

# The James-Stein Estimator and Shrinkage Review

Gary Pai

December 2025

The James-Stein estimator (JSE) [1][2] is a revolutionary concept in statistics that addresses **Stein's Paradox** and utilizes the technique of **shrinkage** to produce a better overall estimator than the standard Maximum Likelihood Estimator (MLE) in multivariate settings.

## 1 Stein's Paradox and Admissibility

When simultaneously estimating the means of three or more independent random variables ( $p \geq 3$ ), the standard sample mean  $\mathbf{X}$  is **inadmissible** under squared-error loss. This means a uniformly better estimator exists.

- **Problem Setting:** We observe a vector of sample means  $\mathbf{X} = (X_1, \dots, X_p)^T$ , where typically  $X_i \sim N(\mu_i, \sigma^2)$ . We want to estimate the true mean vector  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^T$ .
- **Loss Function (Risk):** The quality of an estimator  $\hat{\boldsymbol{\mu}}$  is measured by the Mean Squared Error (MSE), which is the expected loss:

$$R(\boldsymbol{\mu}, \hat{\boldsymbol{\mu}}) = E[\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|^2] = \sum_{i=1}^p E[(\hat{\mu}_i - \mu_i)^2]$$

- **MLE (Inadmissible):** The standard estimator is  $\hat{\boldsymbol{\mu}}_{MLE} = \mathbf{X}$ . Stein showed that for  $p \geq 3$ , there exists an estimator  $\hat{\boldsymbol{\mu}}_{JS}$  such that  $R(\boldsymbol{\mu}, \hat{\boldsymbol{\mu}}_{JS}) < R(\boldsymbol{\mu}, \hat{\boldsymbol{\mu}}_{MLE})$  for all  $\boldsymbol{\mu}$ .

## 2 The James-Stein Estimator (JSE) Formula

The JSE shrinks the vector of estimates  $\mathbf{X}$  towards a central point (often the origin  $\mathbf{0}$ ) to reduce the overall risk.

### 2.1 The General Formula (for known $\sigma^2$ )

The James-Stein estimator for the mean vector  $\boldsymbol{\mu}$  is given by:

$$\hat{\boldsymbol{\mu}}_{JS} = \left(1 - \frac{(p-2)\sigma^2}{\|\mathbf{X}\|^2}\right) \mathbf{X}$$

Where  $\|\mathbf{X}\|^2 = \sum_{i=1}^p X_i^2$ .

### 2.2 The Shrinkage Factor

The term  $c = \left(1 - \frac{(p-2)\sigma^2}{\|\mathbf{X}\|^2}\right)$  is the shrinkage factor. It determines the amount of pulling applied to the individual estimates. This shrinkage works by trading a small amount of **bias** for a large reduction in variance, leading to a lower total MSE.

### 2.3 The Positive-Part JSE (Standardized, $\sigma^2 = 1$ )

A commonly cited and superior version, which prevents over-shrinkage (reversing the sign), is the Positive-Part JSE, often presented with  $\sigma^2 = 1$ :

$$\hat{\boldsymbol{\mu}}_{JS}^+ = \max\left(0, 1 - \frac{p-2}{\|\mathbf{X}\|^2}\right) \mathbf{X}$$

### 3 The Baseball Batting Averages Example

The classic example involves estimating the true batting ability ( $\theta$ ) of  $p = 18$  players based on their initial averages ( $y$ ) in Efron's paper [3]. The JSE combines the individual estimates with the overall grand average ( $\bar{y}$ ).

#### 3.1 The Shrinkage Mechanism

The JSE is often calculated in the form that shrinks toward the grand mean  $\bar{y}$ :

$$\hat{\theta}_{JS,i} = \bar{y} + c \cdot (y_i - \bar{y})$$

Where the shrinkage factor  $c$  is:

$$c = 1 - \frac{(p-3)\sigma^2}{\sum_{i=1}^p (y_i - \bar{y})^2}$$

#### 3.2 Numerical Illustration (Example Player)

For a player with an initial average  $y_i = 0.400$  and a grand average  $\bar{y} = 0.265$ , using an estimated shrinkage factor  $c \approx 0.212$ :

$$\begin{aligned}\hat{\theta}_{JS,i} &= 0.265 + 0.212 \cdot (0.400 - 0.265) \\ &\approx 0.265 + 0.028 \\ &\approx 0.294\end{aligned}$$

The JSE pulls the high initial average of 0.400 down to 0.294, which was found to be much closer to the player's true season-long average. This demonstrates the power of combining information across multiple seemingly independent estimates.

### 4 Modern Relevance

The principle of shrinkage pioneered by James and Stein is fundamental to modern statistical methods and machine learning, particularly in high-dimensional data settings where estimation risk is a major concern. Techniques like **Ridge Regression** and **LASSO** are direct descendants of this concept, utilizing regularization to shrink coefficients toward zero, thereby improving predictive accuracy by managing the bias-variance trade-off.

### References

- [1] Stein, C. (1956). Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, 1, 197–206.
- [2] James, W., & Stein, C. (1961). Estimation with quadratic loss. *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, 1, 361–379.
- [3] Efron, B., & Morris, C. (1975). Data analysis using Stein's estimator and its generalizations. *Journal of the American Statistical Association*, 70(350), 311–319.