

Understanding the James Stein Estimator, as used in Brown 2008. (1)

Model 1

$$\begin{cases} X_i | \theta_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\theta_i, 1) \\ \theta_i \stackrel{\text{ind}}{\sim} \mathcal{N}(0, \tau^2) \end{cases} \quad \theta_i, \tau^2 \text{ unknown}$$

Goal Data $\{X_i\}$. Estimate θ_i via estimator $\hat{\theta}_i = \delta(X_i)$.

Naïve Estimator: MLE

$$\hat{\theta}_i^{(\text{MLE})} = \delta^0(X_i) = X_i$$

Posterior Def $P(\theta_i | X_i) \propto P(X_i | \theta_i) P(\theta_i)$ Bayes' Rule

~~Definition~~

~~Definition~~

$$= \mathcal{N}(X_i | \theta_i, 1) \cdot \mathcal{N}(\theta_i | 0, \tau^2)$$

$$\propto \exp\left(-\frac{1}{2}(X_i - \theta_i)^2\right) \cdot \exp\left(-\frac{\theta_i^2}{2\tau^2}\right)$$

$$\propto \exp\left(-\frac{1}{2}\left(\theta_i^2\left(1 + \frac{1}{\tau^2}\right) - 2\theta_i X_i\right)\right)$$

$$= \exp\left[-\frac{1}{2}\left(\frac{\tau^2 + 1}{\tau^2}\right)\left(\theta_i^2 - \frac{\tau^2}{\tau^2 + 1} 2\theta_i X_i\right)\right]$$

$$\propto \exp\left(-\frac{1}{2\lambda}(\theta_i - \lambda X_i)^2\right)$$

letting $\lambda = \frac{\tau^2}{\tau^2 + 1}$.

Posterior $\theta_i | X_i \sim \mathcal{N}(\lambda X_i, \lambda)$ where $\lambda = \frac{\tau^2}{\tau^2 + 1}$

Bayes Estimator: Posterior Mean

$$\hat{\theta}_i^{(\text{Bayes})} = \delta^*(X_i) = \mathbb{E}(\theta_i | X_i) = \lambda X_i$$

$\hat{\theta}_i^{(\text{Bayes})}$ shrinks $\hat{\theta}_i^{(\text{MLE})}$ towards 0.

Problem τ^2 , and hence λ , is unknown, so we can't use $\hat{\theta}_i^{(\text{Bayes})}$ directly. Need Empirical Bayes

$\lambda = \frac{\tau^2}{1+\tau^2}$ [scribbled out]

Goal Use Empirical Bayes estimator $\hat{\lambda}$, with $\hat{\theta}_i^{(EB)} = \hat{\lambda} X_i$
This is "empirical bayes" because we estimate the hyperparameter λ from the data $\{X_i\}$.

Marginal $X_i \sim \mathcal{N}(0, 1+\tau^2)$

Sum of Squares $S^2 = \|X\|_2^2 = \sum_{i=1}^p X_i^2 \stackrel{d}{=} (\tau^2+1) \chi_p^2$

Lemma $E\left(\frac{1}{\chi_p^2}\right) = \frac{1}{p-2}$ Pf next page

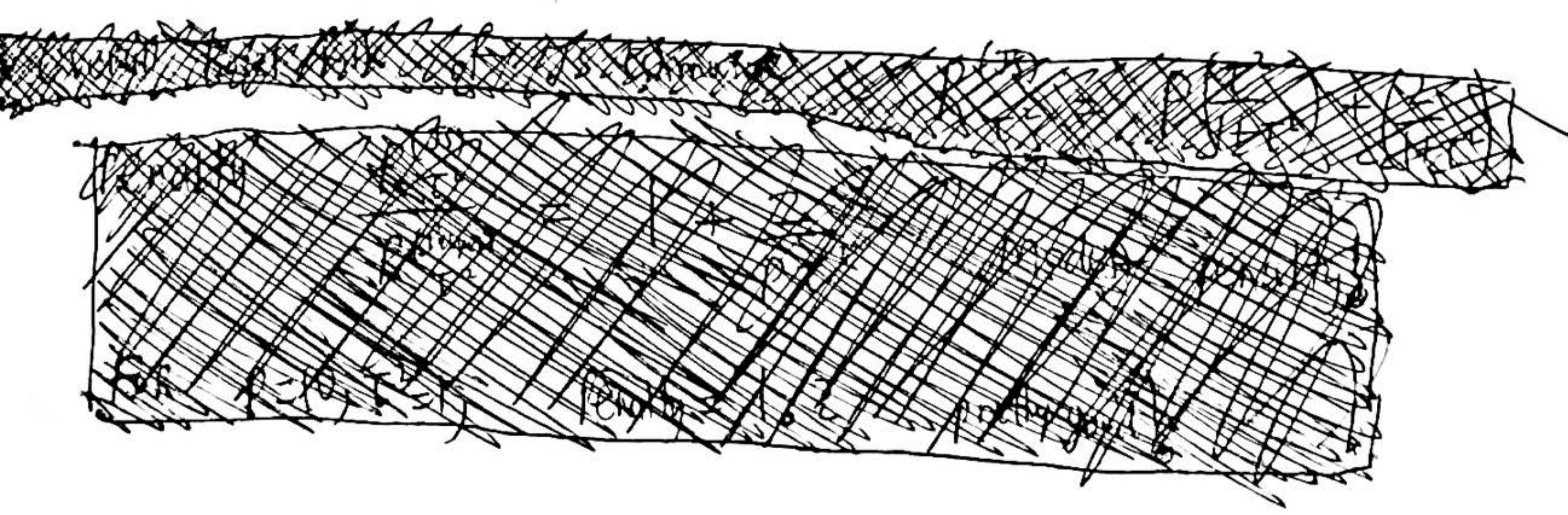
Conclusion $(\tau^2+1) E\left(\frac{1}{S^2}\right) = E\left(\frac{1}{\chi_p^2}\right) = \frac{1}{p-2}$

$\Rightarrow E\left(\frac{p-2}{S^2}\right) = \frac{1}{\tau^2+1} = 1-\lambda$

\Rightarrow Estimator $\hat{\lambda} = 1 - \frac{p-2}{S^2} = 1 - \frac{p-2}{\|X\|_2^2}$

\Rightarrow James Stein Estimator $\hat{\theta}_i^{(JS)} = \hat{\lambda} X_i$

$\hat{\theta}_i^{(JS)} = \left(1 - \frac{p-2}{\|X\|_2^2}\right) X_i$



Pf (lemma)

$$V = \frac{1}{Y}, Y \sim \chi_p^2. \quad F_V(v) = P(V \leq v) = P\left(\frac{1}{Y} \leq v\right) = P\left(Y \geq \frac{1}{v}\right) = 1 - F_Y\left(\frac{1}{v}\right)$$

$$f_V(v) = \frac{d}{dv} F_V(v) = \frac{1}{v^2} f_Y\left(\frac{1}{v}\right).$$

$X = Z^2, Z \sim N(0,1)$, has χ_1^2 distribution, which for $x > 0$ has CDF

$$F_X(x) = P(X \leq x) = P(Z^2 \leq x) = P(|Z| \leq \sqrt{x}) = P(-\sqrt{x} \leq Z \leq \sqrt{x})$$

$$= \int_{-\sqrt{x}}^{\sqrt{x}} \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt = 2 \int_0^{\sqrt{x}} \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt.$$

Then $X \sim \chi_1^2$ has density

$$f_X(x) = \frac{d}{dx} F_X(x) = \frac{2}{\sqrt{2\pi}} e^{-\frac{x}{2}} \cdot \frac{1}{2\sqrt{x}} = \frac{1}{2^{1/2} \Gamma(1/2)} x^{-1/2} e^{-x/2}, \quad x > 0$$

since $\Gamma(1/2) = \sqrt{\pi}$.

Hence $\chi_1^2 \sim \text{Gamma}\left(\frac{1}{2}, \frac{1}{2}\right)$.

MGF of Gamma(α, β) is $M_G(t) = \mathbb{E}e^{Gt} = \left(1 - \frac{t}{\beta}\right)^{-\alpha}$

Hence $M_{\chi_1^2}(t) = (1 - 2t)^{-1/2}$

Hence $M_{\chi_p^2}(t) = \prod_{k=1}^p M_{\chi_1^2}(t) = \prod_{k=1}^p (1 - 2t)^{-1/2} = (1 - 2t)^{-p/2} \Rightarrow \chi_p^2 \sim \text{Gamma}\left(\frac{p}{2}, \frac{1}{2}\right)$

Hence χ_p^2 has density $f_{\chi_p^2}(y) = \frac{(1/2)^{p/2}}{\Gamma(p/2)} y^{\frac{p}{2}-1} e^{-\frac{y}{2}}$.

So, a Random Variable $X \sim \frac{1}{\chi_p^2}$ has density $f(x) = \frac{(1/2)^{p/2}}{\Gamma(p/2)} x^{-\frac{p}{2}-1} e^{-\frac{1}{2x}}, x > 0$.

It is a proper distribution, so $\int_0^{\infty} x^{-\frac{p}{2}-1} e^{-\frac{1}{2x}} dx = \Gamma(p/2) 2^{p/2}$

$$\text{Hence } \mathbb{E}X = \int_0^{\infty} x f(x) dx = \int_0^{\infty} x \frac{(1/2)^{p/2}}{\Gamma(p/2)} x^{-\frac{p}{2}-1} e^{-\frac{1}{2x}} dx$$

$$= \frac{(1/2)^{p/2}}{\Gamma(p/2)} \int_0^{\infty} x^{\frac{(p-2)}{2}-1} e^{-\frac{1}{2x}} dx = \frac{(1/2)^{p/2}}{\Gamma(p/2)} \cdot \Gamma\left(\frac{p-2}{2}\right) 2^{\frac{p-2}{2}}$$

$$= \frac{1}{2} \frac{\Gamma\left(\frac{p}{2}-1\right)}{\Gamma\left(\frac{p}{2}\right)} = \frac{1}{2} \frac{\Gamma\left(\frac{p}{2}-1\right)}{\left(\frac{p}{2}-1\right) \Gamma\left(\frac{p}{2}-1\right)} \quad \text{using } \Gamma(s) = (s-1)\Gamma(s-1).$$

$$= \frac{1}{2} \cdot \frac{2}{p-2} = \frac{1}{p-2} \quad \square$$

Shocking Thm For $p \geq 3$, the James Stein estimator $\hat{\theta}^{(JS)}$ everywhere dominates the MLE $\hat{\theta}^{(MLE)}$ in terms of expected total squared error. That is, for all choices of θ ,

$$\mathbb{E}_{\theta} \|\hat{\theta}^{(JS)} - \theta\|_2^2 < \mathbb{E}_{\theta} \|\hat{\theta}^{(MLE)} - \theta\|_2^2$$

Note This is a frequentist, not a Bayesian, result. $\hat{\theta}^{(JS)}$ is superior no matter what one's prior beliefs about θ may be.

PF Notice $(\hat{\theta}_i - \theta_i)^2 = (X_i - \hat{\theta}_i)^2 - (X_i - \theta_i)^2 + 2(\hat{\theta}_i - \theta_i)(X_i - \theta_i)$.

Sum over $i=1, \dots, p$ and take expectations:

$$\mathbb{E}_{\theta} \|\theta - \hat{\theta}\|_2^2 = \mathbb{E}_{\theta} \|X - \hat{\theta}\|_2^2 - p + 2 \sum_{i=1}^p \text{cov}_{\theta}(\hat{\theta}_i, X_i)$$

cov_{θ} indicates covariance under

$$X \sim \mathcal{N}_p(\theta, I_p)$$

Lemma $\text{cov}_{\theta}(\hat{\theta}_i, X_i) = \mathbb{E}_{\theta} \left(\frac{\partial \hat{\theta}_i}{\partial X_i} \right)$

(see Stein's Lemma) (Just integration by parts)
on next page

Thus $\text{cov}_{\theta}(\hat{\theta}_i, X_i) = \mathbb{E}_{\theta} \frac{\partial}{\partial X_i} \left[\left(1 - \frac{p-2}{\sum_{k=1}^p X_k^2}\right) X_i \right]$

$$= \mathbb{E}_{\theta} \left[\left(1 - \frac{p-2}{\|X\|_2^2}\right) + X_i \left(\frac{(p-2) 2 X_i}{\|X\|_2^4} \right) \right]$$

$$= 1 - (p-2) \mathbb{E} \left(\frac{1}{\|X\|_2^2} \right) + 2(p-2) \mathbb{E} \left(\frac{X_i^2}{\|X\|_2^4} \right)$$

$$\begin{aligned}
 \text{Thus } 2 \sum_{i=1}^p \text{cov}_{\theta}(\hat{\theta}_i, X_i) &= 2 \left[-p(p-2) \mathbb{E} \left(\frac{1}{\|X\|_2^2} \right) + 2(p-2) \mathbb{E} \left(\frac{\|X\|_2^2}{\|X\|_2^4} \right) \right] \\
 &= 2p - 2(p-2)^2 \mathbb{E} \left(\frac{1}{\|X\|_2^2} \right) \\
 &= 2p - 2 \mathbb{E}_{\theta} \left[\frac{(p-2)^2}{\|X\|_2^2} \right]
 \end{aligned}$$

$$\begin{aligned}
 \text{and } \mathbb{E}_{\theta} \|x - \hat{\theta}^{(JS)}\|_2^2 &= \sum_{i=1}^p \mathbb{E}_{\theta} (X_i - \hat{\theta}_i^{(JS)})^2 \\
 &= \sum_{i=1}^p \mathbb{E}_{\theta} (p-2)^2 \frac{X_i^2}{\|X\|_2^4} = \mathbb{E}_{\theta} \left(\frac{(p-2)^2}{\|X\|_2^2} \right)
 \end{aligned}$$

$$\text{Therefore } \mathbb{E}_{\theta} \|\theta - \hat{\theta}^{(JS)}\|_2^2 = \mathbb{E}_{\theta} \left(\frac{(p-2)^2}{\|X\|_2^2} \right) - p + 2p - 2 \mathbb{E}_{\theta} \left(\frac{(p-2)^2}{\|X\|_2^2} \right)$$

$$= p - \mathbb{E}_{\theta} \left(\frac{(p-2)^2}{\|X\|_2^2} \right)$$

$$\leftarrow p = \sum_{i=1}^p 1 = \sum_{i=1}^p \mathbb{E}_{\theta} (X_i - \theta_i)^2 = \mathbb{E}_{\theta} \|\theta - \hat{\theta}^{(MLE)}\|_2^2$$

assuming $p \geq 3$. \square

Stein's lemma Let $X \sim N(\theta, \sigma^2)$, g differentiable function, $\mathbb{E}|g'(x)| < \infty$.

Then
$$\mathbb{E}[g(X)(X-\theta)] = \sigma^2 \mathbb{E}g'(X)$$

Pf

$$\mathbb{E}[g(X)(X-\theta)] = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} g(x)(x-\theta) e^{-\frac{(x-\theta)^2}{2\sigma^2}} dx.$$

$$\left[\begin{array}{l} u = g(x) \quad dv = (x-\theta) e^{-\frac{(x-\theta)^2}{2\sigma^2}} \\ du = g'(x)dx \quad v = -\sigma^2 e^{-\frac{(x-\theta)^2}{2\sigma^2}} \end{array} \right]$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \left[-\sigma^2 g(x) e^{-\frac{(x-\theta)^2}{2\sigma^2}} \right]_{-\infty}^{\infty} + \sigma^2 \int_{-\infty}^{\infty} g'(x) \frac{e^{-\frac{(x-\theta)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma} dx$$

$= 0$ since $\mathbb{E}|g'(x)| < \infty$

$$= \sigma^2 \mathbb{E}g'(X) \quad \square$$

(5)

Model 2
$$\begin{cases} X_i | \theta_i & \text{ind} \sim \mathcal{N}(\theta_i, \sigma_i^2) \\ \theta_i & \text{ind} \sim \mathcal{N}(0, \tau^2) \end{cases} \quad \begin{array}{l} \sigma_i^2 \text{ known} \\ i=1, \dots, p \end{array}$$

Data Transformation
$$\begin{cases} \tilde{X}_i = \frac{X_i}{\sigma_i} \\ \tilde{\theta}_i = \theta_i / \sigma_i \end{cases} \quad \text{so that} \quad \begin{cases} \tilde{X}_i | \tilde{\theta}_i & \text{ind} \sim \mathcal{N}(\tilde{\theta}_i, 1) \\ \tilde{\theta}_i & \text{ind} \sim \mathcal{N}(0, \tau^2 / \sigma_i^2) \end{cases}$$

We know the James Stein estimator in this situation!

$$\hat{\theta}_i^{(JS)} = \left(1 - \frac{p-2}{\|\tilde{X}\|_2^2}\right) \tilde{X}_i = \left(1 - \frac{p-2}{\sum_{k=1}^p X_k^2 / \sigma_k^2}\right) \frac{X_i}{\sigma_i}$$

Define
$$\hat{\theta}_i^{(JS)} = \sigma_i \hat{\tilde{\theta}}_i^{(JS)}$$

Then
$$\hat{\theta}_i^{(JS)} = \left(1 - \frac{p-2}{\sum_{k=1}^p X_k^2 / \sigma_k^2}\right) X_i$$

this is the Battling Avg model used in Brown (2008)

Model 3

$$\begin{cases} X_i | \theta_i & \text{ind} \sim \mathcal{N}(\theta_i, \sigma_i^2) \\ \theta_i & \text{ind} \sim \mathcal{N}(\mu, \tau^2) \end{cases}$$

σ_i^2 known
 $i=1, \dots, p$

Data Transformation

$$\begin{cases} \tilde{X}_i = X_i - \mu \\ \tilde{\theta}_i = \theta_i - \mu \end{cases} \quad \text{so that} \quad \begin{cases} \tilde{X}_i | \tilde{\theta}_i & \text{ind} \sim \mathcal{N}(\tilde{\theta}_i, \sigma_i^2) \\ \tilde{\theta}_i & \text{ind} \sim \mathcal{N}(0, \tau^2) \end{cases}$$

We know the James Stein estimator in this situation!

$$\hat{\theta}_i^{(JS)} = \left(1 - \frac{p-2}{\sum_{k=1}^p \tilde{X}_k^2 / \sigma_k^2} \right) \tilde{X}_i = \left(1 - \frac{p-2}{\sum_{k=1}^p (X_k - \mu)^2 / \sigma_k^2} \right) (X_i - \mu)$$

Define $\hat{\theta}_i^{(JS)} = \mu + \hat{\theta}_i^{(JS)}$

Then

$$\hat{\theta}_i^{(JS)} = \mu + \left(1 - \frac{p-2}{\sum_{k=1}^p (X_k - \mu)^2 / \sigma_k^2} \right) (X_i - \mu)$$

Note $\hat{\theta}_i^{(JS)}$ shrinks the MLE X_i towards the mean μ .

Empirical Bayes Because μ is unknown, we estimate μ from the data as $\hat{\mu}$, and use $\hat{\mu}$ in place of μ in the above $\hat{\theta}_i^{(JS)}$

Model $\begin{cases} X_i | \theta_i \stackrel{\text{ind}}{\sim} N(\theta_i, \sigma_i^2) \\ \theta_i \stackrel{\text{ind}}{\sim} N(\mu, \tau^2) \end{cases}$ σ_i^2 known

MARGINAL $X_i \stackrel{\text{ind}}{\sim} N(\mu, \tau^2 + \sigma_i^2)$

Goal Obtain estimate $\hat{\mu}$ of μ to use in $\hat{\theta}_i^{(JS)}$ [empirical bayes!]

MLE $\hat{\mu}_{MLE}$
 i^{th} likelihood $P(X_i) = \frac{1}{\sqrt{2\pi(\tau^2 + \sigma_i^2)}} e^{-\frac{(X_i - \mu)^2}{2(\tau^2 + \sigma_i^2)}}$

Full log-likelihood $\ell(\mathbf{X} | \mu, \tau^2, \sigma^2) = \sum_{i=1}^p \log P(X_i) = -\frac{1}{2} \sum_{i=1}^p \frac{(X_i - \mu)^2}{\tau^2 + \sigma_i^2} + C(\tau^2, \sigma^2)$

$\frac{\partial \ell}{\partial \mu} = \sum_{i=1}^p \frac{X_i - \mu}{\tau^2 + \sigma_i^2}$

~~MLE~~ $\frac{\partial \ell}{\partial \mu}(\hat{\mu}) = 0 \implies \sum_{i=1}^p \frac{X_i}{\tau^2 + \sigma_i^2} = \hat{\mu} \sum_{i=1}^p \frac{1}{\tau^2 + \sigma_i^2}$

$\implies \hat{\mu}_{MLE} = \frac{\sum X_i / (\tau^2 + \sigma_i^2)}{\sum 1 / (\tau^2 + \sigma_i^2)}$

However, τ^2 is unknown, so we must use an estimate $\hat{\tau}^2$ for τ^2 .

Brown simply uses ~~MLE~~ $\tau^2 = 0$, to make life simple,

Yielding the estimate

$\hat{\mu}_1 = \frac{\sum X_i / \sigma_i^2}{\sum 1 / \sigma_i^2}$

This corresponds to the model $\begin{cases} X_i | \theta_i \sim N(\theta_i, \sigma_i^2) \\ \theta_i \equiv \mu \end{cases}$

which is to assume each batter has a common batting average mean.