



Discussion on “Data Fission: Splitting a Single Data Point”

Zhigen Zhao

To cite this article: Zhigen Zhao (2025) Discussion on “Data Fission: Splitting a Single Data Point”, *Journal of the American Statistical Association*, 120:549, 178-179, DOI: [10.1080/01621459.2024.2416913](https://doi.org/10.1080/01621459.2024.2416913)

To link to this article: <https://doi.org/10.1080/01621459.2024.2416913>



Published online: 14 Apr 2025.



Submit your article to this journal



Article views: 85



View related articles



View Crossmark data

Discussion on “Data Fission: Splitting a Single Data Point”

Zhigen Zhao

Department of Statistics, Operations, and Data Science, Temple University, Philadelphia, PA

I would like to commend the authors for this insightful and thought-provoking paper on “Data FissionSplitting a Single Data Point.” The idea of splitting a data point using random noise to create two parts with specific distributional properties is both innovative and widely applicable, particularly in post-selection inference. The purpose of post-selection inference is to draw from the data and derived empirical Bayes confidence intervals conclusions about the parameters selected through a statistical procedure, denoted $S(\mathbf{y})$. The selection set can be chosen based

on the full dataset, data splitting, or data fission. In this discussion, we would like to comment on the interplay between various

selection rules and post-selection inference methods from the perspective of the Bayes/empirical Bayes framework.

In post-inference problems the main challenge is dealing with “selection bias,” or the “winner’s curse.” For example, consider the normal mean problem discussed in Section 3 of this article,

The interval $M\mathbf{y}_i + (1 - M)\theta \pm z_{\alpha/2} \sqrt{M\sigma^2}$ ensures a good coverage probability for μ_i for any $i \in S(\mathbf{y})$, given this specific choice of θ and σ^2 . When these two hyperparameters are unknown, Hwang and Zhao (2012) and Hwang and Zhao (2013) estimated them through a statistic $\hat{\theta}$ that

EP($\mu_i \in C_{EB}(\mathbf{y})$) $\leq 1 - \alpha$, for any normal prior $N(\theta, \sigma^2)$

In this article, the authors introduce the concept of data fission, where the data is split into $f(\mathbf{y})$ and $g(\mathbf{y})$, with $f(\mathbf{y}_i)$ used for variable selection and $g(\mathbf{y})$ for further inference. For

the normal mean model discussed above the data is split as $f(\mathbf{y}_i) = \mathbf{y}_i + \tau z_i$ and $g(\mathbf{y}) = \mathbf{y}_i - \frac{1}{\tau} z_i$, where $z \sim N(0, \sigma^2)$ is used for inference. The constructed confidence interval is given as

$$y_i \sim N(\mu_i, \sigma^2), i = 1, 2, \dots, p. \quad (1)$$

$$g(\mathbf{y}_i) \pm z_{\alpha/2} \sigma \sqrt{1 + \frac{1}{\tau^2}}.$$

The parameter corresponding to the largest statistic $y_{(n)}$, is of particular interest and is selected for further inference. Directly

using the largest observed value $y_{(n)}$, tends to overestimate this selected parameter leading to low coverage probabilities when $S(\mathbf{y})$ depends solely on when constructing confidence intervals of the form $y_{(n)} \pm z_{\alpha/2} \sigma$. Benjamini and Yekutieli (2005) addressed this issue by adjusting

the confidence interval to account for the number of selected parameters, $R = |S(\mathbf{y})|$, and using a modified confidence coefficient $\alpha R/p$, which tends to be conservative, especially when R is small. In contrast, Dawid (1994) argued that Bayesian inference is immune to selection bias as conditioning on the selection

becomes redundant when the conditioning already applies to the full dataset. In particular, when assuming the prior distribution

$$\mu_i \sim N(\theta, \delta^2), \pi(\theta) \propto 1. \quad (3)$$

it is seen that

$$\mu_i \sim N(\theta, \delta^2), \quad (2)$$

Note that $\mu_i \sim N(\theta, \sigma^2 + \delta^2)$. If we replace θ by $\frac{\sigma^2 + \delta^2}{\sigma^2} z_i$ where $z \sim N(0, \sigma^2)$, then

$$\mu_i | \mathbf{y}, S(\mathbf{y}) = \mu_i | \mathbf{y} \sim N(M\mathbf{y}_i + (1 - M)\theta, M\sigma^2),$$

$$M = \frac{\delta^2}{\sigma^2 + \delta^2}.$$

$$E(\mu_i | \mathbf{y}_i, \theta) \approx y_i + (1 - M)(\theta - y_i),$$

$$E(\mu_i | \mathbf{y}_i, \theta) \approx y_i - \frac{\sigma^2}{\sigma^2 + \delta^2} z_i,$$

which can be viewed as $g(\cdot)$ in the data fission approach when its generalization and potential extensions. I would like to once again congratulate the authors on this exciting paper. I look forward to the authors' rejoinder and to seeing further development of Dawid (1994), it seems that the intervals based on $g(\cdot)$ would provide a valid method for an arbitrary selection rule when assuming the prior distribution (3). In summary,

- When $S(\mathbf{y})$ is chosen arbitrarily, Benjamini and Yekutieli (2005) (BY) method works for any prior distribution, though it requires a longer interval. A much shorter empirical Bayes confidence interval (CI) proposed in Zhao and Hwang (2012) and Hwang and Zhao (2013) guarantees good coverage probabilities for the normal prior (2).
- When $S(\mathbf{y})$ is chosen arbitrarily, the above discussion suggests that the interval centered around $g(\cdot)$ and its empirical Bayesian counterpart could be valid for the class of prior distributions (3) that is broader than (2).
- If $S(\mathbf{Y})$ depends on \mathbf{y} through the $f(\cdot|\mathbf{y})$'s which are independent of $g(\cdot)$'s, then inference based on the $f(\cdot|\mathbf{y})$'s is valid for any prior distribution.

The connection between the proposed approach and the Bayes/empirical Bayes framework offers valuable insights into

Funding

Zhigen Zhao is partially supported by the NSF Grant DMS-2311216.

References

Benjamini, Y., and Yekutieli, D. (2005), "False Discovery Rate-Adjusted Multiple Confidence Intervals for Selected Parameters," *Journal of the American Statistical Association*, 100, 711–733. [\[178, 179\]](#)

Dawid, A. P. (1994), "Selection Paradoxes of Bayesian Inference," *Multivariate Analysis and Its Applications, Lecture Notes-Monograph Series*, 24, 211–220. [\[178, 179\]](#)

Hwang, J. T., and Zhao, Z. (2013), "Empirical Bayes Confidence Intervals for Selected Parameters in High Dimension with Application to Microarray Data Analysis," *Journal of the American Statistical Association*, 108, 607–618. [\[178, 179\]](#)

Zhao, Z., and Hwang, J. T. (2012), "Empirical Bayes False Coverate Rate Controlling Confidence Interval," *Journal of the Royal Statistical Society, Series B*, 74, 871–891. [\[178, 179\]](#)