# Homework 1

## Problem 1
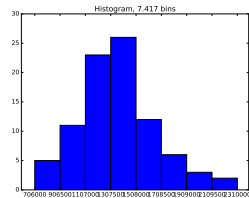


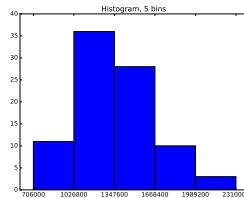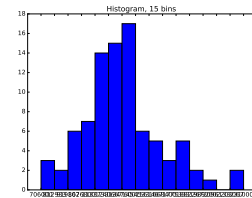Figure 1: Sturges          Figure 2: 5 bins              Figure 3: 15 bins
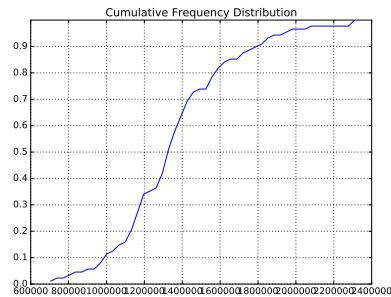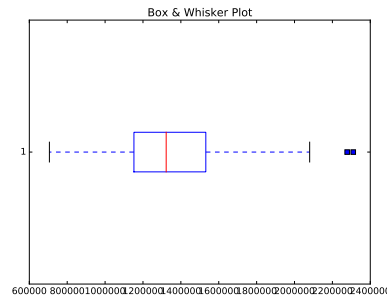


Figure 4: CDF                          Figure 5: Boxplot

The boxplot (figure 5) contains two outliers: 2,278,000 and 2,310,000, both above the 75th-quartile plus the maximum whisker length of $q75 + 1.5 * iqr$. The lower whisker only extends to 706,000 because this is within the maximum whisker length.

Here are the descriptive statistics (question 1c):

    mean, 1352204.54545
    median, 1322500.0
    var, 98171291013.6
    q25, 1152250.0
    q50, 1322500.0
    q75, 1531750.0
    iqr, 379500.0
    skew, 0.661353249864
    qskew, 0.102766798419
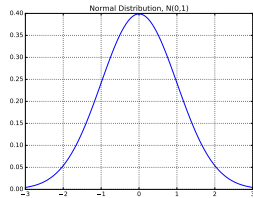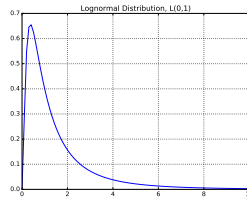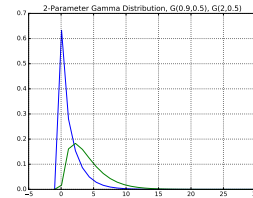
## Problem 2



Figure 6: Normal



Figure 7: Lognormal



Figure 8: Gamma

Below is the python code used to create the plots.

Listing 1: statswrap.py statistics wrapper

```python
from __future__ import division
import matplotlib.pyplot as plt
import numpy
import pandas
import os
import statsmodels.api as sm
import math
from scipy.stats import norm, lognorm, gamma

class Statistics:

    def __init__(self, file, col, output):
        self.filename = os.path.splitext(file)[0]
        self.file = pandas.read_csv(file)
        self.values = self.file[col]
        self.output = output

    def get_stats(self):
        """ mean, median, var, iqr, skew, qskew, quartiles =
        Statistics.get_stats()"""
        """ where 'quartiles' will provide a list [q25, q50, q75]
        """
        mean = self.values.mean() # [0]
        median = self.values.median() # [0]
        var = self.values.var() # [0]
        q25 = self.values.quantile(q=0.25) # [0]
        q50 = self.values.quantile() # [0]
        q75 = self.values.quantile(q=0.75) # [0]
        iqr = q75-q25
        qskew = ((q75-q50)-(q50-q25))/iqr
        skew = self.values.skew() # [0]
        return (mean, median, var, iqr, skew, qskew, [q25, q50,
        q75])
```

```python
31
32    def write_stats(self):
33      mean, median, var, iqr, skew, qskew, quartiles = self.
        get_stats()
34      with open(self.output + 'stats.txt', 'w') as f:
35        f.write('mean, ' + str(mean) + '\n')
36        f.write('median, ' + str(median) + '\n')
37        f.write('var, ' + str(var) + '\n')
38        f.write('q25, ' + str(quartiles[0]) + '\n')
39        f.write('q50, ' + str(quartiles[1]) + '\n')
40        f.write('q75, ' + str(quartiles[2]) + '\n')
41        f.write('iqr, ' + str(iqr) + '\n')
42        f.write('skew, ' + str(skew) + '\n')
43        f.write('qskew, ' + str(qskew))
44
45    def get_sturges(self):
46      n = len(self.file.index)
47      self.nbins = 1+3.3*math.log10(n)
48      smallest = self.values.min() # [0]
49      largest = self.values.max() # [0]
50      binwidth = (largest-smallest)/self.nbins
51      x = smallest.copy()
52      bins = []
53      while x < largest:
54        bins.append(x)
55        x+=binwidth
56      bins.append(x)
57      return bins
58
59    def plot_hist(self, bins):
60      q = self.values.values
61      figfile = ''
62      if bins == 'sturges':
63        figfile = (self.output + 'histogram_sturges.pdf')
64      if bins != 'sturges':
65        self.nbins = bins
66        figfile = (self.output + 'histogram_' + '{0}' + '.pdf').
        format(self.nbins)
67      fig, ax = plt.subplots()
68      counts, bins, patches = ax.hist(q, bins=bins)
69      ax.set_title(('Histogram, {0:.4g} bins').format(self.nbins
        ))
70      ax.set_xticks(bins)
71      fig.savefig(figfile)
72      # plt.show()
73
74    def cumul_freq_dist(self, filename='cumul_freq_dist.pdf'):
75      n = len(self.file.index)
76      sample = self.values.values.ravel()
77      ecdf = sm.distributions.ECDF(sample)
```

```
78    x = numpy.linspace(self.values.min(), self.values.max()) #
        [0]
79    y = ecdf(x)
80    fig, ax = plt.subplots()
81    ax.plot(x,y)
82    ax.set_title('Cumulative Frequency Distribution')
83    ax.set_yticks(numpy.arange(0,1,0.1))
84    plt.grid()
85    fig.savefig(self.output + filename)
86    # plt.show()
87
88  def boxplot(self, filename='boxplot.pdf'):
89    data = self.values.values.ravel()
90    fig, ax = plt.subplots()
91    ax.boxplot(data, 0, 'rs', 0)
92    ax.set_title('Box & Whisker Plot')
93    fig.savefig(self.output + filename)
94    # plt.show()
95
96
97  def norm_dist(self, mean=0, var=1, xstart=-3, xstop=3, xnum
      =100, filename='norm_dist.pdf'):
98    x_axis = numpy.linspace(xstart, xstop, xnum) # xnum steps
99    fig, ax = plt.subplots()
100   ax.plot(x_axis, norm.pdf(x_axis,mean,var))
101   plt.grid()
102   ax.set_title('Normal Distribution, N({0},{1})'.format(mean
      ,var))
103   fig.savefig(self.output + filename)
104   # plt.show()
105
106 def lognorm_dist(self, mean=0, var=1, xstart=0, xstop=10,
      xnum=100, filename='lognorm_dist.pdf'):
107   x_axis = numpy.linspace(xstart, xstop, xnum) # xnum steps
108   fig, ax = plt.subplots()
109   ax.plot(x_axis, lognorm.pdf(x_axis,var,0,numpy.exp(mean)))
110   plt.grid()
111   ax.set_title('Lognormal Distribution, L({0},{1})'.format(
      mean,var))
112   fig.savefig(self.output + filename)
113   # plt.show()
114
115 def gamma_dist(self, klist=[0.9,2], lam=0.5, xstart=-1,
      xstop=10, xnum=30, filename='gamma_dist.pdf'):
116   x_axis = numpy.linspace(xstart, xstop, xnum) # xnum steps
117   try:
118     scale = 1/lam
119   except ZeroDivisionError:
120     scale = 1/0.5
121     print 'ZeroDivisionError: gamma_dist() lambda value set
```

```
           to 0.5 as default '
122        fig , ax = plt.subplots ()
123        for k in klist:
124          ax.plot(x_axis, gamma.pdf(x_axis,k,scale=scale))
125        plt.grid ()
126        ax.set_title('2-Parameter Gamma Distribution, G({0},{1}),
        G({2},{3})'.format(klist[0],lam,klist[1],lam))
127        fig.savefig(self.output + filename)
128        # plt.show()
129
130 if __name__ == "__main__":
131
132     stats = Statistics('hw1given.csv', 'peak_streamflow_cfs', '
        probs_1_2/')
133
134     # Problem 1a
135     sturges = stats.get_sturges ()
136     stats.plot_hist('sturges')
137     stats.plot_hist(5)
138     stats.plot_hist(15)
139
140     # Problem 1b
141     stats.cumul_freq_dist ()
142
143     # Problem 1c
144     stats.write_stats ()
145
146     # Problem 1d
147     stats.boxplot ()
148
149     # Problem 2
150     stats.norm_dist(0,1,-3,3,100)
151     stats.lognorm_dist(0,1,0,10,100)
152     stats.gamma_dist([0.9,2],0.5,-1,30,30)
```