

## 16. Statistics Basic

```
In [1]: 1 %matplotlib inline
        2 import networkx as nx
        3 import numpy as np
        4 import matplotlib.pyplot as plt
        5 import matplotlib.axes as axes
        6 import math
        7 import scipy.optimize as opt
        8 import scipy.stats as stats
        9 import netlab as nl
```

### Confidence Interval

```
In [2]: 1 a = np.random.normal(loc=120, scale=1.5, size=25)
        2 print (a)

[120.52360591 120.65548198 119.22991278 118.09107671 119.13649566
 121.79530688 118.7785582 121.49600567 120.56963384 118.92190098
 118.7548623 118.93454761 118.30938898 121.03437658 119.7735088
 119.33476249 123.31843362 116.57135298 120.27086975 119.96745255
 119.69261621 120.64636377 120.39631637 119.28467209 119.16285545]
```

```
In [3]: 1 n = len(a)
        2 mean = np.mean(a)
        3 sigma = np.std(a)
        4 print ('(%.3f, %.3f, %d)' % (mean, sigma, n))

(119.786, 1.339, 25)
```

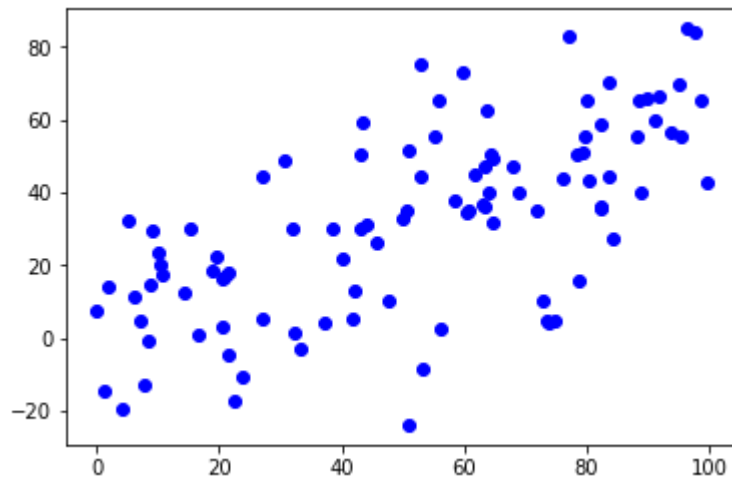
```
In [4]: 1 confidence_level = 95
        2 confidence_coeff = confidence_level / 100.0
        3 confidence_interval = stats.norm.interval(confidence_coeff, loc=mean, scale=
        4 print ('%.3f %.3f' % (confidence_interval[0], confidence_interval[1]))

119.261 120.311
```

### Simple Linear Regression

```
In [5]: 1 data = np.loadtxt('Data\y11.txt')
        2 x = data[:,0]
        3 y = data[:,1]
        4 plt.plot(x, y, 'bo')
```

Out[5]: [<matplotlib.lines.Line2D at 0x250c53ad080>]



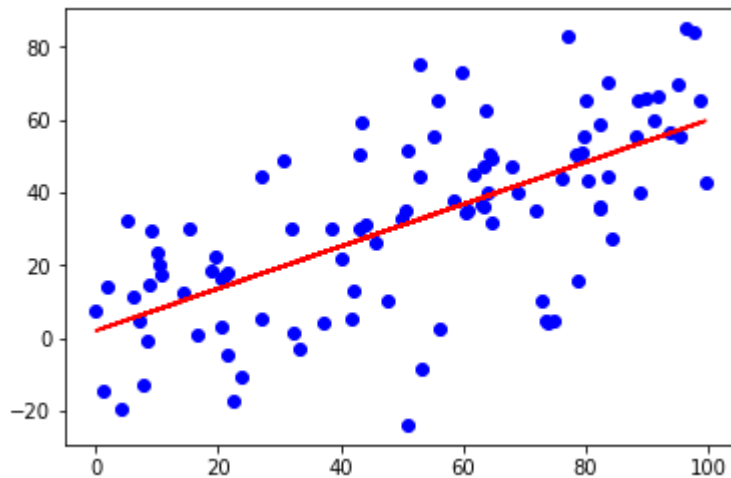
## Linear function

$$f(x) = ax + b$$

```
In [6]: 1 a, b, corr, p, stderr = stats.linregress(x,y)
2 z = a * x + b
3 print ('f(x) = %f x %c %f' % (a, '+' if b >= 0 else '-', abs(b)))
4 print ('corr      = %f' % corr)
5 print ('p        = %f' % p)
6 print ('std err = %f' % stderr)
7 plt.plot(x, y, 'bo')
8 plt.plot(x, z, 'r')
```

```
f(x) = 0.579406 x + 1.958615
corr   = 0.660108
p      = 0.000000
std err = 0.066603
```

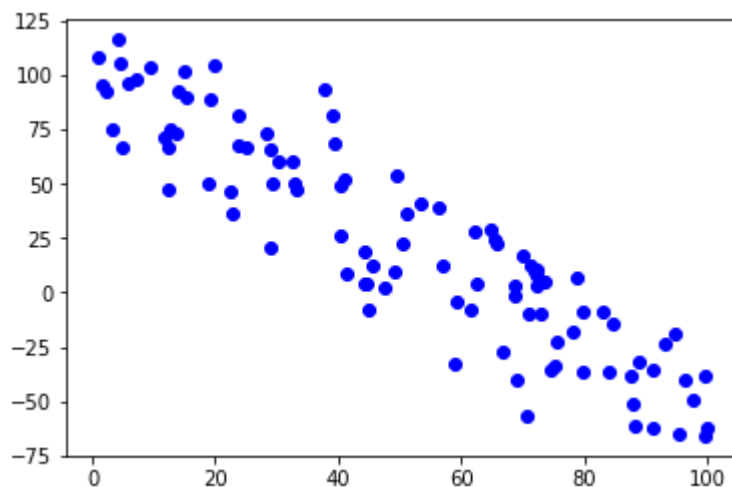
Out[6]: [<matplotlib.lines.Line2D at 0x250c540b5f8>]



```
In [7]: 1
2 stats.linregress??
```

```
In [8]: 1 data = np.loadtxt('Data\y12.txt')
2 x = data[:,0]
3 y = data[:,1]
4 plt.plot(x, y, 'bo')
```

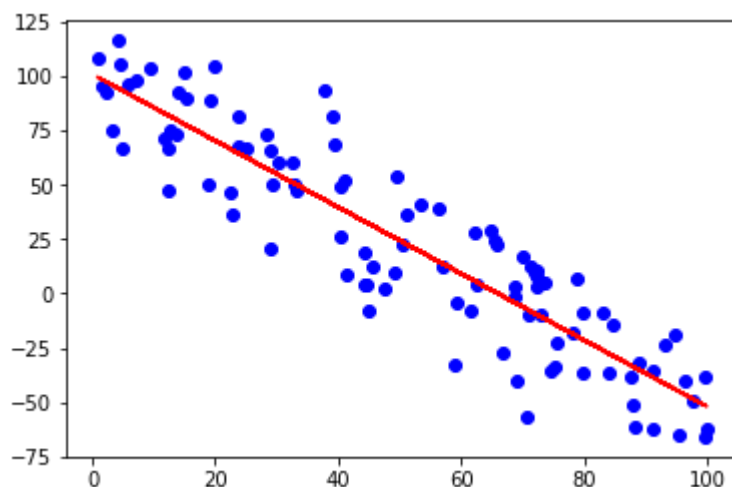
Out[8]: [



```
In [9]: 1 a, b, corr, p, stderr = stats.linregress(x,y)
2 z = a * x + b
3 print ('f(x) = %f x %c %f' % (a, '+' if b >= 0 else '-', abs(b)))
4 print ('corr      = %f' % corr)
5 print ('p        = %f' % p)
6 print ('std err = %f' % stderr)
7 plt.plot(x, y, 'bo')
8 plt.plot(x, z, 'r')
```

```
f(x) = -1.528676 x + 100.858129
corr   = -0.910921
p      = 0.000000
std err = 0.069941
```

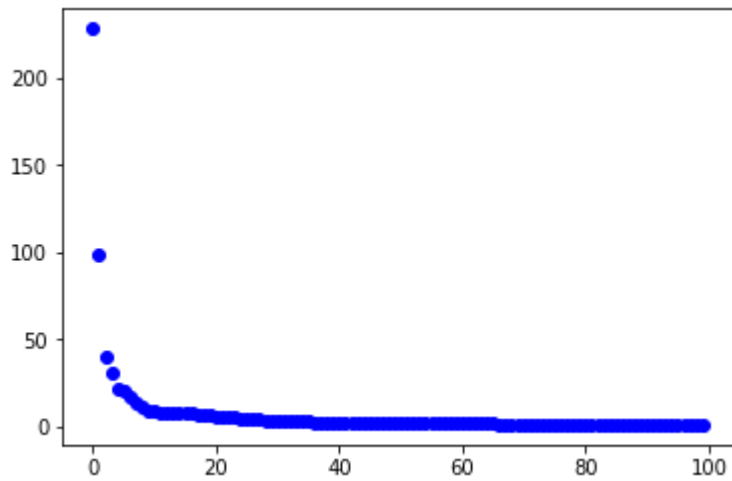
Out[9]: [



## Nonlinear Regression Using Linear Base Functions

```
In [10]: 1 data = np.loadtxt('Data\yp.txt')
          2 x = data[:,0]
          3 y = data[:,1]
          4 plt.plot(x, y, 'bo')
```

Out[10]: [<matplotlib.lines.Line2D at 0x250c6b71d30>]



## Log-log linear function

consider a function as

$$f(x) = bx^a$$

Take logarithm on both sides:

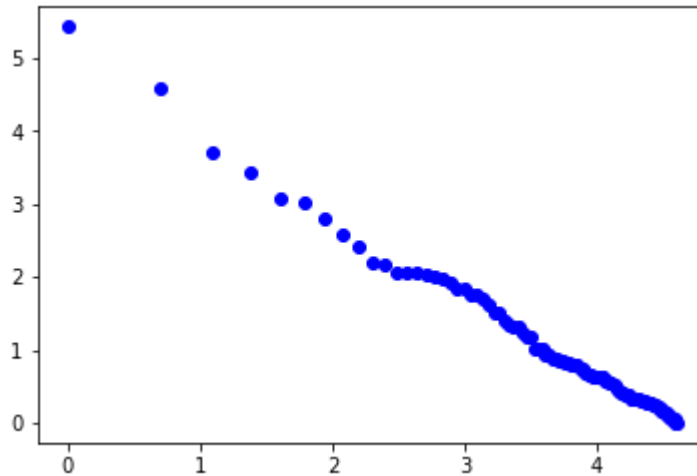
$$\begin{aligned}\log f(x) &= \log b + \log x^a \\ &= \log b + a \log x\end{aligned}$$

By properly substituting:

$$F(x) = B + AX$$

```
In [11]: 1 x = range(1,101)
2 xlog = np.log(x)
3 ylog = np.log(y)
4 plt.plot(xlog, ylog, 'bo')
```

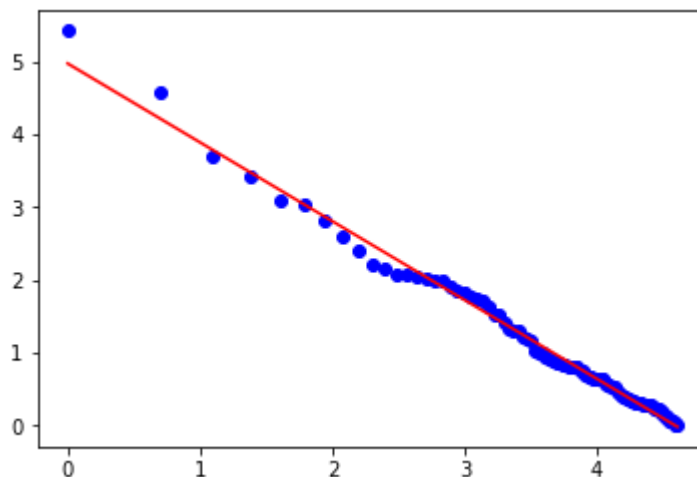
Out[11]: [<matplotlib.lines.Line2D at 0x250c6bca438>]



```
In [12]: 1 a, b, corr, p, stderr = stats.linregress(xlog,ylog)
2 zlog = a * xlog + b
3 print ('corr=', corr)
4 print ('err =', stderr)
5 plt.plot(xlog, ylog, 'bo')
6 plt.plot(xlog, zlog, 'r')
```

```
corr= -0.995468410955407
err = 0.010470269575618493
```

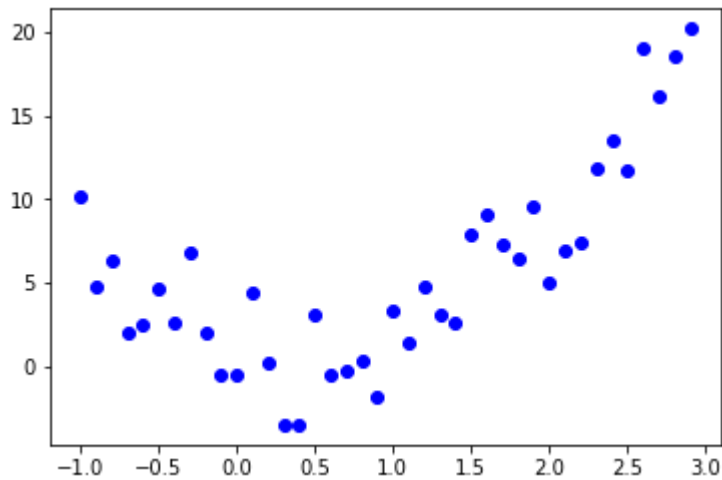
Out[12]: [<matplotlib.lines.Line2D at 0x250c6be8cf8>]



## Nonlinear regression with Optimization

```
In [13]: 1 data = np.loadtxt('Data\y2.txt')
2 x = data[:,0]
3 y = data[:,1]
4 plt.plot(x, y, 'bo')
```

Out[13]: [<matplotlib.lines.Line2D at 0x250c6c76630>]



## Quadratic function

$$f(x) = ax^2 + bx + c$$

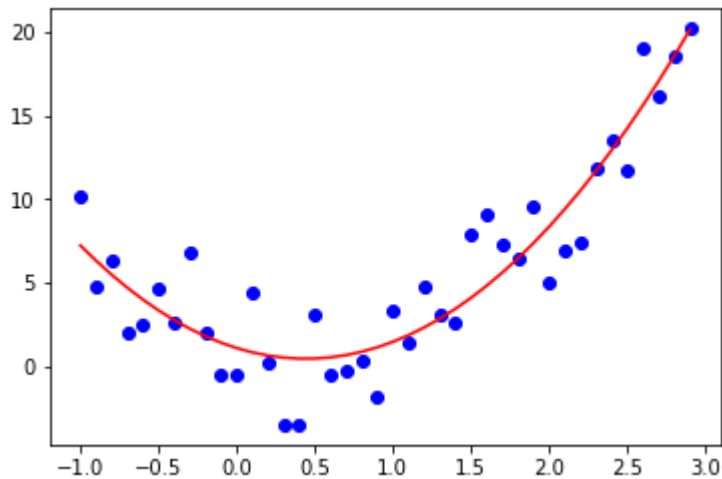
```
In [14]: 1 def quadratic(x,a,b,c):
2         return a * x ** 2 + b * x + c
```

```
In [15]: 1 popt,pcov = opt.curve_fit(quadratic, x, y)
2 print ('pcov=\n', pcov)
3 print ('stdev err=', np.sqrt(np.diag(pcov)))
4 z = quadratic(x, *popt)
```

```
pcov=
[[ 0.10035442 -0.1906734 -0.0431524 ]
 [-0.1906734  0.46905656 -0.01944869]
 [-0.0431524 -0.01944869  0.25720235]]
stdev err= [0.31678766 0.68487704 0.50715121]
```

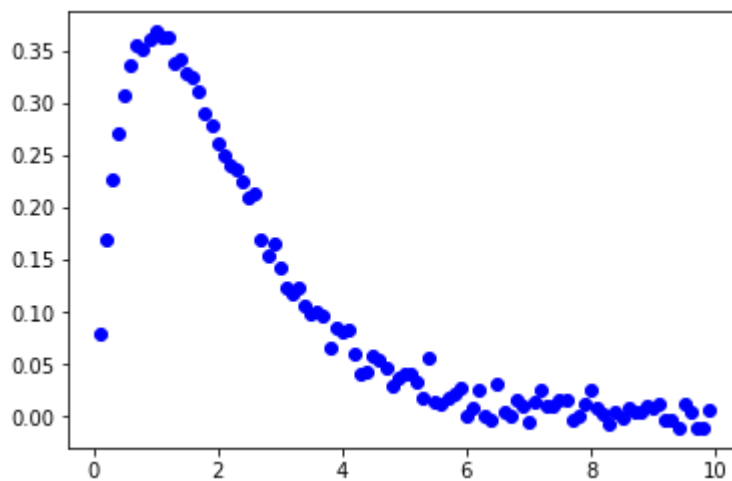
```
In [16]: 1 plt.plot(x, y, 'bo')
         2 plt.plot(x, z, 'r')
```

Out[16]: [<matplotlib.lines.Line2D at 0x250c6c9ed30>]



```
In [17]: 1 data = np.loadtxt('Data\y7.txt')
         2 x = data[:,0]
         3 y = data[:,1]
         4 plt.plot(x, y, 'bo')
```

Out[17]: [<matplotlib.lines.Line2D at 0x250c6d38a90>]



## Gammar function

$$f(x) = ax^b e^{cx}$$

```
In [18]: 1 def gamma(x, a, b, c):
         2     return ( a * x ** b ) * np.exp( c * x )
```

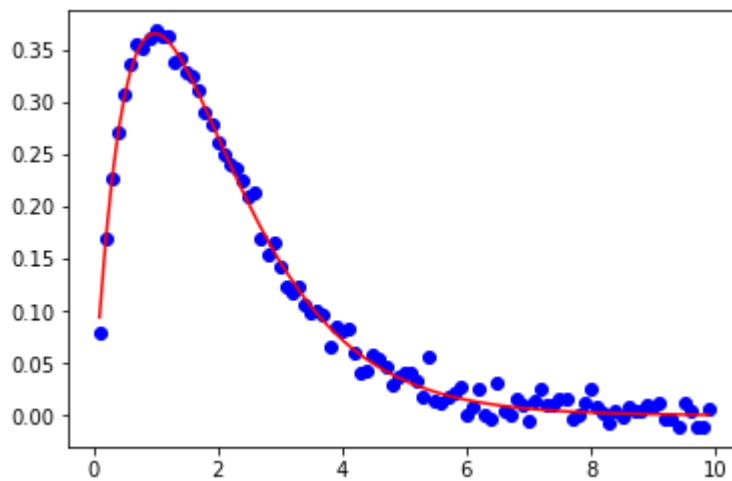


```
In [19]: 1 popt,pcov = opt.curve_fit(gamma, x, y)
2 print ('popt=\n', popt)
3 print ('pcov=\n', pcov)
4 print ('stdev err=', np.sqrt(np.diag(pcov)))
5 z = gamma(x, *popt)
```

```
popt=
[ 0.98708091  0.97783386 -0.99144034]
pcov=
[[ 0.00050611  0.00046711 -0.0003879 ]
 [ 0.00046711  0.00060765 -0.00040746]
 [-0.0003879  -0.00040746  0.00032398]]
stdev err= [0.022497  0.02465048 0.01799935]
```

```
In [20]: 1 plt.plot(x, y, 'bo')
2 plt.plot(x, z, 'r')
```

```
Out[20]: [<matplotlib.lines.Line2D at 0x250c6d66b38>]
```

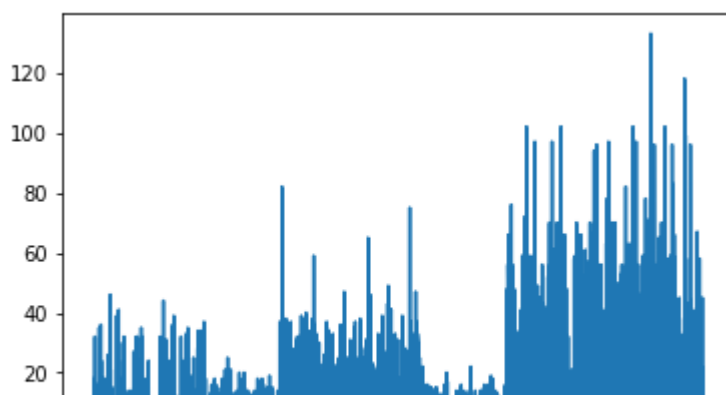


## Degree Fit

```
In [21]: 1 # to show the raw data
2 y = np.loadtxt('Data/degrees.txt', dtype=int)
3 print ('y =', y)
4 print ('size: ', y.size)
5 print ('type: ', y.dtype)
6 print ('shape:', y.shape)
7 print ('ndim: ', y.ndim)
8 plt.plot(y)
```

```
y = [ 1  3  1 ... 45  3 14]
size: 6633
type: int32
shape: (6633,)
ndim: 1
```

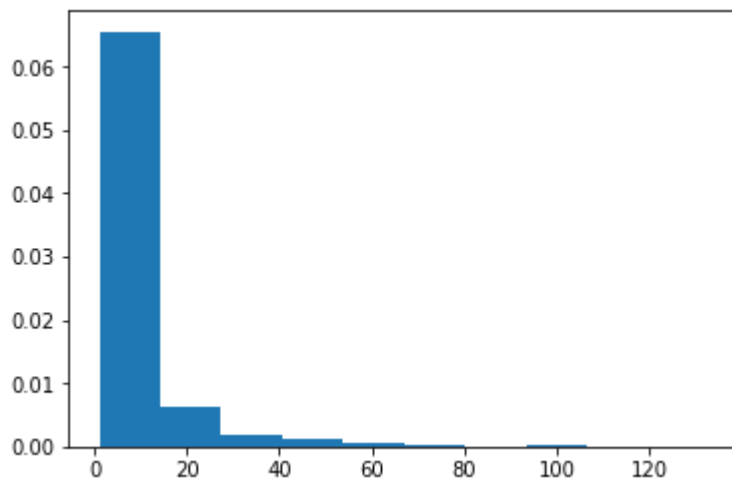
Out[21]: [<matplotlib.lines.Line2D at 0x250c6df7fd0>]



```
In [22]: 1 # to show the histogram
2 plt.hist(y, normed=True)
3 print ('# of points =', len(y))
```

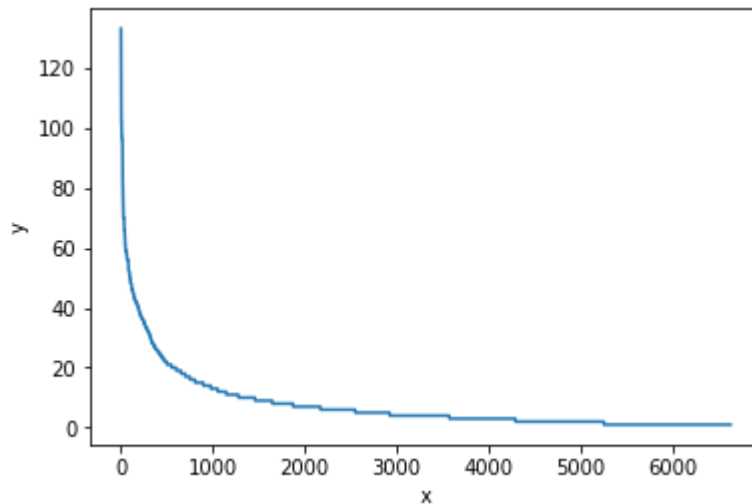
```
# of points = 6633
```

C:\Anaconda3\lib\site-packages\matplotlib\axes\\_axes.py:6462: UserWarning: The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.  
warnings.warn("The 'normed' kwarg is deprecated, and has been "



```
In [23]: 1 y = sorted(y, reverse=True)
2 x = np.arange(1,len(y)+1)
3 # to show the rank
4 plt.plot(x,y)
5 plt.xlabel('x')
6 plt.ylabel('y')
```

Out[23]: Text(0,0.5, 'y')



```
In [24]: 1 def power_law(x,a,c):
2         return c * (x ** -a)
```

```
In [25]: 1 # fail because too big x values
2 popt,pcov = opt.curve_fit(power_law, x, y)
```

C:\Anaconda3\lib\site-packages\ipykernel\_launcher.py:2: RuntimeWarning: overflow encountered in power

C:\Anaconda3\lib\site-packages\ipykernel\_launcher.py:2: RuntimeWarning: overflow encountered in multiply

```
In [26]: 1 def compare(x,y,z):
2         plt.loglog(x,y, marker='o')
3         plt.plot(x,z)
4         plt.show()
```

## Scaling x-axis

The x-axis, or the rank, needs not be an enumerable counter. Instead, you may think the x-axis as a percentile so that it is scaled from 0 to 100.

```
In [27]: 1 delta = 1.0 / float(len(x)) * 100
2 x = np.arange(delta, 100.0+delta, delta)
3
4 print ('\nfitting to power law')
5 # fail in curve_fit due to a bad initial value
6 popt,pcov = opt.curve_fit(power_law, x, y)
7 print ('pcov=\n', pcov)
8 print ('stdev err=', np.sqrt(np.diag(pcov)))
9 z = power_law(x, *popt)
10 compare(x,y,z)
```

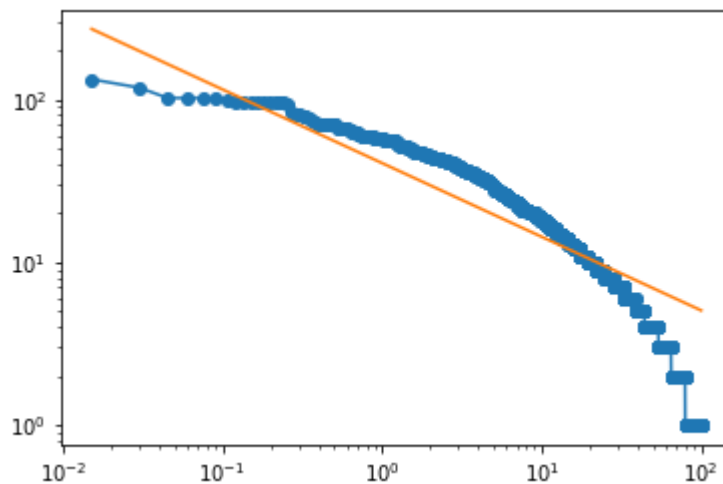
fitting to power law

pcov=

```
[[3.45916823e-06 8.56050347e-05]
```

```
[8.56050347e-05 4.20577177e-02]]
```

stdev err= [0.00185988 0.20507978]



```
In [28]: 1 def power_exp(x,a,b,c):
2         return c * (x ** a) * np.exp(b * x)
```

```
In [29]: 1 print ('\nfitting to power law with arbitrary')
2 popt,pcov = opt.curve_fit(power_exp, x, y)
3 print ('popt=\n', popt)
4 print ('pcov=\n', pcov)
5 print ('stdev err=', np.sqrt(np.diag(pcov)))
6 z = power_exp(x, *popt)
7 compare(x,y,z)
```

fitting to power law with arbitrary

popt=

```
[9.99966593e-01 9.9993868e-01 7.38921955e-46]
```

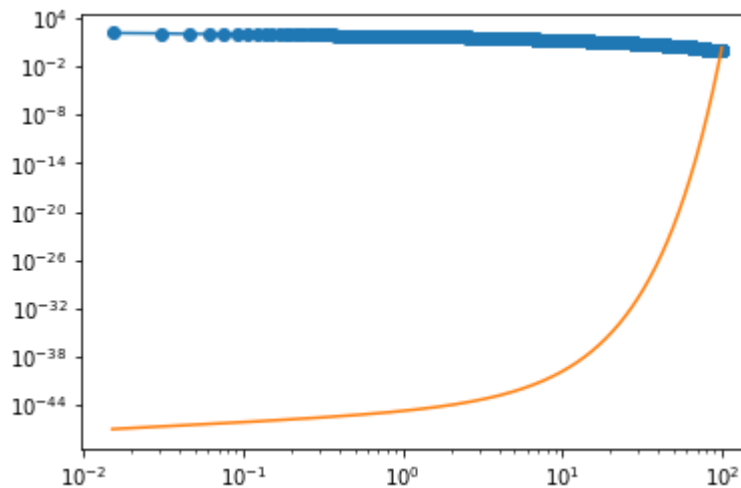
pcov=

```
[[ 2.19813968e+09 -2.22005881e+07 -5.84523088e-36]
```

```
[-2.22005881e+07 2.24225369e+05 5.90347553e-38]
```

```
[-5.84523088e-36 5.90347553e-38 1.55435036e-80]]
```

stdev err= [4.68843223e+04 4.73524412e+02 1.24673588e-40]



## Intial Values

Decide intial values of parameters when they are clear and present. In the above example, we know it is a decreasing function. So,  $a$  and  $b$  should be negative.

```
In [30]: 1 print ('\nfitting to power law with exponential cutoff with an initial value
2 popt,pcov = opt.curve_fit(power_exp, x, y, p0=[-2, -2, 10])
3 print ('popt=\n', popt)
4 print ('pcov=\n', pcov)
5 print ('stdev err=', np.sqrt(np.diag(pcov)))
6 z = power_exp(x, *popt)
7 compare(x,y,z)
```

fitting to power law with exponential cutoff with an initial value

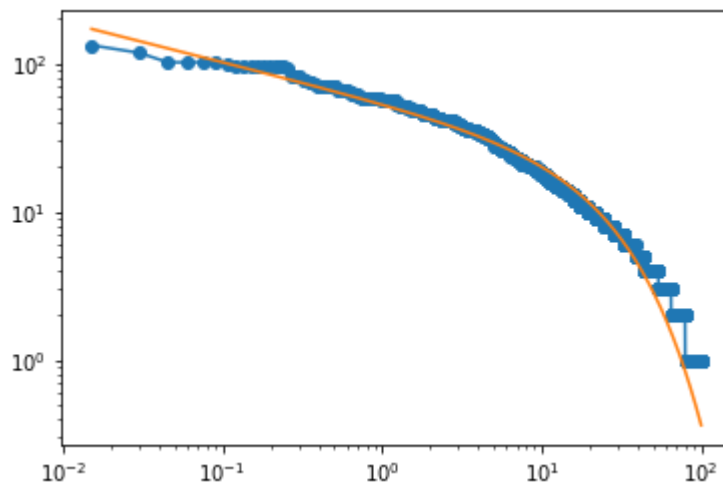
popt=

```
[-2.71764262e-01 -3.76132009e-02  5.47880810e+01]
```

pcov=

```
[[ 9.24605466e-07 -1.35387381e-07  1.18384187e-05]
 [-1.35387381e-07  3.99241869e-08 -8.67196868e-06]
 [ 1.18384187e-05 -8.67196868e-06  7.45344200e-03]]
```

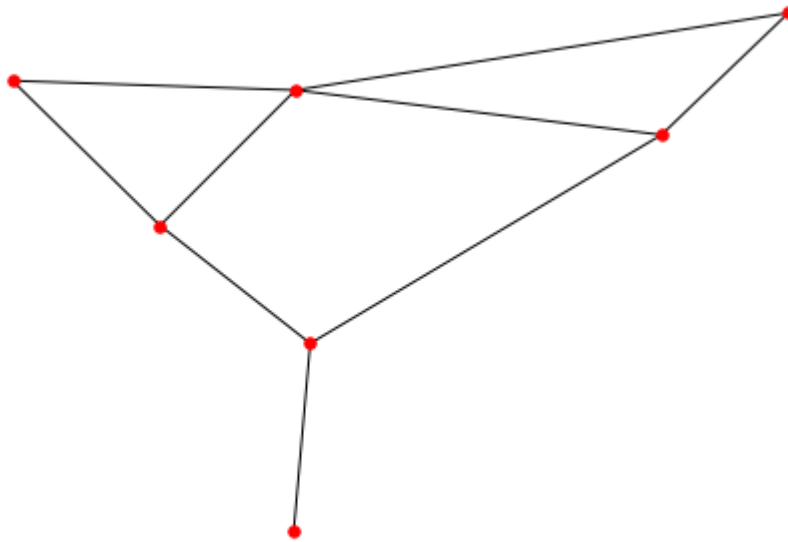
```
stdev err= [0.00096156 0.00019981 0.08633332]
```



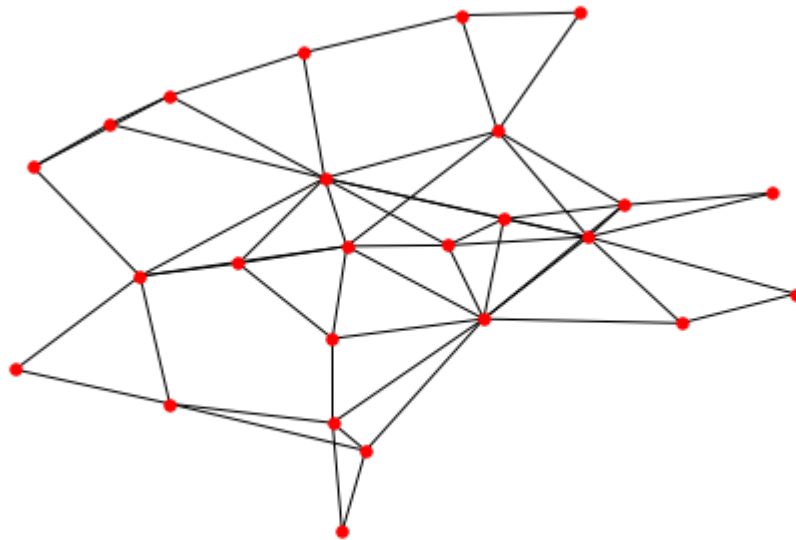
## Visualization of Degree Rank Analysis and Fitting

```
In [41]: 1 symbol = nl.marker()
2
3 def degree_analysis(G, name=None):
4     y = sorted(list(dict(nx.degree(G)).values()), reverse=True)
5     x = range(1, len(y)+1, 1)
6     m0, m1 = next(symbol)
7     popt, pcov = opt.curve_fit(power_law, x, y)
8     print ('nonlinear_regression\npopt=%s\npcov=%s\n' % (popt, pcov))
9     z = power_law(x, *popt)
10
11     plt.loglog(x, y, m0, label=name+' data')
12     plt.loglog(x, z, m1, label=name+' fitted')
13
```

```
In [42]: 1 f1 = 'Graphs/g35.gml'
2 G1 = nx.read_gml(f1)
3 nx.draw(G1, node_size=30)
```



```
In [43]: 1 f2 = 'data/AttMpls.gml'
2 G2 = nx.read_gml(f2)
3 nx.draw(G2, node_size=30)
```

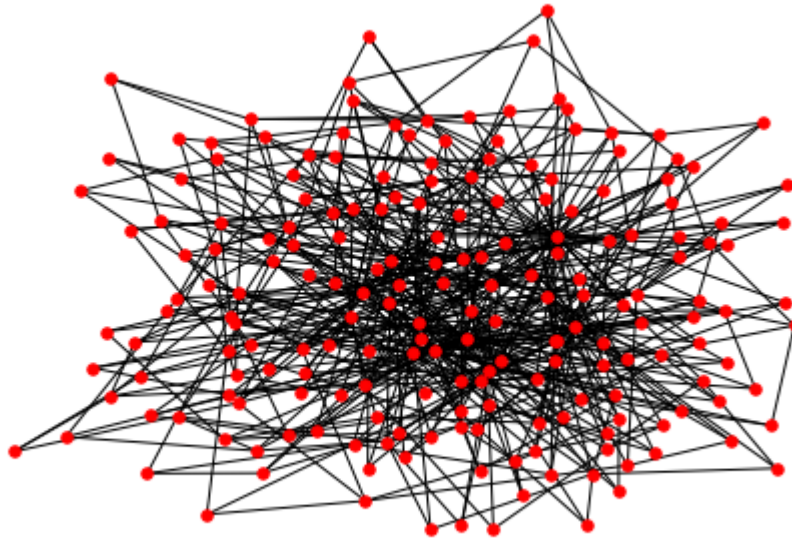


In [44]:

```

1 f3 = 'Graphs/rand200.gml'
2 G3 = nx.read_gml(f3)
3 nx.draw(G3, node_size=30)

```



In [45]:

```

1 degree_analysis(G1, f1)
2 degree_analysis(G2, f2)
3 degree_analysis(G3, f3)
4
5 plt.xlabel('rank')
6 plt.ylabel('degrees')
7 plt.title('Degree Rank Analysis')
8 plt.legend()
9 plt.show()

```

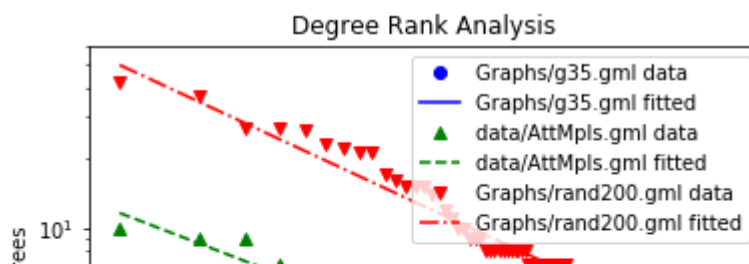
nonlinear\_regression

popt=[ 0.52745048 49.9506148 ]

pcov=

[[3.57145516e-05 3.96954104e-03]

[3.96954104e-03 6.87510514e-01]]



## Advanced Fitting



In [47]:

```

1 def power_cumulative(x,a,c):
2     if a != 1.0:
3         return c/(1-a) * x ** (-a+1)
4     else:
5         return c * np.log(x)
6
7 def draw_xy(x,y):
8     plt.title('Rank of Zipf Sequence')
9     plt.ylabel('#')
10    plt.xlabel('rank(%)')
11    plt.xlim([min(x), max(x)])
12    plt.loglog(x,y,marker='o', label='raw')
13    return
14
15 def fitting_power(x,y):
16    popt, pcov = opt.curve_fit(power_law, x, y)
17    print ('coeff =', popt)
18    z = power_law(x, *popt)
19    plt.plot(x, z, label='power law fit')
20    return
21
22 def fitting_log_bin(x,y):
23    x1 = []
24    y1 = []
25    si = 1.0
26    while si <= len(x):
27        i = int(si)
28        x1.append(x[i])
29        y1.append(y[i])
30        si *= 1.2
31    popt, pcov = opt.curve_fit(power_law, x1, y1)
32    print ('log bin coeff =', popt)
33    z1 = power_law(x1, *popt)
34    plt.plot(x1, z1, marker='s', label='log bin fit')
35    return
36
37 def fitting_cumulative(x,y):
38    y1 = y
39    for i in range(1,len(y1)):
40        y1[i] += y1[i-1]
41
42    popt, pcov = opt.curve_fit(power_cumulative, x, y1, p0=[0.5, 1.0])
43    z1 = power_cumulative(x, *popt)
44    z = np.zeros([len(z1)])
45    z[0] = z1[0]
46    for i in range(1,len(z1)):
47        z[i] = z1[i] - z1[i-1]
48    plt.plot(x,z,label='cum. fit')
49    return

```

In [55]:

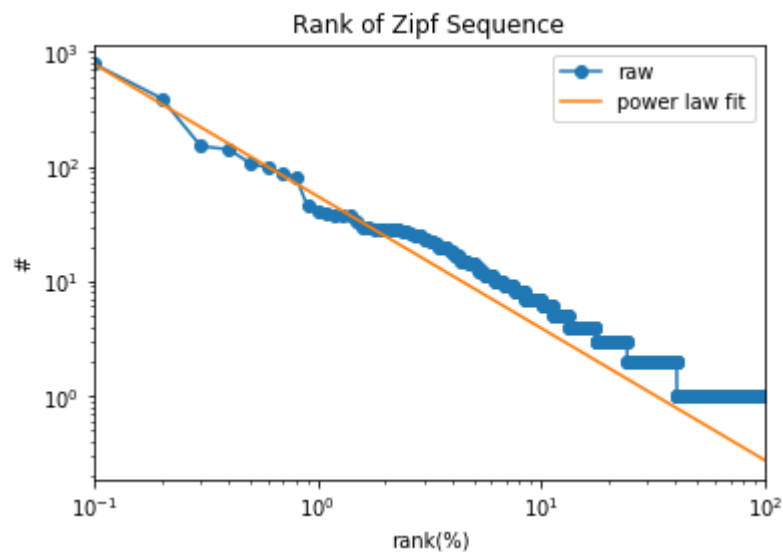
```

1 s = [nx.utils.random_sequence.zipf_rv(alpha=2) for i in range(1000)]
2 y = sorted(s, reverse=True)
3 delta = 1.0 / float(len(y)) * 100
4 x = np.arange(delta,100.0+delta, delta)

```

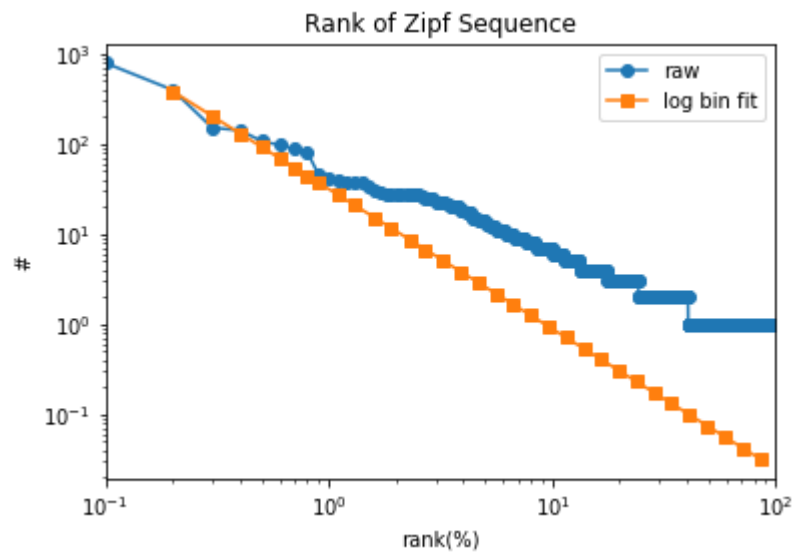
```
In [58]: 1 l0 = draw_xy(x,y)
2         l1 = fitting_power(x,y)
3         plt.legend()
4         plt.show()
```

```
coeff = [ 1.15184945 55.18534468]
```



```
In [59]: 1 l0 = draw_xy(x,y)
2         l2 = fitting_log_bin(x,y)
3         plt.legend()
4         plt.show()
```

```
log bin coeff = [ 1.54956154 31.44460708]
```

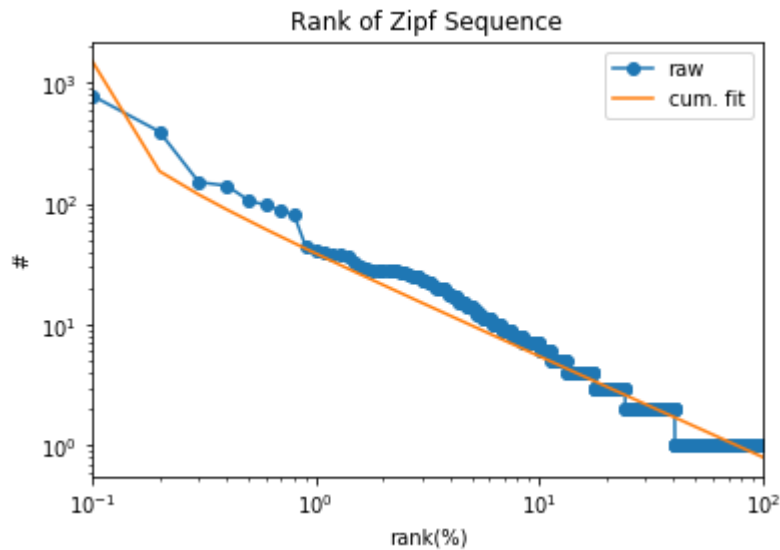


In [60]:

```

1 l0 = draw_xy(x,y)
2 l3 = fitting_cumulative(x,y)
3 plt.legend()
4 plt.show()

```



In [61]:

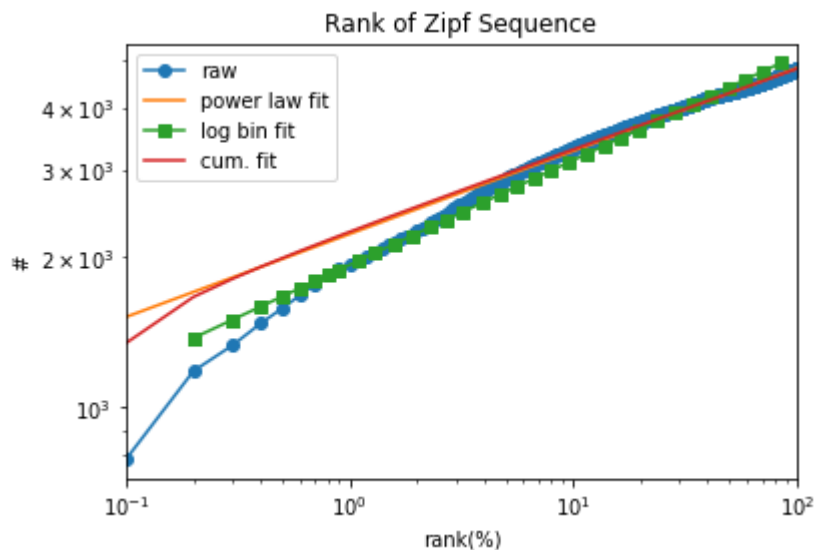
```

1 l0 = draw_xy(x,y)
2 l1 = fitting_power(x,y)
3 l2 = fitting_log_bin(x,y)
4 l3 = fitting_cumulative(x,y)
5 plt.legend()
6 plt.show()

```

```
coeff = [-1.66736396e-01  2.22806881e+03]
```

```
log bin coeff = [-2.10421651e-01  1.93085887e+03]
```



```
In [62]: 1 def rank_vertex_degree(G, mini=False):
2         D = sorted(nx.degree(G).values(),reverse=True) # degree sequence
3         plt.loglog(D,'b-',marker='o')
4         plt.title("Degree Rank")
5         plt.ylabel("degree")
6         plt.xlabel("rank")
7         if mini:
8             mini_draw(G)
9         plt.show()
```

```
In [ ]: 1
```