

Contents

| | | |
|----------|---|-----------|
| 1 | Shrinkage estimators | 1 |
| 2 | Admissible linear shrinkage estimators | 3 |
| 3 | Admissibility of unbiased normal mean estimators | 6 |
| 4 | Motivating the James-Stein estimator | 11 |
| 4.1 | What is wrong with \mathbf{X} ? | 11 |
| 4.2 | An oracle estimator: | 12 |
| 4.3 | Adaptive shrinkage estimation | 13 |
| 5 | Risk of δ_{JS} | 16 |
| 5.1 | Risk bound for δ_{JS} | 16 |
| 5.2 | Stein's identity | 17 |
| 6 | Some oracle inequalities | 21 |
| 6.1 | A simple oracle inequality | 21 |
| 7 | Unknown variance or covariance | 23 |

Much of this content comes from Lehmann and Casella [1998], sections 5.2, 5.4, 5.5, 4.6 and 4.7.

1 Shrinkage estimators

Consider a model $\{p(x|\theta) : \theta \in \Theta\}$ for a random variable X such that

$$E[X|\theta] = \mu(\theta), \quad 0 < \text{Var}[X|\theta] = \sigma^2(\theta) < \infty \quad \forall \theta \in \Theta.$$

A linear estimator $\delta(x)$ for $\mu(\theta)$ is an estimator of the form

$$\delta_{ab}(X) = aX + b.$$

Is δ_{ab} admissible?

Theorem 1 (LC thm 5.2.6). $\delta_{ab}(X) = aX + b$ is inadmissible for $E[X|\theta]$ under squared error loss whenever

1. $a > 1$,
2. $a = 1$ and $b \neq 0$, or
3. $a < 0$.

Proof. The risk of δ_{ab} is

$$\begin{aligned} R(\theta, \delta_{ab}) &= E[(aX + b - \mu)^2 | \theta] \\ &= E[(aX - a\mu - \mu(1 - a) + b)^2 | \theta] \\ &= E[a^2(X - \mu)^2 + (b - \mu(1 - a))^2 + 2a(X - \mu)(b - \mu(1 - a)) | \theta] \\ &= a^2\sigma^2 + (b - \mu(1 - a))^2 \end{aligned}$$

1. If $a > 1$, then $R(\theta, \delta_{ab}) > a^2\sigma^2 > \sigma^2 = R(\theta, X)$, so δ_{ab} is dominated by X .
2. If $a < 0$, then

$$\begin{aligned} R(\theta, \delta_{ab}) &> (b - \mu(1 - a))^2 \\ &= (1 - a)^2(b/(1 - a) - \mu)^2 \\ &= (b/(1 - a) - \mu)^2 = R(\theta, b/(1 - a)) \end{aligned}$$

and so δ_{ab} is dominated by the constant estimator $b/(1 - a)$.

3. If $a = 1$ and $b \neq 0$, then $R(\theta, \delta_{ab}) = \sigma^2 + b^2 > \sigma^2 = R(\theta, X)$, so δ_{ab} is dominated by X .

□

Letting $w = 1 - a$ and $\mu_0 = b/(1 - a)$, the result suggests that if we want to use an admissible linear estimator, it should be of the form

$$\delta(X) = w\mu_0 + (1 - w)X, \quad w \in [0, 1]$$

We call such estimators *linear shrinkage estimators* as they “shrink” the estimate from X towards μ_0 . Intuitively, you can think of μ_0 as your “guess” as to the value of μ , and w as the confidence you have in your guess. Of course, the closer your guess is to the truth, the better your estimator.

If μ_0 represents your guess as to $\mu(\theta)$, it seems natural to require that $\mu_0 \in \mu(\Theta) = \{\mu : \mu = \mu(\theta), \theta \in \Theta\}$, i.e. μ_0 is a possible value of μ .

Lemma 1. *If $\mu(\Theta)$ is convex and $\mu_0 \notin \bar{\mu}(\Theta)$, then $\delta(X) = w\mu_0 + (1 - w)X$ is not admissible.*

Proof.

For the one-dimensional case, suppose $\mu_0 > \mu(\theta) \forall \theta \in \Theta$.

Let $\tilde{\mu}_0 = \sup_{\Theta} \mu(\theta)$, and $\tilde{\delta}(X) = w\tilde{\mu}_0 + (1 - w)X$.

Then $\tilde{\delta}(X)$ dominates $\delta(X)$

(the variances are the same, and the latter has higher bias for all θ).

The proof is similar for the case $\mu_0 < \mu(\theta) \forall \theta \in \Theta$.

□

Exercise: Generalize this result to higher dimensions.

2 Admissible linear shrinkage estimators

We have shown that $\delta(X) = w\mu_0 + (1 - w)X$ is inadmissible for $\mu(\theta) = E[X|\theta]$ if

- $w \notin [0, 1]$ or

- $\mu_0 \notin \mu(\Theta)$.

Restricting attention to $w \in [0, 1]$ and $\mu_0 \in \mu(\Theta)$, it may seem that such estimators should always be admissible, but “always” is almost always too inclusive.

Exercise: Given an example where $w\mu_0 + (1 - w)X$ is not admissible, even with $w \in (0, 1)$ and $\mu_0 \in \mu(\Theta)$.

Linear shrinkage via conjugate priors

What about using a Bayesian argument? Recall,

Theorem. *Any unique Bayes estimator is admissible.*

If we can show that $w\mu_0 + (1 - w)X$ is unique Bayes under some prior, then we will have shown admissibility.

Let $X_1, \dots, X_n \sim$ i.i.d. $p(x|\theta)$, where

$$p(x|\theta) \in \mathcal{P} = \{p(x|\theta) = h(x) \exp(x \cdot \theta - A(\theta)) : \theta \in \mathcal{H}\}$$

Consider estimation of $\mu = E[X|\theta]$ under squared error loss.

Let $\pi(\theta) \propto \exp(n_0\mu_0 \cdot \theta - n_0A(\theta))$ where $n_0 > 0$ and $\mu_0 \in \text{Conv}\{E[X|\theta] : \theta \in \mathcal{H}\}$

Recall that under this prior,

$$E[\mu] \equiv E[E[X|\theta]] = \mu_0.$$

Then $\pi(\theta|x) \propto \exp(n_1\mu_1 \cdot \theta - n_1A(\theta))$, where $n_1 = n_0 + n$ and

$$\begin{aligned} n_1\mu_1 &= n_0\mu_0 + n\bar{x} \\ \mu_1 &= n_0\mu_0/n_1 + n\bar{x}/n_1 \\ &= \frac{n_0}{n_0 + n_1}\mu_0 + \frac{n}{n_0 + n}\bar{x}. \end{aligned}$$

Under this posterior distribution,

$$E[\mu|x] \equiv E[E[X|\theta]|x] = \mu_1.$$

Therefore, the unique Bayes estimator of $\mu = E[X|\theta]$ under squared error loss is

$$\mu_1 = w\mu_0 + (1 - w)\bar{x},$$

and so this linear shrinkage estimator is admissible.

Example (multiple normal means):

Let $\mathbf{X} \sim N_p(\boldsymbol{\theta}, \sigma^2 \mathbf{I})$. First consider the case that σ^2 is known, so that

$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\theta}) &= (2\pi\sigma^2)^{-p/2} \exp(-(\mathbf{x} - \boldsymbol{\theta}) \cdot (\mathbf{x} - \boldsymbol{\theta})/[2\sigma^2]) \\ &\propto_{\boldsymbol{\theta}} \exp(\mathbf{x} \cdot \boldsymbol{\theta}/\sigma^2 - \boldsymbol{\theta} \cdot \boldsymbol{\theta}/[2\sigma^2]). \end{aligned}$$

Consider the normal prior

$$\begin{aligned} \pi(\boldsymbol{\theta}) &= (2\pi\tau^2)^{-p/2} \exp(-(\boldsymbol{\theta} - \boldsymbol{\theta}_0) \cdot (\boldsymbol{\theta} - \boldsymbol{\theta}_0)/[2\tau_0^2]) \\ &\propto_{\boldsymbol{\theta}} \exp(\boldsymbol{\theta}_0 \cdot \boldsymbol{\theta}/\tau_0^2 - \boldsymbol{\theta} \cdot \boldsymbol{\theta}/[2\tau_0^2]). \end{aligned}$$

where τ_0^2 is analogous to $1/n_0$ in the general formulation for exponential families.

The posterior density is

$$\begin{aligned} \pi(\boldsymbol{\theta}|\mathbf{x}) &\propto_{\boldsymbol{\theta}} \exp\{[\boldsymbol{\theta}_0/\tau_0^2 + \mathbf{x}/\sigma^2] \cdot \boldsymbol{\theta} - \boldsymbol{\theta} \cdot \boldsymbol{\theta}/[1/\sigma^2 + 1/\tau_0^2]/2\} \\ &= \exp\{\boldsymbol{\theta}_1 \cdot \boldsymbol{\theta}/\tau_1^2 - \boldsymbol{\theta} \cdot \boldsymbol{\theta}/[2\tau_1^2]\} \end{aligned}$$

where

- $1/\tau_1^2 = 1/\tau_0^2 + 1/\sigma^2$
- $\boldsymbol{\theta}_1 = \frac{1/\tau_0^2}{1/\tau_0^2 + 1/\sigma^2} \boldsymbol{\theta}_0 + \frac{1/\sigma^2}{1/\tau_0^2 + 1/\sigma^2} \mathbf{x} \equiv w\boldsymbol{\theta}_0 + (1 - w)\mathbf{x}$.

So $\{\boldsymbol{\theta}|\mathbf{x}\} \sim N_p(\boldsymbol{\theta}_1, \tau_1^2 \mathbf{I})$, which means that

$$E[\boldsymbol{\theta}|\mathbf{x}] = \boldsymbol{\theta}_1 = w\boldsymbol{\theta}_0 + (1 - w)\mathbf{x}$$

uniquely minimizes the posterior risk under squared error loss. The posterior mean is therefore a unique Bayes estimator and also an admissible estimator of $\boldsymbol{\theta}$. Since this result holds for all $\tau_0^2 > 0$, we have the following:

Lemma 2. For each $w \in (0, 1)$ and $\boldsymbol{\theta}_0 \in \mathbb{R}^p$, the estimator $\delta_{w\boldsymbol{\theta}_0}(\mathbf{x}) = w\boldsymbol{\theta}_0 + (1-w)\mathbf{x}$ is admissible for estimating $\boldsymbol{\theta}$ in the model $\mathbf{X} \sim N_p(\boldsymbol{\theta}, \sigma^2 I)$, $\boldsymbol{\theta} \in \mathbb{R}^p$, where σ^2 is known.

Of course, what we would like is the following lemma:

Lemma 3. For each $w \in (0, 1)$ and $\boldsymbol{\theta}_0 \in \mathbb{R}^p$, the estimator $\delta_{w\boldsymbol{\theta}_0}(\mathbf{X}) = w\boldsymbol{\theta}_0 + (1-w)\mathbf{X}$ is admissible for estimating $\boldsymbol{\theta}$ in the model $\mathbf{X} \sim N_p(\boldsymbol{\theta}, \sigma^2 I)$, $\boldsymbol{\theta} \in \mathbb{R}^p$, $\sigma^2 \in \mathbb{R}^+$.

How can this result be obtained?

Theorem 2. Let $\mathcal{P} = \{p(x|\theta, \psi) : (\theta, \psi) \in \Theta \times \Psi\}$, and for $\psi_0 \in \Psi$, let $\mathcal{P}_{\psi_0} = \{p(x|\theta, \psi_0) : \theta \in \Theta\}$ be a submodel. If δ is admissible for estimating θ under \mathcal{P}_{ψ_0} for each $\psi_0 \in \Psi$, then δ is admissible for estimating θ under \mathcal{P} .

Proof. Suppose δ satisfies the conditions of the theorem but is not admissible.

Then there exists a $\delta' \in \mathcal{D}$ such that

$$\begin{aligned} \forall(\theta, \psi), R((\theta, \psi), \delta') &\leq R((\theta, \psi), \delta) \\ \exists(\theta_0, \psi_0), R((\theta_0, \psi_0), \delta') &< R((\theta_0, \psi_0), \delta). \end{aligned}$$

But this contradicts the assumption that δ is admissible for estimating θ under \mathcal{P}_{ψ_0} . Therefore, no such δ' can exist and so δ is admissible for \mathcal{P} . \square

A corollary to this theorem is the admissibility of $w\boldsymbol{\theta}_0 + (1-w)\mathbf{X}$ in the normal model with unknown variance.

3 Admissibility of unbiased normal mean estimators

Let $\mathbf{X} \sim N_p(\boldsymbol{\theta}, \sigma^2 \mathbf{I})$, $\boldsymbol{\theta} \in \mathbb{R}^p$, $\sigma^2 > 0$.

For estimation of $\boldsymbol{\theta}$ under squared error loss, we have shown that the linear shrinkage estimator $\delta(\mathbf{x}) = w\boldsymbol{\theta}_0 + (1-w)\mathbf{x}$ is

- inadmissible if $w \notin [0, 1]$,
- admissible if $w \in (0, 1)$.

What remains to evaluate is the admissibility for $w \in \{0, 1\}$. Admissibility for $w = 1$ is easy to show - the estimator $\delta_{1\theta_0}(\mathbf{x}) = \boldsymbol{\theta}_0$ beats everything at $\boldsymbol{\theta}_0$ and so can't be dominated. The last and most interesting case is that of $w = 0$, i.e. $\delta_0(\mathbf{X}) = \mathbf{X}$, the unbiased MLE and UMVUE.

Blyth's method

Recall Blyth's method for showing admissibility using a limiting Bayes argument:

Theorem 3 (LC 5.7.13). *Suppose $\Theta \subset \mathbb{R}^p$ is open, and that $R(\theta, \delta)$ is continuous in θ for all $\delta \in \mathcal{D}$. Let δ be an estimator and $\{\pi_n\}$ be a sequence of measures such that for any open ball $B \subset \Theta$,*

$$\frac{R(\pi_n, \delta) - R(\pi_n, \delta_{\pi_n})}{\pi_n(B)} \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Then δ is admissible.

Let's try to use this to show admissibility of $\delta_0(\mathbf{X}) = \mathbf{X}$ in the normal means problem. We begin with the case that $\sigma^2 = 1$ is known.

$$\begin{aligned} \mathbf{X} &\sim N_p(\boldsymbol{\theta}, \mathbf{I}) \\ \boldsymbol{\theta} &\sim N_p(\mathbf{0}, \tau^2 \mathbf{I}) \\ \{\boldsymbol{\theta} | \mathbf{X}\} &\sim N_p\left(\frac{\tau^2}{1+\tau^2} \mathbf{X}, \frac{\tau^2}{1+\tau^2} \mathbf{I}\right) \end{aligned}$$

The unique Bayes estimator is $\delta_{\tau^2} = E[\boldsymbol{\theta} | \mathbf{X}] = \frac{\tau^2}{1+\tau^2} \mathbf{X} \equiv (1 - w) \mathbf{X}$.

To apply the theorem, we need to compute the Bayes risk of δ_{τ^2} and \mathbf{X} under the $N_p(\mathbf{0}, \tau^2 \mathbf{I})$ prior π_{τ^2} . The loss we will use is "scaled" squared error loss, $L(\boldsymbol{\theta}, \mathbf{d}) = \sum (\theta_j - d_j)^2 / p$. Because the risk is the average of the individual MSEs, the Bayes

risk is just the average of the Bayes risks from the p components,

$$\begin{aligned} R(\boldsymbol{\theta}, \delta) &= \mathbb{E}\left[\sum_{j=1}^p (\theta_j - \delta_j)^2\right]/p \\ &= \sum_{j=1}^p \mathbb{E}[(\theta_j - \delta_j)^2]/p, \end{aligned}$$

and so calculating the Bayes risk is similar to calculating the risk in the $p = 1$ problem. For $\delta_{\tau^2}(x) = ax$, where $a = 1 - w = \tau^2/(1 + \tau^2)$, we have

$$\begin{aligned} \mathbb{E}[(aX - \theta)^2] &= \mathbb{E}[(aX - a\theta + (1 - a)\theta)^2] \\ &= a^2\mathbb{E}[(X - \theta)^2] + 2a(1 - a)\mathbb{E}[(X - \theta)\theta] + (1 - a)^2\mathbb{E}[\theta]^2 \\ &= a^2 + (1 - a)^2\tau^2 \\ &= \left(\frac{\tau^2}{1 + \tau^2}\right)^2 + \frac{\tau^2}{(1 + \tau^2)^2} = \frac{\tau^2}{1 + \tau^2} \end{aligned}$$

A more intuitive way to calculate this makes use of the fact that $\delta_{\tau^2}(X) = \mathbb{E}[\theta|X]$, so

$$\begin{aligned} \mathbb{E}[(\theta - \delta_{\tau^2})^2] &= \mathbb{E}[(\theta - \mathbb{E}[\theta|X])^2] \\ &= \mathbb{E}_x[\mathbb{E}_{\theta|x}[(\theta - \mathbb{E}[\theta|X])^2]] \\ &= \mathbb{E}_x[\text{Var}[\theta|X]] \\ &= \mathbb{E}_x[\tau^2/(1 + \tau^2)] = \frac{\tau^2}{1 + \tau^2}. \end{aligned}$$

Similarly,

$$\mathbb{E}[(X - \theta)^2] = \mathbb{E}_\theta[\mathbb{E}_x[\mathbb{E}_{\theta|x}[(X - \theta)^2]]] = \mathbb{E}_\theta[1] = 1.$$

So $R(\pi_{\tau^2}, \delta_{\tau^2}) = \frac{\tau^2}{1 + \tau^2}$ and $R(\pi_{\tau^2}, X) = 1$. Returning to the p -variate case, since the Bayes risk is the arithmetic average of the risks for each of the p components of $\boldsymbol{\theta}$, we have

$$\begin{aligned} R(\pi_{\tau^2}, \delta_{\tau^2}) &= \frac{\tau^2}{1 + \tau^2} \\ R(\pi_{\tau^2}, \mathbf{X}) &= 1. \end{aligned}$$

Note that

$$\begin{aligned}\delta_{\tau^2}(\mathbf{X}) &= \frac{\tau^2}{1+\tau^2} \mathbf{X} \uparrow \mathbf{X} \quad \text{as } \tau^2 \uparrow \infty, \\ R(\pi_{\tau^2}, \delta_{\tau^2}) &\uparrow R(\pi_{\tau^2}, \mathbf{X}) \quad \text{as } \tau^2 \uparrow \infty,\end{aligned}$$

so \mathbf{X} is a “limiting Bayes” estimator, for which the risk difference from the Bayes estimator converges to zero. This is promising - let’s now apply the theorem. Letting B be any open finite ball in \mathbb{R}^p , we need to see if the following limit is zero:

$$\begin{aligned}\lim_{\tau^2 \rightarrow \infty} \frac{R(\pi_{\tau^2}, \mathbf{X}) - R(\pi_{\tau^2}, \delta_{\tau^2})}{\pi_{\tau^2}(B)} &= \lim_{\tau^2 \rightarrow \infty} \frac{(1 - \frac{\tau^2}{1+\tau^2})}{\pi_{\tau^2}(B)} \\ &= \lim_{\tau^2 \rightarrow \infty} [(1 + \tau^2)\pi_{\tau^2}(B)]^{-1}\end{aligned}$$

Now $\pi_{\tau^2}(B) \rightarrow 0$ as $\tau^2 \rightarrow \infty$ for any bounded set B . Therefore, the limit is zero only if

$$\lim_{\tau^2 \rightarrow \infty} \tau^2 \pi_{\tau^2}(B) = \infty.$$

We have

$$\begin{aligned}\tau^2 \pi_{\tau^2}(B) &= \tau^2 \int_B (2\pi\tau^2)^{-p/2} \exp(-\|\boldsymbol{\theta}\|^2/[2\tau^2]) \, d\boldsymbol{\theta} \\ &= (2\pi)^{-p/2} \times (\tau^2)^{1-p/2} \times \int_B \exp(-\|\boldsymbol{\theta}\|^2/[2\tau^2]) \, d\boldsymbol{\theta} \\ &= (2\pi)^{-p/2} \times (a) \times (b).\end{aligned}$$

Now take the limit as $\tau^2 \rightarrow \infty$:

$$\begin{aligned}(b) &\rightarrow \text{Vol}(B) \quad \text{as } \tau^2 \rightarrow \infty \\ (a) &\rightarrow \begin{cases} \infty & \text{if } p = 1 \\ 1 & \text{if } p = 2 \\ 0 & \text{if } p > 2 \end{cases}.\end{aligned}$$

Therefore, the desired limit is achieved for $p = 1$ but not $p > 1$.

- By the theorem, X is admissible for Θ .

- For $p > 1$, this method of showing admissibility does not work.
 - For $p = 2$, \mathbf{X} can be shown to be admissible using Blyth's method with non-normal priors (see LC exercise 5.4.5).
 - For $p > 2$, \mathbf{X} can't be shown to be admissible because it isn't.

Interpreting the failure of Blyth's method:

The admissibility conditions for Blyth's method derived from consideration of the existence of an estimator δ that dominates \mathbf{X} . If such an estimator exists, then by continuity of risks

$$\exists \epsilon > 0 \text{ and an open ball } B \subset \Theta : R(\theta, \mathbf{X}) - R(\theta, \delta) > \epsilon \forall \theta \in B,$$

which implies for each prior π_k that

$$\begin{aligned} R(\pi_k, \mathbf{X}) - R(\pi_k, \delta) &= \int [R(\theta, \mathbf{X}) - R(\theta, \delta)] \pi_k(d\theta) \geq \int_B [R(\theta, \mathbf{X}) - R(\theta, \delta)] \pi_k(d\theta) \\ &\geq \epsilon \pi_k(B). \end{aligned}$$

Integrating with respect to a prior π_k , and comparing to the Bayes risk of the Bayes estimator δ_k under π_k gives

$$R(\pi_k, \mathbf{X}) - R(\pi_k, \delta_k) \geq R(\pi_k, \mathbf{X}) - R(\pi_k, \delta) > \epsilon \pi_k(B) \forall k,$$

as δ_k has Bayes risk less than or equal to that of δ . Could such a δ exist?

Could exist: Suppose B is a ball such that $\pi_k(B)$ goes to zero very fast. Then an estimator (like \mathbf{X}) can have a good limiting Bayes risk and still do poorly on B . This allows for the possibility of domination by another estimator that does better on B .

Couldn't exist: On the other hand, if $R(\pi_k, \mathbf{X}) - R(\pi_k, \delta_k)$ goes to zero very fast (e.g. faster than the probability of any ball B), then in a sense \mathbf{X} would have to be doing well everywhere, and would not be able to be dominated - this is Blyth's method for showing admissibility.

What fails in the admissibility proof for the normal means problem is that for $p > 2$, the probability $\pi_k(B)$ of an open ball B is going to zero much faster than the Bayes risk difference, leaving a large enough "gap" for some other estimator to do better.

4 Motivating the James-Stein estimator

Stein [1956] showed that \mathbf{X} is inadmissible for $\boldsymbol{\theta}$ in the normal means problem when $p > 2$. This was surprising, as \mathbf{X} is the MLE and UMVUE for $\boldsymbol{\theta}$. In this section,

4.1 What is wrong with \mathbf{X} ?

For large p ,

- \mathbf{X} may be close to $\boldsymbol{\theta}$, but
- $\mathbf{X} \cdot \mathbf{X} = \|\mathbf{X}\|^2$ may be far from $\boldsymbol{\theta} \cdot \boldsymbol{\theta} = \|\boldsymbol{\theta}\|^2$.

If $\mathbf{X} \sim N_p(\boldsymbol{\theta}, \mathbf{I})$,

$$\begin{aligned} \mathbb{E}[\|\mathbf{X}\|^2] &= \mathbb{E}\left[\sum_1^p X_j^2\right] \\ &= \sum_1^p (\theta_j^2 + 1) = \|\boldsymbol{\theta}\|^2 + p, \end{aligned}$$

so for large p , the magnitude of the estimator vector \mathbf{X} is expected to be much larger than the magnitude of the estimand vector $\boldsymbol{\theta}$. More insight can be gained as follows: Note that every vector \mathbf{x} can be expressed as

$$\mathbf{x} = s\boldsymbol{\theta} + \mathbf{r}, \text{ for some } s \in \mathbb{R} \text{ and } \mathbf{r} : \boldsymbol{\theta} \cdot \mathbf{r} = 0.$$

Here, the random variable s is the magnitude of the projection of \mathbf{x} in the direction of $\boldsymbol{\theta}$, and \mathbf{r} is the residual vector. Using this decomposition, we can write the squared-error loss of $a\mathbf{x}$ for estimating $\boldsymbol{\theta}$ as

$$\begin{aligned} \|a\mathbf{x} - \boldsymbol{\theta}\|^2 &= (a\mathbf{x} - \boldsymbol{\theta}) \cdot (a\mathbf{x} - \boldsymbol{\theta}) \\ &= ((as - 1)\boldsymbol{\theta} + a\mathbf{r}) \cdot ((as - 1)\boldsymbol{\theta} + a\mathbf{r}) \\ &= (as - 1)^2\|\boldsymbol{\theta}\|^2 + a^2\|\mathbf{r}\|^2 \end{aligned}$$

Now consider replacing \mathbf{x} with $\mathbf{X} \sim N_p(\boldsymbol{\theta}, \mathbf{I})$. The random-variable version of the above equation is then

$$\|a\mathbf{X} - \boldsymbol{\theta}\|^2 = (aS - 1)^2\|\boldsymbol{\theta}\|^2 + a^2\|\mathbf{R}\|^2$$

Exercise: Show that

- $S \sim N(1, \|\boldsymbol{\theta}\|^{-2})$,
- $\|\mathbf{R}\|^2 \sim \chi_{p-1}^2$,
- S and \mathbf{R} are independent.

Now imagine a situation where p is growing but $\|\boldsymbol{\theta}\|^2$ remains fixed. The distribution of $(aS - 1)^2\|\boldsymbol{\theta}\|^2$ remains fixed whereas the distribution of $a^2\|\mathbf{R}\|^2$ blows up. This suggests that if we think $\|\boldsymbol{\theta}\|^2/p$ is small we should use an estimator like $a\mathbf{X}$ with $a < 1$ to control the error that comes from \mathbf{R}^2 . But what should the value of a be?

4.2 An oracle estimator:

Question: Among estimators $a\mathbf{X} : a \in [0, 1]$, which has the smallest risk?

Solution:

$$\begin{aligned} \mathbb{E}[\|a\mathbf{X} - \boldsymbol{\theta}\|^2] &= \mathbb{E}[\|(a\mathbf{X} - a\boldsymbol{\theta}) - (1 - a)\boldsymbol{\theta}\|^2] \\ &= a^2p + (1 - a)^2\|\boldsymbol{\theta}\|^2. \end{aligned}$$

Taking derivatives, the minimizing value \tilde{a} of a satisfies

$$\begin{aligned} 2\tilde{a}p - 2(1 - \tilde{a})\|\boldsymbol{\theta}\|^2 &= 0 \\ \frac{\tilde{a}}{1 - \tilde{a}} &= \frac{\|\boldsymbol{\theta}\|^2}{p} \\ \tilde{a} &= \frac{\|\boldsymbol{\theta}\|^2}{\|\boldsymbol{\theta}\|^2 + p}. \end{aligned}$$

Thus the optimal shrinkage “estimator” is given by $\delta_{\tilde{a}}(\mathbf{X}) = \tilde{a}\mathbf{X}$. This is not really an estimator in the usual sense, because the ideal degree of shrinkage \tilde{a} depends on $\boldsymbol{\theta}$. For this reason, $\tilde{a}\mathbf{X}$ is sometimes called an “oracle estimator:” You would need an oracle to tell you the value of $\|\boldsymbol{\theta}\|^2$ before you could use it. Note that the risk of this estimator is

$$\begin{aligned} E[\|a\mathbf{X} - \boldsymbol{\theta}\|^2] &= \frac{\|\boldsymbol{\theta}\|^4 p + p^2 \|\boldsymbol{\theta}\|^2}{(\|\boldsymbol{\theta}\|^2 + p)^2} \\ &= \frac{p\|\boldsymbol{\theta}\|^2(\|\boldsymbol{\theta}\|^2 + p)}{(\|\boldsymbol{\theta}\|^2 + p)^2} \\ &= \frac{p\|\boldsymbol{\theta}\|^2}{\|\boldsymbol{\theta}\|^2 + p} \end{aligned}$$

and so

$$E[\|a\mathbf{X} - \boldsymbol{\theta}\|^2] = p \frac{\|\boldsymbol{\theta}\|^2}{\|\boldsymbol{\theta}\|^2 + p} < p = E[\|\mathbf{X} - \boldsymbol{\theta}\|^2].$$

The risk differential is large if $\|\boldsymbol{\theta}\|^2$ is small compared to p .

4.3 Adaptive shrinkage estimation

As shown above, the optimal amount of shrinkage \tilde{a} is

$$\tilde{a} = \frac{\|\boldsymbol{\theta}\|^2}{\|\boldsymbol{\theta}\|^2 + p} = \frac{\|\boldsymbol{\theta}\|^2/p}{\|\boldsymbol{\theta}\|^2/p + 1}.$$

Note that $\|\boldsymbol{\theta}\|^2/p$ is the variability of the θ_j values around zero. Can this variability be estimated? Consider the following hierarchical model:

$$\begin{aligned} X_j &= \theta_j + \epsilon_j & \epsilon_1, \dots, \epsilon_p &\stackrel{\text{iid}}{\sim} N(0, 1) \\ & & \theta_1, \dots, \theta_p &\stackrel{\text{iid}}{\sim} N(0, \tau^2) \end{aligned}$$

If you'd like to connect this with some actual inference problem, imagine that each X_j is the sample mean or t -statistic calculated from observations from experiment j , with population mean θ_j .

Suppose you believed this model and knew the value of τ^2 . If you were interested finding an estimator $\delta(\mathbf{X})$ that minimized the the expected squared error $\|\boldsymbol{\theta} - \delta(\mathbf{X})\|^2$ under repeated sampling of

- $\theta_1, \dots, \theta_p$, followed by sampling of
- X_1, \dots, X_p ,

you would want to come up with an estimator $\delta(\mathbf{X})$ that minimized

$$\mathbb{E}[\|\boldsymbol{\theta} - \delta(\mathbf{X})\|^2] = \int \int \|\boldsymbol{\theta} - \delta(\mathbf{x})\|^2 p(d\mathbf{x}|\boldsymbol{\theta}) p(d\boldsymbol{\theta}).$$

Exercise: Show that $\delta_{\tau^2}(\mathbf{X}) = \frac{\tau^2}{1+\tau^2} \mathbf{X}$ minimizes the expected loss.

If we knew τ^2 , then the estimator to use would be $\frac{\tau^2}{1+\tau^2} \mathbf{X}$. We generally don't know τ^2 , but maybe it can be estimated from the data. Under the above model,

$$\begin{aligned} \mathbf{X} &= \boldsymbol{\theta} + \boldsymbol{\epsilon} \\ \boldsymbol{\theta} &\sim N_p(\mathbf{0}, \tau^2 \mathbf{I}) \\ \boldsymbol{\epsilon} &\sim N_p(\mathbf{0}, \mathbf{I}) \\ \text{Cov}(\boldsymbol{\theta}, \boldsymbol{\epsilon}) &= \mathbf{0}. \end{aligned}$$

This means that the distribution of \mathbf{X} marginalized over $\boldsymbol{\theta}$ is

$$\mathbf{X} \sim N_p(\mathbf{0}, (\tau^2 + 1)\mathbf{I}).$$

An unbiased estimator of $\tau^2 + 1$ is clearly $\|\mathbf{X}\|^2/p$, so an unbiased estimator of τ^2 is

$$\hat{\tau}^2 = \frac{\|\mathbf{X}\|^2 - p}{p}.$$

However, we were interested in estimating $\tau^2/(\tau^2 + 1)$, not τ^2 . If $p > 2$, you can use the fact that $\|\mathbf{X}\|^2 \sim \text{gamma}(p/2, 1/[2(\tau^2 + 1)])$ to show that

$$\mathbb{E}[\|\mathbf{X}\|^{-2}] = [(p - 2)(\tau^2 + 1)]^{-1},$$

and so

$$\begin{aligned} \mathbb{E}\left[\frac{p-2}{\|\mathbf{X}\|^2}\right] &= \frac{1}{\tau^2 + 1} \\ \mathbb{E}\left[1 - \frac{p-2}{\|\mathbf{X}\|^2}\right] &= \frac{\tau^2}{\tau^2 + 1}. \end{aligned}$$

Again, $\frac{\tau^2}{\tau^2+1}\mathbf{X}$ would be the optimal estimator in this hierarchical model if we knew τ^2 . If we don't know τ^2 , we might instead consider using

$$\begin{aligned} \delta_{JS}(\mathbf{X}) &= \left(\frac{\widehat{\tau^2}}{\tau^2 + 1}\right)\mathbf{X} \\ &= \left(1 - \frac{p-2}{\|\mathbf{X}\|^2}\right)\mathbf{X}. \end{aligned}$$

This estimator is called the James-Stein estimator. As we will see, it has many interesting properties:

- For large p in the hierarchical normal model, it is almost as good as the oracle estimator $\frac{\tau^2}{\tau^2+1}\mathbf{X}$:

$$\mathbb{E}_{X|\theta}[\|\boldsymbol{\theta} - \delta_{JS}\|^2] \approx \mathbb{E}_{X|\theta}[\|\boldsymbol{\theta} - \frac{\tau^2}{\tau^2+1}\mathbf{X}\|^2]$$

- Even if the hierarchical normal model isn't correct, it still is almost as good as the oracle estimator $\tilde{a}\mathbf{X}$ in the normal means model:

$$\mathbb{E}_{X|\theta}[\|\boldsymbol{\theta} - \delta_{JS}\|^2] \approx \mathbb{E}_{X|\theta}[\|\boldsymbol{\theta} - \tilde{a}\mathbf{X}\|^2]$$

- In the normal means problem, this estimator dominates the unbiased estimator \mathbf{X} if $p > 2$:

$$\mathbb{E}_{X|\theta}[\|\boldsymbol{\theta} - \delta_{JS}\|^2] < \mathbb{E}_{X|\theta}[\|\boldsymbol{\theta} - \mathbf{X}\|^2] \quad \forall \boldsymbol{\theta}$$

We will show this last inequality first.

5 Risk of δ_{JS}

We will show that δ_{JS} dominates \mathbf{X} by showing that the risk $R(\boldsymbol{\theta}, \delta_{JS}) = \mathbb{E}[|\delta_{JS} - \boldsymbol{\theta}|^2]/p$ is uniformly less than 1. This will not be done by computing its risk function directly, but instead by showing that 1 is an upper bound on the risk. This bound will be obtained via an identity that has applications beyond the calculation of $R(\boldsymbol{\theta}, \delta_{JS})$.

5.1 Risk bound for δ_{JS}

We can write the James-Stein estimator as

$$\begin{aligned} \delta_{JS} &= \left(\frac{\widehat{\tau^2}}{1 + \widehat{\tau^2}} \right) \mathbf{x} = \left(1 - \frac{p-2}{\mathbf{x} \cdot \mathbf{x}} \right) \mathbf{x} \\ &= \mathbf{x} - \frac{p-2}{\mathbf{x} \cdot \mathbf{x}} \mathbf{x} \\ &\equiv \mathbf{x} - g(\mathbf{x}). \end{aligned}$$

Under $\mathbf{X} \sim N_p(\boldsymbol{\theta}, \mathbf{I})$,

$$\begin{aligned} \mathbb{E}[|\delta_{JS} - \boldsymbol{\theta}|^2] &= \mathbb{E}[(\mathbf{X} - g(\mathbf{X}) - \boldsymbol{\theta})^2] \\ &= \mathbb{E}[(\mathbf{X} - \boldsymbol{\theta}) - g(\mathbf{X})]^2 \\ &= \mathbb{E}[|\mathbf{X} - \boldsymbol{\theta}|^2] + \mathbb{E}[|g(\mathbf{X})|^2] - 2\mathbb{E}[(\mathbf{X} - \boldsymbol{\theta}) \cdot g(\mathbf{X})]. \end{aligned}$$

where all expectations are with respect to the distribution of \mathbf{X} given $\boldsymbol{\theta}$. The first expectation is p and the second is

$$\mathbb{E}[(p-2)^2 \frac{\mathbf{X} \cdot \mathbf{X}}{(\mathbf{X} \cdot \mathbf{X})^2}] = (p-2)^2 \mathbb{E}[\frac{1}{\mathbf{X} \cdot \mathbf{X}}].$$

The third expectation is more complicated, but in the next subsection we'll derive an identity (Stein's identity) for computing $\mathbb{E}[(\mathbf{X} - \boldsymbol{\theta}) \cdot g(\mathbf{X})]$ that is applicable for arbitrary functions g . Stein's identity as applied to $g(\mathbf{x}) = \frac{p-2}{\mathbf{x} \cdot \mathbf{x}} \mathbf{x}$ gives

$$\mathbb{E}[(\mathbf{X} - \boldsymbol{\theta})g(\mathbf{X})] = \mathbb{E}[\frac{(p-2)^2}{\mathbf{X} \cdot \mathbf{X}}].$$

Using this for the above risk calculation gives

$$\begin{aligned} \mathbb{E}[\|\delta_{JS} - \boldsymbol{\theta}\|^2] &= p + (p-2)^2 \mathbb{E}\left[\frac{1}{\mathbf{X} \cdot \mathbf{X}}\right] - 2(p-2)^2 \mathbb{E}\left[\frac{1}{\mathbf{X} \cdot \mathbf{X}}\right] \\ &= p - \mathbb{E}\left[\frac{(p-2)^2}{\mathbf{X} \cdot \mathbf{X}}\right]. \end{aligned}$$

Note that we haven't actually calculated the risk of δ_{JS} in closed form - our formula depends on the expectation of $1/\mathbf{X} \cdot \mathbf{X}$, which is an inverse-moment of a noncentral χ^2 distribution where the noncentrality parameter depends on $\boldsymbol{\theta}$. However, computing this moment is not necessary to show that δ_{JS} dominates \mathbf{X} : Since $\frac{1}{\mathbf{x} \cdot \mathbf{x}} > 0 \forall \mathbf{x} \in \mathbb{R}^p$, we have

$$\mathbb{E}[\|\delta_{JS} - \boldsymbol{\theta}\|^2] = p - \mathbb{E}\left[\frac{(p-2)^2}{\mathbf{X} \cdot \mathbf{X}}\right] < p = \mathbb{E}[\|\mathbf{X} - \boldsymbol{\theta}\|^2].$$

Since the expectation of $(\mathbf{X} \cdot \mathbf{X})^{-1}$ is complicated, further study of the risk of δ_{JS} is often achieved via a study of its *unbiased risk estimate*. From the above calculation, we see that

$$\mathbb{E}[\|\delta_{JS} - \boldsymbol{\theta}\|^2] = \mathbb{E}\left[p - \frac{(p-2)^2}{\mathbf{X} \cdot \mathbf{X}}\right],$$

and so $p - \frac{(p-2)^2}{\mathbf{X} \cdot \mathbf{X}}$ can be said to be an unbiased estimate of the risk of δ_{JS} .

5.2 Stein's identity

We start with a univariate version of the identity:

Lemma 4 (Stein's identity). *Let $X \sim N(\mu, \sigma^2)$ and let $g(x)$ be such that $\mathbb{E}[|g'|] < \infty$. Then*

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

Proof.

The proof follows from Fubini's theorem and a bit of calculus. Letting $p(x) = \phi([x - \mu]/\sigma)/\sigma$, note that $p'(x) = -(\frac{x-\mu}{\sigma^2})p(x)$. By the fundamental theorem of calculus,

$$p(x) = \int_{-\infty}^x -\left(\frac{y-\mu}{\sigma^2}\right)p(y) dy = \int_x^{\infty} \left(\frac{y-\mu}{\sigma^2}\right)p(y) dy.$$

The expectation we wish to calculate is

$$\begin{aligned} \mathbb{E}[g'(x)] &= \int_{-\infty}^{\infty} g'(x)p(x) dx \\ &= \int_0^{\infty} g'(x)p(x) dx + \int_{-\infty}^0 g'(x)p(x) dx. \end{aligned}$$

Doing the first part, we have

$$\begin{aligned} \int_0^{\infty} g'(x)p(x) dx &= \int_0^{\infty} g'(x) \int_x^{\infty} \left(\frac{y-\mu}{\sigma^2}\right)p(y) dy dx \\ &= \int_{0 < x < y < \infty} g'(x) \left(\frac{y-\mu}{\sigma^2}\right)p(y) dy dx \\ &= \int_0^{\infty} \left(\frac{y-\mu}{\sigma^2}\right)p(y) \int_0^y g'(x) dx dy \\ &= \int_0^{\infty} \left(\frac{y-\mu}{\sigma^2}\right)p(y)[g(y) - g(0)] dy \\ &= \mathbb{E}[(g(X) - g(0))\left(\frac{X-\mu}{\sigma^2}\right)1(X > 0)]. \end{aligned}$$

Similarly, the second part is

$$\int_{-\infty}^0 g'(x)p(x) dx = \mathbb{E}[(g(X) - g(0))\left(\frac{X-\mu}{\sigma^2}\right)1(X < 0)].$$

Adding the two parts gives

$$\begin{aligned} \mathbb{E}[g'(x)] &= \mathbb{E}[(g(X) - g(0))\left(\frac{X-\mu}{\sigma^2}\right)] \\ &= \mathbb{E}[g(X)\left(\frac{X-\mu}{\sigma^2}\right)] - g(0)\mathbb{E}\left[\left(\frac{X-\mu}{\sigma^2}\right)\right] = \mathbb{E}[g(X)\left(\frac{X-\mu}{\sigma^2}\right)], \end{aligned}$$

and so

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

□

Stein's lemma is often alternatively proven with integration by parts. These proofs go roughly as follows: As before,

$$\begin{aligned} p(x) &= (2\pi\sigma^2)^{-1/2} \exp\{-(x - \mu)^2/[2\sigma^2]\} \\ \frac{dp(x)}{dx} &= -\left(\frac{1}{\sigma^2}\right)(x - \mu)p(x). \end{aligned}$$

$$\begin{aligned}
\mathbb{E}[g'(X)] &= \int_{-\infty}^{\infty} g'(x) \times p(x) \, dx \\
&= \lim_{c \rightarrow \infty} \int_{-c}^c g'(x) \times p(x) \, dx \\
&\equiv \lim_{c \rightarrow \infty} \int_{-c}^c v'(x)u(x) \, dx \\
&= \lim_{c \rightarrow \infty} \left(u(x)v(x)|_{-c}^c - \int_{-c}^c v(x)u'(x) \, dx \right) \\
&= \lim_{c \rightarrow \infty} \left(g(x)p(x)|_{-c}^c - \int_{-c}^c g(x)[-(x - \mu)/\sigma^2]p(x) \, dx \right) \\
&= \mathbb{E}[g(X)(X - \mu)/\sigma^2] + \lim_{c \rightarrow \infty} g(x)p(x)|_{-c}^c.
\end{aligned}$$

To complete the proof we have to show that the last limit is zero. This is straightforward to show if $p(x)$ decreases monotonically in $|x|$:

Lemma 5. *Let $p(x)$ be decreasing to zero in $|x|$ and let $\mathbb{E}[|g'(X)|] < \infty$. Then $g(x)p(x) \rightarrow 0$ as $x \rightarrow \pm\infty$.*

Proof. Given $\epsilon > 0 \exists K$ such that

$$\int_K^{\infty} |g'(x)|p(x) \, dx < \epsilon/3.$$

Then for any t sufficiently large,

1. $p(t) < p(K)(\int_0^K |g'(x)|p(x))^{-1} \times \epsilon/3$
2. $p(t)|g(0)| < \epsilon/3$.

From this, we have

$$\begin{aligned}
|g(t)p(t)| &= p(t)|g(0) + \int_0^t g'(x) dx| \\
&= p(t)|g(0) + \int_0^K g'(x) dx + \int_K^t g'(x) dx| \\
&\leq p(t)|g(0)| + p(t) \int_0^K |g'(x)| dx + p(t) + p(t) \int_K^t |g'(x)| dx \\
&\leq p(t)|g(0)| + \frac{p(t)}{p(K)} \int_0^K |g'(x)|p(x) dx + \int_K^t |g'(x)|p(x) dx \\
&\leq \epsilon/3 + \epsilon/3 + \epsilon/3 = \epsilon,
\end{aligned}$$

where the second to last line holds as $p(x)/p(K) < 1$ and $p(x)/p(t) < 1$ on $x \in (0, K)$ and (K, t) respectively, due to $p(x)$ being monotonically decreasing. \square

This identity generalizes to other exponential families. See LC Theorem 1.5.15.

For computing the risk of a vector-valued function $g : \mathbb{R}^p \rightarrow \mathbb{R}^p$, we will need a multivariate version of the above identity.

Lemma 6 (Stein's identity, multivariate version). *Let $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ and let $g(\mathbf{x}) : \mathbb{R}^p \rightarrow \mathbb{R}^p$ such that $E[|dg_i/dx_i|] < \infty$. Then*

$$E[(\mathbf{X} - \boldsymbol{\mu}) \cdot g(\mathbf{X})] = \sigma^2 E[\nabla \cdot g(\mathbf{X})]$$

where $\nabla \cdot g = \sum_{j=1}^p dg_j(\mathbf{x})/dx_j$.

Proof. This is just a corollary of the univariate version:

$$\begin{aligned}
E[dg_p(\mathbf{x})/dx_p] &= \int_{\mathbf{x}_{-p}} \left(\int_{x_p} \frac{dg_p(\mathbf{x})}{dx_p} p(x_p) dx_p \right) \prod_1^{p-1} p(x_j) dx_j \\
&= \int_{\mathbf{x}_{-p}} E_{x_p}[g_p(\mathbf{X})(X_p - \mu_p)] / \sigma^2 \prod_1^{p-1} p(x_j) dx_j \\
&= \frac{1}{\sigma^2} E[g_p(\mathbf{X})(X_p - \mu_p)],
\end{aligned}$$

and similarly for each other $j \in \{1, \dots, p\}$. Therefore

$$\begin{aligned} \mathbb{E}[\nabla \cdot g(\mathbf{X})] &= \sum \mathbb{E}[dg_j(\mathbf{x})/dx_j] \\ &= \sum \mathbb{E}[g_j(\mathbf{X})(X_j - \mu_j)]/\sigma^2 \\ &= \mathbb{E}[g(\mathbf{X}) \cdot (\mathbf{X} - \boldsymbol{\mu})]/\sigma^2. \end{aligned}$$

□

Now we are in a position to apply the lemma to obtain the unbiased risk estimator of δ_{JS} . Recall we needed to calculate $\mathbb{E}[(\mathbf{X} - \boldsymbol{\theta}) \cdot g(\mathbf{X})]$ where $g(\mathbf{x}) = \frac{(p-2)}{\mathbf{x} \cdot \mathbf{x}} \mathbf{x}$. Applying the lemma, we have

$$\begin{aligned} \mathbf{g}(\mathbf{x}) &= (p-2) \left(\frac{x_1}{\mathbf{x} \cdot \mathbf{x}}, \dots, \frac{x_p}{\mathbf{x} \cdot \mathbf{x}} \right) \\ \nabla \cdot \mathbf{g}(\mathbf{x}) &= (p-2) \sum \left(\frac{\mathbf{x} \cdot \mathbf{x} - 2x_j^2}{(\mathbf{x} \cdot \mathbf{x})^2} \right) \\ &= \frac{p-2}{(\mathbf{x} \cdot \mathbf{x})^2} [p\mathbf{x} \cdot \mathbf{x} - 2\mathbf{x} \cdot \mathbf{x}] \\ &= \frac{(p-2)^2}{\mathbf{x} \cdot \mathbf{x}}, \end{aligned}$$

and so

$$\mathbb{E}[(\mathbf{X} - \boldsymbol{\theta}) \cdot g(\mathbf{X})] = \mathbb{E}\left[\frac{(p-2)^2}{\mathbf{X} \cdot \mathbf{X}}\right],$$

as we used above.

6 Some oracle inequalities

6.1 A simple oracle inequality

Recall that if we knew $\|\boldsymbol{\theta}\|^2$, the optimal estimator in the class $\{\delta_a(\mathbf{x}) = a\mathbf{x} : a \in [0, 1]\}$ would be $\delta_{\tilde{a}}$ where

$$\tilde{a} = \frac{\|\boldsymbol{\theta}\|^2/p}{\|\boldsymbol{\theta}\|^2/p + 1}.$$

We also showed that the risk of this estimator is

$$R(\boldsymbol{\theta}, \tilde{\mathbf{a}}\mathbf{X}) = \mathbb{E}[\|\tilde{\mathbf{a}}\mathbf{X} - \boldsymbol{\theta}\|^2]/p = \frac{\|\boldsymbol{\theta}\|^2}{\|\boldsymbol{\theta}\|^2 + p} < 1.$$

Not surprisingly, it turns out that

$$R(\boldsymbol{\theta}, \tilde{\mathbf{a}}\mathbf{X}) \leq R(\boldsymbol{\theta}, \boldsymbol{\delta}_{JS})$$

(use the fact that $\mathbf{X} \cdot \mathbf{X}$ has a noncentral χ^2 distribution, or condition on $\mathbf{X} \cdot \mathbf{X}$).

But how much worse is $\boldsymbol{\delta}_{JS}$ than the oracle estimator $\boldsymbol{\delta}_{\tilde{\mathbf{a}}}$? Recall the risk of $\boldsymbol{\delta}_{JS}$ is

$$R(\boldsymbol{\theta}, \boldsymbol{\delta}_{JS}) = 1 - \frac{(p-2)^2}{p} \mathbb{E}[(\mathbf{X} \cdot \mathbf{X})^{-1}]$$

Since $1/x$ is convex, Jensen's inequality gives

$$\mathbb{E}\left[\frac{1}{\mathbf{X} \cdot \mathbf{X}}\right] \geq \frac{1}{\mathbb{E}[\mathbf{X} \cdot \mathbf{X}]} = \frac{1}{\|\boldsymbol{\theta}\|^2 + p},$$

and so

$$\begin{aligned} R(\boldsymbol{\theta}, \boldsymbol{\delta}_{JS}) &\leq 1 - \frac{(p-2)^2}{p} \frac{1}{\|\boldsymbol{\theta}\|^2 + p} \\ p(R(\boldsymbol{\theta}, \boldsymbol{\delta}_{JS}) - R(\boldsymbol{\theta}, \boldsymbol{\delta}_{\tilde{\mathbf{a}}})) &\leq p - \frac{(p-2)^2}{\|\boldsymbol{\theta}\|^2 + p} - \frac{\|\boldsymbol{\theta}\|^2 p}{\|\boldsymbol{\theta}\|^2 + p} \\ &= \frac{p\|\boldsymbol{\theta}\|^2 + p^2 - p^2 + 4p - 4 - \|\boldsymbol{\theta}\|^2 p}{\|\boldsymbol{\theta}\|^2 + p} \\ &= \frac{4(p-1)}{\|\boldsymbol{\theta}\|^2 + p} \\ &\leq 4 \frac{p-1}{p} \leq 4, \end{aligned}$$

and so

$$R(\boldsymbol{\theta}, \boldsymbol{\delta}_{\tilde{\mathbf{a}}}) \leq R(\boldsymbol{\theta}, \boldsymbol{\delta}_{JS}) \leq R(\boldsymbol{\theta}, \boldsymbol{\delta}_{\tilde{\mathbf{a}}}) + 4/p.$$

Additional work can get you to

$$R(\boldsymbol{\theta}, \boldsymbol{\delta}_{\tilde{\mathbf{a}}}) \leq R(\boldsymbol{\theta}, \boldsymbol{\delta}_{JS}) \leq R(\boldsymbol{\theta}, \boldsymbol{\delta}_{\tilde{\mathbf{a}}}) + 2/p.$$

For more on this, see Johnstone [2002] or Candès [2006].

7 Unknown variance or covariance

Suppose $\mathbf{X} \sim N_p(\boldsymbol{\theta}, \sigma^2 \mathbf{I})$, where σ^2 is known. Letting

- $\tilde{\mathbf{X}} = \mathbf{X}/\sigma$ and
- $\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}/\sigma$,

we have $\tilde{\mathbf{X}} \sim N_p(\tilde{\boldsymbol{\theta}}, \mathbf{I})$. The James-Stein estimator $\tilde{\boldsymbol{\delta}}_{JS}$ of $\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}/\sigma$ is then

$$\begin{aligned}\tilde{\boldsymbol{\delta}}(\tilde{\mathbf{X}}) &= \left(1 - \frac{p-2}{\tilde{\mathbf{X}} \cdot \tilde{\mathbf{X}}}\right) \tilde{\mathbf{X}} \\ &= \left(1 - \frac{\sigma^2(p-2)}{\mathbf{X} \cdot \mathbf{X}}\right) \mathbf{X}/\sigma.\end{aligned}$$

It seems natural then that the JSE of $\boldsymbol{\theta}$ should be σ times the JSE of $\tilde{\boldsymbol{\theta}} = \boldsymbol{\theta}/\sigma$, giving

$$\boldsymbol{\delta}_{JS} = \left(1 - \frac{\sigma^2(p-2)}{\mathbf{X} \cdot \mathbf{X}}\right) \mathbf{X}.$$

Of course, often σ^2 is not known. Is there a version of the JSE in this case? Yes, if you have information about σ^2 . Consider the following hierarchical sampling scheme:

$$\begin{aligned}X_{i,j} &= \theta_j + \epsilon_{i,j}, \quad i = 1, \dots, n, \quad j = 1, \dots, p \\ \{\epsilon_{i,j}\} &\sim \text{i.i.d. } N(0, \sigma^2).\end{aligned}$$

Letting $X_j = \bar{X}_{.j} = \sum_{i=1}^n X_{i,j}/n$, we now basically have the situation described above.

Also note that the data contain information about σ^2 via the pooled sample sum of squares:

$$S = \sum_{j=1}^p (X_{i,j} - \bar{X}_{.j})^2 \sim \sigma^2 \chi_{p(n-1)}^2.$$

Note further that \mathbf{X} and S are statistically independent. For this and similar situations, James and Stein [1961] considered estimators of the following form:

Let $\mathbf{X} \sim N_p(\boldsymbol{\theta}, \sigma^2 \mathbf{I})$ be independent of $S \sim \sigma^2 \chi_k^2$. Define the estimator

$$\boldsymbol{\delta}_c(\mathbf{X}, S) = \left(1 - \frac{cS}{\mathbf{X} \cdot \mathbf{X}}\right) \mathbf{X}.$$

The value of c that minimizes the risk of δ_c for all $\boldsymbol{\theta}$ is $c = (p - 2)/(k + 2)$, resulting in the following estimator:

$$\delta_{JS}(\mathbf{X}, S) = \left(1 - \frac{S}{k+2} \frac{(p-2)}{\mathbf{X} \cdot \mathbf{X}}\right) \mathbf{X}.$$

In particular, note that this estimator dominates \mathbf{X} .

This result also generalizes to the correlated data case: Let $\mathbf{X} \sim N_p(\boldsymbol{\theta}, \Sigma)$ and $\mathbf{S} \sim \text{Wishart}(\Sigma, k)$ be independent. Consider estimators of the form

$$\delta_c(\mathbf{X}, \mathbf{S}) = \left(1 - \frac{c}{\mathbf{X}^T \mathbf{S}^{-1} \mathbf{X}}\right) \mathbf{X}.$$

James and Stein [1961] show that the estimator obtained by setting

$$c = \frac{p - 2}{n - p + 3}$$

minimizes the risk for all values of $\boldsymbol{\theta}$.

See Brown and Han for recent work on these and related problems, including extensions to regression problems.

References

- L.D. Brown and X. Han. Optimal estimation of multidimensional normal means with an unknown variance.
- Emmanuel J. Candès. Modern statistical estimation via oracle inequalities. *Acta Numer.*, 15:257–325, 2006. ISSN 0962-4929. doi: 10.1017/S0962492906230010. URL <http://dx.doi.org.offcampus.lib.washington.edu/10.1017/S0962492906230010>.
- W. James and Charles Stein. Estimation with quadratic loss. In *Proc. 4th Berkeley Sympos. Math. Statist. and Prob., Vol. I*, pages 361–379. Univ. California Press, Berkeley, Calif., 1961.

IM Johnstone. Function estimation and gaussian sequence models. *Unpublished manuscript*, 2002.

E. L. Lehmann and George Casella. *Theory of point estimation*. Springer Texts in Statistics. Springer-Verlag, New York, second edition, 1998. ISBN 0-387-98502-6.

Charles Stein. Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, 1954–1955, vol. I*, pages 197–206, Berkeley and Los Angeles, 1956. University of California Press.