

Post-selection inference: some answers and some questions

Arun Kumar Kuchibhotla

14 September, 2021

Carnegie Mellon University

Table of contents

1. Invalidity of Classical Inference
2. Formulation of the Problem
3. Three Solutions
4. Some Questions

Invalidity of Classical Inference

Example: Invalidity of classical inference under selection

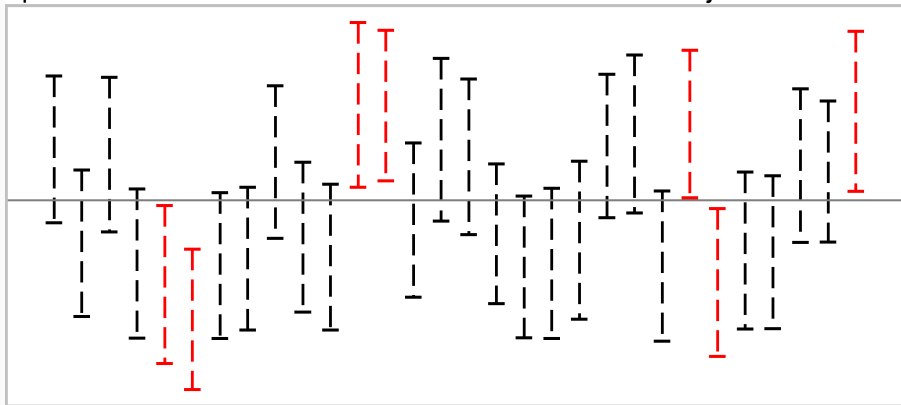
Generate 500 observations from $(X, Y) \sim N(0, I_{p+1})$. ($Y \perp X$)

Select one covariate $X_{\hat{j}}$ that is most correlated with Y .

Coverage of classical **95%** confidence interval

$p = 5$

Unadjusted: **76.9%**



Example: Invalidity of classical inference under selection

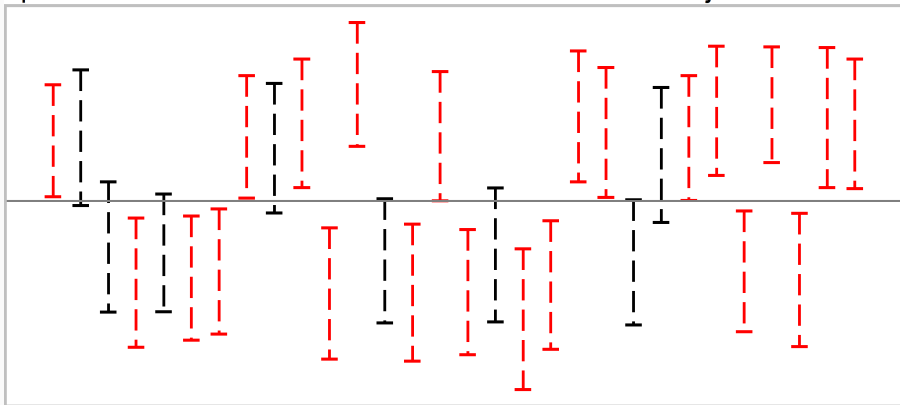
Generate 500 observations from $(X, Y) \sim N(0, I_{p+1})$. ($Y \perp X$)

Select one covariate $X_{\hat{j}}$ that is most correlated with Y .

Coverage of classical **95%** confidence interval

$p = 20$

Unadjusted: **32.6%**



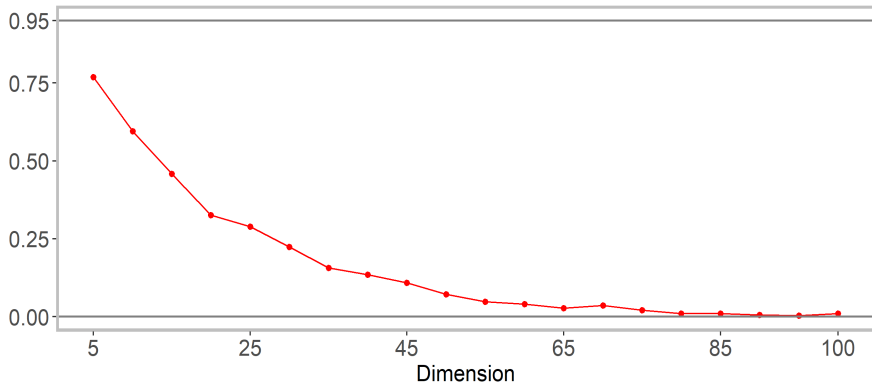
Example: Invalidity of classical inference under selection

Generate 500 observations from $(X, Y) \sim N(0, I_{p+1})$. ($Y \perp X$)

Select one covariate $X_{\hat{j}}$ that is most correlated with Y .

Coverage of classical **95%** confidence interval.

Empirical Coverage



Summary

Unadjusted classical inference can be (very) **misleading**.

Duality of confidence intervals and testing implies that classical tests **may not** control Type I error.

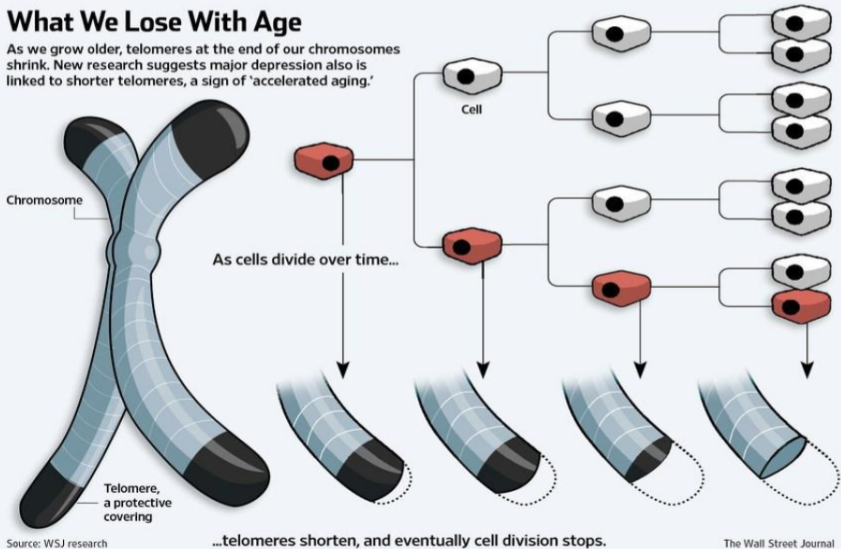
It does not require a pathological selection to invalidate classical inference.

More concerningly, common practice of data exploration is very informal and imprecise.

Case Study 2: Covariate Selection

What We Lose With Age

As we grow older, telomeres at the end of our chromosomes shrink. New research suggests major depression also is linked to shorter telomeres, a sign of 'accelerated aging.'



Source: WSJ research

The Wall Street Journal

Case Study 2: Covariate Selection

nature.com

Scientific Reports, 2019

WGS-based telomere length analysis in Dutch family trios implicates stronger maternal inheritance and a role for *RRM1* gene

“The MLR models were tested by **sequential introduction of predictors and interaction terms.**

• • •

ultimately, from the **three best models** with similar adjusted R squared values the **simplest one was chosen.**”

Formulation of the Problem

There are p hypotheses to start with

$$H_{0,j} : \text{corr}(Y, X_j) = 0, \quad \text{for } 1 \leq j \leq p.$$

Equivalently,

$$H_{0,j} : \beta_j = 0, \quad \text{for } 1 \leq j \leq p,$$

where

$$(\alpha_j, \beta_j) := \underset{(\alpha, \beta)}{\text{argmin}} \mathbb{E} [(Y - \alpha - \beta X_j)^2].$$

There are p hypotheses to start with

$$H_{0,j} : \text{corr}(Y, X_j) = 0, \quad \text{for } 1 \leq j \leq p.$$

Equivalently,

$$H_{0,j} : \beta_j = 0, \quad \text{for } 1 \leq j \leq p,$$

where

$$(\alpha_j, \beta_j) := \underset{(\alpha, \beta)}{\text{argmin}} \mathbb{E} [(Y - \alpha - \beta X_j)^2].$$

Select a $\hat{j} \in \{1, 2, \dots, p\}$ based on the data.

There are p hypotheses to start with

$$H_{0,j} : \text{corr}(Y, X_j) = 0, \quad \text{for } 1 \leq j \leq p.$$

Equivalently,

$$H_{0,j} : \beta_j = 0, \quad \text{for } 1 \leq j \leq p,$$

where

$$(\alpha_j, \beta_j) := \underset{(\alpha, \beta)}{\text{argmin}} \mathbb{E} [(Y - \alpha - \beta X_j)^2].$$

Select a $\hat{j} \in \{1, 2, \dots, p\}$ based on the data.

Test the hypothesis $H_{0,\hat{j}} : \beta_{\hat{j}} = 0$.

Classical (invalid) test:

$$\text{Reject } H_{0,\hat{j}} \text{ if } |t_{\hat{j}}| := \frac{n^{1/2} |\hat{\beta}_{\hat{j}}|}{\hat{\sigma}_{\hat{j}}} \leq 1.96.$$

The General PoSI Problem

For each model $M \subseteq \{1, 2, \dots, p\}$, define the OLS target as

$$\beta_M := \operatorname{argmin}_{\theta \in \mathbb{R}^{|M|}} \mathbb{E} [(Y - X_M^\top \theta)^2].$$

Construct a confidence region $\widehat{\text{CI}}_{\hat{j}, \hat{M}}$ such that

$$\liminf_{n \rightarrow \infty} \mathbb{P} \left(\beta_{\hat{j}, \hat{M}} \in \widehat{\text{CI}}_{\hat{j}, \hat{M}} \right) \geq 1 - \alpha,$$

for any model \hat{M} (of size **at most k**) and $\hat{j} \in \hat{M}$, **irrespective of how it is chosen**.

Simultaneous Inference \Rightarrow Post-selection Inference

$$\mathbb{P} \left(\bigcap_{\substack{|M| \leq k \\ j \in M}} \left\{ \beta_{j \cdot M} \in \widehat{CI}_{j \cdot M} \right\} \right) \leq \inf_{\substack{\widehat{j} \in \widehat{M}, \\ |\widehat{M}| \leq k}} \mathbb{P} \left(\beta_{\widehat{j} \cdot \widehat{M}} \in \widehat{CI}_{\widehat{j} \cdot \widehat{M}} \right).$$

Theorem:

Simultaneous inference is necessary for valid PoSI.

Three Solutions

A (Very) Simple Solution

Apply Bonferroni procedure.

$$\mathbb{P} \left(\bigcap_{\substack{|\mathbb{M}| \leq k \\ j \in \mathbb{M}}} \{ \beta_{j \cdot \mathbb{M}} \in \hat{\mathbb{C}}_{j \cdot \mathbb{M}} \} \right) \geq 1 - \sum_{\substack{|\mathbb{M}| \leq k, \\ j \in \mathbb{M}}} \mathbb{P} \left(\beta_{j \cdot \mathbb{M}} \in \hat{\mathbb{C}}_{j \cdot \mathbb{M}} \right).$$

How many elements in the sum?

$$\sum_{\substack{|\mathbb{M}| \leq k, \\ j \in \mathbb{M}}} 1 = \sum_{s=1}^k s \binom{p}{s} \asymp \left(\frac{ep}{k} \right)^k.$$

Construct $1 - \frac{\alpha}{(ep/k)^k}$ confidence intervals for individual coefficients.

Can be very conservative.

Second Simple Solution

For simultaneous inference, inflate the interval to

$$\widehat{CI}_{j \cdot M} := \left\{ \theta \in \mathbb{R} : \left| \frac{n^{1/2}(\widehat{\beta}_{j \cdot M} - \theta)}{\widehat{\sigma}_{j \cdot M}} \right| \leq K_\alpha \right\},$$

with K_α , the $(1 - \alpha)$ quantile of

$$\max_{|M| \leq k, j \in M} \left| \frac{n^{1/2}(\widehat{\beta}_{j \cdot M} - \beta_{j \cdot M})}{\widehat{\sigma}_{j \cdot M}} \right|.$$

Accounts for dependence.

Disadvantage of these Solutions

Bonferroni Solution:

$$\widehat{CI}_{j \cdot M}^{\text{Bonf}} := \left\{ \theta \in \mathbb{R} : \left| \frac{n^{1/2}(\widehat{\beta}_{j \cdot M} - \theta)}{\widehat{\sigma}_{j \cdot M}} \right| \leq z_{\alpha/(2(ep/k)^k)} \right\}.$$

PoSI Solution:

$$\widehat{CI}_{j \cdot M}^{\text{PoSI}} := \left\{ \theta \in \mathbb{R} : \left| \frac{n^{1/2}(\widehat{\beta}_{j \cdot M} - \theta)}{\widehat{\sigma}_{j \cdot M}} \right| \leq K_{\alpha} \right\}.$$

K_{α} usually grows with largest model size k .

Say, $k = 20$, then

width of intervals for model of size 2

\approx

width of intervals for model of size 20.

The Third Solution

Define

$$\text{(OLS Estimator)} \quad \hat{\beta}_M := \underset{\theta \in \mathbb{R}^{|M|}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (Y_i - X_{i,M}^\top \theta)^2,$$

$$\text{(OLS Target)} \quad \beta_M := \underset{\theta \in \mathbb{R}^{|M|}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \mathbb{E} [(Y_i - X_{i,M}^\top \theta)^2].$$

For any $M \subseteq \{1, 2, \dots, p\}$, consider the confidence region

$$\hat{\text{CI}}_M^{\text{UPoSI}^*} := \left\{ \theta \in \mathbb{R}^{|M|} : \|\hat{\Sigma}_M(\hat{\beta}_M - \theta)\|_\infty \leq C_{xy}(\alpha) + C_{xx}(\alpha) \|\theta\|_1 \right\}.$$

Then for any model \hat{M} chosen based on the data,

$$\mathbb{P} \left(\beta_{\hat{M}} \in \hat{\text{CI}}_{\hat{M}}^{\text{UPoSI}^*} \right) \geq 1 - \alpha,$$

if $C_{xy}(\alpha)$ and $C_{xx}(\alpha)$ denote the $(1 - \alpha)$ joint quantiles of

$$\left\| \frac{1}{n} \sum_{i=1}^n \{X_i Y_i - \mathbb{E}[X_i Y_i]\} \right\|_\infty \quad \text{and} \quad \left\| \frac{1}{n} \sum_{i=1}^n \{X_i X_i^\top - \mathbb{E}[X_i X_i^\top]\} \right\|_\infty.$$

Comparison of Volumes

Reference	$\text{Leb}(\widehat{\text{CI}}_{\widehat{M}})$	Design
Kuchibhotla et al. (2021, AoS)	$(\log p/n)^{ \widehat{M} /2}$	fixed
	$(\widehat{M} \log p/n)^{ \widehat{M} /2}$	random
Berk et al. (2013)	$(k \log(ep/k)/n)^{ \widehat{M} /2}$	fixed/random
Bachoc et al. (2019)		
Kuchibhotla et al. (2021, Econ. Theory)		

Table 1: Volumes of Different PoSI Regions.

Some Questions

Question 1: Optimality of volume

- We have shown that simultaneous inference is necessary and sufficient for post-selection inference in the problem formulation here.
- If arbitrary model selection procedures are allowed and the analyst chooses the covariate subset \widehat{M} , then what is **the best possible** confidence set for $\beta_{\widehat{M}}$ in terms of volume?
- The current best available is $(|\widehat{M}| \log(ep/|\widehat{M}|)/n)^{|\widehat{M}|/2}$. This follows from the third solution above and also the Hierarchical PoSI proposed in Kuchibhotla (2020, PhD Thesis).
- From the literature of sparse high-dimensional linear regression, this volume seems to be the best possible.

What is the smallest volume confidence set for $\beta_{\widehat{M}}$ with an (asymptotic) coverage validity for arbitrary selection of \widehat{M} ?

Question 1: Optimality of volume (Contd.)

- A simpler problem to consider is the normal means setting.
- Suppose $X \sim N(\mu, \Sigma)$ in \mathbb{R}^d with known Σ . Let $\mu = (\mu_1, \dots, \mu_d)^\top$.
- Consider the problem of constructing a confidence interval for $\mu_{\hat{j}}$ for a data-based selection $\hat{j} \in \{1, 2, \dots, d\}$.
- For any fixed $1 \leq j \leq d$, a valid confidence interval for μ_j has width that scales as $\sqrt{\text{Var}(X_j)}$.
- For a data-based selection \hat{j} , one can construct confidence intervals of width of order $\sqrt{\log(\hat{j})} \max_j \sqrt{\text{Var}(X_j)}|_{j=\hat{j}}$.
- It is currently unknown what the smallest volume confidence set for this simpler setting.

Question 2: Dynamic PoSI

- The post-selection inference problem posed is written with respect to a universe.
- For example, in linear regression with covariate selection, we have the universe of selection to be set of all subsets of covariates.
- In practice, this is not how data-drive selection is done. The second model selection method is informed by the output of the first model selection method, i.e., **dynamic selection**.
- For example, one might decide to (univariate) marginal screening first, then depending how many p-values are small might decide where to threshold the p-values for selecting covariates.
- Without any control on either the dynamics or the output of selection, PoSI is *impossible*.
- Dynamic PoSI via randomized output is achieved by adaptive data analysis and stable algorithms. But these methods require specific model selection methods and cannot handle visual selection.

Question 2: Dynamic PoSI (Contd.)

- Is there a generic way to modify the output of an arbitrary model selection strategy to solve dynamic PoSI?
- This includes those selections based on residual plots, QQ plots, and so on.
- At present, a very crude solution is available where by the underlying data is randomized to solve this. But this results in a huge lose of information.
- Furthermore, the methods that are more precise randomize functions of data with the amount of randomization depending on underlying distribution parameters that are seldom available.
- For example, in differential privacy, Laplace mechanism noise depends on the sensitivity of the function of data which is often unavailable.

Question 3: Computable PoSI

- The problem of post-selection inference as posed with respect to a universe of selection has several solutions now.
- Except in the linear regression case, there is no computable solution, i.e., some inference procedure that is not NP-hard.
- This is somewhat similar to the ℓ_0 (sparsity) penalized linear regression.
- The statistics used to compute simultaneous confidence intervals for PoSI at present involve maximum over the universe.
- In the context of covariate selection, this is the maximum over all (sparse) subsets of covariates, which is NP-hard to compute.

Is there a computable solution to PoSI for general linear models?

Question 3: Computable PoSI (Contd.)

- Currently, there is work in progress that provides a randomized algorithm to compute the approximate maximum over the universe of selection.
- These approximations yield valid, albeit conservative confidence intervals for PoSI.
- It is not obvious what would be the best approximation to the maximum for best performance in practice.
- It is also not clear if there is a trade-off between conservativeness and computation.
- Computation is the main hurdle for the practical use of the existing PoSI solutions at present.

Conclusions

Post-selection inference is an important topic of research given the current reproducibility crisis.

There exist several solutions in the current formulation of the PoSI problem.

Their practical use is riddled with several obstacles and resolving them would lead to (more) trustworthy practical data analysis.

Reference: Kuchibhotla et al. (2021) [Valid Post-selection Inference in Model-free Linear Regression](#), Annals of Statistics.

Conclusions

Post-selection inference is an important topic of research given the current reproducibility crisis.

There exist several solutions in the current formulation of the PoSI problem.

Their practical use is riddled with several obstacles and resolving them would lead to (more) trustworthy practical data analysis.

Reference: Kuchibhotla et al. (2021) [Valid Post-selection Inference in Model-free Linear Regression](#), Annals of Statistics.

Thank you!