

# CSC 521 Final Project

## Jonggoo Kang

### Problem1

```
In [1]: import pandas as pd
import numpy as np
import math
import random
import matplotlib.pyplot as plt
import seaborn as sns
import matplotlib inline
import scipy.stats as stats
from scipy.stats import expon
from termcolor import colored
```

It contains recorded accidents over the last 4 years.

The first column indicate where it happened (plant A or plant B).

The second column indicates the day (0 to 4 years ago).

The third column indicates the loss caused by the accidents in dollars.

```
In [2]: col = ['plant','days','losses']
df = pd.read_csv('accidents.csv', header = None, names = col)
print(df.shape)
print(df)
a = colored("By looking at pandas head() function\n", "blue", attrs = ['bold'])
b = colored("By looking at pandas describe() function\n", "blue", attrs = ['bold'])
print(a,df.head())
print(b, df.describe())

(291, 3)

By looking at pandas head() function
  plant  days  losses
0     A     1    3348
1     B     4     181
2     B     7     250
3     A    13    5446
4     A    30   38549

By looking at pandas describe() function
count      291.000000      291.000000
mean       497.202749      9189.185567
std        427.485991      21853.188417
min         1.000000         46.000000
25%        327.000000      987.000000
50%        706.000000      2886.000000
75%       1041.500000      8717.000000
max       1462.000000     272851.000000
```

```
In [3]: df.describe()

Out[3]:
```

	days	losses
count	291.000000	291.000000
mean	697.202749	9189.185567
std	427.485991	21853.188417
min	1.000000	46.000000
25%	327.000000	987.000000
50%	706.000000	2886.000000
75%	1041.500000	8717.000000
max	1462.000000	272851.000000

### Problem1-1. Without any simulation from the data, answer these questions:

1-1-a. The average number of accidents per year in plant A and B

1-1-b. The average loss per accident in plant A and plant B

1-1-c. The average loss in total per year in plant A and plant B

```
In [4]: # 1-1-a. Total number of events for plant A and B
cnt_plantA = df.plant.value_counts()[1]
cnt_plantB = df.plant.value_counts()[0]
# 1-1-a. The average number of accidents per year in plant A and B
avg_plantA_yr = cnt_plantA/4
avg_plantB_yr = cnt_plantB/4

# 1-1-b. The data frame for plant A and plant B
df_plantA = df[df.plant == 'A']
df_plantB = df[df.plant == 'B']
# 1-1-b. the average loss per accident in plant A and plant B
avg_lossA = sum(df_plantA.losses) / len(df_plantA.losses)
avg_lossB = sum(df_plantB.losses) / len(df_plantB.losses)

# 1-1-c. the average loss in total per year in plant A and plant B
avg_lossA_yr = sum(df_plantA.losses) / 4
avg_lossB_yr = sum(df_plantB.losses) / 4

color1_1 = colored(("There are 135 accidents in plant A and 156 accidents in plant B"), "blue",
    attrs = ["bold"])
color1_2 = colored(("1-1. The average number of accidents per year in plant A is"), "blue",at
    trs = ["bold"])
color1_3 = colored("and plant B is", "blue", attrs = ["bold"])
color1_4 = colored("1-2. The average loss per accident in plant A is", "blue", attrs = ["bol
    d"])
color1_5 = colored("1-3. The average loss in total per year in plant A is", "blue", attrs = [
    "bold"])

In [5]: def probl1():
    """This function show the results for probl1"""
    # 1-1-a. sol: total number of count for plant A and plant B are divided by 4 (since 4 year
    s)
    # 1-1-b. sol: sum of losses devided by the number of losses
    # 1-1-c. sol: sum of losses devided by 4 (since 4 years)

    print(color1_1)
    print(df.plant.value_counts())
    print(cnt_plantA, avg_plantA_yr, color_B, avg_plantB_yr)
    print("\tsol: each of total count number for plant A and plant B are divided by 4 (since 4
    years)\n")
    print(color1_3, avg_lossA, color_B, avg_lossB)
    print("\tsol: each total sum of losses for plant A and plant B are divided by the number of
    losses\n")
    print(color1_4, avg_lossA_yr, color_B, avg_lossB_yr)
    print("\tsol: each total sum of losses for plant A and plant B are divided by 4 (since 4 ye
    ars)")
```

therefore,

```
In [6]: probl1()

There are 135 accidents in plant A and 156 accidents in plant B
A      135
B      156
Name: plant, dtype: int64

1-1. The average number of accidents per year in plant A is 33.75 and plant B is 39.0
sol: each of total count number for plant A and plant B are divided by 4 (since 4 ye
ars)

1-2. The average loss per accident in plant A is 17470.155555555557 and plant B is 2022.9615
384615386
sol: each total sum of losses for plant A and plant B are divided by the number of l
osses

1-3. The average loss in total per year in plant A is 589617.75 and plant B is 78895.5
sol: each total sum of losses for plant A and plant B are divided by 4 (since 4 year
s)
```

### Problem1-2. Now assume the time interval between accidents is exponential and the natural log of a loss due to a single accident is a gaussian (aka the loss is lognormal)

1. Implement a simulate once that simulates one year of losses for both plants.
2. Running simulate many. What is the average yearly loss with a relative precision of 10%? Report the bootstrap errors in your result.
- 2-3. How much should the company budget to make sure that it can cover these losses in 90% of the simulated scenarios?

```
In [7]: def E(f,S): return float(sum(f(x) for x in S))/(len(S) or 1)
def mean(X): return E(lambda x:x, X)
def variance(X): return E(lambda x:x**2, X) - E(lambda x:x, X)**2
def sd(X): return math.sqrt(variance(X))
def resample(v):
    return (random.choice(v) for k in range(len(v)))
def bootstrap(scenarios, confidence):
    # len(scenarios) == 1000
    samples = []
    for x in range(100):
        samples.append(mean(resample(scenarios)))
    samples.sort()
    # len(samples) == 100
    i = int((100-confidence)/2)
    j = 99-i
    mu_plus = samples[j]
    mu_minus = samples[i]
    return mu_minus, mu_plus

In [8]: # Day different for plant A and B
diff_plantA = df.plant['days'].diff()
diff_plantB = df.plant['days'].diff()
### The result shows that first row has missing value
### this is because the first column has nothing to subtract
### Therefore, I will add 0 to the first row

# Add 0 to the first row for diff_plantA and diff_plantB
diff_plantA[0] = 0
diff_plantB[0] = 0

# mu for diff_plantA and diff_plantB
mu_diffA = diff_plantA.mean()
mu_diffB = diff_plantB.mean()
#####
# mu_loss for plant A and B
mu_lossA = math.log(avg_lossA)
mu_lossB = math.log(avg_lossB)

# sigma for plant A and B
sig_loglossA = sd(np.log(df[df.plant == 'A'].losses))
sig_loglossB = sd(np.log(df[df.plant == 'B'].losses))

# string variable that I will use the function in the next cell
A = 'A'
B = 'B'

2-1. Implement a simulate once that simulates one year of losses for both plants.
```

```
In [9]: def simulate_once(mu_loss, sig_loss, mu_diff):
    """This function results two outputs:
    the time interval between accidents
    and loss for a single accident"""
    # mu_loss : mu of losses for plant A or plant B
    # sig_loss: sigma of losses for plant A or plant B
    # mu_diff : mu of the time gap between days
    time_once_list = []
    loss_exp_list = []
    while sum(time_once_list) <= 365:
        time_once_list.append(random.exponential(1/mu_diff))
        loss_exp_list.append(np.exp(random.gauss(mu_loss, sig_loss)))

    return time_once_list, loss_exp_list

def simulate_once_result(mu_loss, sig_loss, mu_diff):
    """This function returns "simulate once" outputs for one year"""
    # mu_lossA = log of average losses for plant A
    # sigma_A = log of standard deviation losses for plant A
    # num = the length of dataframe of plant A for one year
    time, loss_exp = simulate_once(mu_loss, sig_loss, mu_diff)
    time_list = sorted(time)
    loss_exp_list = sorted(loss_exp)

    print(colored(("Time Interval"), "blue", attrs = ["bold"]))
    print(time_list)
    y = stats.expon.pdf(time_list, np.mean(time_list), np.std(time_list))
    plt.plot(time_list, y, label = 'Log Loss')
    plt.hist(time_list, normed = True)
    plt.show()

    print(colored(("Losses"), "blue", attrs = ["bold"]))
    print(loss_exp_list)
    y = stats.norm.pdf(loss_exp_list, np.mean(loss_exp_list), np.std(loss_exp_list))
    plt.plot(loss_exp_list, y, label = 'log Loss')
    plt.hist(loss_exp_list, normed = True)
    plt.show()
```

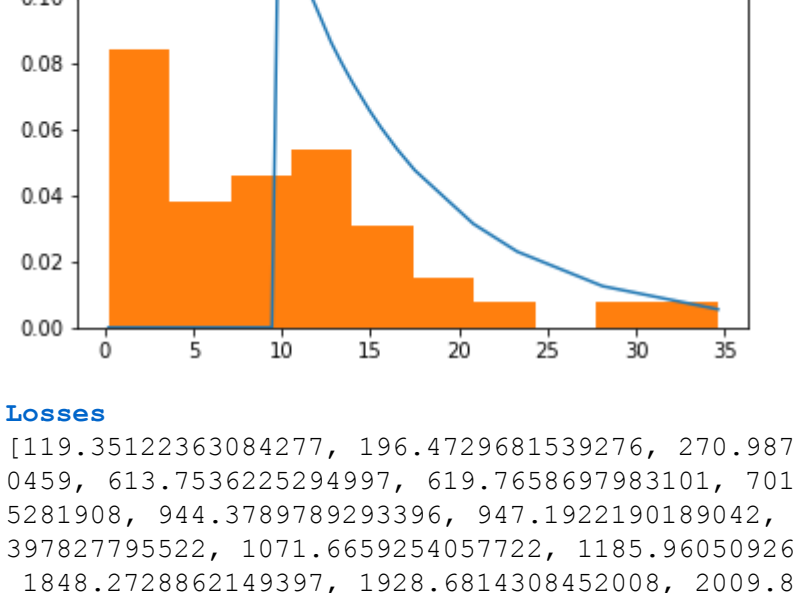
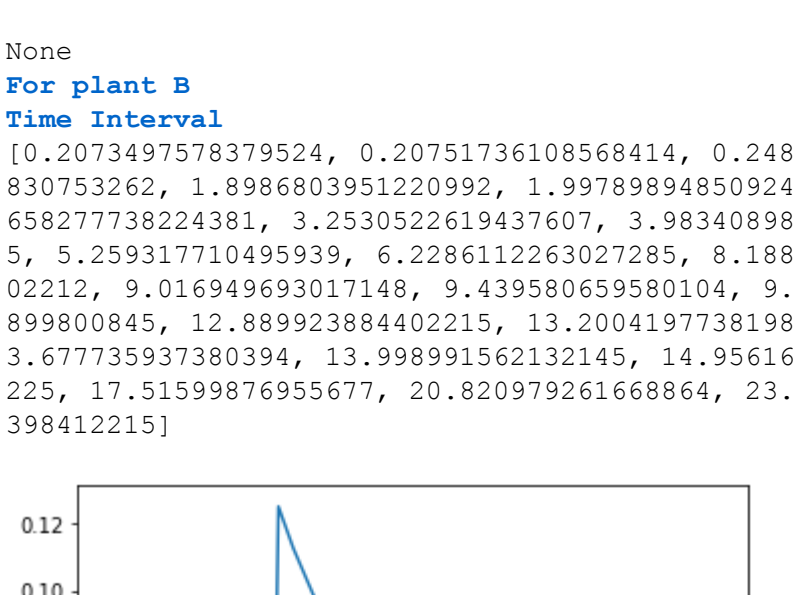
```
In [10]: color2_1 = colored(("2-1. Implement simulate once for one year of losses for both plants."),
    , "red", attrs = ["bold"])
color_forA = colored(("For plant A"), "blue", attrs = ["bold"])
color_forB = colored(("For plant B"), "blue", attrs = ["bold"])

def prob2_1():
    print(color2_1)
    print(color_forA)
    print(simulate_once_result(mu_lossA, sig_loglossA, mu_diffA))
    print(color_forB)
    print(simulate_once_result(mu_lossB, sig_loglossB, mu_diffB))
```

```
In [11]: prob2_1()

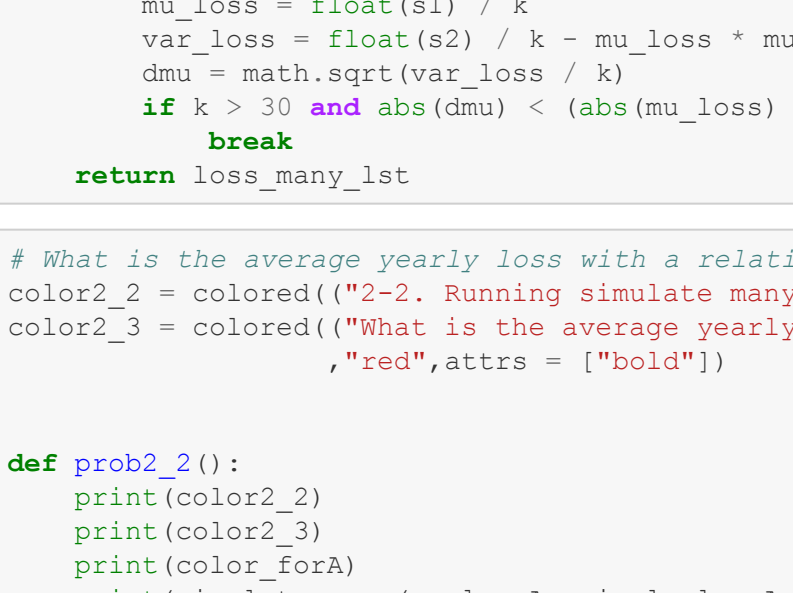
2-1. Implement simulate once for one year of losses for both plants.
For plant A
Time Interval
[0.783428456620441, 1.0595715339468897, 1.1205241092660052, 2.6909227944504632, 2.771867555
7037854, 3.9164488122279716, 981942.4577993618, 1148911.1820191944, 1491161.2110493677, 640399.33
17437744, 736567.843248427, 324098.7989561139, 727922.8648073215, 658158.4057101267, 69017
9.8575305466, 921001.5326630775, 1585486.831382626, 755832.3768965938, 1409991.452840144, 11
98874.549655371, 923403.7816401913, 1040309.27190371, 961174.9490853811, 586695.825214807
729468, 15.27508164639365, 18.76640445387046, 19.11886395952095, 47.2980838864583, 48.7228
936599263, 56.887681875099505]

Losses
[1724.4551142927123, 5197.725940181831, 5709.447941031977, 5858.264513146545, 7583.422224289
0459, 613.7536225294997, 619.765897983101, 701.9324615628221, 726.4439708071459, 932.668013
5281908, 944.3789789293396, 947.1922190189042, 966.3617453405707, 1026.3868769856147, 1051.9
39782775529, 1071.645924697722, 1185.96050262316, 1226.7393065621832, 1707.874064774426,
1468.272862149397, 1528.6814308452068, 2009.8146639971258, 2137.9237125440045, 2329.947143
759308, 2529.419049302382, 2661.83214977381981, 13.32227961685938, 13.388494313472552, 1
3.677535937803794, 13.989896332145, 14.79516946161685, 5282.17014234147, 55338.5731709681,
225.175159987955677, 20.82097926166864, 23.2948308930259975, 28.12048978380772, 34.636889
398412219]
```



```
None

For plant B
Time Interval
[0.2078349758379526, 0.20751736108569414, 0.24833909719763533, 0.5448463892328022, 1.2283652
830735262, 1.956803951220992, 1.9978989485050247, 2.2613378212334094, 2.3632490924301797, 2.
658277738224381, 3.2530522619437607, 3.983408986374237, 4.605862751659509, 4.995226152644618
5, 4.29317710495895, 6.256861226355664, 13915.91059746785, 1226.7393065621832, 1707.874064774426
02212, 9.0169496393017148, 9.439580659580104, 9.806401946946913, 10.650364384477312, 12.774099
89900845, 12.889923884402215, 13.20041977381981, 13.32227961685938, 13.388494313472552, 1
3.677535937803794, 13.989896332145, 14.79516946161685, 5282.17014234147, 55338.5731709681,
225.175159987955677, 20.82097926166864, 23.2948308930259975, 28.12048978380772, 34.636889
398412219]
```



```
None

2-2. Running simulate many. What is the average yearly loss with a relative precision of 10%?

In [13]: def simulate_many(mu, sigma, muTi, rp = 0.10, ns = 1000):
    loss_many_list = []
    time_many_list = []
    s1 = s2 = 0
    for k in range(1, ns):
        time, loss = simulate_once(mu, sigma, muTi)
        time_many_list.append(sum(time))
        loss_many_list.append(sum(loss))
        s1 = sum(loss)
        s2 += sum(losses) * sum(loss)
        mu_loss = float(s1) / k
        var_loss = float(s2) / k - mu_loss * mu_loss
        dmu = math.sqrt(var_loss / k)
        if k > 30 and abs(dmu) < (abs(mu_loss) * rp):
            return loss_many_list

In [13]: # What is the average yearly loss with a relative precision of 10%?
color2_2 = colored(("2-2. Running simulate many"), "red", attrs = ["bold"])
color2_3 = colored(("What is the average yearly loss with a relative precision of 10%?\n"),
    , "red", attrs = ["bold"])

def prob2_2():
    print(color2_2)
    print(color2_3)
    print(color_forA)
    print(simulate_many(mu_lossA, sig_loglossA, mu_diffA))
    print(color_forB)
    print(simulate_many(mu_lossB, sig_loglossB, mu_diffB))

In [14]: prob2_2()

2-2. Running simulate many,
What is the average yearly loss with a relative precision of 10%?
For plant A
[3721525.7398204645, 1025676.3214487784, 928666.2159123345, 1109950.9287458884, 1183638.2156
21488, 53261.52279716, 981942.4577993618, 1148911.1820191944, 1491161.2110493677, 640399.33
17437744, 736567.843248427, 324098.7989561139, 727922.8648073215, 658158.4057101267, 69017
9.8575305466, 921001.5326630775, 1585486.831382626, 755832.3768965938, 1409991.452840144, 11
98874.549655371, 923403.7816401913, 1040309.27190371, 961174.9490853811, 586695.825214807
729468, 15.27508164639365, 18.76640445387046, 19.11886395952095, 47.2980838864583, 48.7228
936599263, 56.887681875099505]
For plant B
[157945.31321776874, 236356.22243638142, 104715.6822680515, 98181.53094967852, 183856.750950
1861, 262216.64170457458, 83181.51123365705, 159575.12373180865, 186648.36143094688, 111194.
96620576625, 126214.40711250513, 162418.4121168092, 135964.2725400062, 178484.61450054045,
112624.07511851931, 155519.66652118257, 149457.83332123014, 193398.76865566542, 168647.5157
253148, 234123.90178912977, 63903.23209600954, 196179.5920665503, 136442.7638386264, 19411
7.61289856583, 221003.45367297914, 154219.95150759025, 186606.30483773185, 184236.865947817
8, 133991.6523030945, 200010.6902279068, 147587.29643957247]
```

2-3. Report the bootstrap errors in your result.

How much should the company budget to make sure that it can cover these losses in 90% of the simulated scenarios?

```
In [15]: color2_4 = colored(("2-3. Report the bootstrap errors in your result."), "red", attrs = ["bold"])
color2_5 = colored(("How much should the company budget to make sure that it can cover these
losses in 90% of the simulated scenarios?"), "red", attrs = ["bold"])

sim_A = simulate_many(mu_lossA, sig_loglossA, mu_diffA)
sim_B = simulate_many(mu_lossB, sig_loglossB, mu_diffB)
sim_list = [sim_A, sim_B]

def prob2_3():
    print(color2_4)
    for i in sim_list:
        if i == sim_A: print(color_forA)
        else: print(color_forB)
        for conf in list(range(90, 100, 10)):
            mu_minus, mu_plus = bootstrap(i, conf)
            print("\tThe confidence interval at (i) is from (i) to (i).format(conf, round(mu_minus,
4)), round(mu_plus, 4))
            if conf == 90:
                print(color2_5)
                print("Therefore, company can budget to make sure that it can cover the losses
in 90%")
                print("with a range from minimunly (i) to maximunly (i)\n".format(round(mu_minus,
4), round(mu_plus, 4)))

In [16]: prob2_3()

2-3. Report the bootstrap errors in your result.
For plant A
The confidence interval at 70 is from 1015236.6844 to 1133091.8762
The confidence interval at 80 is from 994722.0842 to 1165153.0
The confidence interval at 90 is from 963236.359 to 1164626.8115

How much should the company budget to make sure that it can cover these losses in 90% of the
simulated scenarios?
Therefore, company can budget to make sure that it can cover the losses in 90%
with a range from minimunly 963236.359 to maximunly 1164626.8115

For plant B
The confidence interval at 70 is from 151467.0932 to 169551.7947
The confidence interval at 80 is from 147389.2245 to 169070.6512
The confidence interval at 90 is from 148451.7496 to 171798.7104

How much should the company budget to make sure that it can cover these losses in 90% of the
simulated scenarios?
Therefore, company can budget to make sure that it can cover the losses in 90%
with a range from minimunly 148451.7496 to maximunly 171798.7104
```

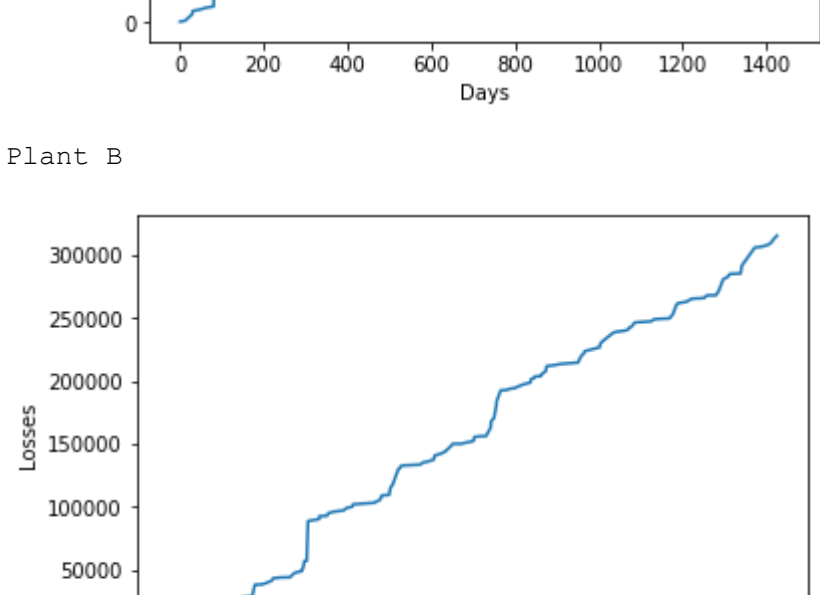
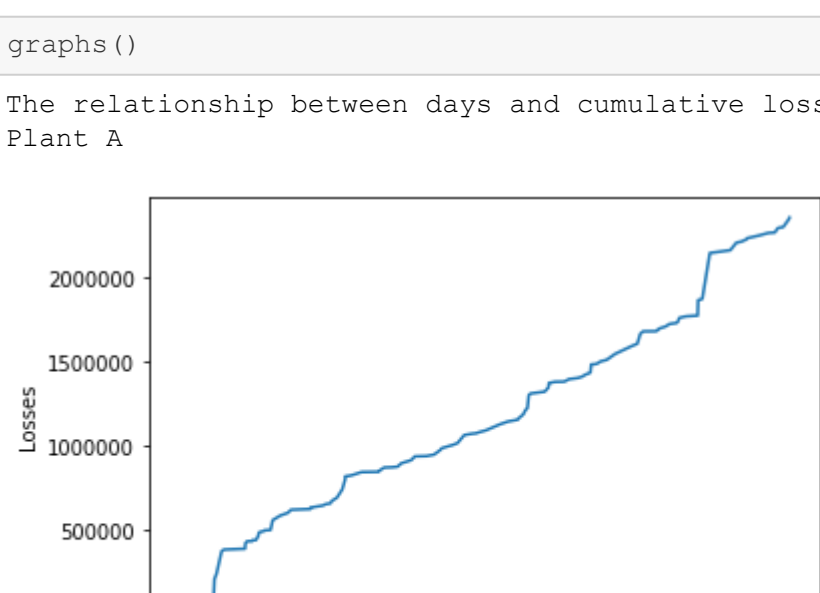
```
In [17]: def graphs():
    days_list = [df_plantA.days, df_plantB.days]
    plant_list = [df_plantA.losses, df_plantB.losses]
    print("The relationship between days and cumulative losses")
    for i, j, k in zip(days_list, plant_list, range(0, 2)):
        if k == 0:
            print("Plant A")
        else:
            print("Plant B")
        plt.plot(np.array(i), np.array(j).cumsum())
        plt.xlabel("Days")
        plt.ylabel("Losses")
        plt.show()

    timeA, logA = simulate_once(mu_lossA, sig_loglossA, mu_diffA)
    timeB, logB = simulate_once(mu_lossB, sig_loglossB, mu_diffB)
    time_list = [timeA, timeB]
    log_list = [logA, logB]

    print("The relationship between cumulative-simulated-time and cumulative-simulated-losses")
    for i, j, k in zip(time_list, log_list, range(0, 2)):
        if k == 0:
            print("Plant A")
        else: print("Plant B")
        plt.plot(np.array(i).cumsum(), np.array(j).cumsum())
        plt.xlabel("Days")
        plt.ylabel("Losses")
        plt.show()

In [18]: graphs()

The relationship between days and cumulative losses
Plant A
```



The relationship between cumulative-simulated-time and cumulative-simulated-losses  
Plant A