

# Doing Bayesian Data Analysis

## §5.1~5.2

Chao-En Yu

National Tsing Hua University

# Reallocation of Credibility

- On a typical day at your location, what is the probability that it is cloudy?
- Suppose you are told it is raining, now what is the probability that it is cloudy?
- We can be pretty sure that
$$p(\text{cloudy}) < p(\text{cloudy}|\text{raining}).$$

# Human Reasoning

- Suppose instead you are told that everyone outside is wearing sunglasses, what is the probability that it is cloudy?
- Most likely, it is true that
$$p(\text{cloudy}) > p(\text{cloudy}|\text{sunglass}).$$
- Notice how we human beings have reasoned in this example.

# Bayes' Rule

- Bayes' rule is merely the mathematical relation between the prior allocation of credibility and the posterior reallocation of credibility conditional on data.

- From the definition of conditional probability,

$$p(c|r) = \frac{p(r, c)}{p(r)}$$

i.e., the probability that  $r$  and  $c$  happen together relative to the probability that  $r$  happens at all.



Only a condition

# Derivation

- Notice that by definition we also have

$$p(r|c) = \frac{p(r, c)}{p(c)} \quad \leftarrow \text{Marginal probability}$$

- Multiplying both sides by  $p(r)$  and  $p(c)$ , respectively, gives

$$\begin{aligned} p(c|r)p(r) &= p(r, c) = p(r|c)p(c) \\ \Rightarrow p(c|r) &= \frac{p(r|c)p(c)}{p(r)} \end{aligned}$$

which is called *Bayes' rule*.

# Bayes' Rule

- We can further re-write the denominator in terms of  $p(r|c)$ :

$$p(c|r) = \frac{p(r|c)p(c)}{\sum_{c^*} p(r|c^*)p(c^*)}$$

- The  $c$  in the numerator is a specific fixed value, whereas the  $c^*$  in the denominator is a variable that takes on all possible values.
- $p(c|r)$  is called posterior;  $p(c)$  is prior;  $p(r|c)$  is likelihood;  $p(r)$  is marginal likelihood (evidence).

# Intuition of Bayes Rule

- Without knowing anything about a person's eye color, all we believe about hair colors is expressed by the marginal probabilities at the bottom.

Eye color	Hair color				Marginal (Eye color)
	Black	Brunette	Red	Blond	
Brown	0.11	0.20	0.04	0.01	0.37
Blue	0.03	0.14	0.03	0.16	0.36
Hazel	0.03	0.09	0.02	0.02	0.16
Green	0.01	0.05	0.02	0.03	0.11
Marginal (hair color)	0.18	0.48	0.12	0.21	1.0

# Intuition of Bayes Rule

- If we are told that the selected person's eyes are blue, we can focus our attention on the Blue row.
- Notice that we have gone from the “prior” (marginal) beliefs about hair color before knowing eye color, to the “posterior” (conditional) beliefs.

Eye color	Hair color				Marginal (Eye color)
	Black	Brunette	Red	Blond	
Blue	$0.03/0.36$ $= 0.08$	$0.14/0.36$ $= 0.39$	$0.03/0.36$ $= 0.08$	$0.16/0.36$ $= 0.45$	$0.36/0.36 = 1.0$

# Prosecutor's Fallacy

- A cab was involved in a hit and run accident at night. Two cab companies, the Green and the Blue, operate in the city. 85% of the cabs in the city are Green and 15% are Blue.
- A witness identified the cab as Blue. The court tested the reliability of the witness under the same circumstances that existed on the night of the accident and concluded that the witness correctly identified each one of the two colors 80% of the time and failed 20% of the time.
- What is the probability that the cab involved in the accident was blue rather green?

$$p(\theta = \ddot{\theta} | T = +) = ?$$

Test result	Disease		Marginal (test result)
	$\theta = \ddot{\theta}$ (present)	$\theta = \ddot{\theta}$ (absent)	
$T = +$	$p(+ \ddot{\theta}) p(\ddot{\theta})$ $= 0.99 \cdot 0.001$	$p(+ \ddot{\theta}) p(\ddot{\theta})$ $= 0.05 \cdot (1 - 0.001)$	$\sum_{\theta} p(+ \theta) p(\theta)$
$T = -$	$p(- \ddot{\theta}) p(\ddot{\theta})$ $= (1 - 0.99) \cdot 0.001$	$p(- \ddot{\theta}) p(\ddot{\theta})$ $= (1 - 0.05) \cdot (1 - 0.001)$	$\sum_{\theta} p(- \theta) p(\theta)$
Marginal (disease)	$p(\ddot{\theta}) = 0.001$	$p(\ddot{\theta}) = 1 - 0.001$	1.0

# Hit Rate vs False Alarm

Test result	Disease	
	$\theta = \ddot{\smile}$ (present)	$\theta = \ddot{\smiley}$ (absent)
$T = +$	True Positive (Hit Rate)	False Positive (False Alarm)
$T = -$	False Negative	True Negative

# Parametric Model

- The key application of Bayes' rule is when the row variable represents data values and the column variable represents parameter values.
- A Bayesian model specifies  $p(\text{data values}|\text{parameter values})$  along with the prior,  $p(\text{parameter values})$
- We use Bayes' rule to convert that to know  $p(\text{parameter values}|\text{data values})$  ← How strongly we should believe in the various parameter values given the data.

# Parameter Inference

- The posterior  $p(\theta|D)$  is obtained by conditionalizing on the row with the observed data value, and that operation is Baye's rule.

	Model parameter			
Data	...	$\theta$ value	...	Marginal
⋮		⋮		⋮
$D$ value	...	$p(D, \theta) = p(D \theta) p(\theta)$	...	$p(D) = \sum_{\theta^*} p(D \theta^*) p(\theta^*)$
⋮		⋮		⋮
Marginal	...	$p(\theta)$	...	

# Caveats

- For continuous parameters, the only change in Bayes' rule is that the marginal likelihood changes from the sum to an integral:

$$p(D) = \int d\theta^* p(D|\theta^*)p(\theta^*)$$

- Suppose we observe some data  $D$  and  $D'$ , does our final belief depend on whether we update with  $D$  first and  $D'$  second, or update with  $D'$  first and  $D$  second?

# Statistical Independence

- When  $p(x, y) = p(x) * p(y), \forall x, y$ , we say that  $x$  and  $y$  are statistically independent.
- By Bayes' rule,  $x$  and  $y$  are statistically independent if  $p(y | x) = p(y), \forall x, y$ .
- Bayesian prefers the second definition for it relates probability to an addition of information, though independence is often a convenient assumption.

# Data-Order Invariance

- The answer is: It depends!
- When the data probabilities are independent,  $p(D, D' | \theta) = p(D | \theta) \cdot p(D' | \theta)$ , then the order of updating has no effect of the final posterior.
- ★ Intuition: If the likelihood function has no dependence on data ordering, then the posterior shouldn't have any dependence on data ordering.

# Proof of Order-Invariance

- Equivalent to assume that datum is equally representative of the underlying process, regardless of when the datum was observed, and regardless of any data observed before or after.

$$\begin{aligned} p(\theta|D', D) &= \frac{p(D', D|\theta) p(\theta)}{\sum_{\theta^*} p(D', D|\theta^*) p(\theta^*)} && \text{Bayes' rule} \\ &= \frac{p(D'|\theta)p(D|\theta) p(\theta)}{\sum_{\theta^*} p(D'|\theta^*)p(D|\theta^*) p(\theta^*)} && \text{by assumption of independence} \\ &= \frac{p(D|\theta)p(D'|\theta) p(\theta)}{\sum_{\theta^*} p(D|\theta^*)p(D'|\theta^*) p(\theta^*)} && \text{multiplication is commutative} \\ &= p(\theta|D, D') && \text{Bayes' rule} \end{aligned}$$

# Doing Bayesian Data Analysis

§4.3~4.4

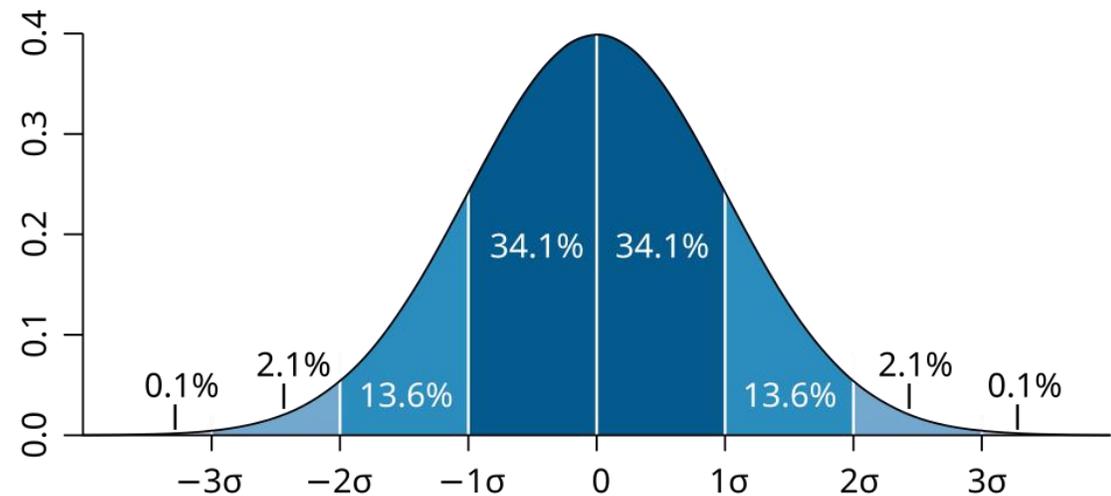
Chao-En Yu

National Tsing Hua University

# Probability Distributions

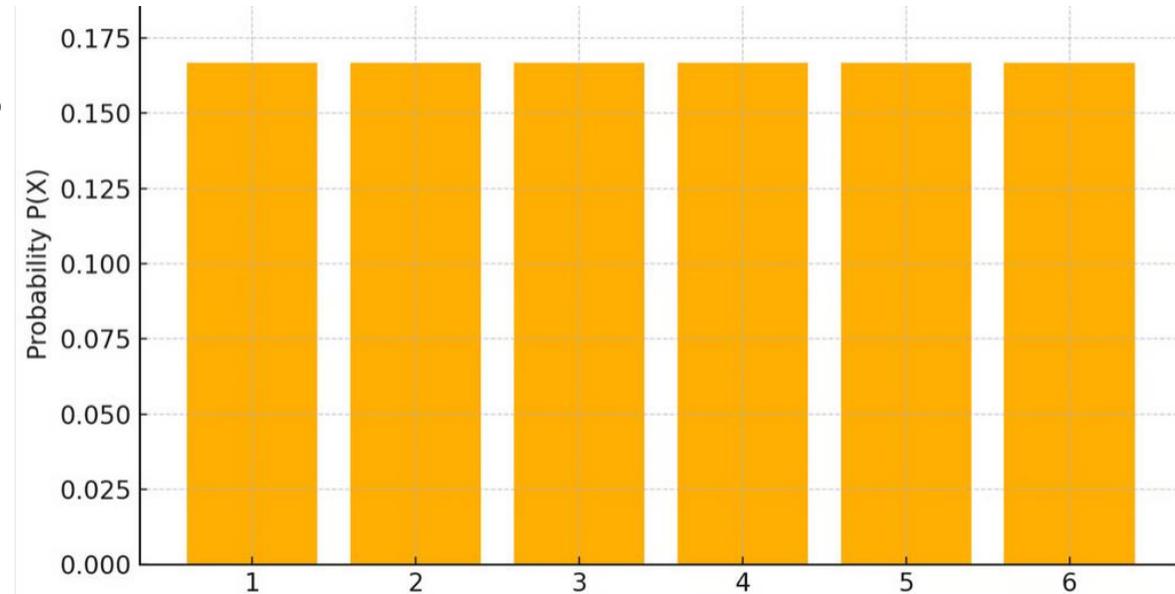
- A probability distribution is a theoretical and mathematical function or table that describes the likelihood of all possible outcomes for a *random variable*. ← Can be discrete or continuous.

- A probability, whether it's outside the head or inside the head, is just a way of assigning numbers to a set of mutually exclusive possibilities.



# Probability Mass Function

- The probability mass function (PMF) is a relationship that defines the probabilities associated with each possible outcome of a discrete random variable. ← the horizontal axis
- PMFs  $P(x)$  are visualized using bar charts (長條圖).
- The vertical axis represents the probability assigned by the PMF.

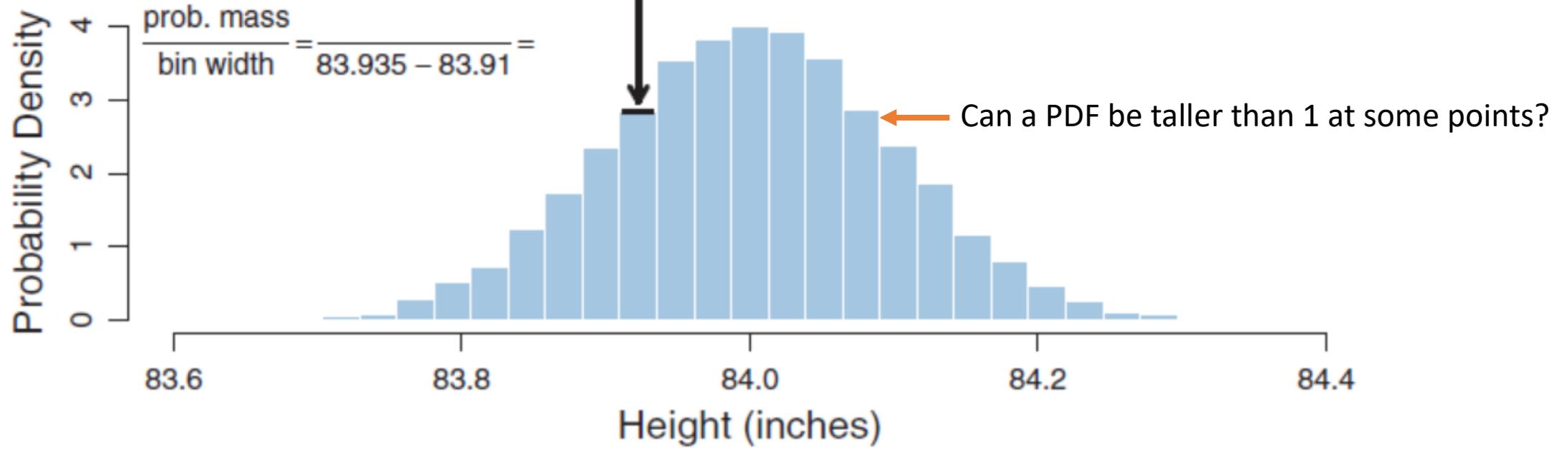
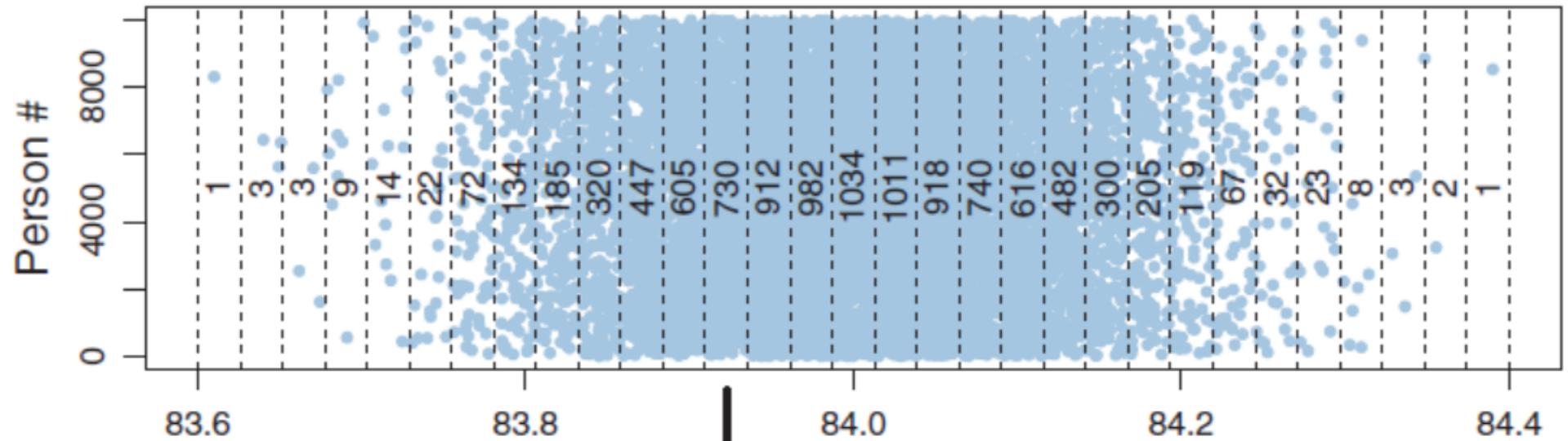


# Probability Density Function

- The probability density function (PDF) defines the probability of a continuous random variable lying between a specific range of values.
- PDFs  $f(x)$  are visualized using histograms (直方圖).
- We make the intervals infinitesimally narrow, and talk about the ratio of the probability mass to the interval width. That ratio is called the *probability density*.

← Is the density a probability? Why?

Total N = 10000



# Properties of PDF

- The probability mass of the  $i$ th interval is denoted  $p([x_i, x_i + \Delta x])$ . Then the sum of those probability masses must be 1:

$$\sum_i p([x_i, x_i + \Delta x]) = 1.$$

- The probability density is the ratio of probability mass over interval width, so we have

$$\sum_i \Delta x \frac{p([x_i, x_i + \Delta x])}{\Delta x} = 1.$$

# Integral of PDF

- Do not confuse the probability density  $p(x)$  with  $p([x_i, x_i + \Delta x])$ , which is the probability mass in an interval. ← Thus, the density is often denoted by  $f(x)$
- The above summation equation becomes an integral:

$$\sum_i \underbrace{\Delta x}_{dx} \underbrace{\frac{p([x_i, x_i + \Delta x])}{\Delta x}}_{p(x)} = 1 \quad \text{that is,} \quad \int dx p(x) = 1.$$

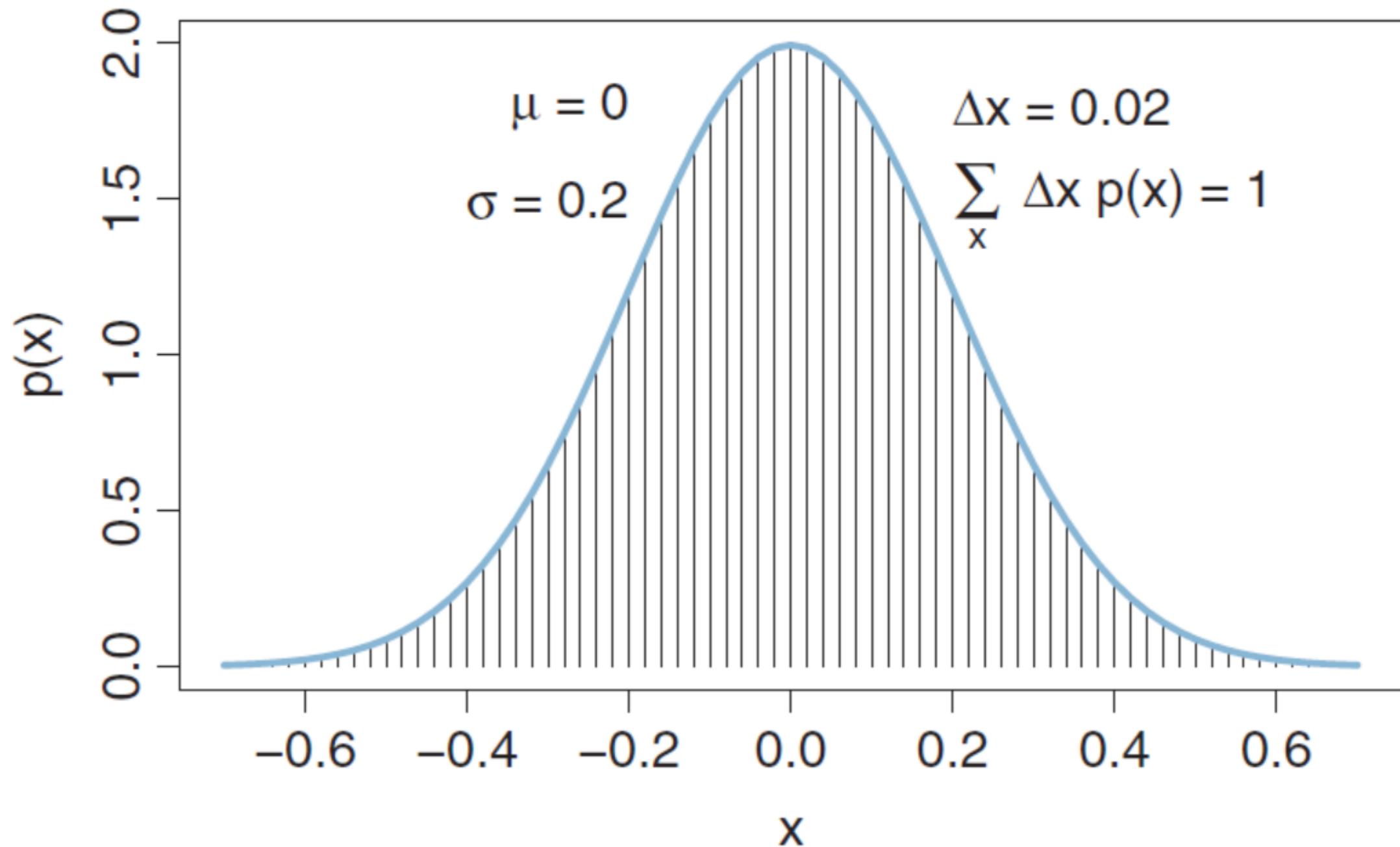
# The Normal Distribution

- Any function that has only nonnegative values and integrates to 1 can be a probability density function.
- The most famous probability density function is the normal distribution, also known as the Gaussian distribution:

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left[\frac{x - \mu}{\sigma}\right]^2\right).$$

- Two parameters: The mean ( $\mu$ ) and the standard deviation ( $\sigma$  governs how wide the bell is).

# Normal Probability Density



# Mean of a Distribution

- The *mean* of a probability distribution (also called the expected value) is an average value in the long run, weighted by the probability of each outcome:

$$E[x] = \sum_x p(x) x \quad \text{or} \quad E[x] = \int dx p(x) x$$

- Ex 1.  $p(x) = 6x(1 - x)$  for  $x \in [0,1]$ .
- Ex 2. Can you verify that  $E[x] = \mu$  for the standard normal distribution ( $\mu = 0, \sigma = 1$ )?

# Variance of a Distribution

- The *variance* of a probability distribution  $\sigma^2$  is a number that represents the dispersion of the distribution away from its mean.

- The definition of variance is the mean squared deviation (MSD) of the  $x$  values from their mean:

$$\text{Var}[x] = \int dx p(x) (x - E[x])^2$$

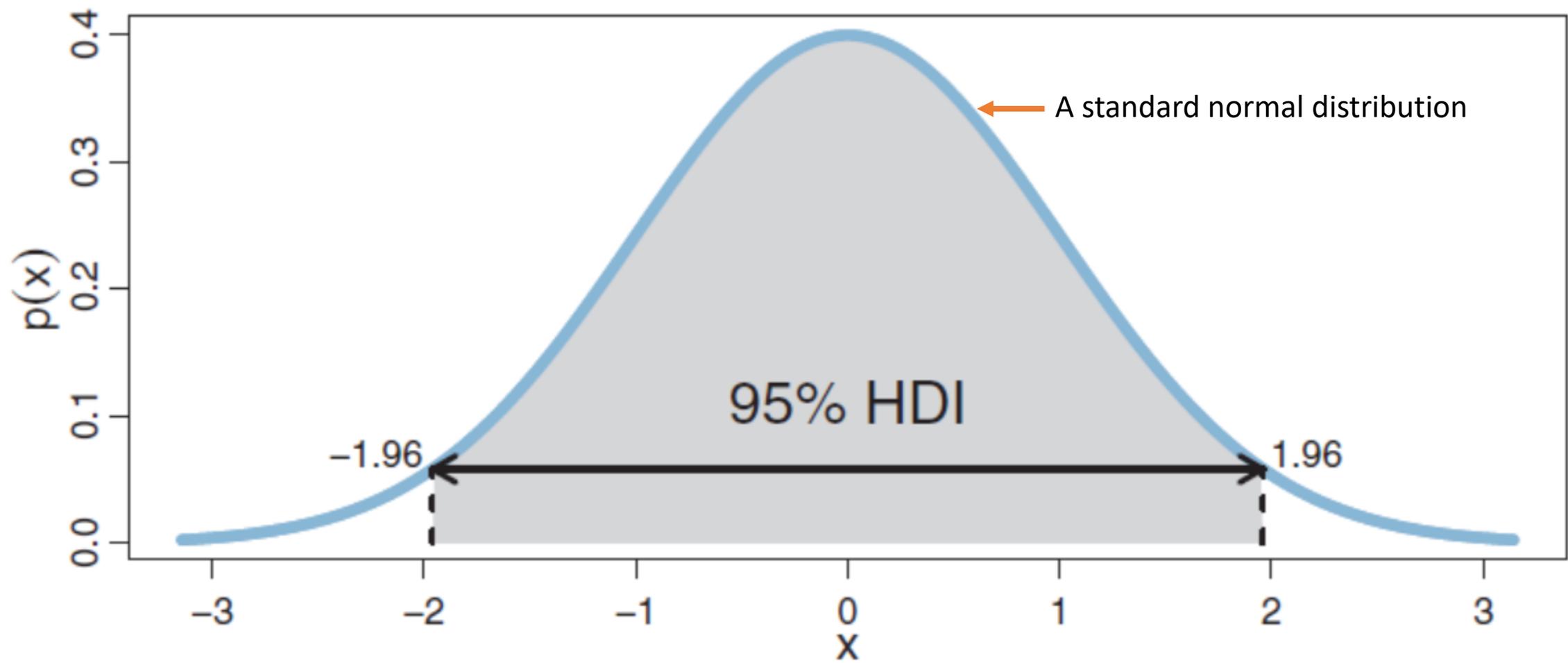
- Exercise 4.5. What is the probability mass under the standard normal curve from  $x = \mu - \sigma$  to  $x = \mu + \sigma$ ?

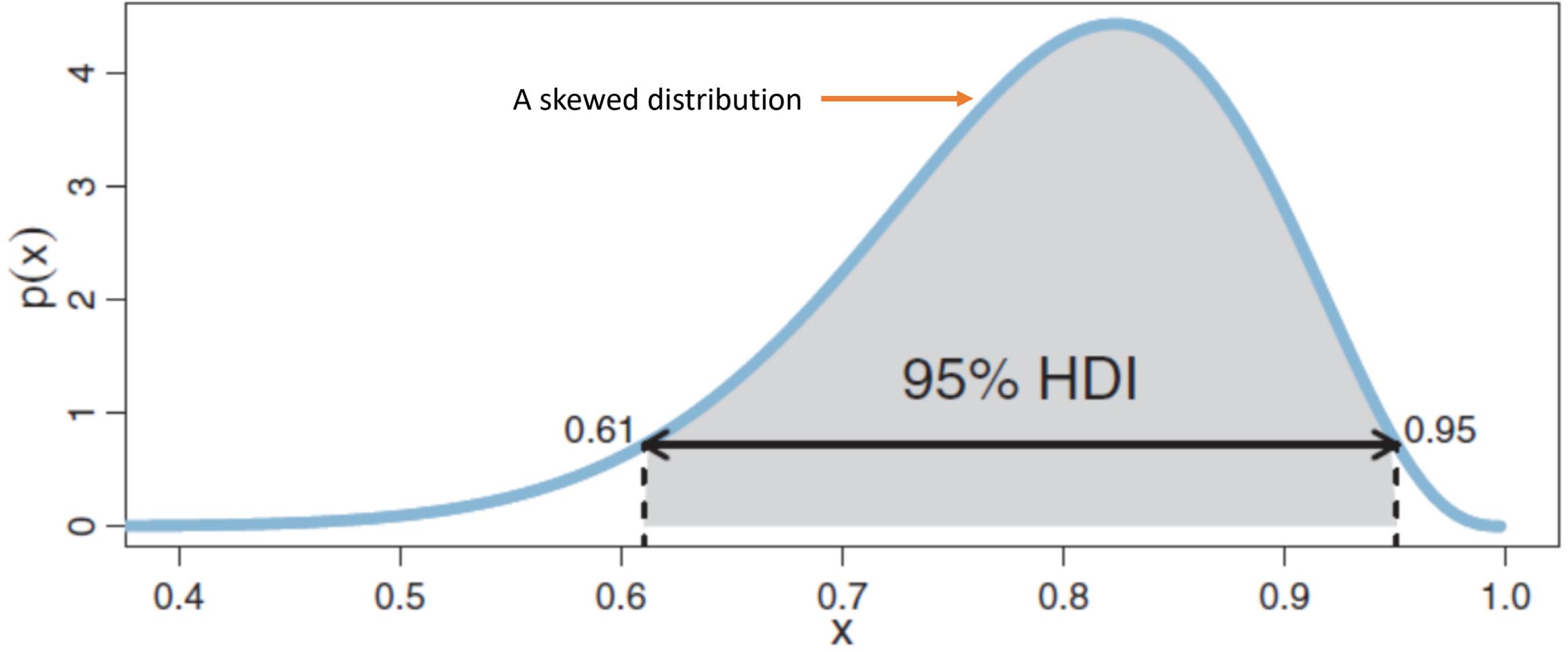
# Standard Deviation

- The square root of the variance, sometimes referred to as root mean squared deviation (RMSD), is called the *standard deviation* of the distribution..
- A probability distribution can refer to probability of measurement values or of parameter values.
- The standard deviation of  $\theta$ , which measures how wide the distribution is, can be thought of as a measure of uncertainty across candidate values.

# Highest Density Interval (HDI)

- Another way of summarizing a distribution is the *highest density interval* (HDI).
- The HDI indicates which points of a distribution are most *credible*. The area under the curve between the HDI limits is, say 0.95 (or  $\geq 0.95$  for p.m.f.s).
- Moreover, every point inside the HDI has higher probability density than any point outside the interval.





# Caveats

- The HDI does not necessarily produce equal-area tails outside the HDI (but confidence interval does).
- With a multimodal distribution, the HDI is split into two subintervals, not good for statistical estimation.
- If the HDI is wide, then beliefs are uncertain. If the HDI is narrow, then beliefs are relatively certain.

