# REVERSING THE STEIN EFFECT

Michael Perlman and Sanjay Chaudhuri

June 18, 2019

# INTRODUCTION

- We identify the *Reverse Stein Effect*.

- A statistician who shrinks his/her data toward a point chosen without reliable knowledge about the underlying value of the parameter to be estimated but based instead upon the observed data will not be protected by the minimax property of shrinkage estimators such as that of James and Stein.

- But instead, he/she will likely incur a greater error than if shrinkage were not used.
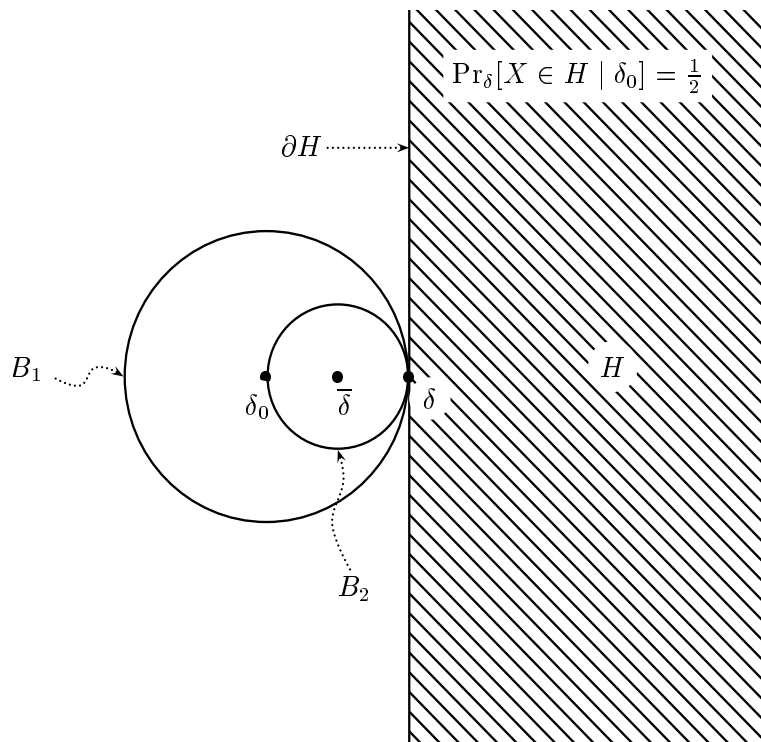
## THE CASE FOR SHRINKAGE: THE STEIN EFFECT

- Suppose, $X = Y + \delta \in \mathbb{R}^p$ is an observed random vector.

- $\delta$ is an unknown location parameter,

- The unobserved random vector $Y$ is absolutely continuous.

- We further assume that, $Y \equiv X - \delta$ is *directionally symmetric*, that is, $\vec{Y} \stackrel{d}{=} -\vec{Y}$, where $\vec{Y} := Y/\|Y\|$ is the unit vector in the direction of $Y$.

- We consider the "shrinkage" estimators for $\delta$ of the form:

$$\hat{\delta}_\gamma \equiv \hat{\delta}_\gamma(X; \delta_0) = \gamma(X - \delta_0) \cdot (X - \delta_0) + \delta_0, \tag{1}$$

where $\gamma \equiv \gamma(X - \delta_0) \in [0, 1)$ and $\delta_0$ is *any* fixed shrinkage target point in $\mathbb{R}^p$.

- The improvement offered by such shrinkage estimators is often referred to as *the Stein Effect*.

# THE STEIN EFFECT (CONTD.)



$\mathrm{Pr}_\delta[X \in H \mid \delta_0] = \frac{1}{2}$

$\partial H$

$B_1$

$\delta_0$

$\bar{\delta}$

$\delta$

$H$

$B_2$

- For fixed $\delta$ and $\delta_0$, let $B_1 \equiv B_1(\|\delta-\delta_0\|; \delta_0) \subset \mathbb{R}^p$ denote the ball of radius $\|\delta - \delta_0\|$ centred at $\delta_0$.
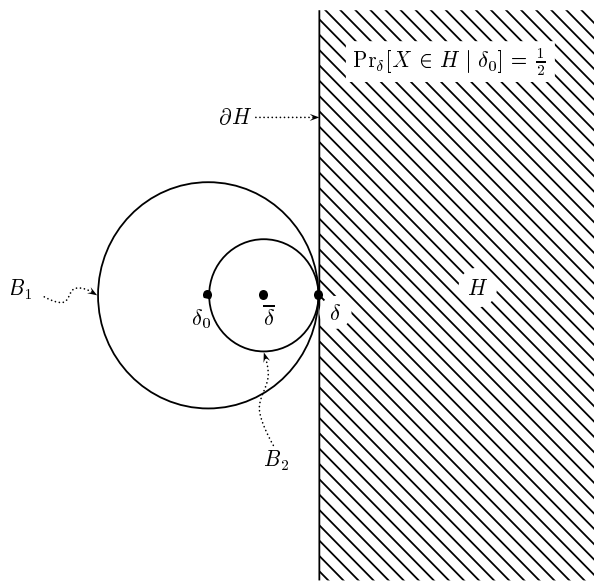- Let $H$ be the halfspace bounded by a hyperplane $\partial H$ tangent to $B_1$ at $\delta$.
- 
$$\mathrm{Pr}_\delta \left[ \|X - \delta_0\| > \|\delta - \delta_0\| \;\middle|\; \delta_0 \right]$$
$$= \mathrm{Pr}_\delta \left[ X \in B_1^c \;\middle|\; \delta_0 \right]$$
$$> \mathrm{Pr}_\delta[ X \in H \mid \delta_0] = \frac{1}{2}, \qquad (2)$$

- The inequality in (2) follows from directional symmetry.

- Furthermore, under somewhat stronger but still general assumptions,

$$\lim_{p \to \infty} \mathrm{Pr}_\delta \left[ \|X - \delta_0\| > \|\delta - \delta_0\| \right] \equiv \lim_{p \to \infty} \mathrm{Pr}_\delta \left[ X \in B_1^c \right] = 1. \qquad (3)$$

3

$\Pr_\delta[X \in H \mid \delta_0] = \frac{1}{2}$

$\partial H$

$B_1$

$\delta_0$  $\overline{\delta}$  $\delta$

$H$

$B_2$

- Thus $\|X - \delta_0\|$ is usually an overestimate of $\|\delta - \delta_0\|$, so an estimator of the form $\gamma(X - \delta_0) \cdot (X - \delta_0)$ for $\delta - \delta_0$ should be preferable to $X - \delta_0$ itself.
- Writing $\delta$ as $(\delta - \delta_0) + \delta_0$ immediately leads to estimators for $\delta$ of the form (1).
- Now, suppose, $\tilde{\gamma} \equiv \tilde{\gamma}(X - \delta_0, \delta - \delta_0)$ is allowed to depend on $\delta$ and $B_2 \equiv B_2(\|\delta - \delta_0\|; \overline{\delta})$ is the ball of radius $\frac{1}{2}\|\delta - \delta_0\|$ centred at $\frac{1}{2}(\delta_0 + \delta) \equiv \overline{\delta}$.
- Then $\{\, X \mid \exists \tilde{\gamma} \in [0, 1) \ni \|\hat{\delta}_{\tilde{\gamma}} - \delta\| < \|X - \delta\|\,\} = B_2^c$.
- Since $B_2^c \supset B_1^c$, $\Pr_\delta[X \in B_2^c \mid \delta_0] > \frac{1}{2}$.

- Furthermore, under the assumptions,

$$\lim_{p \to \infty} \Pr_\delta[X \in B_2^c] = 1. \tag{4}$$

- This shows that if $\delta$ were known then usually *some* shrinkage factor $\tilde{\gamma}$ applied to $X - \delta_0$ will move $X$ closer to $\delta$, again suggesting a search for estimators of the form (1).

# THE STEIN PARADOX

- Assume now that $Y \sim N_p(0, \sigma^2 I_p)$, the multivariate normal distribution with mean 0 and covariance matrix $\sigma^2 I_p$, where $\sigma^2 > 0$ is known, so $X \sim N_p(\delta, \sigma^2 I_p)$.

- In this simple case, the James-Stein (JS) estimator for $\delta$ is given by

$$\hat{\delta}_{JS} \equiv \hat{\delta}_{JS}(X; \delta_0) = \left(1 - \frac{\sigma^2(p-2)}{\|X - \delta_0\|^2}\right)(X - \delta_0) + \delta_0, \qquad (5)$$

  where $\delta_0$ is a fixed but *arbitrary* point in $\mathbb{R}^p$.

- The truncated $\equiv$ "plus-rule" JS estimator

$$\hat{\delta}_{JS}^+ \equiv \hat{\delta}_{JS}^+(X; \delta_0) = \left(1 - \frac{\sigma^2(p-2)}{\|X - \delta_0\|^2}\right)^+ (X - \delta_0) + \delta_0 \qquad (6)$$

  is a shrinkage estimator of the form (1).

- These renowned estimators have the property that when $p \geq 3$, they dominate $X$ under both the mean square error (MSE) and Pitman closeness (PC) criteria.

- For every fixed $\delta, \delta_0 \in \mathbb{R}^p$,

$$\mathsf{E}_\delta[\|\hat{\delta}_{JS}^+(X; \delta_0) - \delta\|^2 \mid \delta_0] < \mathsf{E}_\delta[\|\hat{\delta}_{JS}(X; \delta_0) - \delta\|^2 \mid \delta_0] < \mathsf{E}_\delta[\|X - \delta\|^2] \equiv p\sigma^2, \tag{7}$$

- We can also show:

$$\mathsf{Pr}_\delta[\|\hat{\delta}_{JS}^+(X; \delta_0) - \delta\| < \|X - \delta\| \mid \delta_0] > \mathsf{Pr}_\delta[\|\hat{\delta}_{JS}(X; \delta_0) - \delta\| < \|X - \delta\| \mid \delta_0]$$

$$= \mathsf{Pr}\left[\chi_p^2\left(\frac{\|\delta - \delta_0\|^2}{4\sigma^2}\right) \geq \frac{\|\delta - \delta_0\|^2}{4\sigma^2} + \frac{p-2}{2}\right] > \frac{1}{2} \tag{8}$$

and approaches 1 as $p \to \infty$ if $\frac{\|\delta - \delta_0\|}{\sigma} = o(p)$.

- Note especially that:

(A) the improvements offered by the JS estimators can be great, especially when $p$ is large: if $\delta = \delta_0$ then $\mathsf{MSE}(\hat{\delta}_{JS}^+) < \mathsf{MSE}(\hat{\delta}_{JS}) = 2\sigma^2 \ll p\sigma^2$, and if $\|\delta - \delta_0\| = o(p)$ with $\sigma^2$ fixed then $\mathsf{Pr}_\delta[\|\hat{\delta} - \delta\| < \|X - \delta\|] \to 1$ as $p \to \infty$ for both $\hat{\delta} = \hat{\delta}_{JS}$ and $\hat{\delta}_{JS}^+$;

(B) the MSE and PC dominances of $X$ by $\hat{\delta}_{JS}$ and $\hat{\delta}_{JS}^+$ hold even if the true mean $\delta$ is arbitrarily far from the shrinkage target $\delta_0$.
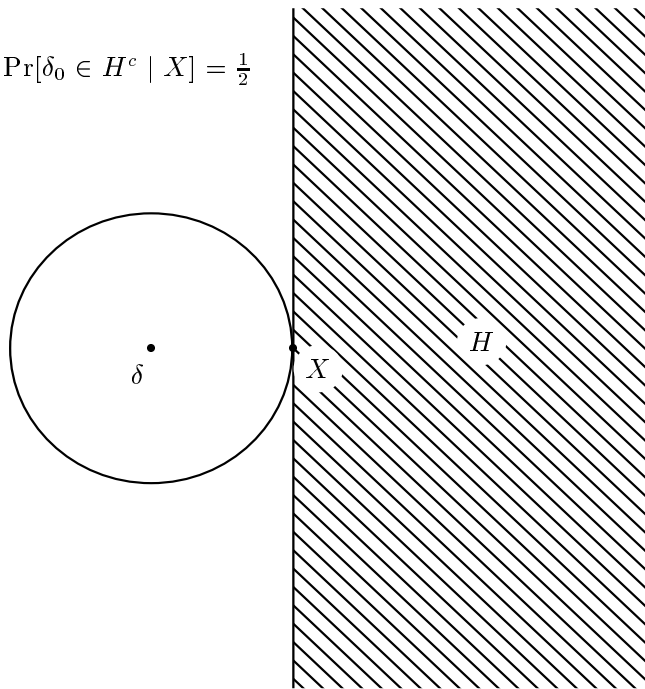
# THE STEIN PARADOX (CONTD.)

- Of the two properties (A) and (B), it is (B) that is most surprising.

- It is not difficult to construct estimators that satisfy (A).

- For example a Bayes estimator w.r. to a normal prior centred at $\delta_0$.

- However, such a Bayes estimator will not satisfy (B), the difference stemming from the fact that the Bayes estimator will have a constant shrinkage factor while the shrinkage factors in (5) and (6) are adaptive.

- When first discovered, the domination of $X$ by the JS estimators was highly surprising, because the estimator $X$ itself is:

  (a) the best unbiased estimator of $\delta$,

  (b) the best translation-invariant estimator of $\delta$,

  (c) the maximum likelihood estimator (MLE) of $\delta$,

  (d) a minimax estimator of $\delta$, and

  (e) an admissible estimator of $\delta$ when $p = 1$ or 2.

  So compelling were these properties of $X$ that its domination by the JS estimators came to be known as *the Stein Paradox*.

# THE REVERSE STEIN EFFECT

$\Pr[\delta_0 \in H^c \mid X] = \frac{1}{2}$
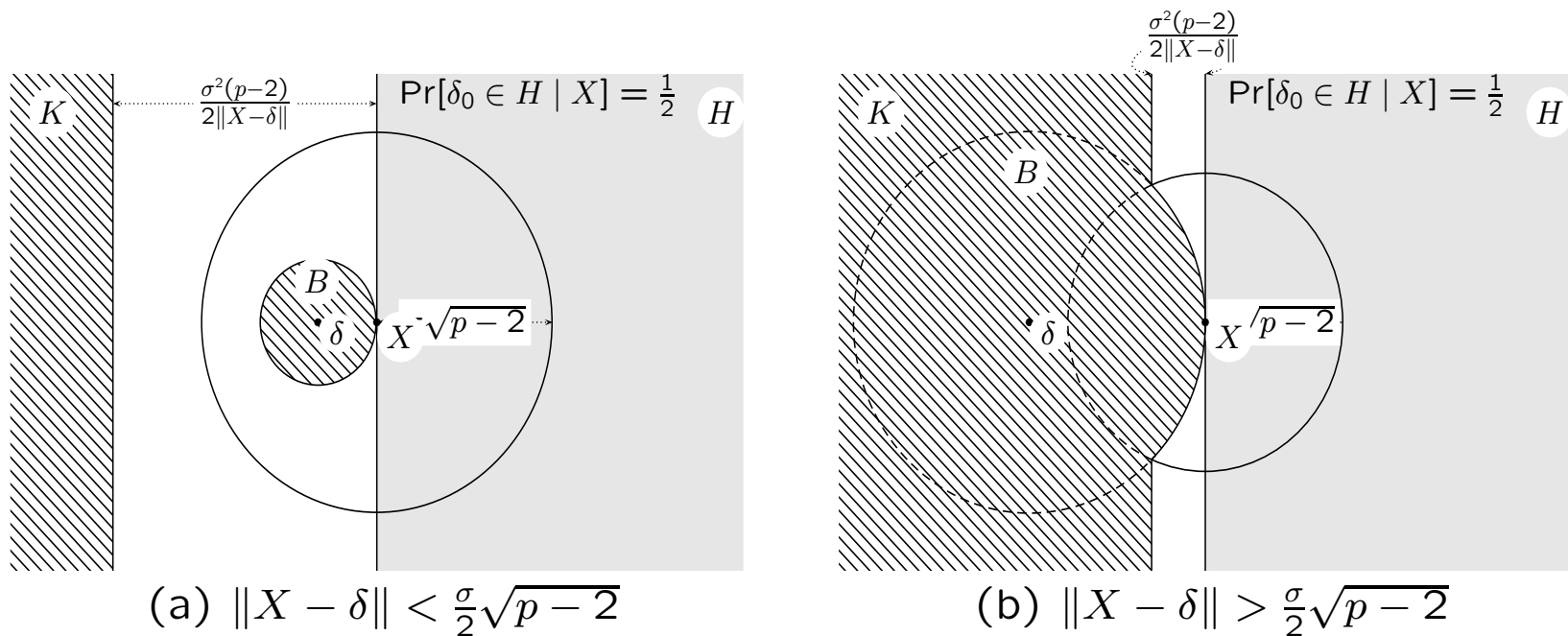
$H$

$\delta$

$X$

- Suppose that the position $X \sim N_p(\delta, \sigma^2 I_p)$ is known and suppose $p \geq 3$.
- If we use the truncated James-Stein estimator $\hat{\delta}_{JS}^+(X; \delta_0)$ to estimate $\delta$ by shrinking $X$ toward a fixed point $\delta_0$, then by (7) and (8), $\hat{\delta}_{JS}^+(X; \delta_0)$ is more likely to be closer to $\delta$ than our present location $X$ is, no matter where $\delta$ is! And what's more, we can shrink $X$ toward any $\delta_0$ that we like!
- Consider the set,

$$\{\, \delta_0 \mid \exists \breve{\gamma} \in [0, 1) \ni \|\hat{\delta}_{\breve{\gamma}} - \delta\| < \|X - \delta\|\}, \quad (9)$$

where $\breve{\gamma} \equiv \breve{\gamma}(X - \delta_0, X - \delta)$ may depend on $\delta$.
- This set is exactly $H^c$.

- Then, since $\Pr[\delta_0 \in H^c | X] = \frac{1}{2}$ by directional symmetry, this shows that shrinkage toward a randomly chosen $\delta_0$ would have at most a 50-50 chance of moving $X$ closer to $\delta$ even when the shrinkage factor is chosen optimally for $\delta$.

- For a representative values of $\delta$, the set of all $\delta_0$ such that the James-Stein shrinkage estimator $\widehat{\delta}_{JS}^{+}(X; \delta_0)$ lies closer to $\delta$ than the observation $X$.



(a) $\|X - \delta\| < \frac{\sigma}{2}\sqrt{p-2}$

(b) $\|X - \delta\| > \frac{\sigma}{2}\sqrt{p-2}$

- We consider the set which is complement to the cross-hatched region in the above Figure.

- By the directional symmetry of $\delta_0$ about $X$, it follows from the Figure that

$$\Pr[\,\|\hat{\delta}^+_{JS}(X;\delta_0) - \delta\| > \|X - \delta\| \mid X\,] \;>\; \Pr[\,\delta_0 \in H \mid X\,] = \frac{1}{2}. \qquad (10)$$

- If $\delta_0$ is actually symmetrically distributed about $X$ then it is easy to see that

$$\mathsf{E}[\,\hat{\delta}^+_{JS}(X;\delta_0) \mid X] = X, \qquad (11)$$

so by Jensen's inequality,

$$\mathsf{E}[\,\|\hat{\delta}^+_{JS}(X;\delta_0) - \delta\|^2 \mid X\,] \;>\; \mathsf{E}[\,\|X - \delta\|^2 \mid X\,] \equiv p\sigma^2 \quad \forall \delta \in \mathbb{R}^p. \qquad (12)$$

- Furthermore, under additional but still general assumptions,

$$\lim_{p\to\infty} \Pr_\delta[\,\|\hat{\delta}^+_{JS}(X;\delta_0) - \delta\| > \|X - \delta\|\,] = 1. \qquad (13)$$

## THE REVERSE STEIN EFFECT (CONTD.)

- Don't (10) and (12) contradict (8) and (7)? For example, under *any* probability distribution for $\delta_0$, (7) yields

$$\mathsf{E}_\delta[\,\|\hat{\delta}_{JS}^+(X;\delta_0) - \delta\|^2\,] \;<\; \mathsf{E}_\delta[\,\|X - \delta\|^2\,] \equiv p\sigma^2 \quad \forall \delta \in \mathbb{R}^p, \qquad (14)$$

  while (12) yields

$$\mathsf{E}_\delta[\,\|\hat{\delta}_{JS}^+(X;\delta_0) - \delta\|^2\,] \;>\; \mathsf{E}_\delta[\,\|X - \delta\|^2\,] \equiv p\sigma^2 \quad \forall \delta \in \mathbb{R}^p. \qquad (15)$$

- So do we shrink or not?

# TO SHRINK OR NOT TO SHRINK – THAT IS THE QUESTION

- Actually, no formal contradiction has occurred.

- The probabilities and expectations appearing in (7), (8), (10), and (12) are conditional probabilities and conditional expectations with different conditioning variables.

- Furthermore, the joint distributions of $(X, \delta_0)$ in (14) and (15) are different, having joint pdfs of the forms $f_\delta(X)f(\delta_0)$ and $f_\delta(X)f(\delta_0|X)$, respectively.

- In the former, $X$ and $\delta_0$ are independent whereas in the latter, $\delta_0$ is dependent.

## THE REVERSE STEIN EFFECT IS REAL

- The Reverse Stein Effect is just as real as the original Stein Effect itself − both are simply manifestations of the strong curvature of spheres in multi-dimensional Euclidean space.

- The figures and the results (10), (12), and (13) show that without some prior knowledge of the location $\delta$, one should not shrink $X$.

- If the shrinkage target $\delta_0$ is chosen without reliable prior information but instead is based upon the data $X$, the minimax/Bayesian robustness property (B) of the JS estimator is lost and no longer guarantees that shrinking is not harmful on average.

- A shrinkage estimator is only as good as, *but no better than*, the prior information upon which it is based. Without reliable *prior*, as opposed to *posterior* (as represented, for example, by "data-dependent" priors), information, *shrinkage is likely to decrease the accuracy of estimation.*

- if the statistical estimation problem is truly invariant under translation then the best invariant estimator should be used, namely $X$ itself.