

Lecture 5: Hierarchical Modeling

Professor Alexander Franks

Announcements

- Reading: Chapter 5 of BDA

Comparing Multiple Related Groups

- Hierarchy of nested populations
- Models which account for this are called *hierarchical* or *multi-level* models

Some examples:

- Patient outcomes within several different hospitals
- People within counties in the United States (e.g. Asthma mortality example)
- Athlete performance in sports
- Genes within a group of animals

Eight schools example

- A study was performed for the Educational Testing Service (ETS) to evaluate the effects of coaching programs on SAT preparation
- Each of eight different schools used a short-term SAT prep coaching program
- Compute the average SAT score in those who did take the program minus those that did not participate in the program
- We observe the average difference varies by school. What accounts for these differences?

Eight schools example

- Interested in “real” differences due to training
- Want to reduce effect of chance variability
- How do we estimate the effect of the program in each of the schools?

Eight schools example

- Consider two extremes:
 - Estimate the effect of the program in every school independently
 - A separate prior distribution for each school effect
 - Or assume the effect is the same in every school
 - Combine all the data
 - A compromise between the above 2 options?

Eight Schools Example

$$y_j \sim N(\theta_j, \sigma_j^2)$$

- y_j is the observed effects of the program in school j
 - Based on a sample of test scores from those in the program and those not in the program
- θ_j are the true *unknown* effects of the program in school j
- Assume variances, σ_j^2 , are *known*
 - e.g. determined by the number of students in the sample

Eight Schools Example

```
1 J <- 8  
2 y = c(28, 8, -3, 7, -1, 1, 18, 12)  
3 sigma <- c(15, 10, 16, 11, 9, 11, 10, 18)
```

- Assuming the effect of the program on each school is identical.
- What are the chances of seeing a value as large as 28?
- As small as -3?

Eight Schools Example

- Assume the effect of the program on each school is identical, i.e. $\theta_j = \mu$
- Assume a flat prior on μ , what is $p(\mu \mid y_1, \dots, y_8, \sigma_1, \dots, \sigma_8)$?

Eight Schools Example

```
1 ## Compute the precision frome each school  
2 prec <- 1/sigma^2  
3  
4 ## global estimate is a weighted vareage  
5 mu_global <- sum(prec * y / sum(prec))  
6 mu_global
```

```
[1] 7.685617
```

Eight Schools Example

- Assume the effect of the program on each school is identical, i.e. $\theta_j = \mu$
- What are the chances of school 1 having an effect large as 28 (given $\sigma_1 = 15$)?
- Y_3 as small as -3 (given $\sigma_3 = 16$)?

Posterior Prediction Under Complete Pooling

```
1 prec <- 1/sigma^2
2
3 ## global estimate is a weighted average
4 mu_global <- sum(prec * y / sum(prec))
5
6 print(sprintf("mu is %f", mu_global))
[1] "mu is 7.685617"
1 1 - pnorm(28, mean=mu_global, sd=sqrt(1/sum(1/sigma^2)+sigma[1]^2))
[1] 0.09560784
1 pnorm(-3, mean=mu_global, sd=sqrt(1/sum(1/sigma^2)+sigma[3]^2))
[1] 0.2587447
```

Eight Schools Example

$$y_j \sim N(\theta_j, \sigma_j^2)$$

- θ_j are the true unknown effects of the program in school j
- y_j is the observed effects of the program in school j
 - Based on a sample of test scores from those in the program and those not in the program
 - Number of people in the sample determine the magnitude of σ_j^2

Eight Schools Example

How do we estimate θ_j ?

- Assume effects are totally independent: $\hat{\theta}_j^{(MLE)} = y_j$ is the MLE
- Assume effects are identical: $\hat{\theta}_j^{(pool)} = \frac{\sum_i \frac{1}{\sigma_i^2} y_i}{\sum \frac{1}{\sigma_i^2}}$
 - Same effect for all schools: estimate using a weighted average of the observed effects

Eight Schools

```
1 theta_j_mle <- y  
2 theta_j_mle
```

```
[1] 28   8  -3   7  -1   1  18  12
```

```
1 theta_j_pooled <- rep(sum(1/sigma^2 * y) / sum(1/sigma^2), J)  
2 theta_j_pooled
```

```
[1] 7.685617 7.685617 7.685617 7.685617 7.685617 7.685617 7.685617 7.685617 7.685617
```

- Compromise: $\hat{\theta}_j^{\text{shrink}} = w\theta_j^{\text{MLE}} + (1 - w)\theta^{\text{pooled}}$

Eight schools example

Add a *shared* normal prior distribution to θ_j

$$\theta_i \sim N(\mu, \tau^2)$$

$$y_j \sim N(\theta_j, \sigma_j^2)$$

- The global mean, μ , is also an unknown parameter. What prior should we choose?
- τ^2 determines how much weight weight we put on the independent estimate vs the pooled estimate.
- A 9-parameter posterior:

$$p(\mu, \theta_1, \dots, \theta_8 \mid y_1, \dots, y_8, \sigma_1, \dots, \sigma_8, \tau^2)$$

Intuition for shrinkage

- $Y_j = \theta_j + \epsilon_j$
 - For simplicity assume $Var(\epsilon_j) = \sigma^2$ for all j
 - θ_j represents true effect in school j (signal)
 - $Var(\theta_j) = \tau^2$ represents how much the true effects vary across schools
 - ϵ_j is sampling variability (noise, chance variation)
- $\hat{\theta}_{MLE} = Y_j$

Intuition for shrinkage

- Consequence: the observed outcomes always have higher variance than the signal, i.e. $\text{Var}(Y_j) > \text{Var}(\theta_j)$
- Intuition: reduce the variance by shrinking estimates to a common mean!
- The variance of the shrunken estimates should be close to τ^2

Eight schools example

Questions:

- Is the training program effective in school j ?
 - What is $P(\theta_j > 0 \mid y)$?
- On average (over all schools) is the training program effective?
 - What is $P(\mu > 0 \mid y)$?
- Will the training program be effective in a new school?
 - What is $P(\mu_{J+1} > 0 \mid y)$?

Eight schools example

Comments:

- The global average, μ , is a parameter so also has uncertainty
- How do we determine how much to shrink, e.g. how do we determine τ^2 ?
- What σ_j^2 were also unknown?

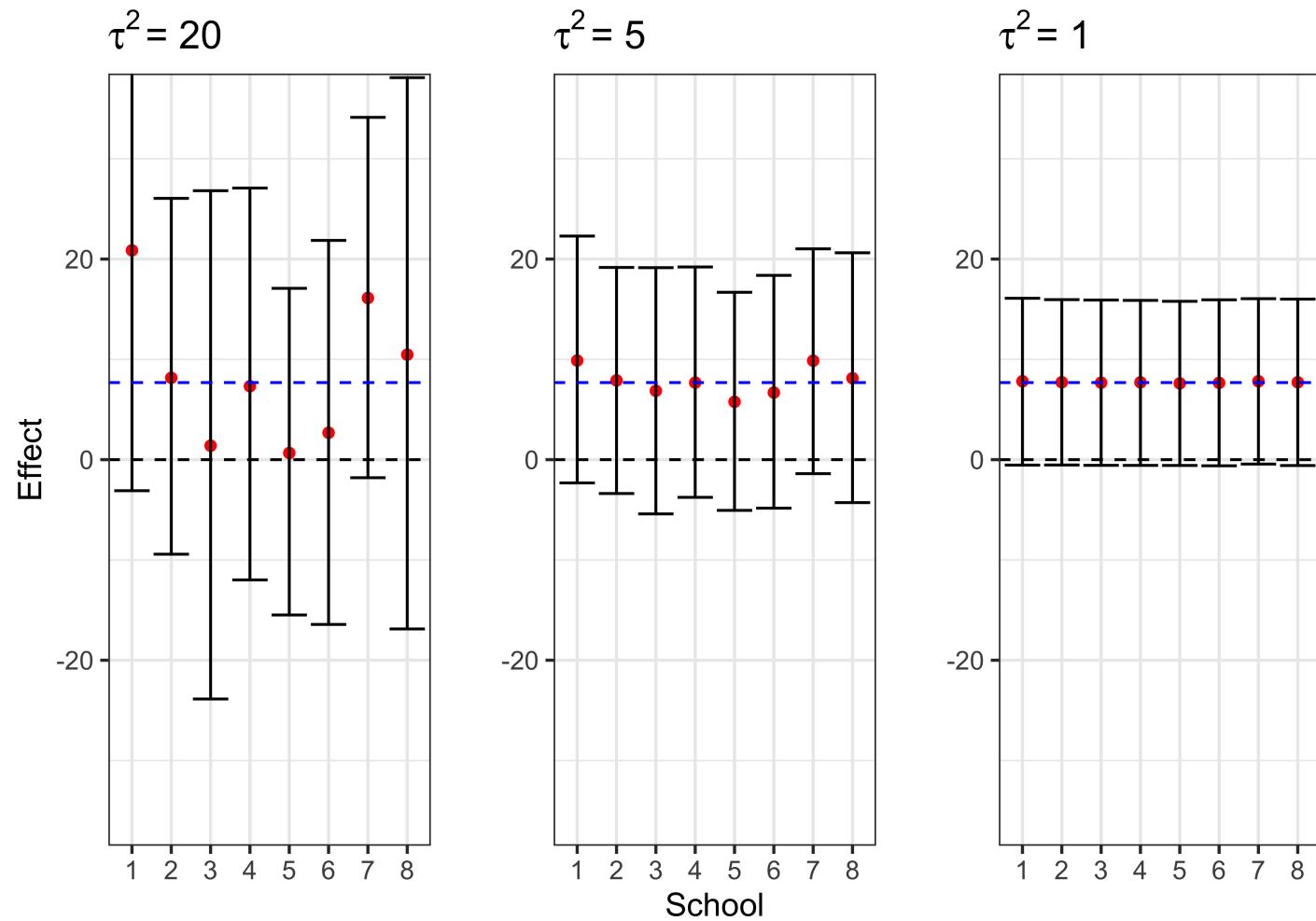
Eight schools example

- If τ^2 is large, the prior for θ_j is not very strong
 - If $\tau^2 \rightarrow \infty$ equivalent to the no pooling model
- If τ^2 is small, we assume a priori that θ_j are very close
 - if $\tau^2 \rightarrow 0$ equivalent to the complete pooling model,
 $\theta_j = \mu$

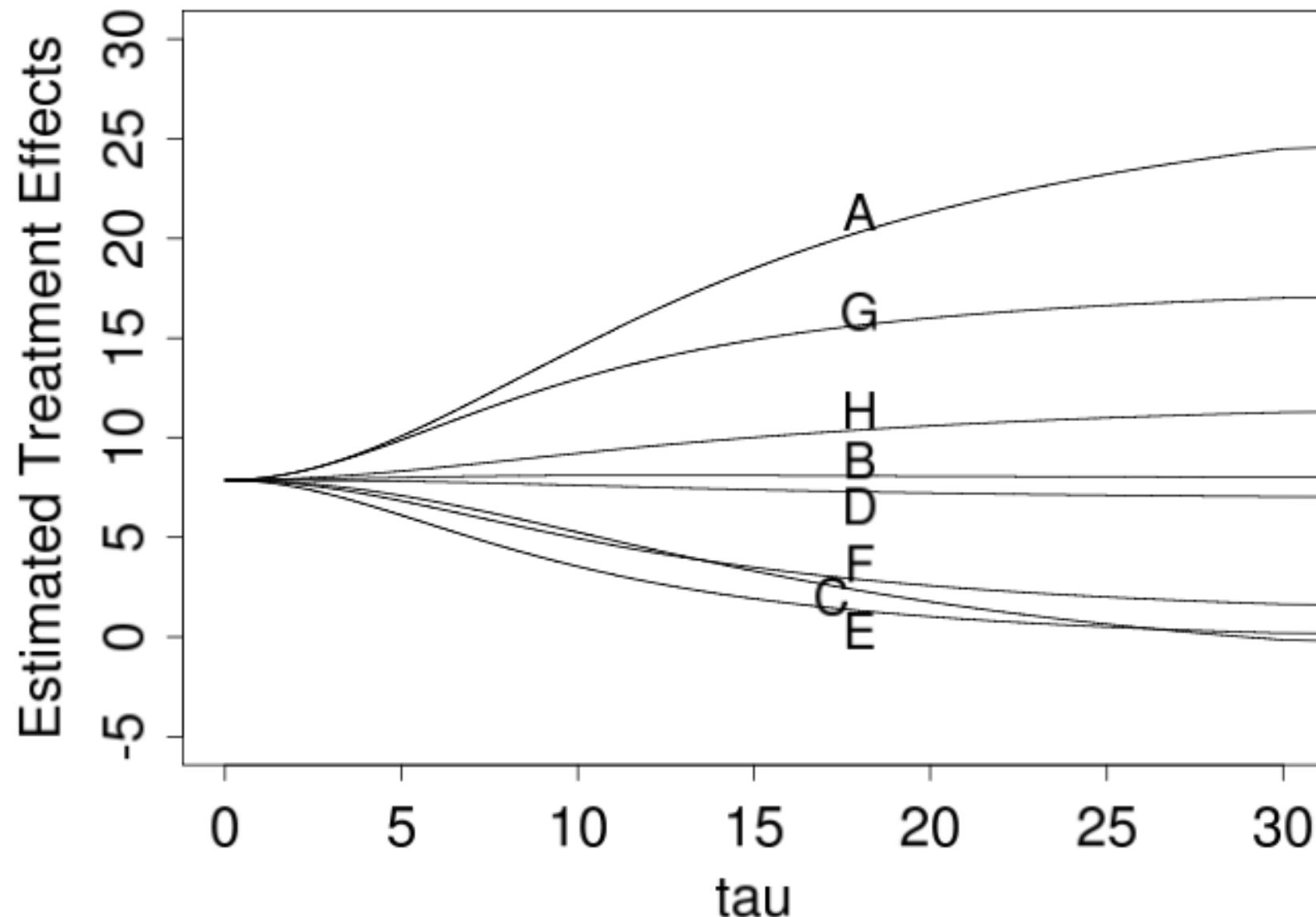
Inference

- Factorize the density into tractable components
 - $p(\mu \mid y_1, \dots, y_8, \tau^2)$
 - $p(\theta_i \mid \mu, y_i, \tau^2)$
- Later: MCMC or other approximate methods

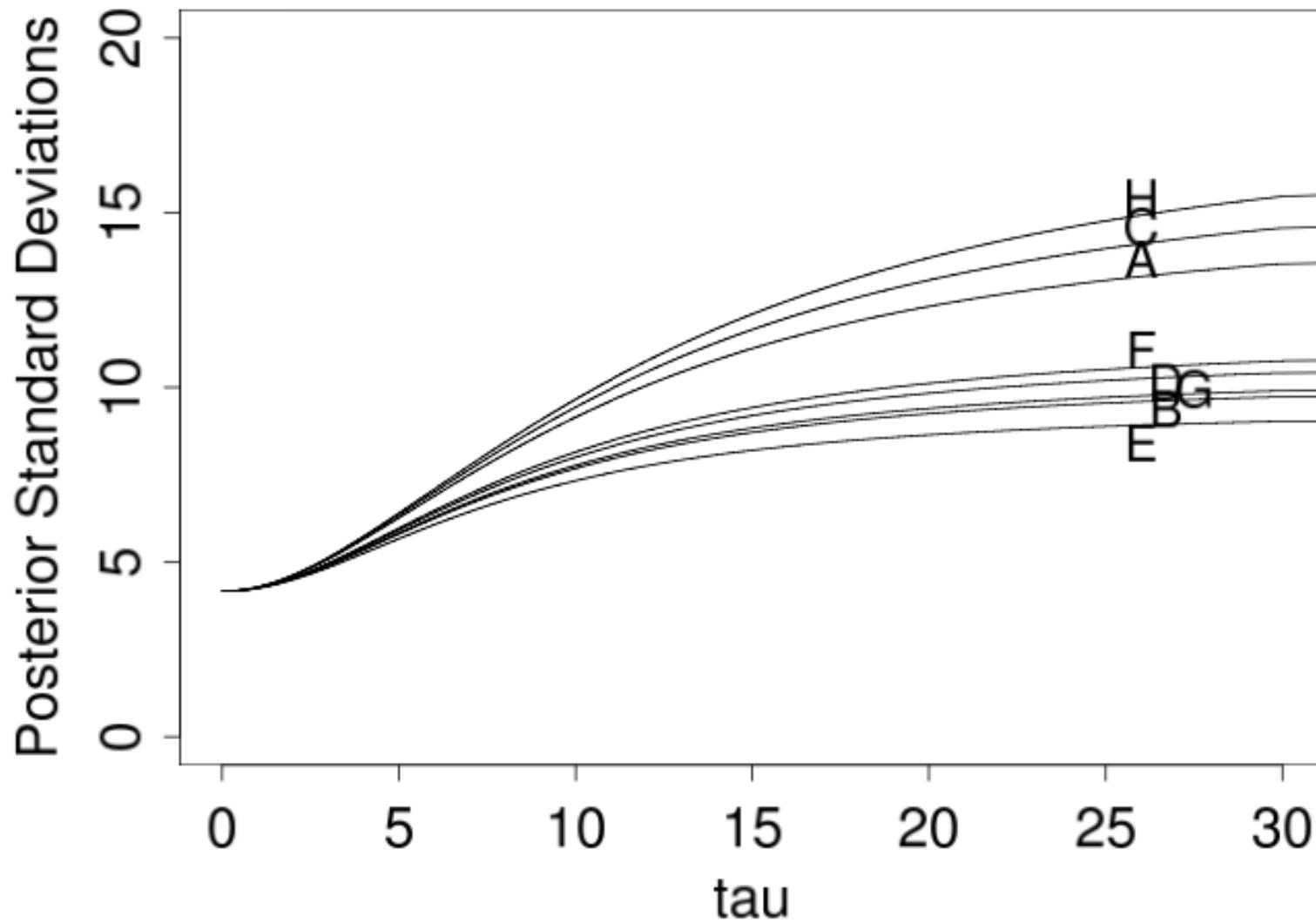
Eight Schools example



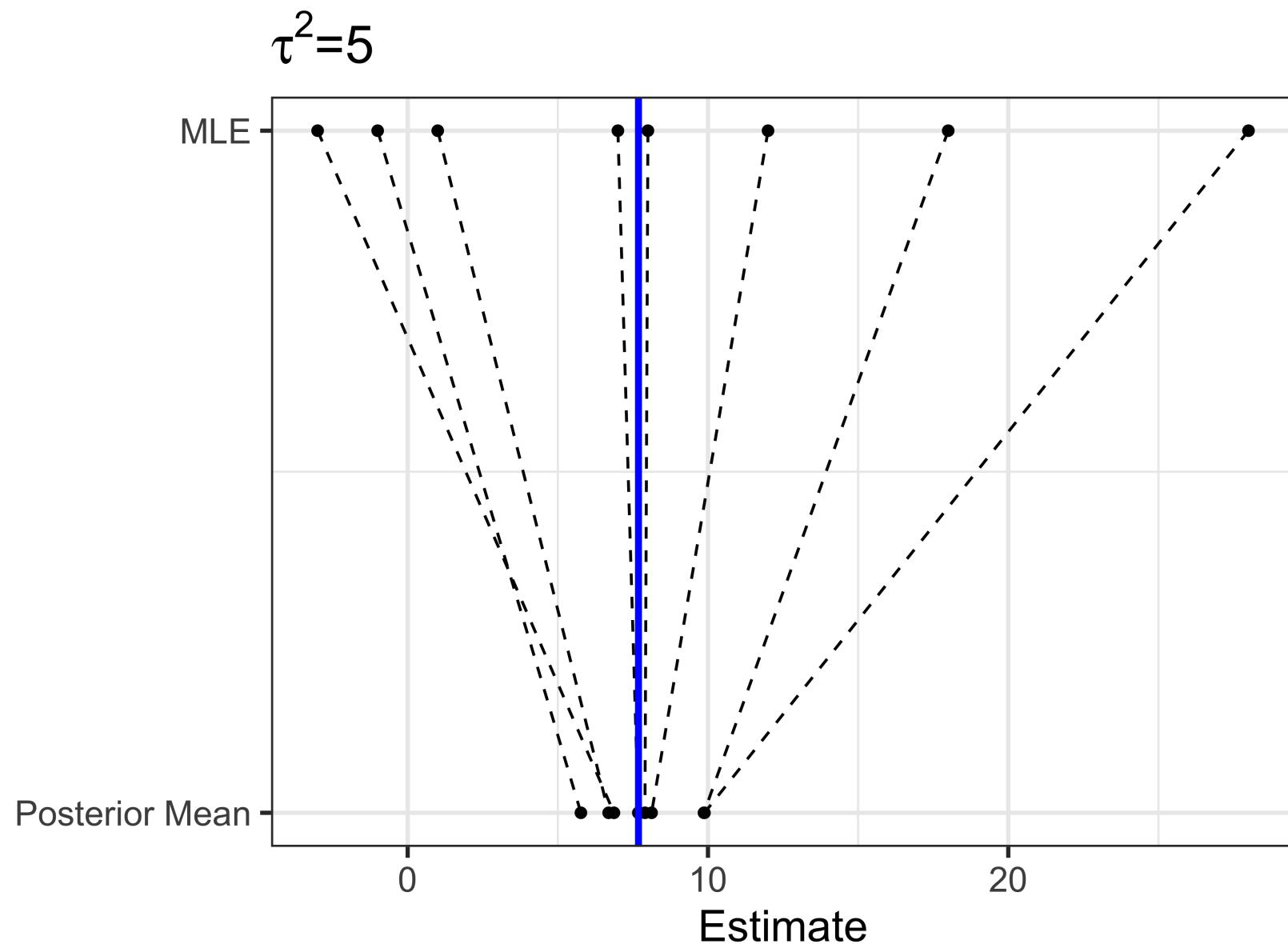
The impact of τ



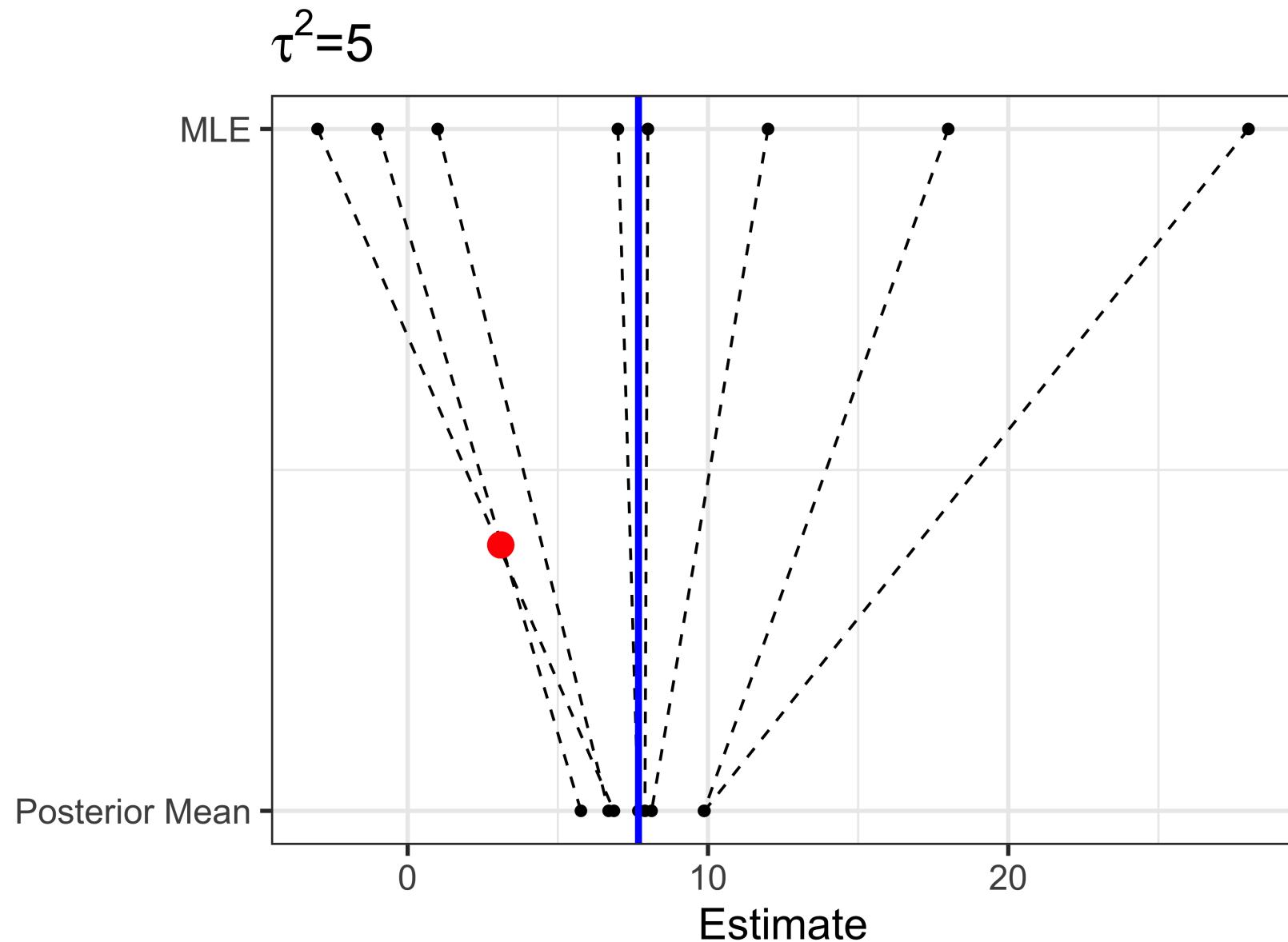
The impact of τ



MLE vs Posterior Mean



MLE vs Posterior Mean



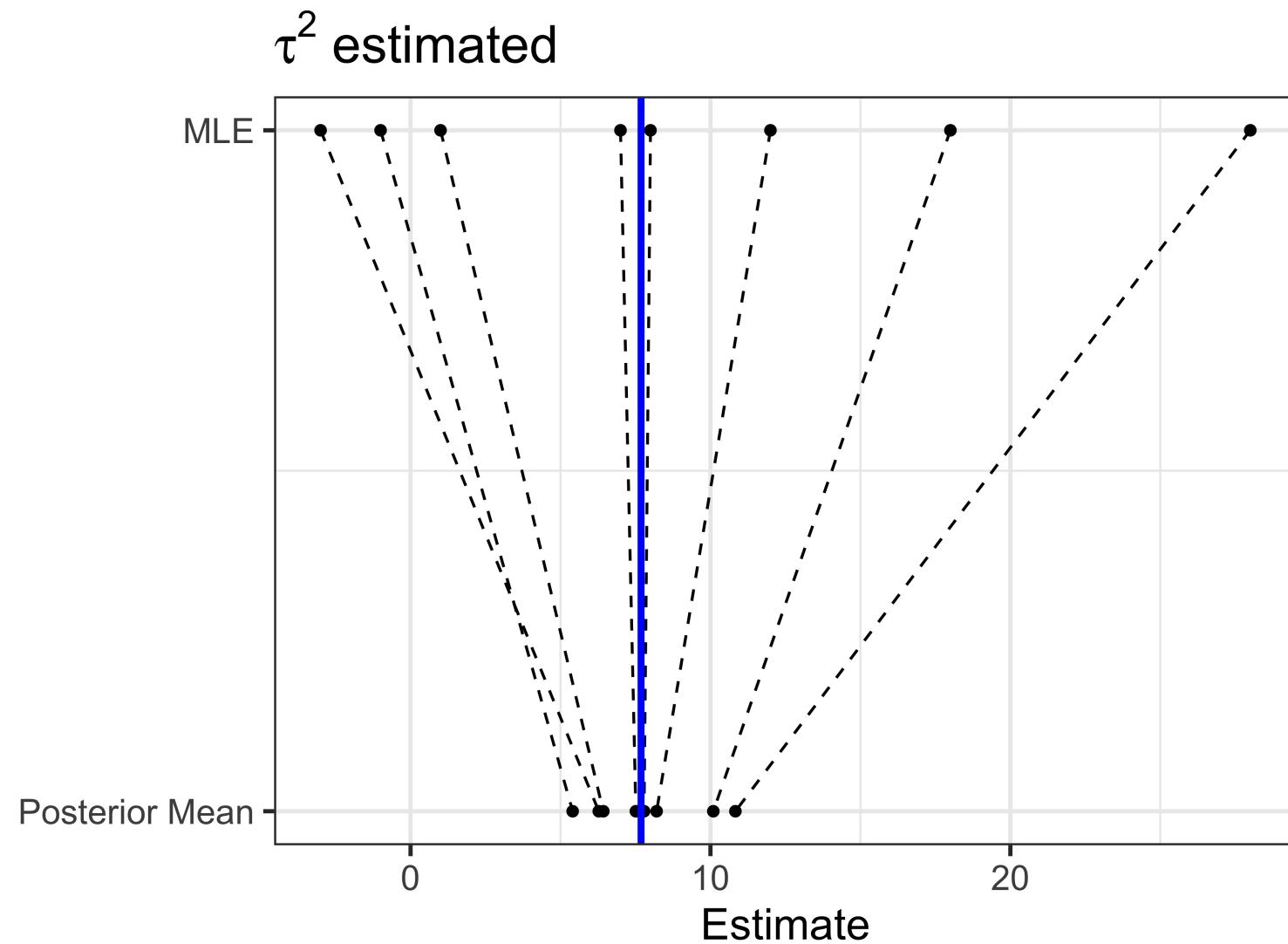
Inference for τ

- Can infer τ (don't need to set τ as a hyperparameter)
- How?

Weak and noninformative Priors on τ

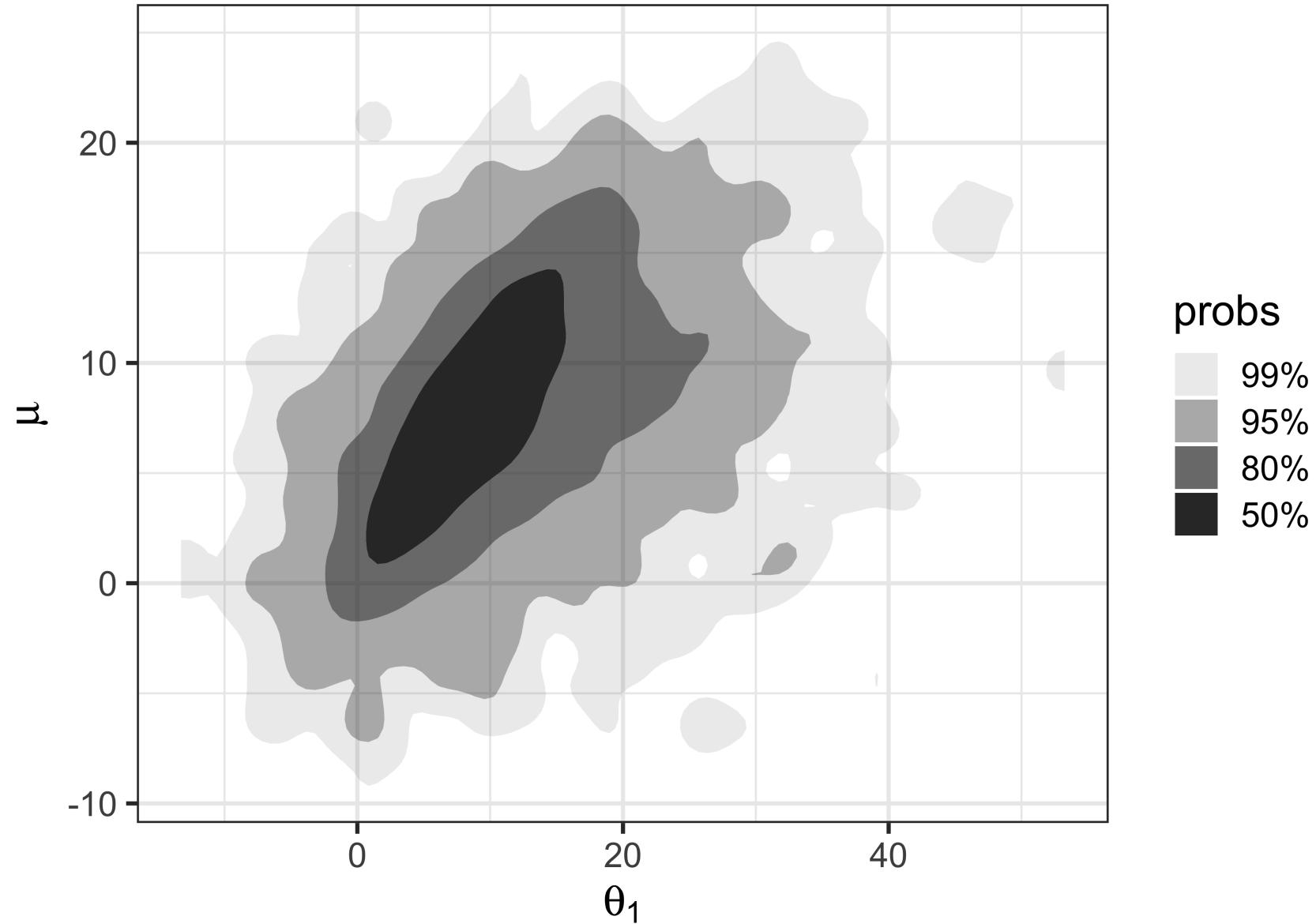
- Consider limits of proper priors
 - Uniform[0, A] as $A \rightarrow \infty$ (ok for $J > 2$)
 - Inverse-Gamma(ϵ, ϵ) as $\epsilon \rightarrow 0$ (improper posterior!)
- Uniform on $\log(\tau)$ (improper posterior)
 - $p(y | \tau) \rightarrow \text{const}$ as $\tau \rightarrow 0$
- Half-Cauchy prior distribution on τ^2
 - Recommended by Gelman et al

MLE vs Posterior Mean

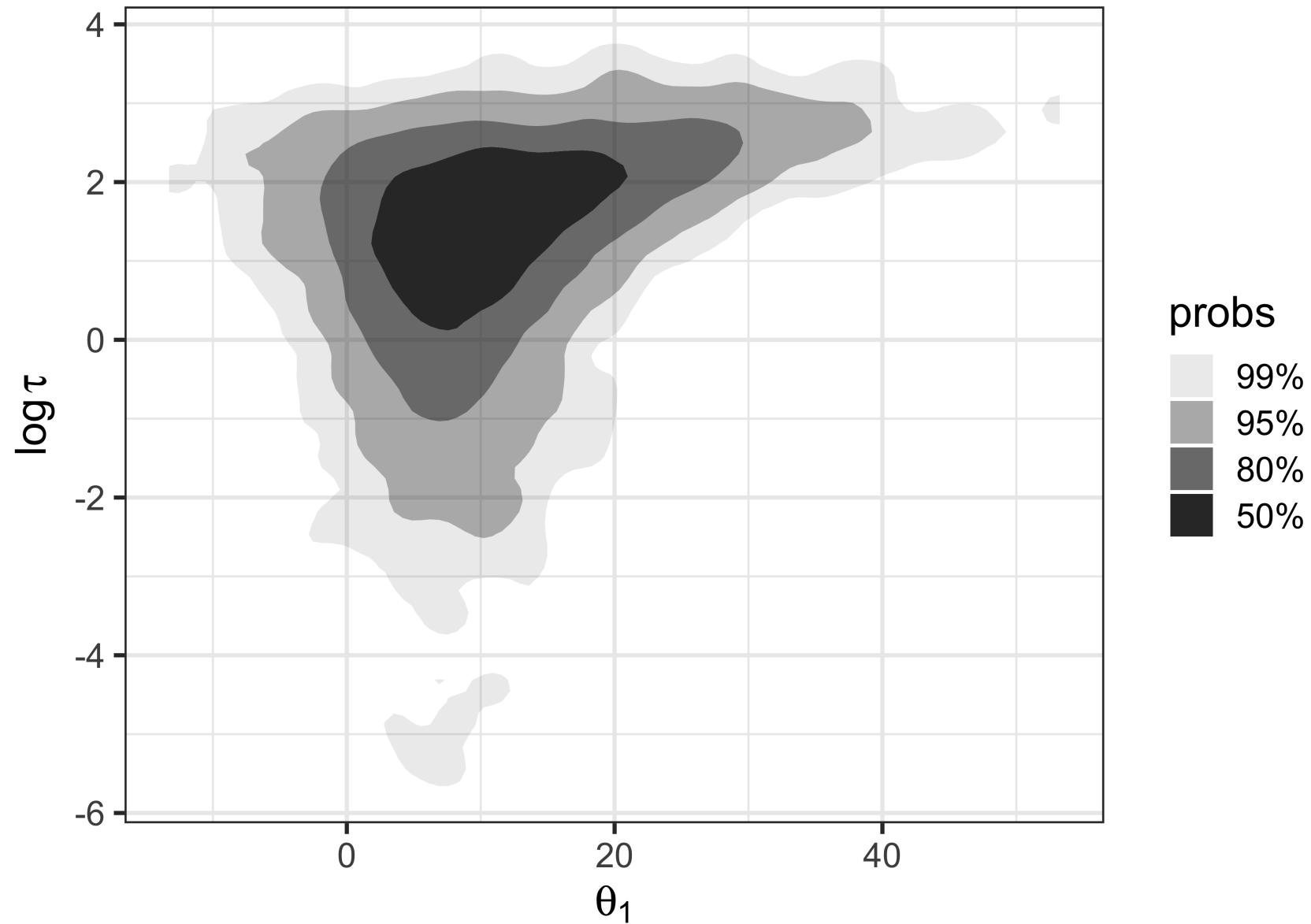


Posterior mean of $\tau = 5.6422886$

Eight Schools Density Plot: θ_1 vs μ



Eight Schools Density Plot: θ_1 vs $\log(\tau)$



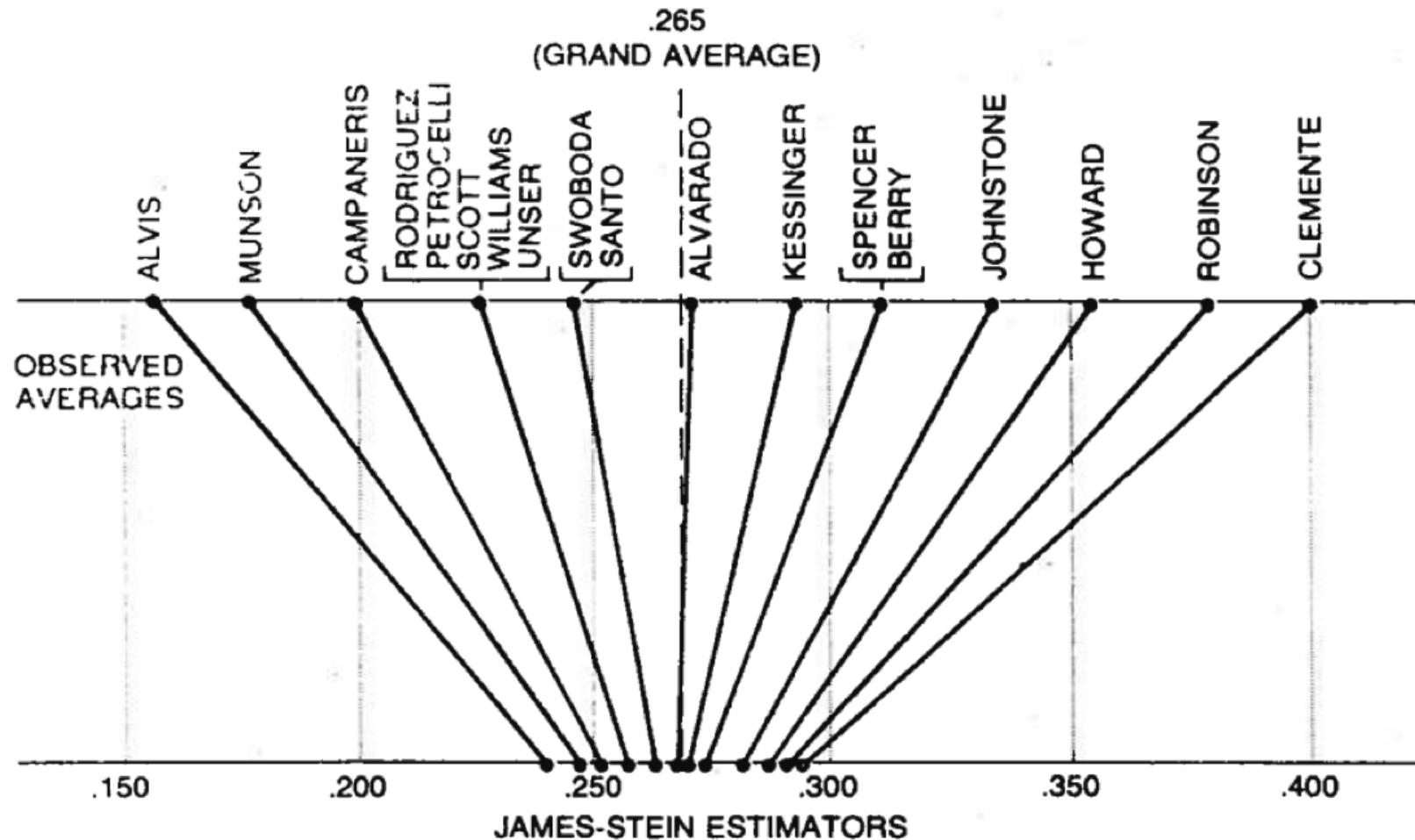
Sampling variance

What if the sampling variances are unknown?
Hierarchical modeling in sports

1. 1970 Batting Averages for 18 Major League Players and Transformed Values X_i, θ_i

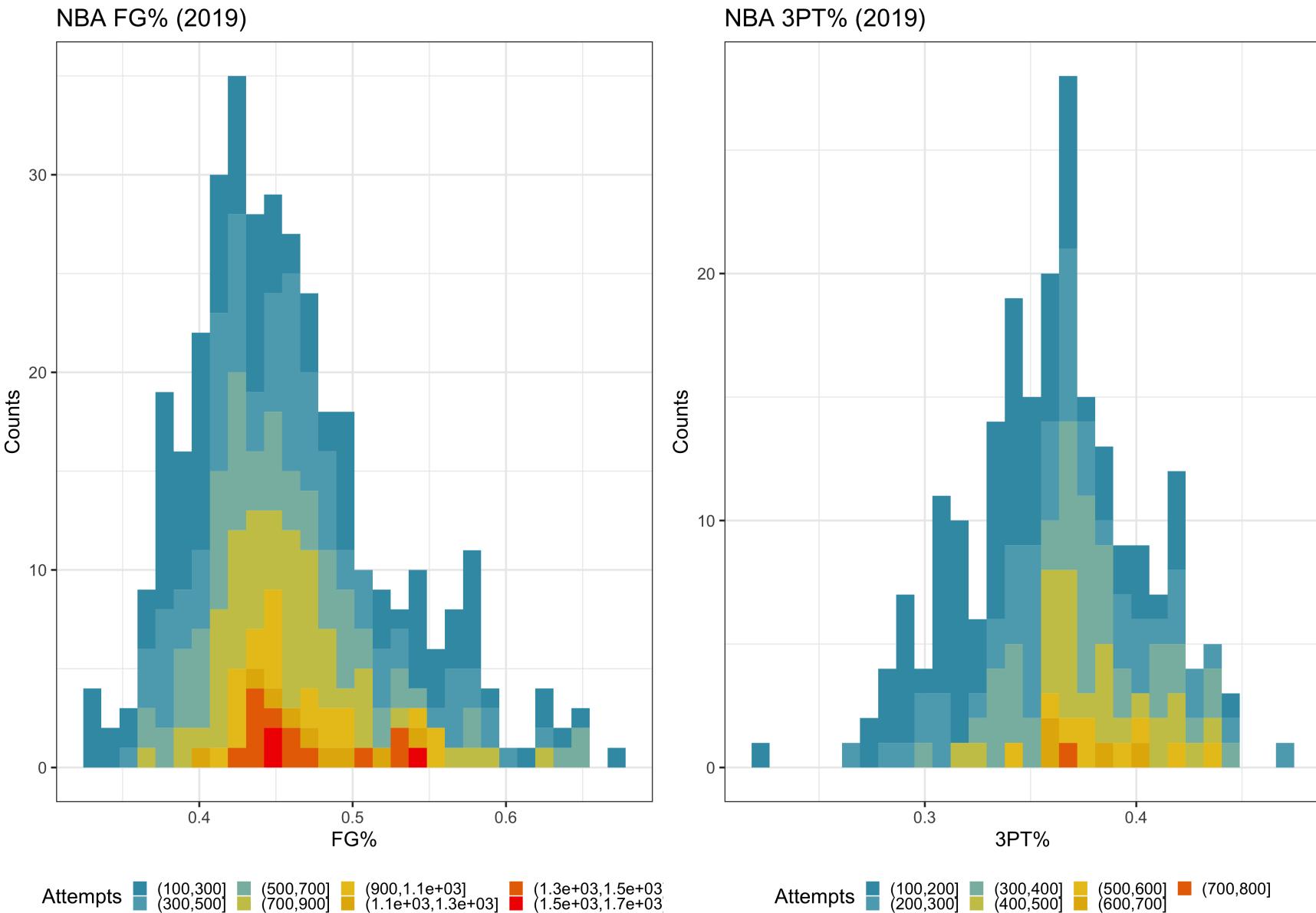
i	Player	$Y_i = \text{batting average for first 45 at bats}$	$p_i = \text{batting average for remainder of season}$	$\text{At bats for remainder of season}$	X_i	θ_i
		(1)	(2)	(3)	(4)	(5)
1	Clemente (Pitts, NL)	.400	.346	367	-1.35	-2.10
2	F. Robinson (Balt, AL)	.378	.298	426	-1.66	-2.79
3	F. Howard (Wash, AL)	.356	.276	521	-1.97	-3.11
4	Johnstone (Cal, AL)	.333	.222	275	-2.28	-3.96
5	Berry (Chi, AL)	.311	.273	418	-2.60	-3.17
6	Spencer (Cal, AL)	.311	.270	466	-2.60	-3.20
7	Kessinger (Chi, NL)	.289	.263	586	-2.92	-3.32
8	L. Alvarado (Bos, AL)	.267	.210	138	-3.26	-4.15
9	Santo (Chi, NL)	.244	.269	510	-3.60	-3.23
10	Swoboda (NY, NL)	.244	.230	200	-3.60	-3.83
11	Unser (Wash, AL)	.222	.264	277	-3.95	-3.30
12	Williams (Chi, AL)	.222	.256	270	-3.95	-3.43
13	Scott (Bos, AL)	.222	.303	435	-3.95	-2.71
14	Petrocelli (Bos, AL)	.222	.264	538	-3.95	-3.30
15	E. Rodriguez (KC, AL)	.222	.226	186	-3.95	-3.89
16	Campaneris (Oak, AL)	.200	.285	558	-4.32	-2.98
17	Munson (NY, AL)	.178	.316	408	-4.70	-2.53
18	Alvis (Mil, NL)	.156	.200	70	-5.10	-4.32

Hierarchical modeling in sports



JAMES-STEIN ESTIMATORS for the 18 baseball players were calculated by “shrinking” the individual batting averages toward the overall “average of the averages.” In this case the grand average is .265 and each of the averages is shrunk about 80 percent of the distance to this value. Thus the theorem on which Stein’s method is based asserts that the true batting abilities are more tightly clustered than the preliminary batting averages would seem to suggest they are.

Basketball Example



Attempts

Attempts

A Beta/Binomial Hierarchical Model

- If we know α and β then the player effects are independent:

$$\begin{aligned} p(\alpha, \beta, \theta_1, \dots, \theta_{50} \mid (y_1, n_1), \dots, (y_{50}, n_{50})) &= \\ p(\alpha, \beta \mid (y_1, n_1), \dots, (y_{50}, n_{50})) \prod_i p(\theta_i \mid (y_i, n_i), \alpha, \beta) \end{aligned}$$

Monte Carlo procedure:

Modeling Three Point Shooting

$$p(\alpha, \beta | y_1, \dots, y_{50})$$

```
1 ## Evaluate the log marginal posterior
2 marginal_posterior <- function(pars) {
3   alpha <- pars[1]
4   beta <- pars[2]
5
6   ## log posterior
7   -5/2*log(alpha+beta) +
8     sum(lgamma(alpha + beta) +
9       lgamma(alpha + y) + lgamma(beta+n-y) -
10      lgamma(alpha) - lgamma(beta) - lgamma(alpha+beta+n))
11 }
12
13 res <- optim(c(50, 50), function(x) -1*marginal_posterior(x))
14 bball_map <- res$par
```

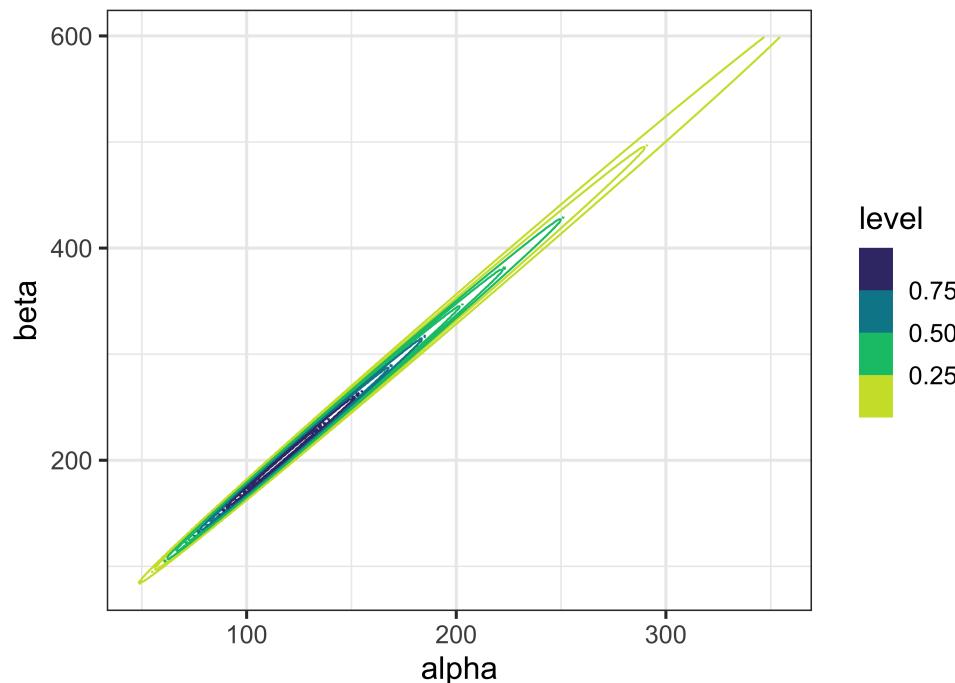
The MAP estimate is (114.4, 196.6)

Empirical Bayes

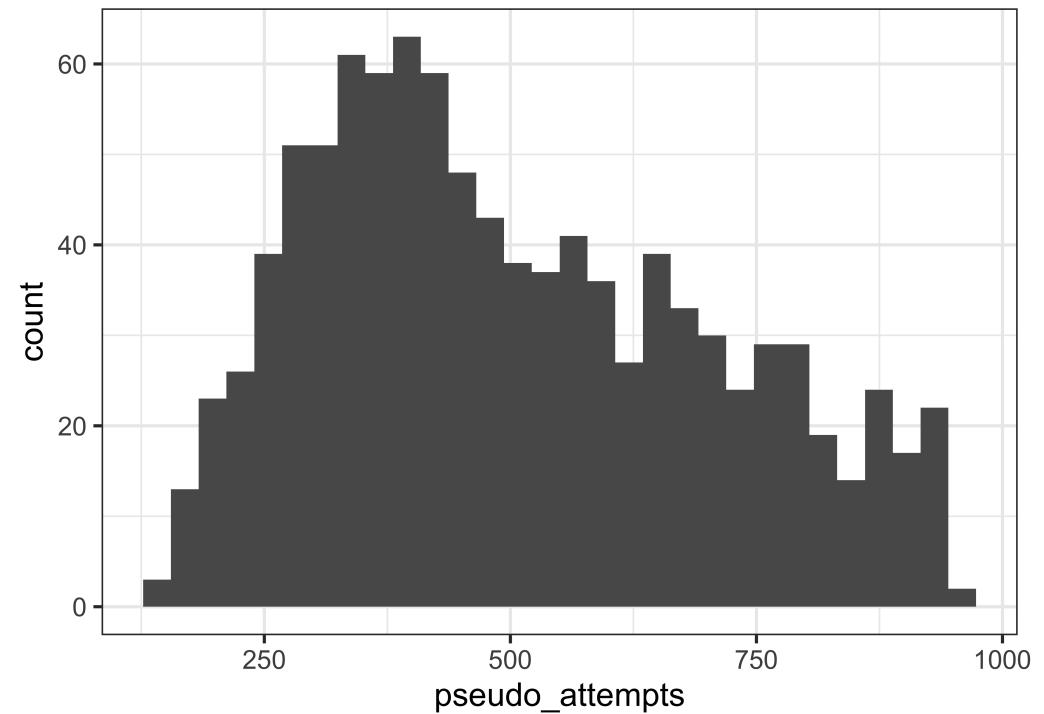
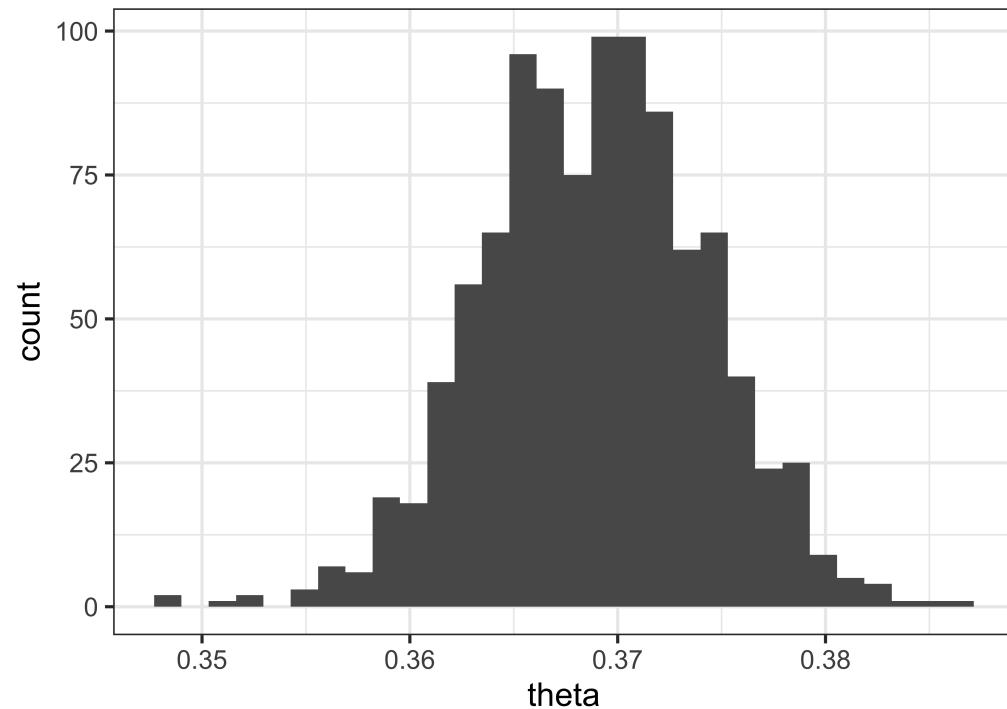
- An **empirical Bayes** solution is to find a MAP estimate of α and β and condition on them as prior parameters.
- Why is this not “true” Bayes?
- When will empirical Bayes be a reasonable approximation to the full Bayes solution?

Modeling Three Point Shooting

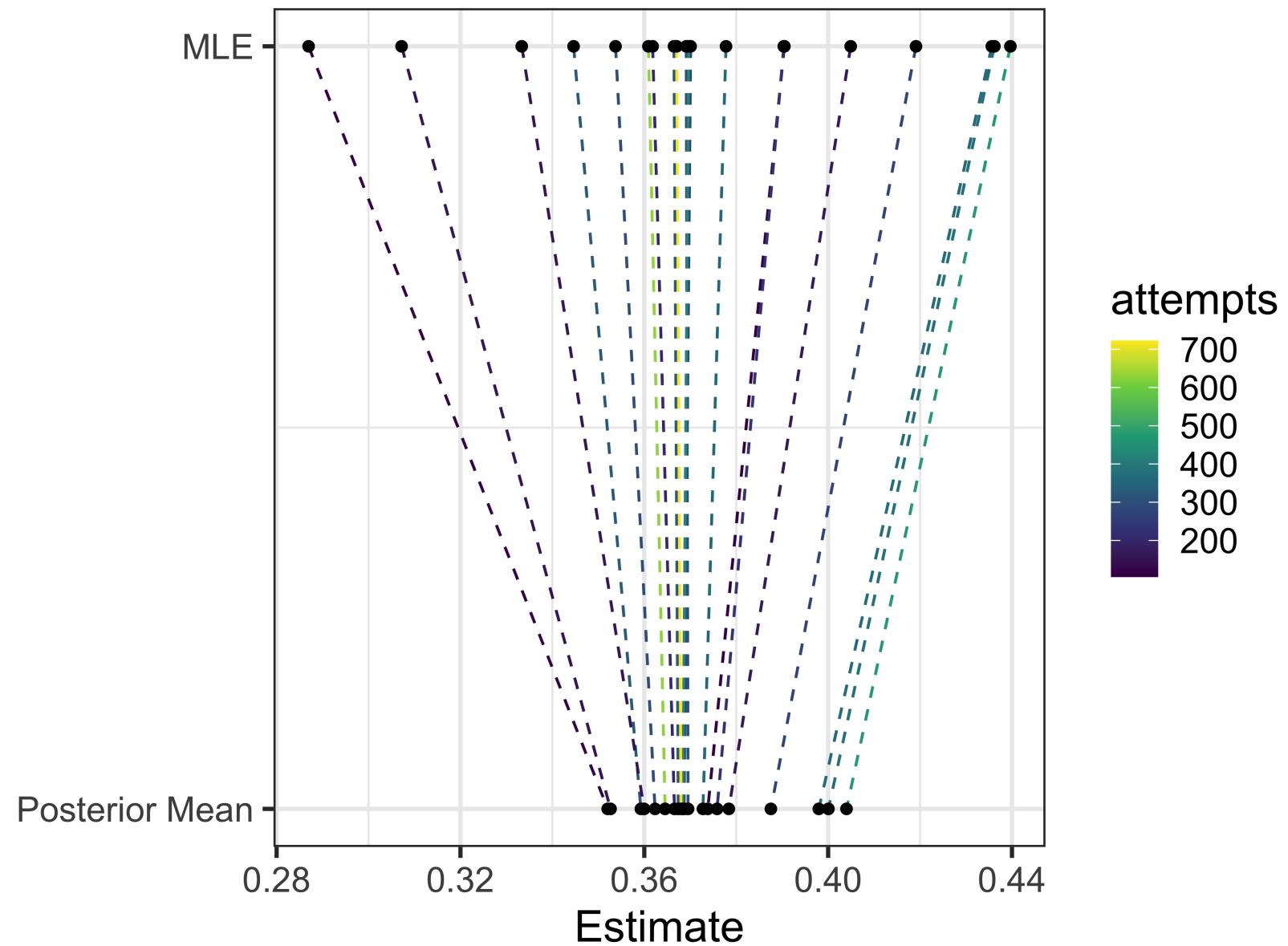
```
1 alpha_beta_grid <- expand.grid(seq(25, 600, by=2), seq(25, 600, by=2))
2 marginal_post <- apply(alpha_beta_grid, 1, marginal_posterior)
3 marginal_post <- marginal_post - max(marginal_post)
4
5 tibble(alpha_beta_grid) %>% rename(alpha=Var1, beta=Var2) %>%
6   ggplot() +
7   geom_contour(aes(x=alpha, y=beta, z=exp(marginal_post), colour=..level..))
8   colorspace::scale_color_binned_sequential("Viridis")
```



Summaries of the Marginal Posterior



Sample theta's

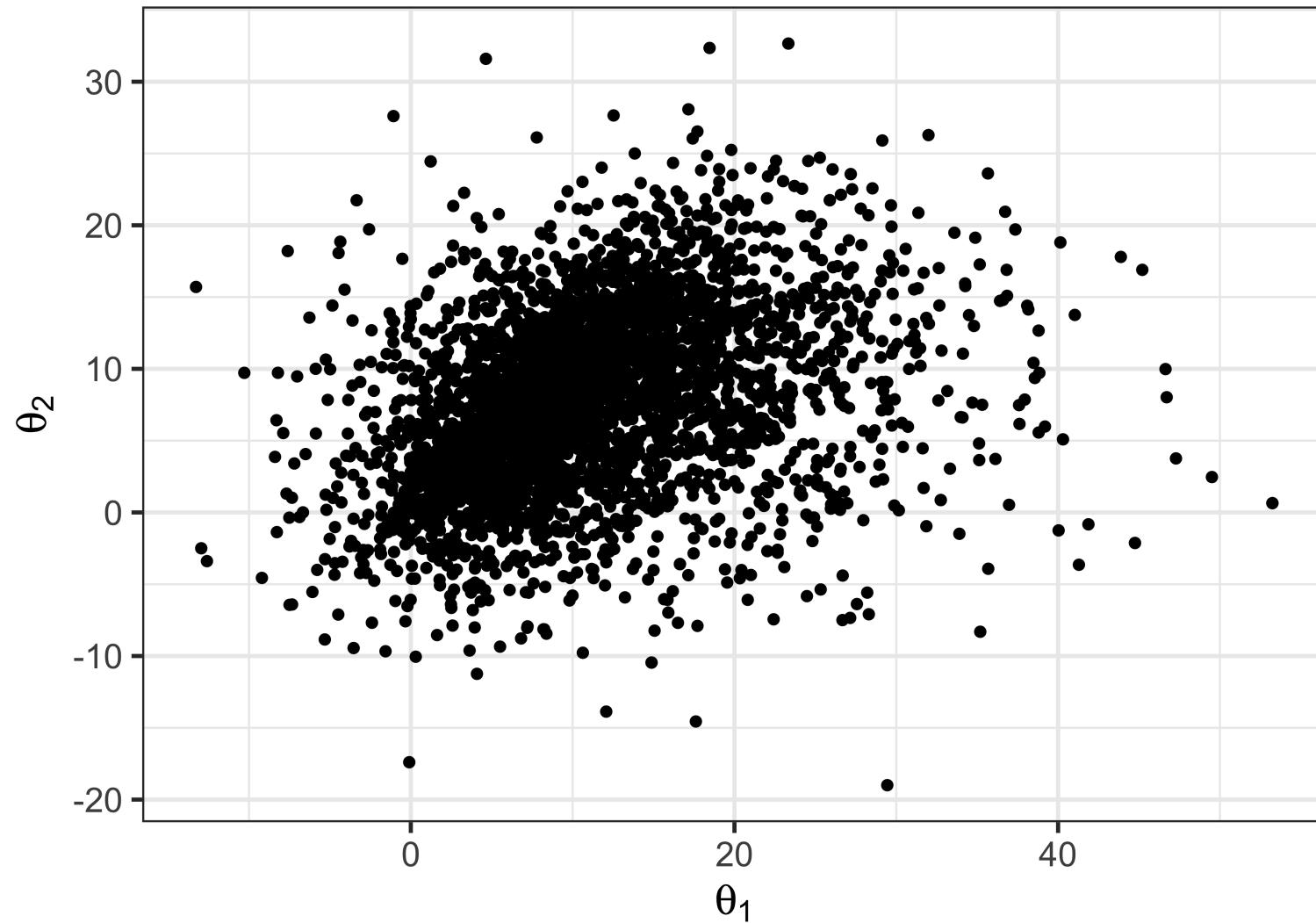


Hierarchical Modeling and Independence

- The shooting skills θ_j are *not* independent!
- E.g. if we know θ_1 is really large, it means θ_j is also more likely to be large
- However, before we see data the distributions of θ_j 's, are indistinguishable
 - Don't know which θ 's will be large and which small

Hierarchical Modeling and Independence

Back to eight schools...



Exchangeability

An exchangeable sequence of random variables is a finite (or infinite) sequence X_1, X_2, \dots, X_n of random variables such that for any finite permutation π of the indices $1, 2, 3, \dots$, the joint probability distribution of the permuted sequence

$$X_{\pi(1)}, X_{\pi(2)}, X_{\pi(3)}, \dots$$

is the same as the joint probability distribution of the original sequence.

Exchangeability

- Consider a set of J experiments in which experiment j has data y_j and parameter θ_j
- The J experiments are related, but no information in the indices that distinguish θ_j
- Assume an exchangeable prior distribution:
 $p(\theta_1, \dots, \theta_J) = p(\theta_{\pi_1}, \dots, \theta_{\pi_J})$ where π is any permutation of the indices $1 \dots J$
- Equivalent to a prior assumption about symmetry among the parameters $(\theta_1, \dots, \theta_J)$

Exchangeability

- Example: All i.i.d random variables are exchangeable
- Example: Bernoulli random variables conditional on the sum
- Example: multivariate normal with common mean μ and equi-correlation

Non-exchangeable random variables

- Time series data (closer in time -> more correlated)
- Spatial data (closer in space -> more correlated)
- Typically, ignorance implies exchangeability
- When exchangeability doesn't hold, can often assume conditional exchangeability

Mixture of i.i.d random variables

- Conditionally i.i.d random variables are exchangeable:

$$p(\theta_1, \dots, \theta_J \mid \theta) = \prod_{j=1}^J p(\theta_j \mid \phi)$$

- Mixtures of i.i.d random variables:

$$p(\theta_1, \dots, \theta_J) = \int \left\{ \prod_{j=1}^J p(\theta_j \mid \phi) \right\} p(\phi) d\phi$$

- If ϕ were known, θ_j would be i.i.d.
- Since ϕ is not known, integrate over uncertainty
- θ_j are a mixture of i.i.d random variables
- Dependent but exchangeable

de Finetti's theorem

- **Theorem:** As $J \rightarrow \infty$ all exchangeable distributions can be expressed as a mixture of independent and identical distributions
- For a finite J , no guarantee that the variables can be represented as a mixture of i.i.d random variables
- But, if variables are exchangeable usually can be modeled as *approximately* a mixture of i.i.d
- Mixture's of i.i.d random variables \rightarrow a hierarchical model is reasonable

