

# Lecture 3: Multiparameter Models

Professor Alexander Franks

1/22/24

- Gridscope → Canvas
  - Hw 2 out tonight, due <sup>next</sup> Sun.  
2/4
  - Hierarchical (Ch. 5)
-

# Bayesian inference in the normal model

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

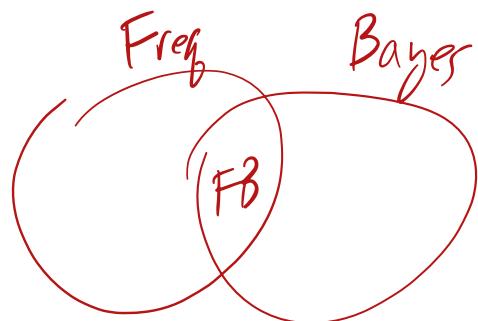
$$\begin{aligned} \sqrt{I(\mu)} &\propto \frac{1}{\sigma^2} \propto \text{const} \\ P(\mu \mid y_1, \dots, y_n) &\propto e^{-\frac{\sum (y_i - \mu)^2}{2\sigma^2}} \times \text{const} \\ &\propto e^{-\frac{n\mu^2 - 2\sum y_i \mu + \text{const}}{2\sigma^2}} \times \frac{1}{n} \end{aligned}$$

$$2 e^{-\frac{(\mu - \bar{y})^2}{2\sigma^2 n}}$$

$$P(\mu | y_1, \dots, y_n, \sigma^2) \sim N(\bar{y}, \frac{\sigma^2}{n})$$

$$\hat{\mu}_{MLE} = \bar{y} \sim N(\mu, \frac{\sigma^2}{n})$$

$$\sqrt{\frac{\sigma^2}{n}} (\mu - \bar{y}) \sim N(0, 1)$$



# Bayesian inference in the normal model

- Assume  $y_1, \dots, y_n \sim N(\mu, \sigma^2)$  with  $\sigma^2$  a known constant
- The normal prior distribution is conjugate for  $\mu$  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  $p(\mu | y_1, \dots, y_n, \sigma^2)$ ?

$\sigma^2$

conjugate

$$P(\mu | y_1, \dots, y_n) \propto L(\mu) P(\mu)$$

$$\propto e^{-\frac{(\mu - \bar{y})^2}{2\sigma^2/n}} e^{-\frac{(\mu - \mu_0)^2}{2\sigma^2/k_0}}$$

$$\propto \exp \left[ -\frac{1}{2} \frac{n+k_0}{\sigma^2} \mu^2 - 2 \left( \frac{n\bar{y}}{\sigma^2} + \frac{k_0}{\sigma^2} \mu_0 \right) \mu + \text{const} \right]$$

$$\propto \exp \left[ -\frac{1}{2} \left( \frac{n+k_0}{\sigma^2} \right) \left( \mu - \frac{\frac{n\bar{y}}{\sigma^2} + \frac{k_0}{\sigma^2} \mu_0}{\frac{(n+k_0)}{\sigma^2}} \right)^2 \right]$$

Mean:  $\frac{n/\sigma^2}{\frac{(n+k_0)}{\sigma^2}} \bar{y} + \frac{\frac{k_0}{\sigma^2} \mu_0}{\frac{n+k_0}{\sigma^2}}$

$$= \frac{n}{n+k_0} \bar{y} + \frac{k_0}{n+k_0} \mu_0$$

$$= w \hat{\mu}_{MLE} + (1-w) \mu_0$$

$$P(\mu | y_1, \dots, y_n) \sim N(w\bar{y} + (1-w)\mu_0, \frac{\sigma^2}{n+k_0})$$

# 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

$$\tau^2 = \frac{\sigma^2}{K_0}$$

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

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

Fisher Weight.

$y_i \stackrel{\text{ind}}{\sim} N(\mu, \sigma_i^2)$ , MLE for  $\mu$ ?

$$\hat{\mu}_{\text{MLE}} = \frac{\sum \frac{1}{\sigma_i^2} y_i}{\sum \frac{1}{\sigma_i^2}}$$

# Known mean, unknown variance

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

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

denom.

$$P(\sigma^2) \propto \frac{1}{\sigma^2} \quad \begin{array}{l} (* P(0) = 1/0 \\ * P(\log \sigma) \propto \text{const} \end{array}$$

7

$$P(\log \sigma) \propto \text{const}$$

$$\sigma^2 = \exp(2 \log \sigma)$$

$$P_{\sigma^2}(\sigma^2) \propto P_{\log \sigma}(\log \sigma) \left| \frac{d \log \sigma}{d \sigma^2} \right|$$

$$\frac{d \log(\sigma^2)}{d \sigma^2} \propto \frac{1}{\sigma^2}$$

# Known mean, unknown variance

- Assume we have  $n$  mean-zero normal observations with variance  $\sigma^2$

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

- The posterior:

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

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

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

$$\text{Inv-Gam}\left(\frac{n}{2}, \frac{\sum d_i^2}{2}\right)$$

Scaled Env- $\chi^2(n, \frac{\sum d_i^2}{n})$

$$\chi_k^2 = \sum_k z_i^2, \quad z_i \sim N(0, 1)$$

$$\sigma^2 \chi_k^2 \quad \chi_k^2 \equiv \text{Gam}(k/2, 1/2)$$

/

# 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:  $p(y | a, b) = \frac{b^a}{\Gamma(a)} y^{-a-1} \exp\left(-\frac{b}{y}\right)$ 
  - 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  $p(\mu, \sigma^2) = p(\mu | \sigma^2)p(\sigma^2)$ .

Specifically:

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

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

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

- $\nu_0$  is interpreted as the prior sample size
- $\sigma_0^2$  is a prior sample variance

$$P(\mu, \sigma^2 | y_1, \dots, y_n) \propto$$

$$\left[ \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mu)^2}{2\sigma^2}} \right] e^{-\frac{(\mu - \mu_0)^2}{2\sigma^2/k_0}} (\sigma^2)^{-\frac{v_0}{2}-1-\frac{v_0\sigma^2}{2\sigma_0^2}}$$

$k_0 \rightarrow 0$  implies what about  $P(\mu | \sigma^2)$

$v_0 \rightarrow 0$  implies what about  $P(\sigma^2)$

$$P(\mu | \sigma^2, y_1, \dots, y_n) P(\sigma^2 | y_1, \dots, y_n)$$

"Conditional post."

"Marginal post."

# Joint inference for the mean and variance

- We write down  $p(\mu, \sigma^2 | y_1, \dots, y_n) \propto L(\mu, \sigma^2)p(\mu, \sigma^2)$ . Now what?
- Estimands of potential interest:
  - $E[\mu | y_1, \dots, y_n]$
  - $E[\sigma | y_1, \dots, y_n]$
  - $E[\sigma/\mu | y_1, \dots, y_n]$  (coefficient of variation)  
 $\underbrace{\sigma}_{CV}$ .

# 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  $\mu$  is  $\underline{\mu_0 = 1.9}$  and Jeffreys prior for  $\sigma^2$ .
- Prior sample sizes: choose  $\kappa_0 = 1$
- $(\bar{y}, s^2) = (1.804, 0.0169)$  are the sufficient statistics

$\underline{\mu_0}$

$p(\sigma^2)$

$\underline{\sigma^2}$

# 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

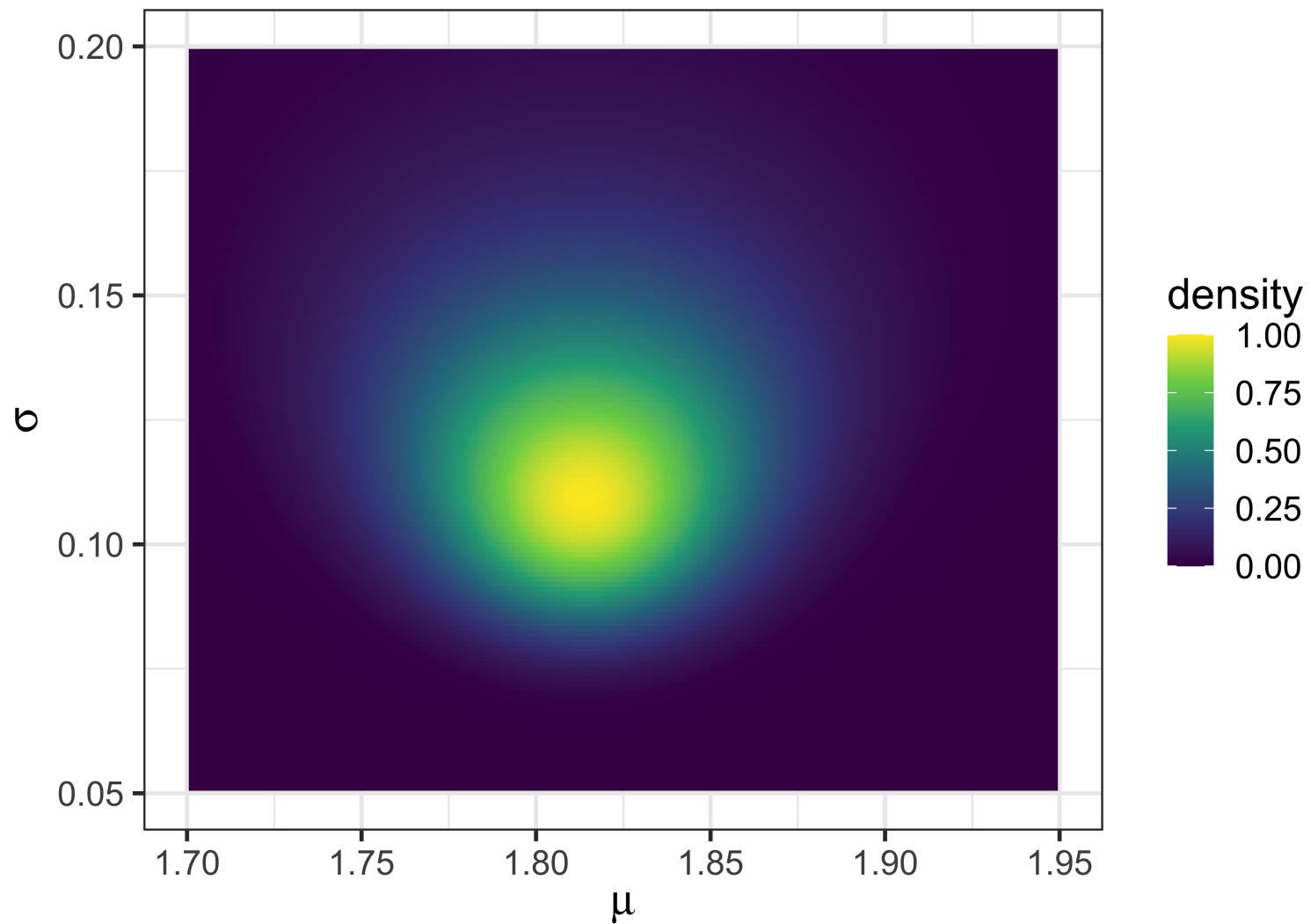
# Grid approximation

```
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")
```

$$y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

# Grid approximation

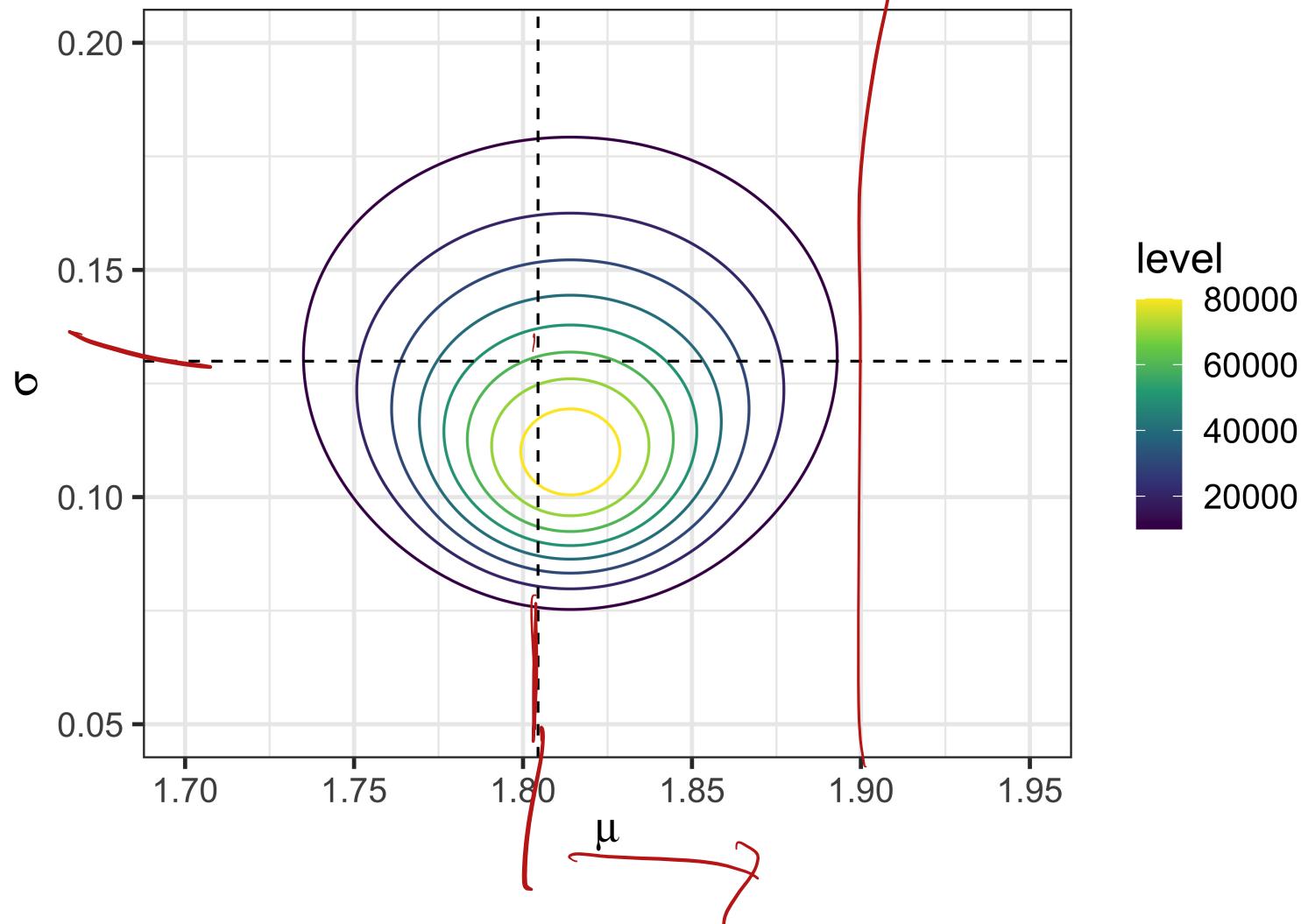
```
1 post_grid %>%
2   mutate(log_density = log_normal_posterior(mu, s)) %>%
3   mutate(density = exp(log_density - max(log_density))) %>%
4   ggplot() +
5   geom_raster(aes(mu, s, fill=density)) +
6   xlim(c(1.7, 1.95)) + ylim(c(0.05, 0.2)) +
7   xlab(expression(mu)) +
8   ylab(expression(sigma)) +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

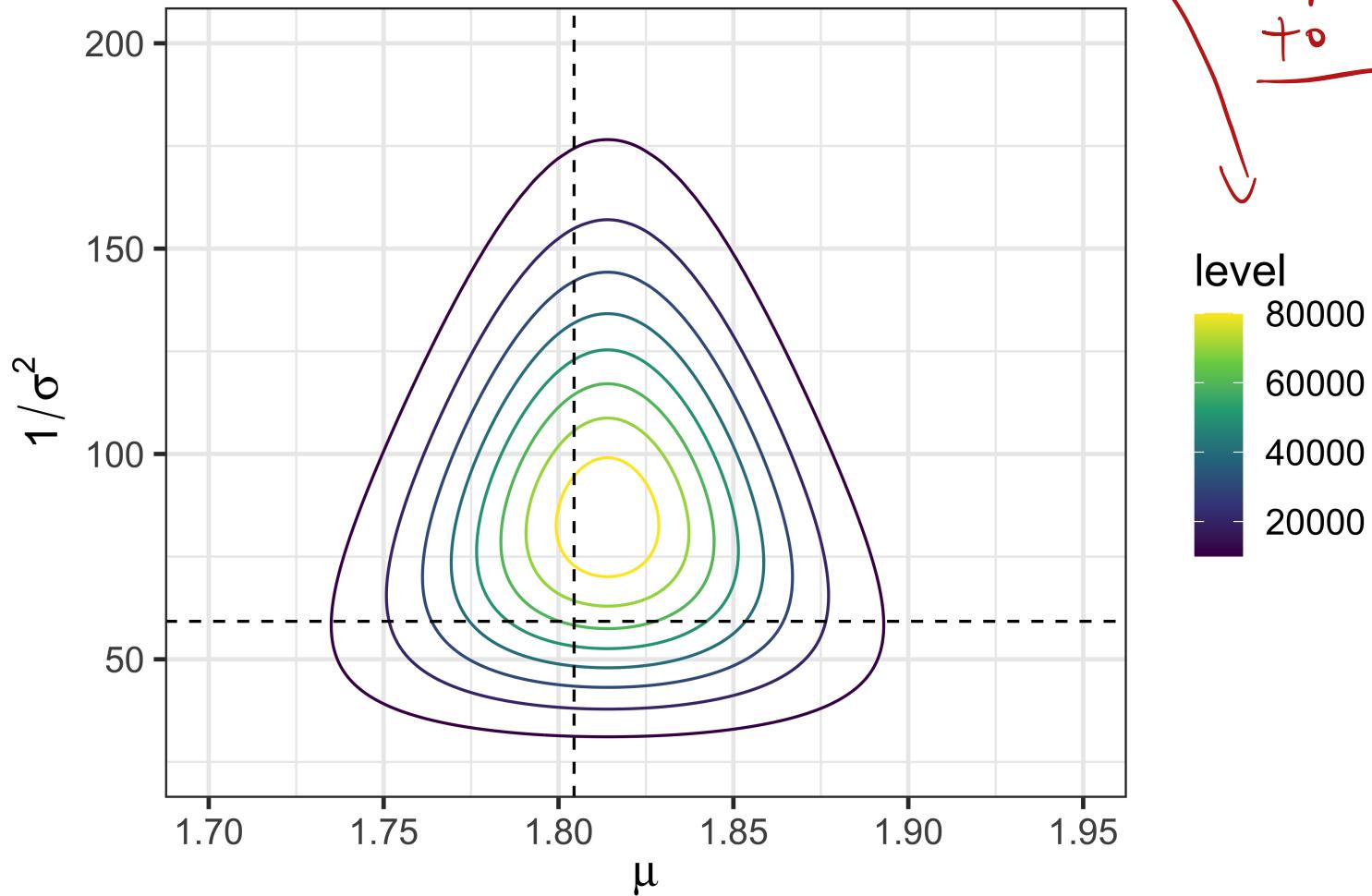


# Contour Plot (Precision)

```
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")
```

$P(\mu, 1/\sigma^2 | \text{data})$

## Posterior Contours



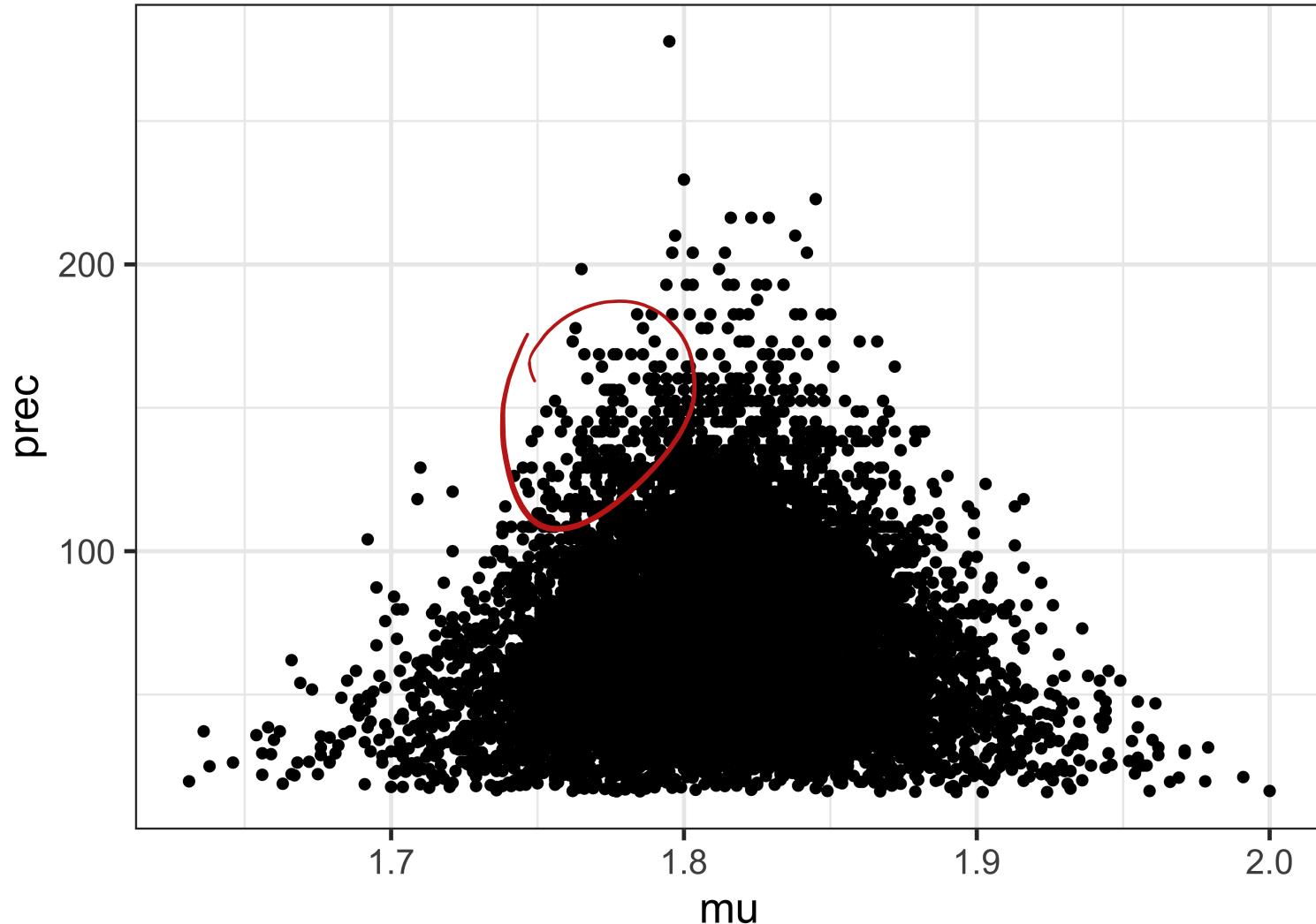
# Sampling from the joint posterior

- 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)
```

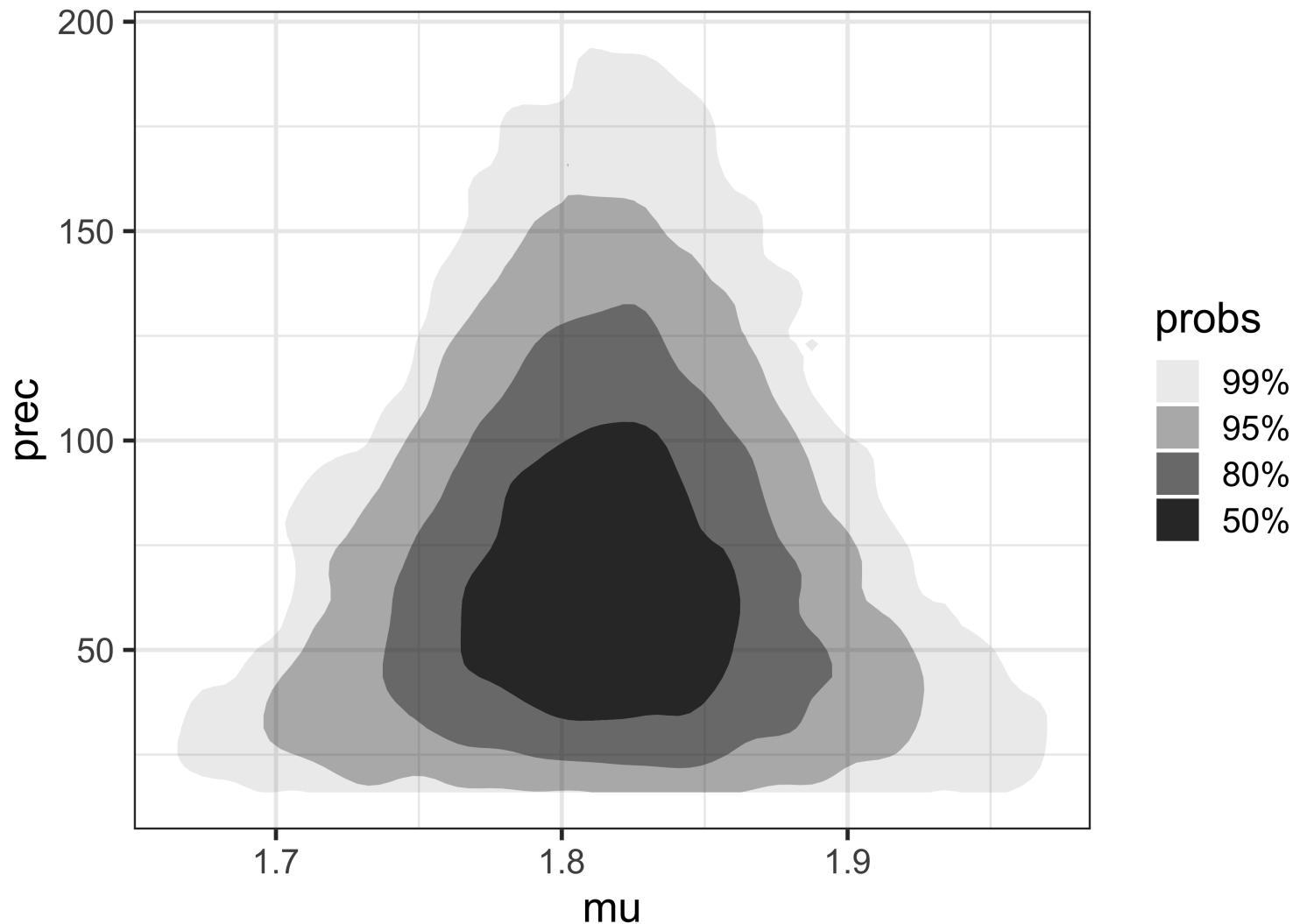
density



# 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

$$\bullet p(\mu, \sigma^2 | y) = p(\mu | \sigma^2, y)p(\sigma^2 | y)$$

- Need to know how to sample from

$$p(\sigma^2 | y) = \int p(\mu, \sigma^2 | y) d\mu \text{ and } p(\mu | \sigma^2, y)$$

- With the proposed conjugate priors, this integral is tractable

$$P(\sigma^2 | y_1, \dots, y_n) \sim \text{Inv-Gamma}\left(\frac{v_n}{2}, \frac{n_n \sigma_n^2}{2}\right)$$

$$v_n = v_0 + n,$$

$$\sigma_n^2 = \frac{1}{v_n} \left[ v_0 \sigma_0^2 + (n-1) s^2 + \frac{k_0 n}{n+k_0} (\bar{y} - \mu_0)^2 \right]$$

$$P(\mu | y_1, \dots, y_n) = \underbrace{\int P(\mu | \sigma^2, y)}_{\downarrow} \underbrace{P(\sigma^2 | y) d\sigma^2}_{\text{Normal}}$$

$$\mu = cZ, \quad Z \sim N(0, 1)$$

$$c \sim \text{Inv-Gamma}$$

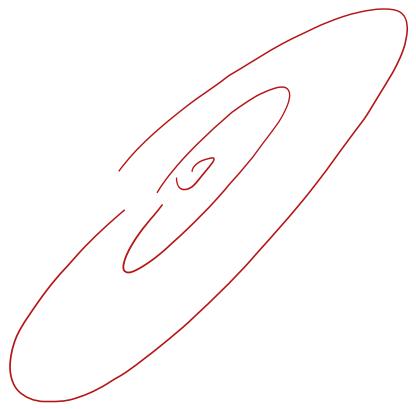
t-distr.

$$y_p = Az, \quad z \sim N(0, I)$$

$$y \sim N(0, \Sigma) \\ "AAT"$$

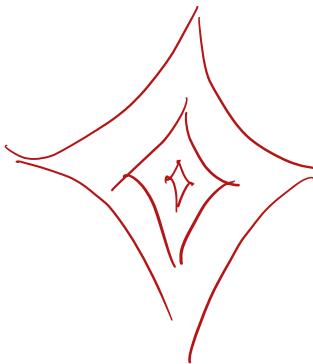
$$y = \sqrt{\frac{v}{c}} Az \sim MVT.$$

$$c \sim \text{Inv-Gamma}$$



still elliptical

$$Y_2 = \begin{pmatrix} c_1 & \\ & c_2 \end{pmatrix} Z$$



# Direct Sampling

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

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

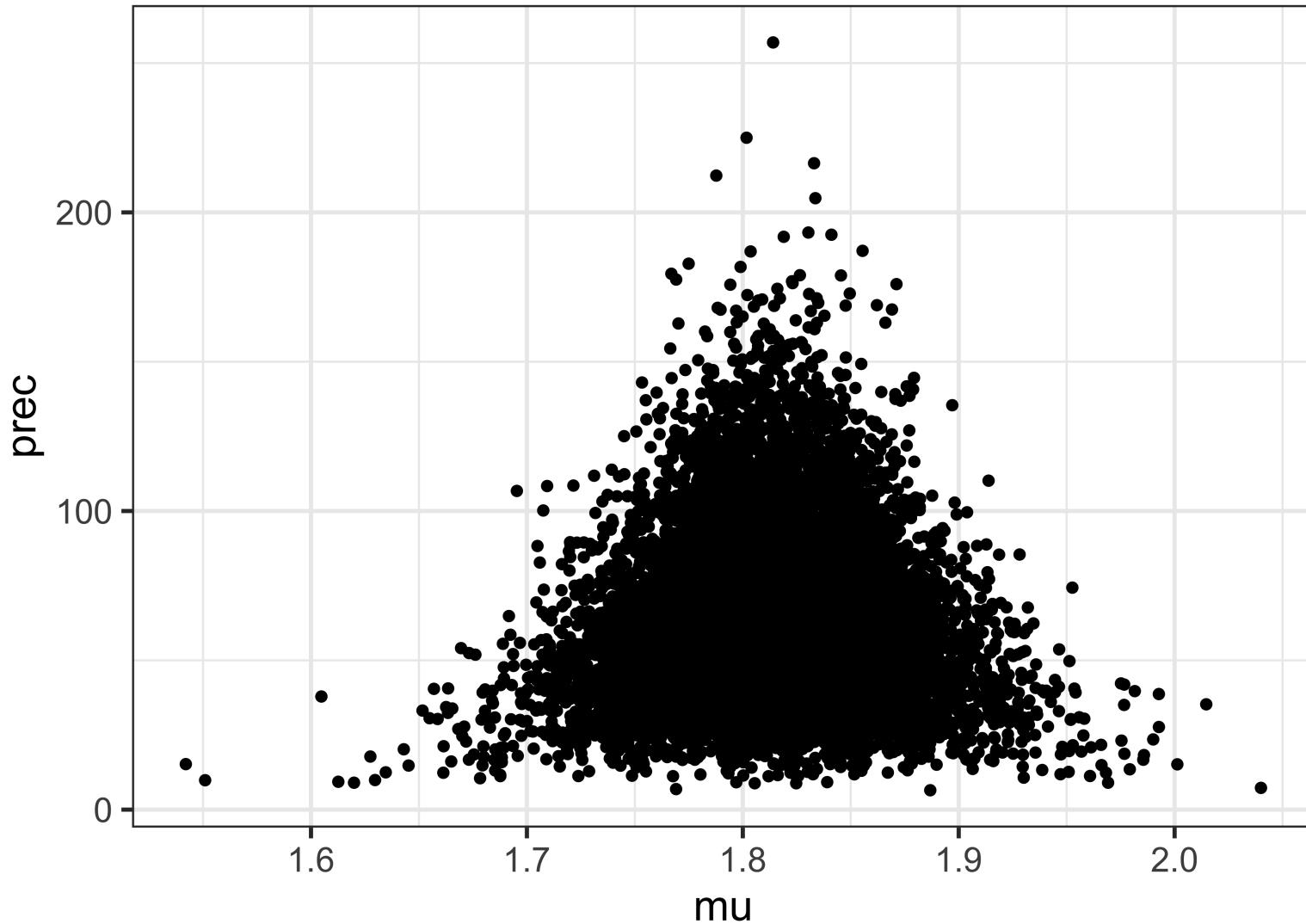
$$\nu_n = \nu_0 + n$$

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

# 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



# The Multivariate Normal Distribution

$$Y_{px1} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_p \end{bmatrix} \sim N_p(\mu, \Sigma)$$

$$P(Y=y | \mu, \Sigma) = (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left(-\frac{1}{2} (y-\mu)' \Sigma^{-1} (y-\mu)\right)$$

$$\mu \in \mathbb{R}^p$$

$\Sigma \in S_p^+$  cone of <sup>sym.</sup> pos. def matrices.

$$a' \Sigma a > 0$$

$p(p+1)/2$  free parameters

Assume  $\Sigma$  is known.

$$P(\mu | \Sigma, y_1, \dots, y_n)$$

$$\mu \sim N_p(\mu_0, \Lambda_0)$$

$$L(\mu) \propto \prod_{i=1}^n (2\pi)^{-p/2} |\Sigma|^{-1/2} \exp\left[-\frac{1}{2}(y_i - \mu)' \Sigma^{-1} (y_i - \mu)\right]$$

$$\propto \exp(-1/2 \mu' A \mu + \mu' b)$$

$$A = n\Sigma^{-1}, \quad b = n\Sigma^{-1}\bar{y}$$

$$P(\mu) \propto \exp(\mu' \Lambda_0^{-1} \mu + \mu' \Lambda_0^{-1} \mu_0)$$

$$P(\mu | y_1, \dots, y_n) \propto \exp(\mu' A_n \mu + \mu' \underbrace{(b + \Lambda_0^{-1} \mu_0)}_{\text{---}})$$

$$A_n = \Sigma^{-1} + \Lambda_0^{-1}, \quad b = n\Sigma^{-1}\bar{y} + \Lambda_0^{-1} \mu_0$$

$$\sim N(\mu_n, \Sigma_n)$$

$$\Sigma_n = (n\Sigma^{-1} + A_0^{-1})^{-1}$$

$$M_n = (n\Sigma^{-1} + A_0^{-1})^{-1} (n\Sigma^{-1}\bar{y} + A_0^{-1}M_0)$$

$\Sigma$  unknown,  $M=0$

$$y_1, \dots, y_n \sim N(0, \Sigma)$$

$y_{ij}$   $i$ th row  
 $j$ th col.

$$S = yy^T = \begin{pmatrix} \sum_i y_{ii}^2 & (\sum_{i=1}^n y_{ii} y_{ik}) \\ \vdots & \ddots & \sum_{i=1}^n y_{ip}^2 \end{pmatrix}$$

"sum of squares matrix"

$$S \sim \text{Wishart}(n, \Sigma)$$

$$P(S) \sim |S|^{-\frac{(n-p-1)}{2}} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} S)}$$

$n > p$ ,  $S$  is pos. def w.p. 1.

$$E[S] = n\Sigma, \quad \frac{YY^T}{n} \rightarrow \Sigma$$

$X \sim \text{Wish}_+$ ,  $X^{-1} \sim \text{Inv-Wish}_+$

Inv-Wish is conjugate for  $\Sigma$

$$L(\Sigma) \propto |\Sigma|^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n \underbrace{y_i' \Sigma^{-1} y_i}_{\text{scalar}}\right)$$

$$\text{tr}(ABC) = \text{tr}(BCA)$$

$$\propto |\Sigma|^{-n/2} \exp\left(-\frac{1}{2} \text{tr}(\Sigma^{-1} S)\right)$$

$$\begin{aligned} & \text{tr}(y' \Sigma^{-1} y) \\ &= \text{tr}(\Sigma^{-1} y y') \end{aligned}$$

$$\Sigma \sim IW(n, \Lambda_0) \rightarrow |\Sigma|^{-\frac{n+p+1}{2}} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} \Lambda_0)}$$

$$P(\Sigma | y_1, \dots, y_n) \sim IW(n+p, S + \Lambda_0)$$

$$\text{Jeffreys} \propto |\Sigma|^{-(p+1)/2}$$

# 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

