# CSE512 Spring 2021 - Machine Learning - Homework 1

Your Name: Sagar Onkar Toshniwal

Solar ID: 113260061

NetID email address: [sagaronkar.toshniwal@stonybrook.edu](mailto:sagaronkar.toshniwal@stonybrook.edu)

Names of people whom you discussed the homework with: None but I have taken reference from Google, Stackoverflow, etc.

# Question 1 – Probability (30 points)

Let X1 and X2 be independent continuous random variables uniformly distributed from 0 to 1. Let X = max(X1; 2X2). Compute:

1. The expectation E(X)

2. The variance V ar(X)

3. The covariance: Cov(X;X1).

Q.1] Let $X_1$ & $X_2$ be Independent random Variables uniformly distributed from 0 to 1. Let $X = max(X_1, 2X_2)$.
Compute :-

ⓐ E(x) Expectation:-

→ Given $X = max(X_1, 2X_2)$

Let $G_x(x) = CDF$ of $X$

$F(x) = CDF$ of $X_1, X_2$

$g(x) = PDF$ of $X$

We know that,

$G_x(x) = Pr(x \leq x)$

$= Pr(max(X_1, 2X_2) \leq x)$

Since $X_1$ & $X_2$ are Independent Random Variables

$= Pr(X_1 \leq x) \, Pr(2X_2 \leq x)$

∴ when $0 \leq x \leq 1$

$= Pr(X_1 \leq x) \, Pr.(X_2 \leq x/2)$

when $1 \leq x \leq 2$

$= Pr(X_2 \leq x/2)$ ... as $X_1 = 1$

$$\therefore \quad G_x(x) = F(x) \cdot F(x/2) \quad \cdots \quad 0 \leq x \leq 1$$
$$= F(x/2) \quad \cdots \quad 1 \leq x \leq 2$$

$$\therefore \quad g(x) = \frac{d}{dx} \left[ G_x(x) \right]$$

$$= \frac{d}{dx} \left[ G_x(x) \right]$$

$$= \frac{d}{dx} \left[ F(x) \cdot F(x/2) \right] \quad \cdots \quad 0 \leq x \leq 1$$

$$\frac{d}{dx} \left[ F(x/2) \right] \quad \cdots \quad 1 \leq x \leq 2$$

$$F(x) = 0 \qquad\qquad F(x/2) = 0 \qquad \cdots \quad x \leq 0$$
$$= x \qquad\qquad\qquad = x/2 \qquad \cdots \quad 0 \leq x \leq 1$$
$$= 1 \qquad\qquad\qquad = x/2 \qquad \cdots \quad 1 \leq x \leq 2$$

$$\therefore \quad g(x) = \frac{d}{dx} \left[ F(x) \cdot F(x/2) \right] \overset{\&}{\ast} \frac{d}{dx} \left[ F(x/2) \right]$$

$$= \frac{d}{dx} \left[ x \cdot \frac{x}{2} \right] \qquad \overset{\&}{\ast} \quad \frac{d}{dx} \left[ x/2 \right] \qquad \cdots \quad \text{for both Intervals.}$$

$$= \frac{1}{2} \frac{d}{dx} \left[ x^2 \right] \overset{\&}{\ast} \qquad \frac{1}{2} \frac{d}{dx} \left[ x \right]$$

$$= x \qquad\qquad\qquad \& \qquad 1$$

$$g(x) = x \qquad \cdots \qquad 0 \leq x \leq 1$$
$$= \frac{1}{2} \qquad \cdots \qquad 1 \leq x \leq 2$$

$\therefore$ ⓐ $\quad E(x) = \displaystyle\int_{-\infty}^{\infty} x \cdot g(x) \, dx$

$\qquad = \displaystyle\int_{-\infty}^{\infty} x \cdot g(x) \, dx$

$\qquad = \displaystyle\int_{0}^{1} x \cdot g(x) \, dx + \int_{1}^{2} x \cdot g(x) \, dx$

$\qquad = \displaystyle\int_{0}^{1} x \cdot x \, dx + \int_{1}^{2} x \cdot \frac{1}{2} \, dx$

$\qquad = \displaystyle\int_{0}^{1} x^2 \, dx + \frac{1}{2} \int_{1}^{2} x \, dx$

$\qquad = \left[ \dfrac{x^3}{3} \right]_0^1 + \dfrac{1}{2} \left[ \dfrac{x^2}{2} \right]_1^2$

$\qquad = \dfrac{1}{3} + \dfrac{3}{4}$

$E(x) \quad = \quad \dfrac{13}{12}$

$\therefore \boxed{E(x) = \dfrac{13}{12}}$

(b)   Var $(x)$

→    Var $(x) = E[x^2] - [E(x)]^2$

$\therefore E[x^2] = \int_{-\infty}^{\infty} x^2 \cdot g(x)\, dx$

$$= \int_{0}^{1} x^2 \cdot x\, dx \quad + \quad \int_{1}^{2} x^2 \cdot \frac{1}{2}\, dx$$

$$= \int_{0}^{1} x^3\, dx \quad + \quad \frac{1}{2}\int_{1}^{2} x^2\, dx$$

$$= \frac{[x^4]_0^1}{4} \quad + \quad \frac{1}{2} \cdot \frac{1}{3} [x^3]_1^2$$

$$= \frac{1}{4} + \frac{7}{6} \quad = \quad \frac{34}{24}$$

$$= \frac{17}{12}$$

$\therefore$ Var $(x) = \quad \dfrac{17}{12} - \left(\dfrac{13}{12}\right)^2$

$$= \quad \frac{17}{12} - \frac{169}{144}$$

$$\boxed{\text{Var } (x) = \quad \frac{35}{144}}$$

③ $Cov(x, X_1)$

→ $Cov(x, X_1) = Cov(X_1, \max(X_1, X_2))$

$= Cov(X_1, X_1) \cdot Pr(x_1 = \max(x, x_2))$

$+$

$Cov(X_1, X_2) \cdot Pr(x = \max(x_1, x_2))$

∵ Since $x_1$ & $x_2$ are Independent.

∴ $Cov(X_1, X_2)$ must be $0$

$= Cov(X_1, X_1) \cdot Pr(x_1 > X_2)$

$= Var(X_1) \cdot Pr(x_1 > x_2)$

Since $Cov(x, X_1) == Var(x_1)$

$= Var(X_1) [(1 - Pr \, x_1 \leq x)]$

$= Var(X_1) [\int_0^1 1 - x \, dx)]$

$Var(X_1) = E(X_1^2) - [E(X_1)]^2$

$= \int_0^1 x^2 \, dx - [\int_0^1 x \, dx]^2$

$= [\frac{x^3}{3}]_0^1 - \{[\frac{x^2}{2}]_0^1\}^2$

$$= \frac{1}{3} - \left(\frac{1}{2}\right)^2$$

$$= \frac{1}{3} - \frac{1}{4}$$

$$= \frac{1}{12}$$

$$\therefore \quad Cov\ (x,\ x_1) = \quad Var\ (x_1) \left[\int_0^1 1 - x\, dx \right]$$

$$= \frac{1}{12} \left[ \int_0^1 1\, dx - \int_0^1 x\, dx \right]$$

$$= \frac{1}{12} \left[ 1 - \left[\frac{x^2}{2}\right]_0^1 \right]$$

$$= \frac{1}{12} \left[ 1 - \frac{1}{2} \right]$$

$$= \frac{1}{12} \times \frac{1}{2}$$

$$= \quad \frac{1}{24}$$

$$\boxed{\therefore \quad Cov\ (x,\ x_1) = \frac{1}{24}}$$

1. E(X) =   13/12
2. Var(X) =  35/144
3. Covv(X,X1) =  1/24

Answer for Question 2.2 :-

(a) (5 points) Report the values of mu0, var0, mu1, var1.

```
[54]: [mu0, var0, mu1, var1] = get_mean_and_variance(X, y)

print ("Mu0 = ",mu0)
print ("Var0 = ",var0)
print ("Mu1 = ",mu1)
print ("Var1 = ",var1)

Mu0 = [63.38 0.4 ]
Var0 = [2.99587347e+02 2.44897959e-01]
Mu1 = [50.51685393 0.47191011]
Var1 = [233.97994033  0.25061893]
```
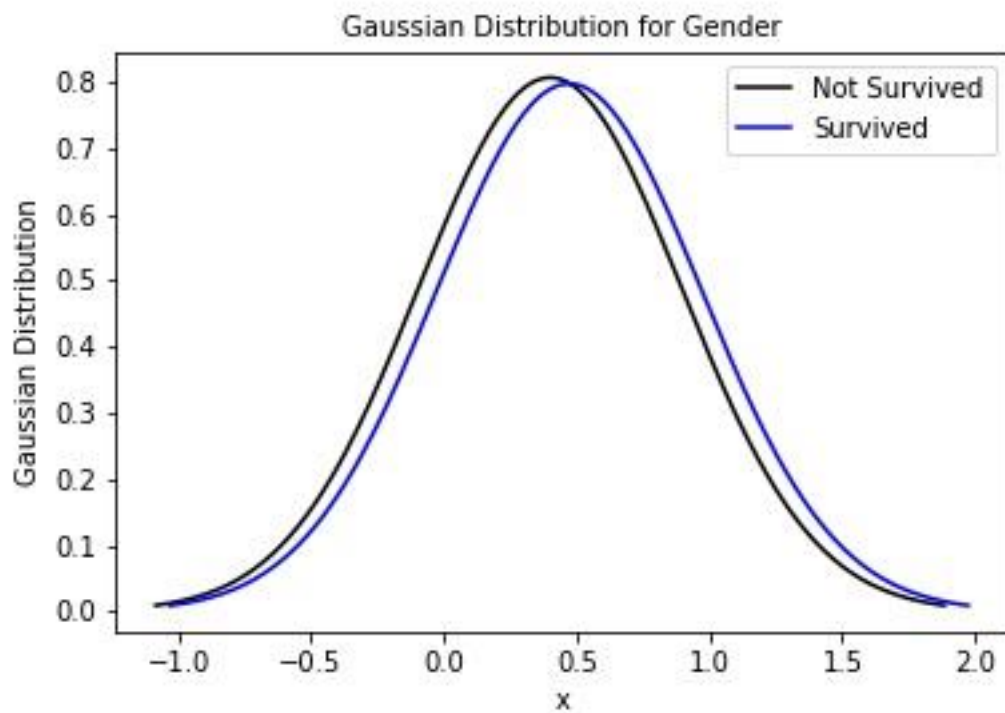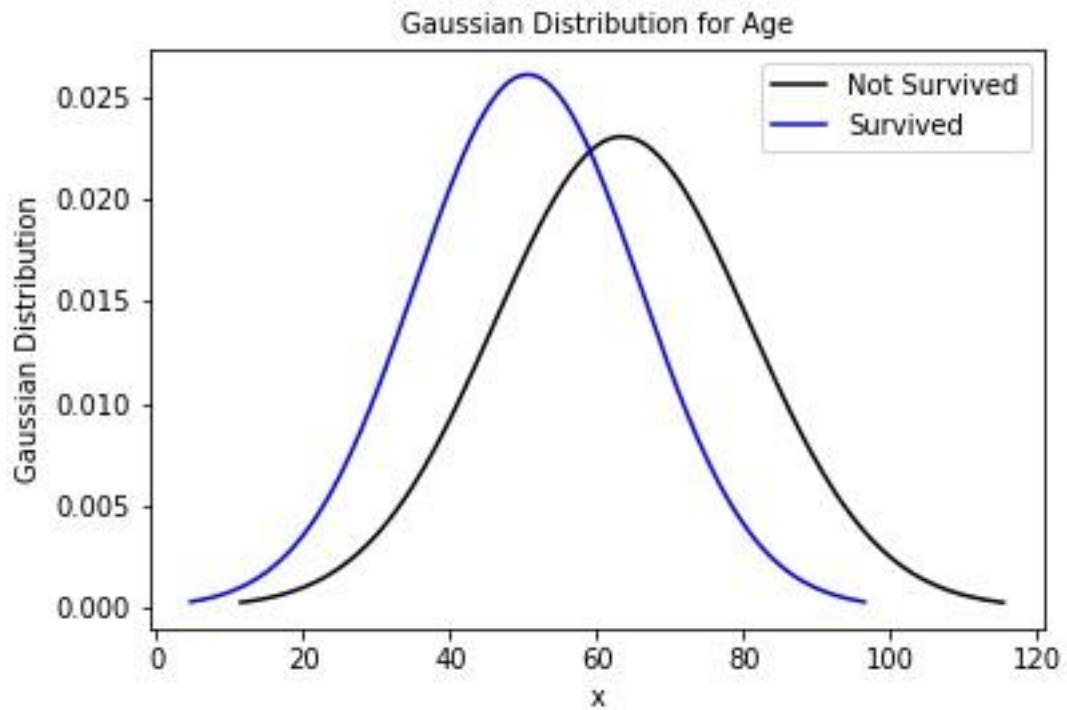
Mu0 = [63.38    0.4 ]
Var0 = [2.99587347e+02   2.44897959e-01]
Mu1 = [50.51685393    0.47191011]
Var1 = [233.97994033     0.25061893]

(b) (5 points) For each feature j, plot the Gaussian distribution with mean mu0[j] and variance var0[j] in black color. On the same graph, plot the Gaussian distribution with mean mu1[j] and variance var1[j] in blue. You can use Python packages matplotlib and scipy.



Gaussian Distribution for Age



Gaussian Distribution for Gender

(c) (5 points) Is it a good idea to approximate gender by a Gaussian distribution? Why or why not?

→ No it is not a good idea to approximate Gender by a Gaussian distribution for the following reason :-

- Gaussian distribution is a type of continuous probability distribution for a real-valued random variable.
- Gender has limited outcomes and is discrete in nature. Gender is similar to a model which has the set of possible outcomes of any single experiment that asks a yes–no question.
- For such discrete data and possible outcomes, **Binomial distribution** is perfectly suited.

## Python code for reference :-

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from scipy.stats import norm
from sklearn.linear_model import LinearRegression
from sklearn import metrics

def get_mean_and_variance(X, y) :
    data = np.append(X,y,1)
    data_0 = data[data[:, -1] == 0, :-1]
    data_1 = data[data[:, -1] == 1, :-1]

    mu0 = np.mean(data_0, axis =0)
    var0 = np.var(data_0, axis =0)
    mu1 = np.mean(data_1, axis =0)
    var1 = np.var(data_1, axis =0)
    return mu0,var0,mu1,var1

df = pd.read_csv("covid19_metadata.csv")
y = df.iloc[:,-1:]
y = y.replace(to_replace = ['Y','N'],value = [1,0])
X = df.iloc[:,:-1]
X.gender = X.gender.replace(to_replace = ['F','M'],value = [1,0])

[mu0, var0, mu1, var1] = get_mean_and_variance(X, y)

print ("mu0 = ",mu0)
print ("var0 = ",var0)
print ("mu1 = ",mu1)
print ("var1 = ",var1)

sigma0 = np.sqrt(var0)
sigma1 = np.sqrt(var1)

fig, (ax1,ax2) = plt.subplots (1,2, figsize = (12,9))
xpts0 = np.linspace(mu0[0] - 3*sigma0[0], mu0[0] + 3*sigma0[0], 200)
xpts1 = np.linspace(mu1[0] - 3*sigma1[0], mu1[0] + 3*sigma1[0], 200)

ax1.plot(xpts0, norm.pdf(xpts0, mu0[0], sigma0[0]),color = 'black')
ax1.plot(xpts1, norm.pdf(xpts1, mu1[0], sigma1[0]),color = 'blue')
ax1.set_title('Gaussian Distribution for Age',fontsize=10)

ax1.set(xlabel='x', ylabel='Gaussian Distribution')
ax1.legend(['Not Survived', 'Survived'])

xpts0 = np.linspace(mu0[1] - 3*sigma0[1], mu0[1] + 3*sigma0[1], 200)
xpts1 = np.linspace(mu1[1] - 3*sigma1[1], mu1[1] + 3*sigma1[1], 200)
ax2.plot(xpts0, norm.pdf(xpts0, mu0[1], sigma0[1]),color = 'black')
ax2.plot(xpts1, norm.pdf(xpts1, mu1[1], sigma1[1]),color = 'blue')
ax2.set_title('Gaussian Distribution for Gender',fontsize=10)

ax2.set(xlabel='x', ylabel='Gaussian Distribution')
ax2.legend(['Not Survived', 'Survived'])
plt.show()
```
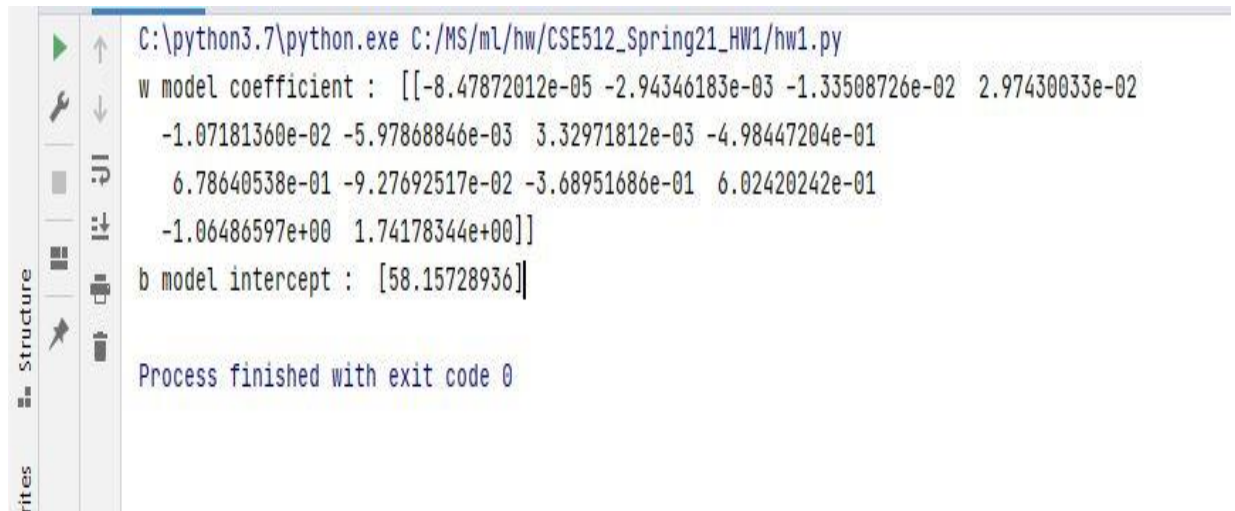
Answer for Question 3.2 :-

(a) (5 points) Report the learned parameters: the weights w and the intercept term b.
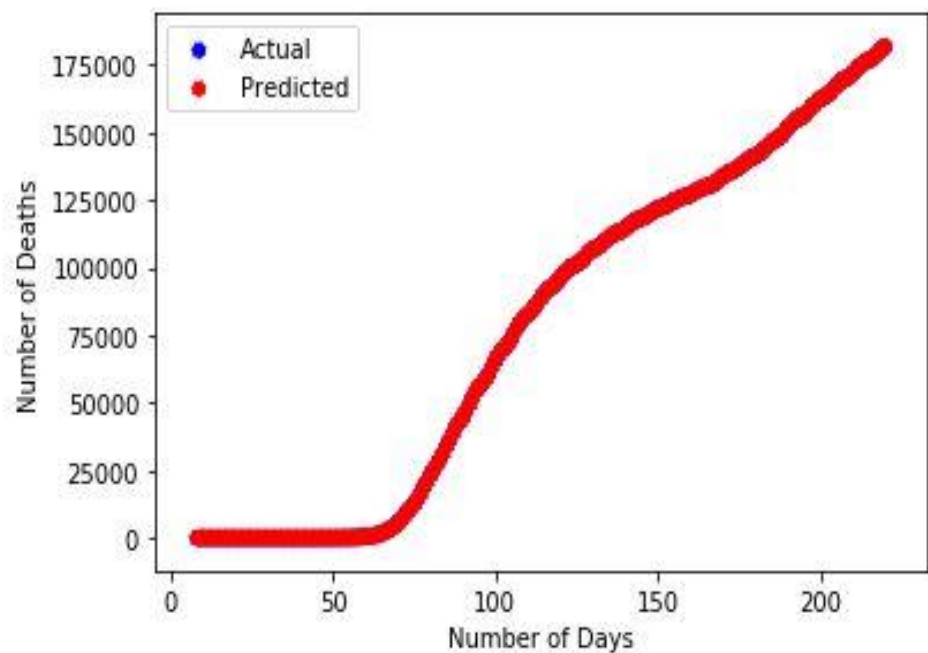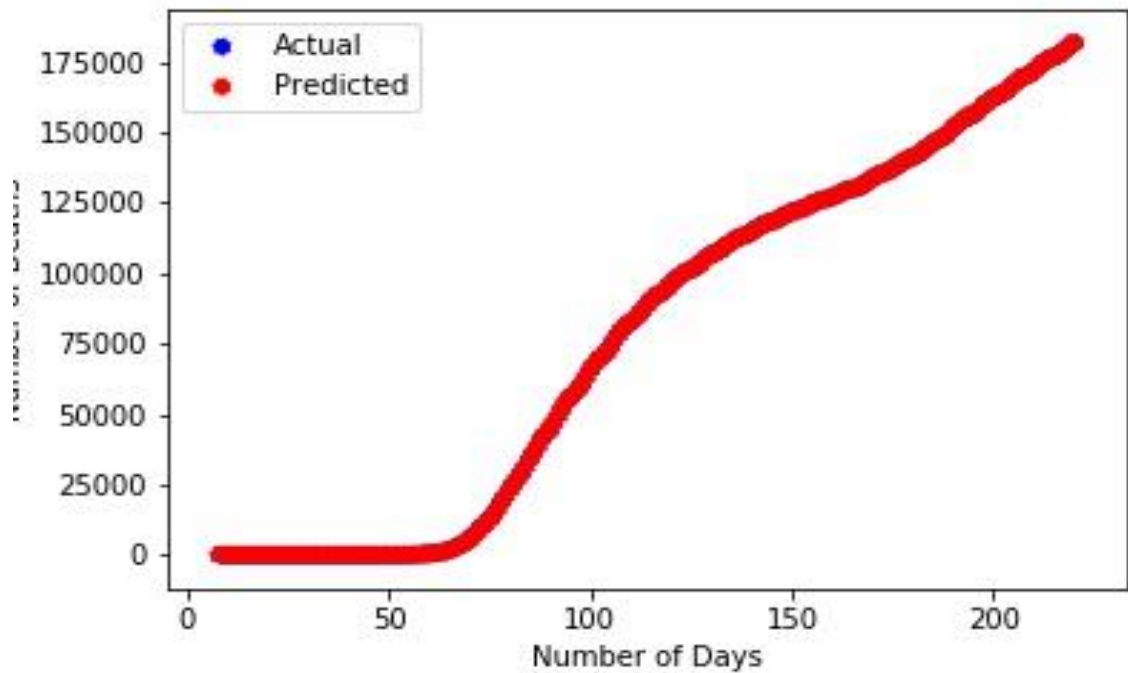
```
      C:\python3.7\python.exe C:/MS/ml/hw/CSE512_Spring21_HW1/hw1.py
      w model coefficient :  [[-8.47872012e-05 -2.94346183e-03 -1.33508726e-02  2.97430033e-02
        -1.07181360e-02 -5.97868846e-03  3.32971812e-03 -4.98447204e-01
         6.78640538e-01 -9.27692517e-02 -3.68951686e-01  6.02420242e-01
        -1.06486597e+00  1.74178344e+00]]
      b model intercept :  [58.15728936]

      Process finished with exit code 0
```

w model coefficient :  **[[-8.47872012e-05**
**-2.94346183e-03**
**-1.33508726e-02**
**2.97430033e-02**
**-1.07181360e-02**
**-5.97868846e-03**
**3.32971812e-03**
**-4.98447204e-01**
 **6.78640538e-01**
**-9.27692517e-02**
**-3.68951686e-01**
**6.02420242e-01**
**-1.06486597e+00**
**1.74178344e+00]]**

b model intercept :  **[58.15728936]**

(b)(5 points) Visualize the actual and predicted death values yt and ^yt (for 8 <= t <= n). Display yt as a function of t and ^yt as a function of t on the same graph. You can use the library matplotlib.pyplot to plot.

(c) (5 points) Use a Gaussian to approximate the distribution of the errors yt - ^yt (for 8 <= t <= n). Report the mean and variance of this Gaussian.

```python
[89]: va = np.var(error,ddof = 1)
      std = np.sqrt(va)


      print ("Mean for Gaussian distribution of error : ", me)
      print ("Variance for Gaussian distribution of error : ", va)
```
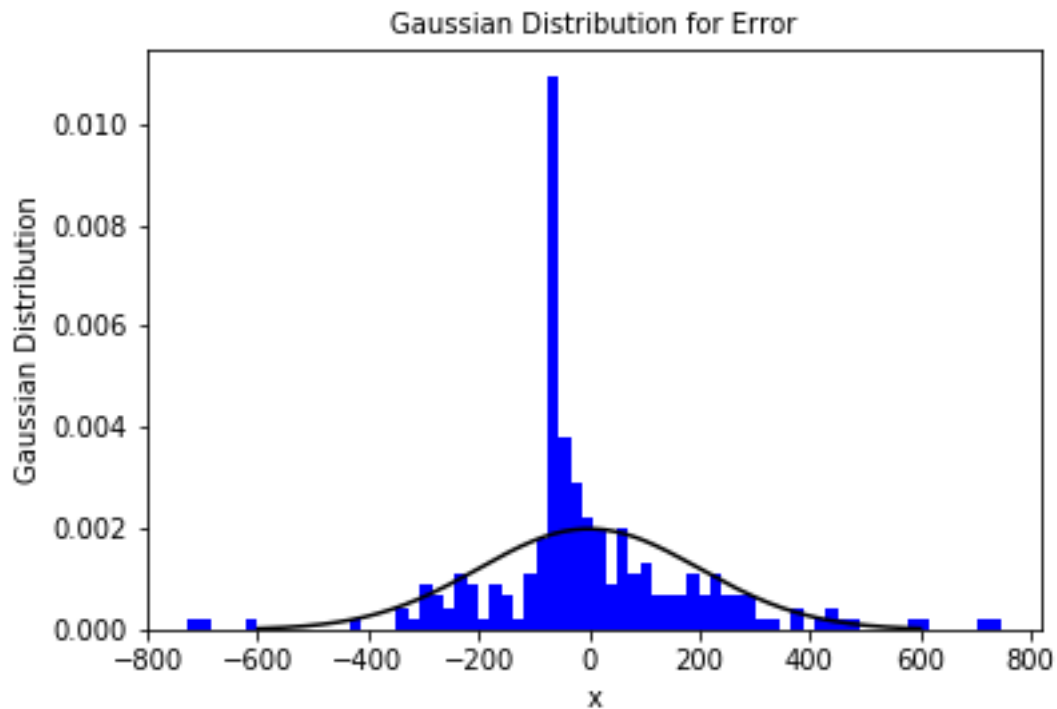
```
Mean for Gaussian distribution of error :  -1.2600558051323527e-11
Variance for Gaussian distribution of error :  39949.11578698032
```

Mean for Gaussian distribution of error :  **-1.2600558051323527e-11**
Variance for Gaussian distribution of error :  **39949.11578698032**

(d) (5 points) Use matplotlib.pyplot.hist to plot the distribution of yt - ^yt (for $8 <= t <= n$). On the same plot, plot the Gaussian function that approximates this distribution. Is Gaussian a good approximation for the distribution of the errors?



Gaussian Distribution for Error

- For approximation of error under Gaussian distribution, each error point should be independent and identically distributed which basically means that the error associated with a time value at a particular time point should be independent of error associated with a different time value at a different time point.
- If this case is satisfied then the Gaussian distribution is a good approximation for the distribution of the errors
- Because in the current given model, the error in the current point is dependent on the previous time point value, the IID (Independent and Identically distributed) assumption is violated
- Hence in our Case, it is not a good Idea to approximate distribution of error using Gaussian function

## Python code for reference :-

```python
def learn_reg_params(x, y) :
    y_copy = y[0, 7:]
    y_copy = np.atleast_2d(y_copy).T
    X = np.zeros((x.shape[1] - 7, 14))
    for t in range(7, x.shape[1]):
        X[t - 7,] = np.concatenate([x[0, t - 7:t], y[0, t - 7:t]])
    model = LinearRegression()
    model.fit(X, y_copy)
    y_pred = model.predict(X)
    print("Learned parameter w = ", model.coef )
    print("Learned parameter b = ", model.intercept )
    error = y_copy - y_pred
    me = np.mean(error)
    var = np.var(error, ddof=1)
    std = np.sqrt(var)

    fig, (ax3, ax4, ax5) = plt.subplots(1, 3, figsize=(12, 9))
    t = np.arange(8, 221)
    ax3.scatter(t, y_copy, marker='o', color='blue', linestyle=':')
    ax3.scatter(t, y_pred, marker='o', color='red', linestyle=':')
    ax3.set(xlabel='days', ylabel='Death')
    ax3.legend(['Actual', 'Predicted'])
    ax3.set_title("Actual and Predicted Death Values")

    xpts0 = np.linspace(me - 3 * std, me + 3 * std, 200)
    ax4.plot(xpts0, norm.pdf(xpts0, me, std), color='black')

    ax4.set title('Gaussian Distribution for Errors', fontsize=10)
    ax4.set(xlabel='X', ylabel='Error')

    xpts0 = np.linspace(me - 3 * std, me + 3 * std, 200)

    ax5.plot(xpts0, norm.pdf(xpts0, me, std), color='black')
    ax5.hist(error, bins=70, color='blue', density=True)
    ax5.set_title('Gaussian Distribution and Histogram plot for Error',
fontsize=10)

    ax5.set(xlabel='X', ylabel='Error')

    plt.show()
    print("w model coefficient : ", model.coef_)
    print("b model intercept : ", model.intercept_)
    return model.coef_,model.intercept_


df1 = pd.read_csv('covid19_time_series.csv', ',')
data1 = df1.to_numpy()
x = data1[:1,1:]
y = data1[1:,1:]
w,b = learn_reg_params(x, y)
```