

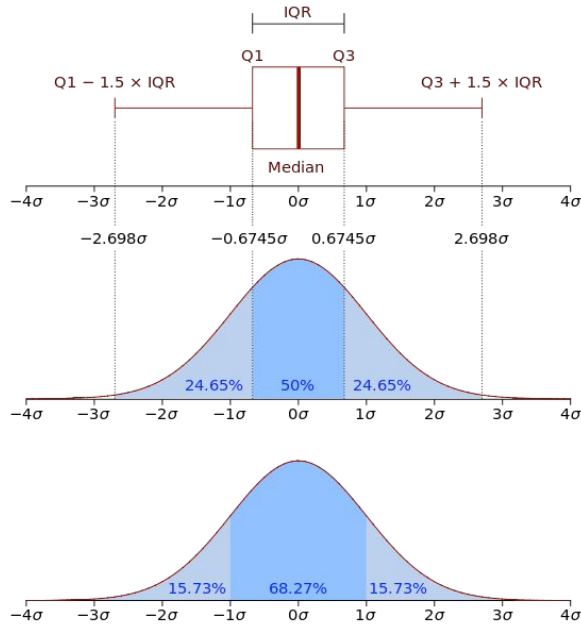
Data Visualization

W7-1

Quiz

- What do you find interesting in today's VotW?
- Explain every element (e.g., box, line, ...) of a box plot. What do they represent and how do we calculate them?
- In histogram, _____ represents the frequency of the data.
- What are the issues when bins are too wide or too narrow in drawing histograms?

Let's construct a box plot

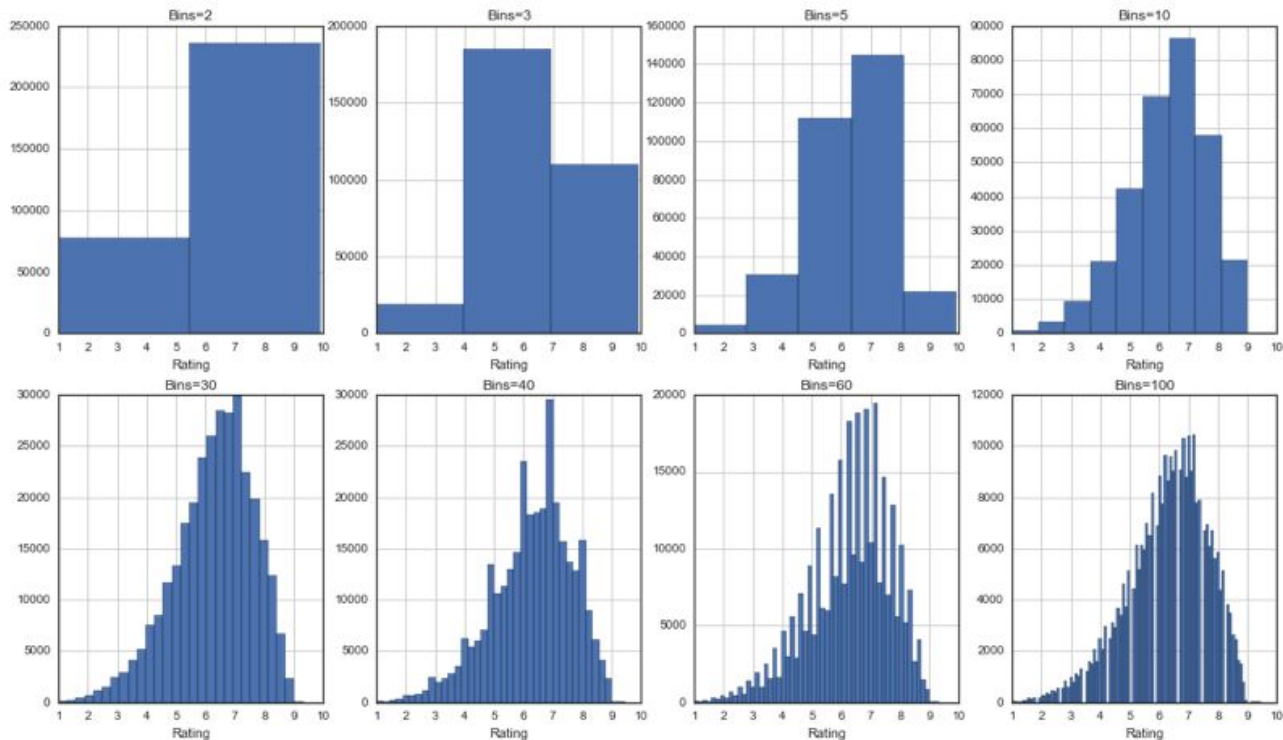


1. We want to know the **median** (the central value in the data).
2. What's the range that capture about the **half of the data points**?
3. How about **most of the data points**? (Many ways to estimate this!)
4. What do we do with the **rest of the data points**?

In histogram, **area**, not
the height, should
represent the frequency!

Too wide bins \rightarrow large bias (inaccurate)

Too narrow bins \rightarrow large variance (overfit to the data/noise)



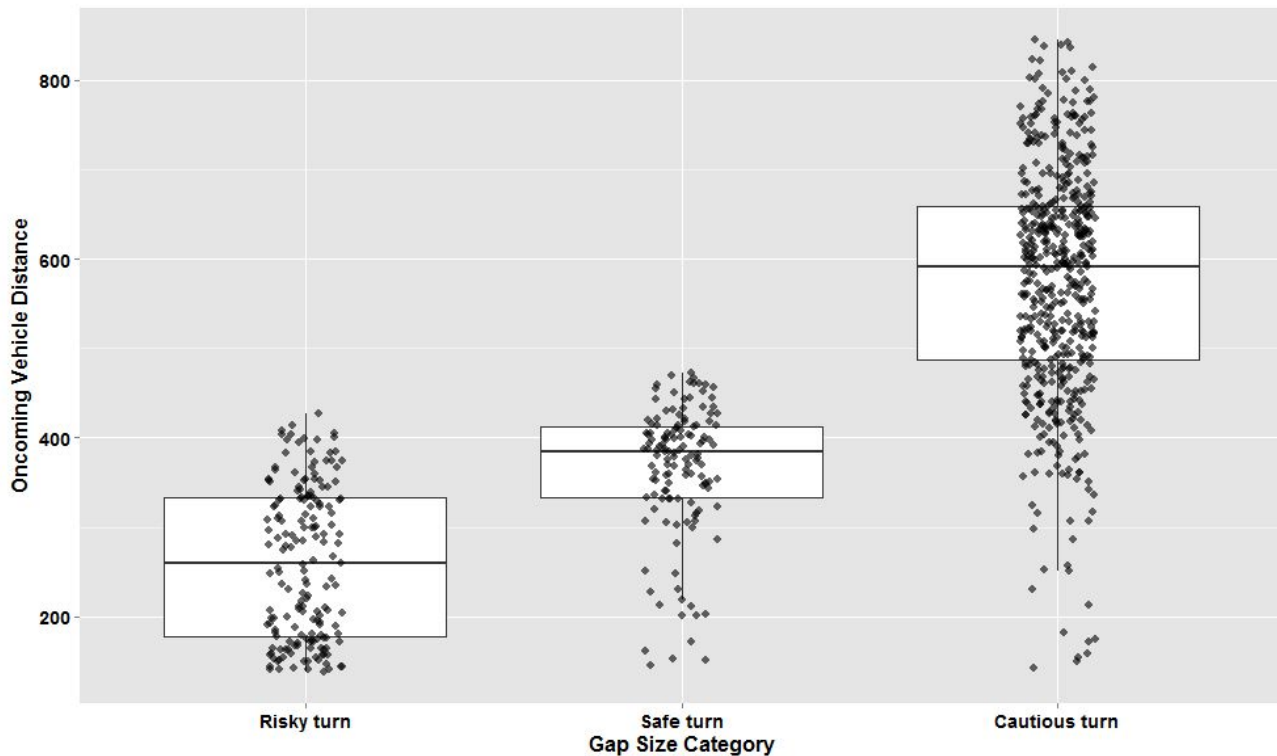
Histogram is cheap and
easy.

Try multiple histograms

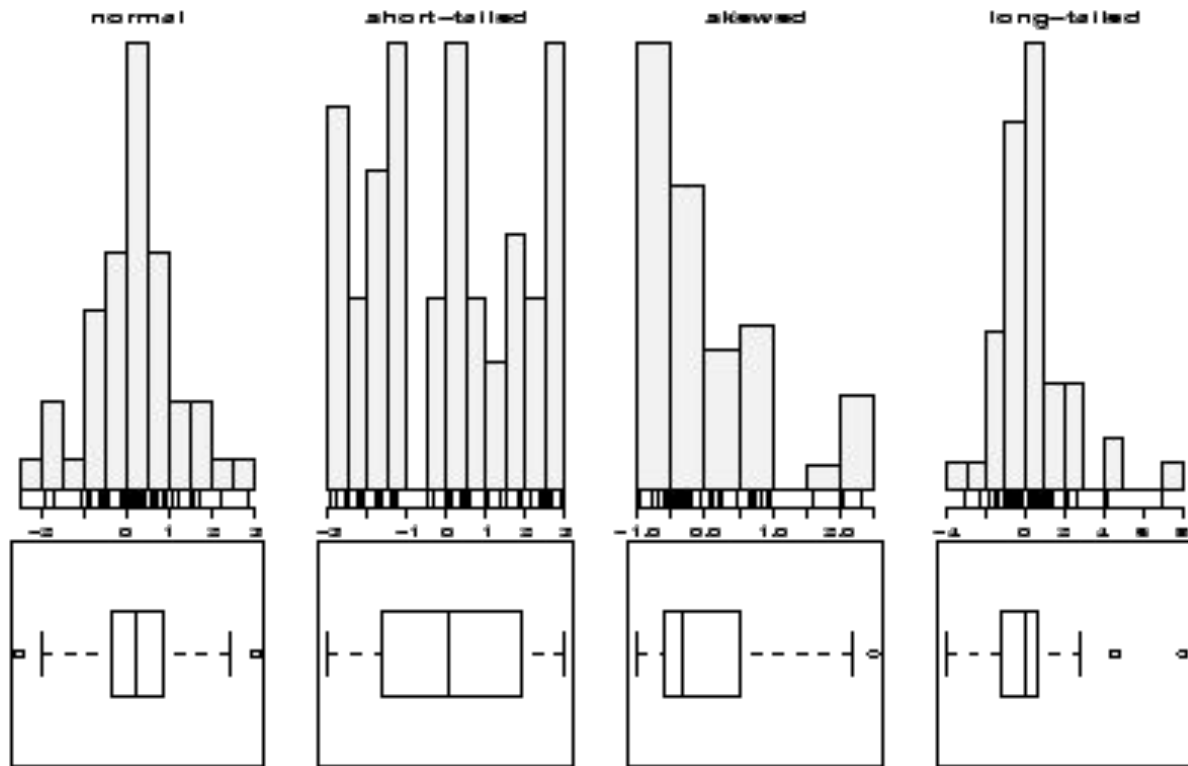
Boxplot? histogram?

Which should I use?

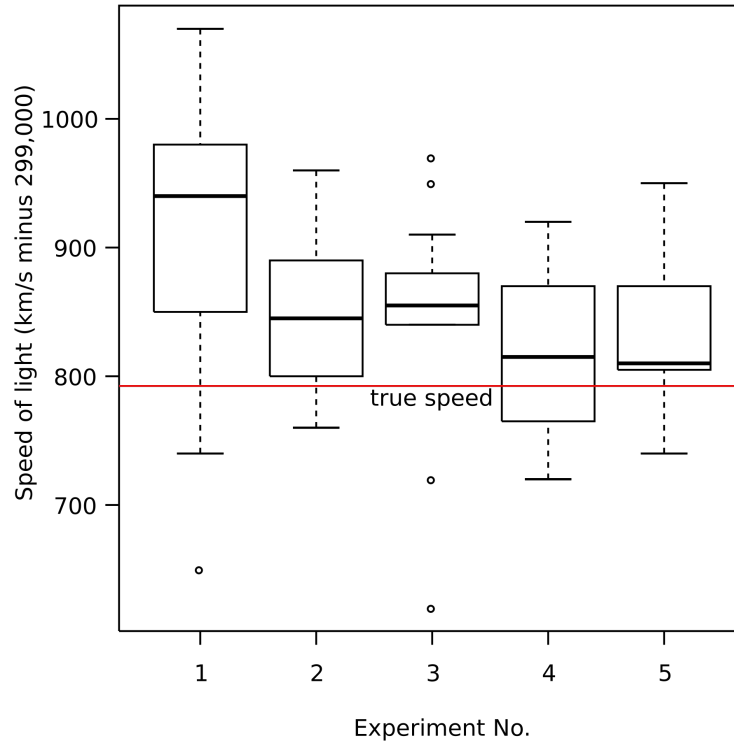
It's easy to combine a boxplot with scatterplot



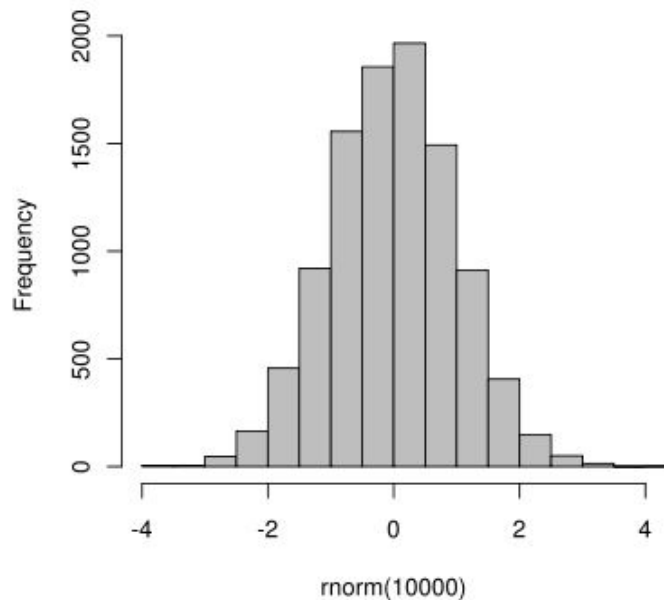
Boxplots hide details



Boxplots are succinct and useful for comparison

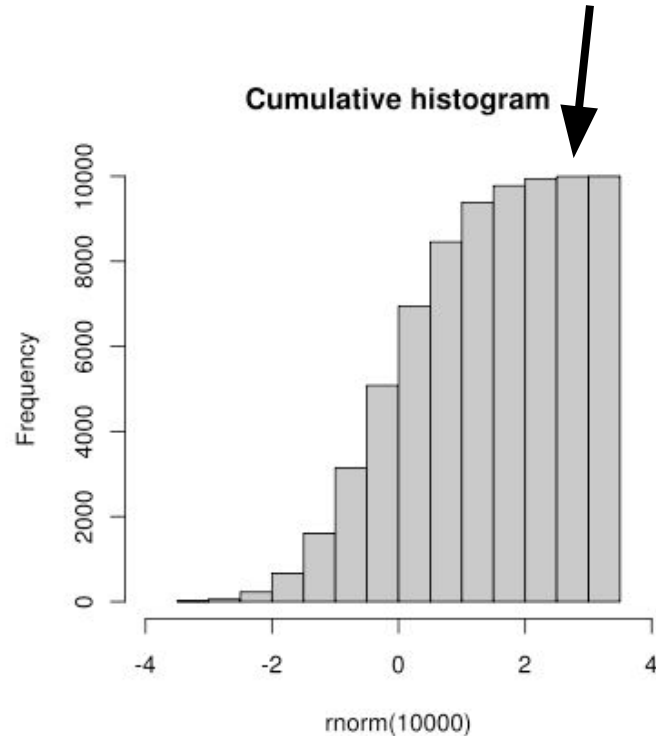
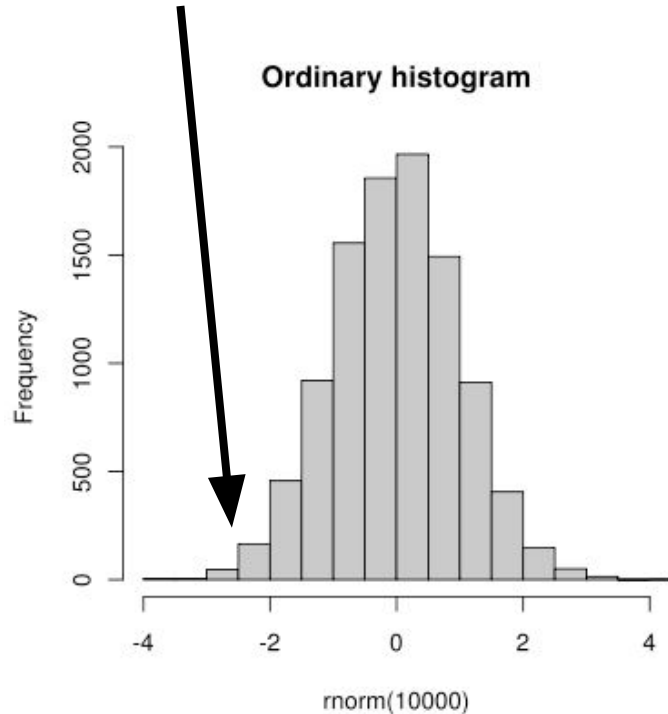


Q: Can we design a **new way to draw histogram** so that you can easily spot the location of median (or even any arbitrary percentile points)?

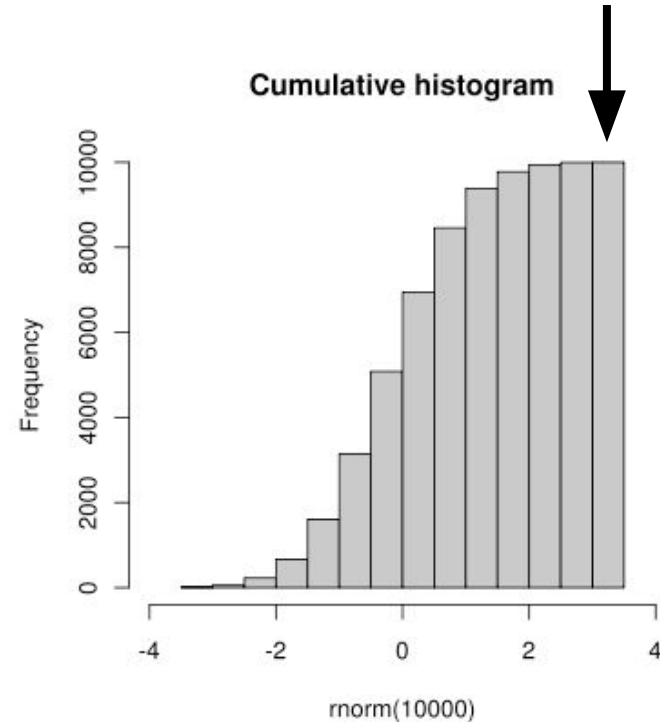
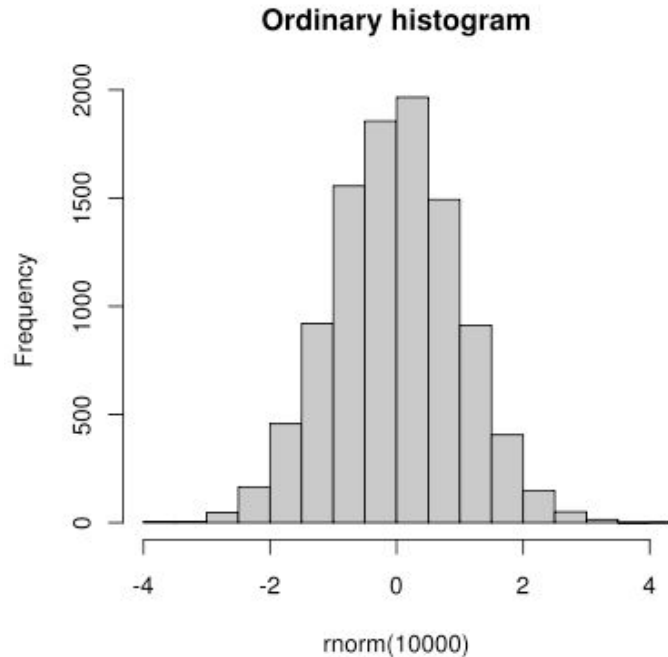


"The number of data points within this bin"

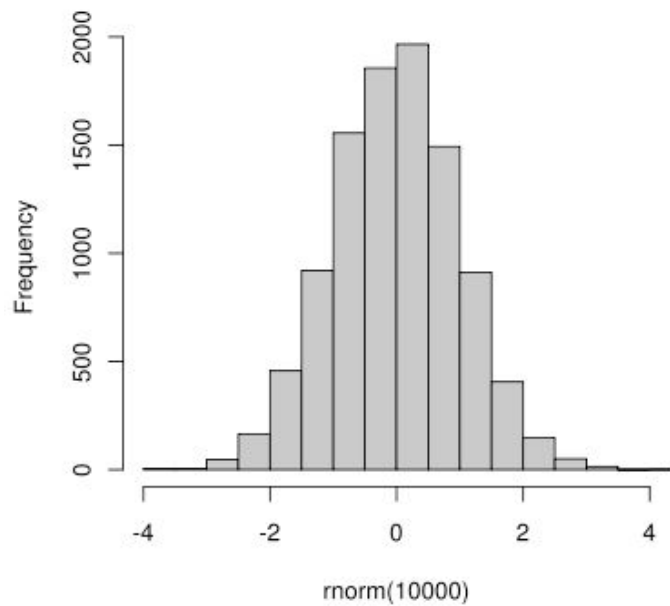
"The number of data points we have seen so far"



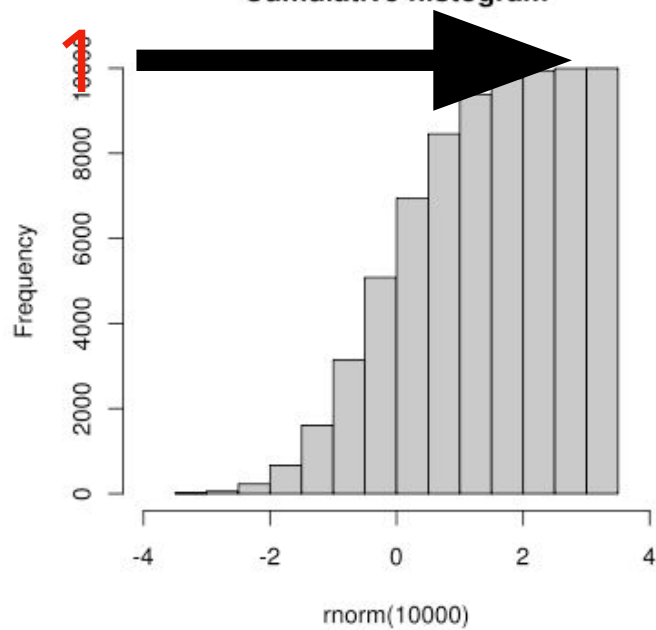
The total number of
data points

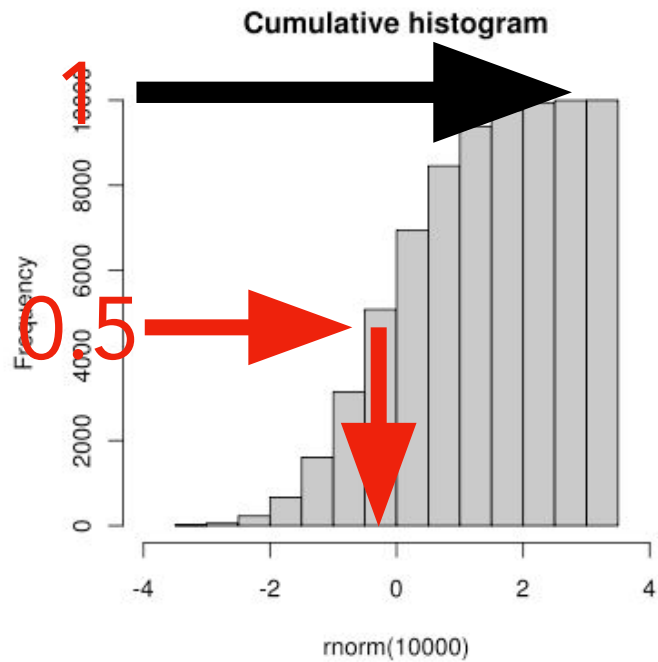
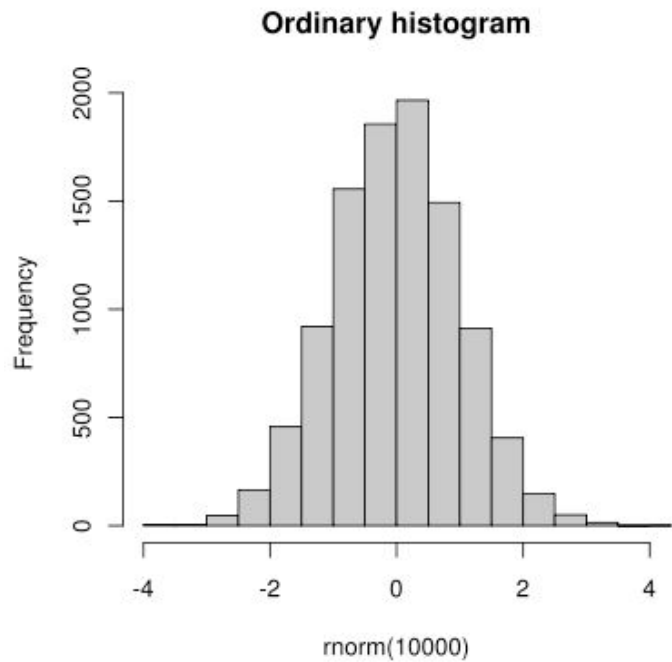


Ordinary histogram



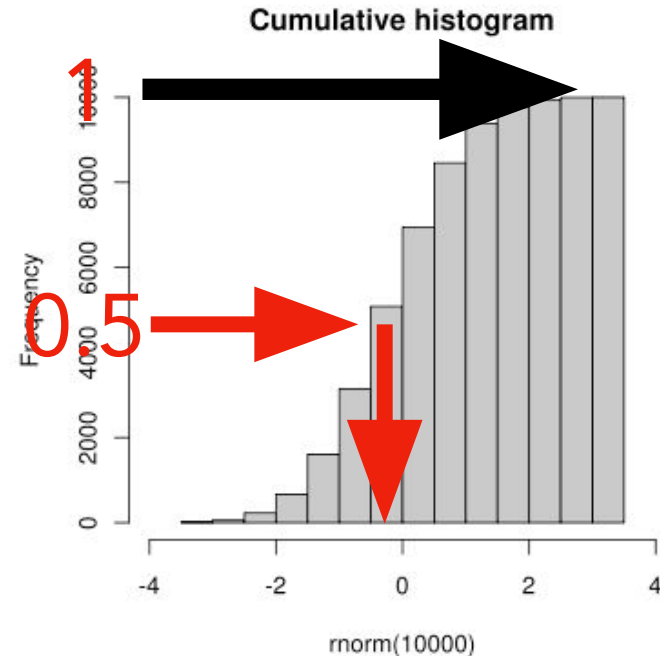
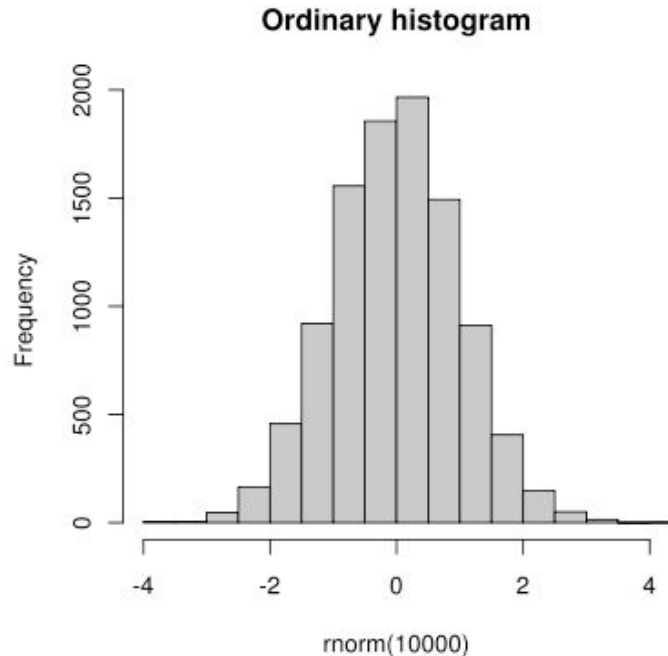
Cumulative histogram

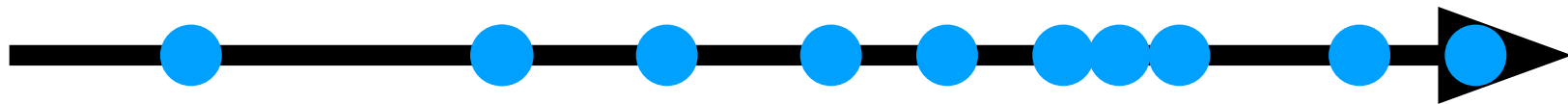


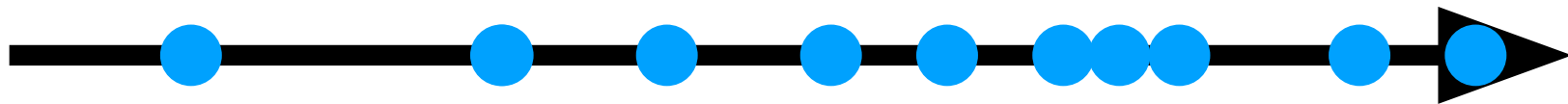


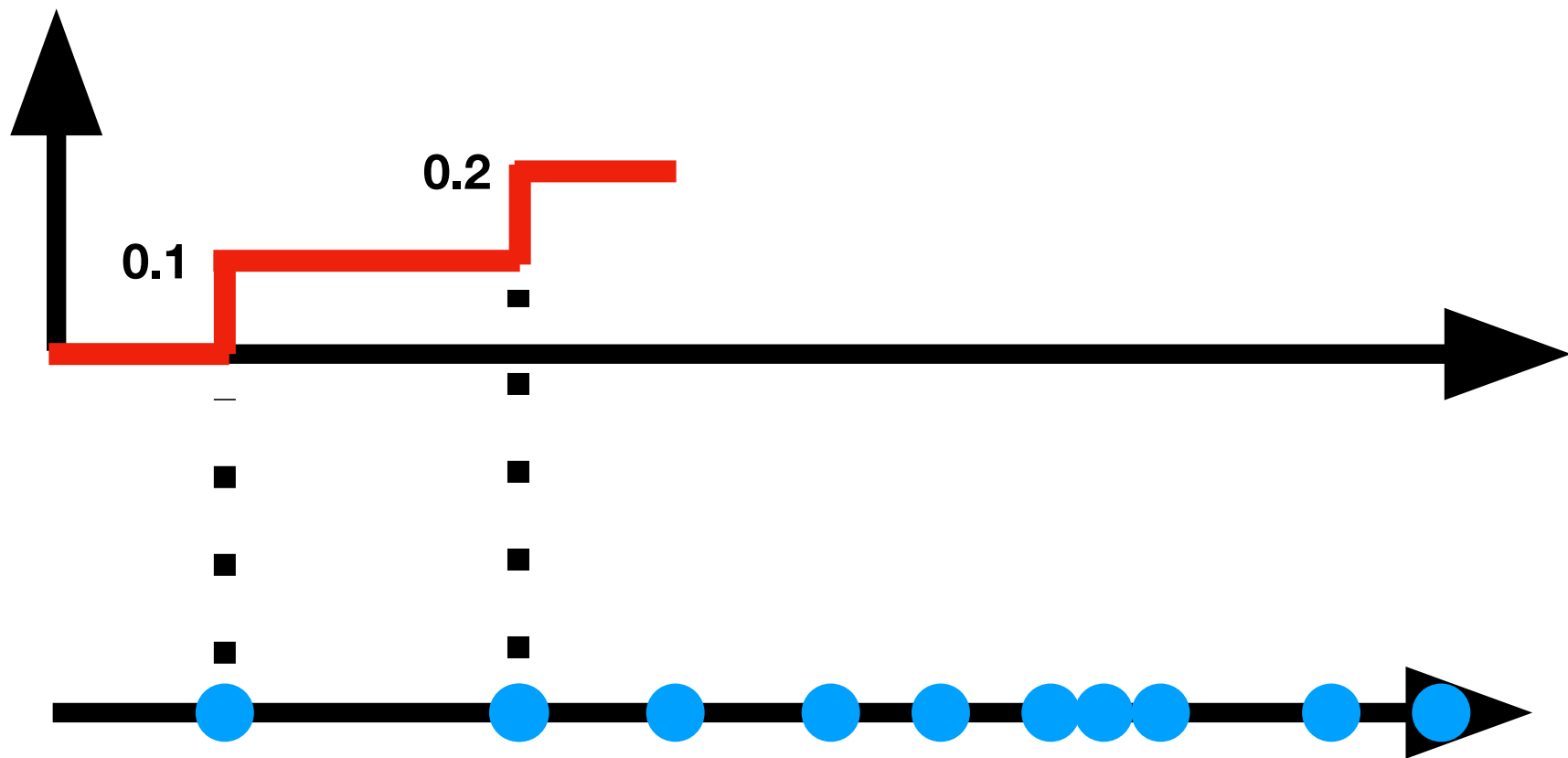
Can we improve it even further???

(to make it have a higher resolution)









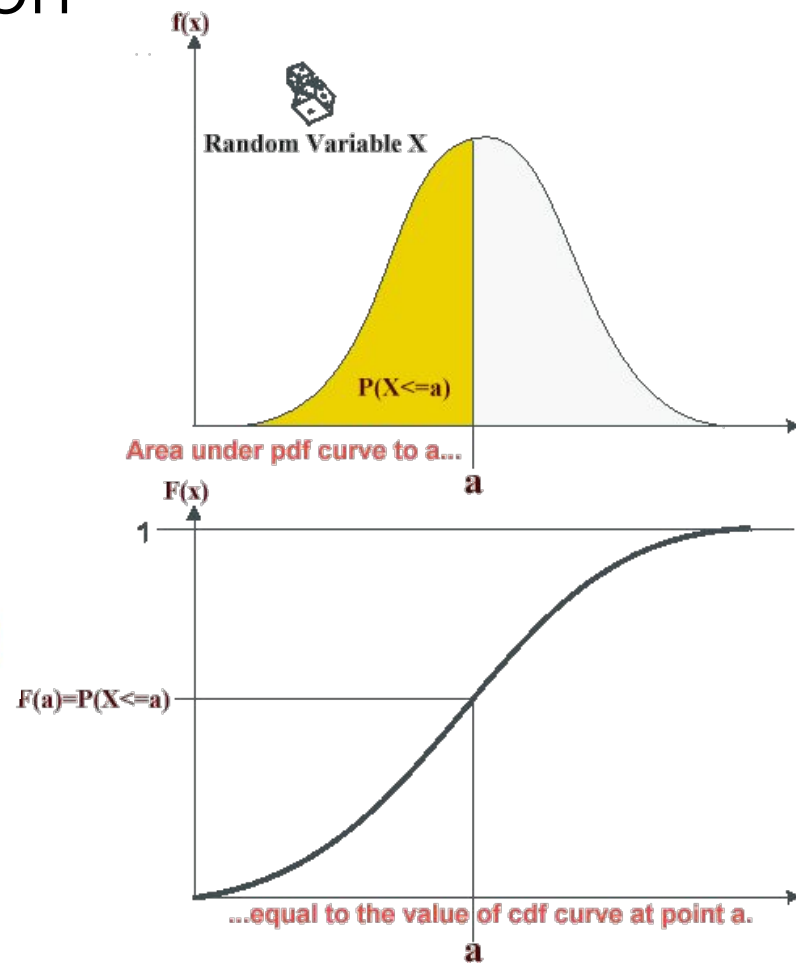
CDF (Cumulative Distribution Function)

$$F_X(x) = P(X \leq x)$$

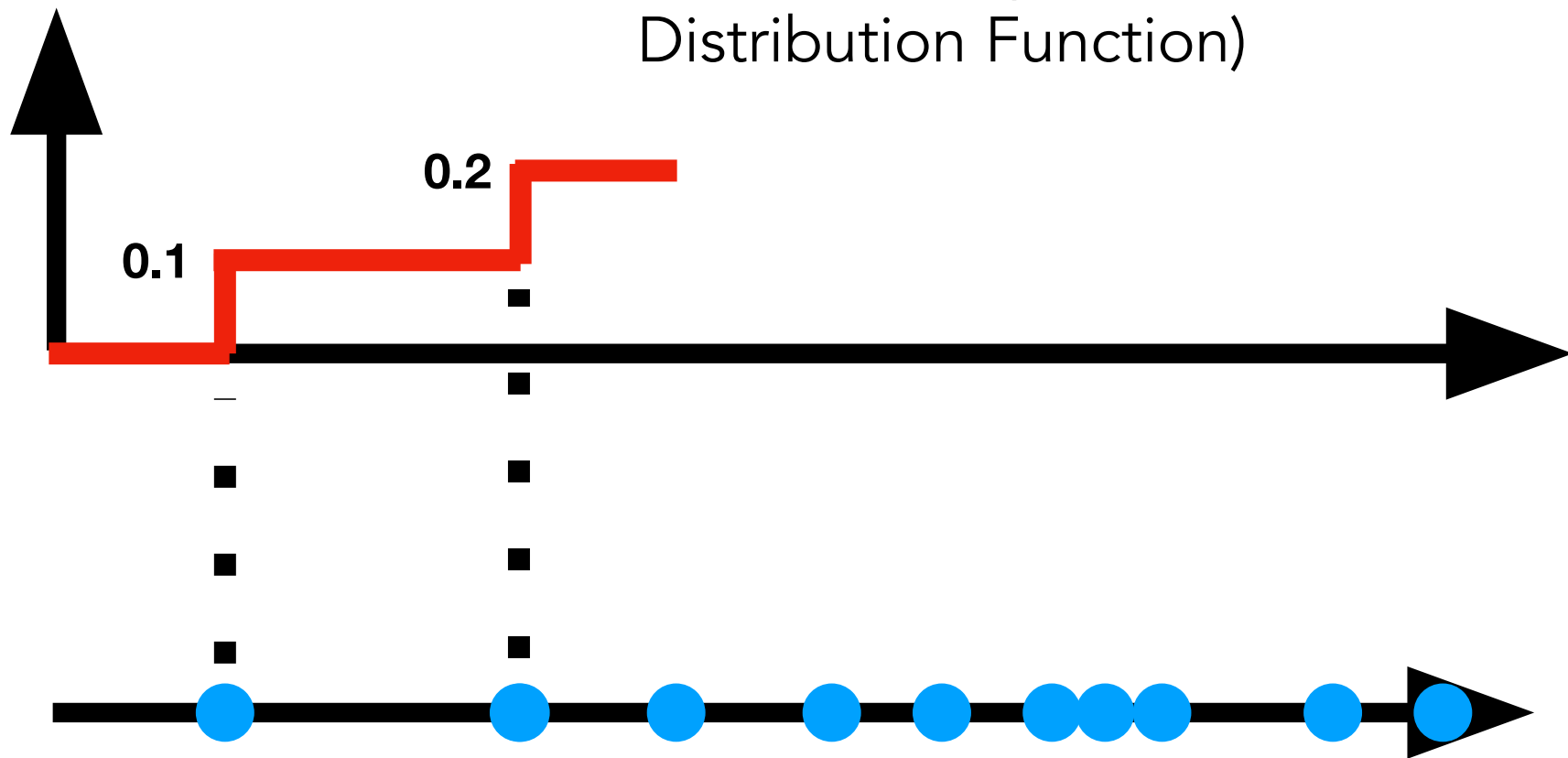
$$F_X(x) = \int_{-\infty}^x \underset{\substack{\uparrow \\ \text{Prob. density function (PDF)}}}{f_X(t)} dt$$

Prob. density function (PDF)

We don't need any bins!!



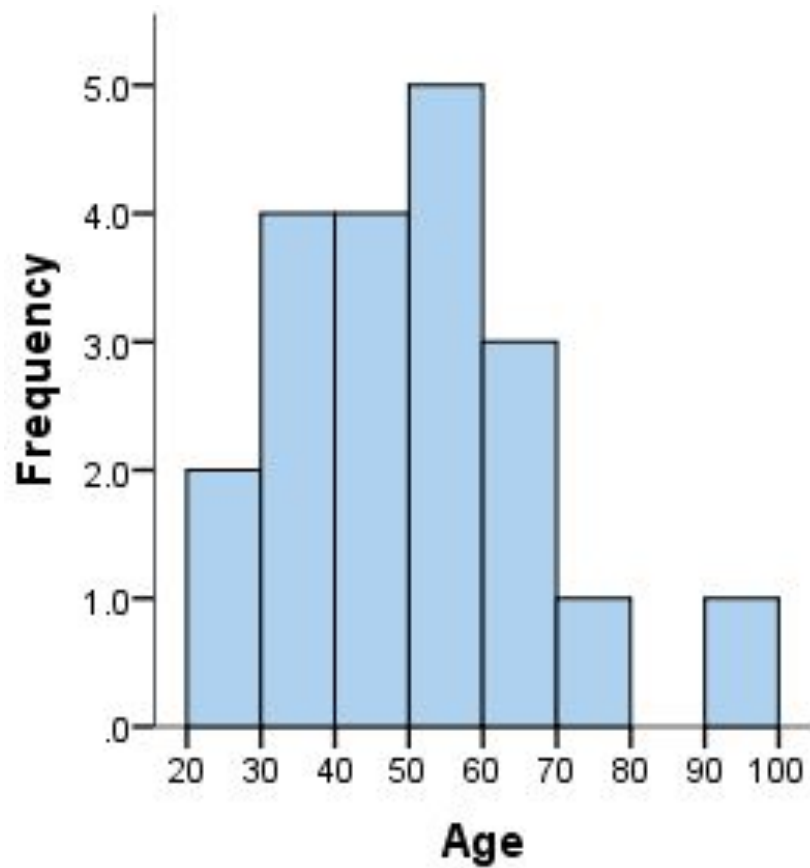
Empirical CDF (Empirical Cumulative Distribution Function)



Let's draw the empirical CDF for the
following data

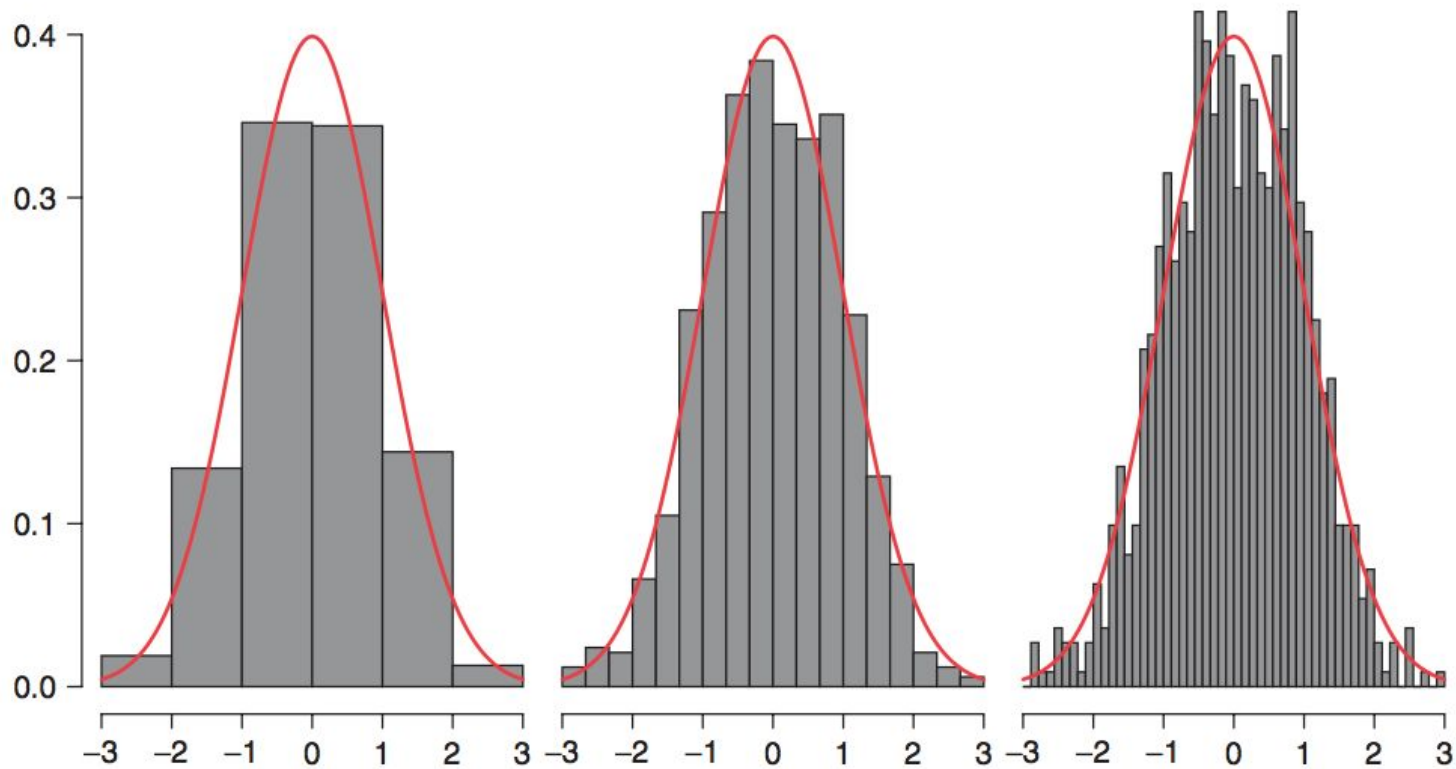
3, 5, 5, 8, 10

Estimation

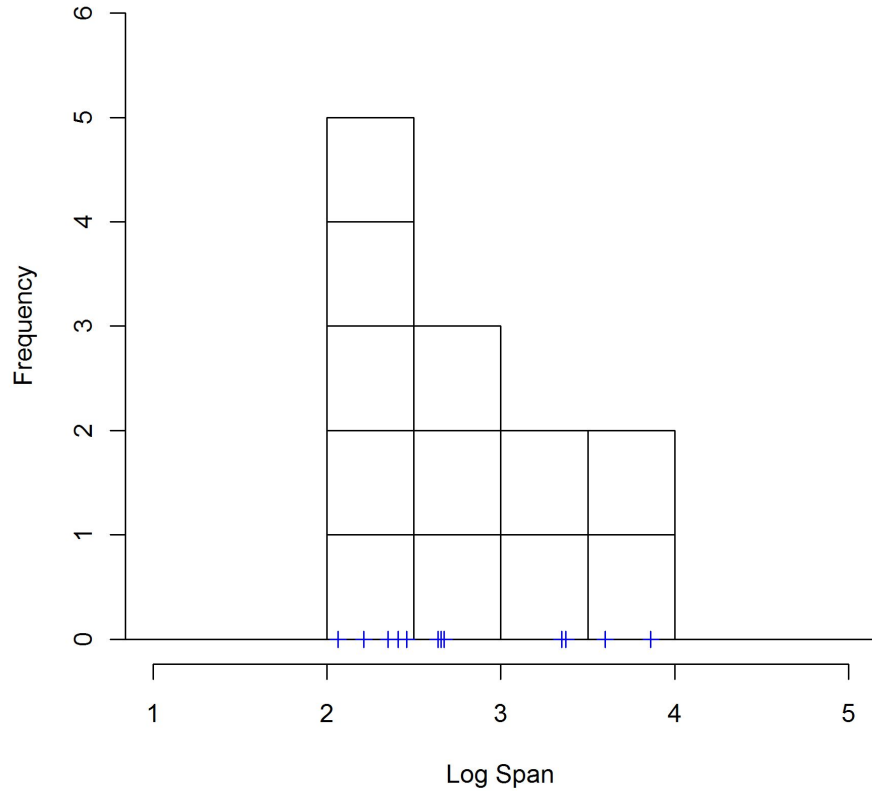


Simple

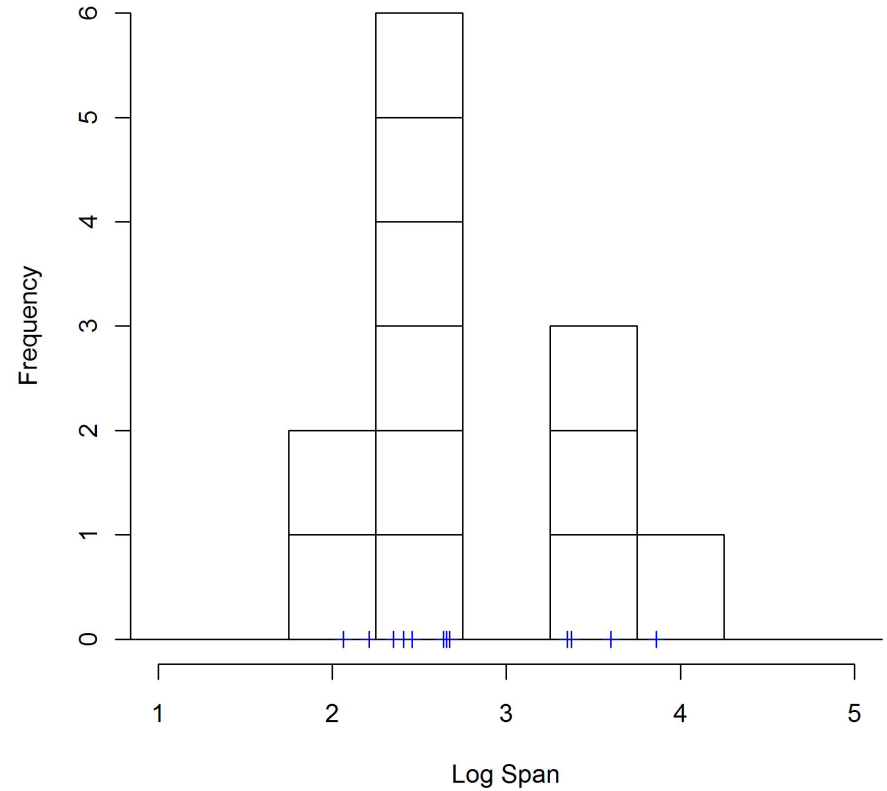
Transparent

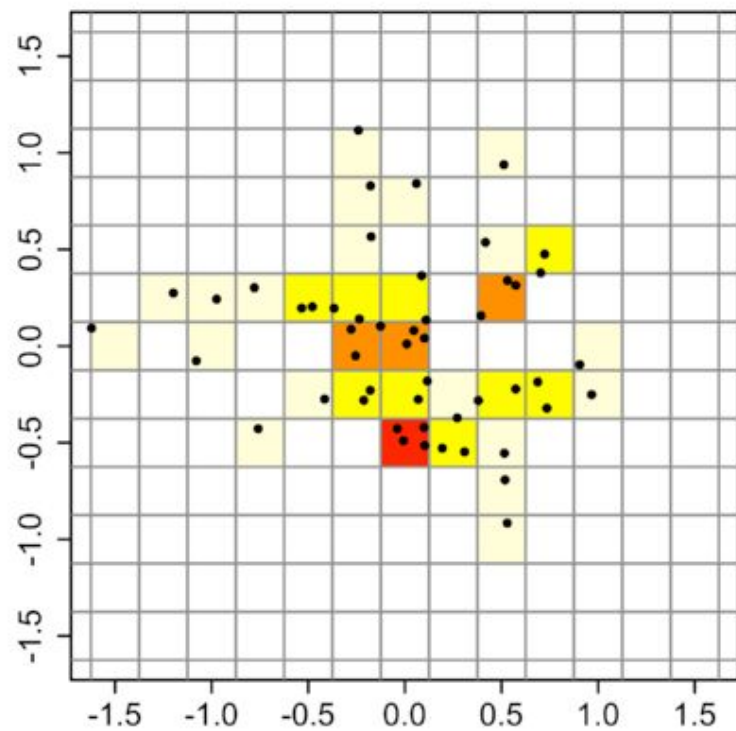
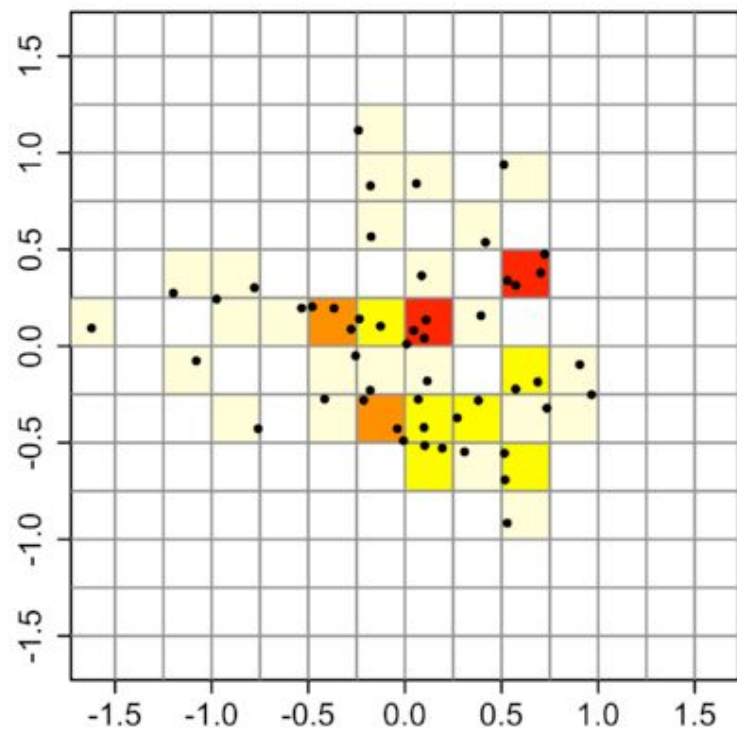


**Histogram with breaks at n.0 and n.5
binwidth=0.5**



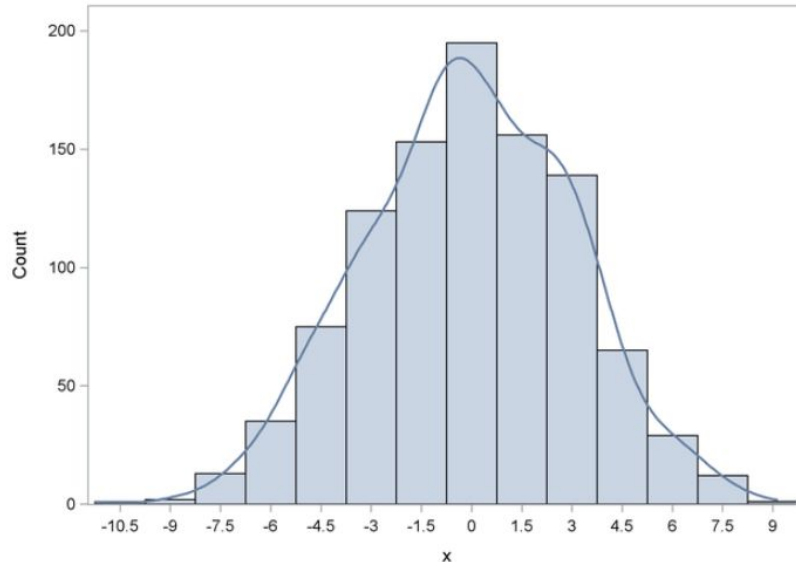
**Histogram with breaks at n.25 and n.75
binwidth=0.5**





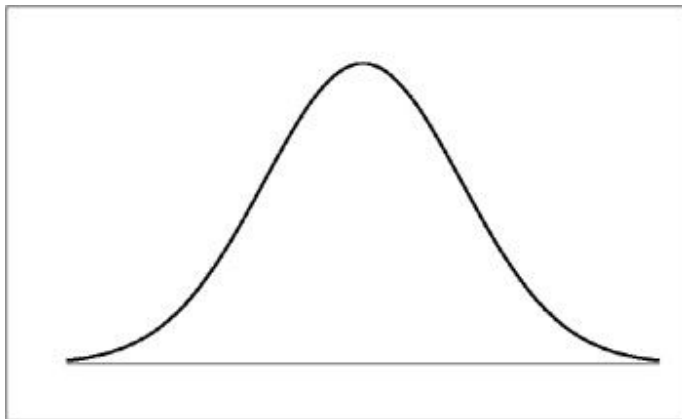
A histogram is a “model” that infers the underlying distribution of the data.

Can we directly infer the underlying (smooth) distribution from data? How can we do that?

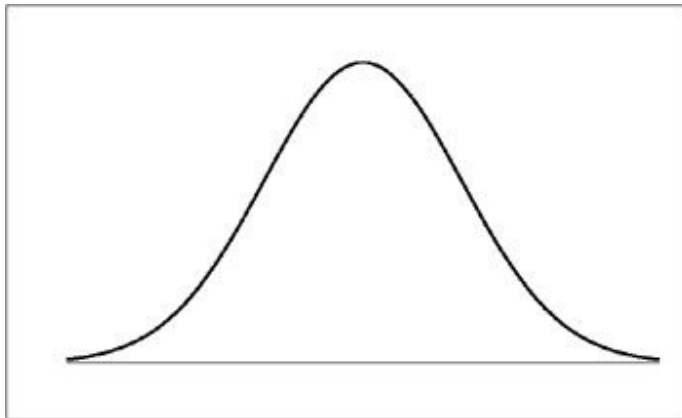


1. Parametric approach

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

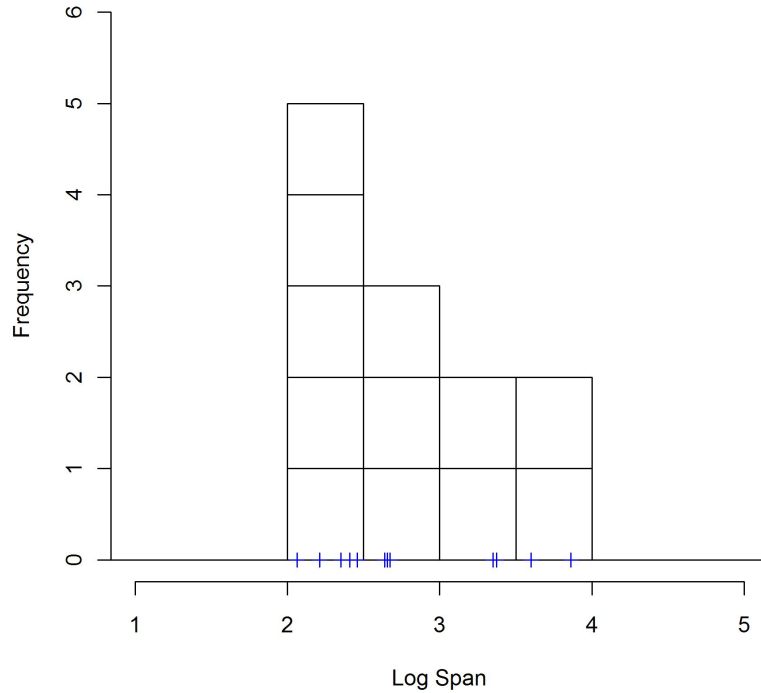


$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

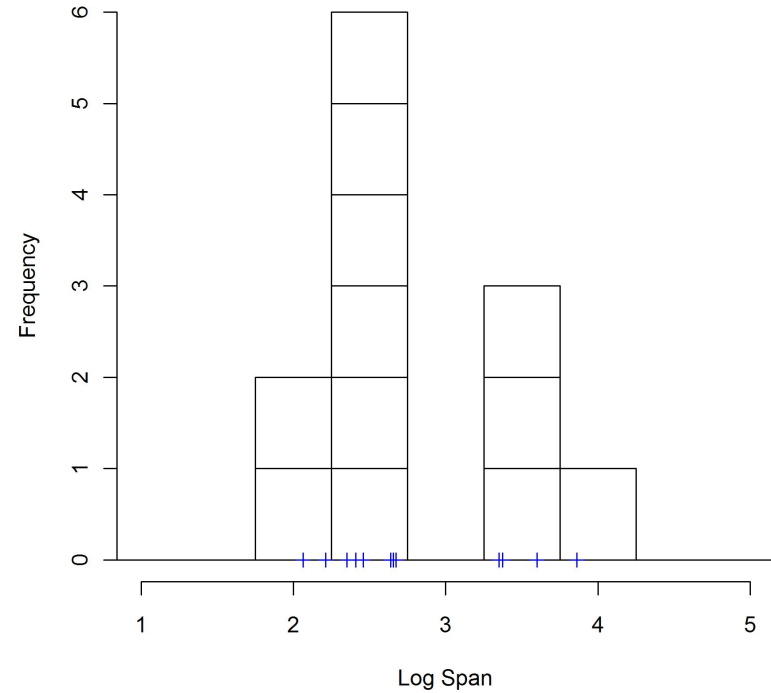


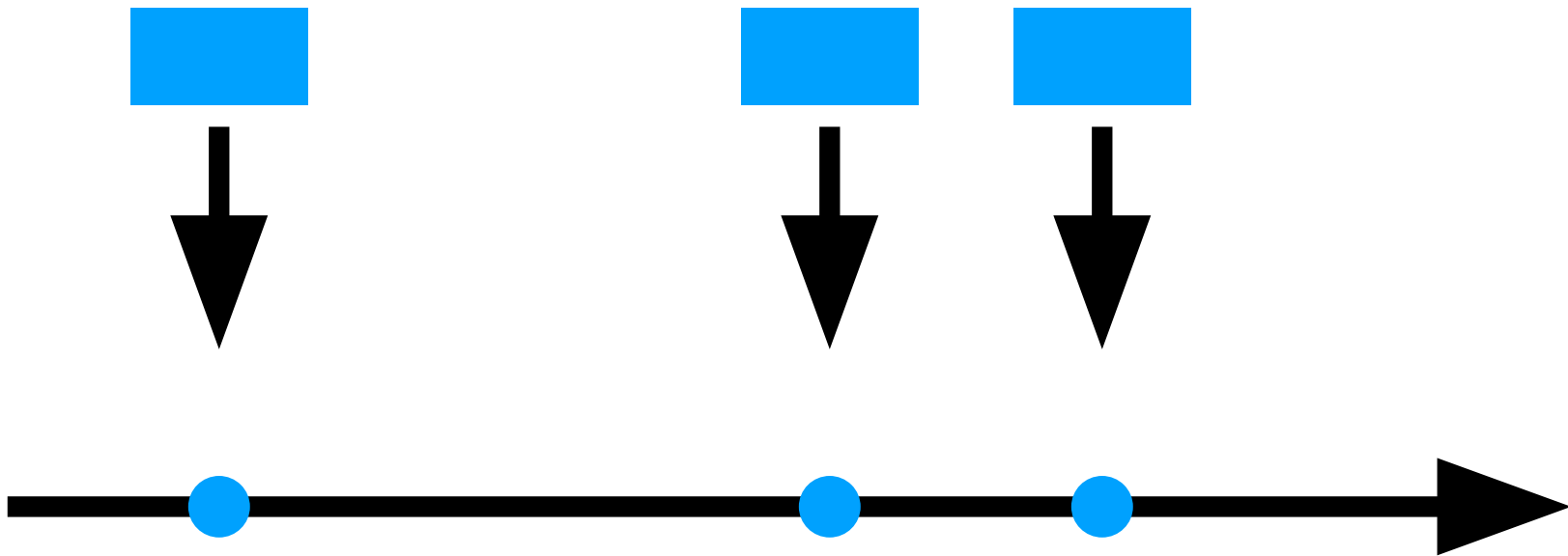
2. Non-parametric approach

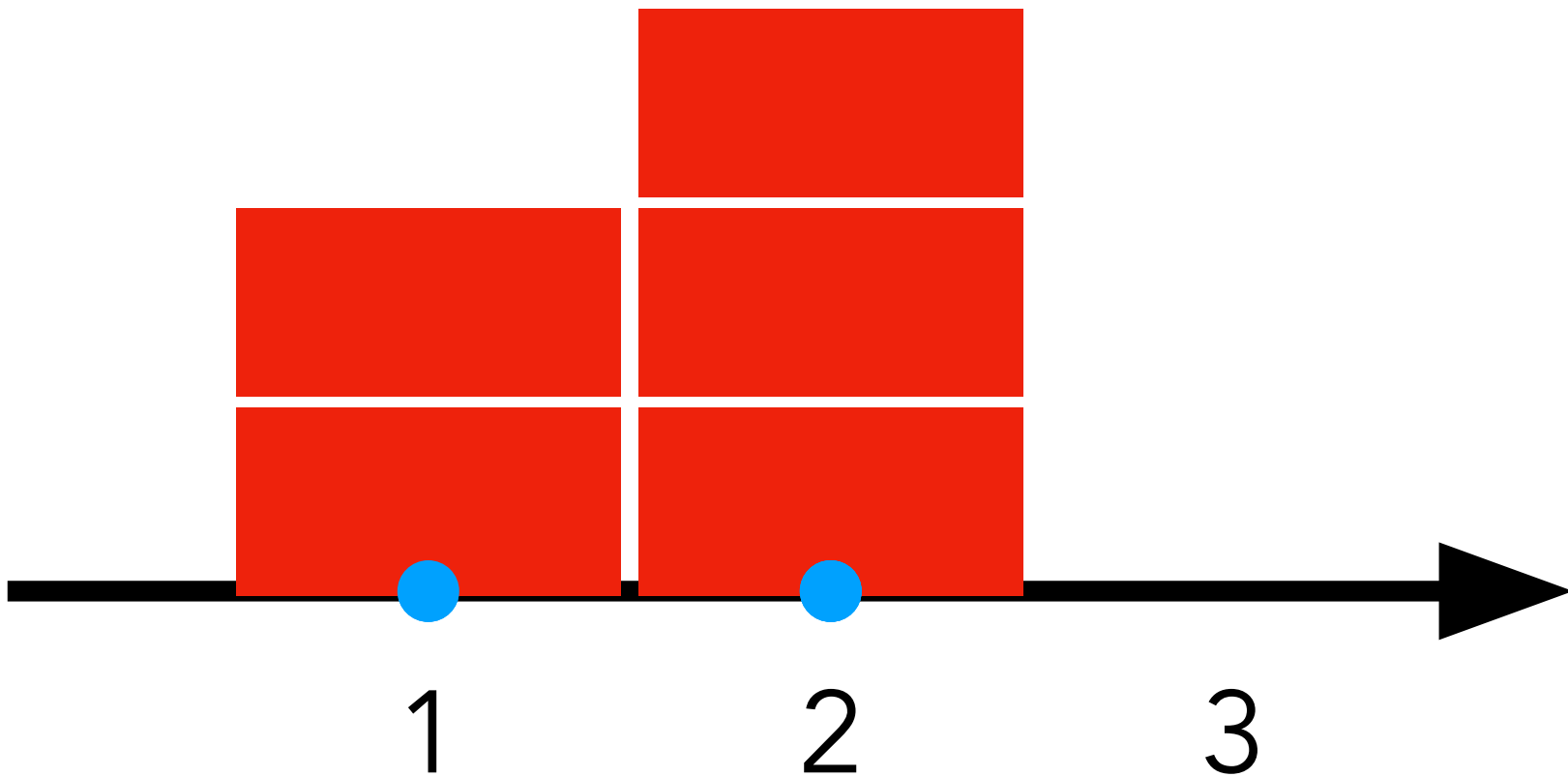
**Histogram with breaks at n.0 and n.5
binwidth=0.5**

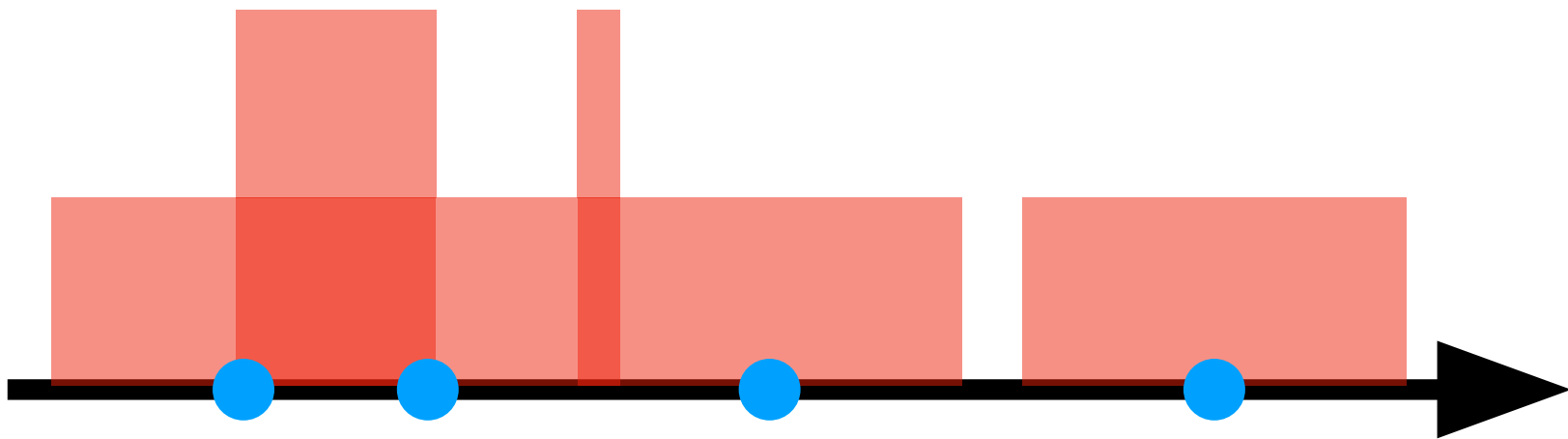


**Histogram with breaks at n.25 and n.75
binwidth=0.5**

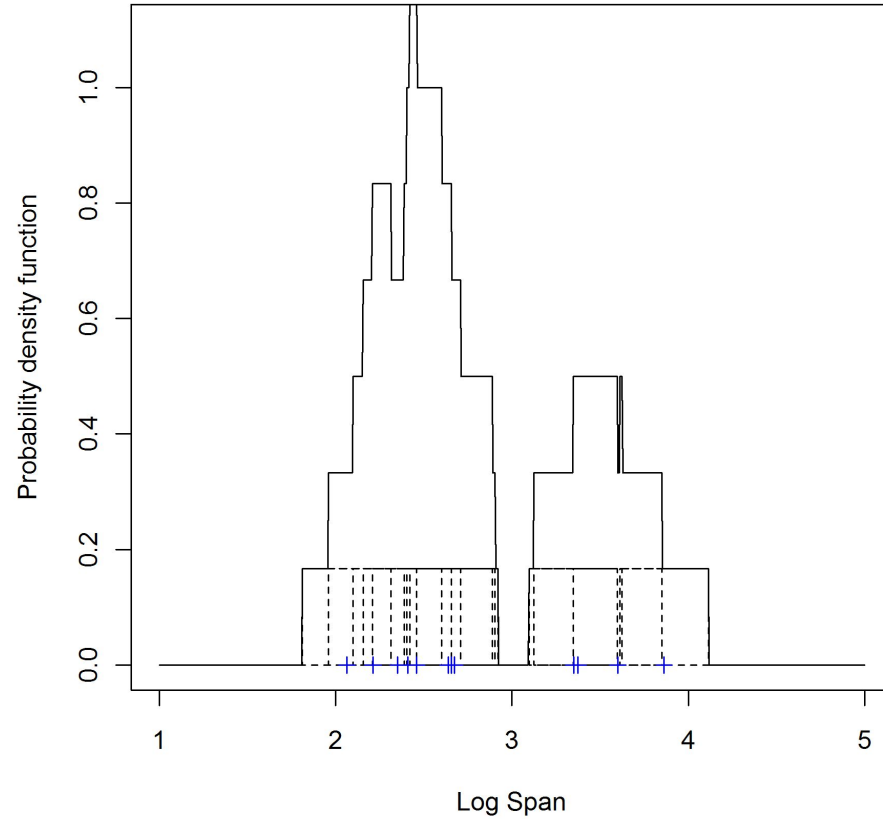




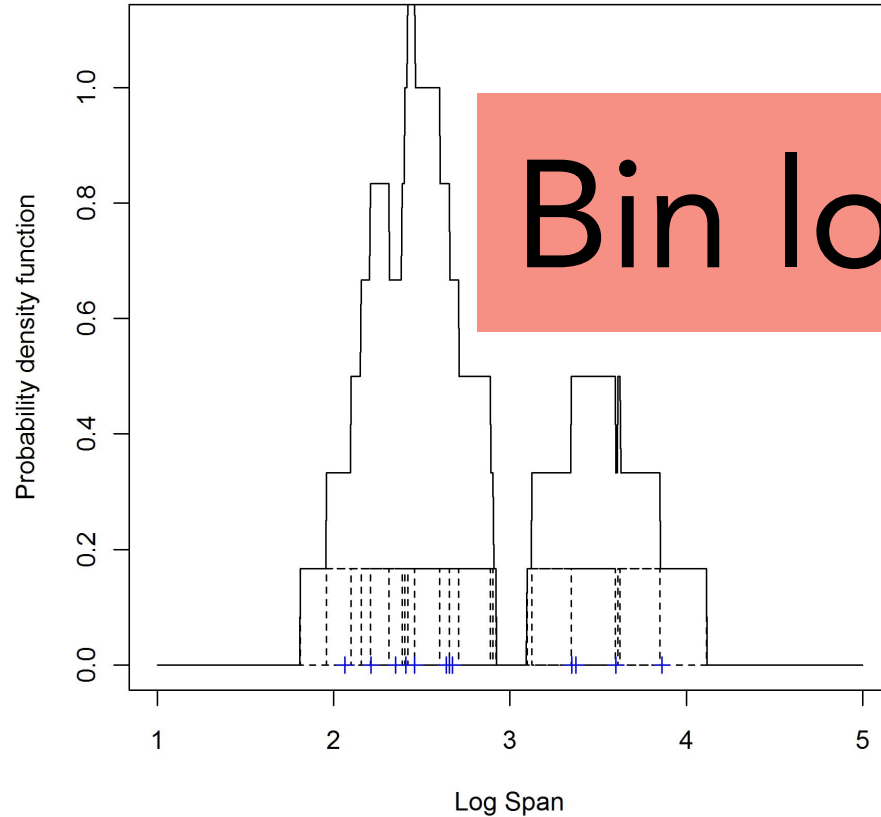


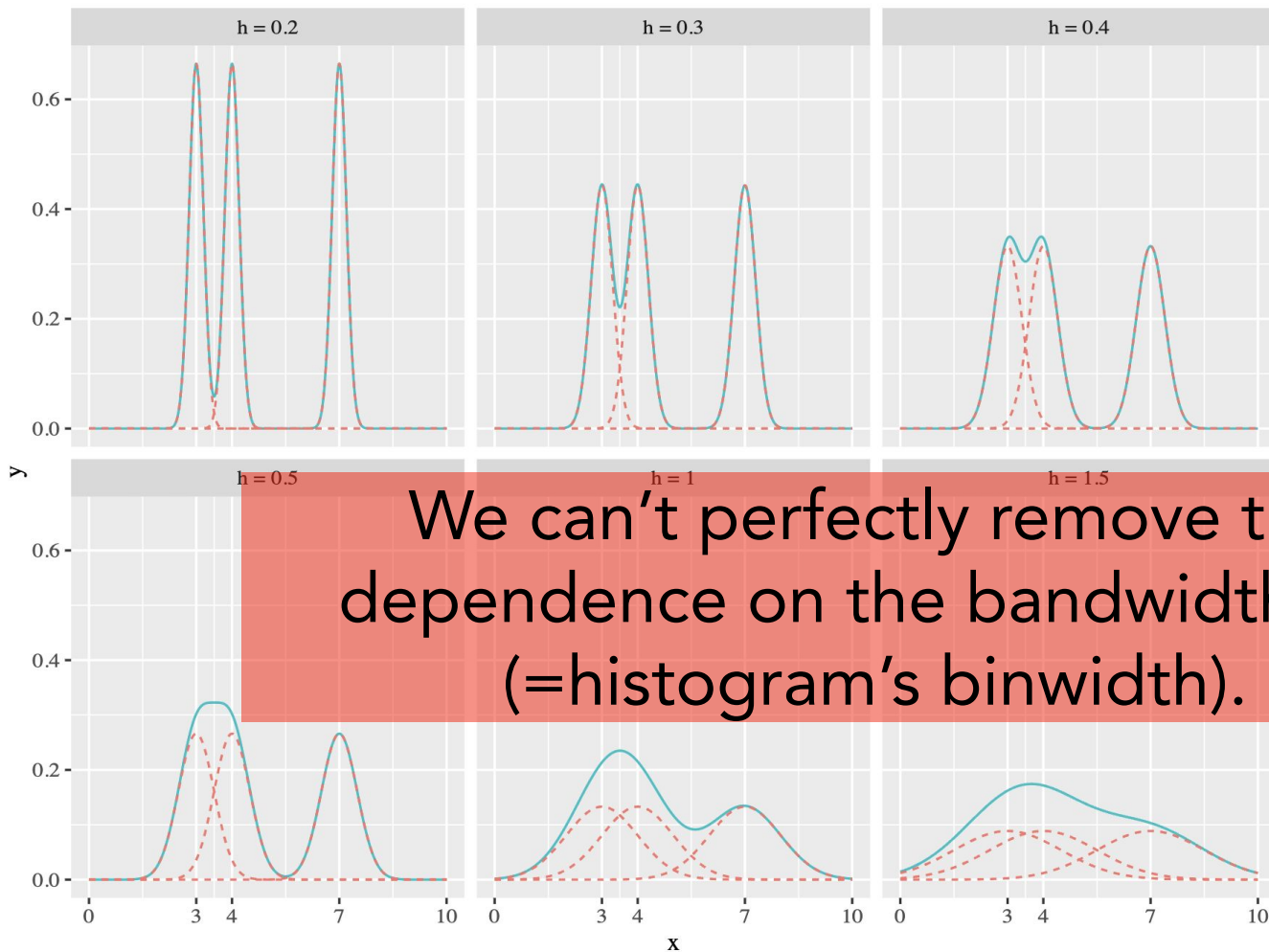


'Histogram' with blocks centred over data points

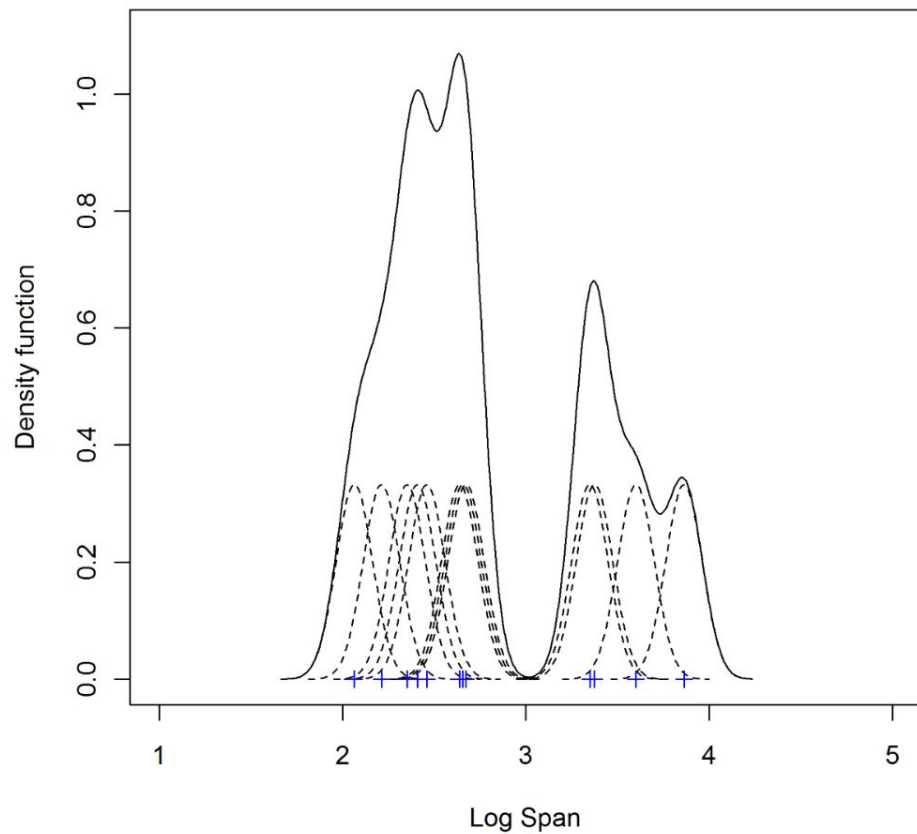


'Histogram' with blocks centred over data points

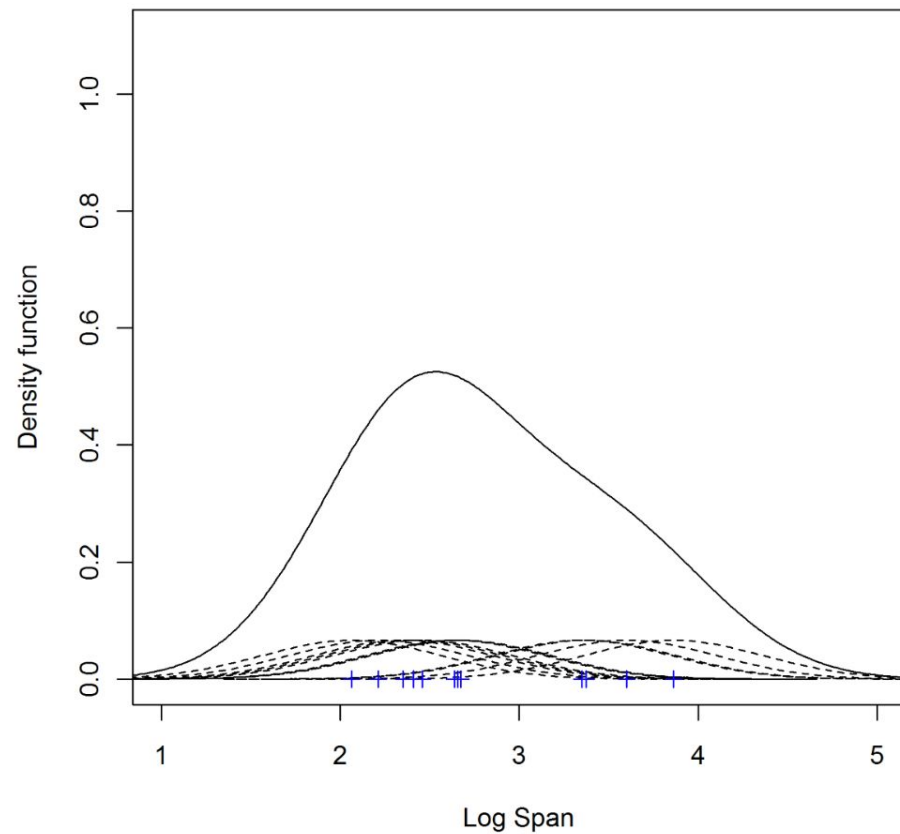




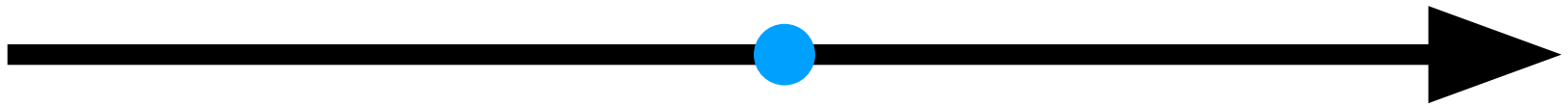
Undersmoothed

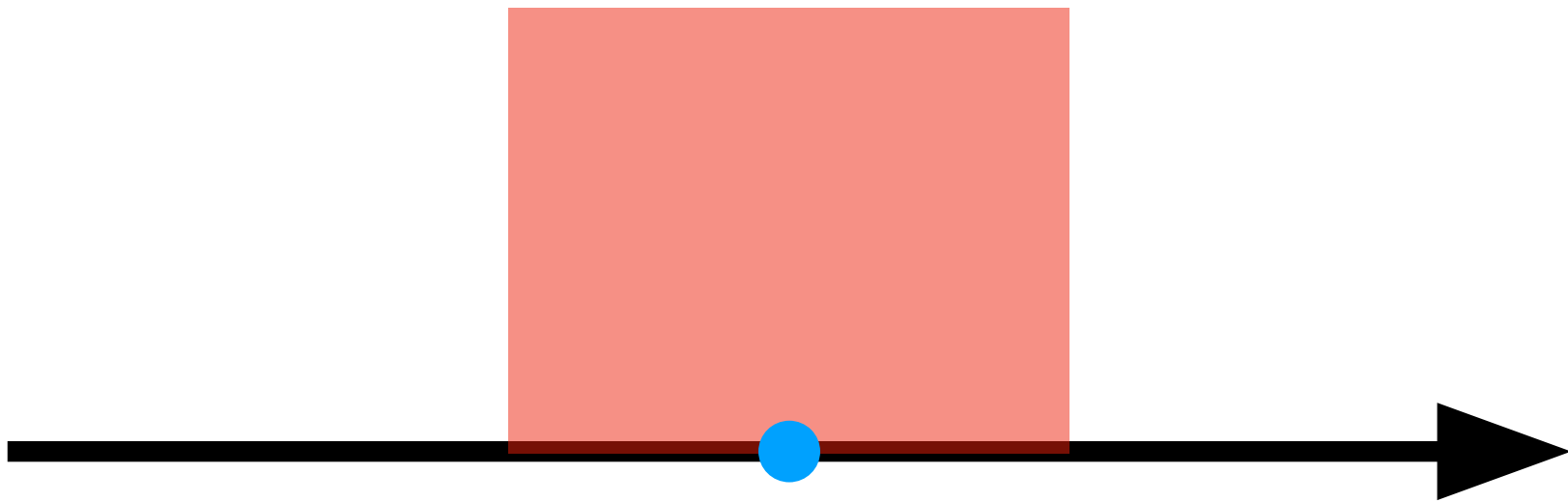


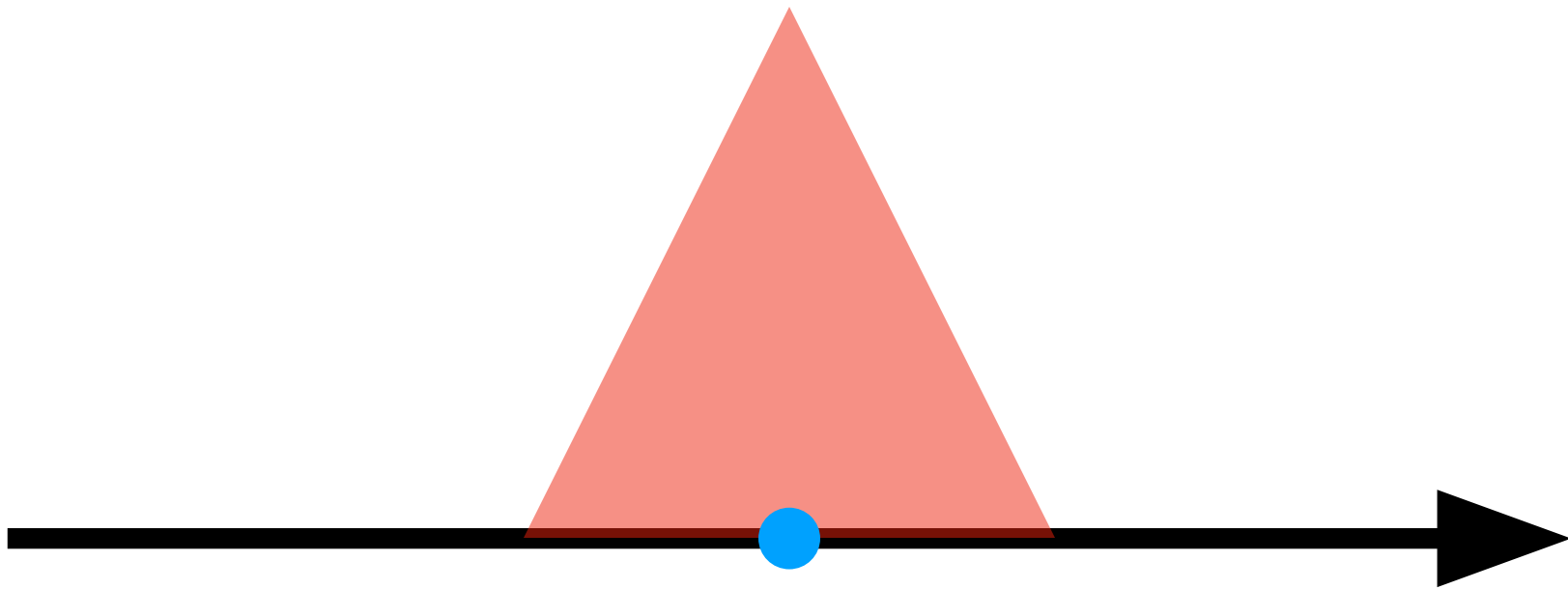
Oversmoothed

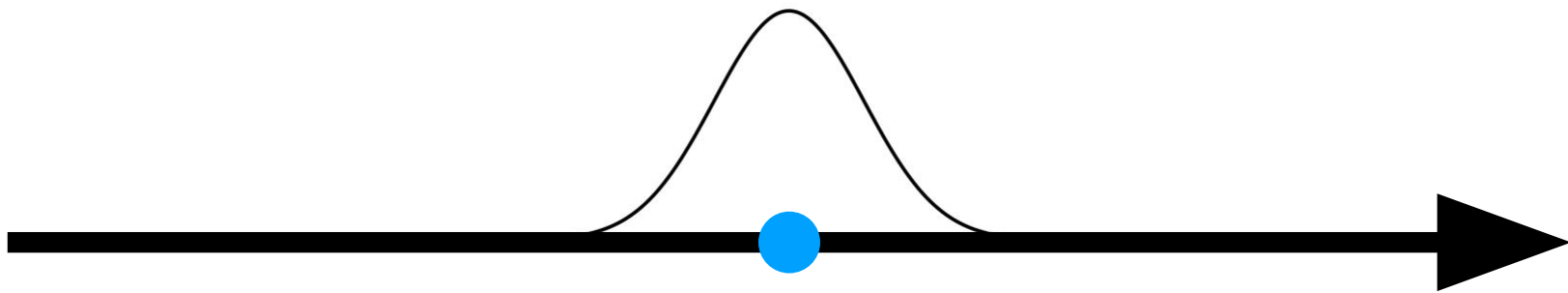


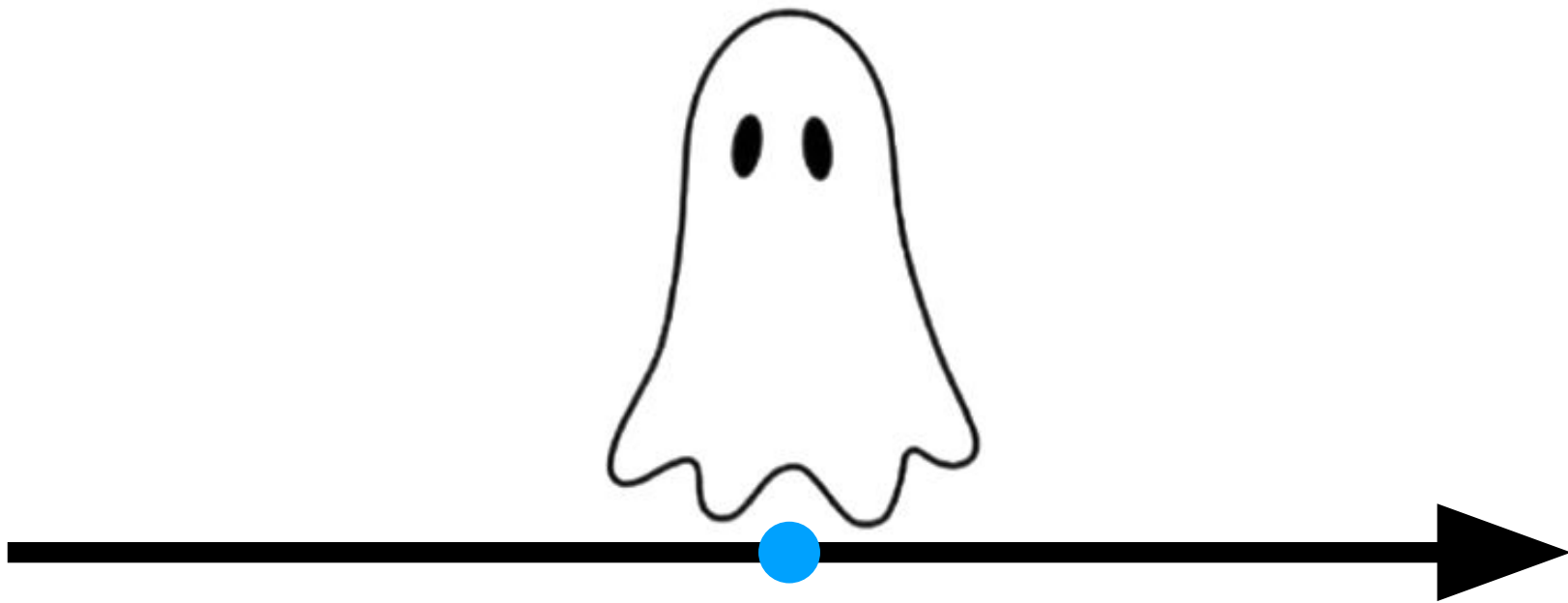
Less arbitrary
More data-driven



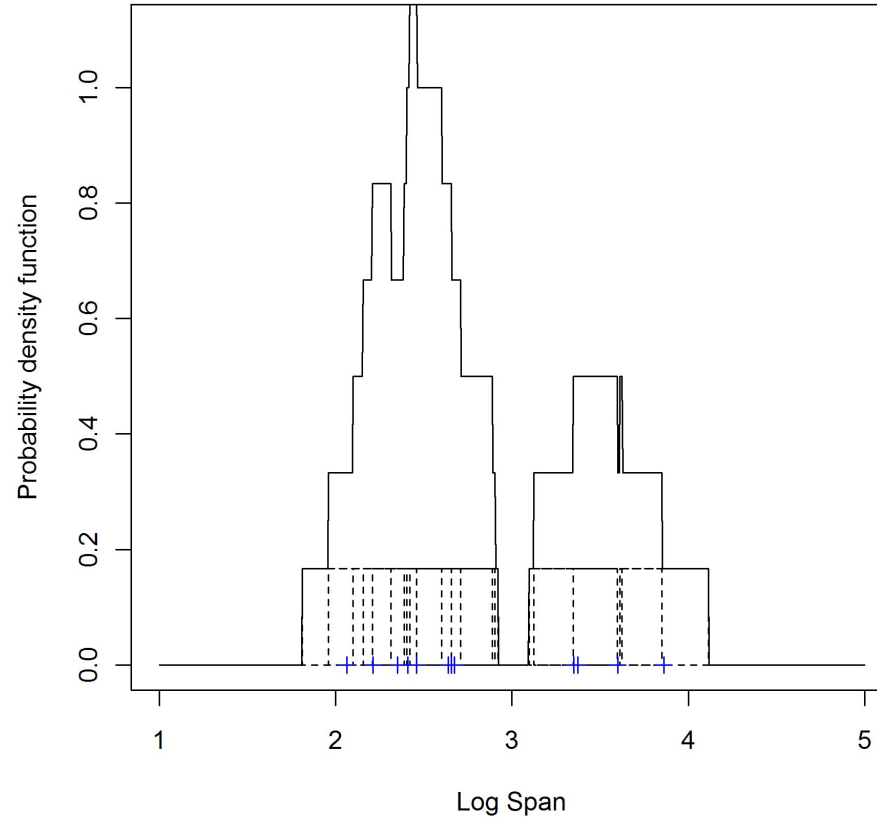


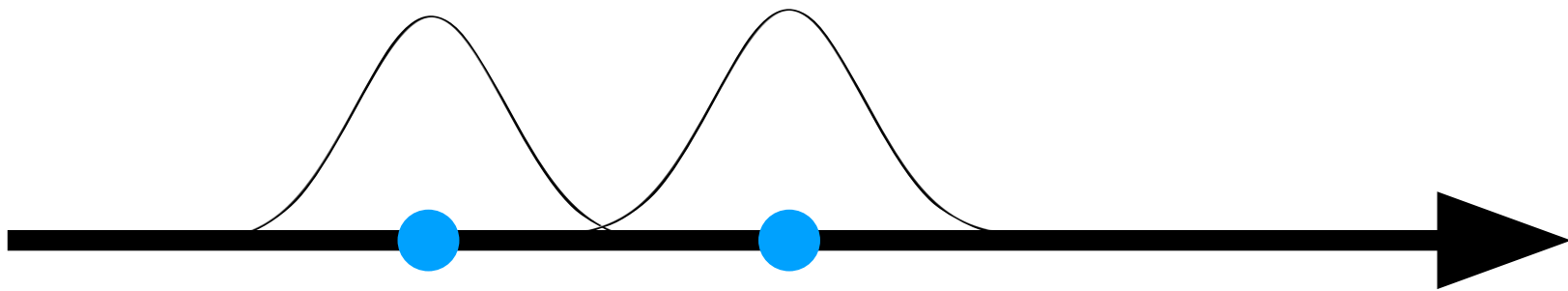




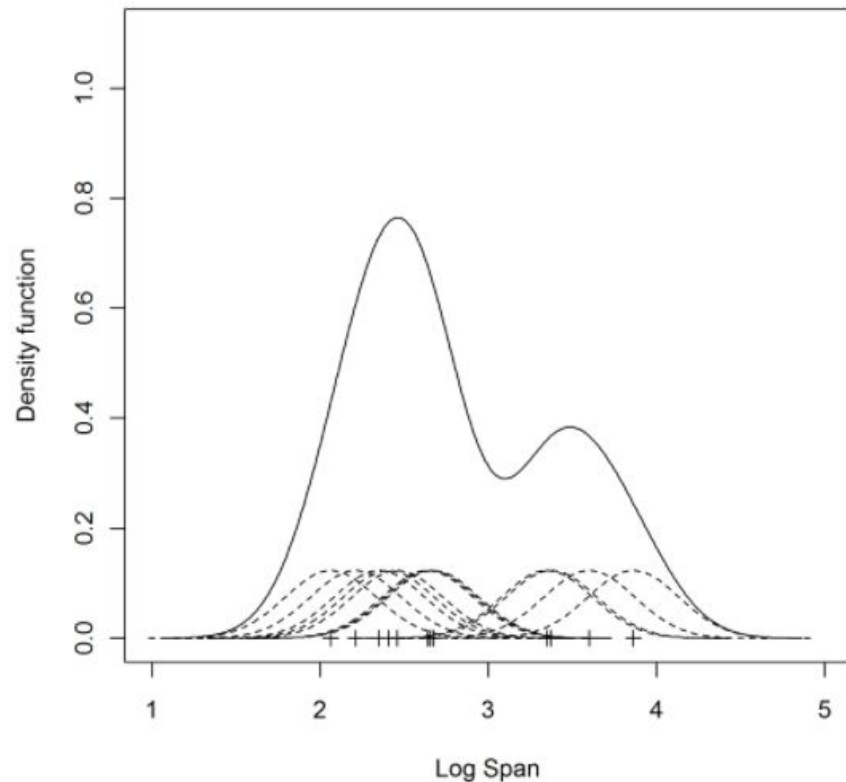


'Histogram' with blocks centred over data points





Kernel Density Estimation (KDE)



"Kernels"

