

# Empirical Bayes, James-Stein

## Outline

- 1) Empirical Bayes
- 2) James-Stein Paradox
- 3) Stein's Lemma
- 4) Stein's unbiased risk estimator (SURE)

# Estimators for Gaussian seq. model

## Gaussian sequence model

$$X \sim N_d(\theta, I_d)$$

Goal: estimate  $\theta \in \mathbb{R}^d$  via  $\delta(x)$

with low  $MSE(\theta; \delta) = \mathbb{E}_\theta \|\theta - \delta(x)\|^2$

Model more general than it might appear:

Ex.  $X_1, \dots, X_n \stackrel{iid}{\sim} (\theta, \sigma^2 I_d)$

$$\Rightarrow Z = \frac{1}{\sigma\sqrt{n}} \sum_{k=1}^n X_k \approx N_d(\theta, I_d)$$

## Estimators

$\delta_0(x) = X$  has much to recommend it

- UMVU
- MLE
- Objective Bayes (flat or Jeffreys prior)

# Linear shrinkage estimator

$$\delta_{\zeta}(x) = (1 - \zeta)X$$

Arises from Bayes:

Simple Bayes:

$$\theta_i \stackrel{\text{iid}}{\sim} N(0, \tau^2)$$

$$\Rightarrow \text{use } \zeta = \frac{1}{1 + \tau^2}$$

Hierarchical Bayes:

$$\tau^2 \sim \lambda_0$$

$$\theta_i | \tau^2 \stackrel{\text{iid}}{\sim} N(0, \tau^2)$$

$$\Rightarrow \text{use } \zeta = \mathbb{E}\left[\frac{1}{1 + \tau^2} \mid x\right] = \hat{\zeta}_{\text{Bayes}}(x)$$

Empirical Bayes:

Estimate  $\zeta$  (est.  $\tau^2$ )

$$X | \tau^2 \stackrel{\text{iid}}{\sim} N(0, 1 + \tau^2)$$

(Suff.)

$$\rightsquigarrow \|X\|^2 \sim (1 + \tau^2) \chi_d^2$$

$$\hat{\zeta}_{\text{MLE}}(x) = d / \|X\|^2$$

$$\hat{\zeta}_{\text{UMVU}}(x) = \frac{d - 2}{\|X\|^2}$$

$$\left( \text{since } \mathbb{E}\left[\frac{1}{Y}\right] = \frac{1}{d-2} \text{ for } Y \sim \chi_d^2 \right)$$

Lemma If  $Y \sim \chi_d^2$ ,  $\mathbb{E} \frac{1}{Y} = \frac{1}{d-2}$

Proof:

# Empirical Bayes

Common situation in hierarchical Bayes models:

$\xi \sim \lambda(\xi)$   $\leftarrow$  one draw  $\Rightarrow$  hard to justify prior  
lots of info  $\Rightarrow$  prior doesn't matter

$\theta_i | \xi \stackrel{\text{iid}}{\sim} \pi_\xi(\theta)$   $\leftarrow$  only  $X_i$  informative  $\Rightarrow$  prior helps  
many draws  $\Rightarrow$  can check fit

$X_i | \xi, \theta \stackrel{\text{ind.}}{\sim} p_{\theta_i}(x)$   $i = 1, \dots, d$

Hybrid approach: treat  $\xi$  as fixed

- Estimate  $\xi$  based on observed data
- Plug in  $\xi$  as though known

Ex.  $\theta_i \sim N(0, \tau^2)$   $\tau^2$  fixed, unknown  
 $X_i | \theta \sim N(\theta_i, 1)$   $i = 1, \dots, d$

Bayes estimator if we knew  $\tau^2$  is

$$\delta_i(X) = (1 - \xi) X_i, \quad \xi = \frac{1}{1 + \tau^2}$$

To estimate  $\xi$ , use  $X \sim N_d(0, \xi^{-1} I_d) = \left(\frac{\xi}{2\pi}\right)^{d/2} e^{-\xi \|X\|^2/2}$

$$\|X\|^2 \sim \xi^{-1} \chi_d^2 \Rightarrow \hat{\xi}_{MLE}^{-1} = d / \|X\|^2$$

Plug in:  $\delta_i(X) = (1 - d / \|X\|^2) X_i$

If  $d$  large, should be near-optimal

# James - Stein Estimator

James & Stein proposed instead ( $d \geq 3$ ):

$$\hat{\sigma}_{SS, i}(X) = \left(1 - \frac{d-2}{\|X\|^2}\right) x_i$$

Emp Bayes Motivation:  $\frac{d-2}{\|X\|^2}$  is UMVUE of  $\sigma^2$

Prop: If  $Y \sim \chi_d^2 = \text{Gamma}\left(\frac{d}{2}, 2\right)$ ,  $d \geq 3$  then

$$\mathbb{E}\left[\frac{1}{Y}\right] = \frac{1}{d-2}$$

Proof: 
$$\mathbb{E}\left[\frac{1}{Y}\right] = \int_0^\infty \frac{1}{y} \frac{1}{2^{d/2} \Gamma(d/2)} \cdot y^{d/2-1} e^{-y/2} dy$$

$$= \frac{2^{(d-2)/2} \Gamma(d/2)}{2^{d/2} \Gamma(d/2)} \int_0^\infty \underbrace{\frac{1}{2^{(d-2)/2} \Gamma(d/2)} y^{(d-2)/2-1} e^{-y/2}}_{\chi_{d-2}^2 \text{ density}} dy$$

Now, use  $\Gamma'(x) = (x-1)\Gamma(x-1) \quad \forall x > 0$

$$\dots = \frac{1}{2} \cdot \frac{1}{(d-2)/2} = \frac{1}{d-2}$$

□

$$\begin{aligned} \zeta \|X\|^2 &\sim \chi_d^2 \Rightarrow \zeta^{-1} \mathbb{E}_\zeta \left[ \frac{1}{\|X\|^2} \right] = \frac{1}{d-2} \\ &\Rightarrow \hat{\sigma} = \frac{d-2}{\|X\|^2} \quad \text{UMVUE} \end{aligned}$$

# James - Stein Paradox

Back to non-Bayesian Gaussian seq. model:

$$X_i \stackrel{\text{iid}}{\sim} N_d(\theta, \sigma^2 I_d), \quad \theta \in \mathbb{R}^d \text{ (fixed)}, \quad \sigma^2 > 0 \text{ known}$$

$i = 1, \dots, n$

Shocking result of James & Stein (1956):

For  $d \geq 3$ , the sample mean  $\bar{X} = \frac{1}{n} \sum X_i$  is inadmissible as an estimator of  $\theta$

under squared error loss:

$$\text{For } \delta_{\text{JS}}(X) = \left(1 - \frac{(d-2)\sigma^2/n}{\|\bar{X}\|^2}\right) \bar{X}$$

$$\text{MSE}(\theta, \delta_{\text{JS}}) < \text{MSE}(\theta, \bar{X}) \quad \forall \theta \in \mathbb{R}^d \text{ (!!!)}$$

$\bar{X}$  is UMVU, Minimax, objective Bayes, ....

Note: Might as well take  $n=1$  (Suff. reduction)  $\Rightarrow \left(1 - \frac{d-2}{\|X\|^2}\right)X$

Note this result holds without assumption of

Bayes model on  $\theta$ : true for  $\theta = (500, -10^6, 4)$

Nothing special about  $\theta$ : for any  $\theta_0 \in \mathbb{R}^d$

$$\delta(X) = \theta_0 + \left(1 - \frac{d-2}{\|X - \theta_0\|^2}\right)(X - \theta_0)$$

also dominates  $X$

Deep implication: shrinkage makes sense even without Bayes justification.

## Linear shrinkage w/o Bayesian assumptions

Gaussian seq. model:  $X \sim N_d(\theta, I_d)$ , fixed  $\theta \in \mathbb{R}^d$

Let  $\delta_\zeta(X) = (1-\zeta)X$ ,  $\zeta$  is tuning parameter

$$\begin{aligned} R(\theta; \delta_\zeta) &= \|\theta - \mathbb{E}\delta_\zeta(X)\|^2 + \sum_i \text{Var}((1-\zeta)X_i) \\ \uparrow \\ \text{(MSE)} & \\ &= \underbrace{\zeta^2 \|\theta\|^2}_{\text{bias}^2} + \underbrace{d(1-\zeta)^2}_{\text{variance}} \end{aligned}$$

What is optimal  $\zeta$ ?

$$\frac{d}{d\zeta} R(\theta; \delta_\zeta) = 2\zeta \|\theta\|^2 - 2(1-\zeta)d$$

$$\Rightarrow \text{minimizer} = \zeta^*(\theta) = \frac{d}{d + \|\theta\|^2} = \frac{1}{1 + \|\theta\|^2/d}$$

$\zeta^*$  always  $> 0$ , but  $\rightarrow 0$  as  $\theta \rightarrow \infty$

What if we estimate  $\zeta^*(\theta)$ ?

How does adaptivity of  $\hat{\zeta}^*(X)$  affect MSE?



# Stein's Lemma

Useful tool for computing / estimating risk in Gaussian estimation problems

Theorem (Stein's Lemma, univariate):

Suppose  $X \sim N(\theta, \sigma^2)$

$h(x): \mathbb{R} \rightarrow \mathbb{R}$  differentiable,  $\mathbb{E} |h'(x)| < \infty$

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

$\stackrel{=}{\text{Cov}}(X, h(X))$

Proof Note we can assume wlog  $h(0) = 0$  (why?)

First assume  $\theta = 0, \sigma^2 = 1$ :

$$\text{Note } \mathbb{E}[Xh(X)] = \int_0^{\infty} x h(x) \phi(x) dx + \int_{-\infty}^0 x h(x) \phi(x) dx$$

$$\begin{aligned} \int_0^{\infty} x h(x) \phi(x) dx &= \int_0^{\infty} x \left[ \int_0^x h(y) dy \right] \phi(x) dx \\ &= \int_0^{\infty} \int_0^{\infty} \mathbb{1}\{y < x\} x h(y) \phi(x) dx dy \\ &= \int_0^{\infty} h(y) \left[ \int_y^{\infty} x \phi(x) dx \right] dy \\ &= \int_0^{\infty} h(y) \phi(y) dy \end{aligned}$$

In the last step we have used:

$$\frac{d}{dx} \left[ \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \right] = -x \cdot \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

Similar argument shows  $\int_{-\infty}^0 x h(x) \phi(x) dx = \int_{-\infty}^0 h(x) \phi(x) dx$

$\Rightarrow$  Result holds for  $\theta = 0, \sigma^2 = 1$

General  $\theta, \sigma^2$ :

write  $X = \theta + \sigma Z, \quad Z \sim N(0, 1)$

$$\begin{aligned} \mathbb{E}[(x-\theta)h(x)] &= \sigma \mathbb{E}[Z h(\theta + \sigma Z)] \\ &= \sigma^2 \mathbb{E}[h'(\theta + \sigma Z)] \\ &= \sigma^2 \mathbb{E}[h'(x)] \end{aligned}$$

# Multivariate Stein's Lemma

Def  $h: \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,  $Dh \in \mathbb{R}^{d \times d}$

$$(Dh(x))_{ij} = \frac{\partial h_i}{\partial x_j}(x)$$

Def (Frobenius norm):  $A \in \mathbb{R}^{d \times d}$

$$\|A\|_F = \left( \sum_{i,j} A_{ij}^2 \right)^{1/2}$$

Theorem (Stein's Lemma, Multivariate):

$$X \sim N_d(\theta, \sigma^2 I_d) \quad \theta \in \mathbb{R}^d$$

$h: \mathbb{R}^d \rightarrow \mathbb{R}^d$  diff'able,  $\mathbb{E} \|Dh(x)\|_F < \infty$

$$\begin{aligned} \text{Then } \mathbb{E} \left[ (x - \theta)' h(x) \right] &= \sigma^2 \mathbb{E} \text{tr}(Dh(x)) \\ &= \sigma^2 \sum_i \mathbb{E} \frac{\partial h_i}{\partial x_i}(x) \end{aligned}$$

Proof

$$\begin{aligned} \mathbb{E} \left[ (x_i - \theta_i) h_i(x) \right] &= \mathbb{E} \left[ \mathbb{E} \left[ (x_i - \theta_i) h_i(x) \mid x_{-i} \right] \right] \\ &= \mathbb{E} \left[ \mathbb{E} \left[ \sigma^2 \frac{\partial h_i}{\partial x_i}(x) \mid x_i \right] \right] \\ &= \sigma^2 \mathbb{E} \frac{\partial h_i}{\partial x_i}(x) \quad \square \end{aligned}$$

# Stein's Unbiased Risk Estimator (SURE)

Can use Stein's Lemma to get unbiased estimator of the MSE of any  $\delta(x)$ :  
apply Stein's Lemma with  $h(x) = X - \delta(x)$

Assume  $\sigma^2 = 1$ :

$$\begin{aligned} R(\theta; \delta) &= \mathbb{E}_{\theta} \left[ \|X - \theta - h(x)\|^2 \right] \\ &= \mathbb{E}_{\theta} \|X - \theta\|^2 + \mathbb{E}_{\theta} \|h(x)\|^2 - 2 \mathbb{E}_{\theta} [(X - \theta)' h(x)] \\ &= d + \mathbb{E}_{\theta} \|h(x)\|^2 - 2 \mathbb{E}_{\theta} \text{tr}(Dh(x)) \end{aligned}$$

$$\Rightarrow \hat{R}(x) = d + \|h(x)\|^2 - 2 \text{tr}(Dh(x))$$

is unbiased for the MSE (estimator b/c only dep. on  $x$ )

Can also compute MSE via  $R = \mathbb{E}_{\theta} \hat{R}$

Ex:  $\delta(x) = X \Rightarrow h(x) = 0, Dh'(x) = 0$   
 $\hat{R} = d = R(\theta; \delta) \quad \forall \theta$

Ex:  $\delta_{\zeta}(x) = (1 - \zeta)x$  for fixed  $\zeta$

$$\Rightarrow h(x) = \zeta x, \quad Dh = \zeta I_d$$

$$\hat{R} = d + \zeta^2 \|x\|^2 - 2\zeta d = (1 - 2\zeta)d + \zeta^2 \|x\|^2$$

# Risk of James-Stein

$$\delta^{JS}(x) = \left(1 - \frac{d-2}{\|x\|^2}\right) X$$

$$\Rightarrow h(x) = \frac{d-2}{\|x\|^2} X$$

$$\|h(x)\|^2 = \frac{(d-2)^2}{\|x\|^2}$$

$$\frac{\partial h_i}{\partial x_i}(x) = \frac{\partial}{\partial x_i} \frac{(d-2)x_i}{\sum_j x_j^2}$$

$$= (d-2) \frac{\|x\|^2 - 2x_i^2}{\|x\|^4}$$

$$\Rightarrow \text{tr}(Dh(x)) = \frac{d-2}{\|x\|^4} \sum_i (\|x\|^2 - 2x_i^2)$$

$$= \frac{(d-2)^2}{\|x\|^2}$$

$$\hat{R} = d + \frac{(d-2)^2}{\|x\|^2} - 2 \frac{(d-2)^2}{\|x\|^2}$$

$$= d - \frac{(d-2)^2}{\|x\|^2}$$

$$R(\theta; \delta_{JS}) = d - \overbrace{(d-2)^2 \mathbb{E}_\theta \left[ \frac{1}{\|x\|^2} \right]}^{> 0}$$

$$< d$$

$$= R(\theta; X)$$

$$\text{If } \theta = 0 \text{ then } \mathbb{E}_\theta \left[ \frac{1}{\|X\|^2} \right] = d-2$$

$$\Rightarrow R(\theta; \delta_{JS}) = d - (d-2) = 2$$

Possibly  $\ll d$  !

$$\theta \rightarrow \infty \text{ then } \mathbb{E}_\theta \left[ \frac{1}{\|X\|^2} \right] \approx \frac{1}{\|\theta\|^2}$$

$$\Rightarrow R(\theta; \delta_{JS}) \approx d - \frac{(d-2)^2}{\|\theta\|^2} \\ \rightarrow d$$

Smaller and smaller advantage but always better.

Note  $\delta_{JS}(X)$  also inadmissible:

$$\delta_{JS+}(X) = \left( 1 - \frac{d-2}{\|X\|^2} \right)_+ X \text{ is strictly better}$$

Practically more useful version:

$$\delta_{JS,2}(X) = \bar{X} + \left( 1 - \frac{d-3}{\|X - \bar{X} \mathbf{1}_d\|^2} \right) (X - \bar{X} \mathbf{1}_d)$$

Dominates  $\delta(X) = X$  for  $d \geq 4$

Taken to logical extreme, suggestion seems dumb:  
should everyone @ Berkeley pool their estimates?

Note  $\mathbb{E} \|\cdot\|^2$  is improved, but  $\mathbb{E}(X_i - \theta_i)^2$  may get worse for individual coordinates.