

On the Distribution of the Adaptive LASSO Estimator (pt I)

Ulrike Schneider
(joint with B.M. Pötscher)

University of Vienna

Stochastics Colloquium, University of Göttingen
January 14, 2009

Outline

- 1 Introduction
- 2 Theoretical results for the adaptive LASSO
- 3 Simulation results
- 4 Impossibility result for the estimation of the cdf.
- 5 Conclusions

Penalized LS (ML) Estimators

Linear regression model

$$\mathbf{y} = \theta_1 \mathbf{x}_{\cdot 1} + \dots + \theta_k \mathbf{x}_{\cdot k} + \varepsilon$$

- response $\mathbf{y} \in \mathbb{R}^n$ (known)
- regressors $\mathbf{x}_{\cdot i} \in \mathbb{R}^n$, $1 \leq i \leq k$ (known)
- errors $\varepsilon \in \mathbb{R}^n$ (unknown)
- parameter vector $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)' \in \mathbb{R}^k$ (unknown)

A penalized least-squares (LS) estimator $\hat{\boldsymbol{\theta}}$ for $\boldsymbol{\theta}$ is given by

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^k} \underbrace{\|\mathbf{y} - X\boldsymbol{\theta}\|^2}_{\text{likelihood or LS -part}} + \underbrace{\lambda_n \rho(\boldsymbol{\theta})}_{\text{penalty}}$$

$\lambda_n > 0$ is a tuning parameter ($\lambda_n = 0$ corresponds to unpenalized/ordinary LS), $X = [\mathbf{x}_{\cdot 1}, \dots, \mathbf{x}_{\cdot k}]$ the $n \times k$ regression matrix.

Penalized LS (ML) Estimators

Linear regression model

$$\mathbf{y} = \theta_1 \mathbf{x}_{.1} + \dots + \theta_k \mathbf{x}_{.k} + \varepsilon$$

- response $\mathbf{y} \in \mathbb{R}^n$ (known)
- regressors $\mathbf{x}_{.i} \in \mathbb{R}^n$, $1 \leq i \leq k$ (known)
- errors $\varepsilon \in \mathbb{R}^n$ (unknown)
- parameter vector $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)' \in \mathbb{R}^k$ (unknown)

A penalized least-squares (LS) estimator $\hat{\boldsymbol{\theta}}$ for $\boldsymbol{\theta}$ is given by

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^k} \underbrace{\|\mathbf{y} - X\boldsymbol{\theta}\|^2}_{\text{likelihood or LS -part}} + \underbrace{\lambda_n p(\boldsymbol{\theta})}_{\text{penalty}}$$

$\lambda_n > 0$ is a tuning parameter ($\lambda_n = 0$ corresponds to unpenalized/ordinary LS), $X = [\mathbf{x}_{.1}, \dots, \mathbf{x}_{.k}]$ the $n \times k$ regression matrix.

Clearly, different penalties give rise to different estimators.

- General class of Bridge-estimators (Frank & Friedman, 1993) using l_γ - type penalties

$$\lambda_n \rho(\boldsymbol{\theta}) = \lambda_n \sum_{i=1}^k |\theta_i|^\gamma$$

$\gamma = 2$: Ridge-estimator (Hoerl & Kennard, 1970)

$\gamma = 1$: LASSO (Tibshirani, 1996).

- Hard- and soft-thresholding estimators.
- Smoothly clipped absolute deviation (SCAD) estimator (Fan & Li, 2001).
- Adaptive LASSO estimator (Zou, 2006).

Clearly, different penalties give rise to different estimators.

- General class of Bridge-estimators (Frank & Friedman, 1993) using l_γ - type penalties

$$\lambda_n p(\boldsymbol{\theta}) = \lambda_n \sum_{i=1}^k |\theta_i|^\gamma$$

$\gamma = 2$: Ridge-estimator (Hoerl & Kennard, 1970)

$\gamma = 1$: LASSO (Tibshirani, 1996).

- Hard- and soft-thresholding estimators.
- Smoothly clipped absolute deviation (SCAD) estimator (Fan & Li, 2001).
- Adaptive LASSO estimator (Zou, 2006).

Clearly, different penalties give rise to different estimators.

- General class of Bridge-estimators (Frank & Friedman, 1993) using l_γ - type penalties

$$\lambda_n p(\boldsymbol{\theta}) = \lambda_n \sum_{i=1}^k |\theta_i|^\gamma$$

$\gamma = 2$: Ridge-estimator (Hoerl & Kennard, 1970)

$\gamma = 1$: LASSO (Tibshirani, 1996).

- Hard- and soft-thresholding estimators.
- Smoothly clipped absolute deviation (SCAD) estimator (Fan & Li, 2001).
- Adaptive LASSO estimator (Zou, 2006).

Clearly, different penalties give rise to different estimators.

- General class of Bridge-estimators (Frank & Friedman, 1993) using l_γ - type penalties

$$\lambda_n p(\boldsymbol{\theta}) = \lambda_n \sum_{i=1}^k |\theta_i|^\gamma$$

$\gamma = 2$: Ridge-estimator (Hoerl & Kennard, 1970)

$\gamma = 1$: LASSO (Tibshirani, 1996).

- Hard- and soft-thresholding estimators.
- Smoothly clipped absolute deviation (SCAD) estimator (Fan & Li, 2001).
- Adaptive LASSO estimator (Zou, 2006).

Clearly, different penalties give rise to different estimators.

- General class of Bridge-estimators (Frank & Friedman, 1993) using l_γ - type penalties

$$\lambda_n p(\boldsymbol{\theta}) = \lambda_n \sum_{i=1}^k |\theta_i|^\gamma$$

$\gamma = 2$: Ridge-estimator (Hoerl & Kennard, 1970)

$\gamma = 1$: LASSO (Tibshirani, 1996).

- Hard- and soft-thresholding estimators.
- Smoothly clipped absolute deviation (SCAD) estimator (Fan & Li, 2001).
- Adaptive LASSO estimator (Zou, 2006).

Relationship to classical PMS-estimators

Bridge-estimators satisfy

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^k} \|y - X\boldsymbol{\theta}\|^2 + \lambda_n \sum_{i=1}^k |\theta_i|^\gamma \quad (0 < \gamma < \infty)$$

For $\gamma \rightarrow 0$, get

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^k} \|y - X\boldsymbol{\theta}\|^2 + \lambda_n \text{card}\{i : \theta_i \neq 0\}$$

which yields a minimum C_p -type procedure such as AIC and BIC.
(l_γ -type penalty with “ $\gamma = 0$ ”)

- For “ $\gamma = 0$ ” procedures are computationally expensive.
- For $\gamma > 0$ (Bridge) estimators are more computationally tractable, especially for $\gamma \geq 1$ (convex objective function).
- For $\gamma \leq 1$, estimators perform model selection

$$P_{n,\theta}(\hat{\theta}_i = 0) > 0 \quad \text{if } \theta_i = 0.$$

Same for SCAD, hard- and soft-thresholding. Phenomenon is more pronounced for smaller γ .

- $\gamma = 1$ (LASSO and adaptive LASSO) as compromise between the wish to detect zeros and computational simplicity. (SCAD leads to a non-convex optimization problem.)

- For “ $\gamma = 0$ ” procedures are computationally expensive.
- For $\gamma > 0$ (Bridge) estimators are more computationally tractable, especially for $\gamma \geq 1$ (convex objective function).
- For $\gamma \leq 1$, estimators perform model selection

$$P_{n,\theta}(\hat{\theta}_i = 0) > 0 \quad \text{if } \theta_i = 0.$$

Same for SCAD, hard- and soft-thresholding. Phenomenon is more pronounced for smaller γ .

- $\gamma = 1$ (LASSO and adaptive LASSO) as compromise between the wish to detect zeros and computational simplicity. (SCAD leads to a non-convex optimization problem.)

- For “ $\gamma = 0$ ” procedures are computationally expensive.
- For $\gamma > 0$ (Bridge) estimators are more computationally tractable, especially for $\gamma \geq 1$ (convex objective function).
- For $\gamma \leq 1$, estimators perform model selection

$$P_{n,\theta}(\hat{\theta}_i = 0) > 0 \quad \text{if } \theta_i = 0.$$

Same for SCAD, hard- and soft-thresholding. Phenomenon is more pronounced for smaller γ .

- $\gamma = 1$ (LASSO and adaptive LASSO) as compromise between the wish to detect zeros and computational simplicity. (SCAD leads to a non-convex optimization problem.)

- For “ $\gamma = 0$ ” procedures are computationally expensive.
- For $\gamma > 0$ (Bridge) estimators are more computationally tractable, especially for $\gamma \geq 1$ (convex objective function).
- For $\gamma \leq 1$, estimators perform model selection

$$P_{n,\theta}(\hat{\theta}_i = 0) > 0 \quad \text{if } \theta_i = 0.$$

Same for SCAD, hard- and soft-thresholding. Phenomenon is more pronounced for smaller γ .

- $\gamma = 1$ (LASSO and adaptive LASSO) as compromise between the wish to detect zeros and computational simplicity. (SCAD leads to a non-convex optimization problem.)

- For “ $\gamma = 0$ ” procedures are computationally expensive.
- For $\gamma > 0$ (Bridge) estimators are more computationally tractable, especially for $\gamma \geq 1$ (convex objective function).
- For $\gamma \leq 1$, estimators perform model selection

$$P_{n,\theta}(\hat{\theta}_i = 0) > 0 \quad \text{if } \theta_i = 0.$$

Same for SCAD, hard- and soft-thresholding. Phenomenon is more pronounced for smaller γ .

- $\gamma = 1$ (LASSO and adaptive LASSO) as compromise between the wish to detect zeros and computational simplicity. (SCAD leads to a non-convex optimization problem.)

The PLS estimator(s) we treat in the following can be viewed to simultaneously perform model selection and parameter estimation.

Some terminology (model selection)

- **Consistent model selection** – Zero coefficients are found with asymptotic probability equal to 1.

$$\lim_{n \rightarrow \infty} P_{n,\theta}(\hat{\theta}_i = 0) = 1 \quad \text{whenever } \theta_i = 0 \quad (1 \leq i \leq k)$$

$$\lim_{n \rightarrow \infty} P_{n,\theta}(\hat{\theta}_i = 0) = 0 \quad \text{whenever } \theta_i \neq 0 \quad (1 \leq i \leq k)$$

An estimator performing consistent model selection is said to have the **sparsity property**.

- **Conservative model selection** – Zero coefficients are found with asymptotic probability less than 1.

$$\lim_{n \rightarrow \infty} P_{n,\theta}(\hat{\theta}_i = 0) < 1 \quad \text{whenever } \theta_i = 0 \quad (1 \leq i \leq k)$$

$$\lim_{n \rightarrow \infty} P_{n,\theta}(\hat{\theta}_i = 0) = 0 \quad \text{whenever } \theta_i \neq 0 \quad (1 \leq i \leq k)$$

Some terminology (model selection)

- **Consistent model selection** – Zero coefficients are found with asymptotic probability equal to 1.

$$\lim_{n \rightarrow \infty} P_{n,\theta}(\hat{\theta}_i = 0) = 1 \quad \text{whenever } \theta_i = 0 \quad (1 \leq i \leq k)$$

$$\lim_{n \rightarrow \infty} P_{n,\theta}(\hat{\theta}_i = 0) = 0 \quad \text{whenever } \theta_i \neq 0 \quad (1 \leq i \leq k)$$

An estimator performing consistent model selection is said to have the **sparsity property**.

- **Conservative model selection** – Zero coefficients are found with asymptotic probability less than 1.

$$\lim_{n \rightarrow \infty} P_{n,\theta}(\hat{\theta}_i = 0) < 1 \quad \text{whenever } \theta_i = 0 \quad (1 \leq i \leq k)$$

$$\lim_{n \rightarrow \infty} P_{n,\theta}(\hat{\theta}_i = 0) = 0 \quad \text{whenever } \theta_i \neq 0 \quad (1 \leq i \leq k)$$

- Consistent vs. conservative model selection can in our context be driven by the asymptotic behavior of the tuning parameters λ_n . Also called “**sparse**ly” vs. “**non-sparse**ly” tuned procedures.
- **Oracle property** – Asymptotic distribution coincides with the one of the **infeasible unpenalized estimator** using the true zero restrictions (with VC-matrix Σ_θ).

$$n^{1/2}(\hat{\theta} - \theta) \rightarrow N(\mathbf{0}, \Sigma_\theta)$$

Seems to suggest that $\hat{\theta}$ performs as well as if we would know the true zero coefficients of θ .

- Consistent vs. conservative model selection can in our context be driven by the asymptotic behavior of the tuning parameters λ_n . Also called “**sparse**ly” vs. “**non-sparse**ly” tuned procedures.
- **Oracle property** – Asymptotic distribution coincides with the one of the **infeasible unpenalized estimator** using the true zero restrictions (with VC-matrix Σ_θ).

$$n^{1/2}(\hat{\theta} - \theta) \rightarrow N(\mathbf{0}, \Sigma_\theta)$$

Seems to suggest that $\hat{\theta}$ performs as well as if we would know the true zero coefficients of θ .

Literature on distributional properties of PLSEs

- Knight & Fu, 2000. Moving-parameter asymptotics for non-sparsely tuned LASSO and Bridge estimators in general.
- Fan & Li, 2001. Fixed-parameter asymptotics for SCAD.
- Zou, 2006. Fixed-parameter asymptotics for sparsely-tuned LASSO and adaptive LASSO.
- Additional papers establishing the **oracle property** for **sparsely-tuned** PLSEs and related estimators within a fixed-parameter framework.

Fan & Li (2002, 2004), Bunea (2004), Bunea & McKeague (2005), Wang & Leng (2007), Li & Liang (2007), Wang, G. Li, & Tsai (2007), Zhang & Li (2007), Wang, R. Li, & Tsai (2007), Zou & Yuan (2008), Zou & Li (2008), Johnson, Lin, & Zeng (2008), ...

This talk is based on

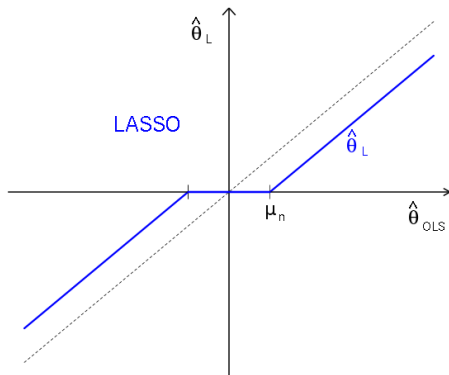
- [Pötscher & Leeb, 2007](#). Finite-sample distribution, moving-parameter asymptotics for hard-thresholding, LASSO, and SCAD. Impossibility result for the estimation of the cdf.
- [Pötscher & Schneider, 2007](#). Analogous results for the adaptive LASSO.
- [Pötscher & Schneider, 2008](#). Finite-sample and asymptotic coverage probabilities of confidence sets for hard-thresholding, LASSO, ad. LASSO.

Definition of the (adaptive) LASSO estimator $\hat{\theta}_{AL}$

LASSO estimator (Tibshirani, 1996)

$$\hat{\theta}_L = \arg \min_{\theta \in \mathbb{R}^k} \|\mathbf{y} - X\theta\|^2 + 2n\mu_n \sum_{i=1}^k |\theta_i| \quad \mu_n > 0$$

Tuning parameter $\lambda_n = 2n\mu_n$. For $k = 1$:

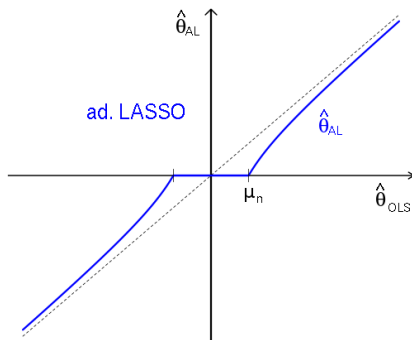


Definition of the (adaptive) LASSO estimator $\hat{\theta}_{AL}$

adaptive LASSO estimator (Zou, 2006)

$$\hat{\theta}_{AL} = \arg \min_{\theta \in \mathbb{R}^k} \|\mathbf{y} - \mathbf{X}\theta\|^2 + 2n\mu_n^2 \sum_{i=1}^k |\theta_i| / |\hat{\theta}_{OLS,j}| \quad \mu_n > 0$$

Tuning parameter $\lambda_n = 2n\mu_n^2$. For $k = 1$:



Two regimes for consistency

In terms of **model selection consistency**, two possible regimes for the tuning parameter μ_n arise.

- 1 The case $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow m$, $0 \leq m < \infty$, corresponds to **conservative** model selection (non-sparsely tuned).
- 2 The case $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$ corresponds to **consistent** model selection (sparsely tuned).

Remark (**estimation consistency**).

If $\mu_n \not\rightarrow 0$, then $\hat{\theta}_{AL}$ is not even consistent for θ . Therefore, $\mu_n \rightarrow 0$ is a “basic condition”.

We will focus on 2 here, also discuss 1 .

Two regimes for consistency

In terms of **model selection consistency**, two possible regimes for the tuning parameter μ_n arise.

- 1 The case $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow m$, $0 \leq m < \infty$, corresponds to **conservative** model selection (non-sparsely tuned).
- 2 The case $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$ corresponds to **consistent** model selection (sparsely tuned).

Remark (**estimation consistency**).

If $\mu_n \not\rightarrow 0$, then $\hat{\theta}_{\text{AL}}$ is not even consistent for θ . Therefore, $\mu_n \rightarrow 0$ is a “basic condition”.

We will focus on 2 here, also discuss 1 .

Two regimes for consistency

In terms of **model selection consistency**, two possible regimes for the tuning parameter μ_n arise.

- 1 The case $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow m$, $0 \leq m < \infty$, corresponds to **conservative** model selection (non-sparsely tuned).
- 2 The case $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$ corresponds to **consistent** model selection (sparsely tuned).

Remark (**estimation consistency**).

If $\mu_n \not\rightarrow 0$, then $\hat{\theta}_{AL}$ is not even consistent for θ . Therefore, $\mu_n \rightarrow 0$ is a “basic condition”.

We will focus on 2 here, also discuss 1 .

Two regimes for consistency

In terms of **model selection consistency**, two possible regimes for the tuning parameter μ_n arise.

- 1 The case $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow m$, $0 \leq m < \infty$, corresponds to **conservative** model selection (non-sparsely tuned).
- 2 The case $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$ corresponds to **consistent** model selection (sparsely tuned).

Remark (**estimation consistency**).

If $\mu_n \not\rightarrow 0$, then $\hat{\theta}_{AL}$ is not even consistent for θ . Therefore, $\mu_n \rightarrow 0$ is a “basic condition”.

We will focus on 2 here, also discuss 1 .

Zou (2006) “oracle property”

Suppose $X'X/n \rightarrow Q > 0$ and $\varepsilon_t \stackrel{\text{iid}}{\sim} (0, \sigma^2)$.

If $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$ and additionally $n^{1/4}\mu_n \rightarrow 0$, then

$$n^{1/2}(\hat{\theta}_{\text{AL}} - \theta) \rightarrow N(\mathbf{0}, \Sigma_{\theta}),$$

where Σ_{θ} is the asymptotic VC-matrix of the restricted LS-estimator based on the **unknown** true zero restrictions.

- Does this theorem provide meaningful insights? Finite-sample distribution?
- Asymptotic behavior under regime ① ?
- What if condition $n^{1/4}\mu_n \rightarrow 0$ is dropped in ② ?
- Pointwise vs. uniform consistency rates?
- Properties of confidence intervals?
- Estimability of finite-sample distribution?

- Does this theorem provide meaningful insights? Finite-sample distribution?
- Asymptotic behavior under regime ① ?
 - What if condition $n^{1/4} \mu_n \rightarrow 0$ is dropped in ② ?
 - Pointwise vs. uniform consistency rates?
 - Properties of confidence intervals?
 - Estimability of finite-sample distribution?

- Does this theorem provide meaningful insights? Finite-sample distribution?
- Asymptotic behavior under regime ① ?
- What if condition $n^{1/4}\mu_n \rightarrow 0$ is dropped in ② ?
- Pointwise vs. uniform consistency rates?
- Properties of confidence intervals?
- Estimability of finite-sample distribution?

- Does this theorem provide meaningful insights? Finite-sample distribution?
- Asymptotic behavior under regime ① ?
- What if condition $n^{1/4}\mu_n \rightarrow 0$ is dropped in ② ?
- Pointwise vs. uniform consistency rates?
- Properties of confidence intervals?
- Estimability of finite-sample distribution?

- Does this theorem provide meaningful insights? Finite-sample distribution?
- Asymptotic behavior under regime ① ?
- What if condition $n^{1/4}\mu_n \rightarrow 0$ is dropped in ② ?
- Pointwise vs. uniform consistency rates?
- Properties of confidence intervals?
- Estimability of finite-sample distribution?

- Does this theorem provide meaningful insights? Finite-sample distribution?
- Asymptotic behavior under regime ① ?
- What if condition $n^{1/4}\mu_n \rightarrow 0$ is dropped in ② ?
- Pointwise vs. uniform consistency rates?
- Properties of confidence intervals?
- Estimability of finite-sample distribution?

- Does this theorem provide meaningful insights? Finite-sample distribution?
- Asymptotic behavior under regime ① ?
- What if condition $n^{1/4}\mu_n \rightarrow 0$ is dropped in ② ?
- Pointwise vs. uniform consistency rates?
- Properties of confidence intervals?
- Estimability of finite-sample distribution?

We answer these questions within a normal linear regression model and address the non-orthogonal case in a simulation study.

Explicit solution in a simple model

- X is non-stochastic ($n \times k$), $rk(X) = k$.
- $\varepsilon \sim N_n(0, \sigma^2 \mathcal{I}_n)$
- For the **theoretical analysis**, assume that σ^2 is known and that $X'X$ is diagonal, in particular $X'X = n\mathcal{I}_k$.
- Remove these assumptions for **simulation results** concerning the finite-sample distribution.

Wlog consider Gaussian location model $y_1, \dots, y_n \stackrel{\text{iid}}{\sim} N(\theta, 1)$.
Then $\hat{\theta}_{\text{OLS}} = \bar{y}$ with $\hat{\theta}_{\text{OLS}} \sim N(\theta, 1/n)$ and

$$\hat{\theta}_{\text{AL}} = \begin{cases} 0 & \text{if } |\bar{y}| \leq \mu_n \\ \bar{y} - \mu_n^2/\bar{y} & \text{if } |\bar{y}| > \mu_n \end{cases}$$

Selects between restricted $\{N(0, 1)\}$ and full model $\{N(\theta, 1) : \theta \in \mathbb{R}\}$

The finite-sample distribution of $\hat{\theta}_{\text{AL}}$

The cdf $F_{n,\theta}(x) = P_{n,\theta}(n^{1/2}(\hat{\theta}_{\text{AL}} - \theta) \leq x)$ of $\hat{\theta}_{\text{AL}}$ is given by

$$\mathbf{1}(n^{1/2}\theta + x \geq 0) \Phi\left(z_{n,\theta}^{(2)}(x)\right) + \mathbf{1}(n^{1/2}\theta + x < 0) \Phi\left(z_{n,\theta}^{(1)}(x)\right).$$

$z_{n,\theta}^{(2)}(x)$ and $z_{n,\theta}^{(1)}(x)$ are $-(n^{1/2}\theta - x)/2 \pm \sqrt{((n^{1/2}\theta + x)/2)^2 + n\mu_n^2}$.

Φ and ϕ the cdf and pdf of $N(0, 1)$, resp.

The finite-sample distribution of $\hat{\theta}_{AL}$

The cdf $F_{n,\theta}(x) = P_{n,\theta}(n^{1/2}(\hat{\theta}_{AL} - \theta) \leq x)$ of $\hat{\theta}_{AL}$ is given by

$$\mathbf{1}(n^{1/2}\theta + x \geq 0) \Phi\left(z_{n,\theta}^{(2)}(x)\right) + \mathbf{1}(n^{1/2}\theta + x < 0) \Phi\left(z_{n,\theta}^{(1)}(x)\right).$$

$z_{n,\theta}^{(2)}(x)$ and $z_{n,\theta}^{(1)}(x)$ are $-(n^{1/2}\theta - x)/2 \pm \sqrt{((n^{1/2}\theta + x)/2)^2 + n\mu_n^2}$.

$dF_{n,\theta}$ is given by

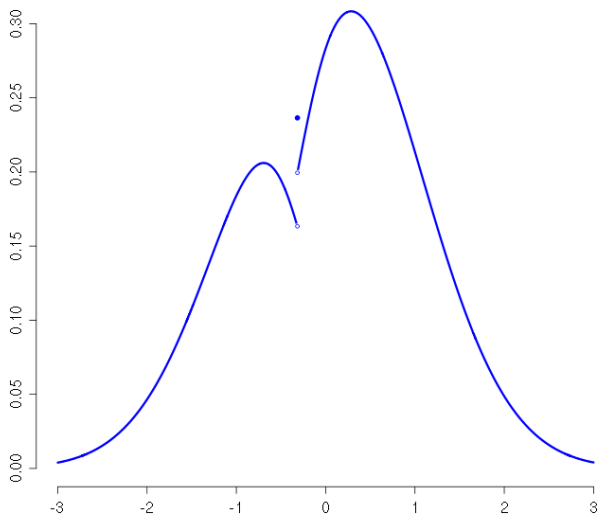
$$\begin{aligned} & \{ \Phi(n^{1/2}(-\theta + \mu_n)) - \Phi(n^{1/2}(-\theta - \mu_n)) \} d\delta_{-n^{1/2}\theta}(x) + \\ & 0.5 \times \{ \mathbf{1}(n^{1/2}\theta + x > 0) \phi\left(z_{n,\theta}^{(2)}(x)\right) (1 + t_{n,\theta}(x)) + \\ & \quad \mathbf{1}(n^{1/2}\theta + x < 0) \phi\left(z_{n,\theta}^{(1)}(x)\right) (1 - t_{n,\theta}(x)) \} dx \end{aligned}$$

where $t_{n,\theta}(x) := (((n^{1/2}\theta + x)/2)^2 + n\mu_n^2)^{-1/2}$.

Φ and ϕ the cdf and pdf of $N(0, 1)$, resp.

The finite-sample distribution of $\hat{\theta}_{AL}$

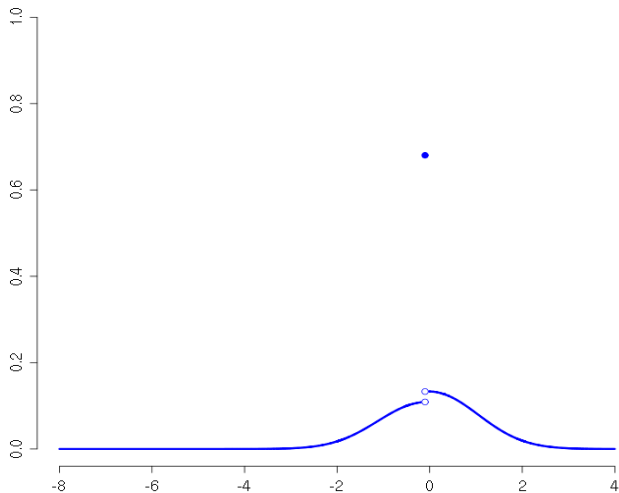
$n = 40, \theta = 0.05, \mu_n = 0.05$



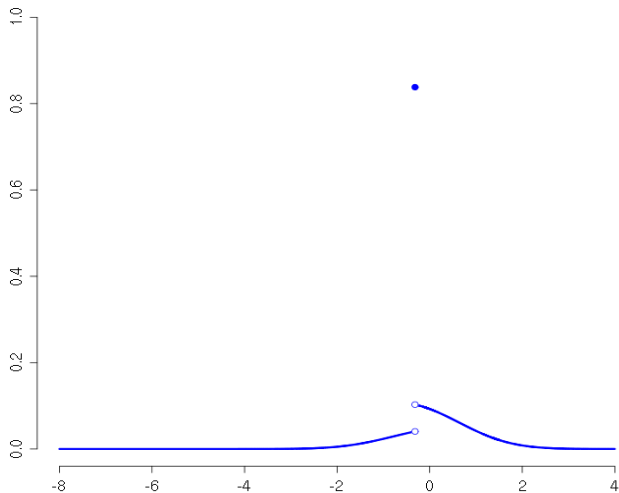
Non-normality??

- Finite-sample distribution is highly non-normal.
- Oracle property predicts normality (asymptotically).

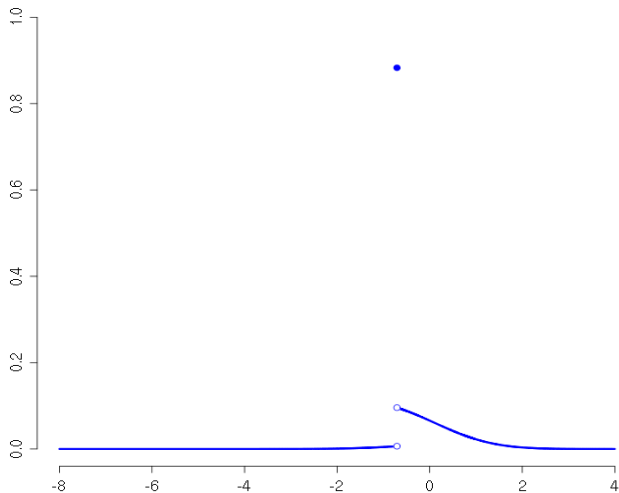
$$n = 1, \quad \mu_n = n^{-1/3} \text{ (consistent case)}$$



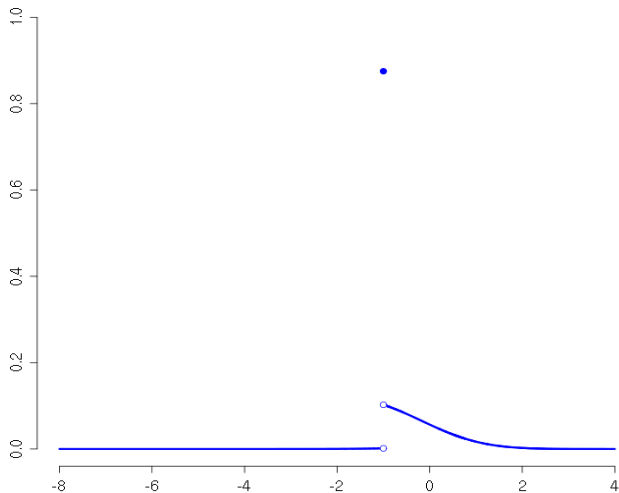
$$n = 10, \quad \mu_n = n^{-1/3} \text{ (consistent case)}$$



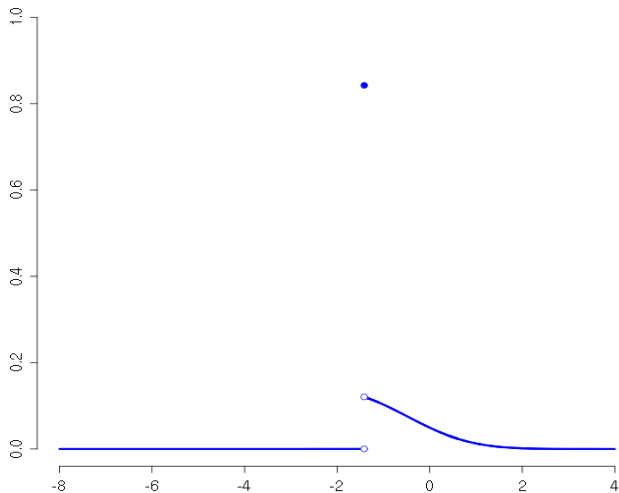
$$n = 50, \quad \mu_n = n^{-1/3} \text{ (consistent case)}$$



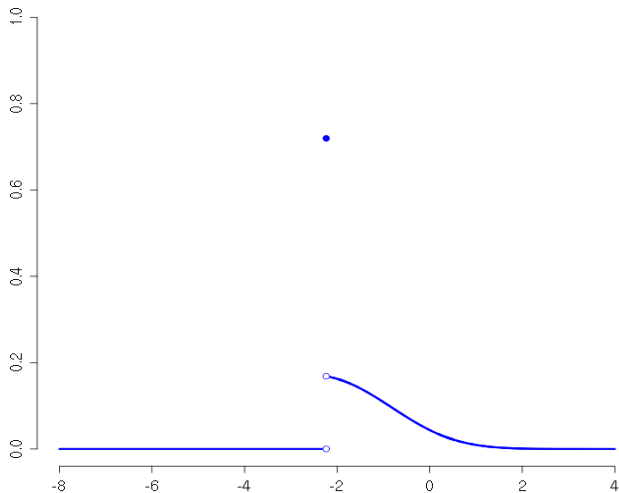
$$n = 100, \quad \mu_n = n^{-1/3} \text{ (consistent case)}$$



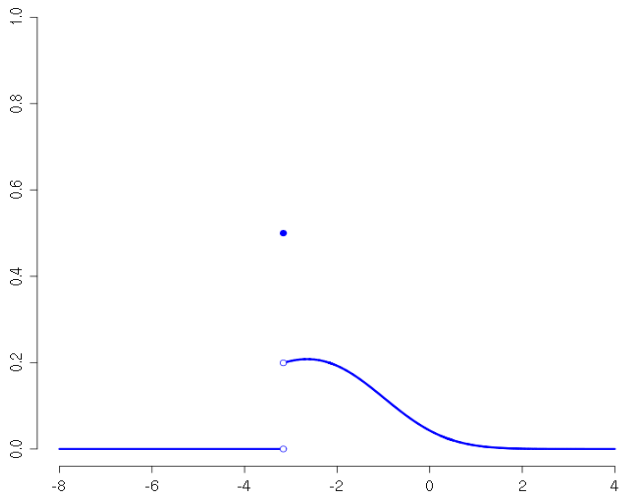
$n = 200$, $\mu_n = n^{-1/3}$ (consistent case)



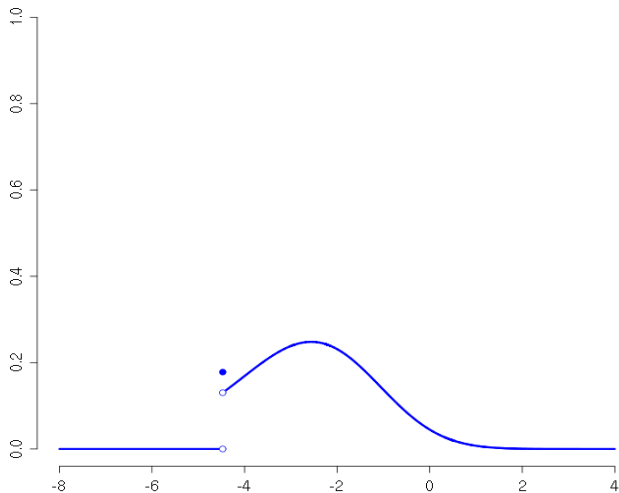
$$n = 500, \quad \mu_n = n^{-1/3} \text{ (consistent case)}$$



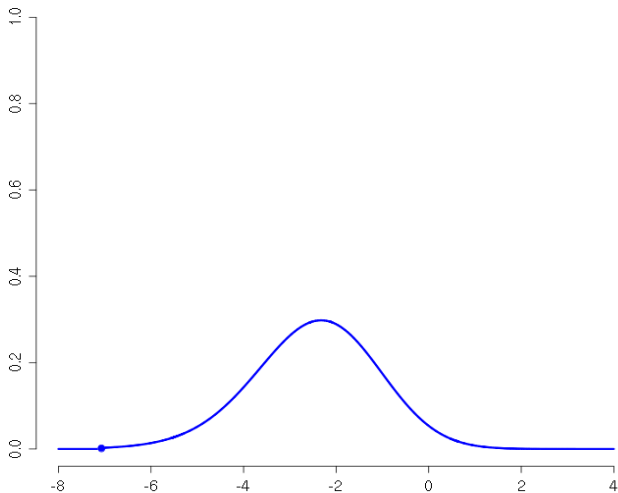
$n = 1000$, $\mu_n = n^{-1/3}$ (consistent case)



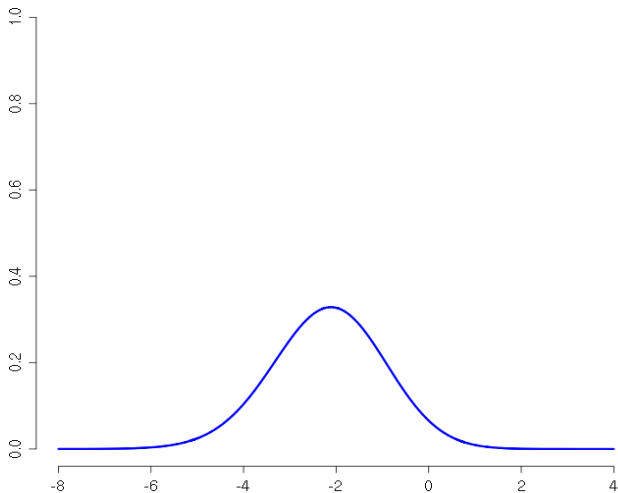
$n = 2000$, $\mu_n = n^{-1/3}$ (consistent case)



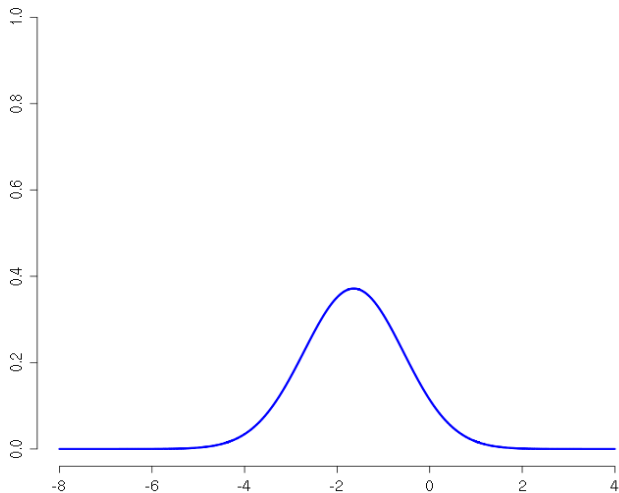
$n = 5000$, $\mu_n = n^{-1/3}$ (consistent case)



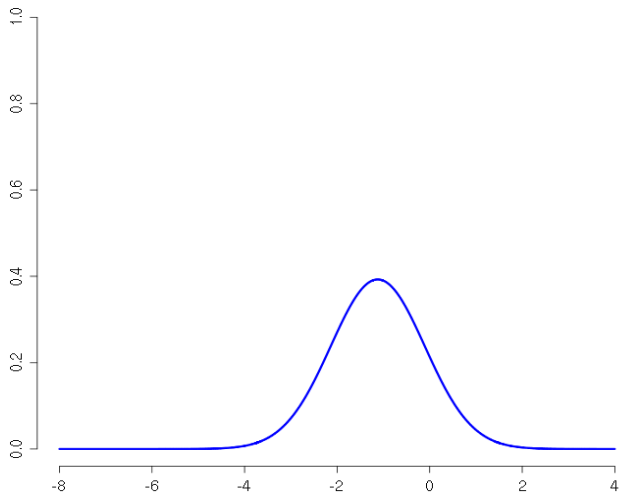
$$n = 10^4, \quad \mu_n = n^{-1/3} \text{ (consistent case)}$$



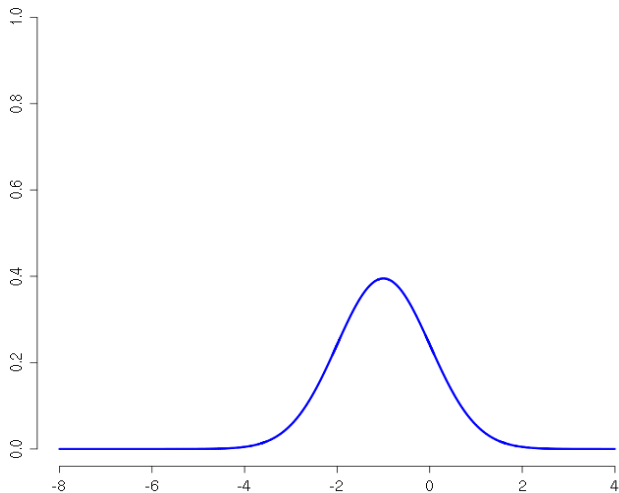
$n = 5 \times 10^4$, $\mu_n = n^{-1/3}$ (consistent case)



$n = 5 \times 10^5$, $\mu_n = n^{-1/3}$ (consistent case)



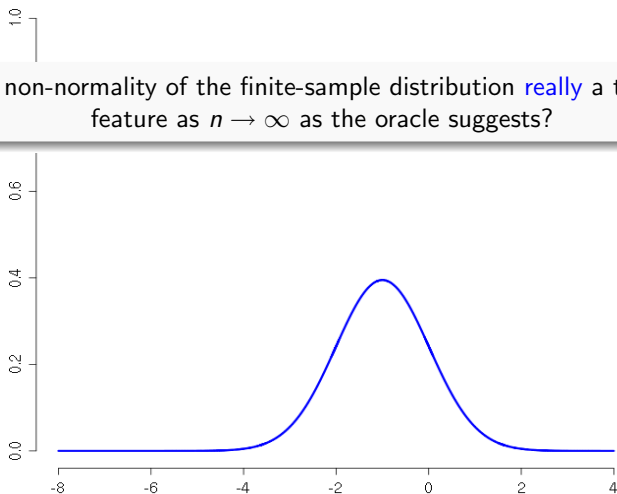
$$n = 10^6, \quad \mu_n = n^{-1/3} \text{ (consistent case)}$$



The Oracle (fixed-parameter asymptotics)

$$n = 10^6, \quad \mu_n = n^{-1/3} \text{ (consistent case)}$$

Is the non-normality of the finite-sample distribution **really** a transient feature as $n \rightarrow \infty$ as the oracle suggests?

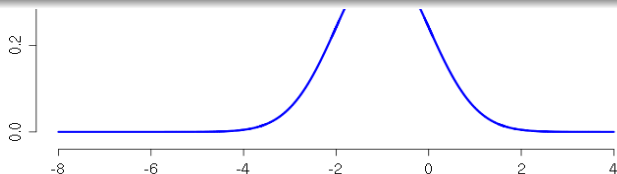


The Oracle (fixed-parameter asymptotics)

$$n = 10^6, \quad \mu_n = n^{-1/3} \text{ (consistent case)}$$

Is the non-normality of the finite-sample distribution **really** a transient feature as $n \rightarrow \infty$ as the oracle suggests?

Need to look at moving-parameter asymptotics!



Moving-parameter asymptotics?

Let underlying parameter θ depend on sample size:

Let $\theta_n \in \mathbb{R}$ be arbitrary, subject only to
 $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$ and $n^{1/2}\theta_n \rightarrow \nu \in \mathbb{R} \cup \{-\infty, \infty\}$.

This is not really a restriction since every subsequence of θ_n contains a further subsequence with these properties. Also note that $\zeta \neq 0$ implies $\nu = \pm\infty$.

2 Consistent case.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$ and $n^{1/2}\theta_n \rightarrow \nu \in \mathbb{R} \cup \{-\infty, \infty\}$. Then F_{n,θ_n} converges weakly to

- If $0 \leq |\zeta| < \infty$: pointmass at $-\nu$
- If $|\zeta| = \infty$: $\Phi(\cdot + \rho/\theta)$ where $n^{1/2}\mu_n^2 \rightarrow \rho$.

Depending on ζ , ν and ρ , three possible limits arise.

- Distribution collapses at a point.
- Total mass escapes to $\pm\infty$.
- Limit distribution is normal (possibly shifted!).

2 Consistent case.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$ and $n^{1/2}\theta_n \rightarrow \nu \in \mathbb{R} \cup \{-\infty, \infty\}$. Then F_{n,θ_n} converges weakly to

- If $0 \leq |\zeta| < \infty$: pointmass at $-\nu$
- If $|\zeta| = \infty$: $\Phi(\cdot + \rho/\theta)$ where $n^{1/2}\mu_n^2 \rightarrow \rho$.

Depending on ζ , ν and ρ , three possible limits arise.

- Distribution collapses at a point.
- Total mass escapes to $\pm\infty$.
- Limit distribution is normal (possibly shifted!).

2 Consistent case.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$ and $n^{1/2}\theta_n \rightarrow \nu \in \mathbb{R} \cup \{-\infty, \infty\}$. Then F_{n,θ_n} converges weakly to

- If $0 \leq |\zeta| < \infty$: pointmass at $-\nu$
- If $|\zeta| = \infty$: $\Phi(\cdot + \rho/\theta)$ where $n^{1/2}\mu_n^2 \rightarrow \rho$.

Depending on ζ , ν and ρ , three possible limits arise.

- Distribution collapses at a point.
- Total mass escapes to $\pm\infty$.
- Limit distribution is normal (possibly shifted!).

2 Consistent case.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$ and $n^{1/2}\theta_n \rightarrow \nu \in \mathbb{R} \cup \{-\infty, \infty\}$. Then F_{n,θ_n} converges weakly to

- If $0 \leq |\zeta| < \infty$: pointmass at $-\nu$
- If $|\zeta| = \infty$: $\Phi(\cdot + \rho/\theta)$ where $n^{1/2}\mu_n^2 \rightarrow \rho$.

Depending on ζ , ν and ρ , three possible limits arise.

- Distribution collapses at a point.
- Total mass escapes to $\pm\infty$.
- Limit distribution is normal (possibly shifted!).

2 Consistent case.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$ and $n^{1/2}\theta_n \rightarrow \nu \in \mathbb{R} \cup \{-\infty, \infty\}$. Then F_{n,θ_n} converges weakly to

- If $0 \leq |\zeta| < \infty$: pointmass at $-\nu$
- If $|\zeta| = \infty$: $\Phi(\cdot + \rho/\theta)$ where $n^{1/2}\mu_n^2 \rightarrow \rho$.

Depending on ζ , ν and ρ , three possible limits arise.

- Distribution collapses at a point.
- Total mass escapes to $\pm\infty$.
- Limit distribution is normal (possibly shifted!).

2 Consistent case.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$ and $n^{1/2}\theta_n \rightarrow \nu \in \mathbb{R} \cup \{-\infty, \infty\}$. Then F_{n,θ_n} converges weakly to

- If $0 \leq |\zeta| < \infty$: pointmass at $-\nu$
- If $|\zeta| = \infty$: $\Phi(\cdot + \rho/\theta)$ where $n^{1/2}\mu_n^2 \rightarrow \rho$.

Depending on ζ , ν and ρ , three possible limits arise.

- Distribution collapses at a point.
- Total mass escapes to $\pm\infty$.
- Limit distribution is normal (possibly shifted!).

Non-normality persists!!

Illustration: collapsing to pointmass

Example 1: $n = 1$, $\zeta = 0$, $\nu = 2$ ($\mu_n = n^{-1/3}$, $\theta_n = 2n^{-1/2}$)

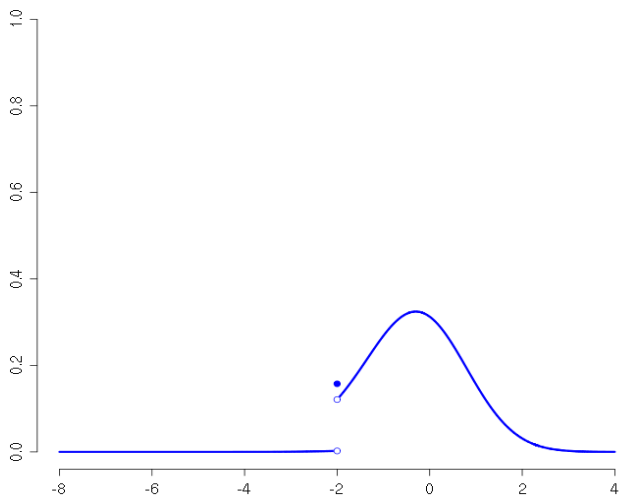


Illustration: collapsing to pointmass

Example 1: $n = 10$, $\zeta = 0$, $\nu = 2$ ($\mu_n = n^{-1/3}$, $\theta_n = 2n^{-1/2}$)

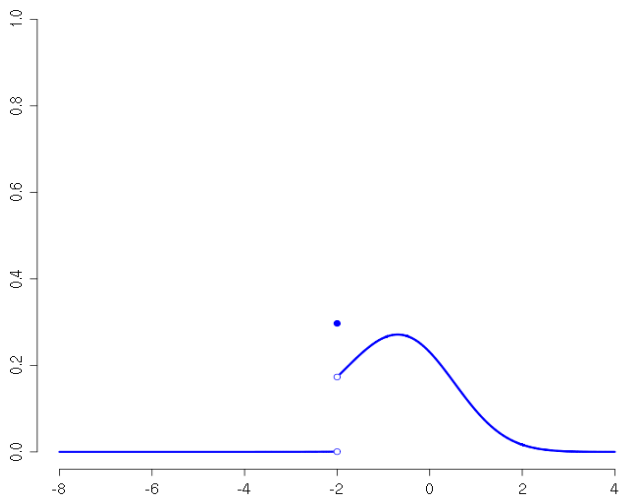


Illustration: collapsing to pointmass

Example 1: $n = 50$, $\zeta = 0$, $\nu = 2$ ($\mu_n = n^{-1/3}$, $\theta_n = 2n^{-1/2}$)

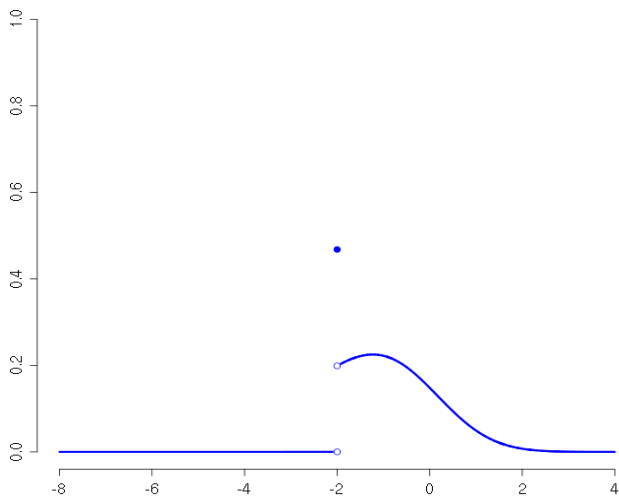


Illustration: collapsing to pointmass

Example 1: $n = 100$, $\zeta = 0$, $\nu = 2$ ($\mu_n = n^{-1/3}$, $\theta_n = 2n^{-1/2}$)

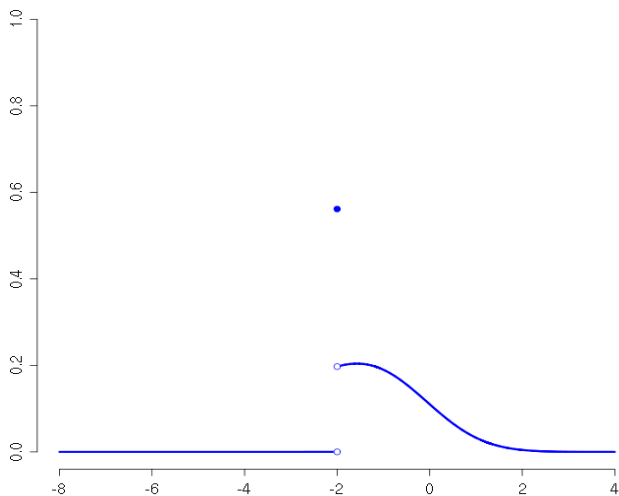


Illustration: collapsing to pointmass

Example 1: $n = 200$, $\zeta = 0$, $\nu = 2$ ($\mu_n = n^{-1/3}$, $\theta_n = 2n^{-1/2}$)

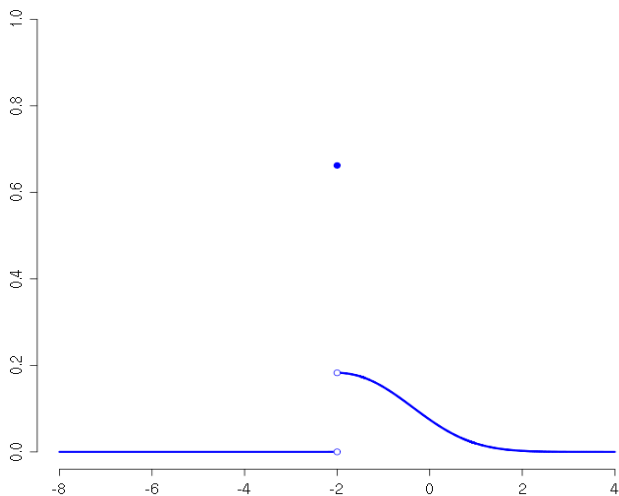


Illustration: collapsing to pointmass

Example 1: $n = 500$, $\zeta = 0$, $\nu = 2$ ($\mu_n = n^{-1/3}$, $\theta_n = 2n^{-1/2}$)

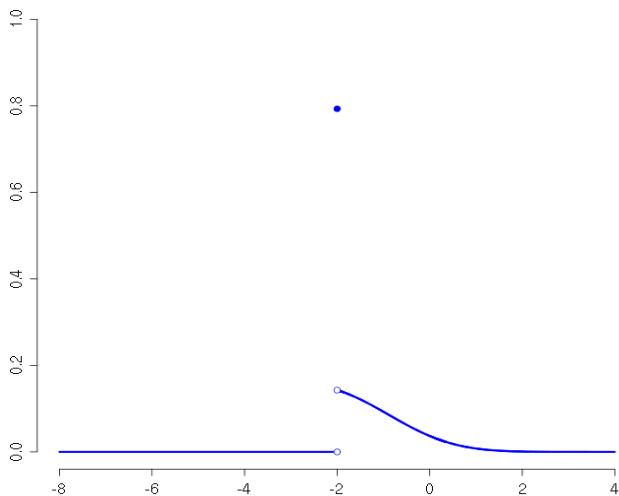


Illustration: collapsing to pointmass

Example 1: $n = 1000$, $\zeta = 0$, $\nu = 2$ ($\mu_n = n^{-1/3}$, $\theta_n = 2n^{-1/2}$)

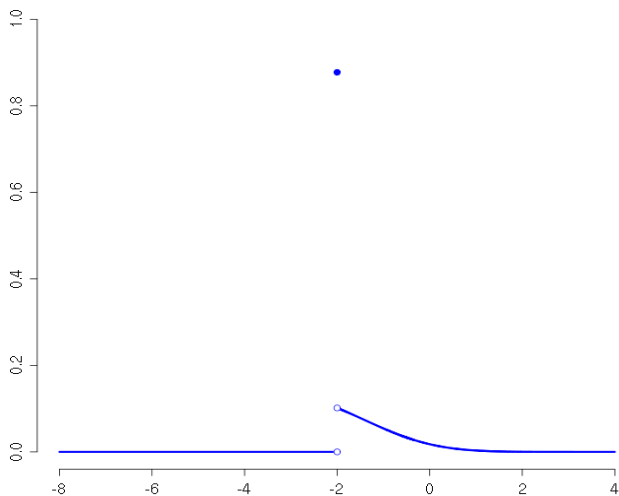


Illustration: collapsing to pointmass

Example 1: $n = 2000$, $\zeta = 0$, $\nu = 2$ ($\mu_n = n^{-1/3}$, $\theta_n = 2n^{-1/2}$)

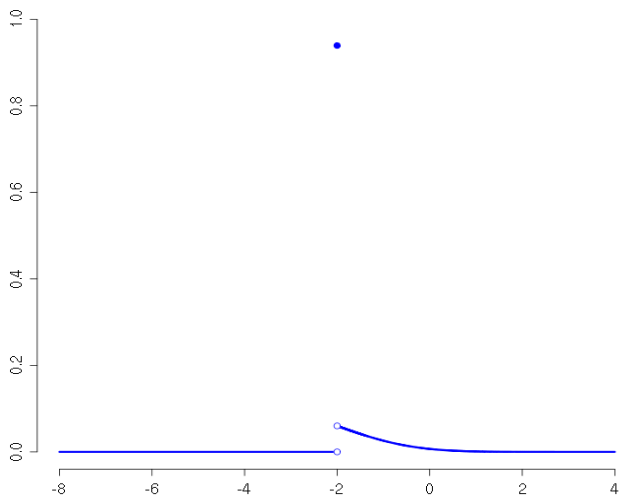


Illustration: collapsing to pointmass

Example 1: $n = 5000$, $\zeta = 0$, $\nu = 2$ ($\mu_n = n^{-1/3}$, $\theta_n = 2n^{-1/2}$)

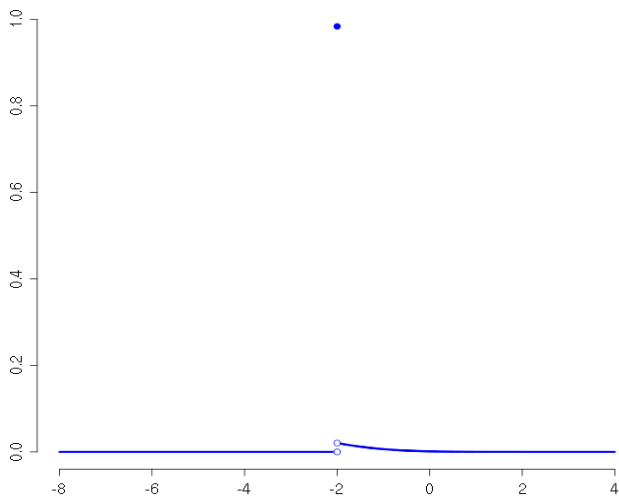


Illustration: collapsing to pointmass

Example 1: $n = 10^4$, $\zeta = 0$, $\nu = 2$ ($\mu_n = n^{-1/3}$, $\theta_n = 2n^{-1/2}$)

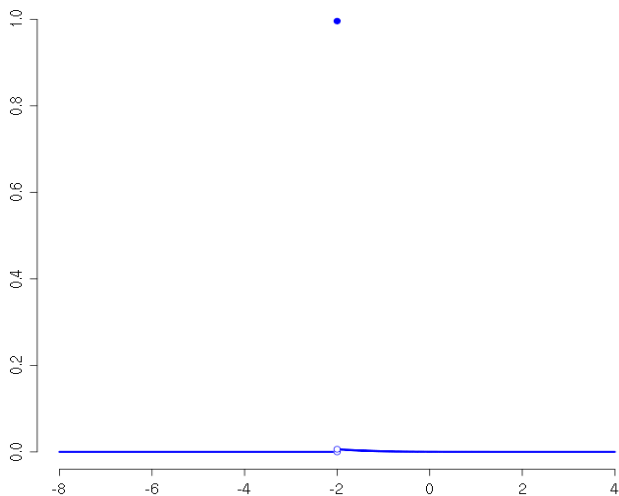


Illustration: collapsing to pointmass

Example 1: $n = 5 \times 10^4$, $\zeta = 0$, $\nu = 2$ ($\mu_n = n^{-1/3}$, $\theta_n = 2n^{-1/2}$)

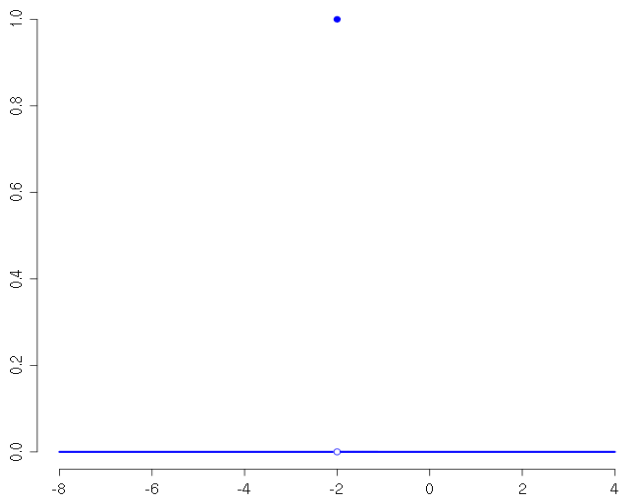
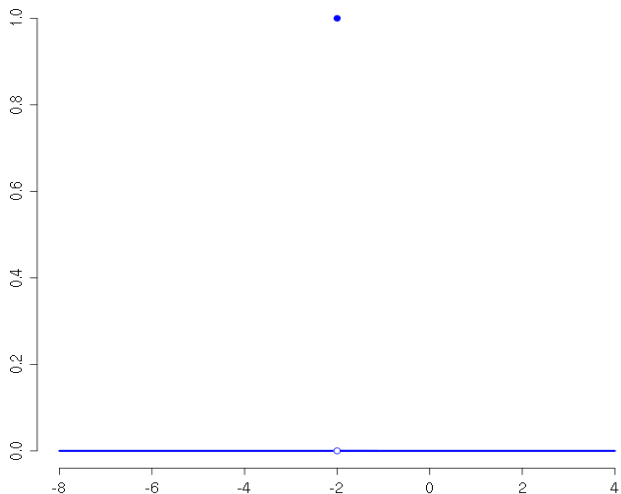


Illustration: collapsing to pointmass

Example 1: $n = 5 \times 10^4$, $\zeta = 0$, $\nu = 2$ ($\mu_n = n^{-1/3}$, $\theta_n = 2n^{-1/2}$)



END

Illustration: mass escaping to $-\infty$

Example 2: $n = 1$, $\zeta = 1$, $\nu = \infty$ ($\mu_n = n^{-1/5}$, $\theta_n = n^{-1/5}$)

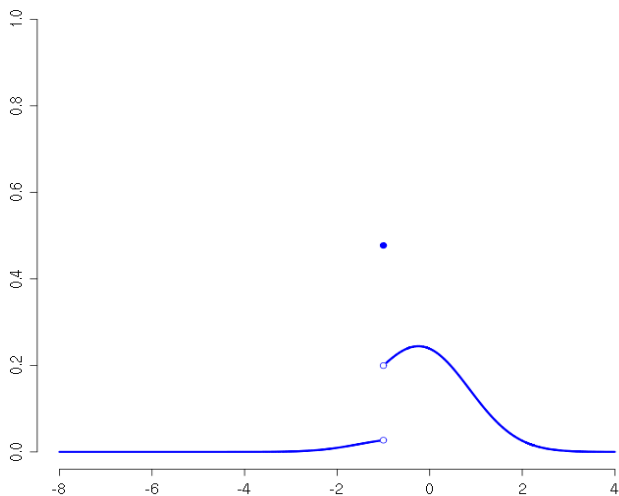


Illustration: mass escaping to $-\infty$

Example 2: $n = 10$, $\zeta = 1$, $\nu = \infty$ ($\mu_n = n^{-1/3}$, $\theta_n = n^{-1/5}$)

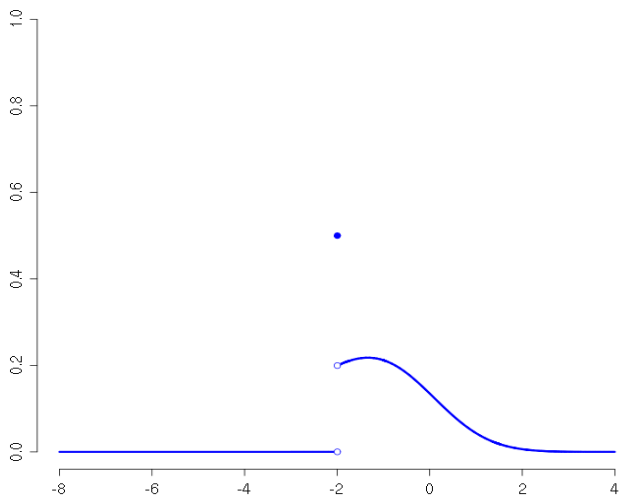


Illustration: mass escaping to $-\infty$

Example 2: $n = 50$, $\zeta = 1$, $\nu = \infty$ ($\mu_n = n^{-1/3}$, $\theta_n = n^{-1/5}$)

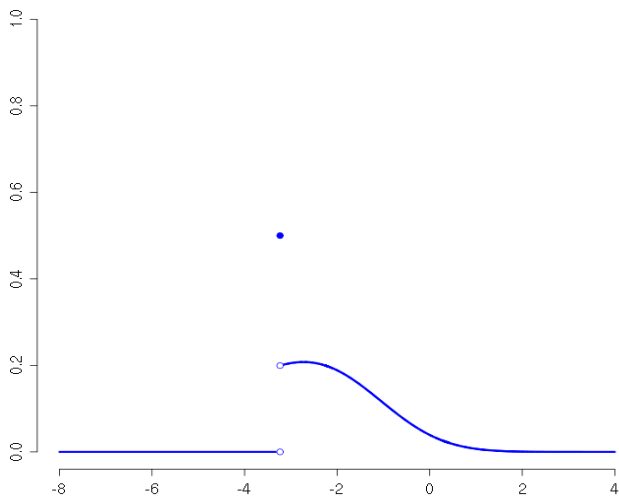


Illustration: mass escaping to $-\infty$

Example 2: $n = 100$, $\zeta = 1$, $\nu = \infty$ ($\mu_n = n^{-1/3}$, $\theta_n = n^{-1/5}$)

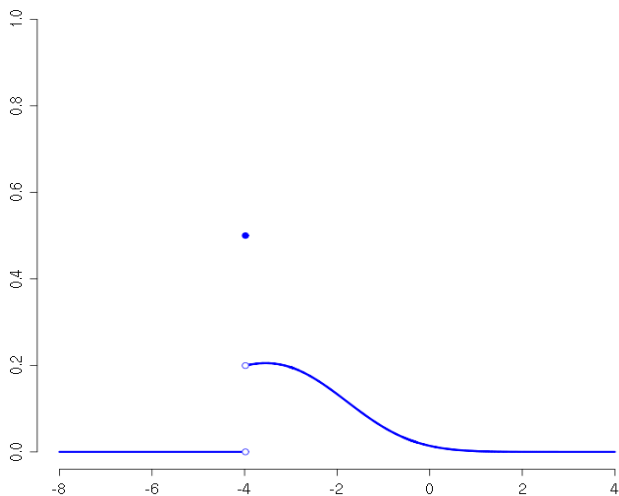


Illustration: mass escaping to $-\infty$

Example 2: $n = 200$, $\zeta = 1$, $\nu = \infty$ ($\mu_n = n^{-1/3}$, $\theta_n = n^{-1/5}$)

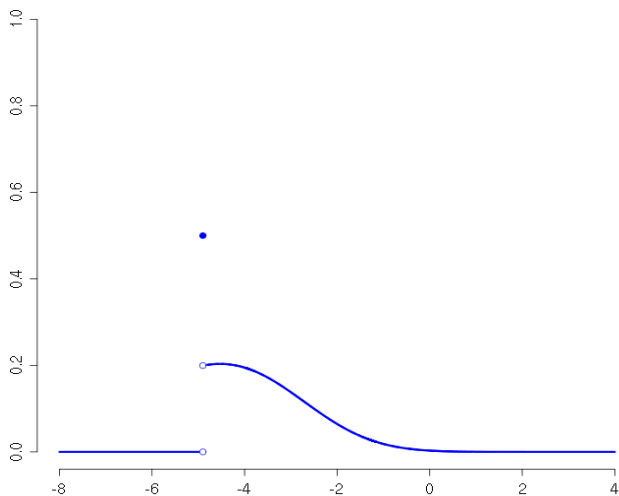


Illustration: mass escaping to $-\infty$

Example 2: $n = 500$, $\zeta = 1$, $\nu = \infty$ ($\mu_n = n^{-1/3}$, $\theta_n = n^{-1/5}$)

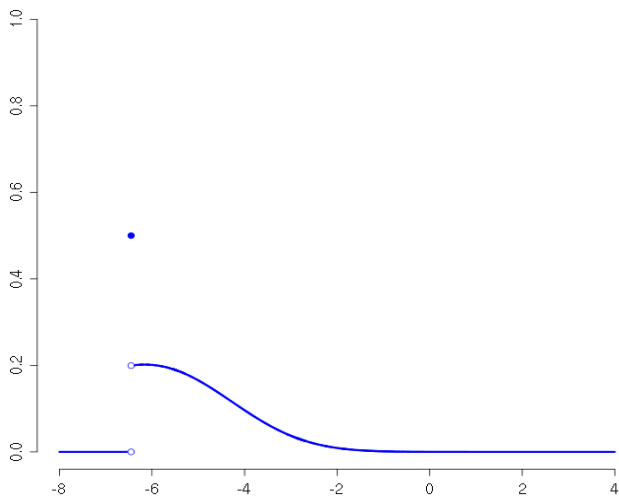


Illustration: mass escaping to $-\infty$

Example 2: $n = 1000$, $\zeta = 1$, $\nu = \infty$ ($\mu_n = n^{-1/3}$, $\theta_n = n^{-1/5}$)

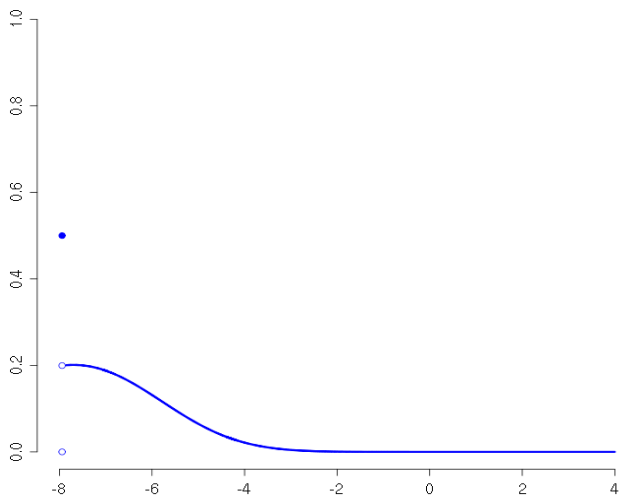


Illustration: mass escaping to $-\infty$

Example 2: $n = 2000$, $\zeta = 1$, $\nu = \infty$ ($\mu_n = n^{-1/3}$, $\theta_n = n^{-1/5}$)

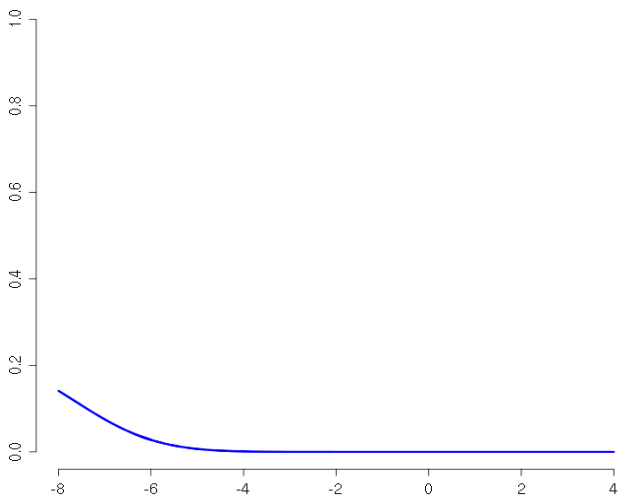


Illustration: mass escaping to $-\infty$

Example 2: $n = 5000$, $\zeta = 1$, $\nu = \infty$ ($\mu_n = n^{-1/3}$, $\theta_n = n^{-1/5}$)

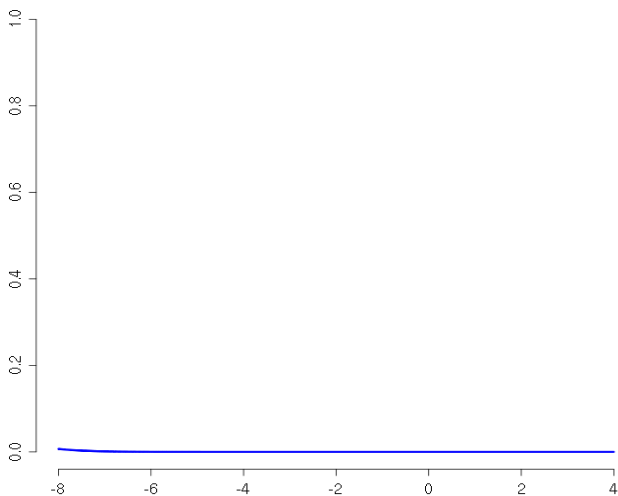


Illustration: mass escaping to $-\infty$

Example 2: $n = 10^4$, $\zeta = 1$, $\nu = \infty$ ($\mu_n = n^{-1/3}$, $\theta_n = n^{-1/5}$)

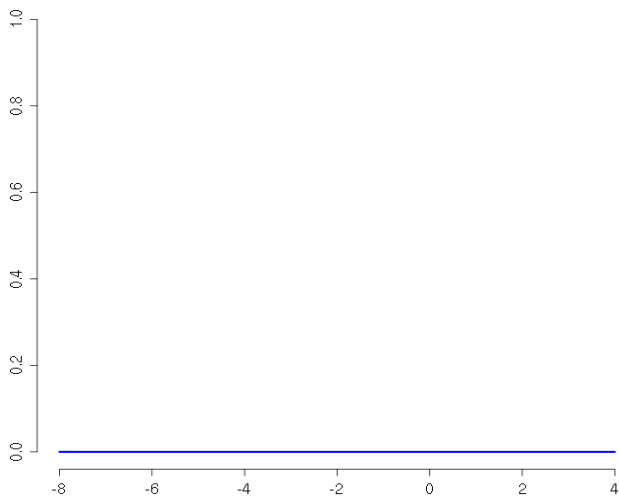
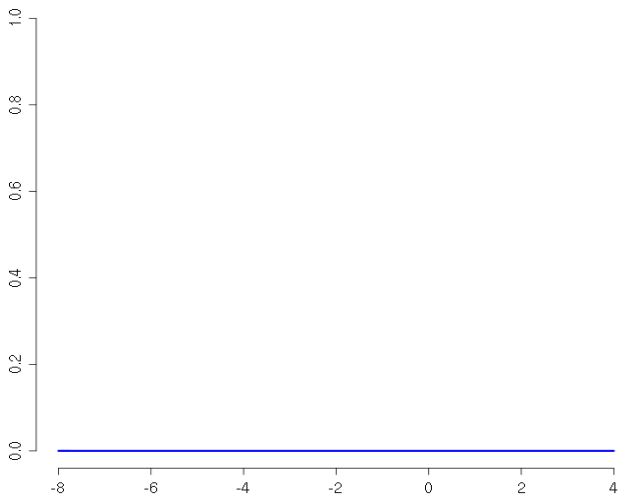


Illustration: mass escaping to $-\infty$

Example 2: $n = 10^4$, $\zeta = 1$, $\nu = \infty$ ($\mu_n = n^{-1/3}$, $\theta_n = n^{-1/5}$)



END

2 Consistent case.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$ and $n^{1/2}\theta_n \rightarrow \nu \in \mathbb{R} \cup \{-\infty, \infty\}$. Then F_{n,θ_n} converges weakly to

- If $0 \leq |\zeta| < \infty$: pointmass at $-\nu$
- If $|\zeta| = \infty$: $\Phi(\cdot + \rho/\theta)$ where $n^{1/2}\mu_n^2 \rightarrow \rho$.

Zou (pointwise case) ? Above theorem implies that

$$F_{n,\theta}(x) \rightarrow \begin{cases} \mathbf{1}(x \geq 0) & \theta = 0 \quad (\implies \zeta, \nu = 0) \\ \Phi(x + \rho/\theta) & \theta \neq 0 \quad (\implies |\zeta| = \infty) \end{cases}$$

Remark: $\rho = 0 \iff n^{1/4}\mu_n \rightarrow 0$.

2 Consistent case.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$ and $n^{1/2}\theta_n \rightarrow \nu \in \mathbb{R} \cup \{-\infty, \infty\}$. Then F_{n,θ_n} converges weakly to

- If $0 \leq |\zeta| < \infty$: pointmass at $-\nu$
- If $|\zeta| = \infty$: $\Phi(\cdot + \rho/\theta)$ where $n^{1/2}\mu_n^2 \rightarrow \rho$.

Zou (pointwise case) ? Above theorem implies that

$$F_{n,\theta}(x) \rightarrow \begin{cases} \mathbf{1}(x \geq 0) & \theta = 0 \quad (\implies \zeta, \nu = 0) \\ \Phi(x + \rho/\theta) & \theta \neq 0 \quad (\implies |\zeta| = \infty) \end{cases}$$

Remark: $\rho = 0 \iff n^{1/4}\mu_n \rightarrow 0$.

- Adaptive LASSO has in a **uniform sense** a **rate of convergence** that is **slower than $n^{1/2}$** .
- The “correct” uniform rate can be shown to be μ_n^{-1} .
- In a moving-parameter framework, the asymptotic distribution of $\mu_n^{-1}(\hat{\theta}_{\text{AL}} - \theta)$ collapses to pointmass.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$. Then

$G_{n,\theta_n} := P(\mu_n^{-1}(\hat{\theta}_{\text{AL}} - \theta) \leq x)$ converges weakly to

- If $|\zeta| < 1$: **pointmass at $-\zeta$**
- If $1 \leq |\zeta| < \infty$: **pointmass at $-1/\zeta$**
- If $|\zeta| = \infty$: **pointmass at 0**

- Adaptive LASSO has in a **uniform sense** a **rate of convergence** that is **slower than $n^{1/2}$** .
- The “correct” uniform rate can be shown to be μ_n^{-1} .
- In a moving-parameter framework, the asymptotic distribution of $\mu_n^{-1}(\hat{\theta}_{\text{AL}} - \theta)$ collapses to pointmass.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$. Then

$G_{n,\theta_n} := P(\mu_n^{-1}(\hat{\theta}_{\text{AL}} - \theta) \leq x)$ converges weakly to

- If $|\zeta| < 1$: **pointmass at $-\zeta$**
- If $1 \leq |\zeta| < \infty$: **pointmass at $-1/\zeta$**
- If $|\zeta| = \infty$: **pointmass at 0**

- Adaptive LASSO has in a **uniform sense** a **rate of convergence** that is **slower than $n^{1/2}$** .
- The “correct” uniform rate can be shown to be μ_n^{-1} .
- In a moving-parameter framework, the asymptotic distribution of $\mu_n^{-1}(\hat{\theta}_{\text{AL}} - \theta)$ collapses to pointmass.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$. Then

$G_{n,\theta_n} := P(\mu_n^{-1}(\hat{\theta}_{\text{AL}} - \theta) \leq x)$ converges weakly to

- If $|\zeta| < 1$: **pointmass at $-\zeta$**
- If $1 \leq |\zeta| < \infty$: **pointmass at $-1/\zeta$**
- If $|\zeta| = \infty$: **pointmass at 0**

- Adaptive LASSO has in a **uniform sense** a **rate of convergence** that is **slower than $n^{1/2}$** .
- The “correct” uniform rate can be shown to be μ_n^{-1} .
- In a moving-parameter framework, the asymptotic distribution of $\mu_n^{-1}(\hat{\theta}_{\text{AL}} - \theta)$ collapses to pointmass.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $\theta_n/\mu_n \rightarrow \zeta \in \mathbb{R} \cup \{-\infty, \infty\}$. Then

$G_{n,\theta_n} := P(\mu_n^{-1}(\hat{\theta}_{\text{AL}} - \theta) \leq x)$ converges weakly to

- If $|\zeta| < 1$: **pointmass at $-\zeta$**
- If $1 \leq |\zeta| < \infty$: **pointmass at $-1/\zeta$**
- If $|\zeta| = \infty$: **pointmass at 0**

Above theorems reflect that

$$\hat{\theta}_{\text{AL}} - \theta = \text{“BIAS”} + \text{“FLUCTUATION”}$$

where

- “BIAS” is $O(n^{-1/2})$ in a **pointwise** sense but is only $O(\mu_n)$ in a **uniform** sense, whereas
- “FLUCTUATION” is always of **order** $n^{-1/2}$.

1 Conservative case.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow m$, $0 \leq m < \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $n^{1/2}\theta_n \rightarrow \nu \in \mathbb{R} \cup \{-\infty, \infty\}$. Then F_{n,θ_n} converges weakly to

- If $\nu \in \mathbb{R}$

$$\mathbf{1}(\nu + x \geq 0) \Phi \left(-(\nu - x)/2 + \sqrt{((\nu + x)/2)^2 + m^2} \right) + \\ \mathbf{1}(\nu + x < 0) \Phi \left(-(\nu - x)/2 - \sqrt{((\nu + x)/2)^2 + m^2} \right)$$

- $\Phi(x)$ if $|\nu| = \infty$.

1 Conservative case.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow m$, $0 \leq m < \infty$. Suppose the true parameter $\theta_n \in \mathbb{R}$ satisfies $n^{1/2}\theta_n \rightarrow \nu \in \mathbb{R} \cup \{-\infty, \infty\}$. Then F_{n,θ_n} converges weakly to

- If $\nu \in \mathbb{R}$

$$\mathbf{1}(\nu + x \geq 0) \Phi \left(-(\nu - x)/2 + \sqrt{((\nu + x)/2)^2 + m^2} \right) + \\ \mathbf{1}(\nu + x < 0) \Phi \left(-(\nu - x)/2 - \sqrt{((\nu + x)/2)^2 + m^2} \right)$$

- $\Phi(x)$ if $|\nu| = \infty$.

Note: Asymptotic distributions are the same as finite-sample distribution, except that $n^{1/2}\theta_n$ and $n^{1/2}\mu_n$ have settled down to their limiting values, capturing finite-sample behavior very well.

- $\hat{\theta}_{AL}$ is now uniformly $n^{1/2}$ -consistent.
- Fixed-parameter asymptotics: previous theorem implies that $F_{n,\theta}(x)$ converges to
 - $\mathbf{1}(x \geq 0) \Phi\left(\frac{x}{2} + \sqrt{\left(\frac{x}{2}\right)^2 + m^2}\right) + \mathbf{1}(x < 0) \Phi\left(\frac{x}{2} - \sqrt{\left(\frac{x}{2}\right)^2 + m^2}\right)$
if $\theta = 0$ ($\nu = 0$)
 - $\Phi(x)$ if $\theta \neq 0$ ($|\nu| = \infty$)
- Fixed-parameter asymptotic distributions are also non-normal, capturing behavior the finite-sample distributions to some extent (no oracle here).

- $\hat{\theta}_{AL}$ is now **uniformly** $n^{1/2}$ -consistent.
- **Fixed-parameter** asymptotics: previous theorem implies that $F_{n,\theta}(x)$ converges to
 - $\mathbf{1}(x \geq 0) \Phi\left(\frac{x}{2} + \sqrt{\left(\frac{x}{2}\right)^2 + m^2}\right) + \mathbf{1}(x < 0) \Phi\left(\frac{x}{2} - \sqrt{\left(\frac{x}{2}\right)^2 + m^2}\right)$
if $\theta = 0$ ($\nu = 0$)
 - $\Phi(x)$ if $\theta \neq 0$ ($|\nu| = \infty$)
- Fixed-parameter asymptotic distributions are also non-normal, capturing behavior the finite-sample distributions to some extent (no oracle here).

- $\hat{\theta}_{AL}$ is now **uniformly** $n^{1/2}$ -consistent.
- **Fixed-parameter** asymptotics: previous theorem implies that $F_{n,\theta}(x)$ converges to
 - $\mathbf{1}(x \geq 0) \Phi\left(\frac{x}{2} + \sqrt{\left(\frac{x}{2}\right)^2 + m^2}\right) + \mathbf{1}(x < 0) \Phi\left(\frac{x}{2} - \sqrt{\left(\frac{x}{2}\right)^2 + m^2}\right)$
if $\theta = 0$ ($\nu = 0$)
 - $\Phi(x)$ if $\theta \neq 0$ ($|\nu| = \infty$)
- Fixed-parameter asymptotic distributions are also non-normal, capturing behavior the finite-sample distributions to some extent (no oracle here).

- $\hat{\theta}_{AL}$ is now **uniformly** $n^{1/2}$ -consistent.
- **Fixed-parameter** asymptotics: previous theorem implies that $F_{n,\theta}(x)$ converges to
 - $\mathbf{1}(x \geq 0) \Phi\left(\frac{x}{2} + \sqrt{\left(\frac{x}{2}\right)^2 + m^2}\right) + \mathbf{1}(x < 0) \Phi\left(\frac{x}{2} - \sqrt{\left(\frac{x}{2}\right)^2 + m^2}\right)$
if $\theta = 0$ ($\nu = 0$)
 - $\Phi(x)$ if $\theta \neq 0$ ($|\nu| = \infty$)
- Fixed-parameter asymptotic distributions are also non-normal, capturing behavior the finite-sample distributions to some extent (no oracle here).

Results are similar for **hard-thresholding**, **soft-thresholding** (**LASSO**), and **SCAD** estimator. (Pötscher & Leeb, 2007).

- Identical results in terms of (uniform) consistency.
- Analogous (asymptotic) distributional results.

Confidence sets based on the adaptive LASSO

Let $C_n = [\hat{\theta}_{AL} - a_n, \hat{\theta}_{AL} + b_n]$.

The infimal coverage probability $\inf_{\theta \in \mathbb{R}} P_{n,\theta}(\hat{\theta}_{AL} \in C_n)$ is given by

$$\Phi(n^{1/2}(a_n - \mu_n)) - \Phi\left(n^{1/2}\left(\frac{a_n - b_n}{2} - \sqrt{\left(\frac{a_n + b_n}{2}\right)^2 + \mu_n^2}\right)\right)$$

if $a_n \leq b_n$ and

$$\Phi\left(n^{1/2}\left(\frac{a_n - b_n}{2} + \sqrt{\left(\frac{a_n + b_n}{2}\right)^2 + \mu_n^2}\right)\right) - \Phi(n^{1/2}(-b_n + \mu_n))$$

if $a_n > b_n$.

Symmetric intervals ($a_n = b_n$) can be shown to be the shortest ones for a given infimal coverage probability δ .

Confidence sets based on PLSEs

- For each $n \in \mathbb{N}$, we have

$$a_{n,H} > a_{n,AL} > a_{n,L} > a_{n,OLS} \quad \text{for a given } \delta > 0$$

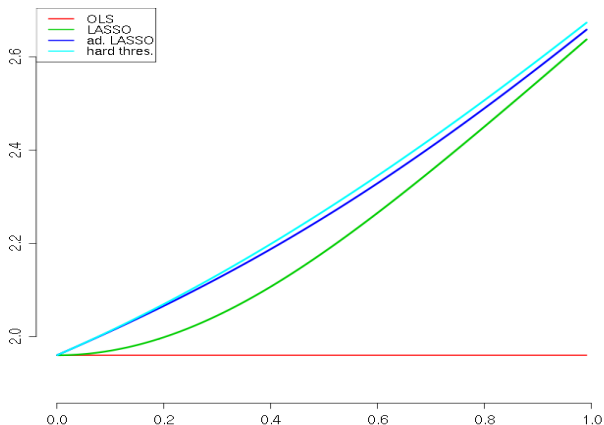
- Asymptotically, the following holds.
- ① **Conservative case.** All quantities are of the same order $n^{-1/2}$.

$$a_{n,H} \sim a_{n,AL} \sim a_{n,L} \sim a_{n,OLS}$$

- ② **Consistent case.** $a_{n,H}$, $a_{n,L}$, and $a_{n,A}$ are one order of magnitude larger than $a_{n,OLS}$.

$$a_H/a_{OLS} \sim a_{AL}/a_{OLS} \sim a_L/a_{OLS} \sim n^{1/2} \mu_n \rightarrow \infty$$

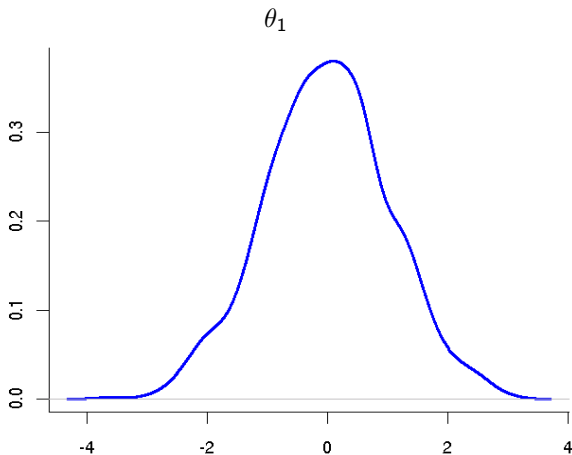
Plot of $n^{1/2}a_n$ against $n^{1/2}\mu_n$ for $\delta = 0.95$.



Simulations - remove orthogonality assumption

$k = 4$, $n = 200$, $\theta = (3, 1.5, 0, 0)' + 2/n^{1/2}(0, 0, 1, 1)'$, $X'X = n\Omega$ with $\Omega_{ij} = 0.5^{|i-j|}$, 1000 simulations

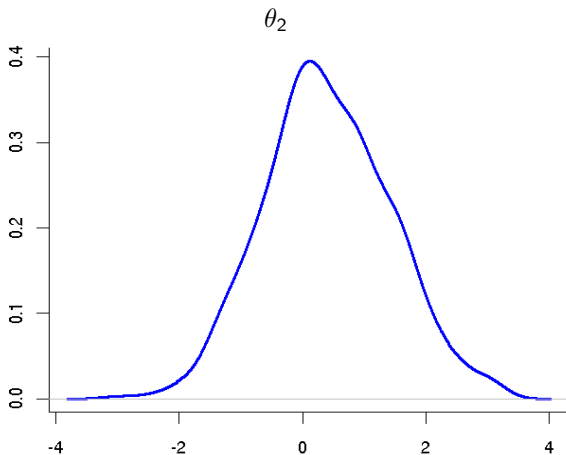
- $\mu_n = n^{-1/3}$



Simulations - remove orthogonality assumption

$k = 4$, $n = 200$, $\theta = (3, 1.5, 0, 0)' + 2/n^{1/2}(0, 0, 1, 1)'$, $X'X = n\Omega$ with $\Omega_{ij} = 0.5^{|i-j|}$, 1000 simulations

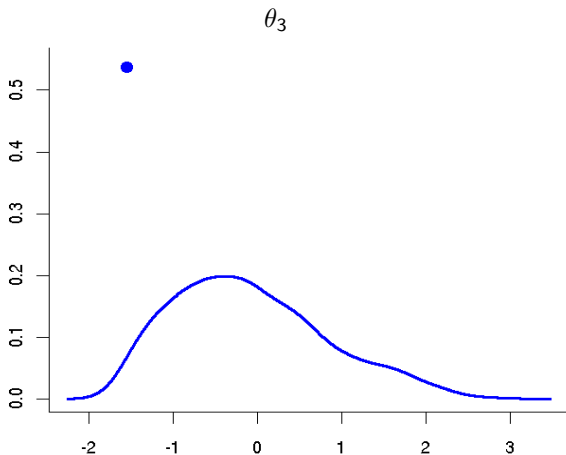
- $\mu_n = n^{-1/3}$



Simulations - remove orthogonality assumption

$k = 4$, $n = 200$, $\theta = (3, 1.5, 0, 0)' + 2/n^{1/2}(0, 0, 1, 1)'$, $X'X = n\Omega$ with $\Omega_{ij} = 0.5^{|i-j|}$, 1000 simulations

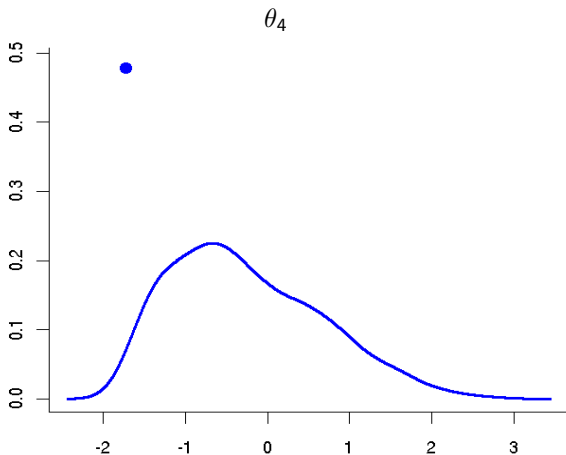
- $\mu_n = n^{-1/3}$



Simulations - remove orthogonality assumption

$k = 4$, $n = 200$, $\theta = (3, 1.5, 0, 0)' + 2/n^{1/2}(0, 0, 1, 1)'$, $X'X = n\Omega$ with $\Omega_{ij} = 0.5^{|i-j|}$, 1000 simulations

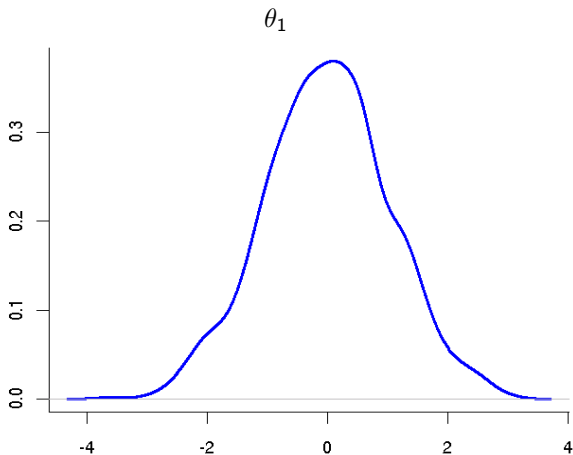
- $\mu_n = n^{-1/3}$



Simulations - remove orthogonality assumption

$k = 4$, $n = 200$, $\theta = (3, 1.5, 0, 0)' + 2/n^{1/2}(0, 0, 1, 1)'$, $X'X = n\Omega$ with $\Omega_{ij} = 0.5^{|i-j|}$, 1000 simulations

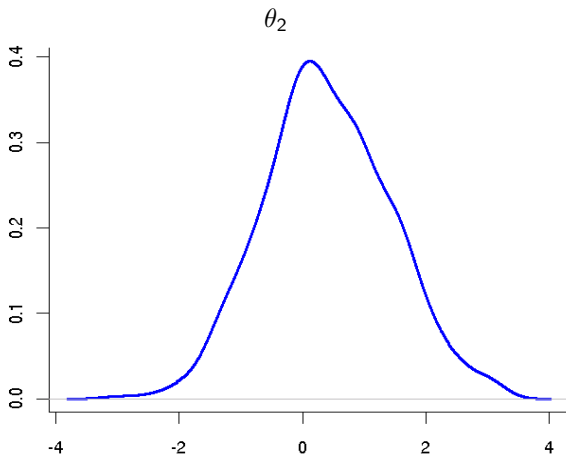
- Choose μ_n through cross-validation.



Simulations - remove orthogonality assumption

$k = 4$, $n = 200$, $\theta = (3, 1.5, 0, 0)' + 2/n^{1/2}(0, 0, 1, 1)'$, $X'X = n\Omega$ with $\Omega_{ij} = 0.5^{|i-j|}$, 1000 simulations

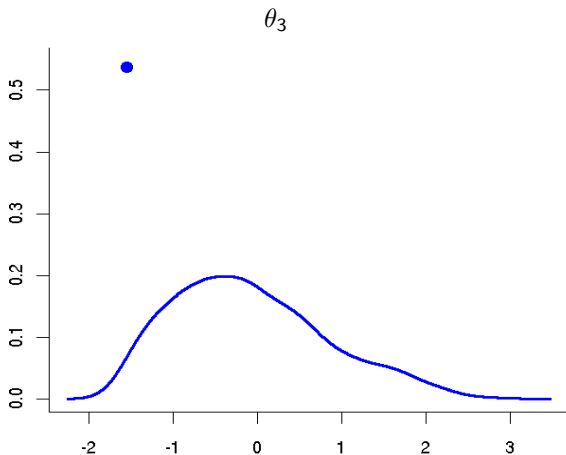
- Choose μ_n through cross-validation.



Simulations - remove orthogonality assumption

$k = 4$, $n = 200$, $\theta = (3, 1.5, 0, 0)' + 2/n^{1/2}(0, 0, 1, 1)'$, $X'X = n\Omega$ with $\Omega_{ij} = 0.5^{|i-j|}$, 1000 simulations

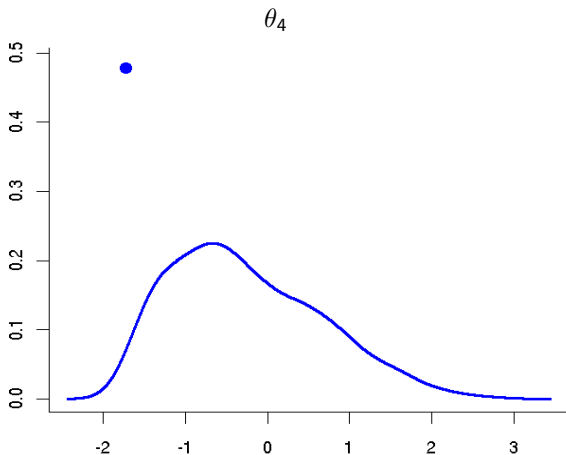
- Choose μ_n through cross-validation.



Simulations - remove orthogonality assumption

$k = 4$, $n = 200$, $\theta = (3, 1.5, 0, 0)' + 2/n^{1/2}(0, 0, 1, 1)'$, $X'X = n\Omega$ with $\Omega_{ij} = 0.5^{|i-j|}$, 1000 simulations

- Choose μ_n through cross-validation.



Estimation of the cdf of $n^{1/2}(\hat{\theta}_{AL} - \theta)$?

Let $F_{n,\theta}$ be the distribution function of $n^{1/2}(\hat{\theta}_{AL} - \theta)$.

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow m$ with $0 \leq m \leq \infty$. Then every consistent estimator $\hat{F}_n(t)$ of $F_{n,\theta}(t)$ satisfies

$$\lim_{n \rightarrow \infty} \sup_{|\theta| < c/n^{1/2}} P_{n,\theta} \left(\left| \hat{F}_n(t) - F_{n,\theta}(t) \right| > \varepsilon \right) \geq \frac{1}{2}$$

for each $\varepsilon < (\Phi(t+m) - \Phi(t-m))/2$ and each $c > |t|$.

In particular, not uniformly consistent estimator for $F_{n,\theta}(t)$ exists!

Analogous result for cdf under μ^{-1} -scaling.

Proof rests on Pötscher & Leeb (2006).

Estimation of the cdf of $n^{1/2}(\hat{\theta}_{AL} - \theta)$?

Finite-sample result:

Let $\mu_n \rightarrow 0$ and $n^{1/2}\mu_n \rightarrow m$ with $0 \leq m \leq \infty$. Then every estimator $\hat{F}_n(t)$ of $F_{n,\theta}(t)$ satisfies

$$\sup_{|\theta| < c/n^{1/2}} P_{n,\theta} \left(\left| \hat{F}_n(t) - F_{n,\theta}(t) \right| > \varepsilon \right) \geq \frac{1}{2}$$

for each $\varepsilon < (\Phi(t+m) - \Phi(t-m))/2$, for each $c > |t|$ and each sample size n . Hence

$$\liminf_{n \rightarrow \infty} \inf_{\hat{F}_n(t)} \sup_{|\theta| < c/n^{1/2}} P_{n,\theta} \left(\left| \hat{F}_n(t) - F_{n,\theta}(t) \right| > \varepsilon \right) = 1$$

for each $\varepsilon < (\Phi(t+m) - \Phi(t-m))/2$ and each $c > |t|$ where the infimum extend over *all* estimators $\hat{F}_n(t)$.

Conclusions

- The **finite-sample distribution** of the adaptive LASSO estimator and other PLSEs are **highly non-normal**.
- **Non-normality persists in large samples**. This can be seen through a “moving-parameter” asymptotic framework.
- Fixed-parameter asymptotics (as underlying the oracle-property) paint a misleading picture of the performance of the estimator due to the **non-uniformity** of these results.
- **Confidence intervals** in the consistent case **are larger by one order of magnitude** compared to the unpenalized estimator.
- The distribution function of the adaptive LASSO estimator and other PLSEs **cannot** be estimated in a **uniformly consistent** manner.
- **NOT a criticism** on PLSEs per se, but relying on fixed-parameter asymptotics in this context is dangerous.

Conclusions

- The **finite-sample distribution** of the adaptive LASSO estimator and other PLSEs are **highly non-normal**.
- **Non-normality persists in large samples**. This can be seen through a “moving-parameter” asymptotic framework.
- Fixed-parameter asymptotics (as underlying the oracle-property) paint a misleading picture of the performance of the estimator due to the **non-uniformity** of these results.
- **Confidence intervals** in the consistent case **are larger by one order of magnitude** compared to the unpenalized estimator.
- The distribution function of the adaptive LASSO estimator and other PLSEs **cannot** be estimated in a **uniformly consistent** manner.
- **NOT a criticism** on PLSEs per se, but relying on fixed-parameter asymptotics in this context is dangerous.

Conclusions

- The **finite-sample distribution** of the adaptive LASSO estimator and other PLSEs are **highly non-normal**.
- **Non-normality persists in large samples**. This can be seen through a “moving-parameter” asymptotic framework.
- Fixed-parameter asymptotics (as underlying the oracle-property) paint a misleading picture of the performance of the estimator due to the **non-uniformity** of these results.
- **Confidence intervals** in the consistent case **are larger by one order of magnitude** compared to the unpenalized estimator.
- The distribution function of the adaptive LASSO estimator and other PLSEs **cannot** be estimated in a **uniformly consistent** manner.
- **NOT a criticism** on PLSEs per se, but relying on fixed-parameter asymptotics in this context is dangerous.

Conclusions

- The **finite-sample distribution** of the adaptive LASSO estimator and other PLSEs are **highly non-normal**.
- **Non-normality persists in large samples**. This can be seen through a “moving-parameter” asymptotic framework.
- Fixed-parameter asymptotics (as underlying the oracle-property) paint a misleading picture of the performance of the estimator due to the **non-uniformity** of these results.
- **Confidence intervals** in the consistent case are **larger by one order of magnitude** compared to the unpenalized estimator.
- The distribution function of the adaptive LASSO estimator and other PLSEs **cannot** be estimated in a **uniformly consistent** manner.
- **NOT a criticism** on PLSEs per se, but relying on fixed-parameter asymptotics in this context is dangerous.

Conclusions

- The **finite-sample distribution** of the adaptive LASSO estimator and other PLSEs are **highly non-normal**.
- **Non-normality persists in large samples**. This can be seen through a “moving-parameter” asymptotic framework.
- Fixed-parameter asymptotics (as underlying the oracle-property) paint a misleading picture of the performance of the estimator due to the **non-uniformity** of these results.
- **Confidence intervals** in the consistent case **are larger by one order of magnitude** compared to the unpenalized estimator.
- The distribution function of the adaptive LASSO estimator and other PLSEs **cannot** be estimated in a **uniformly consistent** manner.
- **NOT a criticism** on PLSEs per se, but relying on fixed-parameter asymptotics in this context is dangerous.

Conclusions

- The **finite-sample distribution** of the adaptive LASSO estimator and other PLSEs are **highly non-normal**.
- **Non-normality persists in large samples**. This can be seen through a “moving-parameter” asymptotic framework.
- Fixed-parameter asymptotics (as underlying the oracle-property) paint a misleading picture of the performance of the estimator due to the **non-uniformity** of these results.
- **Confidence intervals** in the consistent case **are larger by one order of magnitude** compared to the unpenalized estimator.
- The distribution function of the adaptive LASSO estimator and other PLSEs **cannot** be estimated in a **uniformly consistent** manner.
- **NOT a criticism** on PLSEs per se, but relying on fixed-parameter asymptotics in this context is dangerous.

Conclusions

- The **finite-sample distribution** of the adaptive LASSO estimator and other PLSEs are **highly non-normal**.
- **Non-normality persists in large samples**. This can be seen through a “moving-parameter” asymptotic framework.
- Fixed-parameter asymptotics (as underlying the oracle-property) paint a misleading picture of the performance of the estimator due to the **non-uniformity** of these results.
- **Confidence intervals** in the consistent case **are larger by one order of magnitude** compared to the unpenalized estimator.
- The distribution function of the adaptive LASSO estimator and other PLSEs **cannot** be estimated in a **uniformly consistent** manner.
- **NOT a criticism** on PLSEs per se, but relying on fixed-parameter asymptotics in this context is dangerous.

References

-  J. Fan and R. Li. [Variable selection via nonconcave penalized likelihood and its oracle properties](#). *J. Am. Stat. Ass.*, 96:1348–1360, 2001.
-  I. E. Frank and J. H. Friedman. [A statistical view of some chemometrics regression tools \(with discussion\)](#). *Technom.*, 35:109–148, 1993.
-  A.E. Hoerl and R. Kennard. [Ridge regression: Biased estimation for nonorthogonal problems](#). *Technometrics* 12, 55–67.
-  K. Knight and W. Fu. [Asymptotics of lasso-type estimators](#). *Ann. Stat.*, 28:1356–1378, 2000.
-  H. Leeb and B. M. Pötscher. [Performance limits for estimators of the risk or distribution of shrinkage-type estimators, and some general lower risk-bound results](#). *Economet. Theory*, 22:69–97, 2006. [Corrections](#). *Ibidem*, 24:581–583, 2008.
-  B. M. Pötscher. [Confidence sets based on sparse estimators are necessarily large](#). *Manuscript*, 2007. arXiv:0711.1036.
-  B. M. Pötscher and H. Leeb. [On the distribution of penalized maximum likelihood estimators: The LASSO, SCAD, and thresholding](#). *Manuscript*, 2007. arXiv:0711.0660.
-  B. M. Pötscher and U. Schneider. [Confidence sets based on penalized maximum likelihood estimators](#). *Manuscript*, 2008. arXiv:0806.1652.
-  B. M. Pötscher and U. Schneider. [On the distribution of the adaptive lasso estimator](#). *J. Stat. Plan. Inf.*, to appear.
-  R. Tibshirani. [Regression shrinkage and selection via the lasso](#). *J. Roy. Stat. Soc. B*, 58:267–288, 1996.
-  H. Zou. [The adaptive lasso and its oracle properties](#). *J. Am. Stat. Ass.*, 101:1418–1429, 2006.