

## HW 08

This is a set of problems which focus on goodness-of-fit tests for various models of data.

**This homework is due Tuesday 6/26 at midnight. Hand it in with Lab 08 by the due date and time.**

**Be sure to select Run All from the Cell menu before submitting this notebook as your solution!**

```

In [469]: 1 # Jupyter notebook specific
          2 from IPython.display import Image
          3 from IPython.core.display import HTML
          4 from IPython.display import display_html
          5 from IPython.display import display
          6 from IPython.display import Math
          7 from IPython.display import Latex
          8 from IPython.display import HTML
          9
         10 # General useful imports
         11 import numpy as np
         12 from numpy import log, arange, linspace, mean, var, std, exp
         13 from scipy.special import comb
         14 from mpl_toolkits.mplot3d import Axes3D
         15 import matplotlib.pyplot as plt
         16 import matplotlib.mlab as mlab
         17 from numpy.random import seed, random, randint, uniform, choice, binomial, geometric, poisson,
         18 import math
         19 from collections import Counter
         20 import pandas as pd
         21 %matplotlib inline
         22
         23
         24 # Numpy basic stats functions
         25
         26 # https://docs.scipy.org/doc/numpy-1.13.0/reference/routines.statistics.html
         27
         28 X = [1,2,3]
         29
         30 # mean of a list
         31 mean(X)
         32
         33 # population variance
         34 var(X)
         35
         36 # sample variance
         37 var(X, ddof=1)
         38
         39 # population standard deviation
         40 std(X)
         41
         42 # sample standard deviation
         43 std(X, ddof=1)
         44
         45 # Scipy statistical functions
         46
         47 from scipy.stats import norm, expon, chi2
         48
         49 # https://docs.scipy.org/doc/scipy/reference/stats.html
         50
         51 # The following will work with norm replaced by expon or by chi2
         52
         53 # given random variable X (e.g., housing prices) normally distributed with mu = 60, sigma = 4
         54
         55
         56 # pdf of the distribution
         57 norm.pdf(x=50, loc=60, scale=40)
         58 chi2.pdf(x=10, df=5)
         59
         60 #a. Find P(X<50)
         61 norm.cdf(x=50, loc=60, scale=40) # 0.4012936743170763
         62 chi2.cdf(x=10, df=5)
         63
         64 #a. Find P(X>50)
         65 norm.sf(x=50, loc=60, scale=40) # 0.4012936743170763
         66 chi2.sf(x=10, df=5)
         67
         68 #c. Find P(60<X<80)
         69 norm.cdf(x=80, loc=60, scale=40) - norm.cdf(x=60, loc=60, scale=40)
         70
         71 #d. how much top most 5% expensive house cost at least? or find x where P(X>x) = 0.05

```

```

72 norm.isf(q=0.05,loc=60,scale=40)
73
74 #e. how much top most 5% cheapest house cost at least? or find x where P(X<x) = 0.05
75 norm.ppf(q=0.05,loc=60,scale=40)
76
77 #f give the endpoints of the range for the central alpha percent of the distribution
78 norm.interval(alpha=0.3, loc=60, scale=140)
79
80 # Same for exponential distribution
81 beta = 5    # mean of distribution
82
83
84
85 # Utility functions
86
87 # Round to 4 decimal places
88 def round4(x):
89     return round(float(x)+0.00000000001,4)
90
91 def round4List(X):
92     return [round(float(x)+0.00000000001,4) for x in X]
93
94 def show_histogram(freqs,bins,title="Data Histogram"):
95     plt.bar(bins[:-1],freqs,width=(bins[1]-bins[0]),align='edge',tick_label=round4List(bins[:-1]))
96     plt.title(title)
97     plt.ylabel("Frequency")
98     plt.xlabel("Bins [Lo .. Hi]")
99     plt.show()

```

## Problem Zero (Optional)

I suggest that you consolidate the code from lab so that you can take a set of observations and a theoretical distribution (given as a list of probabilities) and print out the results of the test. Use this in problems 1 and 2.

```

In [470]: 1 # Obs = set of observations; Probs = matching set of probabilities for the distribution to test
          2 # los = level of significance
          3
          4 def do_chi2_hypothesis_test(Obs,Probs,los, H0 = "The data follows the given distribution."):
          5     pass
          6
          7

```

Now test your function on Example B.1 on p.285 of Schaum's (same as we went over in class) and confirm that you implemented it correctly.

```

In [471]: 1 Obs = [80,65,70,85]
          2 n = len(Obs)
          3 Uniform = [1/n for k in range(n)]
          4 do_chi2_hypothesis_test(Obs,Uniform,0.05)

```

```

H0 = The data follows the given distribution.
p-value = 0.343030146138
Fail to reject H0 at the 0.05 level of significance.

```

## Problem One (Chi<sup>2</sup> Hypothesis Testing -- Normal Distribution)

For this problem, we would like you to apply the  $\chi^2$  test to data following the normal distribution. The basic ideas of the test are shown in Schaum's Appendix B, Example B.5 on p.288. Before starting this problem, please look at that example.

The  $\chi^2$  test is designed for discrete distributions, so to test a continuous distribution we need to convert it into a discrete one by "binning" the data into some number of intervals. The first part of this problem shows how to do this easily using the numpy histogram function.

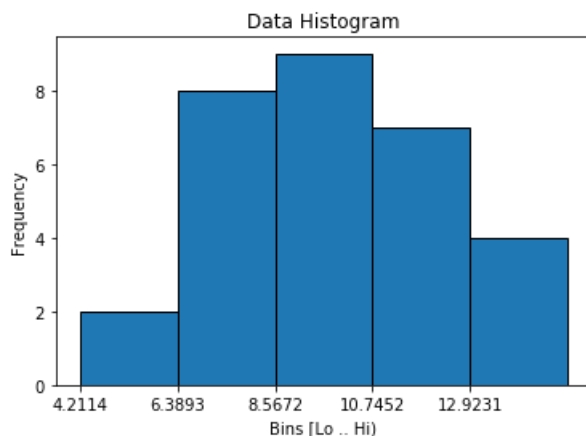
## Part A

In order to bin the observed and theoretical values, we will use the numpy function `histogram`, which takes a list of outcomes and counts the frequencies in `num_bins` bins spread evenly over the range; it returns a list of frequencies for each bin (= how many values from the data set ended up in that bin) and a list of the bin boundaries.

For example, we can separate a set of observed data points into 4 bins as follows:

```
In [472]: 1
2 X1a = [ 11.69690176,    6.94784341,    4.21136807,    9.21746662,    6.84108737,
3         13.08583604,   12.22008482,    7.37578756,    7.19696017,   11.58233656,
4         8.79227342,    9.77640417,   13.29730199,    9.92861583,    6.82766429,
5         10.55806493,   12.74073126,    8.13821793,   13.4815449,    9.67843858,
6         12.49634904,   12.834129,     9.66832643,    6.30537823,    9.78097709,
7         15.10102417,    7.77790787,    8.17102049,   10.51004163,   10.89602475]
8
9 obs,bins = np.histogram(X1a,bins=5)
10
11 print(obs)
12 print(bins)
13
14 show_histogram(obs,bins)
15
```

```
[2 8 9 7 4]
[ 4.21136807  6.38929929  8.56723051 10.74516173 12.92309295
15.10102417]
```



The bins are the ranges (values rounded to 4 decimal places):

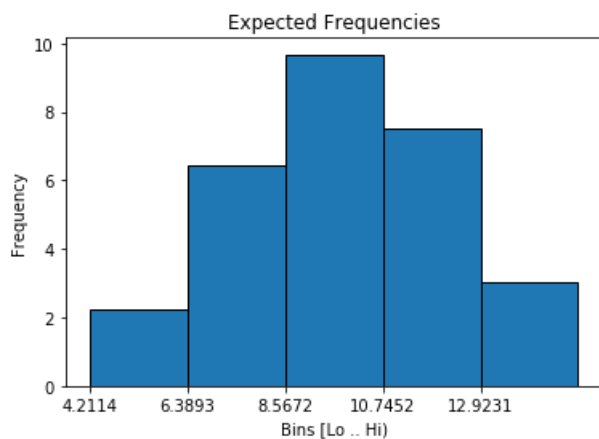
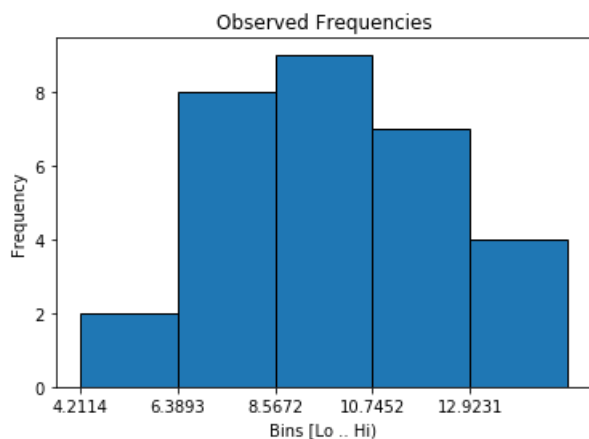
```
[4.2114 .. 6.3893), [6.3893 .. 8.5672), [8.5672 .. 10.7452), [ 10.7452 .. 12.9231), [12.
9231 .. 15.101].
```

Note that the variable `bins` is the list of boundaries between each of the bins, so its length is `num_bins + 1`.

Now complete the following function stub to take a data set (which is potentially from a normal distribution with mean  $\mu$  and standard deviation  $\sigma$ ) and bin the data into a set `Obs` of observed frequencies, and also produce a list `Probs` of expected probabilities for these bins if the data followed a normal distribution with the same mean and standard deviation as the data set.

```
In [473]: 1 def get_obs_and_normal_probs(X,num_bins):
2         pass
3
4         return (obs,probs,bins)
5
```

```
In [474]: 1 # Test your function on the data in X1a
2
3 Obs,Probs,bins = get_obs_and_normal_probs(X1a,5)
4
5 show_histogram(Obs,bins,"Observed Frequencies")
6 # convert probabilities into expected frequencies
7 Exp = [p*len(X1a) for p in Probs]
8 show_histogram(Exp,bins,"Expected Frequencies")
9
10 # Also: try 10 bins or 20 and see what happens!
11
12 # As a rough rule-of-thumb, you want at least approximately 5 values in each bin, so
13 # the largest number of bins would be # of data points / 5. In this case, that would give
14 # us an upper bound of 30/5 = 6 bins. That is why we chose 5.
15
```

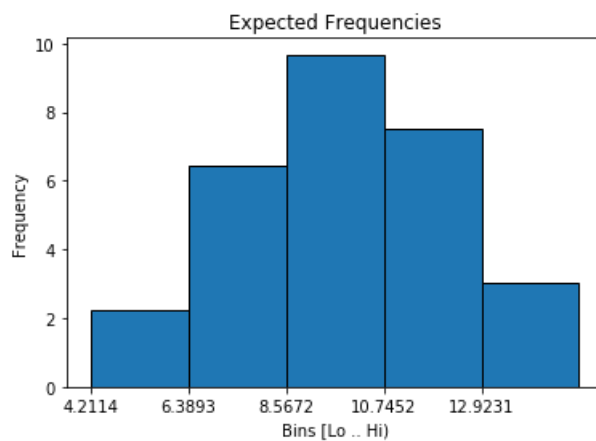
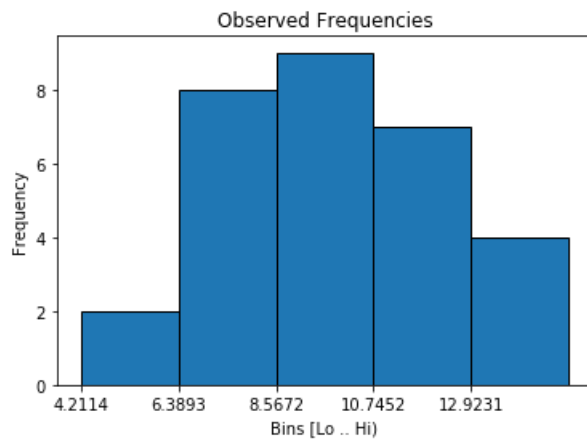


## Part B

Now complete the following stub to do the normal distribution test on the data in X1a, as in the lab, using the code from Problem 0 (if you wish).

```
In [475]: 1 def do_chi2_normal_test(X,num_bins,los):
2     pass
3
```

```
In [476]: 1 X1b = X1a
          2 do_chi2_normal_test(X1b,5,0.1)
          3
          4 # You should get a large p-value, which means that the data is very, very likely to be normal.
```



H0 = The data follows a normal distribution.  
p-value = 0.938696801217  
Fail to reject H0 at the 0.1 level of significance.

## Part C

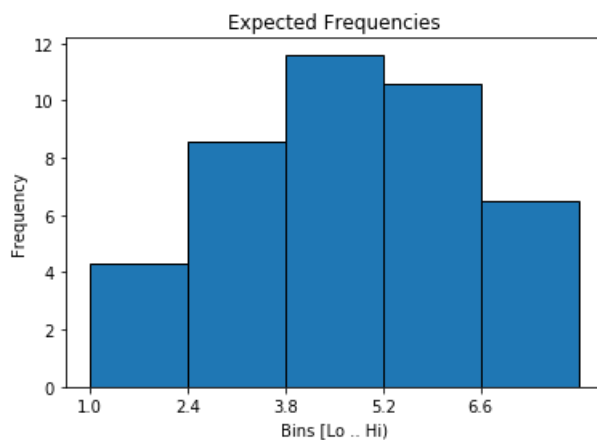
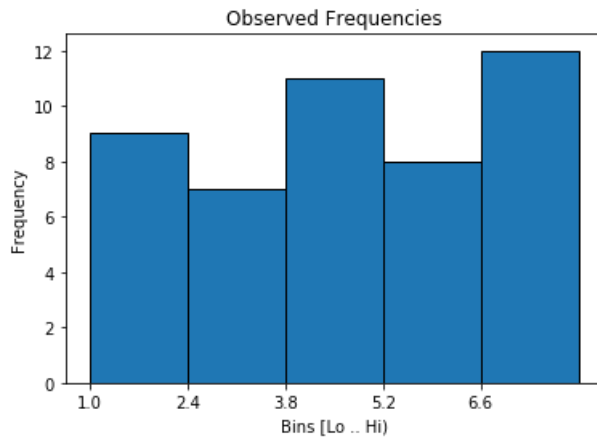
Now test the following data set using 5 bins and  $\alpha = 0.05$

In [477]:

```

1 X1c = [2,5,4,3,6,5,4,2,3,2,3,2,3,2,7,3,2,6,6,6,6,6,7,4,5,6,8,7,1,2,8,8,8,8,8,8,5,4,3,6,5
2 do_chi2_normal_test(X1c,5,0.05)
3

```



$H_0$  = The data follows a normal distribution.  
 p-value = 0.0290361830785  
 Reject  $H_0$  at the 0.05 level of significance.

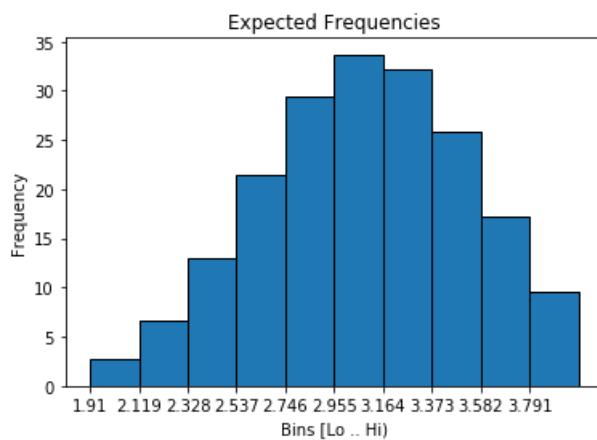
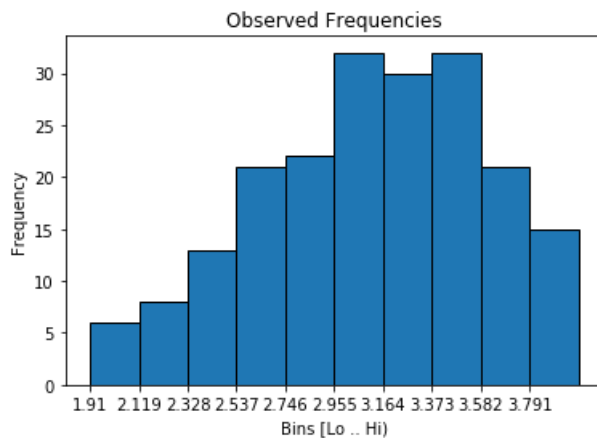
## Problem Two (Chi<sup>2</sup> Normal Test for Real Data)

In this problem you will test a number of real data sets to see how well a normal distribution fits the data.

### Part A

First we will test a small fraction of the BU GPA data to determine if the GPAs follow a normal distribution to a 0.05 level of significance.

```
In [478]: 1 # Test the file "BUDataSmall.csv" to see if the GPA data fits a normal distribution with los =  
2 # Use 10 bins.  
3  
4  
5
```

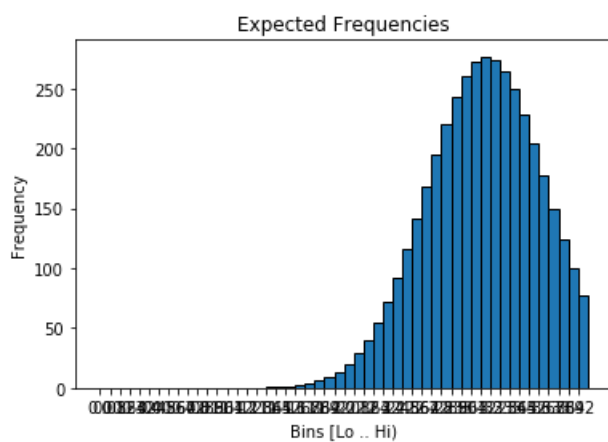
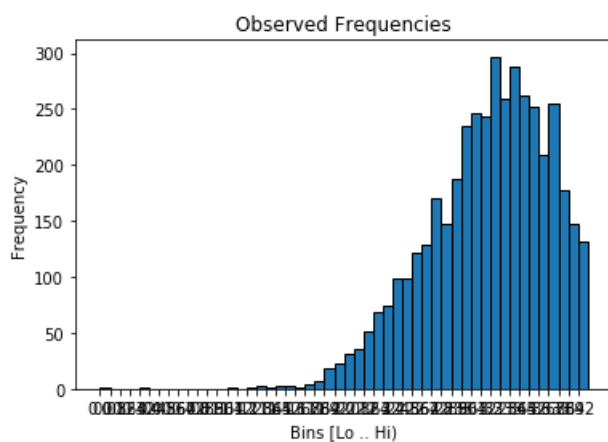


H0 = The data follows a normal distribution.  
p-value = 0.241844264514  
Fail to reject H0 at the 0.05 level of significance.

## Part B



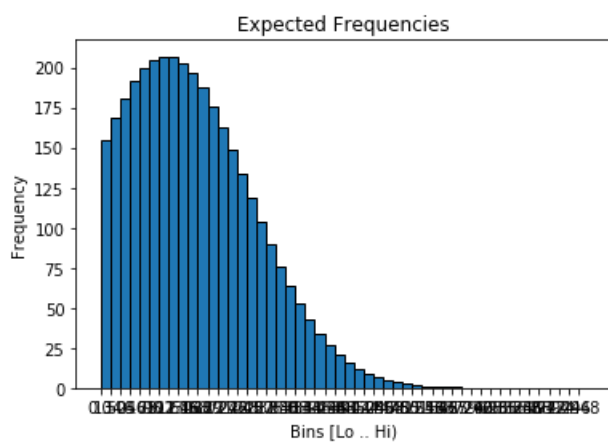
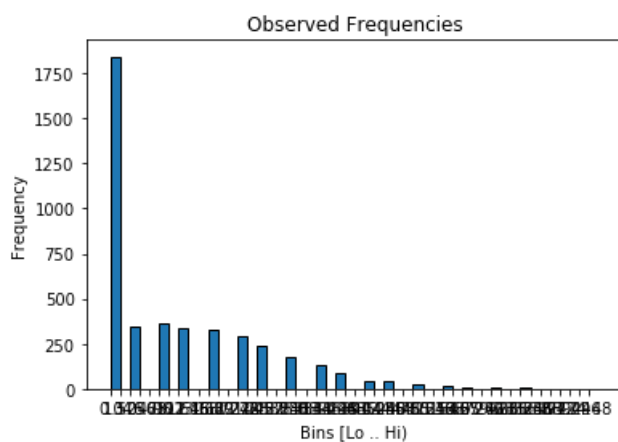
```
In [479]: 1 # Test the file "BUData.csv" to see if the whole data set of GPAs fits a normal distribution w
          2 # Use 50 bins.
          3
          4
          5
```



H0 = The data follows a normal distribution.  
 p-value = 0.0  
 Reject H0 at the 0.05 level of significance.

## Part C

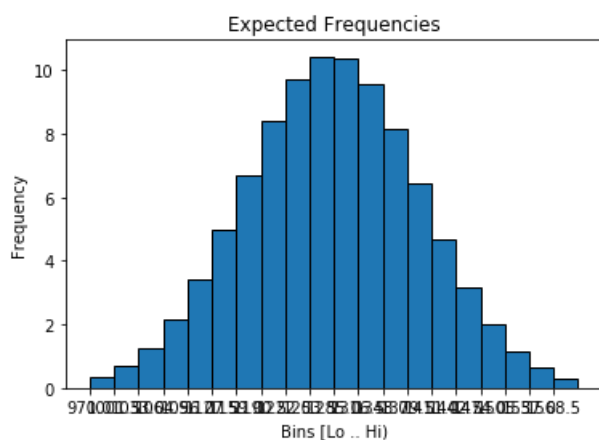
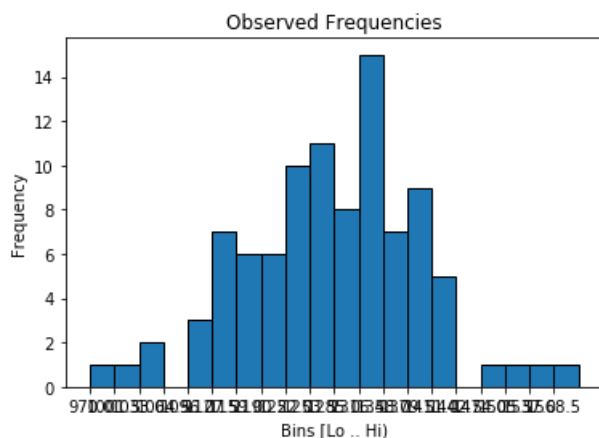
```
In [480]: 1 # Test the file "BUData.csv" to see if the whole data set of GPAs fits a normal distribution w
          2 # Use 50 bins.
          3
          4
          5
```



H0 = The data follows a normal distribution.  
 p-value = 0.0  
 Reject H0 at the 0.05 level of significance.

## Part D

```
In [481]: 1 # Test the file "StudentData3.csv" to see if the SAT_TOTAL data fits a normal distribution with
          2 # Use 20 bins.
          3
          4
          5
```



H0 = The data follows a normal distribution.  
 p-value = 0.660030575093  
 Fail to reject H0 at the 0.05 level of significance.

## Hypothesis Testing using Normal Probability Plots

Basic reference: <http://www.itl.nist.gov/div898/handbook/eda/section3/normprpl.htm>  
 (<http://www.itl.nist.gov/div898/handbook/eda/section3/normprpl.htm>)

From the link <http://www.itl.nist.gov/div898/handbook/eda/section3/eda3676.htm>  
 (<http://www.itl.nist.gov/div898/handbook/eda/section3/eda3676.htm>):

"The test statistic is the correlation coefficient  $\rho$  of the points that make up a normal probability plot. This test statistic is compared with the critical value below. If the test statistic is less than the tabulated value, the null hypothesis that the data came from a population with a normal distribution is rejected."

NOTE: This framework uses an inverse table (compared with the textbook definition of LOS), since you reject when your statistic  $\rho$  is less than the value in the table.

```

In [482]: 1 # Critical Values for the normal probability plot correlation coefficient test for
          2 # normality of a data set, from http://www.itl.nist.gov/div898/handbook/eda/section3/eda3676.
          3
          4 # First column is N = number of samples, second is for 0.01 LOS, third is for 0.05 LOS
          5
          6 CV = [[ 3 , 0.8687 , 0.8790],
          7 [ 4 , 0.8234 , 0.8666],
          8 [ 5 , 0.8240 , 0.8786],
          9 [ 6 , 0.8351 , 0.8880],
         10 [ 7 , 0.8474 , 0.8970],
         11 [ 8 , 0.8590 , 0.9043],
         12 [ 9 , 0.8689 , 0.9115],
         13 [ 10 , 0.8765 , 0.9173],
         14 [ 11 , 0.8838 , 0.9223],
         15 [ 12 , 0.8918 , 0.9267],
         16 [ 13 , 0.8974 , 0.9310],
         17 [ 14 , 0.9029 , 0.9343],
         18 [ 15 , 0.9080 , 0.9376],
         19 [ 16 , 0.9121 , 0.9405],
         20 [ 17 , 0.9160 , 0.9433],
         21 [ 18 , 0.9196 , 0.9452],
         22 [ 19 , 0.9230 , 0.9479],
         23 [ 20 , 0.9256 , 0.9498],
         24 [ 21 , 0.9285 , 0.9515],
         25 [ 22 , 0.9308 , 0.9535],
         26 [ 23 , 0.9334 , 0.9548],
         27 [ 24 , 0.9356 , 0.9564],
         28 [ 25 , 0.9370 , 0.9575],
         29 [ 26 , 0.9393 , 0.9590],
         30 [ 27 , 0.9413 , 0.9600],
         31 [ 28 , 0.9428 , 0.9615],
         32 [ 29 , 0.9441 , 0.9622],
         33 [ 30 , 0.9462 , 0.9634],
         34 [ 31 , 0.9476 , 0.9644],
         35 [ 32 , 0.9490 , 0.9652],
         36 [ 33 , 0.9505 , 0.9661],
         37 [ 34 , 0.9521 , 0.9671],
         38 [ 35 , 0.9530 , 0.9678],
         39 [ 36 , 0.9540 , 0.9686],
         40 [ 37 , 0.9551 , 0.9693],
         41 [ 38 , 0.9555 , 0.9700],
         42 [ 39 , 0.9568 , 0.9704],
         43 [ 40 , 0.9576 , 0.9712],
         44 [ 41 , 0.9589 , 0.9719],
         45 [ 42 , 0.9593 , 0.9723],
         46 [ 43 , 0.9609 , 0.9730],
         47 [ 44 , 0.9611 , 0.9734],
         48 [ 45 , 0.9620 , 0.9739],
         49 [ 46 , 0.9629 , 0.9744],
         50 [ 47 , 0.9637 , 0.9748],
         51 [ 48 , 0.9640 , 0.9753],
         52 [ 49 , 0.9643 , 0.9758],
         53 [ 50 , 0.9654 , 0.9761],
         54 [ 55 , 0.9683 , 0.9781],
         55 [ 60 , 0.9706 , 0.9797],
         56 [ 65 , 0.9723 , 0.9809],
         57 [ 70 , 0.9742 , 0.9822],
         58 [ 75 , 0.9758 , 0.9831],
         59 [ 80 , 0.9771 , 0.9841],
         60 [ 85 , 0.9784 , 0.9850],
         61 [ 90 , 0.9797 , 0.9857],
         62 [ 95 , 0.9804 , 0.9864],
         63 [ 100 , 0.9814 , 0.9869],
         64 [ 110 , 0.9830 , 0.9881],
         65 [ 120 , 0.9841 , 0.9889],
         66 [ 130 , 0.9854 , 0.9897],
         67 [ 140 , 0.9865 , 0.9904],
         68 [ 150 , 0.9871 , 0.9909],
         69 [ 160 , 0.9879 , 0.9915],
         70 [ 170 , 0.9887 , 0.9919],
         71 [ 180 , 0.9891 , 0.9923],

```

```

72 [ 190 , 0.9897 , 0.9927],
73 [ 200 , 0.9903 , 0.9930],
74 [ 210 , 0.9907 , 0.9933],
75 [ 220 , 0.9910 , 0.9936],
76 [ 230 , 0.9914 , 0.9939],
77 [ 240 , 0.9917 , 0.9941],
78 [ 250 , 0.9921 , 0.9943],
79 [ 260 , 0.9924 , 0.9945],
80 [ 270 , 0.9926 , 0.9947],
81 [ 280 , 0.9929 , 0.9949],
82 [ 290 , 0.9931 , 0.9951],
83 [ 300 , 0.9933 , 0.9952],
84 [ 310 , 0.9936 , 0.9954],
85 [ 320 , 0.9937 , 0.9955],
86 [ 330 , 0.9939 , 0.9956],
87 [ 340 , 0.9941 , 0.9957],
88 [ 350 , 0.9942 , 0.9958],
89 [ 360 , 0.9944 , 0.9959],
90 [ 370 , 0.9945 , 0.9960],
91 [ 380 , 0.9947 , 0.9961],
92 [ 390 , 0.9948 , 0.9962],
93 [ 400 , 0.9949 , 0.9963],
94 [ 410 , 0.9950 , 0.9964],
95 [ 420 , 0.9951 , 0.9965],
96 [ 430 , 0.9953 , 0.9966],
97 [ 440 , 0.9954 , 0.9966],
98 [ 450 , 0.9954 , 0.9967],
99 [ 460 , 0.9955 , 0.9968],
100 [ 470 , 0.9956 , 0.9968],
101 [ 480 , 0.9957 , 0.9969],
102 [ 490 , 0.9958 , 0.9969],
103 [ 500 , 0.9959 , 0.9970],
104 [ 525 , 0.9961 , 0.9972],
105 [ 550 , 0.9963 , 0.9973],
106 [ 575 , 0.9964 , 0.9974],
107 [ 600 , 0.9965 , 0.9975],
108 [ 625 , 0.9967 , 0.9976],
109 [ 650 , 0.9968 , 0.9977],
110 [ 675 , 0.9969 , 0.9977],
111 [ 700 , 0.9970 , 0.9978],
112 [ 725 , 0.9971 , 0.9979],
113 [ 750 , 0.9972 , 0.9980],
114 [ 775 , 0.9973 , 0.9980],
115 [ 800 , 0.9974 , 0.9981],
116 [ 825 , 0.9975 , 0.9981],
117 [ 850 , 0.9975 , 0.9982],
118 [ 875 , 0.9976 , 0.9982],
119 [ 900 , 0.9977 , 0.9983],
120 [ 925 , 0.9977 , 0.9983],
121 [ 950 , 0.9978 , 0.9984],
122 [ 975 , 0.9978 , 0.9984],
123 [1000 , 0.9979 , 0.9984]]
124
125
126 # Approximation to uniform order statistic medians due to Filliben:
127 # http://www1.cmc.edu/pages/faculty/MONeill/Math152/Handouts/filliben.pdf
128
129 def U(i,n):
130     if(i == n):
131         return 0.5**(1/n)
132     elif(i == 1):
133         return 1 - 0.5**(1/n)
134     else:
135         return (i - 0.3175)/(n + 0.365)
136
137 # Normal Order Statistic Medians
138
139 def NOSM(i,n,mu,sigma):
140     return norm.ppf(q=U(i,n),loc=mu,scale=sigma)
141
142 def getCVIndex(n,lo,hi):
143     if(lo >= hi):

```

```

144         return lo
145     mid = (lo + hi)//2
146     if(CV[mid][0] == n):
147         return mid
148     elif(CV[mid][0] < n):
149         return getCVIndex(n,mid+1,hi)
150     else:
151         return getCVIndex(n,lo,mid-1)
152
153 # returns the critical value or the next lower value if not in the table
154
155 def getCriticalValue(N,los):
156     N = min(N,1000)
157     i = getCVIndex(N,0,len(CV))
158     if(los == 0.01):
159         return CV[i][1]
160     else:
161         return CV[i][2]
162

```

## Problem Three (Tests using Normal Probability Plots)

In this problem you will write a function similar to the one above, but this time performing a test for fitness of the normal distribution using a normal probability plot. A big difference from the  $\chi^2$  test is that we can do this directly on the data and no "binning" of values is necessary. Thus, it is usually a more sensitive test than  $\chi^2$  for normal data.

The cell above contains all the code you need to do this test, supposing that your test data is in the list Y, here the outline of what you need to do in filling in the stub in the next cell:

```

N = len(Y)

los = level of significance of test.

```

The predicted values on the x-axis are produced as follows:

```

X = [NOSM(i,N,mean(X),std(X,ddof=1)) for i in range(1,N+1)]      # we use the sample s
td dev

```

and the test is performed as follows:

```

If rho(X,Y) < getCriticalValue(N,los) then reject, else fail to reject.

```

To display the probability plot, sort Y and draw a scatter plot of X against the sorted Y. To draw the linear prediction, plot X against X.

## Part One

Complete the following function stub and test it on the data X1a with los = 0.01.

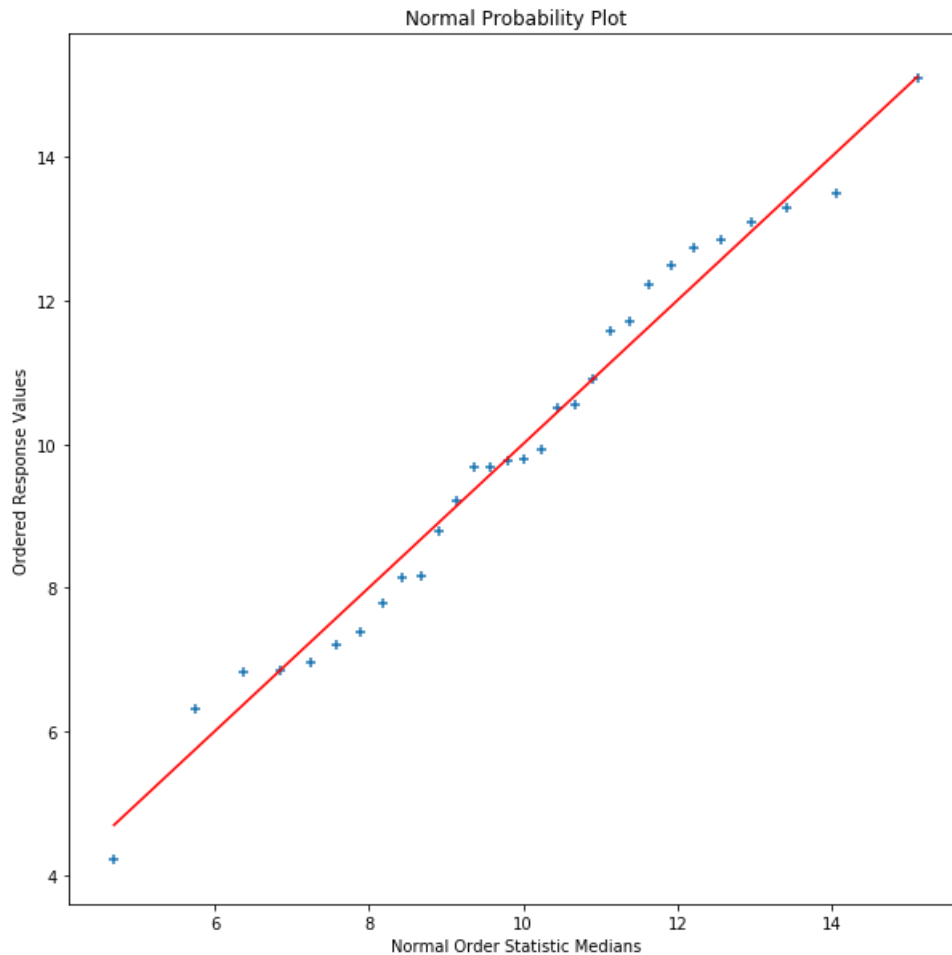
You must draw the probability plot and then print out the result of the test.

```

In [483]: 1 def do_pplot_normal_test(Y,los):
          2     pass
          3
          4
          5

```

```
In [484]: 1 # Test X1a with los = 0.01
          2
          3 do_pplot_normal_test(X1a,0.01)
```



```
rho          = 0.9909
critical value = 0.9462
```

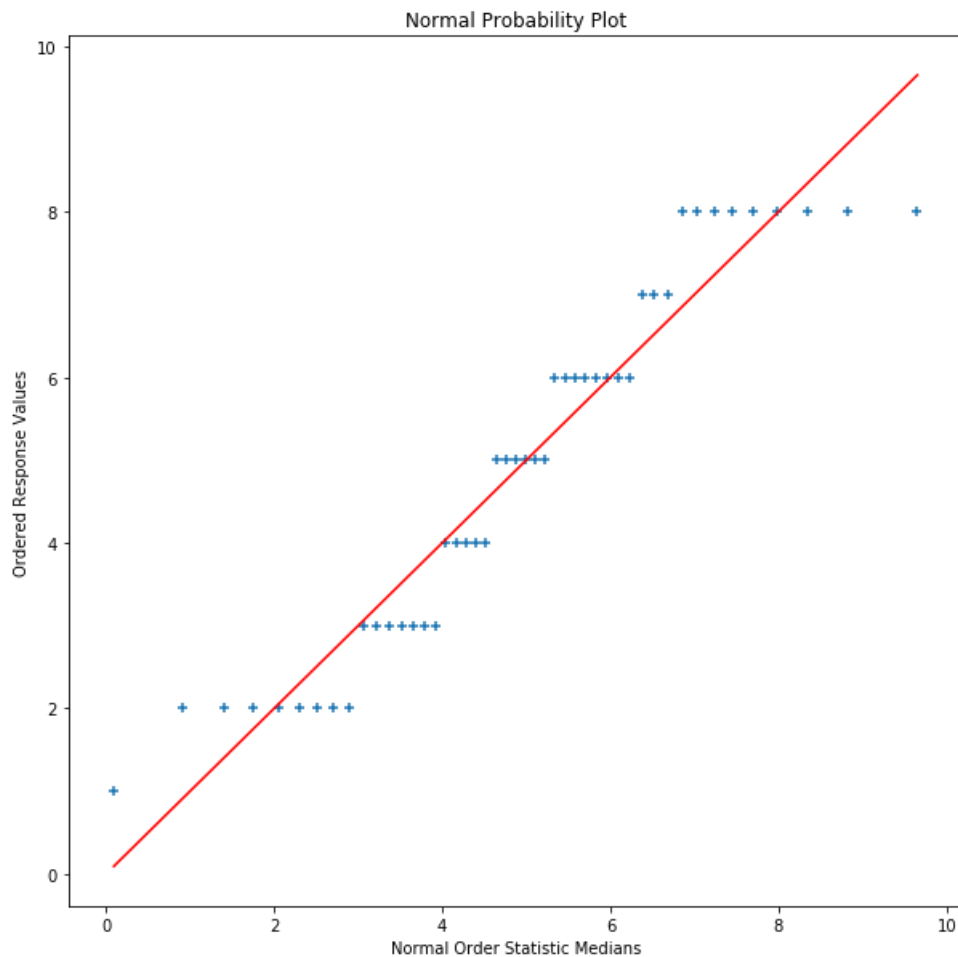
H0 = The samples come from a normal distribution.

Fail to reject at a 0.01 level of significance.

## Part B

Now we would like you to test the data in X1c. It will fail, even at a 0.1 level of significance, because it is fundamentally a discrete distribution (the outcomes are all integers), and the probability plot expects a continuous distribution. Hence, in this case the  $\chi^2$  test is more appropriate.

```
In [485]: 1 # Test  $X_{1c}$  with  $\alpha = 0.1$   
2  
3
```



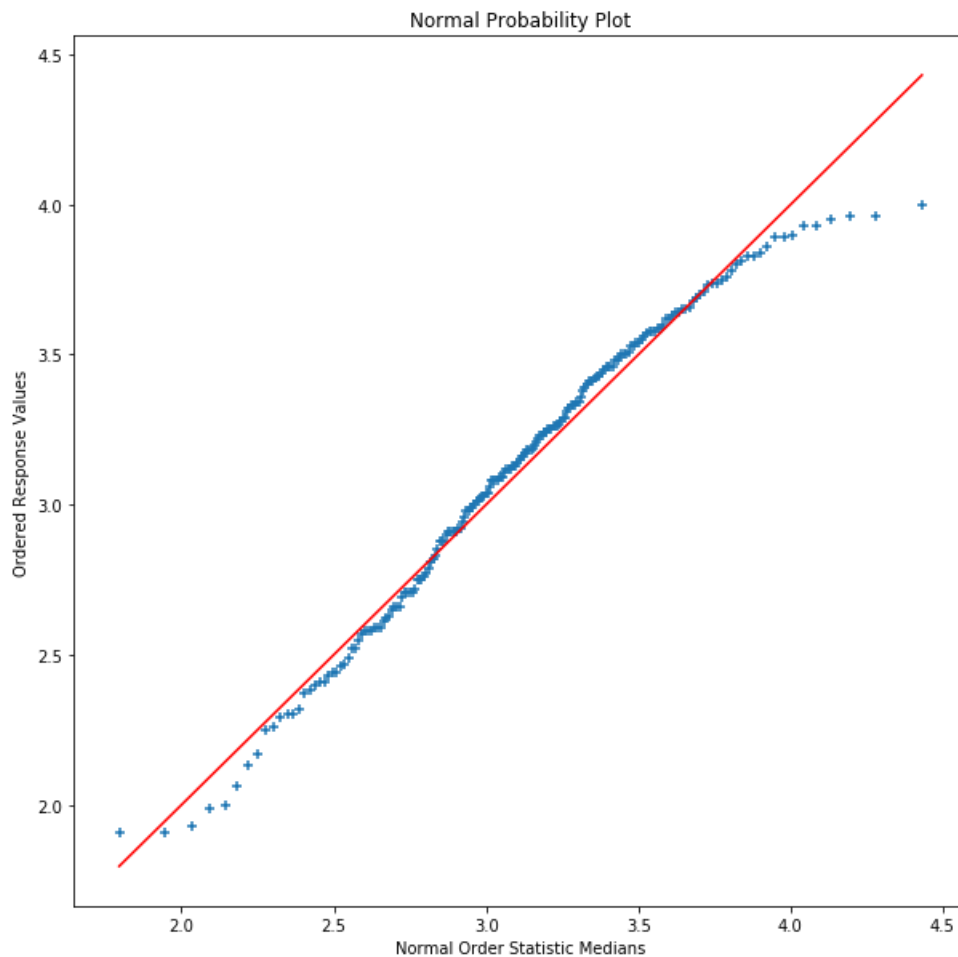
```
rho = 0.9637  
critical value = 0.9748
```

$H_0$  = The samples come from a normal distribution.

Reject at a 0.1 level of significance.



```
In [486]: 1 # Test X2a with los = 0.05  
2  
3
```

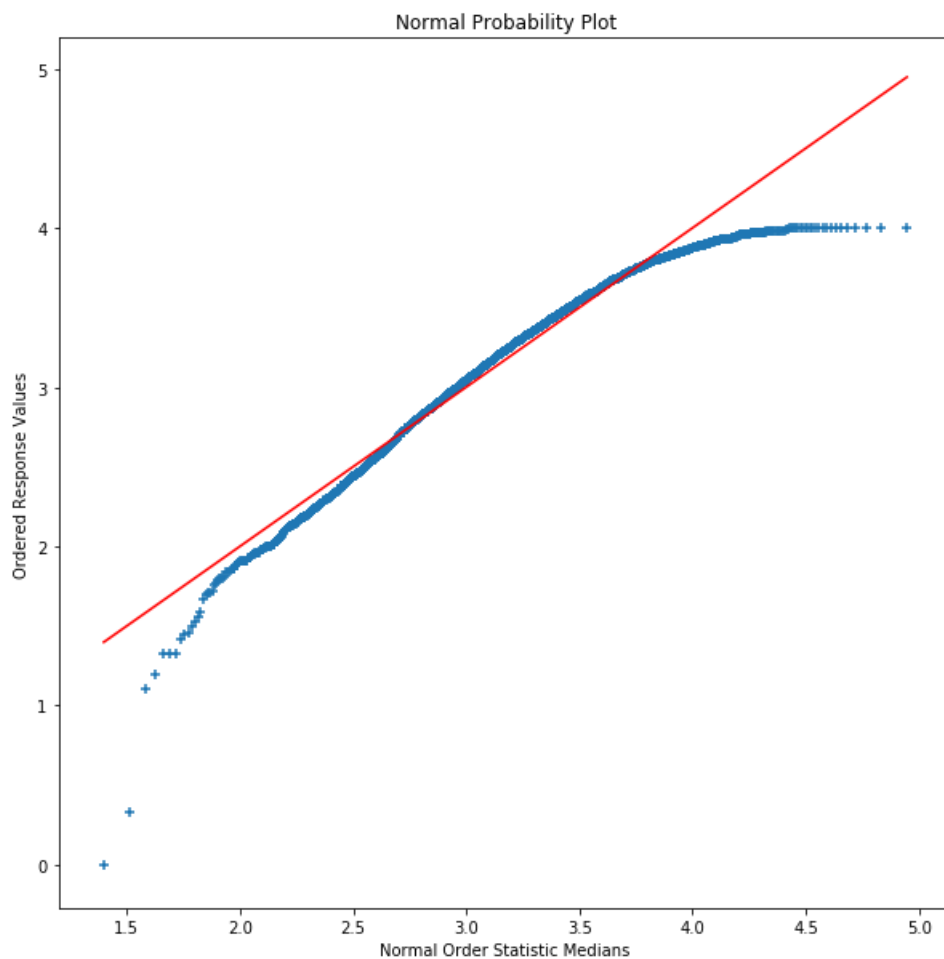


```
rho          = 0.9911  
critical value = 0.993
```

H0 = The samples come from a normal distribution.

Reject at a 0.05 level of significance.

```
In [487]: 1 # Test X2b with los = 0.05  
2  
3
```

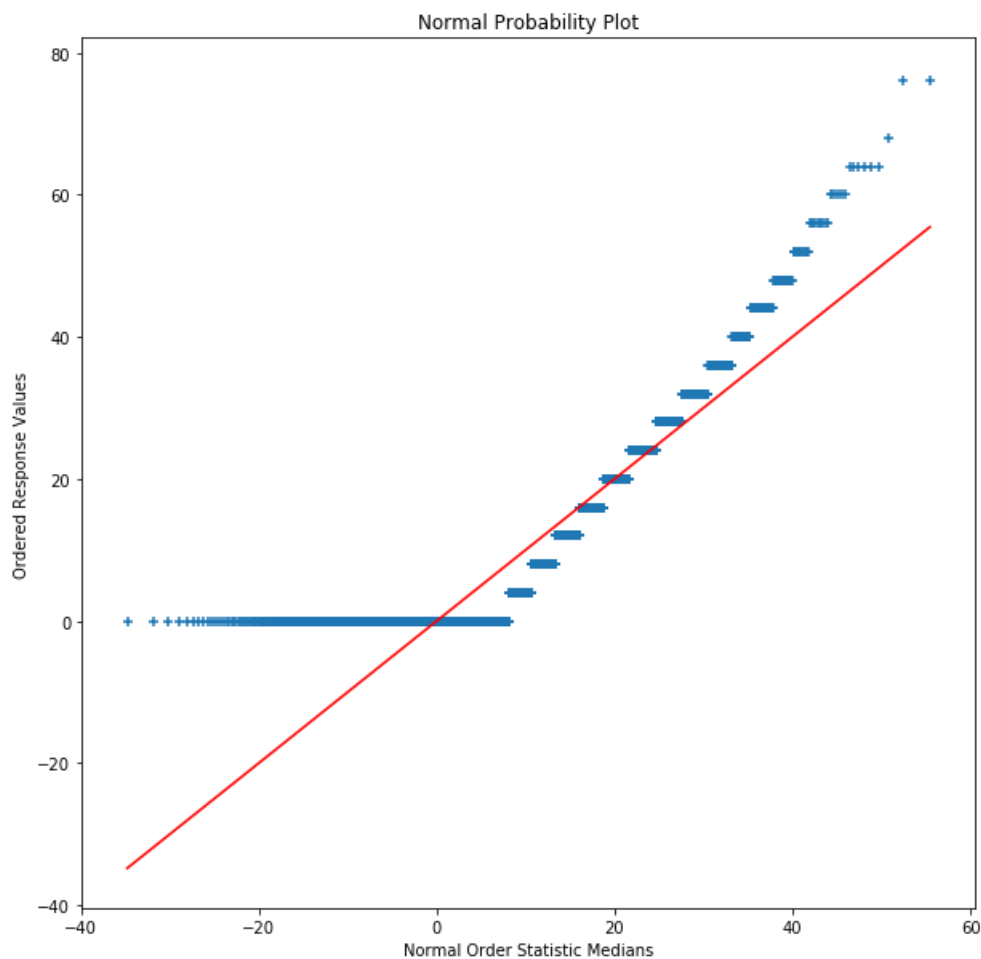


```
rho          = 0.9833  
critical value = 0.9984
```

H0 = The samples come from a normal distribution.

Reject at a 0.05 level of significance.

```
In [488]: 1 # Test X2c with los = 0.05  
2  
3
```

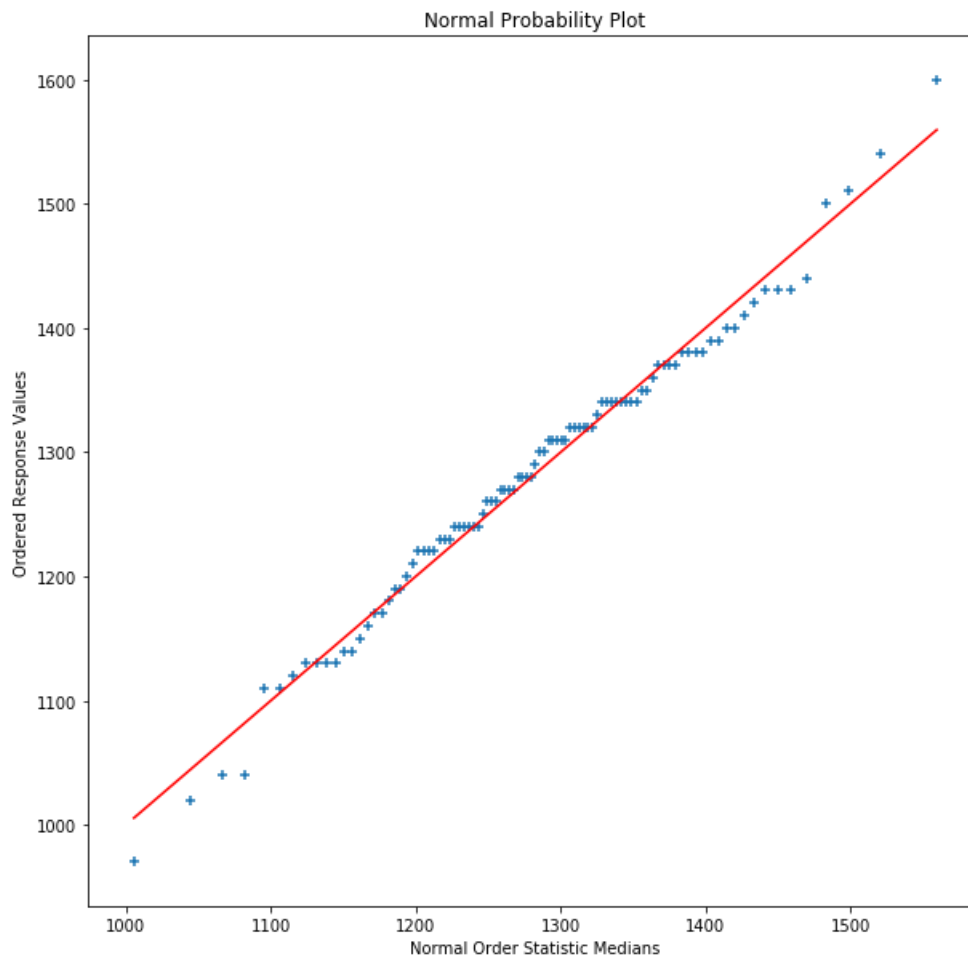


```
rho          = 0.9005  
critical value = 0.9984
```

H0 = The samples come from a normal distribution.

Reject at a 0.5 level of significance.

```
In [489]: 1 # Test X2d with los = 0.05  
          2  
          3
```



```
rho          = 0.993  
critical value = 0.9864
```

$H_0$  = The samples come from a normal distribution.

Fail to reject at a 0.05 level of significance.