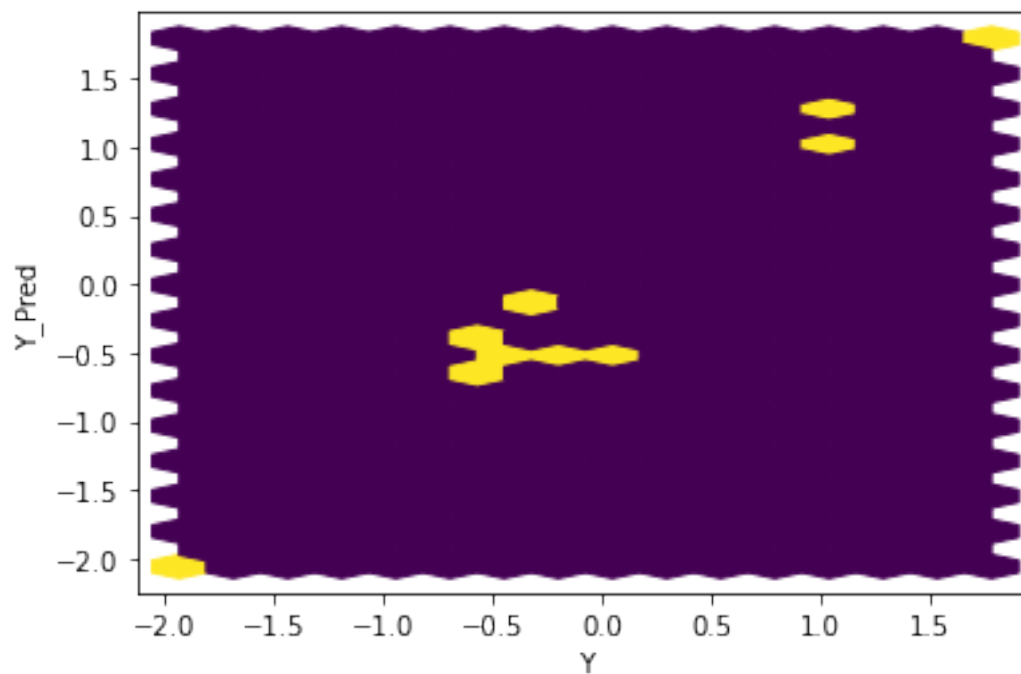
```
[15]: train_test.
       ↪training_GAN_2(discriminator,generator,disc_opt,gen_opt,real_dataset,batch_size,error,crite:
```

Number of epochs needed 524



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

Distribution of Mean Square Error

Mean Square Error: 0.001048027573293151



Distribution of Mean Absolute Error

Mean Absolute Error: 0.02679868476688862

## Manhattan Distance



Mean Manhattan Distance: 0.2679868476688862

## Euclidean Distance

```
Mean Euclidean Distance: 0.09972122825440852
```

# 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)
```

Distribution of Mean Square Error



Mean Square Error: 0.037076254995210105

Distribution of Mean Absolute Error

Mean Absolute Error: 0.12744419066309928
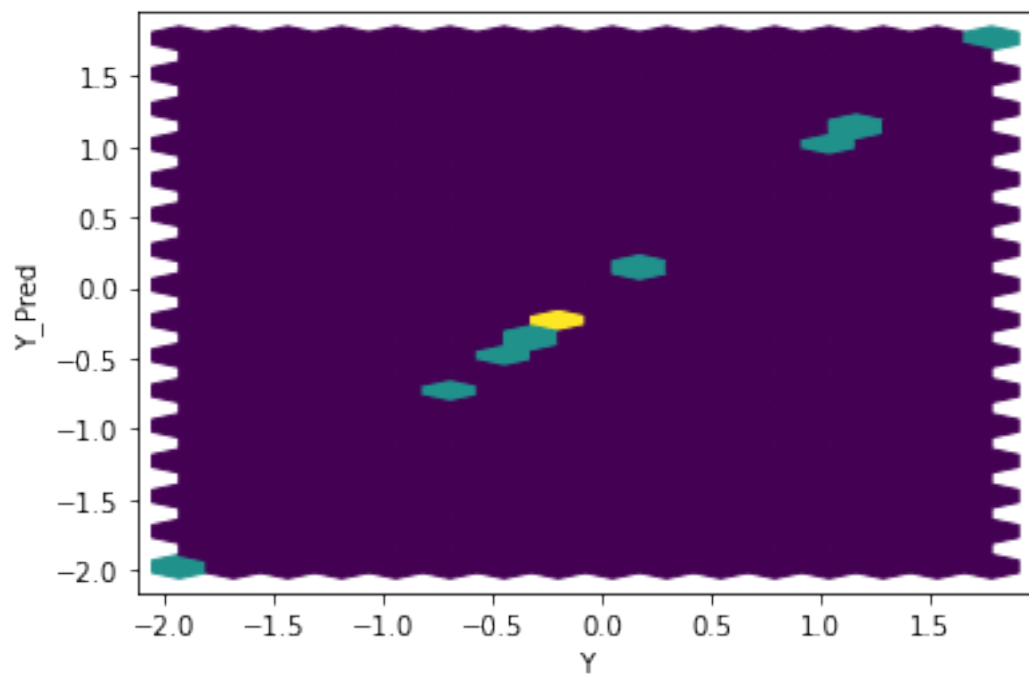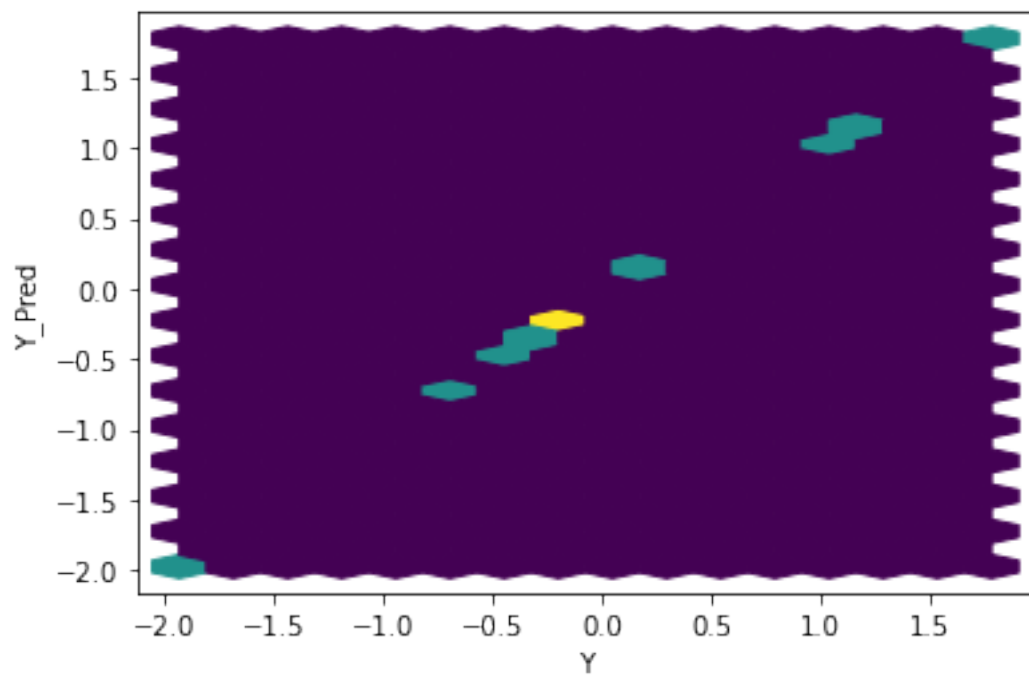Mean Manhattan Distance: 1.2744419066309929


Manhattan Distance

```
Mean Euclidean Distance: 0.6088945945150124
```



Euclidean Distance

## Sanity Checks

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

Discriminator Output for real data



Discriminator Output for noisy data

**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))
```
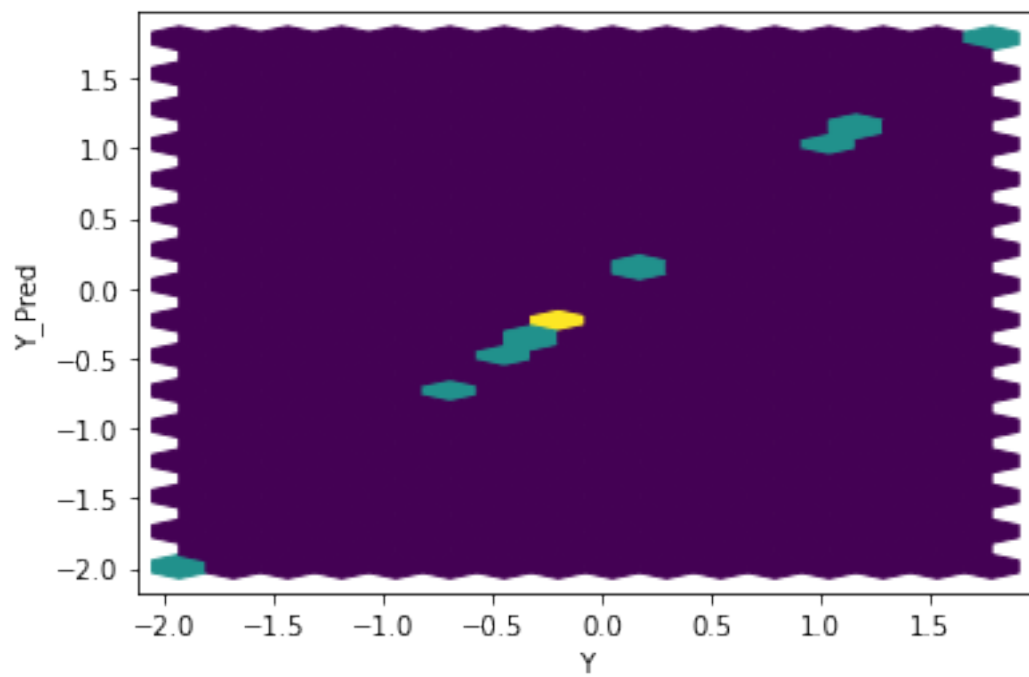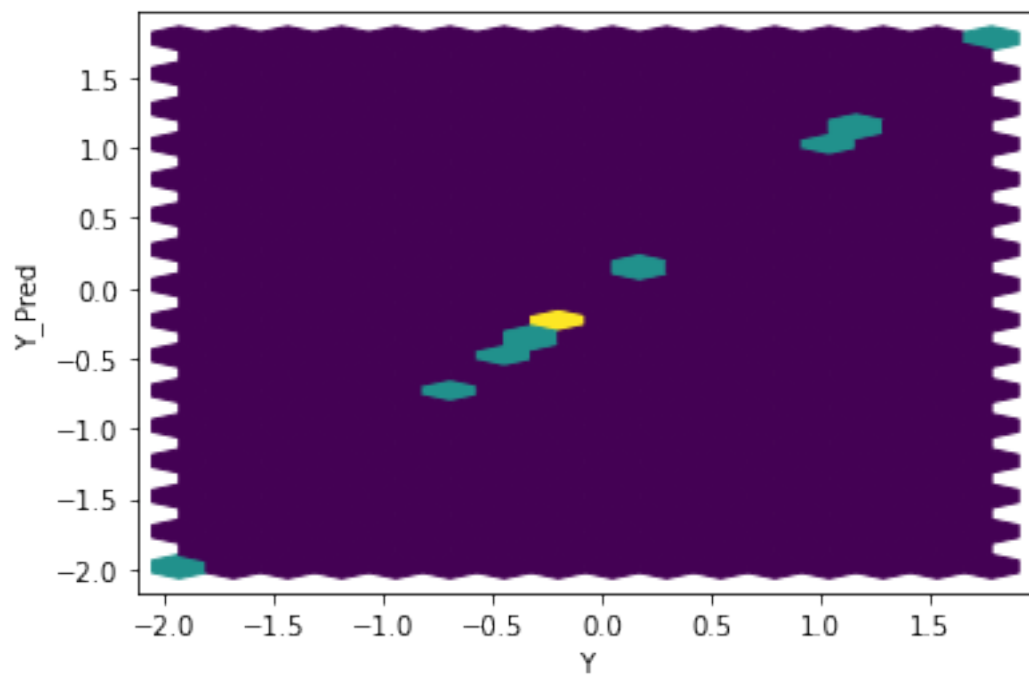
```
[22]: ABC_train_test.
      ↪training_GAN_2(disc,gen,disc_opt,gen_opt,real_dataset,batch_size,␣
      ↪error,criterion,coeff,mean,variance,device)
```
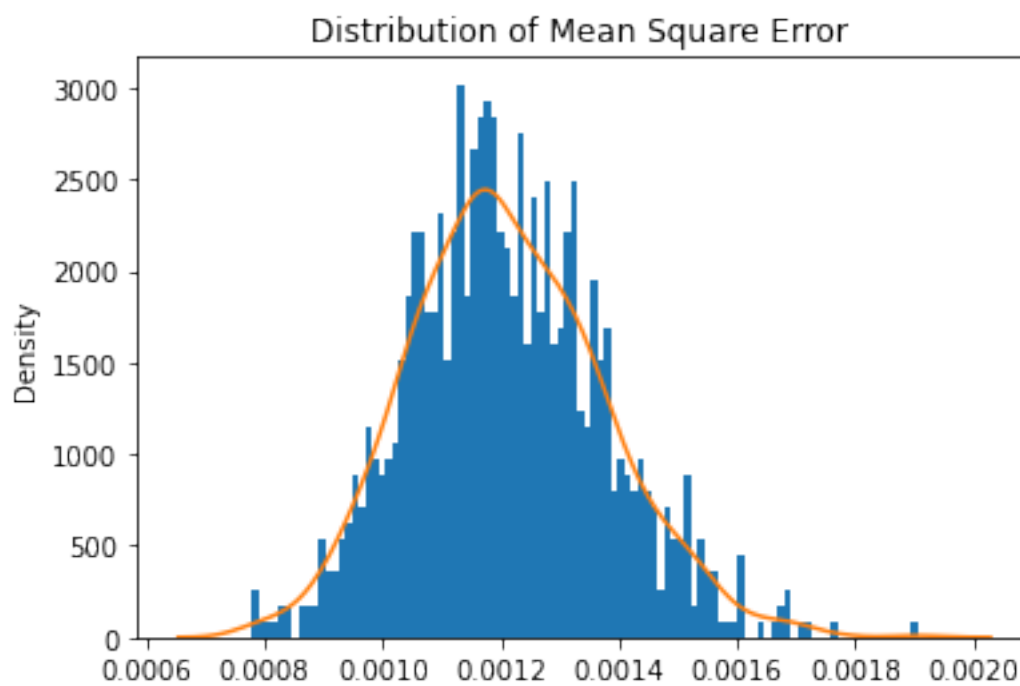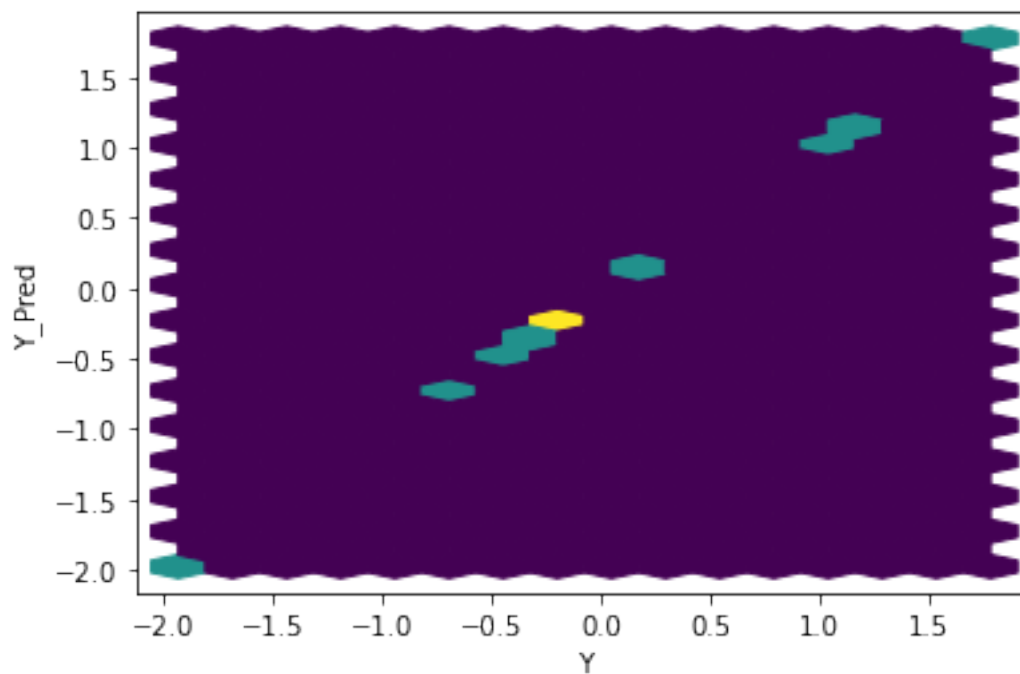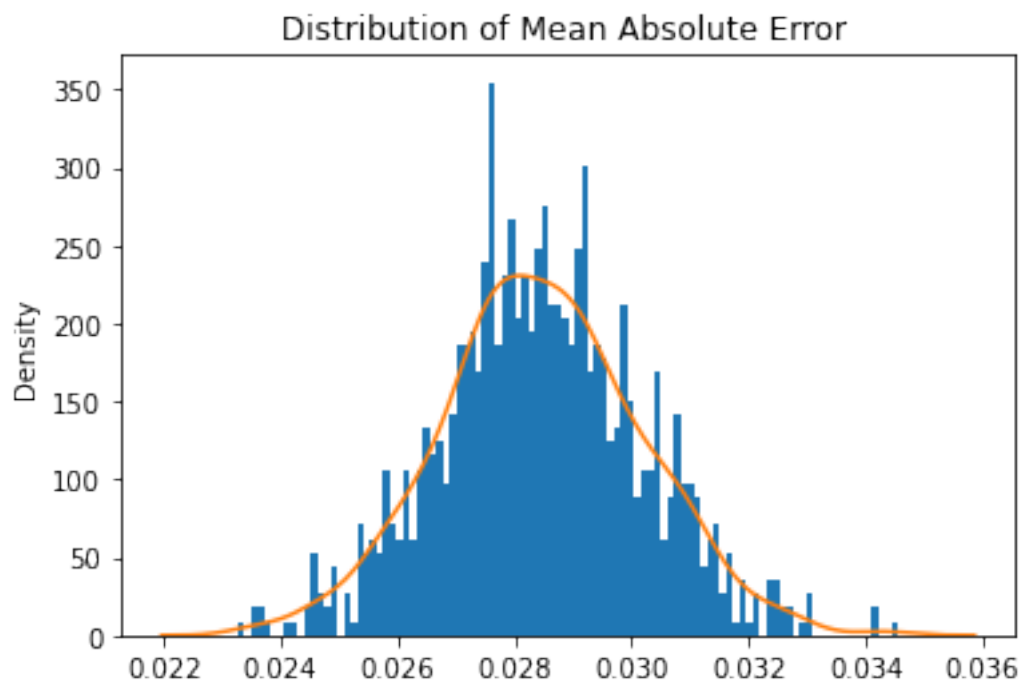
Number of epochs 447



```
[23]: ABC_train_test.test_generator(gen,real_dataset,coeff,mean,variance,device)
```
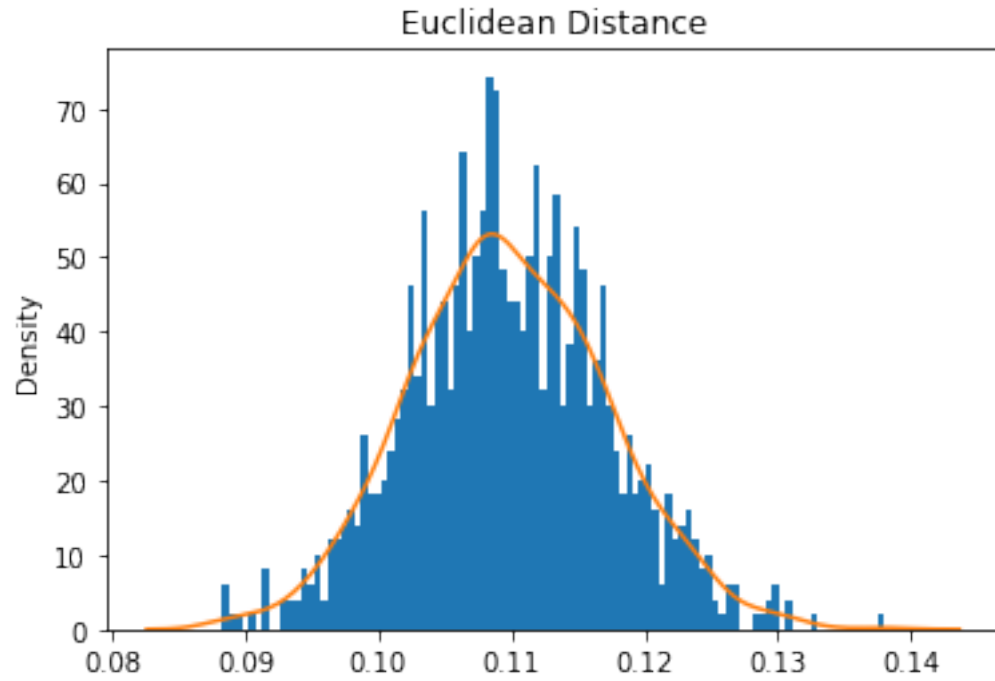
Distribution of Mean Square Error



Mean Square Error: 0.001210072173281629

Distribution of Mean Absolute Error

Mean Absolute Error: 0.028451322716474534
Mean Manhattan Distance: 0.2845132271647453



Manhattan Distance

```
Mean Euclidean Distance: 0.10974866060474416
```



Euclidean Distance

```
[ ]:
```