

Lecture 3: Multiparameter Models

Professor Alexander Franks

1/22/24

Announcements

- Today: Chapter 3
- Next: Chapter 4 (can skim) and Chapter 5

Bayesian inference in the normal model

- Assume $y_1, \dots, y_n \sim N(\mu, \sigma^2)$ with σ^2 a known constant
- Let's start with a Jeffreys' (improper prior): $p(\mu) \propto \text{const}$
- What is the posterior distribution $p(\mu \mid y_1, \dots, y_n, \sigma^2)$?

Bayesian inference in the normal model

- Assume $y_1, \dots, y_n \sim N(\mu, \sigma^2)$ with σ^2 a known constant
- The normal prior distribution is conjugate for μ in the normal sampling model
- Sampling distribution, prior distribution and posterior distribution are all normal.
- Assume the prior is $p(\mu) \sim N(\mu_0, \sigma^2/\kappa_0)$
- What are the parameters of the posterior $p(\mu | y_1, \dots, y_n, \sigma^2)$?

A conjugate prior for the normal likelihood

- The normal distribution is conjugate for the normal likelihood
 - Often called the “normal-normal model”
- $Y_i \sim N(\mu, \sigma^2)$ and $\mu \sim N(\mu_0, \sigma^2/\kappa_0)$ implies that the posterior distribution $p(\mu | y)$ is also normally distributed:

$$\mu | Y \sim N(\mu_n, \tau_n^2)$$

where $\mu_n = \frac{\frac{\kappa_0}{\sigma^2}\mu_0 + \frac{n}{\sigma^2}\bar{y}}{\frac{\kappa_0}{\sigma^2} + \frac{n}{\sigma^2}}$ and $\tau_n^2 = \frac{1}{\frac{\kappa_0}{\sigma^2} + \frac{n}{\sigma^2}}$

The posterior mean and pseudo-counts

$$\begin{aligned}\mu_n &= \frac{\frac{1}{\tau^2}}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}} \mu_0 + \frac{\frac{n}{\sigma^2}}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}} \bar{y} \\ &= (1 - w)\mu_0 + w\bar{y}\end{aligned}$$

$$\text{where } w = \frac{\frac{n}{\sigma^2}}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}}$$

Known mean, unknown variance

- Assume we have n mean-zero normal observations with variance σ^2
- Define $d_i = (y_i - \mu)$ for notational convenience
- What is $p(\sigma^2 \mid \mu, d_1, \dots, d_n)$?

$$L(\sigma^2) \propto (\sigma^2)^{-n/2} \exp \left\{ - \sum_{i=1}^n d_i^2 / (2\sigma^2) \right\}, \sigma^2 > 0$$

Known mean, unknown variance

- Assume we have n mean-zero normal observations with variance σ^2

- Jeffreys prior $p(\sigma^2) \propto \frac{1}{\sigma^2}$ (careful!)

- The posterior:

$$p(\sigma^2 | y) \propto (\sigma^2)^{-n/2-2} \exp \left\{ - \sum_{i=1}^n d_i^2 / (2\sigma^2) \right\}$$

- This distribution called an *inverse-Gamma* distribution.

- If $X \sim \text{Gamma}(a, b)$ then $\frac{1}{X} \sim \text{Inv-Gamma}(a, b)$.

A conjugate prior distribution

- In general, the conjugate prior distribution has the form:

$$p(\sigma^2) \propto (\sigma^2)^{-k_1} e^{\frac{k_2}{\sigma^2}}$$

- Inverse-gamma: $p(y \mid a, b) = \frac{b^a}{\Gamma(a)} y^{-a-1} \exp\left(-\frac{b}{y}\right)$
 - The mean of an inverse-Gamma is $b/(a-1)$

Joint inference for the mean and variance

In the normal model we typically factorize the prior distribution $p(\mu, \sigma) = p(\mu | \sigma)p(\sigma)$.

Specifically:

$$\sigma^2 \sim \text{Inv-Gamma}(\nu_0/2, \nu_0/2\sigma_0^2)$$

$$\mu | \sigma^2 \sim \text{normal}(\mu_0, \sigma^2/\kappa_0)$$

$$Y_1, \dots, Y_n | \mu, \sigma \sim \text{i.i.d. normal}(\mu, \sigma^2)$$

- ν_0 is interpreted as the prior sample size
- σ_0^2 is a prior sample variance

Joint inference for the mean and variance

- We write down $p(\mu, \sigma \mid y_1, \dots, y_n) \propto L(\mu, \sigma)p(\mu, \sigma)$. Now what?
- Estimands of potential interest:
 - $E[\mu \mid y_1, \dots, y_n]$
 - $E[\sigma \mid y_1, \dots, y_n]$
 - $E[\sigma/\mu \mid y_1, \dots, y_n]$ (coefficient of variation)

Posterior inference for arbitrary functions

Method of transformations

1. Find the inverse, $\theta = g^{-1}(\gamma) = \frac{e^\gamma}{1+e^\gamma}$

2. Compute $\frac{dg^{-1}(\gamma)}{d\gamma}$

3. Find

$$p_\gamma(\gamma \mid y_1, \dots, y_n) = \left| \frac{dg^{-1}(\gamma)}{d\gamma} \right| \times p_\theta(g^{-1}(\gamma) \mid y_1, \dots, y_n)$$

Posterior inference for arbitrary functions

For any $\gamma = g(\theta)$ we have

$$E(g(\theta)|y) = \int g(\theta)p(\theta|y)d\theta$$

Examples:

- $E[\gamma | y] = \int \log\left(\frac{\theta}{1-\theta}\right)p(\theta | y)d\theta$
- $Pr(\theta \in R | y) = E(I[\theta \in R]|y)$

$$\int g(\theta)p_\theta(\theta | y_1, \dots, y_n)d\theta = \int \gamma p_\gamma(\gamma | y_1, \dots, y_n)d\gamma$$

Law of the Unconscious Statstician



Monte Carlo Method

- $\bar{\theta} = \sum_{s=1}^S \theta^{(s)} / S \rightarrow \text{E}[\theta | y_1, \dots, y_n]$
- $\sum_{s=1}^S (\theta^{(s)} - \bar{\theta})^2 / (S - 1) \rightarrow \text{Var}[\theta | y_1, \dots, y_n]$
- $\#(\theta^{(s)} \leq c) / S \rightarrow \Pr(\theta \leq c | y_1, \dots, y_n)$
- the α -percentile of $\{\theta^{(1)}, \dots, \theta^{(S)}\} \rightarrow \theta_\alpha$

Almost sure convergence due to the SLLN

Monte Carlo Error

- If posterior samples are independent then:

$$p(\overline{g(\theta)}) \rightarrow N\left(g(\theta), \frac{\text{Var}(g(\theta) \mid y_1, \dots, y_n)}{S}\right)$$

-
- Just because we can write down the (proportional) posterior, doesn't mean we can sample from the distribution!
 - With exponential families and conjugate priors, sampling algorithms are available
 - For more “complicated” distributions we need new tools (e.g. MCMC)

**Back to the normal
model**

Example: midge wing length

- Modeling wing length of different species of midge (small, two-winged flies)
- From prior studies: mean wing length close to 1.9mm.
- Prior mean for μ is $\mu_0 = 1.9$ and Jeffreys prior for σ^2 .
- Prior sample sizes: choose $\kappa_0 = 1$
- $(\bar{y}, s^2) = (1.804, 0.0169)$ are the sufficient statistics

Working with the log posterior

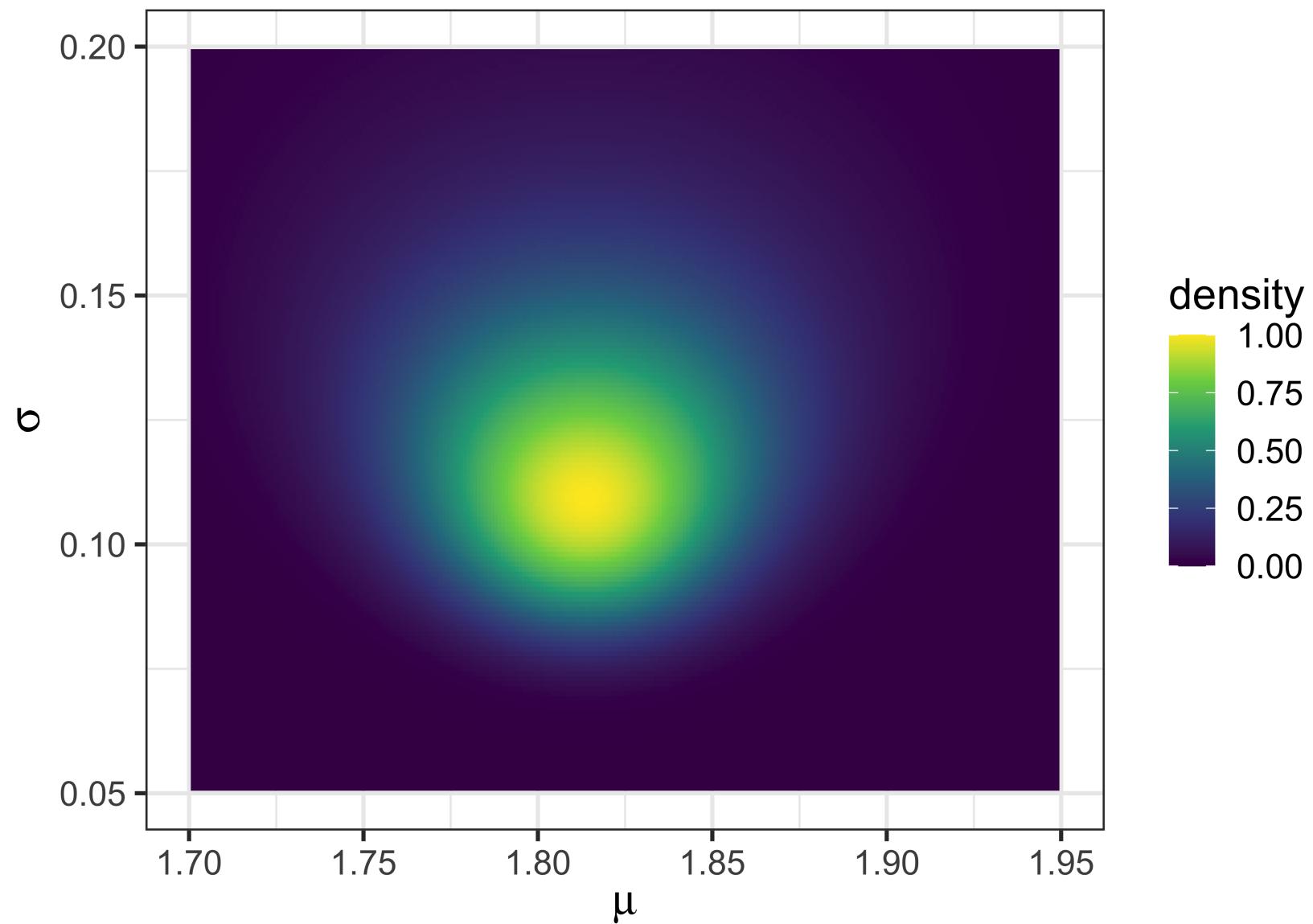
- As always, we will write down $p(\theta \mid y) \propto p(y \mid \theta)p(\theta)$
- In code, we always work with the **log-posterior** for numerical reasons
 - Mathematically it makes no difference, but computationally it is important
 - $L(\theta) \propto \prod p(y_i \mid \theta)$ is very small for moderate sample size (underflow)
 - $\ell(\theta) = \sum \log(p(y_i \mid \theta))$ is numerically stable
- Monte Carlo methods only require that we can evaluate the log posterior

Grid approximation

```
1 log_normal_posterior <- Vectorize(function(mu, sigma) {  
2  
3   ### log likelihood  
4   sum(dnorm(y, mu, sigma, log=TRUE)) +  
5   ## plus log prior  
6   dnorm(mu, mu0, sigma/sqrt(k0), log=TRUE) +  
7   -2*log(sigma)  
8  
9 })  
10  
11  
12 post_grid <- as_tibble(  
13   expand.grid(seq(1.6, 2.0, by=0.001),  
14               seq(0.01, 0.25, by=0.001)))  
15 colnames(post_grid) <- c("mu", "s")
```

Grid approximation

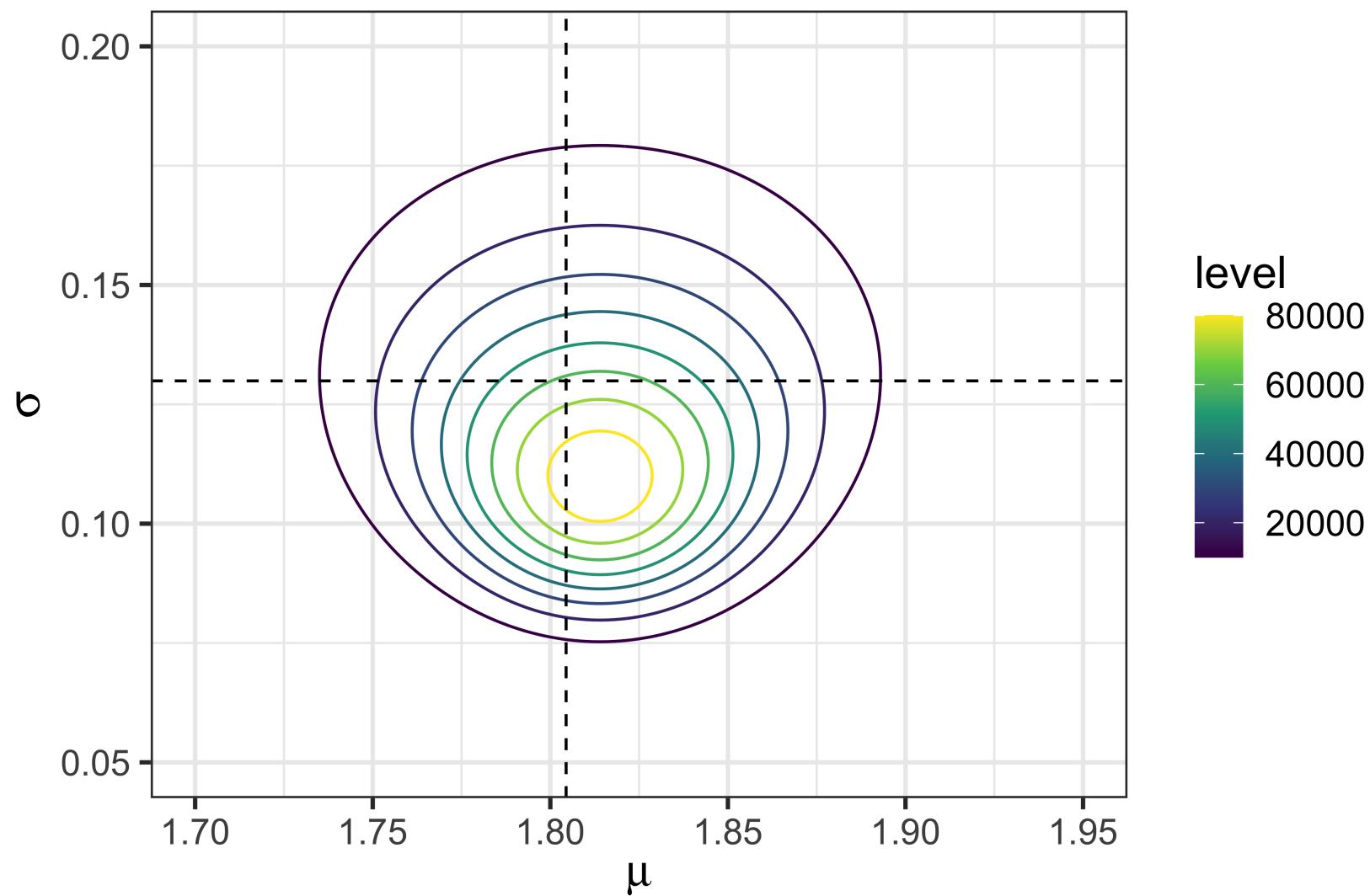
```
1 post_grid %>%
2   mutate(log_density = log_normal_posterior(mu, s)) %>%
3   mutate(density = exp(log_density - max(log_density)))) %>%
4   ggplot() +
5   geom_raster(aes(mu, s, fill=density)) +
6   xlim(c(1.7, 1.95)) + ylim(c(0.05, 0.2)) +
7   xlab(expression(mu)) +
8   ylab(expression(sigma)) +
9   theme_bw(base_size=16
10      ) +
11   scale_fill_continuous(type="viridis")
```



Contour Plot (Standard Deviation)

```
1 post_grid %>%
2   mutate(density = exp(log_normal_posterior(mu, s))) %>%
3   ggplot() +
4   geom_contour(aes(mu, s, z=density, colour=stat(level))) +
5   xlim(c(1.7, 1.95)) + ylim(c(0.05, 0.2)) +
6   xlab(expression(mu)) + ylab(expression(sigma)) +
7   ggtitle("Posterior Contours") +
8   theme_bw(base_size=16) +
9   scale_color_continuous(type="viridis") +
10  geom_hline(yintercept=s, linetype="dashed") +
11  geom_vline(xintercept=ybar, linetype="dashed")
```

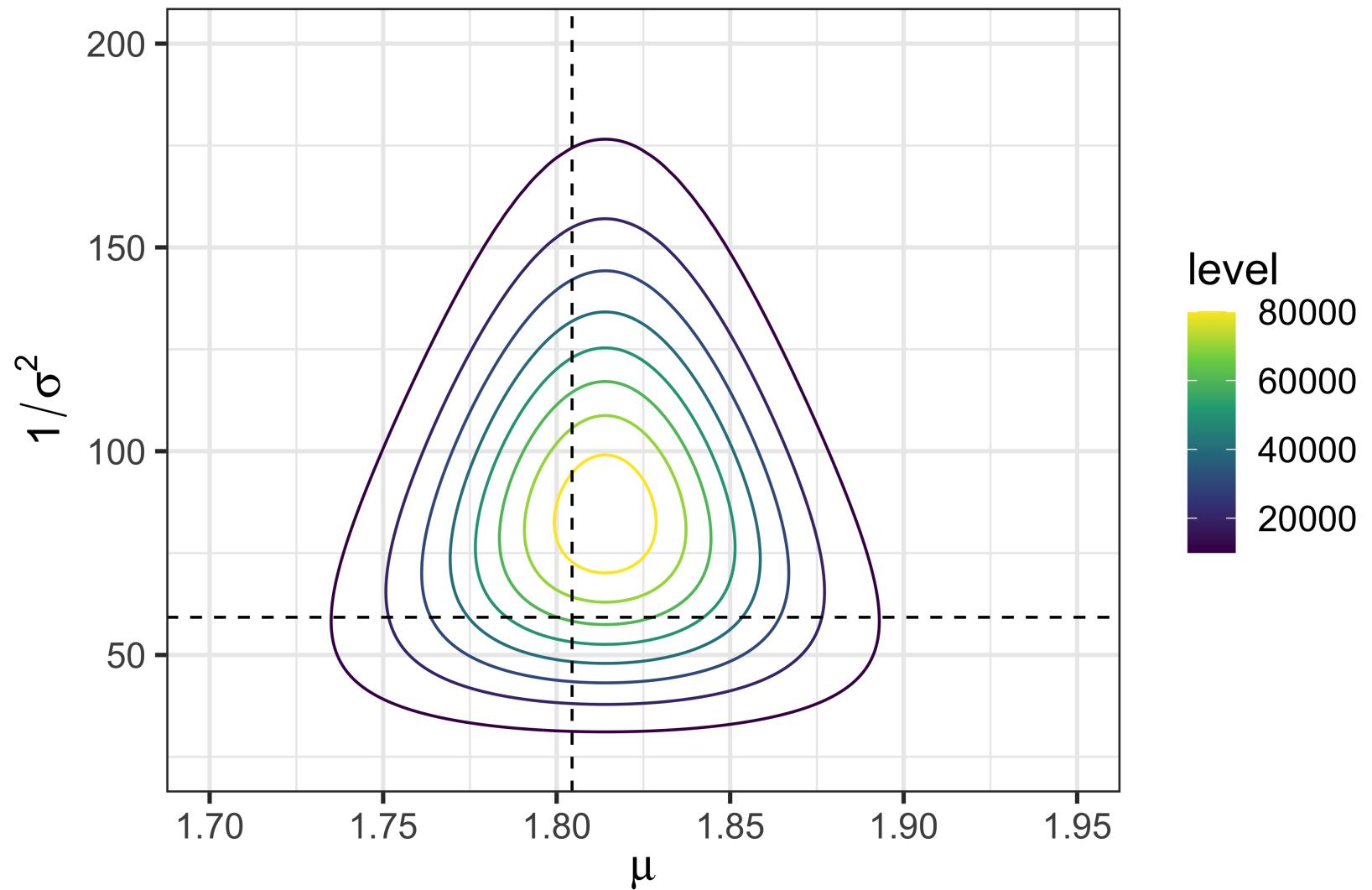
Posterior Contours



Contour Plot (Precision)

```
1 post_grid$prec <- 1/post_grid$s^2
2
3 post_grid %>%
4   mutate(density = exp(log_normal_posterior(mu, sqrt(1/prec)))) %>%
5   ggplot() +
6   geom_contour(aes(mu, prec, z=density, colour=stat(level))) +
7   xlim(c(1.7, 1.95)) + ylim(c(25, 200)) +
8   xlab(expression(mu)) + ylab(expression(1/sigma^2)) +
9   ggtitle("Posterior Contours") +
10  theme_bw(base_size=16) +
11  scale_color_continuous(type="viridis") +
12  geom_hline(yintercept=1/s^2, linetype="dashed") +
13  geom_vline(xintercept=ybar, linetype="dashed")
```

Posterior Contours

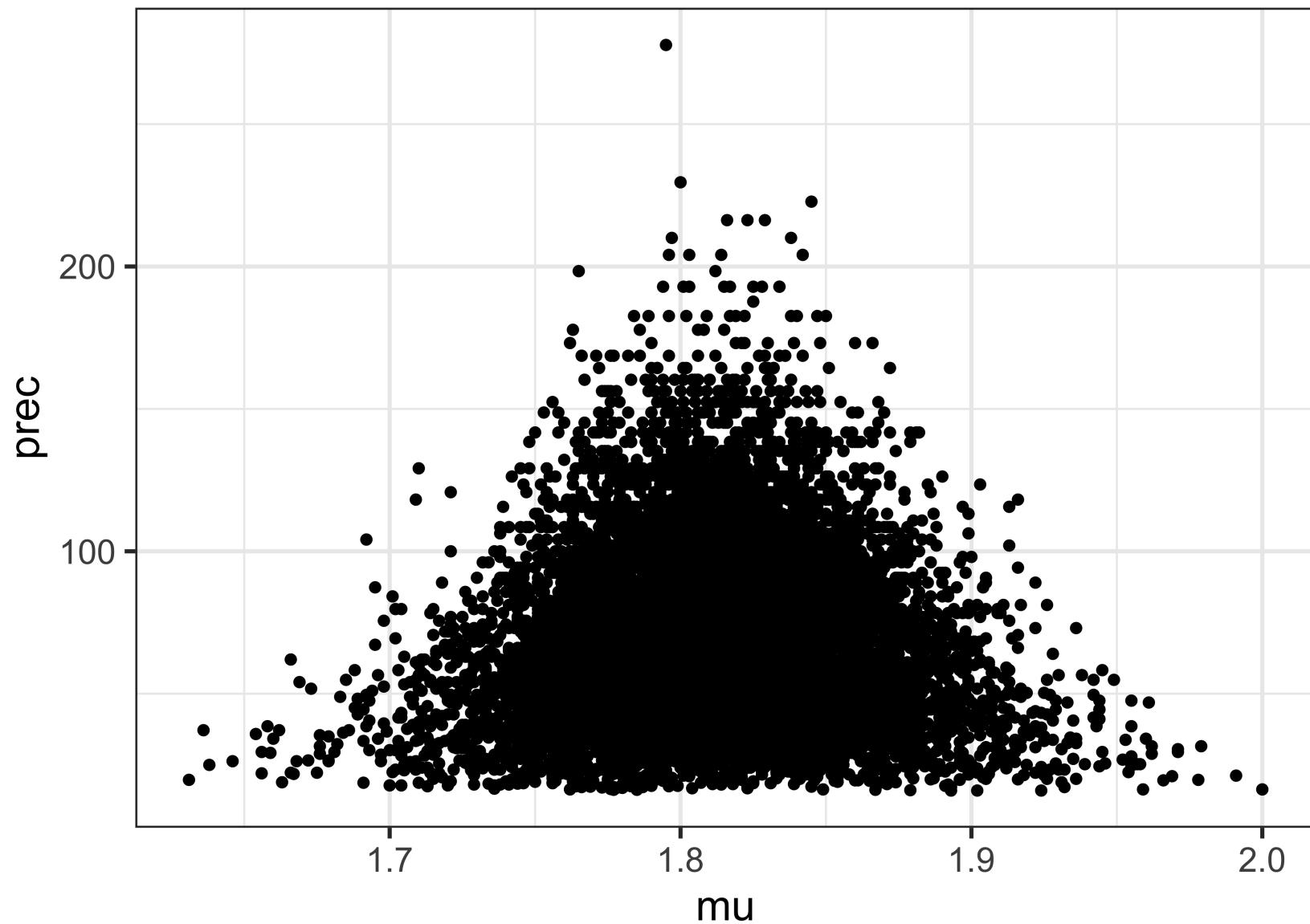


Sampling from the joint posterior

- Contour and raster plots allow us to visualize the posterior (in two dimensions)
 - Grid sampling approximates posterior samples
 - Need to know approximately where the high posterior density is (not easy)
- When we have more than 2 parameters visualization isn't feasible
- How do we summarize the posterior?
 - e.g. posterior means, posterior probabilities, intervals, etc..

Visualizing Posterior Samples

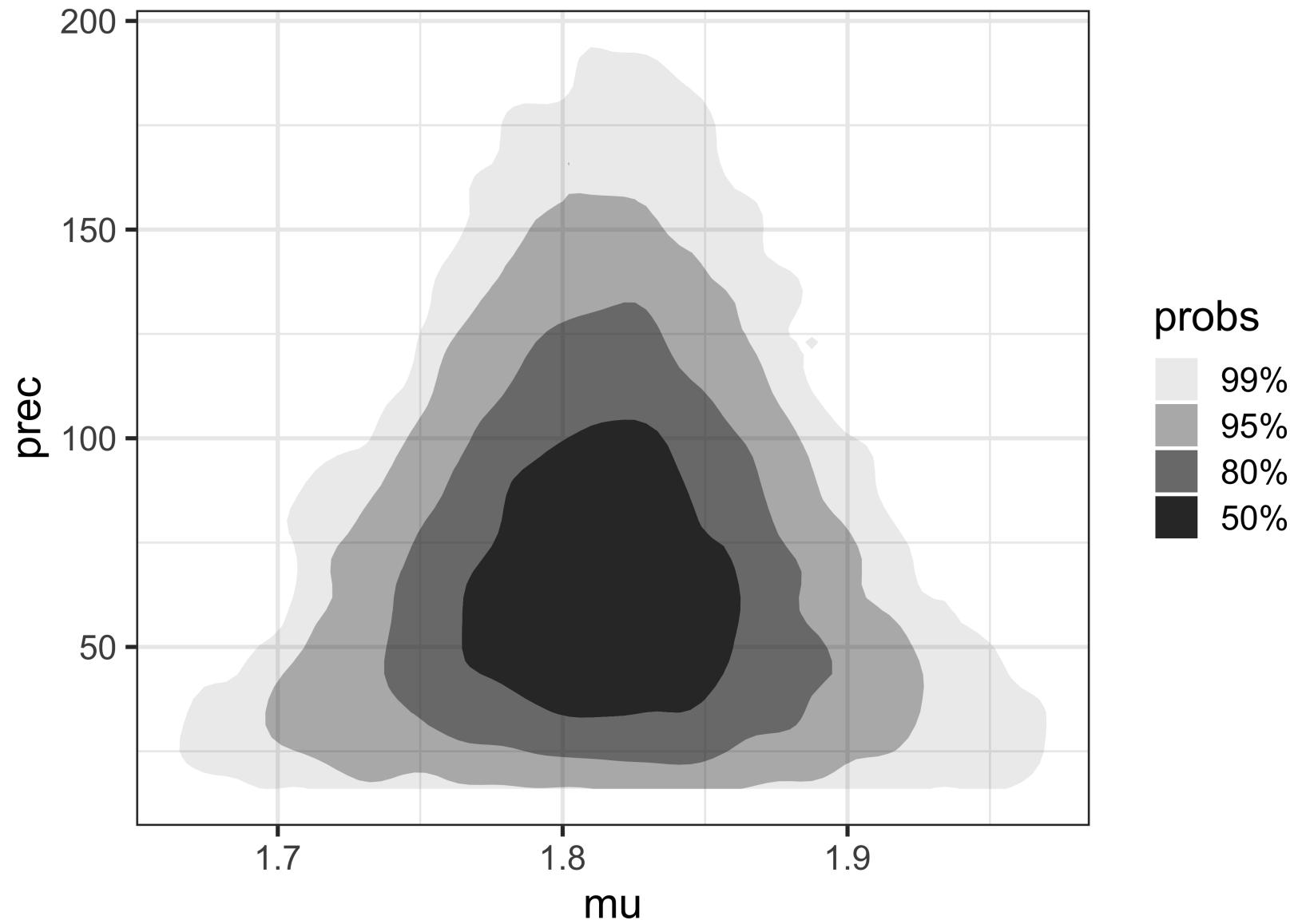
```
1 post_grid %>%
2   mutate(density = exp(log_normal_posterior(mu, sqrt(1/prec)))) %>%
3   mutate(density = density / sum(density)) -> post_grid
4
5 sample_indices <- sample(1:nrow(post_grid), size=10000, replace=TRUE, prob=
6
7 post_grid[sample_indices, ] %>%
8   mutate(prec = 1/s^2) %>%
9   ggplot() +
10  geom_point(aes(x=mu, y=prec)) +
11  theme_bw(base_size=16)
```



Visualizing Posterior Samples

```
1 post_grid %>%
2   mutate(density = exp(log_normal_posterior(mu, sqrt(1/prec)))) %>%
3   mutate(density = density / sum(density)) -> post_grid
4
5 sample_indices <- sample(1:nrow(post_grid), size=10000, replace=TRUE, prob=
6
7 post_grid[sample_indices, ] %>%
8   mutate(prec = 1/s^2) %>%
9   ggplot() +
10  ggdensity::geom_hdr(aes(x=mu, y=prec)) +
11  theme_bw(base_size=16)
```

Visualizing Posterior Samples



Direct Sampling

- $p(\mu, \sigma^2 | y) = p(\mu | \sigma^2, y)p(\sigma^2 | y)$
- Need to know how to sample from
 $p(\sigma^2 | y) = \int p(\mu, \sigma^2 | y) d\mu$ and $p(\mu | \sigma^2, y)$
- With the proposed conjugate priors, this integral is tractable

Direct Sampling

- $p(\mu, \sigma^2 | y) = p(\mu | \sigma^2, y)p(\sigma^2 | y)$
- Need to know how to sample from
 $p(\sigma^2 | y) = \int p(\mu, \sigma^2 | y) d\mu$ and $p(\mu | \sigma^2, y)$
- With the proposed conjugate priors, this integral is tractable

$\{1/\sigma^2 | y_1, \dots, y_n\} \sim \text{gamma}(\nu_n/2, \nu_n \sigma_n^2/2)$, where

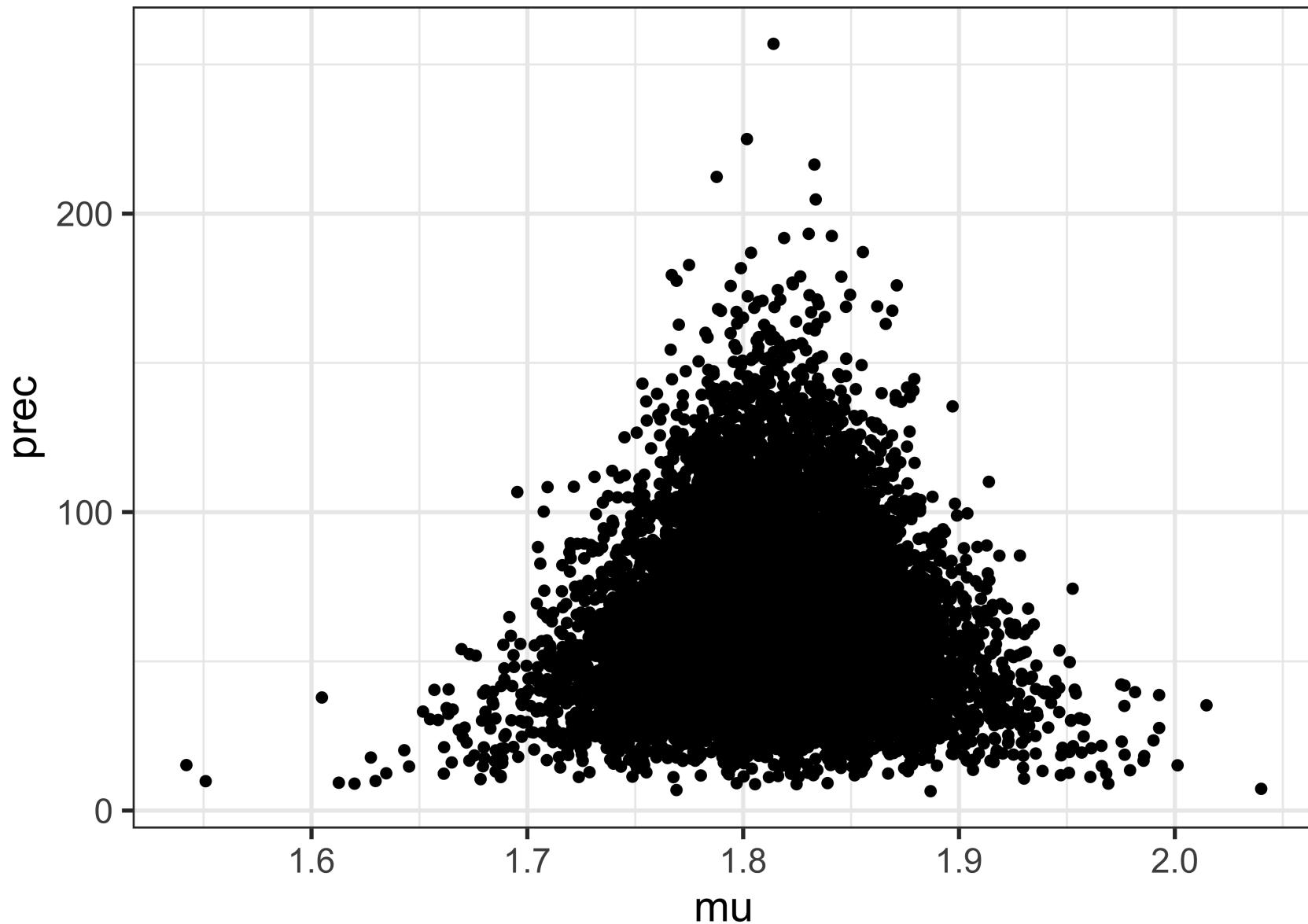
$$\nu_n = \nu_0 + n$$

$$\sigma_n^2 = \frac{1}{\nu_n} \left[\nu_0 \sigma_0^2 + (n - 1)s^2 + \frac{\kappa_0 n}{\kappa_n} (\bar{y} - \mu_0) \right]$$

Direct Sampling

```
1 ## posterior parameters
2 kn <- k0 + n
3 nun <- nu0 + n
4 mun <- (k0 * mu0 + n * ybar) / kn
5 s2n <- (nu0*s20 + (n-1)*s2 + k0*n / kn * (ybar - mu0)^2) / nun
6
7 nsamps <- 10000
8 prec_samps <- rgamma(nsamps, nun/2, nun*s2n/2)
9 mu_samps <- rnorm(nsamps, mun, sqrt((1/prec_samps) / kn))
10
11 tibble(mu=mu_samps, prec=prec_samps) %>%
12   ggplot() +
13   geom_point(aes(x=mu, y=prec)) +
14   #ggdensity::geom_hdr(aes(x=mu, y=prec)) +
15   theme_bw(base_size=16)
```

Direct Sampling



Bayes Estimators

Why the posterior mean?

- Often times we need to make a “decision” by providing a single estimate
- The posterior provides a full distribution over θ , which can be summarized in infinitely many ways
- Specify a *loss function* which describes the cost of estimating $\hat{\theta}$ when the truth is θ

Bayes Estimators

- The *loss function*: $L(\hat{\theta}, \theta)$
 - Squared error: $L(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$
 - Absolute error: $L(\hat{\theta}, \theta) = |\hat{\theta} - \theta|$
- The *Bayes risk* is the posterior expected loss: $E_{\theta|y}[L(\hat{\theta}, \theta)]$
- Summarize the posterior by minimizing the Bayes risk.
- An estimator $\hat{\theta}$ is said to be a Bayes estimator if it minimizes the Bayes risk among all estimators.

Examples Squared error loss

$$\min_{\hat{\theta}} E_{\theta|y} (\hat{\theta} - \theta)^2 = \min_{\hat{\theta}} \int (\hat{\theta} - \theta)^2 p(\theta | y) d\theta$$

The Bias-Variance Tradeoff

- The prior distribution (usually) makes your estimator biased...
- But the prior distribution also (usually) reduces the variance!
- Bayes estimators are (usually) consistent
- Example: compute the frequentist mean and variance of the posterior mean.

Example: IQ scores

- Scoring on IQ tests is designed to yield a $N(100, 15)$ distribution for the general population
- We observe IQ scores for a sample of n individuals from a particular town and estimate μ , the town-specific IQ score
- If we lacked knowledge about the town, a natural choice would be $\mu_0 = 100$
- Suppose the true parameters for this town are $\mu = 112$ and $\sigma = 13$
 - The town is smarter on average than the general population

Example: IQ scores

- What is the mean squared error of the MLE? MSE of the posterior mean?

Example: IQ scores

- What is the mean squared error of the MLE? MSE of the posterior mean?
- $\text{MSE} [\hat{\mu}_{MLE}] = \text{Var} [\hat{\mu}_{MLE}] = \frac{\sigma^2}{n} = \frac{169}{n}$

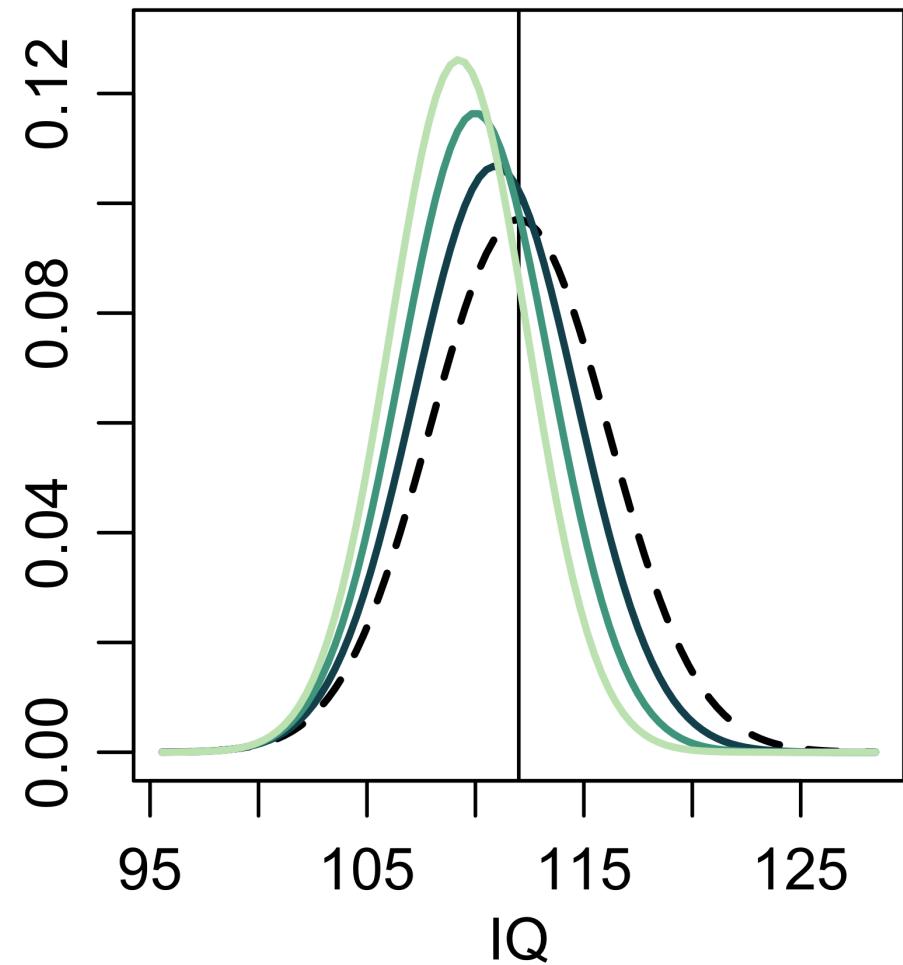
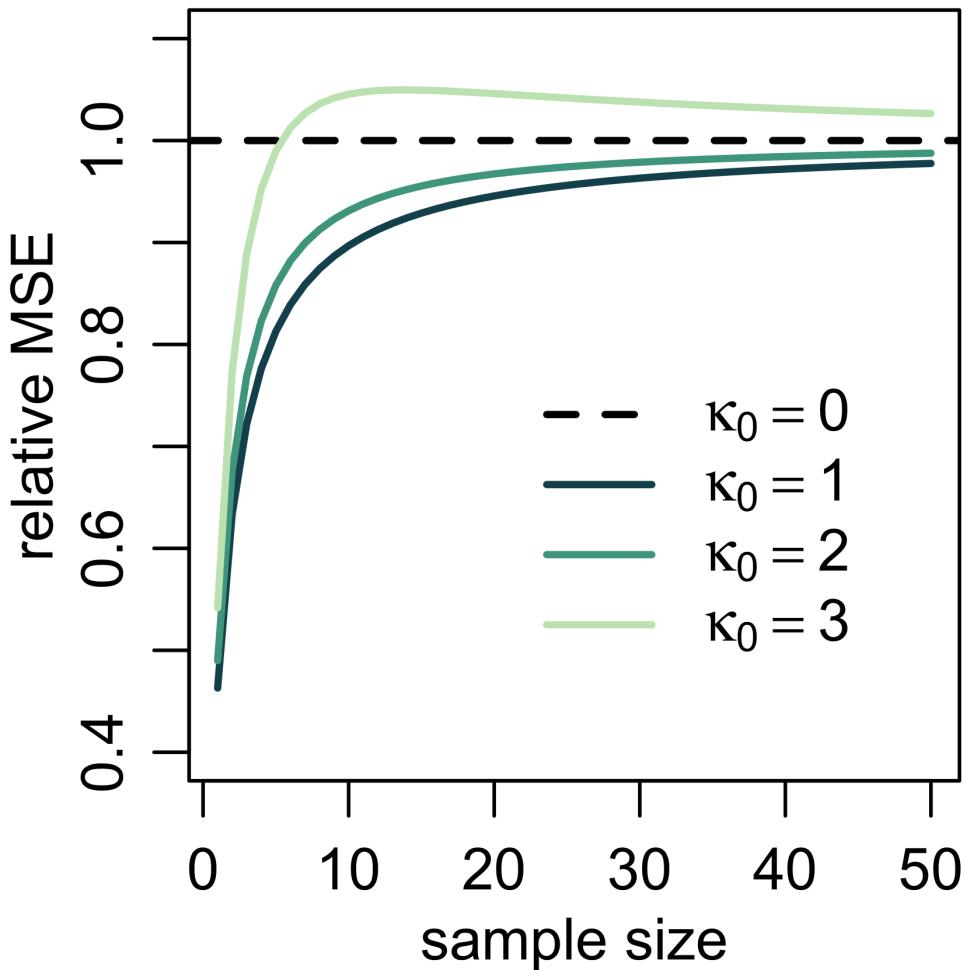
Example: IQ scores

- What is the mean squared error of the MLE? MSE of the posterior mean?
- $\text{MSE} [\hat{\mu}_{MLE}] = \text{Var} [\hat{\mu}_{MLE}] = \frac{\sigma^2}{n} = \frac{169}{n}$
- $\text{MSE} [\hat{\mu}_{PM} | \theta_0] = w^2 \frac{169}{n} + (1 - w)^2 144$

Example: IQ scores

- What is the mean squared error of the MLE? MSE of the posterior mean?
- $\text{MSE} [\hat{\mu}_{MLE}] = \text{Var} [\hat{\mu}_{MLE}] = \frac{\sigma^2}{n} = \frac{169}{n}$
- $\text{MSE} [\hat{\mu}_{PM} | \theta_0] = w^2 \frac{169}{n} + (1 - w)^2 144$
- Reminder: $w = \frac{n}{\kappa_0 + n}$. For what values of n and κ_0 is the MSE smaller for the posterior mean estimator than the maximum likelihood?

Example: IQ scores



The Multivariate Normal Distribution

Dirichlet-Multinomial

Metagenomics example

- Metagenomics is the study of genetic material recovered directly from environmental samples
- Map counts of genetic material to counts of microbial species
- Assume species are sampled with replacement
 - Observed sample is a multinomial distribution
- Total counts isn't meaningful (hard to control how much total sample)
- Relative counts are meaningful

Multinomial Density

- Let $\theta = (\theta_1, \theta_2, \dots, \theta_k)$ with $\sum_i \theta_i = 1$
- If $Y = (Y_1, \dots, Y_K) \sim Mult(n, \theta)$, then:
 - $Y_i \sim Bin(n, \theta_i)$
 - $Y_i + Y_j \sim Bin(n, \theta_i + \theta_j)$
- What is $\hat{\theta}_{MLE}$?

Dirichlet Distribution

