# 36-617: Applied Linear Models

Bayes, Shrinkage, and Multi-Level Models

Brian Junker

132E Baker Hall

brian@stat.cmu.edu

# Announcements

- **No new HW, no quiz this week**
  - ❑ Work on your IDMRAD rough drafts (Due Weds / Grace Fri)
  - ❑ IDMRAD papers submitted to Canvas, not Gradescope
    - ▪ Details will be on Canvas
- **See handout on rough draft IDMRAD papers**
- **Reading: Please at least skim**
  - ❑ Lynch Ch 3: Basics of Bayesian Statistics
    - ▪ Worth reading a little more carefully than a skim
  - ❑ Lynch Ch 4: Modern Model Estimation Part 1: Gibbs Sampling
    - ▪ Read 4.1, 4.2 more carefully; skim rest of chapter
  
  (these are in the week12 folder on canvas!)

# Outline

- **Today:**
  - Shrinkage
  - Review of MLE
  - Crash course in Bayes
  - Normal-Normal Model & Shrinkage
  - MLM's and Shrinkage
- **Project Discussion**
- **After Thanksgiving:**
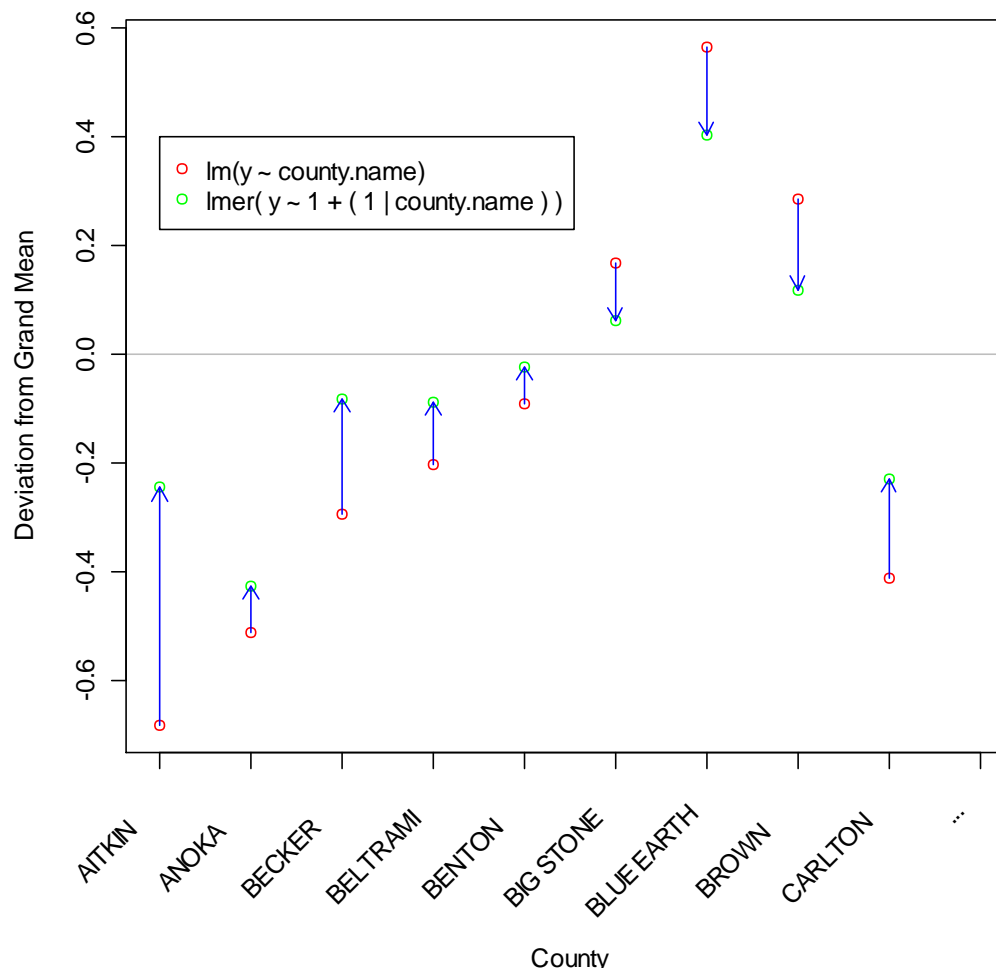  - A little practical Bayes / MCMC for multi-level models

# An MLM phenomenon: Shrinkage

The fitted multilevel model underpredicts high obs's and overpredicts low ones.

The distribution assumptions underlying lmer() "smooth out" extreme observations!

Multi-level models provide more smoothing/shrinkage to groups with smaller sample sizes (since there is less evidence that their values should be different from "grand mean".)

We'll talk about __*why*__ today…

# Methods of Estimation – How can we systematically construct "good" estimators?

- Several methods have proven useful:

  - *Method of Moments (MoM)*: The $k^{th}$ moment of X is $E[X^{\underline{k}}]$. MoM estimators combine unbiased estimates of moments of X.

  - *Least Squares (LS)*: Obtained by minimizing squared error $\sum_{i=1}^{n}(Y_i - E[Y_i])^2$. Ordinary linear regression!

  - *Maximum likelihood (ML)*: The likelihood is the probability of the data we observed. ML estimators (MLE's) choose parameter values that maximize the likelihood.

  - *Bayesian Estimation (Bayes)*: Treat the parameters as random variables, and use Bayes' rule to pick the parameter value most likely, given the data (the "reverse" of ML!)

# Maximum Likelihood Estimators (MLE's)

- Let $X_1, ..., X_n$ be an iid sample from $f_X(x; \theta)$, $x_1, ..., x_n$ are the observed values

- The *likelihood* of the sample is the joint density

$$L(\theta) \quad = \quad f(x_1, \ldots, x_n; \ \theta) \quad = \quad f(x_1; \ \theta) f(x_2; \ \theta) \cdots f(x_n; \ \theta)$$

$$= \quad \prod_{i=1}^{n} f(x_i; \ \theta)$$

- The *maximum likelihood estimate* $\hat{\theta}_{MLE}$ maximizes L($\theta$):

$$L(\hat{\theta}_{MLE}) \geq L(\theta) \quad \forall \ \theta$$

- Strategy: It's usually (but not always) easier to work with the *log likelihood*

$$LL(\theta) = \log L(\theta) = \sum_{i=1}^{n} \log f(x_i; \ \theta) \ .$$

# Pennsylvania, Pre-Midterm Poll

- John Fetterman (D) running for election to the US Senate against Mehmet Oz ®

- In a Suffolk University Poll (October 27-30):
    - 457 of 500 voters expressed a preference for Fetterman or Oz.
    - Of those 457: 233 prefer Fetterman.

- In most polling, weights are attached to each response, to adjust the "representativeness" of the response for things like
    - who is likely to be home when survey worker calls
    - who refuses to answer
    - etc

- We will ignore weights etc and treat the 457 as a simple random sample.

# Possible models for the data

- 457 individual Bernoulli coin flips, $x_i = 1$ for Fetterman, $x_i = 0$ for Oz

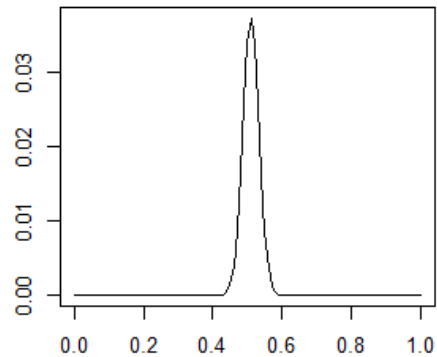$$L_{ber}(p) = \prod_{i=1}^{457} p^{x_i}(1-p)^{1-x_i} = p^{233}(1-p)^{224}$$

- 457 trials, 233 "successes" (Fetterman voters)
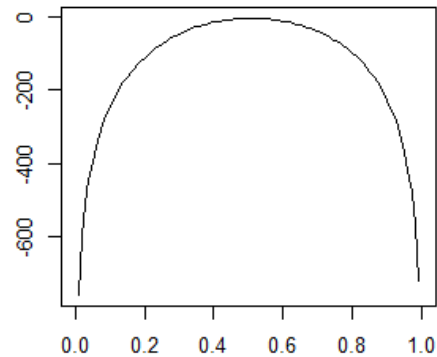
$$L_{bin}(p) = \binom{457}{233} p^{233}(1-p)^{224}$$

- What matters for MLE and SE is *shape*, not *size*!

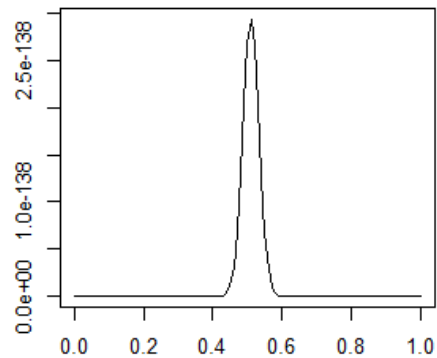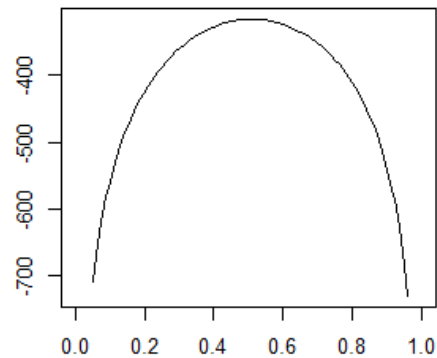# Binomial and Bernoulli Likelihoods

# Finding the MLE…

- If we use the Bernoulli likelihood,

$$
\begin{aligned}
LL_{ber}(p) &= \log L_{ber}(p) \\
= \log p^k(1-p)^{n-k} &= k \log p + (n-k)\log(1-p)
\end{aligned}
$$

- If we use the Binomial likelihood

$$
\begin{aligned}
LL_{bin}(p) &= \log L_{bin}(p) \\
= \log \binom{n}{k} p^k(1-p)^{n-k} &\propto k \log p + (n-k)\log(1-p)
\end{aligned}
$$

- Either way we want to maximize

$$
k \log p + (n-k)\log(1-p)
$$

with k = 233, n=457

# MLE: Point Estimate

■ Differentiating and setting to zero…

$$0 \quad = \quad LL'(p) \quad = \quad \frac{d}{dp}\left[k\log p + (n-k)\log(1-p)\right]$$

$$= \quad \frac{k}{p} - \frac{n-k}{1-p} \quad = \quad \frac{k-pn}{p(1-p)}$$

■ so, clearly,

$$\hat{p} = \frac{k}{n} = \frac{233}{457} = 0.510$$

# Bayes' Rule (a.k.a. Bayes' Theorem)

- A very simple idea with very powerful consequences

- We often start with information like P[A|B] and what we really want is P[B|A].  Bayes' Theorem lets us "turn the conditioning around":

$$\mathbf{P[B|A]} = \frac{P[A\&B]}{P[A]} = \frac{\mathbf{P[A|B]}P[B]}{P[A]}$$

- See https://arbital.com/p/bayes_rule/ for lots of examples and proselytizing.

# Conditional probability & conditional density

- P[A|B] = P[A&B]/P[B]

- P[B] = P[B|A]P[A] + P[B|A$^c$]P[A$^c$]

- P[A & B] = P[B|A]P[A]

- f(x|y) = f(x,y)/f(y)

- $f(y) \quad = \quad \int f(y|x)f(x)dx$

- f(x,y) = f(y|x) f(x)

- Bayes' Theorem:

$$P[B|A] \quad = \quad \frac{P[A\&B]}{P[A]} = \frac{P[A|B]P[B]}{P[A]}$$

$$= \quad \frac{P[A|B]P[B]}{P[A|B]P[B] + P[A|B^c]P[B^c]}$$

- Bayes' Theorem:

$$f(y|x) \quad = \quad \frac{f(x,y)}{f(x)} = \frac{f(x|y)f(y)}{f(x)}$$

$$= \quad \frac{f(x|y)f(y)}{\int f(x|y^*)f(y^*)dy^*}$$

# Bayes' Theorem for Data

- Bayes' Theorem

$$f(y|x) \quad = \quad \frac{f(x,y)}{f(x)} \quad = \quad \frac{f(x|y)f(y)}{f(x)}$$

$$= \quad \frac{f(x|y)f(y)}{\int f(x|y^*)f(y^*)dy^*}$$

Dummy variable of integration

- Let x = data, y = $\theta$ (parameter!); then

$$f(\theta|\text{data}) \quad = \quad \frac{f(\text{data},\theta)}{f(\text{data})} \quad = \quad \frac{f(\text{data}|\theta)f(\theta)}{f(\text{data})}$$

$$= \quad \frac{f(\text{data}|\theta)f(\theta)}{\int f(\text{data}|\theta^*)f(\theta^*)d\theta^*}$$

# Bayes' Theorem for Data

- ■ We call
  - ❑ f($\theta$) the *prior distribution*
  - ❑ f(data|$\theta$) = L($\theta$) the *likelihood*
  - ❑ f($\theta$|data) the *posterior distribution*
- ■ So Bayes' Theorem says

$$f(\theta|\text{data}) \quad = \quad \frac{f(\text{data}|\theta)f(\theta)}{f(\text{data})} \quad \propto \quad f(\text{data}|\theta)f(\theta)$$

- ■ Slogan: (posterior) $\propto$ (likelihood)$\times$(prior)

# Back to 2022 PA pre-midterm poll

- The *likelihood* is the same as before:
$$L(p) \propto p^k (1-p)^{n-k}$$

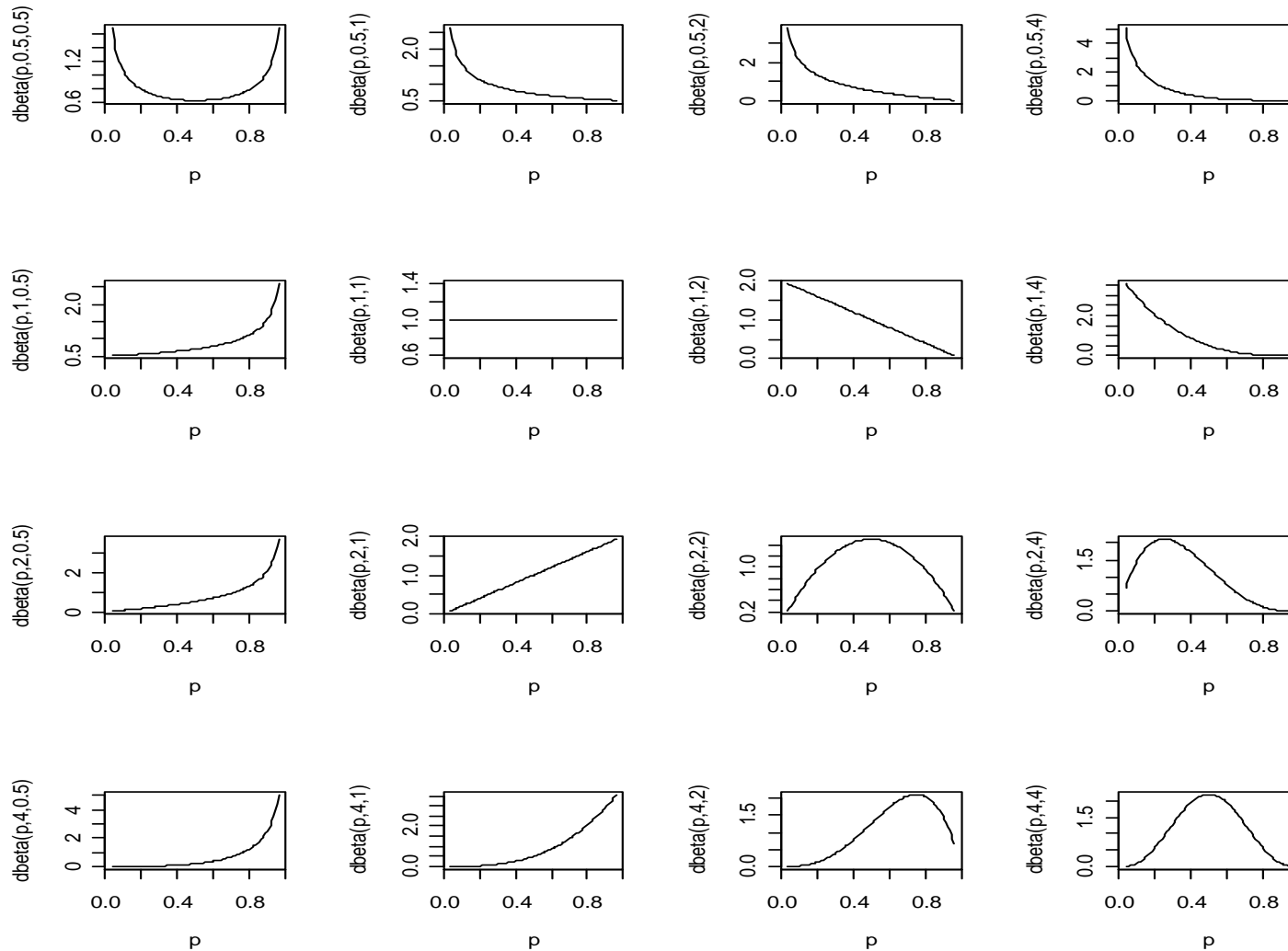- We need a *prior distribution*.  One good choice is a *beta distribution*, with

  ○ Density $\quad f(p|\alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1}(1-p)^{\beta-1}$

  ○ Mean $\quad E[p] = \frac{\alpha}{\alpha+\beta}$

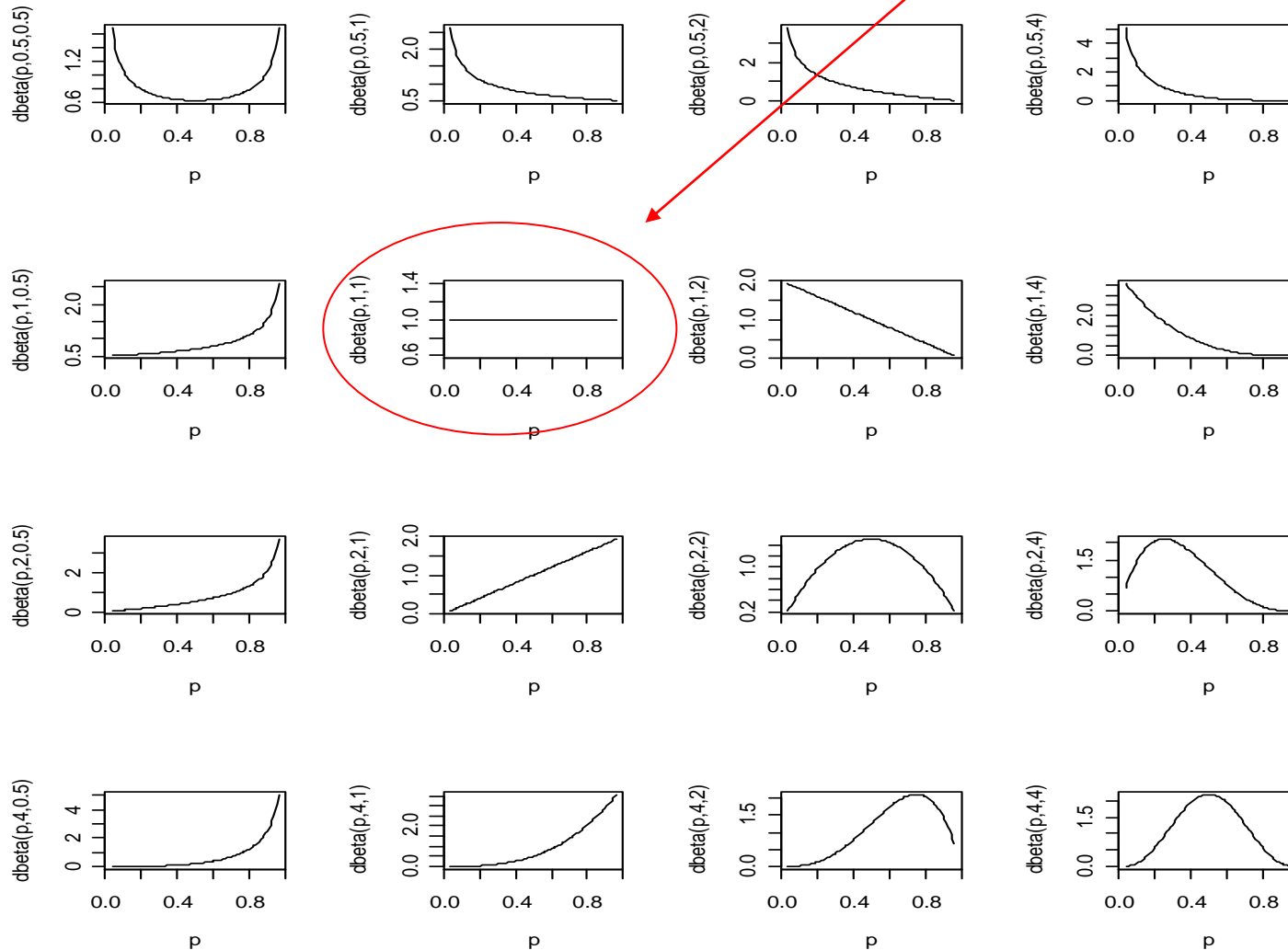  ○ Variance $\quad \text{Var}(p) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$

- Some graphs of beta densities appear on the next slide

# Some Beta Densities

# Some Beta Densities

uniform distribution!

# Choosing prior parameters...

- The *likelihood* is the same as before:

$$L(p) \propto p^k (1-p)^{n-k} = p^{226}(1-p)^{225}$$

- The *prior distribution* is a beta distribution

$$f(p|\alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1}(1-p)^{\beta-1}$$

- $\alpha$ = 1, $\beta$ = 1 gives a uniform distribution – no preference for one p over another!

- Suppose that in a previous poll, 942 prefer Fetterman and 1008 prefer Oz. Could set $\alpha$=942, $\beta$=1008
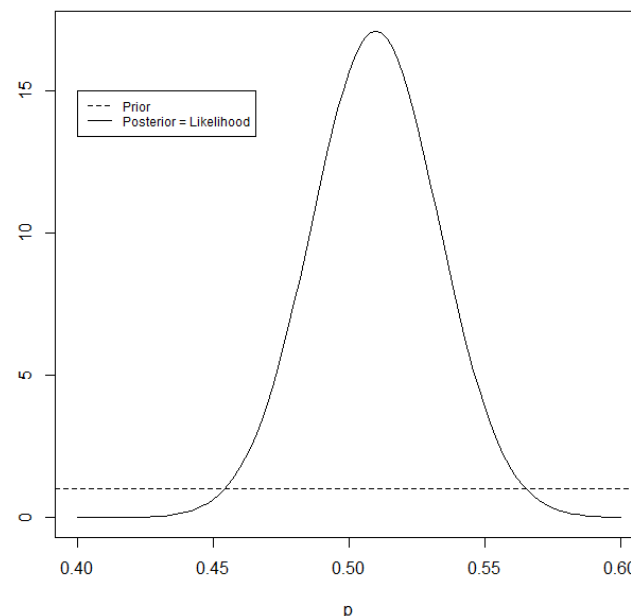
# If $\alpha$=1 and $\beta$=1...

- (posterior) $\propto$ (likelihood)$\times$(prior):

$$f(p|\text{data}) \quad \propto \quad L(p) \times 1 \;=\; p^{233}(1-p)^{224}$$

- Since f(p|data)=L(p),

  posterior mode = MLE

  = 233/457 = 0.5098468

- Since f(p|data) is a beta

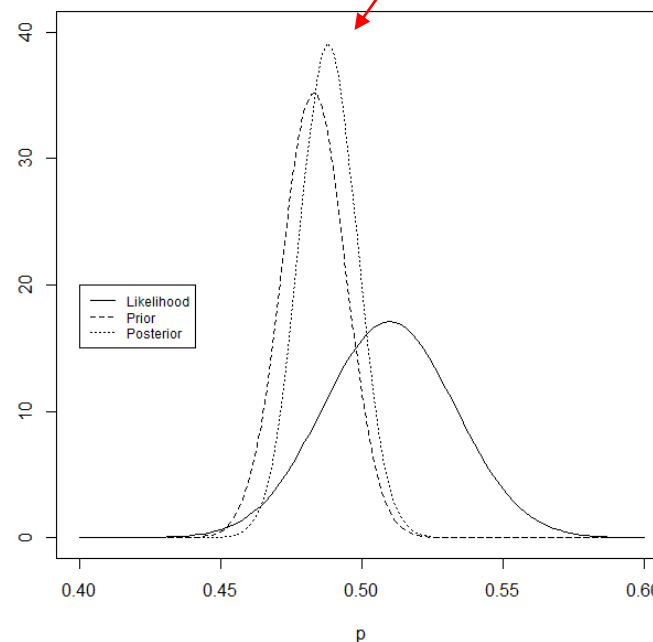  with $\alpha$=234, $\beta$=225,

  E[p|data] = 234/459 = 0.5098039

# If $\alpha$=942, $\beta$=1008…

■ (posterior) $\propto$ (likelihood)$\times$(prior):

$$f(p|\text{data}) \quad \propto \quad L(p) \times p^{941}(1-p)^{1007}$$
$$= \quad p^{1174}(1-p)^{1231}$$

■ Since f(p|data) = beta(p,1175,1232), E[p|data] = 1175/2406 = 0.488 vs MLE=0.510

"shrinkage": posterior between prior & likelihood

# Normal Model: Estimate $\mu$, with $\sigma^2$ Known, One Observation y $\sim$ N($\mu,\sigma^2$)

- For our prior distribution, we'll assume $\mu \sim$ N($\mu_0,\tau_0{}^2$):

$$f(y|\mu) \quad = \quad \frac{1}{\sqrt{2\pi}\sigma}e^{-\frac{1}{2\sigma^2}(y-\mu)^2}$$

$$f(\mu) \quad = \quad \frac{1}{\sqrt{2\pi}\tau_0}e^{-\frac{1}{2\tau_0^2}(\mu-\mu_0)^2}$$

$$f(\mu|y) \quad \propto \quad f(y|\mu)f(\mu) \quad \propto \quad \exp\left\{-\frac{1}{2}\left[\frac{(y-\mu)^2}{\sigma^2}+\frac{(\mu-\mu_0)^2}{\tau_0^2}\right]\right\}$$

- Posterior must be normal for $\mu$ (quadratic in $\mu$!); to identify it, complete the square…

- The exponent of f($\mu$|y) looks like -1/2 times

$$\frac{(y-\mu)^2}{\sigma^2} + \frac{(\mu-\mu_0)^2}{\tau_0^2} = \frac{\tau_0^2+\sigma^2}{\tau_0^2\sigma^2}\left[\mu^2 - \frac{2y\mu\tau_0^2 + 2\mu\mu_0\sigma^2}{\tau_0^2+\sigma^2} + \frac{y^2\tau_0^2 + \mu_0^2\sigma^2}{\tau_0^2+\sigma^2}\right]$$

$$= \frac{\tau_0^2+\sigma^2}{\tau_0^2\sigma^2}\left[\left(\mu - \frac{y\tau_0^2 + \mu_0\sigma^2}{\tau_0^2+\sigma^2}\right)^2 + \mathsf{junk}(y,\sigma^2,\mu_0,\tau_0^2)\right]$$

$$= \frac{1}{\tau_1^2}(\mu-\mu_1)^2 + (\mathsf{known\ junk})$$

so that $\mu$|y $\sim$ N($\mu_1$, $\tau_1{}^2$), where

$$\tau_1^2 = \frac{\tau_0^2\sigma^2}{\tau_0^2+\sigma^2} = \frac{1}{1/\sigma^2 + 1/\tau_0^2}$$

$$\mu_1 = \frac{y\tau_0^2 + \mu_0\sigma^2}{\tau_0^2+\sigma^2} = \left(\frac{\tau_0^2}{\tau_0^2+\sigma^2}\right)y + \left(\frac{\sigma^2}{\tau_0^2+\sigma^2}\right)\mu_0$$

- The exponent of $f(\mu|y)$ looks like -1/2 times

$$
\begin{aligned}
\frac{(y-\mu)^2}{\sigma^2} + \frac{(\mu-\mu_0)^2}{\tau_0^2} &= \frac{\tau_0^2+\sigma^2}{\tau_0^2\sigma^2}\left[\mu^2 - \frac{2y\mu\tau_0^2 + 2\mu\mu_0\sigma^2}{\tau_0^2+\sigma^2} + \frac{y^2\tau_0^2 + \mu_0^2\sigma^2}{\tau_0^2+\sigma^2}\right] \\
&= \frac{\tau_0^2+\sigma^2}{\tau_0^2\sigma^2}\left[\left(\mu - \frac{y\tau_0^2 + \mu_0\sigma^2}{\tau_0^2+\sigma^2}\right)^2 + \mathsf{junk}(y,\sigma^2,\mu_0,\tau_0^2)\right] \\
&= \frac{1}{\tau_1^2}(\mu-\mu_1)^2 + (\mathsf{known\ junk})
\end{aligned}
$$

so that $\mu|y \sim N(\mu_1, \tau_1{}^2)$, where

$$
\begin{aligned}
\tau_1^2 &= \frac{\tau_0^2\sigma^2}{\tau_0^2+\sigma^2} = \frac{1}{1/\sigma^2 + 1/\tau_0^2} \\
\mu_1 &= \frac{y\tau_0^2 + \mu_0\sigma^2}{\tau_0^2+\sigma^2} = \left(\frac{\tau_0^2}{\tau_0^2+\sigma^2}\right)y + \left(\frac{\sigma^2}{\tau_0^2+\sigma^2}\right)\mu_0
\end{aligned}
$$

# n Observations $y_i \sim N(\mu, \sigma^2)$

- **Since**

$$
\begin{aligned}
p(y_1, \ldots, y_n | \mu) &= N(y_1, \ldots, y_n | \mu, \sigma^2) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y_i - \mu)^2} \\
&\propto N(\bar{y} | \mu, \sigma^2/n) \equiv \frac{1}{\sqrt{2\pi\sigma^2/n}} e^{-\frac{1}{2\sigma^2/n}(\bar{y} - \mu)^2}
\end{aligned}
$$

we can apply the results for one observation

- $p(y_1, \ldots, y_n | \mu) \propto N(\bar{y} | \mu, \sigma_n^2),\ \sigma_n^2 = \sigma^2/n$
- $p(\mu) = N(\mu | \mu_o, \tau_o^2)$
- $p(\mu | data) = N(\mu | \mu_n, \tau_n^2)$ where

$$
\begin{aligned}
\tau_n^2 &= \frac{1}{1/\sigma_n^2 + 1/\tau_0^2} = \frac{1}{n/\sigma^2 + 1/\tau_0^2} \\
\mu_n &= \frac{\bar{y}/\sigma_n^2 + \mu_0/\tau_0^2}{1/\sigma_n^2 + 1/\tau_0^2} = \left(\frac{\tau_0^2}{\tau_0^2 + \sigma^2/n}\right)\bar{y} + \left(\frac{\sigma^2/n}{\tau_0^2 + \sigma^2/n}\right)\mu_0
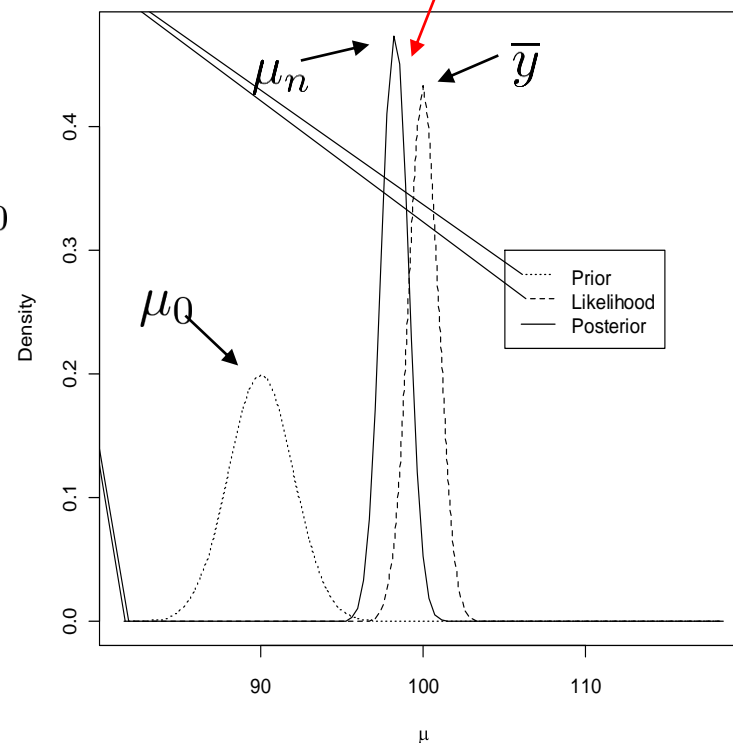\end{aligned}
$$

# Normal Mean, Example

- Suppose we know $\sigma$=12, we look at n=169 IQ scores, and we find $\bar{y}$ = 100.

- We use as prior N($\mu_0$, $\tau_0^2$) with $\mu_0$=90, $\tau_0^2$ = 4

- Shrinkage determined by

$$\mu_n = \left( \frac{\tau_0^2}{\tau_0^2 + \sigma^2/n} \right) \bar{y} + \left( \frac{\sigma^2/n}{\tau_0^2 + \sigma^2/n} \right) \mu_0$$

- $\dfrac{\tau_0^2}{\tau_0^2 + \sigma^2/n}$ is the *reliability*

- n larger $\Rightarrow$

  reliability larger $\Rightarrow$

  less shrinkage

# Minnesota Radon Example

- **Emphasize Distribution Structure**

$$\text{Level 2:} \quad \mu_j \quad \overset{iid}{\sim} \quad N(\mu_0, \tau_0^2)$$

$$\text{Level 1:} \quad y_i \quad \overset{indep}{\sim} \quad N(\mu_{j[i]}, \sigma^2)$$

- **Emphasize Bayesian point of view**

$$\text{Prior:} \quad \mu_j \quad \overset{iid}{\sim} \quad N(\mu_0, \tau_0^2)$$

$$\text{Likelihood:} \quad y_i \quad \overset{indep}{\sim} \quad N(\mu_{j[i]}, \sigma^2)$$

- **Emphasize two-stage (multistage) sampling**

$$\mu_0$$

Mean radon across MN

$$\mu_1 \quad \mu_2 \quad \cdots \quad \mu_J$$

County-level differences from grand mean

$$y_1 y_2 \cdots y_{n_1} \quad y_{n_1+1} \cdots y_{n_2} \quad\quad y_{n_{J-1}+1} \cdots y_n$$

individual house levels

# Minnesota Radon Example

$$\mu_0$$

Mean radon across MN

$$\mu_1 \qquad \mu_2 \qquad \cdots \qquad \mu_I$$

County-level differences
from grand mean

$$y_{11} y_{12} \cdots y_{1n_1} \quad y_{21} y_{22} \cdots y_{2n_2} \qquad y_{I1} y_{I2} \cdots y_{In_I}$$

individual house levels

- In each county i with $n_i$ houses, the posterior mean radon level $E[\mu_i | y_{i1}, \ldots, y_{in_i}]$ will be

$$\mu_i^{post} = \left( \frac{\tau_0^2}{\tau_0^2 + \sigma^2/n_i} \right) \bar{y}_i + \left( \frac{\sigma^2/n_i}{\tau_0^2 + \sigma^2/n_i} \right) \mu_0$$

  - When $n_i$ large, $\mu_i^{post} \approx \bar{y}_i$
  - When $n_i$ small, $\mu_i^{post} \approx \mu_o$

# Minnesota Radon Example

- **In the figure, the grand mean is $\mu_0$**

- **In each county i with $n_i$ houses, posterior mean is**

$$\mu_i^{post} = \left(\frac{\tau_0^2}{\tau_0^2 + \sigma^2/n_i}\right)\bar{y}_i$$

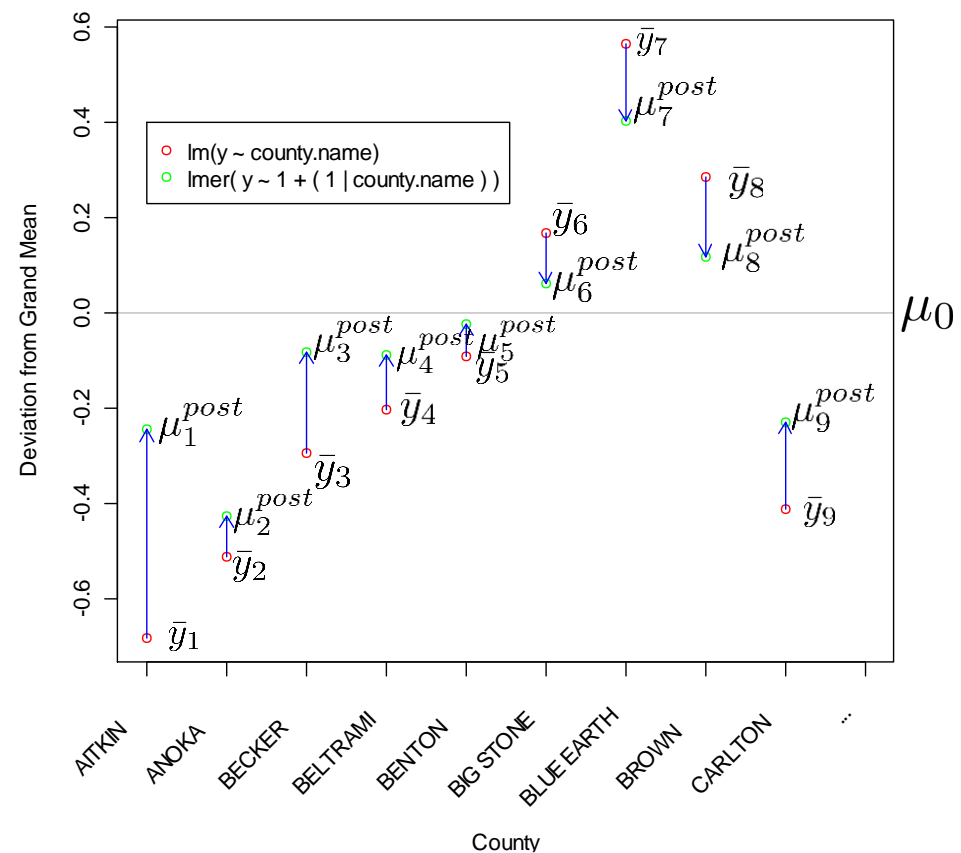$$+ \left(\frac{\sigma^2/n_i}{\tau_0^2 + \sigma^2/n_i}\right)\mu_0$$

- ☐ When $n_i$ large, $\mu_i^{post} \approx \bar{y}_i$
- ☐ When $n_i$ small, $\mu_i^{post} \approx \mu_0$

# MLM's and Shrinkage

- The random effect "estimates" that `lmer` produces with `ranef()` are a form of posterior means E[η|data] for each η.

- The posterior means E[η|data] are always shrunk toward the prior mean 0, so that the random effects α are always shrunk toward the corresponding fixed effects β.

- The Bayesian pov not only provides insights, but also

  ❑ Novel ways to expand the multi-level model framework

  ❑ Simulation-based methods of estimation (MCMC with `jags`, Hamiltonian MC with `stan`, etc.)

*(we will take a look at some of this next week!)*

# Summary

- Today:
  - Shrinkage
  - Review of MLE
  - Crash course in Bayes
  - Normal-Normal Model & Shrinkage
  - MLM's and Shrinkage
- Project discussion
- After Thanksgiving:
  - A little practical Bayes / MCMC for multi-level models