

“PoSI” — Valid Post-Selection Inference

Andreas Buja

joint work with

Richard Berk, Lawrence Brown, Kai Zhang, Linda Zhao

Department of Statistics, The Wharton School
University of Pennsylvania
Philadelphia, USA

UF 2014/01/18

Larger Problem: Non-Reproducible Empirical Findings

- Indicators of a problem

(from: Berger, 2012, “Reproducibility of Science: P-values and Multiplicity”)

Larger Problem: Non-Reproducible Empirical Findings

- Indicators of a problem

(from: Berger, 2012, “Reproducibility of Science: P-values and Multiplicity”)

- ▶ Bayer Healthcare reviewed 67 in-house attempts at replicating findings in published research:
 - < 1/4 were viewed as replicated.
- ▶ Arrowsmith (2011, Nat. Rev. Drug Discovery 10):
 - Increasing failure rate in Phase II drug trials

Larger Problem: Non-Reproducible Empirical Findings

- Indicators of a problem

(from: Berger, 2012, “Reproducibility of Science: P-values and Multiplicity”)

- ▶ Bayer Healthcare reviewed 67 in-house attempts at replicating findings in published research:
 - < 1/4 were viewed as replicated.
- ▶ Arrowsmith (2011, Nat. Rev. Drug Discovery 10):
 - Increasing failure rate in Phase II drug trials
- ▶ Ioannidis (2005, PLOS Medicine):
 - “Why Most Published Research Findings Are False”

Larger Problem: Non-Reproducible Empirical Findings

- Indicators of a problem

(from: Berger, 2012, “Reproducibility of Science: P-values and Multiplicity”)

- ▶ Bayer Healthcare reviewed 67 in-house attempts at replicating findings in published research:
 - < 1/4 were viewed as replicated.
- ▶ Arrowsmith (2011, Nat. Rev. Drug Discovery 10):
 - Increasing failure rate in Phase II drug trials
- ▶ Ioannidis (2005, PLOS Medicine):
 - “Why Most Published Research Findings Are False”
- ▶ Simmons, Nelson, Simonsohn (2011, Psychol.Sci):
 - “False-Positive Psychology: Undisclosed Flexibility in Data Collection and Analysis Allows Presenting Anything as Significant,”

Larger Problem: Non-Reproducible Empirical Findings

- Indicators of a problem

(from: Berger, 2012, “Reproducibility of Science: P-values and Multiplicity”)

- ▶ Bayer Healthcare reviewed 67 in-house attempts at replicating findings in published research:
 - < 1/4 were viewed as replicated.
- ▶ Arrowsmith (2011, Nat. Rev. Drug Discovery 10):
 - Increasing failure rate in Phase II drug trials
- ▶ Ioannidis (2005, PLOS Medicine):
 - “Why Most Published Research Findings Are False”
- ▶ Simmons, Nelson, Simonsohn (2011, Psychol.Sci):
 - “False-Positive Psychology: Undisclosed Flexibility in Data Collection and Analysis Allows Presenting Anything as Significant,”

- Many potential causes – two major ones:

Larger Problem: Non-Reproducible Empirical Findings

- Indicators of a problem

(from: Berger, 2012, “Reproducibility of Science: P-values and Multiplicity”)

- ▶ Bayer Healthcare reviewed 67 in-house attempts at replicating findings in published research:
 - < 1/4 were viewed as replicated.
- ▶ Arrowsmith (2011, Nat. Rev. Drug Discovery 10):
 - Increasing failure rate in Phase II drug trials
- ▶ Ioannidis (2005, PLOS Medicine):
 - “Why Most Published Research Findings Are False”
- ▶ Simmons, Nelson, Simonsohn (2011, Psychol.Sci):
 - “False-Positive Psychology: Undisclosed Flexibility in Data Collection and Analysis Allows Presenting Anything as Significant,”

- Many potential causes – two major ones:

- ▶ publication bias: “file drawer problem” (Rosenthal 1979)
- ▶ statistical biases: “researcher degrees of freedom” (SNS 2011)

Statistical Biases – one among several

- Hypothesis: A statistical bias is due to
an absence of accounting for model/variable selection.

Statistical Biases – one among several

- Hypothesis: A statistical bias is due to
an absence of accounting for model/variable selection.
- Model selection is done on several levels:

Statistical Biases – one among several

- Hypothesis: A statistical bias is due to
an absence of accounting for model/variable selection.
- Model selection is done on several levels:
 - ▶ formal selection: AIC, BIC, Lasso, ...

Statistical Biases – one among several

- Hypothesis: A statistical bias is due to
an absence of accounting for model/variable selection.
- Model selection is done on several levels:
 - ▶ formal selection: AIC, BIC, Lasso, ...
 - ▶ informal selection: residual plots, influence diagnostics, ...

Statistical Biases – one among several

- Hypothesis: A statistical bias is due to
an absence of accounting for model/variable selection.
- Model selection is done on several levels:
 - ▶ **formal selection**: AIC, BIC, Lasso, ...
 - ▶ **informal selection**: residual plots, influence diagnostics, ...
 - ▶ **post hoc selection**: “The effect size is too small in relation to the cost of data collection to warrant inclusion of this predictor.”

Statistical Biases – one among several

- Hypothesis: A statistical bias is due to
an absence of accounting for model/variable selection.
- Model selection is done on several levels:
 - ▶ **formal selection**: AIC, BIC, Lasso, ...
 - ▶ **informal selection**: residual plots, influence diagnostics, ...
 - ▶ **post hoc selection**: “The effect size is too small in relation to the cost of data collection to warrant inclusion of this predictor.”
- Suspicions:

Statistical Biases – one among several

- Hypothesis: A statistical bias is due to
an absence of accounting for model/variable selection.
- Model selection is done on several levels:
 - ▶ **formal selection**: AIC, BIC, Lasso, ...
 - ▶ **informal selection**: residual plots, influence diagnostics, ...
 - ▶ **post hoc selection**: “The effect size is too small in relation to the cost of data collection to warrant inclusion of this predictor.”
- Suspensions:
 - ▶ All three modes of model selection may be used in much empirical research.

Statistical Biases – one among several

- Hypothesis: A statistical bias is due to
an absence of accounting for model/variable selection.
- Model selection is done on several levels:
 - ▶ **formal selection**: AIC, BIC, Lasso, ...
 - ▶ **informal selection**: residual plots, influence diagnostics, ...
 - ▶ **post hoc selection**: “The effect size is too small in relation to the cost of data collection to warrant inclusion of this predictor.”
- Suspensions:
 - ▶ All three modes of model selection may be used in much empirical research.
 - ▶ Ironically, the most thorough and competent data analysts may also be the ones who produce the most spurious findings.

Statistical Biases – one among several

- Hypothesis: A statistical bias is due to
an absence of accounting for model/variable selection.
- Model selection is done on several levels:
 - ▶ **formal selection**: AIC, BIC, Lasso, ...
 - ▶ **informal selection**: residual plots, influence diagnostics, ...
 - ▶ **post hoc selection**: “The effect size is too small in relation to the cost of data collection to warrant inclusion of this predictor.”
- Suspicions:
 - ▶ All three modes of model selection may be used in much empirical research.
 - ▶ Ironically, the most thorough and competent data analysts may also be the ones who produce the most spurious findings.
 - ▶ If we develop valid post-selection inference for “adaptive Lasso”, say, it won’t solve the problem because few empirical researchers would commit themselves **a priori** to **one formal** selection method and nothing else.
⇒ “Meta-Selection Problem”

The Problem of Post-Selection Inference

How can Variable Selection **invalidate** Conventional Inference?

The Problem of Post-Selection Inference

How can Variable Selection **invalidate** Conventional Inference?

- Conventional inference after variable selection ignores the fact that the model was obtained through a **stochastic selection process**.

The Problem of Post-Selection Inference

How can Variable Selection **invalidate** Conventional Inference?

- Conventional inference after variable selection ignores the fact that the model was obtained through a **stochastic selection process**.
- Stochastic variable selection **distorts sampling distributions** of the post-selection parameter estimates: Most selection procedures search for **strong**, hence **highly significant looking** predictors.

The Problem of Post-Selection Inference

How can Variable Selection **invalidate** Conventional Inference?

- Conventional inference after variable selection ignores the fact that the model was obtained through a **stochastic selection process**.
- Stochastic variable selection **distorts sampling distributions** of the post-selection parameter estimates: Most selection procedures search for **strong**, hence **highly significant looking** predictors.
- Some forms of the problem has been known for decades:
Koopmans (1949); Buehler and Feddersen (1963); Brown (1967); and Olshen (1973); Sen (1979); Sen and Saleh (1987); Dijkstra and Veldkamp (1988); Arabatzis et al. (1989); Hurvich and Tsai (1990); Regal and Hook (1991); Pötscher (1991); Chiou and Han (1995a,b); Giles (1992); Giles and Srivastava (1993); Kabaila (1998); Brockwell and Gordeon (2001); Leeb and Pötscher (2003; 2005; 2006a; 2006b; 2008a; 2008b); Kabaila (2005); Kabaila and Leeb (2006); Berk, Brown and Zhao (2009); Kabaila (2009).

Example: Length of Criminal Sentence

Question: What covariates predict length of a criminal sentence best?

A small empirical study:

- $N = 250$ observations.

Example: Length of Criminal Sentence

Question: What covariates predict length of a criminal sentence best?

A small empirical study:

- $N = 250$ observations.
- Response: log-length of sentences

Example: Length of Criminal Sentence

Question: What covariates predict length of a criminal sentence best?

A small empirical study:

- $N = 250$ observations.
- Response: log-length of sentences
- $p = 11$ covariates (predictors, explanatory variables):
 - ▶ race
 - ▶ gender
 - ▶ initial age
 - ▶ marital status
 - ▶ employment status
 - ▶ seriousness of crime
 - ▶ psychological problems
 - ▶ education
 - ▶ drug related
 - ▶ alcohol usage
 - ▶ prior record

Example: Length of Criminal Sentence

Question: What covariates predict length of a criminal sentence best?

A small empirical study:

- $N = 250$ observations.
- Response: log-length of sentences
- $p = 11$ covariates (predictors, explanatory variables):
 - ▶ race
 - ▶ gender
 - ▶ initial age
 - ▶ marital status
 - ▶ employment status
 - ▶ seriousness of crime
 - ▶ psychological problems
 - ▶ education
 - ▶ drug related
 - ▶ alcohol usage
 - ▶ prior record
- What variables should be included?

Example: Length of Criminal Sentence (contd.)

- All-subset search with BIC chooses a model \hat{M} with seven variables:
 - ▶ initial age
 - ▶ gender
 - ▶ employment status
 - ▶ seriousness of crime
 - ▶ drugs related
 - ▶ alcohol usage
 - ▶ prior records

Example: Length of Criminal Sentence (contd.)

- All-subset search with BIC chooses a model \hat{M} with seven variables:

- ▶ initial age
- ▶ gender
- ▶ employment status
- ▶ seriousness of crime
- ▶ drugs related
- ▶ alcohol usage
- ▶ prior records

- t -statistics of selected covariates, in descending order:

- ▶ $|t_{\text{alcohol}}| = 3.95;$
- ▶ $|t_{\text{prior records}}| = 3.59;$
- ▶ $|t_{\text{seriousness}}| = 3.57;$
- ▶ $|t_{\text{drugs}}| = 3.31;$
- ▶ $|t_{\text{employment}}| = 3.04;$
- ▶ $|t_{\text{initial age}}| = 2.56;$
- ▶ $|t_{\text{gender}}| = 2.33.$

Example: Length of Criminal Sentence (contd.)

- All-subset search with BIC chooses a model \hat{M} with seven variables:
 - ▶ initial age
 - ▶ gender
 - ▶ employment status
 - ▶ seriousness of crime
 - ▶ drugs related
 - ▶ alcohol usage
 - ▶ prior records
- t -statistics of selected covariates, in descending order:
 - ▶ $|t_{\text{alcohol}}| = 3.95;$
 - ▶ $|t_{\text{prior records}}| = 3.59;$
 - ▶ $|t_{\text{seriousness}}| = 3.57;$
 - ▶ $|t_{\text{drugs}}| = 3.31;$
 - ▶ $|t_{\text{employment}}| = 3.04;$
 - ▶ $|t_{\text{initial age}}| = 2.56;$
 - ▶ $|t_{\text{gender}}| = 2.33.$
- Can we use the cutoff $t_{.975, 250-8} = 1.97$?

Linear Model Inference and Variable Selection

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon$$

- \mathbf{X} = fixed design matrix, $N \times p$, $N > p$, full rank.
- $\epsilon \sim \mathcal{N}_N(\mathbf{0}, \sigma^2 \mathbf{I}_N)$

In textbooks:

- 1 Variables selected
- 2 Data seen
- 3 Inference produced

In common practice:

- 1 Data seen
- 2 Variables selected
- 3 Inference produced

Linear Model Inference and Variable Selection

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon$$

- \mathbf{X} = fixed design matrix, $N \times p$, $N > p$, full rank.
- $\epsilon \sim \mathcal{N}_N(\mathbf{0}, \sigma^2 \mathbf{I}_N)$

In textbooks:

- 1 Variables selected
- 2 Data seen
- 3 Inference produced

In common practice:

- 1 Data seen
- 2 Variables selected
- 3 Inference produced

Is this inference valid?

Evidence from a Simulation

Generate Y from the following linear model:

$$Y = \beta x + \sum_{j=1}^{10} \gamma_j z_j + \epsilon,$$

where $p = 11$, $N = 250$, and $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ iid.

► More Details

Evidence from a Simulation

Generate Y from the following linear model:

$$Y = \beta x + \sum_{j=1}^{10} \gamma_j z_j + \epsilon,$$

where $p = 11$, $N = 250$, and $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ iid.

► More Details

- For simplicity: “Protect” x and select only among z_1, \dots, z_{10} ; interest is in inference for β .

Evidence from a Simulation

Generate Y from the following linear model:

$$Y = \beta x + \sum_{j=1}^{10} \gamma_j z_j + \epsilon,$$

where $p = 11$, $N = 250$, and $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ iid.

► More Details

- For simplicity: “Protect” x and select only among z_1, \dots, z_{10} ; interest is in inference for β .
- Model selection: All-subset search with BIC among z_1, \dots, z_{10} ; always including x .

Evidence from a Simulation

Generate Y from the following linear model:

$$Y = \beta x + \sum_{j=1}^{10} \gamma_j z_j + \epsilon,$$

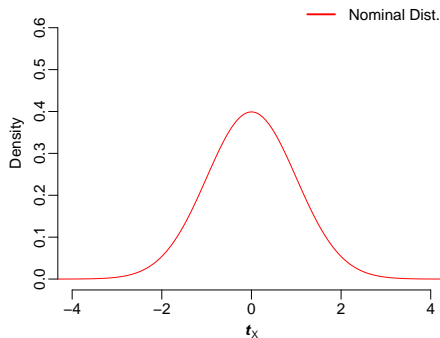
where $p = 11$, $N = 250$, and $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ iid.

► More Details

- For simplicity: “Protect” x and select only among z_1, \dots, z_{10} ; interest is in inference for β .
- Model selection: All-subset search with BIC among z_1, \dots, z_{10} ; always including x .
- Proper coverage of a 95% CI on the slope β of x under the chosen model requires that the t -statistic is about $\mathcal{N}(0, 1)$ distributed.

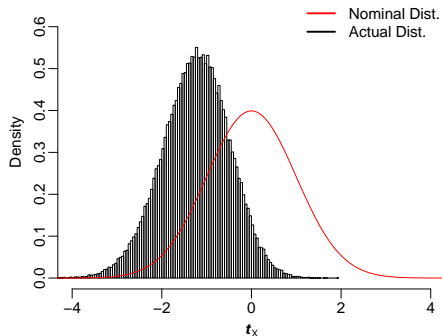
Evidence from a Simulation (contd.)

Marginal Distribution of Post-Selection t -statistics:



Evidence from a Simulation (contd.)

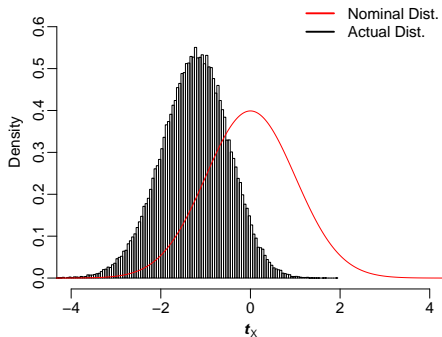
Marginal Distribution of Post-Selection t -statistics:



- The overall coverage probability of the conventional post-selection CI is $83.5\% < 95\%$.

Evidence from a Simulation (contd.)

Marginal Distribution of Post-Selection t -statistics:



- The overall coverage probability of the conventional post-selection CI is **83.5% < 95%**.
- For $p = 30$, the coverage probability can be as low as **39%**.

The PoSI Procedure — Rough Outline

- We propose to construct **Post Selection Inference** (PoSI) with guarantees for the coverage of CIs and Type I errors of tests.

The PoSI Procedure — Rough Outline

- We propose to construct **Post Selection Inference** (PoSI) with guarantees for the coverage of CIs and Type I errors of tests.
- We **widen** CIs and retention intervals to achieve correct/conservative post-selection coverage probabilities. This is the **price** we have to pay.

The PoSI Procedure — Rough Outline

- We propose to construct **Post Selection Inference** (PoSI) with guarantees for the coverage of CIs and Type I errors of tests.
- We **widen** CIs and retention intervals to achieve correct/conservative post-selection coverage probabilities. This is the **price** we have to pay.
- The approach is a reduction of PoSI to **simultaneous inference**.

The PoSI Procedure — Rough Outline

- We propose to construct **Post Selection Inference** (PoSI) with guarantees for the coverage of CIs and Type I errors of tests.
- We **widen** CIs and retention intervals to achieve correct/conservative post-selection coverage probabilities. This is the **price** we have to pay.
- The approach is a reduction of PoSI to **simultaneous inference**.
- Simultaneity is across **all submodels** and **all slopes** in them.

The PoSI Procedure — Rough Outline

- We propose to construct **Post Selection Inference** (PoSI) with guarantees for the coverage of CIs and Type I errors of tests.
- We **widen** CIs and retention intervals to achieve correct/conservative post-selection coverage probabilities. This is the **price** we have to pay.
- The approach is a reduction of PoSI to **simultaneous inference**.
- Simultaneity is across **all submodels** and **all slopes** in them.
- As a result, we obtain

valid PoSI for all variable selection procedures!

The PoSI Procedure — Rough Outline

- We propose to construct **Post Selection Inference** (PoSI) with guarantees for the coverage of CIs and Type I errors of tests.
- We **widen** CIs and retention intervals to achieve correct/conservative post-selection coverage probabilities. This is the **price** we have to pay.
- The approach is a reduction of PoSI to **simultaneous inference**.
- Simultaneity is across **all submodels** and **all slopes** in them.
- As a result, we obtain

valid PoSI for all variable selection procedures!

- But first we need some preliminaries on

Targets of Inference and Inference in Wrong Models

PoSI — What's in a word?

<http://www.thefreedictionary.com/posies>

Translations

posy [ˈpəʊzi] *N* → ramillete *m*

posy [ˈpəʊzi] *n* → petit bouquet *m*

posy
n → Sträußchen *nt*

posy [ˈpəʊzi] *n* → mazzolino (di fiori)

posy
n **posy** [ˈpəʊzi]

a small bunch of flowers a *posy* of primroses. ruiker بقة kitka kytic̨ka lille buket der Strauß μπουκετάκι ramillete (vāike) lillekimp نسته گل kukkavihko petit bouquet (de fleurs) פּוֹסֵי פִּרְחִים पुष्पगुच्छ
kitica cvijeća kis csokor seikat bunga blómvöndur mazzolino 花束 꽃다발 puokštele puķu pušķītis
sejambak bunga boeket liten bukett, blomst bukiecik ramalhetę buchetał бучет(ик) цветов
kytička šopek buketič [] buketт ചക്കചെടി küçük çiçek demeti 花束 букетик kwitib جيوٹا گندسہ chùm
hoa nhỏ 花束

Submodels — Notation, Parameters, Assumptions

Submodels — Notation, Parameters, Assumptions

- Denote a submodel by the integers $M = \{j_1, j_2, \dots, j_m\}$ for the predictors:

$$\mathbf{X}_M = (\mathbf{x}_{j_1}, \mathbf{x}_{j_2}, \dots, \mathbf{x}_{j_m}) \in \mathbb{R}^{N \times m}.$$

Submodels — Notation, Parameters, Assumptions

- Denote a submodel by the integers $M = \{j_1, j_2, \dots, j_m\}$ for the predictors:

$$\mathbf{X}_M = (\mathbf{x}_{j_1}, \mathbf{x}_{j_2}, \dots, \mathbf{x}_{j_m}) \in \mathbb{R}^{N \times m}.$$

- The LS estimators in the submodel M are

$$\hat{\beta}_M = (\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{X}_M^T \mathbf{Y} \in \mathbb{R}^m$$

Submodels — Notation, Parameters, Assumptions

- Denote a submodel by the integers $M = \{j_1, j_2, \dots, j_m\}$ for the predictors:

$$\mathbf{X}_M = (\mathbf{x}_{j_1}, \mathbf{x}_{j_2}, \dots, \mathbf{x}_{j_m}) \in \mathbb{R}^{N \times m}.$$

- The LS estimators in the submodel M are

$$\hat{\beta}_M = (\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{X}_M^T \mathbf{Y} \in \mathbb{R}^m$$

- What does $\hat{\beta}_M$ estimate, **not** assuming the truth of M ?

Submodels — Notation, Parameters, Assumptions

- Denote a submodel by the integers $M = \{j_1, j_2, \dots, j_m\}$ for the predictors:

$$\mathbf{X}_M = (\mathbf{X}_{j_1}, \mathbf{X}_{j_2}, \dots, \mathbf{X}_{j_m}) \in \mathbb{R}^{N \times m}.$$

- The LS estimators in the submodel M are

$$\hat{\beta}_M = (\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{X}_M^T \mathbf{Y} \in \mathbb{R}^m$$

- What does $\hat{\beta}_M$ estimate, **not** assuming the truth of M ?

A: Its expectation — i.e., we ask for unbiasedness.

$$\mu := \mathbf{E}[\mathbf{Y}] \in \mathbb{R}^N \quad \text{arbitrary!!}$$

$$\beta_M := \mathbf{E}[\hat{\beta}_M] = (\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{X}_M^T \mu$$

- Once again: We do **not** assume that the submodel is correct, i.e., we allow $\mu \neq \mathbf{X}_M \beta_M$! But $\mathbf{X}_M \beta_M$ is the best approximation to μ .

Submodels versus the Full Model — Confusions

- Abbreviate $\hat{\beta} := \hat{\beta}_{M_F}$ and $\beta := \beta_{M_F}$, $M_F = \{1, 2, \dots, p\}$ = full model.

Questions:

- ▶ How do submodel estimates $\hat{\beta}_M$ relate to full-model estimates $\hat{\beta}$?
- ▶ How do submodel parameters β_M relate to full-model parameters β ?

Is $\hat{\beta}_M$ a subset of $\hat{\beta}$ and β_M a subset of β ?

Submodels versus the Full Model — Confusions

- Abbreviate $\hat{\beta} := \hat{\beta}_{M_F}$ and $\beta := \beta_{M_F}$, $M_F = \{1, 2, \dots, p\}$ = full model.
Questions:

- ▶ How do submodel estimates $\hat{\beta}_M$ relate to full-model estimates $\hat{\beta}$?
- ▶ How do submodel parameters β_M relate to full-model parameters β ?

Is $\hat{\beta}_M$ a subset of $\hat{\beta}$ and β_M a subset of β ?

- Answer: Unless \mathbf{X} is an orthogonal design,
 - ▶ $\hat{\beta}_M$ is **not** a subset of $\hat{\beta}$ and
 - ▶ β_M is **not** a subset of β .

Submodels versus the Full Model — Confusions

- Abbreviate $\hat{\beta} := \hat{\beta}_{M_F}$ and $\beta := \beta_{M_F}$, $M_F = \{1, 2, \dots, p\}$ = full model.
Questions:

- ▶ How do submodel estimates $\hat{\beta}_M$ relate to full-model estimates $\hat{\beta}$?
- ▶ How do submodel parameters β_M relate to full-model parameters β ?

Is $\hat{\beta}_M$ a subset of $\hat{\beta}$ and β_M a subset of β ?

- Answer: Unless \mathbf{X} is an orthogonal design,
 - ▶ $\hat{\beta}_M$ is **not** a subset of $\hat{\beta}$ and
 - ▶ β_M is **not** a subset of β .
- Reason: Slopes, both estimates and parameters, depend on what the other predictors are — **in value and in meaning**.

Submodels versus the Full Model — Confusions

- Abbreviate $\hat{\beta} := \hat{\beta}_{M_F}$ and $\beta := \beta_{M_F}$, $M_F = \{1, 2, \dots, p\}$ = full model.
Questions:

- ▶ How do submodel estimates $\hat{\beta}_M$ relate to full-model estimates $\hat{\beta}$?
- ▶ How do submodel parameters β_M relate to full-model parameters β ?

Is $\hat{\beta}_M$ a subset of $\hat{\beta}$ and β_M a subset of β ?

- Answer: Unless \mathbf{X} is an orthogonal design,
 - ▶ $\hat{\beta}_M$ is **not** a subset of $\hat{\beta}$ and
 - ▶ β_M is **not** a subset of β .
- Reason: Slopes, both estimates and parameters, depend on what the other predictors are — **in value and in meaning**.
- Message: $\hat{\beta}_M$ does **not** estimate full-model parameters!
(Exception: The full model is causal or “data generating” ...
Submodel estimates suffer then from “omitted variables bias.”)

‘Slopes depend on ...’ — An Illustration

- Survey of potential purchasers of a new high-tech gizmo:
 - ▶ Response: “LoP” = Likelihood of Purchase (self-reported on a Likert scale)
 - ▶ Predictor 1: Age
 - ▶ Predictor 2: Income

‘Slopes depend on ...’ — An Illustration

- Survey of potential purchasers of a new high-tech gizmo:
 - ▶ Response: “LoP” = Likelihood of Purchase (self-reported on a Likert scale)
 - ▶ Predictor 1: Age
 - ▶ Predictor 2: Income
- Expectation: Younger customers have higher LoP, that is, $\beta_{\text{Age}} < 0$.

‘Slopes depend on ...’ — An Illustration

- Survey of potential purchasers of a new high-tech gizmo:
 - ▶ Response: “LoP” = Likelihood of Purchase (self-reported on a Likert scale)
 - ▶ Predictor 1: Age
 - ▶ Predictor 2: Income
- Expectation: Younger customers have higher LoP, that is, $\beta_{\text{Age}} < 0$.
- Outcome of the analysis:

‘Slopes depend on ...’ — An Illustration

- Survey of potential purchasers of a new high-tech gizmo:
 - ▶ Response: “LoP” = Likelihood of Purchase (self-reported on a Likert scale)
 - ▶ Predictor 1: Age
 - ▶ Predictor 2: Income
- Expectation: Younger customers have higher LoP, that is, $\beta_{\text{Age}} < 0$.
- Outcome of the analysis:
 - ▶ Against expectations, a regression of LoP on Age alone indicates that older customers have higher LoP: $\beta_{\text{Age}} > 0$

‘Slopes depend on ...’ — An Illustration

- Survey of potential purchasers of a new high-tech gizmo:
 - ▶ Response: “LoP” = Likelihood of Purchase (self-reported on a Likert scale)
 - ▶ Predictor 1: Age
 - ▶ Predictor 2: Income
- Expectation: Younger customers have higher LoP, that is, $\beta_{\text{Age}} < 0$.
- Outcome of the analysis:
 - ▶ Against expectations, a regression of LoP on Age alone indicates that older customers have higher LoP: $\beta_{\text{Age}} > 0$
 - ▶ But a regression of LoP on Age **and** Income indicates that, **adjusted** for Income, younger customers have higher LoP: $\beta_{\text{Age} \bullet \text{Income}} < 0$

‘Slopes depend on ...’ — An Illustration

- Survey of potential purchasers of a new high-tech gizmo:
 - ▶ Response: “LoP” = Likelihood of Purchase (self-reported on a Likert scale)
 - ▶ Predictor 1: Age
 - ▶ Predictor 2: Income
- Expectation: Younger customers have higher LoP, that is, $\beta_{\text{Age}} < 0$.
- Outcome of the analysis:
 - ▶ Against expectations, a regression of LoP on Age alone indicates that older customers have higher LoP: $\beta_{\text{Age}} > 0$
 - ▶ But a regression of LoP on Age **and** Income indicates that, **adjusted** for Income, younger customers have higher LoP: $\beta_{\text{Age} \bullet \text{Income}} < 0$
 - ▶ Enabling factor: (partial) collinearity between Age and Income.

‘Slopes depend on ...’ — An Illustration

- Survey of potential purchasers of a new high-tech gizmo:
 - ▶ Response: “LoP” = Likelihood of Purchase (self-reported on a Likert scale)
 - ▶ Predictor 1: Age
 - ▶ Predictor 2: Income
- Expectation: Younger customers have higher LoP, that is, $\beta_{\text{Age}} < 0$.
- Outcome of the analysis:
 - ▶ Against expectations, a regression of LoP on Age alone indicates that older customers have higher LoP: $\beta_{\text{Age}} > 0$
 - ▶ But a regression of LoP on Age **and** Income indicates that, **adjusted** for Income, younger customers have higher LoP: $\beta_{\text{Age} \bullet \text{Income}} < 0$
 - ▶ Enabling factor: (partial) collinearity between Age and Income.
- A case of Simpson’s paradox: $\beta_{\text{Age}} > 0 > \beta_{\text{Age} \bullet \text{Income}}$.
 - ▶ The **marginal** and the **Income-adjusted** slope have very **different values** and **different meanings**.

Adjustment, Estimates, Parameters, t -Statistics

- Notation and facts for the components of $\hat{\beta}_M$ and β_M , assuming $j \in M$:
 - ▶ Let $\mathbf{X}_{j \cdot M}$ be the predictor \mathbf{X}_j adjusted for the other predictors in M :

$$\mathbf{X}_{j \cdot M} := (\mathbf{I} - \mathbf{H}_{M \setminus \{j\}}) \mathbf{X}_j \perp \mathbf{X}_k \forall k \in M \setminus \{j\}.$$

Adjustment, Estimates, Parameters, t -Statistics

- Notation and facts for the components of $\hat{\beta}_M$ and β_M , assuming $j \in M$:
 - ▶ Let $\mathbf{X}_{j \bullet M}$ be the predictor \mathbf{X}_j adjusted for the other predictors in M :

$$\mathbf{X}_{j \bullet M} := (\mathbf{I} - \mathbf{H}_{M \setminus \{j\}}) \mathbf{X}_j \perp \mathbf{X}_k \quad \forall k \in M \setminus \{j\}.$$

- ▶ Let $\hat{\beta}_{j \bullet M}$ be the slope estimate and $\beta_{j \bullet M}$ be the parameter for \mathbf{X}_j in M :

$$\hat{\beta}_{j \bullet M} := \frac{\mathbf{X}_{j \bullet M}^T \mathbf{Y}}{\|\mathbf{X}_{j \bullet M}\|^2}, \quad \beta_{j \bullet M} := \frac{\mathbf{X}_{j \bullet M}^T \mathbf{E}[\mathbf{Y}]}{\|\mathbf{X}_{j \bullet M}\|^2}.$$

Adjustment, Estimates, Parameters, t -Statistics

- Notation and facts for the components of $\hat{\beta}_M$ and β_M , assuming $j \in M$:
 - ▶ Let $\mathbf{X}_{j \bullet M}$ be the predictor \mathbf{X}_j adjusted for the other predictors in M :

$$\mathbf{X}_{j \bullet M} := (\mathbf{I} - \mathbf{H}_{M \setminus \{j\}}) \mathbf{X}_j \perp \mathbf{X}_k \quad \forall k \in M \setminus \{j\}.$$

- ▶ Let $\hat{\beta}_{j \bullet M}$ be the slope estimate and $\beta_{j \bullet M}$ be the parameter for \mathbf{X}_j in M :

$$\hat{\beta}_{j \bullet M} := \frac{\mathbf{X}_{j \bullet M}^T \mathbf{Y}}{\|\mathbf{X}_{j \bullet M}\|^2}, \quad \beta_{j \bullet M} := \frac{\mathbf{X}_{j \bullet M}^T \mathbf{E}[\mathbf{Y}]}{\|\mathbf{X}_{j \bullet M}\|^2}.$$

- ▶ Let $t_{j \bullet M}$ be the t -statistic for $\hat{\beta}_{j \bullet M}$ and $\beta_{j \bullet M}$:

$$t_{j \bullet M} := \frac{\hat{\beta}_{j \bullet M} - \beta_{j \bullet M}}{\hat{\sigma} / \|\mathbf{X}_{j \bullet M}\|} = \frac{\mathbf{X}_{j \bullet M}^T (\mathbf{Y} - \mathbf{E}[\mathbf{Y}])}{\|\mathbf{X}_{j \bullet M}\| \hat{\sigma}}.$$

Parameters One More Time

Parameters One More Time

- Once more: If the predictors are partly collinear (non-orthogonal) then

$$M \neq M' \Rightarrow \beta_{j \bullet M} \neq \beta_{j \bullet M'} \quad \text{in value and in meaning.}$$

Motto: A difference in adjustment implies a difference in parameters.

Parameters One More Time

- Once more: If the predictors are partly collinear (non-orthogonal) then

$$M \neq M' \Rightarrow \beta_{j \bullet M} \neq \beta_{j \bullet M'} \quad \text{in value and in meaning.}$$

Motto: A difference in adjustment implies a difference in parameters.

- It follows that there are up to $p2^{p-1}$ different parameters $\beta_{j \bullet M}$!

Parameters One More Time

- Once more: If the predictors are partly collinear (non-orthogonal) then

$$M \neq M' \Rightarrow \beta_{j \bullet M} \neq \beta_{j \bullet M'} \quad \text{in value and in meaning.}$$

Motto: A difference in adjustment implies a difference in parameters.

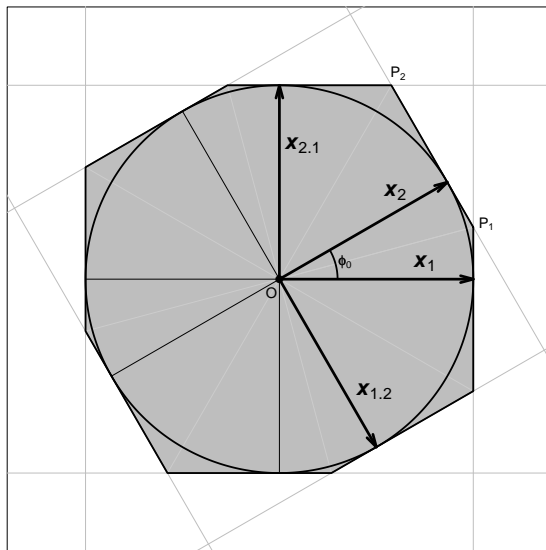
- It follows that there are up to $p2^{p-1}$ different parameters $\beta_{j \bullet M}$!
- However, they are intrinsically p -dimensional:

$$\beta_M = (\mathbf{X}_M^T \mathbf{X}_M)^{-1} \mathbf{X}_M^T \mathbf{X} \beta$$

where \mathbf{X} and β are from the full model.

- Hence each $\beta_{j \bullet M}$ is a lin. comb. of the full model parameters β_1, \dots, β_p .

Geometry of Adjustment



Column space
of \mathbf{X} for $p=2$
predictors,
partly collinear

Error Estimates $\hat{\sigma}^2$

Error Estimates $\hat{\sigma}^2$

- Important: To enable simultaneous inference for all $t_{j\bullet M}$,

Error Estimates $\hat{\sigma}^2$

- Important: To enable simultaneous inference for all $t_{j \bullet M}$,
 - ▶ do **not** use the error estimate $\hat{\sigma}_M^2 := \|\mathbf{Y} - \mathbf{X}_M \hat{\beta}_M\|^2 / (n - m)$ in M ;
(the selected model M may well be 1st order wrong;)
 - ▶ instead, **for all models** M use $\hat{\sigma}^2 = \hat{\sigma}_{Full}^2$ from the full model;

Error Estimates $\hat{\sigma}^2$

- Important: To enable simultaneous inference for all $t_{j \bullet M}$,
 - ▶ do **not** use the error estimate $\hat{\sigma}_M^2 := \|\mathbf{Y} - \mathbf{X}_M \hat{\beta}_M\|^2 / (n - m)$ in M ;
(the selected model M may well be 1st order wrong;)
 - ▶ instead, **for all models** M use $\hat{\sigma}^2 = \hat{\sigma}_{Full}^2$ from the full model;
- $\implies t_{j \bullet M}$ will have a t -distribution **with the same dfs** $\forall M, \forall j \in M$.

Error Estimates $\hat{\sigma}^2$

- Important: To enable simultaneous inference for all $t_{j \bullet M}$,
 - ▶ do **not** use the error estimate $\hat{\sigma}_M^2 := \|\mathbf{Y} - \mathbf{X}_M \hat{\beta}_M\|^2 / (n - m)$ in M ;
(the selected model M may well be 1st order wrong;)
 - ▶ instead, **for all models** M use $\hat{\sigma}^2 = \hat{\sigma}_{Full}^2$ from the full model; $\implies t_{j \bullet M}$ will have a t -distribution **with the same dfs** $\forall M, \forall j \in M$.
- What if even the full model is 1st order wrong?

Error Estimates $\hat{\sigma}^2$

- Important: To enable simultaneous inference for all $t_{j \cdot M}$,
 - ▶ do **not** use the error estimate $\hat{\sigma}_M^2 := \|\mathbf{Y} - \mathbf{X}_M \hat{\beta}_M\|^2 / (n - m)$ in M ;
(the selected model M may well be 1st order wrong;)
 - ▶ instead, **for all models M** use $\hat{\sigma}^2 = \hat{\sigma}_{Full}^2$ from the full model; $\implies t_{j \cdot M}$ will have a t -distribution **with the same dfs $\forall M, \forall j \in M$.**
- What if even the full model is 1st order wrong?
Answer: $\hat{\sigma}_{Full}^2$ will be inflated and inference will be conservative.
But better estimates are available if ...

Error Estimates $\hat{\sigma}^2$

- Important: To enable simultaneous inference for all $t_{j \bullet M}$,
 - ▶ do **not** use the error estimate $\hat{\sigma}_M^2 := \|\mathbf{Y} - \mathbf{X}_M \hat{\beta}_M\|^2 / (n - m)$ in M ;
(the selected model M may well be 1st order wrong;)
 - ▶ instead, **for all models** M use $\hat{\sigma}^2 = \hat{\sigma}_{Full}^2$ from the full model; $\implies t_{j \bullet M}$ will have a t -distribution **with the same dfs** $\forall M, \forall j \in M$.

- What if even the full model is 1st order wrong?

Answer: $\hat{\sigma}_{Full}^2$ will be inflated and inference will be conservative.

But better estimates are available if ...

- ▶ exact replicates exist: use $\hat{\sigma}^2$ from the 1-way ANOVA of replicates;
- ▶ a larger than the full model can be assumed 1st order correct: use $\hat{\sigma}_{Large}^2$;
- ▶ a previous dataset provided a valid estimate: use $\hat{\sigma}_{previous}^2$;
- ▶ nonparametric estimates are available: use $\hat{\sigma}_{nonpar}^2$ (Hall and Carroll 1989).

Error Estimates $\hat{\sigma}^2$

- Important: To enable simultaneous inference for all $t_{j \bullet M}$,
 - ▶ do **not** use the error estimate $\hat{\sigma}_M^2 := \|\mathbf{Y} - \mathbf{X}_M \hat{\beta}_M\|^2 / (n - m)$ in M ;
(the selected model M may well be 1st order wrong;)
 - ▶ instead, **for all models M** use $\hat{\sigma}^2 = \hat{\sigma}_{Full}^2$ from the full model;
- $\implies t_{j \bullet M}$ will have a t -distribution **with the same dfs $\forall M, \forall j \in M$.**

- What if even the full model is 1st order wrong?

Answer: $\hat{\sigma}_{Full}^2$ will be inflated and inference will be conservative.

But better estimates are available if ...

- ▶ exact replicates exist: use $\hat{\sigma}^2$ from the 1-way ANOVA of replicates;
- ▶ a larger than the full model can be assumed 1st order correct: use $\hat{\sigma}_{Large}^2$;
- ▶ a previous dataset provided a valid estimate: use $\hat{\sigma}_{previous}^2$;
- ▶ nonparametric estimates are available: use $\hat{\sigma}_{nonpar}^2$ (Hall and Carroll 1989).

PS: In the fashionable $p > N$ literature, what is their $\hat{\sigma}^2$?

Statistical Inference under First Order Incorrectness

Statistical Inference under First Order Incorrectness

- Statistical inference, one parameter at a time:

If $r = \text{dfs}$ in $\hat{\sigma}^2$ and $K = t_{1-\alpha/2, r}$, then the confidence intervals

$$CI_{j \cdot M}(K) := [\hat{\beta}_{j \cdot M} \pm K \hat{\sigma} / \|\mathbf{X}_{j \cdot M}\|]$$

satisfy each

$$\mathbf{P}[\beta_{j \cdot M} \in CI_{j \cdot M}(K)] = 1 - \alpha.$$

Statistical Inference under First Order Incorrectness

- Statistical inference, one parameter at a time:

If $r = \text{dfs in } \hat{\sigma}^2$ and $K = t_{1-\alpha/2, r}$, then the confidence intervals

$$\text{CI}_{j \cdot \mathbf{M}}(K) := [\hat{\beta}_{j \cdot \mathbf{M}} \pm K \hat{\sigma} / \|\mathbf{X}_{j \cdot \mathbf{M}}\|]$$

satisfy each

$$\mathbf{P}[\beta_{j \cdot \mathbf{M}} \in \text{CI}_{j \cdot \mathbf{M}}(K)] = 1 - \alpha.$$

- Achieved so far:

$$\mathbf{Y} = \boldsymbol{\mu} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}_N(\mathbf{0}, \sigma^2 \mathbf{I})$$

Statistical Inference under First Order Incorrectness

- Statistical inference, one parameter at a time:

If $r = \text{dfs}$ in $\hat{\sigma}^2$ and $K = t_{1-\alpha/2, r}$, then the confidence intervals

$$\text{CI}_{j \cdot M}(K) := [\hat{\beta}_{j \cdot M} \pm K \hat{\sigma} / \|\mathbf{X}_{j \cdot M}\|]$$

satisfy each

$$\mathbf{P}[\beta_{j \cdot M} \in \text{CI}_{j \cdot M}(K)] = 1 - \alpha.$$

- Achieved so far:

$$\mathbf{Y} = \boldsymbol{\mu} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}_N(\mathbf{0}, \sigma^2 \mathbf{I})$$

- ▶ No assumption is made that the submodels are 1st order correct;

Statistical Inference under First Order Incorrectness

- Statistical inference, one parameter at a time:

If $r = \text{dfs}$ in $\hat{\sigma}^2$ and $K = t_{1-\alpha/2, r}$, then the confidence intervals

$$\text{CI}_{j \cdot M}(K) := [\hat{\beta}_{j \cdot M} \pm K \hat{\sigma} / \|\mathbf{X}_{j \cdot M}\|]$$

satisfy each

$$\mathbf{P}[\beta_{j \cdot M} \in \text{CI}_{j \cdot M}(K)] = 1 - \alpha.$$

- Achieved so far:

$$\mathbf{Y} = \boldsymbol{\mu} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}_N(\mathbf{0}, \sigma^2 \mathbf{I})$$

- ▶ No assumption is made that the submodels are 1st order correct;
- ▶ Even the full model may be 1st order incorrect if a valid $\hat{\sigma}^2$ is otherwise available;

Statistical Inference under First Order Incorrectness

- Statistical inference, one parameter at a time:

If $r = \text{dfs}$ in $\hat{\sigma}^2$ and $K = t_{1-\alpha/2, r}$, then the confidence intervals

$$\text{CI}_{j \cdot \mathbf{M}}(K) := [\hat{\beta}_{j \cdot \mathbf{M}} \pm K \hat{\sigma} / \|\mathbf{X}_{j \cdot \mathbf{M}}\|]$$

satisfy each

$$\mathbf{P}[\beta_{j \cdot \mathbf{M}} \in \text{CI}_{j \cdot \mathbf{M}}(K)] = 1 - \alpha.$$

- Achieved so far:

$$\mathbf{Y} = \boldsymbol{\mu} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}_N(\mathbf{0}, \sigma^2 \mathbf{I})$$

- ▶ No assumption is made that the submodels are 1st order correct;
- ▶ Even the full model may be 1st order incorrect if a valid $\hat{\sigma}^2$ is otherwise available;
- ▶ A single error estimate opens up the possibility of simultaneous inference across submodels.

Variable Selection

Variable Selection

- What is a variable selection procedure?

Variable Selection

- What is a variable selection procedure?

A map $\mathbf{Y} \mapsto \hat{\mathbf{M}} = \hat{\mathbf{M}}(\mathbf{Y}), \quad \mathbb{R}^N \rightarrow \mathcal{P}(\{1, \dots, p\})$

Variable Selection

- What is a variable selection procedure?

A map $\mathbf{Y} \mapsto \hat{\mathbf{M}} = \hat{\mathbf{M}}(\mathbf{Y}), \quad \mathbb{R}^N \rightarrow \mathcal{P}(\{1, \dots, p\})$

- ▶ $\hat{\mathbf{M}}$ divides the response space \mathbb{R}^N into up to 2^p subsets.

Variable Selection

- What is a variable selection procedure?

A map $\mathbf{Y} \mapsto \hat{\mathbf{M}} = \hat{\mathbf{M}}(\mathbf{Y}), \quad \mathbb{R}^N \rightarrow \mathcal{P}(\{1, \dots, p\})$

- ▶ $\hat{\mathbf{M}}$ divides the response space \mathbb{R}^N into up to 2^p subsets.
- ▶ In a fixed-predictor framework, selection purely based on \mathbf{X} does not invalidate inference (example: deselect predictors based on VIF, \mathbf{H} , ...).

Variable Selection

- What is a variable selection procedure?

A map $\mathbf{Y} \mapsto \hat{\mathbf{M}} = \hat{\mathbf{M}}(\mathbf{Y}), \quad \mathbb{R}^N \rightarrow \mathcal{P}(\{1, \dots, p\})$

- ▶ $\hat{\mathbf{M}}$ divides the response space \mathbb{R}^N into up to 2^p subsets.
 - ▶ In a fixed-predictor framework, selection purely based on \mathbf{X} does not invalidate inference (example: deselect predictors based on VIF, \mathbf{H} , ...).
- Facing up to post-selection inference: Confusers!

Variable Selection

- What is a variable selection procedure?

A map $\mathbf{Y} \mapsto \hat{\mathbf{M}} = \hat{\mathbf{M}}(\mathbf{Y}), \quad \mathbb{R}^N \rightarrow \mathcal{P}(\{1, \dots, p\})$

- ▶ $\hat{\mathbf{M}}$ divides the response space \mathbb{R}^N into up to 2^p subsets.
 - ▶ In a fixed-predictor framework, selection purely based on \mathbf{X} does not invalidate inference (example: deselect predictors based on VIF, \mathbf{H} , ...).
- Facing up to post-selection inference: Confusers!
 - ▶ Target of Inference: the vector $\beta_{\hat{\mathbf{M}}(\mathbf{Y})}$, its components $\beta_{j \cdot \hat{\mathbf{M}}(\mathbf{Y})}$ for $j \in \hat{\mathbf{M}}(\mathbf{Y})$.

Variable Selection

- What is a variable selection procedure?

A map $\mathbf{Y} \mapsto \hat{\mathbf{M}} = \hat{\mathbf{M}}(\mathbf{Y}), \quad \mathbb{R}^N \rightarrow \mathcal{P}(\{1, \dots, p\})$

- ▶ $\hat{\mathbf{M}}$ divides the response space \mathbb{R}^N into up to 2^p subsets.
 - ▶ In a fixed-predictor framework, selection purely based on \mathbf{X} does not invalidate inference (example: deselect predictors based on VIF, \mathbf{H} , ...).
- Facing up to post-selection inference: Confusers!
 - ▶ Target of Inference: the vector $\beta_{\hat{\mathbf{M}}(\mathbf{Y})}$, its components $\beta_{j \cdot \hat{\mathbf{M}}(\mathbf{Y})}$ for $j \in \hat{\mathbf{M}}(\mathbf{Y})$.
 - ▶ The target of inference is **random**.

Variable Selection

- What is a variable selection procedure?

A map $\mathbf{Y} \mapsto \hat{\mathbf{M}} = \hat{\mathbf{M}}(\mathbf{Y}), \quad \mathbb{R}^N \rightarrow \mathcal{P}(\{1, \dots, p\})$

- ▶ $\hat{\mathbf{M}}$ divides the response space \mathbb{R}^N into up to 2^p subsets.
 - ▶ In a fixed-predictor framework, selection purely based on \mathbf{X} does not invalidate inference (example: deselect predictors based on VIF, \mathbf{H} , ...).
- Facing up to post-selection inference: Confusers!
 - ▶ Target of Inference: the vector $\beta_{\hat{\mathbf{M}}(\mathbf{Y})}$, its components $\beta_{j \cdot \hat{\mathbf{M}}(\mathbf{Y})}$ for $j \in \hat{\mathbf{M}}(\mathbf{Y})$.
 - ▶ The target of inference is **random**.
 - ▶ The target of inference has a **random dimension**: $\beta_{\hat{\mathbf{M}}(\mathbf{Y})} \in \mathbb{R}^{|\hat{\mathbf{M}}(\mathbf{Y})|}$

Variable Selection

- What is a variable selection procedure?

A map $\mathbf{Y} \mapsto \hat{\mathbf{M}} = \hat{\mathbf{M}}(\mathbf{Y}), \quad \mathbb{R}^N \rightarrow \mathcal{P}(\{1, \dots, p\})$

- ▶ $\hat{\mathbf{M}}$ divides the response space \mathbb{R}^N into up to 2^p subsets.
- ▶ In a fixed-predictor framework, selection purely based on \mathbf{X} does not invalidate inference (example: deselect predictors based on VIF, \mathbf{H} , ...).
- Facing up to post-selection inference: Confusers!
 - ▶ Target of Inference: the vector $\beta_{\hat{\mathbf{M}}(\mathbf{Y})}$, its components $\beta_{j \bullet \hat{\mathbf{M}}(\mathbf{Y})}$ for $j \in \hat{\mathbf{M}}(\mathbf{Y})$.
 - ▶ The target of inference is **random**.
 - ▶ The target of inference has a **random dimension**: $\beta_{\hat{\mathbf{M}}(\mathbf{Y})} \in \mathbb{R}^{|\hat{\mathbf{M}}(\mathbf{Y})|}$
 - ▶ Conditional on $j \in \hat{\mathbf{M}}$, the target component $\beta_{j \bullet \hat{\mathbf{M}}(\mathbf{Y})}$ has **random meanings**.

Variable Selection

- What is a variable selection procedure?

A map $\mathbf{Y} \mapsto \hat{\mathbf{M}} = \hat{\mathbf{M}}(\mathbf{Y}), \quad \mathbb{R}^N \rightarrow \mathcal{P}(\{1, \dots, p\})$

- ▶ $\hat{\mathbf{M}}$ divides the response space \mathbb{R}^N into up to 2^p subsets.
- ▶ In a fixed-predictor framework, selection purely based on \mathbf{X} does not invalidate inference (example: deselect predictors based on VIF, \mathbf{H} , ...).

- Facing up to post-selection inference: Confusers!

- ▶ Target of Inference: the vector $\beta_{\hat{\mathbf{M}}(\mathbf{Y})}$, its components $\beta_{j \cdot \hat{\mathbf{M}}(\mathbf{Y})}$ for $j \in \hat{\mathbf{M}}(\mathbf{Y})$.
- ▶ The target of inference is **random**.
- ▶ The target of inference has a **random dimension**: $\beta_{\hat{\mathbf{M}}(\mathbf{Y})} \in \mathbb{R}^{|\hat{\mathbf{M}}(\mathbf{Y})|}$
- ▶ Conditional on $j \in \hat{\mathbf{M}}$, the target component $\beta_{j \cdot \hat{\mathbf{M}}(\mathbf{Y})}$ has **random meanings**.
- ▶ When $j \notin \hat{\mathbf{M}}$ both $\beta_{j \cdot \hat{\mathbf{M}}}$ and $\hat{\beta}_{j \cdot \hat{\mathbf{M}}}$ are **undefined**.

Variable Selection

- What is a variable selection procedure?

A map $\mathbf{Y} \mapsto \hat{\mathbf{M}} = \hat{\mathbf{M}}(\mathbf{Y}), \quad \mathbb{R}^N \rightarrow \mathcal{P}(\{1, \dots, p\})$

- ▶ $\hat{\mathbf{M}}$ divides the response space \mathbb{R}^N into up to 2^p subsets.
- ▶ In a fixed-predictor framework, selection purely based on \mathbf{X} does not invalidate inference (example: deselect predictors based on VIF, \mathbf{H} , ...).

- Facing up to post-selection inference: Confusers!

- ▶ Target of Inference: the vector $\beta_{\hat{\mathbf{M}}(\mathbf{Y})}$, its components $\beta_{j \cdot \hat{\mathbf{M}}(\mathbf{Y})}$ for $j \in \hat{\mathbf{M}}(\mathbf{Y})$.
- ▶ The target of inference is **random**.
- ▶ The target of inference has a **random dimension**: $\beta_{\hat{\mathbf{M}}(\mathbf{Y})} \in \mathbb{R}^{|\hat{\mathbf{M}}(\mathbf{Y})|}$
- ▶ Conditional on $j \in \hat{\mathbf{M}}$, the target component $\beta_{j \cdot \hat{\mathbf{M}}(\mathbf{Y})}$ has **random meanings**.
- ▶ When $j \notin \hat{\mathbf{M}}$ both $\beta_{j \cdot \hat{\mathbf{M}}}$ and $\hat{\beta}_{j \cdot \hat{\mathbf{M}}}$ are **undefined**.
- ▶ Hence the coverage probability $\mathbf{P}[\beta_{j \cdot \hat{\mathbf{M}}} \in \text{CI}_{j \cdot \hat{\mathbf{M}}}(K)]$ is **undefined**.

Universal Post-Selection Inference

Universal Post-Selection Inference

- Candidates for meaningful coverage probabilities:

Universal Post-Selection Inference

- Candidates for meaningful coverage probabilities:
 - ▶ $\mathbf{P}[j \in \hat{M} \ \& \ \beta_{j \cdot \hat{M}} \in CI_{j \cdot \hat{M}}(K)] \quad (\leq \mathbf{P}[j \in \hat{M}])$

Universal Post-Selection Inference

- Candidates for meaningful coverage probabilities:

- ▶ $\mathbf{P}[j \in \hat{M} \ \& \ \beta_{j \cdot \hat{M}} \in \text{CI}_{j \cdot \hat{M}}(K)]$ ($\leq \mathbf{P}[j \in \hat{M}]$)
- ▶ $\mathbf{P}[\beta_{j \cdot \hat{M}} \in \text{CI}_{j \cdot \hat{M}}(K) \mid j \in \hat{M}]$ ($\mathbf{P}[j \in \hat{M}] = ???$)

Universal Post-Selection Inference

- Candidates for meaningful coverage probabilities:

- ▶ $\mathbf{P}[j \in \hat{M} \ \& \ \beta_{j \cdot \hat{M}} \in \text{CI}_{j \cdot \hat{M}}(K)]$ ($\leq \mathbf{P}[j \in \hat{M}]$)
- ▶ $\mathbf{P}[\beta_{j \cdot \hat{M}} \in \text{CI}_{j \cdot \hat{M}}(K) \mid j \in \hat{M}]$ ($\mathbf{P}[j \in \hat{M}] = ???$)
- ▶ $\mathbf{P}[\forall j \in \hat{M} : \beta_{j \cdot \hat{M}} \in \text{CI}_{j \cdot \hat{M}}(K)]$

All are meaningful; the last will be our choice.

Universal Post-Selection Inference

- Candidates for meaningful coverage probabilities:

- ▶ $\mathbf{P}[j \in \hat{M} \ \& \ \beta_{j \cdot \hat{M}} \in CI_{j \cdot \hat{M}}(K)]$ ($\leq \mathbf{P}[j \in \hat{M}]$)
- ▶ $\mathbf{P}[\beta_{j \cdot \hat{M}} \in CI_{j \cdot \hat{M}}(K) \mid j \in \hat{M}]$ ($\mathbf{P}[j \in \hat{M}] = ???$)
- ▶ $\mathbf{P}[\forall j \in \hat{M} : \beta_{j \cdot \hat{M}} \in CI_{j \cdot \hat{M}}(K)]$

All are meaningful; the last will be our choice.

- Overcoming the next difficulty:

- ▶ Problem: None of the above coverage probabilities are known or can be estimated for most selection procedures \hat{M} .

Universal Post-Selection Inference

- Candidates for meaningful coverage probabilities:

- ▶ $\mathbf{P}[j \in \hat{M} \ \& \ \beta_{j \cdot \hat{M}} \in CI_{j \cdot \hat{M}}(K)]$ ($\leq \mathbf{P}[j \in \hat{M}]$)
- ▶ $\mathbf{P}[\beta_{j \cdot \hat{M}} \in CI_{j \cdot \hat{M}}(K) \mid j \in \hat{M}]$ ($\mathbf{P}[j \in \hat{M}] = ???$)
- ▶ $\mathbf{P}[\forall j \in \hat{M} : \beta_{j \cdot \hat{M}} \in CI_{j \cdot \hat{M}}(K)]$

All are meaningful; the last will be our choice.

- Overcoming the next difficulty:

- ▶ Problem: None of the above coverage probabilities are known or can be estimated for most selection procedures \hat{M} .
- ▶ Solution: Ask for more!

Universal Post-Selection Inference

- Candidates for meaningful coverage probabilities:

- ▶ $\mathbf{P}[j \in \hat{M} \ \& \ \beta_{j \cdot \hat{M}} \in \text{CI}_{j \cdot \hat{M}}(K)]$ ($\leq \mathbf{P}[j \in \hat{M}]$)
- ▶ $\mathbf{P}[\beta_{j \cdot \hat{M}} \in \text{CI}_{j \cdot \hat{M}}(K) \mid j \in \hat{M}]$ ($\mathbf{P}[j \in \hat{M}] = ???$)
- ▶ $\mathbf{P}[\forall j \in \hat{M} : \beta_{j \cdot \hat{M}} \in \text{CI}_{j \cdot \hat{M}}(K)]$

All are meaningful; the last will be our choice.

- Overcoming the next difficulty:

- ▶ Problem: None of the above coverage probabilities are known or can be estimated for most selection procedures \hat{M} .
- ▶ Solution: Ask for more!

Universal Post-Selection Inference **for all selection procedures** is doable.

Reduction to Simultaneous Inference

Reduction to Simultaneous Inference

Lemma

For any variable selection procedure $\hat{M} = \hat{M}(\mathbf{Y})$, we have the following “significant triviality bound”:

$$\max_{j \in \hat{M}} |t_{j \cdot \hat{M}}| \leq \max_{\mathbf{M}} \max_{j \in \mathbf{M}} |t_{j \cdot \mathbf{M}}| \quad \forall \mathbf{Y}, \boldsymbol{\mu} \in \mathbb{R}^N.$$

Reduction to Simultaneous Inference

Lemma

For any variable selection procedure $\hat{M} = \hat{M}(\mathbf{Y})$, we have the following “significant triviality bound”:

$$\max_{j \in \hat{M}} |t_{j \cdot \hat{M}}| \leq \max_{\mathbf{M}} \max_{j \in \mathbf{M}} |t_{j \cdot \mathbf{M}}| \quad \forall \mathbf{Y}, \boldsymbol{\mu} \in \mathbb{R}^N.$$

Theorem

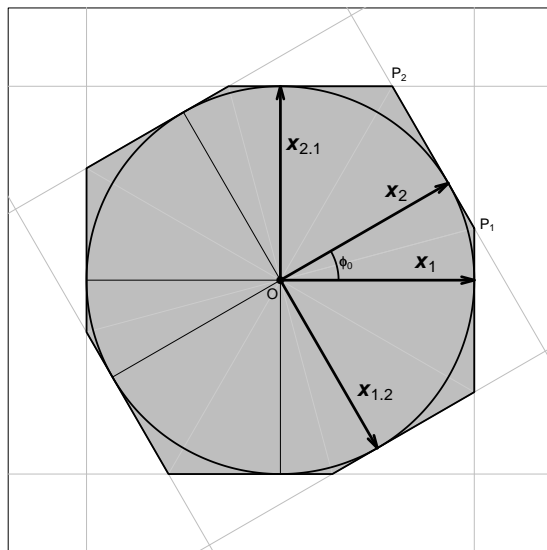
Let K be the $1 - \alpha$ quantile of the “max-max- $|t|$ ” statistic of the lemma:

$$\mathbf{P} \left[\max_{\mathbf{M}} \max_{j \in \mathbf{M}} |t_{j \cdot \mathbf{M}}| \leq K \right] \stackrel{(\geq)}{=} 1 - \alpha.$$

Then we have the following universal PoSI guarantee:

$$\mathbf{P} \left[\beta_{j \cdot \hat{M}} \in CI_{j \cdot \hat{M}}(K) \quad \forall j \in \hat{M} \right] \geq 1 - \alpha \quad \forall \hat{M}.$$

PoSI Geometry — Simultaneity



PoSI polytope
= intersection
of all t -bands.

How Conservative is PoSI?

- Is there a model selection procedure that requires full PoSI protection?

How Conservative is PoSI?

- Is there a model selection procedure that requires full PoSI protection?
- Consider \hat{M} defined as follows:

$$\hat{M} := \operatorname{argmax}_M \max_{j \in M} |t_{j \cdot M}|$$

How Conservative is PoSI?

- Is there a model selection procedure that requires full PoSI protection?
- Consider \hat{M} defined as follows:

$$\hat{M} := \operatorname{argmax}_M \max_{j \in M} |t_{j \cdot M}|$$

A polite name: “Single Predictor Adjusted Regression” =: **SPAR**

How Conservative is PoSI?

- Is there a model selection procedure that requires full PoSI protection?
- Consider \hat{M} defined as follows:

$$\hat{M} := \operatorname{argmax}_M \max_{j \in M} |t_{j \cdot M}|$$

A polite name: “Single Predictor Adjusted Regression” =: **SPAR**

A crude name: “Significance Hunting”

a special case of “p-hacking” (Simmons, Nelson, Simonsohn 2011)

How Conservative is PoSI?

- Is there a model selection procedure that requires full PoSI protection?
- Consider \hat{M} defined as follows:

$$\hat{M} := \operatorname{argmax}_M \max_{j \in M} |t_{j \cdot M}|$$

A polite name: “Single Predictor Adjusted Regression” =: **SPAR**

A crude name: “Significance Hunting”

a special case of “p-hacking” (Simmons, Nelson, Simonsohn 2011)

- SPAR requires the full PoSI protection — by construction!

How Conservative is PoSI?

- Is there a model selection procedure that requires full PoSI protection?
- Consider \hat{M} defined as follows:

$$\hat{M} := \operatorname{argmax}_M \max_{j \in M} |t_{j \cdot M}|$$

A polite name: “Single Predictor Adjusted Regression” =: **SPAR**

A crude name: “Significance Hunting”

a special case of “p-hacking” (Simmons, Nelson, Simonsohn 2011)

- SPAR requires the full PoSI protection — by construction!
- How realistic is SPAR in describing real data analysts behaviors?
 - ▶ It ignores the goodness of fit of the selected model.
 - ▶ It looks for the minimal achievable p-value / strongest “effect”.

Computing PoSI

Computing PoSI

- The simultaneity challenge: there are $p2^{p-1}$ statistics $|t_{j \bullet M}|$.

Computing PoSI

- The simultaneity challenge: there are $p2^{p-1}$ statistics $|t_{j \bullet M}|$.

p	1	2	3	4	5	6	7	8	9	10
$\# t_{j \bullet M} $	1	4	12	32	80	192	448	1,024	2,304	5,120
p	11	12	13	14	15	16	17	18	19	20
$\# t_{j \bullet M} $	11,264	24,576	53,248	114,688	245,760	524,288	1,114,112	2,359,296	4,980,736	10,485,760

Computing PoSI

- The simultaneity challenge: there are $p2^{p-1}$ statistics $|t_{j \bullet M}|$.

p	1	2	3	4	5	6	7	8	9	10
$\# t_{j \bullet M} $	1	4	12	32	80	192	448	1,024	2,304	5,120
p	11	12	13	14	15	16	17	18	19	20
$\# t_{j \bullet M} $	11,264	24,576	53,248	114,688	245,760	524,288	1,114,112	2,359,296	4,980,736	10,485,760

- Monte Carlo-approximation of K_{PoSI} in R, brute force, for $p \approx 20$.

Computing PoSI

- The simultaneity challenge: there are $p2^{p-1}$ statistics $|t_{j \bullet M}|$.

p	1	2	3	4	5	6	7	8	9	10
$\# t_{j \bullet M} $	1	4	12	32	80	192	448	1,024	2,304	5,120
p	11	12	13	14	15	16	17	18	19	20
$\# t_{j \bullet M} $	11,264	24,576	53,248	114,688	245,760	524,288	1,114,112	2,359,296	4,980,736	10,485,760

- Monte Carlo-approximation of K_{PoSI} in R, brute force, for $p \lesssim 20$.
- Computations are specific to a design \mathbf{X} : $K_{\text{PoSI}} = K_{\text{PoSI}}(\mathbf{X}, \alpha, df)$

Computing PoSI

- The simultaneity challenge: there are $p2^{p-1}$ statistics $|t_{j \bullet M}|$.

p	1	2	3	4	5	6	7	8	9	10
$\# t_{j \bullet M} $	1	4	12	32	80	192	448	1,024	2,304	5,120
p	11	12	13	14	15	16	17	18	19	20
$\# t_{j \bullet M} $	11,264	24,576	53,248	114,688	245,760	524,288	1,114,112	2,359,296	4,980,736	10,485,760

- Monte Carlo-approximation of K_{PoSI} in R, brute force, for $p \lesssim 20$.
- Computations are specific to a design \mathbf{X} : $K_{\text{PoSI}} = K_{\text{PoSI}}(\mathbf{X}, \alpha, df)$
- Computations depend only on the inner product matrix $\mathbf{X}^T \mathbf{X}$.

Computing PoSI

- The simultaneity challenge: there are $p2^{p-1}$ statistics $|t_{j \bullet M}|$.

p	1	2	3	4	5	6	7	8	9	10
$\# t_{j \bullet M} $	1	4	12	32	80	192	448	1,024	2,304	5,120
p	11	12	13	14	15	16	17	18	19	20
$\# t_{j \bullet M} $	11,264	24,576	53,248	114,688	245,760	524,288	1,114,112	2,359,296	4,980,736	10,485,760

- Monte Carlo-approximation of K_{PoSI} in R, brute force, for $p \lesssim 20$.
- Computations are specific to a design \mathbf{X} : $K_{\text{PoSI}} = K_{\text{PoSI}}(\mathbf{X}, \alpha, df)$
- Computations depend only on the inner product matrix $\mathbf{X}^T \mathbf{X}$.
 \Rightarrow The limiting factor is p (N may only matter for $\hat{\sigma}^2$).

Computing PoSI

- The simultaneity challenge: there are $p2^{p-1}$ statistics $|t_{j \bullet M}|$.

p	1	2	3	4	5	6	7	8	9	10
$\# t_{j \bullet M} $	1	4	12	32	80	192	448	1,024	2,304	5,120
p	11	12	13	14	15	16	17	18	19	20
$\# t_{j \bullet M} $	11,264	24,576	53,248	114,688	245,760	524,288	1,114,112	2,359,296	4,980,736	10,485,760

- Monte Carlo-approximation of K_{PoSI} in R, brute force, for $p \lesssim 20$.
- Computations are specific to a design \mathbf{X} : $K_{\text{PoSI}} = K_{\text{PoSI}}(\mathbf{X}, \alpha, df)$
- Computations depend only on the inner product matrix $\mathbf{X}^T \mathbf{X}$.
 \Rightarrow The limiting factor is p (N may only matter for $\hat{\sigma}^2$).
- One Monte Carlo computation is good for any α and any error df .

Computing PoSI

- The simultaneity challenge: there are $p2^{p-1}$ statistics $|t_{j \bullet M}|$.

p	1	2	3	4	5	6	7	8	9	10
$\# t_{j \bullet M} $	1	4	12	32	80	192	448	1,024	2,304	5,120
p	11	12	13	14	15	16	17	18	19	20
$\# t_{j \bullet M} $	11,264	24,576	53,248	114,688	245,760	524,288	1,114,112	2,359,296	4,980,736	10,485,760

- Monte Carlo-approximation of K_{PoSI} in R, brute force, for $p \lesssim 20$.
- Computations are specific to a design \mathbf{X} : $K_{\text{PoSI}} = K_{\text{PoSI}}(\mathbf{X}, \alpha, df)$
- Computations depend only on the inner product matrix $\mathbf{X}^T \mathbf{X}$.
 \Rightarrow The limiting factor is p (N may only matter for $\hat{\sigma}^2$).
- One Monte Carlo computation is good for any α and any error df .
- Computations of universal upper bounds:

$$K_{\text{univ}}(p, \alpha, df) \geq K_{\text{PoSI}}(\mathbf{X}, \alpha, df) \quad \forall \mathbf{X} \dots \times p.$$

Scheffé Protection Yields Valid PoSI

Scheffé Protection Yields Valid PoSI

- **Scheffé Simultaneous Inference** is based on the statistic

$$\sup_{\mathbf{x} \in \text{col}(\mathbf{X}) \setminus \{\mathbf{0}\}} \frac{|\mathbf{x}^T (\mathbf{Y} - \mathbf{E}[\mathbf{Y}])|}{\|\mathbf{x}\| \hat{\sigma}} \sim \sqrt{p F_{p, df}}.$$

Scheffé Protection Yields Valid PoSI

- **Scheffé Simultaneous Inference** is based on the statistic

$$\sup_{\mathbf{x} \in \text{col}(\mathbf{X}) \setminus \{\mathbf{0}\}} \frac{|\mathbf{x}^T(\mathbf{Y} - \mathbf{E}[\mathbf{Y}])|}{\|\mathbf{x}\| \hat{\sigma}} \sim \sqrt{p F_{p, df}}.$$

- ▶ The Scheffé method provides sim. inference **for all** linear “contrasts”.

Scheffé Protection Yields Valid PoSI

- **Scheffé Simultaneous Inference** is based on the statistic

$$\sup_{\mathbf{x} \in \text{col}(\mathbf{X}) \setminus \{\mathbf{0}\}} \frac{|\mathbf{x}^T(\mathbf{Y} - \mathbf{E}[\mathbf{Y}])|}{\|\mathbf{x}\| \hat{\sigma}} \sim \sqrt{p F_{p, df}}.$$

- ▶ The Scheffé method provides sim. inference **for all** linear “contrasts”.
- ▶ The Scheffé constant is $K_{\text{Sch}} = K_{\text{Sch}}(p, \alpha, df) = \sqrt{p F_{p, df; 1-\alpha}}$.

Scheffé Protection Yields Valid PoSI

- **Scheffé Simultaneous Inference** is based on the statistic

$$\sup_{\mathbf{x} \in \text{col}(\mathbf{X}) \setminus \{0\}} \frac{|\mathbf{x}^T (\mathbf{Y} - \mathbf{E}[\mathbf{Y}])|}{\|\mathbf{x}\| \hat{\sigma}} \sim \sqrt{p F_{p, df}}.$$

- ▶ The Scheffé method provides sim. inference **for all** linear “contrasts”.
- ▶ The Scheffé constant is $K_{\text{Sch}} = K_{\text{Sch}}(p, \alpha, df) = \sqrt{p F_{p, df; 1-\alpha}}$.

- Compare: **PoSI Simultaneous Inference** is based on the statistic

$$\max_M \max_{j \in M} \frac{|\mathbf{X}_{j \bullet M}^T (\mathbf{Y} - \mathbf{E}[\mathbf{Y}])|}{\|\mathbf{X}_{j \bullet M}\| \hat{\sigma}}$$

Scheffé Protection Yields Valid PoSI

- **Scheffé Simultaneous Inference** is based on the statistic

$$\sup_{\mathbf{x} \in \text{col}(\mathbf{X}) \setminus \{\mathbf{0}\}} \frac{|\mathbf{x}^T (\mathbf{Y} - \mathbf{E}[\mathbf{Y}])|}{\|\mathbf{x}\| \hat{\sigma}} \sim \sqrt{p F_{p, df}}.$$

- ▶ The Scheffé method provides sim. inference **for all** linear “contrasts”.
- ▶ The Scheffé constant is $K_{\text{Sch}} = K_{\text{Sch}}(p, \alpha, df) = \sqrt{p F_{p, df; 1-\alpha}}$.

- Compare: **PoSI Simultaneous Inference** is based on the statistic

$$\max_M \max_{j \in M} \frac{|\mathbf{X}_{j \cdot M}^T (\mathbf{Y} - \mathbf{E}[\mathbf{Y}])|}{\|\mathbf{X}_{j \cdot M}\| \hat{\sigma}}$$

The PoSI contrasts are a subset of the Scheffé contrasts, hence:

Scheffé Protection Yields Valid PoSI

- **Scheffé Simultaneous Inference** is based on the statistic

$$\sup_{\mathbf{x} \in \text{col}(\mathbf{X}) \setminus \{0\}} \frac{|\mathbf{x}^T (\mathbf{Y} - \mathbf{E}[\mathbf{Y}])|}{\|\mathbf{x}\| \hat{\sigma}} \sim \sqrt{p F_{p, df}}.$$

- ▶ The Scheffé method provides sim. inference **for all** linear “contrasts”.
- ▶ The Scheffé constant is $K_{\text{Sch}} = K_{\text{Sch}}(p, \alpha, df) = \sqrt{p F_{p, df; 1-\alpha}}$.

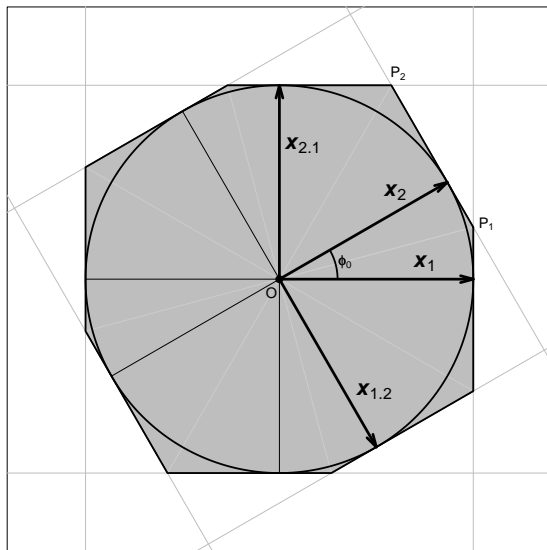
- Compare: **PoSI Simultaneous Inference** is based on the statistic

$$\max_M \max_{j \in M} \frac{|\mathbf{X}_{j \bullet M}^T (\mathbf{Y} - \mathbf{E}[\mathbf{Y}])|}{\|\mathbf{X}_{j \bullet M}\| \hat{\sigma}}$$

The PoSI contrasts are a subset of the Scheffé contrasts, hence:

- ▶ Scheffé statistic \geq PoSI statistic
- ▶ $K_{\text{Sch}} \geq K_{\text{PoSI}}$
- ▶ Scheffé yields universally valid conservative PoSI.

The Scheffé Ball and the PoSI Polytope



Circle =
Scheffé Ball

The PoSI
polytope is
tangent to
the ball.

Orthogonal Designs have the Smallest PoSI Constants

Orthogonal Designs have the Smallest PoSI Constants

- In orthogonal designs there is no adjustment:

$$\mathbf{X}_{j \bullet M}^{\text{orth}} = \mathbf{X}_j^{\text{orth}} \quad \forall M, j (\ni M)$$

Orthogonal Designs have the Smallest PoSI Constants

- In orthogonal designs there is no adjustment:

$$\mathbf{X}_{j \bullet M}^{\text{orth}} = \mathbf{X}_j^{\text{orth}} \quad \forall M, j (\ni M)$$

- The PoSI statistic simplifies to $\max_{j=1 \dots p} |t_{j \bullet \{j\}}|$,
hence the PoSI guarantee reduces to

simultaneity for p orthogonal contrasts.

Orthogonal Designs have the Smallest PoSI Constants

- In orthogonal designs there is no adjustment:

$$\mathbf{X}_{j \bullet M}^{\text{orth}} = \mathbf{X}_j^{\text{orth}} \quad \forall M, j (\ni M)$$

- The PoSI statistic simplifies to $\max_{j=1 \dots p} |t_{j \bullet \{j\}}|$,
hence the PoSI guarantee reduces to

simultaneity for p orthogonal contrasts.

- The PoSI constant for orthogonal designs is uniformly smallest:

$$K_{\text{orth}}(p, \alpha, df) \leq K_{\text{PoSI}}(\mathbf{X}_{\dots \times p}, \alpha, df) \quad \forall p, \alpha, df, \mathbf{X}_{\dots \times p}$$

PoSI Asymptotics

PoSI Asymptotics

- Natural asymptotics for the PoSI constant $K_{\text{PoSI}}(\mathbf{X}_{\dots \times p}, \alpha, df)$ are in terms of design sequences $p \mapsto \mathbf{X}_{\dots \times p}$ as $p \uparrow \infty$ and $df = \infty$, i.e., σ known.

PoSI Asymptotics

- Natural asymptotics for the PoSI constant $K_{\text{PoSI}}(\mathbf{X}_{\dots \times p}, \alpha, df)$ are in terms of design sequences $p \mapsto \mathbf{X}_{\dots \times p}$ as $p \uparrow \infty$ and $df = \infty$, i.e., σ known.
- The Scheffé constant has the following rate in p :

$$K_{\text{Sch}}(p, \alpha) = \sqrt{\chi_{p; 1-\alpha}^2} \sim \sqrt{p}.$$

PoSI Asymptotics

- Natural asymptotics for the PoSI constant $K_{\text{PoSI}}(\mathbf{X}_{\dots \times p}, \alpha, df)$ are in terms of design sequences $p \mapsto \mathbf{X}_{\dots \times p}$ as $p \uparrow \infty$ and $df = \infty$, i.e., σ known.

- The Scheffé constant has the following rate in p :

$$K_{\text{Sch}}(p, \alpha) = \sqrt{\chi_{p; 1-\alpha}^2} \sim \sqrt{p}.$$

- ▶ This represents an upper bound on the PoSI rate.

PoSI Asymptotics

- Natural asymptotics for the PoSI constant $K_{\text{PoSI}}(\mathbf{X}_{\dots \times p}, \alpha, df)$ are in terms of design sequences $p \mapsto \mathbf{X}_{\dots \times p}$ as $p \uparrow \infty$ and $df = \infty$, i.e., σ known.

- The Scheffé constant has the following rate in p :

$$K_{\text{Sch}}(p, \alpha) = \sqrt{\chi_{p; 1-\alpha}^2} \sim \sqrt{p}.$$

- ▶ This represents an upper bound on the PoSI rate.
- ▶ We know a sharper rate bound to be $0.866\ldots\sqrt{p}$.

PoSI Asymptotics

- Natural asymptotics for the PoSI constant $K_{\text{PoSI}}(\mathbf{X}_{\dots \times p}, \alpha, df)$ are in terms of design sequences $p \mapsto \mathbf{X}_{\dots \times p}$ as $p \uparrow \infty$ and $df = \infty$, i.e., σ known.
- The Scheffé constant has the following rate in p :

$$K_{\text{Sch}}(p, \alpha) = \sqrt{\chi_{p;1-\alpha}^2} \sim \sqrt{p}.$$

- ▶ This represents an upper bound on the PoSI rate.
- ▶ We know a sharper rate bound to be $0.866\ldots\sqrt{p}$.
- ▶ We know of design sequences that reach $0.78\ldots\sqrt{p}$.

PoSI Asymptotics

- Natural asymptotics for the PoSI constant $K_{\text{PoSI}}(\mathbf{X}_{\dots \times p}, \alpha, df)$ are in terms of design sequences $p \mapsto \mathbf{X}_{\dots \times p}$ as $p \uparrow \infty$ and $df = \infty$, i.e., σ known.

- The Scheffé constant has the following rate in p :

$$K_{\text{Sch}}(p, \alpha) = \sqrt{\chi_{p; 1-\alpha}^2} \sim \sqrt{p}.$$

- ▶ This represents an upper bound on the PoSI rate.
- ▶ We know a sharper rate bound to be $0.866 \dots \sqrt{p}$.
- ▶ We know of design sequences that reach $0.78 \dots \sqrt{p}$.

- The lowest rate is achieved by orthogonal designs with a rate

$$K_{\text{orth}}(p, \alpha) \sim \sqrt{2 \log p}.$$

PoSI Asymptotics

- Natural asymptotics for the PoSI constant $K_{\text{PoSI}}(\mathbf{X}_{\dots \times p}, \alpha, df)$ are in terms of design sequences $p \mapsto \mathbf{X}_{\dots \times p}$ as $p \uparrow \infty$ and $df = \infty$, i.e., σ known.

- The Scheffé constant has the following rate in p :

$$K_{\text{Sch}}(p, \alpha) = \sqrt{\chi_{p; 1-\alpha}^2} \sim \sqrt{p}.$$

- ▶ This represents an upper bound on the PoSI rate.
- ▶ We know a sharper rate bound to be $0.866\ldots\sqrt{p}$.
- ▶ We know of design sequences that reach $0.78\ldots\sqrt{p}$.

- The lowest rate is achieved by orthogonal designs with a rate

$$K_{\text{orth}}(p, \alpha) \sim \sqrt{2 \log p}.$$

- Hence there is a wide range of rates for the PoSI constants:

$$\sqrt{2 \log p} \lesssim K_{\text{PoSI}}(\mathbf{X}_{\dots \times p}, \alpha) \lesssim \sqrt{p}$$

PoSI Asymptotics

- Natural asymptotics for the PoSI constant $K_{\text{PoSI}}(\mathbf{X}_{\dots \times p}, \alpha, df)$ are in terms of design sequences $p \mapsto \mathbf{X}_{\dots \times p}$ as $p \uparrow \infty$ and $df = \infty$, i.e., σ known.

- The Scheffé constant has the following rate in p :

$$K_{\text{Sch}}(p, \alpha) = \sqrt{\chi_{p; 1-\alpha}^2} \sim \sqrt{p}.$$

- ▶ This represents an upper bound on the PoSI rate.
- ▶ We know a sharper rate bound to be $0.866\ldots\sqrt{p}$.
- ▶ We know of design sequences that reach $0.78\ldots\sqrt{p}$.

- The lowest rate is achieved by orthogonal designs with a rate

$$K_{\text{orth}}(p, \alpha) \sim \sqrt{2 \log p}.$$

- Hence there is a wide range of rates for the PoSI constants:

$$\sqrt{2 \log p} \lesssim K_{\text{PoSI}}(\mathbf{X}_{\dots \times p}, \alpha) \lesssim \sqrt{p}$$

- Under all circumstances, K should not be $t_{df; 1-\alpha/2} = \mathcal{O}(1)!$

Worst Case PoSI Asymptotics — Some Details

- Comments on upper bound $K_{\text{PoSI}} \lesssim 0.866... \sqrt{p}$:

Worst Case PoSI Asymptotics — Some Details

- Comments on upper bound $K_{\text{PoSI}} \lesssim 0.866... \sqrt{p}$:
 - ▶ Ignores the PoSI structure, i.e., the many orthogonalities from adjustment.

Worst Case PoSI Asymptotics — Some Details

- Comments on upper bound $K_{\text{PoSI}} \lesssim 0.866... \sqrt{p}$:
 - ▶ Ignores the PoSI structure, i.e., the many orthogonalities from adjustment.
 - ▶ Based purely on the growth rate: $|\{\mathbf{X}_{j:\mathbf{M}} : j \in \mathbf{M}\}| = p2^{p-1} \sim 2^p$

Worst Case PoSI Asymptotics — Some Details

- Comments on upper bound $K_{\text{PoSI}} \lesssim 0.866... \sqrt{p}$:
 - ▶ Ignores the PoSI structure, i.e., the many orthogonalities from adjustment.
 - ▶ Based purely on the growth rate: $|\{\mathbf{X}_{j:\mathbf{M}} : j \in \mathbf{M}\}| = p 2^{p-1} \sim 2^p$
 - ▶ Bound is achieved by random selection of $\sim 2^p$ random vectors:
 $\mathbf{X}_1, \dots, \mathbf{X}_{p2^{p-1}} \sim U(S^{p-1})$ i.i.d. versus $\{\mathbf{X}_{j:\mathbf{M}} : j \in \mathbf{M}\}$

Worst Case PoSI Asymptotics — Some Details

- Comments on upper bound $K_{\text{PoSI}} \lesssim 0.866... \sqrt{p}$:
 - ▶ Ignores the PoSI structure, i.e., the many orthogonalities from adjustment.
 - ▶ Based purely on the growth rate: $|\{\mathbf{X}_{j:\mathbf{M}} : j \in \mathbf{M}\}| = p 2^{p-1} \sim 2^p$
 - ▶ Bound is achieved by random selection of $\sim 2^p$ random vectors:
 $\mathbf{X}_1, \dots, \mathbf{X}_{p 2^{p-1}} \sim U(S^{p-1})$ i.i.d. versus $\{\mathbf{X}_{j:\mathbf{M}} : j \in \mathbf{M}\}$
 - ▶ Reduction to radial problem: $\max_{j:\mathbf{M}} \langle \mathbf{U}, \mathbf{X}_{j:\mathbf{M}} \rangle$ ($\mathbf{U} \sim U(S^{p-1})$).
 \Rightarrow Wyner's (1967) bounds on sphere packing apply.

Worst Case PoSI Asymptotics — Some Details

- Comments on upper bound $K_{\text{PoSI}} \lesssim 0.866... \sqrt{p}$:
 - ▶ Ignores the PoSI structure, i.e., the many orthogonalities from adjustment.
 - ▶ Based purely on the growth rate: $|\{\mathbf{X}_{j:\mathbf{M}} : j \in \mathbf{M}\}| = p2^{p-1} \sim 2^p$
 - ▶ Bound is achieved by random selection of $\sim 2^p$ random vectors:
 $\mathbf{X}_1, \dots, \mathbf{X}_{p2^{p-1}} \sim U(S^{p-1})$ i.i.d. versus $\{\mathbf{X}_{j:\mathbf{M}} : j \in \mathbf{M}\}$
 - ▶ Reduction to radial problem: $\max_{j:\mathbf{M}} \langle \mathbf{U}, \mathbf{X}_{j:\mathbf{M}} \rangle$ ($\mathbf{U} \sim U(S^{p-1})$).
 \Rightarrow Wyner's (1967) bounds on sphere packing apply.
- Comments on lower bound $K_{\text{PoSI}} \gtrsim 0.78 \sqrt{p}$:
 - ▶ Best lower bound known to date is found by construction of an example.

$$\left| \begin{array}{cccccc} 1 & 0 & 0 & 0 & \dots & c \\ 0 & 1 & 0 & 0 & \dots & c \\ 0 & 0 & 1 & 0 & \dots & c \\ 0 & 0 & 0 & 1 & \dots & c \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & \sqrt{1 - (p-1)c^2} \end{array} \right|$$

- ▶ This is not be the ultimate worst case yet.

Example: Length of Criminal Sentence (contd.)

- Reminder: t -statistics of selected covariates, in descending order:

- ▶ $|t_{\text{alcohol}}| = 3.95;$

- ▶ $|t_{\text{prior records}}| = 3.59;$

- ▶ $|t_{\text{seriousness}}| = 3.57;$

- ▶ $|t_{\text{drugs}}| = 3.31;$

- ▶ $|t_{\text{employment}}| = 3.04;$

- ▶ $|t_{\text{initial age}}| = 2.56;$

- ▶ $|t_{\text{gender}}| = 2.33.$

Example: Length of Criminal Sentence (contd.)

- Reminder: t -statistics of selected covariates, in descending order:
 - ▶ $|t_{\text{alcohol}}| = 3.95;$
 - ▶ $|t_{\text{prior records}}| = 3.59;$
 - ▶ $|t_{\text{seriousness}}| = 3.57;$
 - ▶ $|t_{\text{drugs}}| = 3.31;$
 - ▶ $|t_{\text{employment}}| = 3.04;$
 - ▶ $|t_{\text{initial age}}| = 2.56;$
 - ▶ $|t_{\text{gender}}| = 2.33.$
- The PoSI constant is $K_{\text{PoSI}} \approx 3.1$, hence we would claim significance for the four variables on the left.
- For comparison, the Scheffé constant is $K_{\text{Sch}} \approx 4.5$, leaving us with no significant predictors at all.
- Similarly, Bonferroni with $\alpha/(p2^{p-1})$ yields $K_{\text{Bonf}} \approx 4.7$.

Conclusions

- Valid universal post-selection inference is possible.

Conclusions

- Valid universal post-selection inference is possible.
- Necessary buy-ins:
 - ▶ Each submodel has its own slope parameters.
 - ▶ Use one $\hat{\sigma}$ you believe in for all $t_{j\bullet M}$.
 - ▶ Valid inference for “wrong” models is meaningful.

Conclusions

- Valid universal post-selection inference is possible.
- Necessary buy-ins:
 - ▶ Each submodel has its own slope parameters.
 - ▶ Use one $\hat{\sigma}$ you believe in for all $t_{j\bullet M}$.
 - ▶ Valid inference for “wrong” models is meaningful.
- PoSI is not procedure-specific, hence is conservative. However:
 - ▶ PoSI is valid even for selection that is informal and post-hoc.
 - ▶ PoSI is necessary for selection based on “significance hunting”.

Conclusions

- Valid universal post-selection inference is possible.
- Necessary buy-ins:
 - ▶ Each submodel has its own slope parameters.
 - ▶ Use one $\hat{\sigma}$ you believe in for all $t_{j \bullet M}$.
 - ▶ Valid inference for “wrong” models is meaningful.
- PoSI is not procedure-specific, hence is conservative. However:
 - ▶ PoSI is valid even for selection that is informal and post-hoc.
 - ▶ PoSI is necessary for selection based on “significance hunting”.
- Asymptotics in p suggests strong dependence of K_{PoSI} on design \mathbf{X} .

Conclusions

- Valid universal post-selection inference is possible.
- Necessary buy-ins:
 - ▶ Each submodel has its own slope parameters.
 - ▶ Use one $\hat{\sigma}$ you believe in for all $t_{j \bullet M}$.
 - ▶ Valid inference for “wrong” models is meaningful.
- PoSI is not procedure-specific, hence is conservative. However:
 - ▶ PoSI is valid even for selection that is informal and post-hoc.
 - ▶ PoSI is necessary for selection based on “significance hunting”.
- Asymptotics in p suggests strong dependence of K_{PoSI} on design \mathbf{X} .
- Challenges:
 - ▶ PoSI under heteroskedasticity, random \mathbf{X} , general misspecification, ...
 - ▶ Understanding the design geometry that drives $K_{\text{PoSI}}(\mathbf{X})$.
 - ▶ Computations of K_{PoSI} for large p .

Conclusions

- Valid universal post-selection inference is possible.
- Necessary buy-ins:
 - ▶ Each submodel has its own slope parameters.
 - ▶ Use one $\hat{\sigma}$ you believe in for all $t_{j\bullet M}$.
 - ▶ Valid inference for “wrong” models is meaningful.
- PoSI is not procedure-specific, hence is conservative. However:
 - ▶ PoSI is valid even for selection that is informal and post-hoc.
 - ▶ PoSI is necessary for selection based on “significance hunting”.
- Asymptotics in p suggests strong dependence of K_{PoSI} on design \mathbf{X} .
- Challenges:
 - ▶ PoSI under heteroskedasticity, random \mathbf{X} , general misspecification, ...
 - ▶ Understanding the design geometry that drives $K_{\text{PoSI}}(\mathbf{X})$.
 - ▶ Computations of K_{PoSI} for large p .

THANK YOU!

Details of Simulation

- Number of simulated datasets: 100,000.
- $\beta = 0, \gamma_j = 4, \forall j = 1, \dots, 10$.
- Vectors \mathbf{x} and \mathbf{z}_j 's are standardized to have zero mean and unit variance.
 - ▶ The correlation between \mathbf{x} and each \mathbf{z}_j is 0.7, $\forall j = 1, \dots, 10$.
 - ▶ The correlation between \mathbf{z}_{j_1} and \mathbf{z}_{j_2} is 0.5, $\forall j_1, j_2 = 1, \dots, 10$.

▶ Back

PoSI1 — SPAR1: Focus on One Predictor of Interest

- Sometimes there is one focal predictor of interest, X_p .

PoSI1 — SPAR1: Focus on One Predictor of Interest

- Sometimes there is one focal predictor of interest, \mathbf{X}_p .
- Inference is desired only for $\beta_{p \cdot M}$ (test statistics: $t_{p \cdot M}$).

PoSI1 — SPAR1: Focus on One Predictor of Interest

- Sometimes there is one focal predictor of interest, \mathbf{X}_p .
- Inference is desired only for $\beta_{p \bullet M}$ (test statistics: $t_{p \bullet M}$).
- Search only models M that contain p : $p \in M$ ($\# = 2^{p-1}$).
Purpose: Boost the statistical significance of $\hat{\beta}_{p \bullet M}$.

PoSI1 — SPAR1: Focus on One Predictor of Interest

- Sometimes there is one focal predictor of interest, \mathbf{X}_p .
- Inference is desired only for $\beta_{p \bullet M}$ (test statistics: $t_{p \bullet M}$).
- Search only models M that contain p : $p \in M$ ($\# = 2^{p-1}$).
Purpose: Boost the statistical significance of $\hat{\beta}_{p \bullet M}$.
- **PoSI1** produces a constant K_{PoSI1} whose intervals $\text{CI}_{p \bullet M}(K_{\text{PoSI1}})$ are valid after any variable selection procedure \hat{M} that is subject to $p \in \hat{M}$.

PoSI1 — SPAR1: Focus on One Predictor of Interest

- Sometimes there is one focal predictor of interest, \mathbf{X}_p .
- Inference is desired only for $\beta_{p \bullet \mathbf{M}}$ (test statistics: $t_{p \bullet \mathbf{M}}$).
- Search only models \mathbf{M} that contain p : $p \in \mathbf{M}$ ($\# = 2^{p-1}$).
Purpose: Boost the statistical significance of $\hat{\beta}_{p \bullet \mathbf{M}}$.
- **PoSI1** produces a constant K_{PoSI1} whose intervals $\text{CI}_{p \bullet \mathbf{M}}(K_{\text{PoSI1}})$ are valid after any variable selection procedure $\hat{\mathbf{M}}$ that is subject to $p \in \hat{\mathbf{M}}$.
- A selection procedure that requires the full PoSI1 protection:
SPAR1, defined by $\hat{\mathbf{M}} := \operatorname{argmax}_{\mathbf{M} \ni p} |t_{p \bullet \mathbf{M}}|$.

PoSI1 — SPAR1: Focus on One Predictor of Interest

- Sometimes there is one focal predictor of interest, \mathbf{X}_p .
- Inference is desired only for $\beta_{p \bullet \mathbf{M}}$ (test statistics: $t_{p \bullet \mathbf{M}}$).
- Search only models \mathbf{M} that contain p : $p \in \mathbf{M}$ ($\# = 2^{p-1}$).
Purpose: Boost the statistical significance of $\hat{\beta}_{p \bullet \mathbf{M}}$.
- PoSI1 produces a constant K_{PoSI1} whose intervals $\text{CI}_{p \bullet \mathbf{M}}(K_{\text{PoSI1}})$ are valid after any variable selection procedure $\hat{\mathbf{M}}$ that is subject to $p \in \hat{\mathbf{M}}$.
- A selection procedure that requires the full PoSI1 protection:
SPAR1, defined by $\hat{\mathbf{M}} := \operatorname{argmax}_{\mathbf{M} \ni p} |t_{p \bullet \mathbf{M}}|$.

Conclusions:

- PoSI1 is more appropriate for some situations than full PoSI.
- Trivially, $K_{\text{PoSI1}} < K_{\text{PoSI}}$, but sometimes not by much!

Ways to Limit the Size of the PoSI Problem

- The full universe of models for full PoSI: all non-singular submodels
 - ▶ $\mathcal{M}_{\text{all}} = \{M : M \subset \{1, 2, \dots, p\}, 0 < |M| \leq \min(n, p), \text{rank}(\mathbf{X}_M) = |M|\}.$
- Useful sub-universes:

Ways to Limit the Size of the PoSI Problem

- The full universe of models for full PoSI: all non-singular submodels
 - ▶ $\mathcal{M}_{\text{all}} = \{M : M \subset \{1, 2, \dots, p\}, 0 < |M| \leq \min(n, p), \text{rank}(\mathbf{X}_M) = |M|\}.$
- Useful sub-universes:
 - ▶ Protect one or more predictors, as in PoSI1: $\mathcal{M} = \{M : p \in M\}.$

Ways to Limit the Size of the PoSI Problem

- The full universe of models for full PoSI: all non-singular submodels
 - ▶ $\mathcal{M}_{\text{all}} = \{M : M \subset \{1, 2, \dots, p\}, 0 < |M| \leq \min(n, p), \text{rank}(\mathbf{X}_M) = |M|\}.$
- Useful sub-universes:
 - ▶ Protect one or more predictors, as in PoSI1: $\mathcal{M} = \{M : p \in M\}.$
 - ▶ Sparsity, i.e., submodels of size m' or less: $\mathcal{M} = \{M : |M| \leq m'\}.$

Ways to Limit the Size of the PoSI Problem

- The full universe of models for full PoSI: all non-singular submodels
 - ▶ $\mathcal{M}_{\text{all}} = \{M : M \subset \{1, 2, \dots, p\}, 0 < |M| \leq \min(n, p), \text{rank}(\mathbf{X}_M) = |M|\}.$
- Useful sub-universes:
 - ▶ Protect one or more predictors, as in PoSI1: $\mathcal{M} = \{M : p \in M\}.$
 - ▶ Sparsity, i.e., submodels of size m' or less: $\mathcal{M} = \{M : |M| \leq m'\}.$
 - ▶ Richness, i.e., drop fewer than m' predictors from the full model: $\mathcal{M} = \{M : |M| \geq p - m'\}.$

Ways to Limit the Size of the PoSI Problem

- The full universe of models for full PoSI: all non-singular submodels
 - ▶ $\mathcal{M}_{\text{all}} = \{M : M \subset \{1, 2, \dots, p\}, 0 < |M| \leq \min(n, p), \text{rank}(\mathbf{X}_M) = |M|\}.$
- Useful sub-universes:
 - ▶ Protect one or more predictors, as in PoSI1: $\mathcal{M} = \{M : p \in M\}.$
 - ▶ Sparsity, i.e., submodels of size m' or less: $\mathcal{M} = \{M : |M| \leq m'\}.$
 - ▶ Richness, i.e., drop fewer than m' predictors from the full model: $\mathcal{M} = \{M : |M| \geq p - m'\}.$
 - ▶ Nested sets of models, as in polynomial regression, AR models, ANOVA.

PoSI Significance: Strong Error Control

For each $j \in M$, consider the t -test statistic

$$t_{0,j \cdot \mathbf{M}} = \frac{\hat{\beta}_{j \cdot \mathbf{M}} - 0}{\hat{\sigma}_{\cdot} \cdot ((\mathbf{X}_{\mathbf{M}}^T \mathbf{X}_{\mathbf{M}})^{-1})_{jj}^{\frac{1}{2}}}.$$

Theorem

Let H_1 be the random set of true alternatives in \hat{M} ,
and \hat{H}_1 the random set of rejections in \hat{M} :

$$\hat{H}_1 = \{(j, \hat{M}) : j \in \hat{M}, |t_{0,j,\hat{M}}| > K\} \quad \text{and} \quad H_1 = \{(j, \hat{M}) : j \in \hat{M}, \beta_{j,\hat{M}} \neq 0\}.$$

Then

$$\mathbb{P}(\hat{H}_1 \subset H_1) \geq 1 - \alpha.$$

If we repeat the sampling process many times, the probability that all PoSI rejections are correct is at least $1 - \alpha$, no matter how the model is selected.