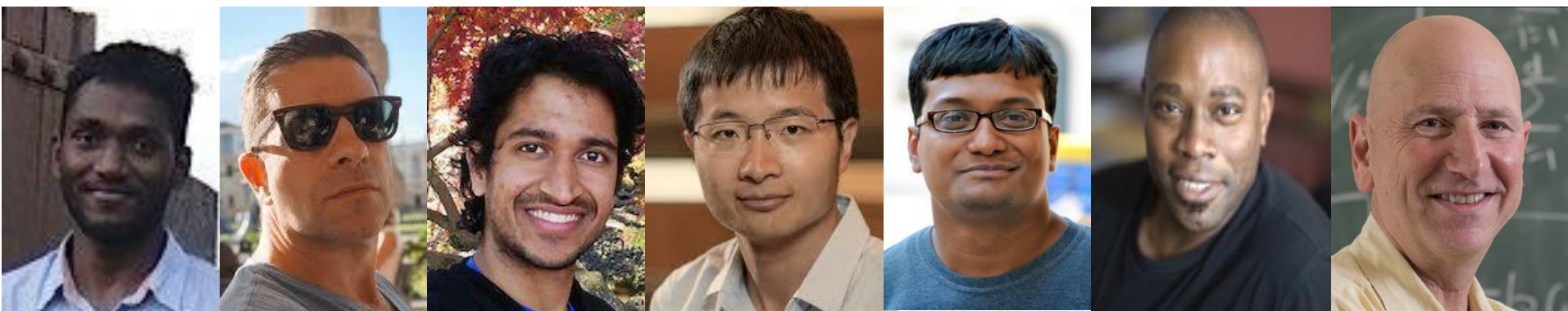




# Arun Kumar Kuchibhotla

University of Pennsylvania

<https://arun-kuchibhotla.github.io/>



# Overview of my research interests

Robust Statistics

Single Index Models

Nonparametric  
Regression

Post-selection  
Inference

Misspecification  
Analysis

Concentration Inequalities  
and  
High-dimensional CLT

# Valid Post-selection Inference

Why and How

**Arun Kumar Kuchibhotla**  
University of Pennsylvania

Larry Brown



Andreas Buja



Junhui Cai



Linda Zhao



Ed George



# Roadmap

Data Snooping: Effects and Examples

Formulation of the Problem

Solution for Covariate Selection

Example and Conclusions

# Roadmap

## Data Snooping: Effects and Examples

- Effects illustrated with stepwise selection.
- Data snooping in textbooks and practice.

## Formulation of the Problem

- The Problem & literature review for covariate selection.

## Solution for Covariate Selection

- Key contributions
- Simulations & main components of the theory.

## Example and Conclusions

- Real data example, Extensions & Summary

# Roadmap

## Data Snooping: Effects and Examples

- Effects illustrated with stepwise selection.
- Data snooping in textbooks and practice.

Formulation of the Problem

Solution for Covariate Selection

Example and Conclusions

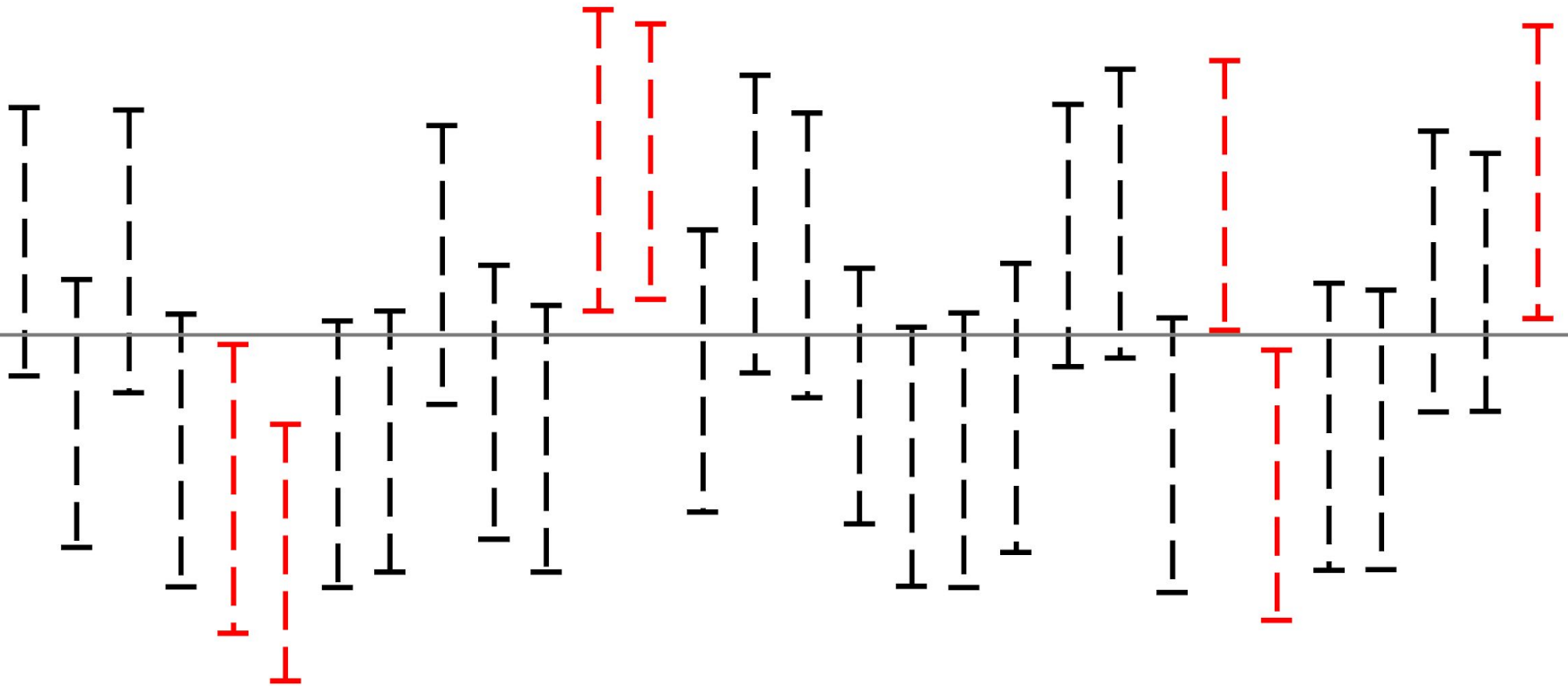
# Effects of Data Snooping: Stepwise Selection

$(X, Y) \sim N(0, I_{p+1}) \implies$  500 observations

$X_{\hat{j}}$  most correlated with  $Y \implies$  95% CI for slope

$p = 5$

Unadjusted: **76.9%**



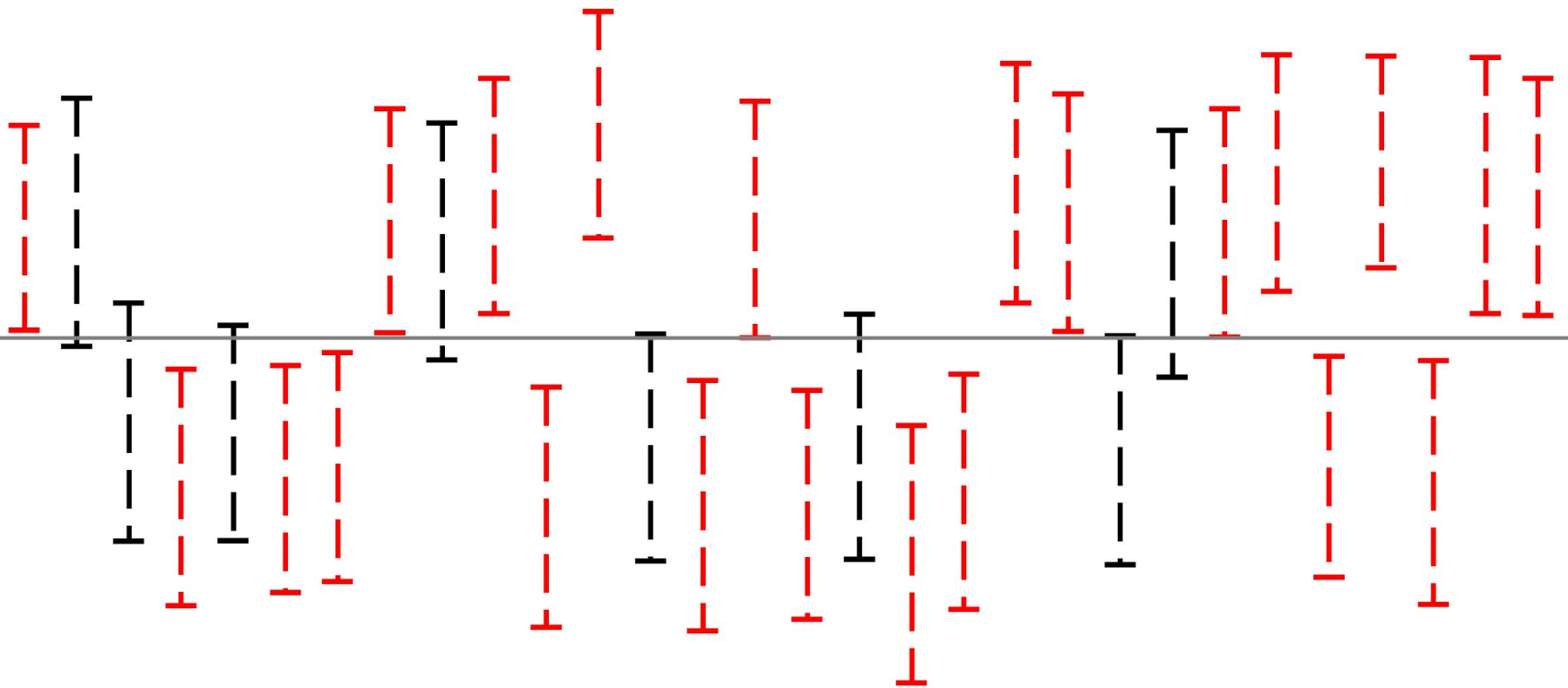
# Effects of Data Snooping: Stepwise Selection

$(X, Y) \sim N(0, I_{p+1}) \implies$  500 observations

$X_{\hat{j}}$  most correlated with  $Y \implies$  95% CI for slope

$p = 20$

Unadjusted: **32.6%**



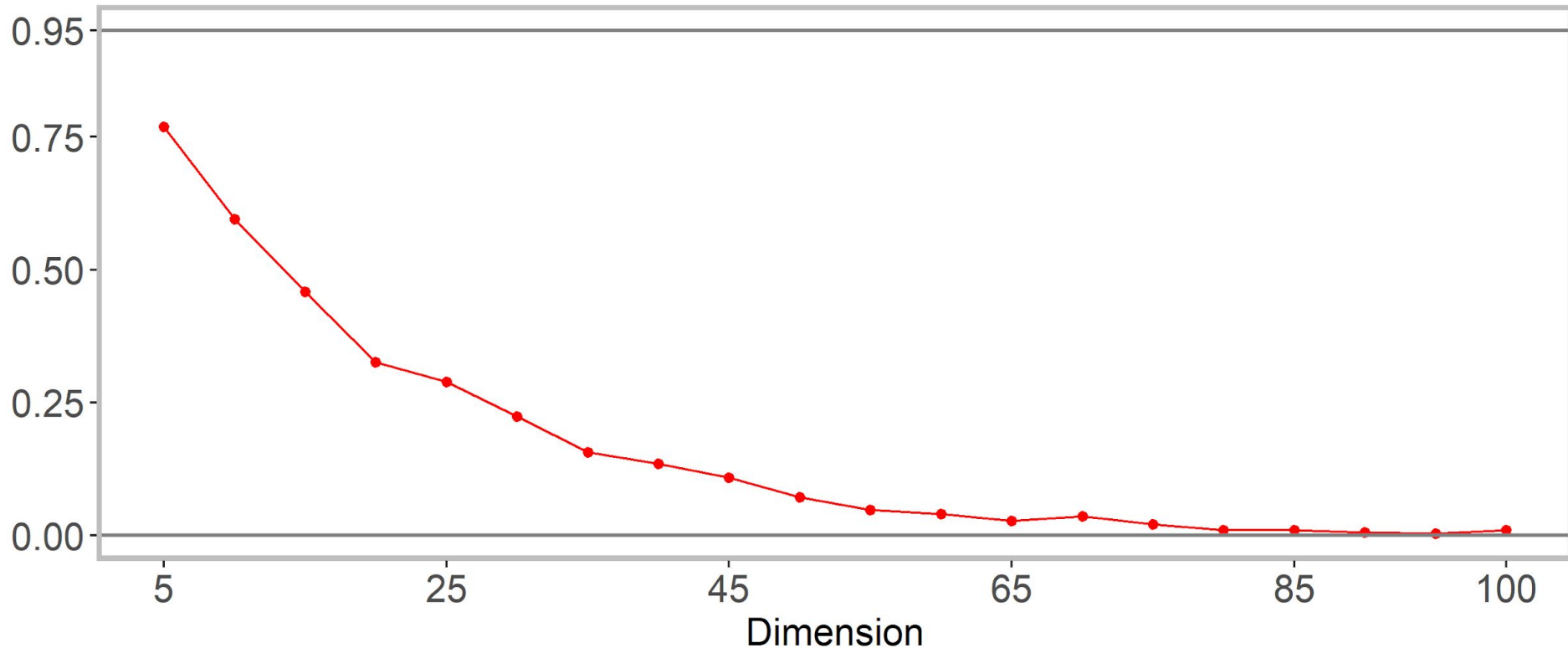


# Effects of Data Snooping: Stepwise Selection

$(X, Y) \sim N(0, I_{p+1}) \implies 500$  observations

$X_{\hat{j}}$  most correlated with  $Y \implies 95\%$  CI for slope

Empirical Coverage



# Some Notes

❖ Unadjusted inference after data snooping can be (very) **misleading**.

❖ Data snooping contributes to the

## **Replicability Crisis**

- Inability to replicate conclusions in future studies.
- 95% CIs should imply correct conclusions in 95% of studies.

❖ More concerningly, common practice of data snooping is more **informal** and **imprecise** than the example shown.

# Case Study 1: Covariate Selection

**British  
Medical  
Journal  
2005**

Postdischarge mortality in children with acute infectious diseases: derivation of postdischarge mortality prediction models

variate imputation using chained equations.<sup>12</sup> Following univariate analysis, candidate models were generated using a stepwise selection procedure minimising Akaike's Information Criterion (AIC). This method is

# Case Study 1: Covariate Selection

**British  
Medical  
Journal  
2005**

Postdischarge mortality in children with acute infectious diseases: derivation of postdischarge mortality prediction models

variate imputation using chained equations.<sup>12</sup> Following univariate analysis, candidate models were generated using a stepwise selection procedure minimising Akaike's Information Criterion (AIC). This method is

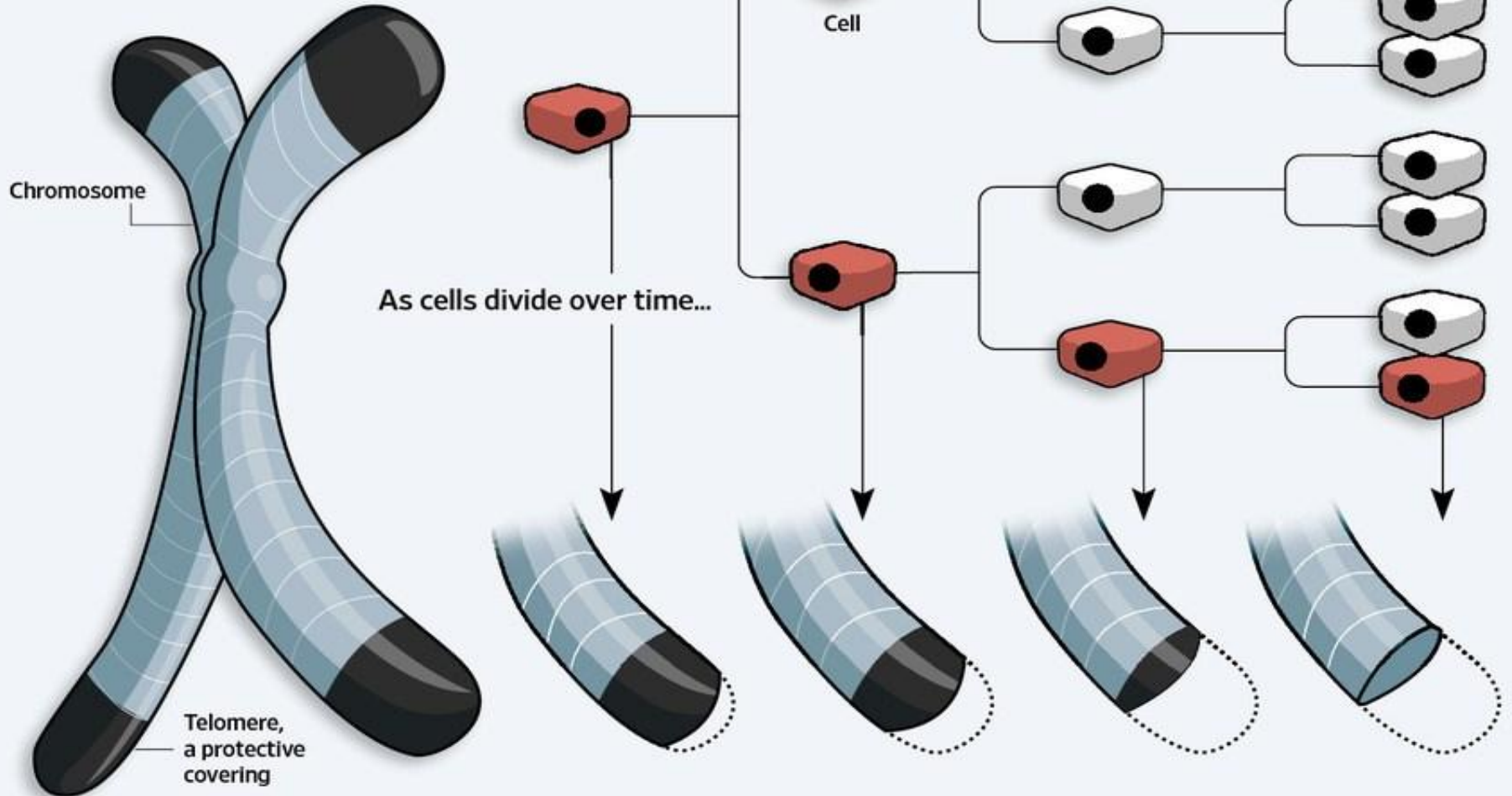
• • •

final selection of a model was judged on model parsimony (the simpler the better), availability of the predictors (with respect to minimal resources and cost), and the attained sensitivity (with at least 50% specificity). All

# 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.'



# 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.**”

# Case Study 3: Transformations

Harrison and Rubinfeld (1978)<sup>★</sup> write

to determine the best fitting functional form. Comparing models with either median value of owner-occupied homes ( $MV$ ) or  $\text{Log}(MV)$  as the dependent variable, we found that the semilog version provided a slightly better fit. Using  $\text{Log}(MV)$  as the dependent variable, we concentrated on estimating a nonlinear term in  $NOX$ ; i.e., we included  $NOX^p$  in the equation, where  $p$  is an unknown

The statistical fit in the equation was best when  $p$  was set equal to 2.0,

<sup>★</sup>H & R (1978) Hedonic housing prices and the demand for clean air.

# Roadmap

## Data Snooping: Effects and Examples

- Effects illustrated with stepwise selection.
- Data snooping in textbooks and practice.

## Formulation of the Problem

- The Problem & literature review for covariate selection.

## Solution for Covariate Selection

## Example and Conclusions



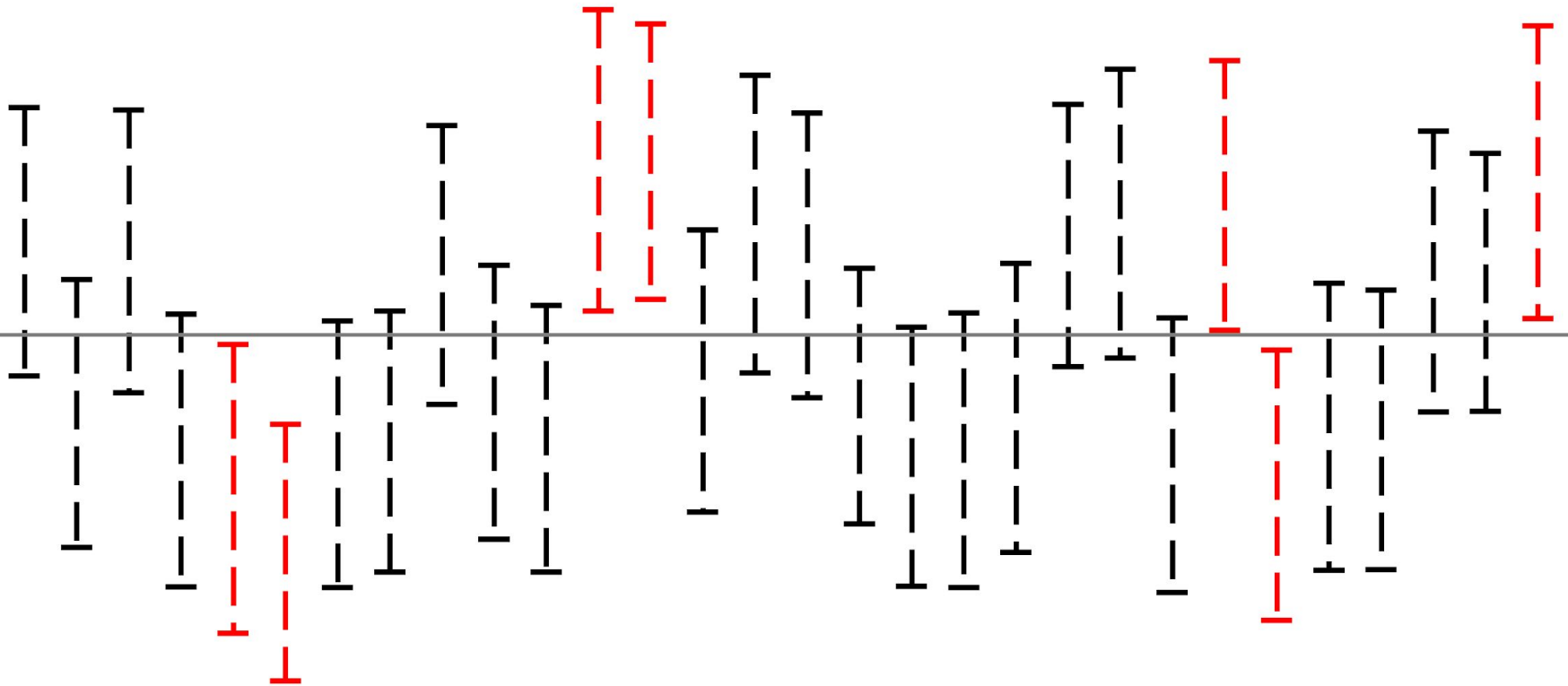
# Effects of Data Snooping: Stepwise Selection

$(X, Y) \sim N(0, I_{p+1}) \implies$  500 observations

$X_{\hat{j}}$  most correlated with  $Y \implies$  95% CI for slope

$p = 5$

Unadjusted: **76.9%**



# Post-selection Inference: Problem 1

There are  $p$  covariates and for each  $1 \leq j \leq p$

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

Want a valid CI for the **parameter/target**  $\beta_{\hat{j}}$ :

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

irrespective of how  $\hat{j}$  is chosen based on data.

# Post-selection Inference: Problem 2

For each  $M \subseteq \{1, 2, \dots, p\}$ ,

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

Want a valid CI for  $\beta_{\hat{j} \cdot \hat{M}}$ , a **coordinate** of  $\beta_{\hat{M}}$ :

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

irrespective of how  $\hat{M}$  with size  $\leq k$  and  $\hat{j} \in \hat{M}$  are chosen based on data.

# Solution 0: Sample Splitting

$\mathcal{D}_1 := \{(X_1, Y_1), \dots, (X_n, Y_n)\}$

Variable selection,  
Transformations etc.

$\{(X'_1, Y'_1), \dots, (X'_n, Y'_n)\} =: \mathcal{D}_2$

Finally, use once  
for inference.

- ✓ Allows **arbitrary exploration** in  $\mathcal{D}_1$ .
- ✗ **Cannot revise** the model after using  $\mathcal{D}_2$ .
- ✗ **Invalidity** when using multiple splits.
- ✗ Applicable only for **independent** data.

# Recap: Post-selection Inference

For each  $M \subseteq \{1, 2, \dots, p\}$ ,

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

Want a valid CI for  $\beta_{\hat{j} \cdot \hat{M}}$ , a **coordinate** of  $\beta_{\hat{M}}$ :

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

irrespective of how  $\hat{M}$  with size  $\leq k$  and  $\hat{j} \in \hat{M}$  are chosen based on data.

# Literature Review

Buehler and Feddersen (1963); Olshen (1973); Sen (1979); Rencher and Pun (1980); Freedman (1983); Sen and Saleh (1987); Dijkstra and Veldkamp (1988); Hurvich and Tsai (1990); Potscher (1991); Pfeiffer, Redd and Carroll (2017).



Cox (1965); Kabaila (1998); Hjort and Claeskens (2003); Claeskens and Carroll (2007); Berk et al. (2013); Lee, Sun, Sun and Taylor (2016); Tibshirani, Taylor, Lockhart and Tibshirani (2016); Bachoc, Preinerstorfer and Steinberger (2019); Rinaldo, Wasserman, G'Sell and Lei (2019).

# Roadmap

## Data Snooping: Effects and Examples

- Effects illustrated with stepwise selection.
- Data snooping in textbooks and practice.

## Formulation of the Problem

- The Problem & literature review for covariate selection.

## Solution for Covariate Selection

- Key contributions
- Simulations & main components of the theory.

## Example and Conclusions

# A Guiding Principle

Simultaneous Inference  $\Rightarrow$  Post-selection Inference  
FWER Control

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

- Simultaneity implies valid CIs for arbitrary selection  $\widehat{\mathbf{M}}$
- Simultaneity implies *infinite* revisions of a selection.
- Simultaneity also guarantees validity if multiple models are reported.



# A Key Result

Simultaneous Inference  
FWER Control  $\iff$  Post-selection  
Inference


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

**Theorem:** Simultaneous inference is *necessary* for valid Post-selection inference.

K et al. (2019) Valid Post-selection Inference in Model-free Linear Regression.  
*Annals of Statistics (Forthcoming)*.

# Solution 1: Uniform Adjustment

The classical interval for  $\beta_{j \cdot M}$  is

$$\left\{ \theta : \left| \frac{\sqrt{n}(\hat{\beta}_{j \cdot M} - \theta)}{\hat{\sigma}_{j \cdot M}} \right| \leq z_{\alpha/2} \right\}.$$


For *simultaneity*, inflate the confidence regions:

$$\widehat{\text{CI}}_{j \cdot M}^{\text{PoSI}} := \left\{ \theta : \left| \frac{\sqrt{n}(\hat{\beta}_{j \cdot M} - \theta)}{\hat{\sigma}_{j \cdot M}} \right| \leq K_{\alpha} \right\},$$

$$K_{\alpha} = (1 - \alpha) \text{ quantile of } \max_{|\mathbf{M}| \leq k, j \in \mathbf{M}} \left| \frac{\sqrt{n}(\hat{\beta}_{j \cdot M} - \beta_{j \cdot M})}{\hat{\sigma}_{j \cdot M}} \right|$$

# Our Contributions

We show

- in an **assumption-lean** setting,
- for **independent** and weakly **dependent** obs.,
- for  $p \gg n$  and maximal model size  $k = k_n$ ,
- for fixed or **random** covariates,

$K_\alpha$  can be estimated using bootstrap.

In the worst case,

$$K_\alpha \asymp \sqrt{k \operatorname{Log}(p/k)}. \quad (\operatorname{Log}(x) = 1 + \log x)$$

Berk et al. (2013): **Homoscedastic Gaussian** response, **fixed X**.

Bachoc et al (2019): **assumption-lean** but **fixed X** and **fixed p**.

# Simulation Examples: Revisiting Stepwise Selection

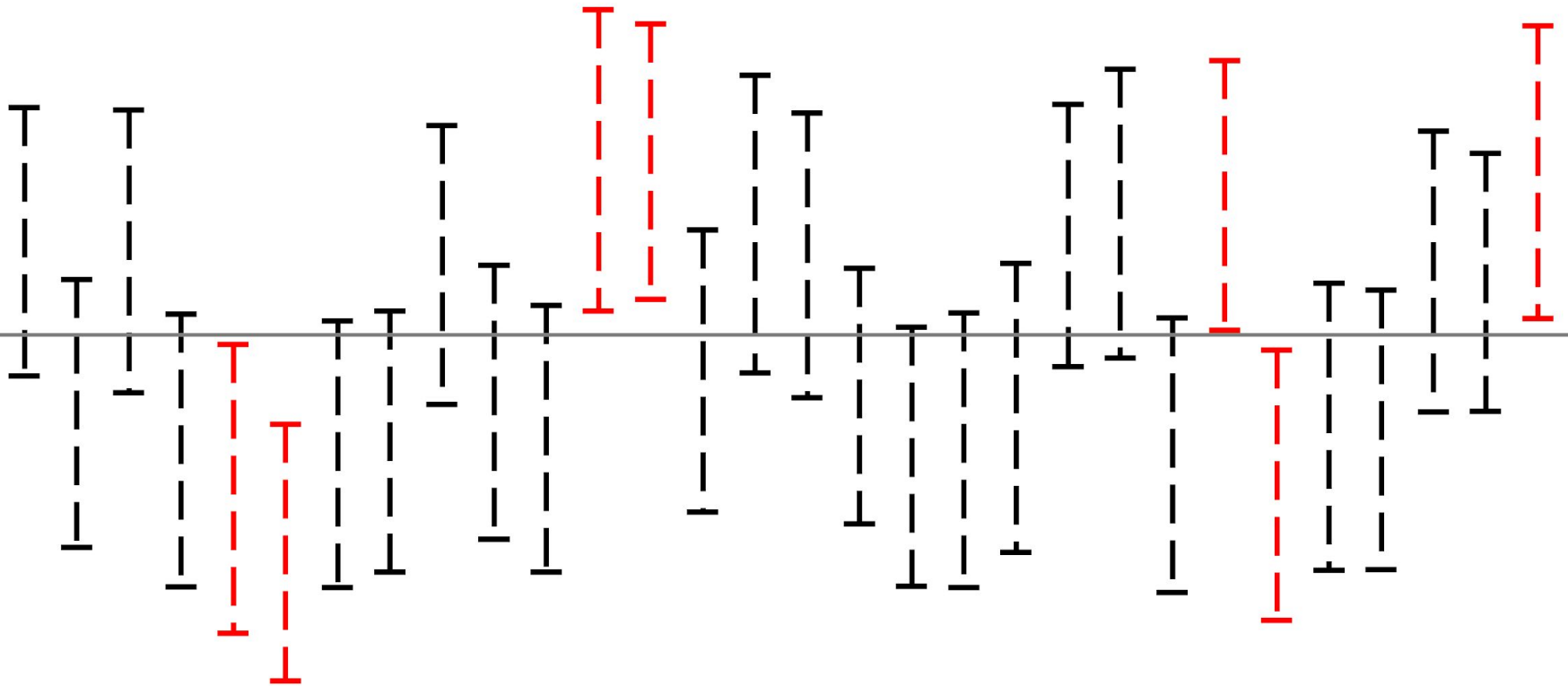
# Effects of Data Snooping: Stepwise Selection

$(X, Y) \sim N(0, I_{p+1}) \implies$  500 observations

$X_{\hat{j}}$  most correlated with  $Y \implies$  95% CI for slope

$p = 5$

Unadjusted: **76.9%**

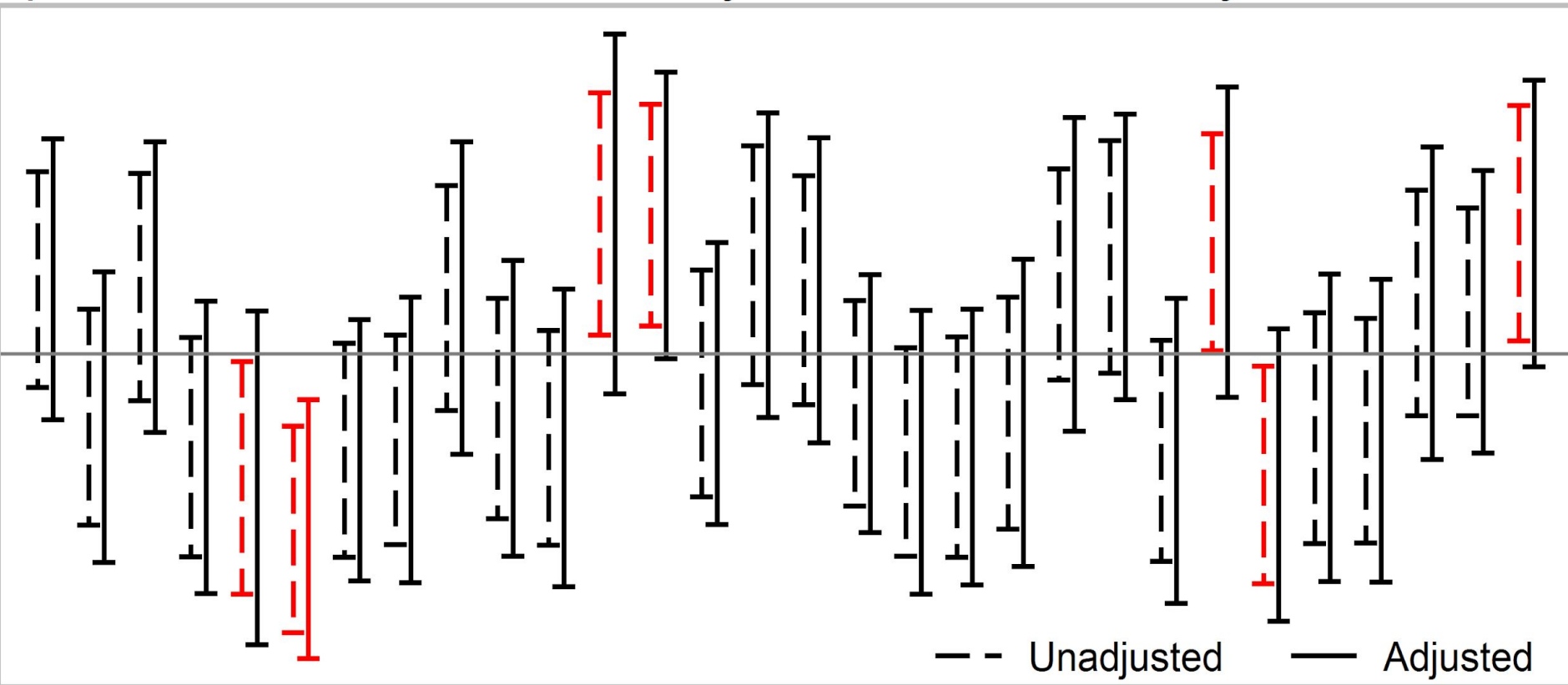


# Effects of Data Snooping: Stepwise Selection

$(X, Y) \sim N(0, I_{p+1}) \Rightarrow$  500 observations

$X_{\hat{j}}$  most correlated with  $Y \Rightarrow$  95% CI for slope

$p = 5$                       Unadjusted: **76.9%**                      Adjusted: **95.2%**



# Effects of Data Snooping: Stepwise Selection

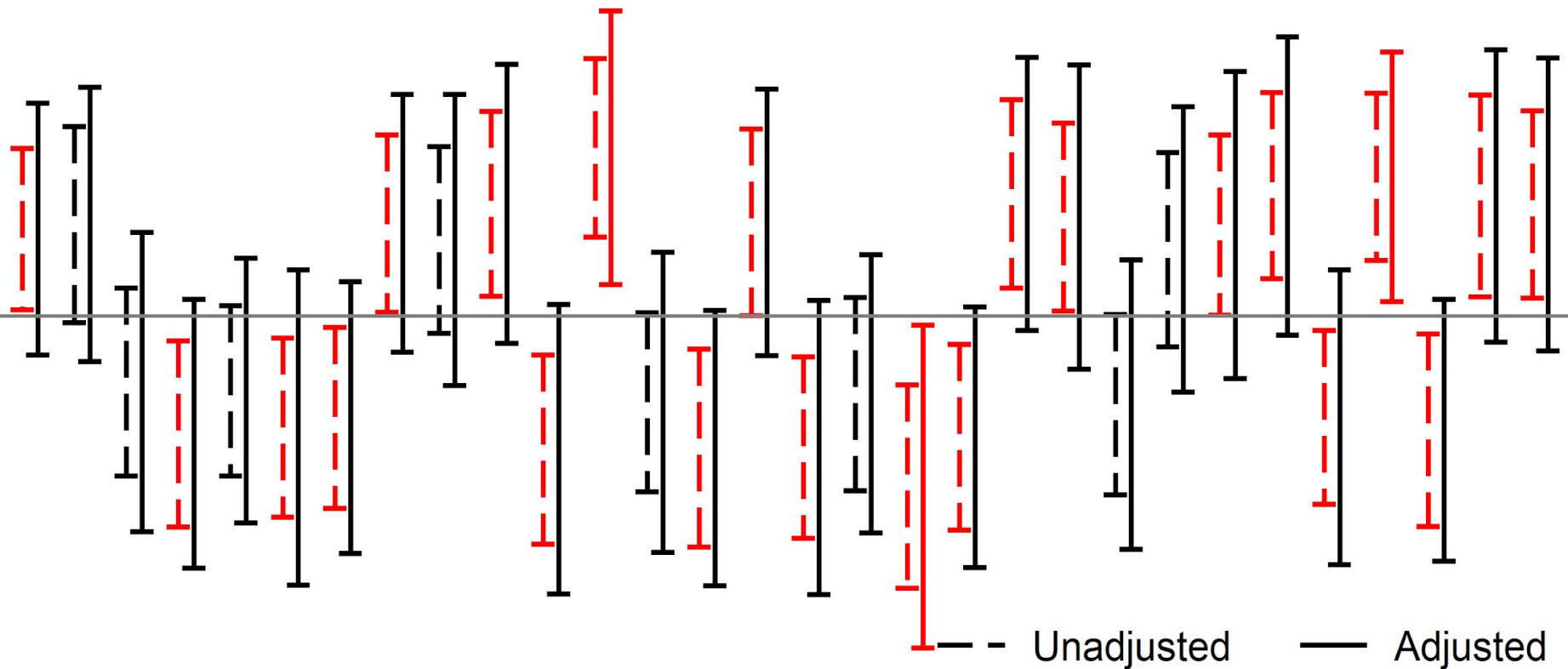
$(X, Y) \sim N(0, I_{p+1}) \Rightarrow$  500 observations

$X_{\hat{j}}$  most correlated with  $Y \Rightarrow$  95% CI for slope

$p = 20$

Unadjusted: **32.6%**

Adjusted: **93.4%**



# Effects of Data Snooping: Stepwise Selection

$$(X, Y) \sim N(0, I_{p+1})$$

$\Rightarrow$

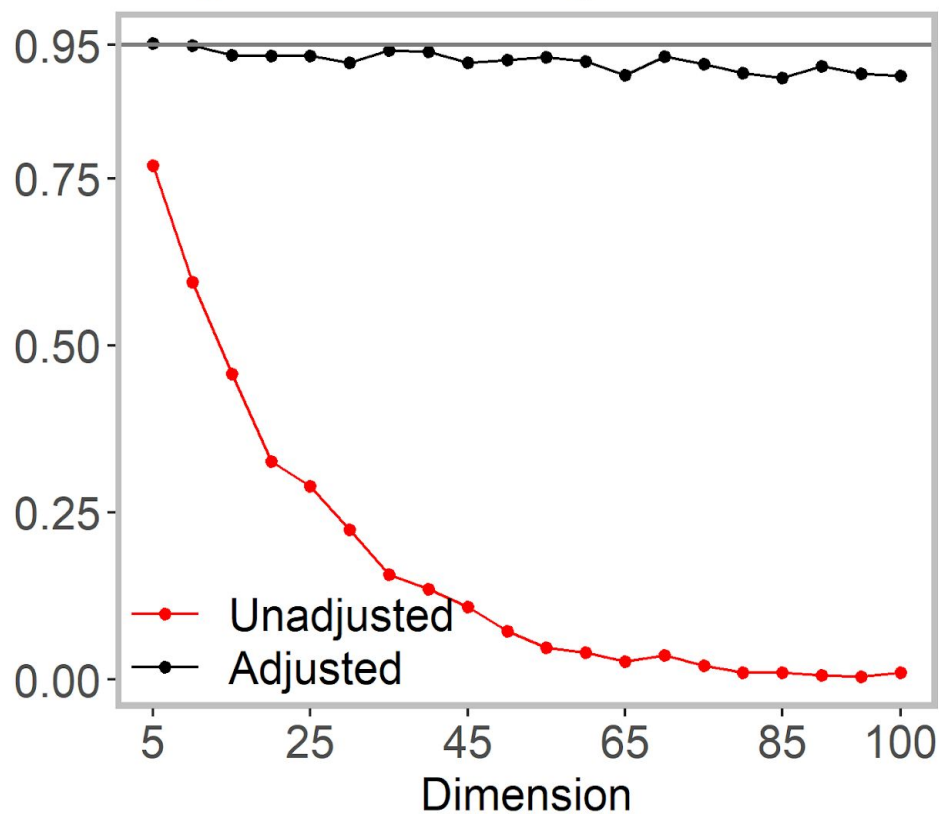
500 observations

$X_{\hat{j}}$  most correlated with  $Y$

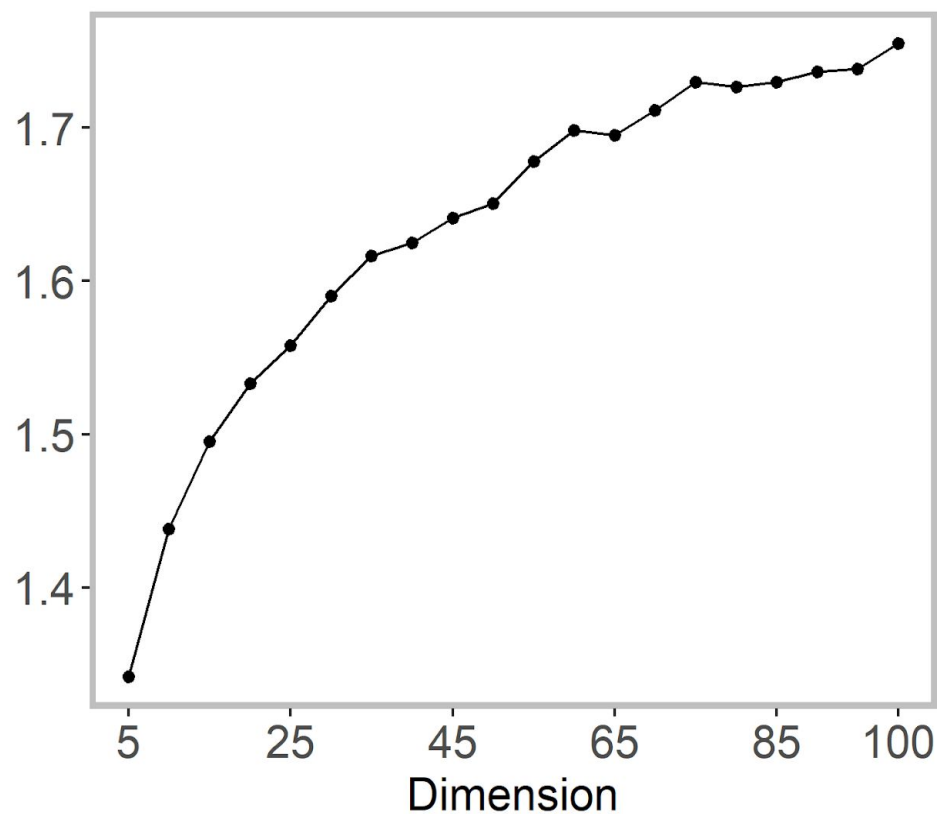
$\Rightarrow$

95% CI for slope

## Empirical Coverage



## Ratio of Widths





# Key Steps in the Proof

# Uniform Linear Representation

For independent, sub-Gaussian data  $(X_i, Y_i), 1 \leq i \leq n$

$$\max_{\substack{|\mathbf{M}| \leq k, \\ j \in \mathbf{M}}} \left| \frac{\sqrt{n}(\hat{\beta}_{j \cdot \mathbf{M}} - \beta_{j \cdot \mathbf{M}})}{\hat{\sigma}_{j \cdot \mathbf{M}}} - \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_{j \cdot \mathbf{M}}(X_i, Y_i) \right| = O_p \left( \frac{k \text{Log}(p/k)}{\sqrt{n}} \right).$$

- Doesn't require any **parametric model** assumptions.
- A **finite sample** result. Allows for diverging  $p, k$ .
- **Extends beyond** independent & sub-Gaussian data.

# Three Line Proof

$$\max_{\substack{|\mathbf{M}| \leq k, \\ j \in \mathbf{M}}} \left| \frac{\sqrt{n}(\hat{\beta}_{j \cdot \mathbf{M}} - \beta_{j \cdot \mathbf{M}})}{\hat{\sigma}_{j \cdot \mathbf{M}}} - \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_{j \cdot \mathbf{M}}(X_i, Y_i) \right| = O_p \left( \frac{k \text{Log}(p/k)}{\sqrt{n}} \right).$$

implies

$$\max_{\substack{|\mathbf{M}| \leq k, \\ j \in \mathbf{M}}} \left| \frac{\sqrt{n}(\hat{\beta}_{j \cdot \mathbf{M}} - \beta_{j \cdot \mathbf{M}})}{\hat{\sigma}_{j \cdot \mathbf{M}}} \right| \xrightarrow[\text{by triangle ineq.}]{\text{Close in Probability}} \max_{\substack{|\mathbf{M}| \leq k, \\ j \in \mathbf{M}}} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_{j \cdot \mathbf{M}}(X_i, Y_i) \right|$$

# Three Line Proof

$$\max_{\substack{|\mathbf{M}| \leq k, \\ j \in \mathbf{M}}} \left| \frac{\sqrt{n}(\hat{\beta}_{j \cdot \mathbf{M}} - \beta_{j \cdot \mathbf{M}})}{\hat{\sigma}_{j \cdot \mathbf{M}}} - \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_{j \cdot \mathbf{M}}(X_i, Y_i) \right| = O_p \left( \frac{k \text{Log}(p/k)}{\sqrt{n}} \right).$$

implies

$$\max_{\substack{|\mathbf{M}| \leq k, \\ j \in \mathbf{M}}} \left| \frac{\sqrt{n}(\hat{\beta}_{j \cdot \mathbf{M}} - \beta_{j \cdot \mathbf{M}})}{\hat{\sigma}_{j \cdot \mathbf{M}}} \right| \xrightarrow[\text{by triangle ineq.}]{\text{Close in Probability}} \max_{\substack{|\mathbf{M}| \leq k, \\ j \in \mathbf{M}}} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_{j \cdot \mathbf{M}}(X_i, Y_i) \right|$$

$$(k \text{Log}(p/k))^5 = o(n)$$

Close in  
Distribution  
by HDCLT

$$\max_{|\mathbf{M}| \leq k, j \in \mathbf{M}} |G_{j \cdot \mathbf{M}}|$$

K et al. (2018) High-dimensional CLT: Improvements,  
Non-uniform Extensions and Large Deviations.  
arXiv:1806.06153

# Three Line Proof

$$\max_{\substack{|\mathbf{M}| \leq k, \\ j \in \mathbf{M}}} \left| \frac{\sqrt{n}(\hat{\beta}_{j \cdot \mathbf{M}} - \beta_{j \cdot \mathbf{M}})}{\hat{\sigma}_{j \cdot \mathbf{M}}} - \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_{j \cdot \mathbf{M}}(X_i, Y_i) \right| = O_p \left( \frac{k \text{Log}(p/k)}{\sqrt{n}} \right).$$

implies

$$\max_{\substack{|\mathbf{M}| \leq k, \\ j \in \mathbf{M}}} \left| \frac{\sqrt{n}(\hat{\beta}_{j \cdot \mathbf{M}} - \beta_{j \cdot \mathbf{M}})}{\hat{\sigma}_{j \cdot \mathbf{M}}} \right| \xrightarrow[\text{by triangle ineq.}]{\text{Close in Probability}} \max_{\substack{|\mathbf{M}| \leq k, \\ j \in \mathbf{M}}} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_{j \cdot \mathbf{M}}(X_i, Y_i) \right|$$

Close in  
Distribution  
by triangle ineq.

Close in  
Distribution  
by HDCLT

Justifies bootstrap!

$$\max_{|\mathbf{M}| \leq k, j \in \mathbf{M}} |G_{j \cdot \mathbf{M}}|$$

# Roadmap

## Data Snooping: Effects and Examples

- Effects illustrated with stepwise selection.
- Data snooping in textbooks and practice.

## Formulation of the Problem

- The Problem & literature review for covariate selection.

## Solution for Covariate Selection

- Key contributions
- Simulations & main components of the theory.

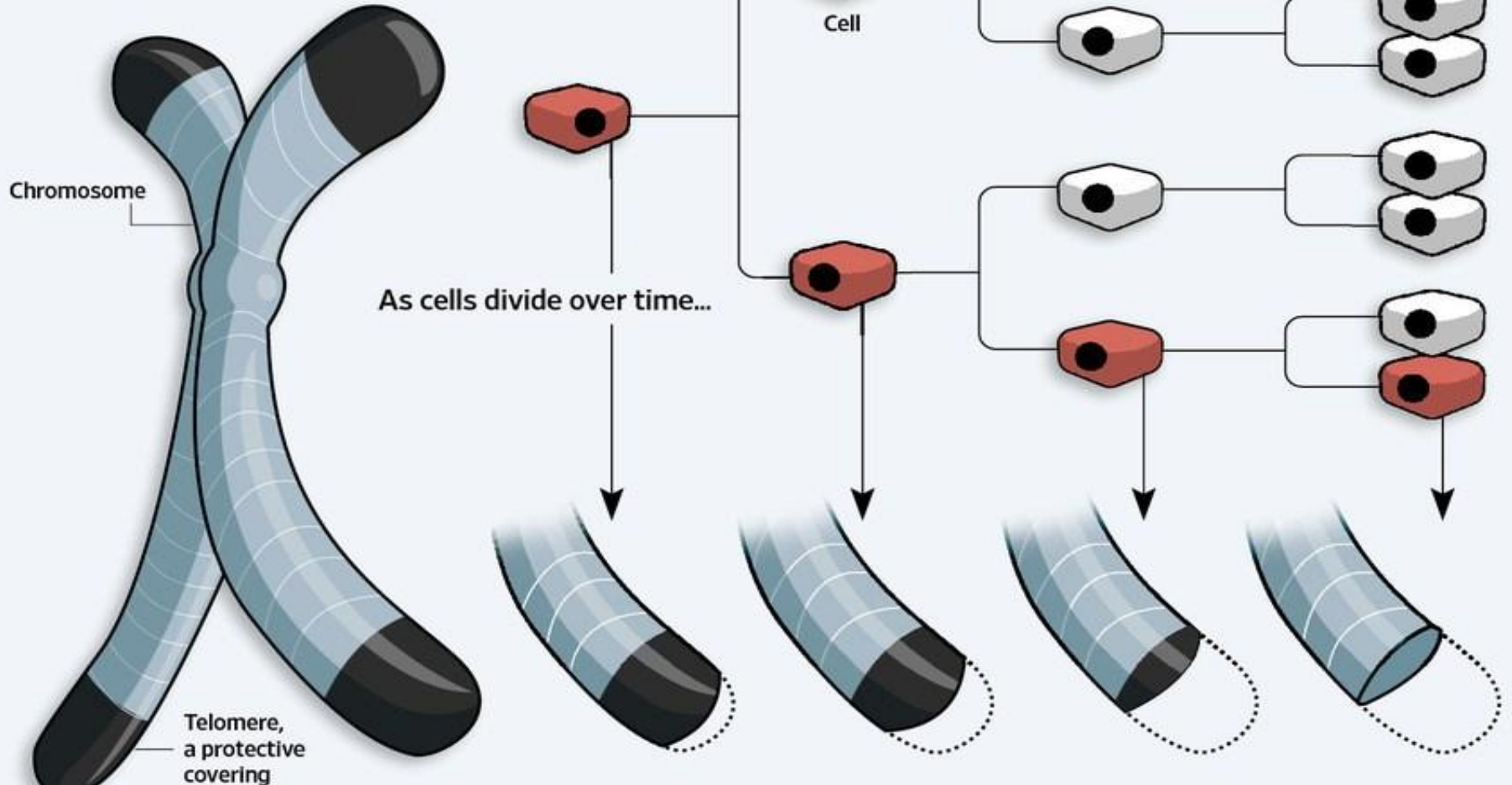
## Example and Conclusions

- Real Data Example, Extensions & Summary.

# Telomere Length Analysis

## 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.'



# 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.**”



# Telomere Length Analysis

TL inheritance patterns based on 246 families.

❖ *Dependent Variable*: **MTL** (Mean telomere length)

❖ *Child Variables*:

➤ **Sex**

➤ **Age**

❖ *Parental Variables*:

➤ **mMTL** (mother MTL)

➤ **fMTL** (father MTL)

➤ **MAC** (mother's age at conception)

➤ **PAC** (father's age at conception)

(Additionally, 15 interaction variables were considered.)

# Adjusted Inference: Telomere Length Analysis

Covariate	Unadjusted	Adjusted
AGE	✓	✗
mMTL	✓	✓
fMTL	✓	✓
MAC	✓	✗
PAC	✗	✗

✓ : Significant at 5% level

✗ : Insignificant at 5% level

# Summary and Conclusions

- ❖ Data snooping contributes to replicability crisis.
- ❖ Inference is possible after data snooping.

Classical Framework	New Framework
Fix the test & model	Fix a universe
Collect the data	Collect the data

- ❖ The framework allows for
  - Misspecified models; Random covariates;
  - Dependent data; High-dimensional features;
  - Variable Transformations;
  - M-estimators: **logistic/Poisson/Quantile/Cox.**

# References

Kuchibhotla A., Brown L., Buja A., Cai J., George E., Zhao L. (2019)  
Valid Post-selection Inference in Model-free Linear Regression.  
*Annals of Statistics (Forthcoming)*.

Kuchibhotla A., Brown L., Buja A., George E., Zhao L. (2018)  
A Model Free Perspective for Linear Regression: Uniform-in-model Bounds  
for Post Selection Inference. *arXiv:1802.05801*

Kuchibhotla A. (2018)  
Deterministic Inequalities for Smooth M-estimators. *arXiv:1809.05172*

Kuchibhotla A., Mukherjee S., and Banerjee D. (2018)  
High-dimensional CLT: Improvements, Non-uniform Extensions and Large  
Deviations. *arXiv:1806.06153*

Kuchibhotla A., Brown L., Buja A., Cai J. (2019)  
All of Linear Regression. *arXiv:1910.06386*

Thank you for your attention

# Post-selection for Transformations

For each  $g \in \mathcal{G} \subseteq L_2(Y)$ ,

$$\beta_g := \operatorname{argmin}_{\theta \in \mathbb{R}^p} \mathbb{E}[(g(Y) - X^\top \theta)^2].$$

Want a valid CI for  $\beta_{1.\hat{g}}$ , a coordinate of  $\beta_{\hat{g}}$

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

irrespective of how  $\hat{g} \in \mathcal{G}$  is chosen based on data.

# Solution for Transformations

Inflate the classical intervals,

$$\left\{ \theta \in \mathbb{R} : \left| \frac{\sqrt{n}(\hat{\beta}_{1.g} - \theta)}{\hat{\sigma}_{1.g}} \right| \leq K_\alpha \right\},$$

with  $K_\alpha$  being the  $(1 - \alpha)$  quantile of

$$\max_{g \in \mathcal{G}} \left| \frac{\sqrt{n}(\hat{\beta}_{1.g} - \beta_{1.g})}{\hat{\sigma}_{1.g}} \right|,$$

Bootstrap applies with **validity** for “nice” function classes  $\mathcal{G}$ , for example, Box-Cox family.

# Transformations: Boston Housing Data

This dataset has 506 census tracts with 13 features.

- ❖ *Dependent Variable*: **MV**, Median value of house.
- ❖ *Covariate of Interest*: **NOX**, Nitrogen Oxide Conc.
- ❖ *Structural Variables*: **RM** (No. of rms), **AGE** (% of homes bfr 1940).
- ❖ *Neighborhood Variables*: **CRIM** (Crime rate), **ZN** (% of res. land zoned for lots > 25K ft<sup>2</sup>), **INDUS** (% non-retail business acres per twn), **RIVER** (Charles river dummy), **TAX** (Property tax rate), **PTRATIO** (Pupil-teacher ratio), **B** (Racial diversity), **LSTAT** (% of lwr socio-econ. status of population).
- ❖ *Accessibility Variables*: **DIS** (Distance to Employment Ctr.), **RAD** (Distance to Radial Highway).

# Transformations: Boston Housing Data

Harrison and Rubinfeld (1978)<sup>★</sup> write

to determine the best fitting functional form. Comparing models with either median value of owner-occupied homes ( $MV$ ) or  $\text{Log}(MV)$  as the dependent variable, we found that the semilog version provided a slightly better fit. Using  $\text{Log}(MV)$  as the dependent variable, we concentrated on estimating a nonlinear term in  $NOX$ ; i.e., we included  $NOX^p$  in the equation, where  $p$  is an unknown

The statistical fit in the equation was best when  $p$  was set equal to 2.0,

<sup>★</sup>H & R (1978) Hedonic housing prices and the demand for clean air.



# Adjusted Inference: Boston Housing

Covariate	Unadjusted	Adjusted	Covariate	Unadjusted	Adjusted
NOX <sup>2</sup>	✓	✓	TAX	✓	✓
RM	✓	✓	PTRATIO	✓	✓
AGE	✗	✗	B	✓	✗
CRIM	✓	✓	LSTAT	✓	✓
ZN	✗	✗	DIS	✓	✓
INDUS	✗	✗	RAD	✓	✓
RIVER	✓	✗			

✓ : Significant at 5% level

✗ : Insignificant at 5% level

# Implications of PoSI for Applications

Similar to Boston housing data and TL data,

How do conclusions in applied data analysis change when **exploration is accounted for?**

# High-dimensional CLT

Anderson, Hall and Titterington (1998, JSPI):

Let  $\mathcal{R}$  be the class of all rectangles in  $\mathbb{R}^p$ .

A proof will be given in Section 3.2. It follows from this result and the Theorem that if (2.1), (2.6) and (2.7) hold,

$$\sup_{B \in \mathcal{R}} \left| P(S \in B) - \int_B \phi(x) dx \right| = O\{n^{-1/2}(\log p)^{3/2}\}.$$

# Maximal Inequalities

- High-dimensional rates are sensitive to tails.
- What is the order of

$$\mathbb{E} \left[ \max_{1 \leq j \leq d} \left| \frac{1}{n} \sum_{i=1}^n X_{i,j} \right| \right], \quad \text{subject to}$$

$$\max_{1 \leq j \leq d} \text{Var}(X_{1,j}) \leq A^2 \quad \text{and} \quad \mathbb{E} [\|X_1\|_\infty^q] \leq B^q?$$

# Multiple Testing under Dependence

- Multiple testing often requires independence.
- FWER is possible under arbitrary dependence.
- What about FDR control?
- How does BH procedure behave?



# Arun Kumar Kuchibhotla

University of Pennsylvania

Webpage: <https://arun-kuchibhotla.github.io/>

Email: [arunku@upenn.edu](mailto:arunku@upenn.edu)

# Case Study 3: Model Building

## Modeling Home Prices Using Realtor Data

Journal of  
Statistics  
Education  
2008

Iain Pardoe

Lundquist College of Business, University of Oregon

- 76 Oregon homes and 12 features.
- **Try a linear model** with 12 predictors “as is”.
- Residuals imply **non-linearity**.  $\text{Age} \rightarrow \text{Age}^2$ .
- Bath and Bed also have **high p-values**, so add an interaction  $\text{Bath} \times \text{Bed}$  to the model.
- Price is skewed suggesting a **log-transformation**.

# Case Study 3: Transformations

In the context of curve fitting to bivariate data, Stine and Foster (2014)<sup>★</sup> on page 515, write

“Picking a transformation requires practice, and you **may need to try several to find one** that is interpretable and captures the pattern in the data.”

<sup>★</sup>Stine and Foster, *Statistics for Business: Decision Making and Analysis*.



# Extensions

- The solution applies to most other estimation problems.
- Examples include **logistic/Poisson** regression, **quantile** regression and **Cox** regression.
- Solution also applies to **transformation of variables**.