

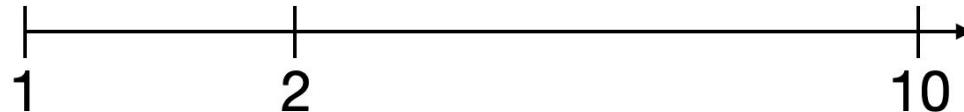
# Data Visualization

## W9-1

# Quiz

- What do you find interesting in today's VotW?
- Explain the key differences between interpolation and regression.
- Can you find 4 and 5? Then explain how.

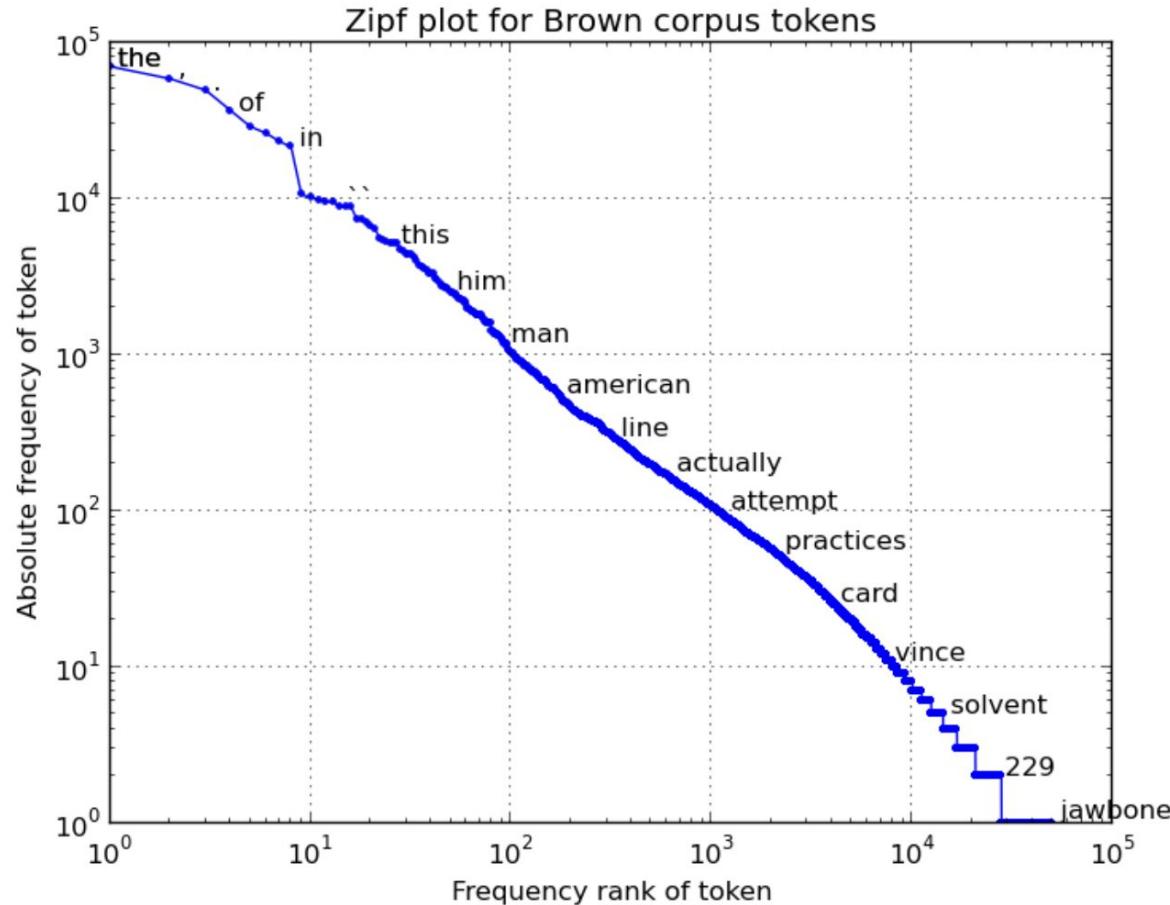
Logscale

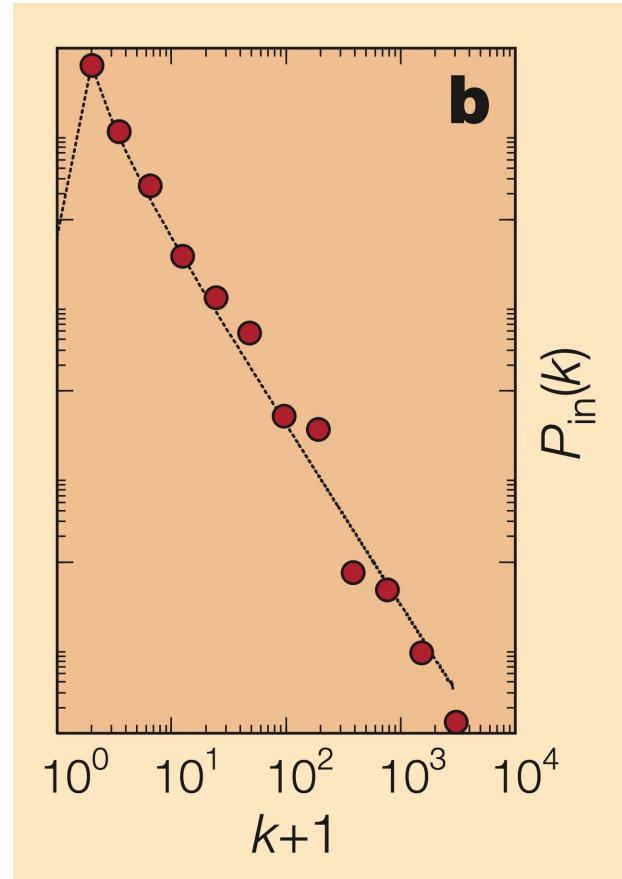
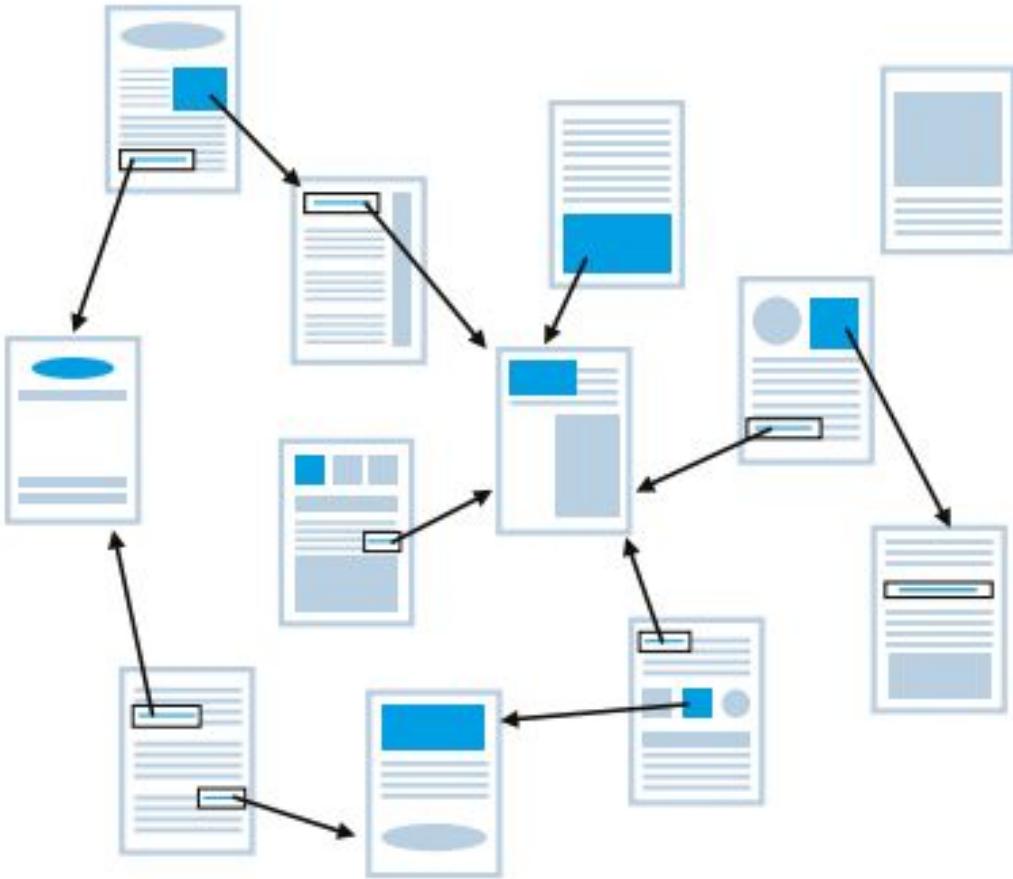


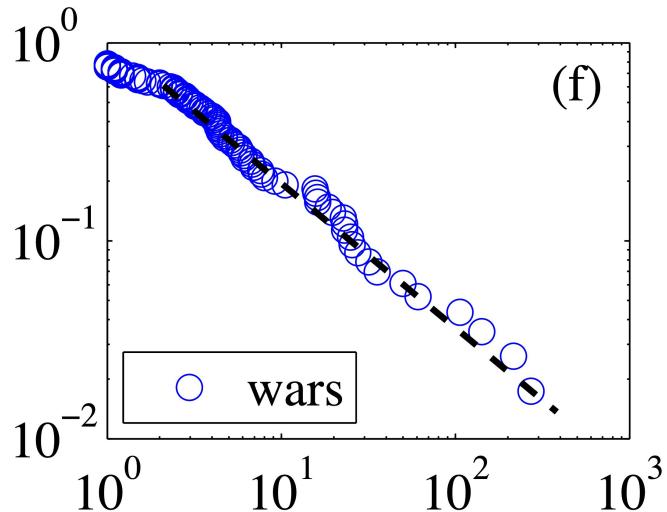
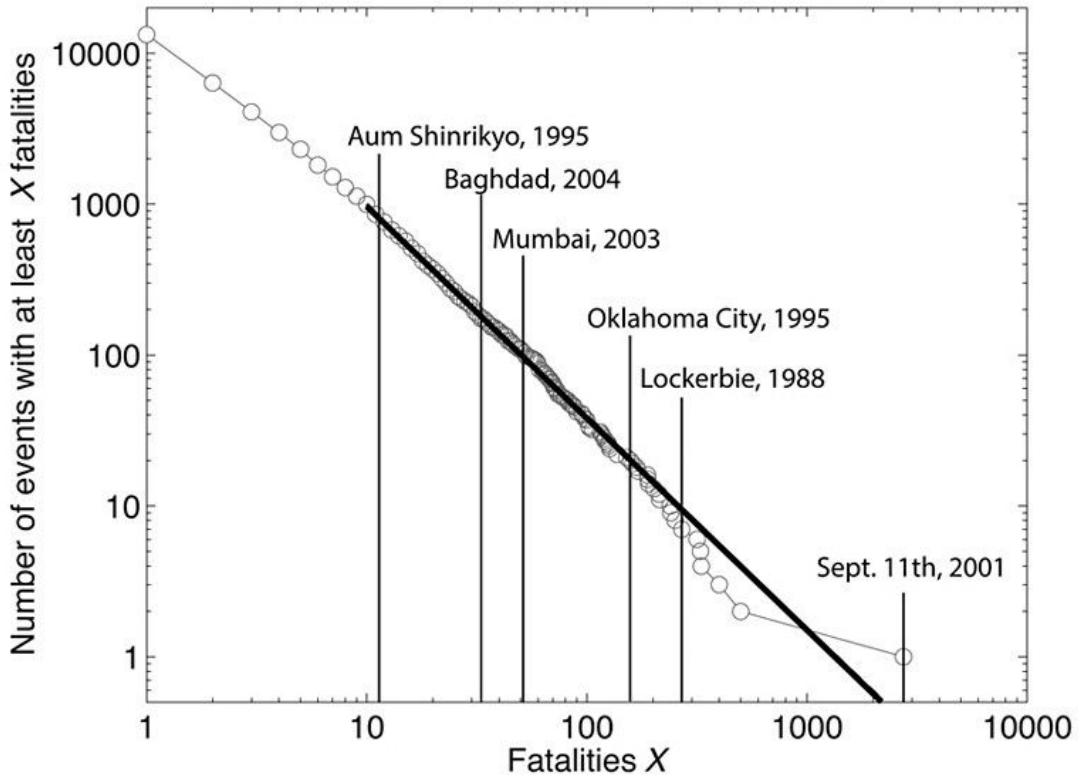
# 'Power-law' distributions

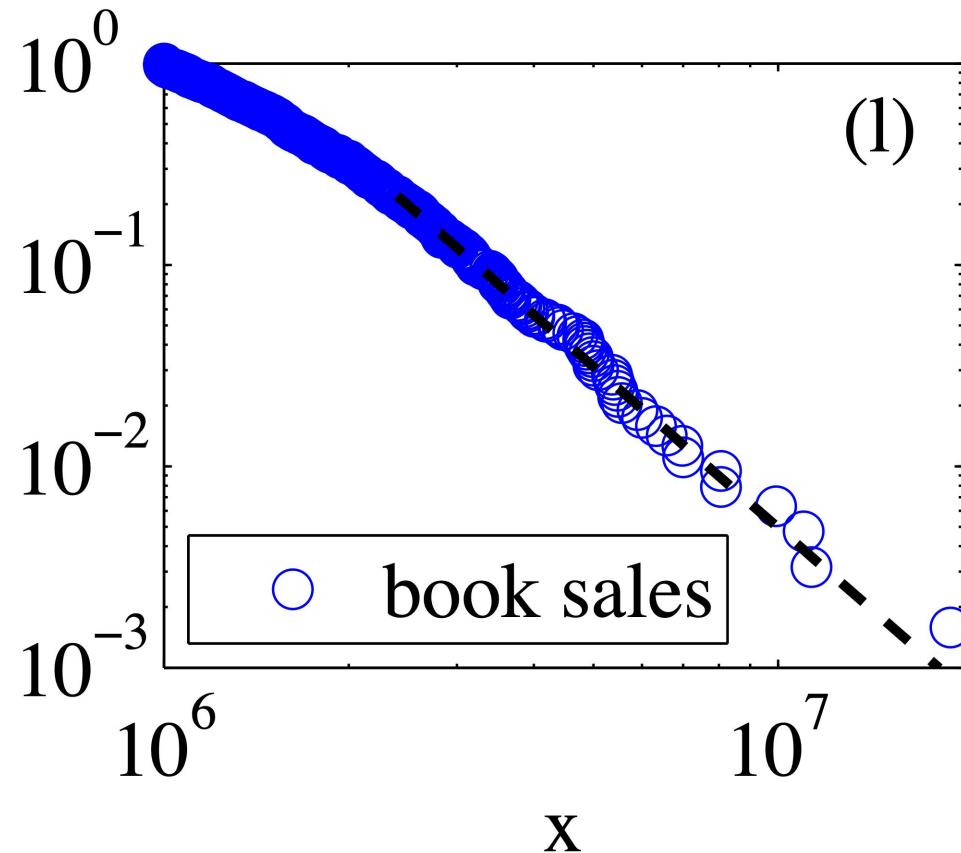
$$f(x) = cx^{-\alpha}$$

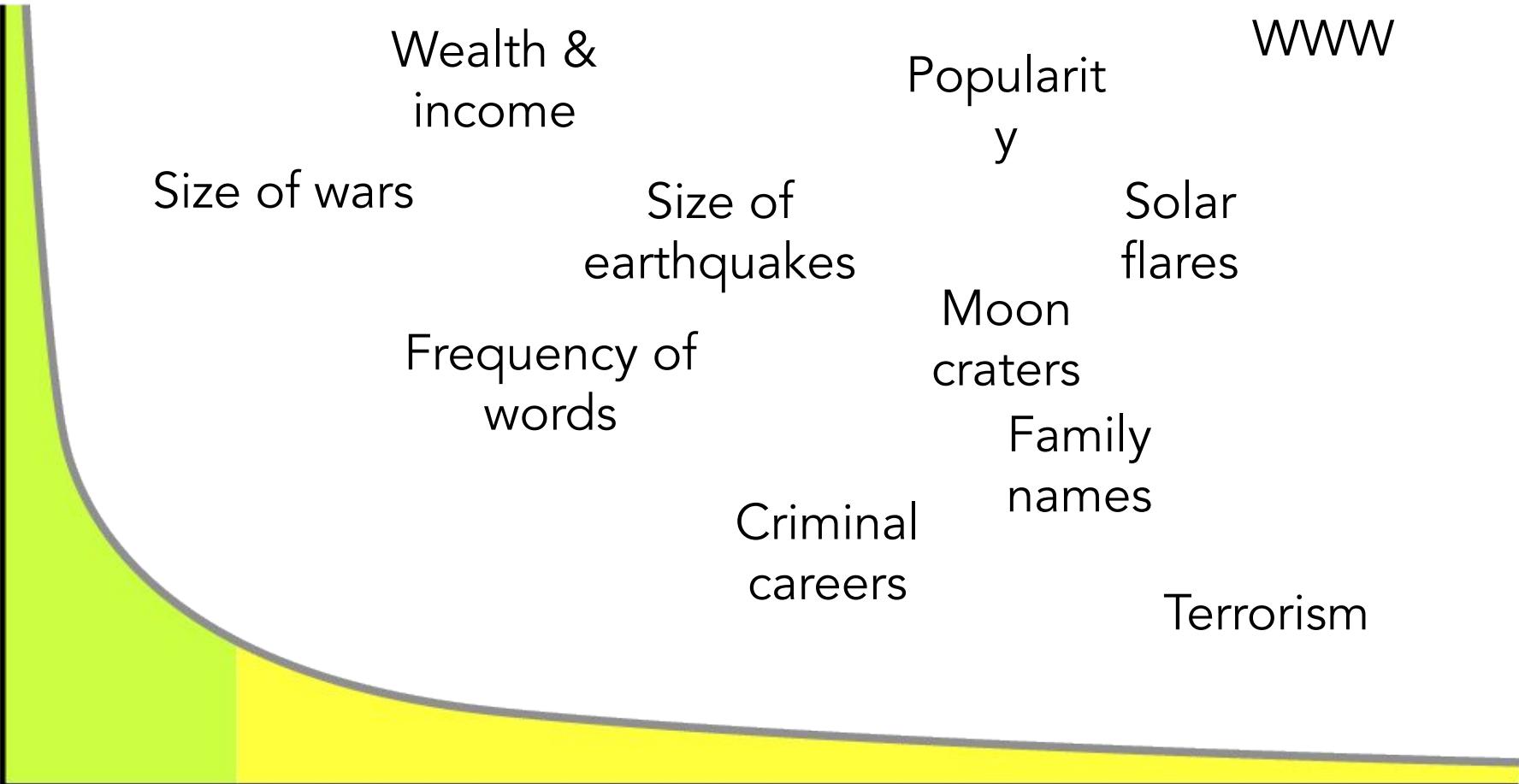
$$\begin{aligned}\log f(x) &= \log(cx^{-\alpha}) \\ &= \log c - \alpha \log x\end{aligned}$$



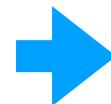
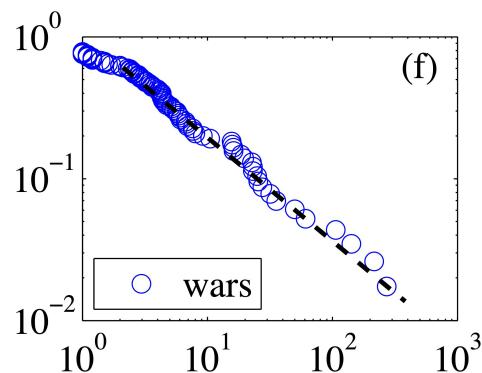
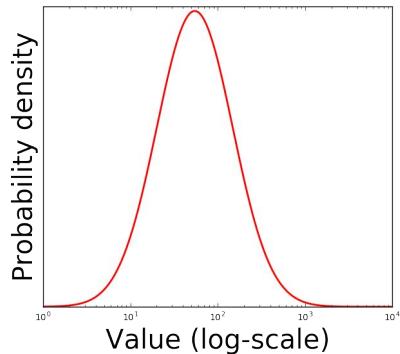
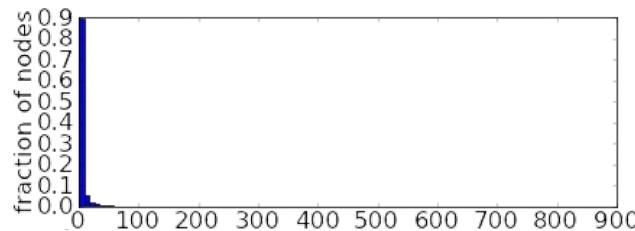
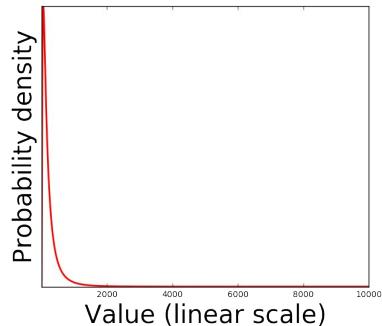




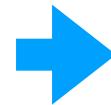




# Have a fat tail?

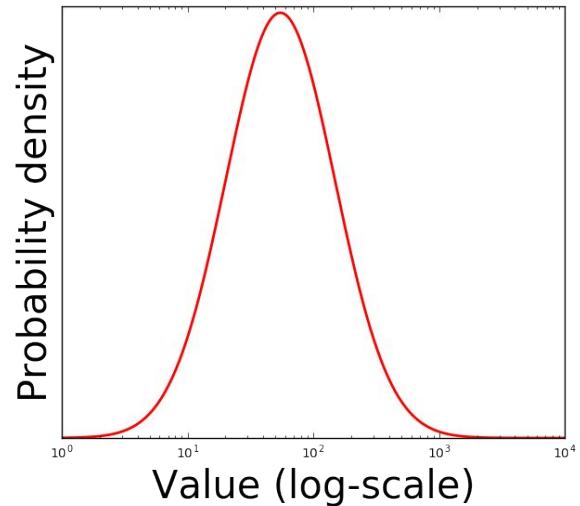
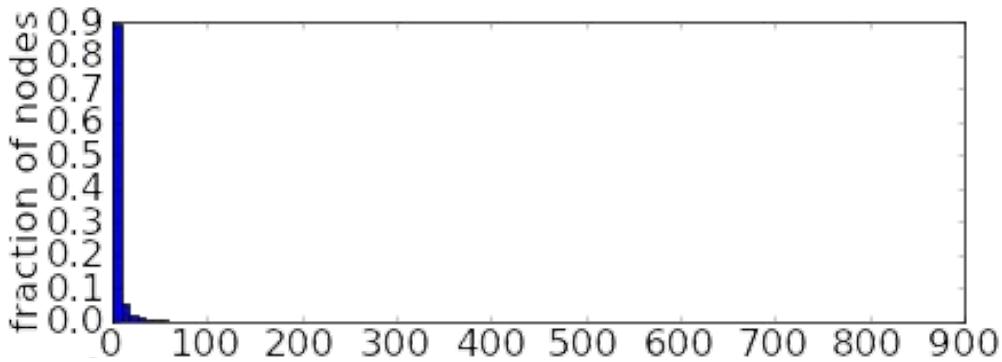


Linear scale:  
Can't see any!



log scale:  
Useful!

# Can we draw a useful histogram in log-scale?



Histogram in  
log-scale

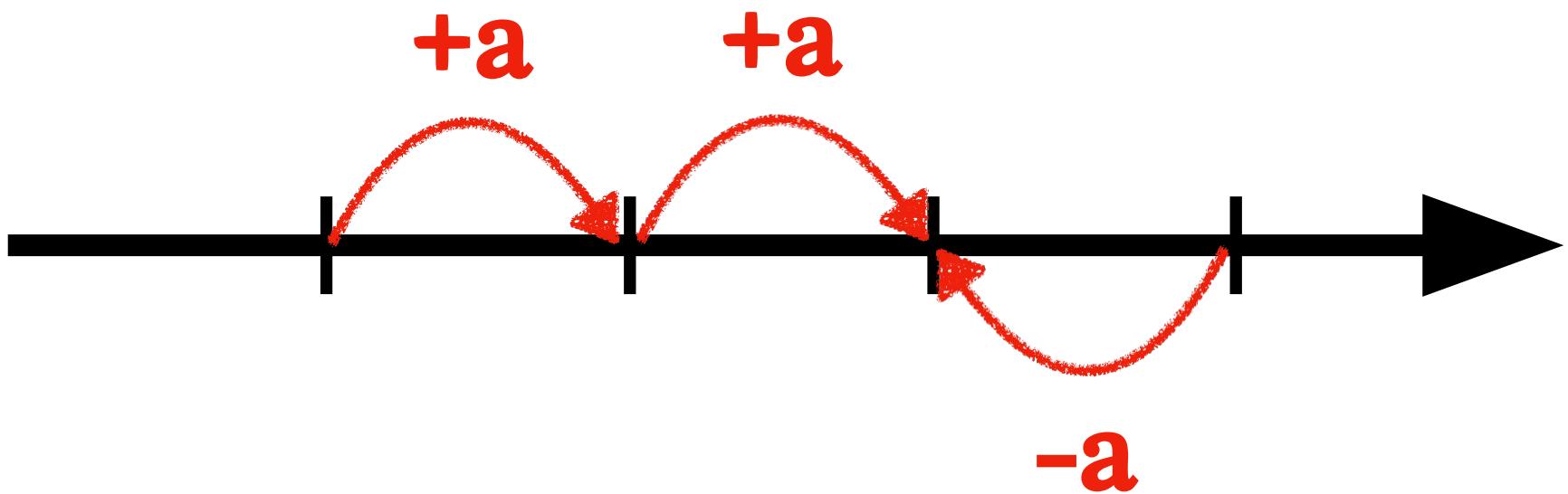
How should we create bins  
in log-scale?

Constant width in  
log-scale?

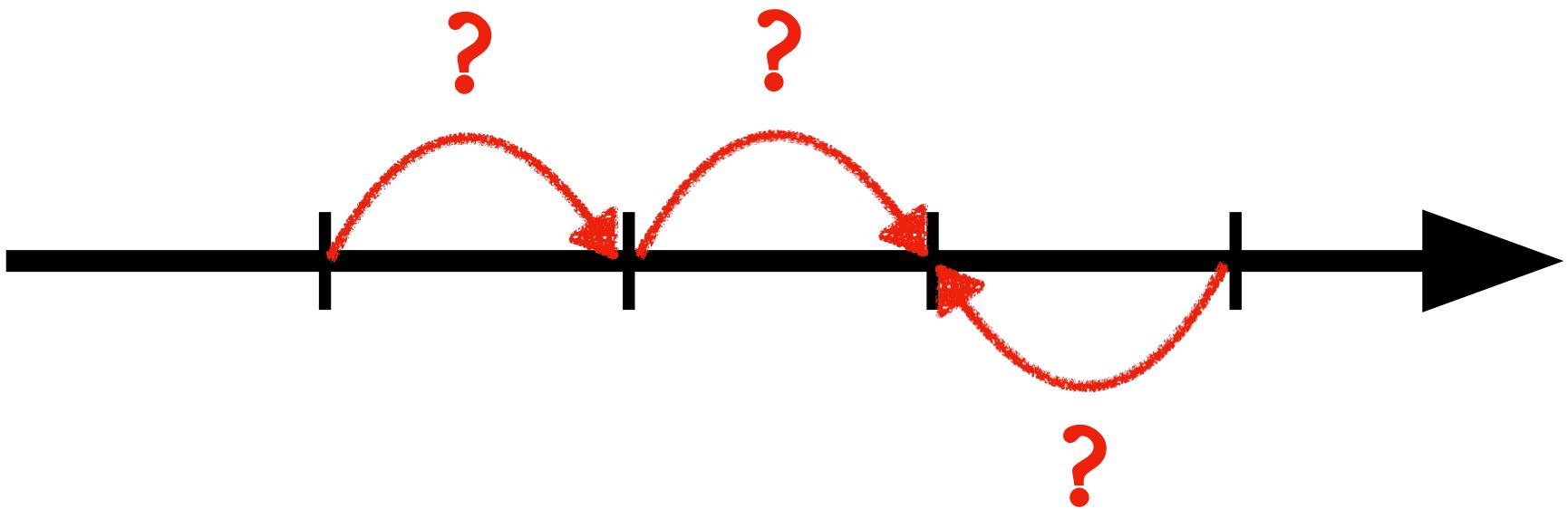
First bin: [1, 3)

What are the next several  
bins? (same width bins in  
log-scale)

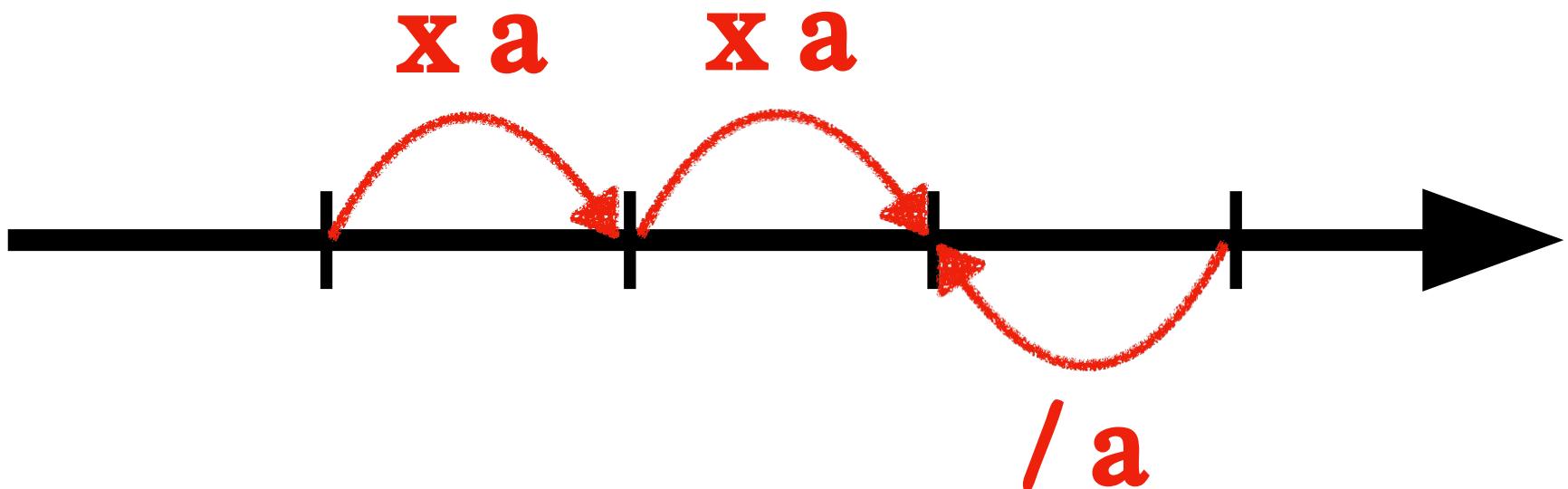
# In linear scale



# In log scale?

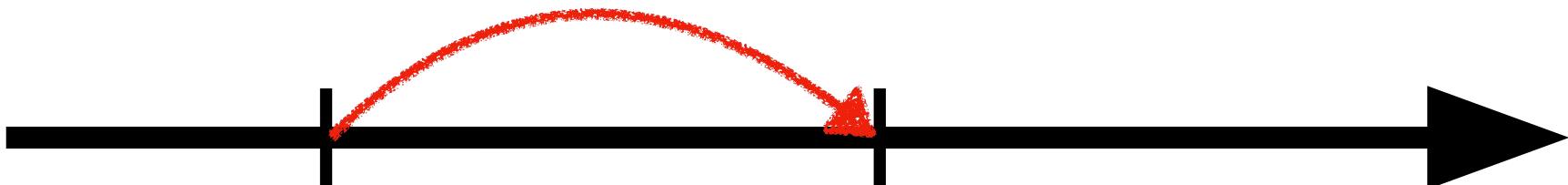


# In log scale?



(log-scale)

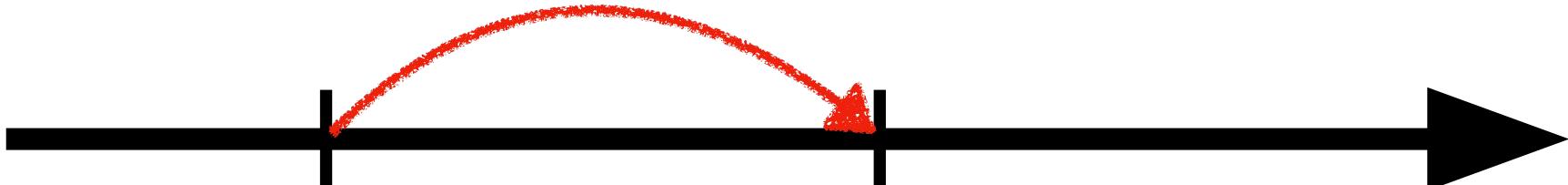
x 2



I

2

x 2

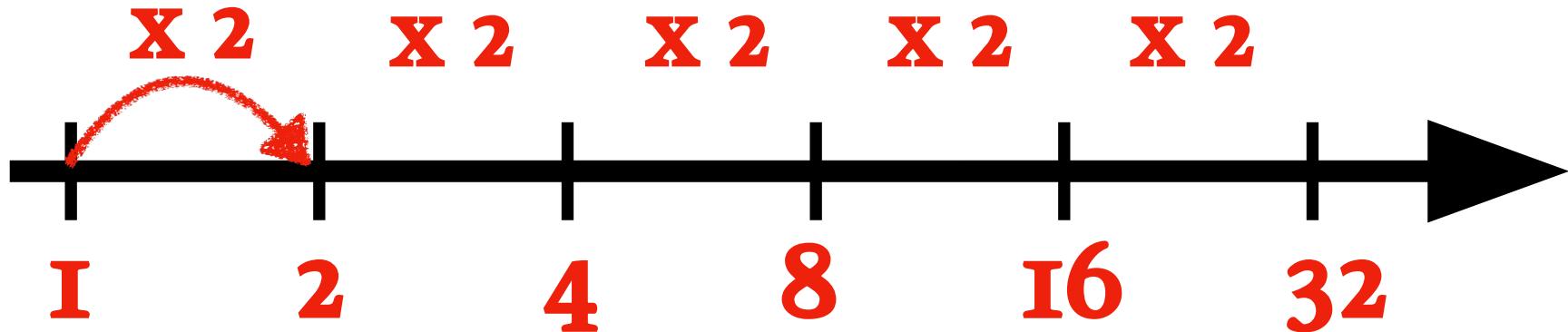


1,000,000

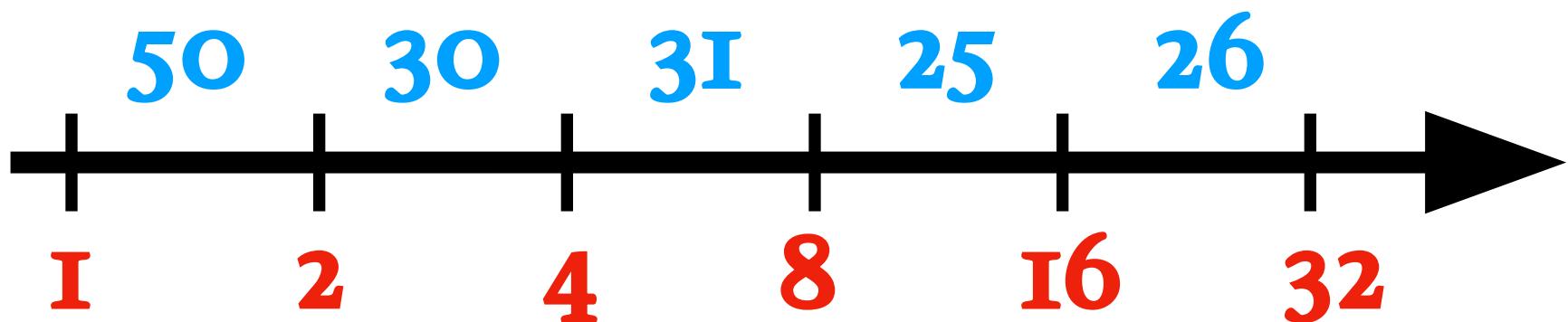
2,000,000

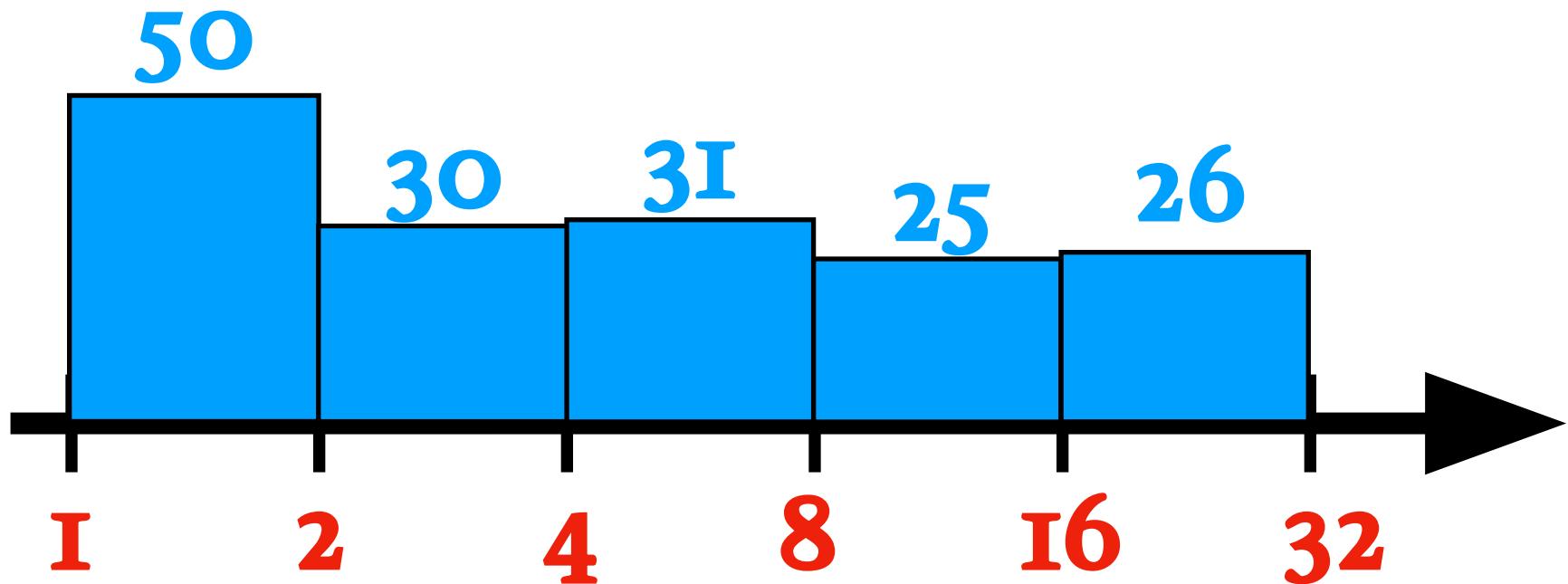
(You can choose  
other widths too)

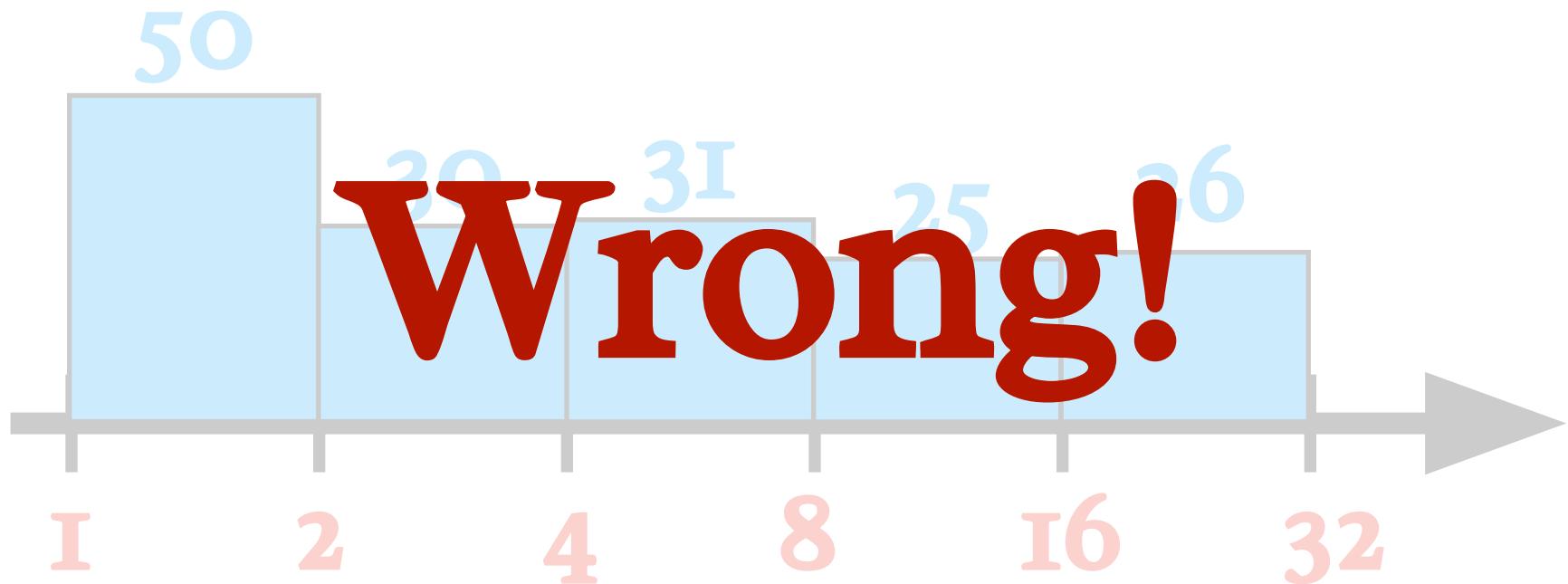
(log-scale)



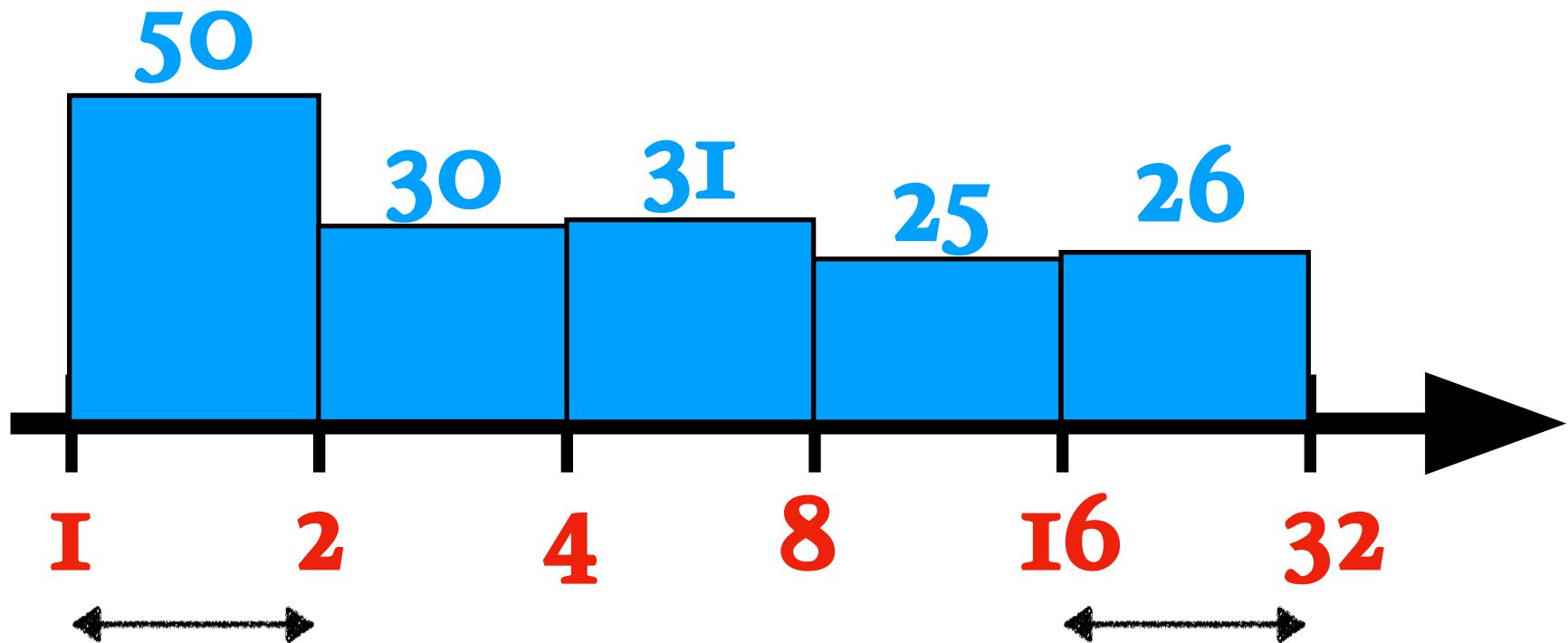
“Log-binning”  
(logarithmic  
binning)



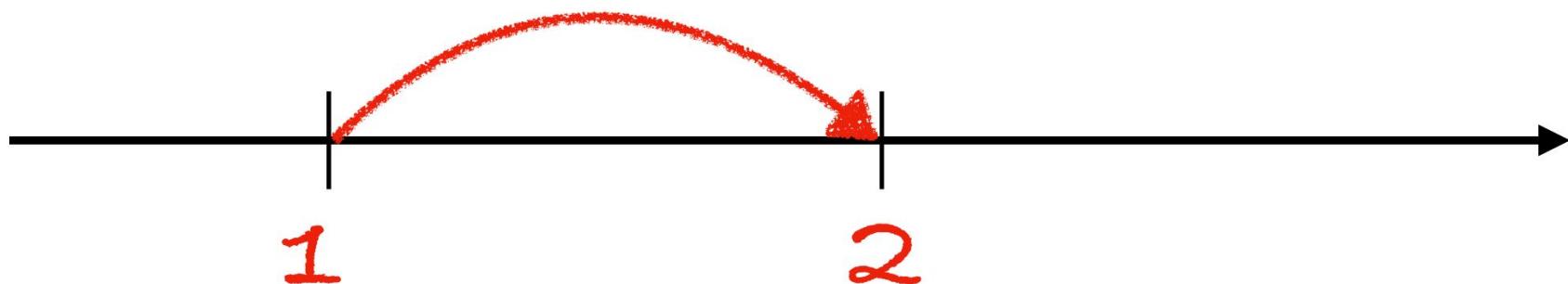




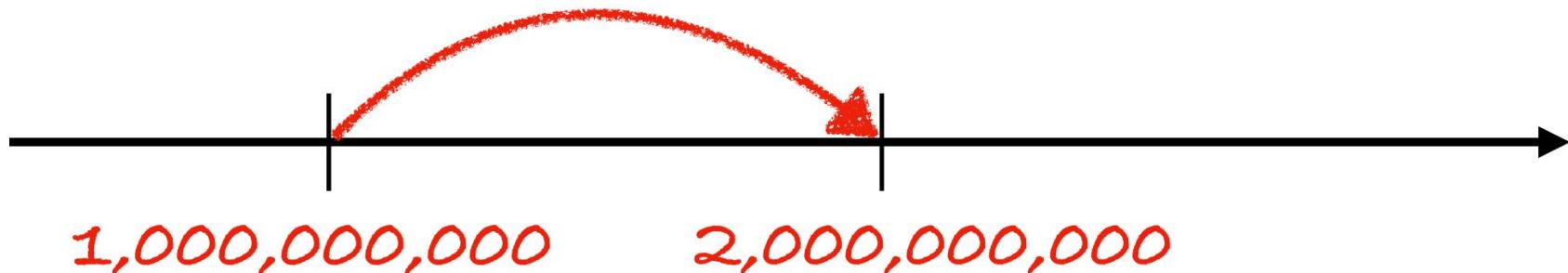
“Area, *not the* height,  
*represents the frequency!*”

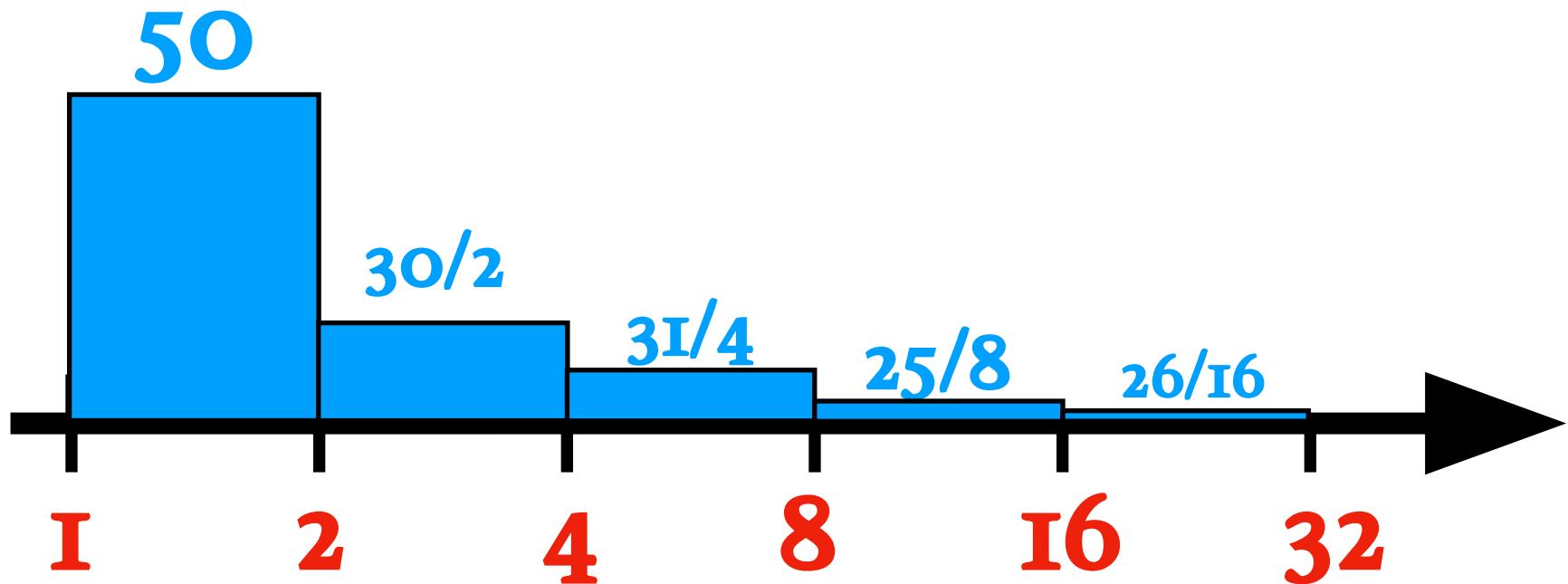


$x_2, +1$



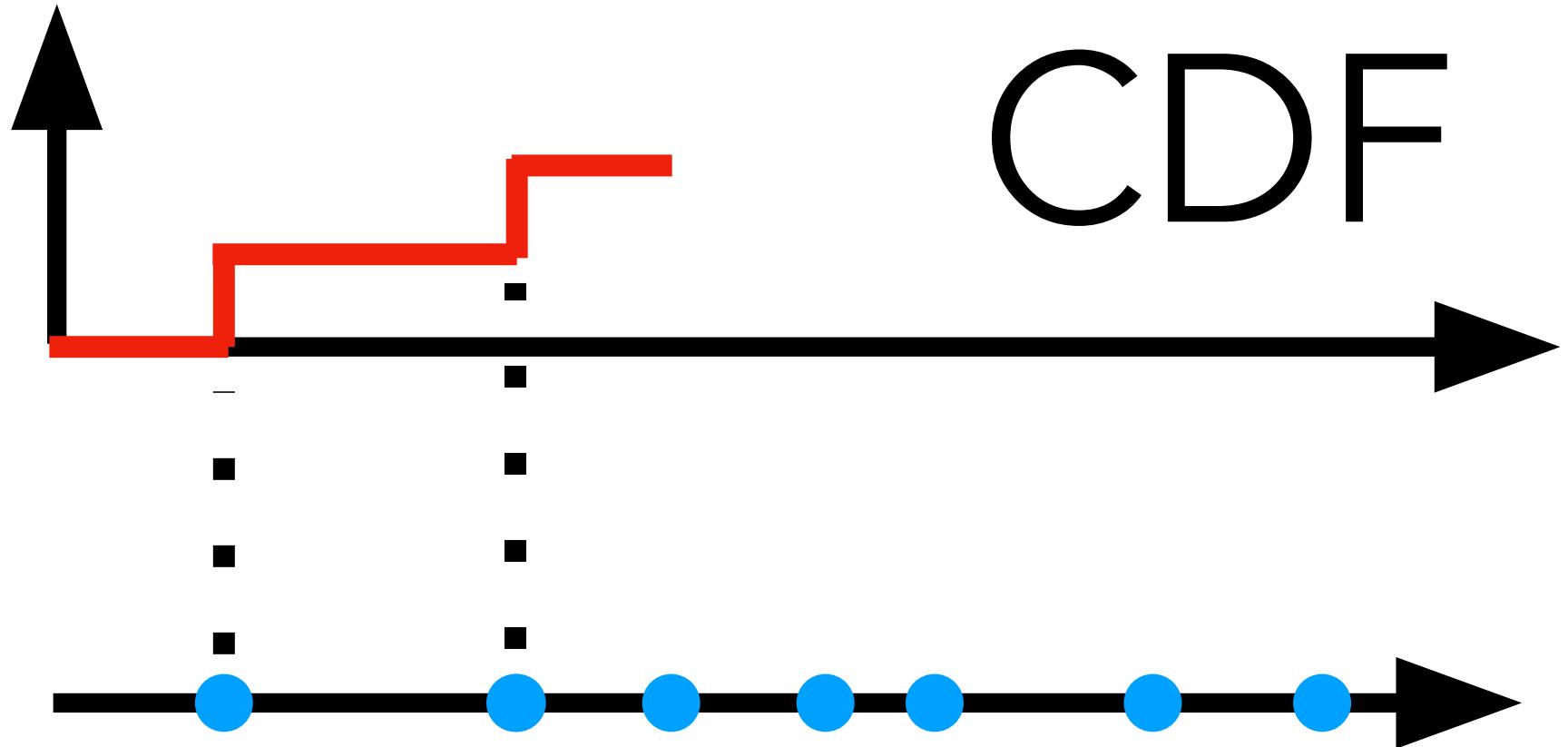
$x_2, +1,000,000,000$

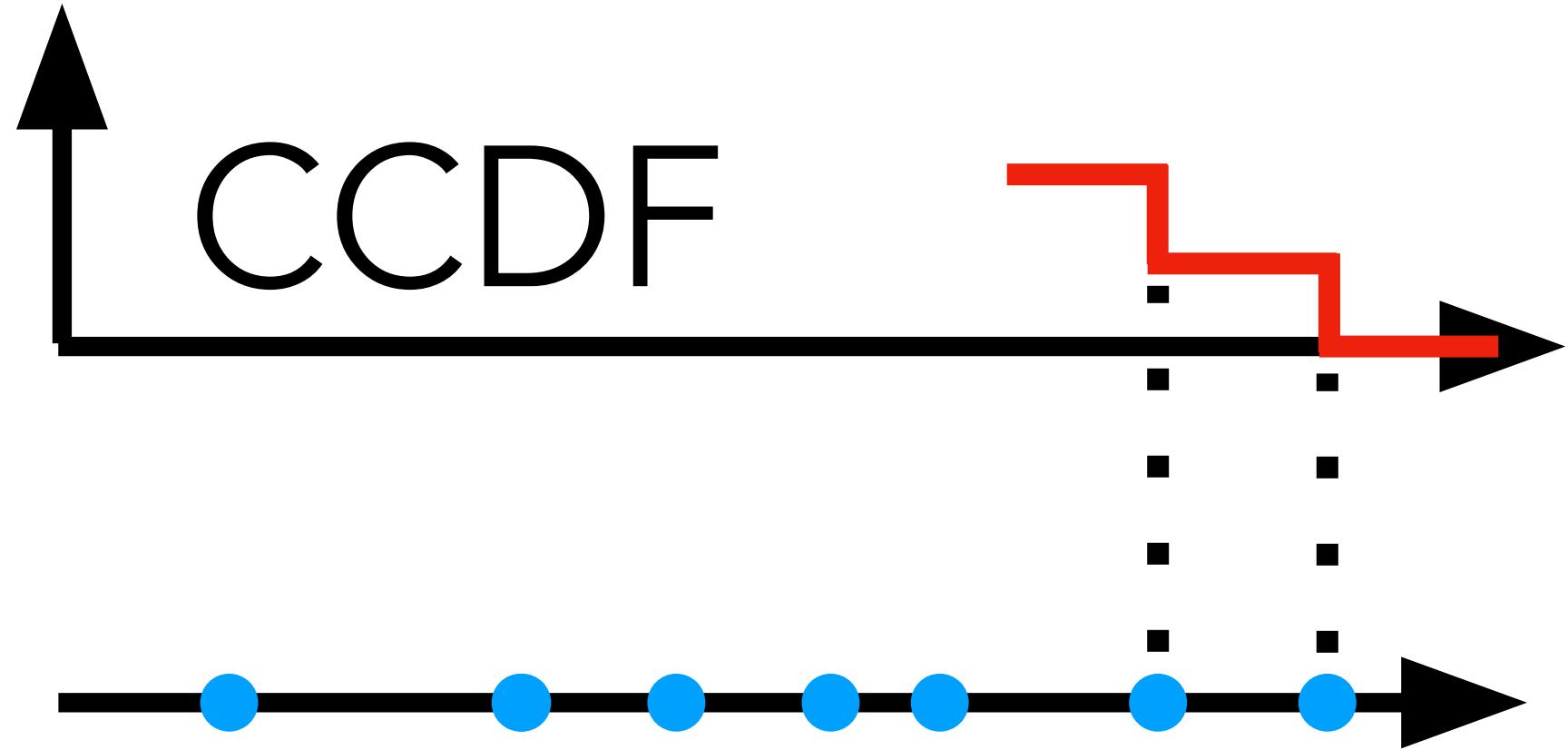


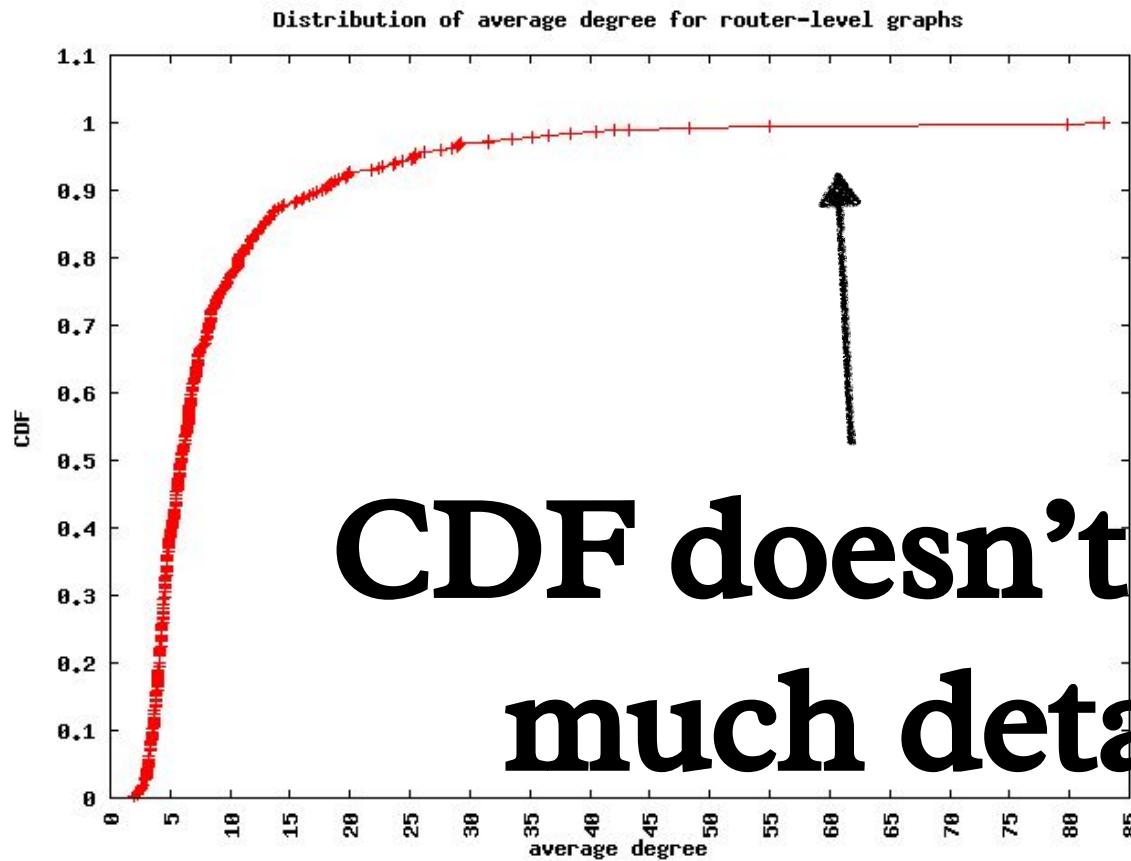


A nice  
alternative

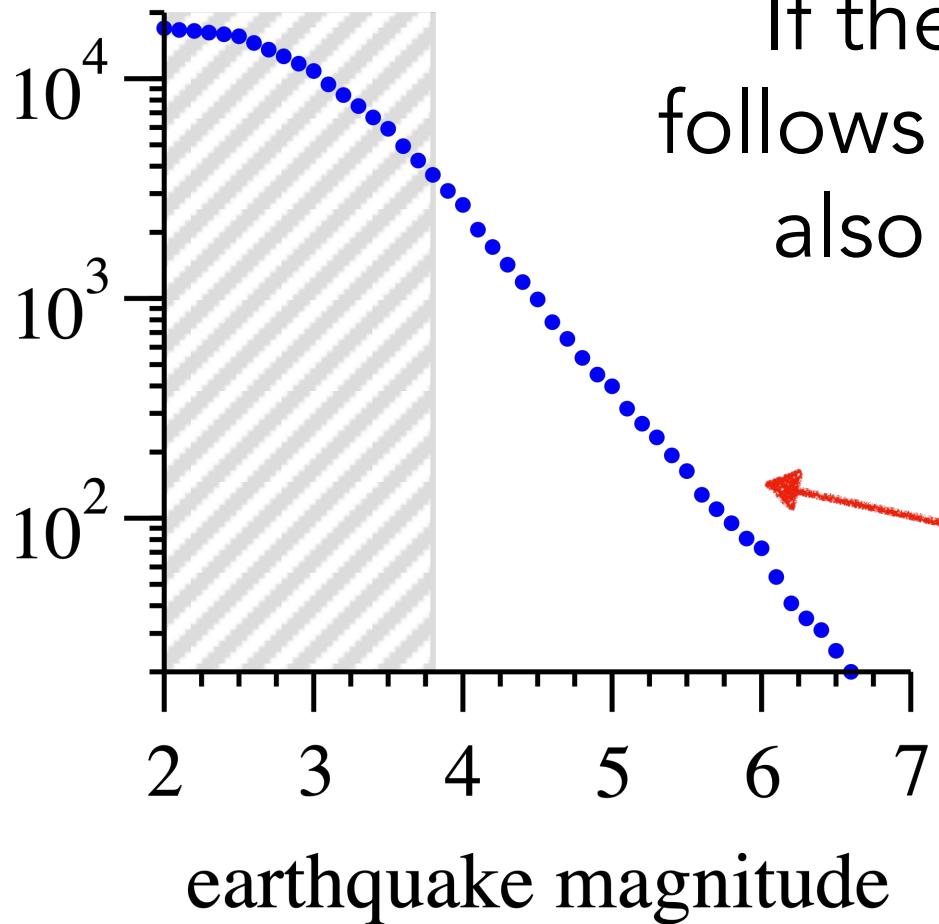
CDF







If the original distribution follows a power-law, the CCDF also follows a power-law.



**more useful**

$$CCDF(x) = P(X > x) \quad (\text{c.f. } CDF(x) = P(X \leq x))$$

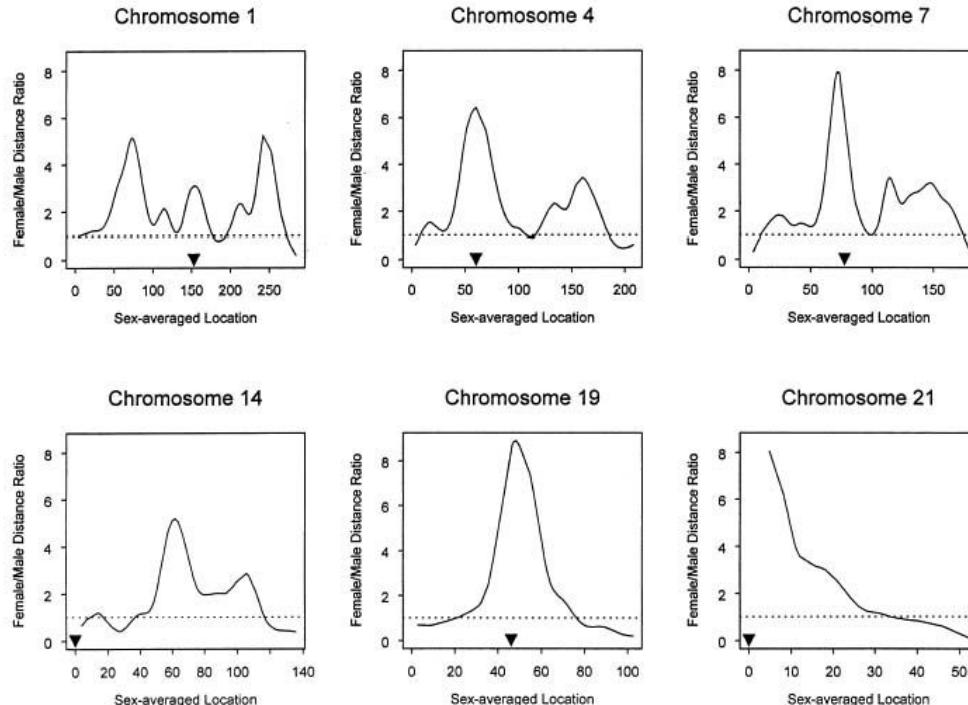
If  $f(x) = cx^{-\alpha}$ ,

$$\begin{aligned} CCDF(x) &= \int_x^{\infty} ct^{-\alpha} dt \\ &= \left[ \frac{c}{1-\alpha} t^{1-\alpha} \right]_x^{\infty} \\ &= \frac{c}{1-\alpha} x^{1-\alpha} = Cx^{-(\alpha-1)} \end{aligned}$$

When visualizing heavy-tailed distributions such as power-law distribution,

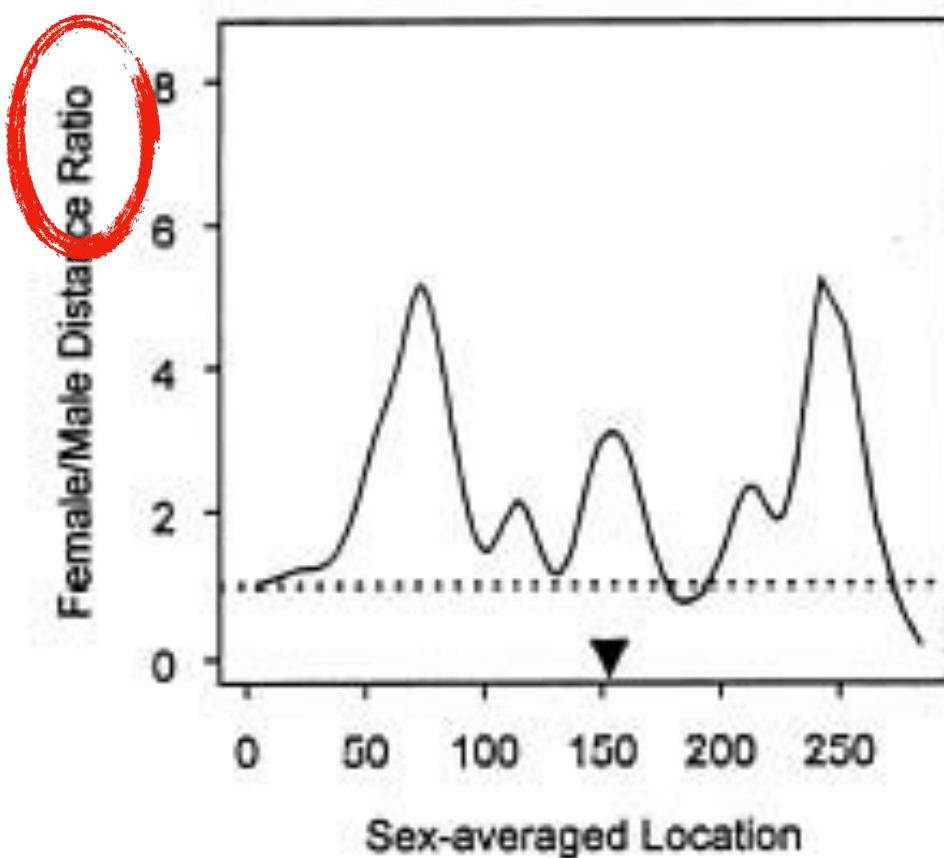
CCDF is often much more useful than CDF.

# What's wrong with these plots?

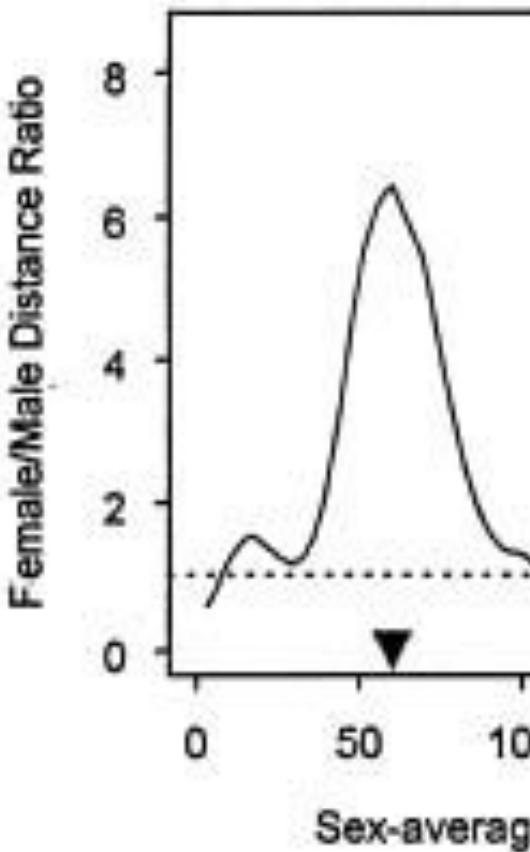


**Figure 1** Plots of the female:male genetic-distance ratio against sex-averaged genetic location (in cM) along six selected chromosomes. Approximate locations of the centromeres are indicated by the triangles. The dashed lines correspond to equal female and male distances.

# Chromosome 1



# Chrom



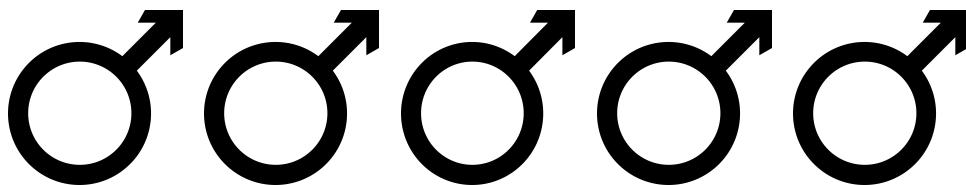
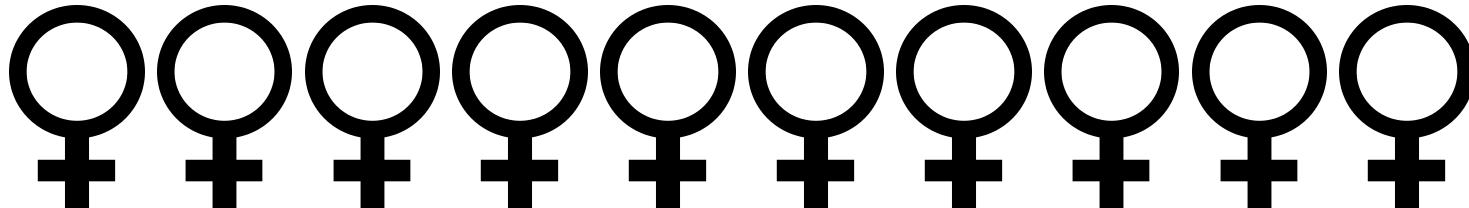
What is a ratio?

$A : B$

Numerator

Denominator

$A / B$



10:5

2:1

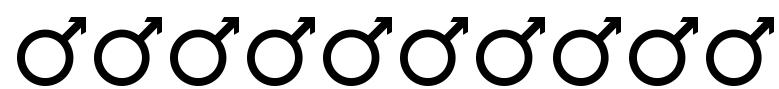
$10/5 = 2$



10:5

2:1

$$10/5 = 2$$



10:5

5:10

2:1

1:2

$$10/5 = 2$$

$$5/10 = 0.5$$

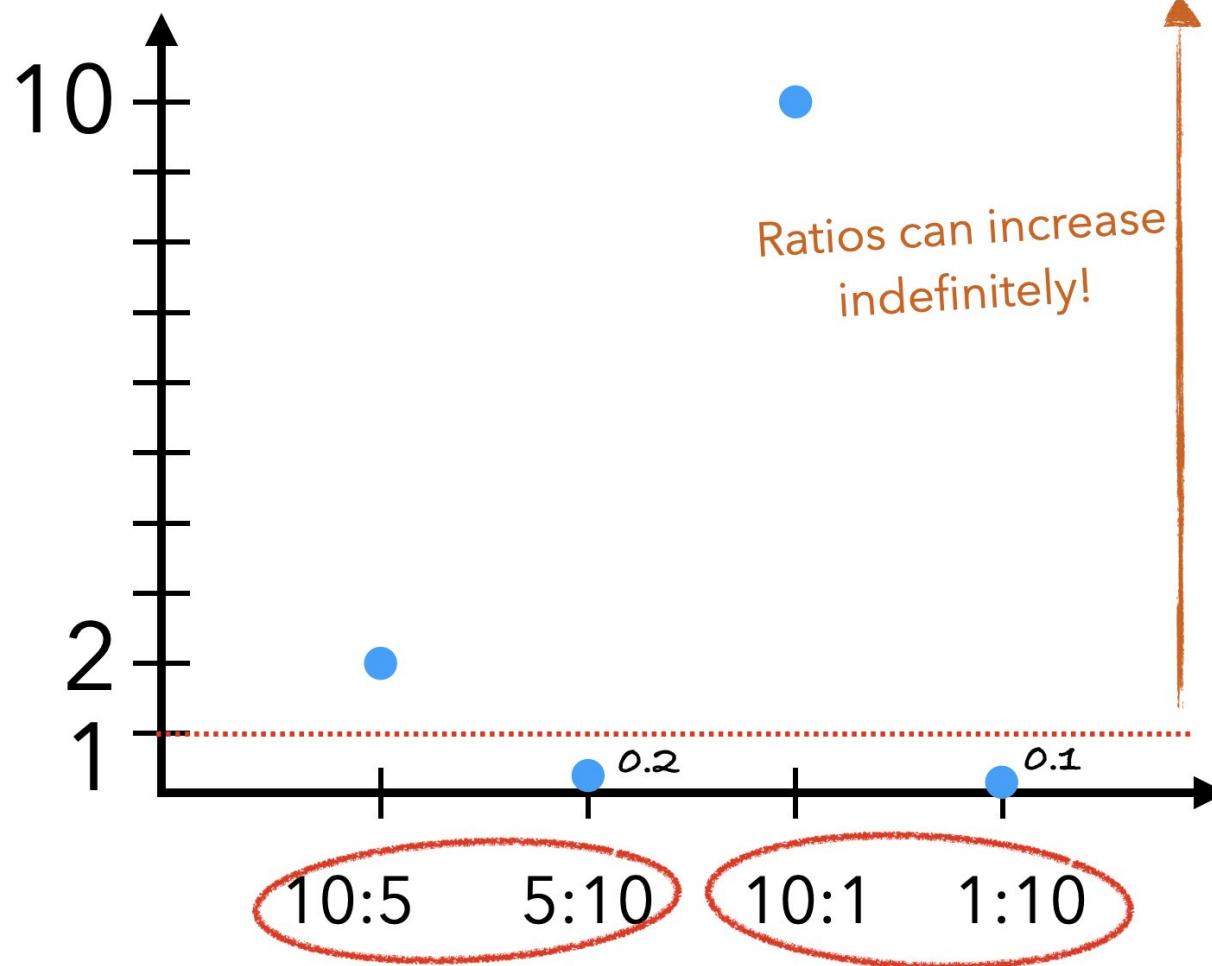


10:1

1:10

10

0.1



Log scale!

$$\log \frac{A}{B} = \log A - \log B$$

$$\log \frac{10}{5} = \underline{\log 2}$$

$$\log \frac{5}{10} = \log 1 - \log 2 = \underline{-\log 2}$$

$$\log \frac{10}{1} = \underline{\log 10}$$

$$\log \frac{1}{10} = \log 1 - \log 10 = -\underline{\log 10}$$

