

## Module 3: Introduction to the Normal Distribution

Rebecca C. Steorts

# Announcements

Thank you for your feedback on the course so far! If you have not yet filled out the survey, I would greatly appreciate your feedback.  
Response rate: 24/40.

# What you appreciated/liked

1. The pace of the lecture.
2. The interactive lectures and use of the ipad.
3. The pre-recorded videos.
4. Most students like the homework and length.
5. Most students like the organization and structure of the course.  
If you did not, there is a video explaining how to navigate the course webpage that I went over on the first day of class.
6. Students appreciated how available myself + the TA Team for answering questions. (If you're not reaching out; please reach out). We are here to help you!
7. Students appreciated the focus on learning over grading.
8. Students said the course was on par or much better than their virtual lectures.
9. Students appreciated that there is a wealth of information available to them.

# Student Requests

- ▶ Students asked for a cheat sheet regarding the course. Done!
- ▶ Extra OH. Done!
- ▶ Students would like to see expectations regarding the graded homeworks. Done!
- ▶ We're moving to Piazza. (See updates to course webpage, Sakai, syllabus).
- ▶ ipad notes have been pre-recorded, so these are already there:)
- ▶ I've posted all my (live) written ipad notes in each github module in case you want to look back over these for your exam preparation. You also have the typed notes. If you want to see a different derivation, you can look in my typed lecture notes or in Hoff.

# Student Requests

- ▶ “I’m overwhelmed by the reading, video recordings, and amount of material.”

If you’re getting overwhelmed, please talk with me 1-1, so I can talk with you regarding making a plan and strategy that works best with you.

- ▶ A few students would prefer videos on Panopto. I unfortunately will not have time to move all videos to Panopto and make this extremely organized this semester. I really appreciate the feedback though.

## Other Requests

- ▶ Two students would like different due dates of assignments.
- ▶ Two students did not like the course webpage. (More feedback would be helpful on what you don't like about it).

Since most students didn't ask for a change, I'm going to keep these dates the same. If you have a special circumstance that falls under the flexibility teaching policy that I offer, please email me by today to talk about options.

## Other Requests

- ▶ The labs now just have one file (Rmd and pdf). Please check that you can render these as you need LaTeX installed.
- ▶ When class connects to a lab, I've noted in the online lectures.
- ▶ Pop me an email if you have any clarification, and I'm always happy to modify the lecture notes to improve the clarity.
- ▶ If there is a major change moving forward, I will post an announcement on Piazza.

# Agenda

- ▶ The normal distribution
- ▶ The variance versus precision
- ▶ The re-parameterized normal distribution
- ▶ Common properties
- ▶ The normal-uniform model
- ▶ The normal-normal model



# Normal distribution

The normal distribution  $\mathcal{N}(\mu, \sigma^2)$

- ▶ with mean  $\mu \in \mathbb{R}$  and variance  $\sigma^2 > 0$  - (standard deviation  $\sigma = \sqrt{\sigma^2}$ ) has p.d.f.

$$\mathcal{N}(x \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

for  $x \in \mathbb{R}$ .

It is often more convenient to write the p.d.f. in terms of the **precision**, or inverse variance,  $\lambda = 1/\sigma^2$  rather than the variance.

# Re-parameterized Normal

In this parametrization, the p.d.f. is

$$\mathcal{N}(x \mid \mu, \lambda^{-1}) = \sqrt{\frac{\lambda}{2\pi}} \exp \left( -\frac{1}{2} \lambda (x - \mu)^2 \right)$$

since  $\sigma^2 = 1/\lambda = \lambda^{-1}$ .

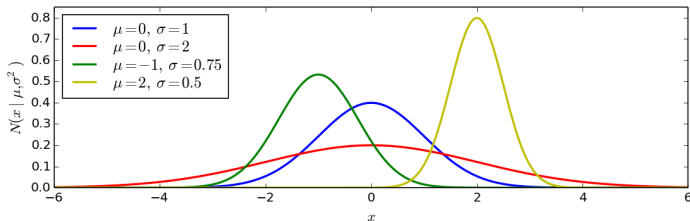


Figure 1: Normal distribution with various choices of  $\mu$  and  $\sigma$ .

# Normality?

- ▶ The central limit theorem (CLT) states that the sum of a large number of independent random variables tends to be approximately normally distributed.
- ▶ Real world data often appears approximately normal.

# Normality?

- ▶ Human heights and other body measurements,
- ▶ Cumulative hydrologic measures such as annual rainfall or monthly river discharge,
- ▶ Errors in astronomical or physical observations,
- ▶ Diffusion of a substance in a liquid or gas.
- ▶ Some things are products of many independent variables (rather than sums), and in such cases the logarithm will be approximately normal since it is a sum of many independent variables

Example: stock market indices, due to the effect of compound interest.

# Properties of the Normal distribution

- ▶ Mean, median, and mode are all the same ( $\mu$ )
- ▶ Symmetric about the mean
- ▶ 95% probability within  $\pm 1.96\sigma$  of the mean (roughly,  $\pm 2\sigma$ )
- ▶ If  $X \sim \mathcal{N}(\mu, \sigma^2)$  and  $Y \sim \mathcal{N}(m, s^2)$  independently, then

$$aX + bY \sim \mathcal{N}(a\mu + bm, a^2\sigma^2 + b^2s^2). \quad (1)$$

- ▶ Careful: `rnorm`, `dnorm`, `pnorm`, and `qnorm` in R take the mean and **standard deviation**  $\sigma$  as arguments (not mean and variance  $\sigma^2$ ). For example, `rnorm(n,m,s)` generates  $n$  normal random variables from  $\mathcal{N}(m, s^2)$ .

# Normal-Uniform

$$X_1, \dots, X_n \mid \theta \stackrel{\text{iid}}{\sim} \mathcal{N}(\theta, \sigma^2).$$

Assume the prior on  $\theta$  is constant over the real line. We can write this as  $p(\theta) \propto 1$ , where  $-\infty < \theta < \infty$ .

Derive the posterior distribution.<sup>1</sup>

---

<sup>1</sup>Please observe that the posterior is not conjugate to the prior in this situation, so it's very important to make sure that the posterior is a proper distribution.

## Solution

$$p(\theta \mid x_{1:n}) \propto \mathcal{N}(x_{1:n} \mid \theta, \sigma^2) \times 1 \quad (2)$$

$$\propto \left(\frac{\ell}{2\pi}\right)^{n/2} \exp\left(-\frac{1}{2}\ell \sum_i (x_i - \theta)^2\right)$$

$$\propto \exp\left(-\frac{1}{2}\ell \sum_i (x_i - \bar{x} + \bar{x} - \theta)^2\right)$$

$$\propto \exp\left(-\frac{1}{2}\ell \sum_i (x_i - \bar{x})^2\right) \exp\left(-\frac{1}{2}\ell \sum_i (\bar{x} - \theta)^2\right)$$

$$\propto \exp\left(-\frac{1}{2}\ell \sum_i (\bar{x} - \theta)^2\right) \quad (3)$$

$$= \exp\left(-\frac{n\ell}{2}(\theta - \bar{x})^2\right) \quad (4)$$

This implies that

$$\theta \mid x_{1:n} \sim N(\bar{x}, (n\ell)^{-1})$$

## Normal-Normal

$$X_1, \dots, X_n \mid \theta \stackrel{\text{iid}}{\sim} \mathcal{N}(\theta, \lambda^{-1}).$$

Assume the precision  $\lambda = 1/\sigma^2$  is known and fixed, and  $\theta$  is given a  $\mathcal{N}(\mu_0, \lambda_0^{-1})$  prior:

$$\boldsymbol{\theta} \sim \mathcal{N}(\mu_0, \lambda_0^{-1})$$

i.e.,  $p(\theta) = \mathcal{N}(\theta \mid \mu_0, \lambda_0^{-1})$ . This is sometimes referred to as a **Normal–Normal** model.



## Posterior derivation

We begin with the **likelihood** of the normal distribution.

For any  $x$  and  $\ell$ ,

$$\begin{aligned}\mathcal{N}(x \mid \theta, \ell^{-1}) &= \sqrt{\frac{\ell}{2\pi}} \exp\left(-\frac{1}{2}\ell(x - \theta)^2\right) \\ &\propto_{\theta} \exp\left(-\frac{1}{2}\ell(x^2 - 2x\theta + \theta^2)\right) \\ &\propto_{\theta} \exp\left(\ell x \theta - \frac{1}{2}\ell\theta^2\right).\end{aligned}\tag{5}$$

Note: we drop the **constant term** and we will do this often when working with the normal distribution.

## Posterior derivation (continued)

We now consider the **prior** distribution on  $\theta$ .

Due to the symmetry of the normal p.d.f.,

$$\mathcal{N}(\theta \mid \mu_0, \lambda_0^{-1}) = \mathcal{N}(\mu_0 \mid \theta, \lambda_0^{-1}) \propto_{\theta} \exp(\lambda_0 \mu_0 \theta - \frac{1}{2} \lambda_0 \theta^2), \quad (6)$$

where  $x = \mu_0$  and  $\ell = \lambda_0$ .

## Posterior derivation (continued)

Let

$$L = \lambda_0 + n\lambda \quad \text{and} \quad M = \frac{\lambda_0\mu_0 + \lambda \sum_{i=1}^n x_i}{\lambda_0 + n\lambda}.$$

$$\begin{aligned} p(\theta|x_{1:n}) &\propto p(\theta)p(x_{1:n}|\theta) \\ &= \mathcal{N}(\theta \mid \mu_0, \lambda_0^{-1}) \prod_{i=1}^n \mathcal{N}(x_i \mid \theta, \lambda^{-1}) \\ &\stackrel{(a)}{\propto} \exp(\lambda_0\mu_0\theta - \tfrac{1}{2}\lambda_0\theta^2) \exp(\lambda(\sum x_i)\theta - \tfrac{1}{2}n\lambda\theta^2) \\ &= \exp\left((\lambda_0\mu_0 + \lambda \sum x_i)\theta - \tfrac{1}{2}(\lambda_0 + n\lambda)\theta^2\right) \\ &= \exp(LM\theta - \tfrac{1}{2}L\theta^2) \\ &\stackrel{(b)}{\propto} \mathcal{N}(M \mid \theta, L^{-1}) = \mathcal{N}(\theta \mid M, L^{-1}), \end{aligned}$$

where step (a) uses Equations 5 and 6, and step (b) uses Equation 5 with  $x = M$  and  $\ell = L$ .

## Posterior derivation (continued)

Recall

$$L = \lambda_0 + n\lambda \quad \text{and} \quad M = \frac{\lambda_0 \mu_0 + \lambda \sum_{i=1}^n x_i}{\lambda_0 + n\lambda}.$$

It turns out that the posterior is

$$\boldsymbol{\theta} | \mathbf{x}_{1:n} \sim \mathcal{N}(M, L^{-1}) \quad (7)$$

i.e.,  $p(\theta | \mathbf{x}_{1:n}) = \mathcal{N}(\theta | M, L^{-1})$ .

Thus, the normal distribution is, itself, a conjugate prior for the mean of a normal distribution with known precision.

# Heights of Adult Humans

- ▶ Heights tend to be normally distributed because there are many independent genetic and environmental factors which contribute additively to overall height
- ▶ This leads to a normal distribution due to the central limit theorem.

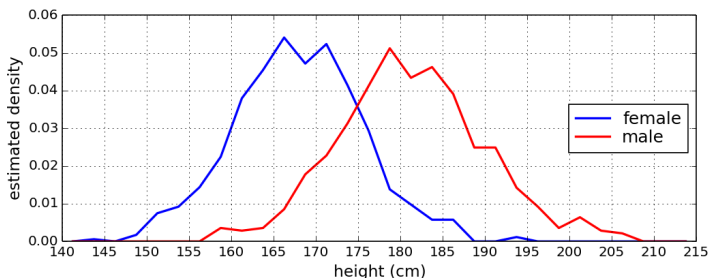


Figure 2: Estimated densities of the heights of Dutch women and Dutch men based on a sample of 695 women and 562 men.

# Heights of Adult Humans

- ▶ Consider combined distribution of heights (pooling females and males together). Would this be normal?
- ▶ It is thought that such data is bimodal (having two maxima). Is it really bimodal? (See, Schilling et al. (2002))

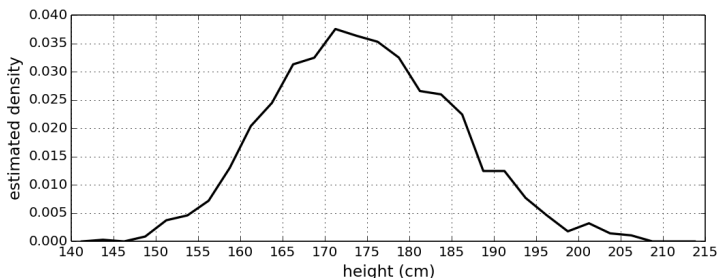


Figure 3: Estimated density for Dutch women and men together, assuming there is an equal proportion of women and men in the population.

# Heights of Adult Humans, Combined

At a glance, while the heights of women and men separately do appear to be roughly normally distributed, the combined distribution does not look bimodal. How could we test whether it is bimodal in a more precise way?

# Our Assumptions

- ▶ Assume female heights and male heights are each normally distributed.
- ▶ Assume both female heights and male heights have different means but the same standard deviation.
- ▶ Assume that there is an equal proportion of women and men in the population.
- ▶ Then, it is known that the combined distribution is bimodal if and only if the difference between the means is greater than twice the standard deviation (Helguerro, 1904).



## Model

In mathematical notation: Assume the female heights are

$$X_1, \dots, X_k \stackrel{\text{iid}}{\sim} \mathcal{N}(\theta_f, \sigma^2),$$

where  $k = 695$ , the male heights are

$$Y_1, \dots, Y_\ell \stackrel{\text{iid}}{\sim} \mathcal{N}(\theta_m, \sigma^2),$$

where  $\ell = 562$ , and the p.d.f. of the combined distribution of heights is

$$\frac{1}{2}\mathcal{N}(x \mid \theta_f, \sigma^2) + \frac{1}{2}\mathcal{N}(x \mid \theta_m, \sigma^2).$$

(This is an example of what is called a two-component **mixture** distribution.)

# Model

Let's put independent normal priors on  $\theta_f$  and  $\theta_m$ :

$$p(\theta_f, \theta_m) = p(\theta_f)p(\theta_m) = \mathcal{N}(\theta_f \mid \mu_{0,f}, \sigma_0^2)\mathcal{N}(\theta_m \mid \mu_{0,m}, \sigma_0^2).$$

- ▶ Assume  $\sigma^2$  is known.
- ▶ For the purposes of this example, let's use  $\sigma = 8$  centimeters (about 3 inches).
- ▶ Based on common knowledge of typical human heights, let's choose the prior parameters (a.k.a. hyperparameters) as follows:

$\mu_{0,f}$	(mean of prior on female mean ht)	165 cm ( $\approx 5$ ft, 5 in)
$\mu_{0,m}$	(mean of prior on male mean ht)	178 cm ( $\approx 5$ ft, 10 in)
$\sigma_0$	(std. dev. of priors on mean ht)	15 cm ( $\approx 6$ in)

# Bimodal Fact

It is known (Helguerro, 1904) that the combined distribution is bimodal if and only if

$$|\theta_f - \theta_m| > 2\sigma.$$

So, to address our question of interest (“Is human height bimodal?”), we would like to compute the posterior probability that this is the case, i.e., we want to know

$$\mathbb{P}(\text{bimodal} \mid \text{data}) = \mathbb{P}(|\boldsymbol{\theta}_f - \boldsymbol{\theta}_m| > 2\sigma \mid x_{1:k}, y_{1:\ell}).$$

# Results

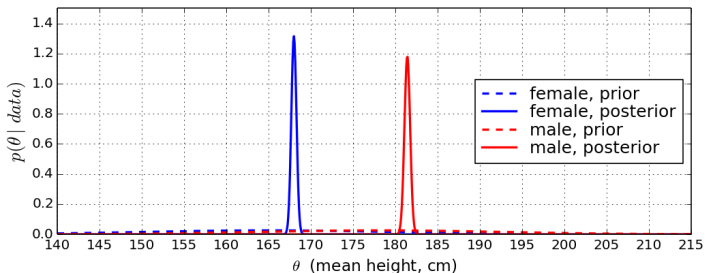


Figure 4: Priors and posteriors for the mean heights of Dutch women and men.

## Results (continued)

We can compute the posteriors for  $\theta_f$  and  $\theta_m$  using Equation 7 for each of them, independently. Figure 4 shows the priors and posteriors.

- ▶ Sample means:  $\bar{x} = 168.0$  cm (5 feet 6.1 inches) for females, and  $\bar{y} = 181.4$  cm (5 feet 11.4 inches) for males.
- ▶ Posterior means:  $M_f = 168.0$  cm for females, and  $M_m = 181.4$  cm for males. (Essentially identical to the sample means, due to the relatively large sample size and relatively weak prior.)
- ▶ Posterior standard deviations:  $1/\sqrt{L_f} = 0.30$  cm and  $1/\sqrt{L_m} = 0.34$  cm.

## Results (continued)

By Equation 1 (a linear combination of independent normals is normal),

$$\boldsymbol{\theta}_m - \boldsymbol{\theta}_f \mid x_{1:k}, y_{1:\ell} \sim \mathcal{N}(M_m - M_f, L_m^{-1} + L_f^{-1}) = \mathcal{N}(13.4, 0.45^2)$$

so we can compute  $\mathbb{P}(\text{bimodal} \mid \text{data})$  using the normal c.d.f.  $\Phi$ :

$$\begin{aligned}\mathbb{P}(\text{bimodal} \mid \text{data}) &= \mathbb{P}(|\boldsymbol{\theta}_m - \boldsymbol{\theta}_f| > 2\sigma \mid x_{1:k}, y_{1:\ell}) \\ &= \Phi(-2\sigma \mid 13.4, 0.45^2) + (1 - \Phi(2\sigma \mid 13.4, 0.45^2)) \\ &= 6.1 \times 10^{-9}.\end{aligned}$$

Intuitive interpretation: The posteriors are about 13 or 14 centimeters apart, which is under the  $2\sigma = 16$  threshold for bimodality, and they are sufficiently concentrated that the posterior probability of bimodality is essentially zero.

# Takeaways

- ▶ We reviewed the univariate normal distribution and properties of it (such as symmetry about the mean).
- ▶ We discussed types of data that empirically have a normal distribution to motivate its usage.
- ▶ We derived the Normal-Uniform model, where the uniform distribution is the first improper prior we have seen in this course.
- ▶ We derived the Normal-Normal model.
- ▶ We considered an application of human heights from Schilling et al. (2002), where we showed that the posterior probability of bimodality was zero (regarding a plot of combined heights).

## Detailed Takeaways for Exam

- ▶ Working with the univariate normal distribution
- ▶ Knowing the precision and variance relationship
- ▶ Knowing the shape of the normal distribution
- ▶ Knowing the CLT
- ▶ Knowing what type of data is approximately normal and what data is not
- ▶ Knowing properties of the normal that we will often work with such as symmetry about the mean
- ▶ Being able to derive the posterior of the normal-uniform model
- ▶ Understanding that a non-informative prior is highly informative and why.
- ▶ Being able to derive the normal-normal posterior distribution
- ▶ Knowing how to write out the posterior mean and variance intuitively



## Detailed Takeaways for Exam (Continued)

- ▶ You should be able to work with a two-component mixture model (as in the case study)
- ▶ Suppose I provide you with a case study such as the one on human heights. You should be able to explain why it makes sense to model the data as normal and what parameters would be reasonable given the description.
- ▶ Given as case study, you should be able to incorporate a pilot (or prior study) into your prior distribution and back up how you incorporate it.
- ▶ You should be able to update the posterior accordingly.
- ▶ You be able to state benefits of your model specification and weaknesses.
- ▶ You should be able to state a better model that would be more realistic once we have covered the normal-normal-gamma model (module 4).
- ▶ You should be able to discuss sensitivity of the posterior analysis for any hyper-parameters.

## Exercise

Recall that Recall

$$L = \lambda_0 + n\lambda \quad \text{and} \quad M = \frac{\lambda_0 \mu_0 + \lambda \sum_{i=1}^n x_i}{\lambda_0 + n\lambda}.$$

What happens to the posterior mean and the precision as  $n \rightarrow \infty$ .

## Exercise

Let's consider the posterior mean.

$$\begin{aligned} M &= \frac{\lambda_0 \mu_0 + \lambda \sum_{i=1}^n x_i}{\lambda_0 + n\lambda} \\ &= \frac{\lambda_0 \mu_0 + \lambda n \bar{x}}{\lambda_0 + n\lambda} \\ &= \frac{\lambda_0 \mu_0 / n + \lambda \bar{x}}{\lambda_0 / n + \lambda} \end{aligned}$$

As  $n \rightarrow \infty$ ,

$$M \rightarrow \lambda \bar{x} / \lambda = \bar{x}.$$

## Exercise

Let's consider the posterior variance.

$$L^{-1} = \frac{1}{\lambda_0 + n\lambda}$$

As  $n \rightarrow \infty$ ,

$$L^{-1} = \frac{1}{\lambda_0 + n\lambda} = \frac{1/n}{\lambda_0/n + \lambda} \approx \frac{1}{n\lambda} \rightarrow 0$$

# Interpretation

What can we learn from this example?

If our sample size is large enough (for any application or real world example), then

1. the posterior mean will tend to the the sample mean.
2. the posterior variance will tend to 0.