

Introduction to Machine Learning

Week 5 - Normal/Gaussian Distribution - unknown mean and variance

Spring 2025

Instructor: Dr. Patrick Skeba

We walked through all the details of the normal-normal model...

- Our goal was to learn the unknown mean, μ , given N observations, \mathbf{x} , and **ASSUMING** the likelihood noise, σ , was known.
- We compared the MLE to the Bayesian approach.
- Discussed asymptotic behavior, when the prior standard deviation, τ_0 , approaches infinity and when the sample size, N , approaches infinity.

Normal-normal model with known σ

$$p(\mu|\mathbf{x}, \sigma) \propto \prod_{n=1}^N \{\text{normal}(x_n|\mu, \sigma)\} \times \text{normal}(\mu|\mu_0, \tau_0)$$

But now...let's relax the assumption on a known σ ...

- Perhaps instead of weighing a 25 dumbbell, I want to weigh myself.
- I no longer feel as confident in being able to assume a fixed value for σ .
- We now have TWO unknown parameters to learn: μ and σ

The posterior is now a **JOINT** distribution between μ and σ conditioned on \mathbf{x}

- We denote the joint posterior distribution as:

$$p(\mu, \sigma \mid \mathbf{x})$$

Does the likelihood change?

- The measurement process is essentially the same. I'm still weighing something.
- So a Gaussian still seems appropriate as the likelihood for this case:

$$p(\mathbf{x} \mid \mu, \sigma) = \prod_{n=1}^N (\text{normal}(x_n \mid \mu, \sigma))$$

Does the likelihood change?

- The measurement process is essentially the same. I'm still weighing something.
- So a Gaussian still seems appropriate as the likelihood for this case:

$$p(\mathbf{x} \mid \mu, \sigma) = \prod_{n=1}^N (\text{normal}(x_n \mid \mu, \sigma))$$

Observations are conditionally independent GIVEN the unknown mean and unknown noise

What about the prior?

- We now have to encode our belief not just about μ or not just about σ ...
- We have to encode our prior belief about their **joint distribution!**

$$p(\mu, \sigma)$$

The joint posterior is proportional to the likelihood times the prior

- Bayesian formulation for the unknown mean and unknown likelihood noise:

$$p(\mu, \sigma | \mathbf{x}) \propto \prod_{n=1}^N \{\text{normal}(x_n | \mu, \sigma)\} \times p(\mu, \sigma)$$

How do we handle the joint prior?

- Most statistics textbooks focus on factoring the prior as:

$$p(\mu|\sigma)p(\sigma)$$

- Allows making use of a conjugate prior on σ .
- Makes use of the math presented last time, and leads to analytic solutions for the marginal posterior on σ .

However, let's consider a different formulation

- We will assume a-priori the parameters are independent:

$$p(\mu, \sigma) = p(\mu)p(\sigma)$$

- Represents that knowing something about σ does not tell me anything about μ .

The un-normalized posterior distribution can then be written:

$$p(\mu, \sigma | \mathbf{x}) \propto \prod_{n=1}^N \{\text{normal}(x_n | \mu, \sigma)\} \times p(\mu)p(\sigma)$$

- Even though a-priori the parameters are assumed independent, the posterior may have a relationship between the parameters via the observations!

Independent prior specification

- We will continue to use a **normal prior** on the unknown mean:

$$p(\mu \mid \mu_0, \tau_0) = \text{normal}(\mu \mid \mu_0, \tau_0)$$

- For σ let's use a **uniform prior** defined between a lower bound, l , and an upper bound, u :

$$p(\sigma \mid l, u) = \text{uniform}(\sigma \mid l, u)$$

The un-normalized joint posterior:

$$p(\mu, \sigma | \mathbf{x}) \propto \prod_{n=1}^N \{\text{normal}(x_n | \mu, \sigma)\} \times \text{normal}(\mu | \mu_0, \tau_0) \times \text{uniform}(\sigma | l, u)$$

The joint posterior distribution does **NOT** have an analytic or closed form expression!

Let's start putting some numbers to this problem

- For the prior on my unknown mean weight...
 - I will use $\mu_0 = 250$ and $\tau_0 = 2$.
 - Thus, a-priori there is a $\approx 99\%$ probability that I weigh less than 255 pounds!
- For the prior on the unknown likelihood noise:
 - Set the lower bound to $l = 0.5$ and the upper bound to $u = 5.5$.
 - The units on the bounds are pounds.

Why those bounds on σ ?

- In this example, σ is the sampling noise or measurement error...remember the likelihood!

$$x_n \mid \mu, \sigma \sim \text{normal}(x_n \mid \mu, \sigma)$$

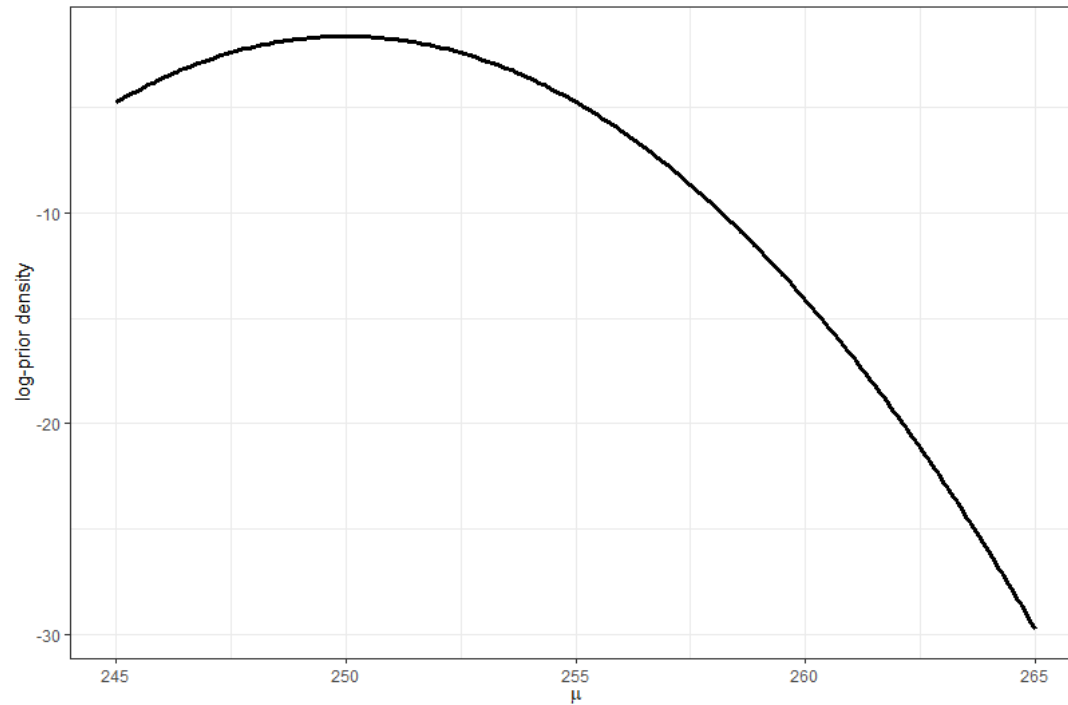
- σ is the lack of repeatability of a measurement.

Why those bounds on σ ?

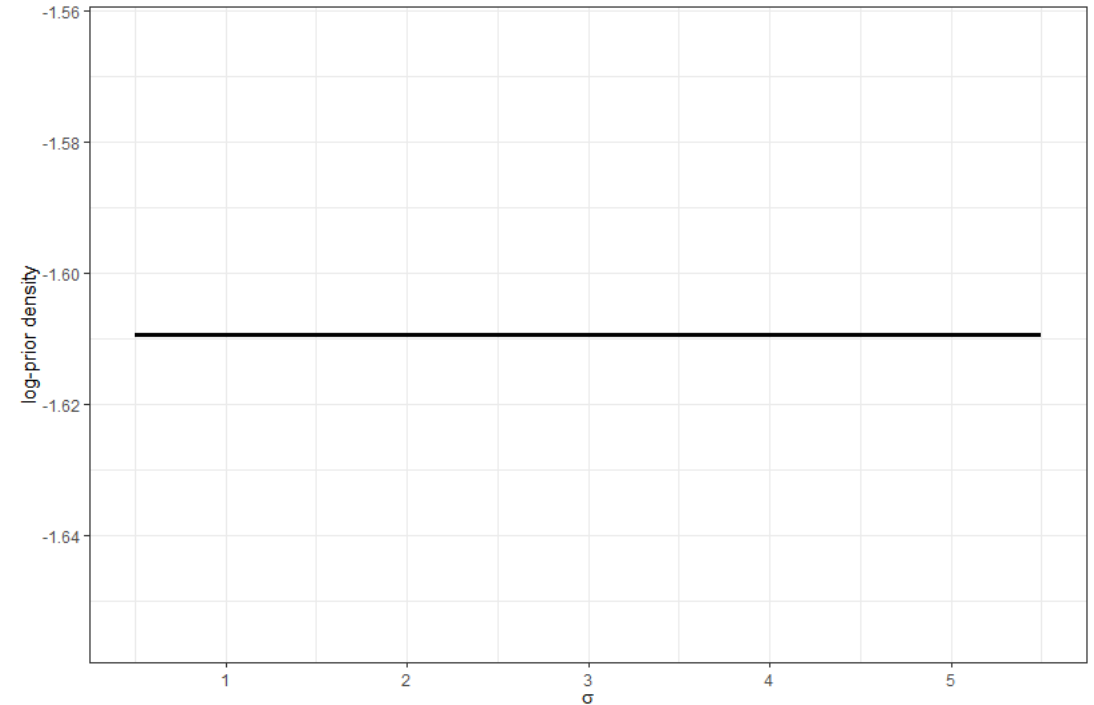
- We are not considering day-to-day variation in my weight.
- The noise represents variability in repeated measurements. For example, I step on the scale 10 times in a row.
- If I weigh 250 pounds, $\sigma = 5$ pounds means that there's $\approx 95\%$ chance that the scale would read between 240 and 260 across those 10 measurements!

What do the marginal log-prior densities look like?

a-priori μ is centered at $\mu_0 = 250$ with $\tau_0 = 2$

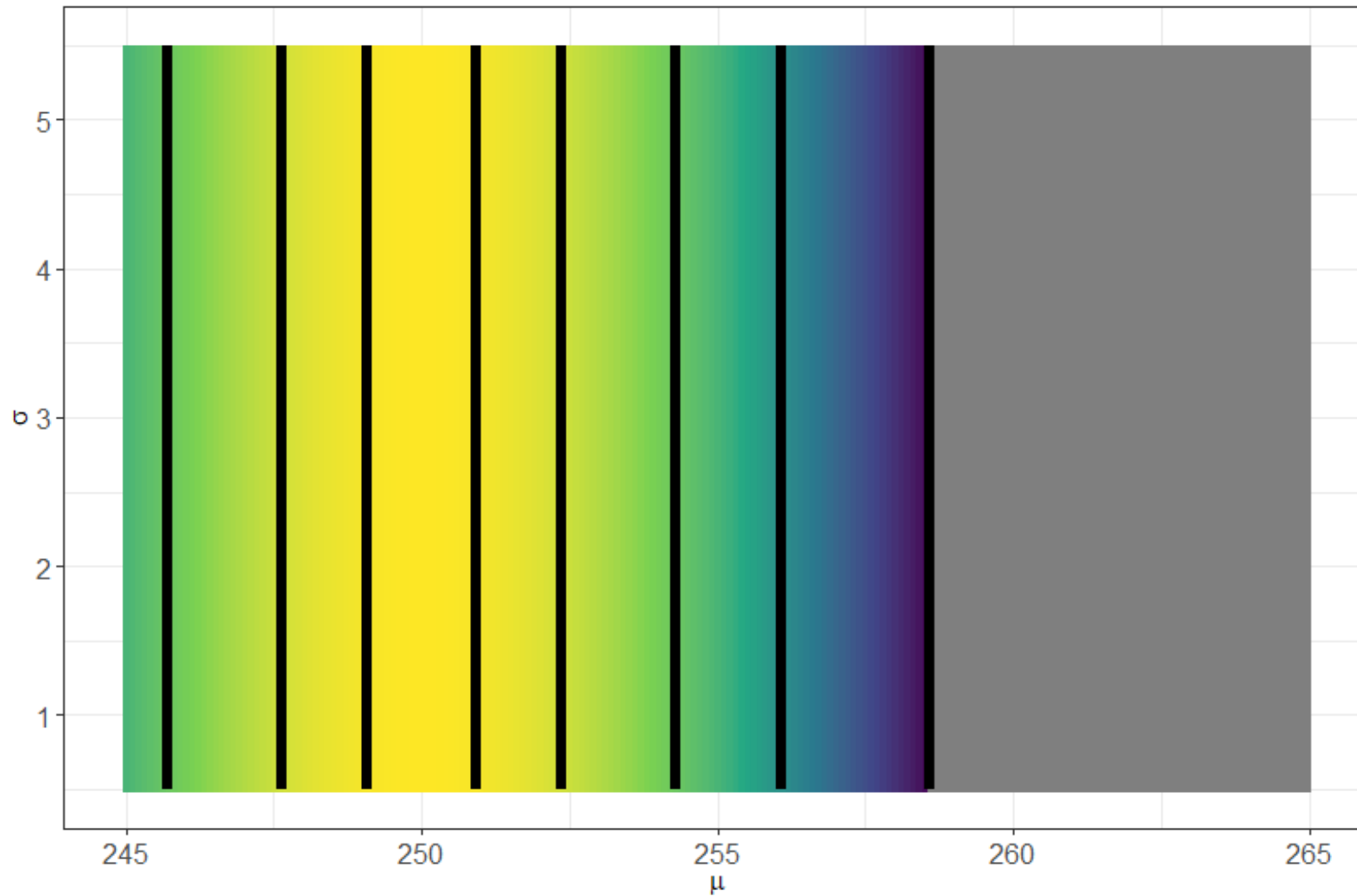


a-priori σ favors no interval in particular
between $l = 0.5$ and $u = 5.5$



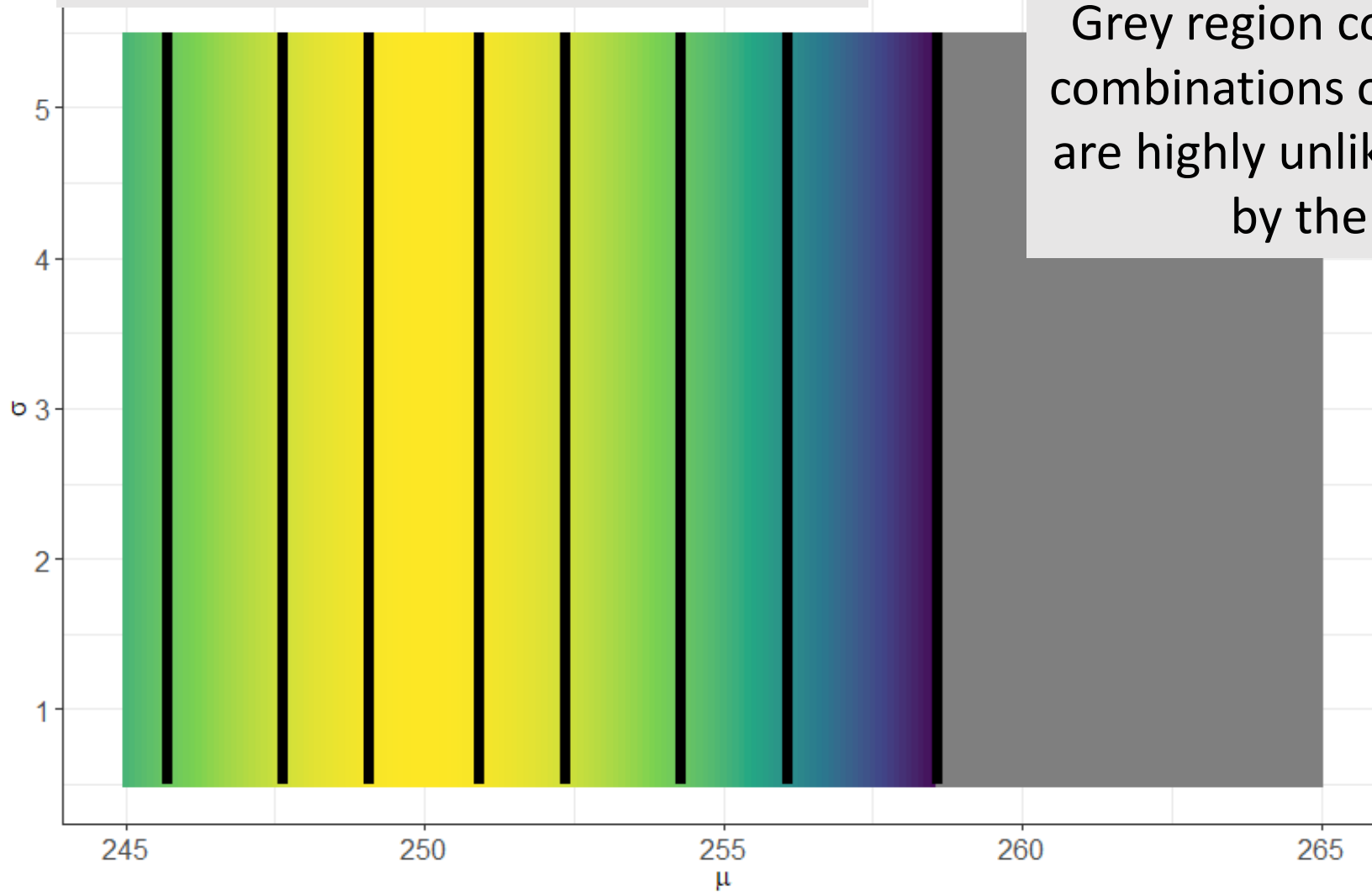
But, σ is constrained to exist ONLY between the
assumed lower and upper bounds!

Putting the two together reveals the log-prior surface



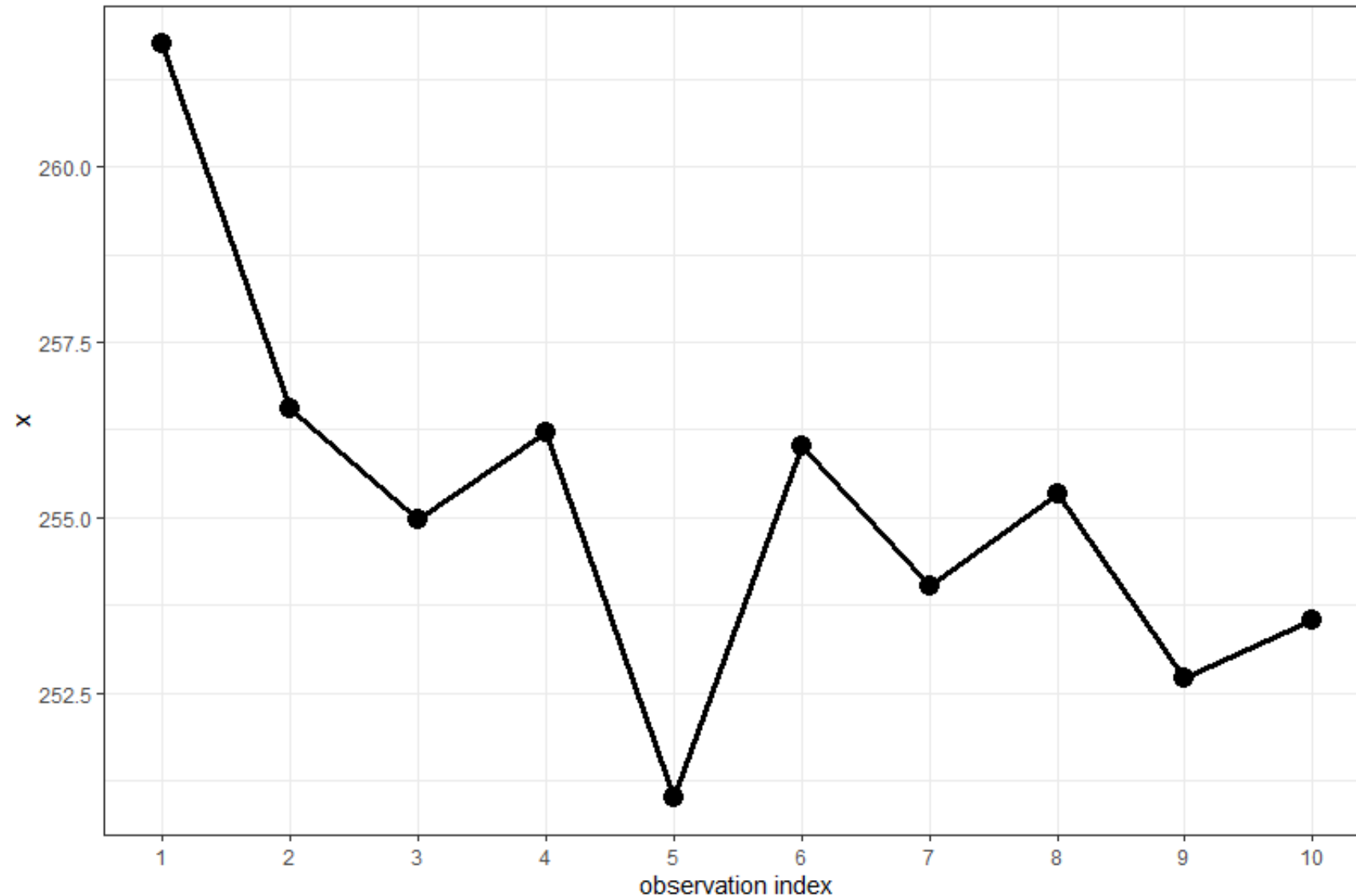
Putting the log-prior surface

Vertical contours are displayed because of the uniform prior on σ !



Grey region corresponds to combinations of μ and σ that are highly unlikely, as viewed by the prior!

Let's take a look at 10 measurements



Without an analytic expression...how can we proceed?

- Since we have just two unknowns, we can **visualize** the posterior, or rather the log-posterior, surface!
- We know how to write down the un-normalized log-posterior:

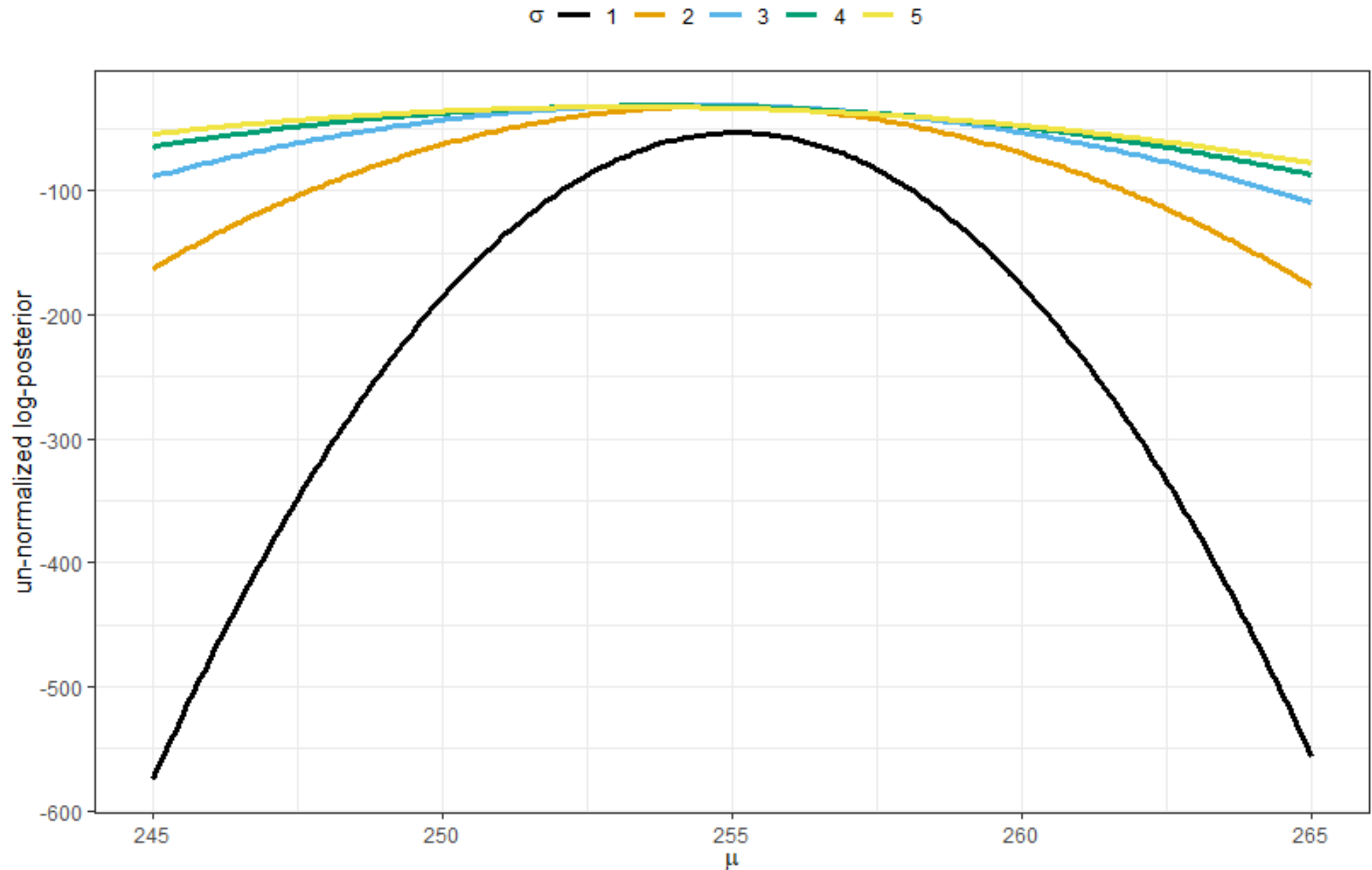
$$\log[p(\mu, \sigma \mid \mathbf{x})] \propto \sum_{n=1}^N (\log[\text{normal}(x_n \mid \mu, \sigma)]) + \log[\text{normal}(\mu \mid \mu_0, \tau_0)] + \log[\text{uniform}(\sigma \mid l, u)]$$

- If we specify candidate values for (μ, σ) we can calculate the un-normalized log-posterior!

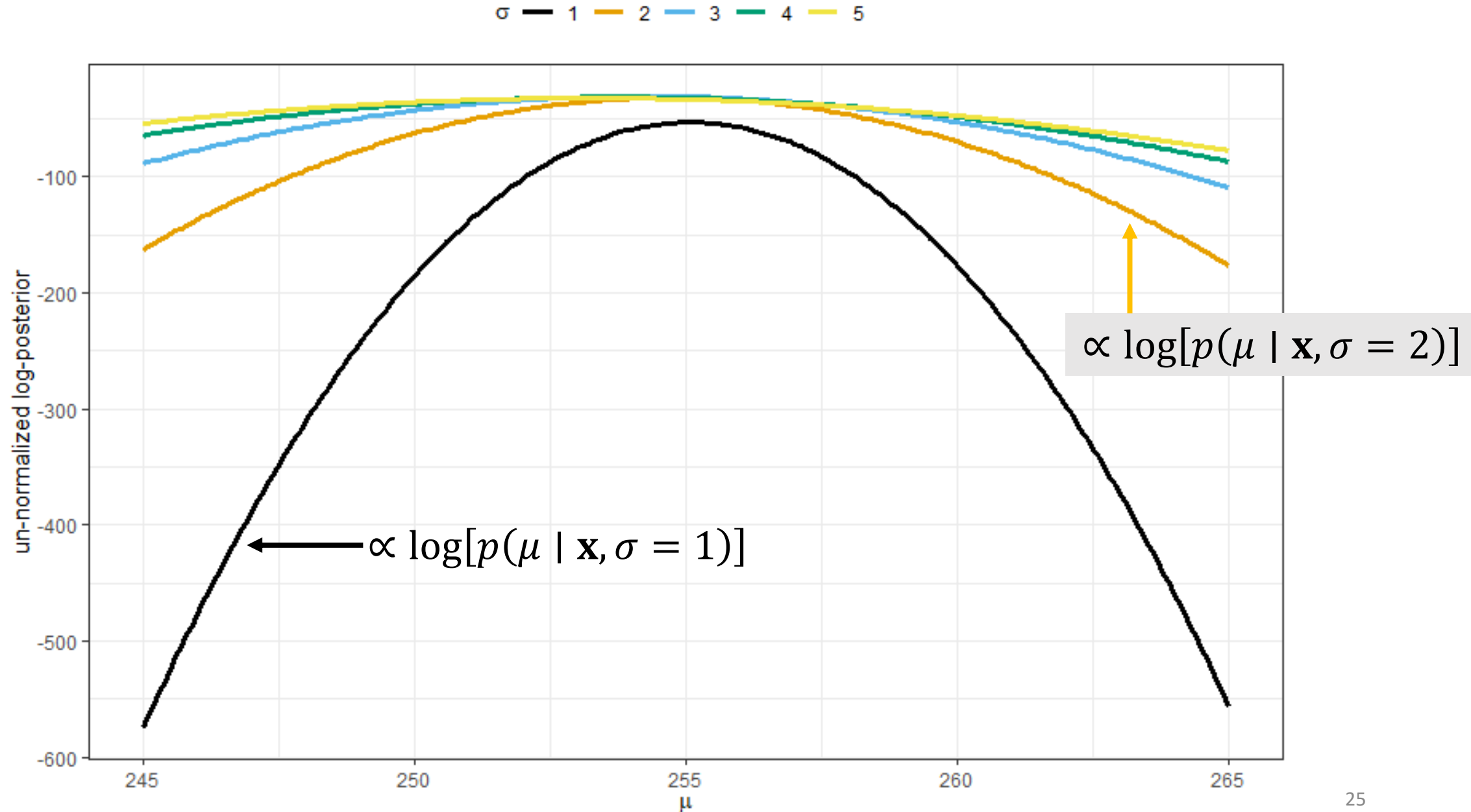
Before visualizing the surface, let's focus on one parameter at a time

- First, plot the log-posterior with respect to μ at a few values of σ .
- Then we will plot the log-posterior with respect to σ at a few values of μ .
- Lastly, we will examine the log-posterior surface with respect to both μ and σ .

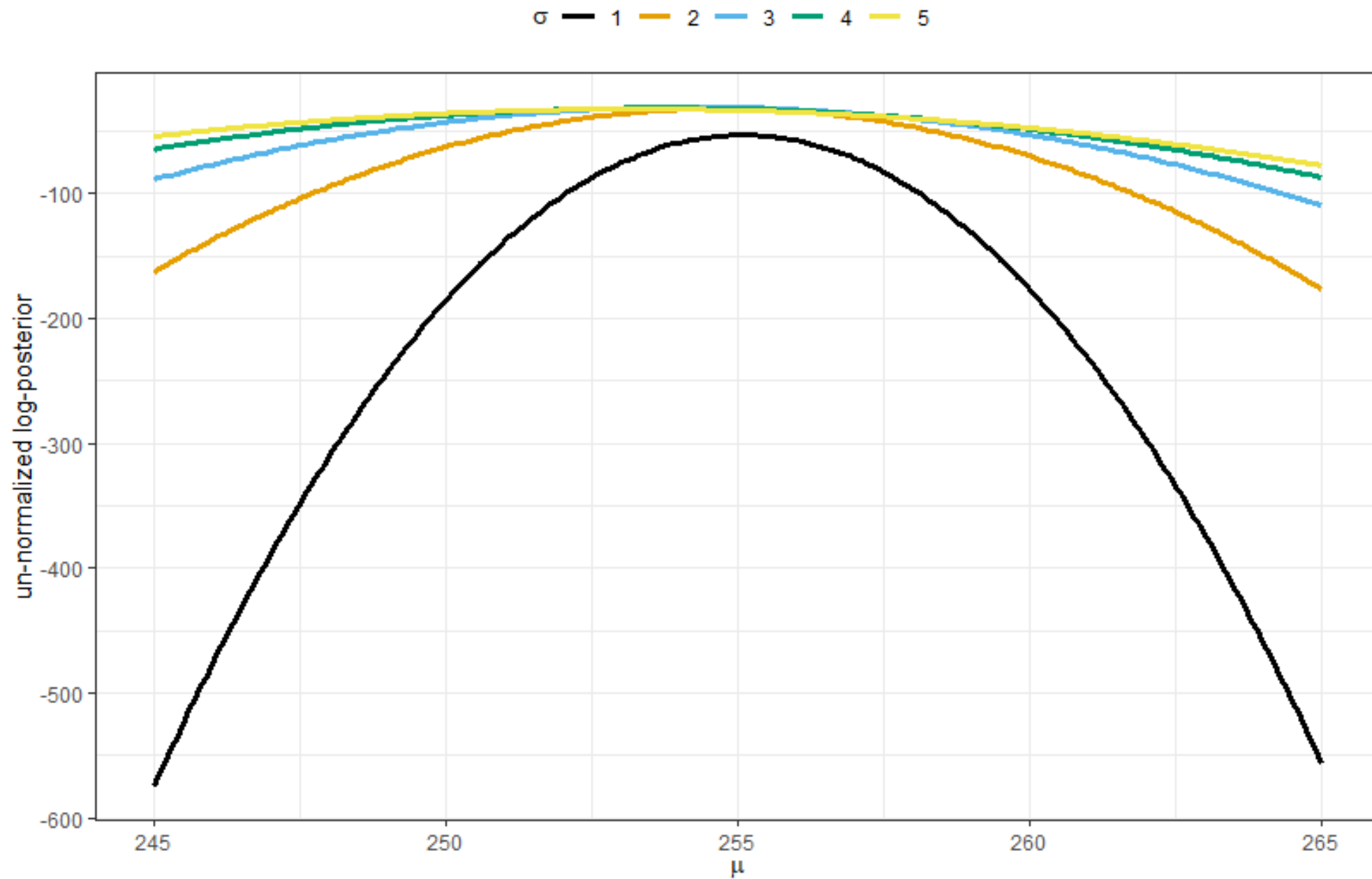
The log-posterior with respect to μ given specific values of σ should look familiar...



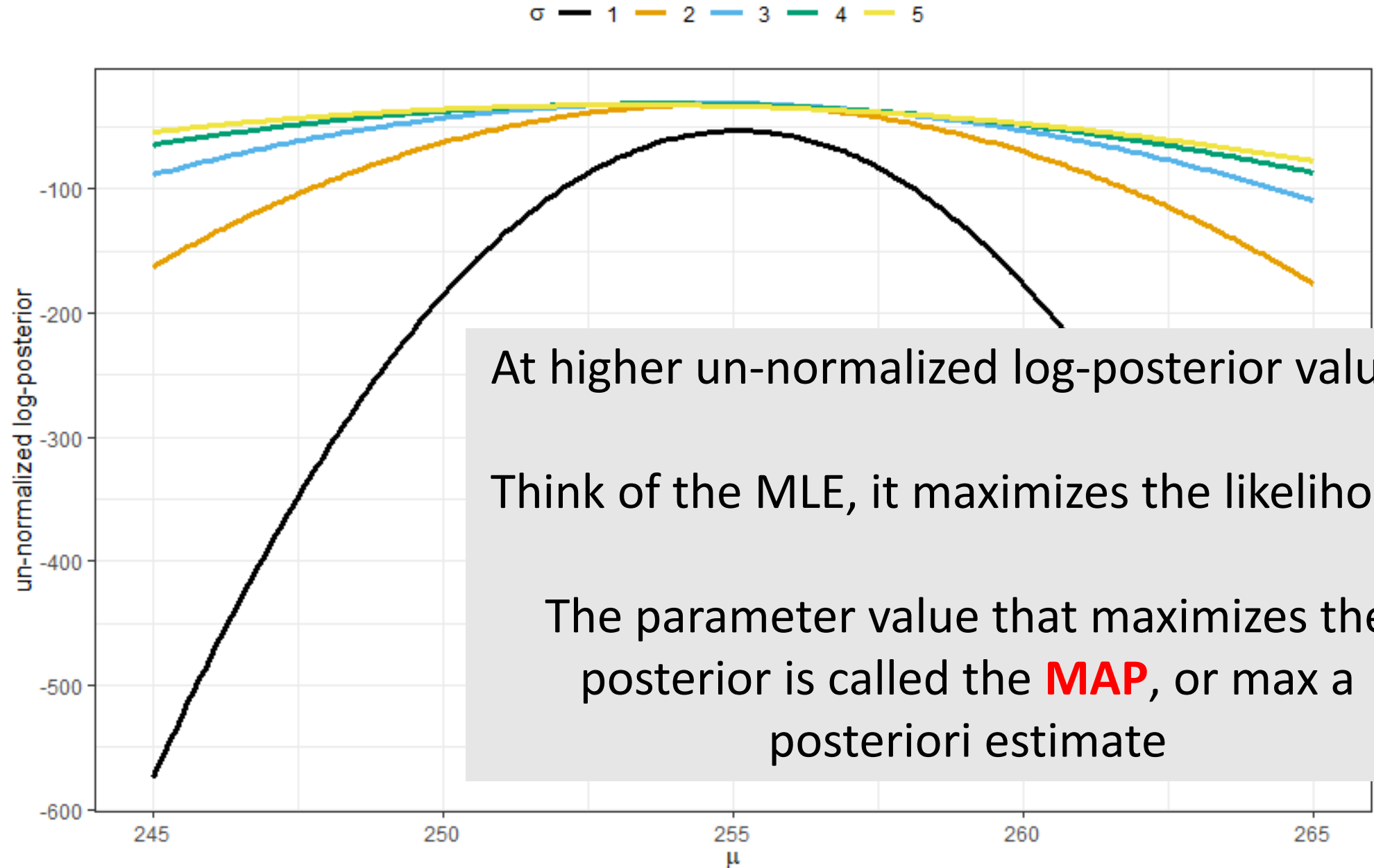
The log-posterior with respect to μ given specific values of σ should look familiar...



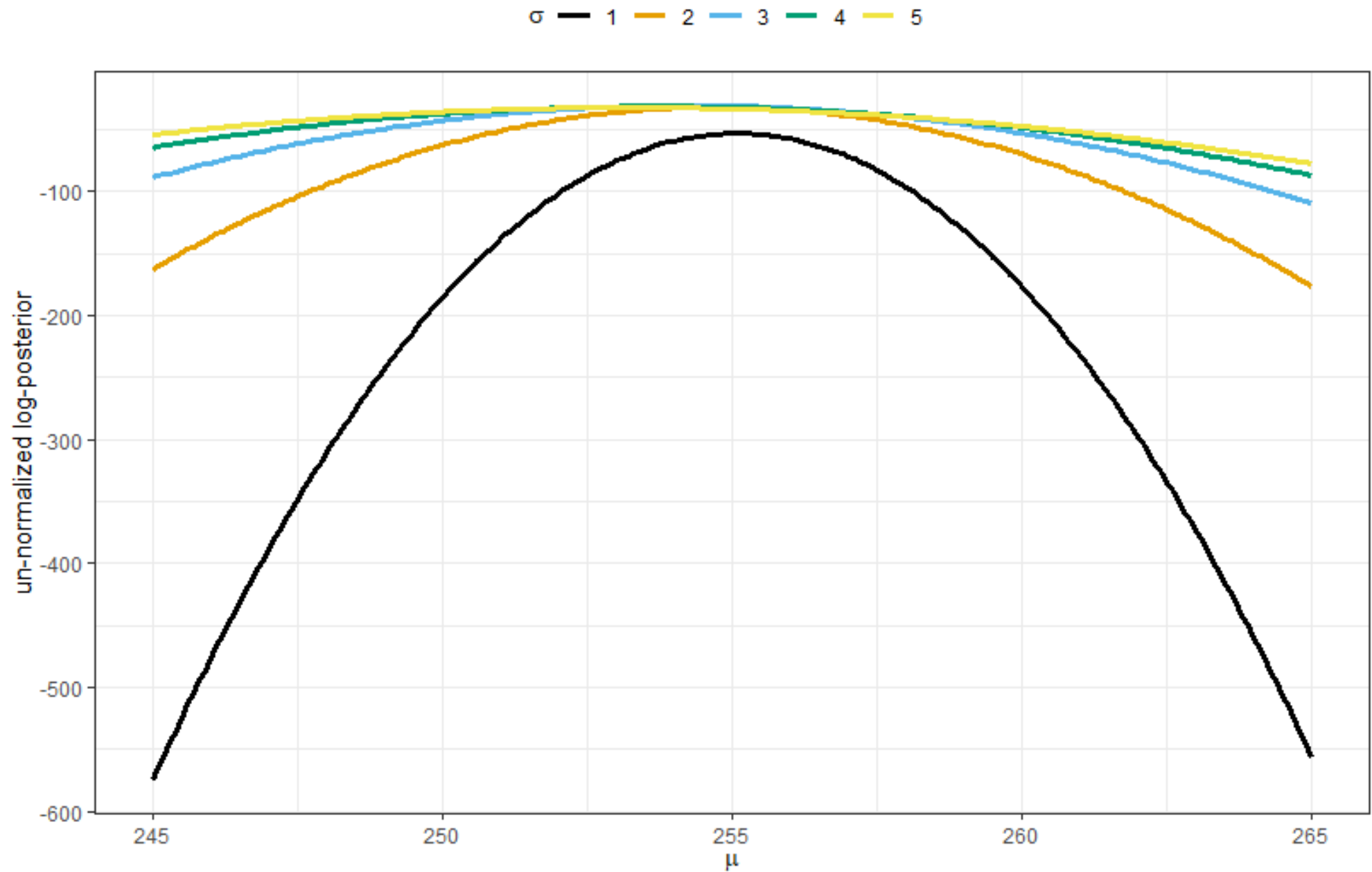
Posterior plausible values of μ are located where?



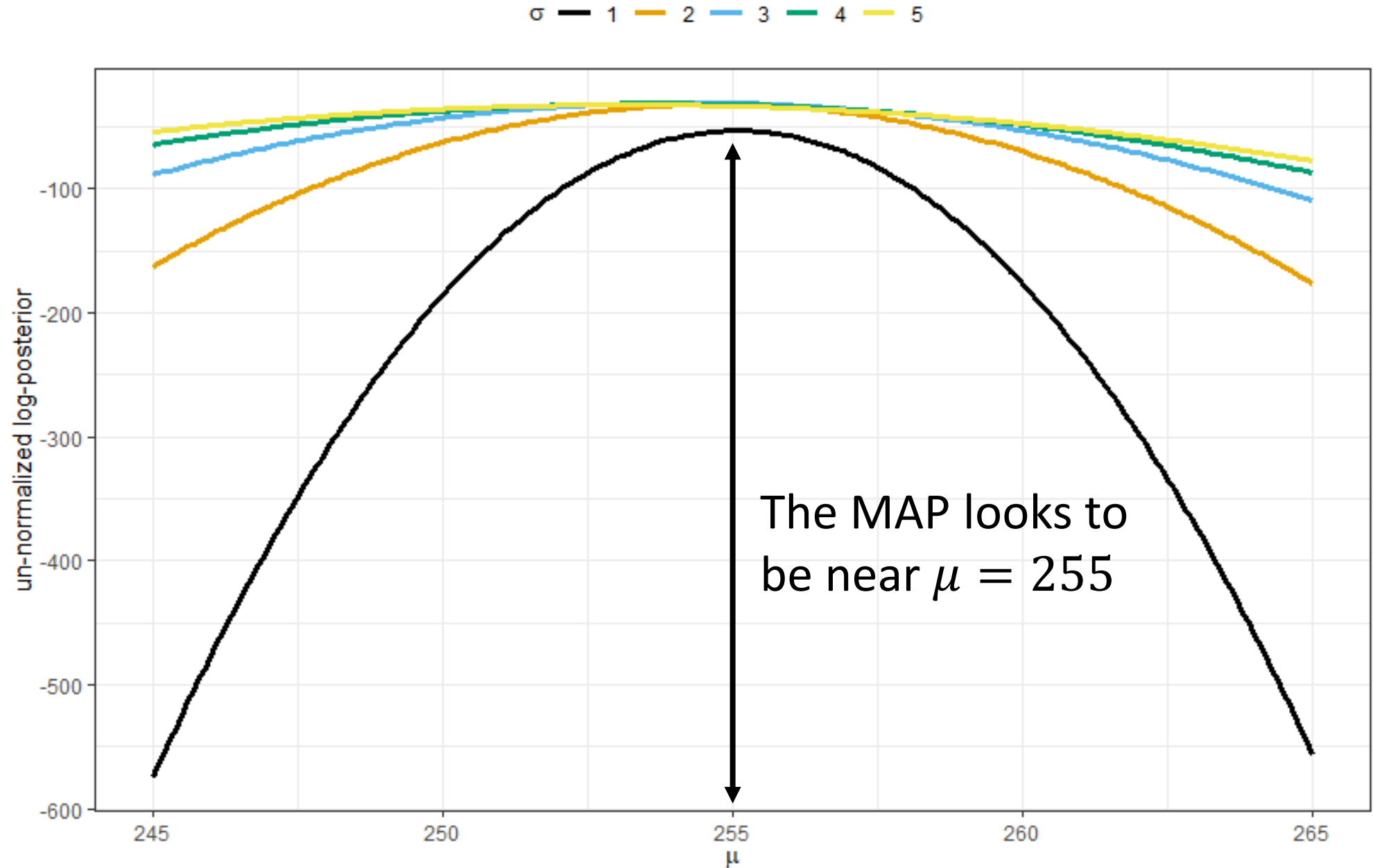
Posterior plausible values of μ are located where?



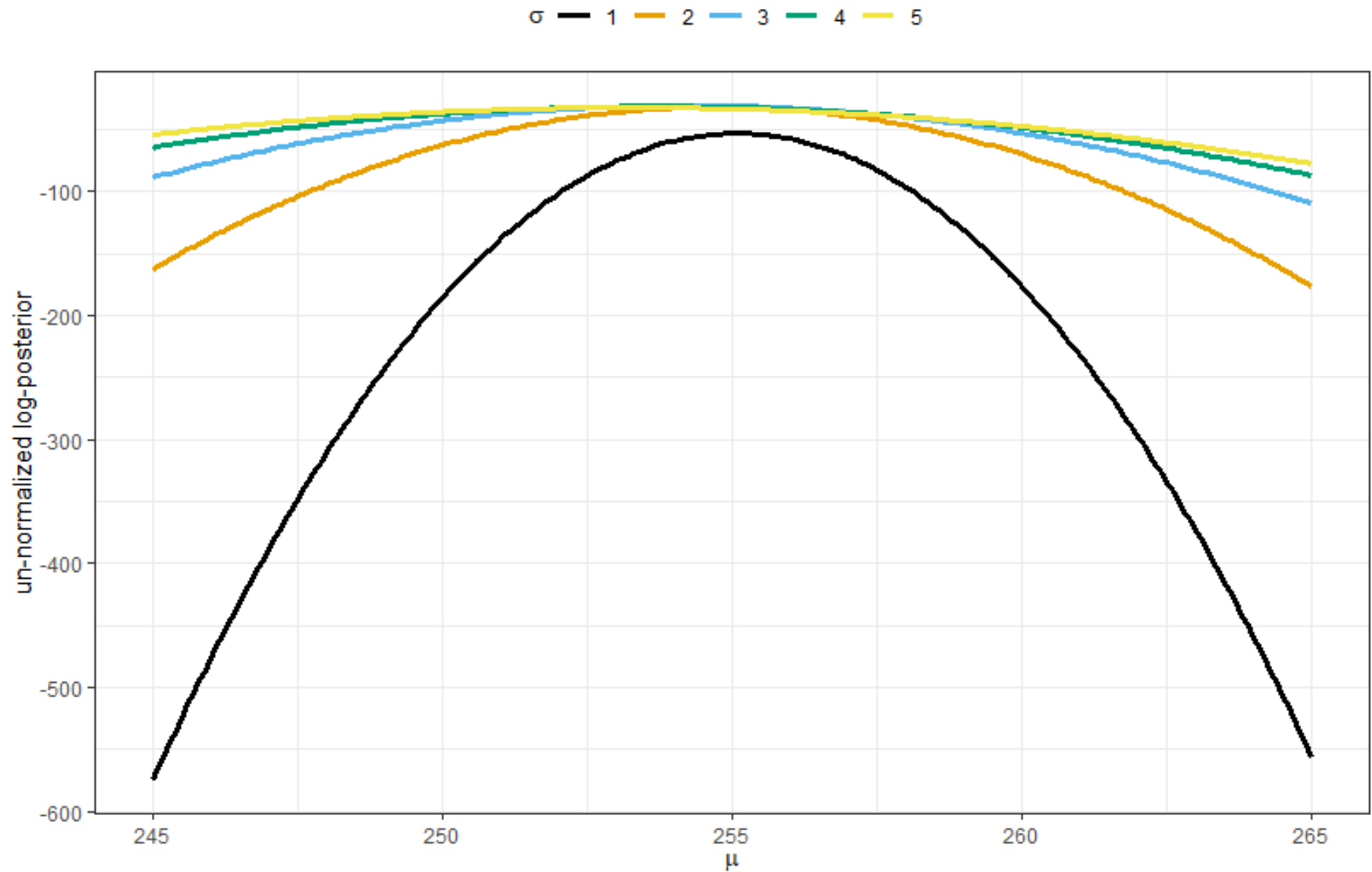
For $\sigma = 1$, what is the MAP on μ ?



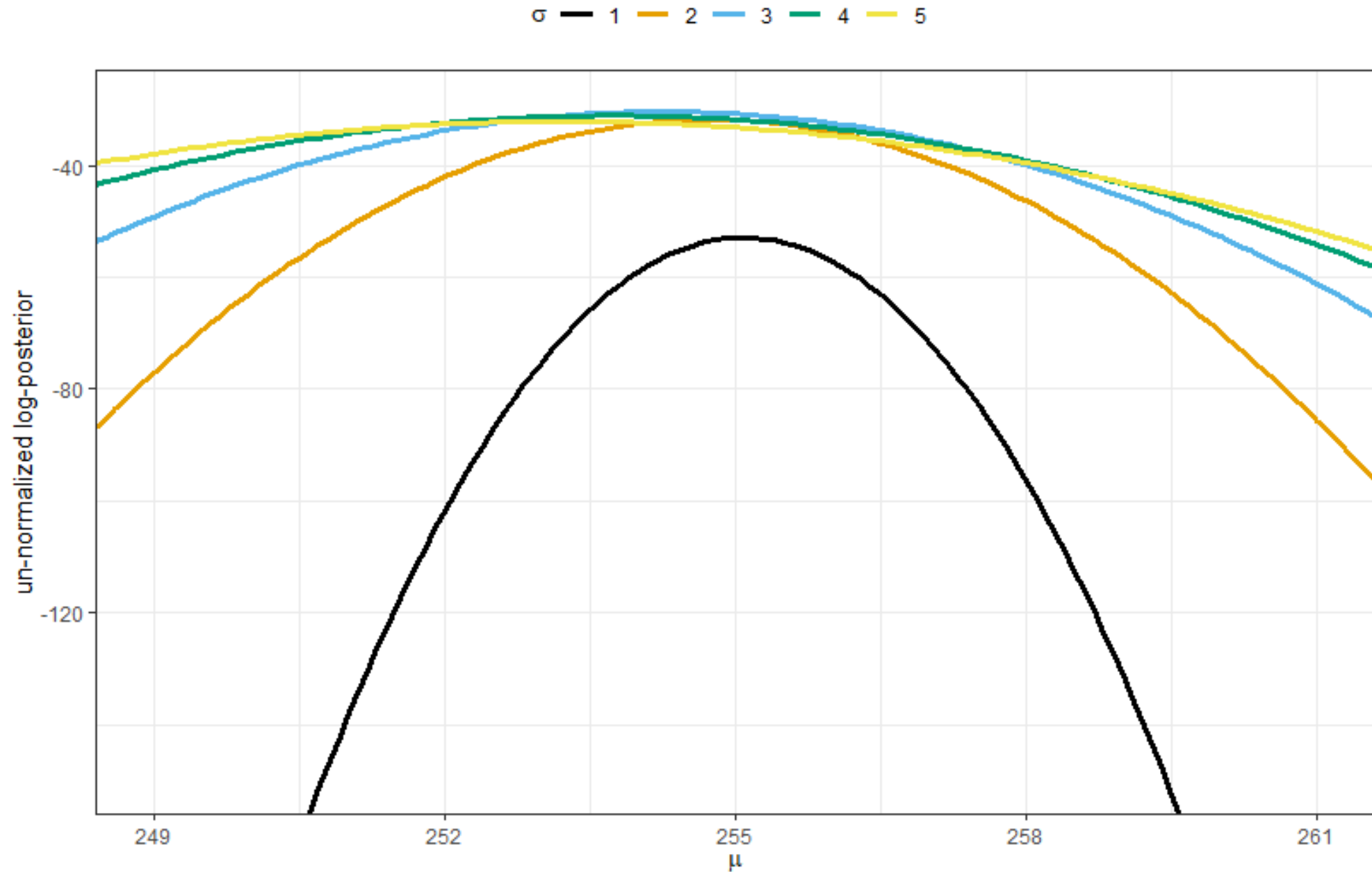
For $\sigma = 1$, what is the MAP on μ ?



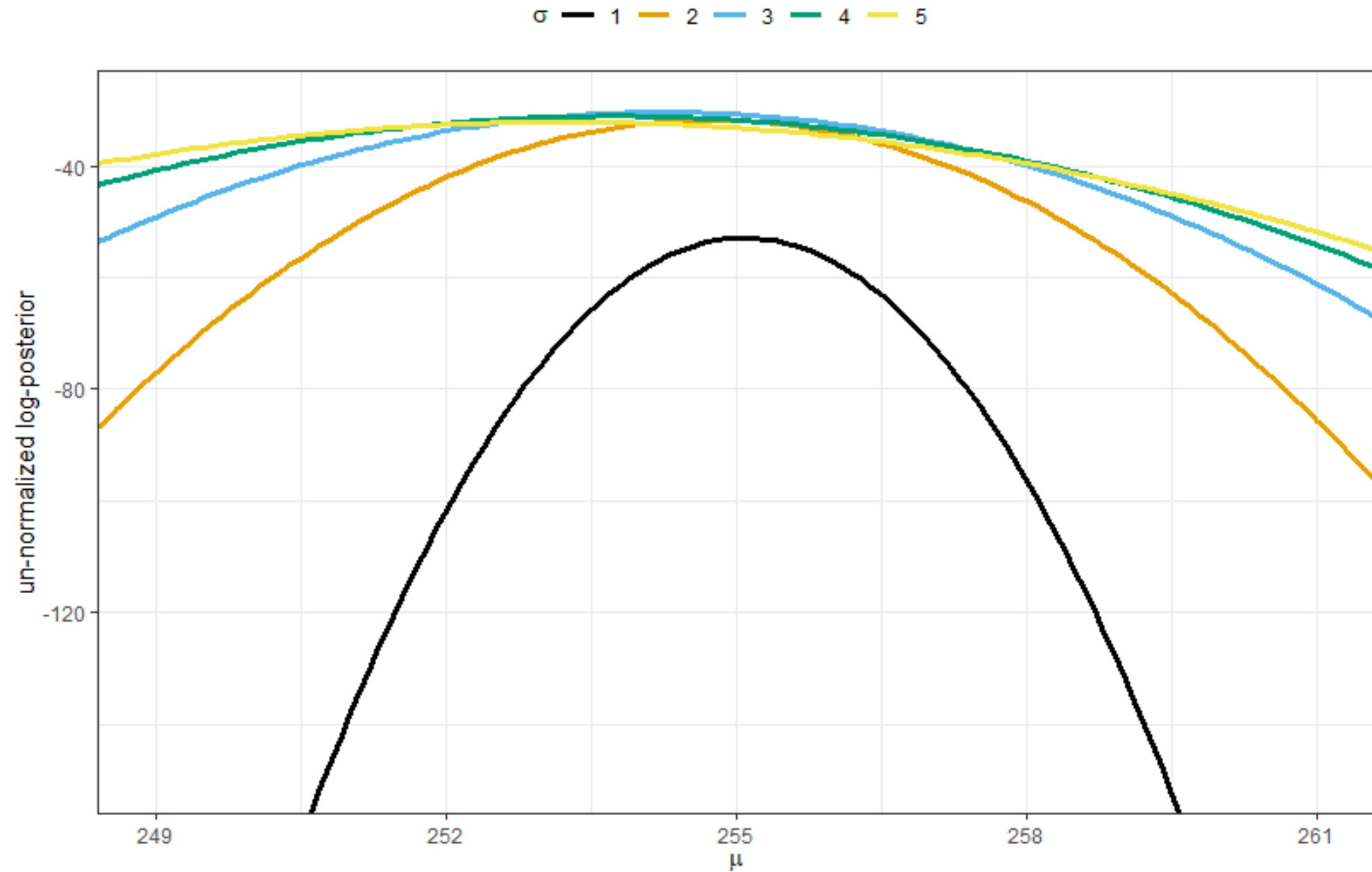
But...is $\sigma = 1$ more plausible than other values of σ ?



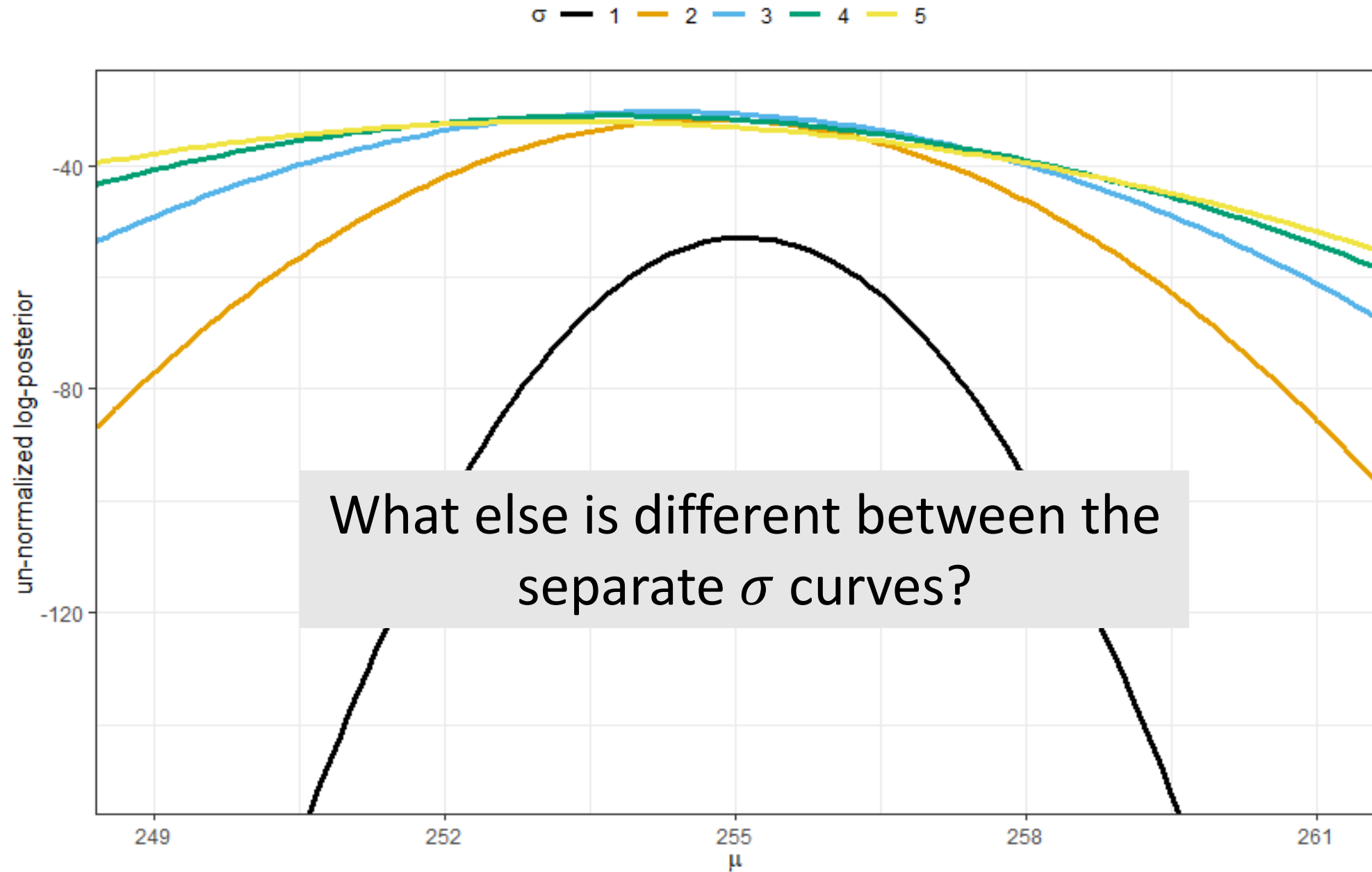
Zooming in, we see that the other curves have higher log-posterior values!



$\sigma > 1$ appears to be more plausible!

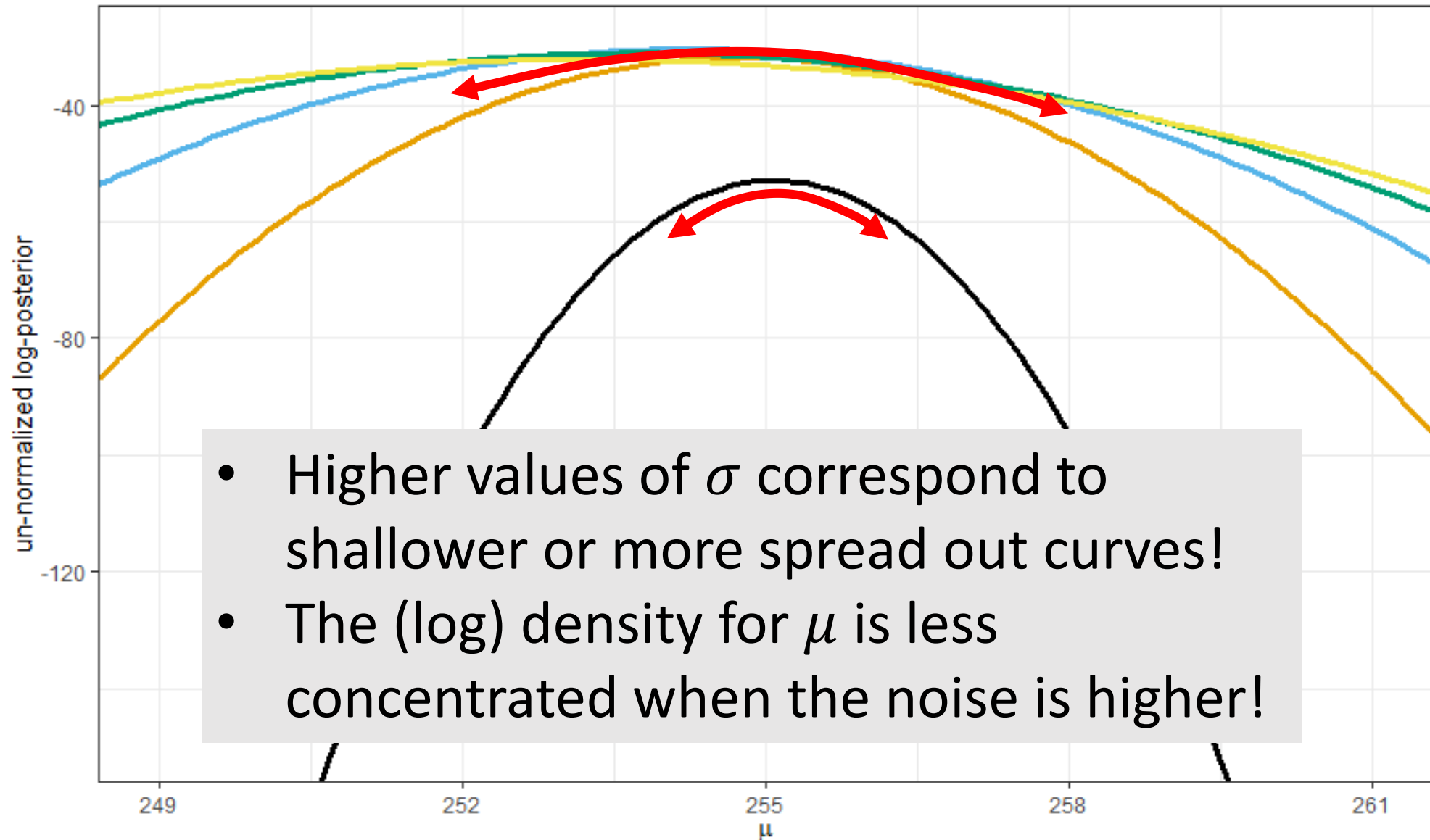


$\sigma > 1$ appears to be more plausible!

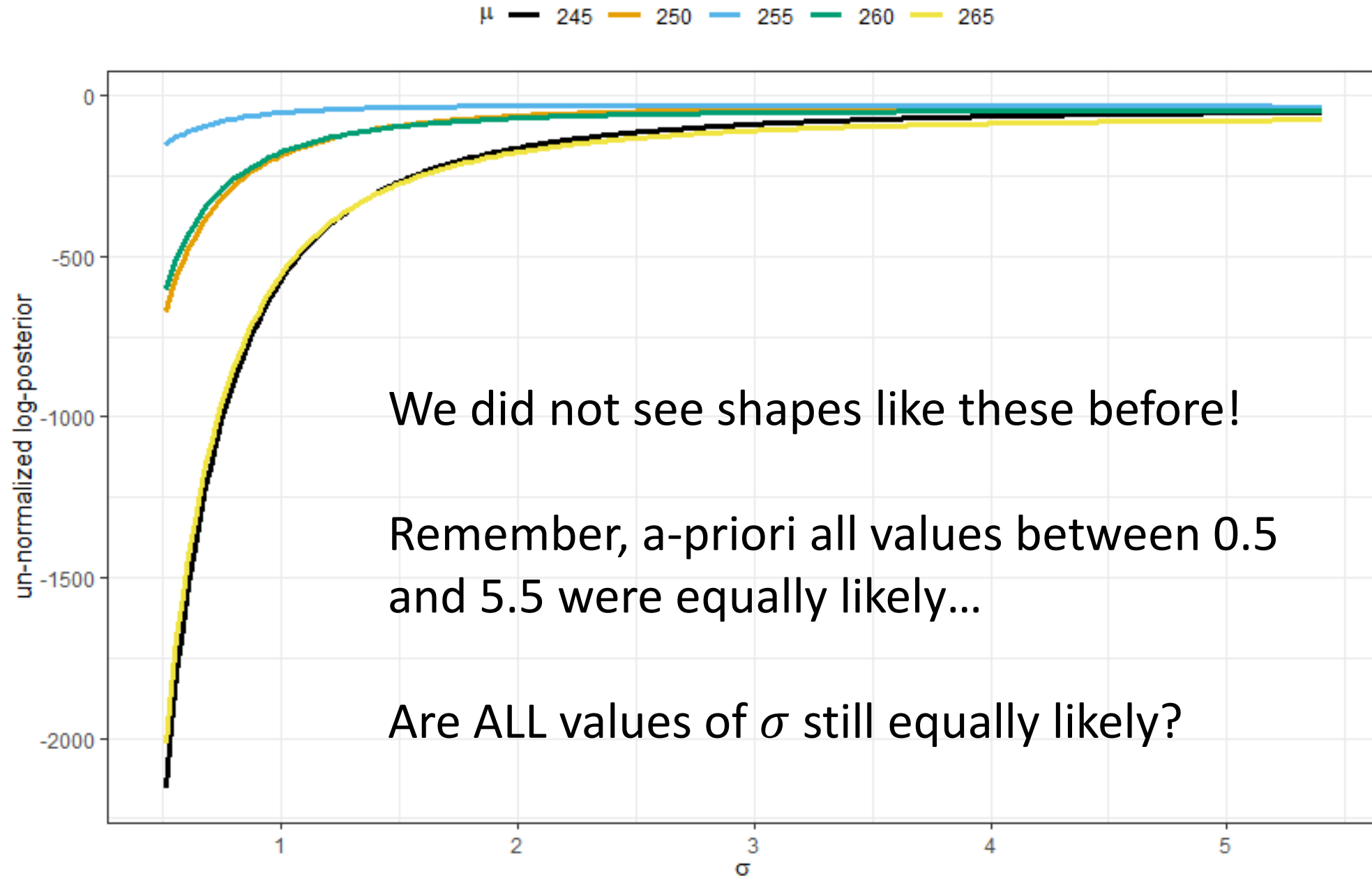


Look at the curvature!

σ — 1 — 2 — 3 — 4 — 5



Now consider the log-posterior with respect σ

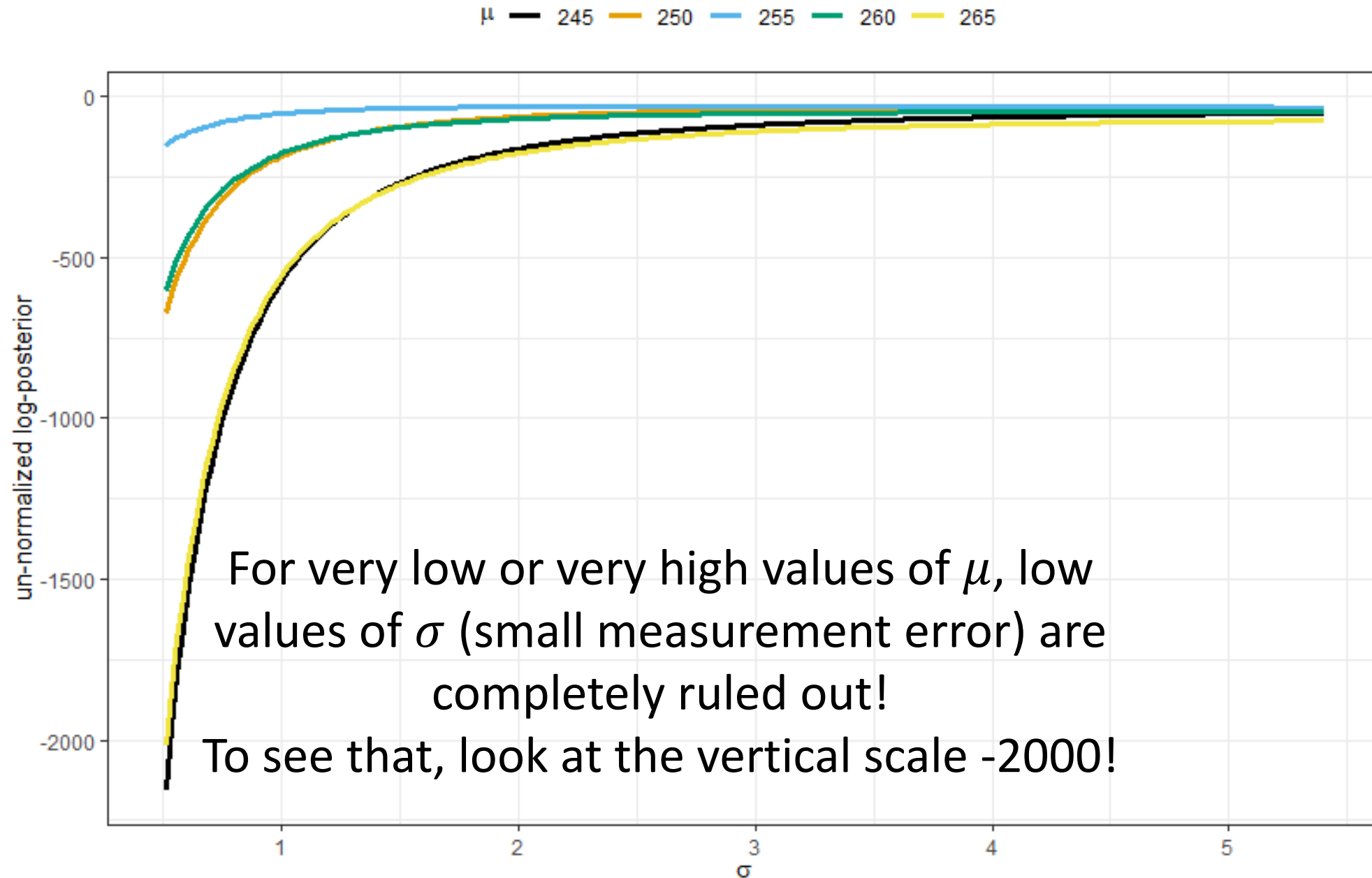


We did not see shapes like these before!

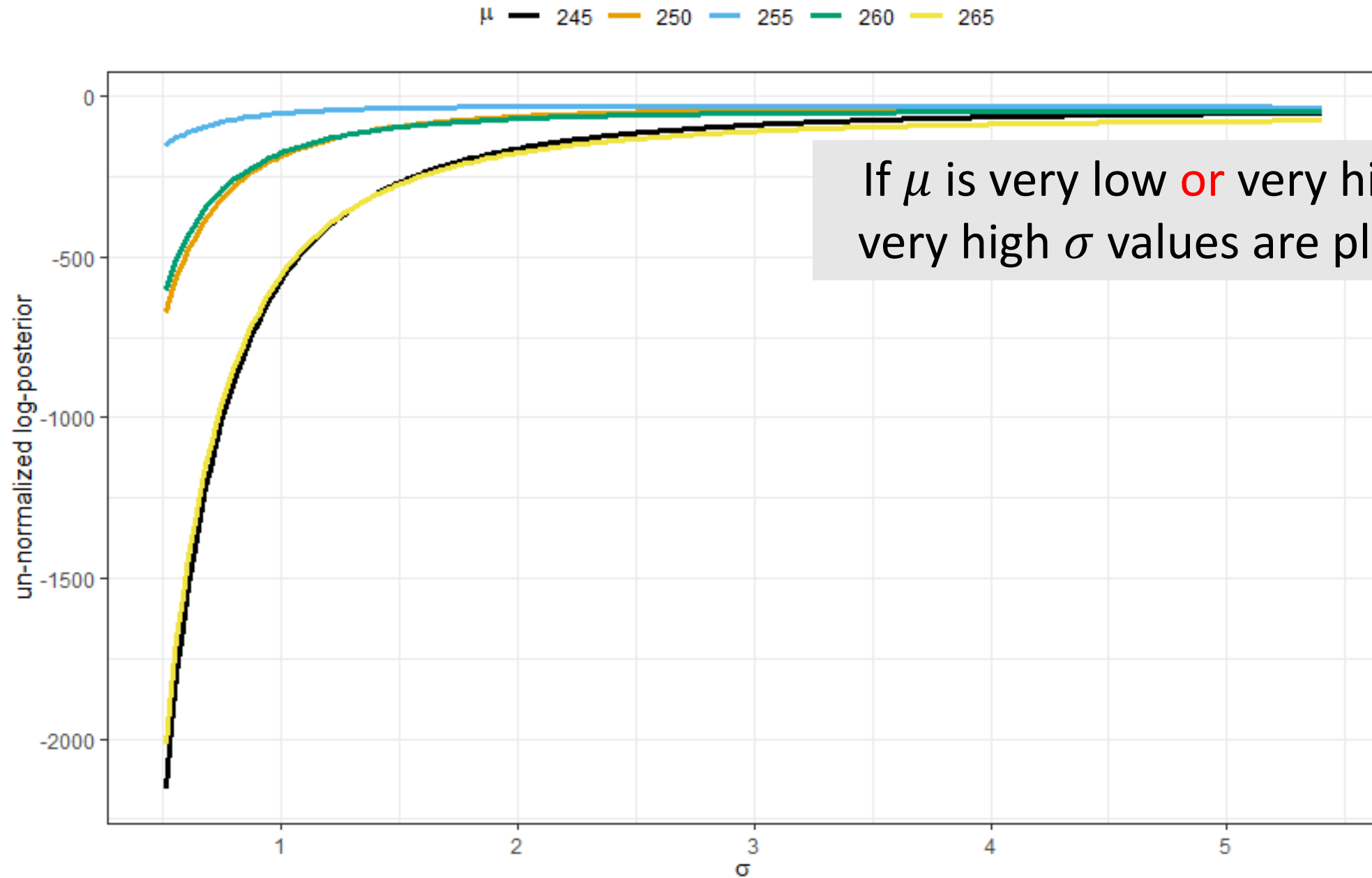
Remember, a-priori all values between 0.5 and 5.5 were equally likely...

Are ALL values of σ still equally likely?

Are all values of σ equally likely? The answer depends on μ .



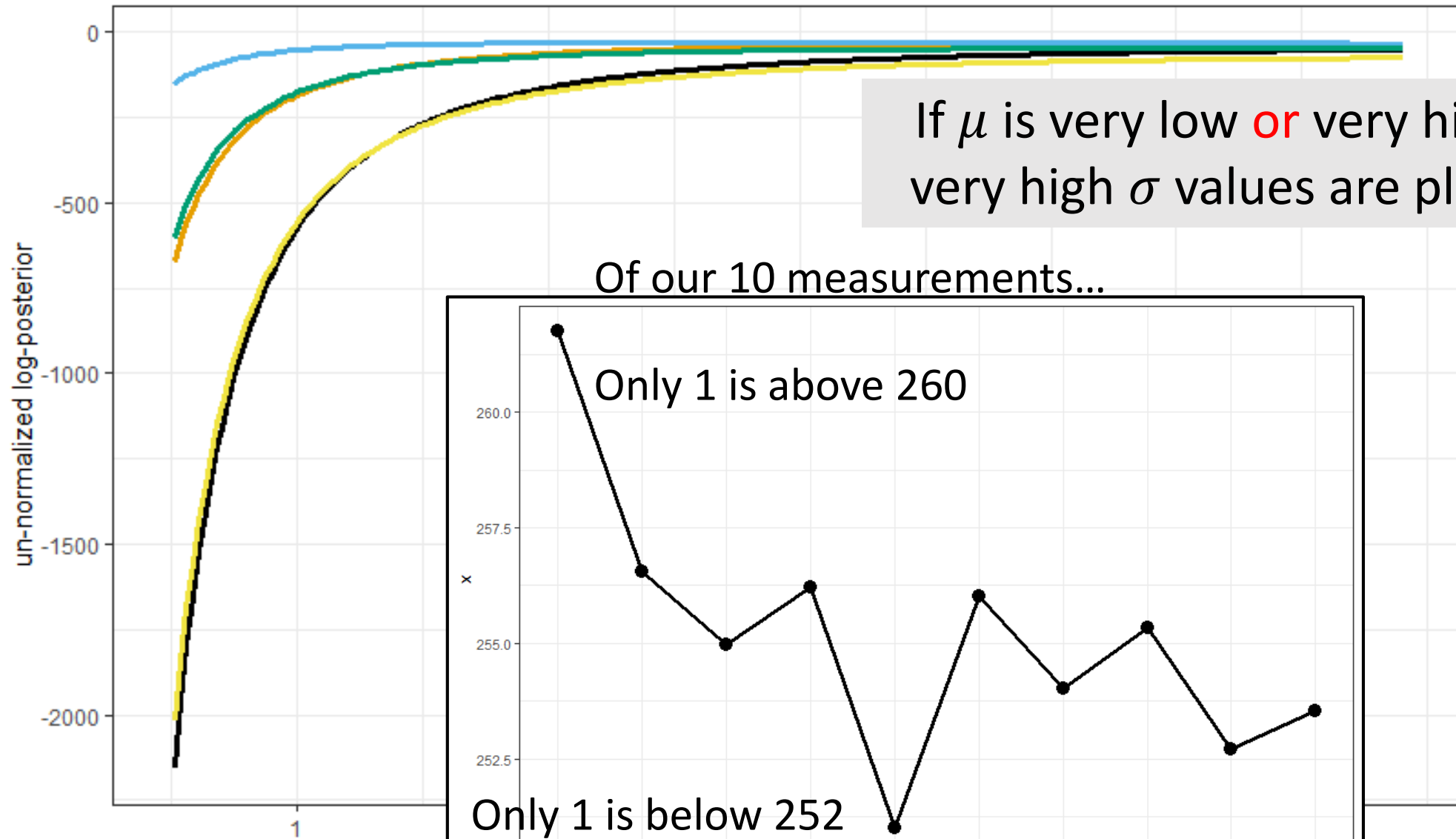
Are all values of σ equally likely? The answer depends on μ .



If μ is very low **or** very high only very high σ values are plausible!

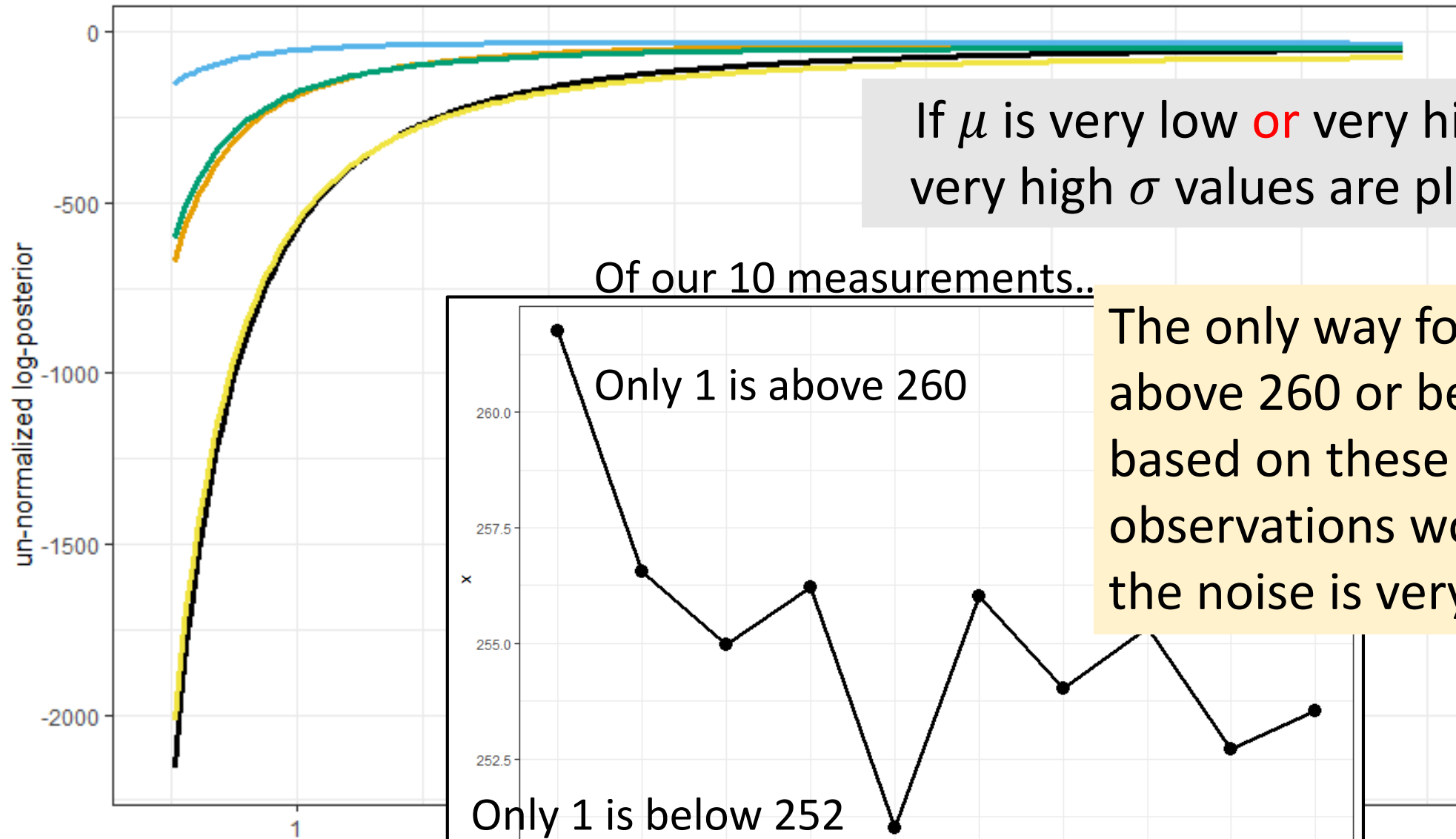
Are all values of σ equally likely? The answer depends on μ .

μ — 245 — 250 — 255 — 260 — 265

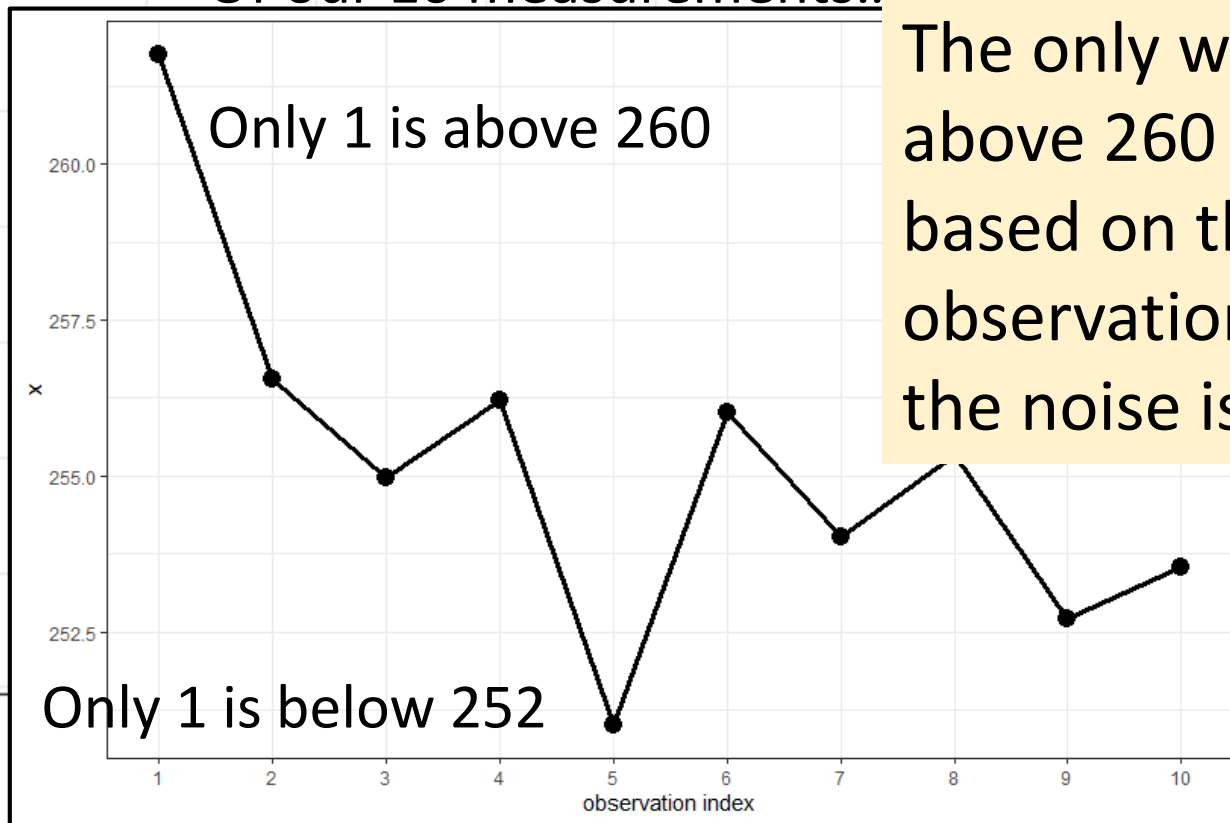


Are all values of σ equally likely? The answer depends on μ .

μ — 245 — 250 — 255 — 260 — 265

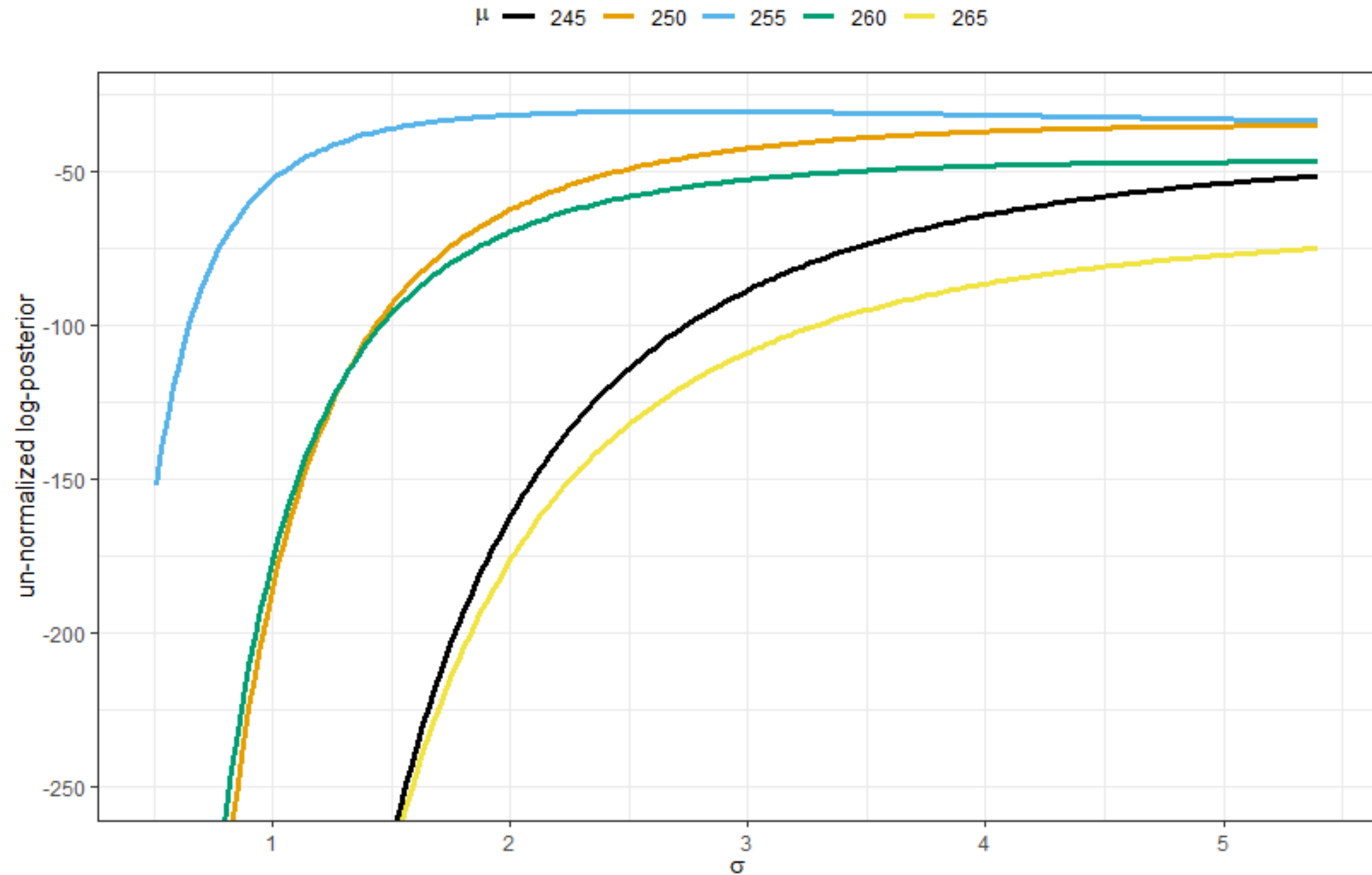


Of our 10 measurements..

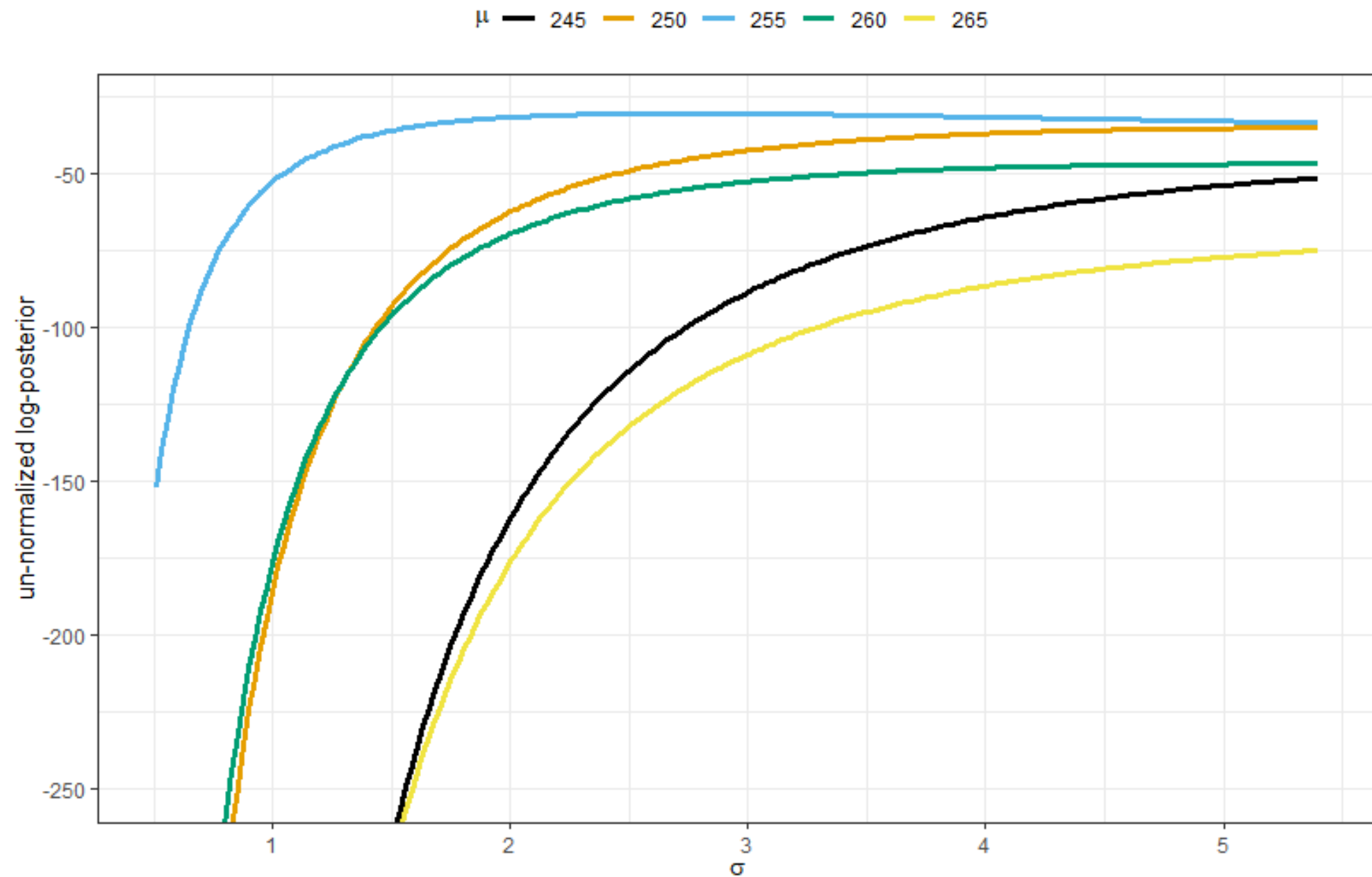


The only way for μ to be above 260 or below 250 based on these 10 observations would be if the noise is very high!

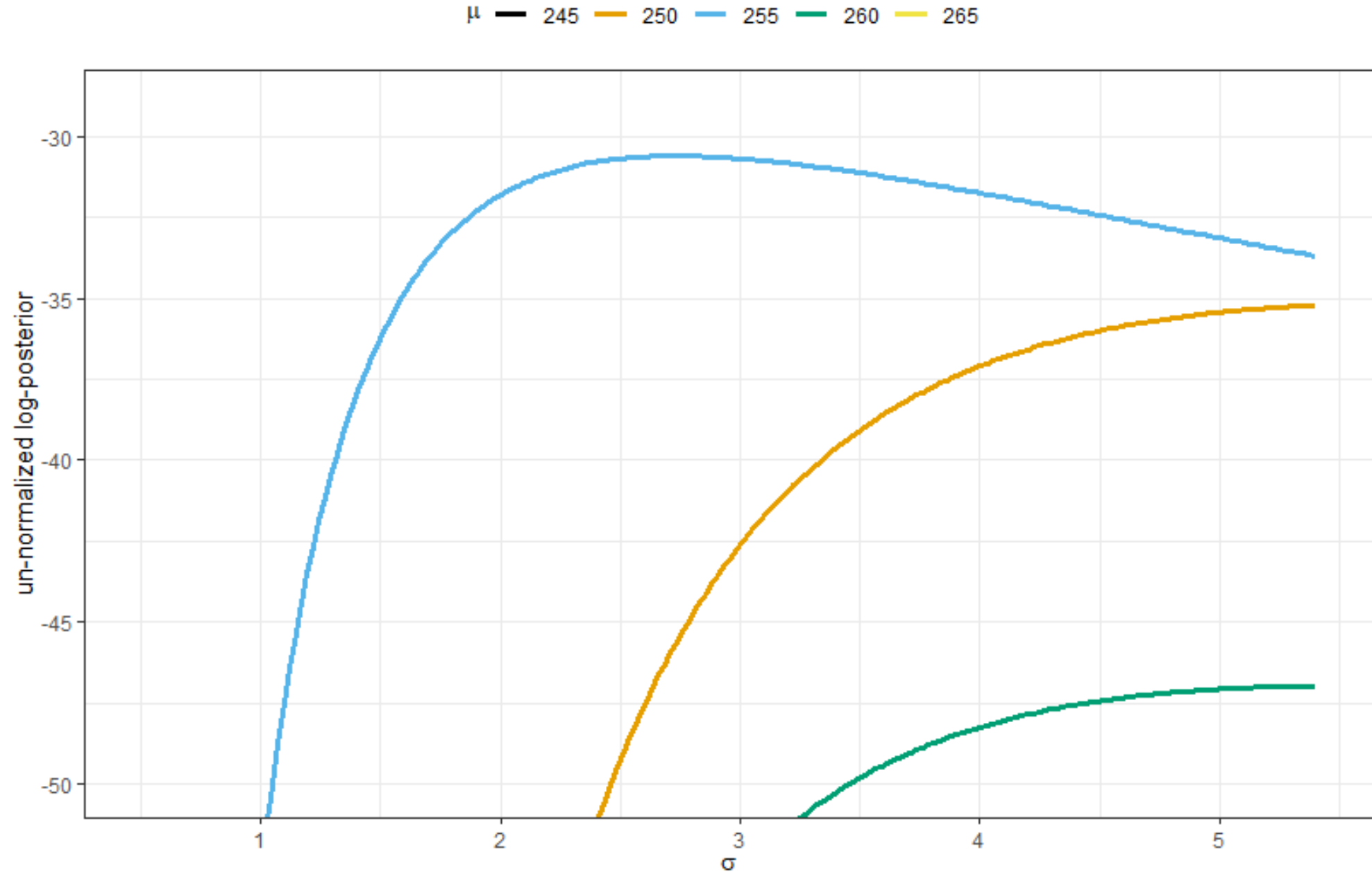
Zooming in we can see the very low or very high μ values are not as plausible as others!



σ does not seem to have as sharply defined MAP as we saw on μ ...

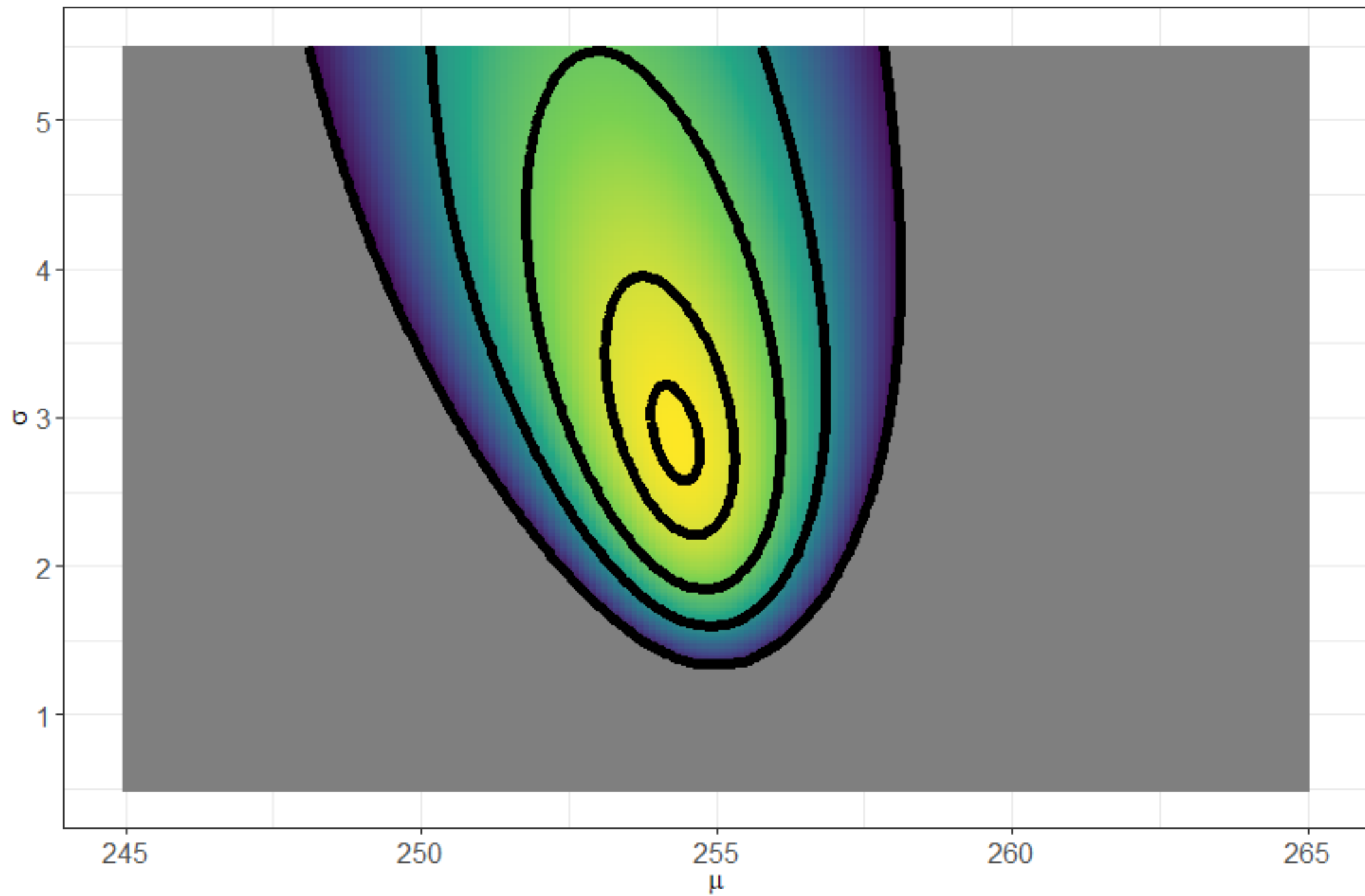


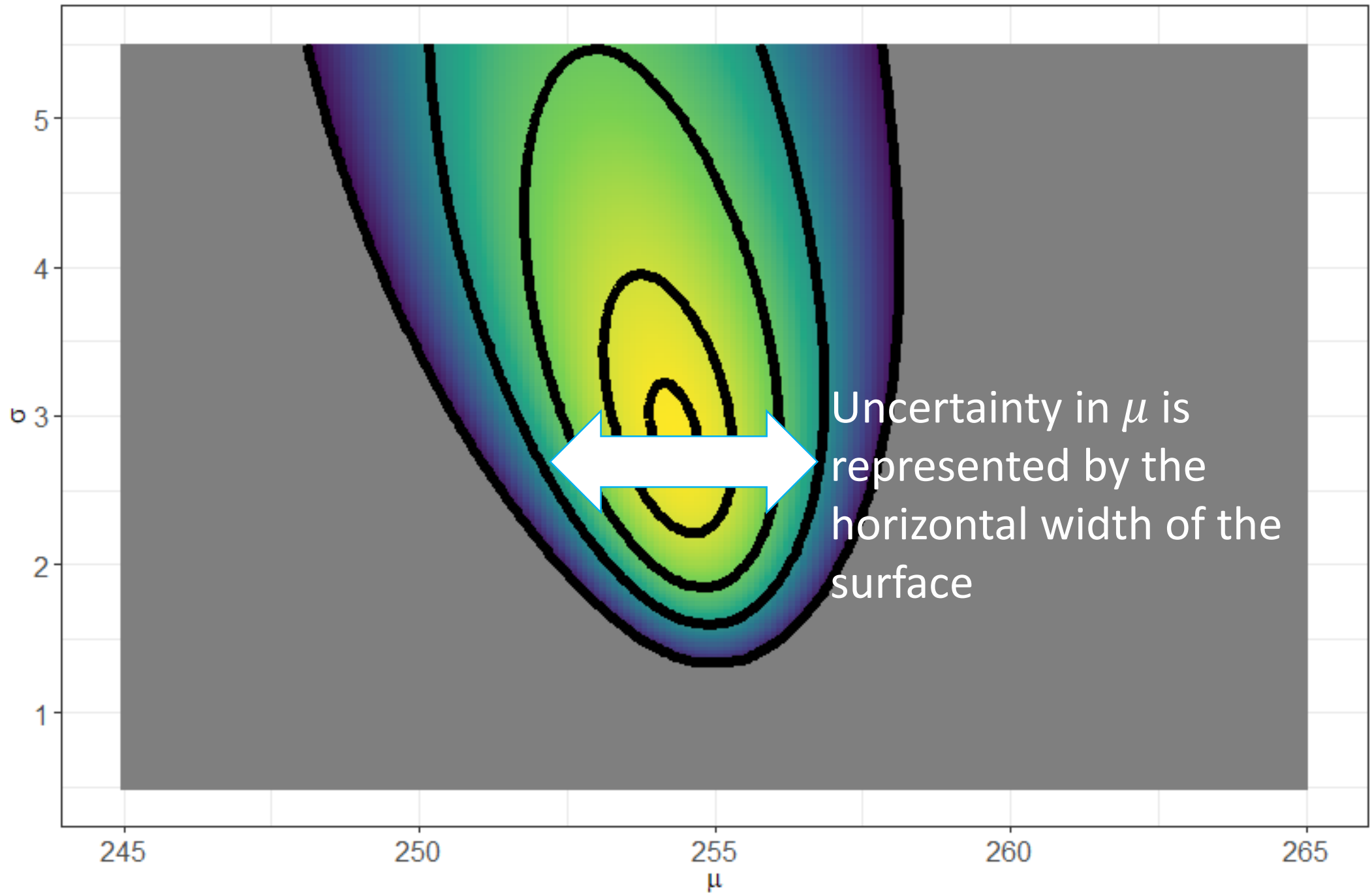
BUT...remember the vertical scale! Zoom in further to reveal the MAP!

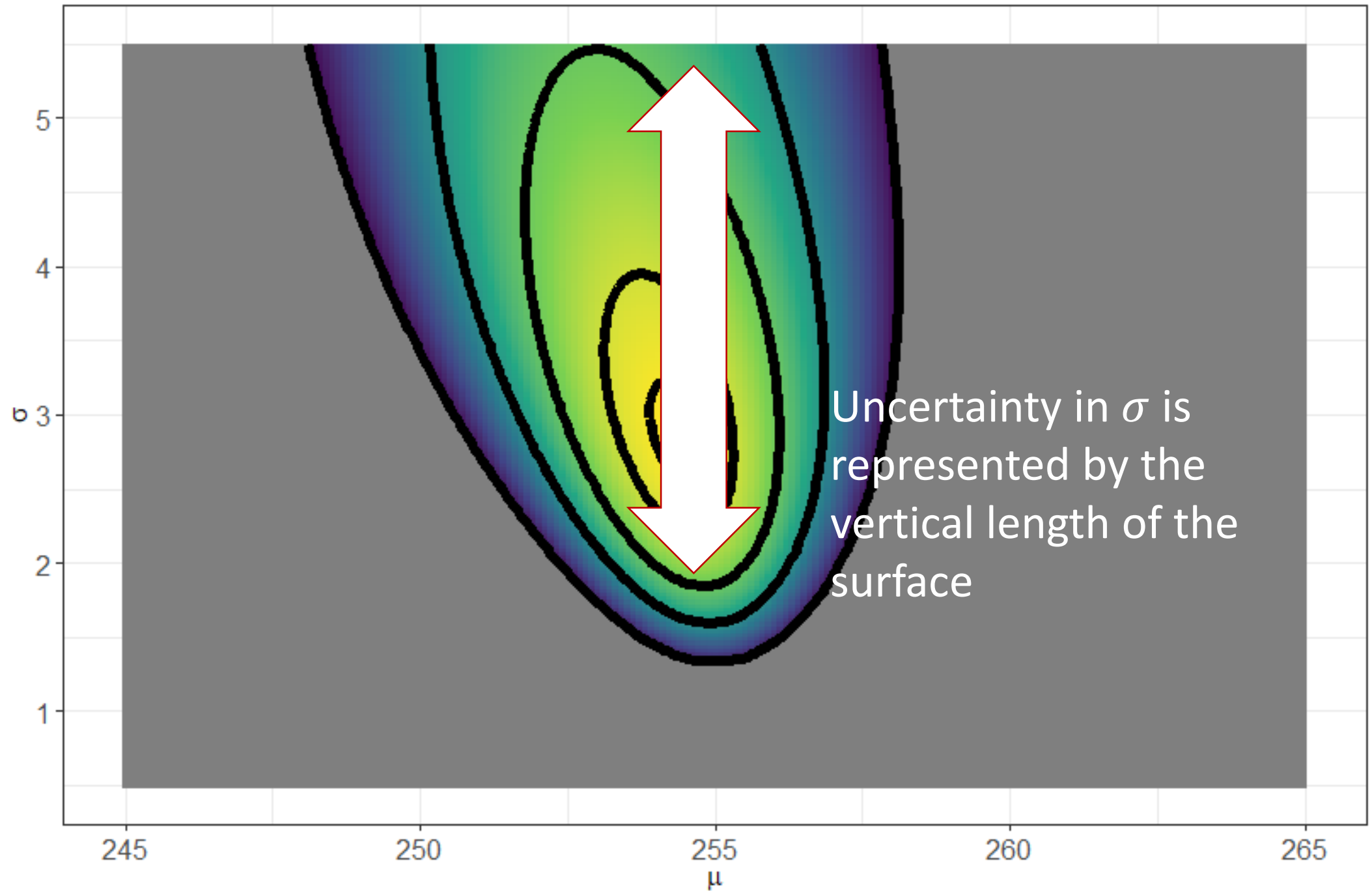


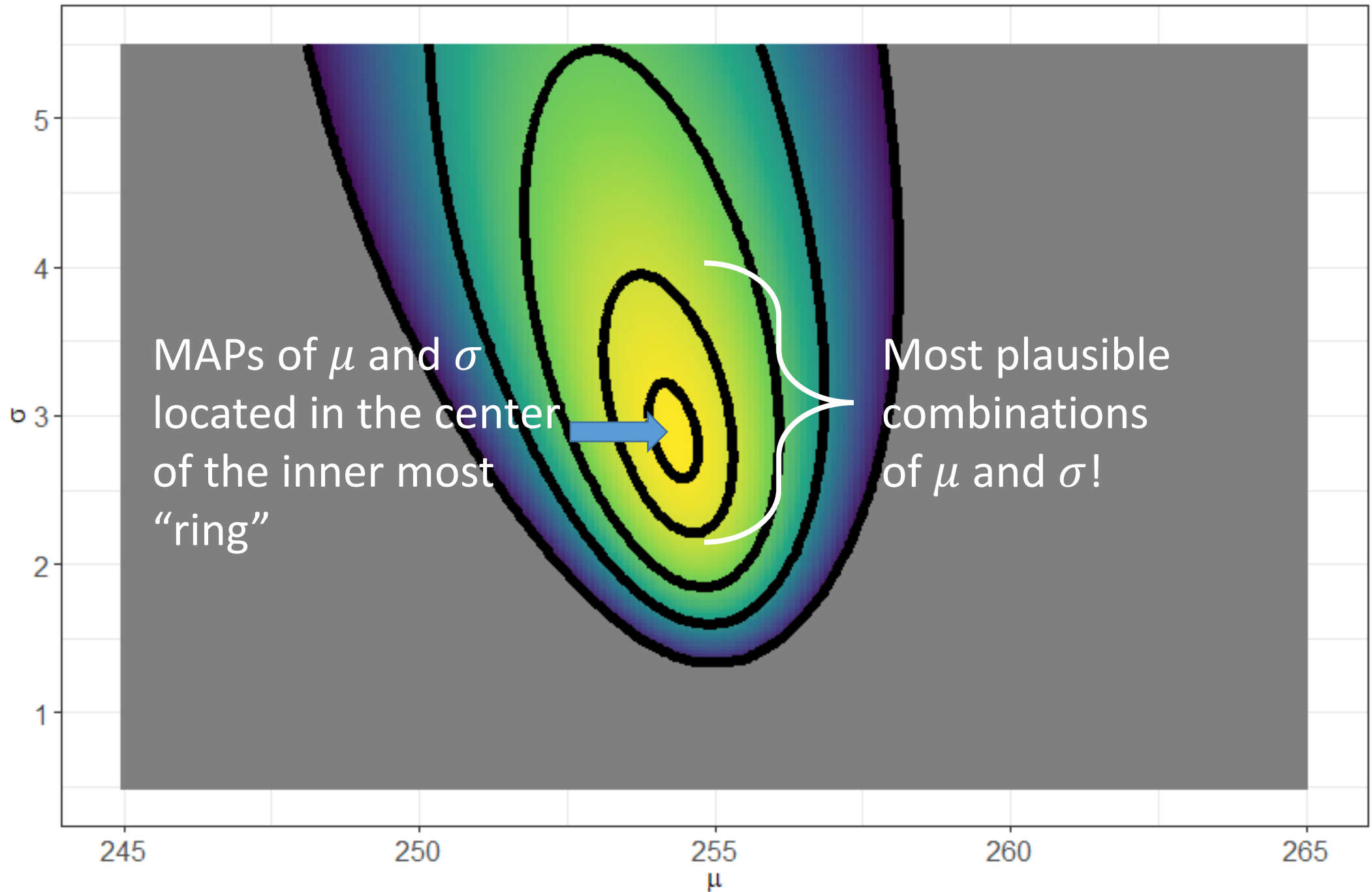
After all of that...let's now look at the surface

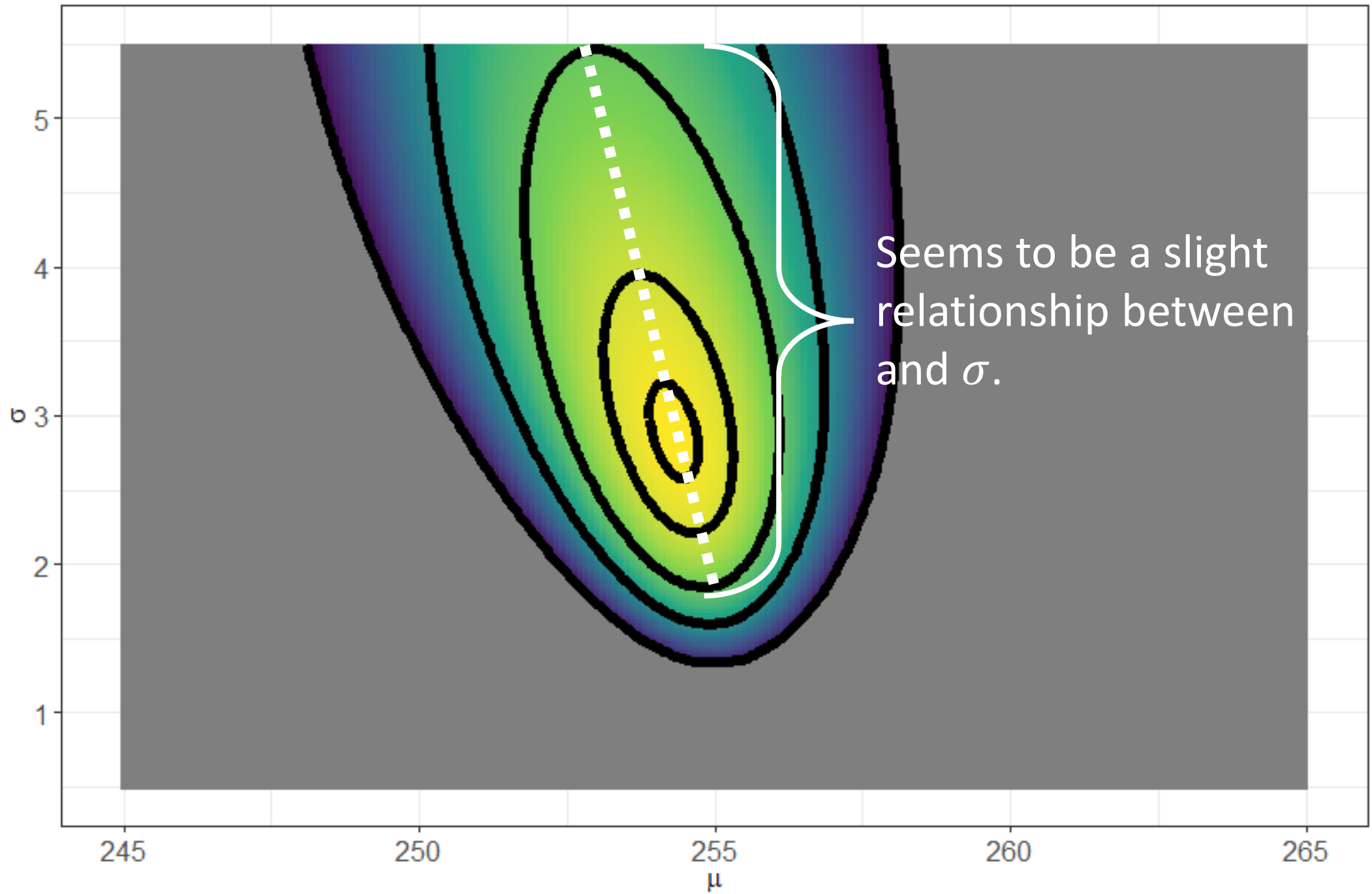
- The log-posterior surface with respect to μ and σ will be displayed similar to how the log-prior surface was displayed.
- Fill represents the value un-normalized log-posterior, bright yellow are higher values and dark blue are lower values.
- **Combinations that are very implausible are greyed out.**
- Contour lines represent the top 10%, 50%, 90%, 99%, and 99.99% of combinations.

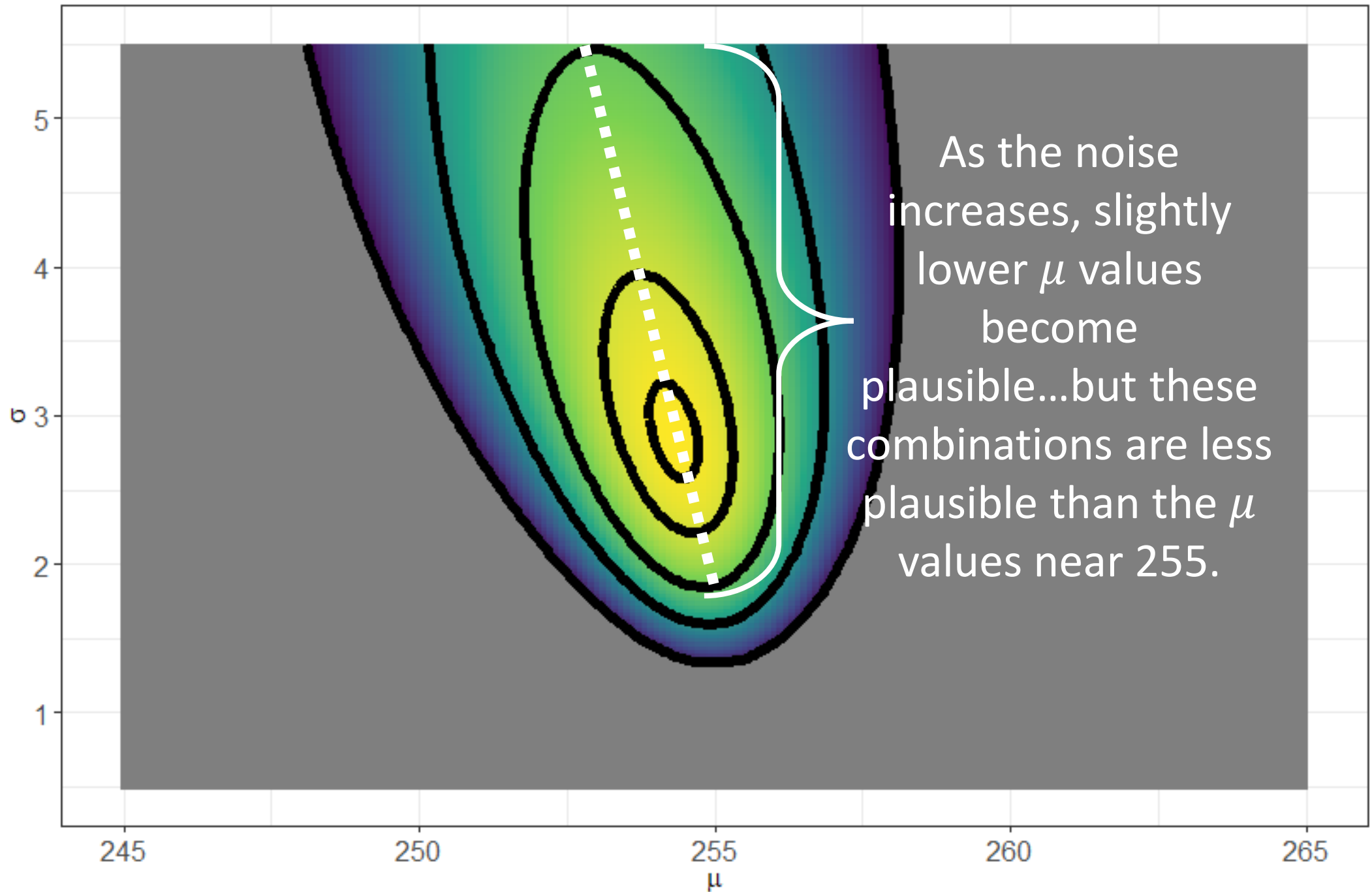






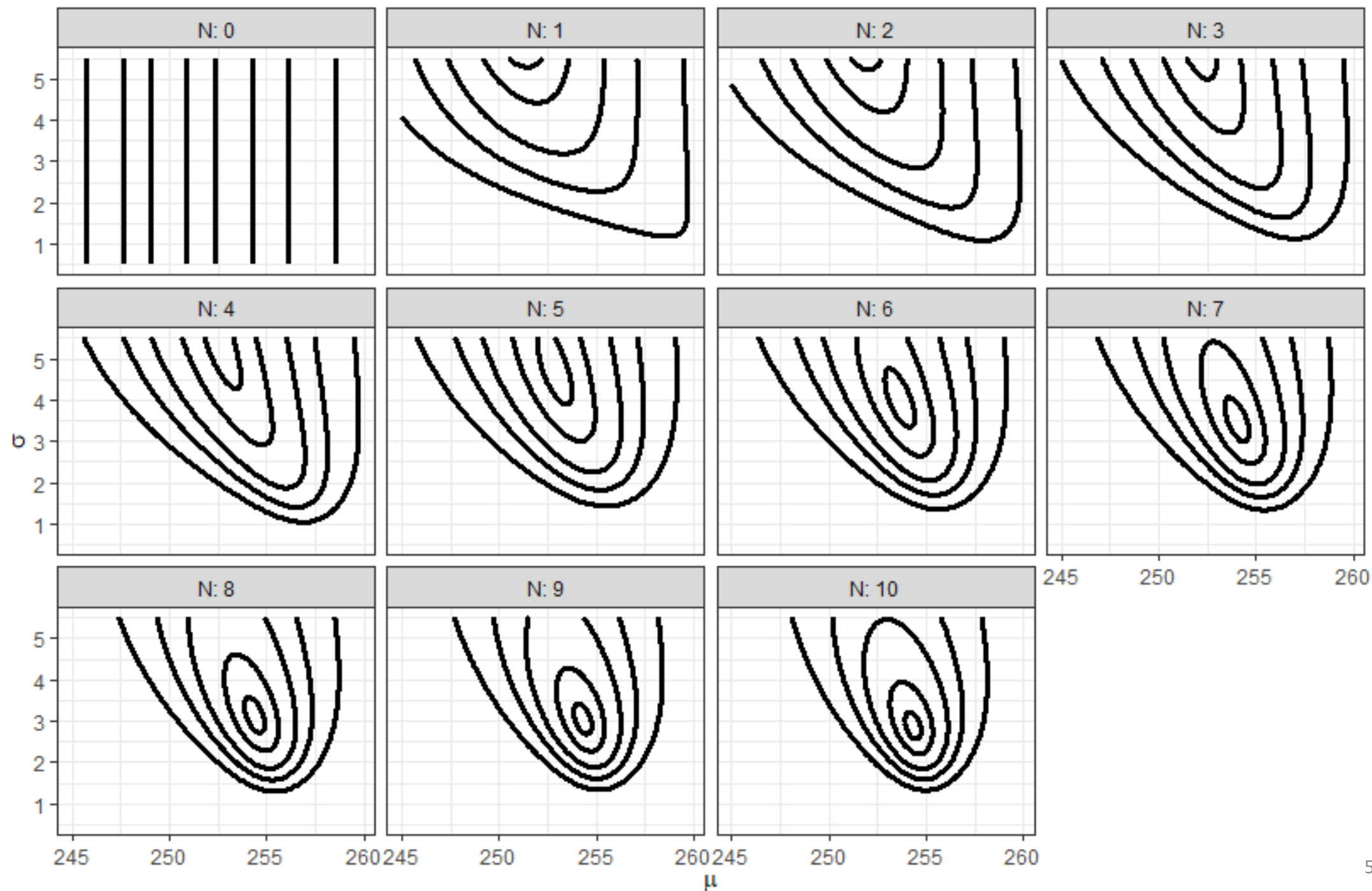




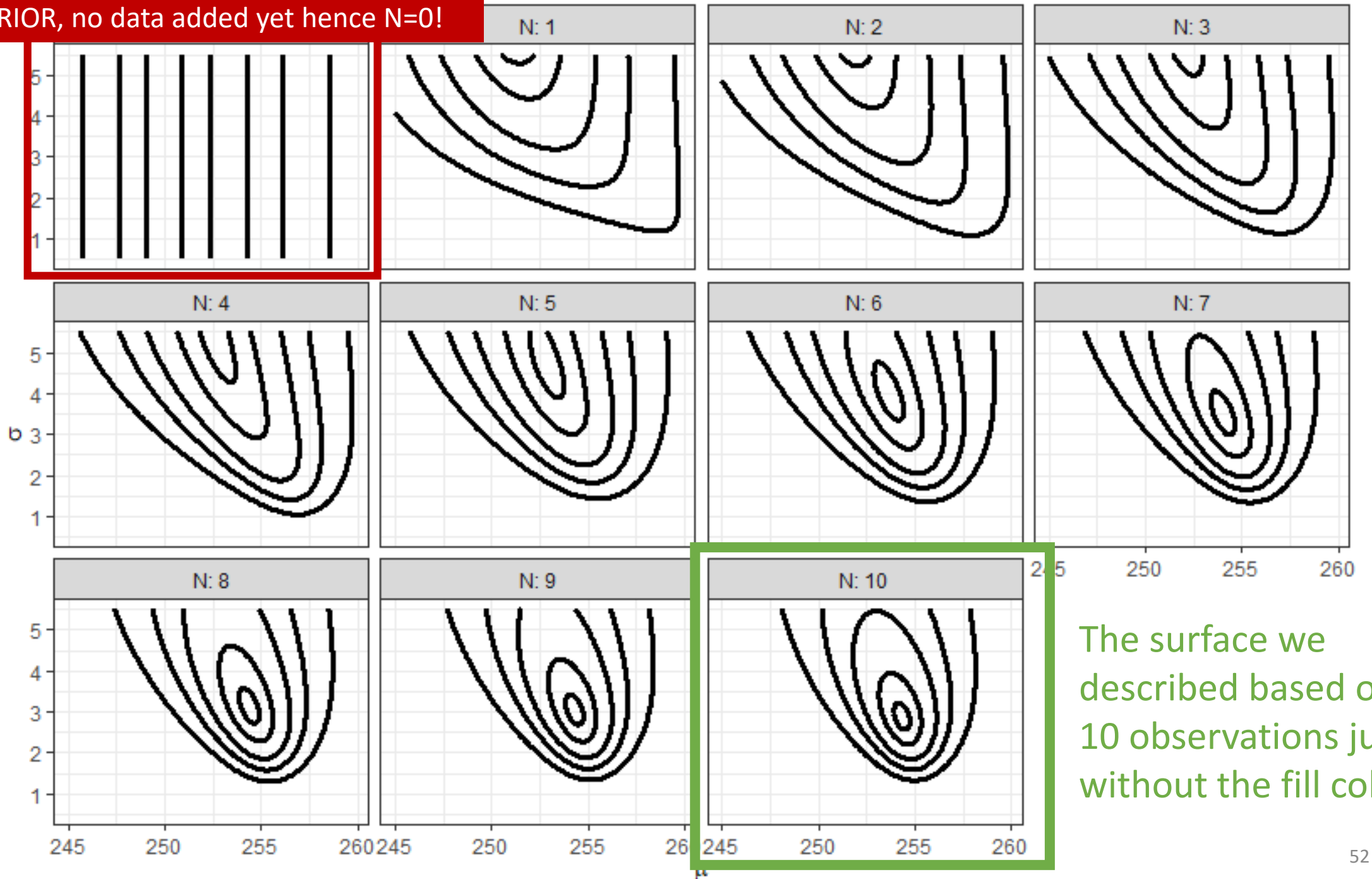


We just visualized the joint posterior distribution of 2 unknown parameters!

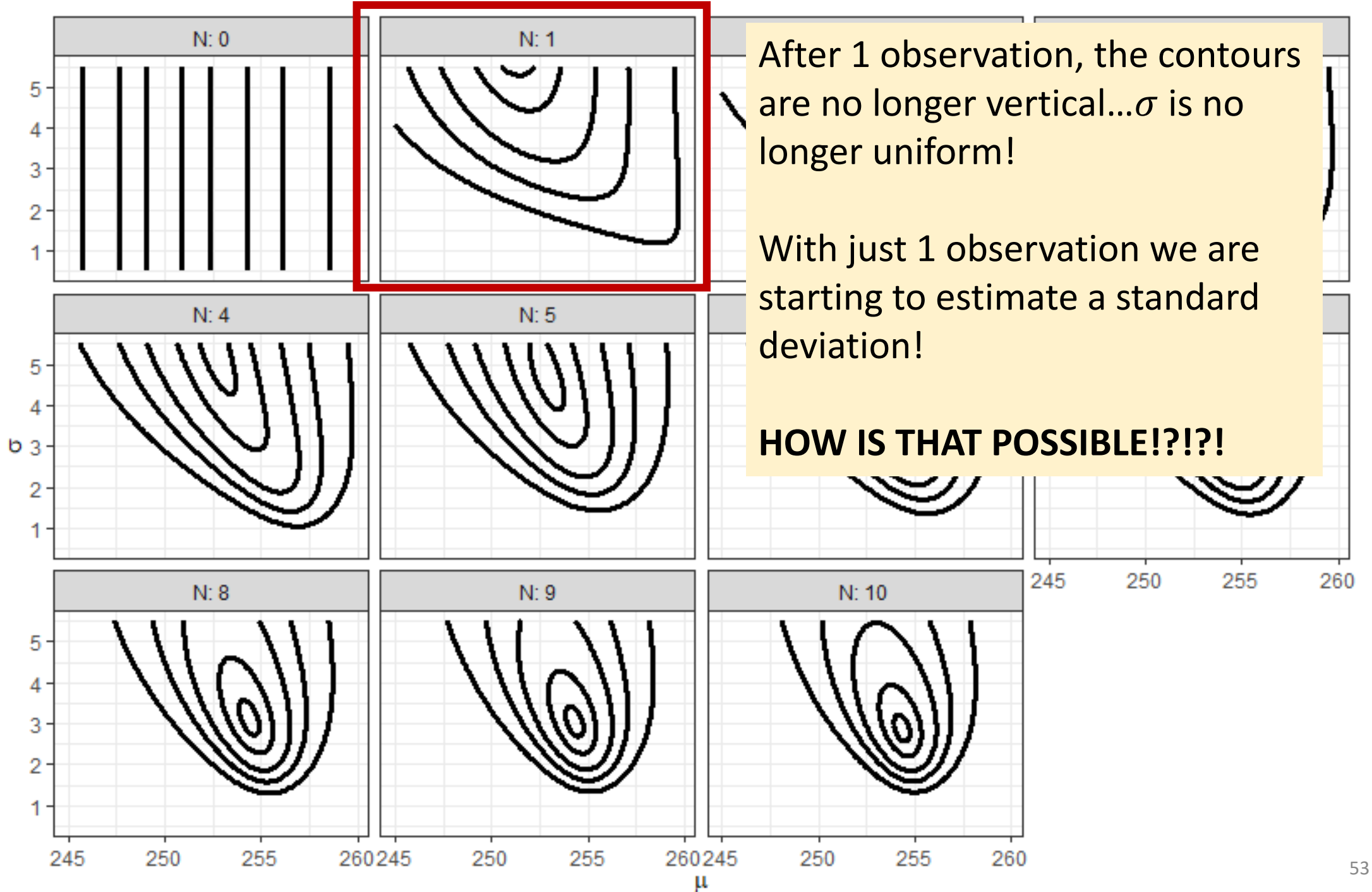
- It was based on our prior beliefs and 10 observations.
- How did the prior morph or evolve into the posterior we just visualized?
- **Let's now step through how the log-posterior surface changes as we sequentially add each data point.**
- The following figures show the contours, but do not fill in the color.
- The contours represent the same iso-probability contours as the previous figure.

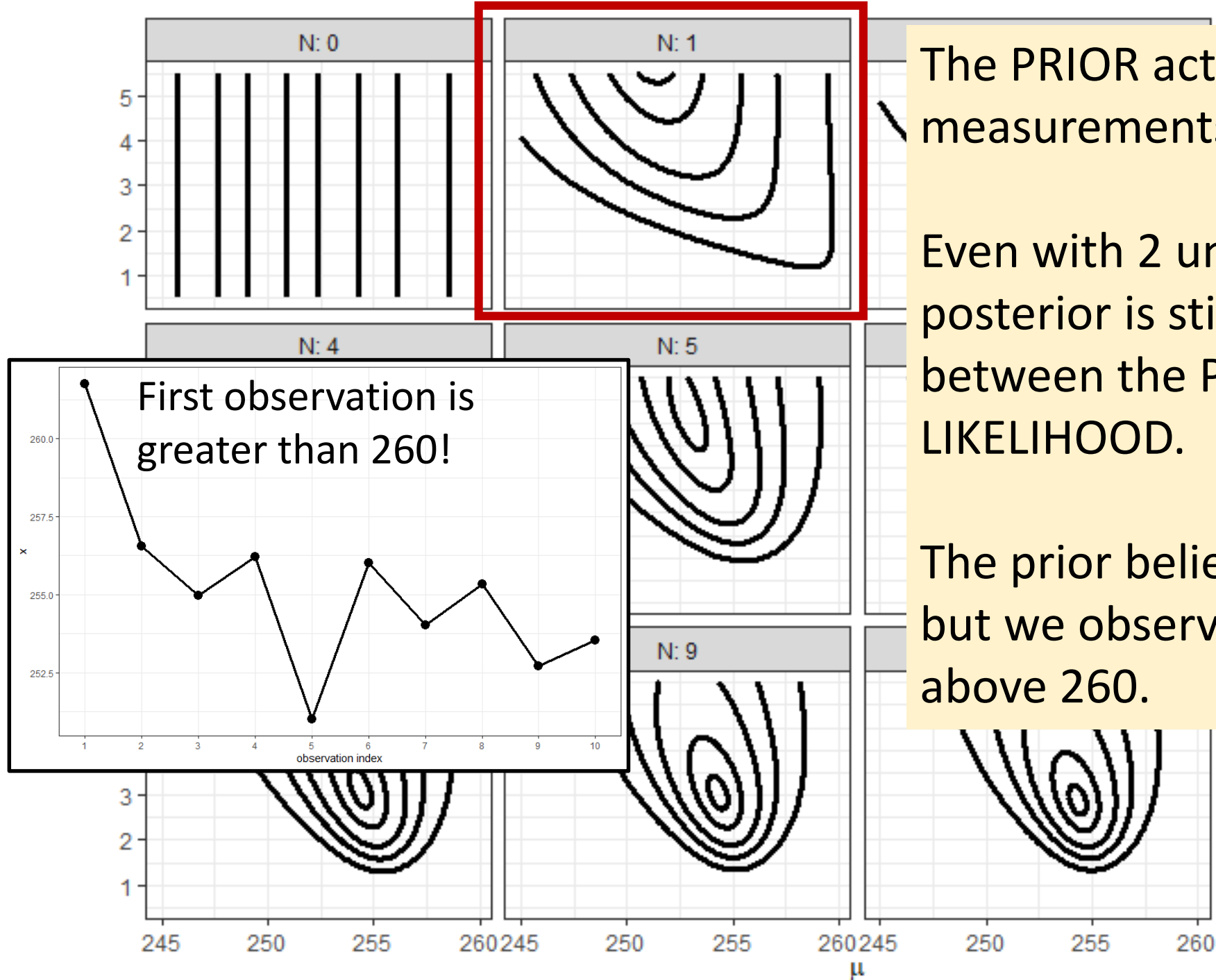


PRIOR, no data added yet hence $N=0$!



The surface we described based on 10 observations just without the fill color.



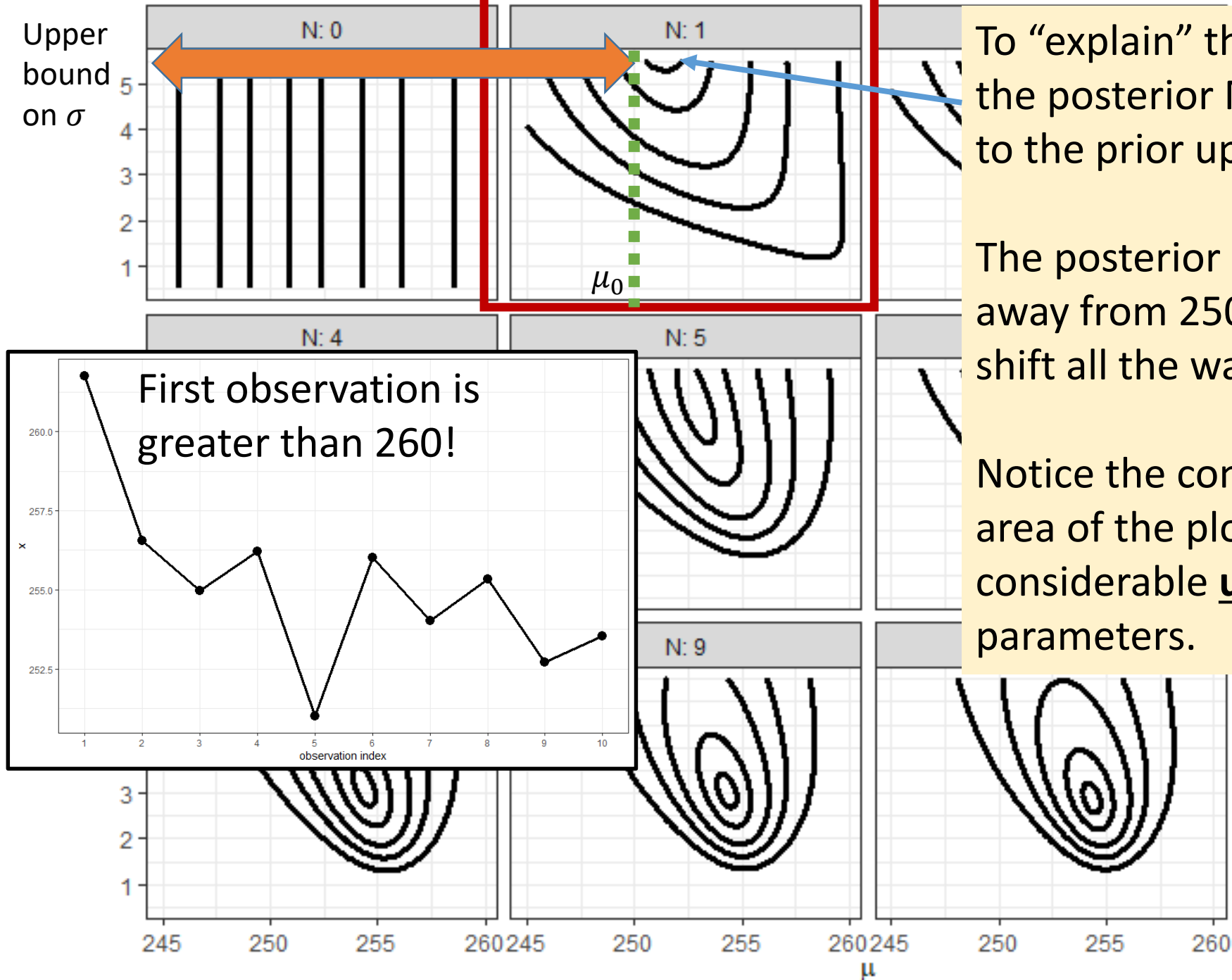


The PRIOR acts as “past measurements”!

Even with 2 unknowns, the posterior is still a **compromise** between the PRIOR and the LIKELIHOOD.

The prior believes μ is near 250, but we observed a measurement above 260.

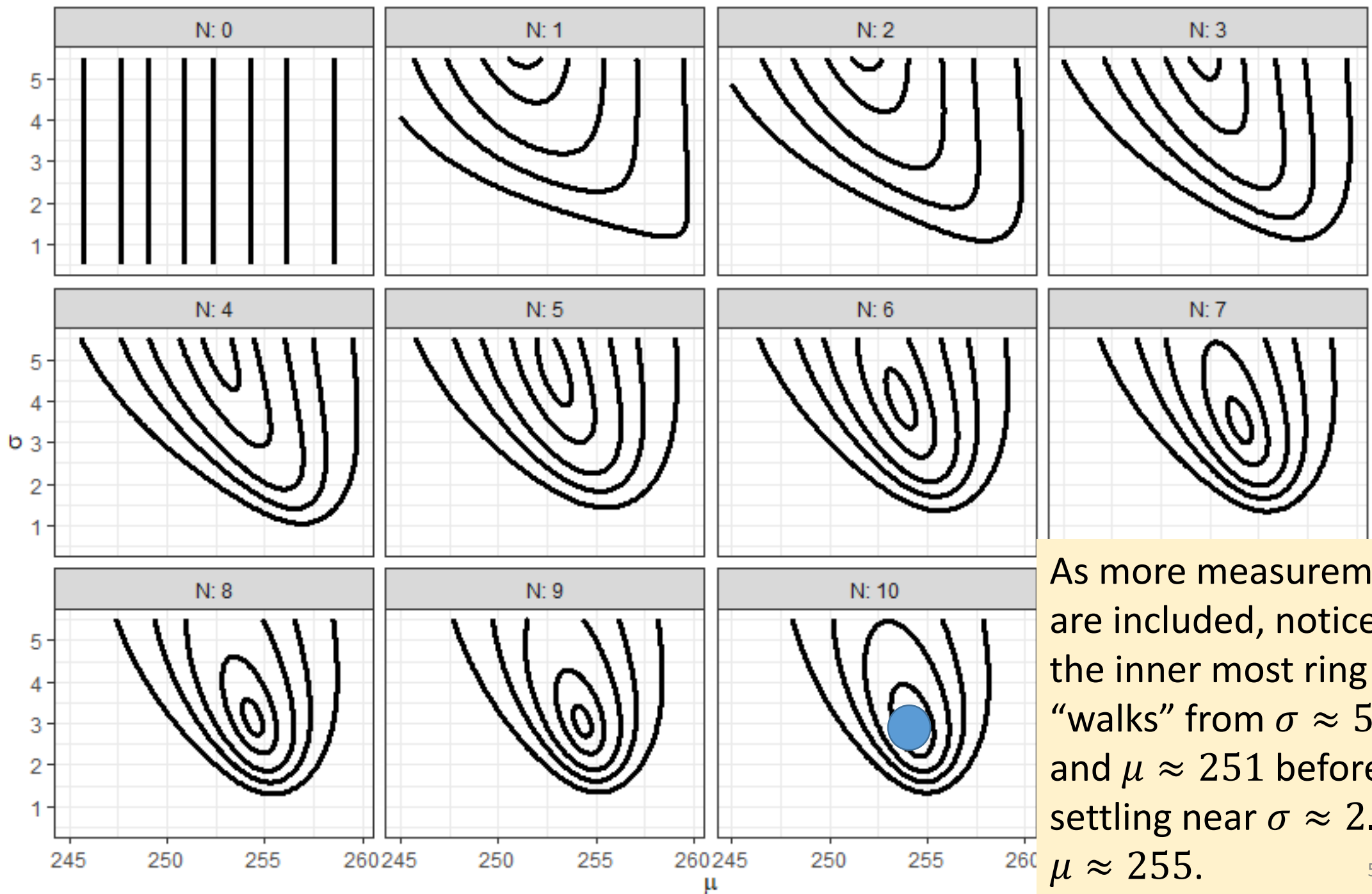
Upper bound on σ



To “explain” the first observation, the posterior MAP on σ is pushed to the prior upper bound.

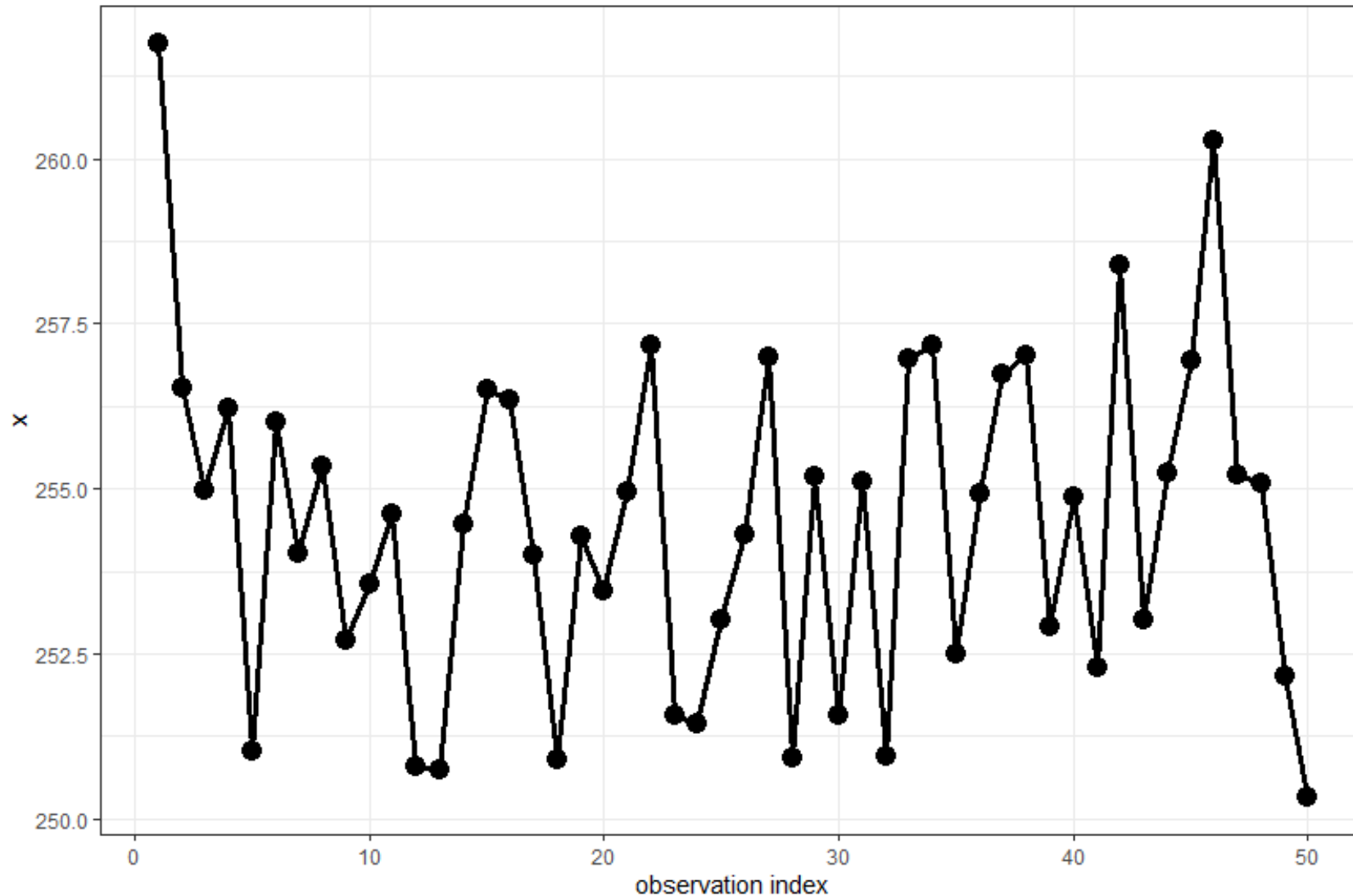
The posterior MAP on μ shifts away from 250, but it does not shift all the way up to 260.

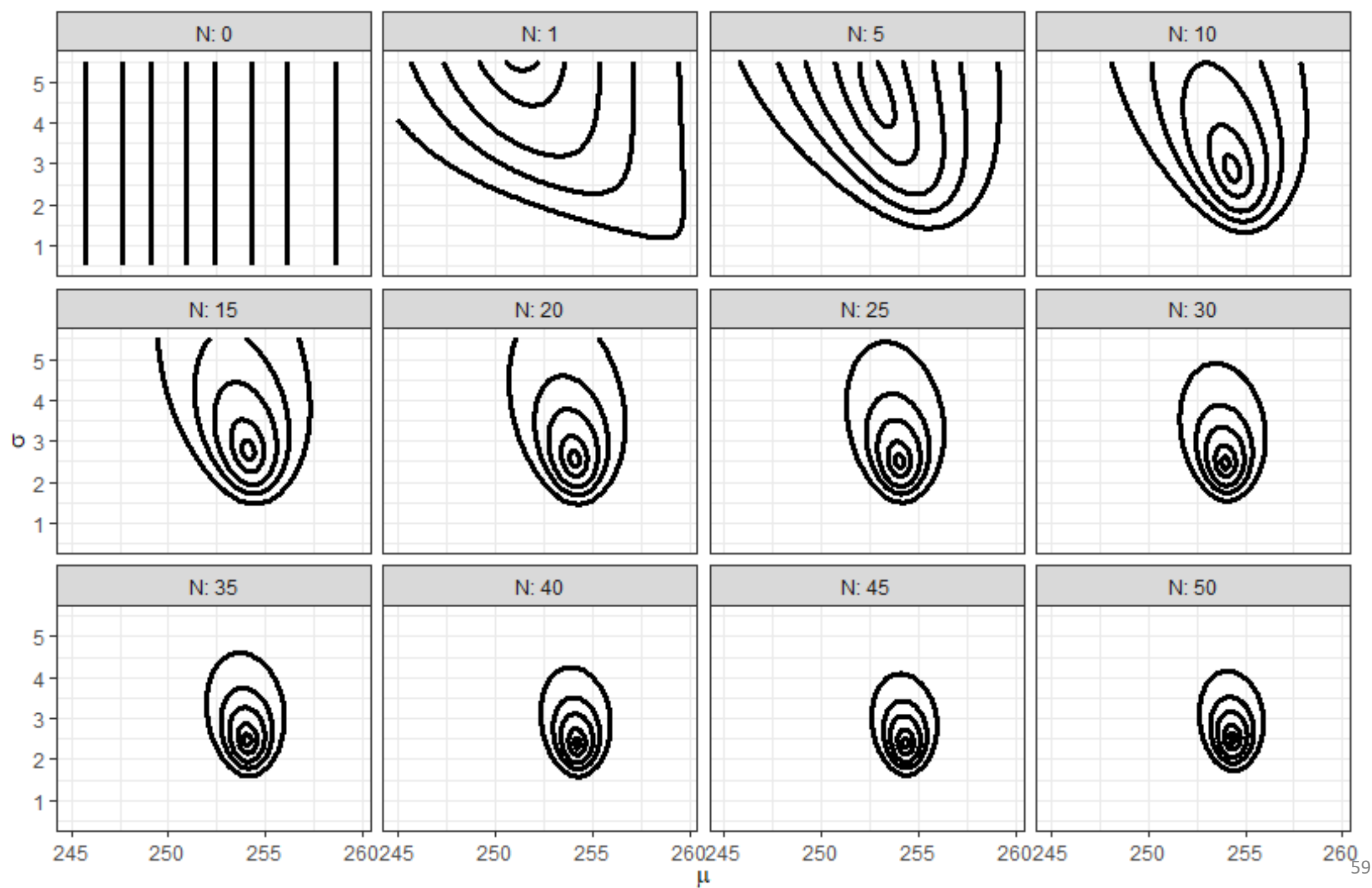
Notice the contours cover a large area of the plot, representing considerable **uncertainty** in both parameters.



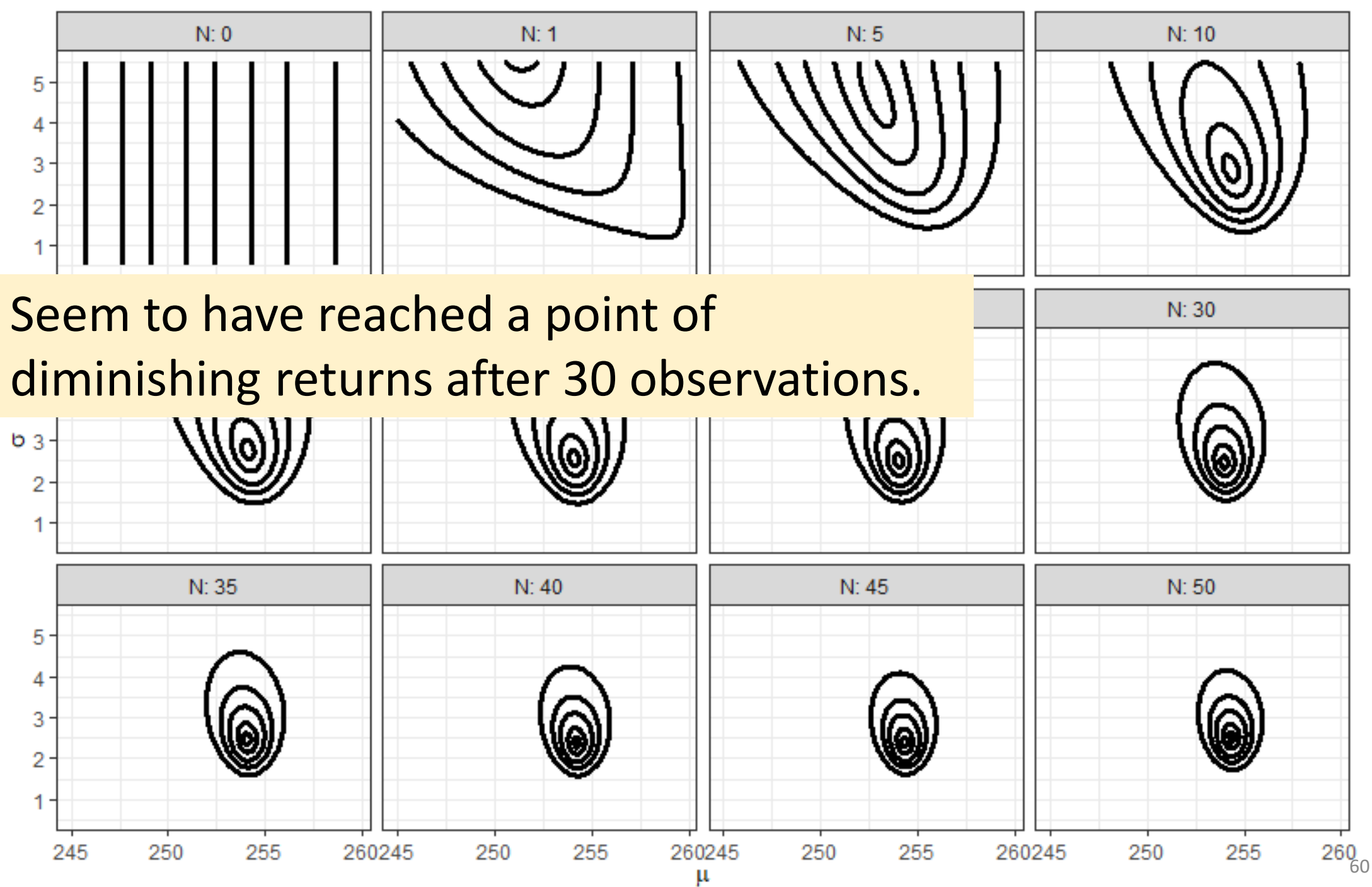
As more measurements are included, notice how the inner most ring “walks” from $\sigma \approx 5.5$ and $\mu \approx 251$ before settling near $\sigma \approx 2.5$ and $\mu \approx 255$.

Next, consider up to 50 measurements

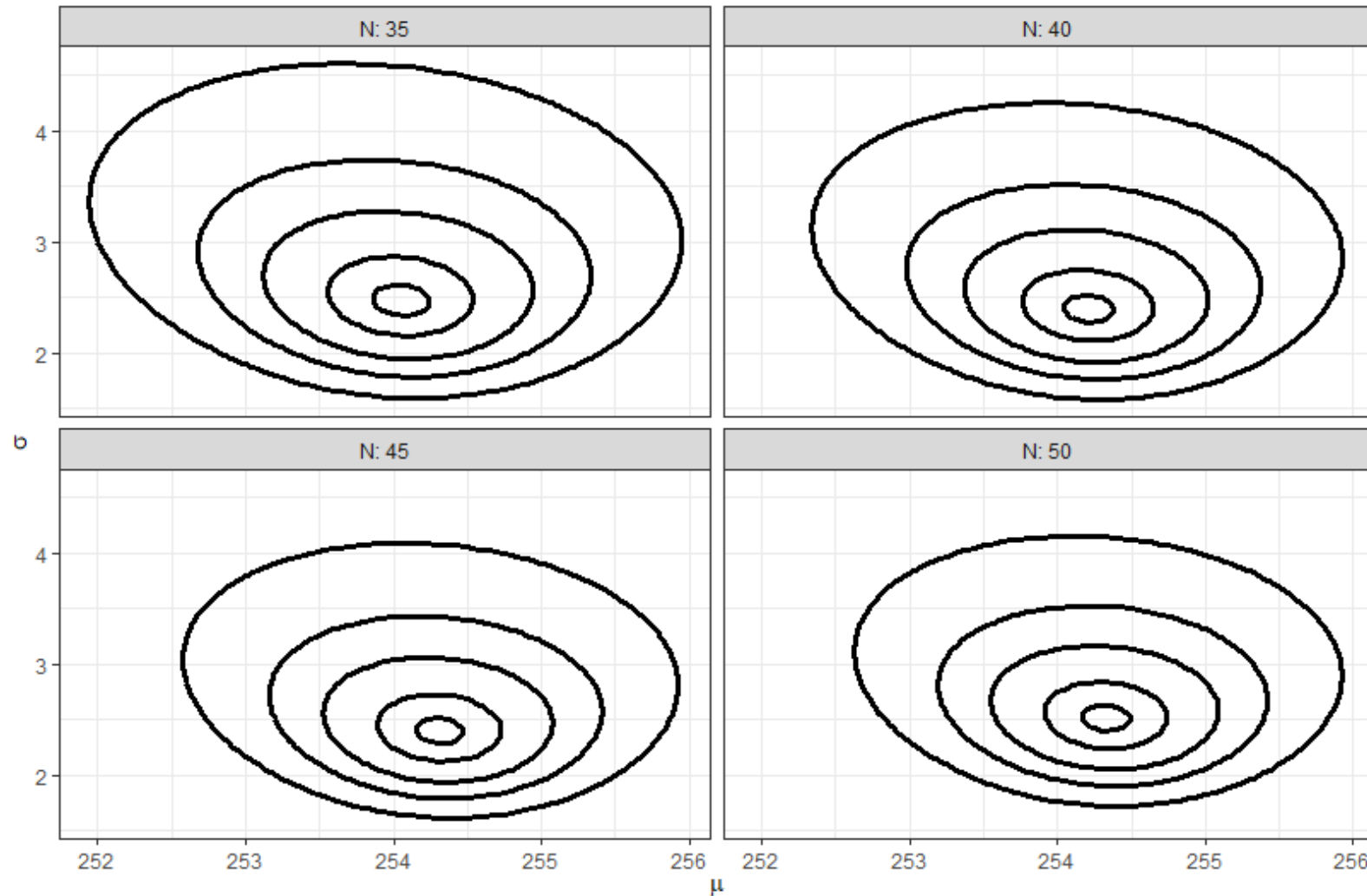




Seem to have reached a point of diminishing returns after 30 observations.



How would you describe the contour shapes?



What's happening to the surface as the sample size increases?

- The contours are starting to become...more regular...or more *normal* looking...around the posterior mode.
- What if we could approximate the posterior as a Gaussian centered around the mode?

Laplace, Quadratic, or Normal approximation

- Approximate the joint posterior distribution with a **Multivariate Normal** (MVN) centered on the **MAP**.
- Benefits:
 - Straightforward to implement.
 - Relatively fast to execute.
 - Scales to a moderate number of variables.
- Cons:
 - Let's see with an example later...

First things first...what's a MVN?

- Generalization of the Gaussian distribution to more than 1 dimension.
- Each dimension (variable) is a Gaussian and each subset of variables are MVN.

MVN density function

- There are D elements to the vector of variables:

$$\mathbf{x} = \{x_1, x_2, \dots, x_d, \dots, x_D\}$$

- **IMPORTANT:** D refers to the number of variables, NOT the number of observations!

$$p(x_1, x_2, \dots, x_D | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Break down the terms in the density

$$p(x_1, x_2, \dots, x_D, \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right\}$$

Vector of means
associated with each
element in the x-vector:
 $\mu = \{\mu_1, \mu_2, \dots, \mu_d, \dots, \mu_D\}$

Break down the terms in the density

$$p(x_1, x_2, \dots, x_D | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

$D \times D$ (variance-) covariance matrix between all elements of the \mathbf{x} -vector.

Off-diagonal elements store the covariance between the variables.

Main-diagonal elements store the variance of each element in the \mathbf{x} -vector

Break down the terms in the density

$$p(x_1, x_2, \dots, x_D | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Determinant of the
covariance matrix.

Break down the terms in the density

$$p(x_1, x_2, \dots, x_D | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Inverse of the
covariance matrix.

Break down the terms in the density

$$p(x_1, x_2, \dots, x_D | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Transpose of the
 $(\mathbf{x} - \boldsymbol{\mu})$ vector

Break down the terms in the density

$$p(x_1, x_2, \dots, x_D | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

What's this?

Break down the terms in the density

$$p(x_1, x_2, \dots, x_D | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Multidimensional generalization of
the 1-D Gaussian term:

$$\left(\frac{x - \mu}{\sigma} \right)^2$$

Break down the terms in the density

$$p(x_1, x_2, \dots, x_D | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

$\sqrt{(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}$ is a generalized distance known as the Mahalanobis distance

Bivariate Gaussian – 2D case

- $D = 2$ the vector of elements becomes: $\mathbf{x} = \{x_1, x_2\}$
- Define the correlation coefficient between the two variables as, ρ .
- The mean vector and covariance matrix are:

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$



$$\text{cov}(x_1, x_2) = \text{cov}(x_2, x_1) = \rho\sigma_1\sigma_2$$

Bivariate Gaussian – marginal distributions

- Each variable has a marginal Gaussian distribution:

$$x_1 | \mu_1, \sigma_1 \sim \text{normal}(x_1 | \mu_1, \sigma_1)$$

$$x_2 | \mu_2, \sigma_2 \sim \text{normal}(x_2 | \mu_2, \sigma_2)$$

Holds for higher dimensions!

Bivariate Gaussian – conditional distribution

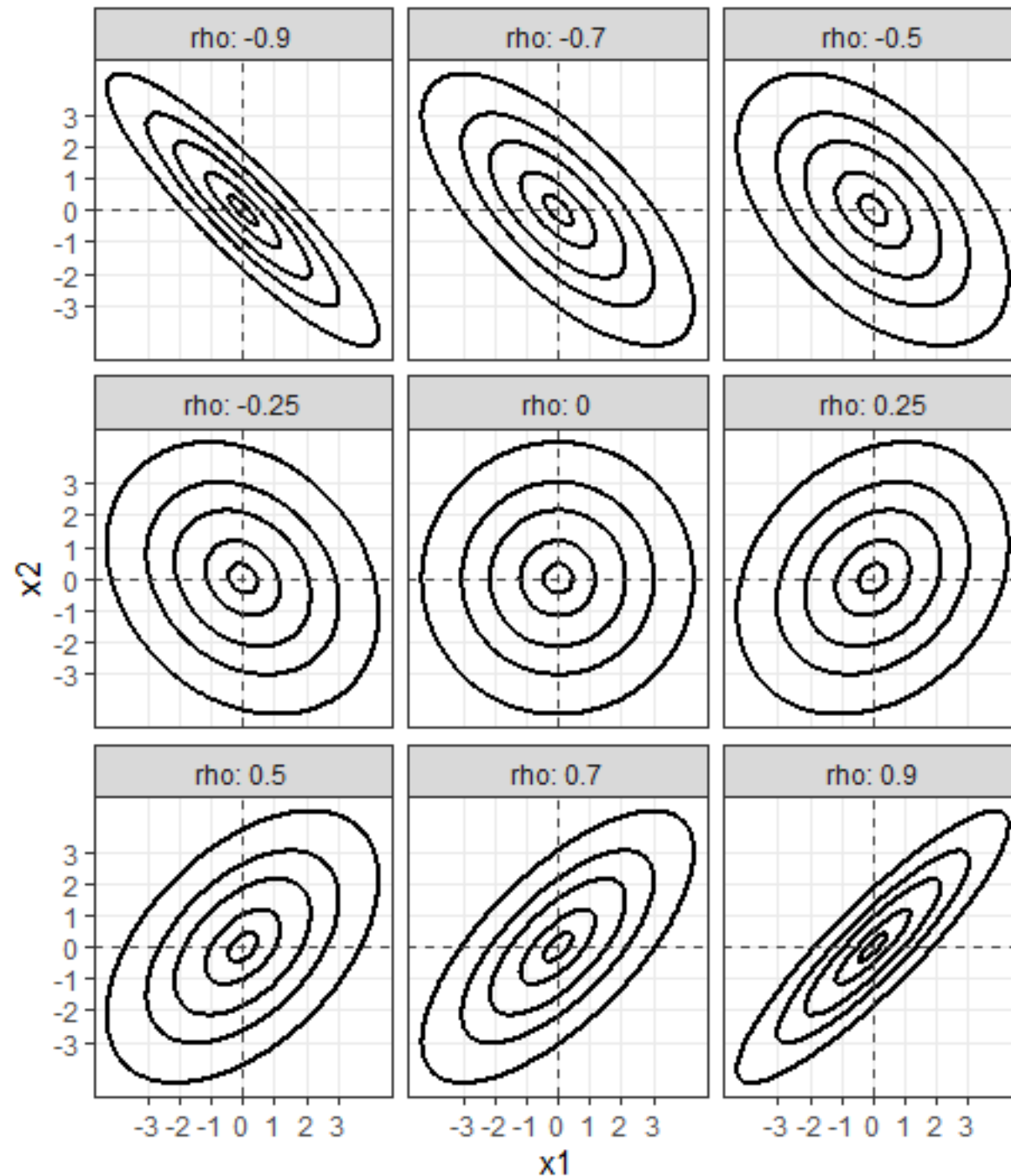
- The conditional distribution of one variable given the other...is also a Gaussian!

$$x_1 | x_2, \boldsymbol{\mu}, \boldsymbol{\Sigma} \sim \mathcal{N} \left(\mu_1 + \frac{\sigma_1}{\sigma_2} \rho (x_2 - \mu_2), (1 - \rho^2) \sigma_1^2 \right)$$

Holds for higher dimensions!

When ρ is near zero, the bivariate Gaussian density looks like circles.

As $|\rho|$ increases away from zero the circles look more and more like ellipses.

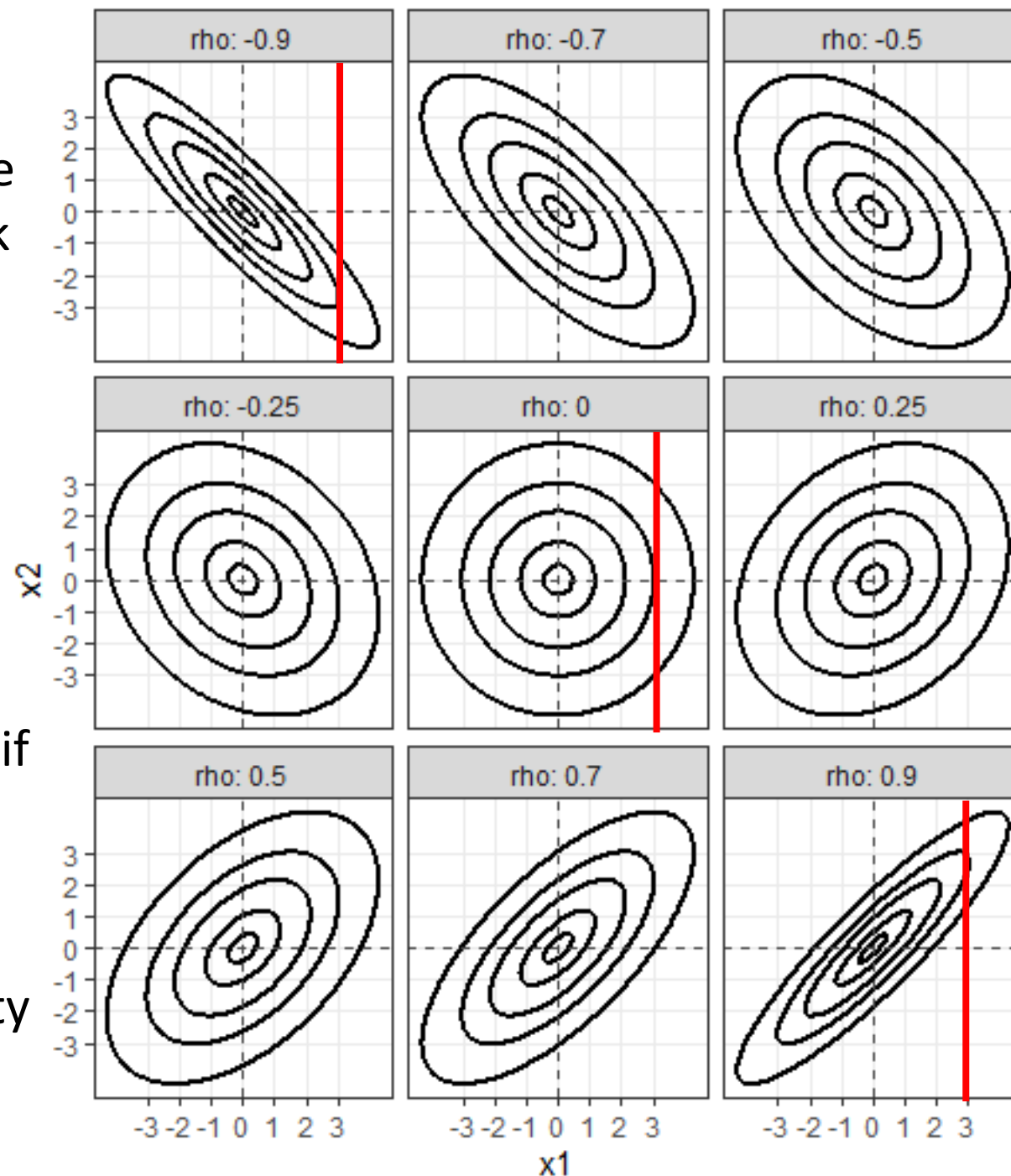


As $|\rho|$ increases away from 0, knowing information about one variable changes what we think about the other!

Consider if $x_1 = 3$, the conditional mean of $x_2 \mid x_1$ depends on the correlation coefficient!

Compare the uncertainty in x_2 if $x_1 = 3$ as ρ changes.

When $|\rho| > 0$ specifying one variable reduces the uncertainty in the value of the other.



You might be asking...what's the point of the prior...

- In this toy example, we know the measurement noise, σ !
- On top of that, the noise is relatively small:

$$\frac{1}{25} = 0.04$$

- In real problems...we will not know σ and it may not be small!