

Confounder Adjustment in Multiple Hypothesis Testing

Qingyuan Zhao

Department of Statistics, Stanford University

January 28, 2016

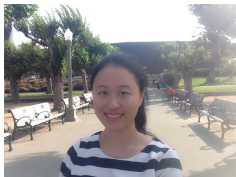
Slides are available at <http://web.stanford.edu/~qyzhao/>.

Collaborators

Confounder
Adjustment

Qingyuan
Zhao

Jingshu Wang



Trevor Hastie



Art Owen



Introduction

Background

Motivating
Examples

Previous Work

Model and
Inference

Model and
Identifiability

Estimation
Hypothesis Tests

Numerical
Examples

Summary

Microarray experiments

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating

Examples

Previous Work

Model and
Inference

Model and
Identifiability

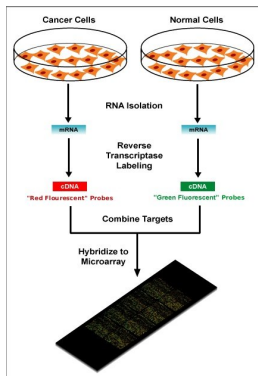
Estimation

Hypothesis Tests

Numerical

Examples

Summary



- Responses: normalized gene expression level.
- Primary variables (variables of interest): treatment, disease status, etc.
- Control covariates: age, gender, batch, date, etc.

Microarray data analysis

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating

Examples

Previous Work

Model and
Inference

Model and
Identifiability

Estimation

Hypothesis Tests

Numerical
Examples

Summary

Biologist: “Which genes are (causally) related to this disease?”
Statistician: “Let me run some analysis.”

Two common practices

- ① **Sparse regression:** regress the primary variable on the genes. More common for SNP data and predictive tasks.
- ② **Association tests/screening (this talk):** for each gene, perform a significance test of correlation with the primary variable.

Statistician: “Here a short list of candidate genes with **false discovery rate (FDR) $\leq 20\%$** .”

Biologist: “Good, let me validate these discoveries.”

Concerns

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

J. P. Ioannidis. [Why most published research findings are false.](#)
Chance, 18(4):40–47, 2005

Two major challenges to reproducibility in genetic screening:

- ① **Correlated tests:** Is the FDR still controlled? If not, can we correct the statistical analysis?
 - Well studied in the last 15 years [Benjamini and Yekutieli, 2001, Storey et al., 2004, Efron, 2007, Fan et al., 2012].
- ② **Confounded tests (this talk):** the individual association tests are biased in presence of unobserved confounders. Can we still provide a good candidate list?
 - Equally long history [e.g. Alter et al., 2000, Price et al., 2006]. Still many open questions.

Confounding

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

Brief history

- Fisher [1935] first uses the term in experiment designs.
- Kish [1959] first uses its modern meaning:
A mixing of effects of unobserved extraneous factors (called confounders) with the effect of interest.
- Huge literature, but mostly in causal inference.

Aliases for confounders in genetic screening:

- “systematic ancestry differences” [Price et al., 2006].
- “batch effects” (widely used by biologists).
- “surrogate variables” [Leek and Storey, 2007, 2008].
- “unwanted variation” [Gagnon-Bartsch and Speed, 2012].
- “latent effects” [Sun et al., 2012].

Example 1: gender study

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

Which genes are more expressed in male/female?

A microarray experiment by Vawter et al. [2004]:

- Postmortem samples from the brains of 10 individuals.
- For each individual, 3 samples from different cortices.
- Each sample is sent to 3 different labs for analysis.
- Two different microarray platforms are used by the labs.

In total, $10 \times 3 \times 3 = 90$ samples.

This example was first used by Gagnon-Bartsch and Speed [2012] to demonstrate the importance to “remove unwanted variation”.

Screening

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

Notation

- \mathbf{Y} : $n \times p$ matrix of gene expression.
 - \mathbf{X} : $n \times 1$ vector of gender.
-
- Simplest association test:
Regress each column of \mathbf{Y} (gene) on \mathbf{X} .
 - In R, run `summary(lm(Y~X))`.
 - Equivalent to a two-sample t -test with equal variance.

Histogram of t-statistics

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating
Examples

Previous Work

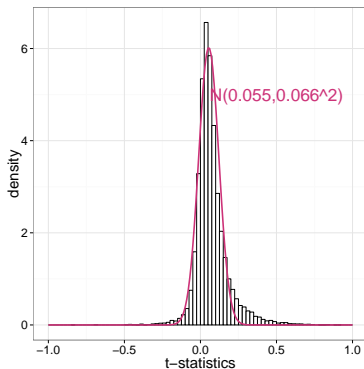
Model and
Inference

Model and
Identifiability

Estimation
Hypothesis Tests

Numerical
Examples

Summary



Skewed and very underdispersed.

What happened?

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating
Examples

Previous Work

Model and
Inference

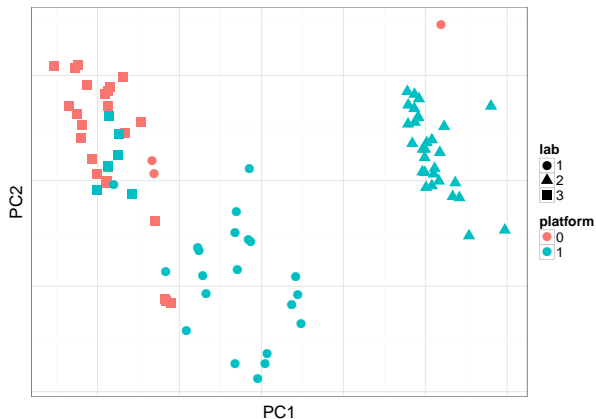
Model and
Identifiability

Estimation

Hypothesis Tests

Numerical
Examples

Summary



Association test

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating
Examples

Previous Work

Model and
Inference

Model and
Identifiability

Estimation
Hypothesis Tests

Numerical
Examples

Summary

Notation

- **Y**: $n \times p$ matrix of gene expression.
- **X**: $n \times 1$ vector of gender.
- **Z**: $n \times d$ matrix of control covariates (lab and platform).

- Modified association test:

Regress each column of **Y** (gene) on **X** and **Z**.

- In R, run `summary(lm(Y~X+Z))`.
- Report the significance of the coefficients of **X**.

Histogram of t-statistics

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating
Examples

Previous Work

Model and
Inference

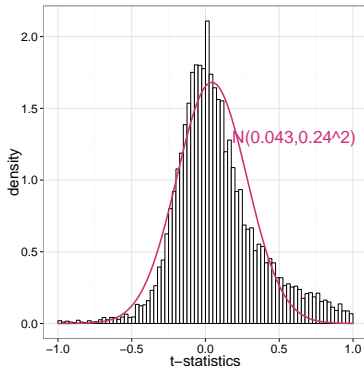
Model and
Identifiability

Estimation

Hypothesis Tests

Numerical
Examples

Summary



Better, but still problematic.

Reasonable guess: **there are more unobserved confounders!**

Example 2: COPD study

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating
Examples

Previous Work

Model and
Inference

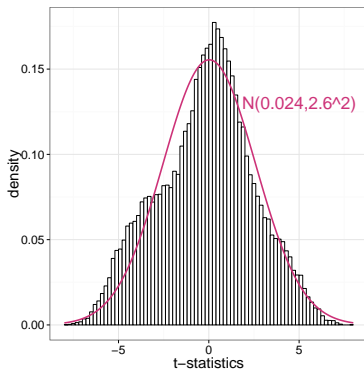
Model and
Identifiability

Estimation
Hypothesis Tests

Numerical
Examples

Summary

- COPD = chronic obstructive pulmonary disease.
- Singh et al. [2011] tried to find genes associated with the severity of COPD (moderate or severe).



Overdispersed and skewed.

Example 3: Mutual fund selection

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating
Examples

Previous Work

Model and
Inference

Model and
Identifiability

Estimation
Hypothesis Tests

Numerical
Examples

Summary

Barras et al. [2010] used the following model to select mutual funds:

$$Y_{it} = \alpha_i + \gamma_i^T \mathbf{Z}_t + e_{it}, \quad i = 1, \dots, n, t = 1, \dots, p.$$

- Y_{it} : observed log-return of fund i at time t .
- α_i : risk-adjusted return (Goal: find funds with positive α).
- \mathbf{Z}_t : systematic risk factors.

They assumed:

- α is sparse (Berk and Green equilibrium);
- No unobserved risk factors (is that possible/necessary?).

Idea 0: Remove the largest principal component(s)

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating
Examples

Previous Work

Model and
Inference

Model and
Identifiability

Estimation
Hypothesis Tests

Numerical
Examples

Summary

EIGENSTRAT [Price et al., 2006]

Regression model:

$$\mathbf{Y}_{n \times p} = \mathbf{X}_{n \times 1} \boldsymbol{\beta}_{p \times 1}^T + \mathbf{Z}_{n \times r} \boldsymbol{\Gamma}_{p \times r}^T + \mathbf{E}_{n \times p},$$

where \mathbf{Z} is the first r PC(s) of \mathbf{Y} .

- Motivation: in SNP, the largest PC(s) usually correspond to ancestry difference.
- Weakness: can easily remove true signals.

Idea 1: Use control genes

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating
Examples

Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

Same regression model:

$$\mathbf{Y}_{n \times p} = \mathbf{X}_{n \times 1} \beta_{p \times 1}^T + \mathbf{Z}_{n \times r} \mathbf{\Gamma}_{p \times r}^T + \mathbf{E}_{n \times p},$$

RUV2 [Gagnon-Bartsch and Speed, 2012]

If we know $\beta_c = \mathbf{0}$ (negative controls),

- ① Run PCA on $col_c(\mathbf{Y})$ to obtain \mathbf{Z} .
- ② Run the regression for $col_{-c}(\mathbf{Y})$.

- Example: bacterial RNAs (spike-in controls).
- Limited to the availability and number of negative controls.

Idea 2: Sparsity

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating

Examples

Previous Work

Model and
Inference

Model and
Identifiability

Estimation

Hypothesis Tests

Numerical

Examples

Summary

Same regression model:

$$\mathbf{Y}_{n \times p} = \mathbf{X}_{n \times 1} \boldsymbol{\beta}_{p \times 1}^T + \mathbf{Z}_{n \times r} \boldsymbol{\Gamma}_{p \times r}^T + \mathbf{E}_{n \times p},$$

Idea: If $\boldsymbol{\beta}$ contains actual effects, it should be a sparse vector.

SVA [Leek and Storey, 2008]

Iterate between

- 1 Weighted PCA on \mathbf{Y} (based on how likely $\boldsymbol{\beta} = \mathbf{0}$).
- 2 Regress \mathbf{Y} on \mathbf{X} and the estimated PCs.

- Does not always converge.

Idea 2: Sparsity

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating
Examples

Previous Work

Model and
Inference

Model and
Identifiability

Estimation
Hypothesis Tests

Numerical
Examples

Summary

Same regression model:

$$\mathbf{Y}_{n \times p} = \mathbf{X}_{n \times 1} \boldsymbol{\beta}_{p \times 1}^T + \mathbf{Z}_{n \times r} \boldsymbol{\Gamma}_{p \times r}^T + \mathbf{E}_{n \times p},$$

Idea: If $\boldsymbol{\beta}$ contains actual effects, it should be a sparse vector.

LEAPP [Sun, Zhang, and Owen, 2012]

- 1 Run PCA on the residuals of $\mathbf{Y} \sim \mathbf{X}$.
- 2 Run a sparse regression.

Our contributions: a unifying framework

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

Missing in previous methods:

- Explicit assumptions on the latent variables.
- Model identification conditions.
- Theoretical guarantees.
- Multiple primary and secondary covariates.
- Practical guidelines: when is confounder adjustment necessary/useful?

Statistical model for confounding

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

- Linear model for the responses (e.g. gene expression)

$$\mathbf{Y}_{n \times p} = \mathbf{X}_{n \times 1} \beta_{p \times 1}^T + \mathbf{Z}_{n \times r} \mathbf{\Gamma}_{p \times r}^T + \mathbf{E}_{n \times p},$$

- \mathbf{X} : primary variable (disease, treatment, gender, etc.);
 - \mathbf{Z} : unobserved confounders;
 - β : primary effects that we are interested in.
- Missing in the literature: dependence of \mathbf{Z} and \mathbf{X}

$$\mathbf{Z}_{n \times r} = \mathbf{X}_{n \times 1} \alpha_{r \times 1}^T + \mathbf{W}_{n \times r},$$

- Additional distributional assumptions:

$$X_i \stackrel{\text{i.i.d.}}{\sim} \text{mean } 0, \text{ variance } 1, \quad i = 1, \dots, n,$$

$$\mathbf{E} \stackrel{\text{i.i.d.}}{\sim} \mathbf{N}(\mathbf{0}, \mathbf{\Sigma}), \quad \mathbf{E} \perp (\mathbf{X}, \mathbf{Z}), \quad \mathbf{\Sigma} = \text{diag}(\{\sigma_j^2\}_{j=1}^p),$$

$$\mathbf{W} \stackrel{\text{i.i.d.}}{\sim} \mathbf{N}(\mathbf{0}, \mathbf{I}_r), \quad \mathbf{W} \perp \mathbf{X}.$$

Marginal effects and direct effects

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

The model can be rewritten as

$$\mathbf{Y}_{n \times p} = \mathbf{X}_{n \times 1} (\boldsymbol{\beta}_{p \times 1} + \boldsymbol{\Gamma}_{p \times r} \boldsymbol{\alpha}_{r \times 1})^T + (\mathbf{W}\boldsymbol{\Gamma} + \mathbf{E}),$$

which gives the population identity

$$\boldsymbol{\tau}_{p \times 1} = \boldsymbol{\beta} + \boldsymbol{\Gamma}\boldsymbol{\alpha}.$$

- $\boldsymbol{\tau}$: marginal effects.
- $\boldsymbol{\beta}$: direct effects (more meaningful).

COPD data: marginal effects vs. direct effects

Confounder
Adjustment

Qingyuan
Zhao

Introduction

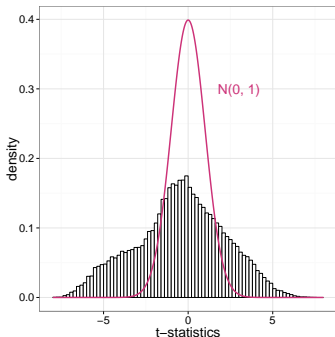
Background
Motivating
Examples
Previous Work

Model and
Inference

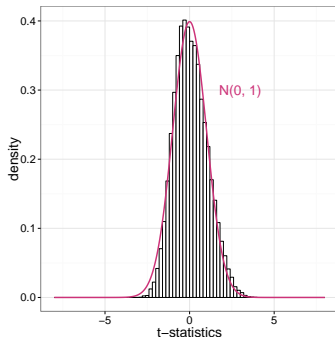
Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary



(a) Before adjustment
(t -statistics for $\tau_j = 0$).



(b) After adjustment
(t -statistics for $\beta_j = 0$).

Identifiability of β

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

To identify α and β from

$$\tau_{p \times 1} = \beta_{p \times 1} + \Gamma \alpha_{r \times 1},$$

there are p equations but $p + r$ parameters.

Proposition [Wang, Z., Hastie, and Owen, 2015]

Suppose Γ can be identified. β is identifiable under either of the two following conditions:

- 1 Negative control: for a known negative control set \mathcal{C} ,

$$\beta_{\mathcal{C}} = \mathbf{0}, |\mathcal{C}| \geq r, \text{rank}(\Gamma_{\mathcal{C}}) = r.$$

- 2 Sparsity: $\|\beta\|_0 \leq \lfloor (p - r)/2 \rfloor$ (the maximum breakdown point),

$$\text{rank}(\Gamma_{\mathcal{C}}) = r, \forall \mathcal{C} \subset \{1, \dots, p\} \text{ such that } |\mathcal{C}| = r.$$

Rotation

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

Householder transformation

$$\mathbf{X}_{n \times 1} = \mathbf{Q}\mathbf{R}$$

where $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is orthogonal with $\mathbf{R} = (\|\mathbf{X}\|_2, 0, \dots, 0)^T$.

- For simplicity, assume $\|\mathbf{X}\|_2 = \sqrt{n}$.
- Can be easily extended to multiple variables \mathbf{X} .

Rotation (LEAPP)

Left-Multiply \mathbf{Q}^T to $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta}^T + \mathbf{Z}\boldsymbol{\Gamma}^T + \mathbf{E}$, we get

$$\text{row}_1(\mathbf{Q}^T \mathbf{Y}) \sim N(\sqrt{n}(\boldsymbol{\beta} + \boldsymbol{\Gamma}\boldsymbol{\alpha}), \boldsymbol{\Gamma}\boldsymbol{\Gamma}^T + \boldsymbol{\Sigma}),$$

$$\text{row}_{-1}(\mathbf{Q}^T \mathbf{Y}) \stackrel{\text{i.i.d.}}{\sim} N(\mathbf{0}, \boldsymbol{\Gamma}\boldsymbol{\Gamma}^T + \boldsymbol{\Sigma}).$$

Two-step estimation

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

- 1 Run factor analysis for

$$\text{row}_1(\mathbf{Q}^T \mathbf{Y}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{\Gamma} \mathbf{\Gamma}^T + \mathbf{\Sigma})$$

to obtain $\hat{\mathbf{\Gamma}}$ and $\hat{\mathbf{\Sigma}}$. Identifiability follows from classical results in factor analysis [e.g. Anderson and Rubin, 1956].

- 2 Run linear regression for the marginal effects

The diagram shows the equation
$$\frac{\text{row}_1(\mathbf{Q}^T \mathbf{Y})_{p \times 1}}{\sqrt{n}} = \hat{\mathbf{\Gamma}}_{p \times r} \boldsymbol{\alpha}_{r \times 1} + \beta_{p \times 1} + \tilde{\mathbf{E}}_1 / \sqrt{n}$$
 with colored boxes and arrows. The left side is in a pink box. The first term on the right, $\hat{\mathbf{\Gamma}}_{p \times r}$, is in a blue box with a blue arrow pointing to it from the label 'design matrix' below. The second term, $\boldsymbol{\alpha}_{r \times 1}$, is in a green box with a green arrow pointing to it from the label 'coefficients' below. The remaining terms, $\beta_{p \times 1} + \tilde{\mathbf{E}}_1 / \sqrt{n}$, are not boxed. A red arrow points from the label 'response' below to the pink box on the left.

response

design matrix

coefficients

How accurate is $\hat{\Gamma}$?

Confounder
Adjustment

Qingyuan
Zhao

Introduction
Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

Assumptions

- High-dimensional data: $n \rightarrow \infty, p \rightarrow \infty$.
- Assume that the factors are strong enough:
$$\lim_{p \rightarrow \infty} \frac{1}{p} \mathbf{\Gamma}^T \mathbf{\Sigma}^{-1} \mathbf{\Gamma} \text{ exists and is positive definite.}$$
- Consistent estimate of r [Bai and Ng, 2002].

Theoretical Results for MLE

- Consistent estimate of $\mathbf{\Gamma}$ and $\mathbf{\Sigma}$ [Bai and Li, 2012] and
$$\sqrt{n}(\hat{\mathbf{\Gamma}}_j - \mathbf{\Gamma}_j) \xrightarrow{d} N(\mathbf{0}, \sigma_j^2 \mathbf{I}_r), \quad \sqrt{n}(\hat{\sigma}_j - \sigma_j) \xrightarrow{d} N(0, 2\sigma_j^4),$$
- Uniform consistency if $n^k/p \rightarrow \infty$ for some $k > 0$ [Wang, Z., Hastie, and Owen, 2015].

Strategy 1: Estimate β via negative controls

Recall the marginal effects are

$$\frac{\tilde{\mathbf{Y}}_{p \times 1}^T}{\sqrt{n}} = \mathbf{\Gamma}_{p \times r} \boldsymbol{\alpha}_{r \times 1} + \beta_{p \times 1} + \tilde{\mathbf{E}}_1 / \sqrt{n}$$

response design matrix coefficients

In the negative control scenario, we know $\beta_c = \mathbf{0}$.

Generalized Least Squares (GLS) estimator

$$\hat{\boldsymbol{\alpha}}^{\text{NC}} = (\hat{\mathbf{\Gamma}}_c^T \hat{\mathbf{\Sigma}}_c^{-1} \hat{\mathbf{\Gamma}}_c)^{-1} \hat{\mathbf{\Gamma}}_c^T \hat{\mathbf{\Sigma}}_c^{-1} \tilde{\mathbf{Y}}_{1,c}^T / \|\mathbf{X}\|_2$$

$$\hat{\boldsymbol{\beta}}_{-c}^{\text{NC}} = \tilde{\mathbf{Y}}_{1,-c}^T / \|\mathbf{X}\|_2 - \hat{\mathbf{\Gamma}}_{-c} \hat{\boldsymbol{\alpha}}^{\text{NC}}$$

Note: RUV4 [Gagnon-Bartsch et al., 2013] = Ordinary Least Squares (OLS).

Asymptotic distribution of $\hat{\beta}^{\text{NC}}$

Confounder
Adjustment

Qingyuan
Zhao

Introduction
Background
Motivating
Examples
Previous Work

Model and
Inference
Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

Theorem (Wang, Z., Hastie, and Owen [2015])

Under the assumptions of uniform convergence of $\hat{\Sigma}$ and $\hat{\Gamma}$ and $\lim_{p \rightarrow \infty} \frac{1}{|\mathcal{C}|} \Gamma_{\mathcal{C}}^T \Sigma_{\mathcal{C}}^{-1} \Gamma_{\mathcal{C}} \succ \mathbf{0}$, then for any finite index set \mathcal{S} such that $\mathcal{S} \cap \mathcal{C} = \emptyset$:

- ① *If the number of negative controls $|\mathcal{C}| \rightarrow \infty$,*

$$\sqrt{n}(\hat{\beta}_S^{\text{NC}} - \beta_S) \xrightarrow{d} N(\mathbf{0}, (1 + \|\alpha\|_2^2) \Sigma_S)$$

- ② *If $\lim_{p \rightarrow \infty} |\mathcal{C}| < \infty$,*

$$\sqrt{n}(\hat{\beta}_S^{\text{NC}} - \beta_S) \xrightarrow{d} N(\mathbf{0}, (1 + \|\alpha\|_2^2)(\Sigma_S + \Delta_S))$$

$$\text{where } \Delta_S = \lim_{p \rightarrow \infty} \Gamma_S (\Gamma_{\mathcal{C}}^T \Sigma_{\mathcal{C}}^{-1} \Gamma_{\mathcal{C}})^{-1} \Gamma_S^T.$$

Strategy 2: Estimate β via sparsity

Recall

$$\frac{\tilde{\mathbf{Y}}_{p \times 1}^T}{\sqrt{n}} = \mathbf{\Gamma}_{p \times r} \boldsymbol{\alpha}_{r \times 1} + \beta_{p \times 1} + \tilde{\mathbf{E}}_1 / \sqrt{n}$$

response design matrix coefficients

Idea: if $\|\beta\|_0 \ll p$, $\beta_j \neq 0$ is an outlier in this regression.

Robust regression estimator (simplification of LEAPP)

$$\hat{\alpha}^{\text{RR}} = \arg \min \sum_{j=1}^p \rho \left(\frac{\tilde{Y}_{1j} / \sqrt{n} - \hat{\Gamma}_j^T \alpha}{\hat{\sigma}_j} \right)$$
$$\hat{\beta}^{\text{RR}} = \tilde{\mathbf{Y}}_1^T / \sqrt{n} - \hat{\Gamma} \hat{\alpha}^{\text{RR}}$$

Confounder
Adjustment

Qingyuan
Zhao

Introduction
Background
Motivating
Examples
Previous Work

Model and
Inference
Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples
Summary

Asymptotic distribution of $\hat{\beta}^{\text{RR}}$

Confounder
Adjustment

Qingyuan
Zhao

Introduction
Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

Assumptions on the loss function $\rho(x)$

The derivatives ρ' , ρ'' and ρ''' exist and are bounded.
 $\rho(0) = \rho'(0) = 0$, $\rho''(0) > 0$ and $\rho'(x) \cdot x \geq 0$.
(e.g. Tukey's bisquare)

Theorem (Wang, Z., Hastie, and Owen [2015])

Under the assumptions of uniform convergence of $\hat{\Sigma}$ and $\hat{\Gamma}$ and the above assumption of the loss function, if $\min(\|\beta\|_0, \|\beta\|_1)\sqrt{n}/p \rightarrow 0$, then for any finite index set S :

$$\sqrt{n}(\hat{\beta}_S^{\text{RR}} - \beta_S) \xrightarrow{d} N(\mathbf{0}, (1 + \|\alpha\|_2^2)\Sigma_S).$$

Oracle efficiency

Confounder
Adjustment

Qingyuan
Zhao

Introduction
Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

In either the sparsity or negative control scenario ($|\mathcal{C}| \rightarrow \infty$):

$$\sqrt{n}(\hat{\beta}_S - \beta_S) \xrightarrow{d} N(\mathbf{0}, (1 + \|\alpha\|_2^2)\Sigma_S)$$

Oracle estimator

Consider the model

$$\mathbf{Y} = \mathbf{X}\beta^T + \mathbf{Z}\Gamma^T + \mathbf{E}.$$

If \mathbf{Z} were observed, the oracle OLS estimator would be

$$\sqrt{n}(\hat{\beta}_S^{\text{OLS}} - \beta_S) \sim N(\mathbf{0}, (1 + \|\alpha\|_2^2)\Sigma_S).$$

$\hat{\beta}_S$ is as efficient asymptotically as the oracle estimator!

Significance test for confounding

Confounder
Adjustment

Qingyuan
Zhao

Introduction
Background
Motivating
Examples
Previous Work

Model and
Inference
Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

Theorem (Wang, Z., Hastie, and Owen [2015])

Under the above assumptions for oracle efficiency and the null hypothesis that $H_{0,\alpha} : \alpha = \mathbf{0}$, we have

$$n \cdot \hat{\alpha}^T \hat{\alpha} \xrightarrow{d} \chi_r^2$$

where χ_r^2 is the chi-square distribution with r degree of freedom.

Recipes

- 1 Graphical diagnostics: the histogram of test statistics.
- 2 Positive controls: e.g. X/Y genes for gender.
- 3 Asymptotic χ^2 test. If significant, check $\hat{\Gamma}$.

Multiple hypothesis testing

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating
Examples

Previous Work

Model and
Inference

Model and
Identifiability

Estimation
Hypothesis Tests

Numerical
Examples

Summary

Two-sided asymptotic z-tests

Test $H_{j0} : \beta_j = 0$ vs. $H_{j1} : \beta_j \neq 0$ for $j = 1, \dots, p$.

$$t_j = \frac{\sqrt{n}\hat{\beta}_j}{\hat{\sigma}_j\sqrt{1 + \|\hat{\alpha}\|^2}}, \quad P_j = 2(1 - \Phi(|t_j|)).$$

Theorem (Wang, Z., Hastie, and Owen [2015])

Under the assumptions for oracle efficiency, the overall type I error and the familywise error rate (FWER) can be asymptotically controlled.

FDR control: ongoing work.

Simulation: $n = 100$, $p = 5000$ and $r = 10$

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background

Motivating

Examples

Previous Work

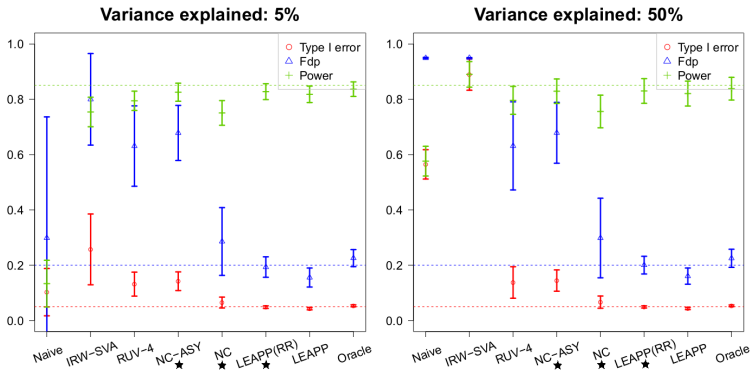
Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

- Sparsity: $\|\beta\|_0/p = 0.05$; NC: $|\mathcal{C}| = 30$.
- $\mathbf{\Gamma}$ uniform from orthogonal matrices; $\sigma_i^2 \stackrel{i.i.d.}{\sim} \text{InvGamma}(3, 2)$.
- Variance of \mathbf{X} explained by \mathbf{Z} : $\max_{\rho} \text{corr}(X_i, \rho^T \mathbf{Z}_i) = \frac{\|\alpha\|_2}{1 + \|\alpha\|_2}$.



COPD data: severity as primary variable

Confounder
Adjustment

Qingyuan
Zhao

Introduction

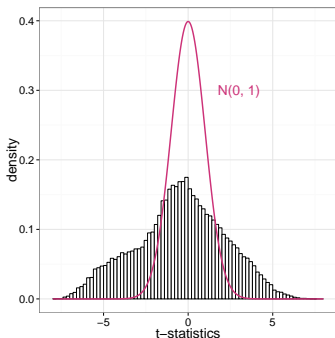
Background
Motivating
Examples
Previous Work

Model and
Inference

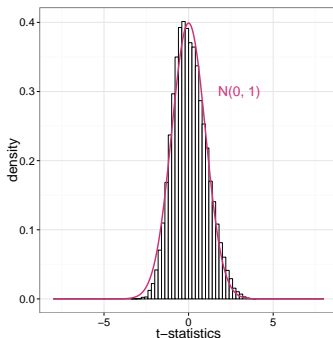
Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary



(a) Naive linear regression.

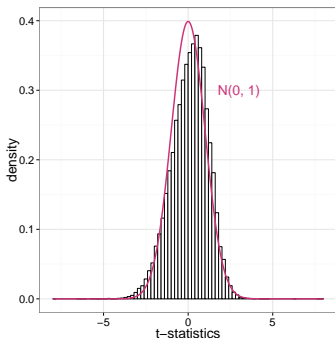


(b) After adjustment.

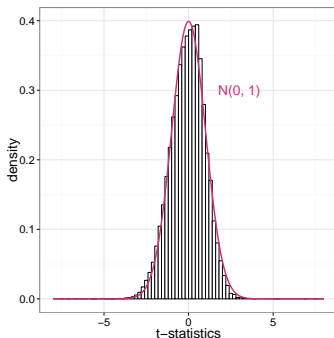
- $\hat{r} = 1$ [Onatski, 2010].
- $\hat{\alpha} \approx 0.98$, variance explained is approximately 22%.
- Test of confounding: p-value ≈ 0 .

COPD data: gender as primary variable

Genes associated with gender should come from X/Y chromosomes (positive controls).



(a) Naive linear regression.



(b) After adjustment.

- $\hat{\alpha} \approx -0.27$, variance explained is approximately 3%.
- Test of confounding: $p\text{-value} \approx 1.2 \times 10^{-3}$.

COPD data: gender as primary variable

Confounder
Adjustment

Qingyuan
Zhao

Introduction

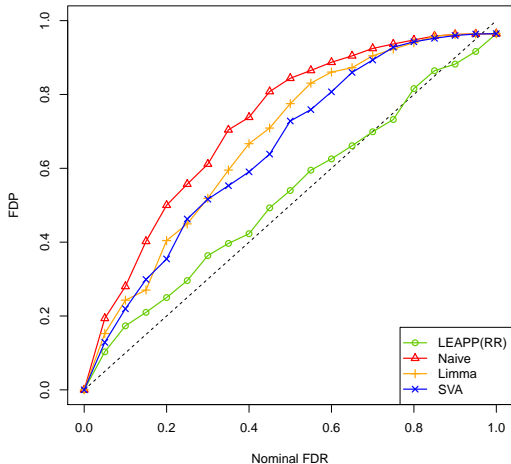
Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary



COPD data: gender as primary variable

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

Method	X/Y Genes in Top 100
LEAPP(RR)	69
Naive	58
Limma	58
SVA	68

Mutual fund selection (preliminary results)

Confounder
Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

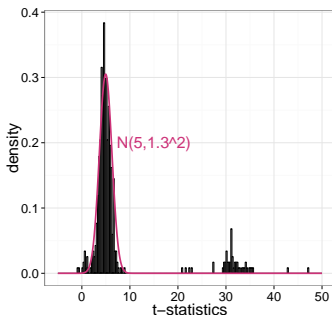
Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

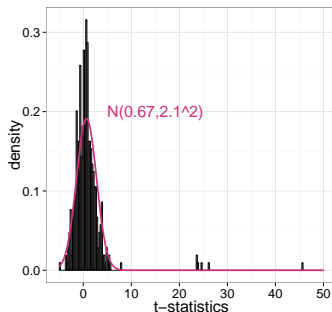
Numerical
Examples

Summary

- $p = 469$ mutual funds with monthly returns available in CRSP database in Jan. 1980 – Dec. 2000 ($n = 240$).
- Apply the RR procedure with $r = 6$ without adjusting for any observed systematic risk factor.



(a) Naive linear regression.



(b) After adjustment.

Summary

Confounder
Adjustment

Qingyuan
Zhao

Introduction
Background
Motivating
Examples
Previous Work

Model and
Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical
Examples

Summary

Recap

- Linear model with unobserved confounding factors.
- Identification conditions: negative control and sparsity.
- Two-step estimation of the primary effects.
- Asymptotic distributions and oracle efficiency.
- Hypothesis tests for confounding and the primary effects.

Open problems

- Correlated noise: approximate factor models.
- Weak factors: random matrix theory.
- Non-Gaussian data: RNA-seq, GWAS.
- Beyond linearity?

Resources

Confounder Adjustment

Qingyuan
Zhao

Introduction

Background
Motivating
Examples
Previous Work

Model and Inference

Model and
Identifiability
Estimation
Hypothesis Tests

Numerical Examples

Summary

- J. Wang, Z., T. Hastie, and A. B. Owen. **Confounder adjustment in multiple hypothesis testing.** *under revision for Annals of Statistics*, 2015.
 - Available on arXiv.
- Software: cate on CRAN.
(<https://cran.r-project.org/web/packages/cate/index.html>)
 - Package vignette available online.
 - Unified interface for existing packages sva, ruv, leapp.
 - We also support formula:

```
results <- cate(~ gender | . - gender - 1, data, ...)
```