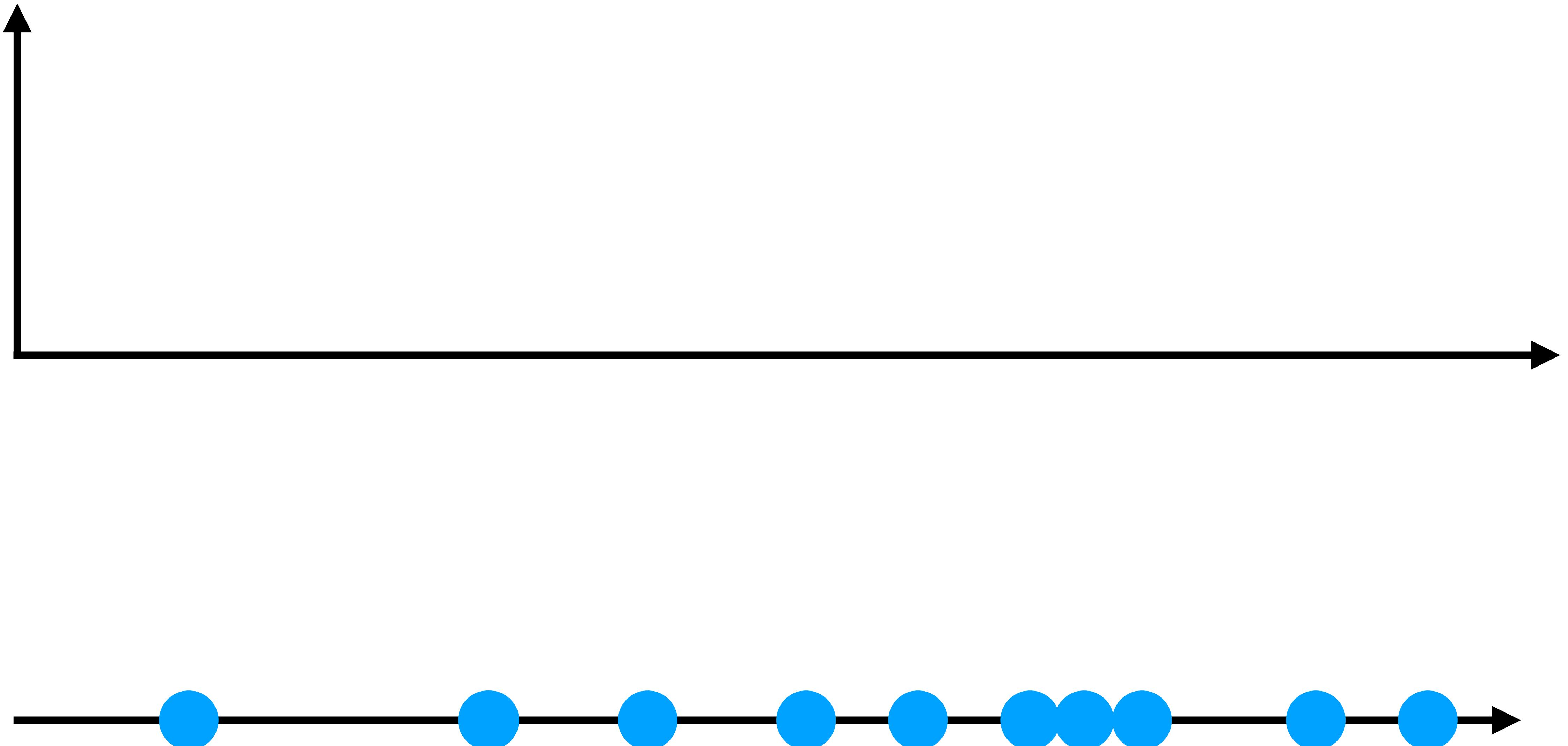
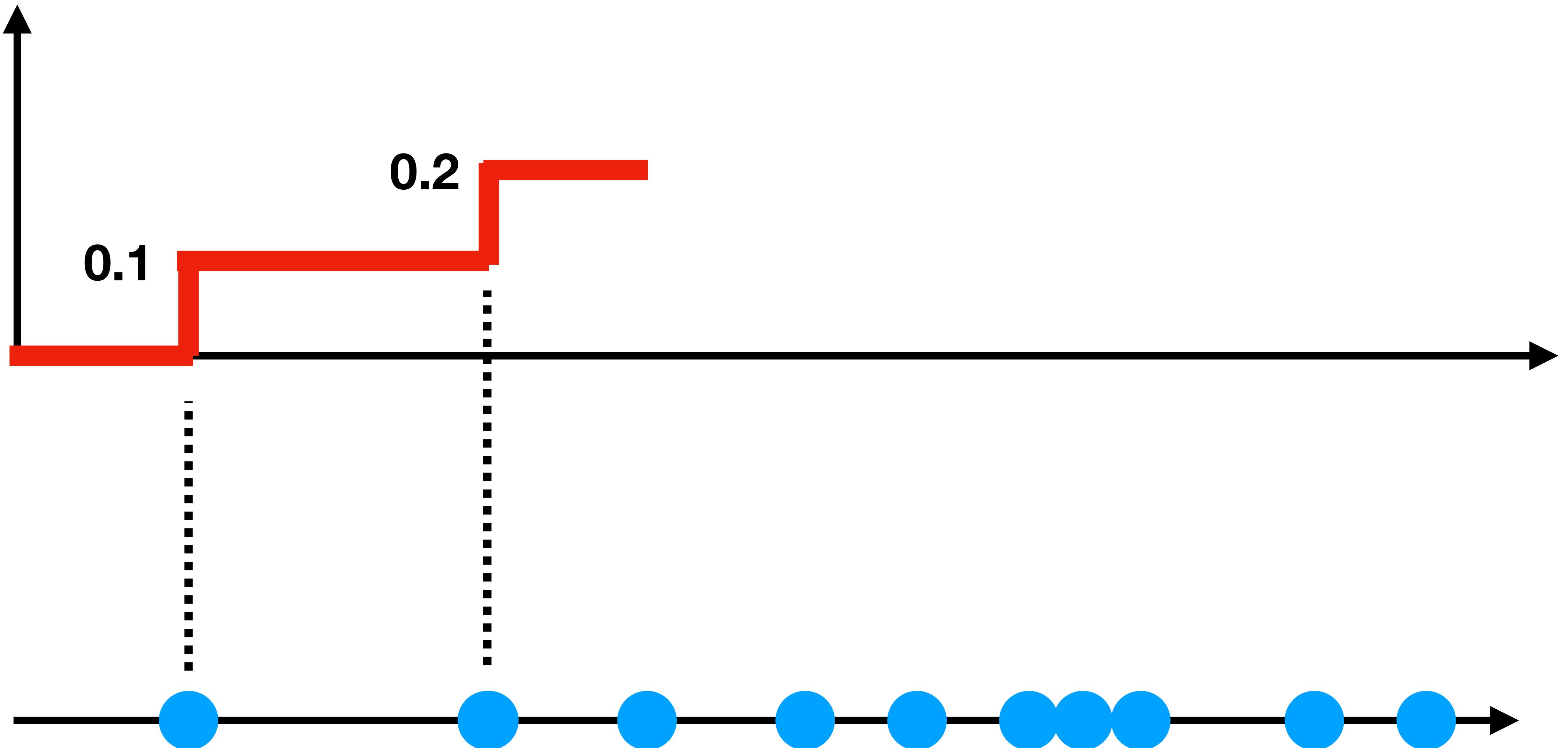


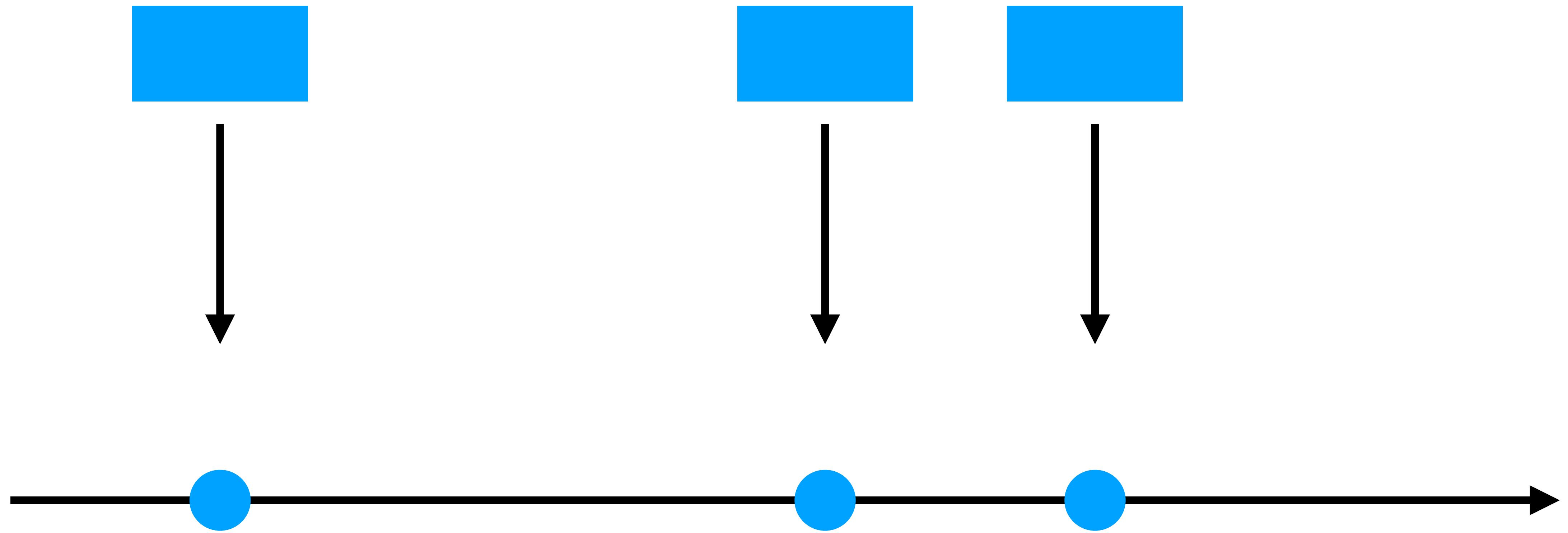
Quiz

- Draw an empirical CDF plot using the following data points:
 $\{10, 1, 2, 5, 8, 10, 5\}$
- Explain how KDE (Kernel Density Estimation) works.
- Explain the choices that you have when creating a KDE.

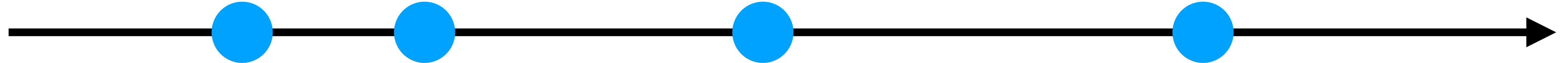


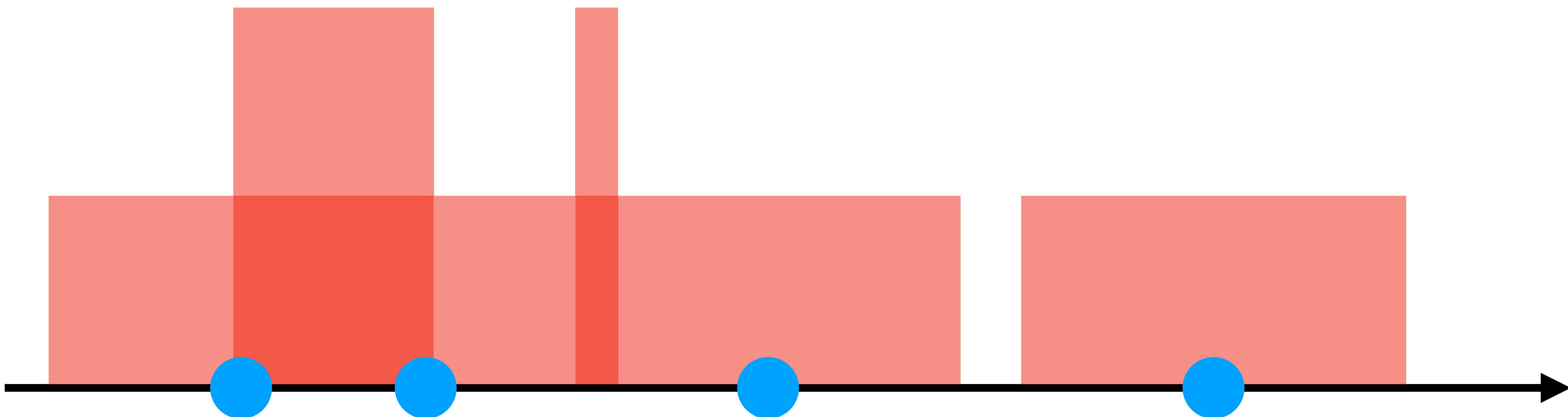




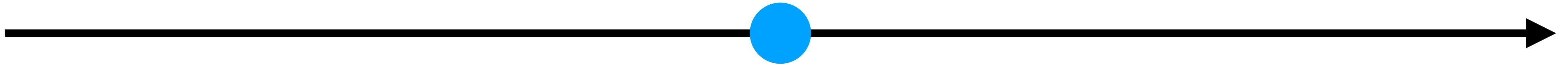


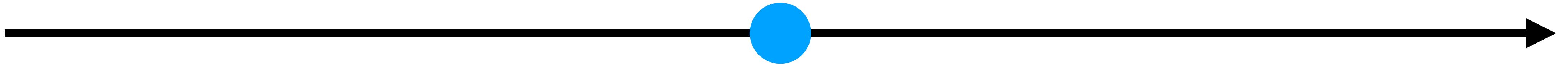


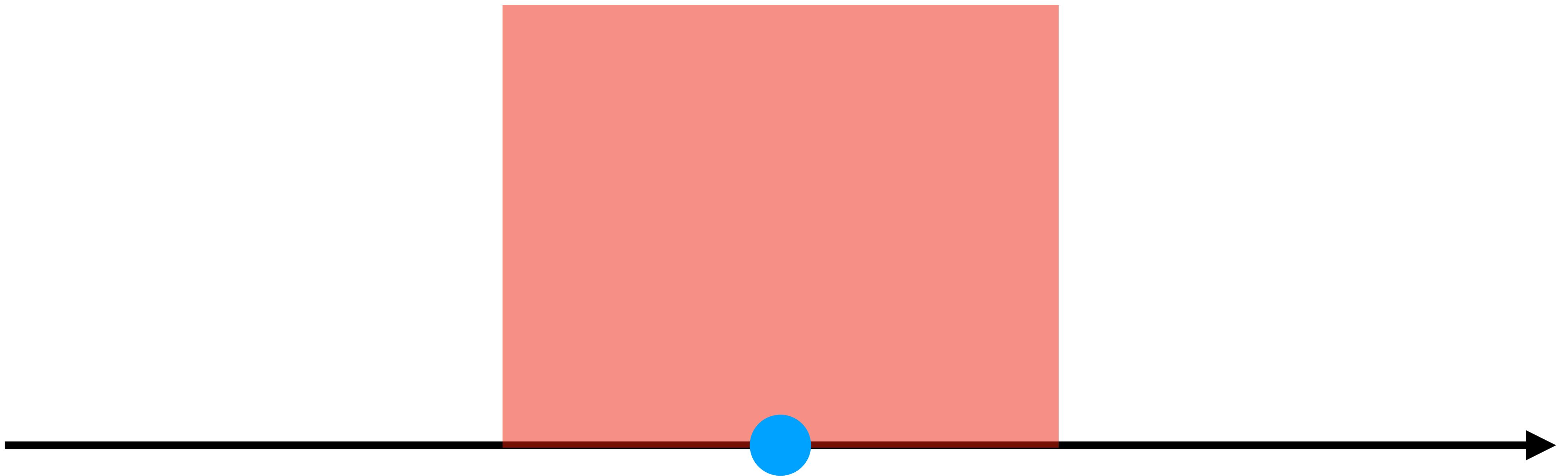


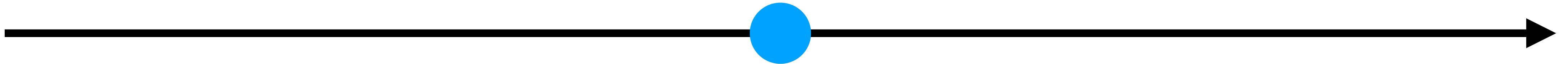


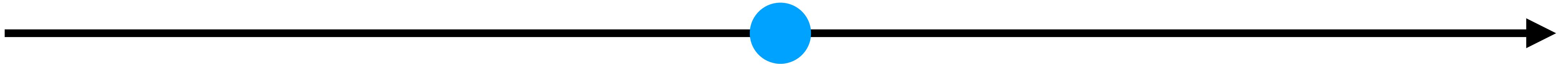


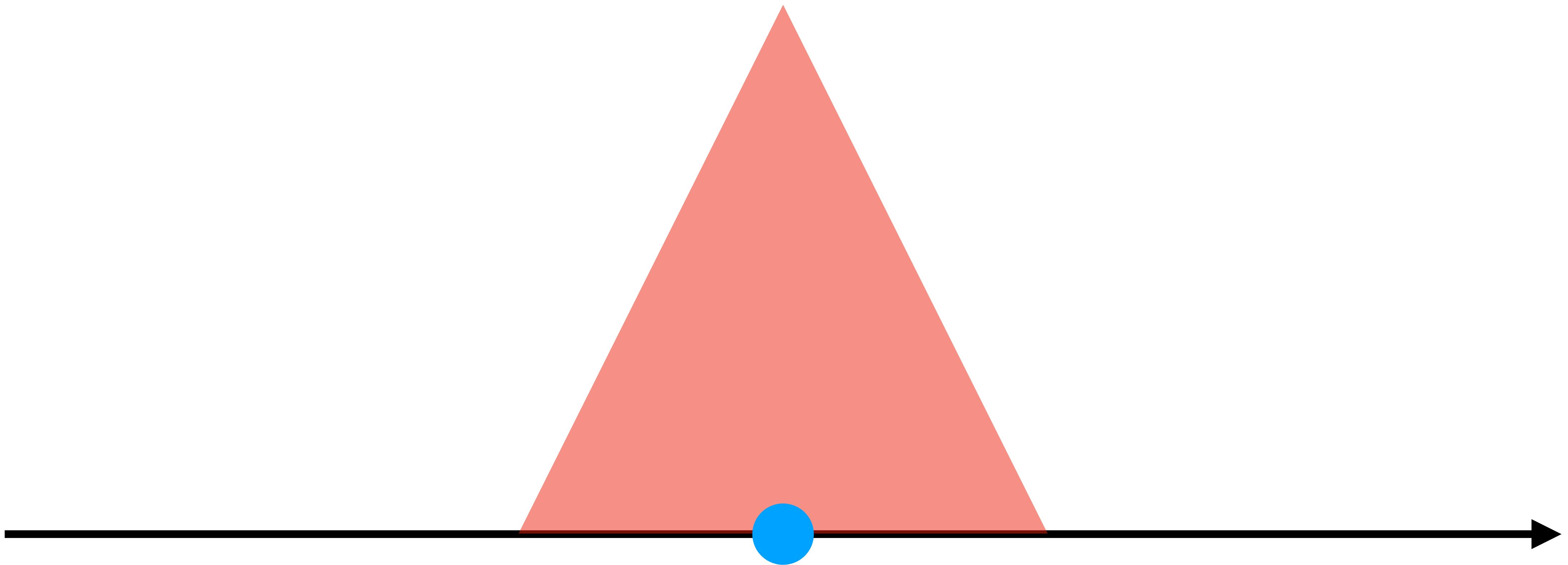


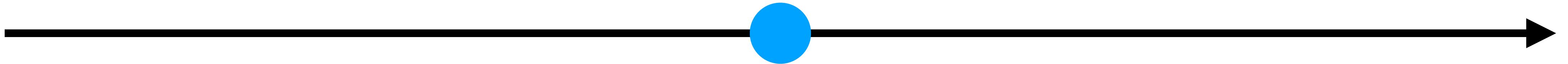


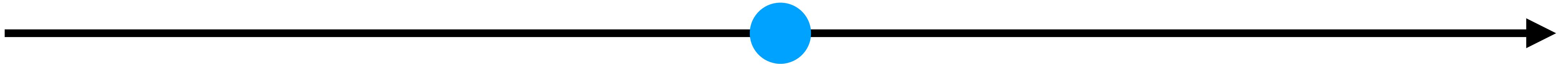


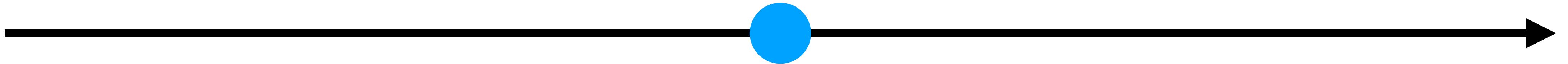


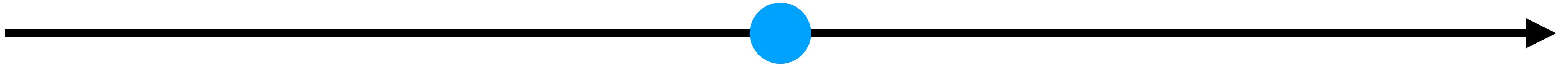


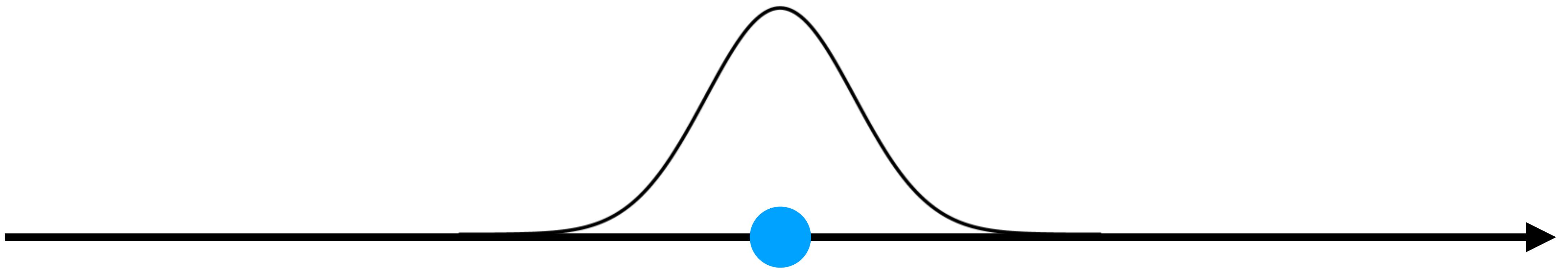


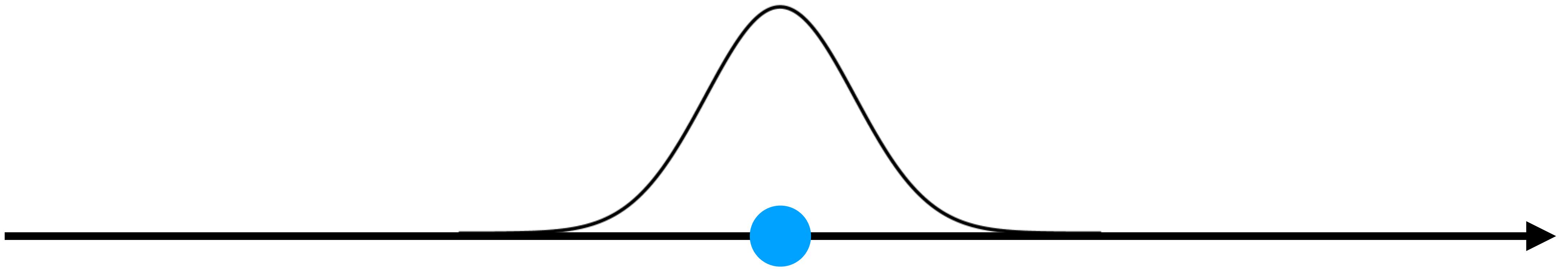


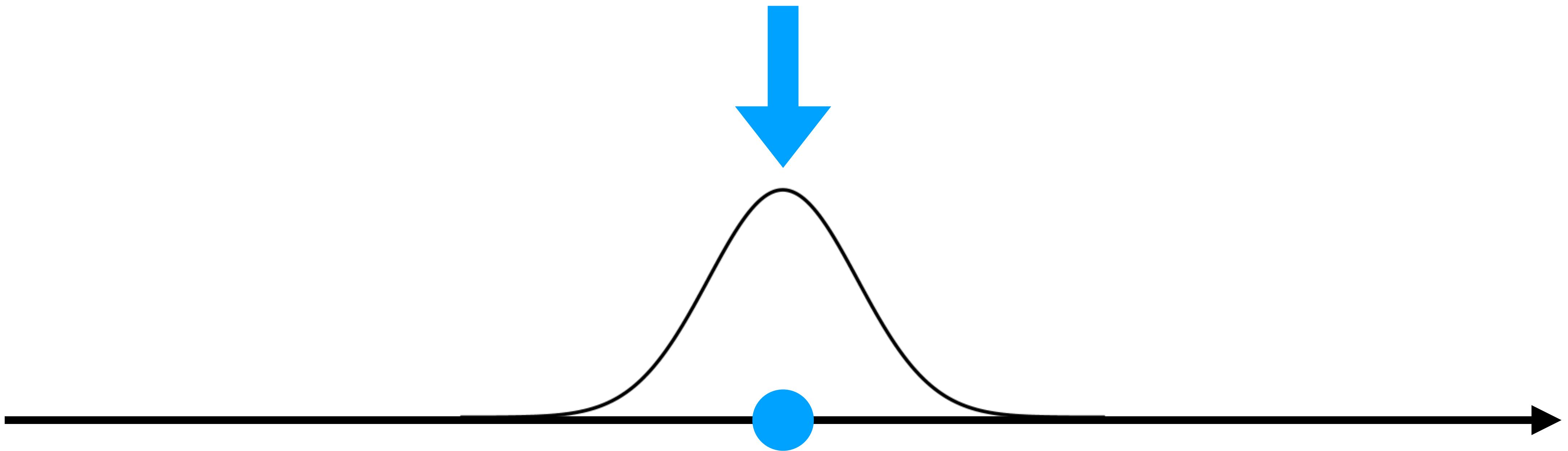


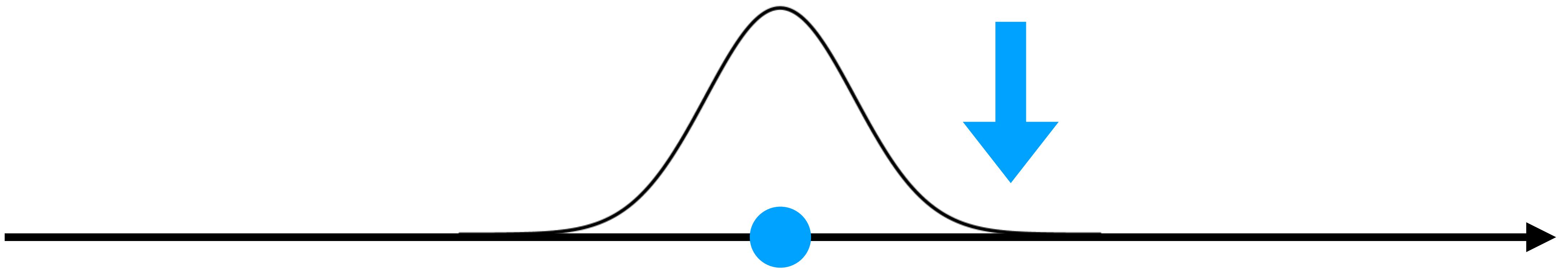




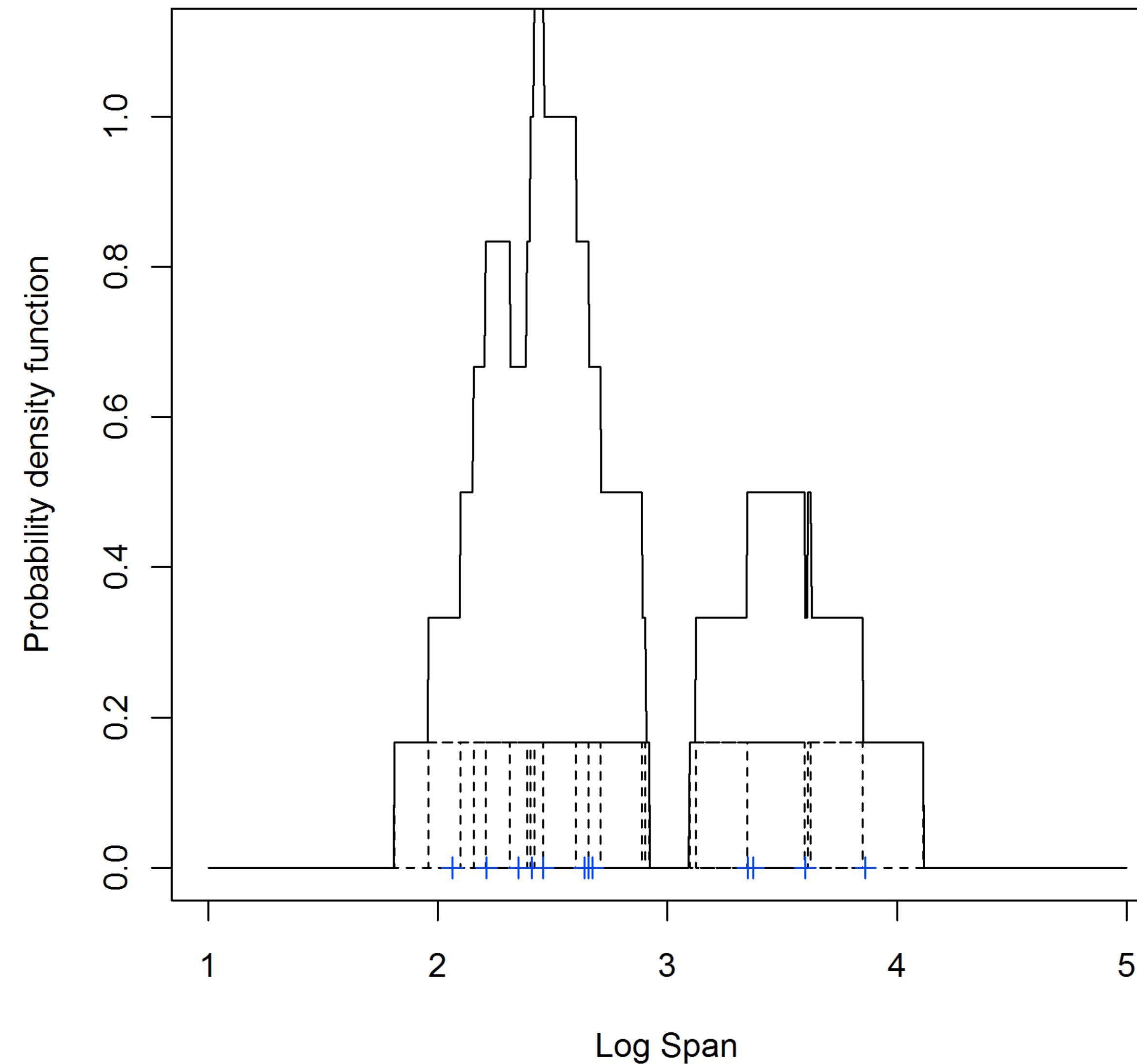


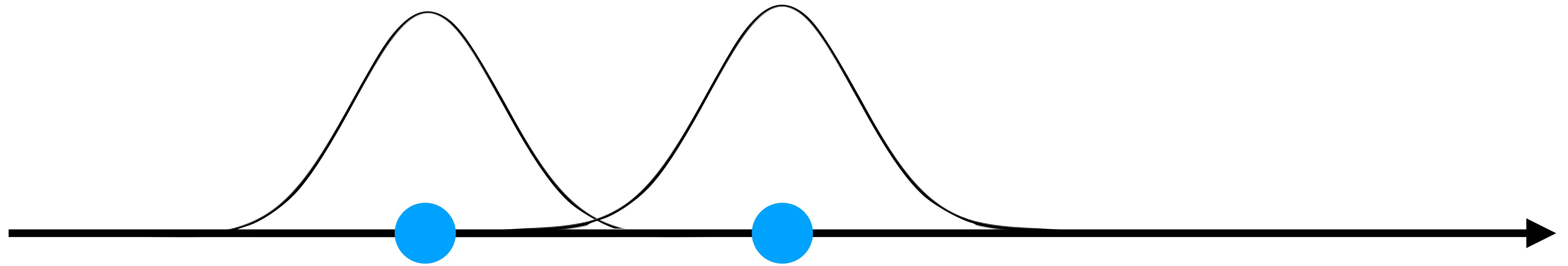




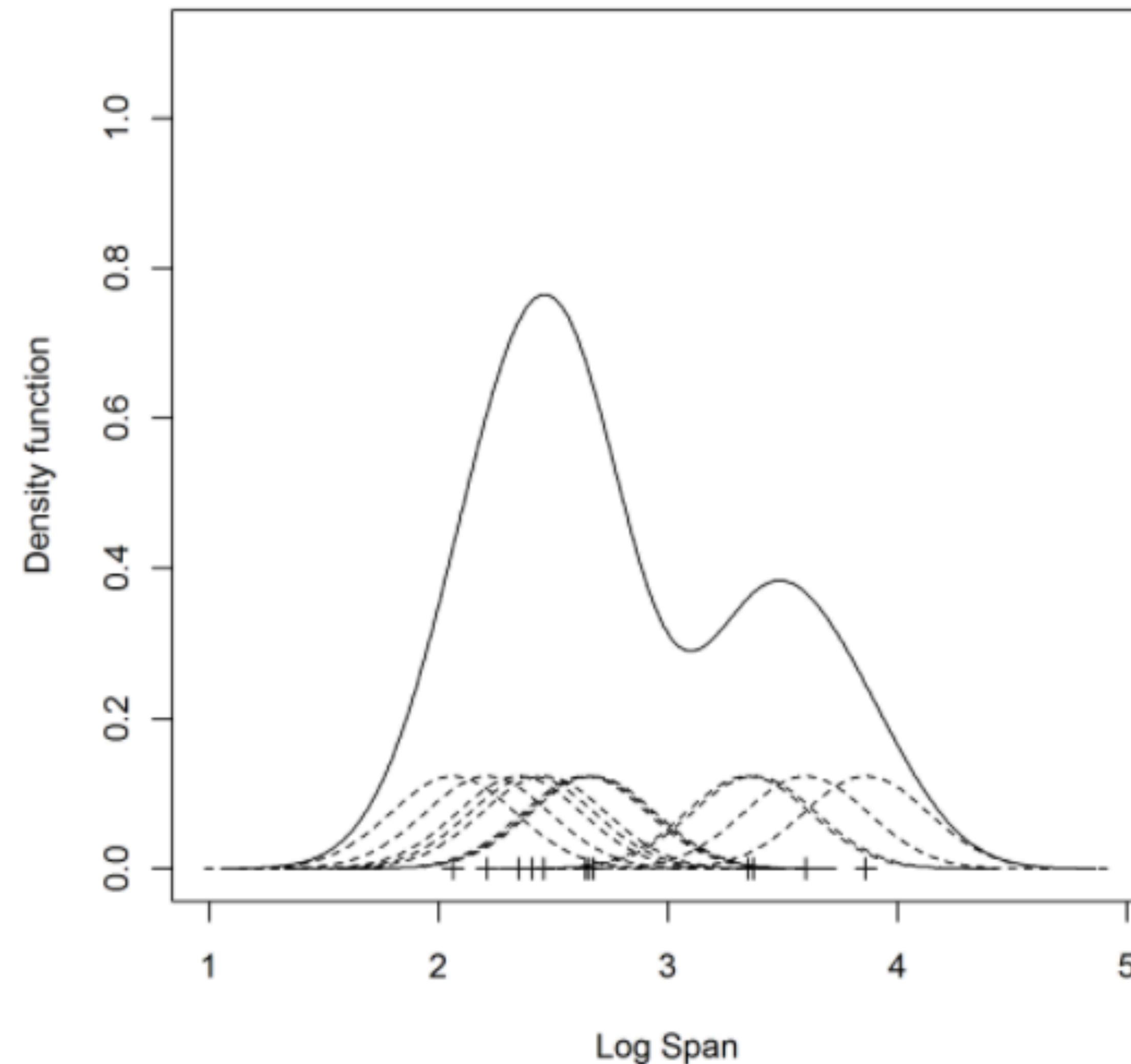


'Histogram' with blocks centred over data points

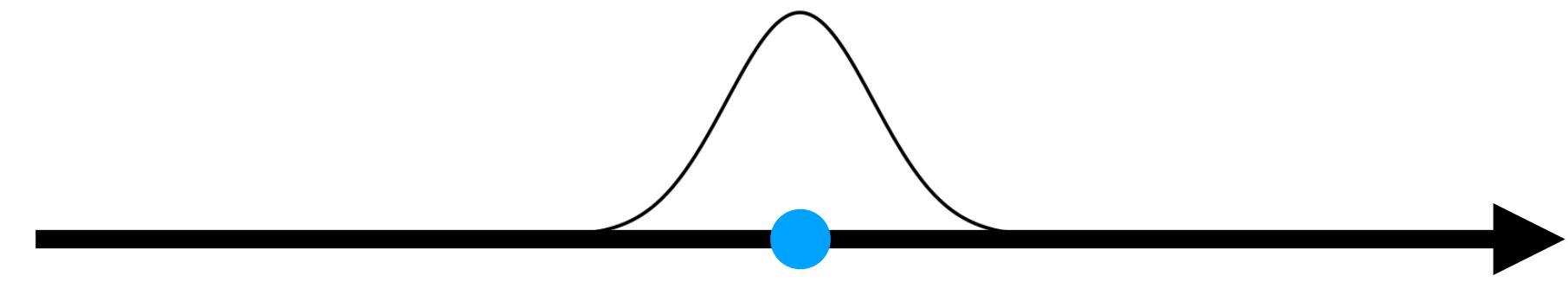
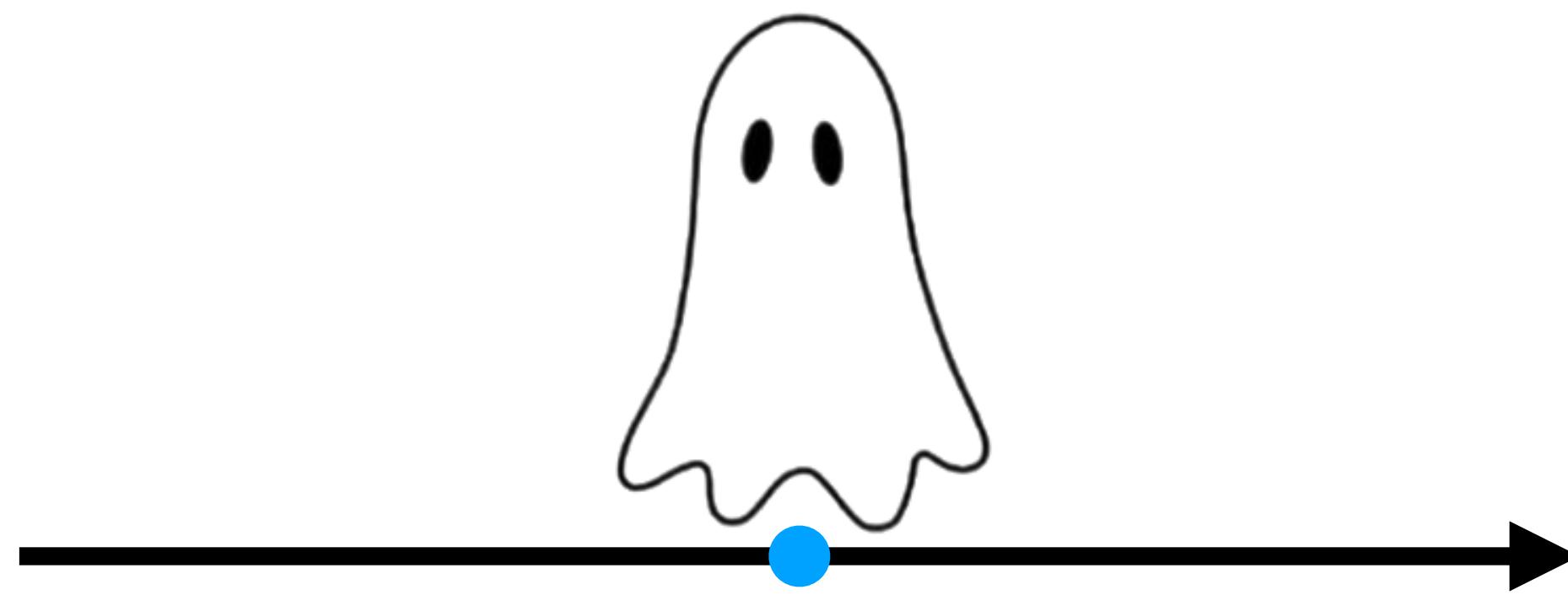
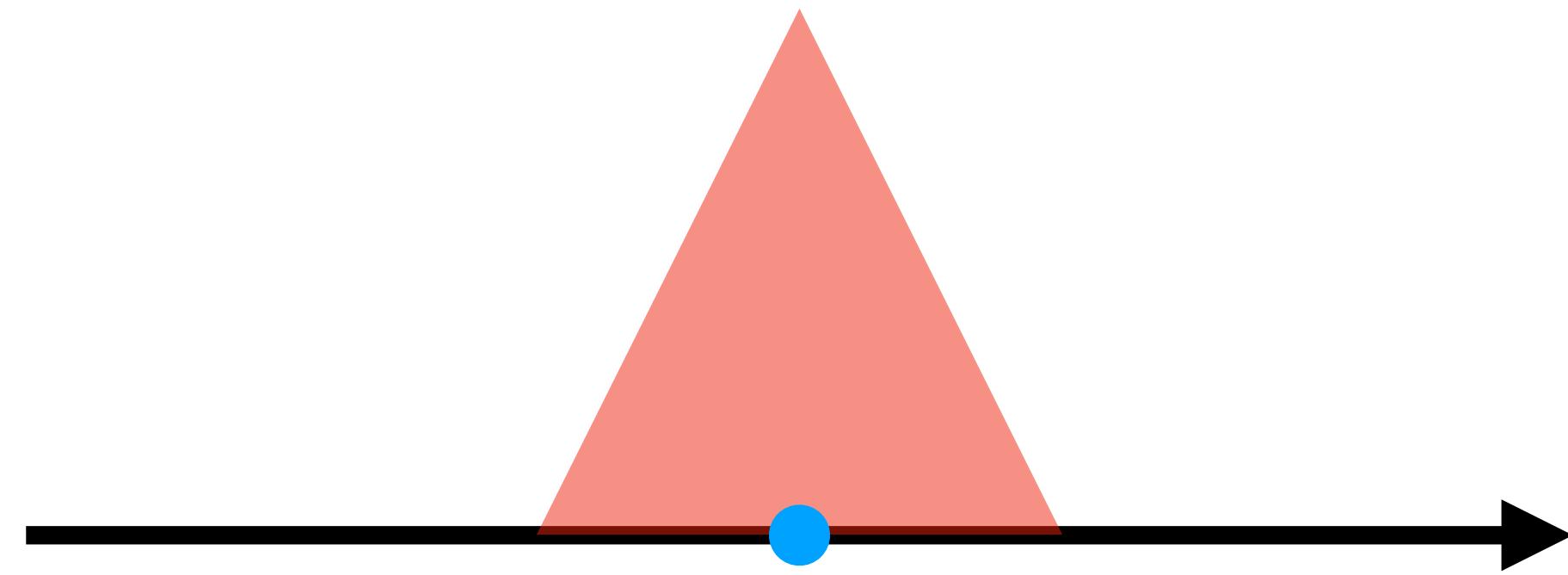
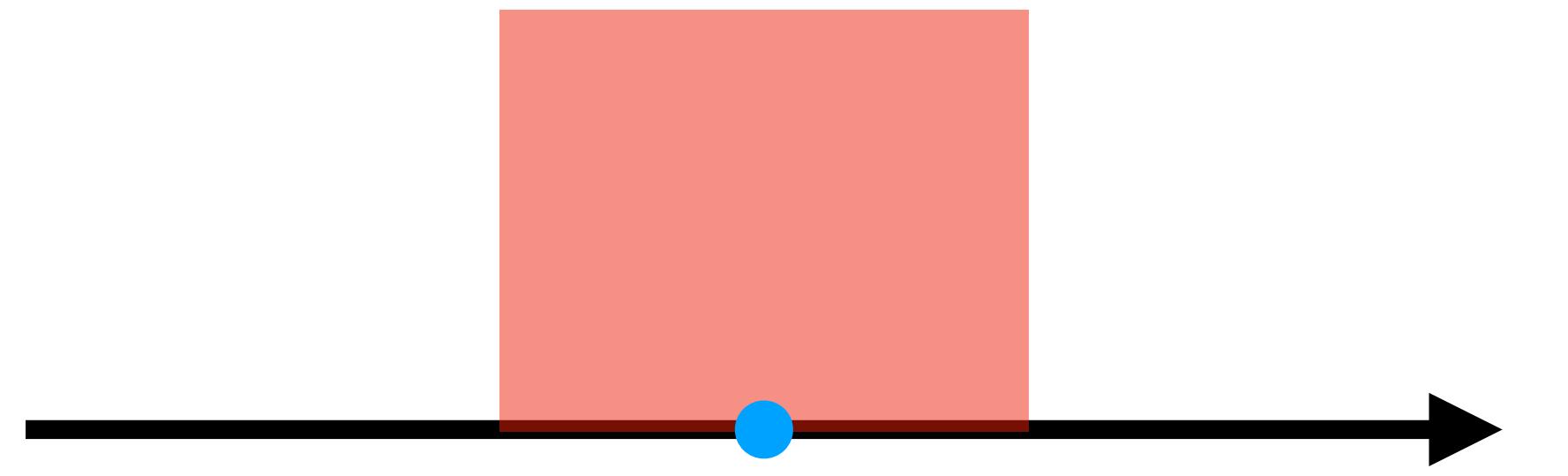




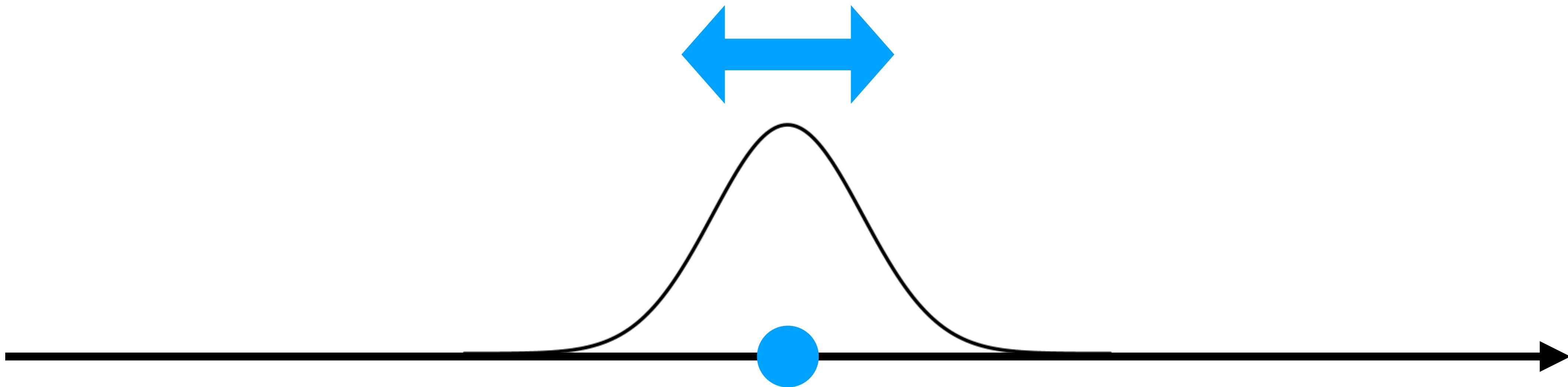
Kernel Density Estimation (KDE)

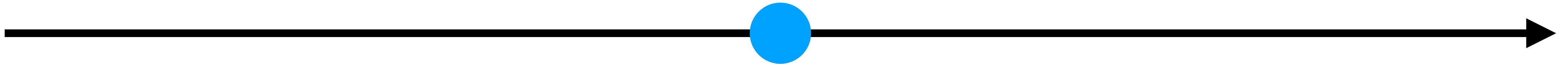


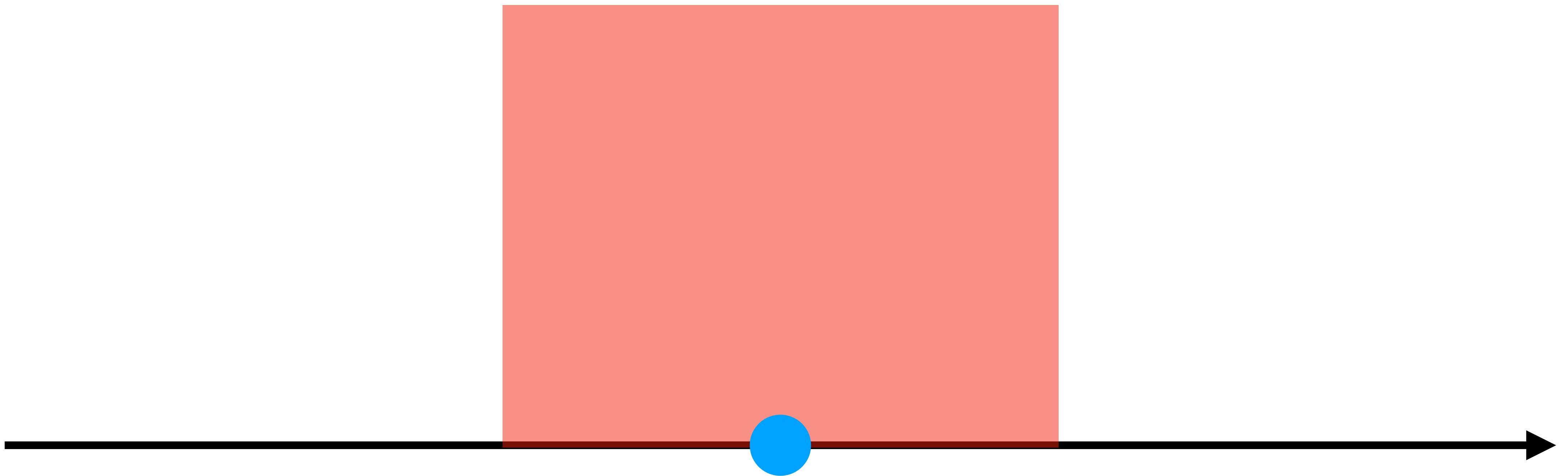
"Kernels"

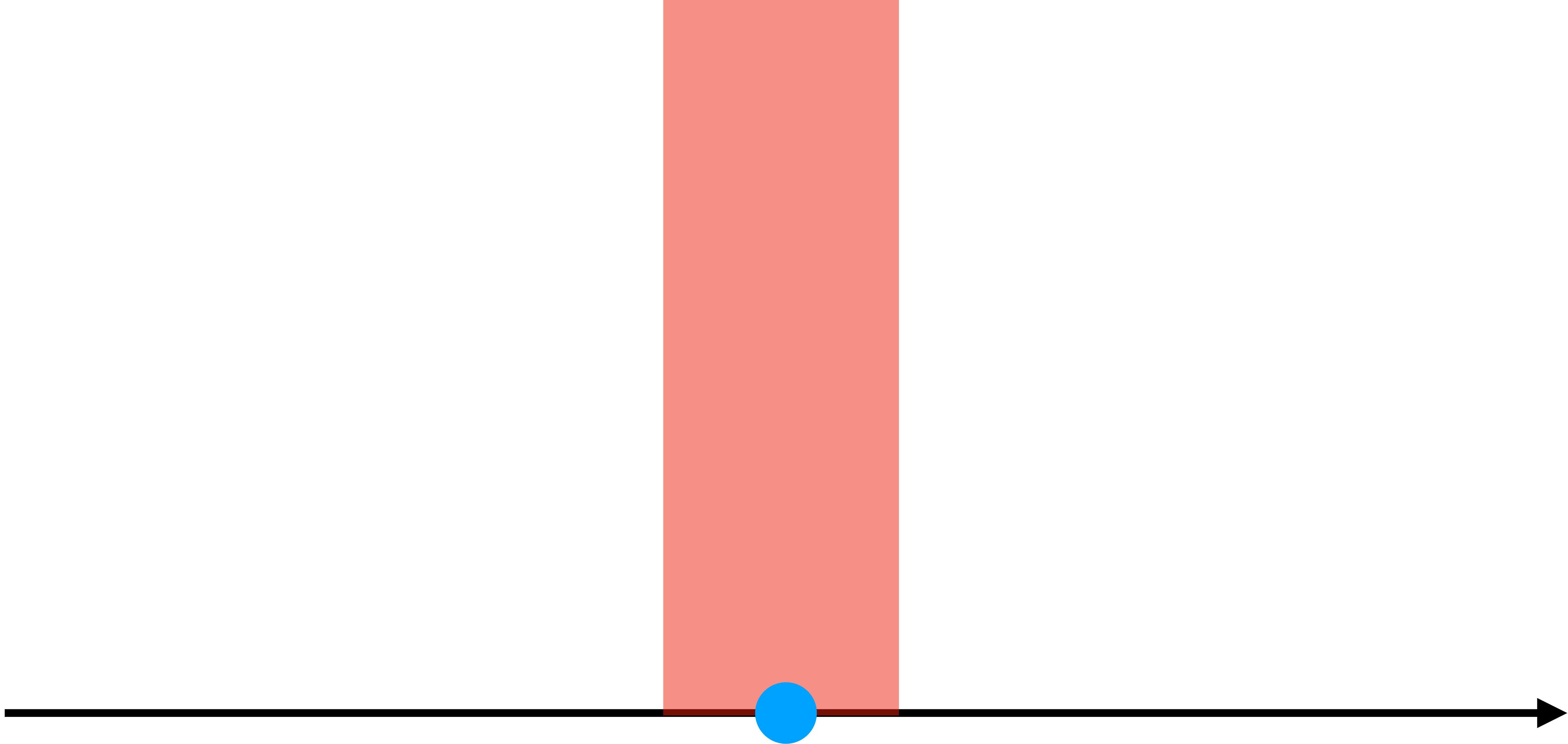


Bandwidth









Interpolation & Extrapolation

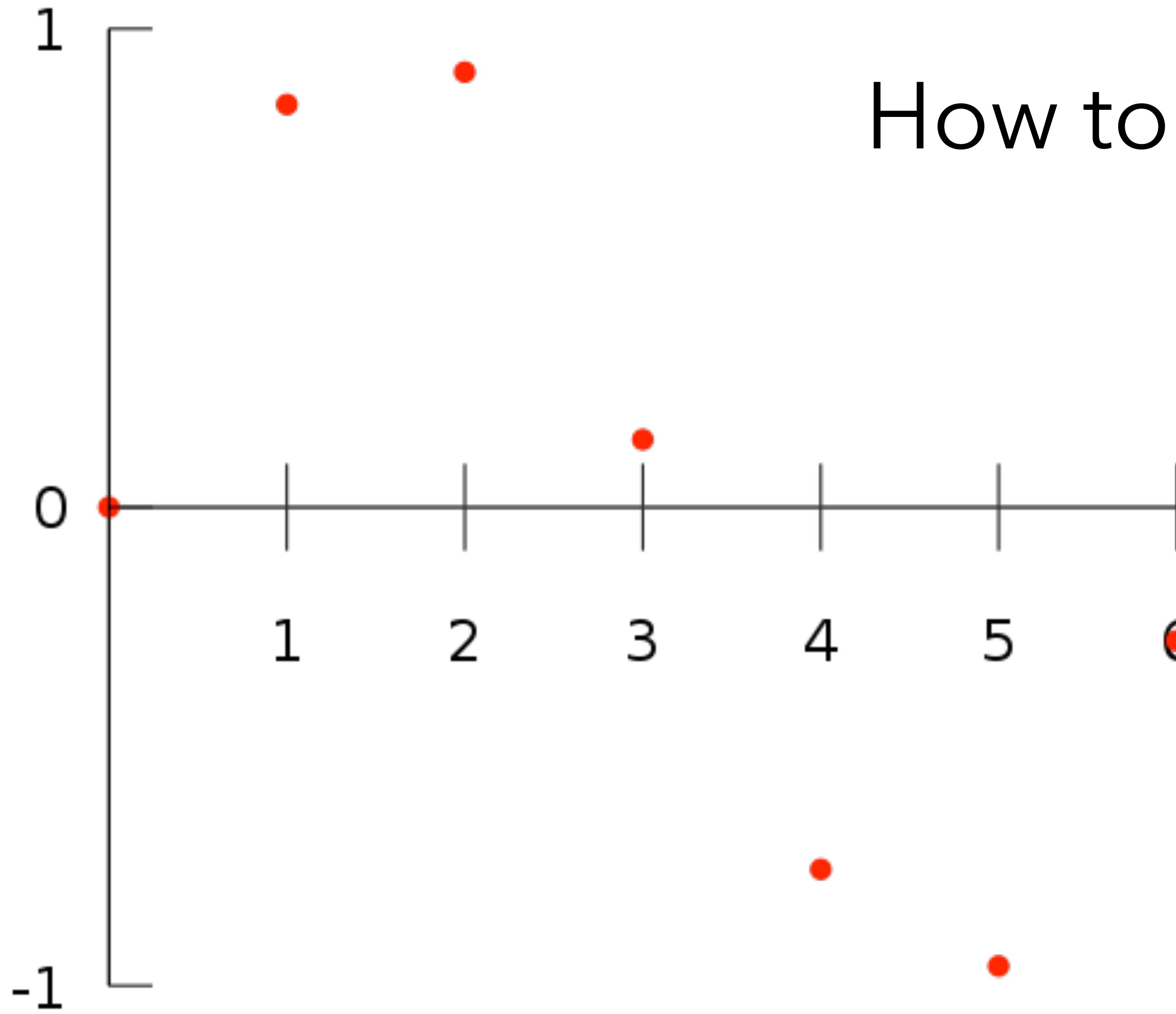
Interpolation

How can we fill the gaps between
(or beyond) data points?

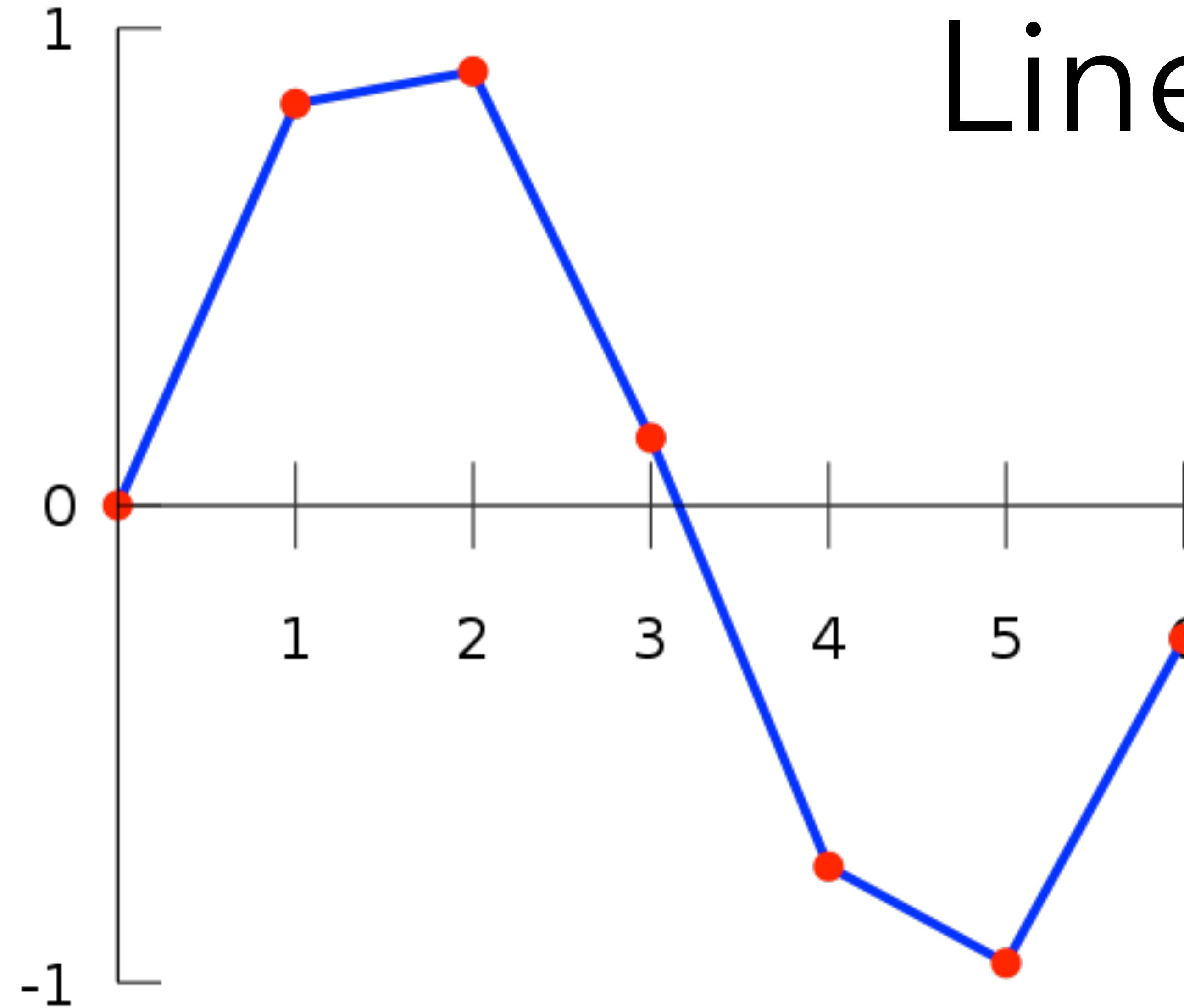
Interpolation

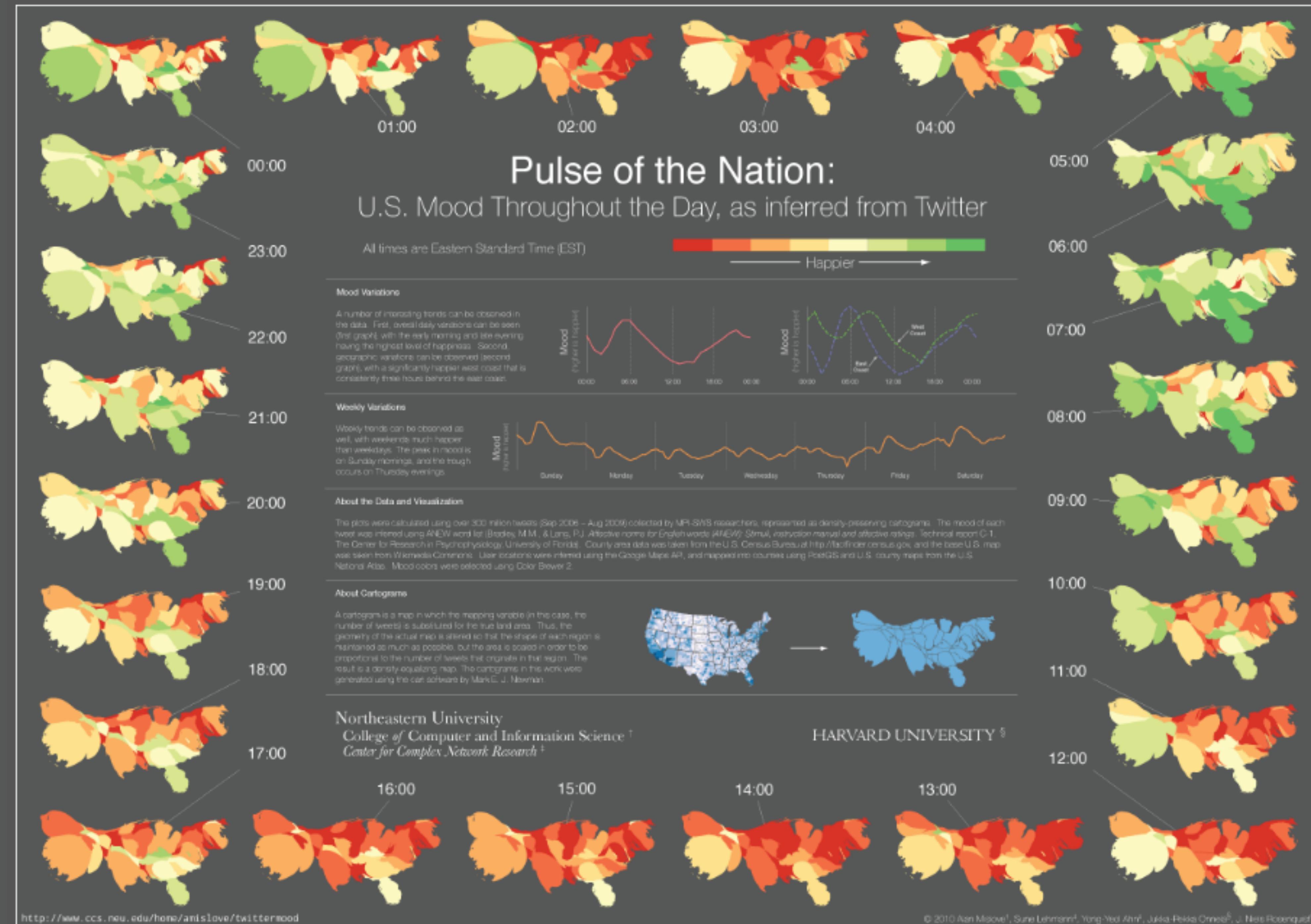
Let's connect the dots.

How to interpolate?

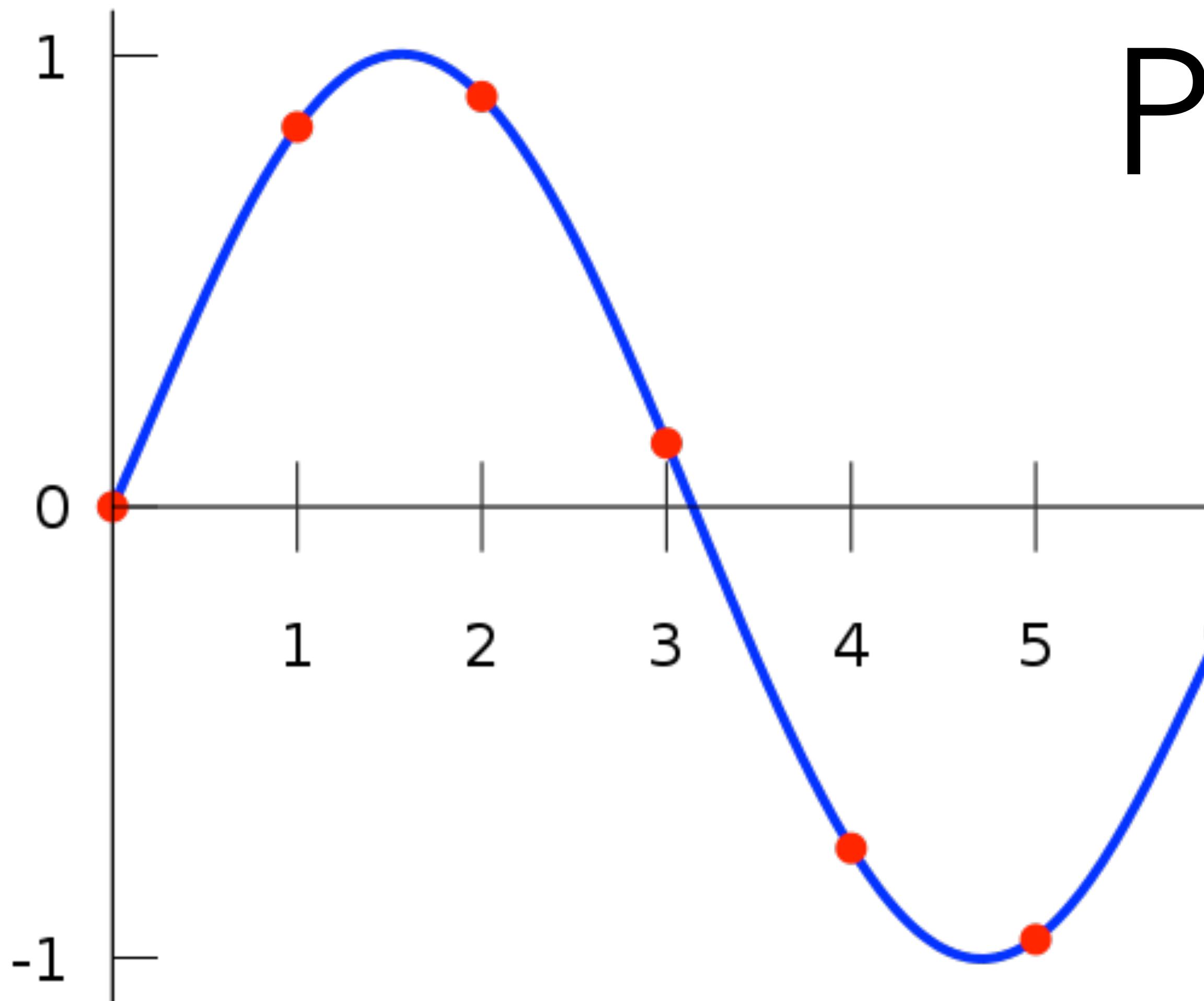


Linear



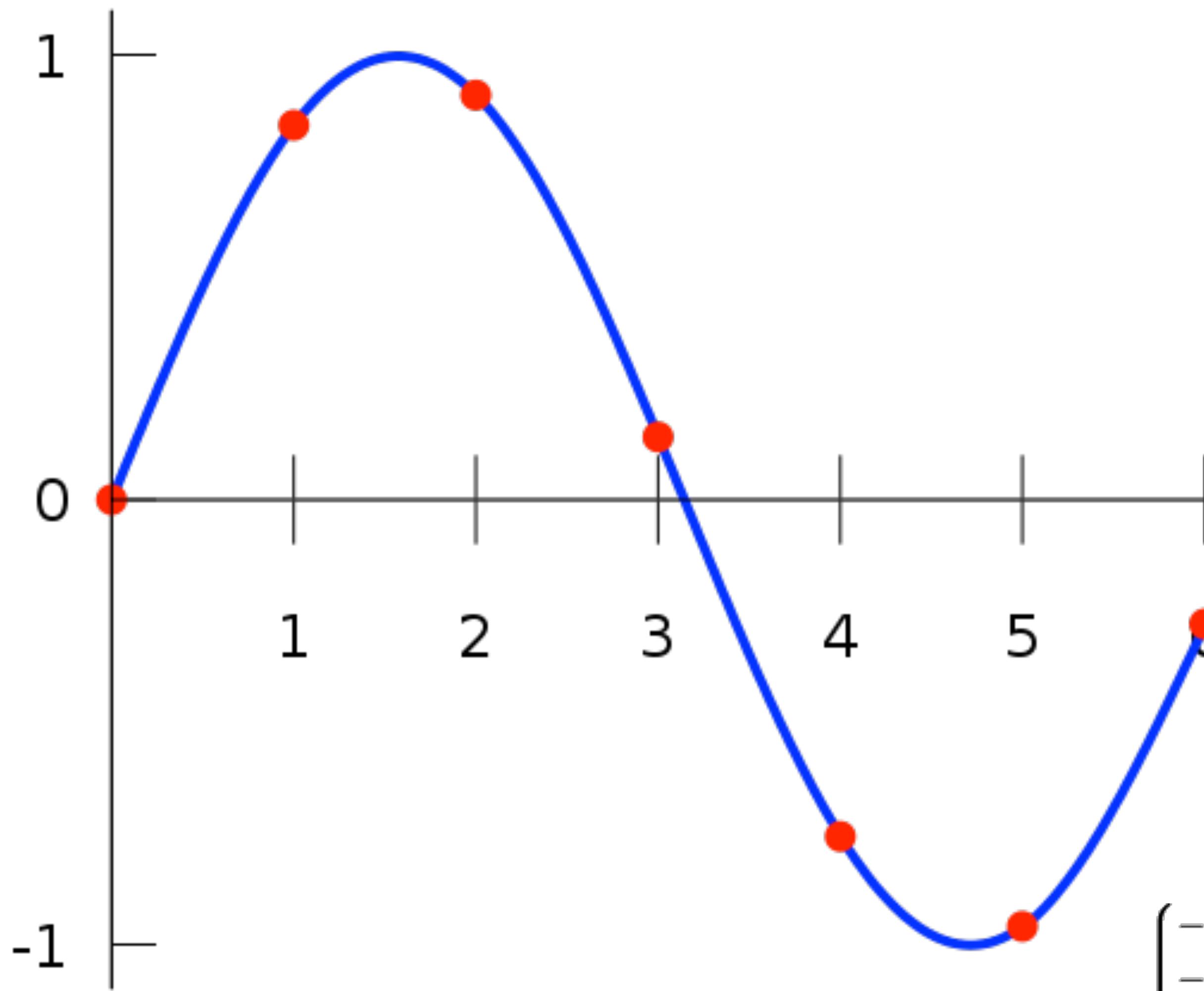


Polynomial



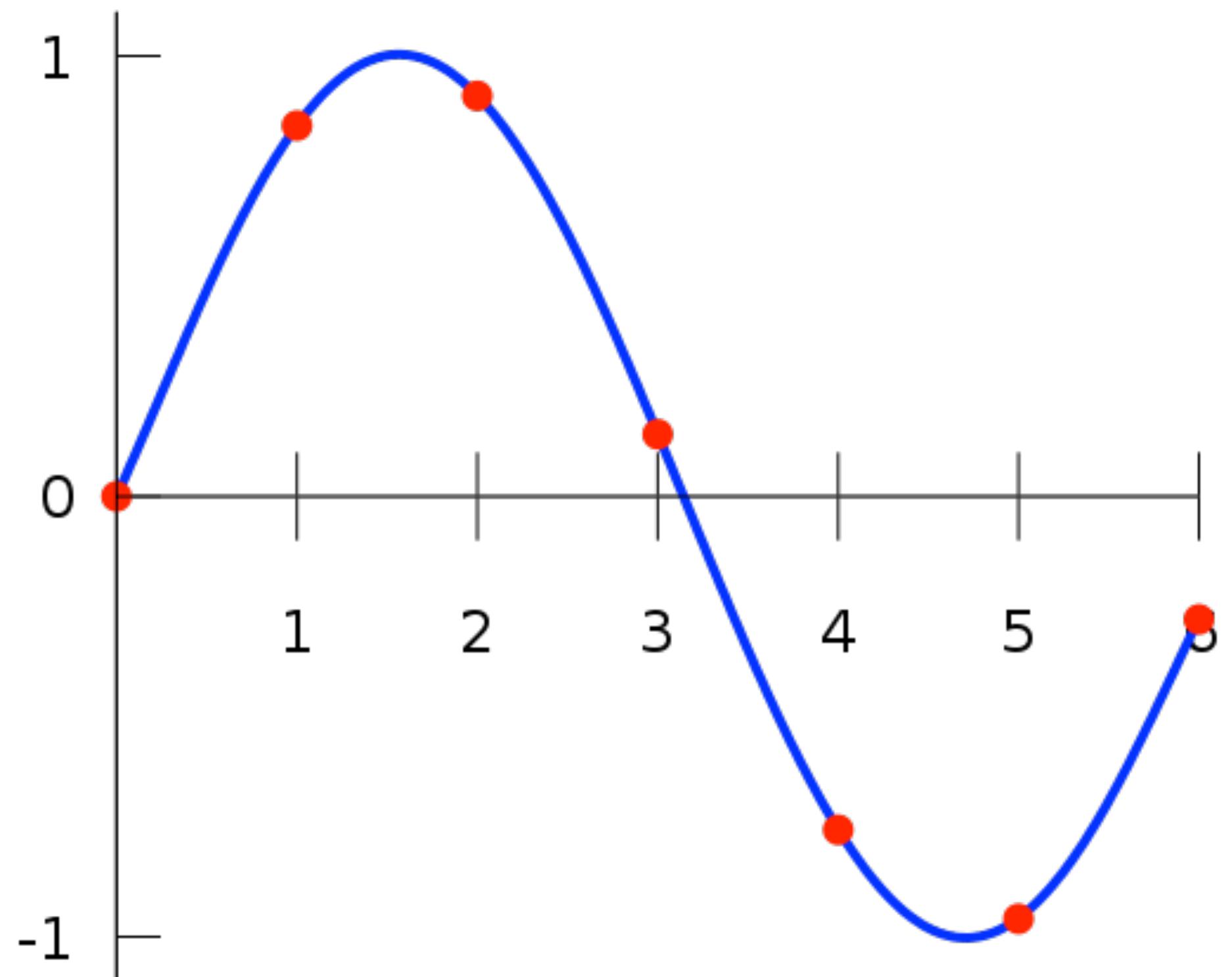
$$p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_2 x^2 + a_1 x + a_0.$$

Splines

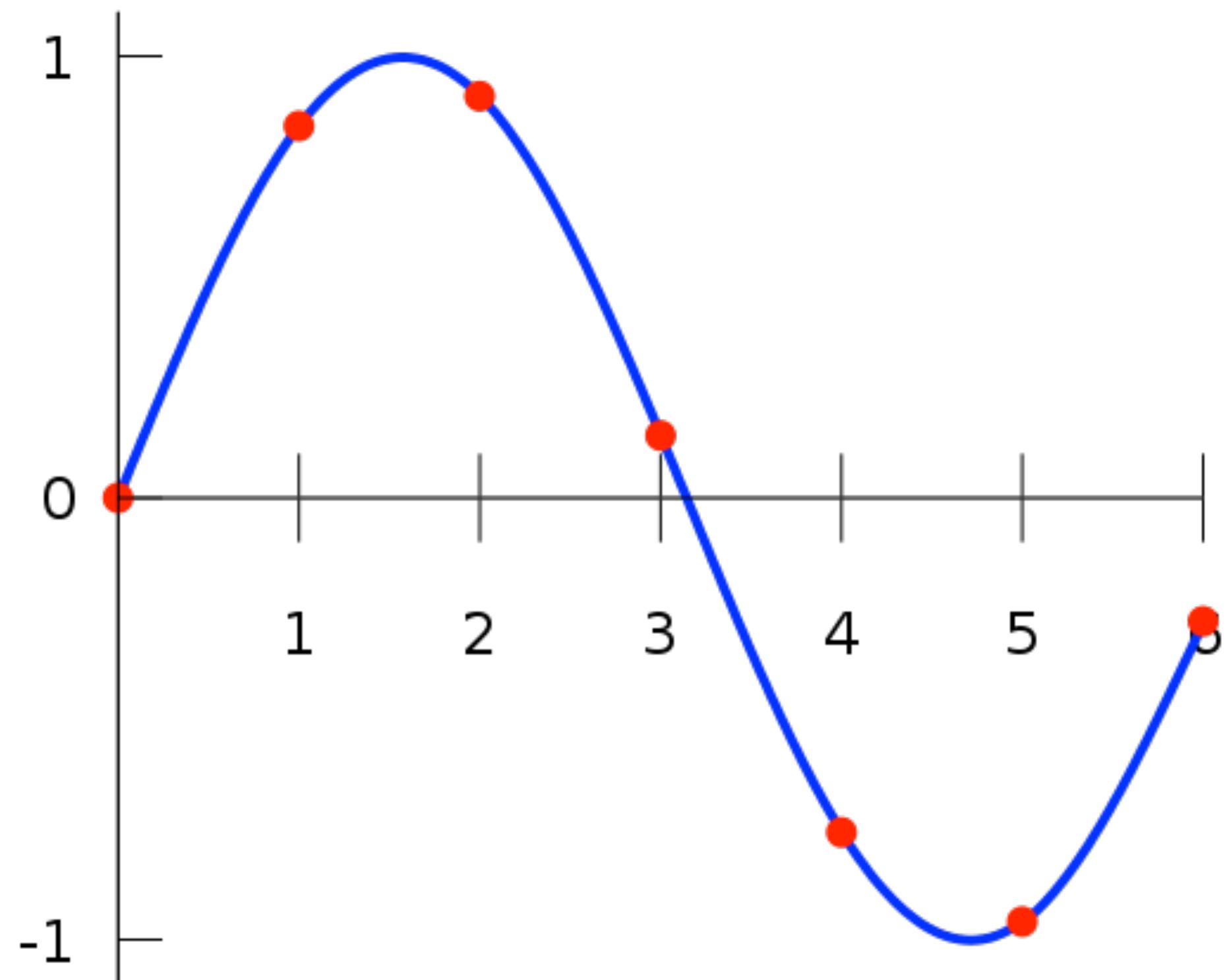


$$f(x) = \begin{cases} -0.1522x^3 + 0.9937x, & \text{if } x \in [0, 1], \\ -0.01258x^3 - 0.4189x^2 + 1.4126x - 0.1396, & \text{if } x \in [1, 2], \\ 0.1403x^3 - 1.3359x^2 + 3.2467x - 1.3623, & \text{if } x \in [2, 3], \\ 0.1579x^3 - 1.4945x^2 + 3.7225x - 1.8381, & \text{if } x \in [3, 4], \\ 0.05375x^3 - 0.2450x^2 - 1.2756x + 4.8259, & \text{if } x \in [4, 5], \\ -0.1871x^3 + 3.3673x^2 - 19.3370x + 34.9282, & \text{if } x \in [5, 6]. \end{cases}$$

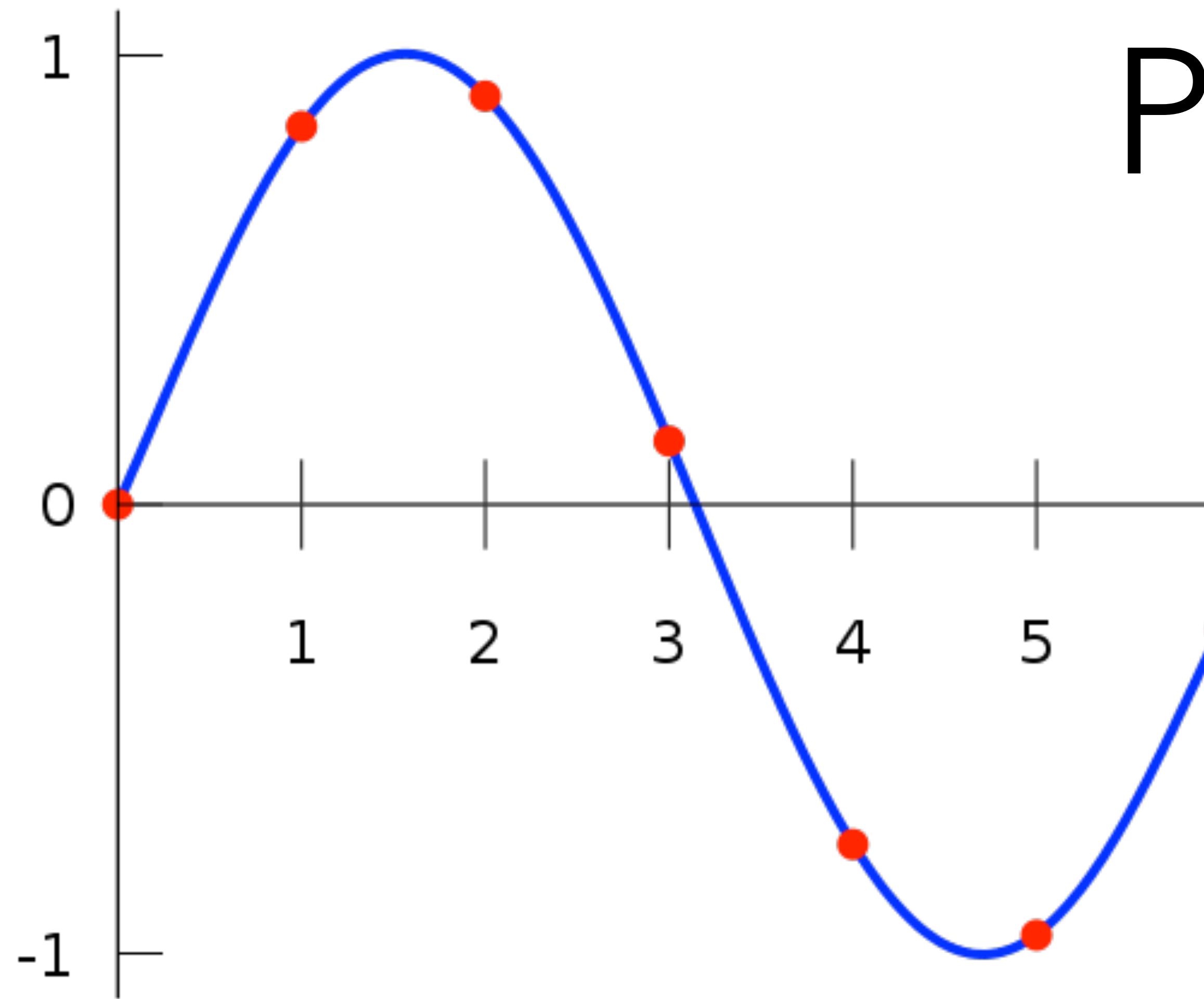
???



VS.

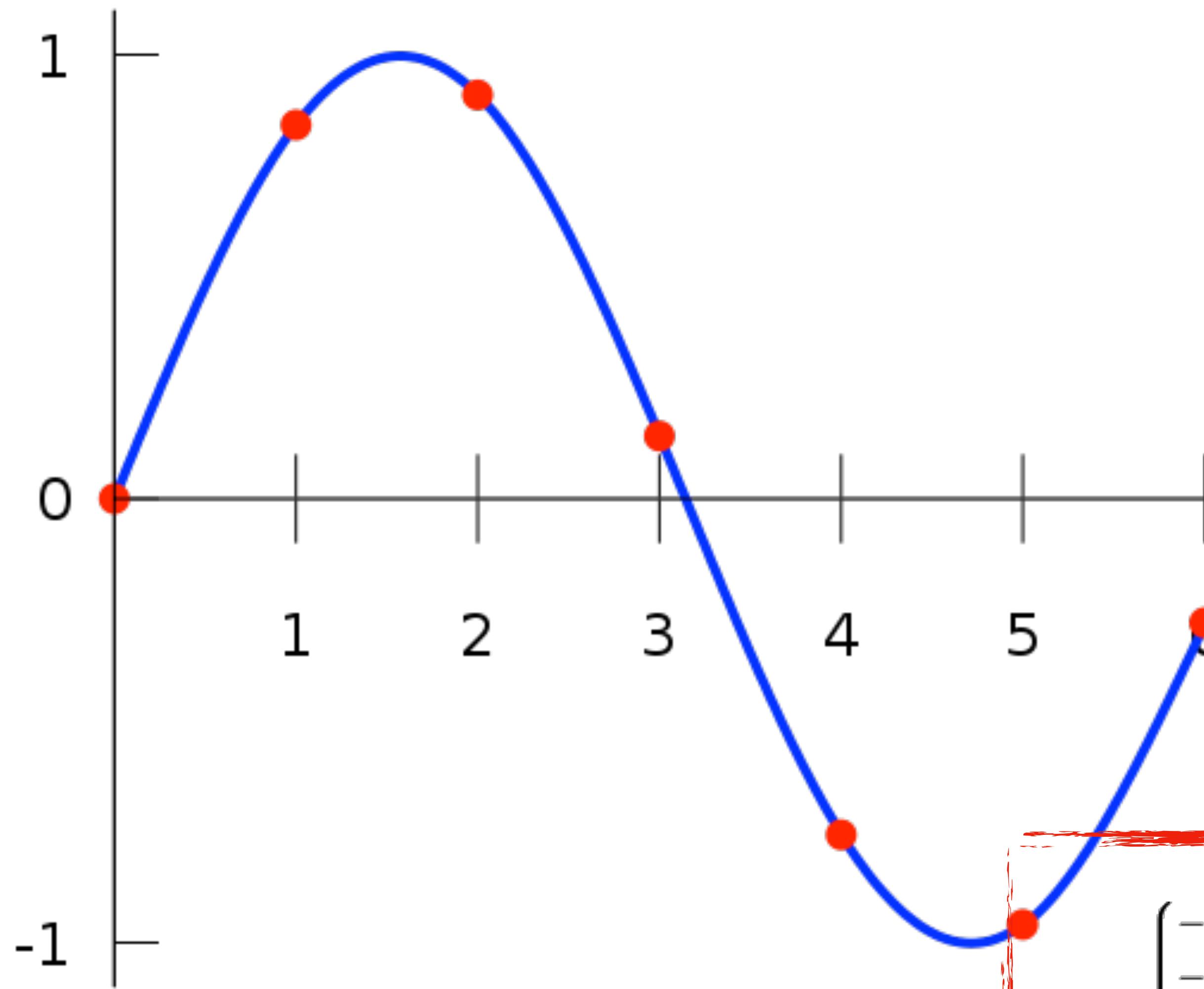


Polynomial



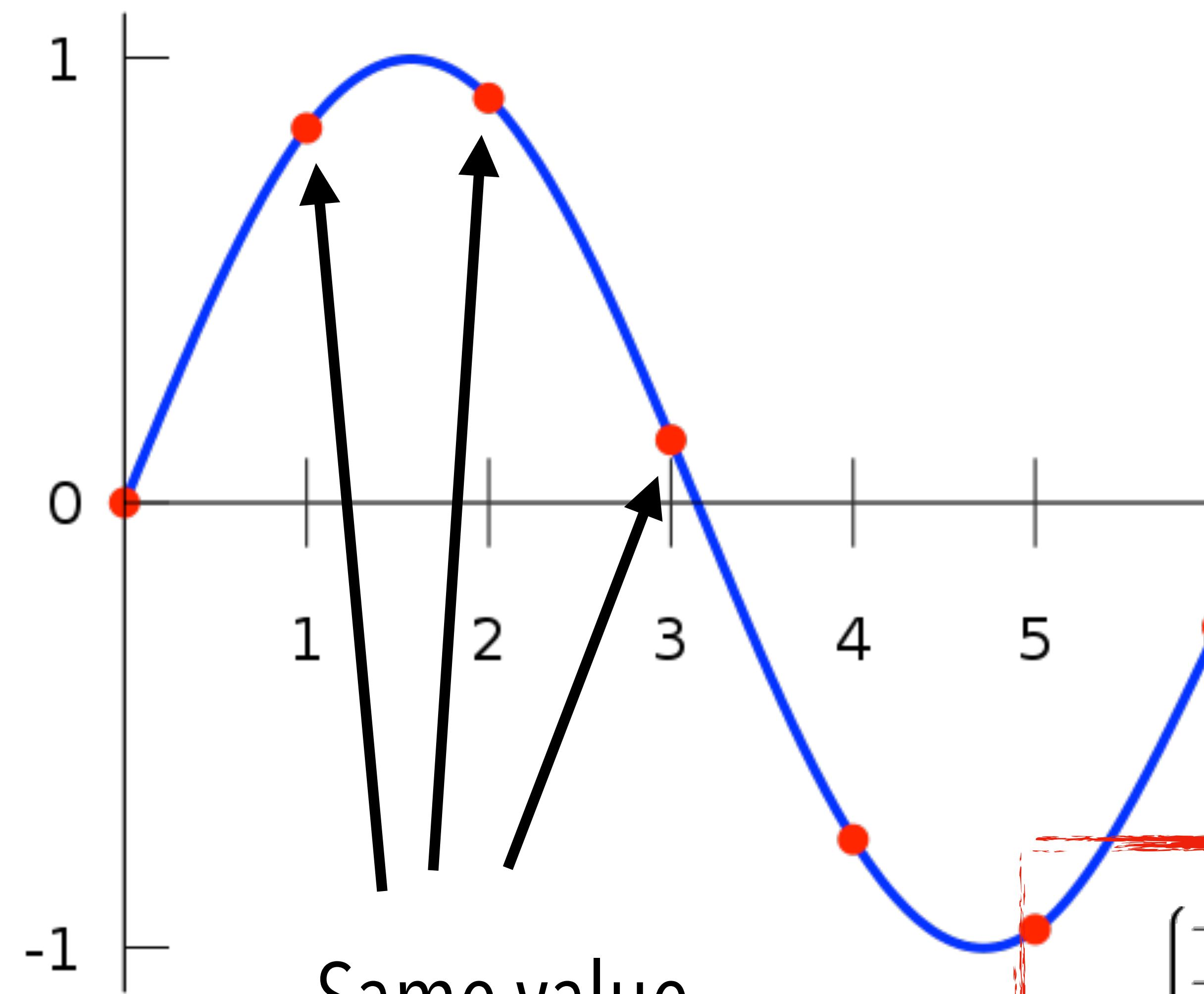
$$p(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_2 x^2 + a_1 x + a_0.$$

Splines



$$f(x) = \begin{cases} -0.1522x^3 + 0.9937x, & \text{if } x \in [0, 1], \\ -0.01258x^3 - 0.4189x^2 + 1.4126x - 0.1396, & \text{if } x \in [1, 2], \\ 0.1403x^3 - 1.3359x^2 + 3.2467x - 1.3623, & \text{if } x \in [2, 3], \\ 0.1579x^3 - 1.4945x^2 + 3.7225x - 1.8381, & \text{if } x \in [3, 4], \\ 0.05375x^3 - 0.2450x^2 - 1.2756x + 4.8259, & \text{if } x \in [4, 5], \\ -0.1871x^3 + 3.3673x^2 - 19.3370x + 34.9282, & \text{if } x \in [5, 6]. \end{cases}$$

Splines



Same value
Same f' and f''

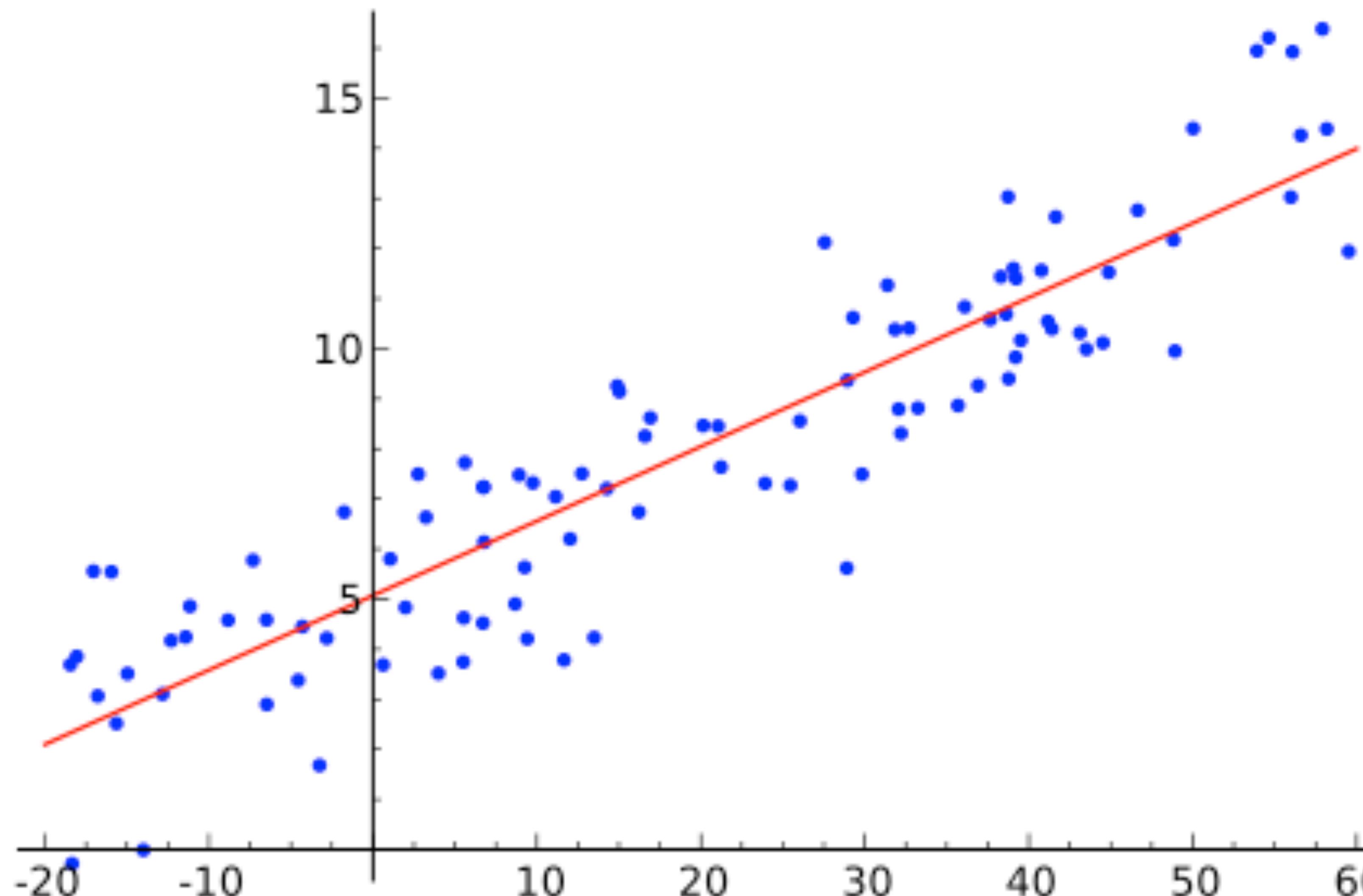
$$f(x) = \begin{cases} -0.1522x^3 + 0.9937x, & \text{if } x \in [0, 1], \\ -0.01258x^3 - 0.4189x^2 + 1.4126x - 0.1396, & \text{if } x \in [1, 2], \\ 0.1403x^3 - 1.3359x^2 + 3.2467x - 1.3623, & \text{if } x \in [2, 3], \\ 0.1579x^3 - 1.4945x^2 + 3.7225x - 1.8381, & \text{if } x \in [3, 4], \\ 0.05375x^3 - 0.2450x^2 - 1.2756x + 4.8259, & \text{if } x \in [4, 5], \\ -0.1871x^3 + 3.3673x^2 - 19.3370x + 34.9282, & \text{if } x \in [5, 6]. \end{cases}$$

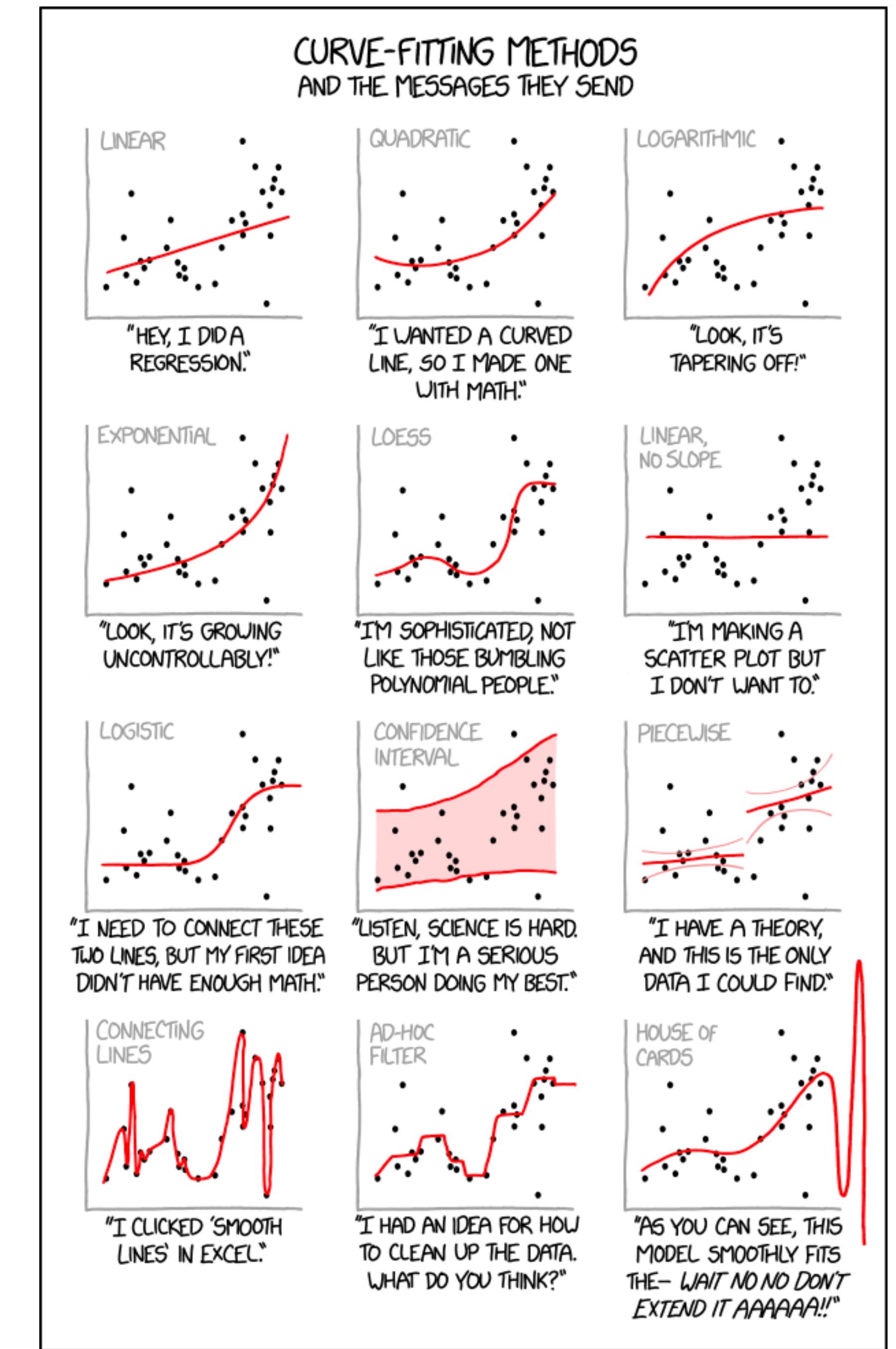
Ok. Interpolation “connects” the dots.

But what if we have some noise?

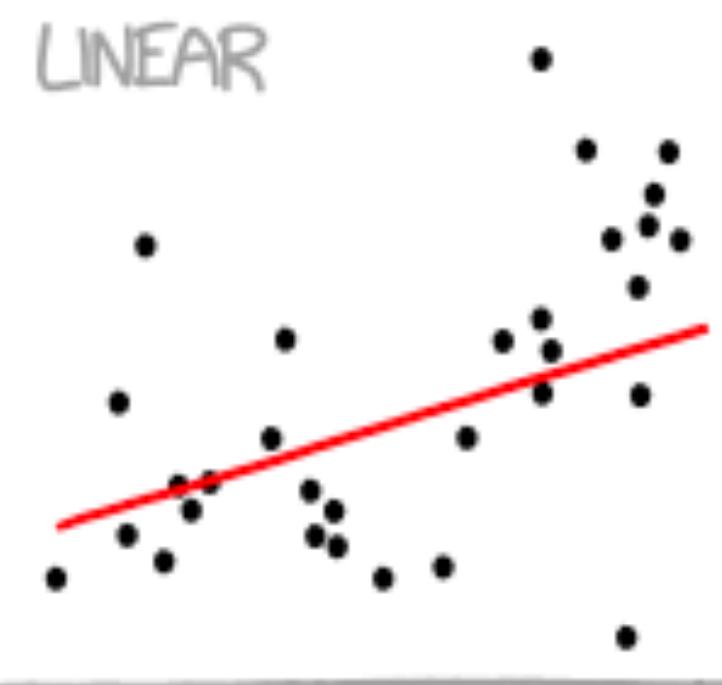
How can we visualize the **trends**?

Regression: a parametric approach

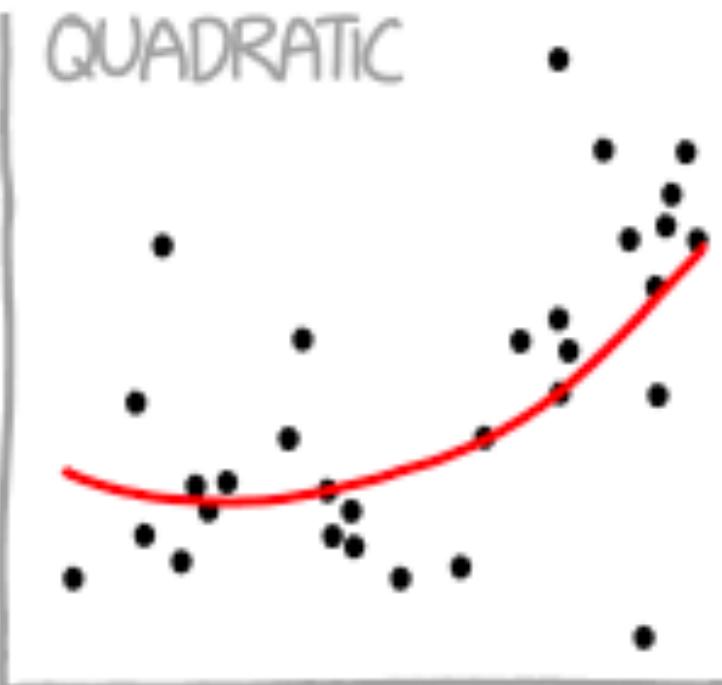




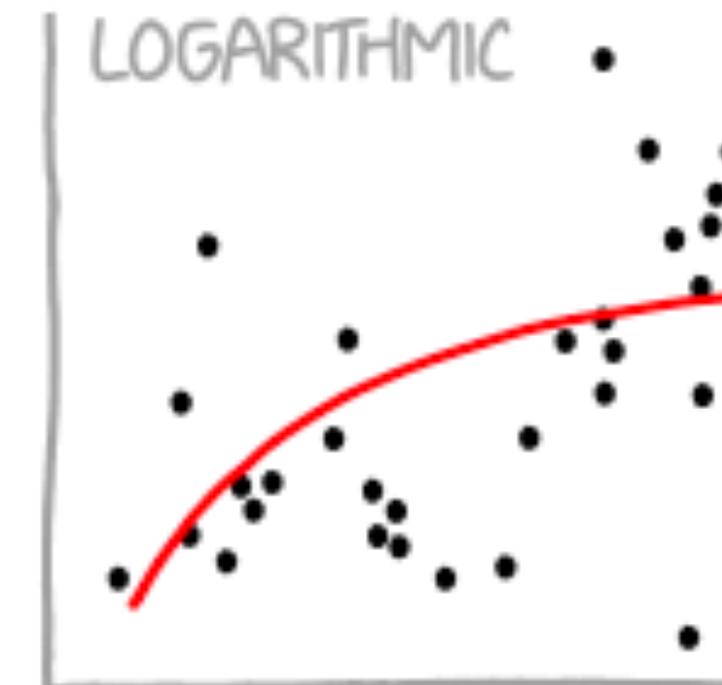
CURVE-FITTING METHODS AND THE MESSAGES THEY SEND



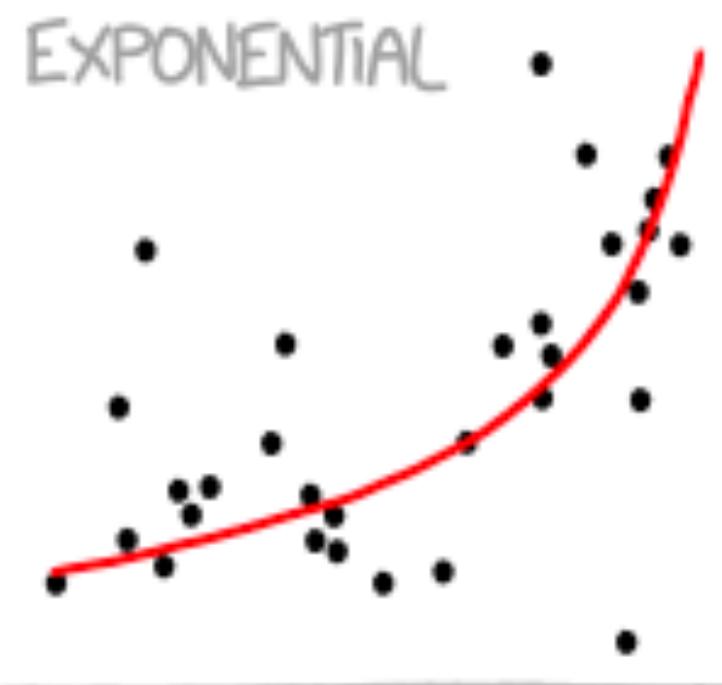
"HEY, I DID A
REGRESSION."



"I WANTED A CURVED
LINE, SO I MADE ONE
WITH MATH."



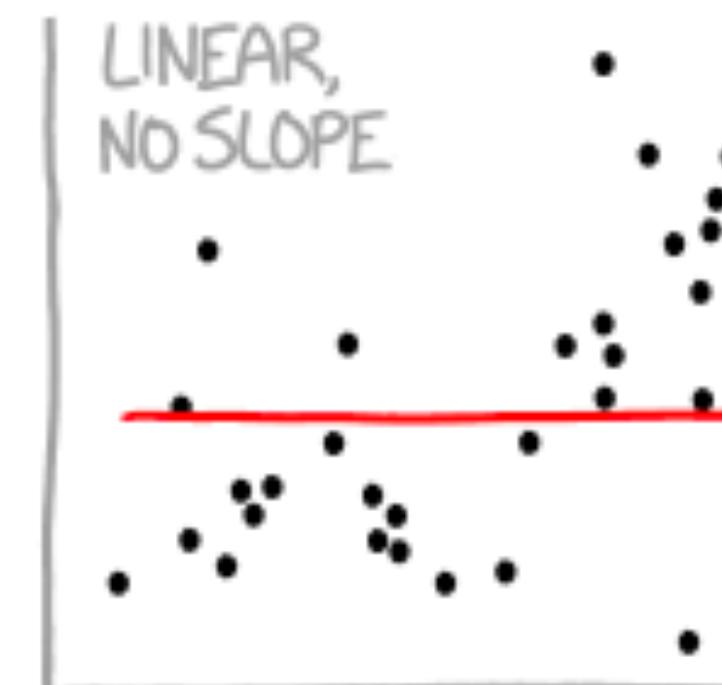
"LOOK, IT'S
TAPERING OFF!"



"LOOK, IT'S GROWING
UNCONTROLLABLY!"



"I'M SOPHISTICATED, NOT
LIKE THOSE BUMBLING
POLYNOMIAL PEOPLE."

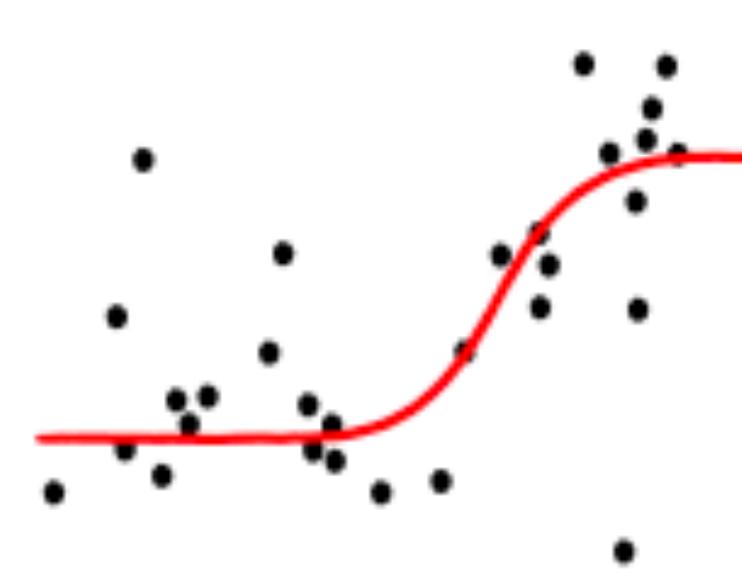


"I'M MAKING A
SCATTER PLOT BUT
I DON'T WANT TO."



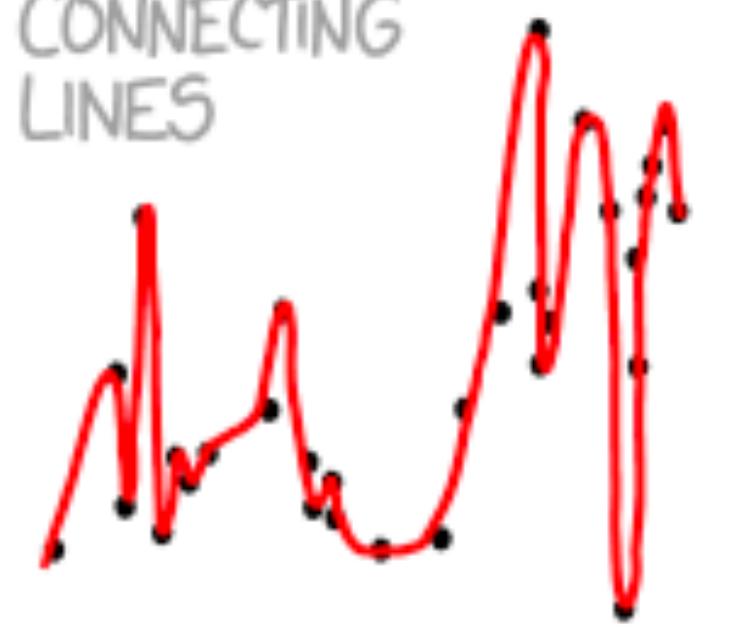
"LOOK, IT'S GROWING UNCONTROLLABLY!"

LOGISTIC



"I NEED TO CONNECT THESE TWO LINES, BUT MY FIRST IDEA DIDN'T HAVE ENOUGH MATH."

CONNECTING LINES



"I CLICKED 'SMOOTH LINES' IN EXCEL."

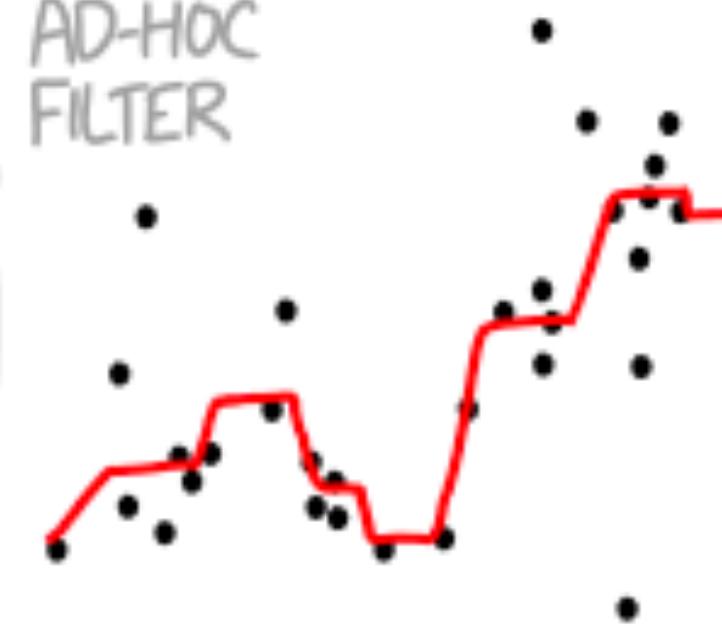
"IM SOPHISTICATED, NOT LIKE THOSE BUMBLING POLYNOMIAL PEOPLE."

CONFIDENCE INTERVAL



"LISTEN, SCIENCE IS HARD. BUT I'M A SERIOUS PERSON DOING MY BEST."

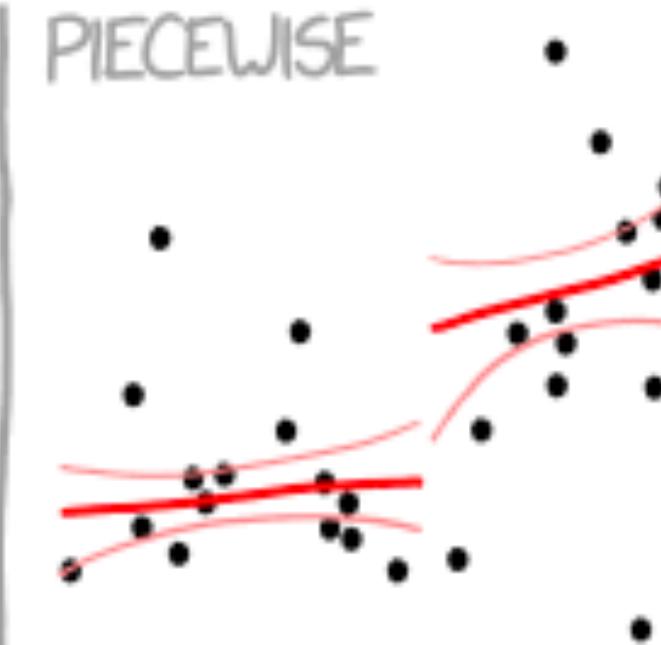
AD-HOC FILTER



"I HAD AN IDEA FOR HOW TO CLEAN UP THE DATA. WHAT DO YOU THINK?"

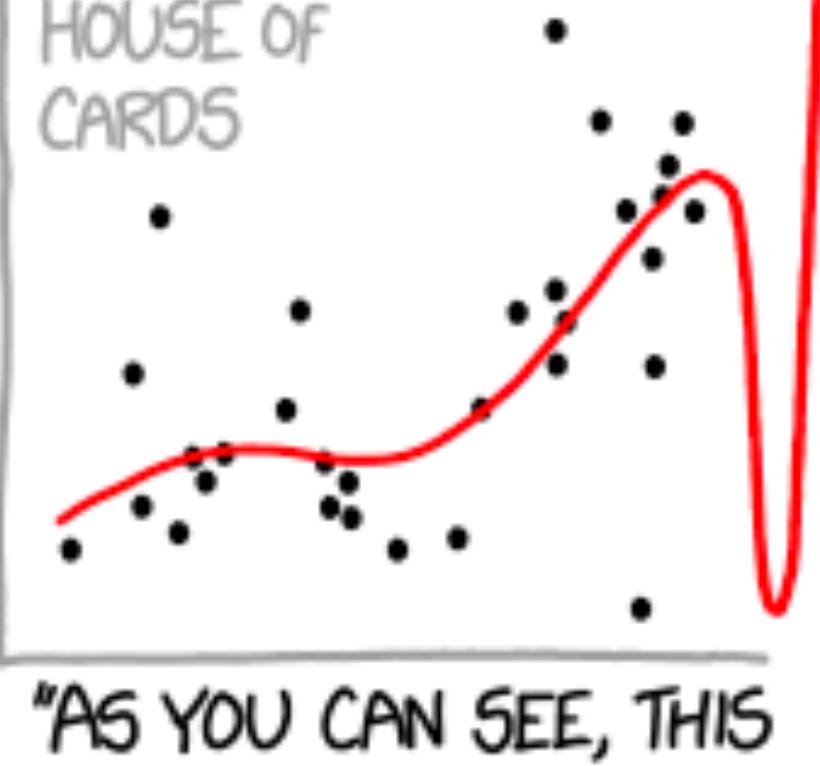
"IM MAKING A SCATTER PLOT BUT I DON'T WANT TO."

PIECEWISE



"I HAVE A THEORY, AND THIS IS THE ONLY DATA I COULD FIND."

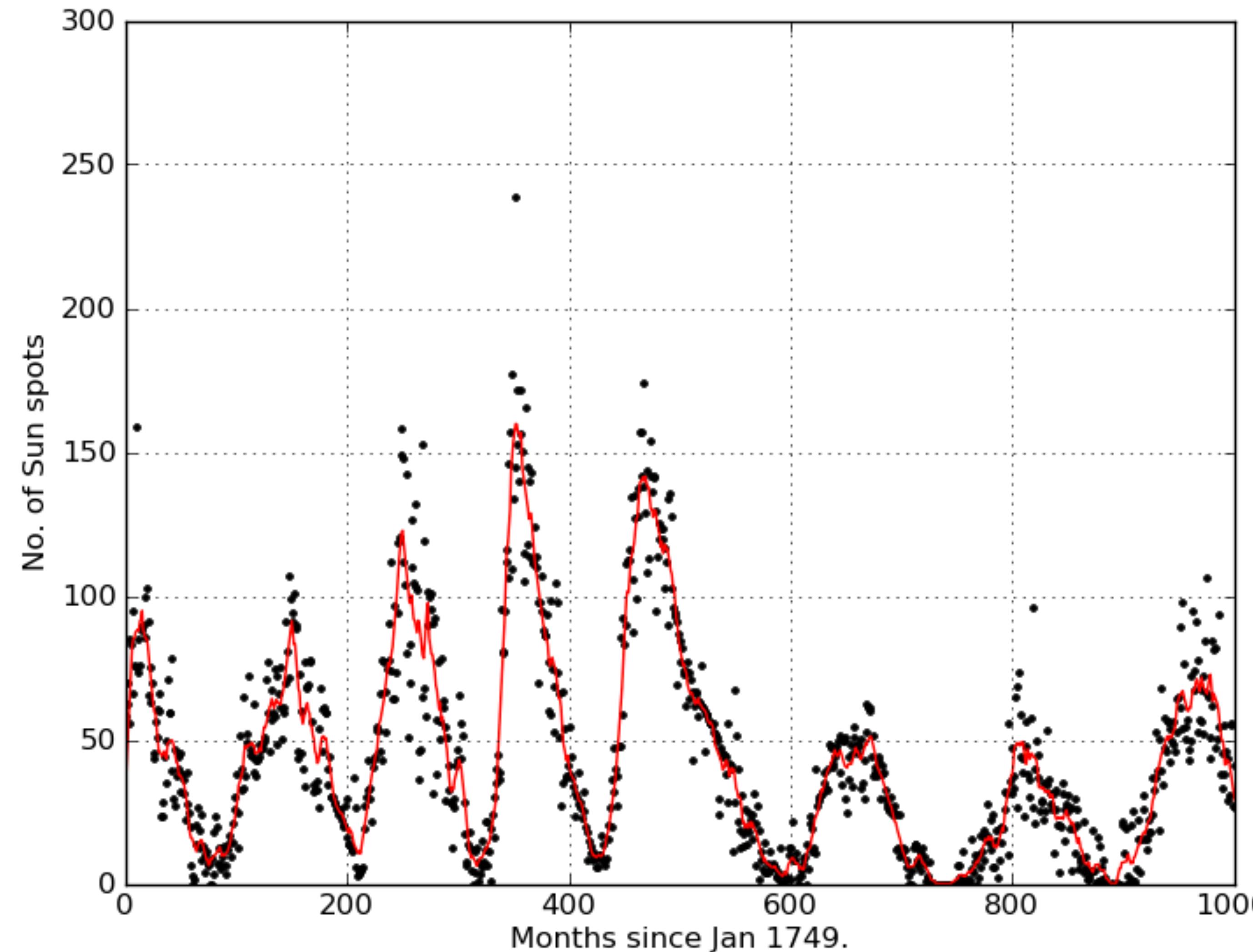
HOUSE OF CARDS

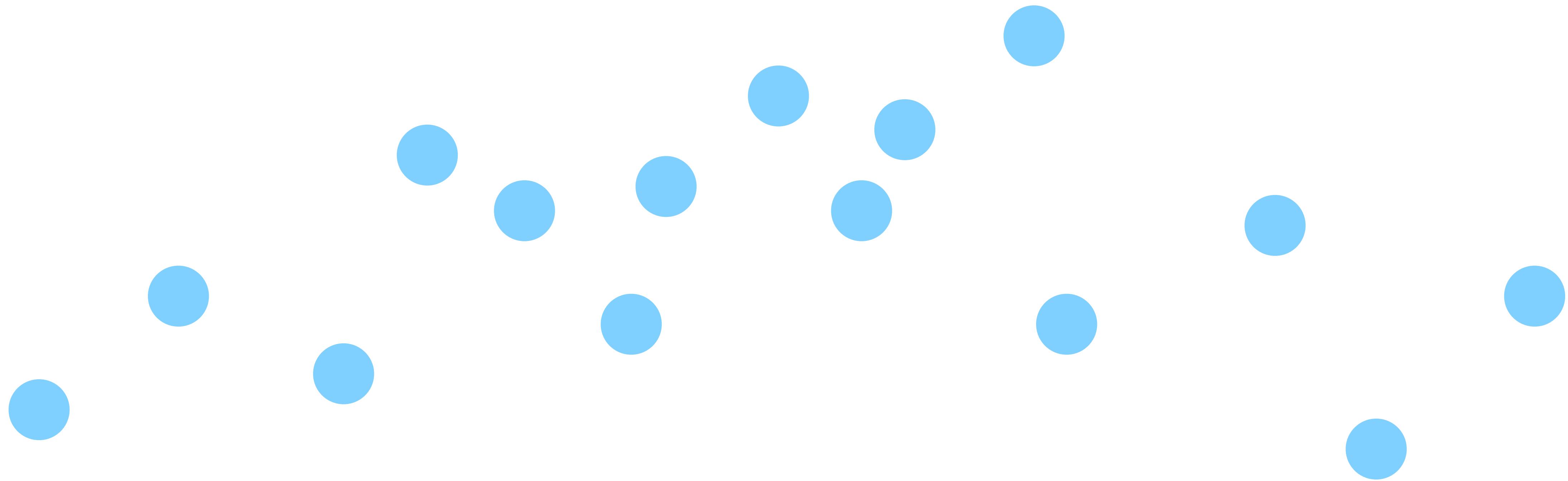


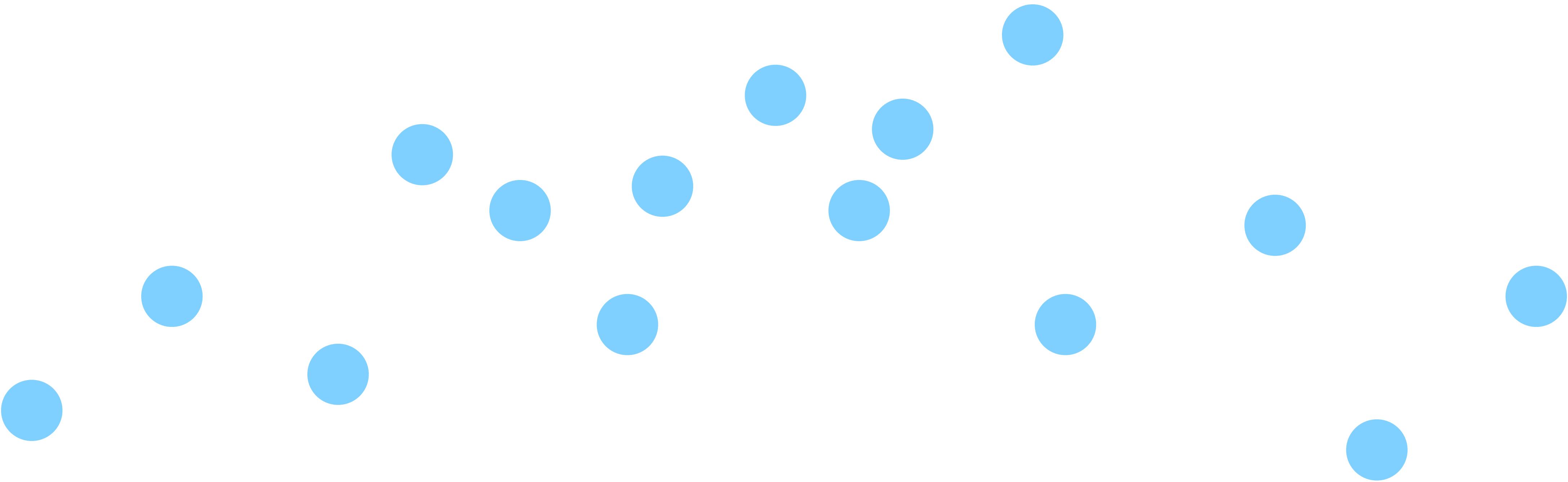
"AS YOU CAN SEE, THIS MODEL SMOOTHLY FITS THE- WAIT NO NO DONT EXTEND IT AAAAAAA!!"

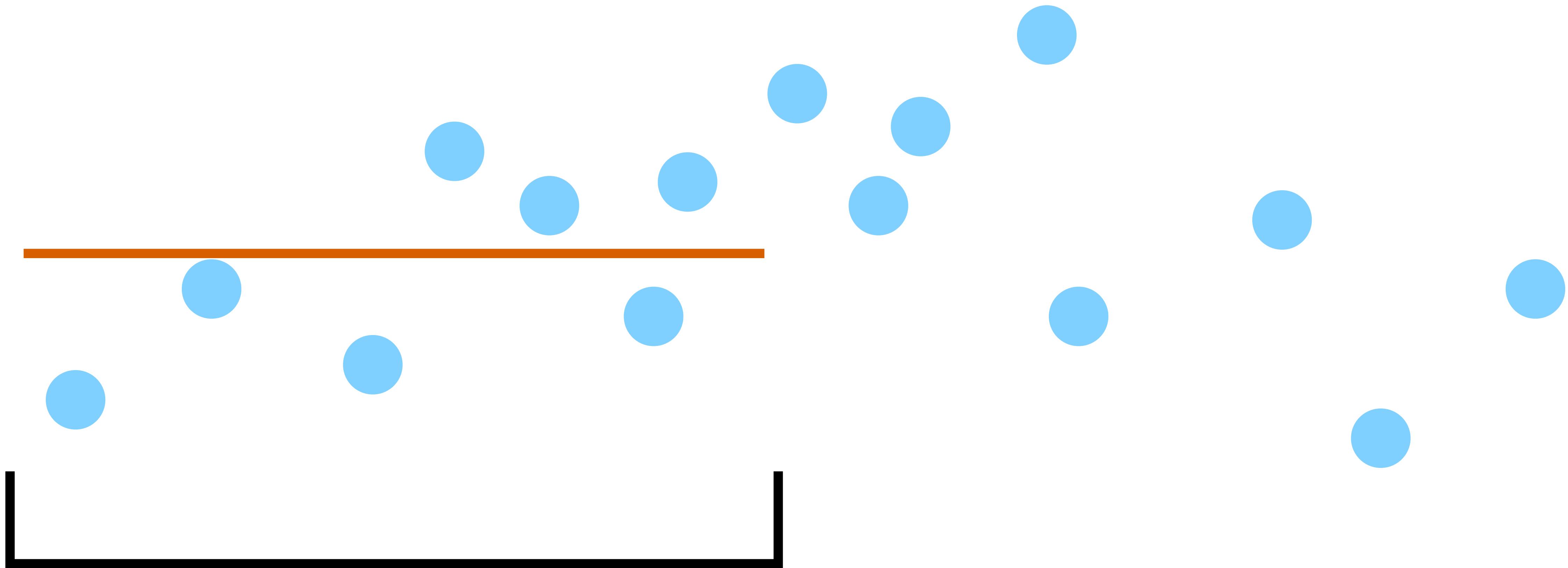
What would be non-parametric approach towards regression (trend detection)?

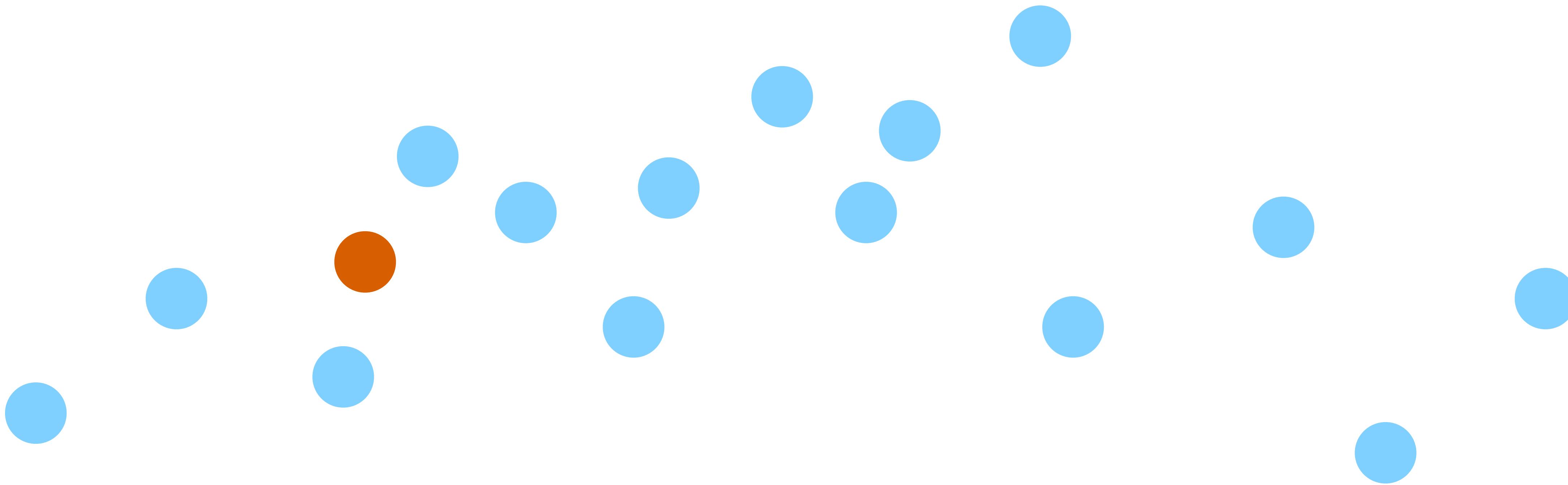
Moving average: a non-parametric approach

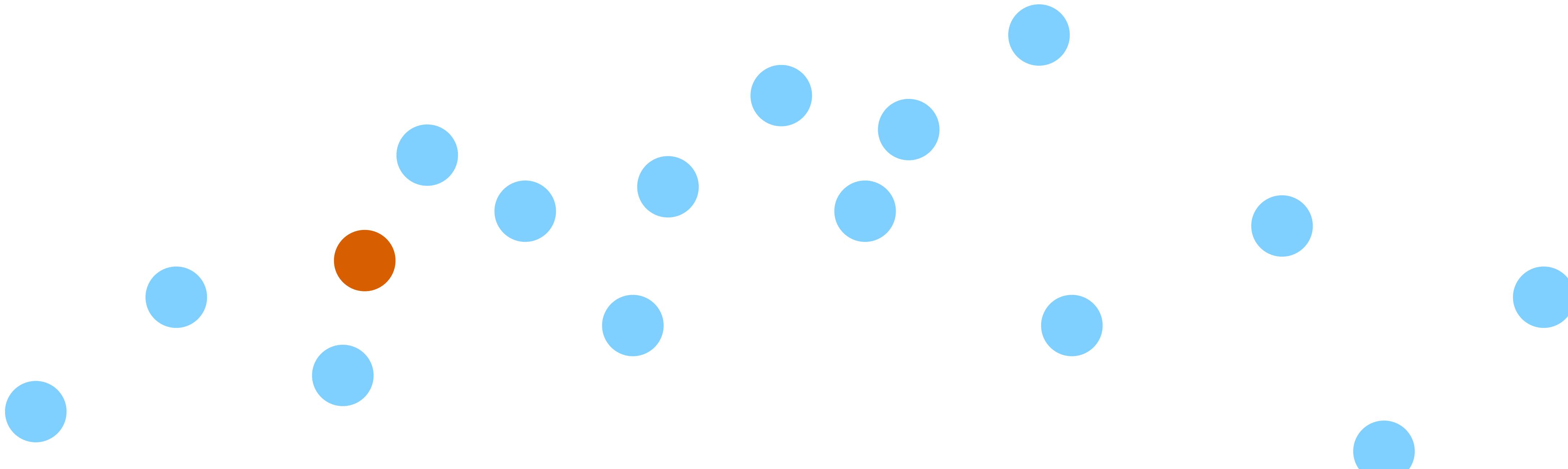
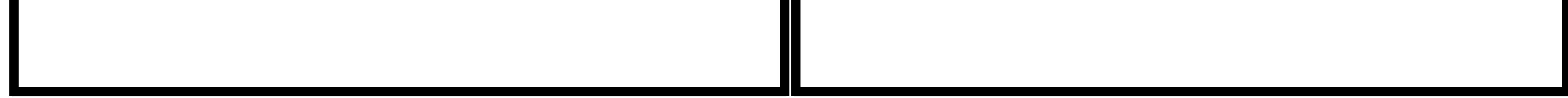


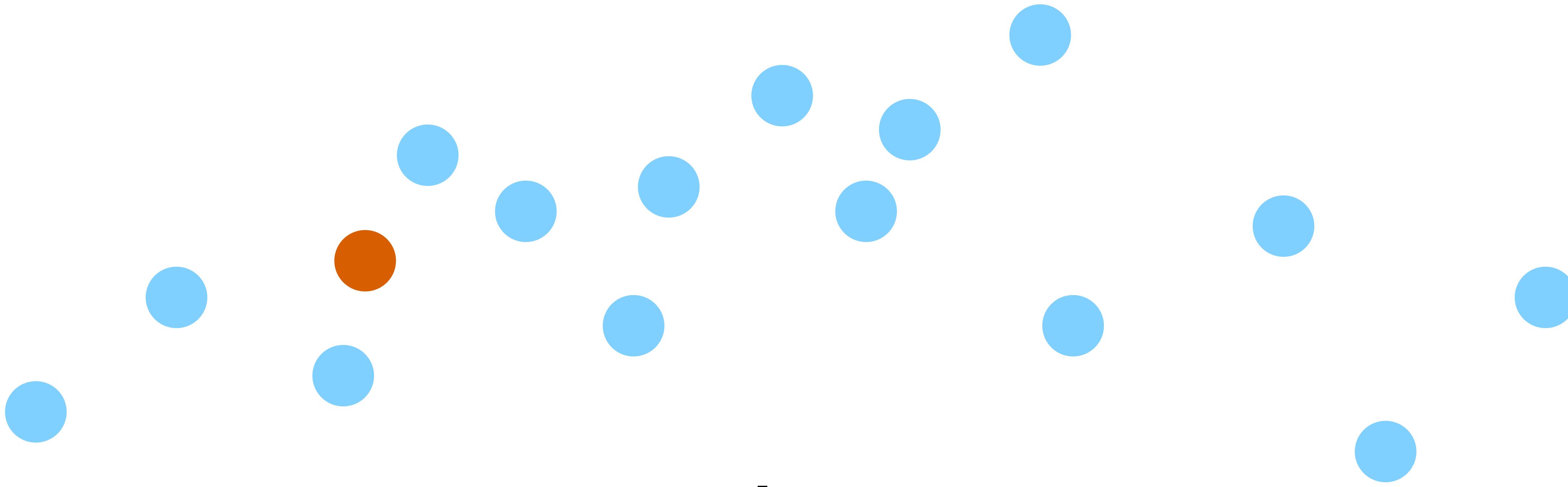


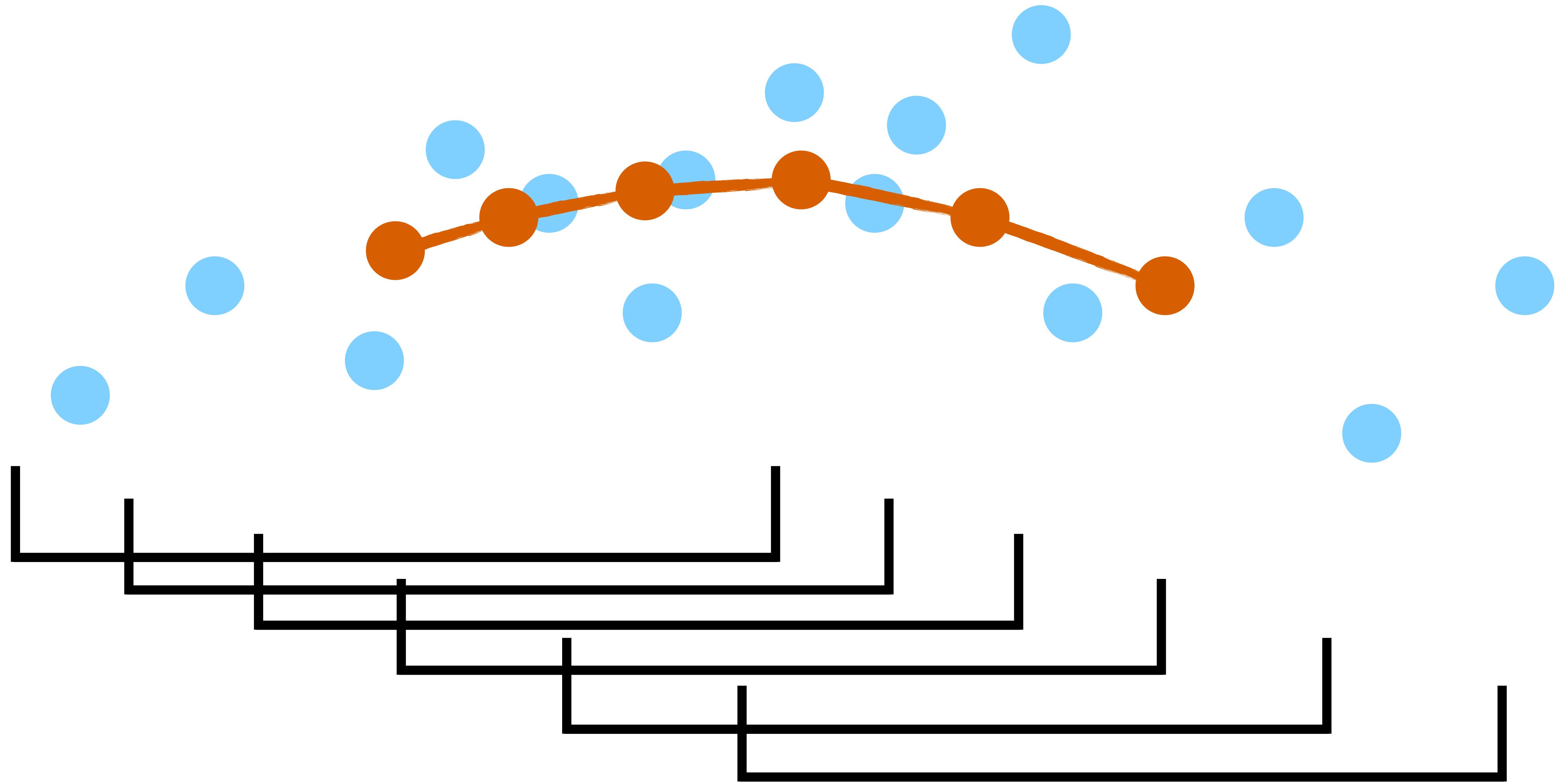










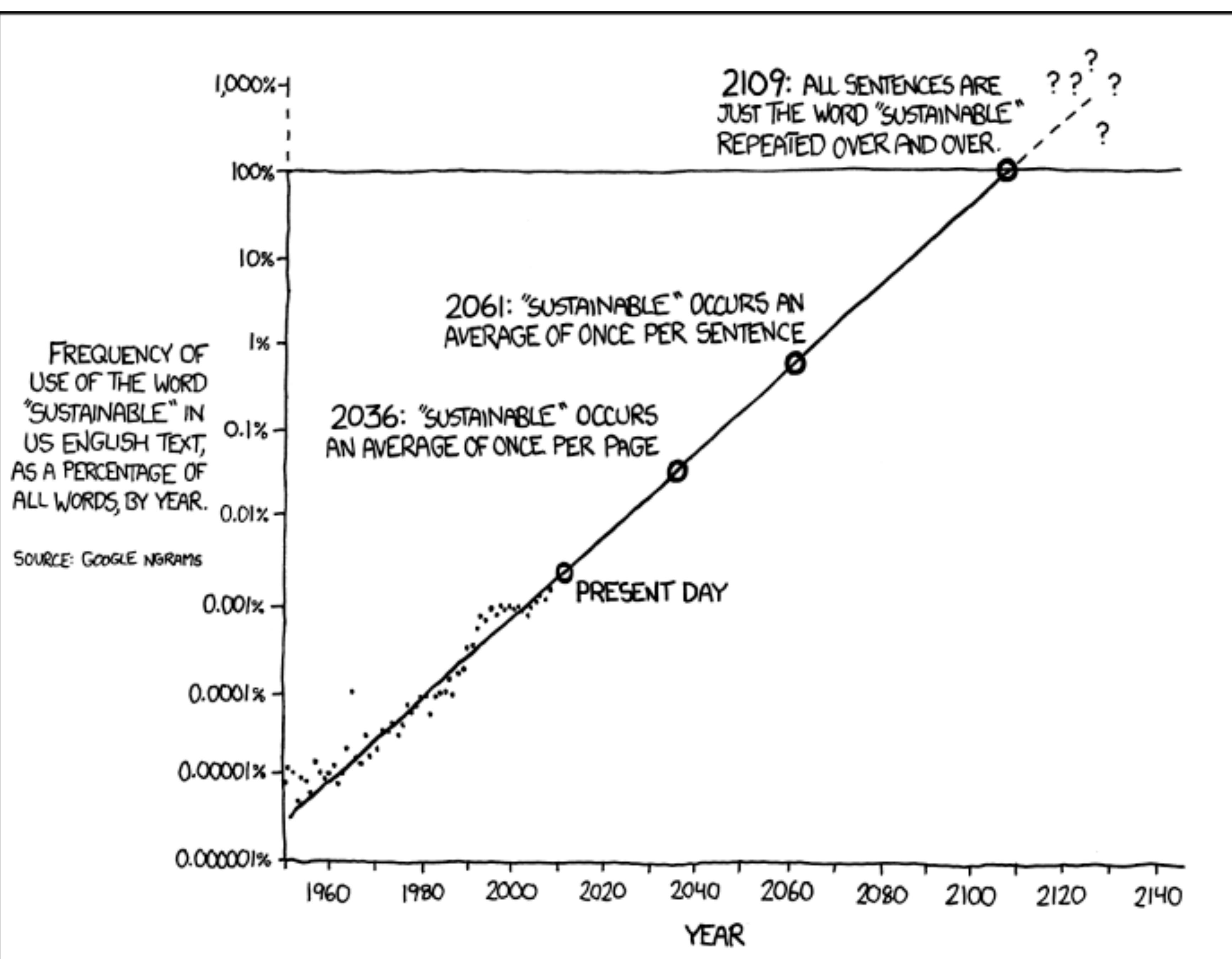


Extrapolation

Let's extend the trend we know

MY HOBBY: EXTRAPOLATING



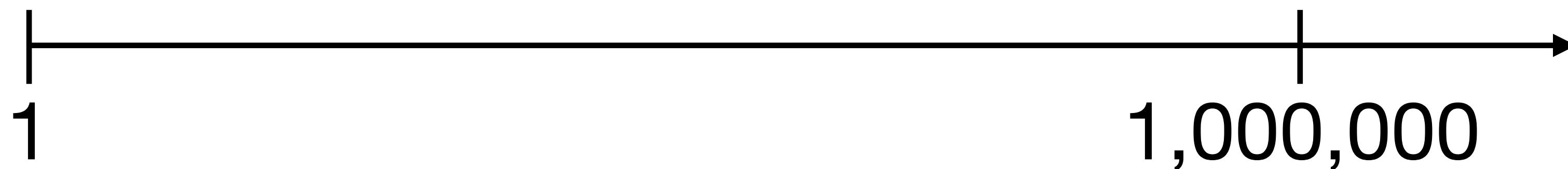


THE WORD "SUSTAINABLE" IS UNSUSTAINABLE.

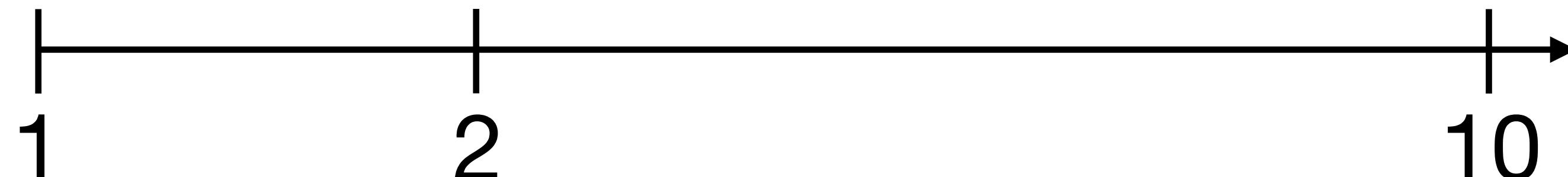
xkcd

Logarithm

Where is 1,000?



Where is 5?

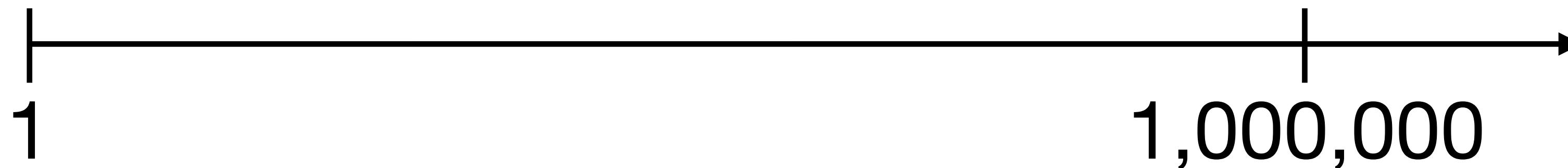


Where is 1,000?



(if linear scale)

Where is 1,000?



How about in **log scale**?

$$1,000 / 1 = 1,000$$

$$1,000 / 1 = 1,000$$

$$1,000,000 / 1,000 = 1,000$$

1,000 / 1 = 1,000

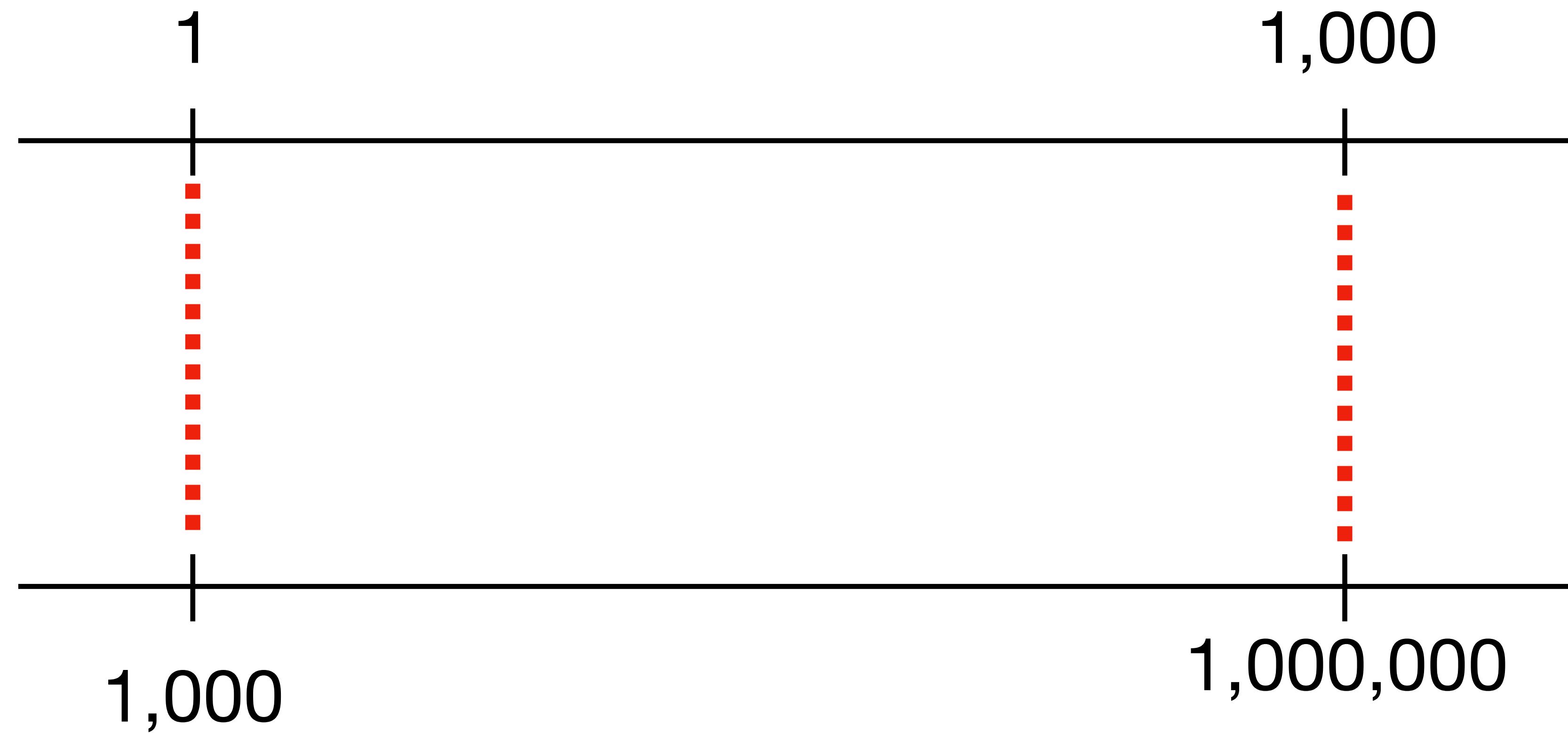
1,000 / 1 = 1,000

1,000,000 / 1,000 = 1,000

In log-scale,



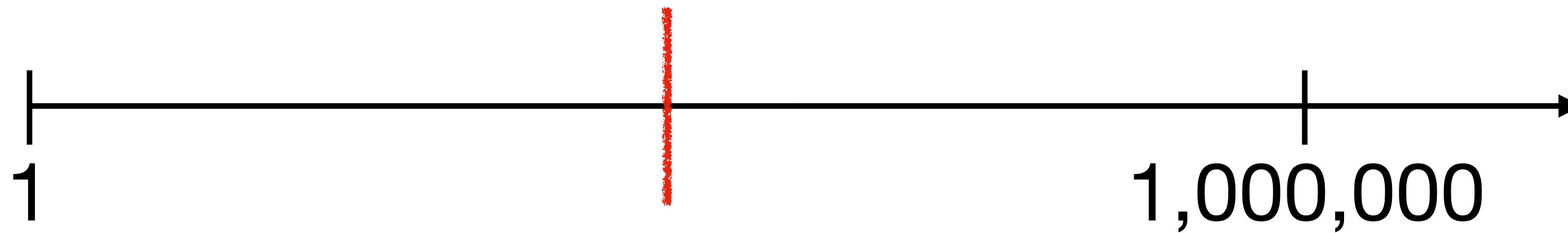
In log-scale,



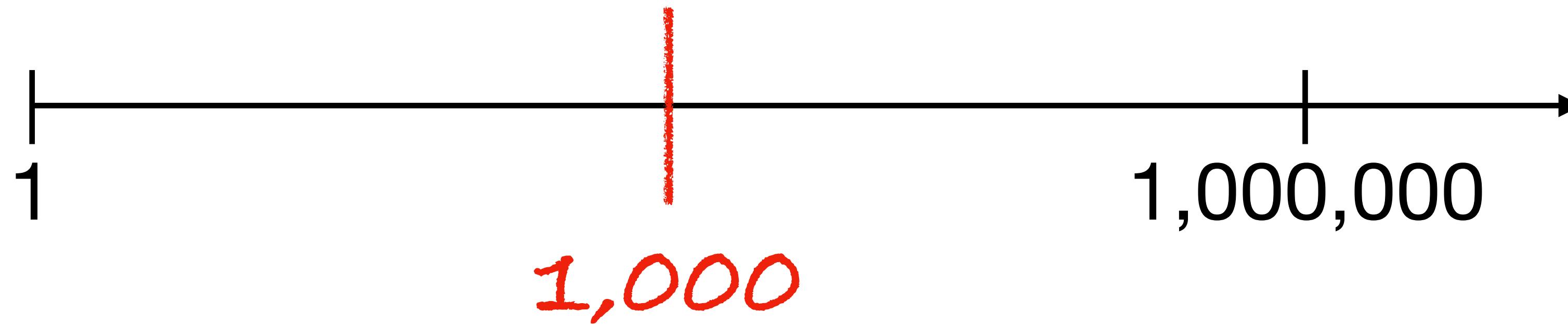
Where is 1,000?



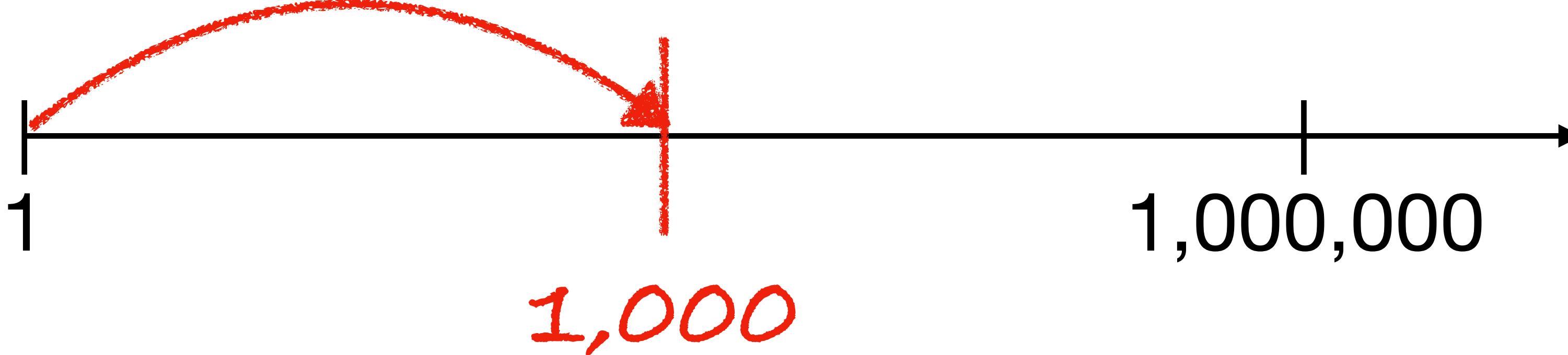
Where is 1,000?



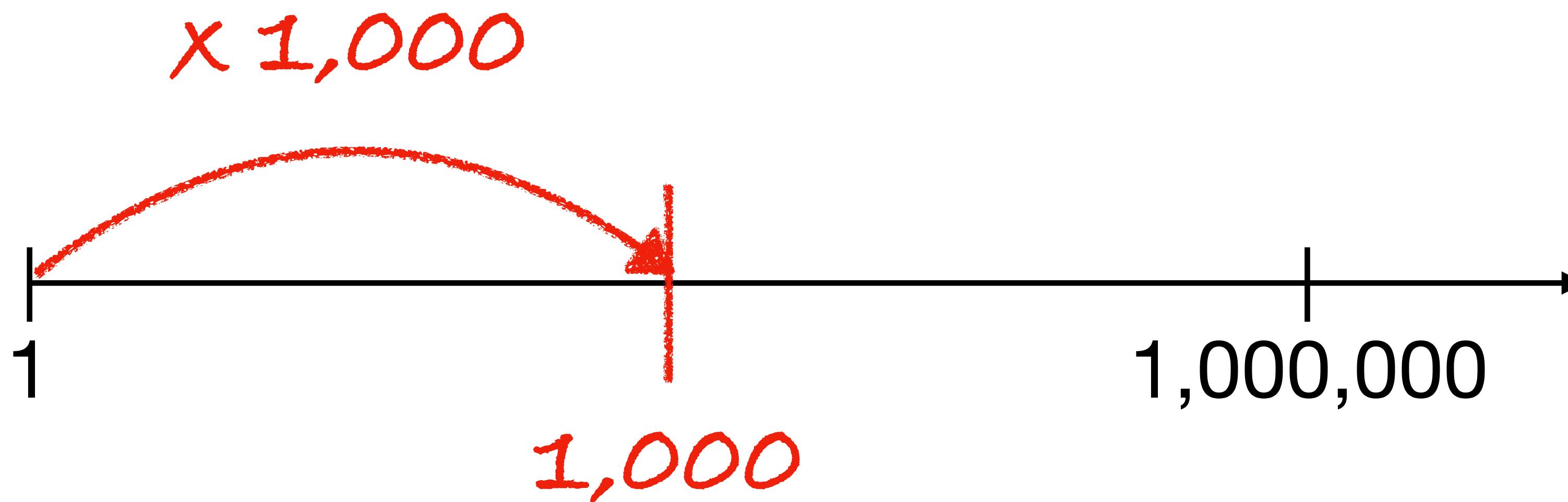
Where is 1,000?



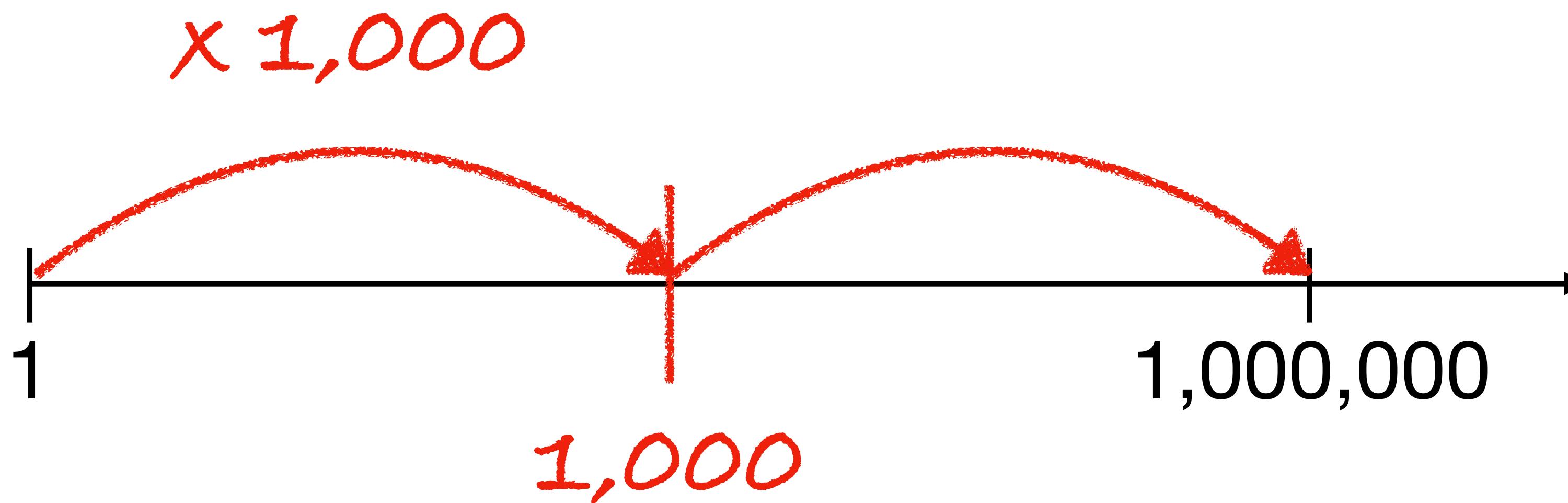
Where is 1,000?



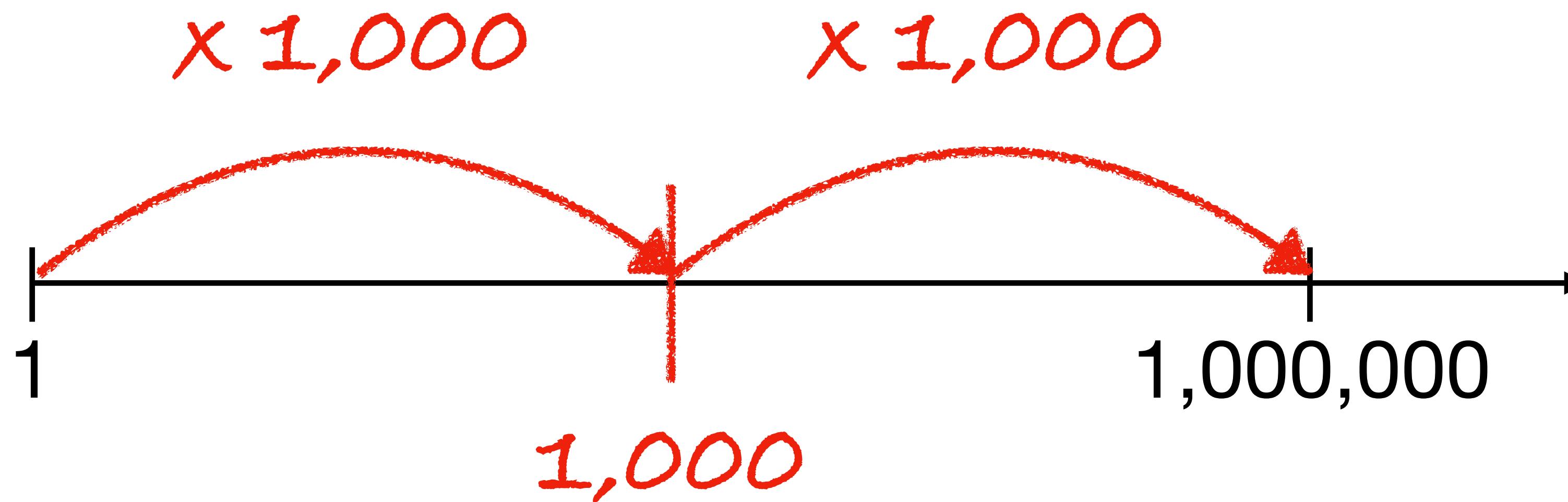
Where is 1,000?



Where is 1,000?



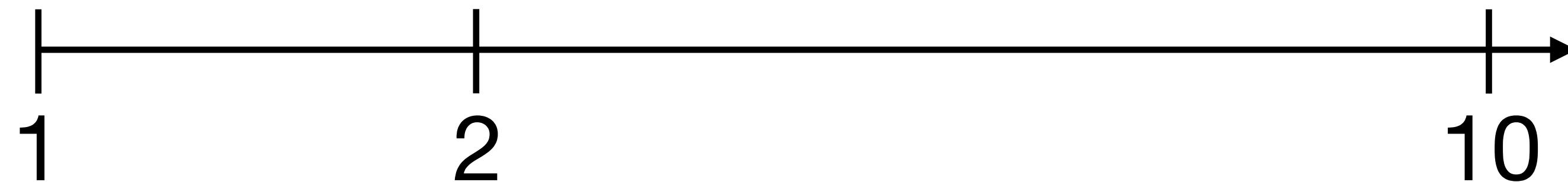
Where is 1,000?



Where is 5?

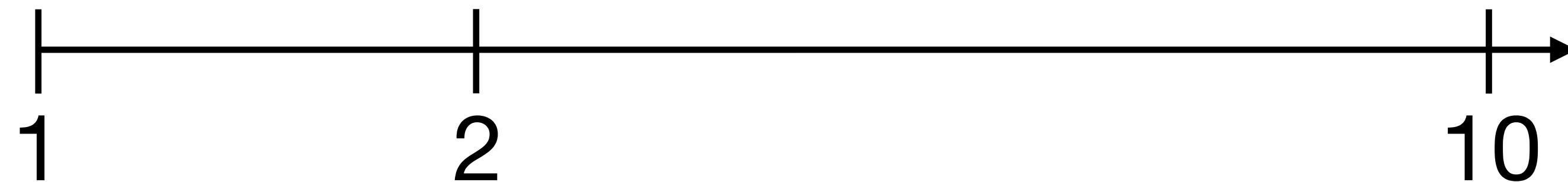


Where is 5?



Where is 5?

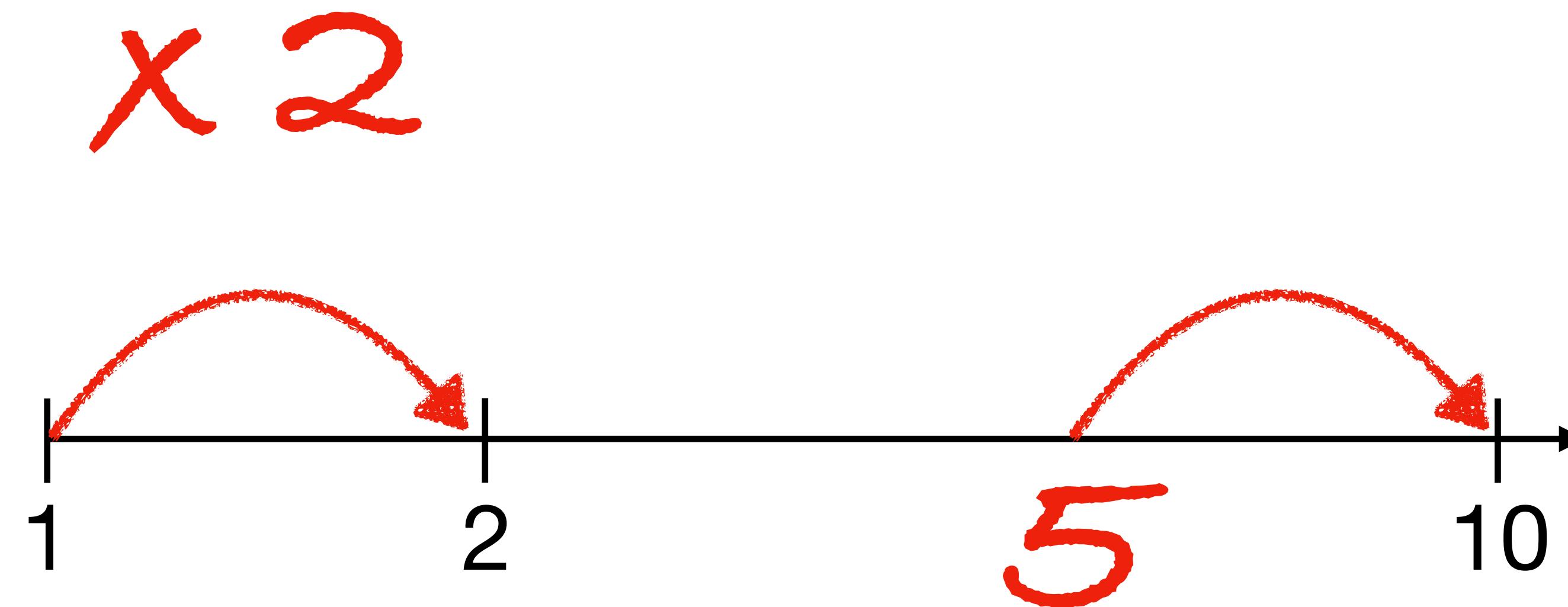
x_2



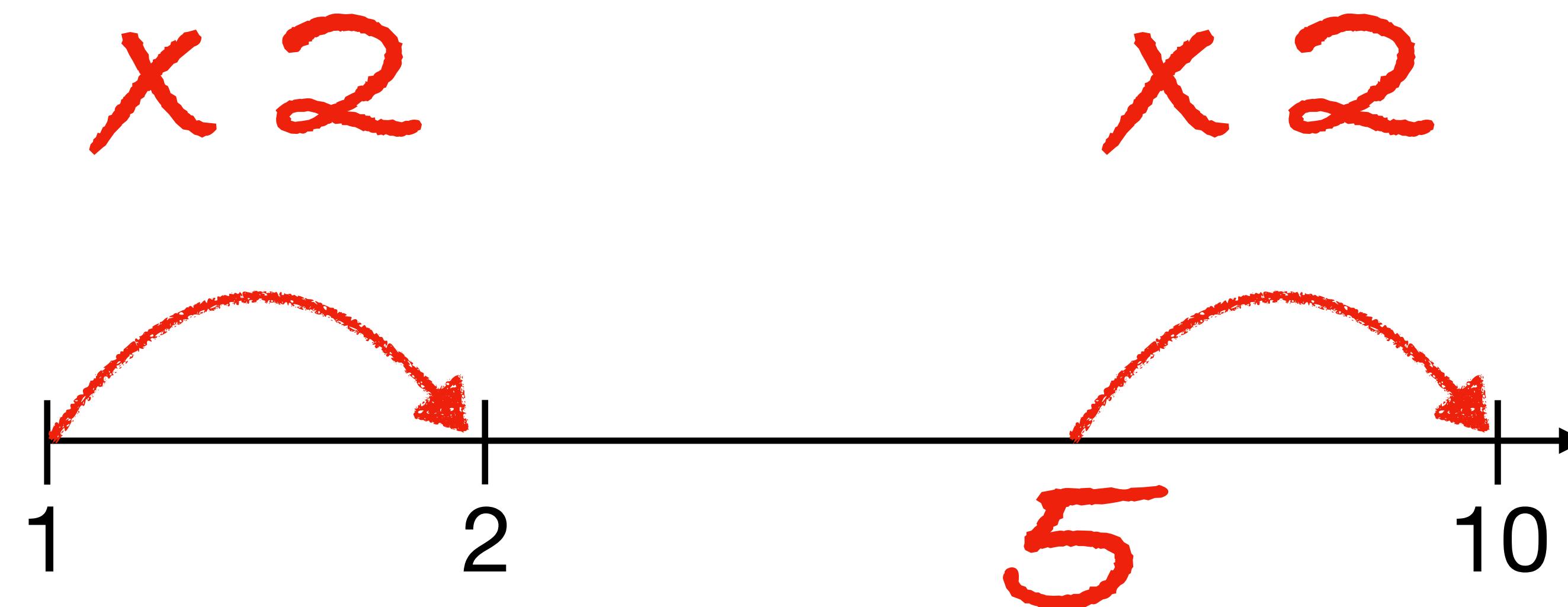
Where is 5?



Where is 5?



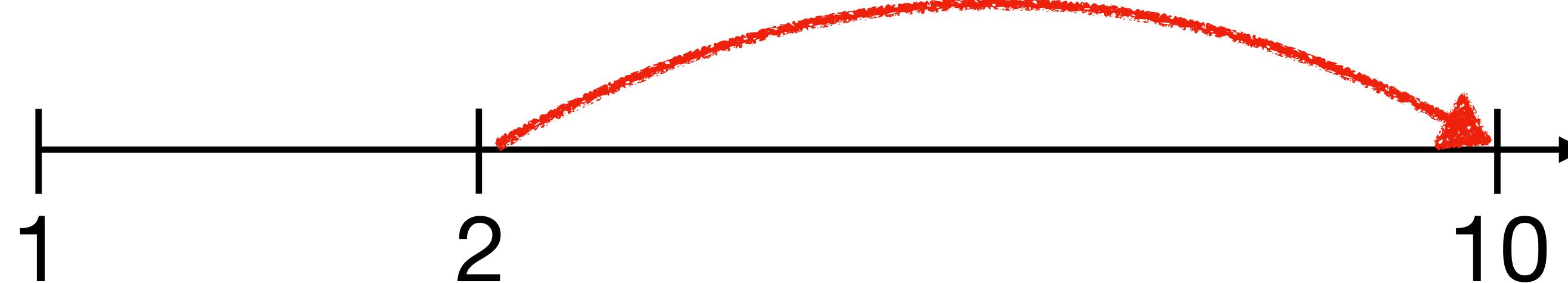
Where is 5?



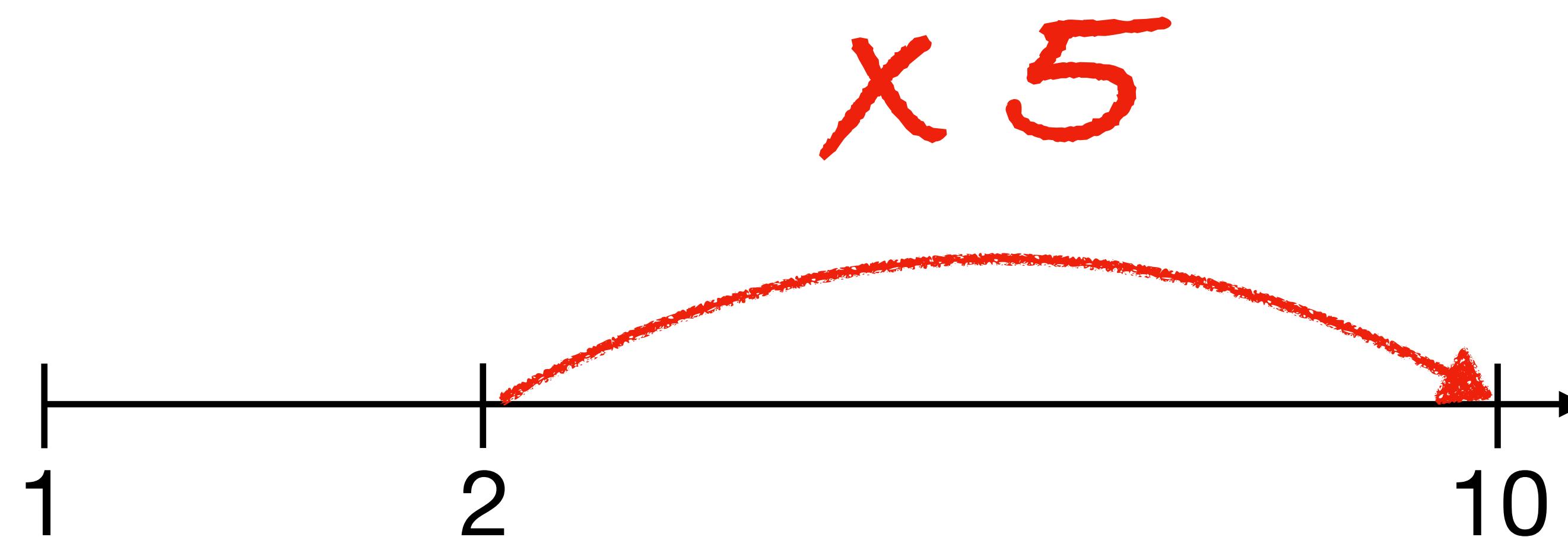
Where is 5?



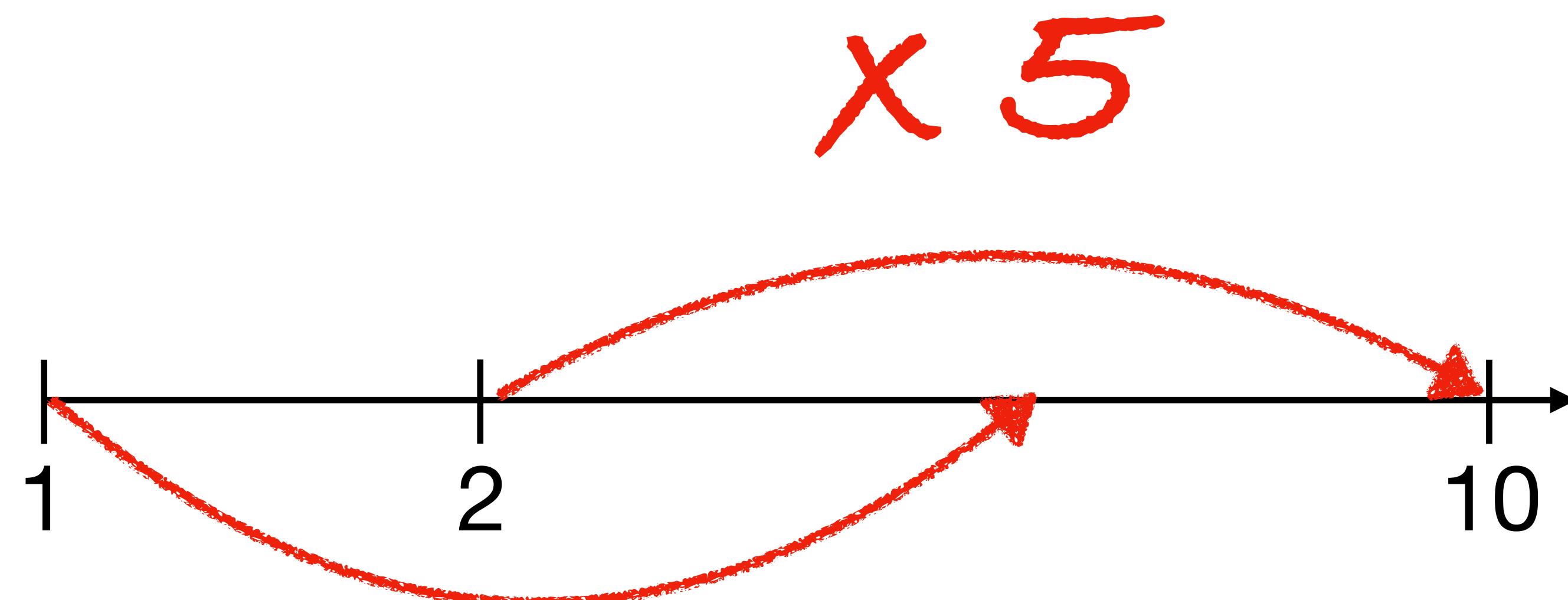
Where is 5?



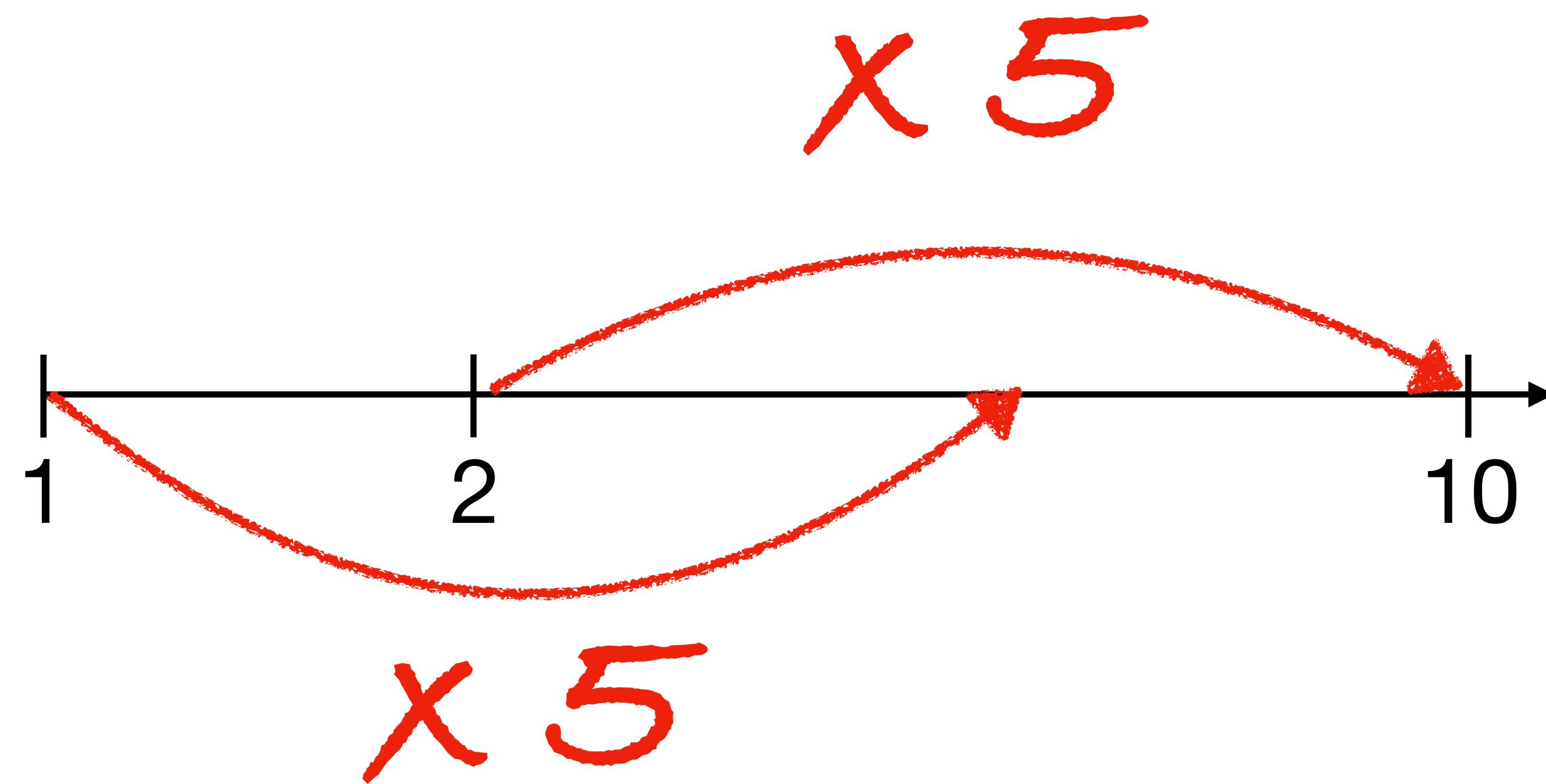
Where is 5?



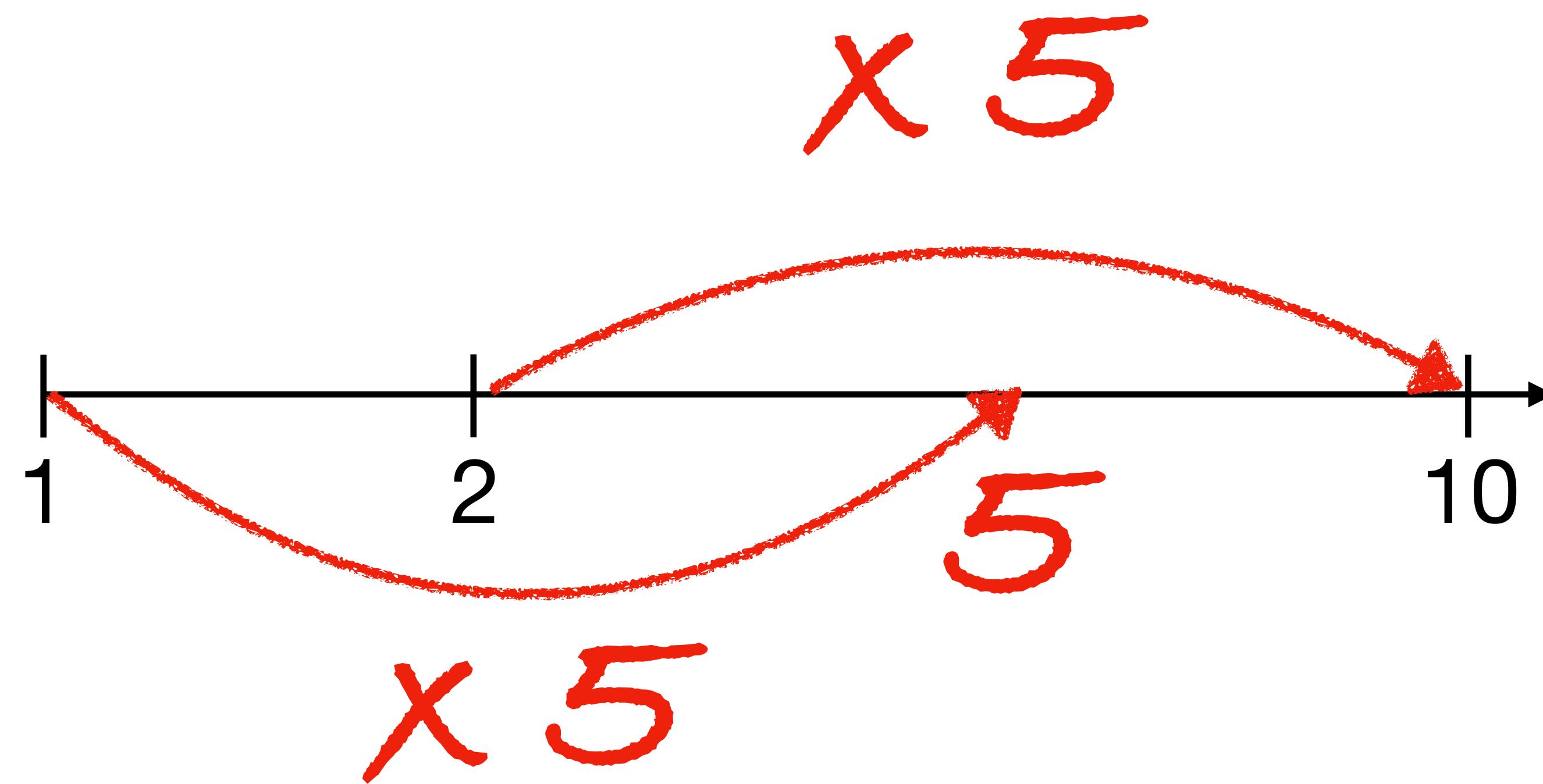
Where is 5?



Where is 5?



Where is 5?



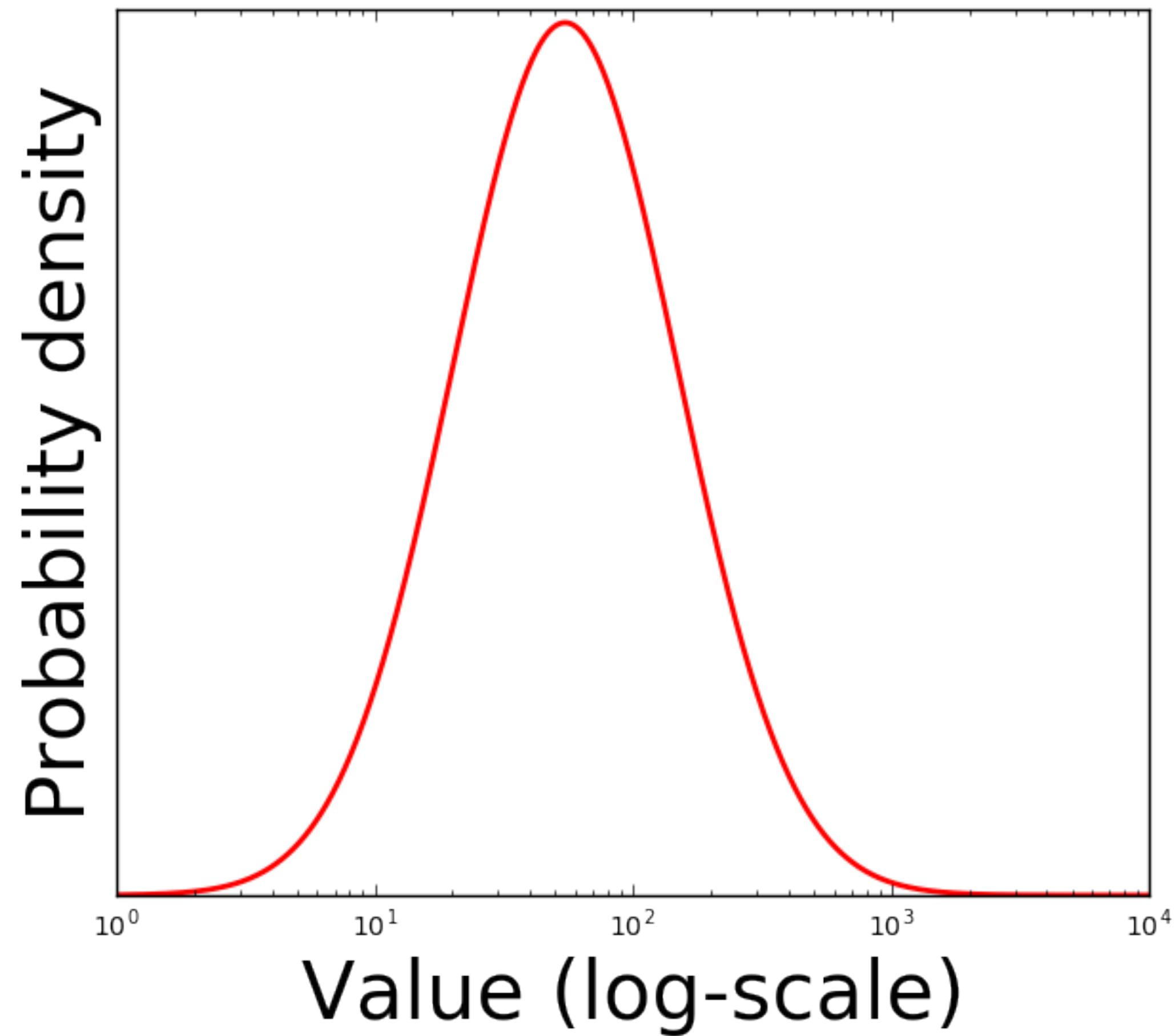
Distance \sim multiplication

Why logscale?

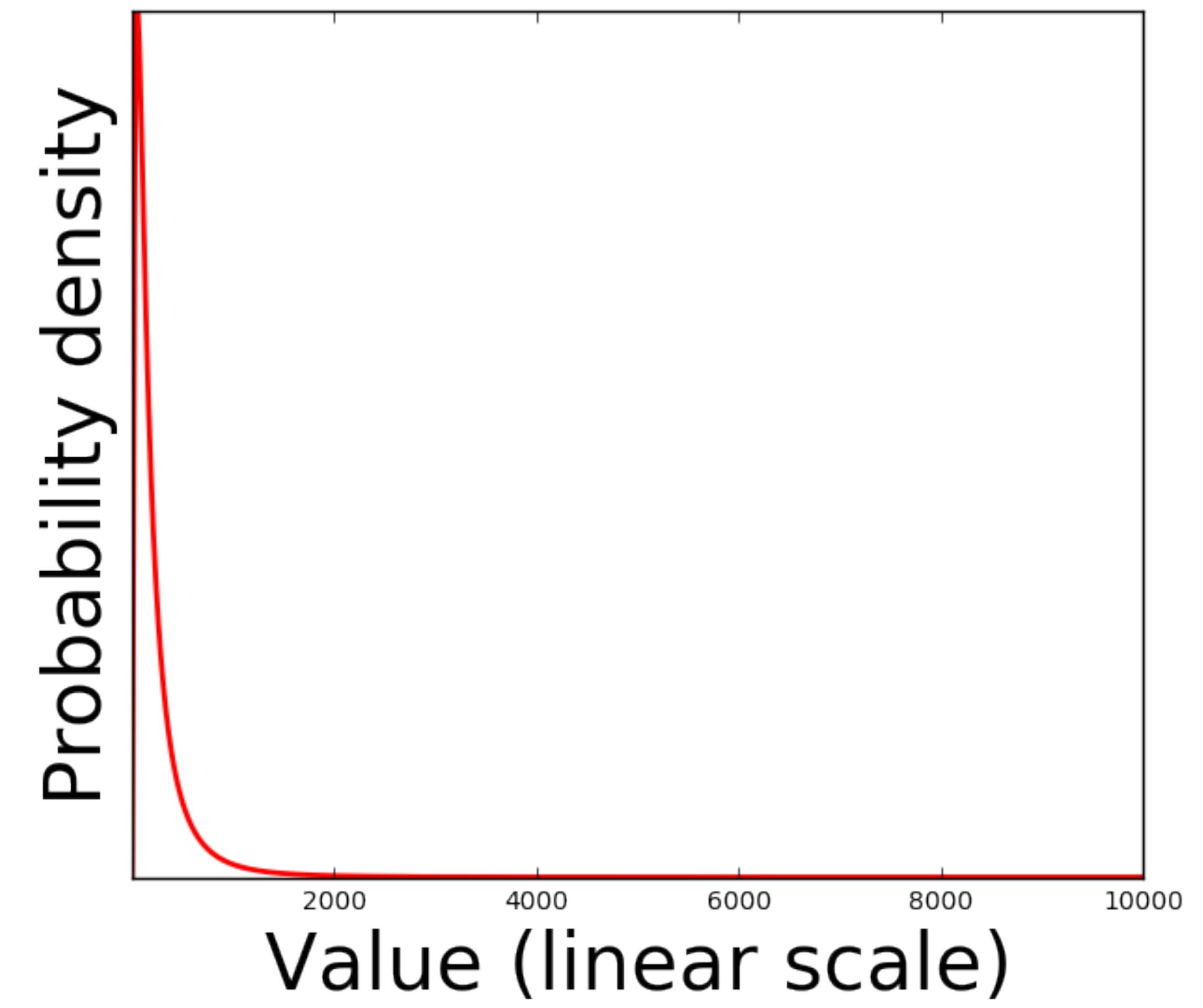
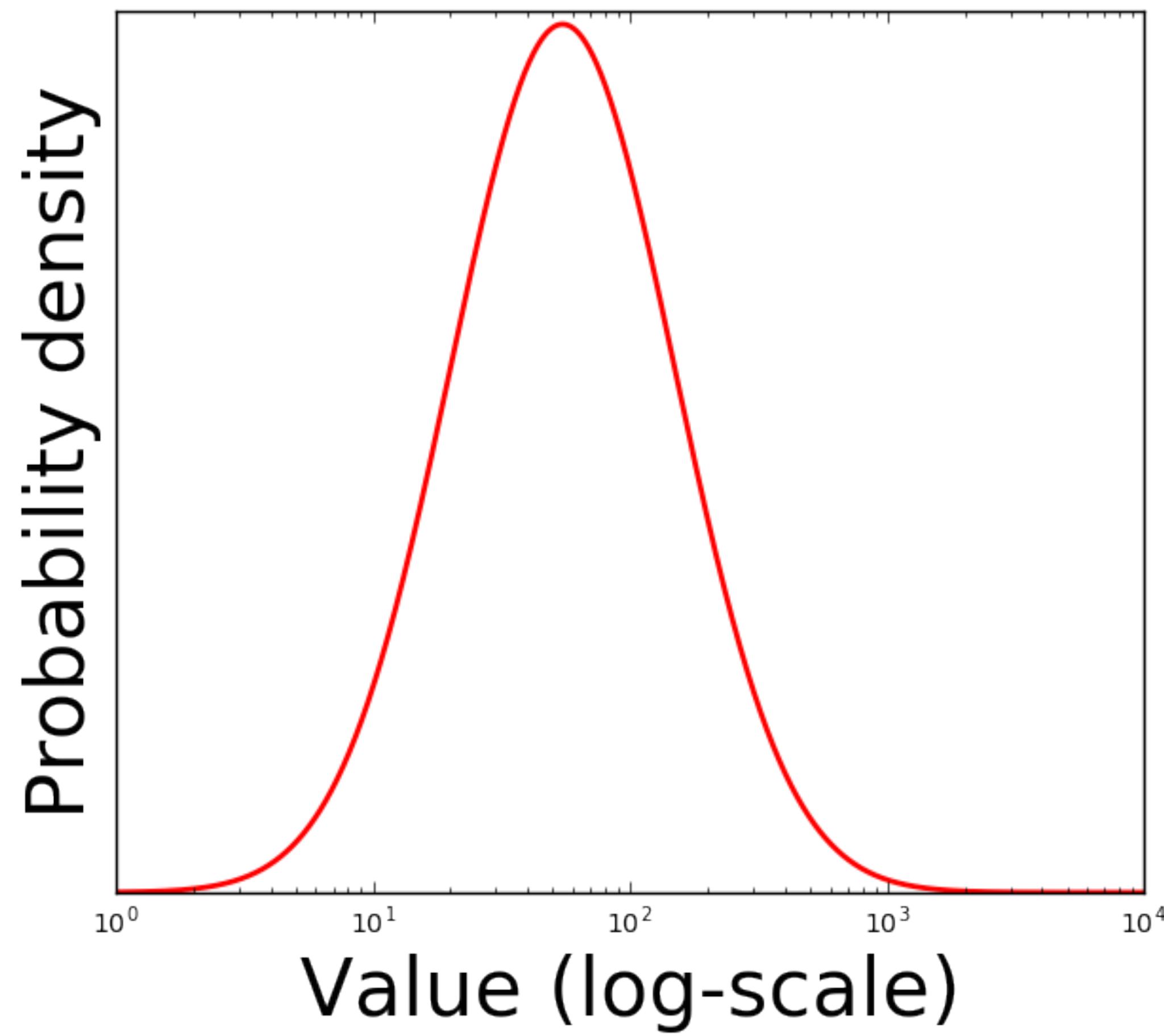
Many datasets span
many orders of magnitude.

Multiplicative process in linear scale
= Additive process in log-scale

Log-normal distribution



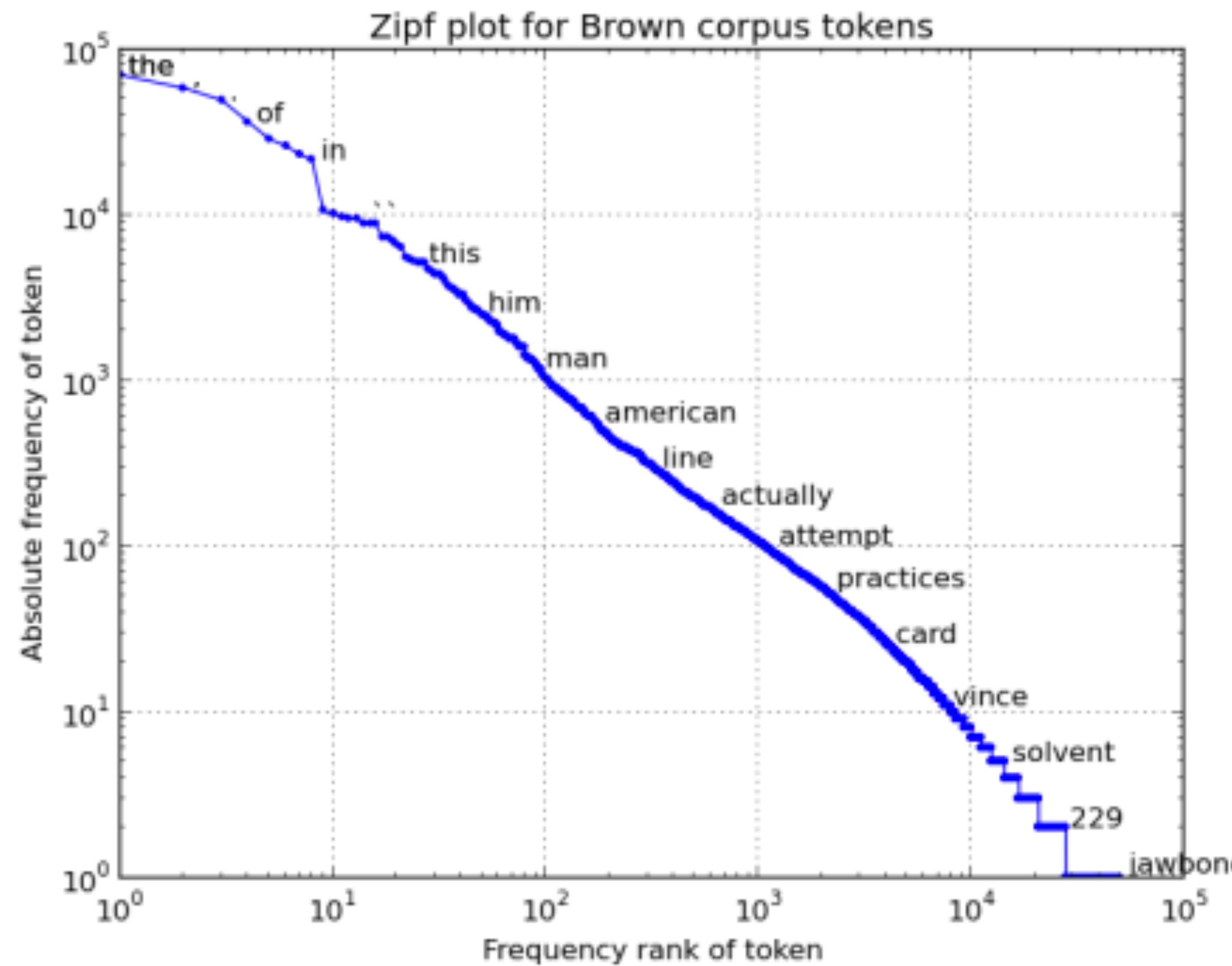
Log-normal distribution

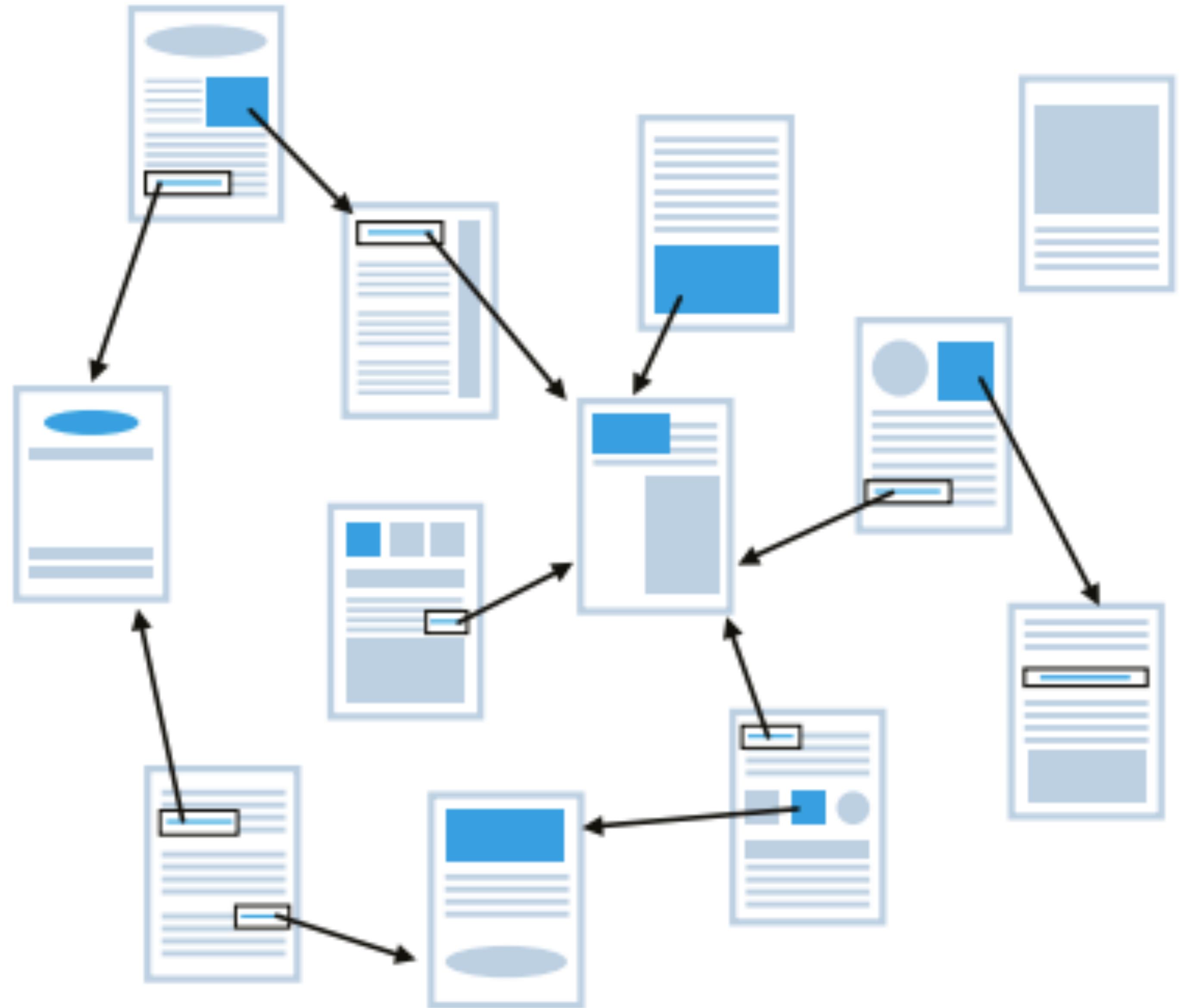


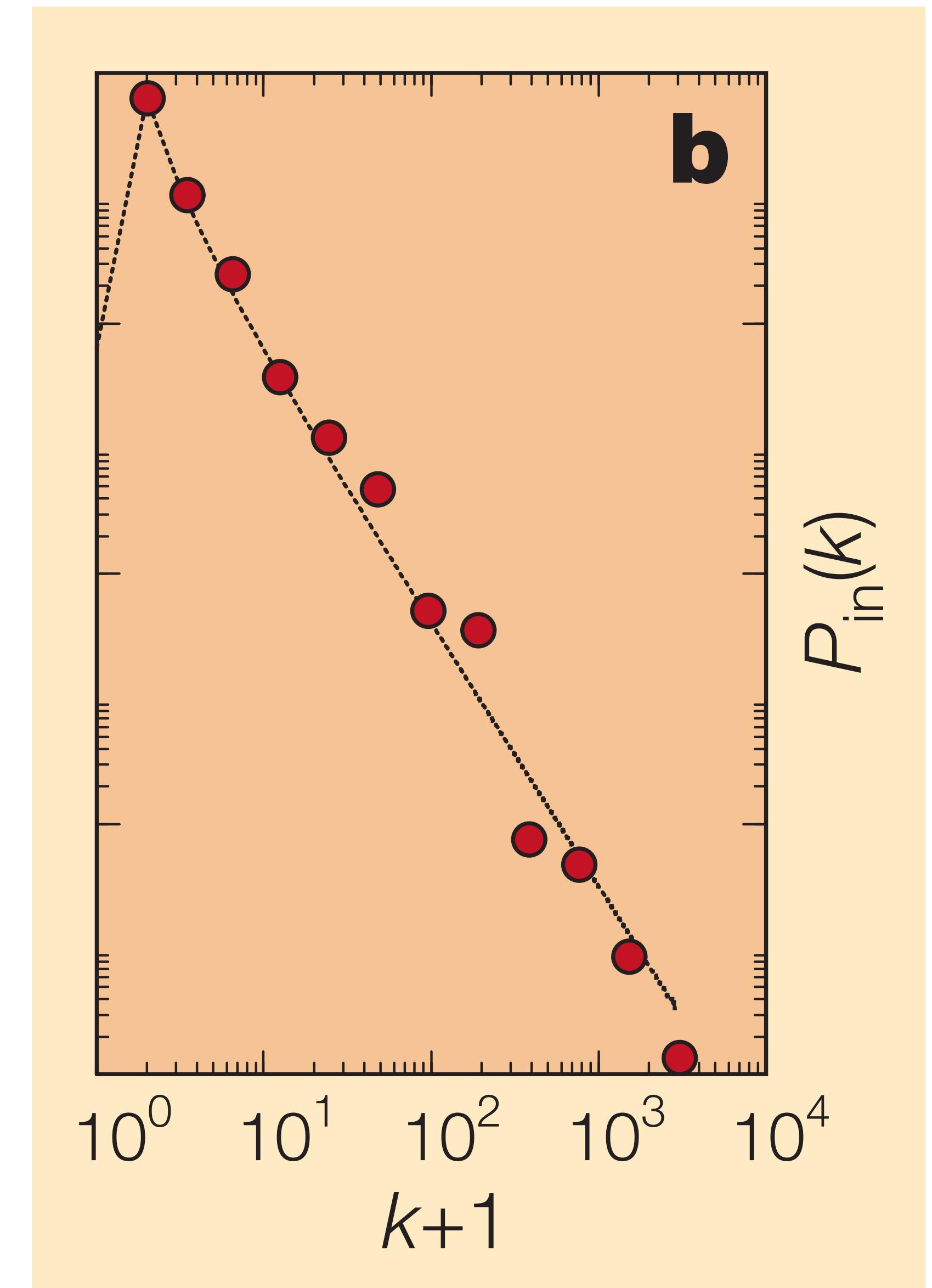
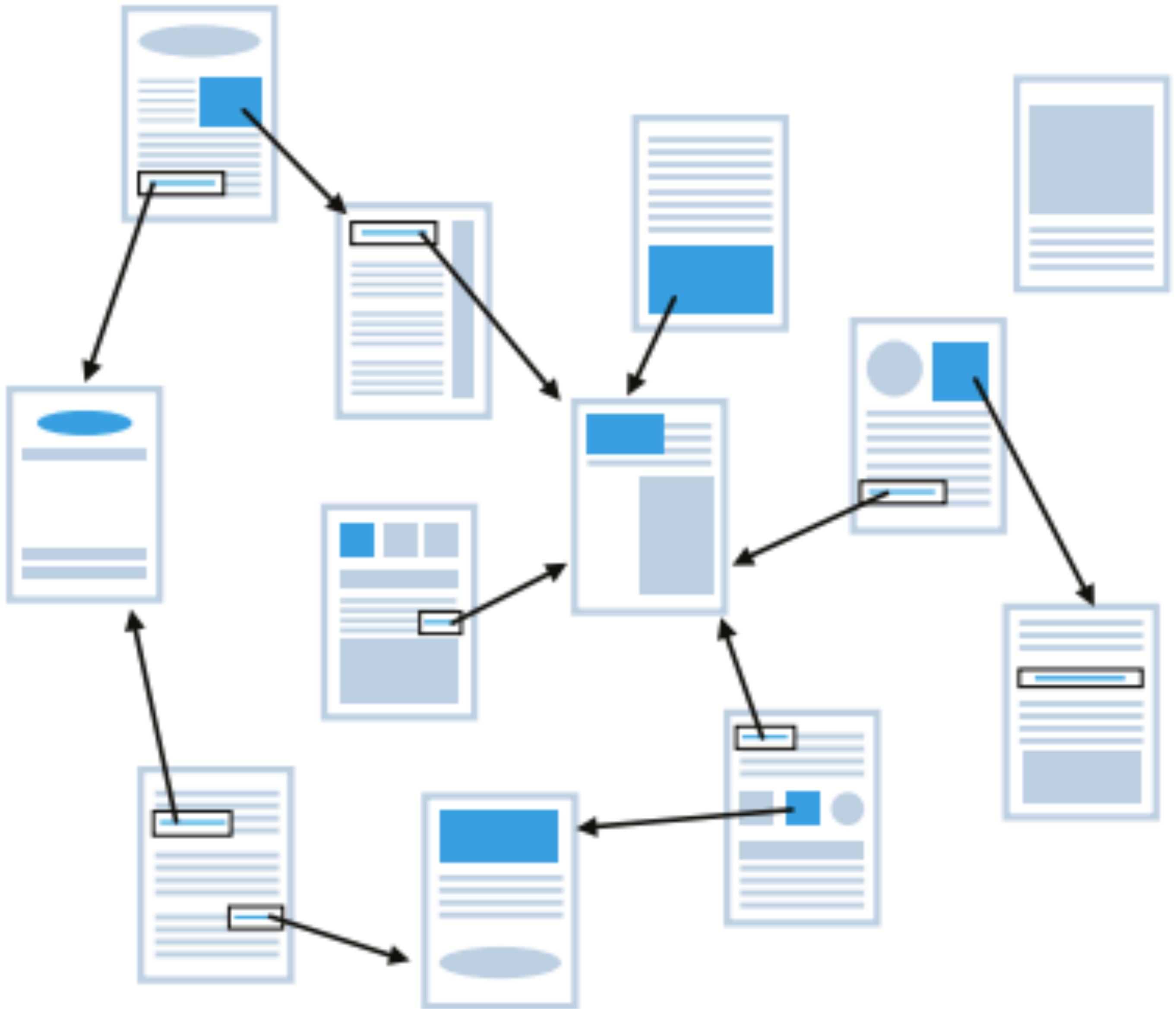
'Power-law' distributions

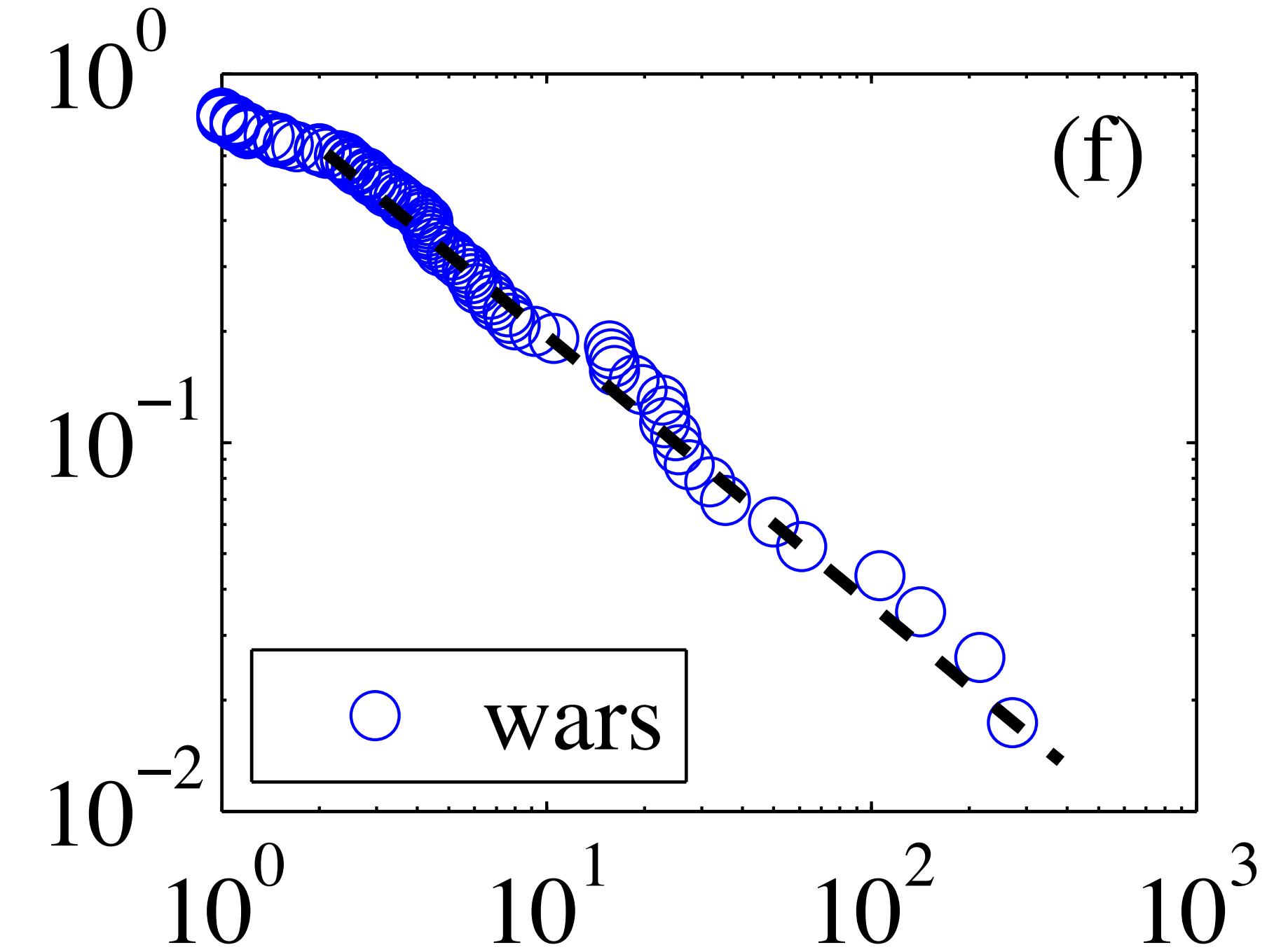
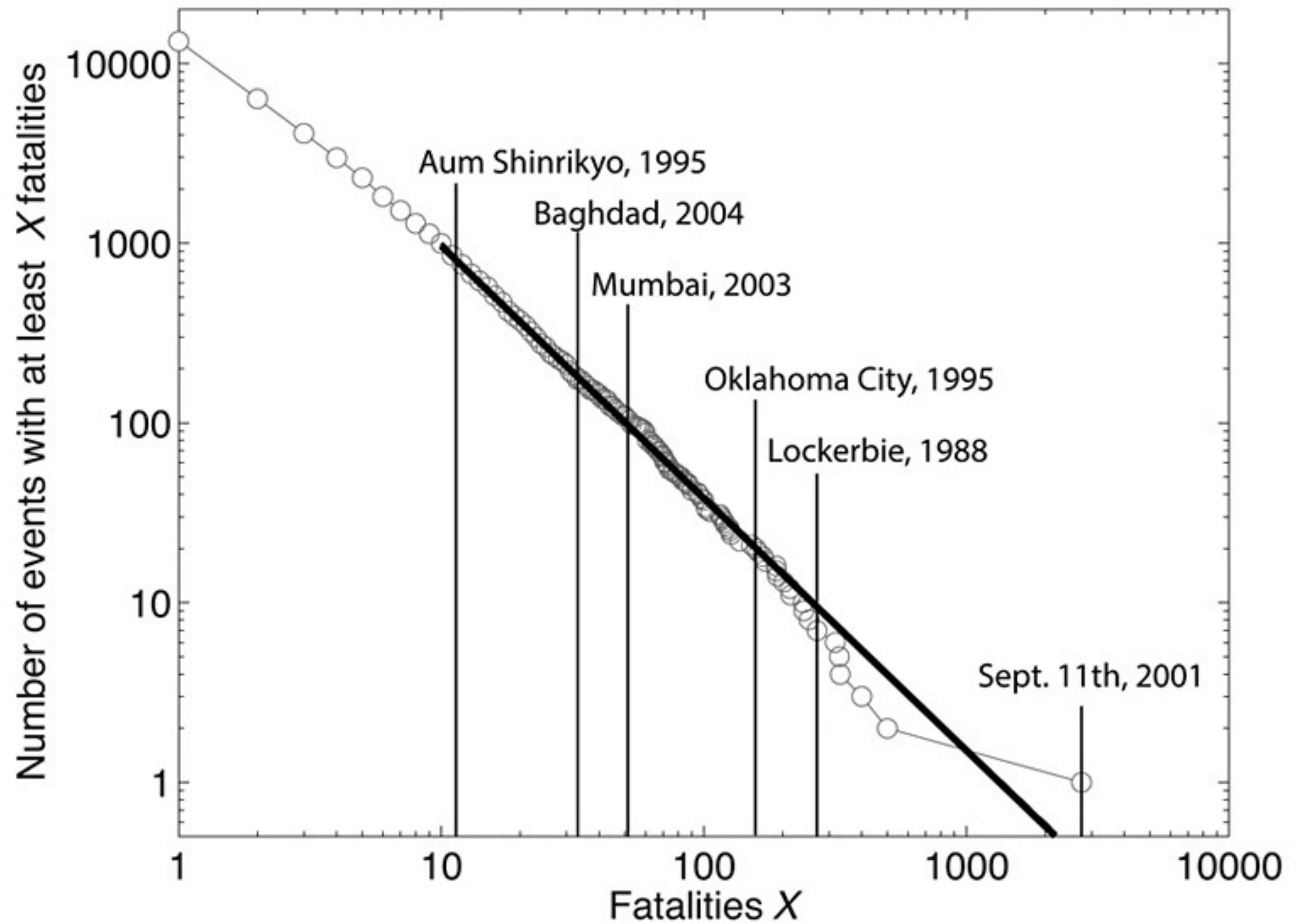
$$f(x) = cx^{-\alpha}$$

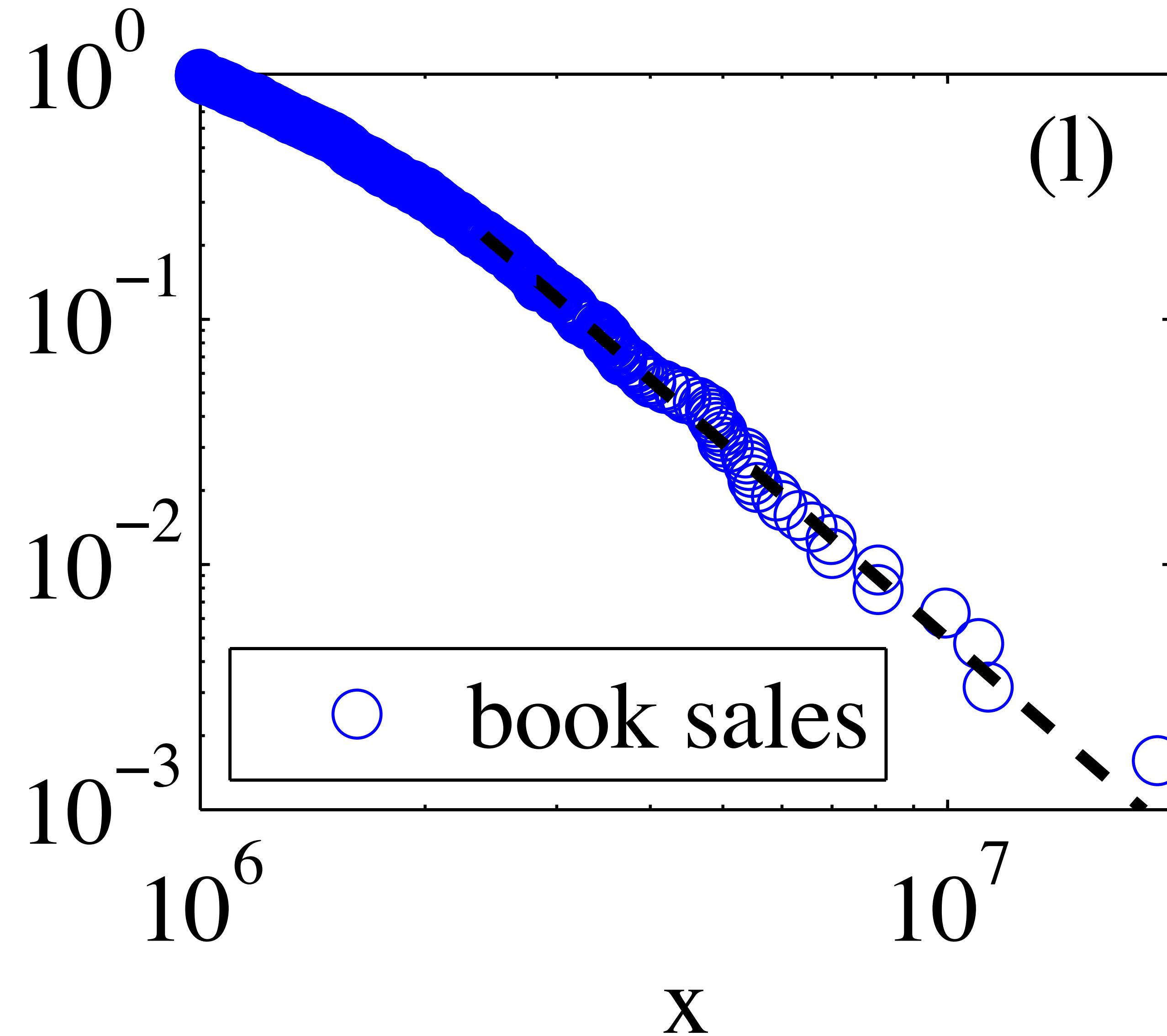
progressing in all directions
but amidst collisions preserves
a regularity in all its movements
and regulated by a law of gravitation
which is the law of the motion of
all matter in the universe.



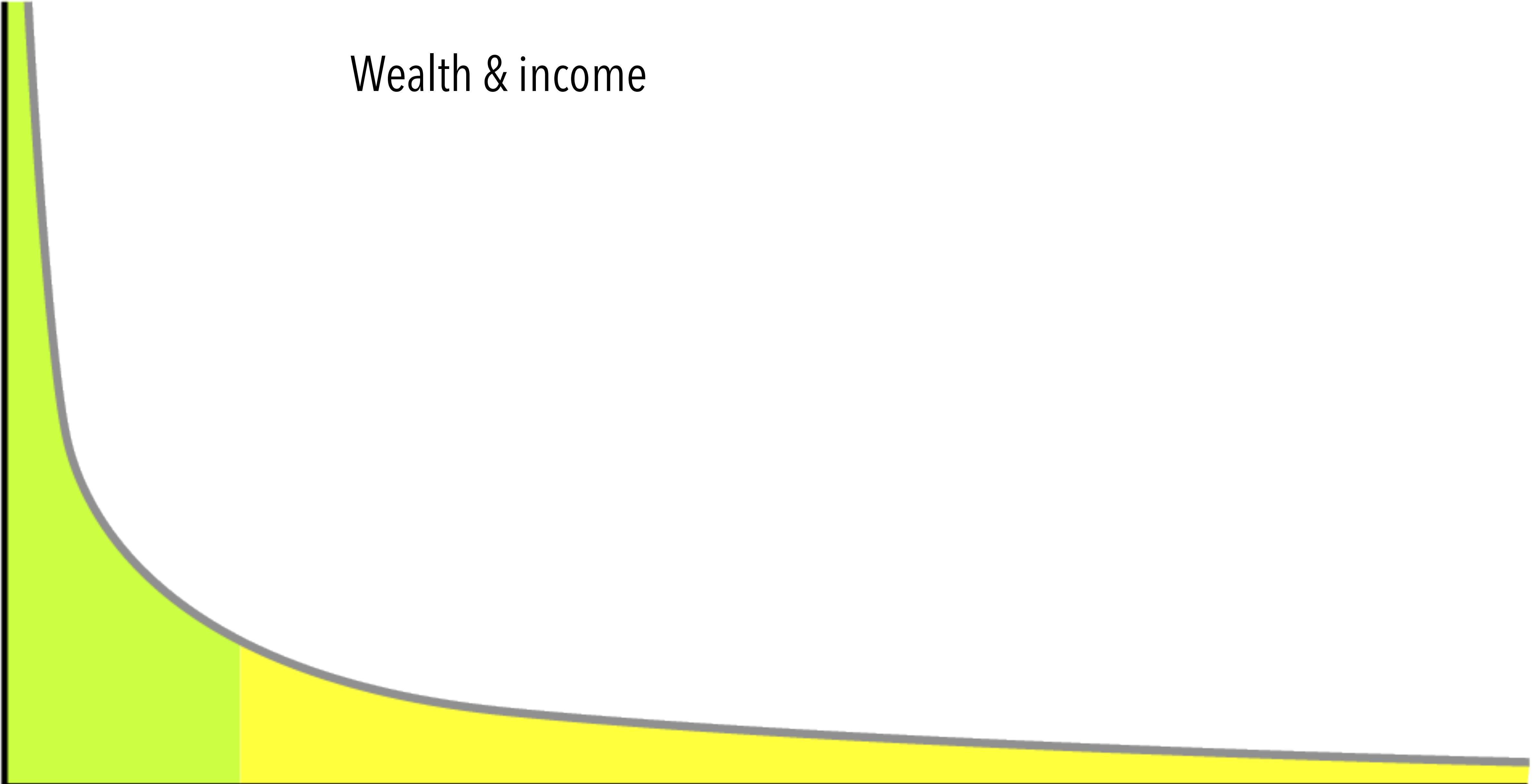




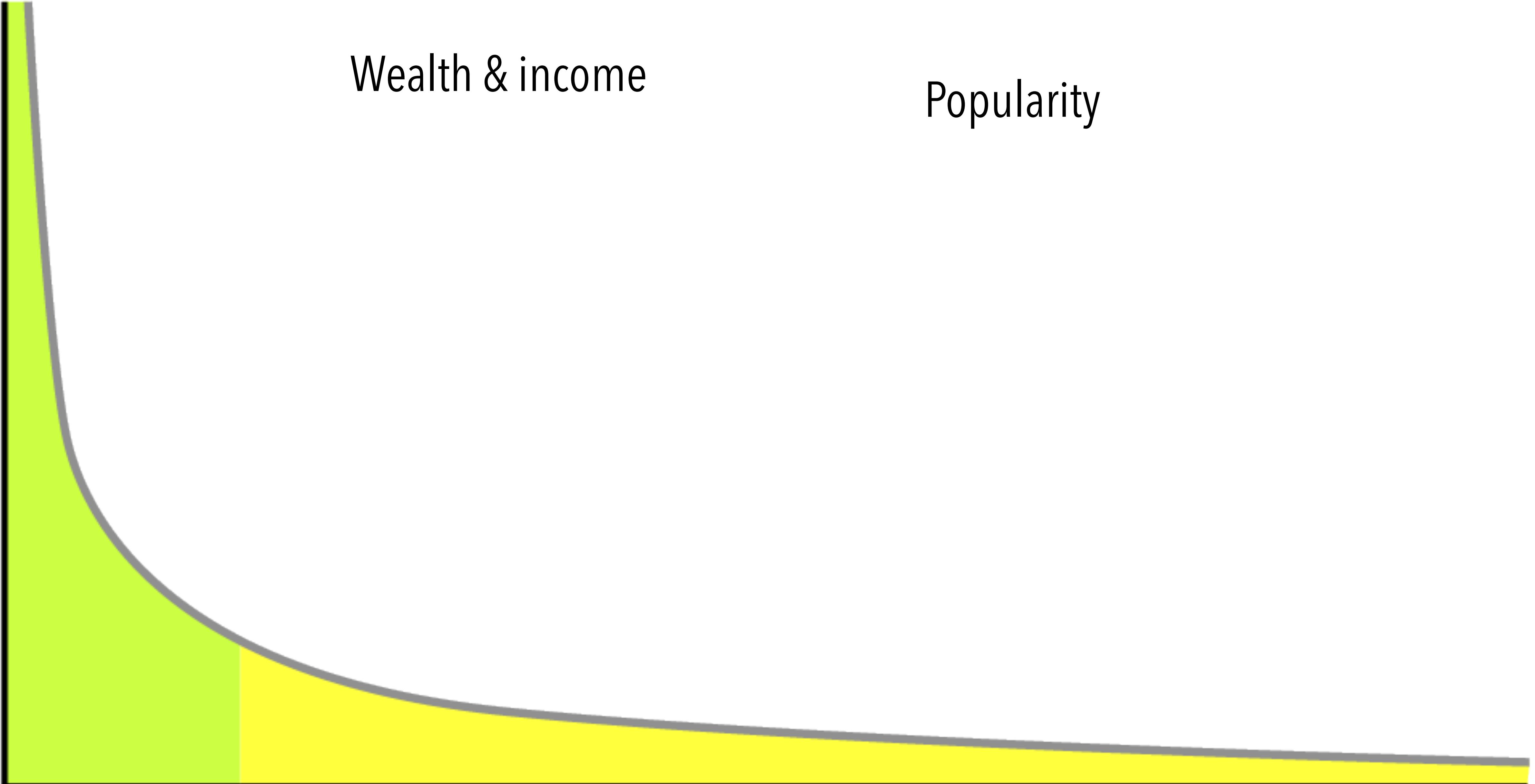






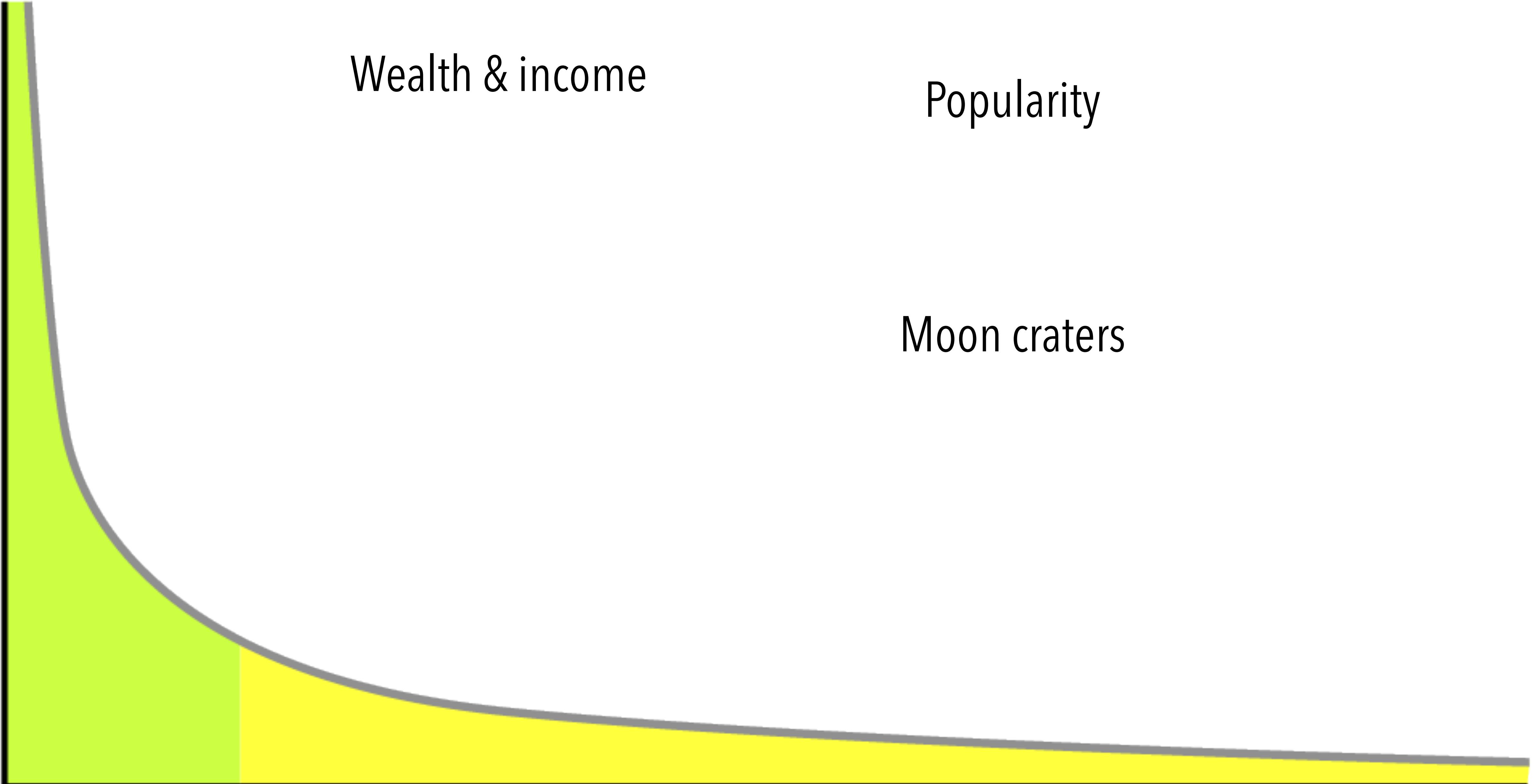


Wealth & income



Wealth & income

Popularity

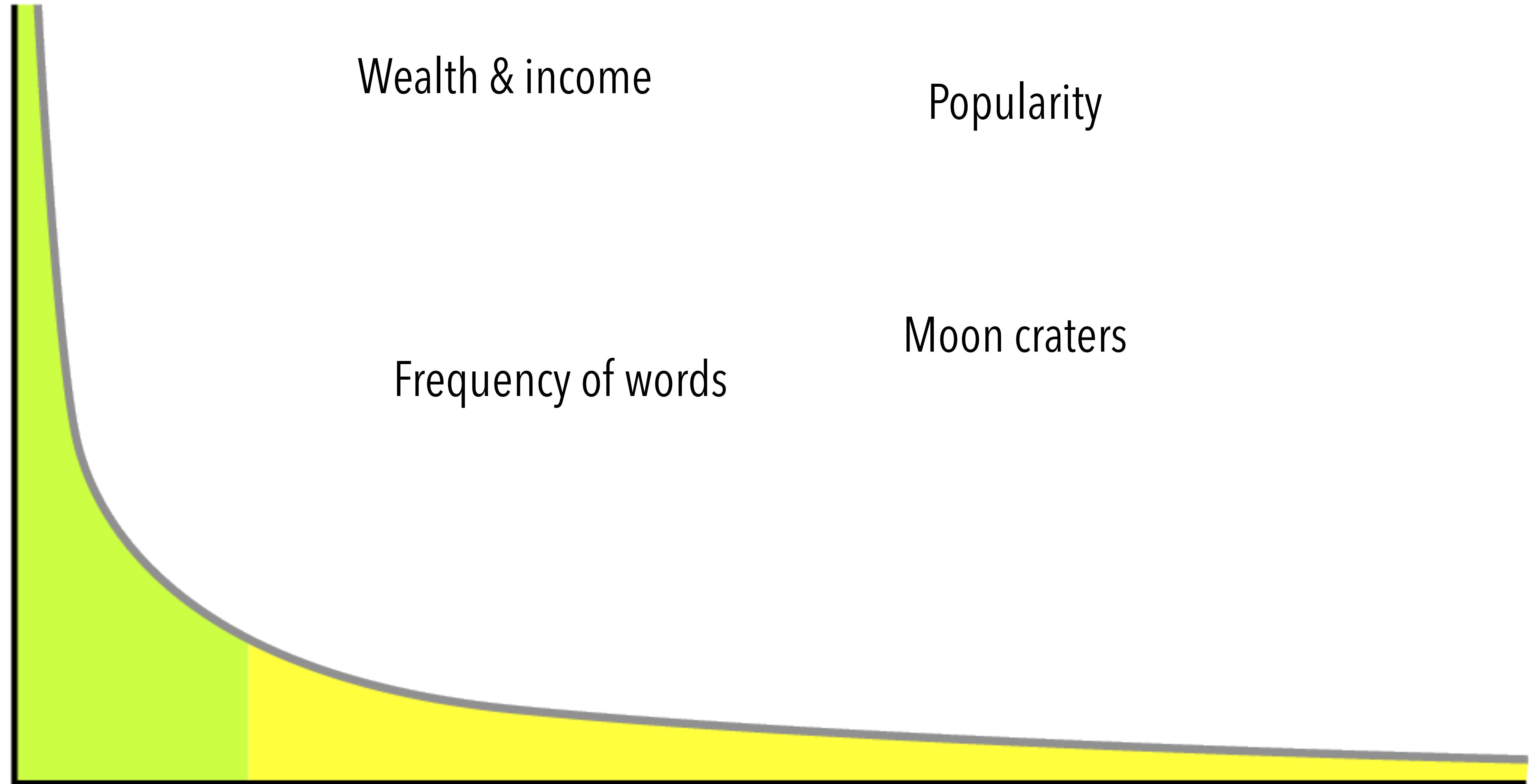


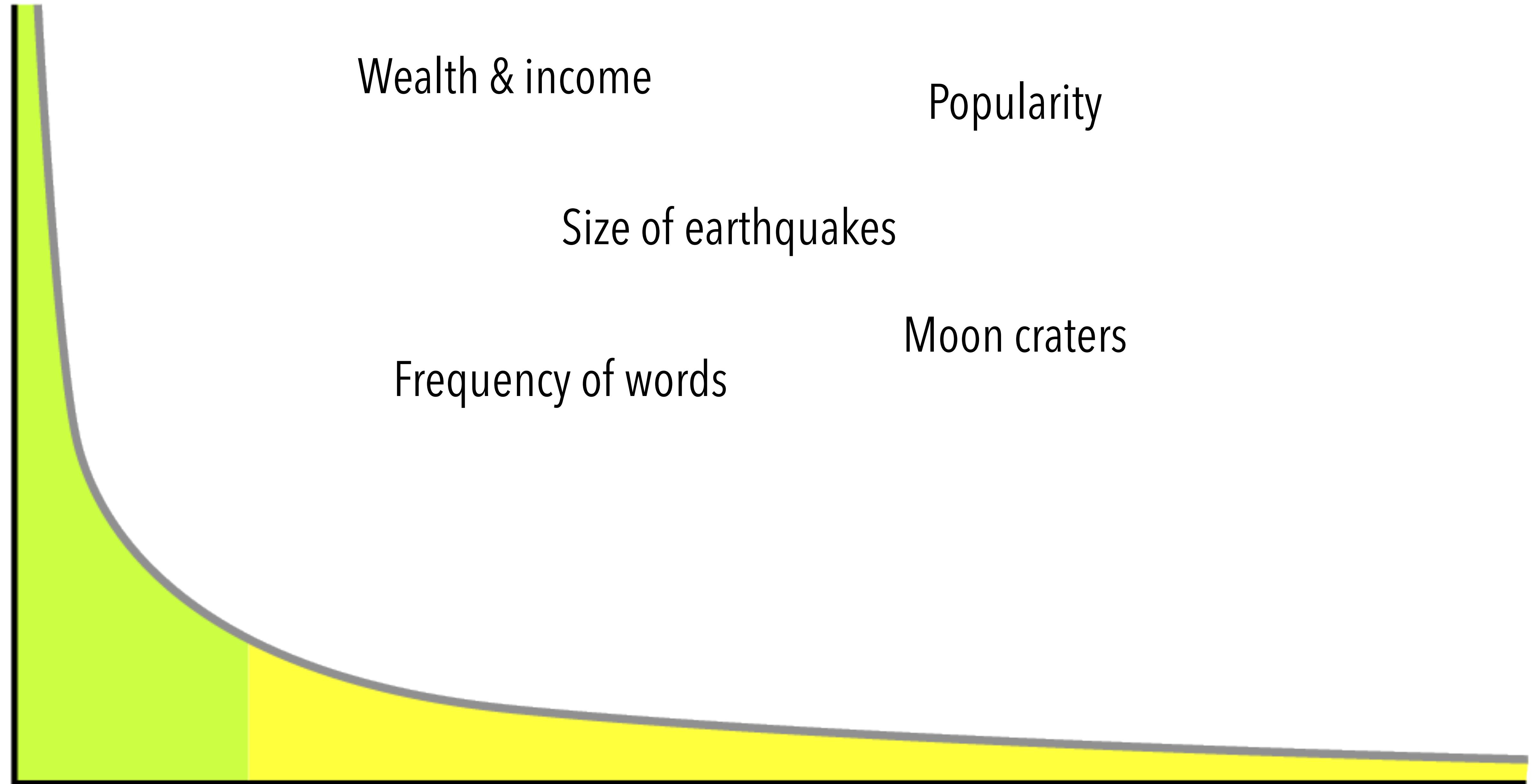
A yellow bell curve graph showing the distribution of wealth and income. The curve starts at a very high frequency on the left and tapers off towards the right. The peak of the curve is located on the far left.

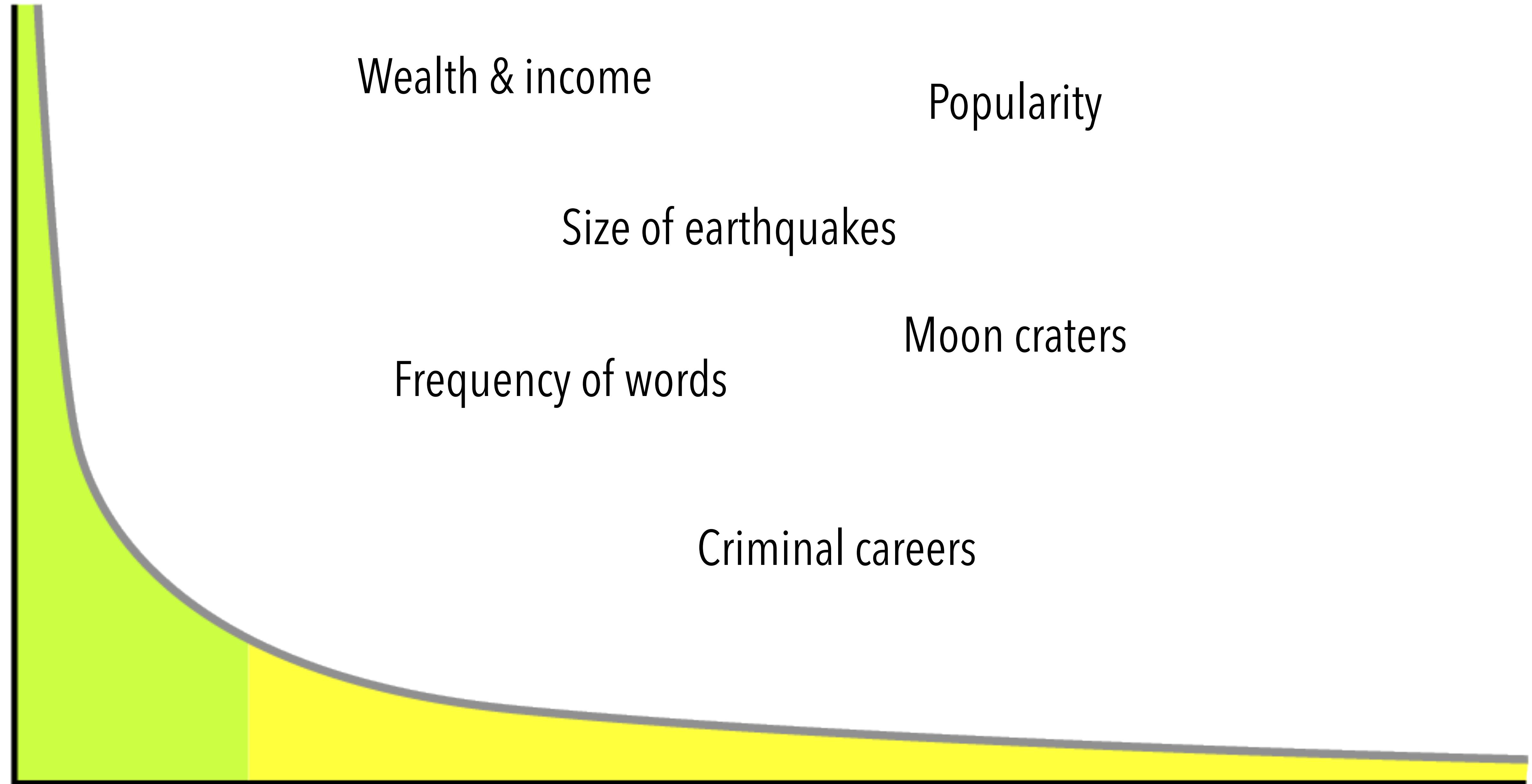
Wealth & income

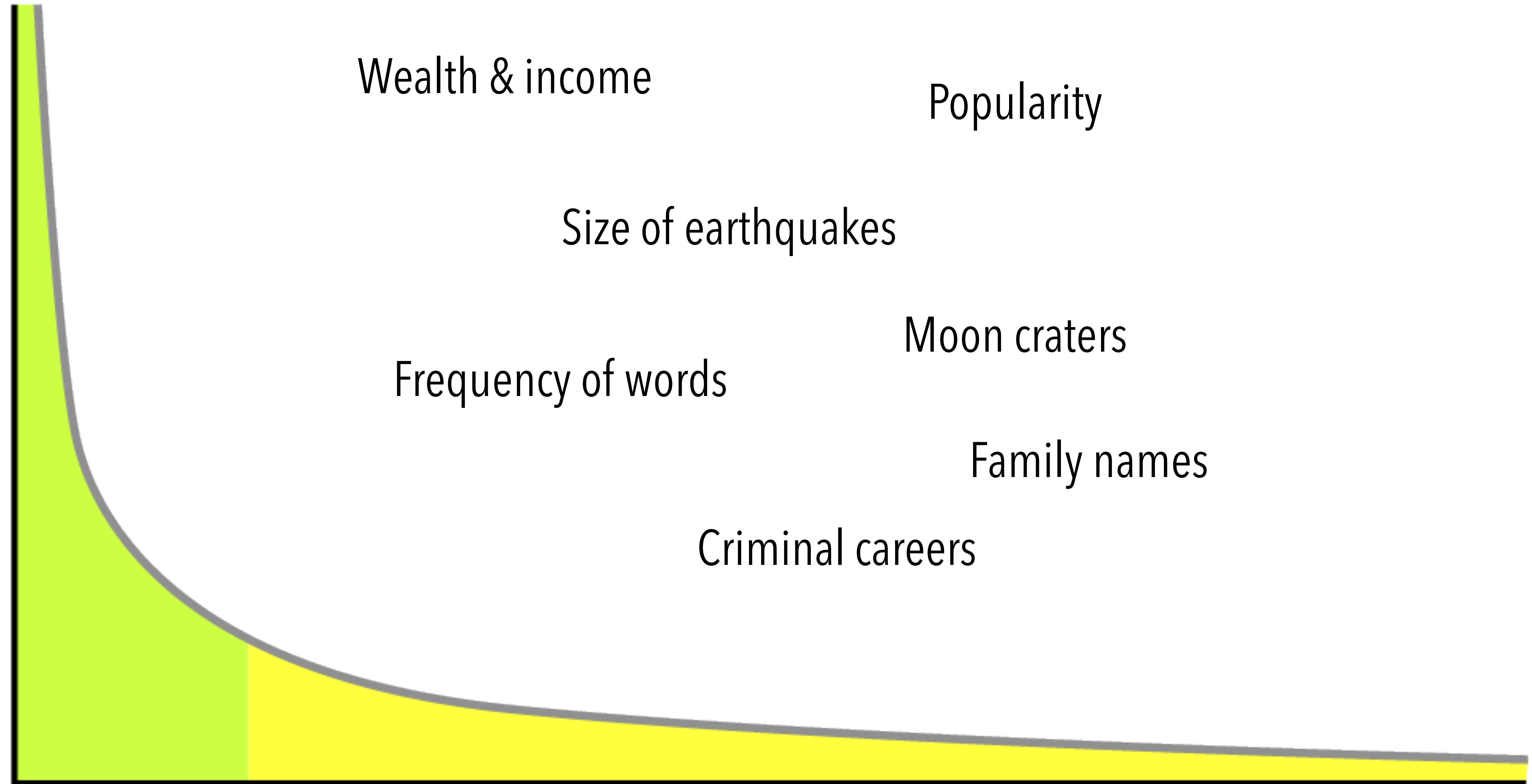
Popularity

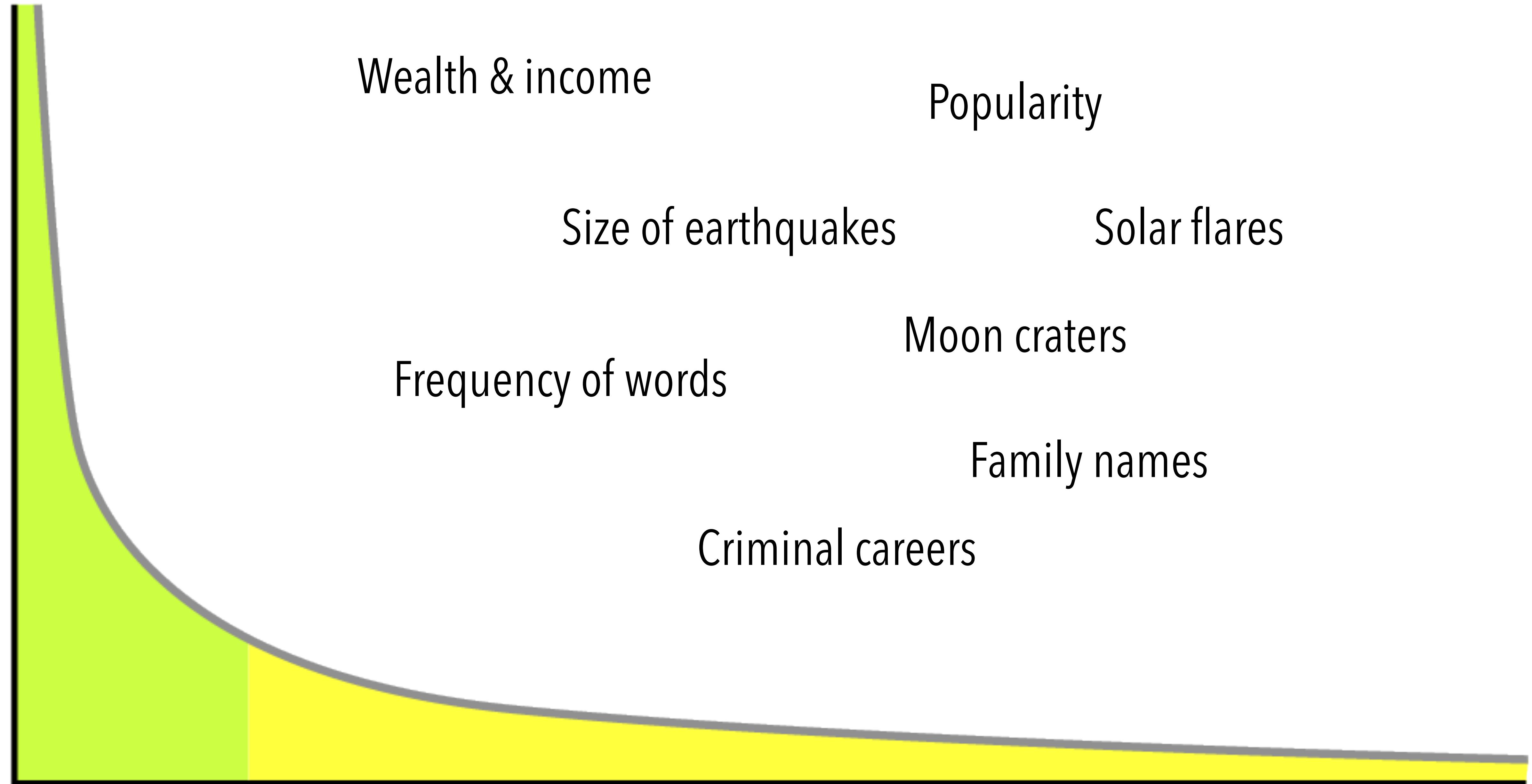
Moon craters

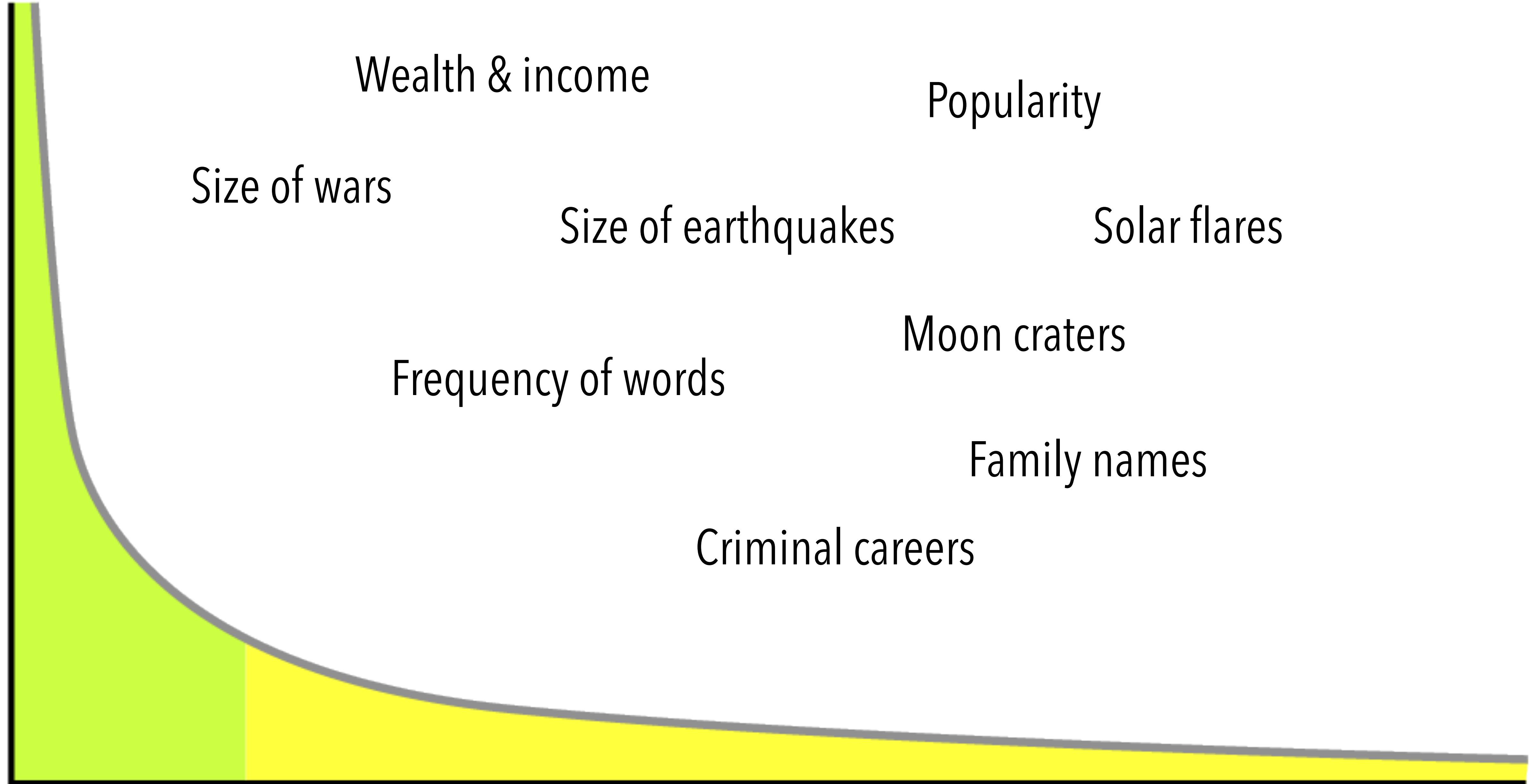


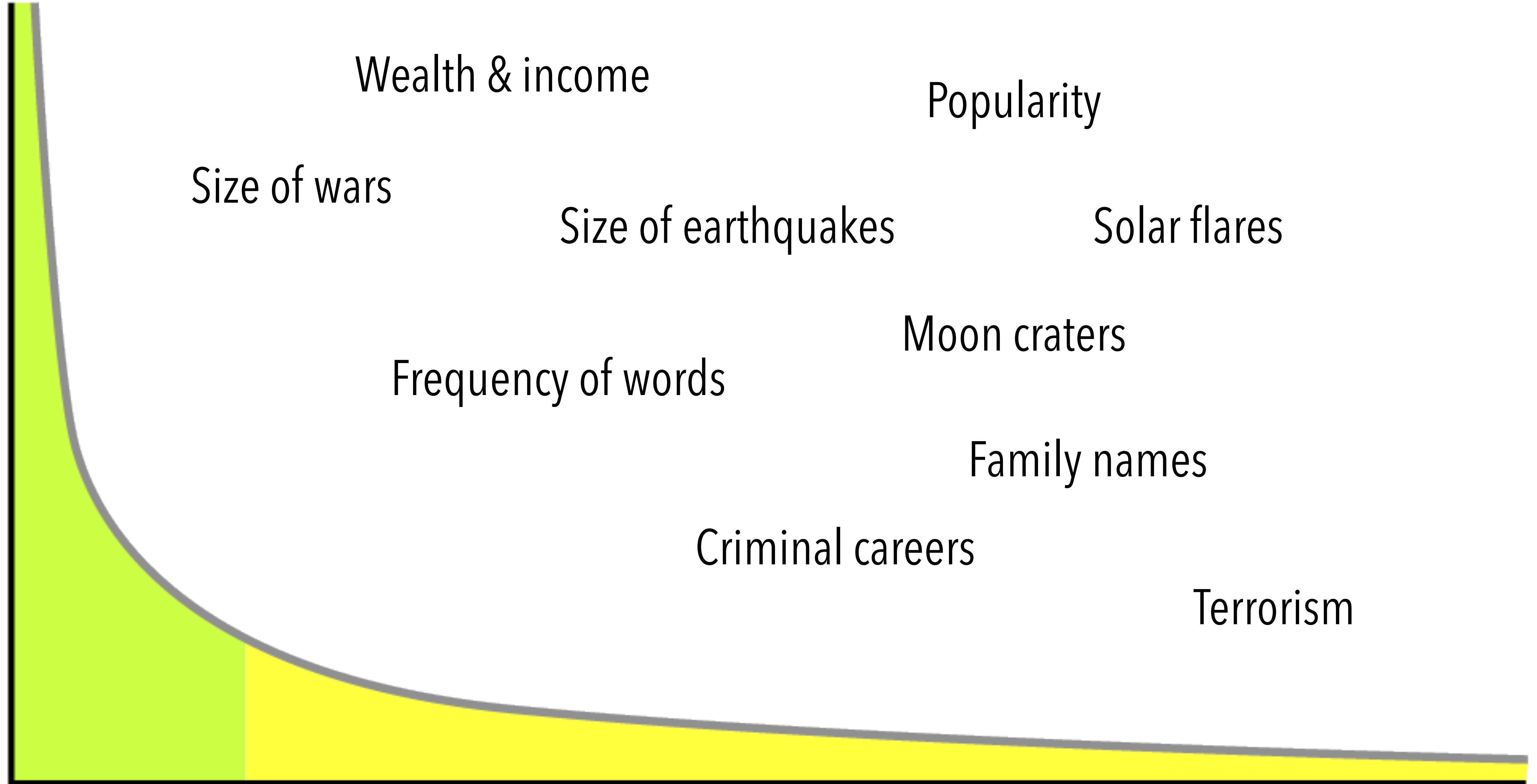


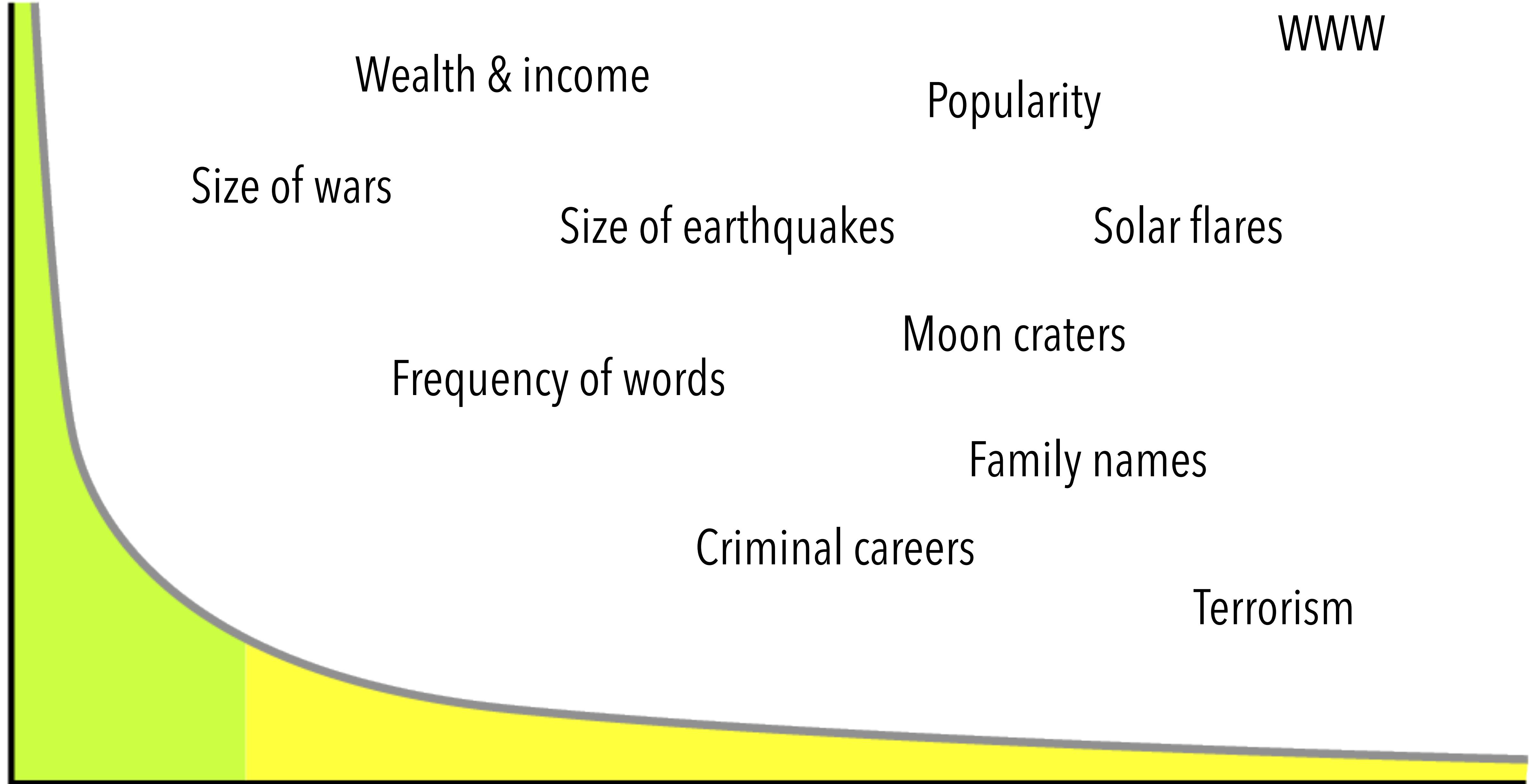


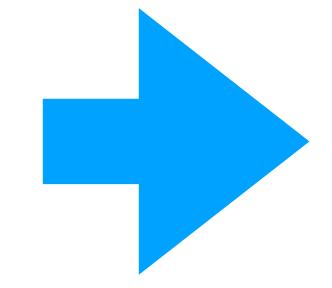
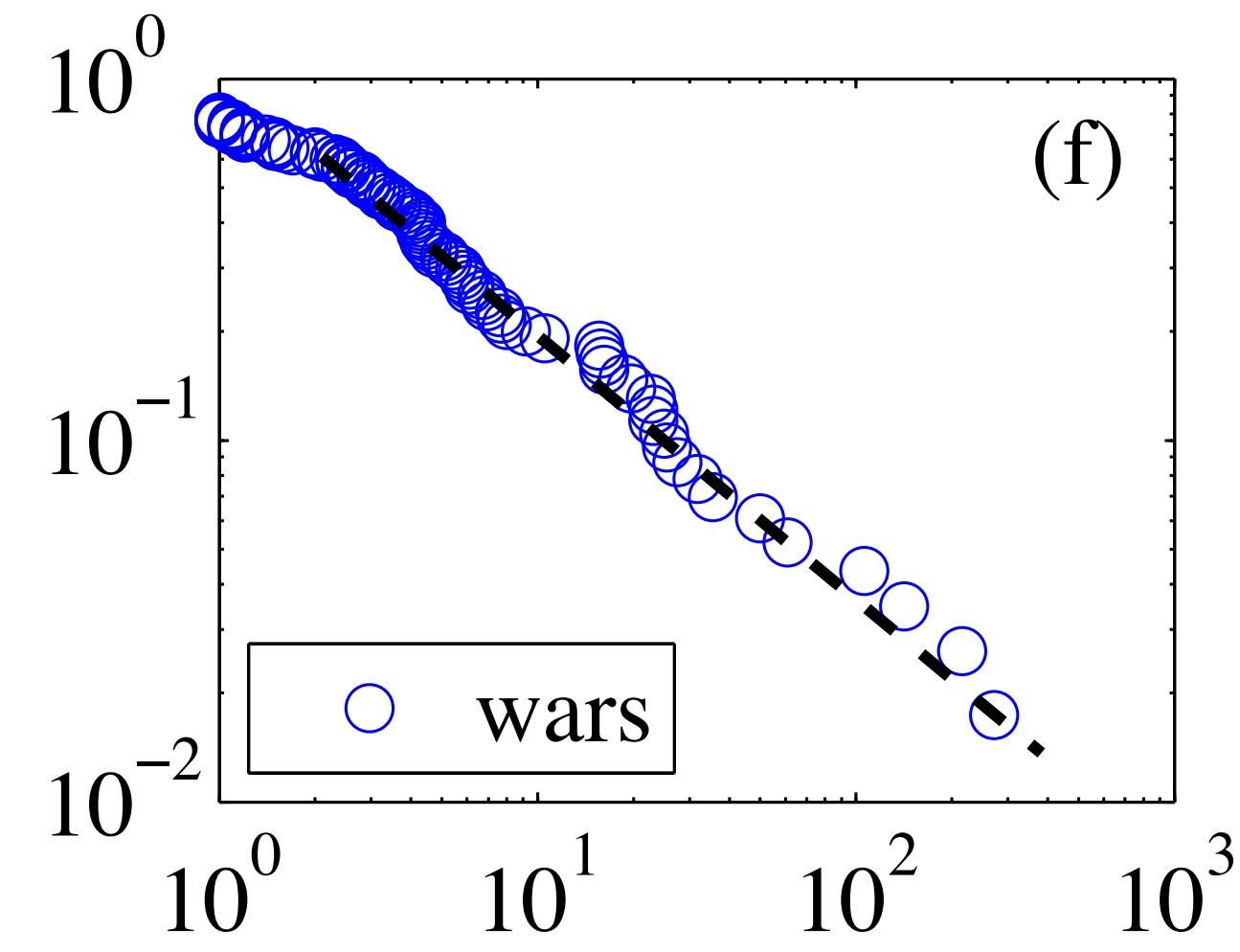
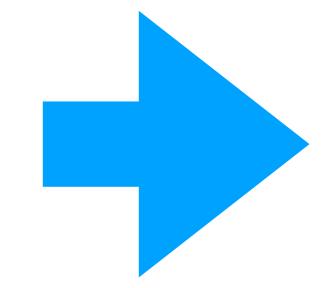
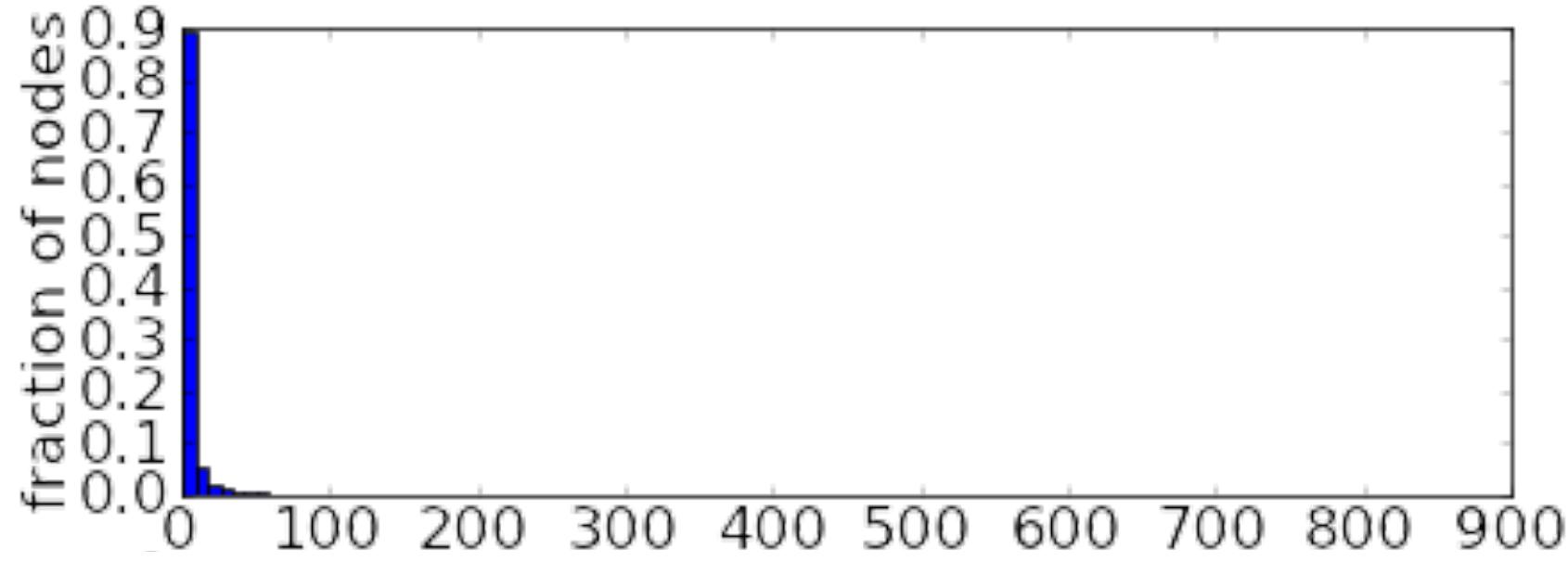
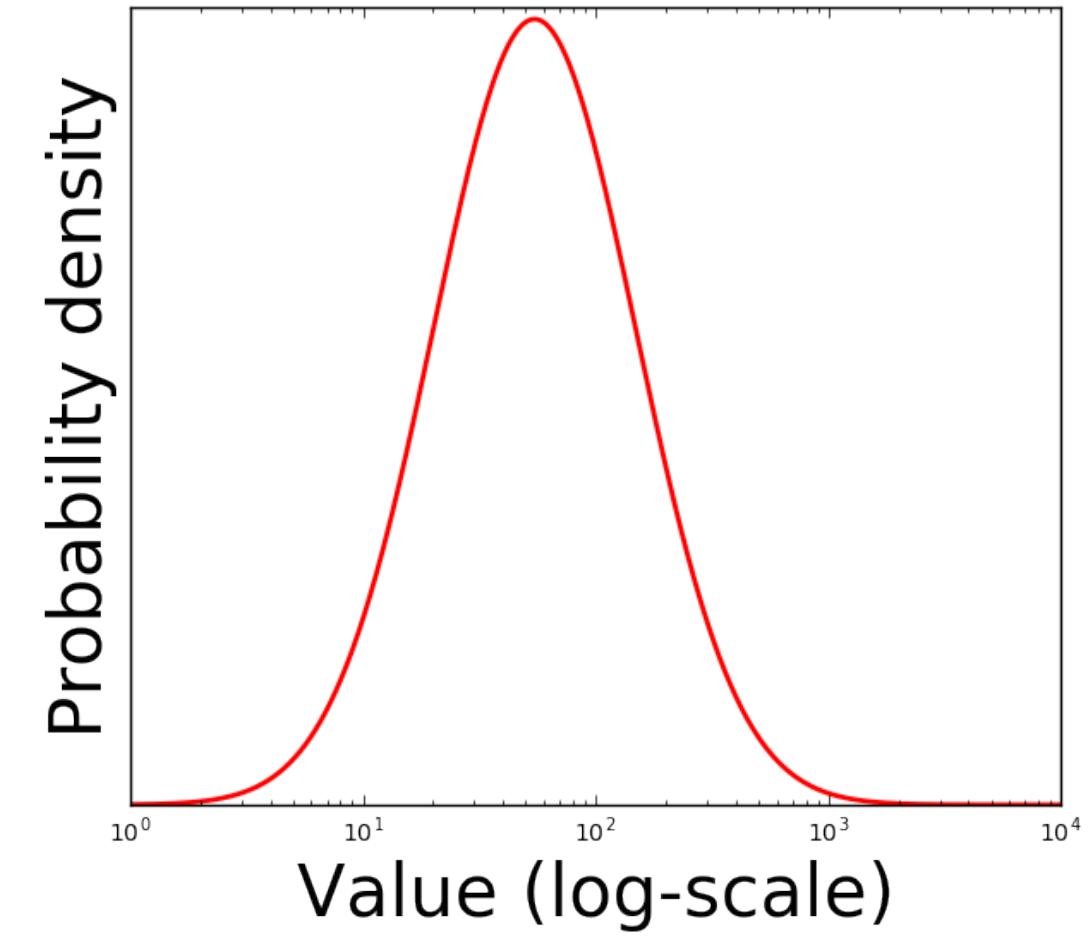


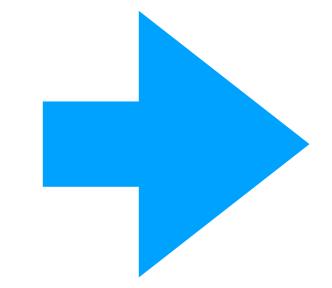
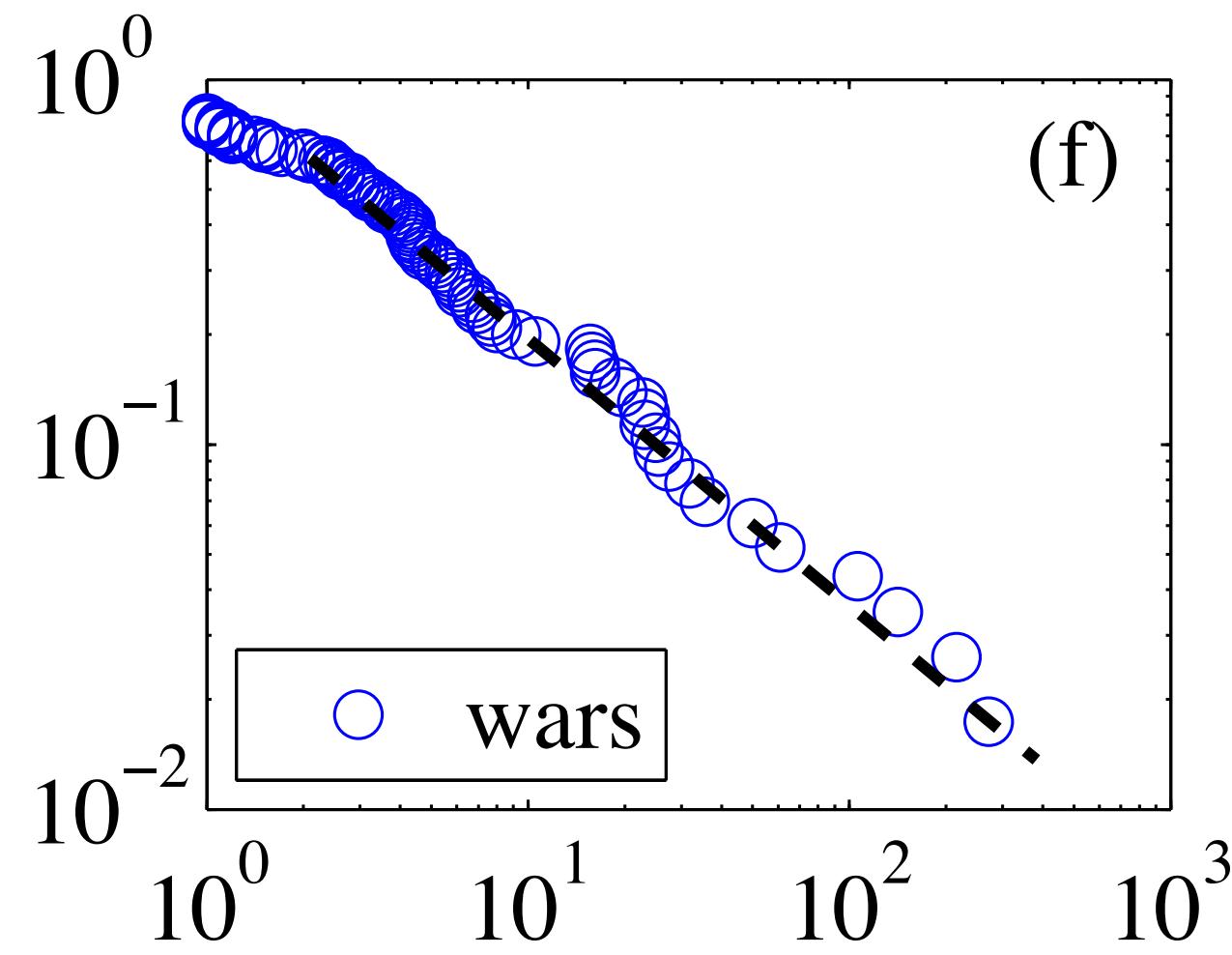
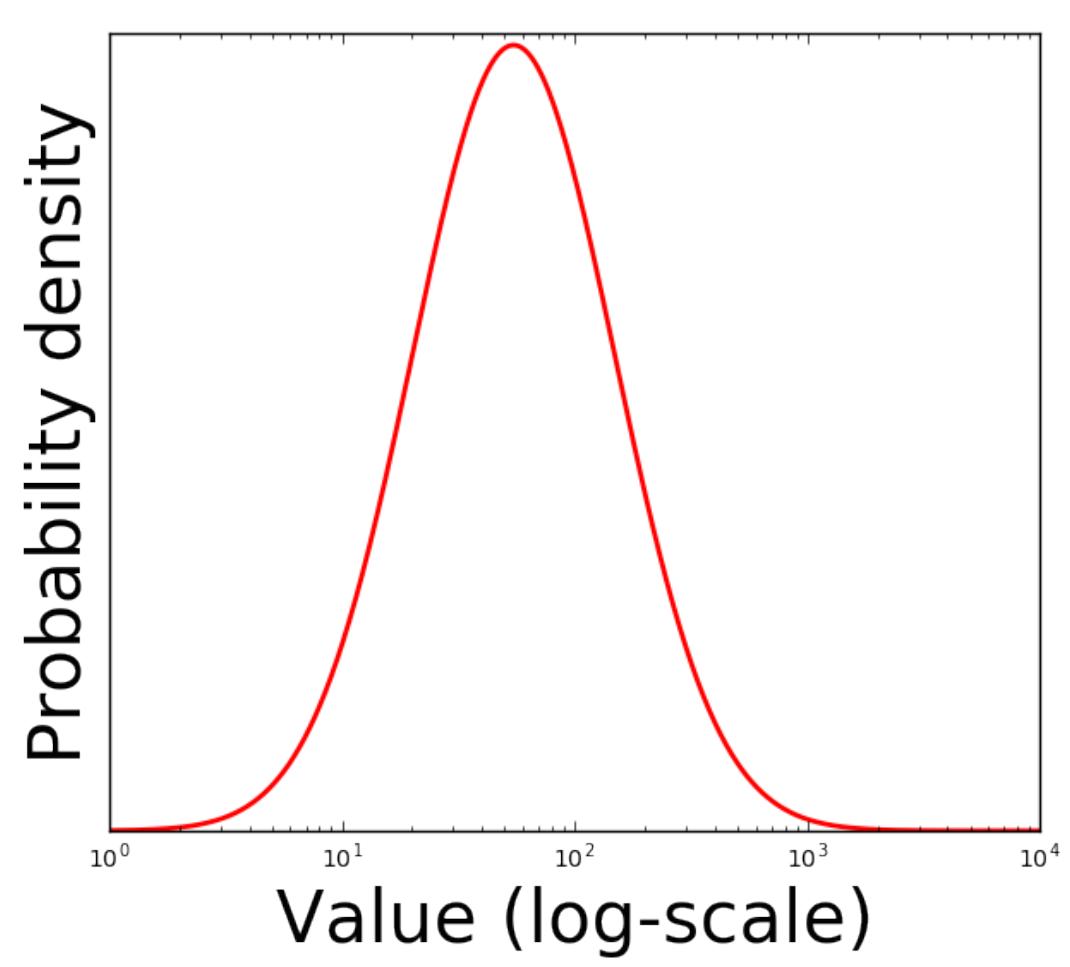
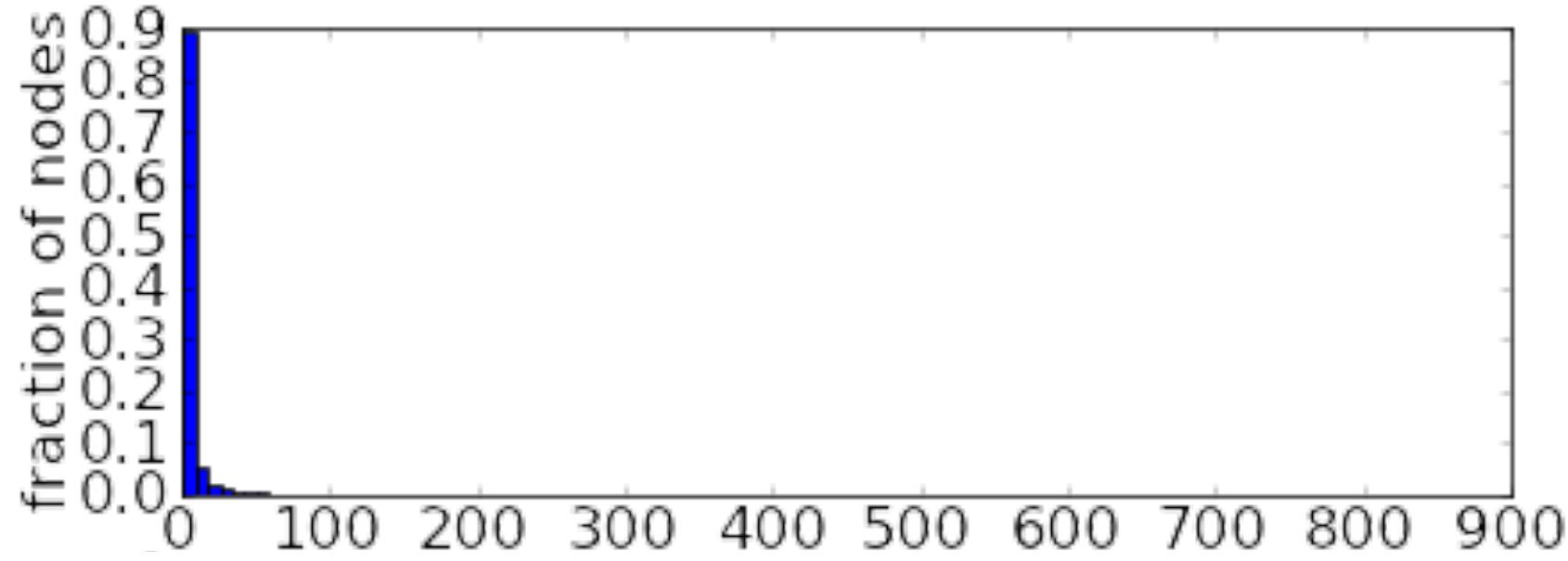




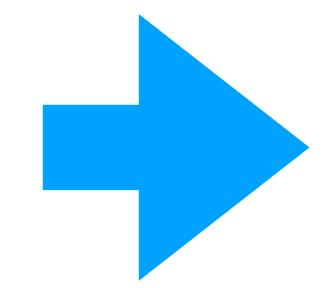


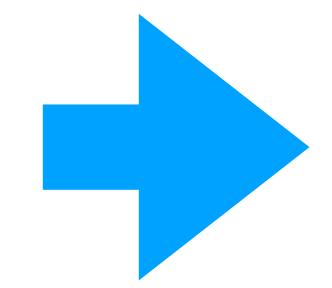
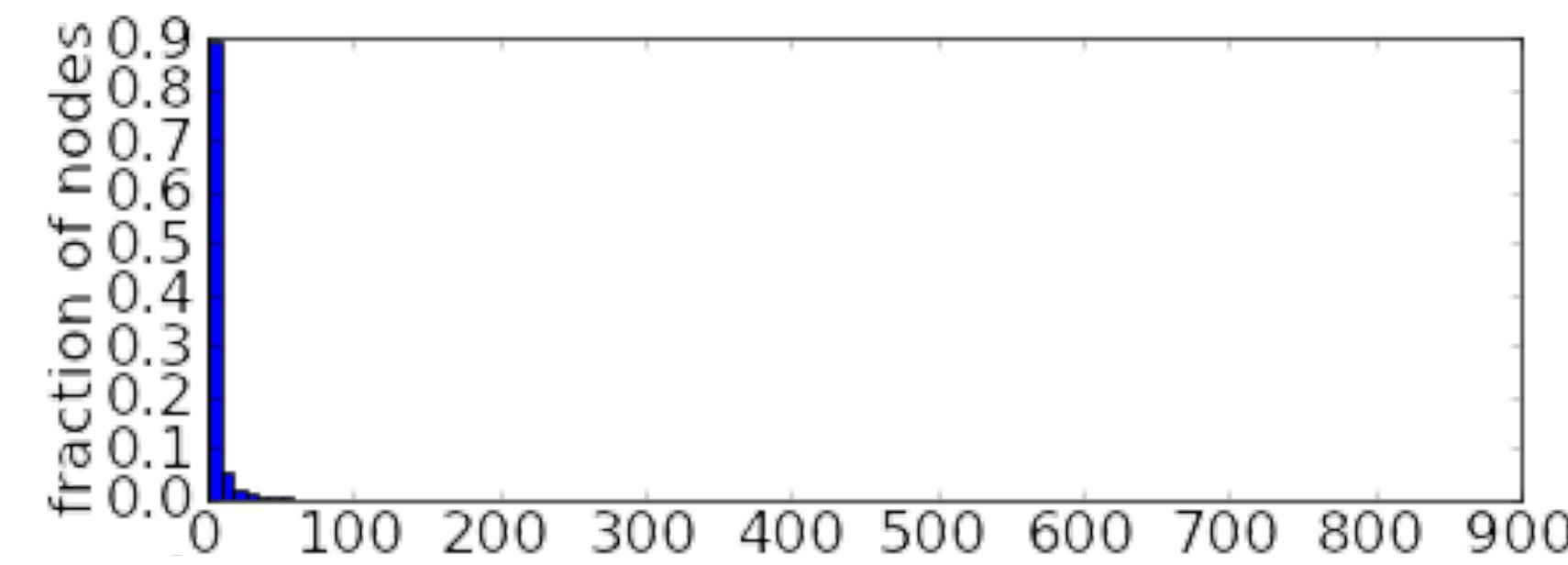
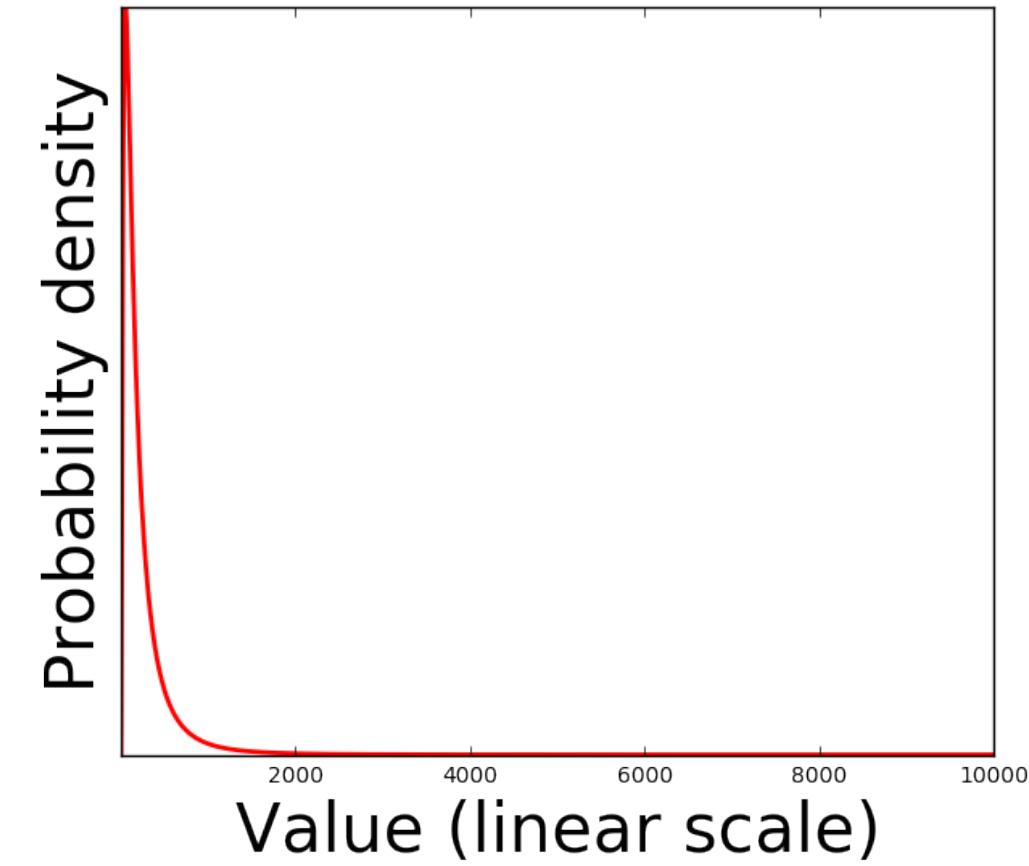




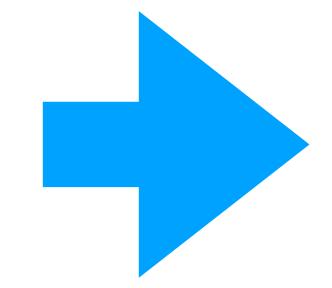
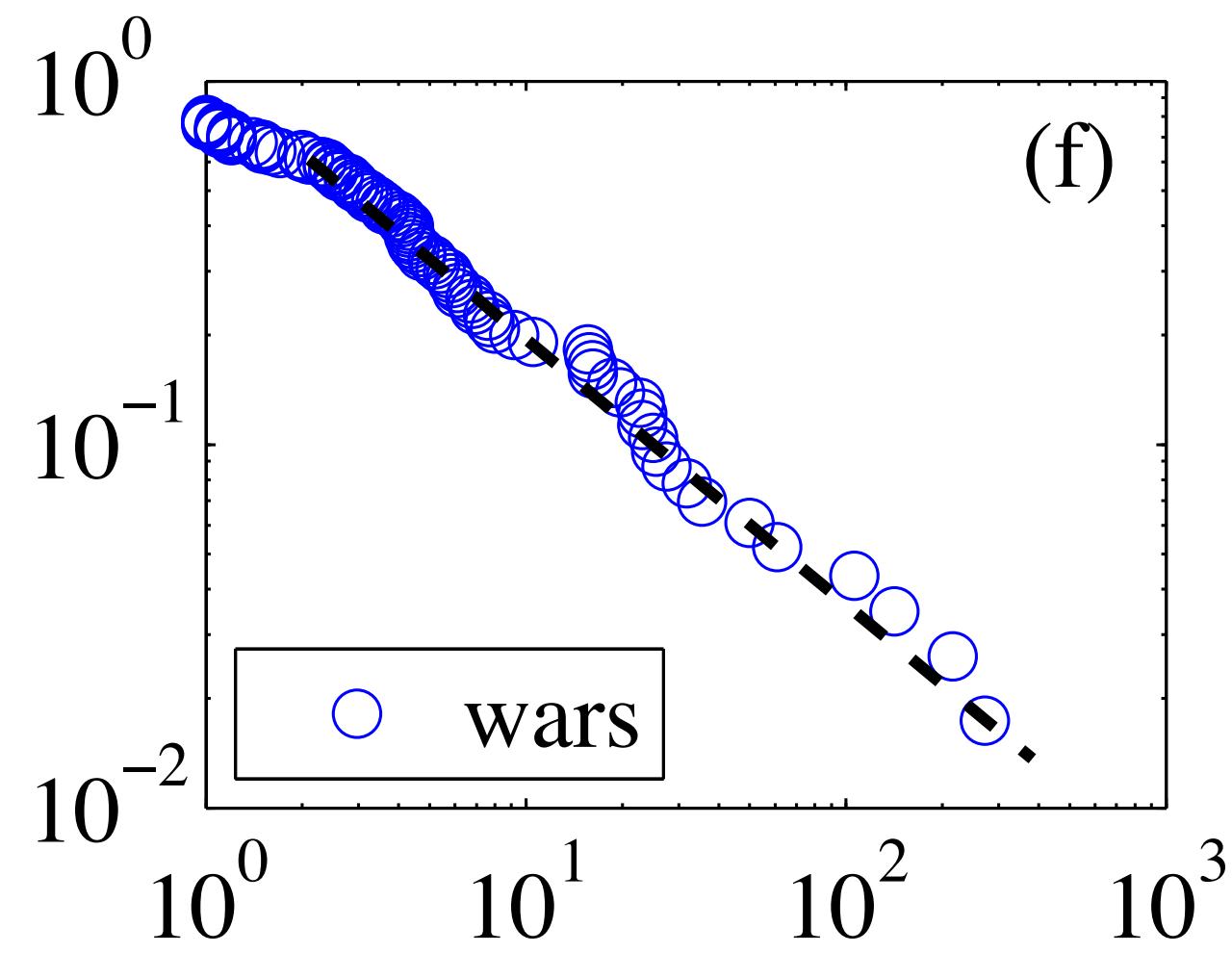
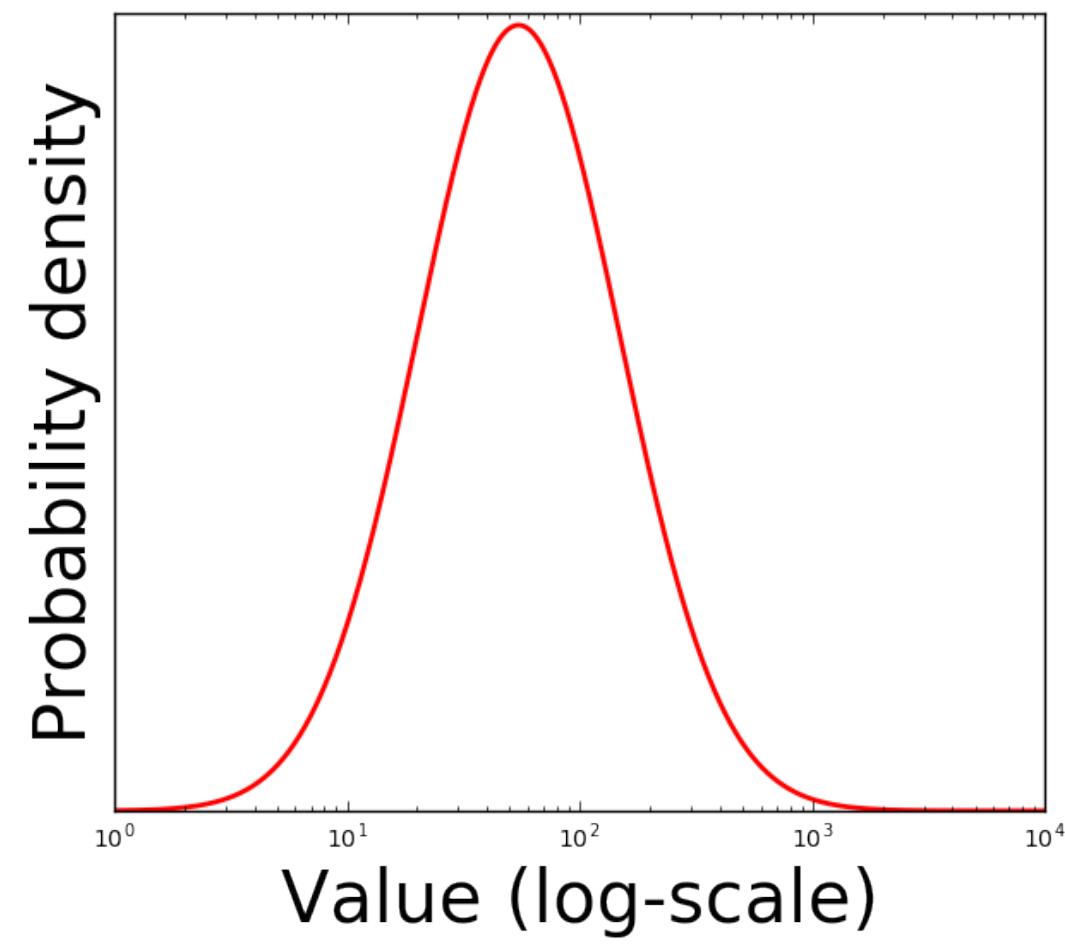


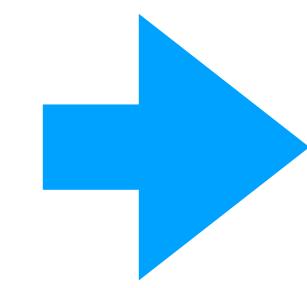
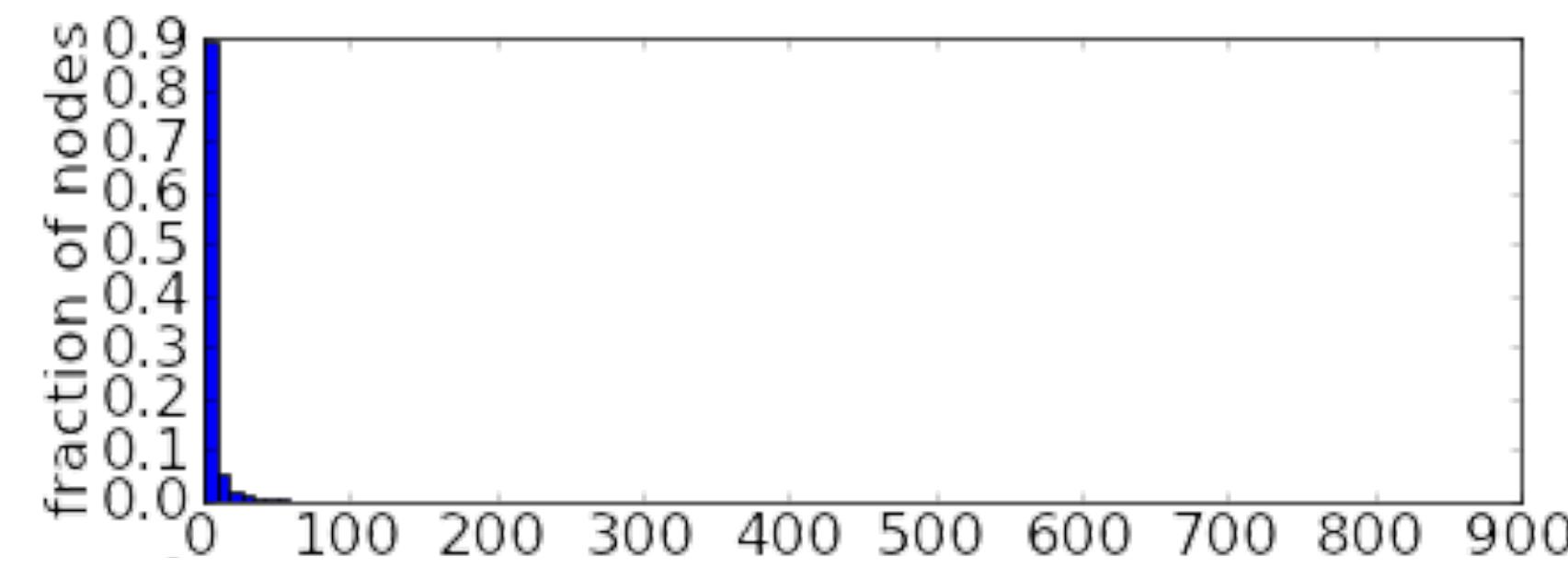
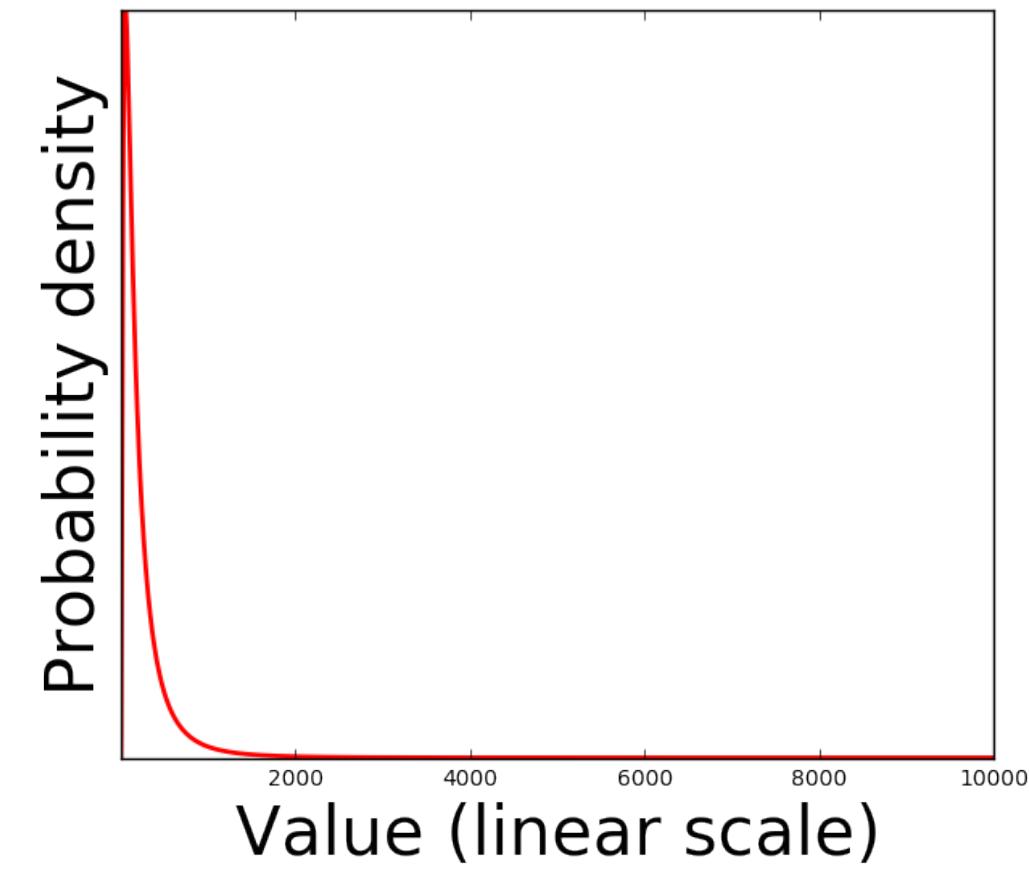
Linear scale:
Can't see any!



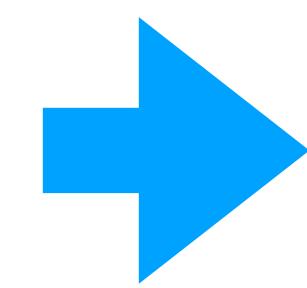
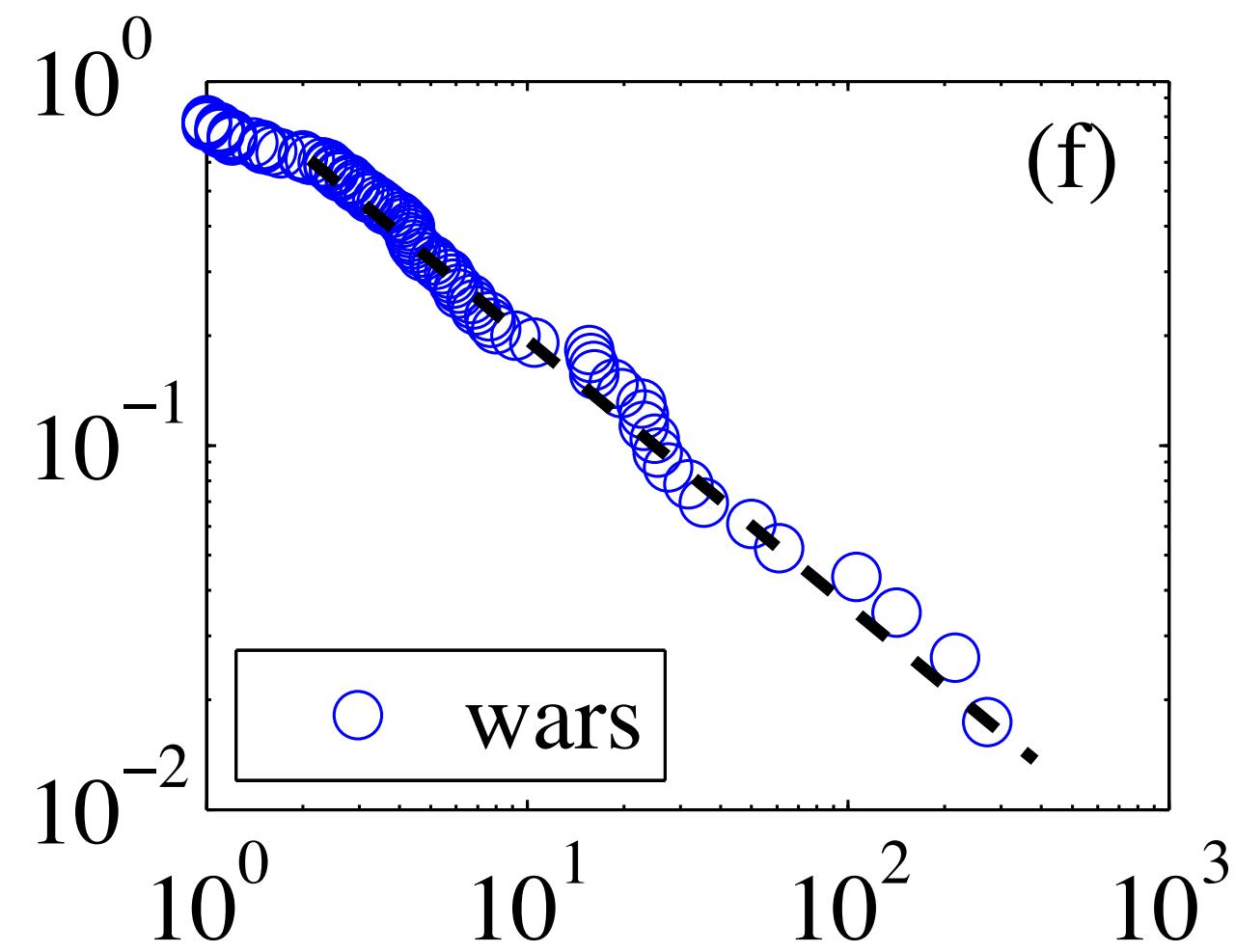
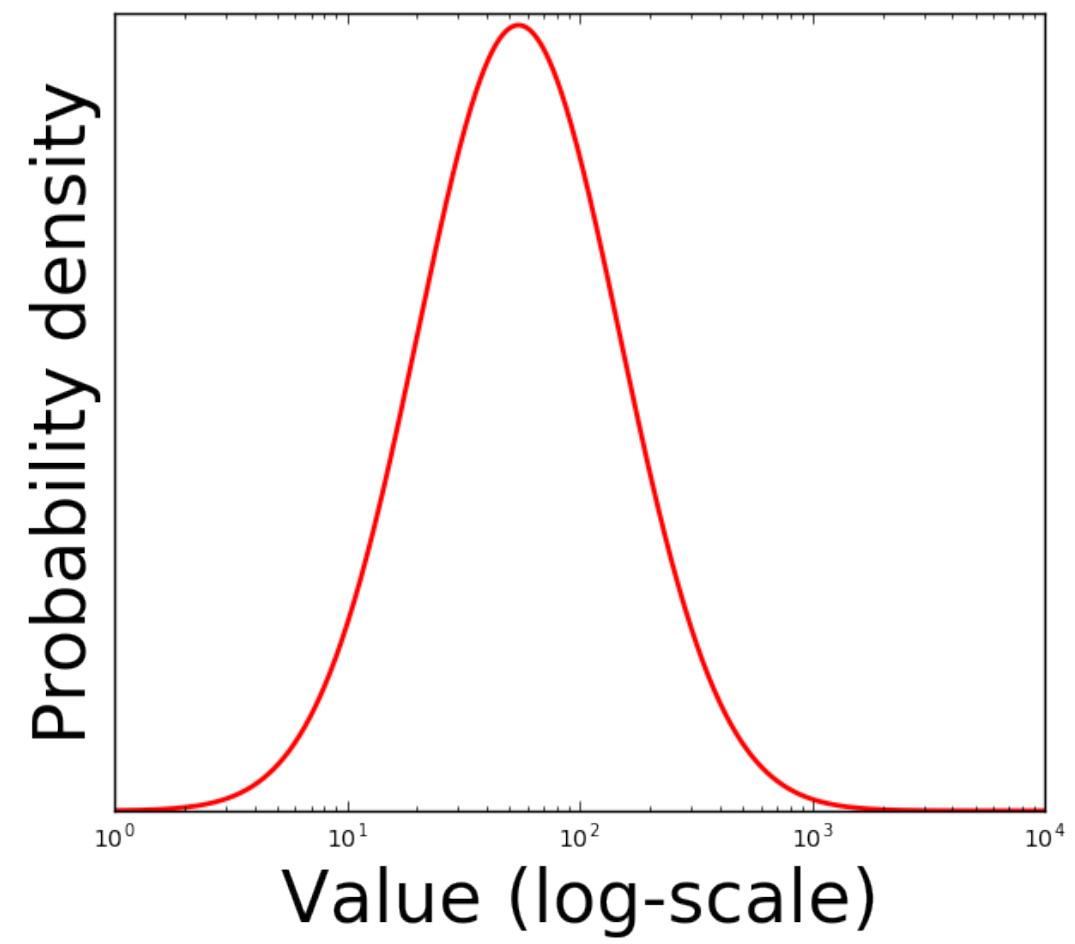


Linear scale:
Can't see any!



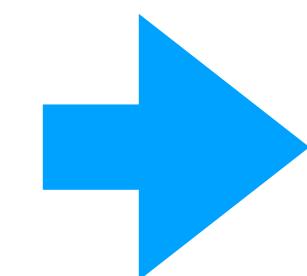
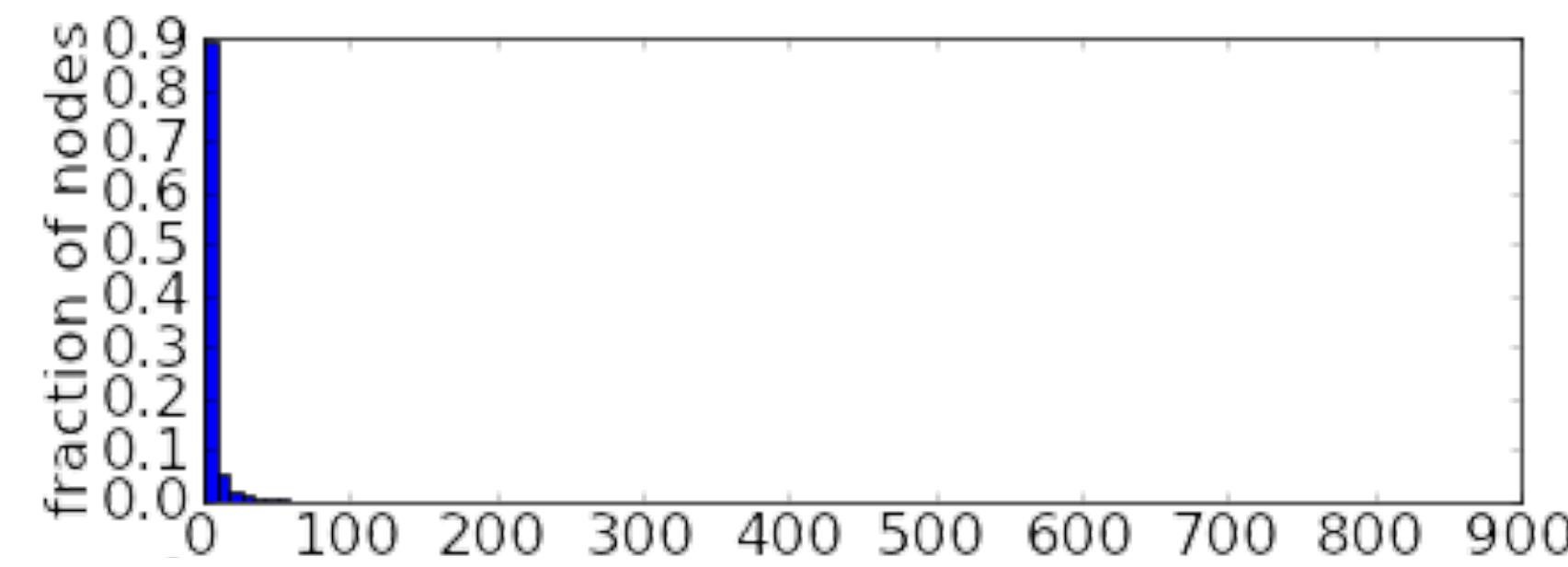
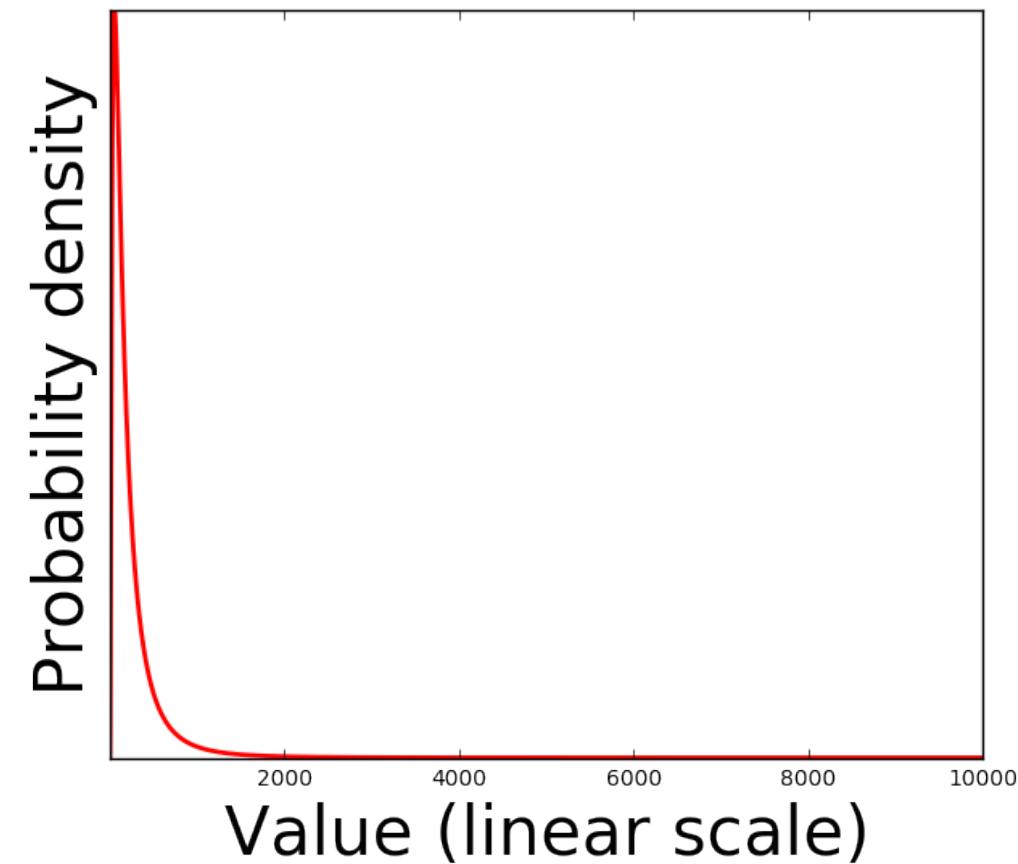


Linear scale:
Can't see any!

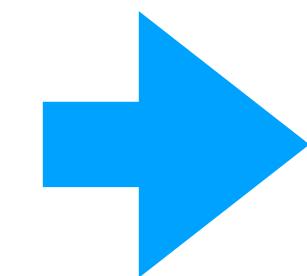
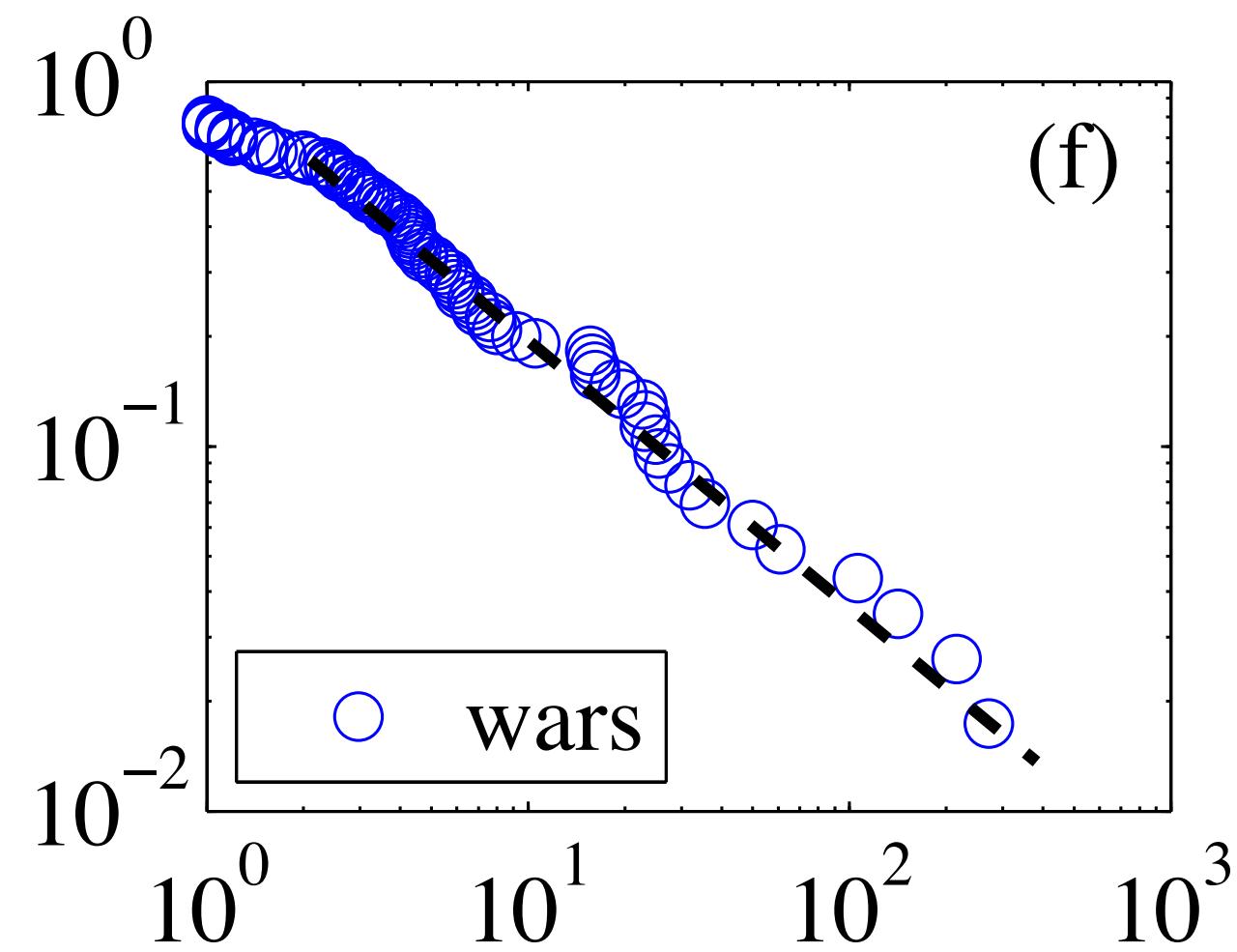
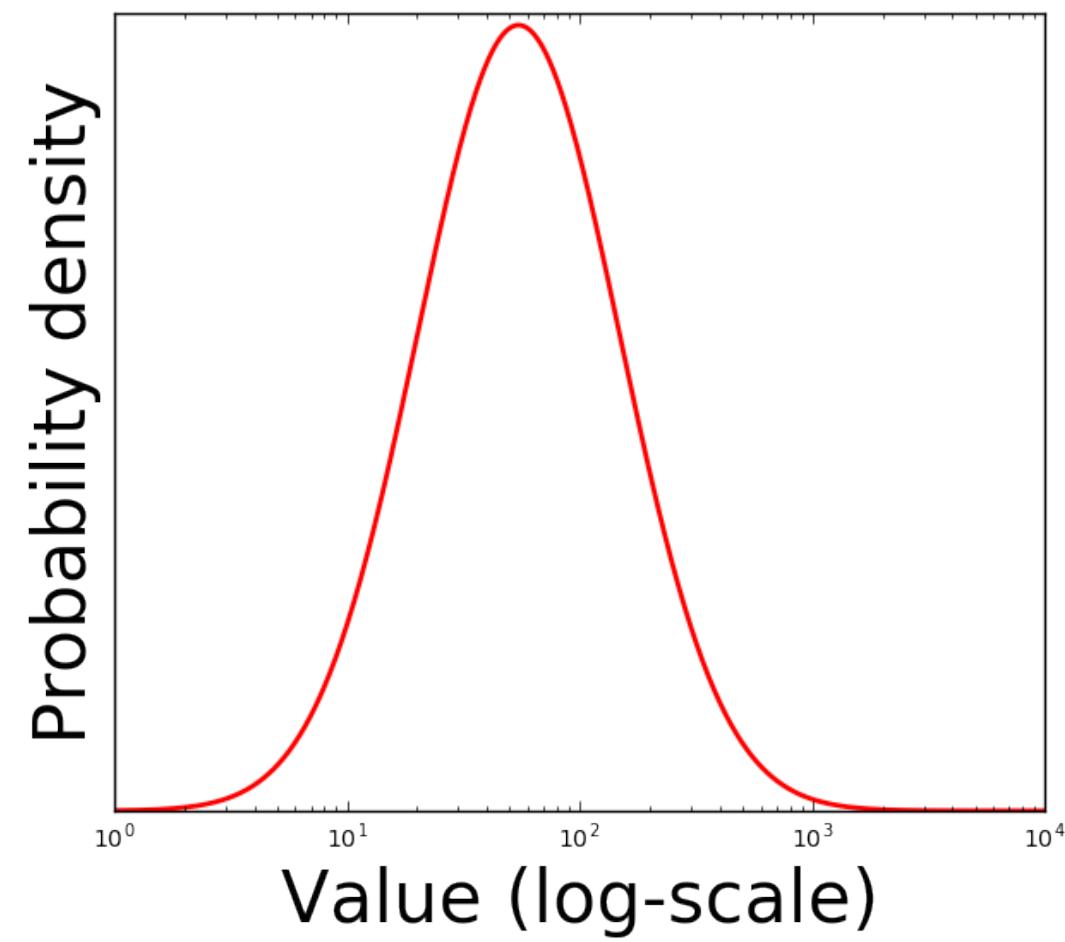


log scale:
Usefull!

Have a fat tail?

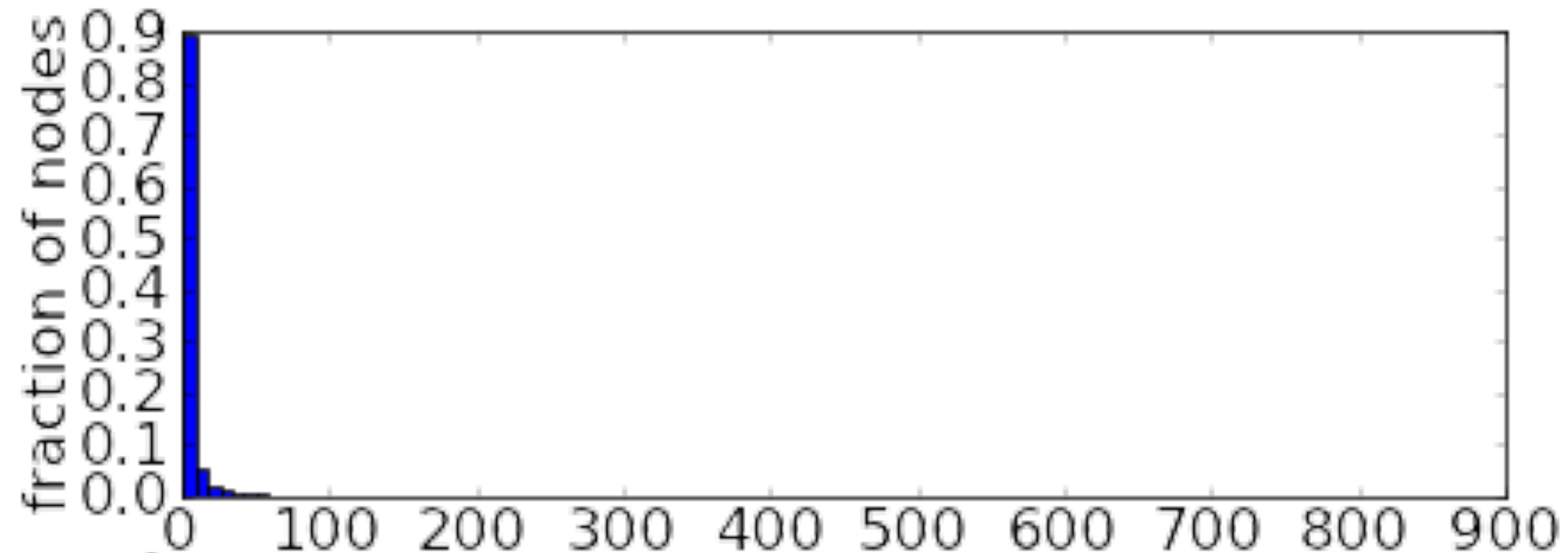


Linear scale:
Can't see any!



log scale:
Usefull!

What would you do?



Histogram in log-scale

How should we create bins
in log-scale?

Constant width in log-scale?

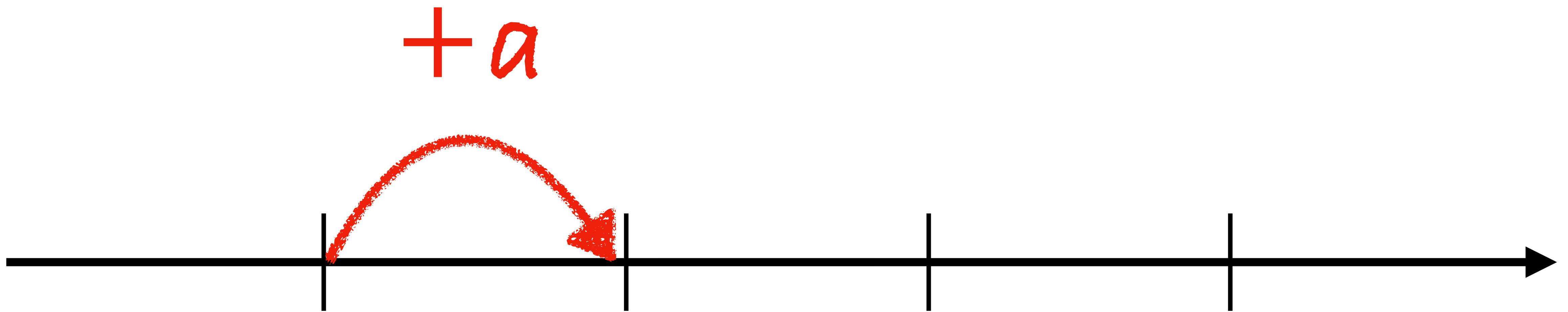
First bin: [1, 3)

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

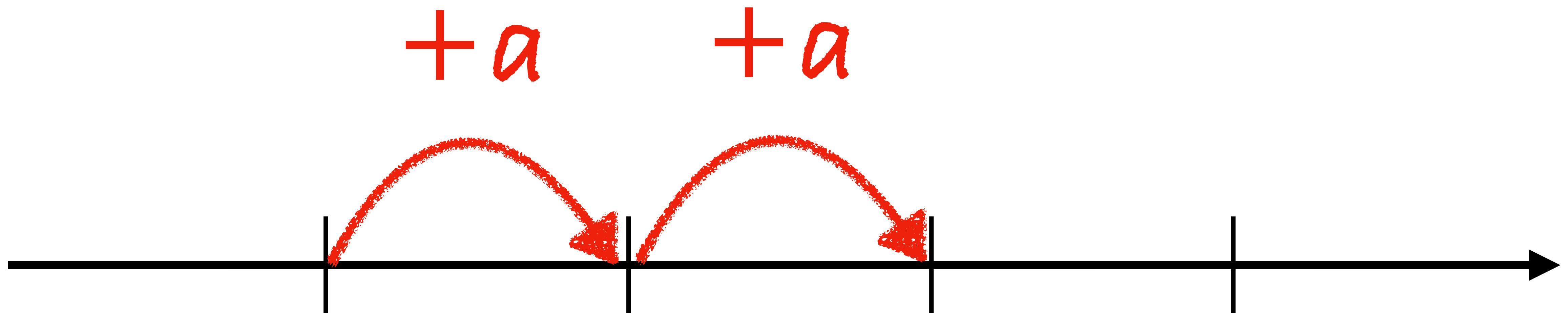
In linear scale



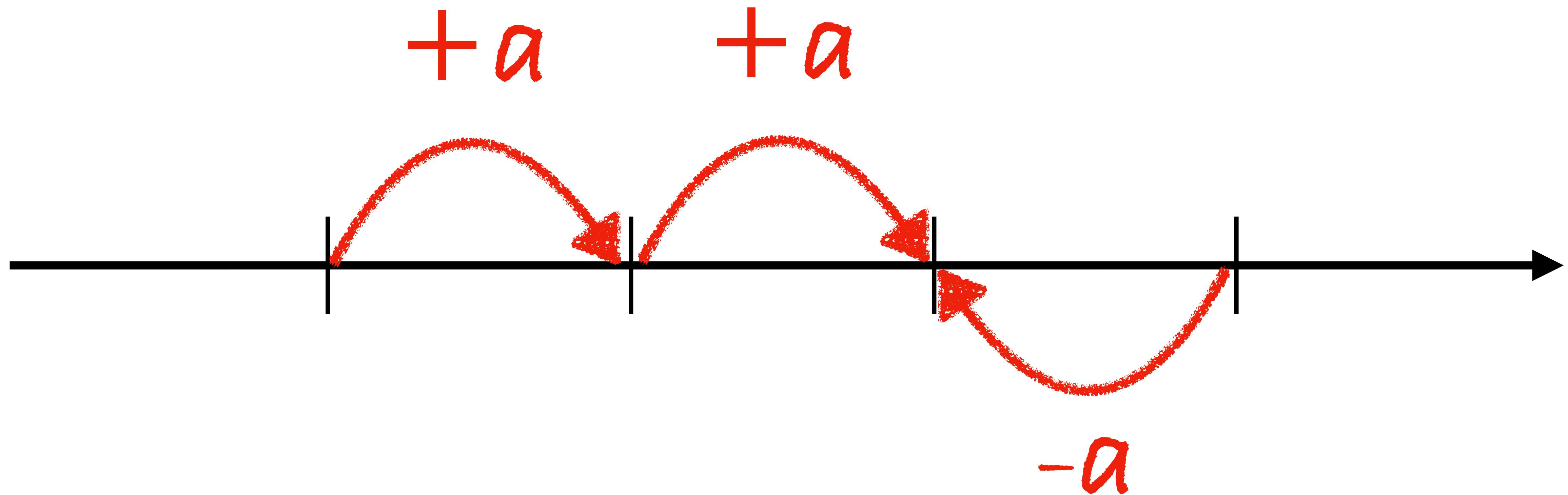
In linear scale



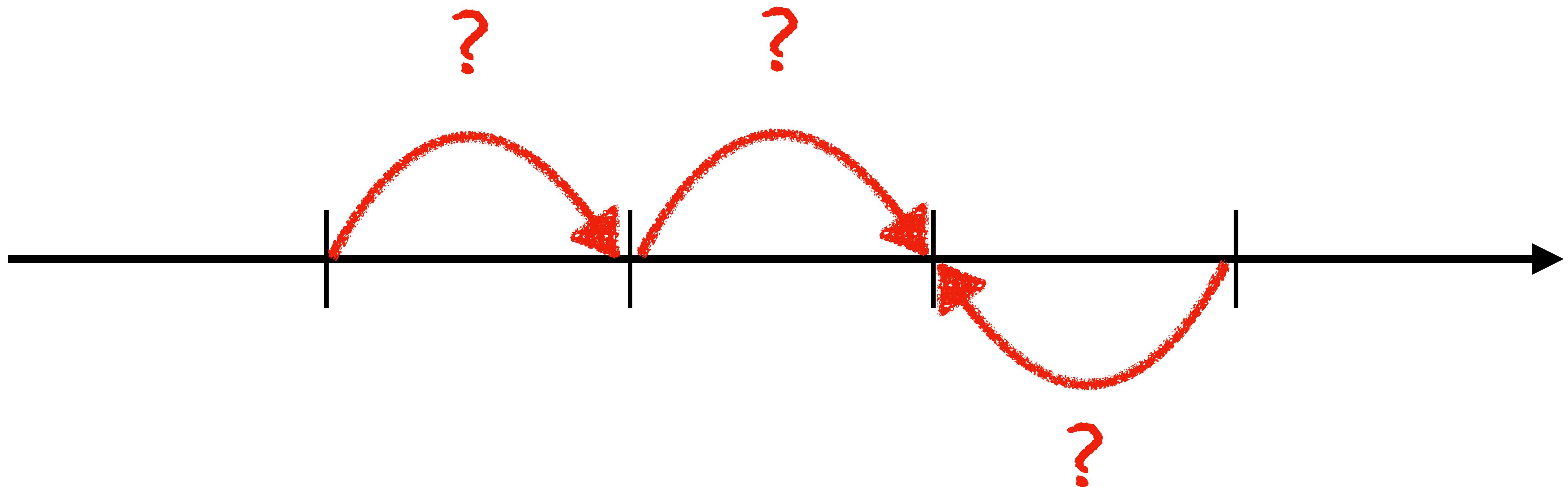
In linear scale



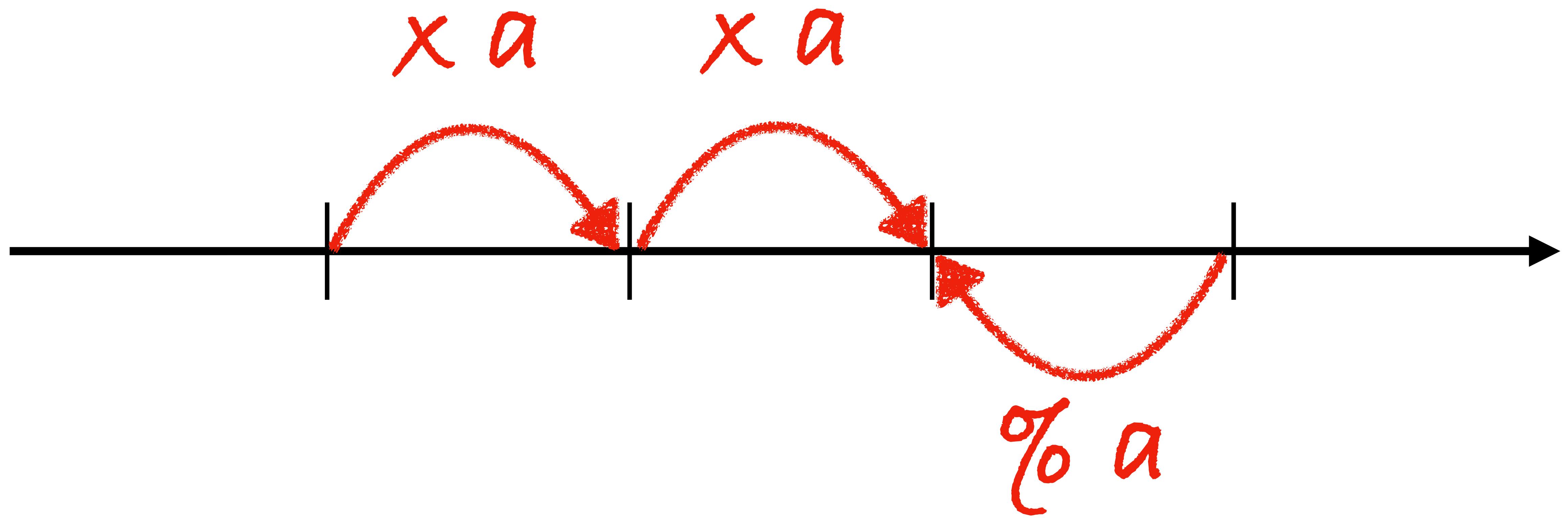
In linear scale



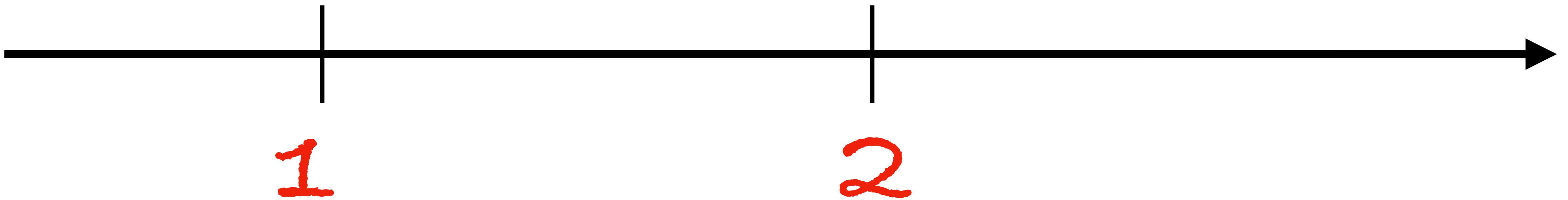
In log scale?



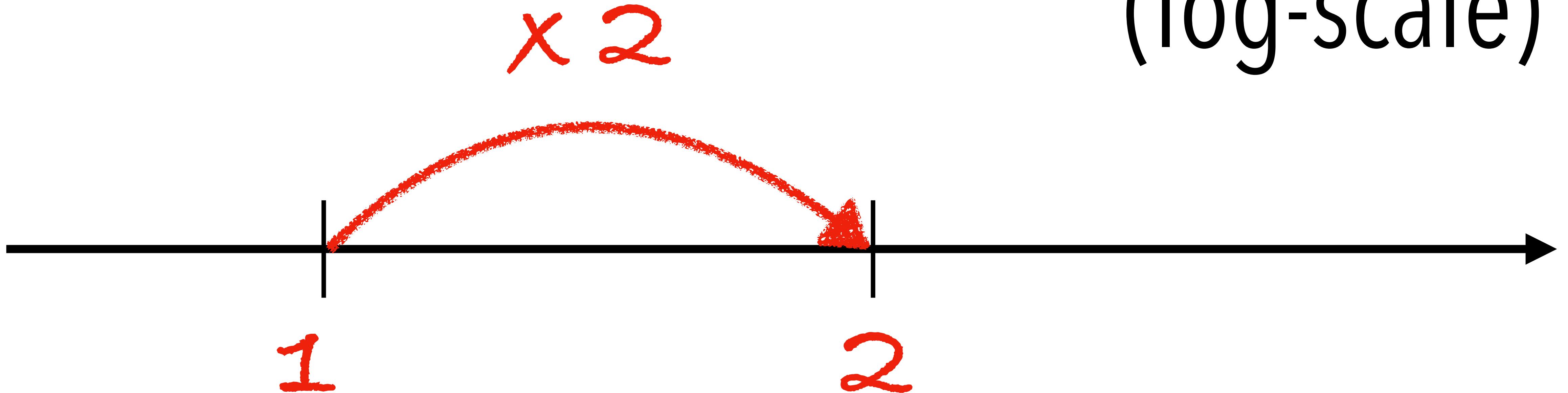
In log scale?



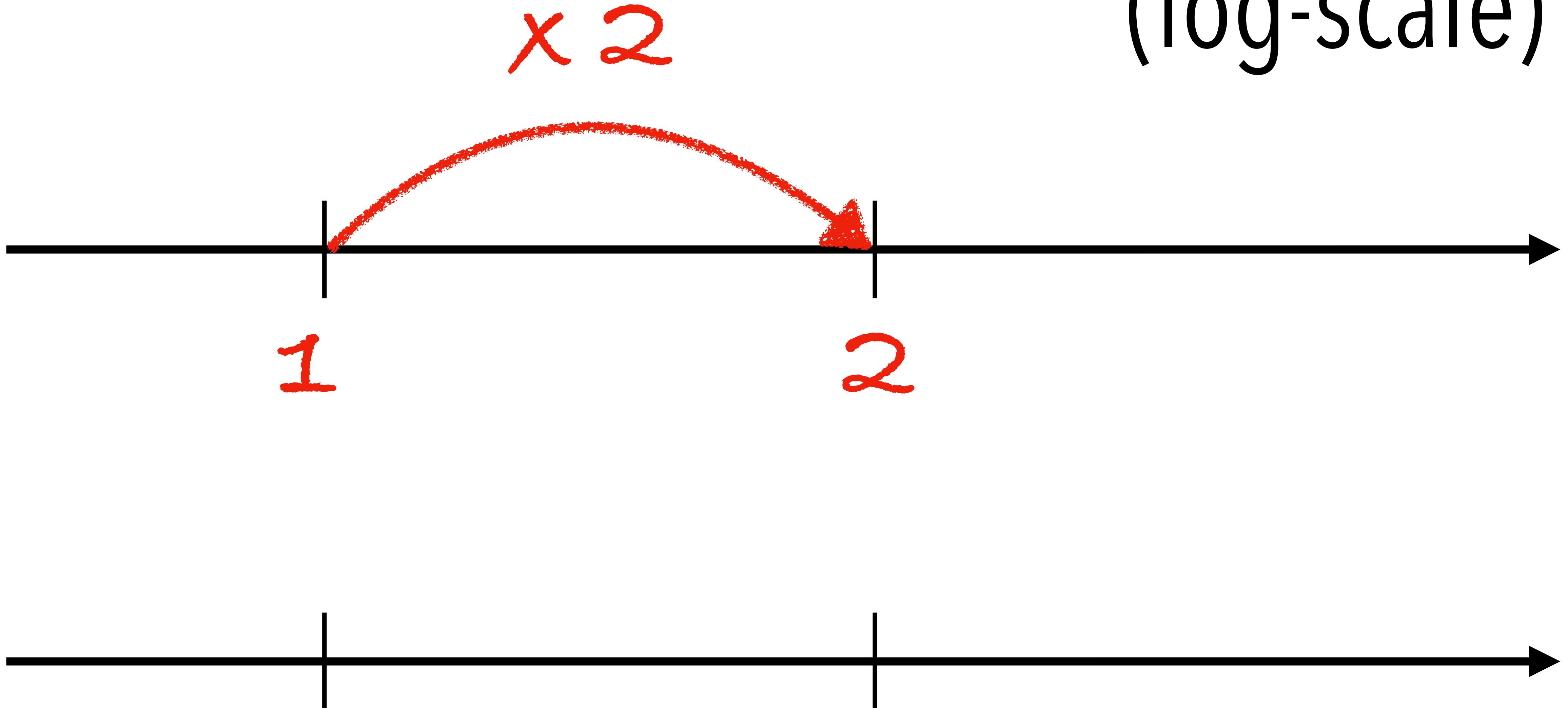
(log-scale)



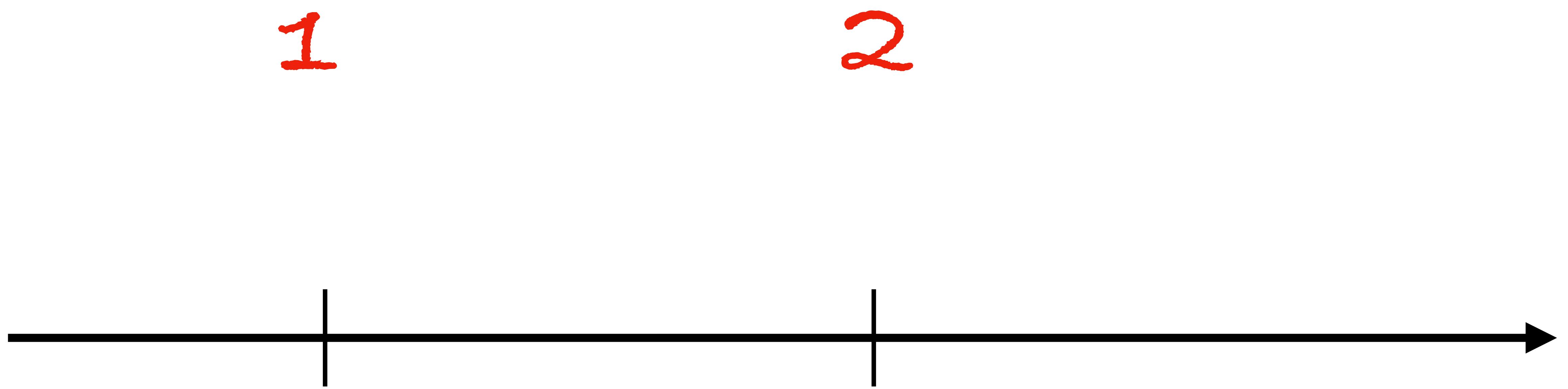
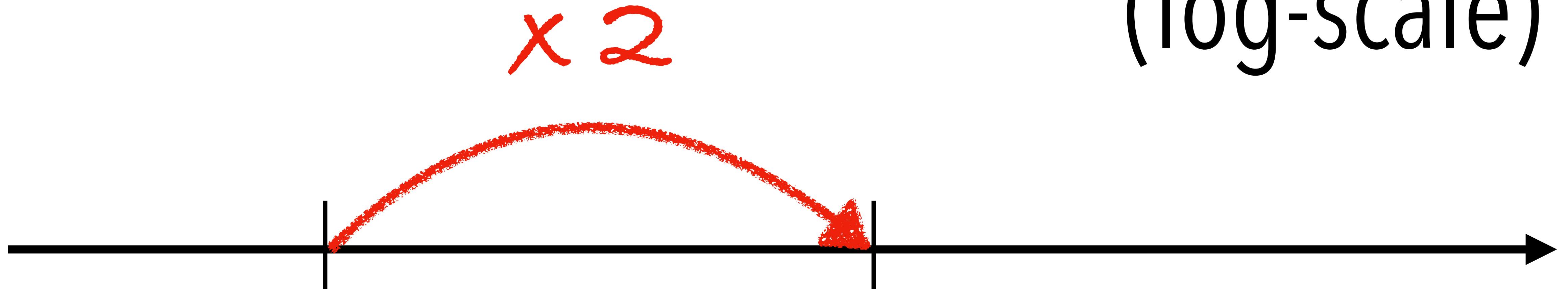
(log-scale)



(log-scale)

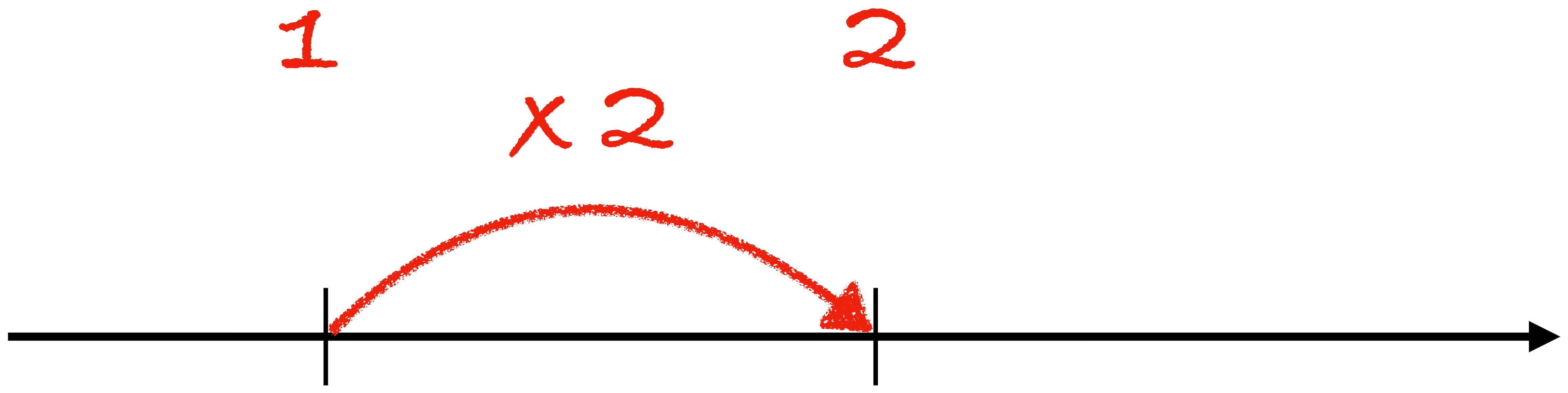


(log-scale)



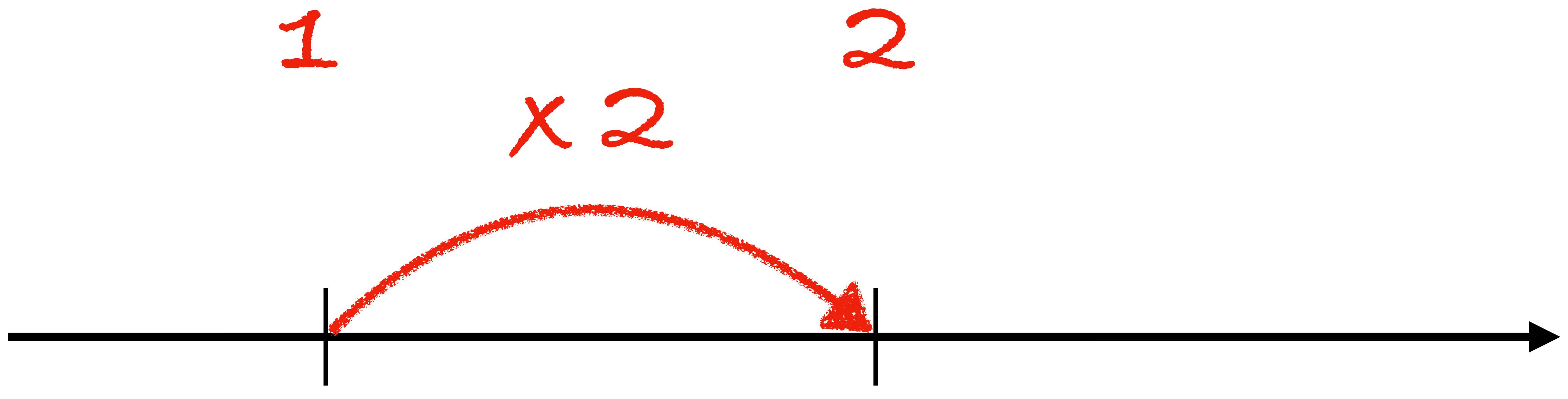
1,000,000

(log-scale)



1,000,000

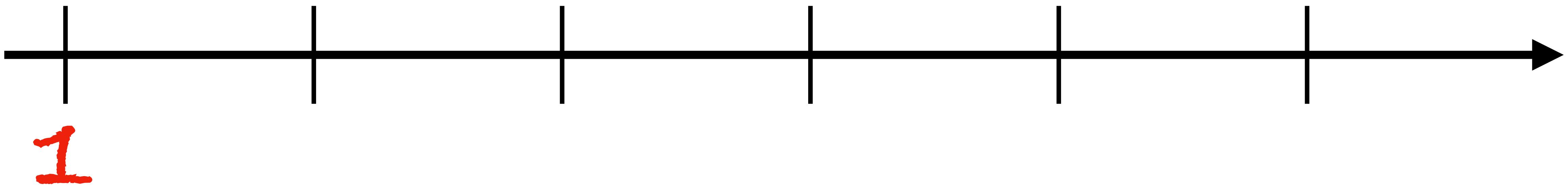
(log-scale)



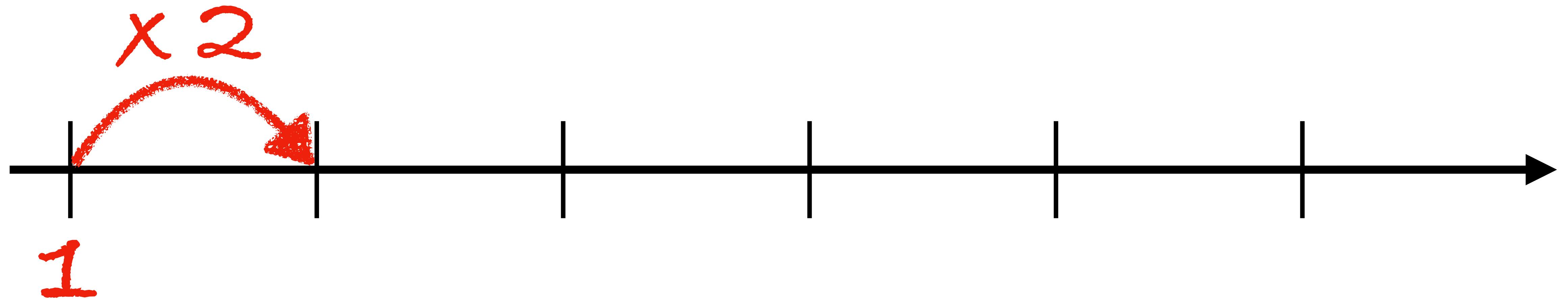
1,000,000

2,000,000

(log-scale)

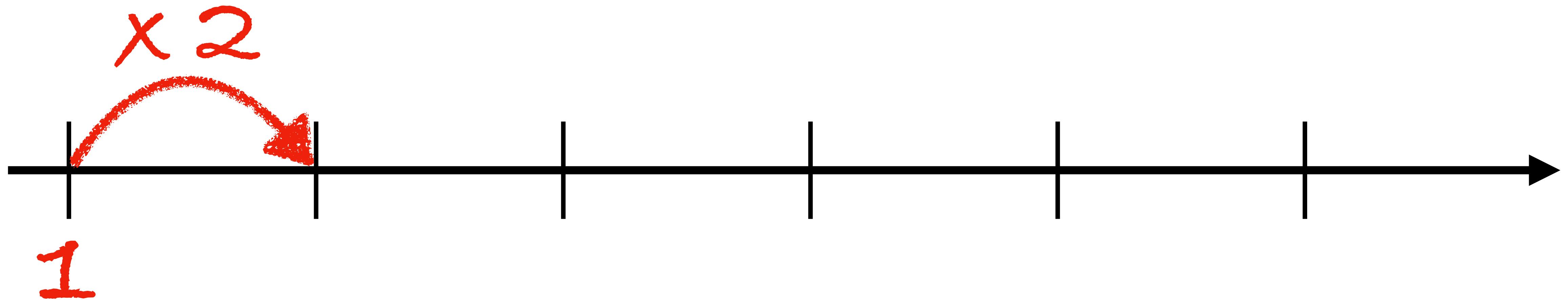


(log-scale)



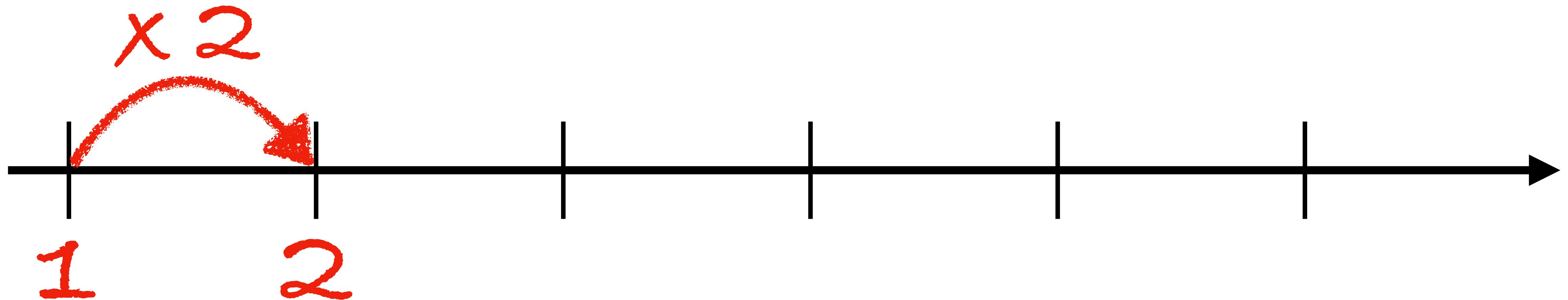
(You can choose
other widths too)

(log-scale)

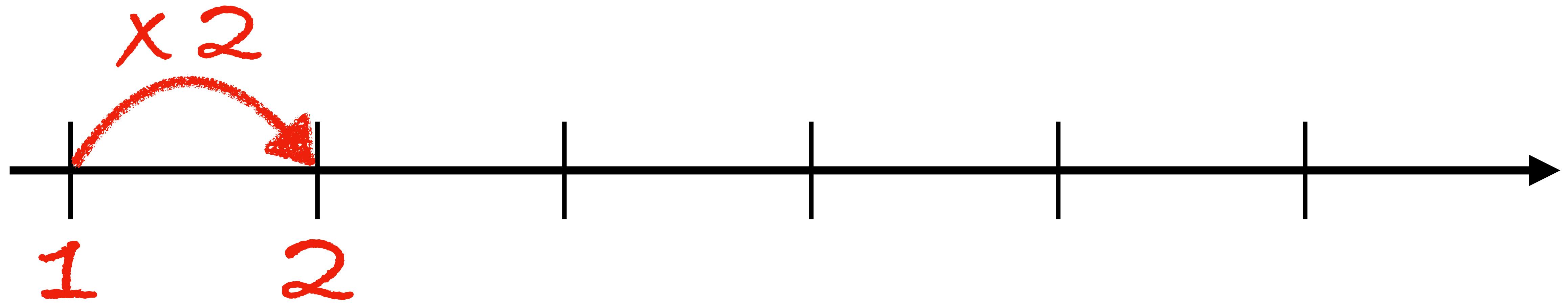


(You can choose
other widths too)

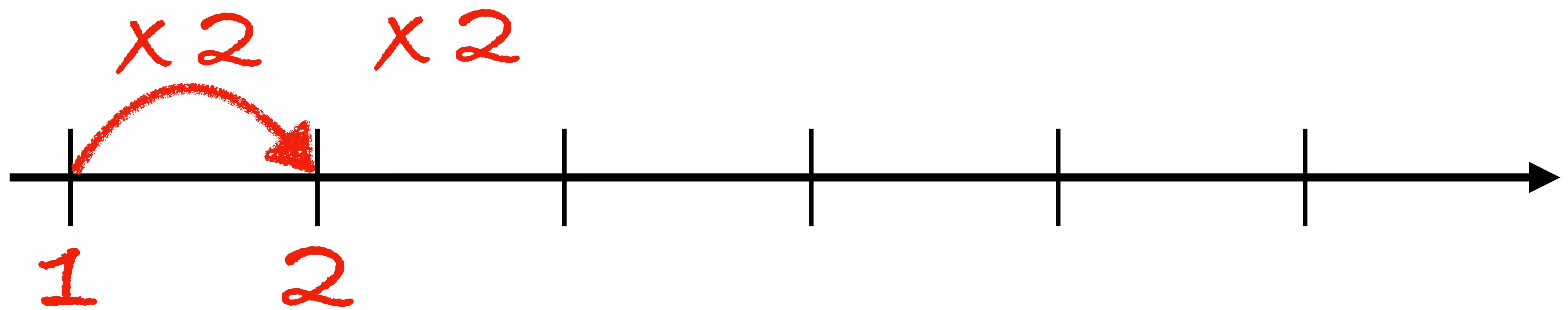
(log-scale)



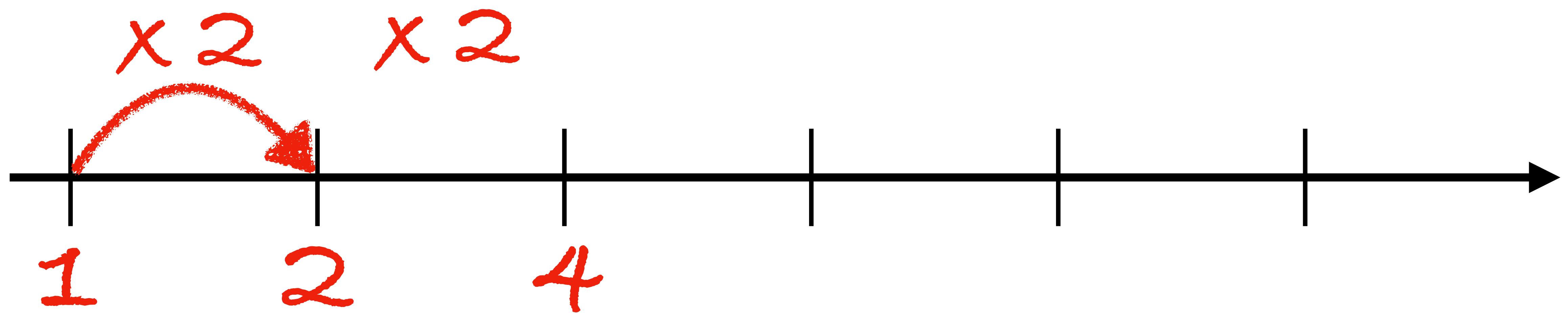
(log-scale)



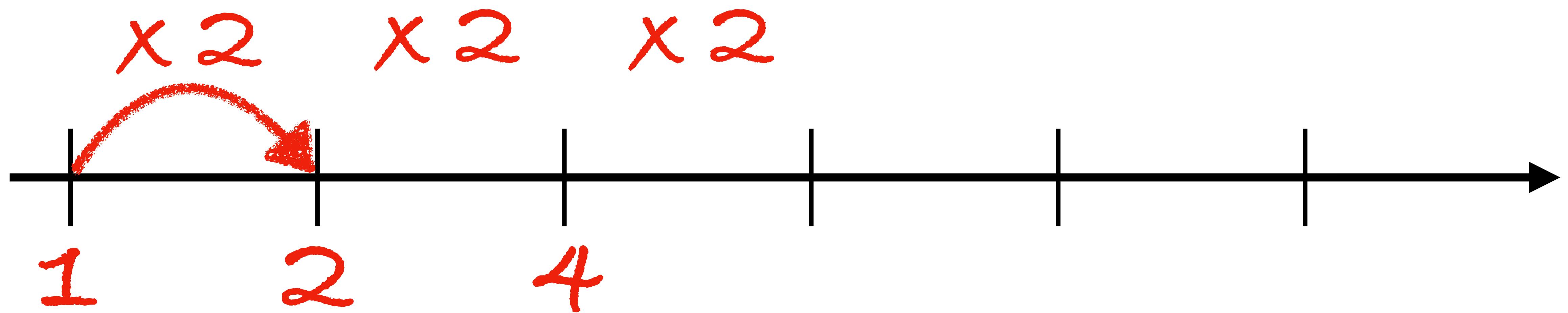
(log-scale)



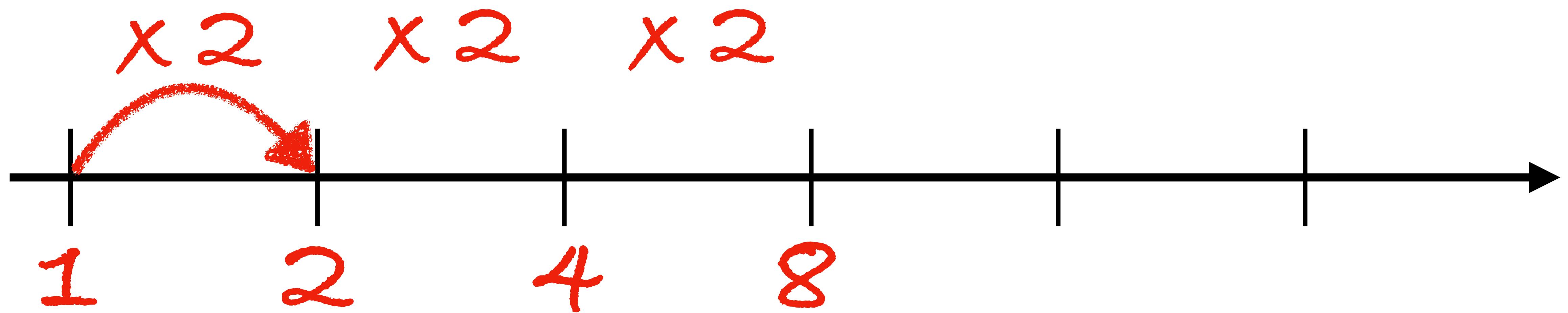
(log-scale)



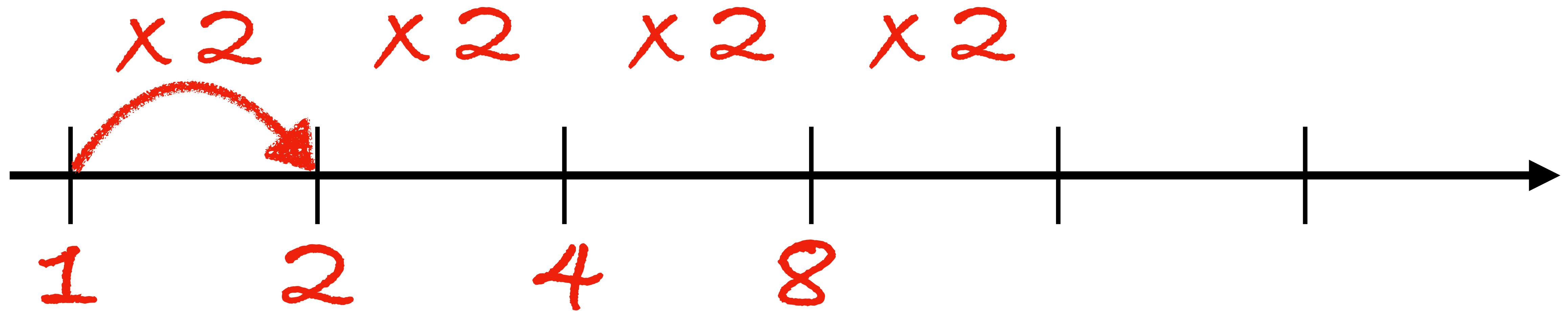
(log-scale)



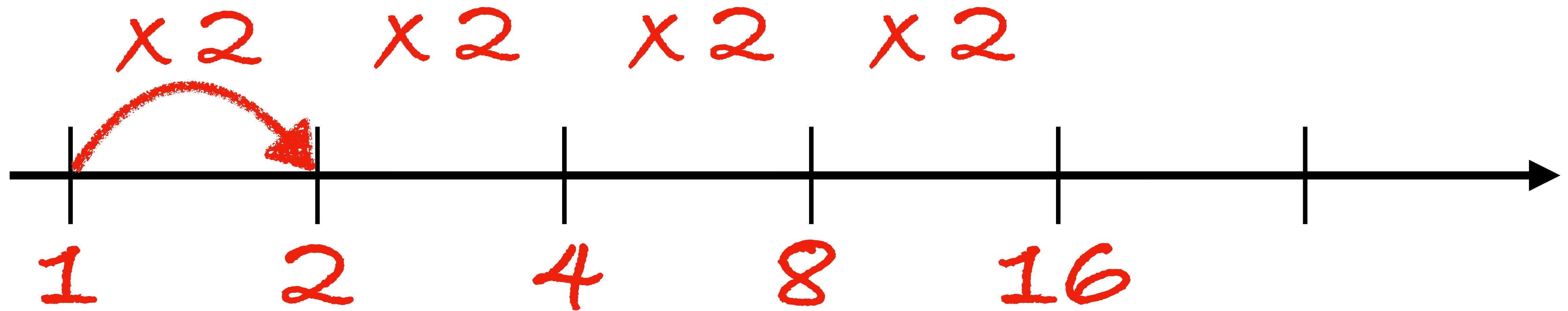
(log-scale)



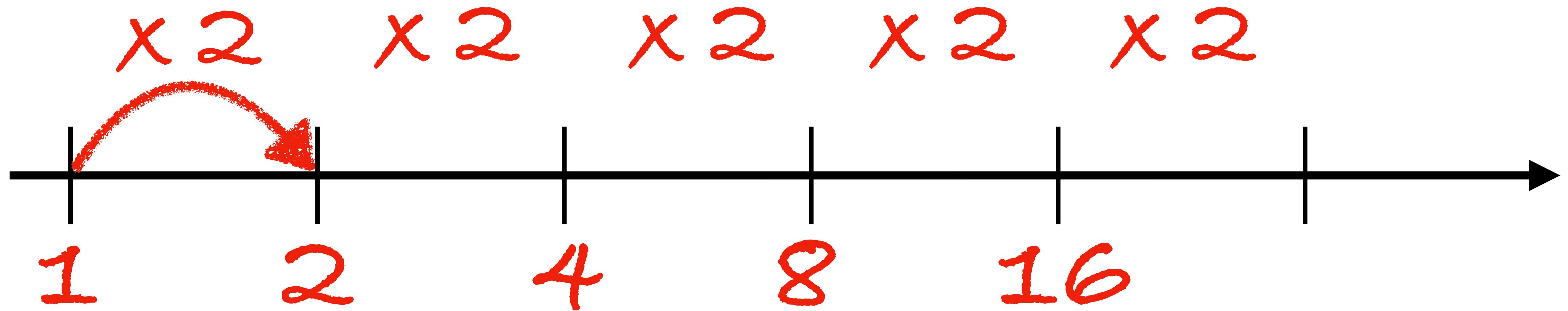
(log-scale)



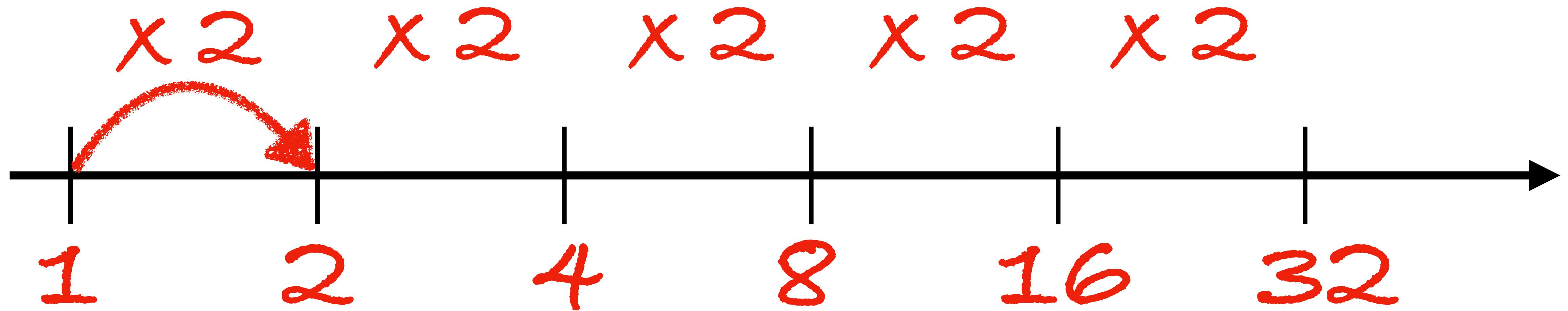
(log-scale)



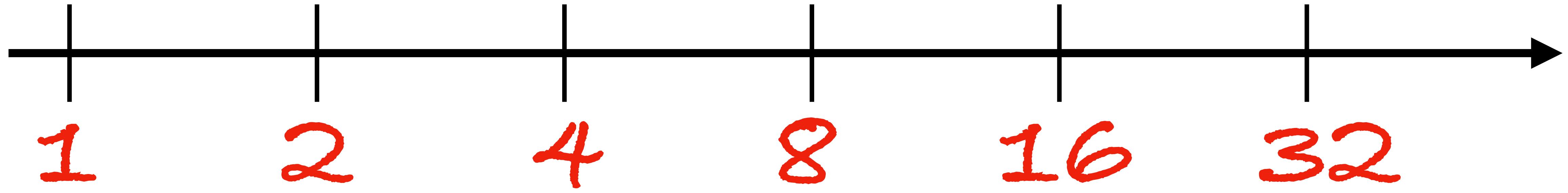
(log-scale)



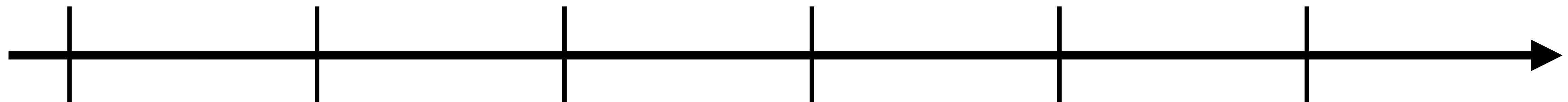
(log-scale)



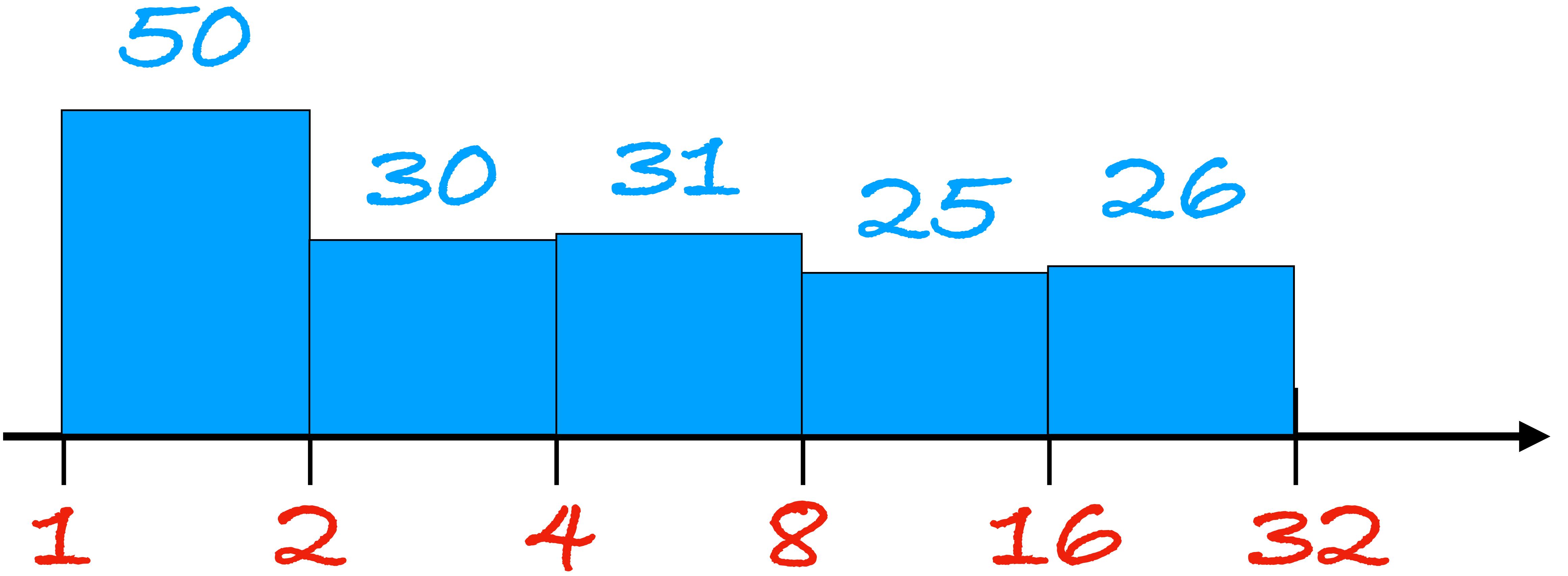
“Log-binning”
(logarithmic binning)

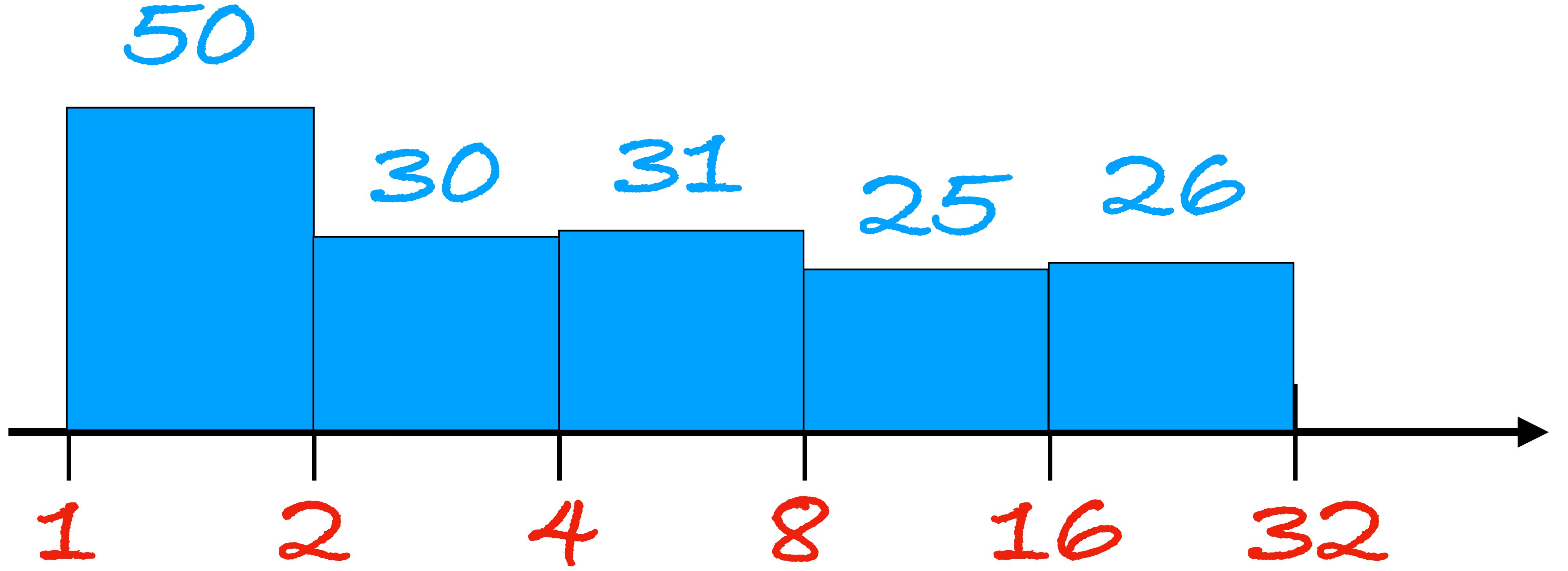


50 30 31 25 26



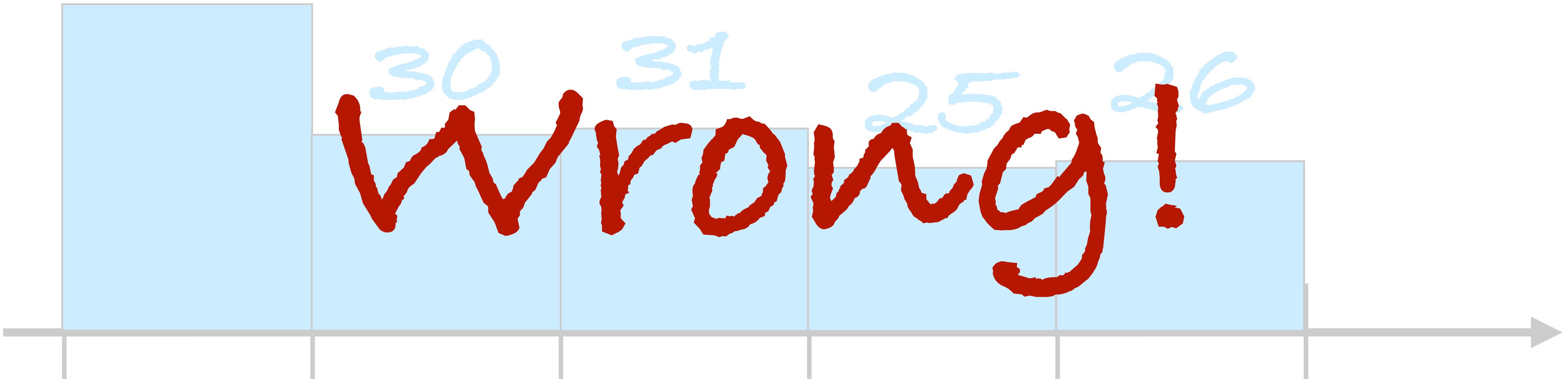
1 2 4 8 16 32





50

30 31
Wrong!
25 36



1

2

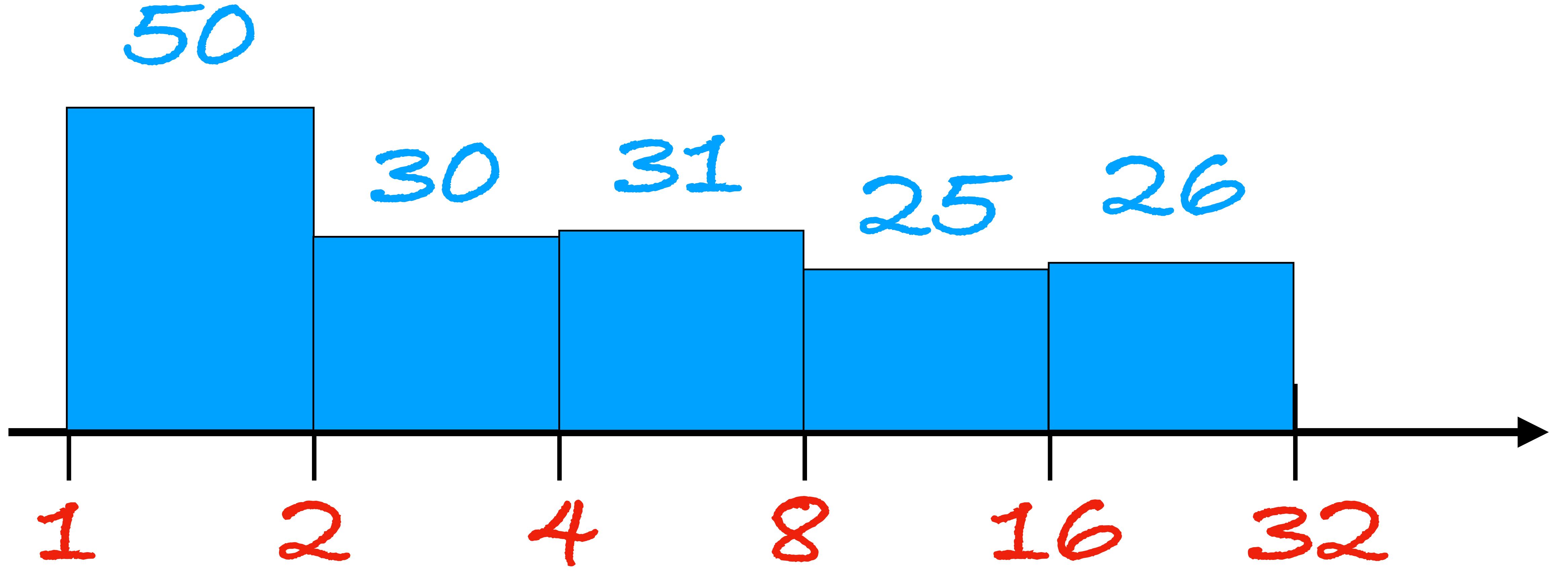
4

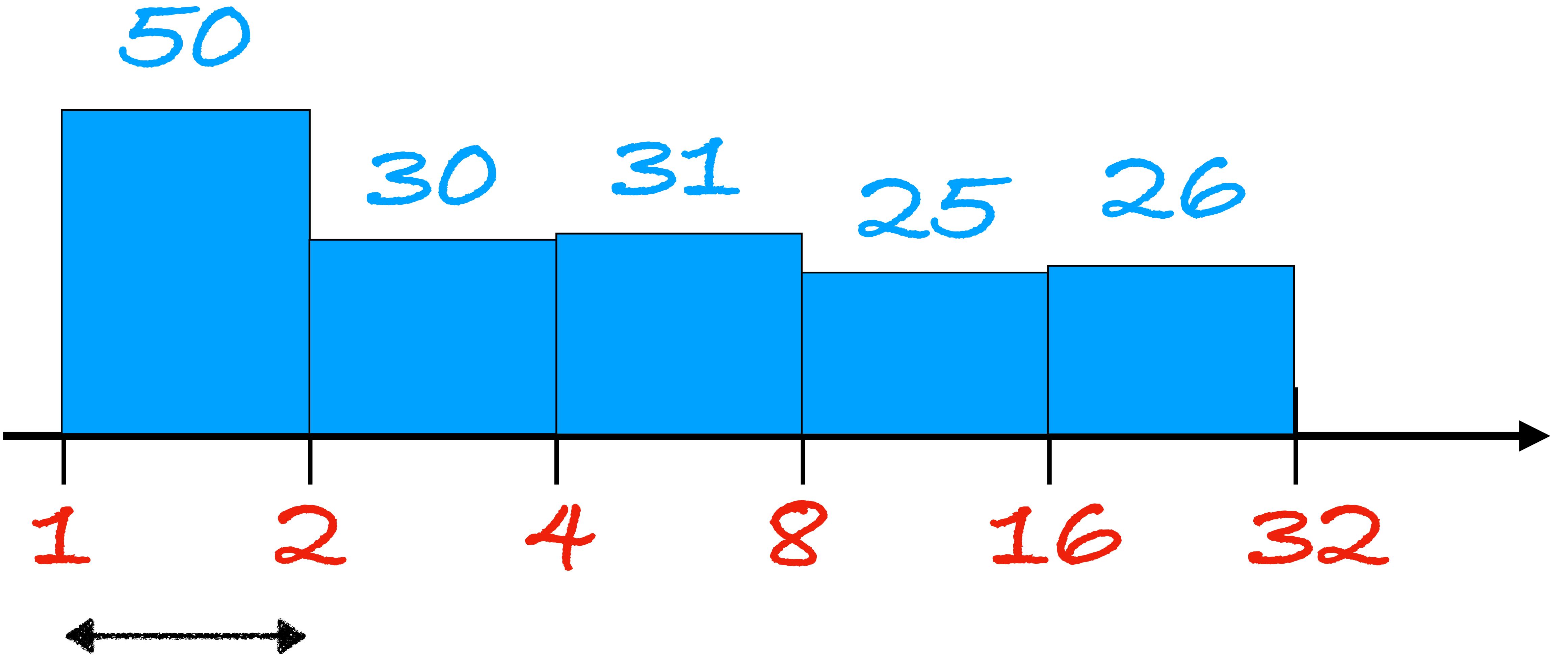
8

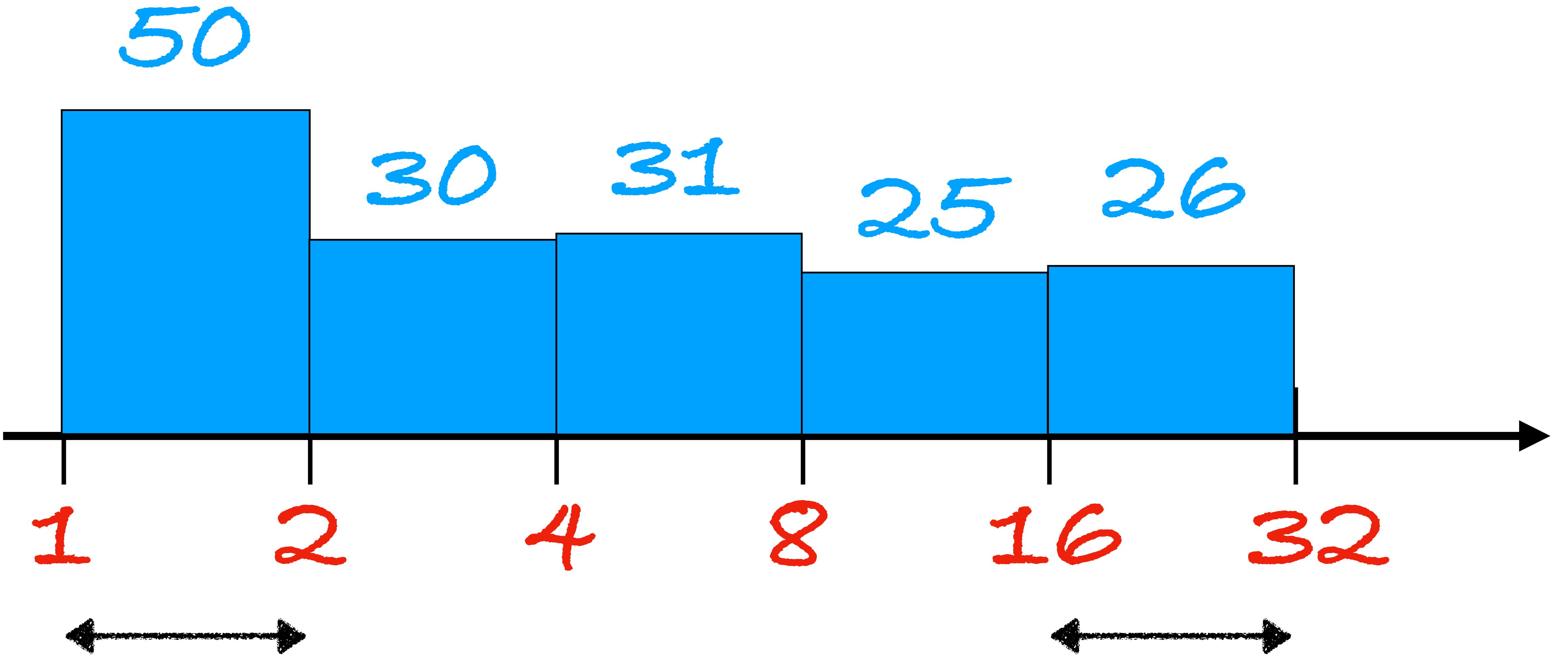
16

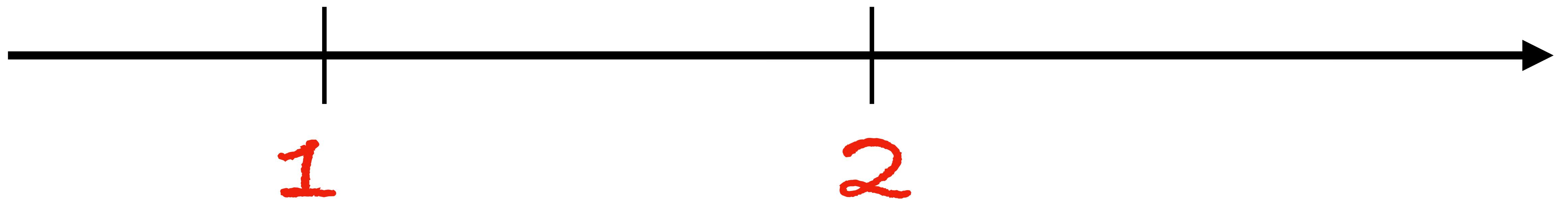
32

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

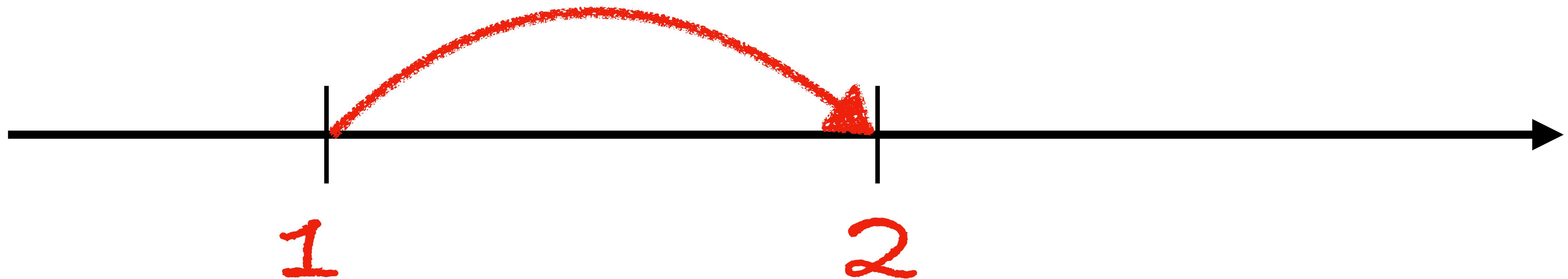




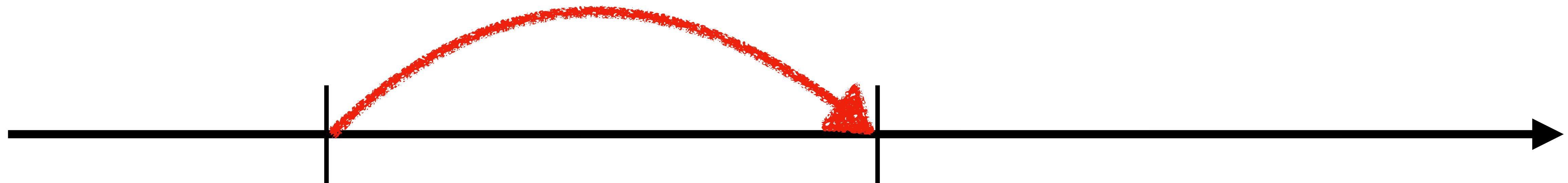




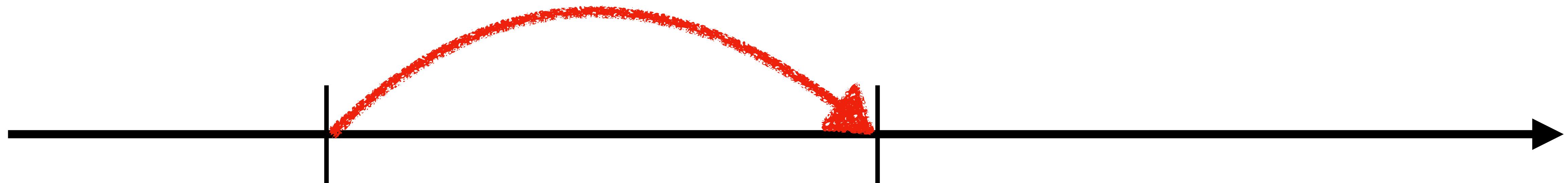
$x_2, +1$



$x_2, +1$



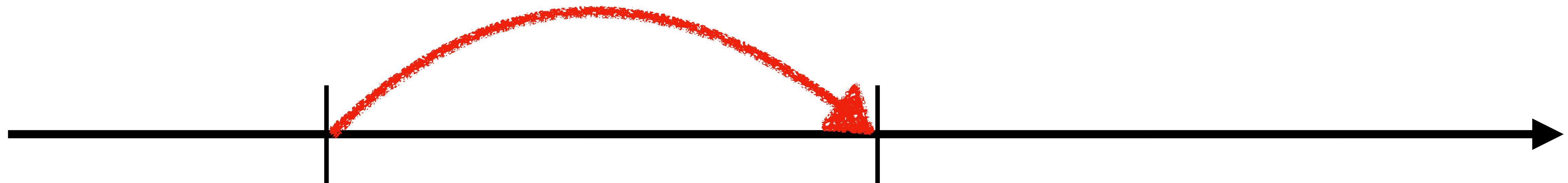
$x_2, +1$



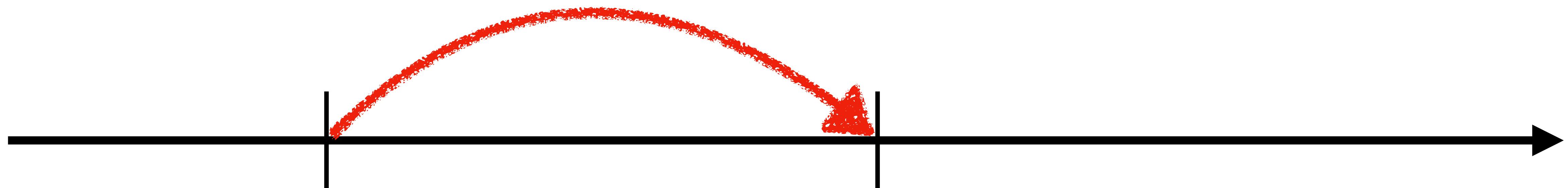
1,000,000,000



$x_2, +1$

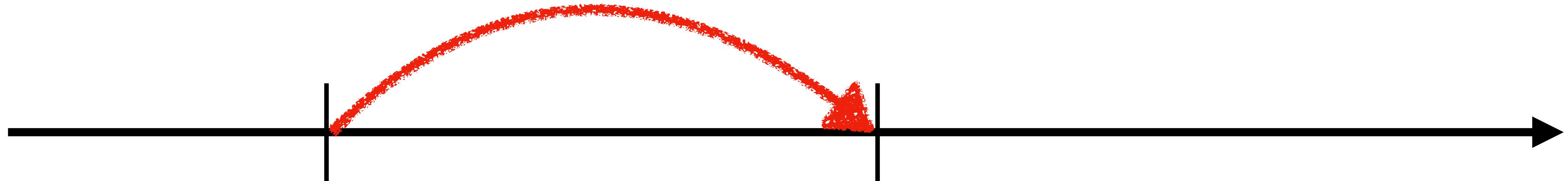


$x_2, +1,000,000,000$

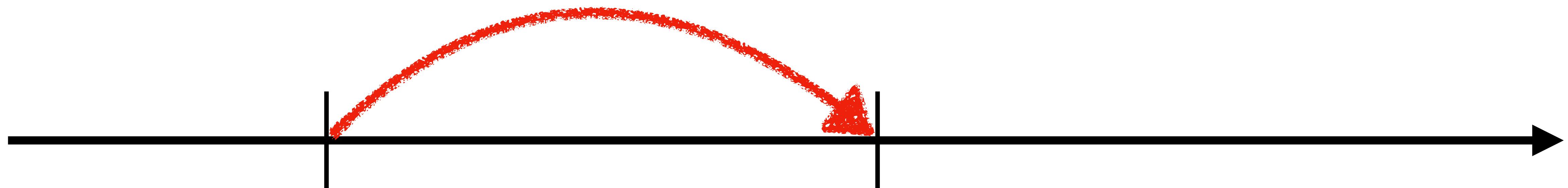


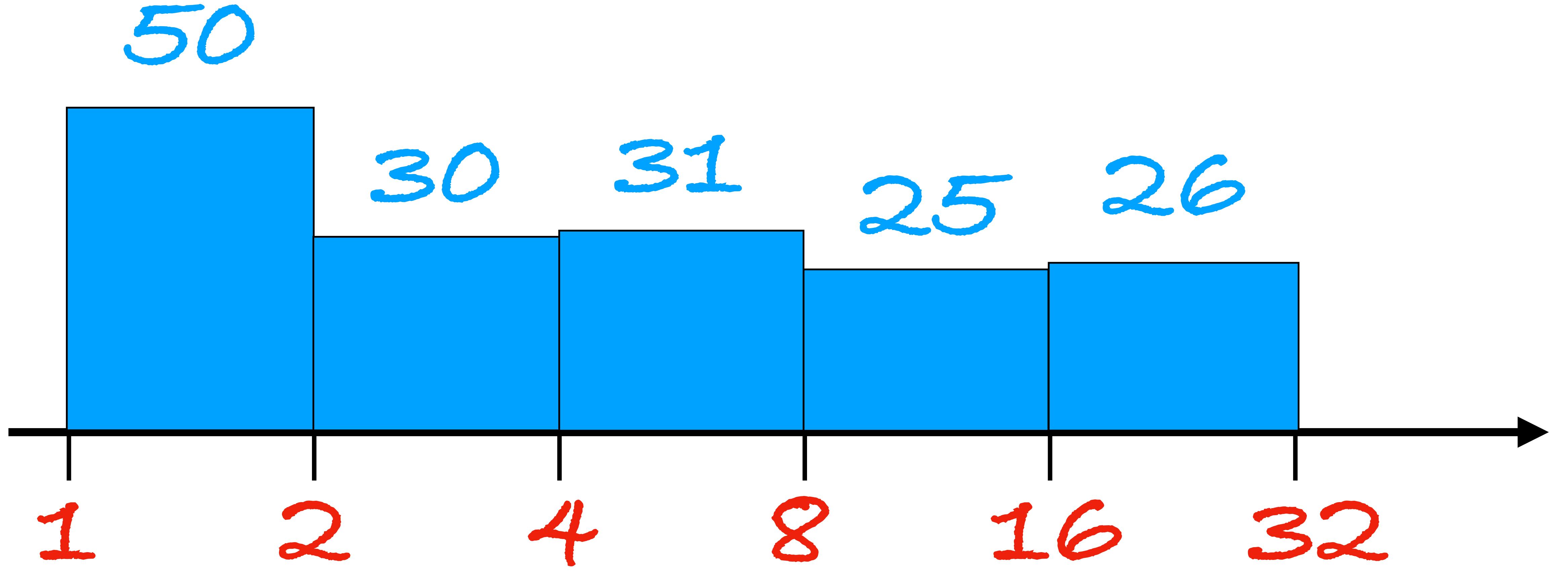
1,000,000,000

$x_2, +1$

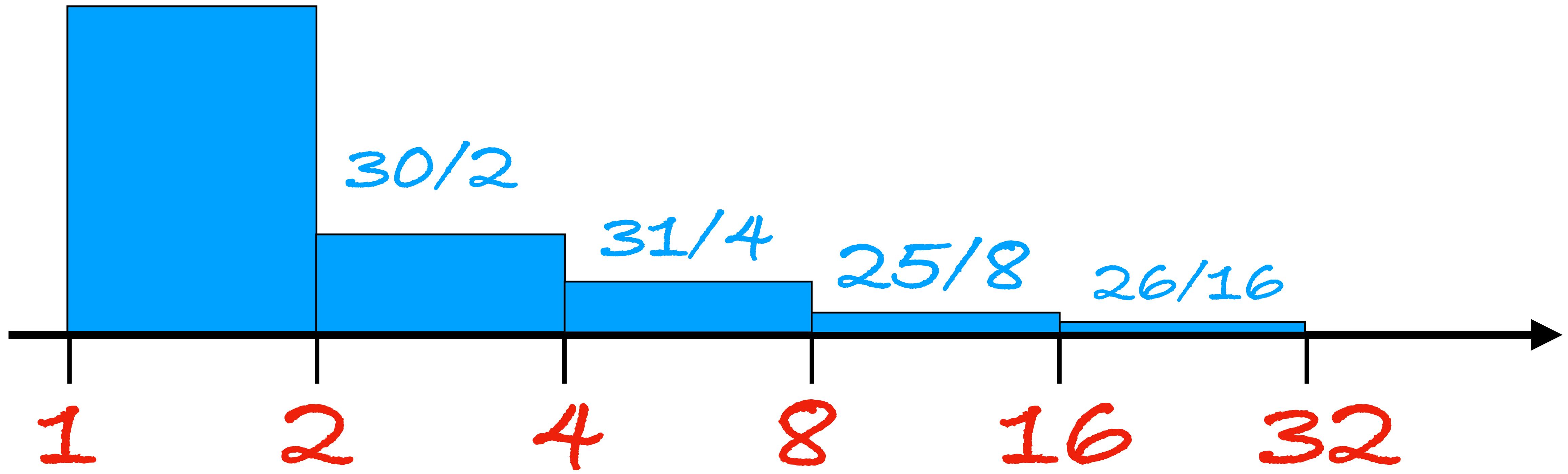


$x_2, +1,000,000,000$

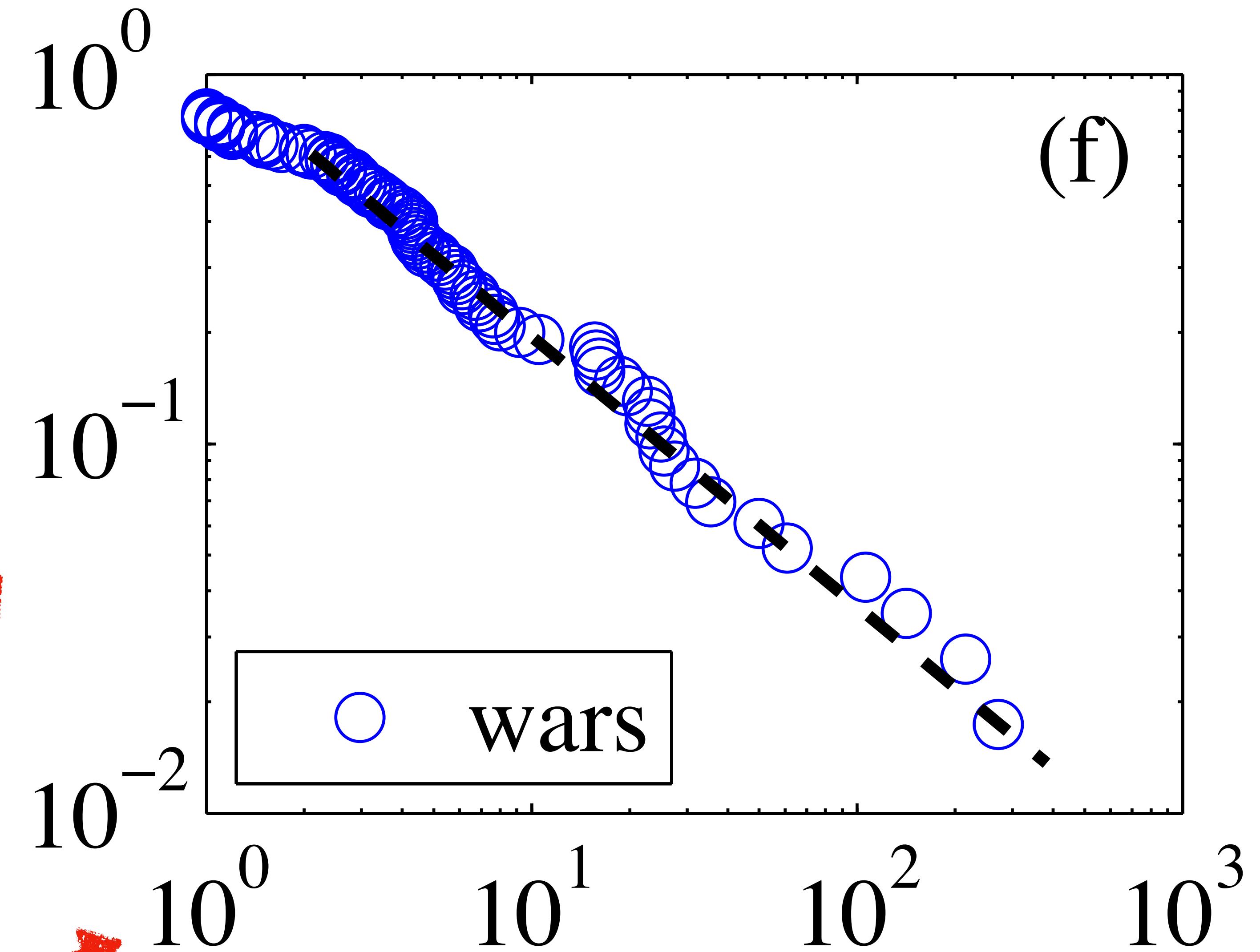




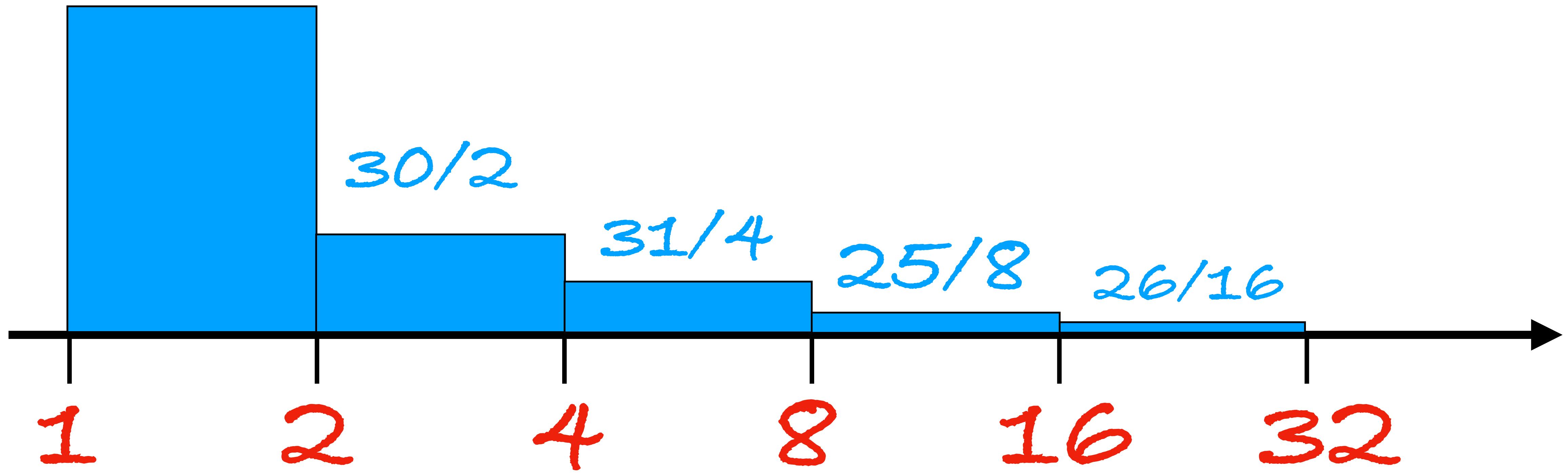
50



Both in
log-scale

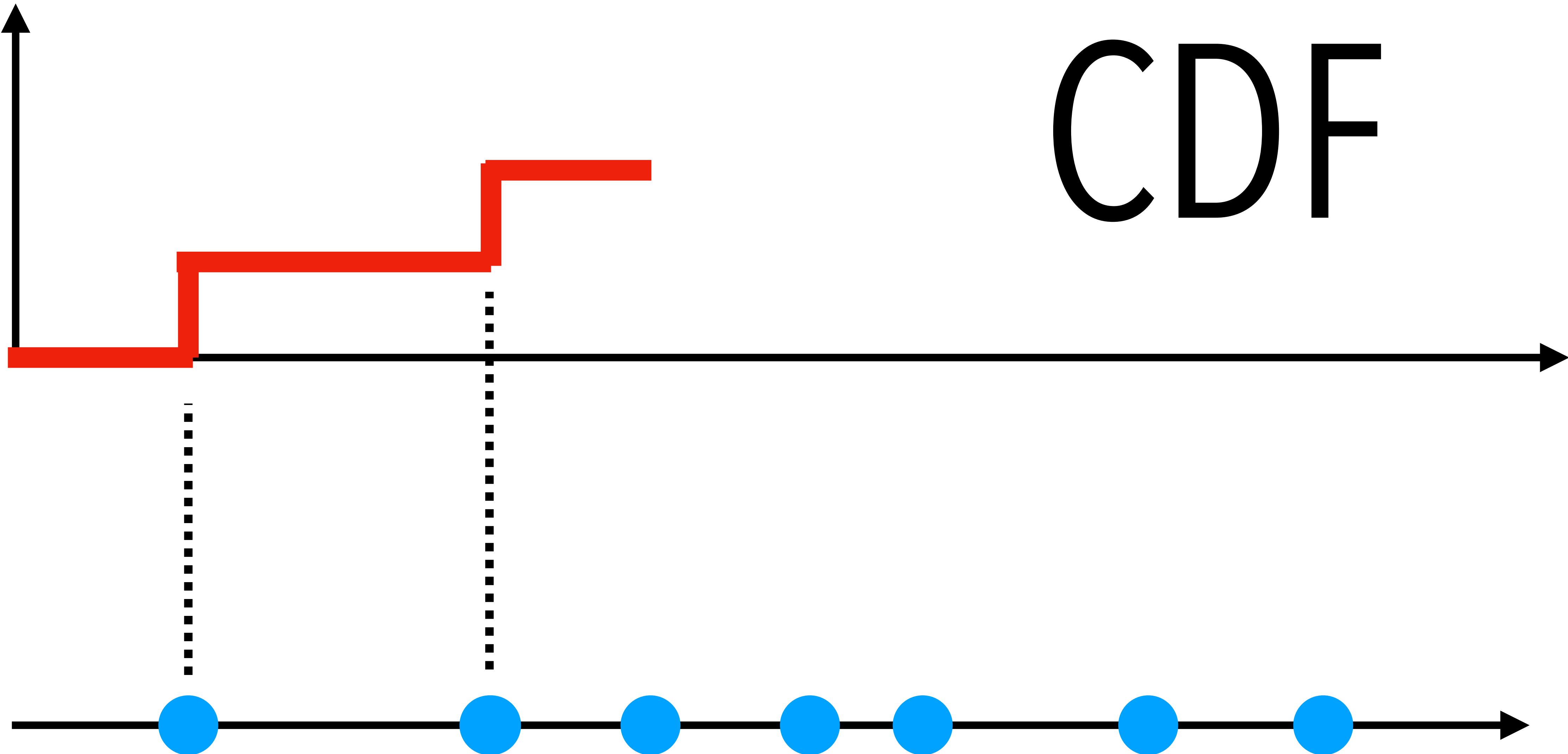


50



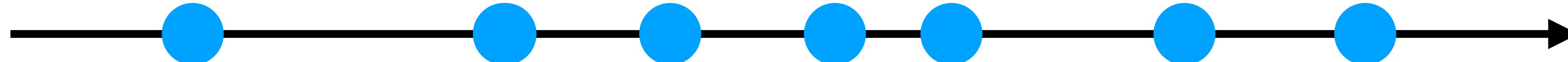
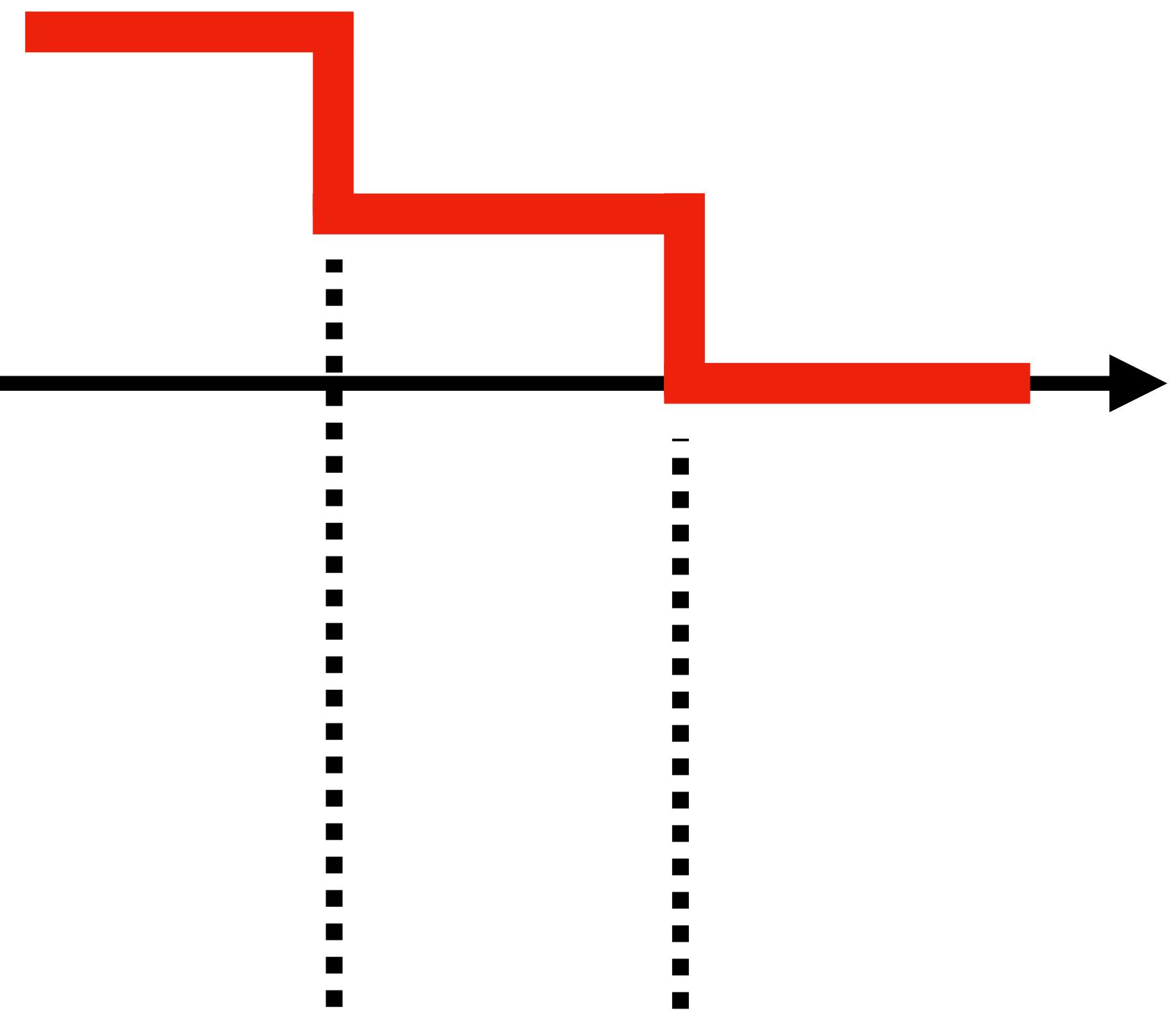
A nice alternative

CDF

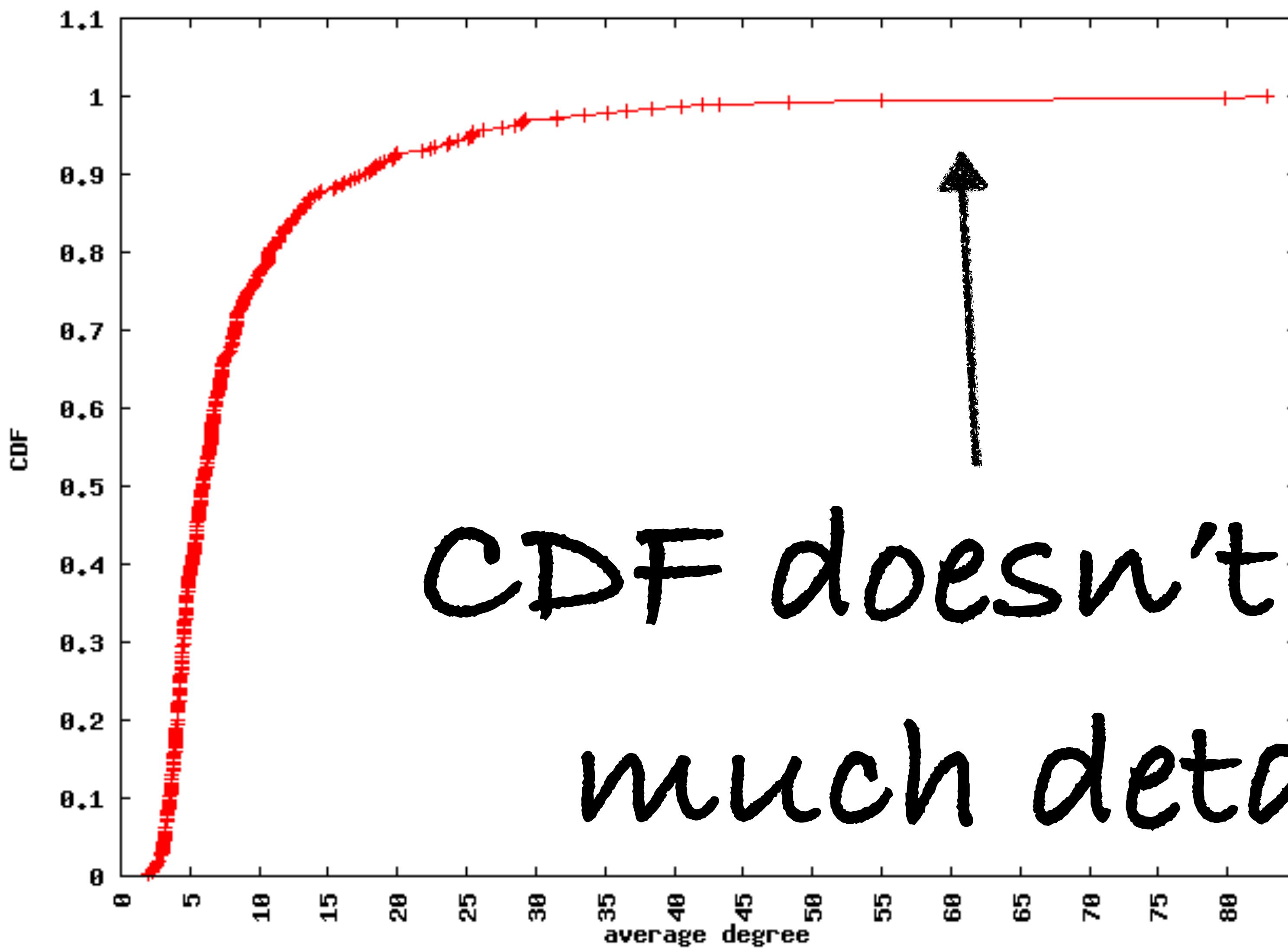


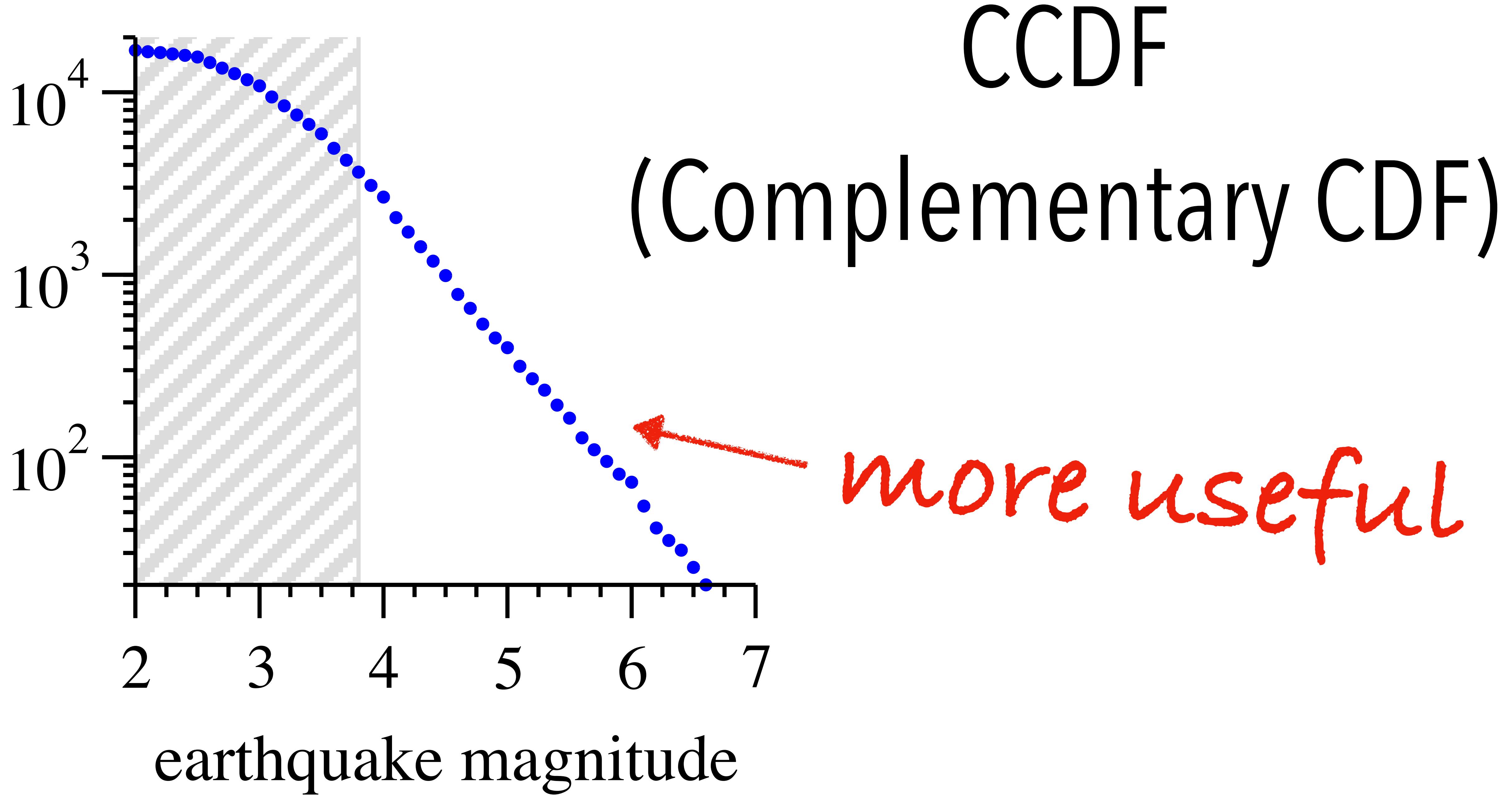


CCDF



Distribution of average degree for router-level graphs





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

CCDF is often much more useful than **CDF**.

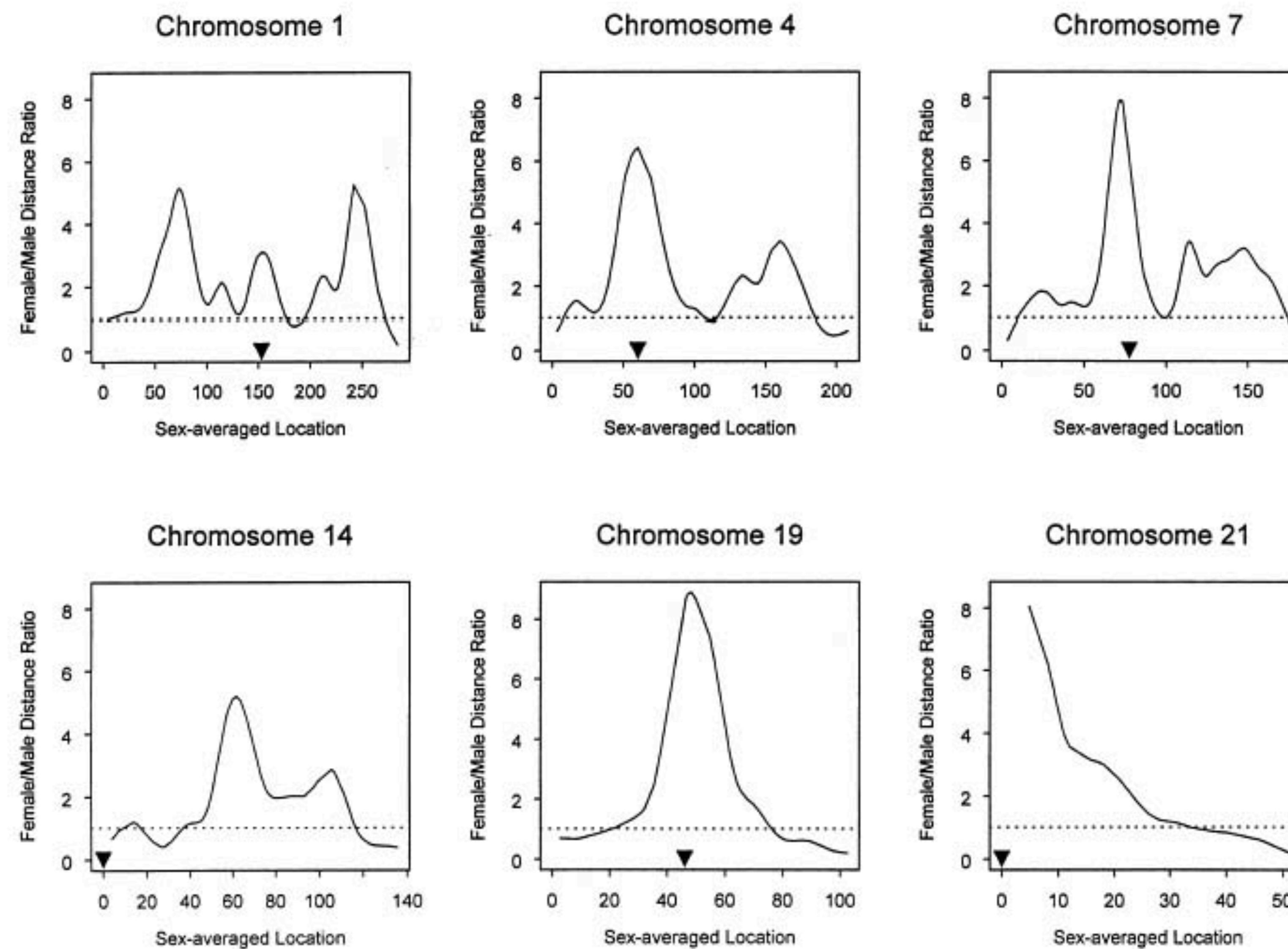
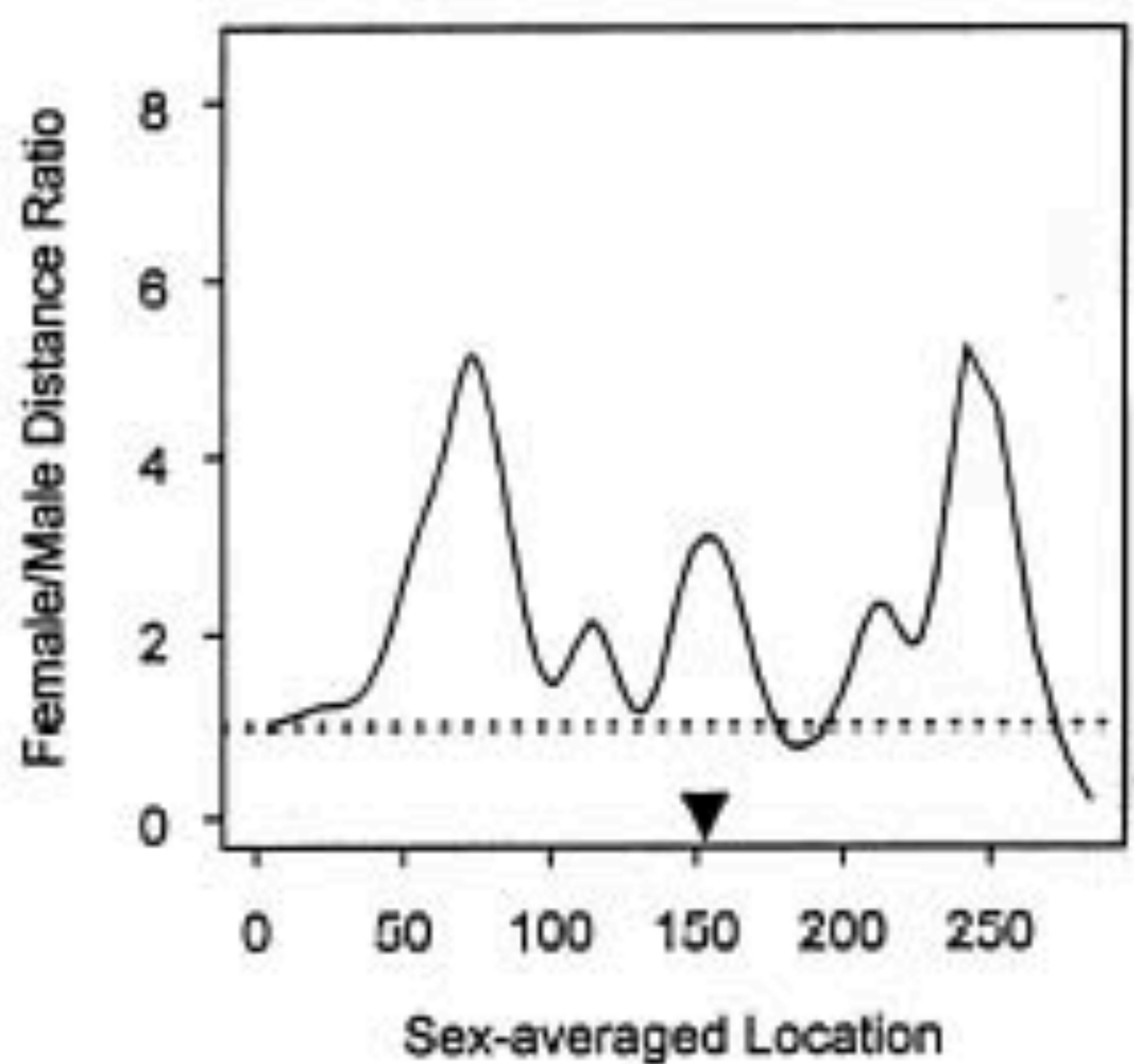
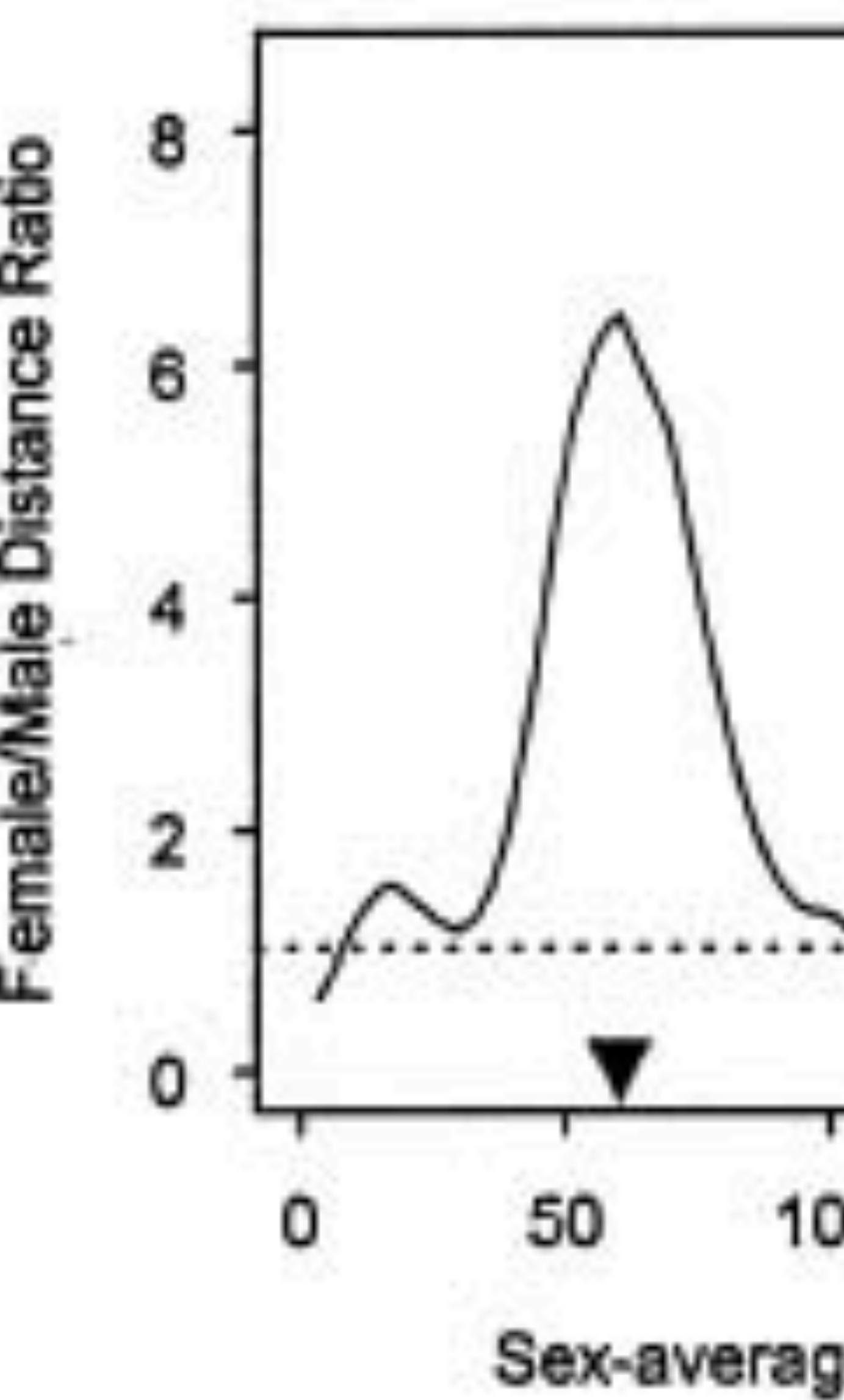


Figure 1 Plots of the female:male genetic-distance ratio against sex-averaged genetic location (in cM) along six selected chromosomes. Approximate locations of the centromeres are indicated by the triangles. The dashed lines correspond to equal female and male distances.

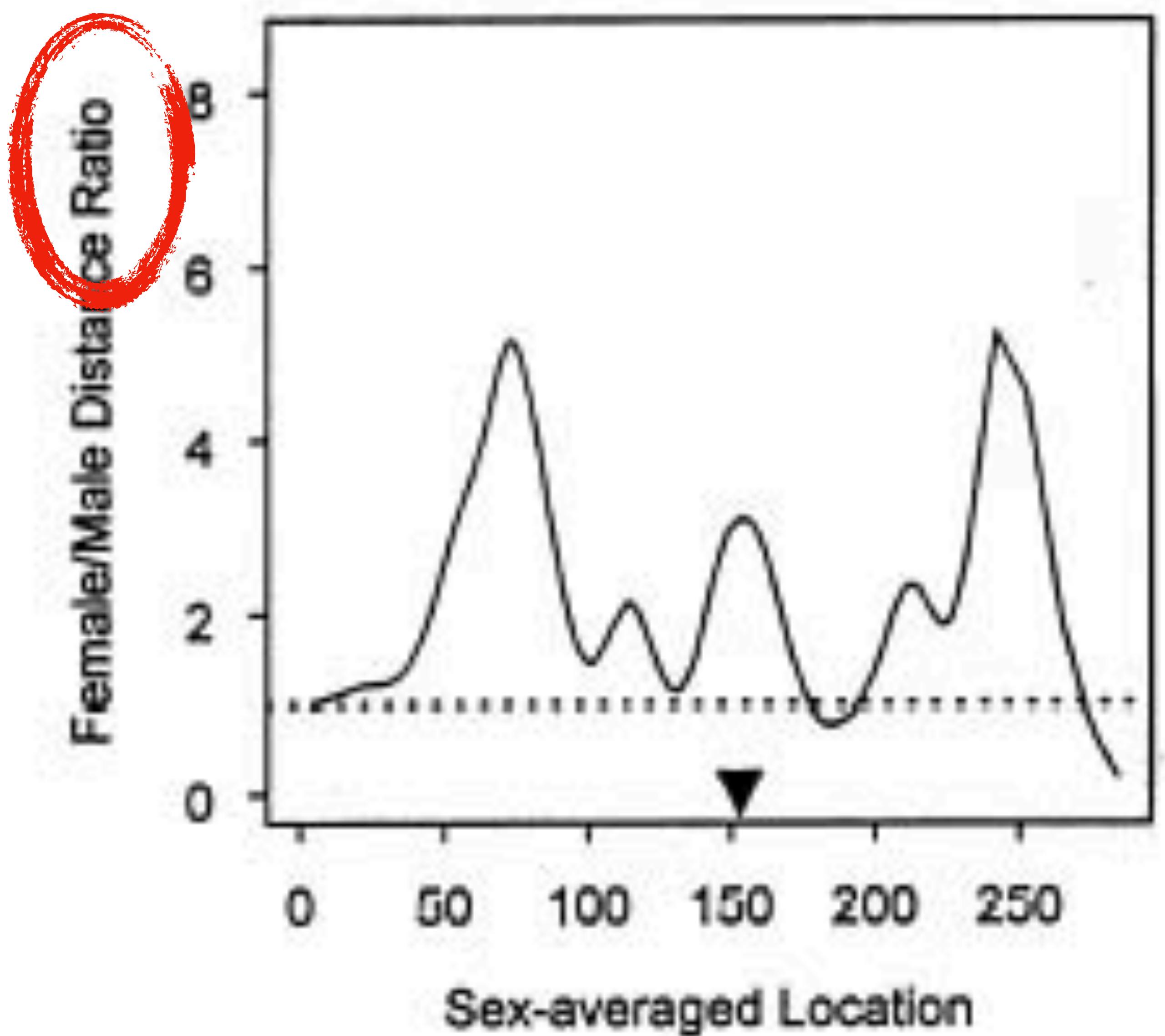
Chromosome 1



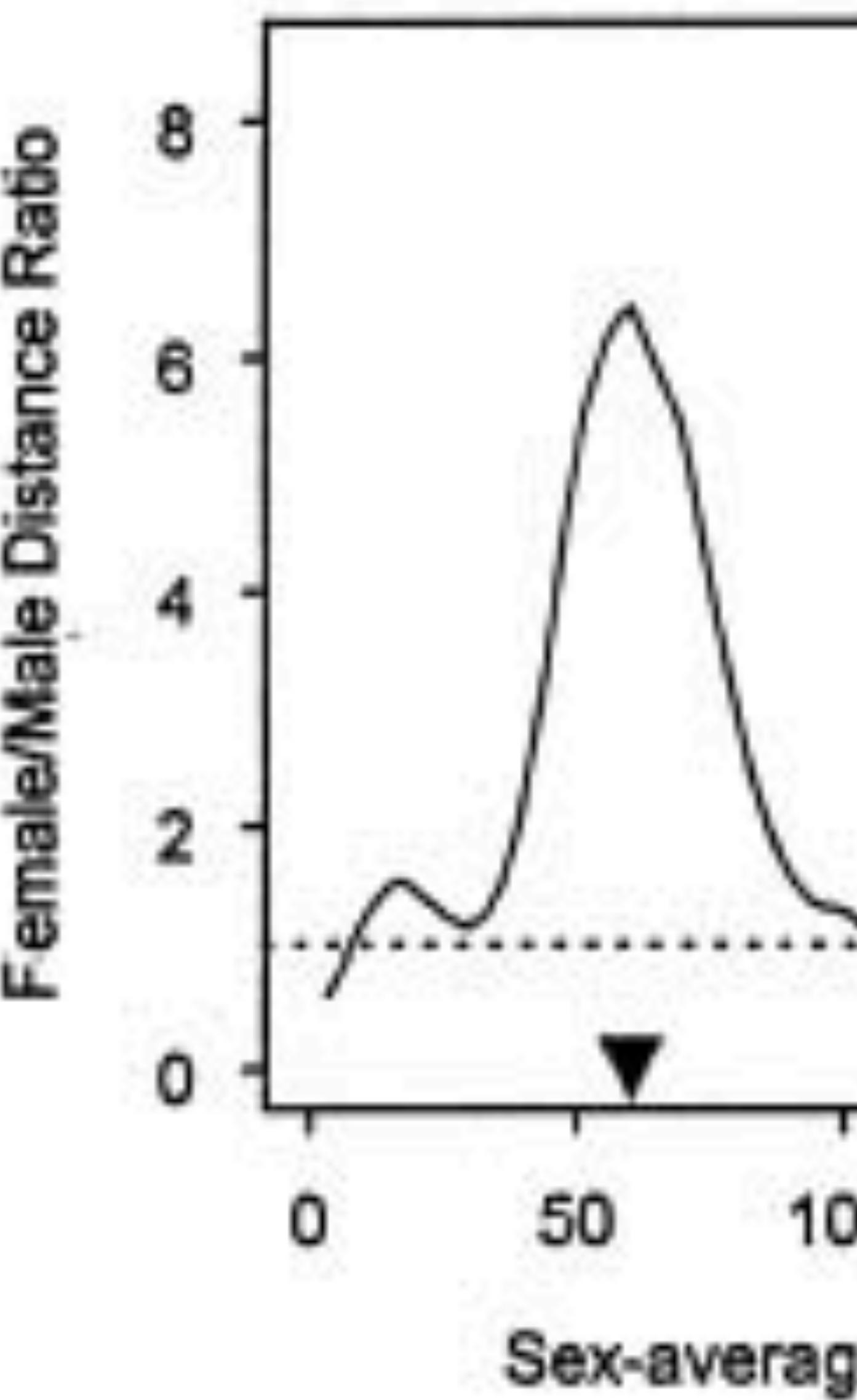
Chromosom



Chromosome 1



Chromosom



What's wrong with this plot?

What is a ratio?

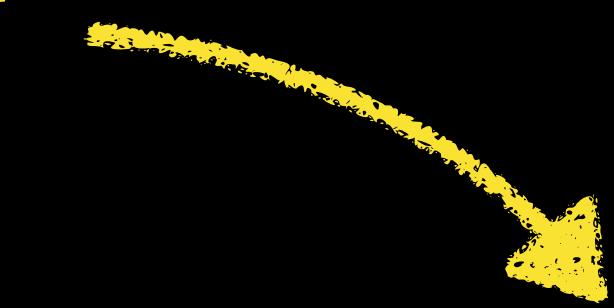
A : B

$A : B$

A / B

$A : B$

Numerator

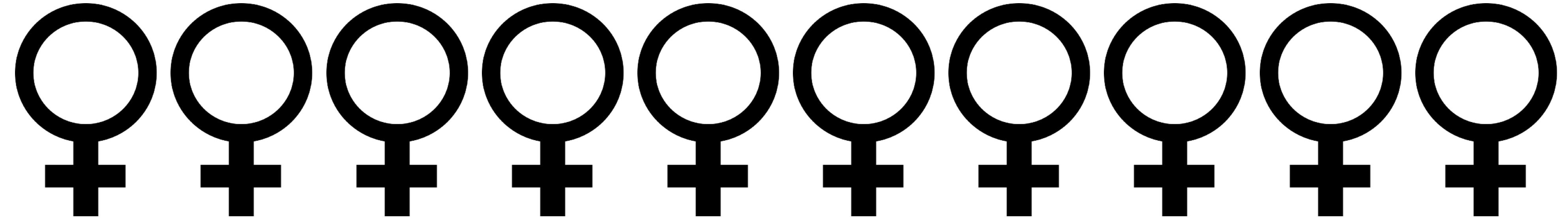
 A / B

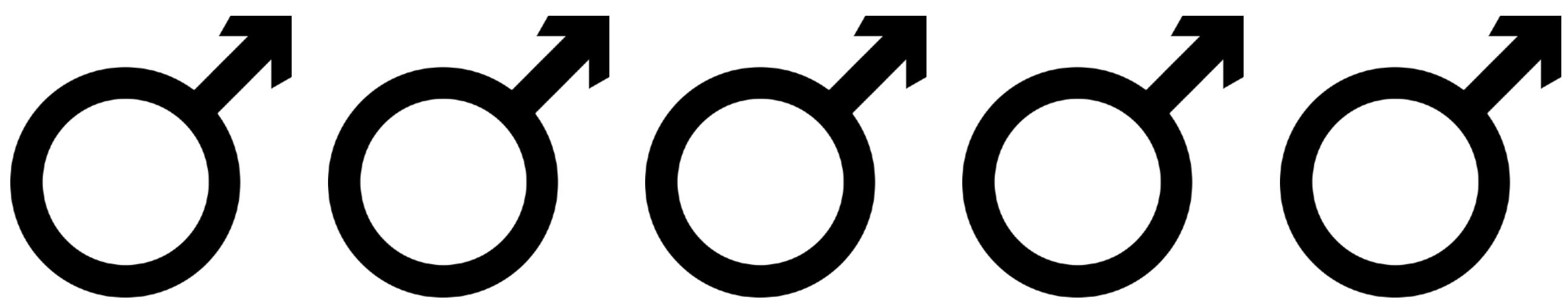
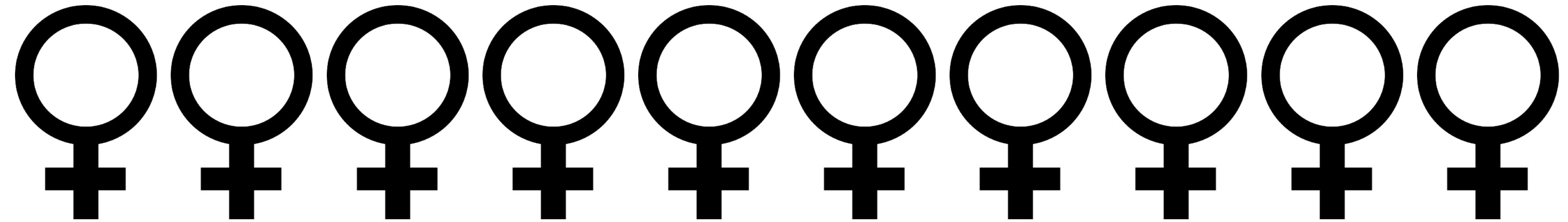
$A : B$

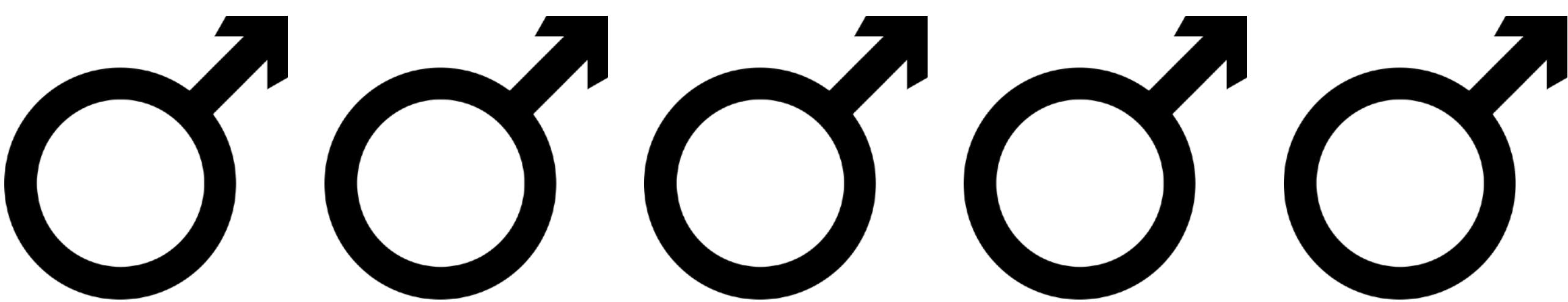
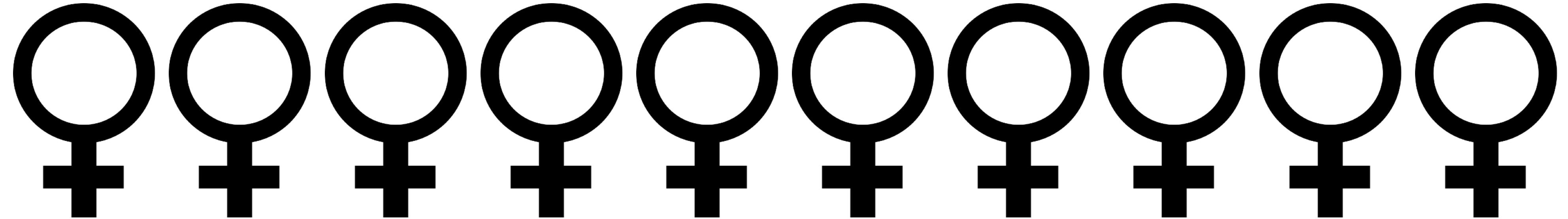
Numerator

Denominator

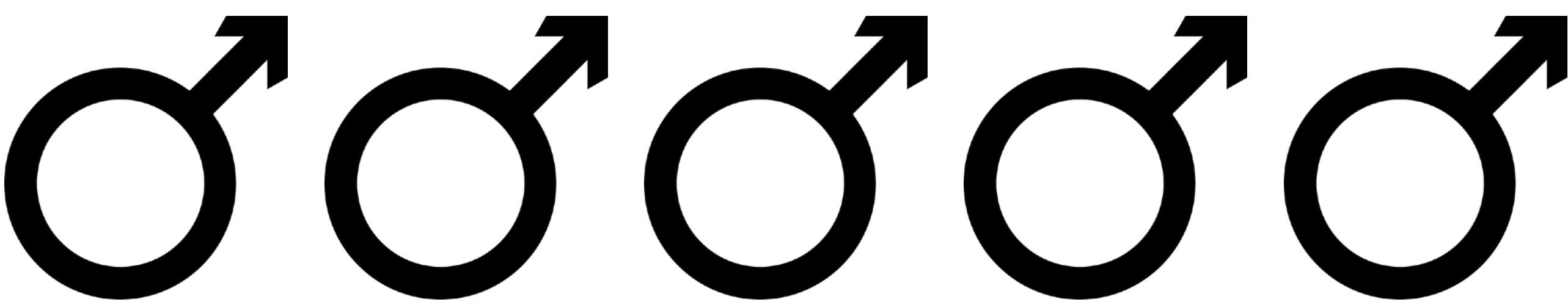
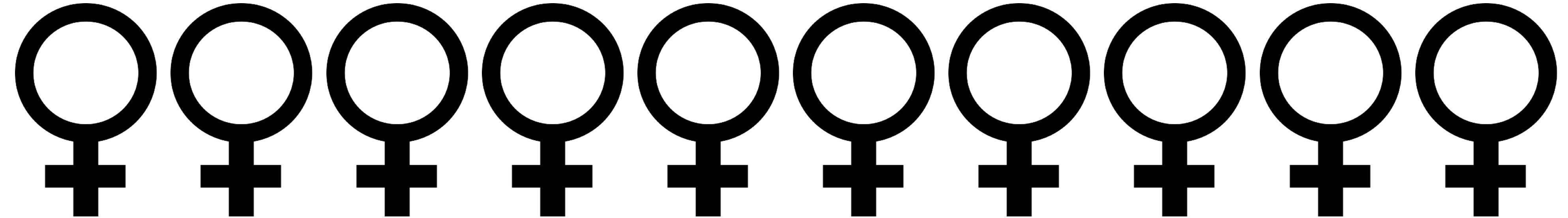
 A / B





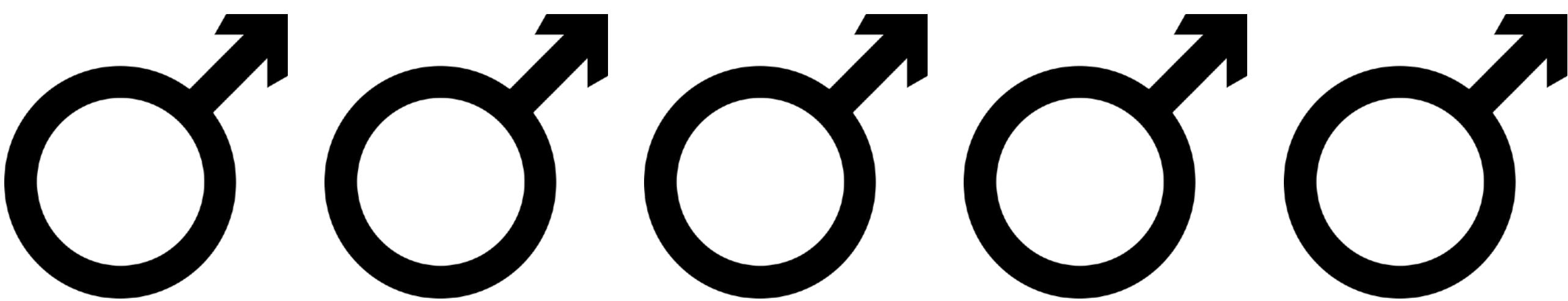
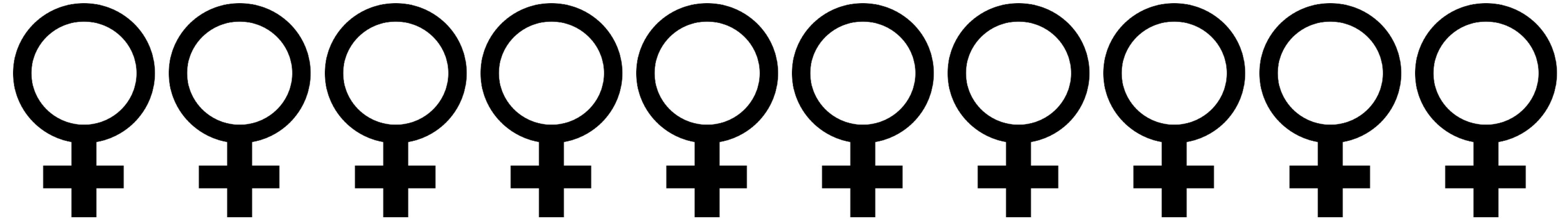


10:5



10:5

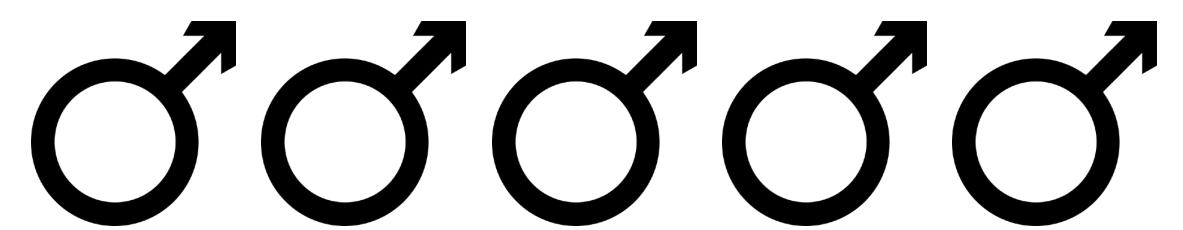
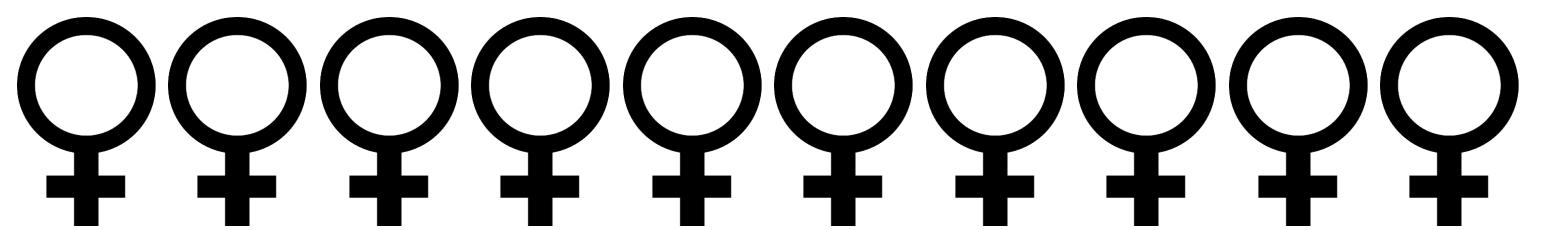
2:1



10:5

2:1

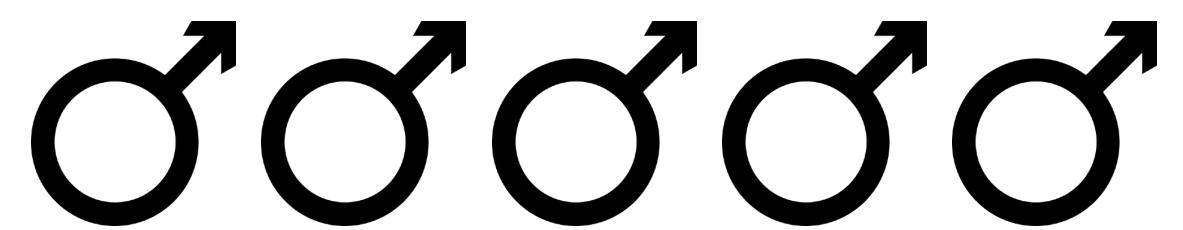
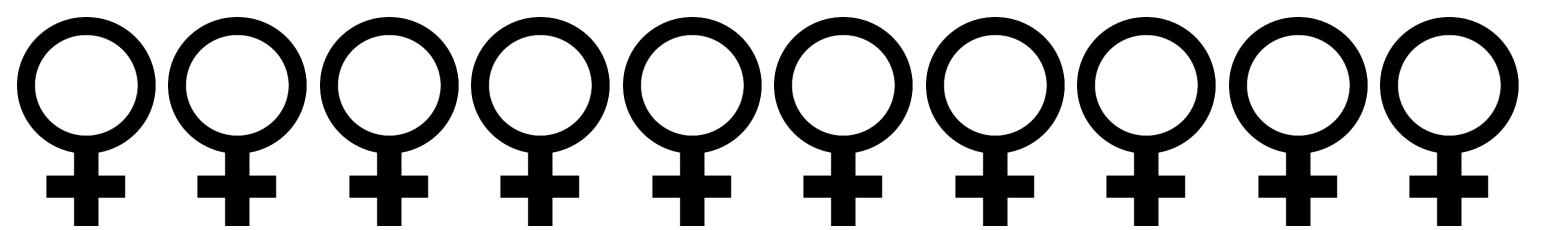
10/5 = 2



10:5

2:1

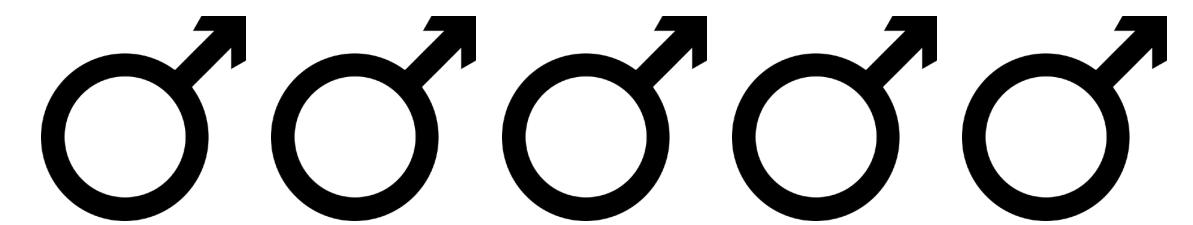
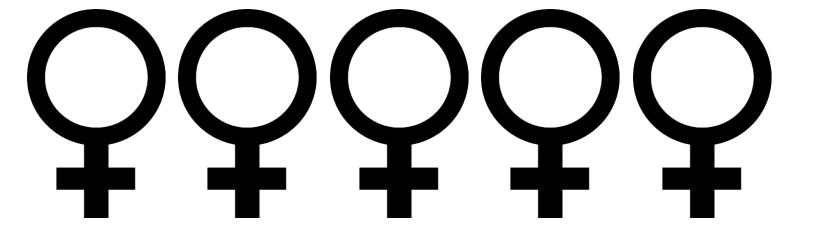
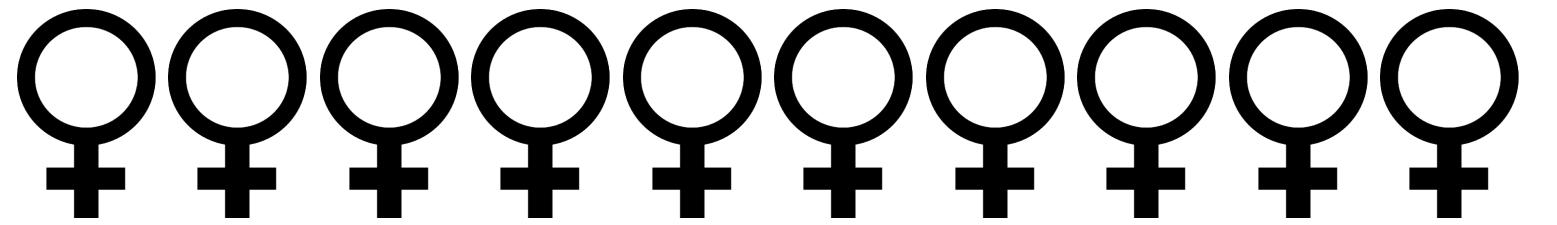
$$10/5 = 2$$



10:5

2:1

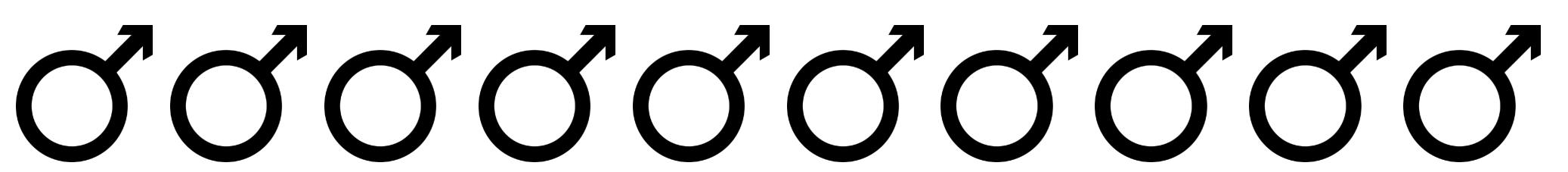
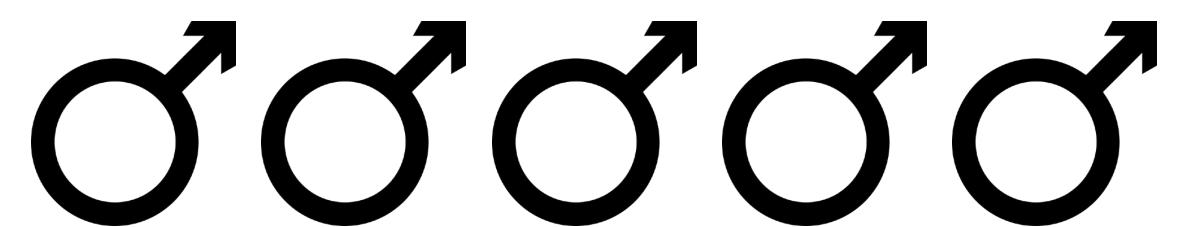
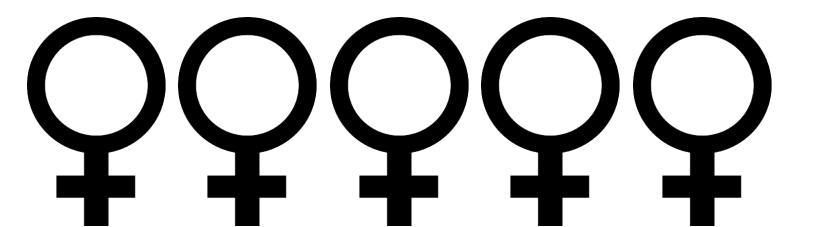
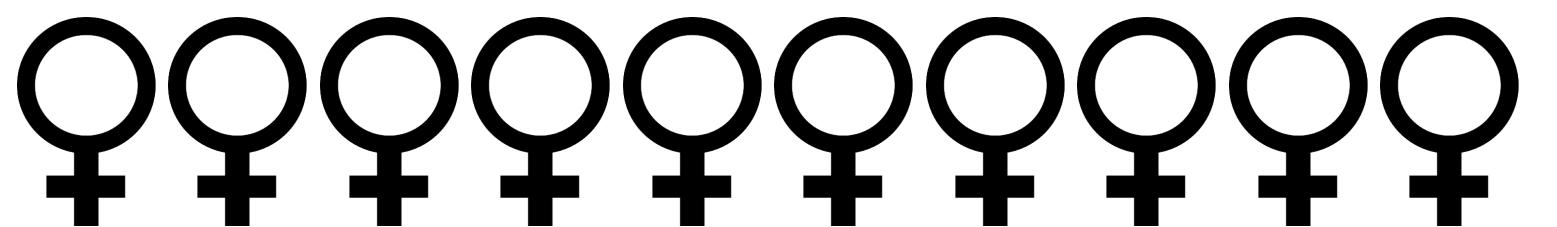
$$10/5 = 2$$



10:5

2:1

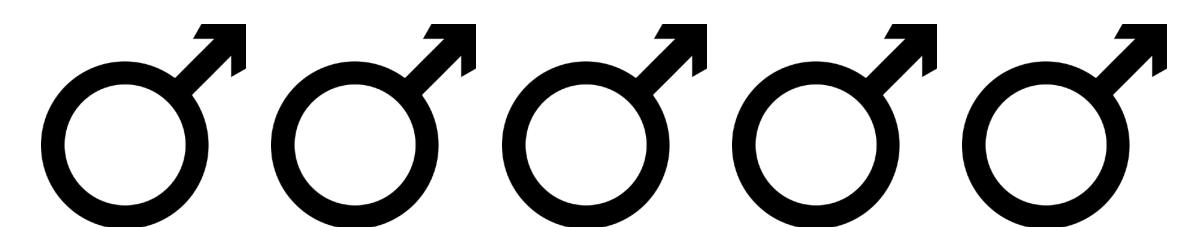
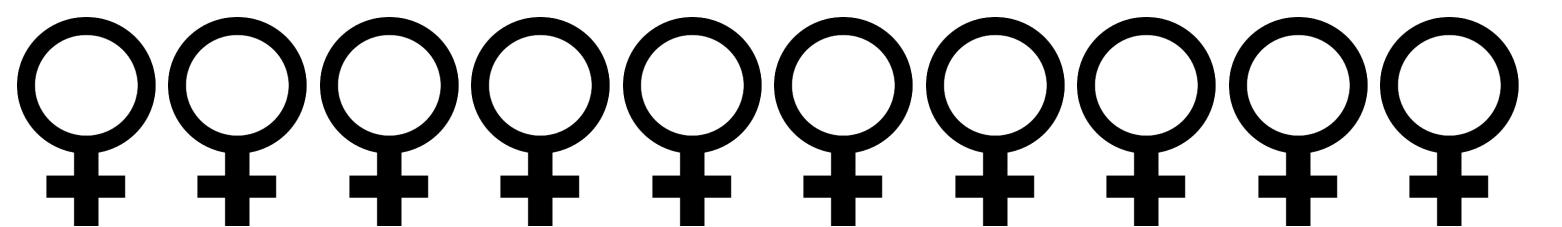
$$10/5 = 2$$



10:5

2:1

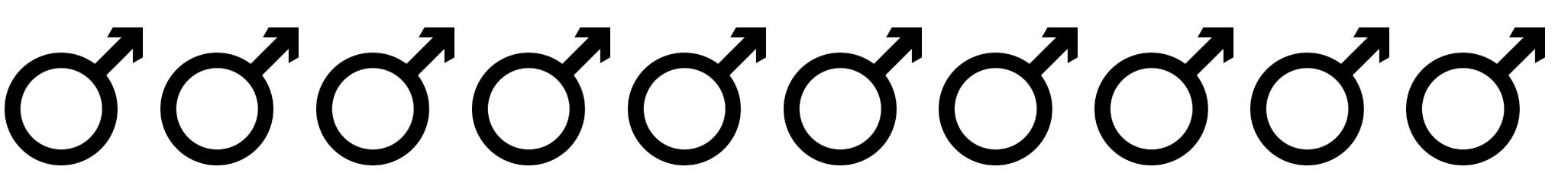
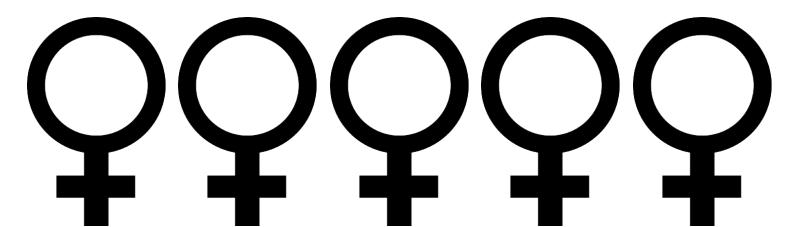
$$10/5 = 2$$



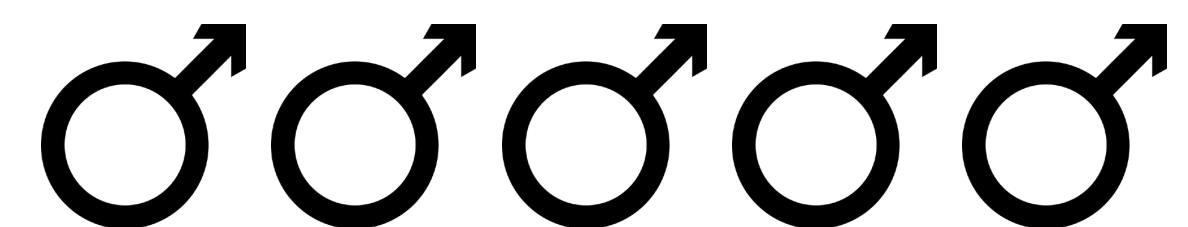
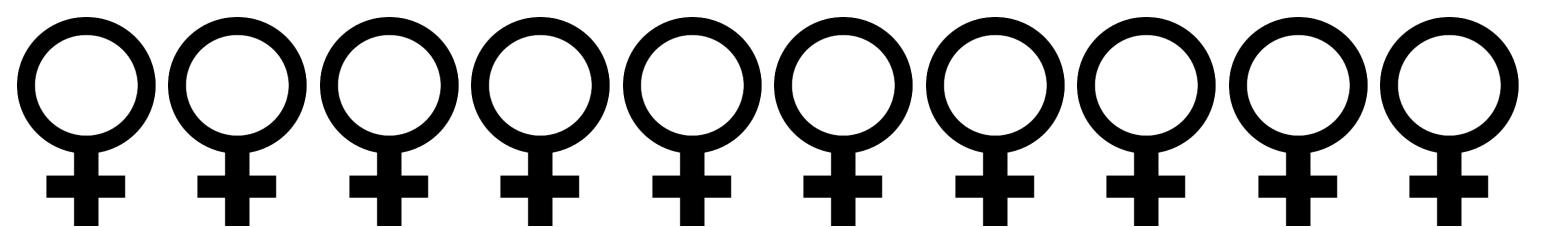
10:5

2:1

$$10/5 = 2$$



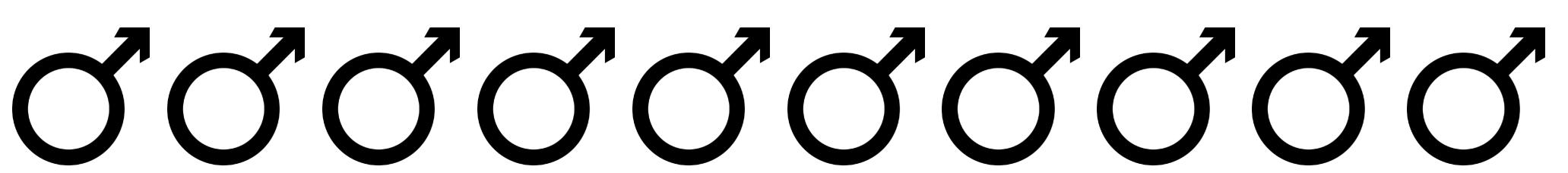
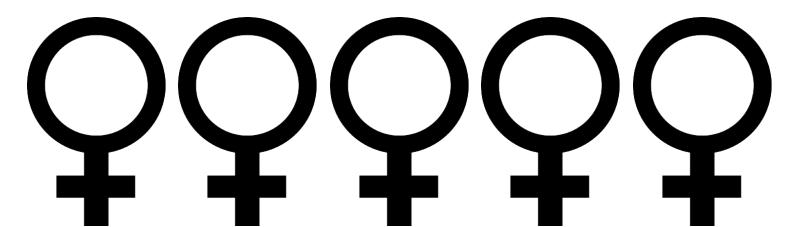
5:10



10:5

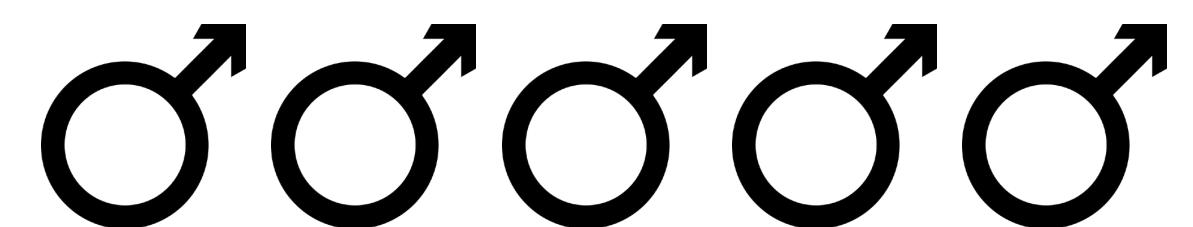
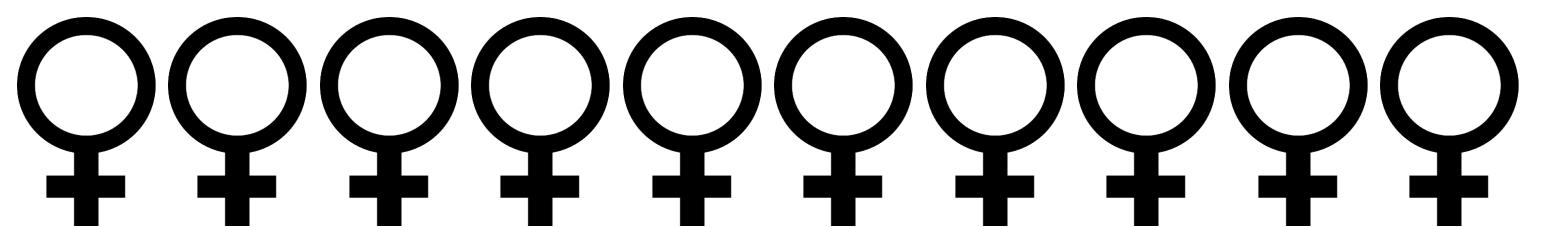
2:1

$$10/5 = 2$$



5:10

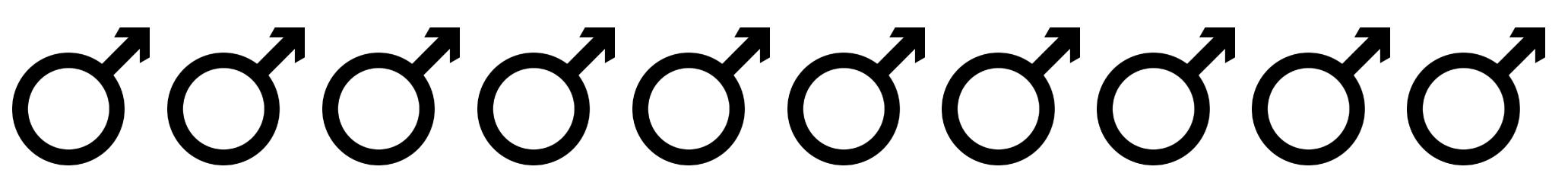
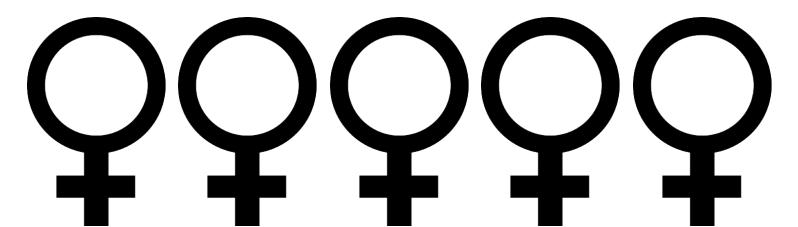
1:2



10:5

2:1

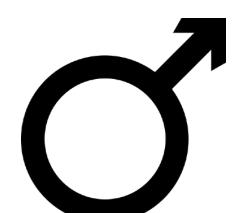
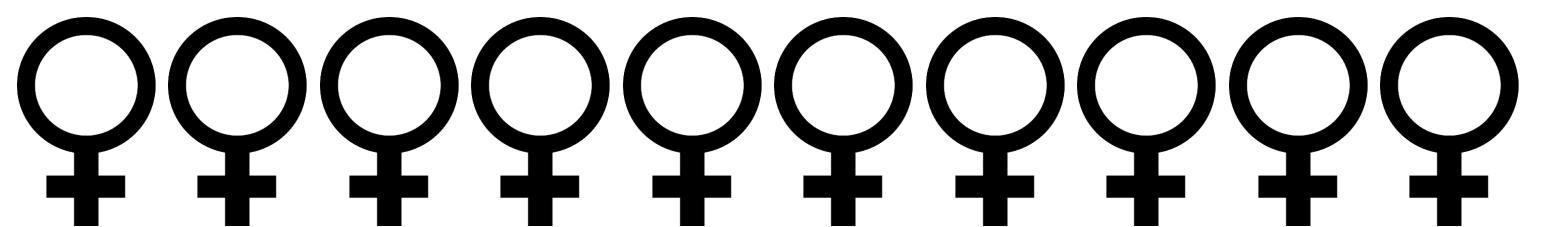
$10/5 = 2$



5:10

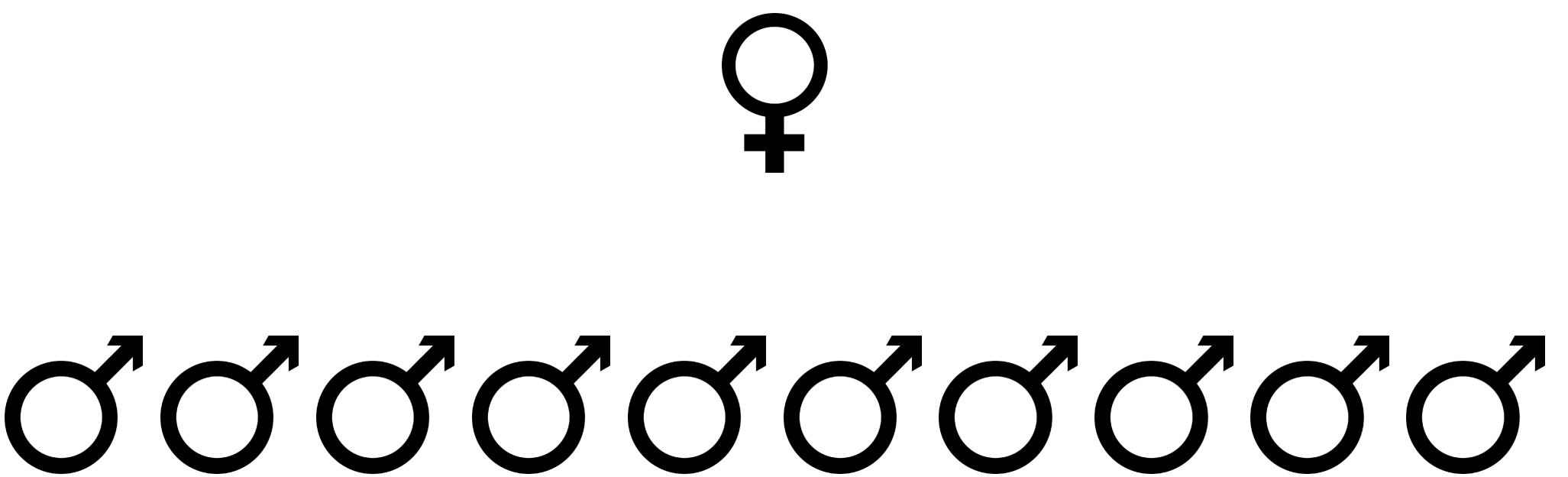
1:2

$5/10 = 0.5$



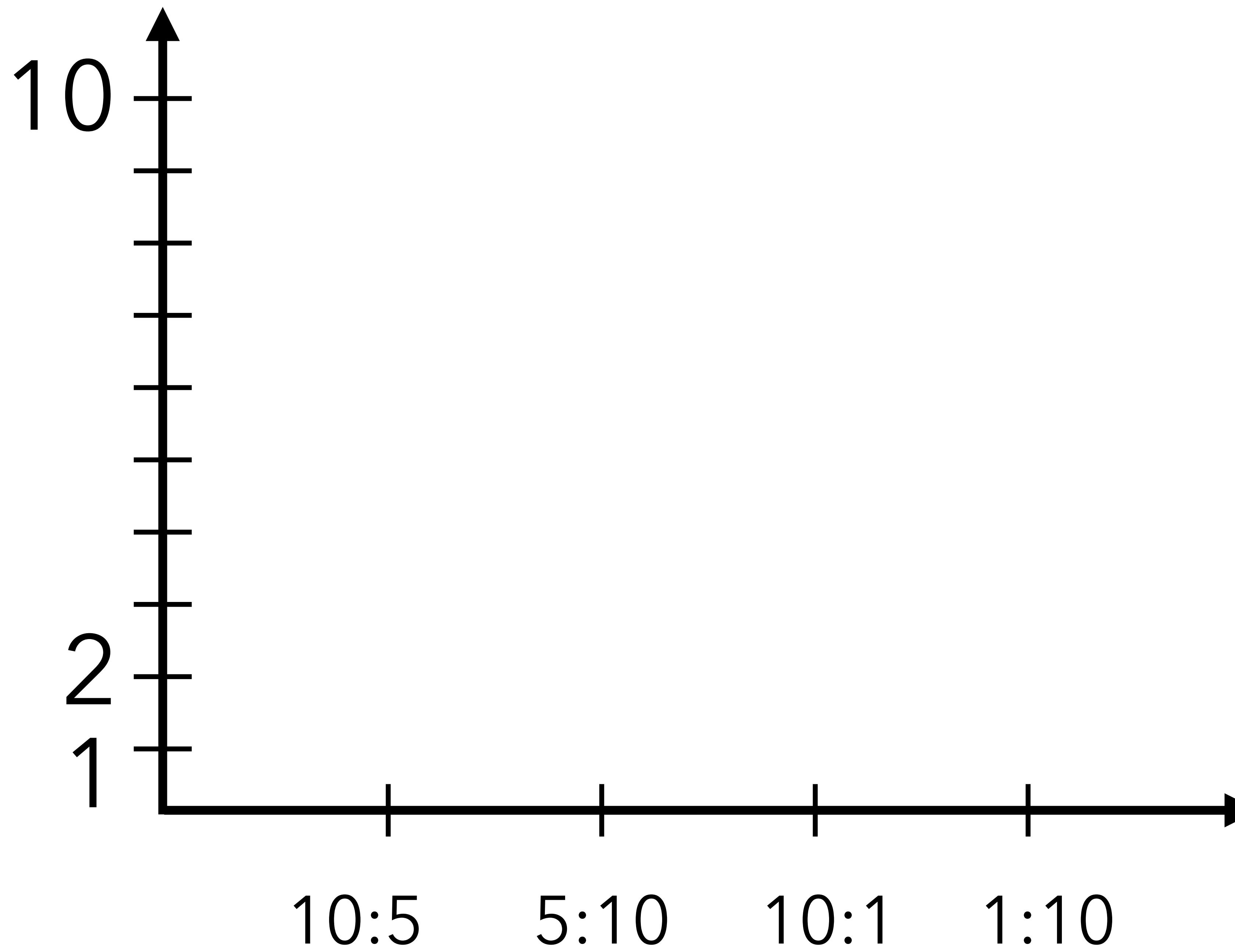
10:1

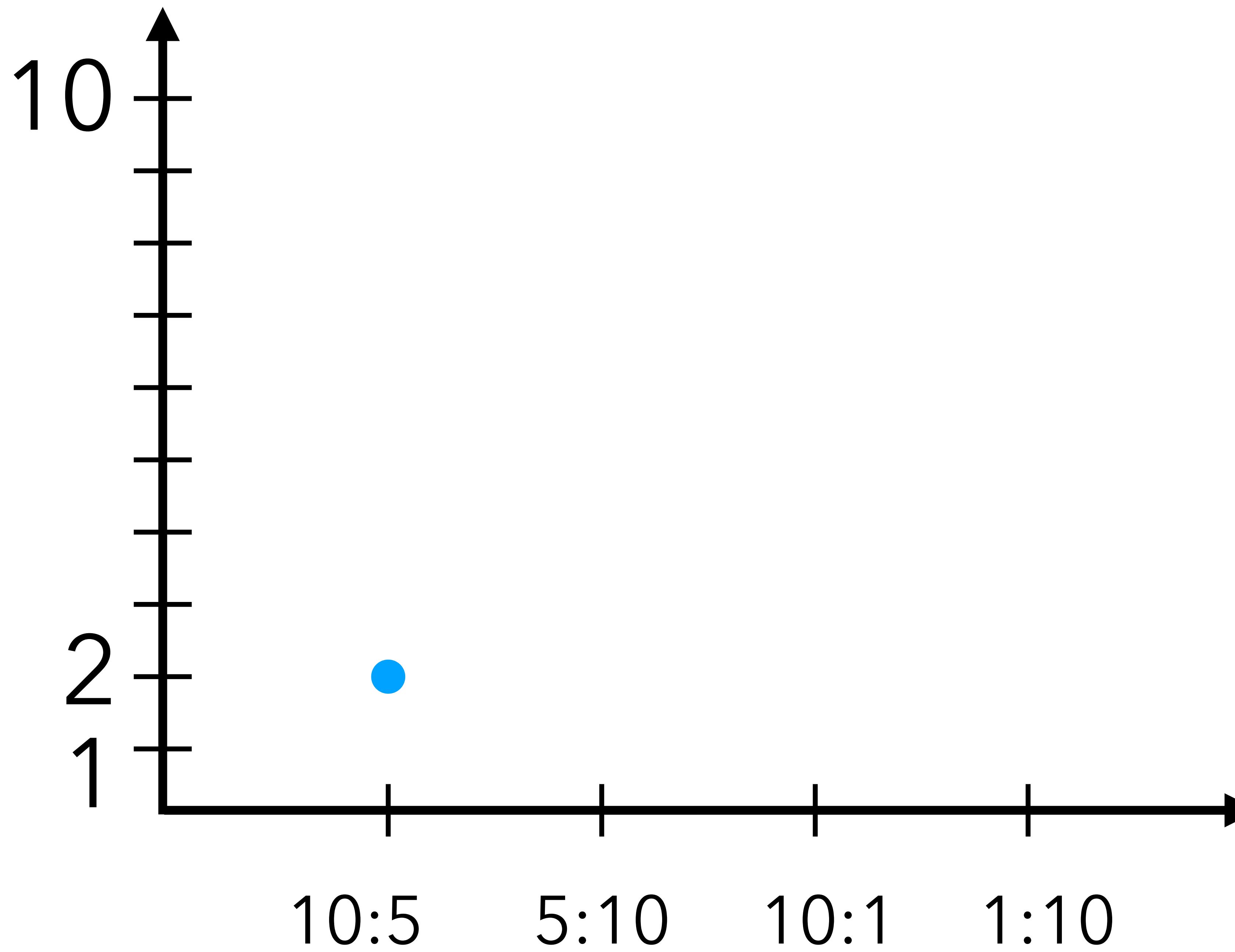
10

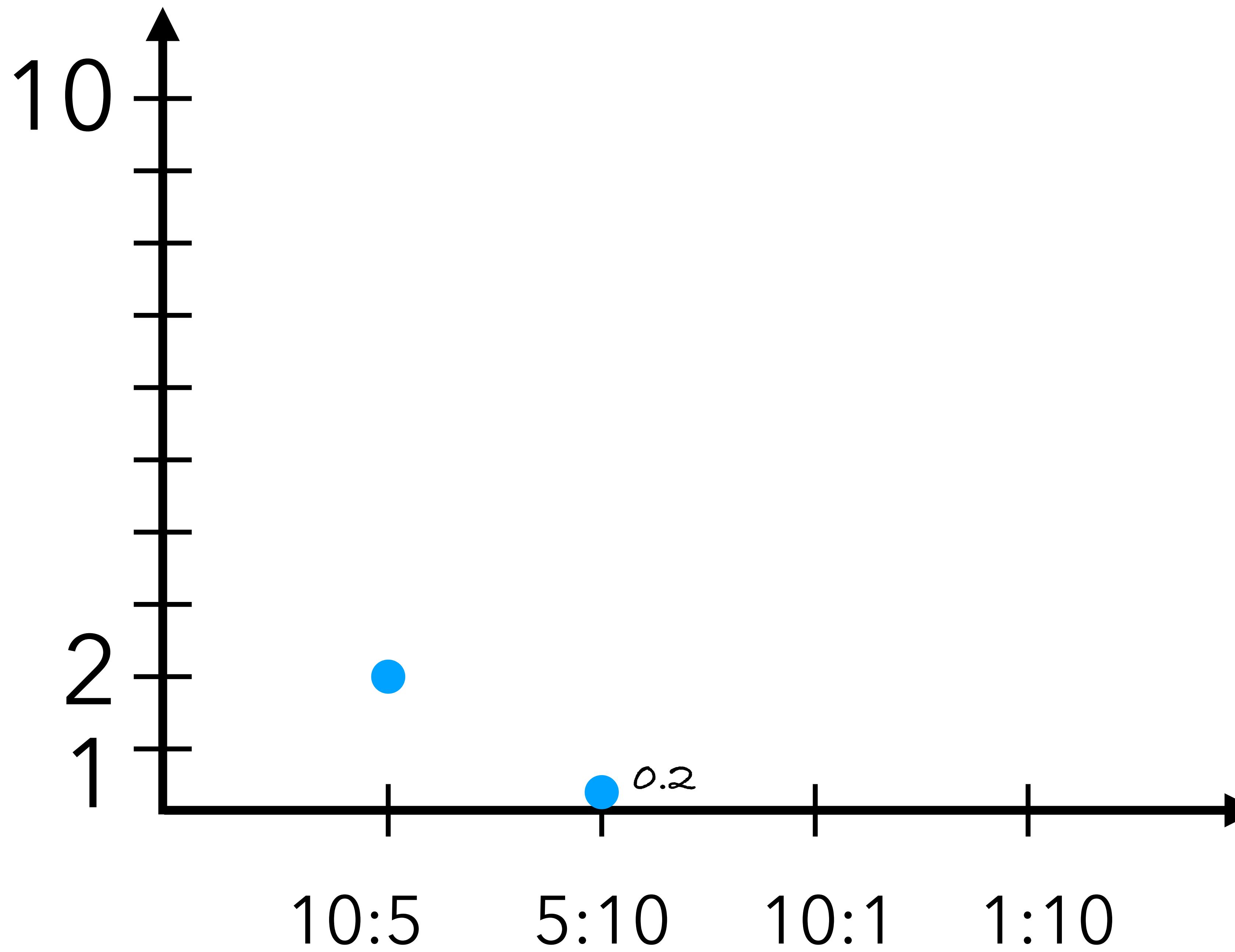


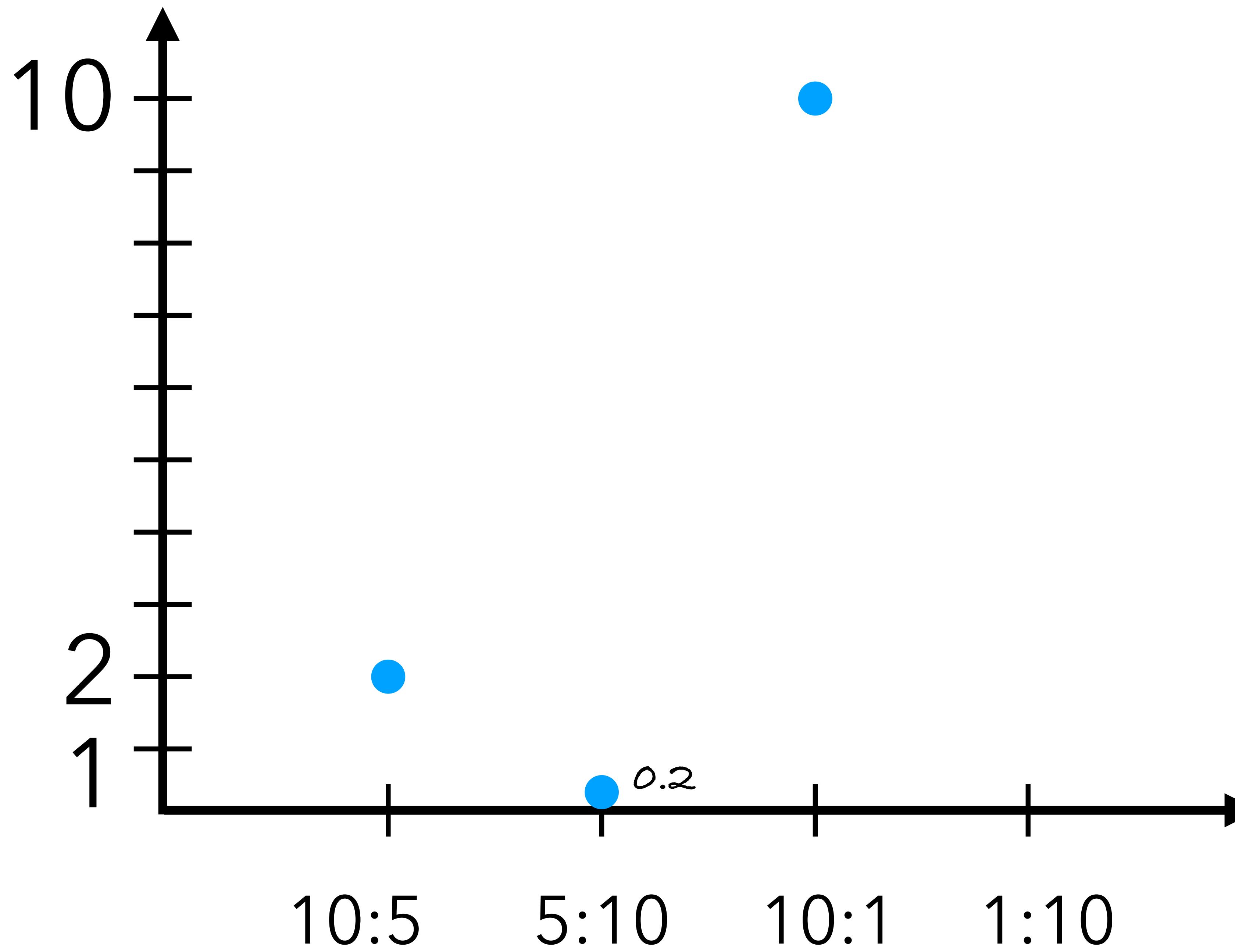
1:10

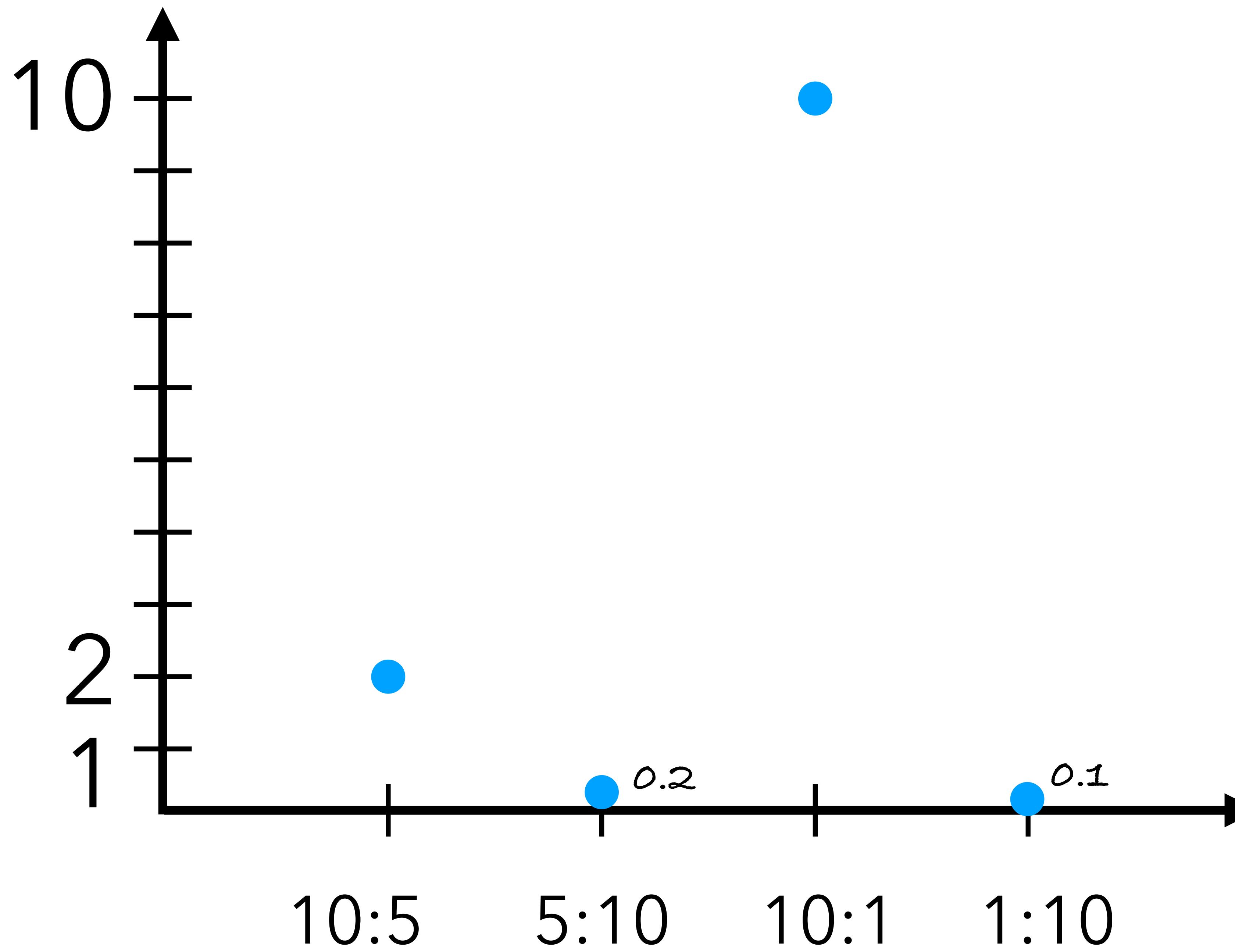
0.1

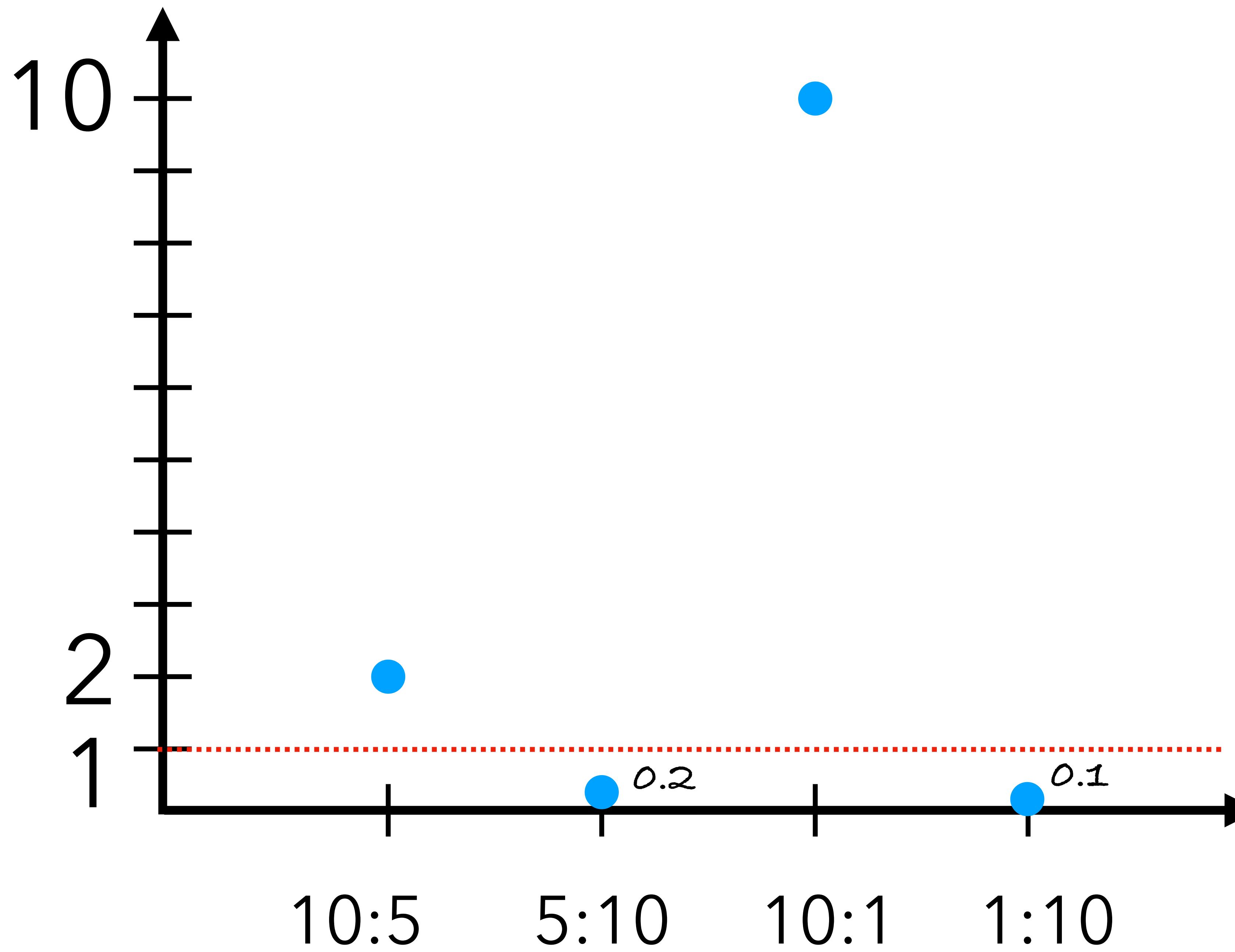


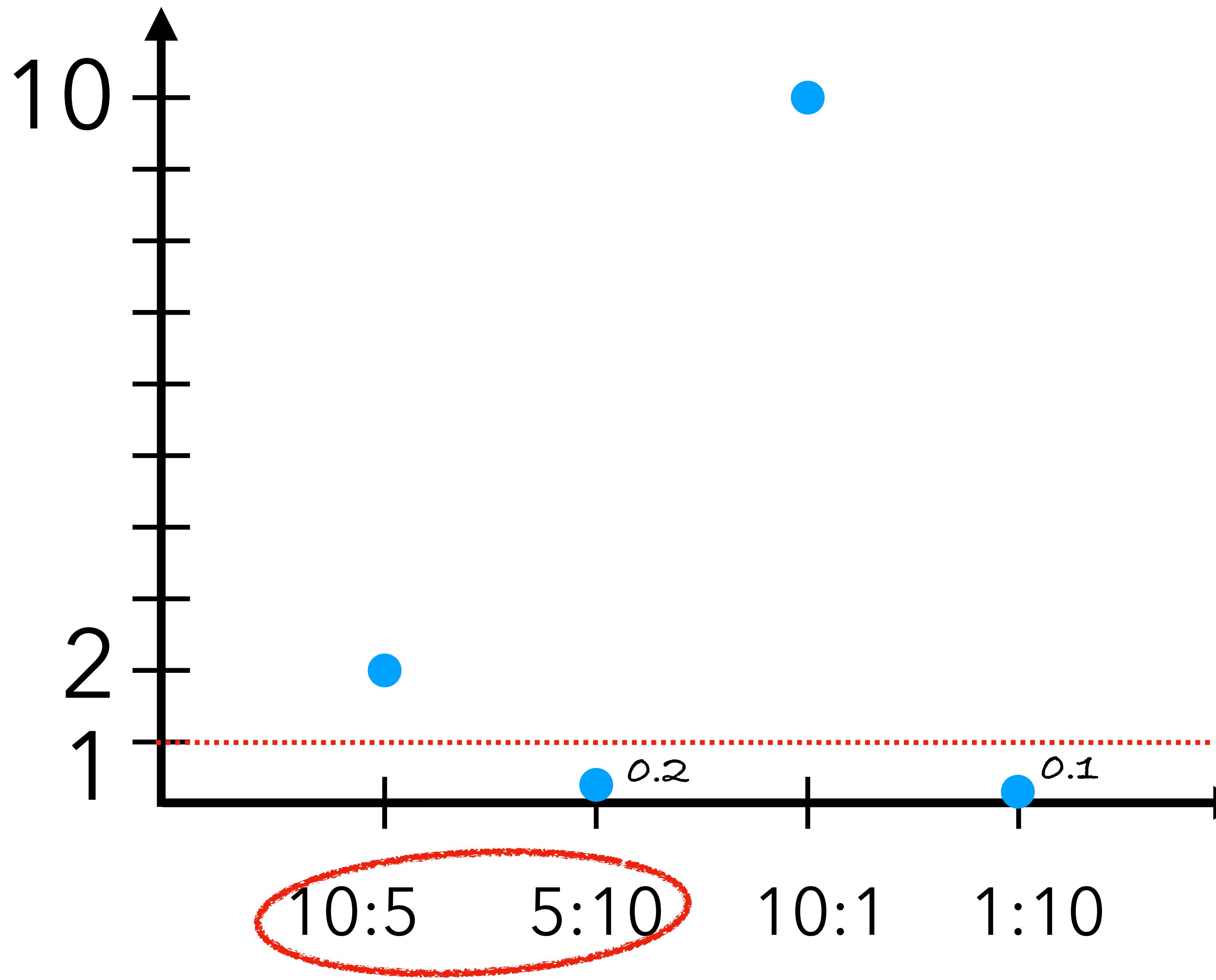


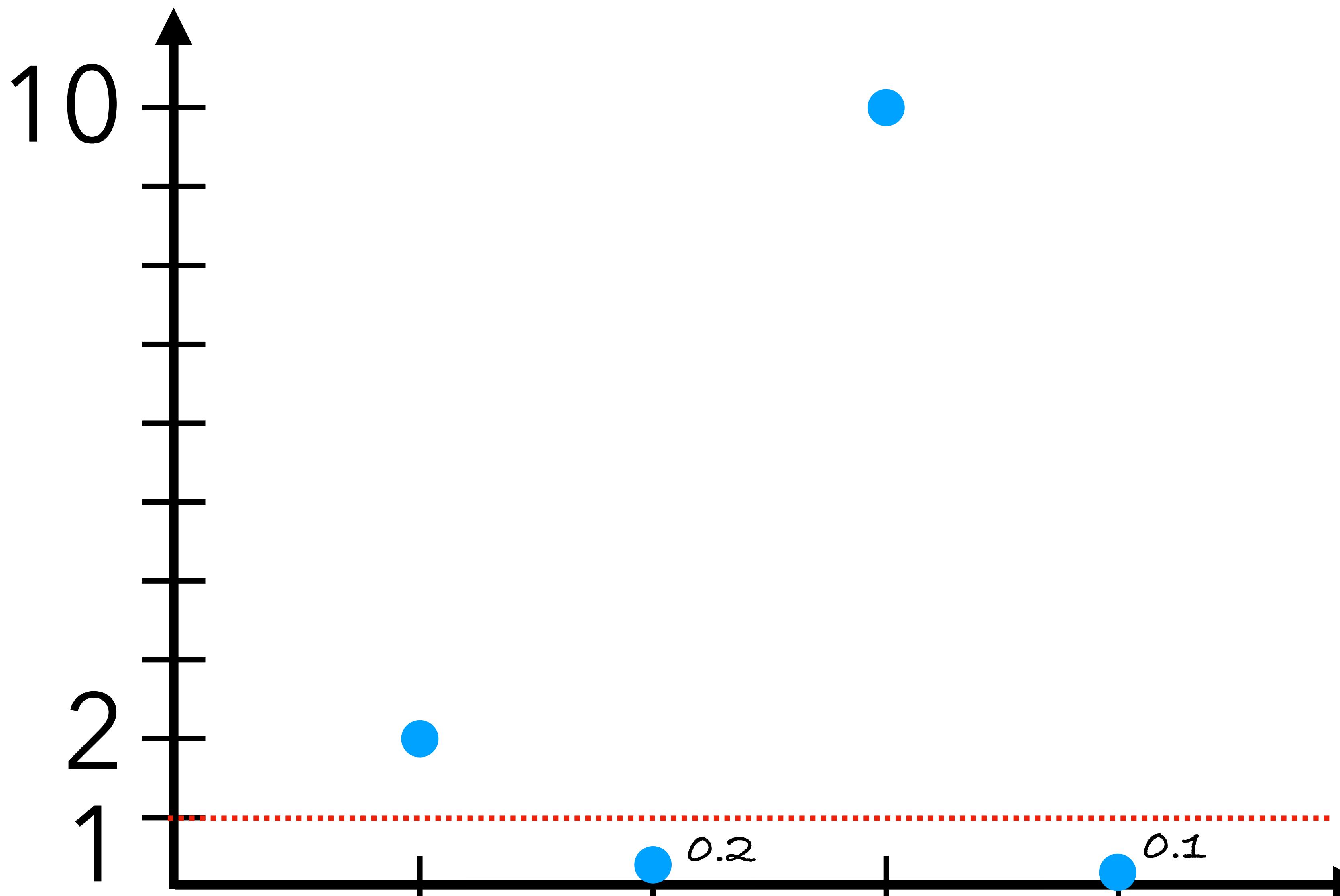










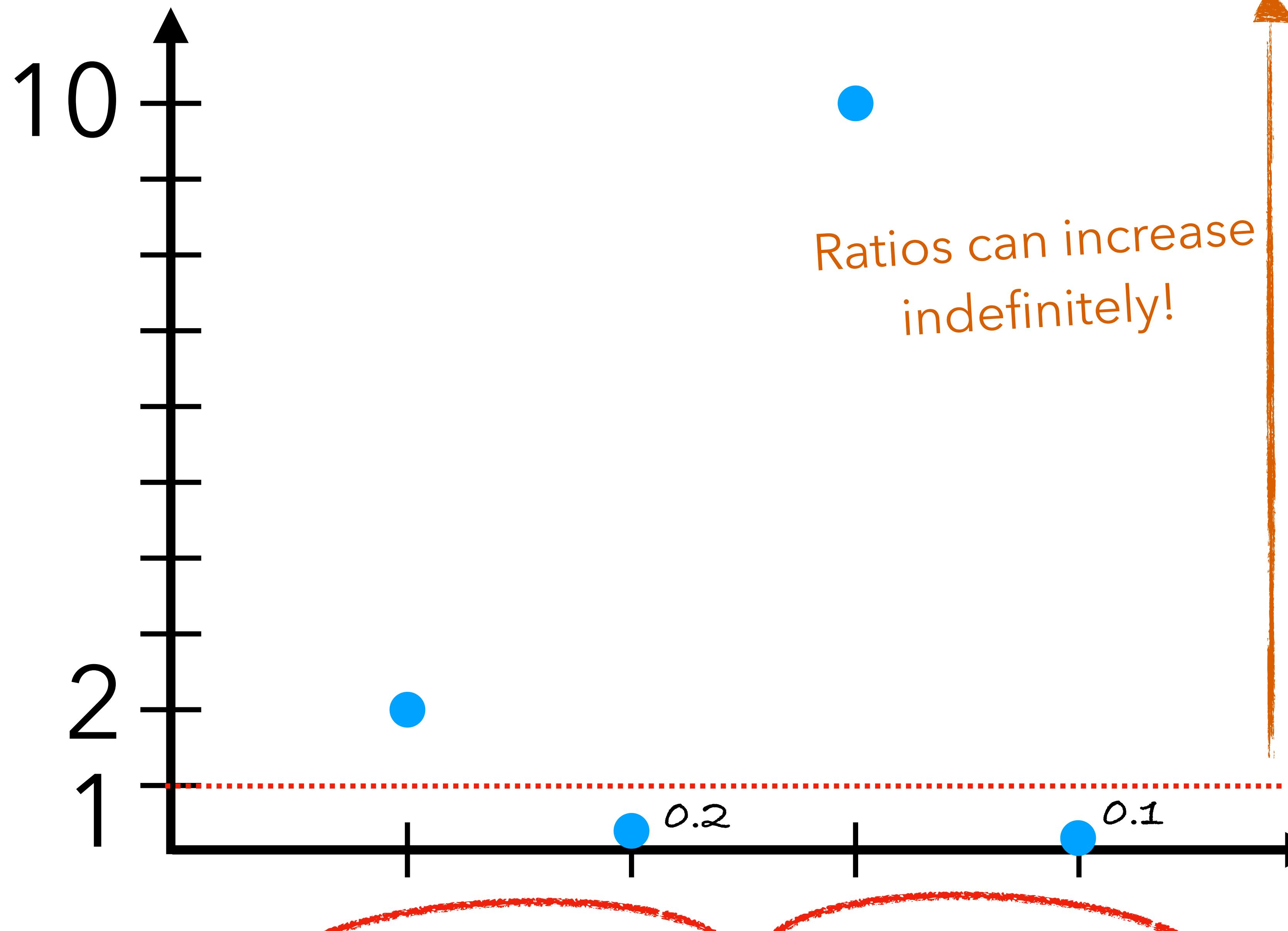


10:5

5:10

10:1

1:10



Ratios can increase
indefinitely!

Log scale!

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

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

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

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

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

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

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

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

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

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

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

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

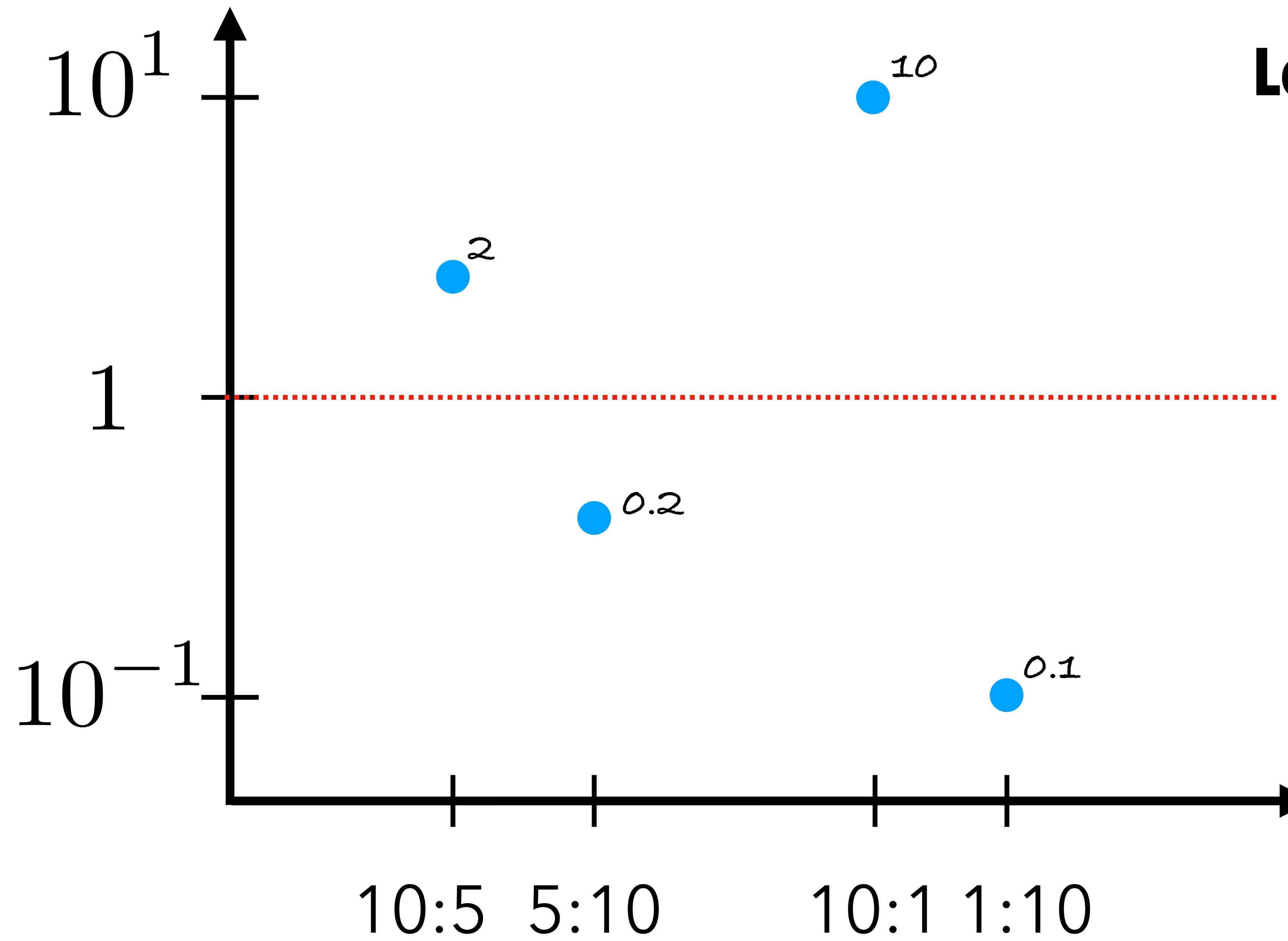
$$\log \frac{10}{1} = \underline{\log 10}$$

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

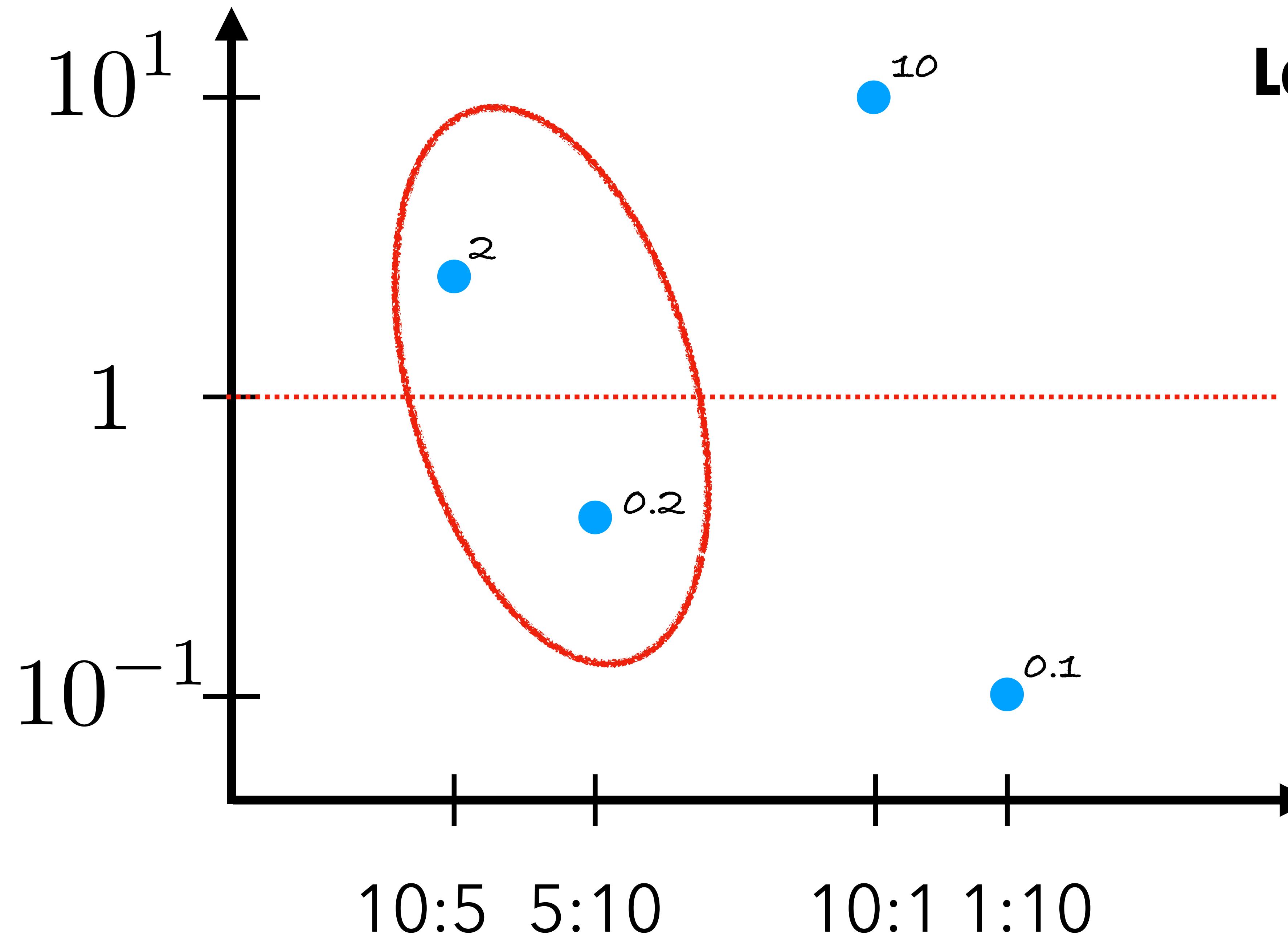
$$\log \frac{10}{1} = \underline{\log 10}$$

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

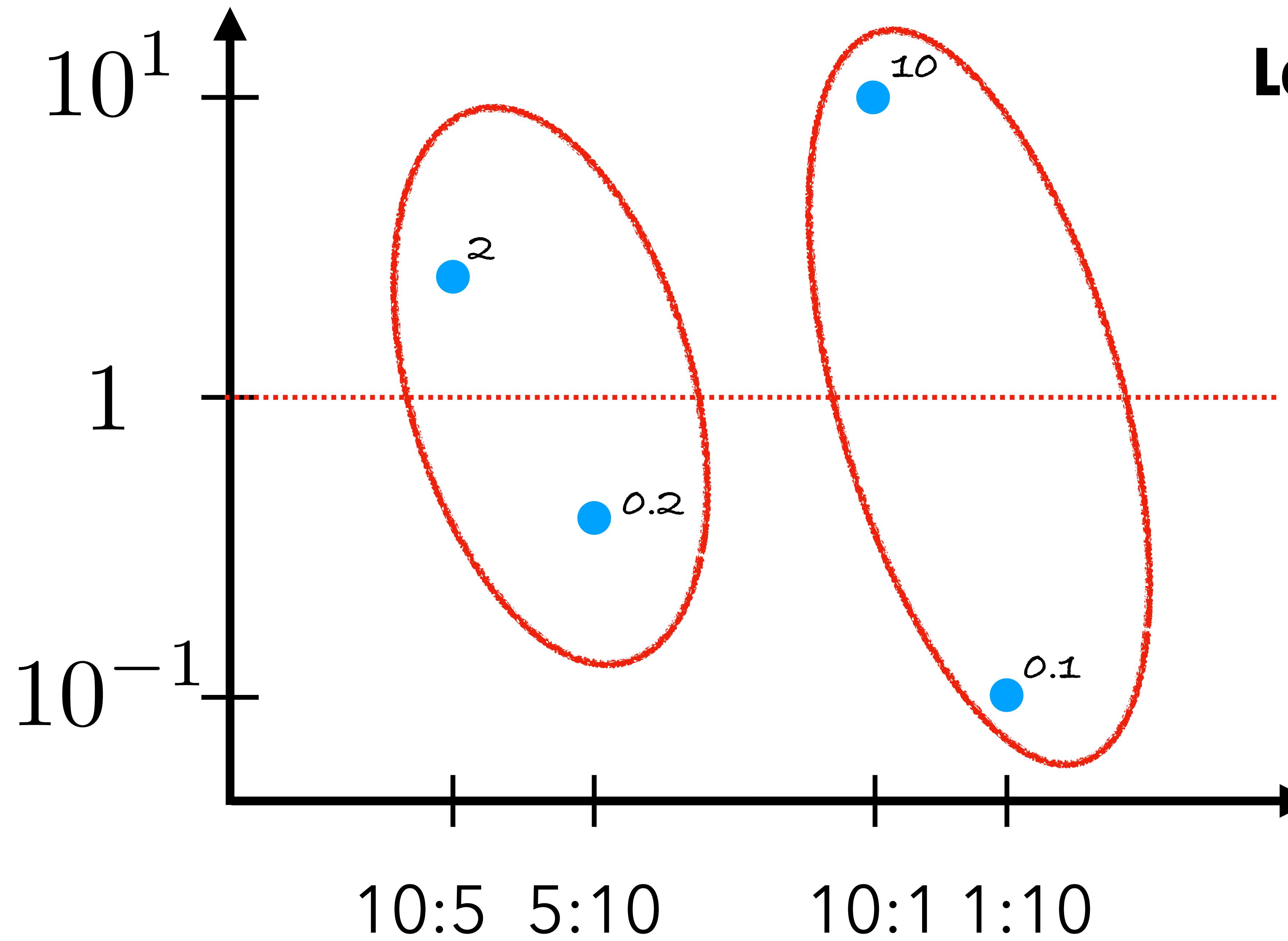
Log scale



Log scale



Log scale



Log scale

