

Functional Estimation in High Dimensional Problems

Arin Madenci and Tudor Cristea-Platon

BST235

December 2019

- 1 Introduction
- 2 Scenario #1: $d/n \rightarrow 0$
- 3 Scenario #2: $d/n \rightarrow \rho \leq 1$, Σ known
- 4 Scenario #3: $d/n \rightarrow \rho \leq 1$, Σ unknown (estimable)
- 5 Scenario #4: $d/n \rightarrow \rho > 1$, Σ known
- 6 Scenario #5: $d/n \rightarrow \rho > 1$, Σ unknown (estimable)

Introduction

To understand the behavior of models and performance limits of model-fitting procedures, **residual variance** and **proportion of explained variance** are critical estimands.

Well-known applications, such as:

- Information about scale of an estimator's risk under ℓ^2 loss
- Computation of model selection statistics
- Regression diagnostics (e.g., goodness-of-fit testing)
- Signal-to-noise ratio

Motivation

Estimators of residual variance and proportion of explained variance may or may not be reliable in high-dimensional conditions.

We will evaluate several estimators:

- “Plug-in” estimator (OLS)
- Dicker 2013
- “EigenPrism,” Janson, Barber, and Candes 2017

Motivation

After briefly reviewing notation, assumptions, and definitions, we will compare their performance under a variety of conditions of d (number of parameters) and n (number of observations):

- $\frac{d}{n} \rightarrow 0$
- $\frac{d}{n} \rightarrow \rho \leq 1$, with Σ known
- $\frac{d}{n} \rightarrow \rho \leq 1$, with Σ unknown
- $\frac{d}{n} \rightarrow \rho > 1$, with Σ known
- $\frac{d}{n} \rightarrow \rho > 1$, with Σ unknown

Notation

- Linear model $y_i = x_i^T \beta + \epsilon_i$, $i = 1, \dots, n$
- $y = (y_1, \dots, y_n)^T \in \mathbb{R}^n$ denotes the n -dimensional vector of observed outcomes
- $X = (x_1, \dots, x_n)^T \in \mathbb{R}^{n \times d}$ denotes the $n \times d$ matrix of observed predictors, where $x_1 = (x_{11}, \dots, x_{1d})^T, \dots, x_n = (x_{n1}, \dots, x_{nd})^T \in \mathbb{R}^d$ are d -dimensional predictors
- $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T \in \mathbb{R}^n$ is a vector of unobserved i.i.d. errors with $\mathbb{E}(\epsilon_i) = 0$ and $\text{var}(\epsilon_i) = \sigma^2 > 0$
- $\beta = (\beta_1, \dots, \beta_d)^T \in \mathbb{R}^d$ is an unknown d -dimensional parameter

Assumptions

- Random predictors x_i with mean $\mathbb{E}(x_i) = 0$ and $d \times d$ positive definite covariance matrix $\text{cov}(x_i) = \Sigma$
- $x_i \perp \epsilon$
- $\epsilon_1, \dots, \epsilon_n \sim N(0, \sigma^2)$
- $x_1, \dots, x_n \sim N(0, \Sigma)$
- Note: $\mathbb{E}(x_i) = 0$, without loss of generality; for cases in which $\mathbb{E}(x_i) \neq 0$, include intercept term, center data, and replace n with $n - 1$.

Definitions: Signal and Noise

- τ^2 : ℓ^2 -signal strength
 - ▶ Let $\tau^2 = \beta^T \Sigma \beta = \|\Sigma^{1/2} \beta\|^2$
 - ▶ Where $\|\cdot\|$ is the ℓ^2 -norm
- σ^2 : residual variance (noise)
 - ▶ $\sigma^2 = \text{Var}(\epsilon_i) = \text{Var}(\mathbb{E}(y_i|x_i))$

Returning to our question of interest:

How can we identify effective estimators of these quantities (i.e., τ^2 and σ^2) in high-dimensional linear models with large d and n ... and especially for $d > n$?

Scenario #1: $d/n \rightarrow 0$

For estimation of σ^2 :

- $\hat{\sigma}_{MLE}^2 = \frac{1}{n} \|P_{C(\mathbb{X})^\perp}(\mathbf{y})\|^2$ is the MLE of σ^2 .
- More specifically, $\hat{\sigma}_{MLE}^2 = \frac{\boldsymbol{\epsilon}^T P_{C(\mathbb{X})^\perp} \boldsymbol{\epsilon}}{n} = \frac{\sigma^2}{n} \left(\frac{\boldsymbol{\epsilon}}{\sigma}\right)^T P_{C(\mathbb{X})^\perp} \left(\frac{\boldsymbol{\epsilon}}{\sigma}\right)$.
- We note that $\left. \frac{\sigma^2}{n} \left(\frac{\boldsymbol{\epsilon}}{\sigma}\right)^T P_{C(\mathbb{X})^\perp} \left(\frac{\boldsymbol{\epsilon}}{\sigma}\right) \right| \mathbb{X} \sim \chi_{n-\text{rank}(\mathbb{X})}^2(0)$.
- From previous work¹, we know that for $d \leq n$, $P(\text{rank}(\mathbb{X}) = d) = 1$ as \mathbf{x}_i is continuous.

¹Eaton ML and Perlman MD. The Non-Singularity of Generalized Sample Covariance Matrices. *The Annals of Statistics* 1973;1:710–7.

Scenario #1: $d/n \rightarrow 0$

Given that $\frac{\sigma^2}{n} \left(\frac{\epsilon}{\sigma} \right)^T P_{C(\mathbb{X})^\perp} \left(\frac{\epsilon}{\sigma} \right) \Big| \mathbb{X} \sim \chi_{n-d}^2(0)$

This implies:

$$\mathbb{E}(\hat{\sigma}_{MLE}^2) = \sigma^2 \frac{n-d}{n} \rightarrow \sigma^2 \quad \text{as } n \rightarrow \infty \quad \text{s.t. } n-d \rightarrow \infty$$

$$\mathbb{V}ar(\hat{\sigma}_{MLE}^2) = \frac{\sigma^4}{n^2} 2(n-d) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

Result: $\hat{\sigma}_{MLE}^2$ and $\hat{\tau}_{MLE}^2$ (not shown) unbiased for $d/n \rightarrow 0$.

Scenario #1: $d/n \rightarrow 0$

For MLE estimation of τ^2 :

$$\begin{aligned}\mathbb{E}[\hat{\tau}_{MLE}^2] &= \mathbb{E}[\|\hat{\beta}_{MLE}\|_2^2] \\ &= \mathbb{E}[\|(X^T X)^{-1} X^T Y\|_2^2] \\ &= \mathbb{E}[\|(X^T X)^{-1} X^T (X \hat{\beta}_{MLE} + \epsilon)\|_2^2] \\ &= \mathbb{E}[\|\hat{\beta}_{MLE} + (X^T X)^{-1} X^T \epsilon\|_2^2] \\ &= \mathbb{E}[\|\hat{\beta}_{MLE}\|_2^2] + \mathbb{E}[\|(X^T X)^{-1} X^T \epsilon\|_2^2]\end{aligned}$$

Scenario #1: $d/n \rightarrow 0$

Bias of $\hat{\tau}_{MLE}^2$:

$$\begin{aligned}\mathbb{E}|| (X^T X)^{-1} X^T \epsilon ||_2^2 &= \mathbb{E}[\epsilon^T X (X^T X)^{-2} X^T \epsilon] \\ &= \mathbb{E}[tr(X (X^T X)^{-2} X^T \sigma^2 I)] \\ &= \sigma^2 \mathbb{E}[tr((X^T X)^{-1})] \\ &= \frac{\sigma^2}{n} \mathbb{E}[tr\left(\frac{X^T X}{n}\right)^{-1}] \\ &= \frac{\sigma^2}{n} tr\left(\frac{\Sigma^{-1}}{n - d - 1}\right)\end{aligned}$$

For $\Sigma = I$:

$$= \sigma^2 \frac{d}{n(n - d - 1)}$$

Numerical simulation: $\sigma^2 = \tau^2 = 1$

- $x_1, \dots, x_n \in \mathbb{R}^d \sim N(0, I)^\dagger$
- $\beta^* \in \mathbb{R}^d$
 - ▶ $\beta_1^*, \dots, \beta_{\lfloor d/2 \rfloor}^* \sim \text{unif}(0, 1)$
 - ▶ $\beta_{\lfloor d/2 \rfloor + 1}^*, \dots, \beta_d^* \sim N(0, 1)$
- $\beta = \beta^* (\beta^{*T} \beta^*)^{-1/2}$ s.t. $\tau^2 = 1$
- $Y = X\beta + \epsilon$, $\epsilon \sim N(0, 1)$ s.t. $\sigma^2 = 1$
- Monte Carlo simulation with 300 replications

[†] Note: assume $\Sigma = I$ without loss of generality (for $\Sigma \neq I$, replace (X, β) with $(X\Sigma^{-1/2}, \Sigma^{1/2}\beta)$).

Numerical simulation for σ^2 : $d/n \rightarrow 0$

Plug-in estimator (OLS) performs well, as expected:

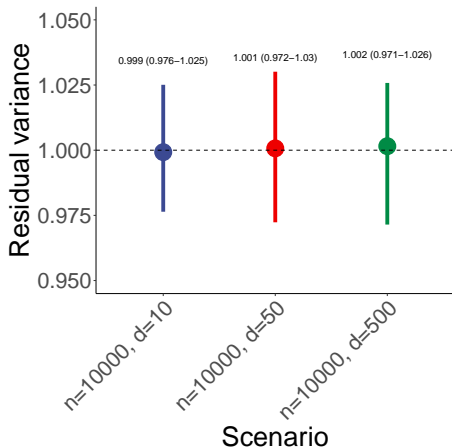


Figure: Plug-in estimator for σ^2 in scenarios with $d/n \rightarrow 0$

Numerical simulation for τ^2 : $d/n \rightarrow 0$

Again, plug-in estimator (OLS) performs well

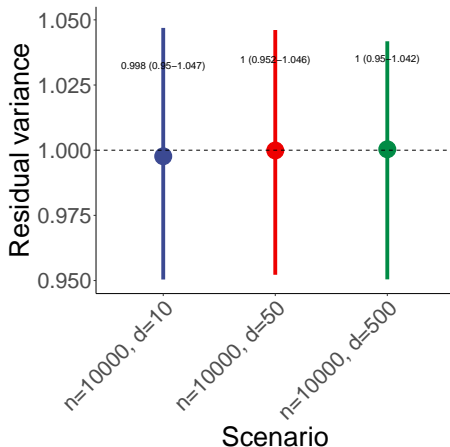


Figure: Plug-in estimator for τ^2 in scenarios with $d/n \rightarrow 0$

Numerical simulation: $d/n \rightarrow 0$

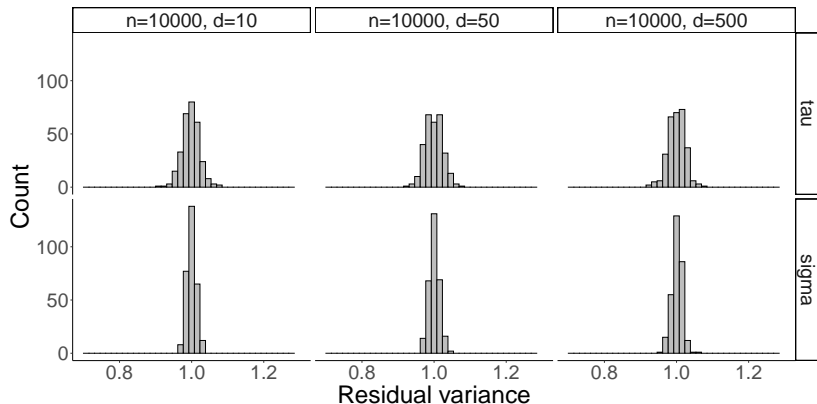


Figure: Histograms of plug-in estimator for σ^2 and τ^2 in scenarios with $d/n \rightarrow 0$

- 1 Introduction
- 2 Scenario #1: $d/n \rightarrow 0$
- 3 Scenario #2: $d/n \rightarrow \rho \leq 1$, Σ known
- 4 Scenario #3: $d/n \rightarrow \rho \leq 1$, Σ unknown (estimable)
- 5 Scenario #4: $d/n \rightarrow \rho > 1$, Σ known
- 6 Scenario #5: $d/n \rightarrow \rho > 1$, Σ unknown (estimable)

Scenario #2: $d/n \rightarrow \rho \leq 1$

Plug-in estimator:

$$\mathbb{E}[\hat{\sigma}_{MLE}^2] = \sigma^2 \frac{d}{n(n-d-1)}$$

Dicker (2014) proposes the following method of moments estimators:

- $\hat{\sigma}^2 = \frac{d+n+1}{n(n+1)} \|y\|_2^2 - \frac{1}{n(n+1)} \|X^T y\|_2^2$
- $\hat{\tau}^2 = -\frac{d}{n(n+1)} \|y\|_2^2 + \frac{1}{n(n+1)} \|X^T y\|_2^2$

Scenario #2: $d/n \rightarrow \rho \leq 1$ (known Σ)

Sketch of proof:²

$$\begin{aligned}\|y\|_2^2 &\sim (\sigma^2 + \tau^2)\chi_n^2 \implies \mathbb{E}\left(\frac{1}{n}\|y\|_2^2\right) = \tau^2 + \sigma^2 \\ \mathbb{E}\left(\frac{1}{n^2}\|X^T y\|_2^2\right) &= \frac{d+n+1}{n}\tau^2 + \frac{d}{n}\sigma^2\end{aligned}$$

- Above identities can be written as linear combinations of τ^2 and σ^2 .
- Unbiased estimators of τ^2 and σ^2 may be found by taking linear combinations of $n^{-1}\|y\|_2^2$ and $n^{-2}\|X^T y\|_2^2$.
- Unbiased estimators in the case that $\Sigma = I$.

²Full reference available in Dicker 2014.

Numerical simulation for σ^2 : $d/n \rightarrow \rho \leq 1$ (known Σ)

Is it so bad to use our plug-in estimator for scenarios in which $d/n \rightarrow \rho \leq 1$ (as opposed to $d/n \rightarrow 0$)?

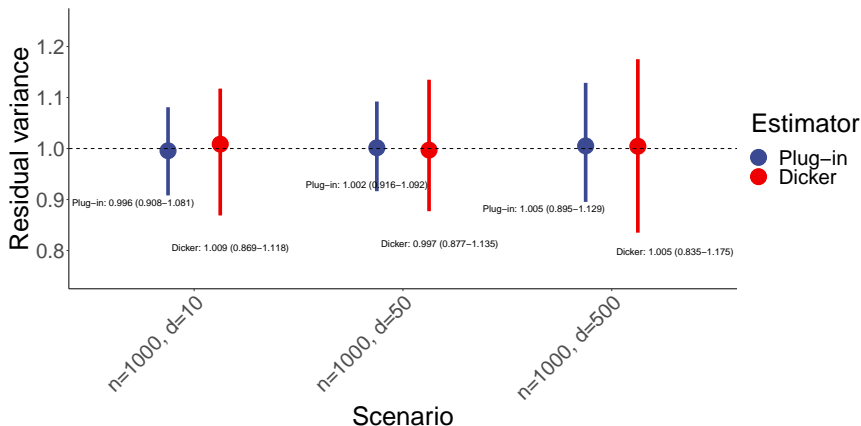


Figure: Plug-in and Dicker estimators for σ^2 in scenarios with $d/n \rightarrow \rho \leq 1$

Numerical simulation for τ^2 : $d/n \rightarrow \rho \leq 1$ (known Σ)

Is it so bad to use our plug-in estimator for scenarios in which $d/n \rightarrow \rho \leq 1$ (as opposed to $d/n \rightarrow 0$)?

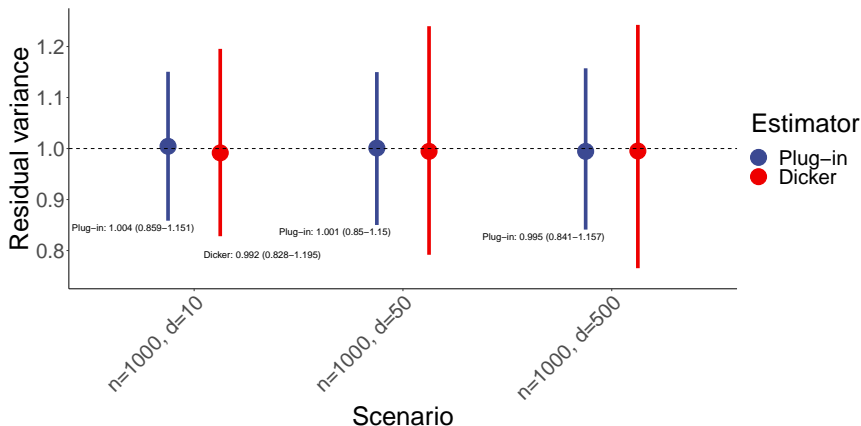


Figure: Plug-in and Dicker estimators for τ^2 in scenarios with $d/n \rightarrow \rho \leq 1$

Numerical simulation: $d/n \rightarrow \rho \leq 1$ (known Σ)

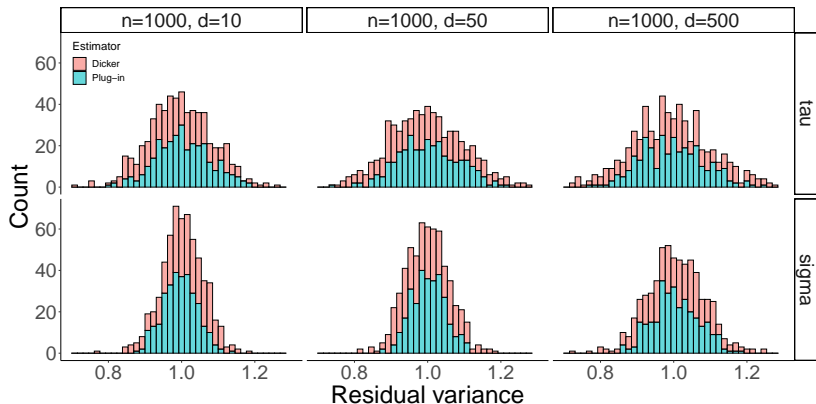


Figure: Histograms of plug-in and Dicker estimators for σ^2 and τ^2 in scenarios with $d/n \rightarrow \rho \leq 1$ (known Σ)

- 1 Introduction
- 2 Scenario #1: $d/n \rightarrow 0$
- 3 Scenario #2: $d/n \rightarrow \rho \leq 1$, Σ known
- 4 Scenario #3: $d/n \rightarrow \rho \leq 1$, Σ unknown (estimable)
- 5 Scenario #4: $d/n \rightarrow \rho > 1$, Σ known
- 6 Scenario #5: $d/n \rightarrow \rho > 1$, Σ unknown (estimable)

Scenario #3: $d/n \rightarrow \rho \leq 1$, Σ unknown (estimable)

With unknown Σ :

- Plug-in estimator can no longer be used
- Dicker estimator must be modified

Take the unknown- Σ analogs of known- Σ Dicker estimators $\hat{\sigma}^2$ and $\hat{\tau}^2$ above, as follows (using some positive definite estimator, $\hat{\Sigma}$):

$$\hat{\sigma}^2(\hat{\Sigma}) = \frac{d+n+1}{n(n+1)} \|y\|_2^2 - \frac{1}{n(n+1)} \|\hat{\Sigma}^{-1/2} X^T y\|_2^2$$
$$\hat{\tau}^2(\hat{\Sigma}) = -\frac{d}{n(n+1)} \|y\|_2^2 + \frac{1}{n(n+1)} \|\hat{\Sigma}^{-1/2} X^T y\|_2^2$$

Scenario #3: $d/n \rightarrow \rho \leq 1$, Σ unknown

We will consider two options for estimating Σ :

① Empiric Σ : $\hat{\Sigma} = \frac{1}{n} X^T X$

② Banded Σ : $\tilde{\Sigma}_{k,p} = B_k(\hat{\Sigma}_p)$

$$\begin{bmatrix} a_{11} & a_{12} & 0 & \cdots & \cdots & \cdots & \cdots & 0 \\ a_{21} & a_{22} & a_{23} & \ddots & & & & \vdots \\ 0 & a_{32} & a_{33} & a_{34} & \ddots & & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & & a_{n-2n-3} & a_{n-2n-2} & a_{n-2n-1} & 0 \\ \vdots & & & & \ddots & a_{n-1n-2} & a_{n-1n-1} & a_{n-1n} \\ 0 & \cdots & \cdots & \cdots & \cdots & 0 & a_{nn-1} & a_{nn} \end{bmatrix}$$

Where $B_k(M) = [m_{ij} \mathbb{1}(|i-j| \leq k)]$

Numerical simulation for σ^2 : $d/n \rightarrow \rho \leq 1$, Σ unknown

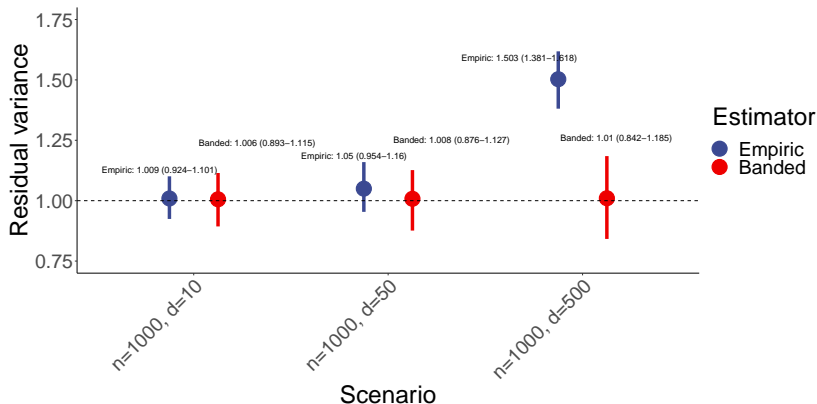


Figure: Numerical simulation for σ^2

Note: # bands, $k = 5$ for simulations.

Numerical simulation for τ^2 : $d/n \rightarrow \rho \leq 1$, Σ unknown

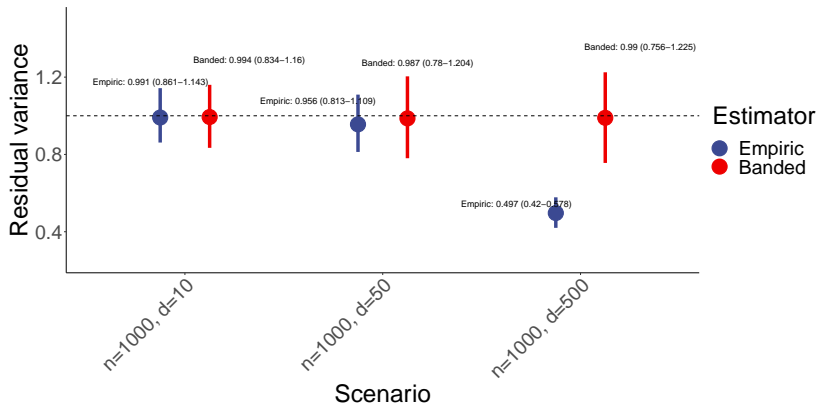


Figure: Numerical simulation for τ^2

Numerical simulation for τ^2 : $d/n \rightarrow \rho \leq 1$, Σ unknown

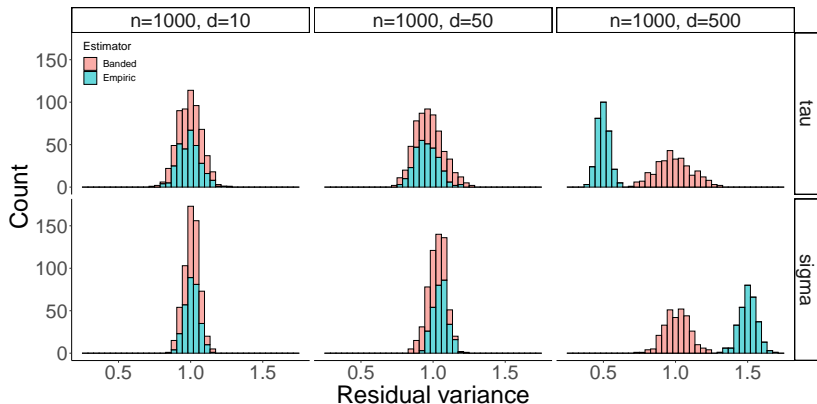


Figure: Histograms of Dicker estimator for σ^2 and τ^2 using empiric and banded estimators of Σ , in scenarios with $d/n \rightarrow \rho \leq 1$ (unknown Σ)

- 1 Introduction
- 2 Scenario #1: $d/n \rightarrow 0$
- 3 Scenario #2: $d/n \rightarrow \rho \leq 1$, Σ known
- 4 Scenario #3: $d/n \rightarrow \rho \leq 1$, Σ unknown (estimable)
- 5 Scenario #4: $d/n \rightarrow \rho > 1$, Σ known
- 6 Scenario #5: $d/n \rightarrow \rho > 1$, Σ unknown (estimable)

$$d/n \rightarrow \rho > 1, \Sigma \text{ known}$$

Introduce another possible estimator for σ^2 and τ^2 : EigenPrism (Janson, Barber, and Candès):

- Goal to develop estimators unbiased for σ^2 and τ^2 , which are asymptotically normally distributed, and with estimable tight bound on variance.
- No need for knowledge of the noise-level (σ^2) or any assumption on the structure of the coefficient vector β (e.g. sparsity).

$d/n \rightarrow \rho > 1$, Σ known

Sketch of proof for EigenPrism estimators:

$$\begin{aligned}\mathbb{E}\left(\sum_{i=1}^n w_i z_i^2 \mid d\right) &= \sum_{i=1}^n w_i \left(\lambda_i \tau^2 + \sigma^2\right) \\ &= \tau^2 \sum_{i=1}^n w_i \lambda_i + \sigma^2 \sum_{i=1}^n w_i\end{aligned}$$

Unbiased estimator for τ^2 , when constraining $\sum_{i=1}^n w_i = 0$ and $\sum_{i=1}^n w_i \lambda_i = 1$

$d/n \rightarrow \rho > 1$, Σ known

$$\begin{aligned}\mathbb{V}ar\left(\sum_{i=1}^n w_i z_i^2 | d\right) &= \tau^2 \sum_{i=1}^n w_i \lambda_i + \sigma^2 \sum_{i=1}^n w_i \\ &\leq 2(\theta^2 + \sigma^2)^2 \cdot \max\left(\sum_{i=1}^n w_i^2, \sum_{i=1}^n (w_i \lambda_i)^2\right)\end{aligned}$$

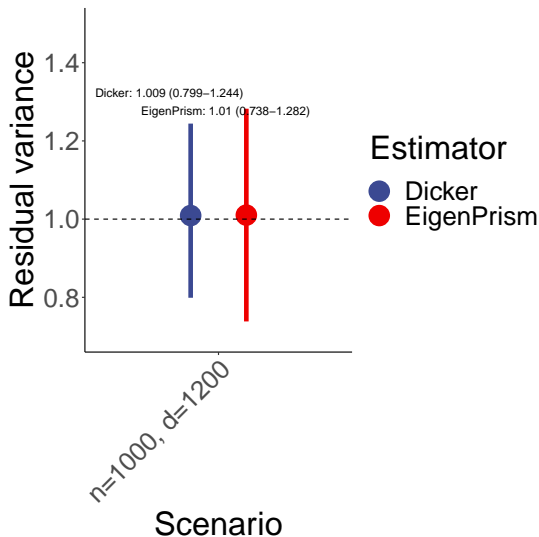
$\mathcal{P}_1 = \operatorname{argmin}_{w \in \mathbb{R}^n} \max\left(\sum_{i=1}^n w_i^2, \sum_{i=1}^n (w_i \lambda_i)^2\right),$
s.t. $\sum_{i=1}^n w_i = 0, \quad \sum_{i=1}^n w_i \lambda_i = 1$, and with solution w^*

EigenPrism:

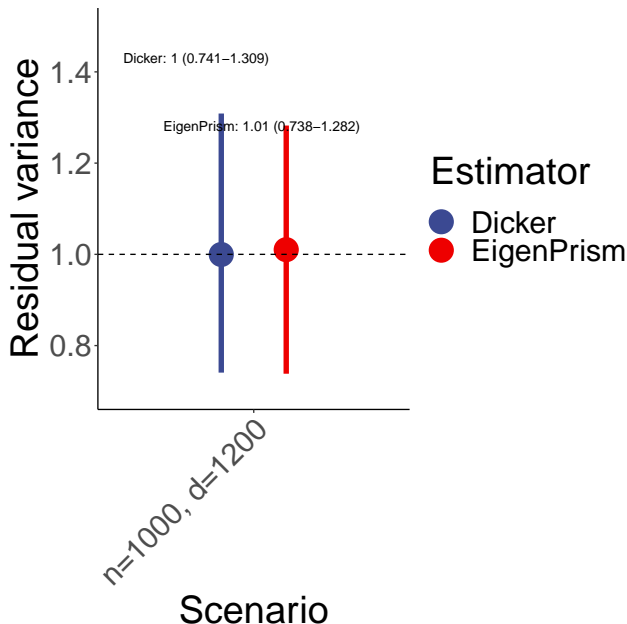
- $\mathbb{E}(\sum_{i=1}^n w_i^* z_i^2 | d) = \hat{\tau}^2$
- $\text{SD}(\sum_{i=1}^n w_i^* z_i^2 | d) \lesssim \sqrt{2\text{val}(\mathcal{P}_1)} \frac{\|y\|_2^2}{n}$

Numerical simulation for σ^2 : $d/n \rightarrow \rho > 1$, Σ known

With $d > n$, cannot use plug-in estimator ($X^T X$ non-invertible). Dicker estimator performs well.



Numerical simulation for τ^2 : $d/n \rightarrow \rho > 1$, Σ known



Numerical simulation for τ^2 : $d/n \rightarrow \rho > 1$, Σ known

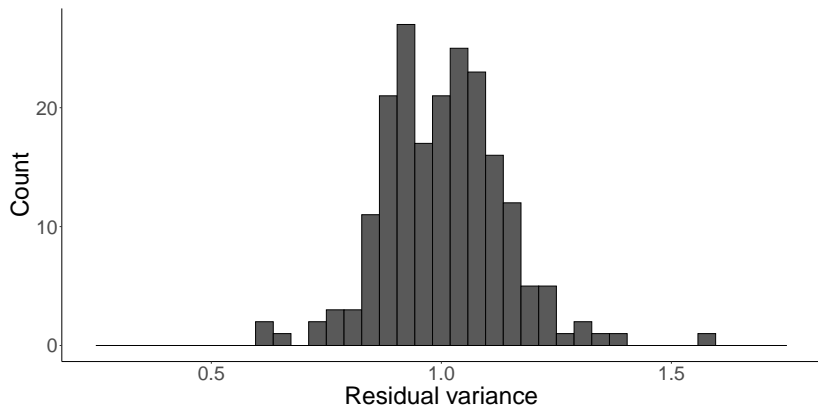


Figure: Histograms of Dicker estimator for σ^2 and τ^2 using empiric and banded estimators of Σ , $d = 1200$, $n = 1000$ (known Σ)

- 1 Introduction
- 2 Scenario #1: $d/n \rightarrow 0$
- 3 Scenario #2: $d/n \rightarrow \rho \leq 1$, Σ known
- 4 Scenario #3: $d/n \rightarrow \rho \leq 1$, Σ unknown (estimable)
- 5 Scenario #4: $d/n \rightarrow \rho > 1$, Σ known
- 6 Scenario #5: $d/n \rightarrow \rho > 1$, Σ unknown (estimable)

Numerical simulation for σ^2 : $d/n \rightarrow \rho > 1$, Σ unknown

- Empiric covariance matrix, $n^{-1}X^T X$ no longer norm-consistent for Σ .
- Generally not possible to find norm-consistent estimator for Σ , without further information/assumptions.

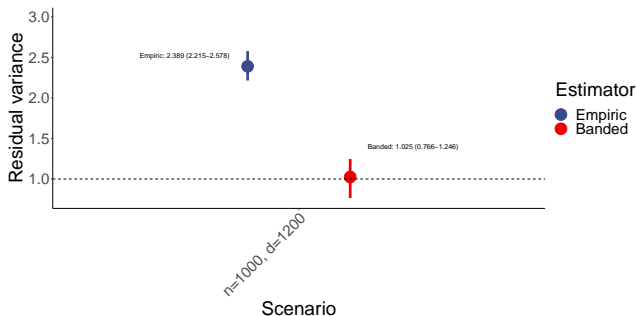


Figure: Numerical simulation for σ^2

Note: Empiric estimator now requires Moore-Penrose generalized inverse of Σ .

Numerical simulation for τ^2 : $d/n \rightarrow \rho > 1$, Σ unknown

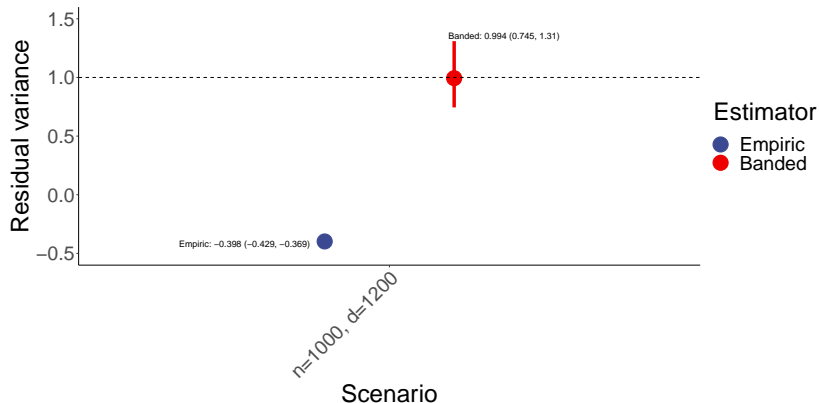







Figure: Numerical simulation for τ^2

References

-  Eaton ML and Perlman MD. The Non-Singularity of Generalized Sample Covariance Matrices. [The Annals of Statistics](#) 1973;1:710–7.
-  Dicker LH. Optimal equivariant prediction for high-dimensional linear models with arbitrary predictor covariance. [Electronic Journal of Statistics](#) 2013;7:1806–34.
-  Dicker LH. Variance estimation in high-dimensional linear models. [Biometrika](#) 2014;101:269–84.
-  Janson L, Barber RF, and Candès E. EigenPrism: inference for high dimensional signal-to-noise ratios. [Journal of the Royal Statistical Society: Series B \(Statistical Methodology\)](#) 2017;79:1037–65.
-  Bickel PJ and Levina E. Regularized estimation of large covariance matrices. [The Annals of Statistics](#) 2008;36:199–227.