



Adaptive Sparse Estimation With Side Information

Trambak Banerjee, Gourab Mukherjee, and Wenguang Sun

Department of Data Sciences and Operations, University of Southern California, Los Angeles, CA

ABSTRACT

The article considers the problem of estimating a high-dimensional sparse parameter in the presence of side information that encodes the sparsity structure. We develop a general framework that involves first using an auxiliary sequence to capture the side information, and then incorporating the auxiliary sequence in inference to reduce the estimation risk. The proposed method, which carries out adaptive Stein's unbiased risk estimate-thresholding using side information (ASUS), is shown to have robust performance and enjoy optimality properties. We develop new theories to characterize regimes in which ASUS far outperforms competitive shrinkage estimators, and establish precise conditions under which ASUS is asymptotically optimal. Simulation studies are conducted to show that ASUS substantially improves the performance of existing methods in many settings. The methodology is applied for analysis of data from single cell virology studies and microarray time course experiments. Supplementary materials for this article are available online.

ARTICLE HISTORY

Received November 2018
Accepted October 2019

KEYWORDS

Adaptive shrinkage estimation; Higher order minimax risk; Inference with side information; Sparsity; SURE shrinkage; Two-sample inference

1. Introduction

The recent technological advancements have made it possible to collect vast amounts of data with various types of side information such as domain knowledge, expert insights, covariates in the primary data, and secondary data from related studies. In a wide range of fields including genomics, neuroimaging, and signal processing, incorporating side information promises to yield more accurate and meaningful results. However, few analytical tools are available for extracting and combining information from different data sources in high-dimensional data analysis. This article aims to develop new theory and methodology for leveraging side information to improve the efficiency in estimating a high-dimensional sparse parameter. We study the following closely related issues: (i) how to properly extract or construct an auxiliary sequence to capture useful sparsity information; (ii) how to combine the auxiliary sequence with the primary summary statistics to develop more efficient estimators; and (iii) how to assess the relevance and usefulness of the side information, as well as the robustness and optimality of the proposed method.

1.1. Motivating Applications


Sparsity is an essential phenomenon that arises frequently in modern scientific studies. In a range of data-intensive application fields such as genomics and neuroimaging, only a small fraction of data contain useful signals. The detection, estimation, and testing of a high-dimensional sparse object have many important applications and have been extensively studied in the literature (Donoho and Jin 2004; Johnstone and Silverman 2004; Abramovich et al. 2006). For instance, in the RNA-seq study

that will be analyzed in Section 4.3, the goal is to estimate the true expression levels of $n = 53,216$ genes for the virus strain VZV, which is the causative agent of varicella (chickenpox) and zoster (shingles) in humans (Zerboni et al. 2014). The parameter of interest (the population mean vector of gene expression) is sparse as it is known that very few genes in the generic RNA-seq kits express themselves in these single-cell virology studies (Sen et al. 2018). The accurate identification and estimation of nonzero large effects is helpful for the discovery of novel genetic biomarkers, which constitutes a key step in the development of new treatments and personalized medicine (Erickson and Sabatti 2005; Matsui 2013; Holland et al. 2016). Another example arises from microarray time-course (MTC) experiments that will be analyzed in Section E of the supplementary materials. The goal is to identify genes that exhibit a specific pattern of differential expression over time. The temporal pattern, which can be revealed by estimating the differences in expression levels of genes between two time points, would help gain insights into the mechanisms of the underlying biological processes (Calvano et al. 2005; Sun and Wei 2011). After baseline removal, the parameter of interest is the difference between two mean vectors that are both individually sparse.

In practice, the intrinsic sparsity structure of the high-dimensional parameter is often captured by side information, which can be obtained as either summary statistics from secondary data sources or can be constructed as a covariate sequence from the original data. For instance, in the RNA-seq data, expression levels corresponding to other four experimental conditions (C1, C2, C3, and C4) are also available for the same n genes through related studies conducted in the lab. The heat map in Figure 1 shows that the sparse structure of the mean

CONTACT Wenguang Sun  wenguan@marshall.usc.edu  Department of Data Sciences and Operations, 401W Bridge Hall, Marshall School of Business, University of Southern California, Los Angeles, CA 90089.

Color versions of one or more of the figures in the article can be found online at www.tandfonline.com/r/JASA.

 Supplementary materials for this article are available online. Please go to www.tandfonline.com/r/JASA.

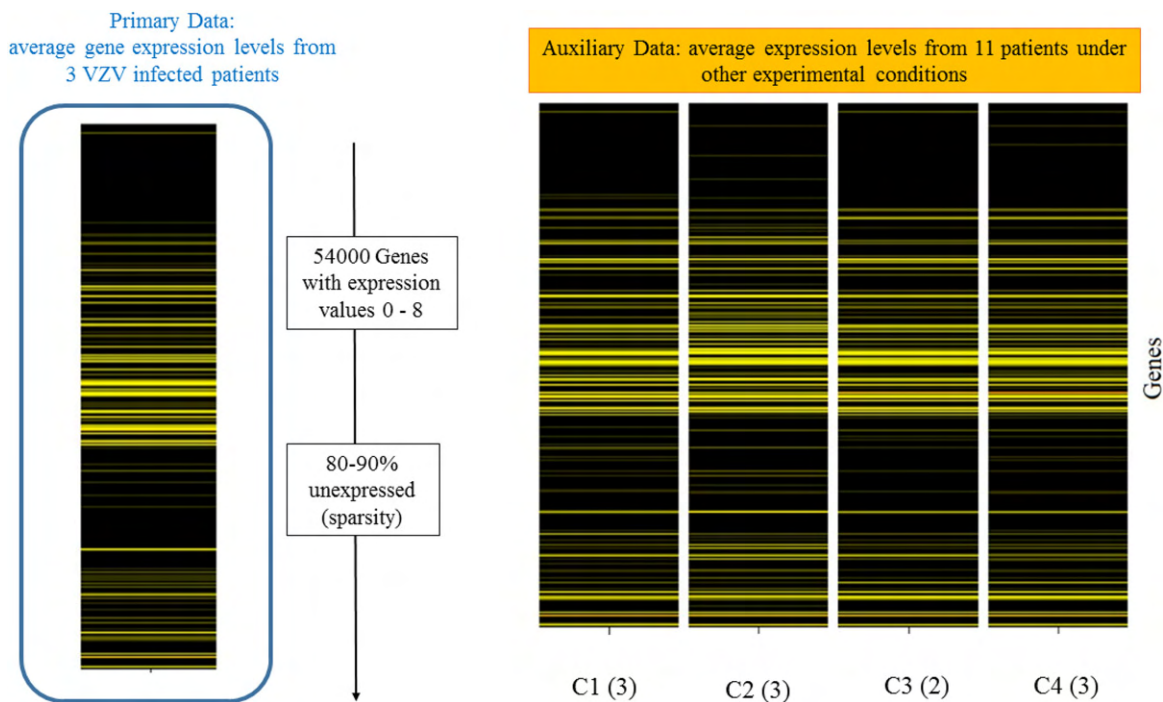


Figure 1. Heat map showing the average expression levels in the RNA-seq study. Left panel: VZV; right panel from top to bottom: C1, C2, C3, and C4, where the number of replicates (patients) is shown in parenthesis. We can see that 80–90% of the genes under the VZV condition are unexpressed (black), and the same sparse structure seems to be roughly maintained in the other four experimental conditions. Useful side information on sparsity can be extracted from secondary data (C1–C4) and be combined with the primary data (VZV) to construct more efficient estimators.

transcription levels of the genes for VZV is roughly maintained by the same set of genes in subjects from the other four conditions. The common structural information shared by both cases (VZV) and controls (C1–C4) can be exploited to construct more efficient estimation procedures. In the two-sample sparse estimation problem considered in the MTC study (analyzed in Section E of the supplementary materials), we illustrate that a covariate sequence can be constructed from the original data matrix to assist inference by capturing the sparseness of the mean difference. Intuitively, incorporating side information promises to improve the efficiency of existing methods and interpretability of results. However, in conventional practice, such useful auxiliary data have been largely ignored in analysis.

1.2. ASUS: A General Framework for Leveraging Side Information

In this article, we develop a general integrative framework for sparse estimation that is capable of handling side information that may be extracted from (i) prior or domain-specific knowledge, (ii) covariate sequence based on the same (original) data, or (iii) summary statistics based on secondary data sources. Let $\theta = (\theta_1, \dots, \theta_n)$ be an unknown high-dimensional sparse parameter. Our study focuses on the class of nonlinear thresholding estimators (see Mallat 2008, chap. 11; Johnstone 2015, chaps. 8 and 13), which have been widely used in the sparse case where many coordinates of θ are small or zero.

The proposed estimation framework involves two steps: first constructing an auxiliary sequence $\mathbf{S} = (S_i : 1 \leq i \leq n)$ to capture the sparse structure, and second combining \mathbf{S} with the primary statistics, denoted $\mathbf{Y} = (Y_i : 1 \leq i \leq n)$,

via a group-wise adaptive thresholding algorithm. Our idea is that the coordinates of θ become nonexchangeable in light of side information. To reflect this heterogeneity, we divide all coordinates into K groups based on S_i . The side information is then incorporated in our estimation procedure by applying soft-thresholding estimators separately, thereby fine tuning the group-wise thresholds to capture the varied sparsity levels across groups. The optimal grouping and thresholds are chosen adaptively via a data-driven approach, which employs the Stein’s unbiased risk estimate (SURE) criterion to minimize the total estimation risk. The proposed method, which carries out adaptive SURE-thresholding using side information (ASUS), is shown to have robust performance and enjoy optimality properties. ASUS is simple and intuitive, but nevertheless provides a general framework for information pooling in sparse estimation problems. Concretely, since ASUS does not rely on any functional relationships between \mathbf{S} and θ , it is robust and effective in leveraging side information in a wide range of scenarios. In Section 2.2, we demonstrate that this flexible framework can be applied to various sparse estimation problems.

The amount of efficiency gain of ASUS depends on two factors: (i) the usefulness of the side information and (ii) the effectiveness in utilizing the side information. To understand the first issue, we formulate in Section 3 a hierarchical model to assess the informativeness of an auxiliary sequence. Our theoretical analysis characterizes the conditions under which methods ignoring side information are suboptimal compared to an “oracle” with perfect knowledge on sparsity structure. To investigate the second issue, Section 3 establishes precise conditions under which ASUS is asymptotically optimal, in the sense that its maximal risk is close to the theoretical limit that is

attained by the oracle. Finally, we carry out a theoretical analysis on the robustness of ASUS; our results show that pooling non-informative side information would not harm the performance of data combination procedures. Our asymptotic results are built upon the elegant higher order minimax risk evaluations developed by Johnstone (1994).

1.3. Connections With Existing Work and Our Contributions

ASUS is a nonlinear shrinkage estimator that incorporates relevant side information by choosing data-adaptive thresholds to reflect the varied sparsity levels across groups. We use the SURE criterion for simultaneous tuning of the grouping and shrinkage parameters. Our methodology is related to Xie, Kou, and Brown (2012), Tan (2015), and Weinstein et al. (2018), which utilized SURE to devise algorithms reflecting optimal shrinkage directions. However, these works are developed for different purposes (addressing the heteroscedasticity issue in the data) and do not cover the sparse case.

The notion of side information in estimation has been explored in several research fields. In information theory for instance, sparse source coding with side information is a well-studied problem (Wyner 1975; Cover and Thomas 2012; Watanabe, Kuzuoka, and Tan 2015). However, these methodologies focus on very different goals and cannot be directly applied to our problem. In the statistical literature, the use of side information in sparse estimation problems has been mainly limited to regression settings where the side information must be in the form of a linear function of θ (Ke, Jin, and Fan 2014; Kou and Yang 2015). By contrast, our estimation framework utilizes a more flexible scheme that does not require the specification of any functional relationship between θ and the side information. The proposed ASUS algorithm is simple and intuitive but nevertheless enjoys strong numerical and theoretical properties. Our simulation studies show that it can substantially outperform competitive methods in many settings. ASUS is a robust data combination procedure in the sense that asymptotically it would not under-perform methods ignoring side information when the auxiliary data are non-informative (see Theorem 4).

The proposed research makes several new theoretical contributions. First, we develop general principles for constructing and pooling the side information, which guarantees proper information extraction and robust performance of ASUS. Second, we formulate a theoretical framework to assess the usefulness of side information. Third, we establish precise conditions under which ASUS is asymptotically optimal. Finally, we extend the sparse minimax decision theory of Johnstone (2015), which provides the foundation for a range of sparse inference problems (Abramovich et al. 2006; Abramovich, Grinshtein, and Pensky 2007; Cai, Low, and Ma 2014; Tibshirani 2014; Collier, Comminges, and Tsybakov 2017), to derive new high-order characterizations of the maximal risk of soft-thresholding estimators.

1.4. Organization of the Article

Section 2 describes the proposed ASUS procedure. Section 3 presents theoretical analyses. The numerical performances of

ASUS are investigated using both simulated and real data in Section 4. Section 5 concludes with a discussion. Additional numerical results and proofs are given in the supplementary materials.

2. Adaptive Sparse Estimation With Side Information

This section first describes the model and assumptions (Section 2.1), then discusses how to construct the auxiliary sequence (Section 2.2), and finally proposes the methodology (Section 2.3).

2.1. Model and Assumptions

To conduct a systematic study of the influence of side information for estimating θ , we consider a hierarchical model that relates the primary and auxiliary datasets through a latent vector $\xi = (\xi_1, \dots, \xi_n)$, which represents the noiseless side information that encodes the sparsity information of θ . The latent vector ξ cannot be observed directly but may be partially revealed by an auxiliary sequence (noisy side information) $S = (S_1, \dots, S_n)$. For instance, in the RNA-seq example, the parameter of interest is the population mean of the gene expression levels for diseased patients, and the latent variable ξ_i may represent the quantitative outcome of a complex gene regulation process that determines whether gene i expresses itself under the influence of a certain experimental condition. The primary and secondary data, respectively, correspond to gene expression levels for the patients from the concerned (i.e., VZV infected) and other related groups. The primary and auxiliary statistics Y_i and S_i for gene i can be constructed based on the corresponding sample means.

For n parallel units, the summary statistic Y_i for the i th unit is modeled by

$$Y_i = \theta_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma_i^2), \tag{1}$$

where, by convention, σ_i^2 are assumed to be known or can be well estimated from data (e.g., Brown and Greenshtein 2009; Xie, Kou, and Brown 2012; Weinstein et al. 2018). We further assume that both θ and S are related to the latent vector ξ through some unknown real-valued functions h_θ and h_s :

$$\theta_i = h_\theta(\xi_i, \eta_{1i}), \tag{2}$$

$$S_i = h_s(\xi_i, \eta_{2i}), \tag{3}$$

where η_{1i} and η_{2i} follow some unspecified priors, and represent independent random perturbations that are independent of ξ_i ; concrete examples for Models (1)–(3) are discussed in Section 2.2.

Remark 1. The above model can be conceptualized as a Bayesian hierarchical model:

$$Y_i | (\theta_i, S_i) \sim N(\theta_i, \sigma_i^2), \quad (\theta_i, S_i) | \xi_i \sim f_1(\theta | \xi_i) f_2(S | \xi_i), \\ \xi_i \stackrel{iid}{\sim} f_3(\xi),$$

where f_1, f_2, f_3 are unknown densities. In Equations (2) and (3), ξ_i is a random quantity and independent of η_{1i} and η_{2i} . As a special case of Equation (2), we can write $\theta_i = h_\theta(\xi_i)$ without the

random perturbations η_{1i} . Our theory is mainly stated in terms of random ξ_i 's for ease of presentation. However, we note that our theoretical results still hold even when ξ_i is deterministic because the theory in Section 2.3 is derived conditional on ξ_i , and the proof in Section 3 is built upon an empirical density function (10).

The hierarchical models (1)–(3) provide a general and flexible framework for our methodological and theoretical developments. In particular, it covers a wide range of scenarios by allowing the strength of the side information to vary from completely non-informative (e.g., when ξ_i is useless, or when S_i and ξ_i are independent for all i) to perfectly informative (e.g., when $\theta_i = \xi_i$ and $S_i = \xi_i$ for all i). In Section 3, the usefulness of the latent vector ξ is investigated via Equation (2), and the informativeness of the auxiliary sequence \mathbf{S} is characterized by Equations (2) and (3).

2.2. Constructing the Auxiliary Sequence: Principles and Examples

A key step in our methodological development is to properly extract side information using an auxiliary sequence. The sequence \mathbf{S} can be constructed from various data sources including the following three basic settings: (i) prior or domain-specific knowledge; (ii) covariates or discard data in the same primary dataset; or (iii) secondary data from related studies. We stress that our estimation framework is valid for all three settings as long as \mathbf{S} fulfills the following two fundamental principles.

The first principle is *informativeness*, which requires that S_i should be chosen or constructed in a way to encode the sparse structure effectively. The second principle is *conditional independence*, which requires that S_i must be conditionally independent of Y_i given the latent variable ξ_i . The conditional independence assumption, which is implied by Models (1)–(3), ensures proper shrinkage directions and plays a key role in establishing the robustness of ASUS. Examples 1–4 present specific instances of auxiliary sequences fulfilling such principles, wherein the auxiliary sequences may either be readily available from distinct but related experiments or can be carefully constructed from the same (original) data to capture important structural information that is discarded by conventional practice.

Example 1 (Prioritized subset analysis (PSA, Li et al. 2008)). In genome wide association studies, prior data and domain knowledge such as known gene functions or interactions may be used to construct an auxiliary sequence \mathbf{S} that can prioritize the discovery of SNPs in certain genomic regions. Typically, the primary dataset can be summarized as a vector $\mathbf{Y} = (Y_1, \dots, Y_n)$, where Y_i are either taken as differential allele frequencies between diseased and control groups, or z -values based on χ^2 -tests assessing the association between the allele frequency and the disease status. Let $\mathbf{S} = (S_1, \dots, S_n) \in \{-1, 1\}^n$ be an auxiliary sequence, where $S_i = 1$ if SNP i is in the prioritized subset and $S_i = -1$ otherwise. \mathbf{S} can be viewed as perturbations of the true state sequence $\xi = (\xi_1, \dots, \xi_n)$, where $\xi_i = 1$ if SNP i is associated with the disease and $\xi_i = -1$ otherwise. The informativeness and independence

principles are fulfilled when (i) the prioritized subset contains SNPs that are more likely to hold disease susceptible variants and (ii) the perturbations of ξ are random (hence, Y_i and S_i are conditionally independent given ξ_i). Both (i) and (ii) seem reasonable assumptions in PSA studies.

Example 2 (One-sample inference). In the RNA-seq study, let the primary data be $\{Y_{ij} : i = 1, \dots, n; j = 1, \dots, k_y\}$ that record the expression levels of n genes from k_y subjects infected by VZV. The primary statistics are $\mathbf{Y} = (\bar{Y}_1, \dots, \bar{Y}_n)$, where $\bar{Y}_i = k_y^{-1} \sum_{j=1}^{k_y} Y_{ij}$. Let the secondary data be $\{X_{ij} : i = 1, \dots, n; j = 1, \dots, k_x\}$ that record the expression levels of the same n genes for k_x subjects but under different conditions C1–C4. The auxiliary sequence can be constructed as $\mathbf{S} = (S_1, \dots, S_n) = (|\bar{X}_1|, \dots, |\bar{X}_n|)$, where $\bar{X}_i = k_x^{-1} \sum_{j=1}^{k_x} X_{ij}$. Thus although we record the expression levels of the same set of n genes, in the case of the primary data the genes are infected with the VZV virus whereas for the secondary data the expression levels are recorded under the influence of agents that are different from that of the VZV virus. The latent state ξ_i represents whether gene i expresses itself under any of the conditions. Now we check whether the two information extraction principles are fulfilled. First, the informativeness principle holds since, as demonstrated by the heat map in Figure 1, inactive genes under VZV are likely to remain inactive under the other conditions. The sparse structure is captured by the auxiliary sequence, where a small S_i signifies an inactive gene. Second, Section 2.1 has explained how the RNA-seq data may be sensibly conceptualized via Models (1)–(3), where \bar{Y}_i and S_i are conditionally independent given the latent variable ξ_i , fulfilling the second principle.

Example 3 (Two-sample inference). Consider the MTC study discussed in the introduction (and analyzed in Section E of the supplementary materials). Let $\{Y_{ij,t_d} : i = 1, \dots, n; j = 1, \dots, k_i; d = 0, 1, 2\}$ record the expression levels of n genes from k_i subjects at time points t_0 (baseline), t_1 and t_2 . Let $\bar{Y}_{i,d} = k_i^{-1} \sum_{j=1}^{k_i} (Y_{ij,t_d} - Y_{ij,t_0})$ be the average expression levels of gene i at time point t_d after baseline adjustment, $d = 1, 2$. Denote $\mu_{i,d} = E(\bar{Y}_{i,d})$ and $\boldsymbol{\mu}_d = (\mu_{i,d} : 1 \leq i \leq n)$. Then both $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$ are individually sparse. The parameter of interest is $\theta_i = \mu_{i,1} - \mu_{i,2}$, which can be estimated by the primary statistic $Y_i = \bar{Y}_{i,1} - \bar{Y}_{i,2}$. Denote the union support $\mathcal{U} = \{i : \mu_{i,1} \neq 0 \text{ or } \mu_{i,2} \neq 0\}$. Then \mathcal{U} can be exploited to screen out zero effects since if $i \notin \mathcal{U}$, we must have $\theta_i = 0$. Consider the sequence $S_i = |\bar{Y}_{i,1} + \kappa_i \bar{Y}_{i,2}|$, where $\kappa_i = \hat{\sigma}_{i,1}/\hat{\sigma}_{i,2}$ and $\hat{\sigma}_{i,d}^2 = (k_i - 1)^{-1} \sum_{j=1}^{k_i} (Y_{ij,t_d} - Y_{ij,t_0} - \bar{Y}_{i,d})^2$. Then the auxiliary sequence is informative since a large S_i provides strong evidence that $i \in \mathcal{U}$. The union support encodes the sparse structure of $\boldsymbol{\theta}$. Moreover, Y_i and S_i are asymptotically independent with our choice of κ_i (Cai, Sun, and Wang 2018, Proposition 6). Hence, both principles are fulfilled.

Example 4 (Estimation under the ANOVA setting). This example is an extension of Example 3 to multi-sample inference. Consider m conditions $d = 1, \dots, m, m \geq 2$. The parameter of interest is $\boldsymbol{\theta}_{n \times 1} = \Gamma \mathbf{a}$, where $\Gamma_{n \times m} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_m)$, $\mu_{i,d} = \mathbb{E}(\bar{Y}_{i,d})$, and $\mathbf{a}_{m \times 1}$ is a vector of known weights. Here, $\boldsymbol{\theta}$ may represent a weighted average of true transcription levels of n genes across

m time points. Let $\mathbf{D}_i = (\bar{Y}_{i,1}, \dots, \bar{Y}_{i,m})$ be the vector of average expression level of gene i for the m time points after baseline adjustment and denote $\mathbb{D}_{n \times m} = (\mathbf{D}_1, \dots, \mathbf{D}_n)^T$. To estimate θ , our proposed framework suggests using the usual unbiased estimator $\mathbf{Y} = \mathbb{D}\mathbf{a}$ as the primary statistic, and $\mathbf{S} = \mathbb{D}\mathbf{b}$ as the auxiliary sequence for some weights \mathbf{b} . The informativeness principle from Example 3 continues to hold under this setting. To fulfil the independence principle, we choose \mathbf{b} such that $\text{cov}(\mathbf{Y}, \mathbf{S}) = \mathbf{0}$.

In Examples 3 and 4, the auxiliary sequence \mathbf{S} is constructed from the same original data matrix. We give some intuitions to explain why \mathbf{S} is useful. The conventional practice reduces the original data into a vector of summary statistics \mathbf{Y} . However, this data reduction step often causes significant loss of information and thus leads to suboptimal procedures. Specifically, the information on the sparseness of the union support \mathcal{U} is lost in the data reduction step. The key idea in Example 3 is that the auxiliary sequence \mathbf{S} captures the structural information on sparsity, which is discarded by conventional practice. Therefore by incorporating \mathbf{S} into the inferential process we can improve the efficiency of existing methods. Note that \mathbf{Y} is not a sufficient statistic for estimating θ , the minimax estimation error based on (\mathbf{Y}, \mathbf{S}) can greatly improve the performance of all estimators that are based on \mathbf{Y} alone; a rigorous theoretical analysis is carried out in the proof of Theorem 2. To summarize, the above examples illustrate that the side information can be either “external” (Examples 1 and 2) or “internal” (Examples 3 and 4). The key in the proposed estimation framework, which we discuss next, is to construct a proper auxiliary sequence that fulfills the two fundamental principles. We shall develop a unified estimation framework that is capable of handling both internal and external side information.

We conclude this section with two remarks. First, the conditional independence assumption can be relaxed; the methodology would work as long as Y_i and S_i are conditionally *uncorrelated* (cf. Proposition 1). Second, we do not require Y_i or θ_i to be related to S_i through any functional forms; hence, classical regression techniques (even nonparametric models) cannot be applied in the above scenarios. We aim to develop a general information pooling strategy that does not involve any prescribed functional relationships; a methodology in this spirit is described next.

2.3. The ASUS Estimator and Its Risk Properties

Let \mathbf{Y} and \mathbf{S} denote the primary statistics and auxiliary sequence obeying Models (1)–(3). Let $\eta_t(\cdot)$ be a soft-thresholding operator such that

$$\eta_t(Y_i) = \begin{cases} -Y_i\sigma_i^{-1}, & \text{if } |Y_i\sigma_i^{-1}| \leq t; \\ -t \text{ sign}(Y_i\sigma_i^{-1}), & \text{otherwise.} \end{cases}$$

The proposed ASUS estimator operates in two steps: first constructing K groups using \mathbf{S} , and second applying soft-thresholding within each group using \mathbf{Y} . The construction of the groups relies only on \mathbf{S} . The tuning parameters for both grouping and shrinkage are determined using the SURE criterion.

Procedure 1. For $k = 1, \dots, K$ and $\boldsymbol{\tau} = \{\tau_1 < \dots < \tau_{K-1}\}$, denote $\widehat{\mathcal{I}}_k^\tau = \{i : \tau_{k-1} < S_i \leq \tau_k\}$ with $\tau_0 = -\infty, \tau_K = \infty$. Consider the following class of shrinkage estimators:

$$\hat{\theta}_i^{\text{SI}}(\mathcal{T}) := Y_i + \sigma_i \eta_{t_k}(Y_i) \quad \text{if } i \in \widehat{\mathcal{I}}_k^\tau, \tag{4}$$

where $\mathcal{T} = \{\tau_1, \dots, \tau_{K-1}, t_1, \dots, t_K\}$ and each of the threshold hyper-parameters t_1, \dots, t_K varies in $[0, t_n]$ with $t_n = (2 \log n)^{1/2}$. Thus, the set of all possible hyper-parameter \mathcal{T} values is $\mathcal{H}_n = \mathbf{R}_+^{K-1} \times [0, t_n]^K$. Define the SURE function

$$S(\mathcal{T}, \mathbf{Y}, \mathbf{S}) = n^{-1} \left[\sum_{i=1}^n \sigma_i^2 + \sum_{k=1}^K \sum_{i \in \widehat{\mathcal{I}}_k^\tau} \left\{ \sigma_i^2 (|Y_i \sigma_i^{-1}| \wedge t_k)^2 - 2\sigma_i^2 I(|Y_i \sigma_i^{-1}| \leq t_k) \right\} \right]. \tag{5}$$

Let $\hat{\mathcal{T}} = \arg \min_{\mathcal{T} \in \mathcal{H}_n} S(\mathcal{T}, \mathbf{Y}, \mathbf{S})$. Then, the ASUS estimator is given by $\hat{\theta}_i^{\text{SI}}(\hat{\mathcal{T}})$.

Remark 2. When θ is very sparse, the empirical fluctuations in the SURE function would have nonnegligible effects on thresholding procedures. We suggest choosing t_1, \dots, t_k for a given grouping by implementing a hybrid scheme that is similar to the SureShrink estimator of Donoho and Johnstone (1995), for example, setting $t_k = t_n$ if $|\widehat{\mathcal{I}}_k^\tau|^{-1} \sum_{i \in \widehat{\mathcal{I}}_k^\tau} (Y_i^2 / \sigma_i^2) \wedge t_n^2 - 1 \leq n^{-1/2} \log^{3/2} n$.

We present a toy example to illustrate why ASUS works. Consider the two-sample inference problem described by Example 3 in Section 2.2. Let $\theta_i = \mu_{i,1} - \mu_{i,2}$ and $\bar{Y}_{i,d} \sim N(\mu_{i,d}, 0.25)$, where $d = 1, 2, i = 1, \dots, n$, and $n = 10^4$. For $\boldsymbol{\mu}_1$, we generate the first 20% of its coordinates randomly from Unif(4, 6), the next 20% randomly from Unif(2, 3), and set the remaining coordinates to 0. For $\boldsymbol{\mu}_2$, the first 20% are from Unif(1, 2), the next 20% from Unif(1, 6), and the remaining 0. Finally, we let $\bar{Y}_i = \bar{Y}_{i,1} - \bar{Y}_{i,2}$ and $S_i = |\bar{Y}_{i,1} + \bar{Y}_{i,2}|$. The left panel in Figure 2 presents the histogram of $\mathbf{Y} = (\bar{Y}_i : 1 \leq i \leq n)$, where the lighter shade corresponds to \bar{Y}_i with $\theta_i = 0$. The SureShrink estimator in Donoho and Johnstone (1995) chooses threshold $t = 0.6$ for all observations, resulting in an MSE of 0.338. Imagine that an oracle has the perfect knowledge about the two groups ($\theta_i = 0$ vs. $\theta_i \neq 0$). In group 0, SureShrink chooses $t_0 = 4.2$, whereas in group 1, SureShrink chooses $t_0 = 0.15$. The total MSE is reduced to 0.20 by adopting varied thresholds for the two groups. In practice, the groups cannot be identified perfectly but can be partially revealed by the auxiliary statistic $S_i = |\bar{Y}_{i,1} + \bar{Y}_{i,2}|$, where a small S_i signifies a possible zero effect. Our simulation studies in Section 4 show that by exploiting the side information in S_i , ASUS achieves substantial gain in performance over conventional methods.

Let $l_n(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}) = n^{-1} \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2^2$ denote the squared error loss of estimating $\boldsymbol{\theta}$ using $\hat{\boldsymbol{\theta}}$. For each member $\hat{\boldsymbol{\theta}}^{\text{SI}}(\mathcal{T})$ in our class of estimators, $\mathcal{T} \in \mathcal{H}_n$, denote its risk by $r_n(\mathcal{T}; \boldsymbol{\theta}) = \mathbb{E} \left[l_n \left\{ \boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{\text{SI}}(\mathcal{T}) \right\} \right]$, where the expectation is taken with respect to the joint distribution of (Y_i, S_i) . The next proposition shows that (5) provides an unbiased estimate of the true risk.

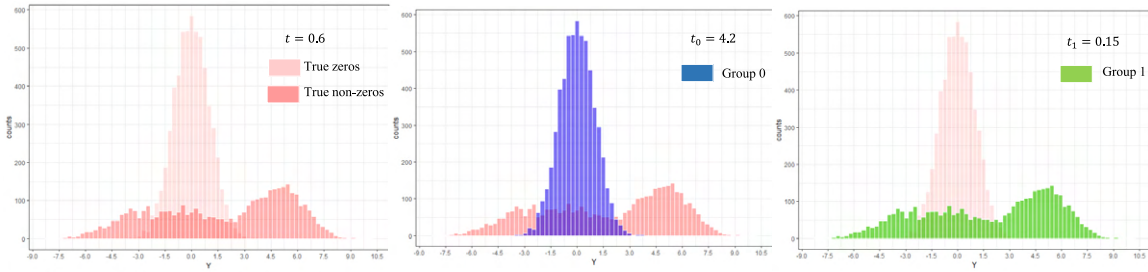


Figure 2. Toy example depicting ASUS. Left: SureShrink estimator at $t = 0.6$. Middle: ASUS with group 0 and $t_0 = 4.2$. Right: ASUS with group 1 and $t_1 = 0.15$.

Proposition 1. Consider Models (1)–(3). Then given ξ_i , the pair $\{(Y_i - \theta_i)\eta_{t_k}(Y_i), I(i \in \hat{\mathcal{I}}_k^t)\}$ are uncorrelated. It follows that $r_n(\mathcal{T}; \theta) = \mathbb{E}\{S(\mathcal{T}, \mathbf{Y}, \mathbf{S})\}$.

Next we study the large-sample behavior of the proposed SURE criterion. As in Xie, Kou, and Brown (2012), we impose the following assumption on the fourth moment of the noise distributions:

$$(A1) \quad \overline{\lim}_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \sigma_i^4 < \infty.$$

The following theorem shows that the risk estimate $S(\mathcal{T}, \mathbf{Y}, \mathbf{S})$ is uniformly close to the true risk as well as the loss, justifying our proposed hyper-parameter estimate $\hat{\mathcal{T}}$. Compared to Xie, Kou, and Brown (2012, Theorem 3.1) and Brown, Mukherjee, and Weinstein (2018, Theorem 4.1), we obtain explicit rates of convergence by tracking the empirical fluctuations in the SURE function through sharper concentration inequalities.

Theorem 1. Under Assumption A1, with $c_n = n^{1/2}(\log n)^{-\delta}$ for any $\delta > 3/2$, we have

- (a) $\lim_{n \rightarrow \infty} c_n \mathbb{E} \left\{ \sup_{\mathcal{T} \in \mathcal{H}_n} \left| S(\mathcal{T}, \mathbf{Y}, \mathbf{S}) - r_n(\mathcal{T}; \theta) \right| \right\} = 0,$
- (b) $\lim_{n \rightarrow \infty} c_n \mathbb{E} \left[\sup_{\mathcal{T} \in \mathcal{H}_n} \left| S(\mathcal{T}, \mathbf{Y}, \mathbf{S}) - l_n(\theta, \hat{\theta}^{\text{SI}}(\mathcal{T})) \right| \right] = 0,$

where the expectation is with respect to the joint distribution of \mathbf{Y}, \mathbf{S} .

Define \mathcal{T}^{OL} as the minimizer of the true loss function: $\mathcal{T}^{\text{OL}} = \arg \min_{\mathcal{T} \in \mathcal{H}_n} l_n(\theta, \hat{\theta}^{\text{SI}}(\mathcal{T}))$. \mathcal{T}^{OL} is referred to as the oracle loss hyper-parameter as it involves the knowledge of θ . It provides the theoretical limit that one can reach if allowed to minimize the true loss. Let $\hat{\theta}^{\text{SI}}(\mathcal{T}^{\text{OL}})$ be the corresponding oracle loss estimator. The following corollary establishes the asymptotic optimality of $\hat{\mathcal{T}}$.

Corollary 1. Under Assumption A1, if $\lim_{n \rightarrow \infty} c_n n^{-1/2} \log^\delta n = 0$ for any $\delta > 3/2$, then

- (a) The loss of $\hat{\theta}^{\text{SI}}(\hat{\mathcal{T}})$ converges in probability to the loss of $\hat{\theta}^{\text{SI}}(\mathcal{T}^{\text{OL}})$:

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[l_n \left\{ \theta, \hat{\theta}^{\text{SI}}(\hat{\mathcal{T}}) \right\} \geq l_n \left\{ \theta, \hat{\theta}^{\text{SI}}(\mathcal{T}^{\text{OL}}) \right\} + c_n^{-1} \epsilon \right] = 0 \text{ for any } \epsilon > 0.$$

- (b) The risk of $\hat{\theta}^{\text{SI}}(\hat{\mathcal{T}})$ converges to the risk of the oracle loss estimator:

$$\lim_{n \rightarrow \infty} c_n \mathbb{E} \left[l_n \left\{ \theta, \hat{\theta}^{\text{SI}}(\hat{\mathcal{T}}) \right\} - l_n \left\{ \theta, \hat{\theta}^{\text{SI}}(\mathcal{T}^{\text{OL}}) \right\} \right] = 0.$$

2.4. Approximating the Bayes Rule by ASUS

This section discusses a Bayes setup and illustrates how ASUS may be conceptualized as an approximation to the Bayes oracle estimator.

Consider a hierarchical model where θ_i has an unspecified prior and $Y_i \stackrel{\text{ind.}}{\sim} N(\theta_i, \sigma_i^2)$ with σ_i^2 known. In the absence of any auxiliary sequence \mathbf{S} and when σ_i are all equal to, say σ , the optimal estimator is

$$\delta_i^\pi = E(\theta_i | y_i) = y_i + \sigma^2 \frac{f'(y_i)}{f(y_i)}, \tag{6}$$

which is known as Tweedie’s formula (Efron 2011). When the marginal densities $f(y_i)$ are unknown, (6) can be implemented in an empirical Bayes (EB) framework. For example, Brown and Greenshtein (2009) used kernel methods to estimate unknown densities and showed that the resulting EB estimator is asymptotically optimal under mild conditions. Under the sparse setting, an effective approach to incorporate the sparsity structure is to consider, for example, spike-and-slab priors (Johnstone and Silverman 2004). In decision theory, it has been established that the posterior median is minimax optimal under spike-and-slab priors (see Johnstone and Silverman 2004, Theorem 1). Hence, the soft-threshold estimators can be viewed as good surrogates to the Bayes rule under sparsity. When the sparsity level is unknown, the threshold should be chosen adaptively using a data-driven method.

For a given pair of primary and auxiliary statistics (Y_i, S_i) , the Bayes oracle estimator is

$$\delta_i^\pi = E(\theta_i | Y_i, S_i). \tag{7}$$

Conditionally on S_i , a Tweedie’s formula for Equation (7) can be written which would require estimating the conditional marginal densities $f(y_i | s_i)$ and its derivatives. ASUS can be viewed as a two-step approximation to the oracle estimator (7). The first step involves using the auxiliary sequence to divide the n coordinates into K groups: $\delta_i^\pi \approx \hat{\delta}_k(Y_i) = E(\theta_i | Y_i, i \in G_k) = E(\theta_i | Y_i, S_i^* = k)$, which can be viewed as a discrete approximation to the oracle rule (7) by discretizing S_i as a categorical variable S_i^* taking values $k = 1, \dots, K$. The second step

involves setting thresholds for separate groups to incorporate the updated structural information from the auxiliary sequence. This step makes sense because under the sparse regime, it is natural to use the class of soft-thresholding estimators as a convenient surrogate to the Bayes rule, and ideally the threshold should be set differently to reflect the varied sparsity levels across the groups. Finally, the optimal grouping and optimal thresholds are chosen by minimizing a SURE criterion.

This Bayesian interpretation reveals that ASUS may suffer from information loss in the discretization step. However, fully utilizing the auxiliary data by modeling \mathbf{S} as a continuous variable is practically impossible under the ASUS framework since the search algorithm cannot deal with a diverging number of groups. Moreover, directly implementing (7) using bivariate Tweedie approaches is highly nontrivial and requires further research. ASUS, thus, seems to provide a simple, feasible yet effective framework to incorporate the side information.

3. Theoretical Analysis

This section studies the theoretical properties of ASUS under the important setting where θ is sparse. By contrast, the results of Section 2.3 hold for any sequence θ . To simplify the presentation, we focus on a class of thresholding estimators that utilize two groups. The two-group model provides a natural choice for some important applications such as the PSA and RNA-seq study, but the proposed ASUS framework can handle more groups. The major goal of our theoretical analysis is to gain insights on sparse inference with side information, for which the simple two-group setup helps in two ways. First, it leads to a concise and intuitive characterization of the potential influence of side information on simultaneous estimation. Second, it enables us to develop precise conditions under which ASUS is asymptotically optimal.

3.1. Asymptotic Setup

Consider hierarchical models (1)–(3). We begin by considering an oracle estimator $\tilde{\theta}_n^{\text{SI}}(\mathcal{T}_n^{\text{OR}})$ that directly uses the noiseless side information ξ :

$$\tilde{\theta}_{i,n}^{\text{SI}}(\mathcal{T}_n^{\text{OR}}) := \begin{cases} Y_i + \sigma_i \eta_{t_1^*}(Y_i) & \text{if } i \in \mathcal{I}_{1,n}^{\tau^*}, \\ Y_i + \sigma_i \eta_{t_2^*}(Y_i) & \text{if } i \in \mathcal{I}_{2,n}^{\tau^*}, \end{cases} \quad (8)$$

where $\mathcal{I}_{1,n}^{\tau} = \{i : \xi_i \leq \tau\}$, $\mathcal{I}_{2,n}^{\tau} = \{i : \xi_i > \tau\}$, and

$$\mathcal{T}_n^{\text{OR}} := (\tau_n^*, t_{1,n}^*, t_{2,n}^*) = \arg \min_{\mathcal{T} \in \mathbb{R} \times [0, t_n] \times [0, t_n]} \mathbb{E} l_n \left\{ \theta, \tilde{\theta}_n^{\text{SI}}(\mathcal{T}) \right\}. \quad (9)$$

Remark 3. Both the oracle estimator $\tilde{\theta}_n^{\text{SI}}(\mathcal{T}_n^{\text{OR}})$ and the oracle loss estimator $\hat{\theta}_n^{\text{SI}}(\mathcal{T}^{\text{OL}})$ assume the knowledge of θ . However, they are different in that the former creates groups based on ξ , whereas the latter uses \mathbf{S} . The purposes of introducing these two oracle estimators are different: $\hat{\theta}_n^{\text{SI}}(\mathcal{T}^{\text{OL}})$ is used to assess the effectiveness of the SURE criterion; by contrast, $\tilde{\theta}_n^{\text{SI}}(\mathcal{T}_n^{\text{OR}})$ is employed to evaluate the usefulness of the noiseless side information, that is, the maximal improvement in performance that can be achieved by incorporating ξ .

Denote $\pi_{1,n} = n^{-1} \sum_{i=1}^n \mathcal{I}(\xi_i \leq \tau_n^*)$ and $\pi_{2,n} = 1 - \pi_{1,n}$. Intuitively, the optimal partition τ_n^* (within the class of thresholding procedures utilizing two groups) is chosen to maximize the “discrepancy” between the two groups. For units in group $\mathcal{I}_{k,n}^{\tau^*}$, the mixture density of θ_i is given by

$$g_{k,n}(\theta) = (1 - p_{k,n}) \delta_0 + p_{k,n} h_{k,n}(\theta), \quad k = 1, 2, \quad (10)$$

where δ_0 is a Dirac delta function (null effects), $h_{k,n}$ is the (alternative) empirical density of nonnull effects. Following Remark 1, our theory developed based on the empirical density (10) can handle both random and deterministic models; this can be more clearly seen in our proofs of the theorems. Here, $p_{k,n}$ is the conditional proportion of nonnull effects for a given group and may be conceptualized as the probability that a randomly selected unit in group $\mathcal{I}_{k,n}^{\tau^*}$ is a nonnull effect.

We consider an asymptotic setup based on the sparse estimation framework in Johnstone (2015, chap. 8.6), which has been widely used in high-dimensional sparse inference (Johnstone and Silverman 1997; Donoho and Johnstone 1998; Abramovich et al. 2006; Mukherjee and Johnstone 2015; Cai and Sun 2017). Let $p_{1,n} = n^{-\alpha}$ and $p_{2,n} = n^{-\beta}$ for some $0 < \alpha < \beta \leq 1$. Define $\rho_n = \pi_{1,n}^{-1} \pi_{2,n}$. Consider the following parameter space

$$\Theta_n(\alpha, \beta, \rho_n) = \left\{ \theta \in \mathbb{R}^n : \|\theta\|_0 \leq n(n^{-\alpha} + \rho_n n^{-\beta}) / (1 + \rho_n) \right\}.$$

The maximal risk of ASUS over $\Theta_n(\alpha, \beta, \rho_n)$ is

$$\mathcal{R}_n^{\text{AS}}(\alpha, \beta, \rho_n) = \sup_{\theta \in \Theta_n(\alpha, \beta, \rho_n)} r_n(\hat{\mathcal{T}}, \theta).$$

Correspondingly, over the same parameter space $\Theta_n(\alpha, \beta, \rho_n)$, we let $\mathcal{R}_n^{\text{OS}}(\alpha, \beta, \rho_n)$ denote the maximal risk of the oracle procedure $\tilde{\theta}_n^{\text{SI}}(\mathcal{T}_n^{\text{OR}})$, and $\mathcal{R}_n^{\text{NS}}(\alpha, \beta, \rho_n)$ the minimax risk of all soft thresholding estimators without side information.

The risk difference $\mathcal{R}_n^{\text{NS}} - \mathcal{R}_n^{\text{OS}}$ is a key quantity that will be used in later analysis as the benchmark decision theoretical improvement due to incorporation of side information. Specifically, the noiseless side information ξ is *useful* if it provides nonnegligible improvement on the risk:

$$\lim_{n \rightarrow \infty} n(\mathcal{R}_n^{\text{NS}} - \mathcal{R}_n^{\text{OS}}) = \infty. \quad (11)$$

Moreover, the ASUS estimator is *asymptotically optimal* if its risk improvement over $\mathcal{R}_n^{\text{NS}}(\alpha, \beta, \rho_n)$ is asymptotically equal to that of the oracle:

$$\text{RI}_n = \frac{\mathcal{R}_n^{\text{NS}} - \mathcal{R}_n^{\text{AS}}}{\mathcal{R}_n^{\text{NS}} - \mathcal{R}_n^{\text{OS}}} \rightarrow 1 \text{ as } n \rightarrow \infty. \quad (12)$$

3.2. Usefulness of Side Information

We focus on Model (10), a hypothetical model based on the oracle partition τ_n^* . We state a few conditions that are needed in later analysis; some are essential for characterizing the situations where the side information is useful, that is, the oracle estimator $\tilde{\theta}_n^{\text{SI}}(\mathcal{T}_n^{\text{OR}})$ would provide nonnegligible efficiency gain over competitive estimators.

$$(A2.1) \quad \lim_{n \rightarrow \infty} \rho_n n^{-\gamma_0} = 0 \text{ for some } \gamma_0 < \beta - \alpha.$$

(A2.2) For some $\nu < 1/2$ and $k_n = \log n$, $\lim_{n \rightarrow \infty} k_n^\nu (1 - \pi_{1,n}) = \infty$.

(A2.3) For some $\nu < 1/2$, $\lim_{n \rightarrow \infty} n^\nu \pi_{1,n} p_{1,n} = \infty$.

(A2.4) Let $\bar{\sigma}_n^2 = n^{-1} \sum_{i=1}^n \sigma_i^2$ and $0 < \underline{\lim}_{n \rightarrow \infty} \bar{\sigma}_n^2 \leq \overline{\lim}_{n \rightarrow \infty} \bar{\sigma}_n^2 < \infty$.

Remark 4. (A2.1) implies $\pi_{2,n} p_{2,n} / (\pi_{1,n} p_{1,n}) \rightarrow 0$, which ensures that the oracle partition is effective in the sense that the two resulting groups have different sparsity levels. The asymmetric condition can be easily flipped for generalization. (A2.2) is a mild condition which allows $\pi_{1,n}$ to approach 1 but at a controlled rate. (A2.3) prevents the trivial setting where ASUS reduces to the SureShrink procedure with universal threshold $\sqrt{2 \log n}$, that is, the side information would not have any influence in the estimation process. See Lemma 3 (Section B of the supplementary materials) which shows that if $\overline{\lim}_{n \rightarrow \infty} n^{1/2} \pi_{1,n} p_{1,n} < \infty$, then ASUS reduces to the SureShrink procedure, that is, there is no need for creating groups. (A2.4) is a mild condition that is satisfied in most real life applications.

Now we study the usefulness of the noiseless side information. Following the theory in Johnstone (1994), the next theorem explicitly evaluates the risk difference $\mathcal{R}_n^{\text{NS}} - \mathcal{R}_n^{\text{OS}}$ up to higher order terms. The analysis overcomes the crudeness of the first-order asymptotics for evaluating thresholding rules as pointed out by Bickel (1983) and Johnstone (1994).

Theorem 2. Consider the oracle estimator defined in (8) and (9). Under Assumption A2.1, with $k_n = \log n$, for all $\nu < 1$, we have,

$$\begin{aligned} \mathcal{R}_n^{\text{NS}}(\alpha, \beta, \rho_n) - \mathcal{R}_n^{\text{OS}}(\alpha, \beta, \rho_n) \\ = \pi_{1,n} p_{1,n} \bar{\sigma}_n^2 \left\{ \log \pi_{1,n}^{-1} (2 - 3\alpha^{-1} k_n^{-1}) + O(k_n^{-\nu}) \right\}. \end{aligned}$$

It follows from (A2.3) that $\lim_{n \rightarrow \infty} n(\mathcal{R}_n^{\text{NS}} - \mathcal{R}_n^{\text{OS}}) = \infty$, establishing (11).

3.3. Asymptotic Optimality of ASUS

To evaluate the efficiency of ASUS, we need to compare the segmentation used by ASUS with that used by the oracle estimator. For a given segmentation hyper-parameter τ , define

$$\tilde{q}_{i,n}^{jk}(\tau) := \mathbb{P}_n(\hat{I}_i^j | I_i^k) \quad \text{for } j, k \in \{1, 2\}, i = 1, \dots, n,$$

where $\hat{I}_i^1 = \{S_i \leq \tau\}$, $I_i^1 = \{\xi_i \leq \tau_n^*\}$, $\hat{I}_i^2 = \mathbf{R} \setminus \hat{I}_i^1$, $I_i^2 = \mathbf{R} \setminus I_i^1$, and the probability operator \mathbb{P}_n is based on Model (10). Let

$$\begin{aligned} q_{i,n}^{jk}(\tau) &= \tilde{q}_{i,n}^{jk}(\tau) \\ \text{if } \inf_{\tau \in \mathbf{R}} \pi_{2,n} \tilde{q}_n^{12}(\tau) + \pi_{1,n} \tilde{q}_n^{21}(\tau) &< \inf_{\tau \in \mathbf{R}} \pi_{1,n} \tilde{q}_n^{11}(\tau) + \pi_{2,n} \tilde{q}_n^{22}(\tau) \end{aligned}$$

and otherwise $q_{i,n}^{jk}(\tau) = \tilde{q}_{i,n}^{kk}(\tau)$ and $q_{i,n}^{kk}(\tau) = 1 - q_{i,n}^{jk}(\tau)$ for $j \neq k$. Denote the weighted average

$$q_n^{jk}(\tau) = \frac{\sum_{i=1}^n q_{i,n}^{jk}(\tau) \sigma_i^2}{\sum_{i=1}^n \sigma_i^2}, \quad j, k \in \{1, 2\}.$$

Viewing the data-driven grouping step of ASUS as a classification procedure with the oracle segmentation corresponding

to the true states, we can conceptualize $q_n^{21}(\tau_n)$ and $q_n^{12}(\tau_n)$ as misclassification rates. Define the efficiency ratio

$$\mathcal{E}_n = \frac{\mathcal{R}_n^{\text{NS}} - \mathcal{R}_n^{\text{OS}}}{\mathcal{R}_n^{\text{AS}} - \mathcal{R}_n^{\text{OS}}}. \quad (13)$$

For notational simplicity, the dependence of this ratio on α, β, ρ_n is not explicitly marked. It follows from (12) that $\text{RI}_n = 1 - \mathcal{E}_n^{-1}$. Hence, a larger \mathcal{E}_n signifies better performance of ASUS. In particular, $\mathcal{E}_n \rightarrow \infty$ implies the asymptotic optimality of ASUS. The poly-log rates in the following theorem are sharp.

Theorem 3. Assume (A2.1)–(A2.4) hold. Let $k_n = \log n$. If there exists a sequence $\{\tau_n\}_{n \geq 1}$ such that

$$\lim_{n \rightarrow \infty} k_n^2 q_n^{21}(\tau_n) = 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} \rho_n q_n^{12}(\tau_n) = 0, \quad (14)$$

then ASUS is asymptotically optimal. In particular, for all $\nu < 1$ we have

$$\underline{\lim}_{n \rightarrow \infty} k_n^{-\nu} \mathcal{E}_n \geq 2 \underline{\lim}_{n \rightarrow \infty} \log \pi_{1,n}^{-1}. \quad (15)$$

Next we present two hierarchical models, respectively, with sub-Gaussian (SG) and sub-Exponential (SEXP) tails, under which the misclassification rates can be adequately controlled. Let $S_i | \xi_i$ be independent random variables with $\mu_i := \mu_i(\xi_i)$ and $(v_i(\xi_i), b_i(\xi_i))$ such that $\mathbb{E}\{\exp(\lambda(S_i - \mu_i))\} \leq \exp(\nu_i^2 \lambda^2 / 2)$ for all i and all $|\lambda| \leq b_i^{-1}$. Let $\overline{\lim}_i b_i < \infty$, $\overline{\lim}_i v_i < \infty$, and $\bar{b}_n = \sup_{1 \leq i \leq n} \max(2v_i^2, b_i)$. When $b_i = 0$, the distribution of S_i has sub-Gaussian tails. For two partitions A and B of the set $\{1, \dots, n\}$, define the ℓ_1 distance between the two sets $\{\mu_i : i \in A\}$ and $\{\mu_i : i \in B\}$ by $\text{dist}(A, B) = \inf\{|x - y| : x \in A, y \in B\}$. Let $c_n = \bar{b}_n(2 \log k_n + \log \rho_n)$. The following lemma provides a sufficient condition under which the requirements on misclassification rates (14) are satisfied. The proof of the lemma follows directly from the standard bounds for sub-Gaussian and sub-Exponential tails.

Lemma 1. Let $I_{1,n}^* = \{i : \xi_i \leq \tau_n^*\}$ and $I_{2,n}^* = \{1, \dots, n\} \setminus I_{1,n}^*$. The requirements on misclassification rates given by (14) are satisfied if

$$\underline{\lim}_{n \rightarrow \infty} c_n^{-\gamma} \text{dist}(I_{1,n}^*, I_{2,n}^*) > \gamma,$$

where γ is 1/2 if $\sup_i b_i = 0$ and 1 otherwise.

3.4. Robustness of ASUS

This section carries out a theoretical analysis to address the concern whether the performance of data combination procedures would deteriorate when pooling non-informative auxiliary data. We first characterize asymptotic regimes under which auxiliary data are non-informative (while the attention is confined to the prescribed class of two-group ASUS estimators), and then show that under such regimes, ASUS is robust in performance in the sense that it does not under-perform standard soft-thresholding methods.

Theorem 4. Suppose (A2.1)–(A2.4) hold. Let $\rho_n = n^{\nu_0}$ and $k_n = \log n$.

- (a) Consider the following situations: (i) $\lim_{n \rightarrow \infty} k_n^{-1} \rho_n q_n^{21}(\tau_n) = \infty$; and (ii) $\lim_{n \rightarrow \infty} n \rho_n q_n^{21}(\tau_n) = 0$ but $\lim_{n \rightarrow \infty} k_n^{-1} \rho_n q_n^{12}(\tau_n) = \infty$. If for all sequence $\{\tau_n\}_{n \geq 1}$ either (i) or (ii) holds, then we must have $\lim_{n \rightarrow \infty} \mathcal{E}_n = 1$. Hence, the auxiliary data are non-informative.
- (b) We always have $\lim_{n \rightarrow \infty} \mathcal{E}_n \geq 1$. Thus, even when pooling non-informative auxiliary data ASUS would be at least as efficient as competing soft thresholding based methods that do not use auxiliary data.

Our next result characterizes the performance of soft-thresholding estimators, where their efficacies are measured by the ratio of their respective maximal risks with respect to that of the oracle. The subsequent analysis is carried out using the ratios $\mathcal{R}_n^{AS} / \mathcal{R}_n^{OS}$ and $\mathcal{R}_n^{NS} / \mathcal{R}_n^{OS}$, instead of the ratios of the risk differences (e.g., RI_n and \mathcal{E}_n). In this metric, we see that any optimally tuned soft-thresholding procedure is robust; but the improvement due to the incorporation of the side information can be observed in the varied convergence rates. Concretely, we show that the maximal risk of any soft thresholding scheme lies within a constant multiple of the oracle risk \mathcal{R}_n^{OS} irrespective of the informativeness of the side information. Particularly, if $\lim_{n \rightarrow \infty} \pi_{1,n} > 0$, then $\lim_{n \rightarrow \infty} k_n^\nu (\mathcal{R}_n^{NS} / \mathcal{R}_n^{OS} - 1) = 0$ for all $\nu < 1$. By contrast, $\mathcal{R}_n^{AS} / \mathcal{R}_n^{OS}$ tends to 1 at a faster rate under the conditions of [Theorem 3](#).

Lemma 2. Let $c_n = \log \pi_{1,n}^{-1} / \{\alpha k_n - 1.5 \log(2\alpha k_n) + 2.5 + \log \phi(0)\}$ and $k_n = \log n$. For any $\nu < 1$, under Assumptions (A2.1)–(A2.4), we have

$$\begin{aligned} \lim_{n \rightarrow \infty} k_n^{2\nu} \{ \mathcal{R}_n^{NS} / \mathcal{R}_n^{OS} - \min(1 + c_n, \beta/\alpha) \} &= 0; \\ \overline{\lim}_{n \rightarrow \infty} k_n^{2\nu} \{ \mathcal{R}_n^{AS} / \mathcal{R}_n^{OS} - \min(1 + c_n, \beta/\alpha) \} &\leq 0. \end{aligned}$$

Under the conditions of [Theorem 3](#), if there exists $\delta > 0$ such that $\lim_{n \rightarrow \infty} k_n^\delta \log \pi_{1,n}^{-1} = \infty$, then

$$\begin{aligned} \lim_{n \rightarrow \infty} k_n^{1+\delta} (\mathcal{R}_n^{NS} / \mathcal{R}_n^{OS} - 1) &= \infty \quad \text{and} \\ \lim_{n \rightarrow \infty} k_n^{2\nu} (\mathcal{R}_n^{AS} / \mathcal{R}_n^{OS} - 1) &= 0. \end{aligned}$$

Hence, the risk of ASUS approaches the oracle risk at a faster rate.

4. Numerical Results

In this section, we compare the performance of ASUS against several competing methods, including (i) the SureShrink (SS) estimator in Donoho and Johnstone (1995), (ii) the extended James Stein estimator (EJS) discussed in Brown (2008), (iii) the empirical Bayes thresholding (EBT) in Johnstone and Silverman (2004), and (iv) the auxiliary screening (Aux-Scr) procedure using simulated data in [Section 4.2](#) and a real dataset in [Section 4.3](#). The ‘‘Aux-Scr’’ method is motivated by a comment for a reviewer. The idea is to first utilize \mathbf{S} to conduct a preliminary screening of the data, then discard coordinates that appear to contain little information, and finally apply soft-thresholding estimators on remaining coordinates. A detailed description of the Aux-Scr method is provided in [Section A](#) of the supplementary materials. More simulation results and an additional

real data analysis are provided in [Sections D and E](#) of the supplementary materials. Our numerical results suggest that ASUS enjoys superior numerical performance and the efficiency gain over competitive estimators is substantial in many settings.

4.1. Implementation and R-Package asus

The R-package `asus` has been developed to implement our proposed methodology. In this section, we provide some implementation details upon which our package has been built.

Our scheme for choosing \mathcal{T} involves minimizing $S(\mathcal{T}, \mathbf{Y}, \mathbf{S})$ with respect to \mathcal{T} . In particular, the optimal \mathcal{T} is given by

$$\hat{\mathcal{T}} = \underset{\tau \in \Delta_n, t_1, \dots, t_K \in [0, t_n]}{\text{arg min}} \quad S(\mathcal{T}, \mathbf{Y}, \mathbf{S}), \quad (16)$$

where Δ_n is a collection of $K - 1$ dimensional distinct points spanning \mathbf{R}_+^{K-1} and t_n denotes the universal threshold of $\sqrt{2 \log n}$. To solve this minimization problem, we proceed as follows: Let $S_{(1)}, S_{(n)}$ be the smallest and largest S_i , respectively. Consider a set of m_n equi-spaced points spanning $(S_{(1)}, S_{(n)})$ and take Δ_n to be a $\binom{m_n}{K-1} \times K - 1$ matrix where each row is a $K - 1$ dimensional sorted vector constructed out of the m_n points. For each τ^j in the j th row of Δ_n , determine $\{t_1^j, \dots, t_K^j\}$ by minimizing the SURE function for the K groups $\hat{\mathcal{T}}_k^\tau$. This step can easily be carried out via the hybrid scheme discussed in Donoho and Johnstone (1995). Using [Proposition 1](#), we compute $S(\mathcal{T}, \mathbf{Y}, \mathbf{S})$ at $\mathcal{T} = \{\tau^j, t_1^j, \dots, t_K^j\}$, and repeat this process for $j = 1, \dots, \binom{m_n}{K-1}$ to find $\hat{\mathcal{T}}$ using [Equation \(16\)](#). For choosing an appropriate K , the procedure discussed above can be repeated for each candidate value of K and an estimate of K may be taken to be the one that minimizes the SURE estimate of risk of ASUS over the candidate values of K . In [Section F](#) of the supplementary materials, we present a simple example that demonstrates this procedure for choosing K . Our practical recommendation is to take $m_n = 50 \log n$ and $K = 2$ which is computationally inexpensive and tends to provide substantial reduction in overall risk against the competing estimators in both simulations and real data examples we considered.

4.2. Simulation

This section presents results from two simulation studies, respectively, investigating the performances of ASUS in one-sample and two-sample estimation problems. To reveal the usefulness of side information and investigate the effectiveness of ASUS, we also include the oracle estimator $\hat{\theta}^{SI}(\mathcal{T}_n^{OR})$ in the comparison. The MSE of the oracle estimator (OR), which provides the lowest attainable risk, serves as a benchmark for assessing the performance of various methods. The R code that reproduces our simulation results can be downloaded from the following link: <https://github.com/trambakbanerjee/ASUS>.

4.2.1. One-Sample Estimation With Side Information

We generate our data based on hierarchical models (1)–(3), where we fix $n = 5000$, $K = 2$, and take $h_\theta(\xi_i, \eta_{1i}) = \xi_i + \eta_{1i}$. We simulate η_{1i} from a sparse mixture model $(1 - n^{-1/2})\delta_0 + n^{-1/2}N(2, 0.01)$. The latent vector $\boldsymbol{\xi}$ is simulated under the following two scenarios:

$$(S1) \xi \sim \left(\underbrace{\text{Unif}(6, 7)}_{\text{sample size} = 50}, \underbrace{\text{Unif}(2, 3)}_{\text{sample size} = 200}, \underbrace{0, \dots, 0}_{\text{sample size} = n - 250} \right),$$

$$(S2) \xi \sim \left(\underbrace{\text{Unif}(4, 8)}_{\text{sample size} = 200}, \underbrace{\text{Unif}(1, 3)}_{\text{sample size} = 800}, \underbrace{0, \dots, 0}_{\text{sample size} = n - 10^3} \right)$$

with $Y_i \sim N(\theta_i, 1)$. In practice, we only observe an auxiliary sequence \mathcal{S} , which can be viewed as a noisy version of ξ . To assess the impact of noise on the performance of ASUS, we consider four different settings. In settings 1 and 2, we simulate m samples of $\eta_2 = (\eta_{21}, \dots, \eta_{2n})$ from two different distributions and generate auxiliary sequences \mathcal{S}_1 and \mathcal{S}_2 as follows:

- (1) $\eta_{2i}^{(1)} \stackrel{iid}{\sim} \text{Laplace}(0, 4)$ with $\mathcal{S}_1 = |\xi + \bar{\eta}_2^{(1)}|$,
- (2) $\eta_{2i}^{(2)} \stackrel{iid}{\sim} \chi_{10}^2$ with $\mathcal{S}_2 = |\xi + \bar{\eta}_2^{(2)}|$,

where $\bar{\eta}_2^{(k)}$ is the average of $\eta_2^{(k)}$ over the m samples. For settings 3 and 4, we first introduce perturbations in the latent variable vector ξ and then generate auxiliary sequences $\mathcal{S}_3, \mathcal{S}_4$ as follows:

- (3) $\tilde{\xi}_i = \xi_i I_{\xi_i \neq 0} + \text{LogN}(0, 5/\sqrt{m}) I_{\xi_i = 0}$ with $\mathcal{S}_3 = |\tilde{\xi} + \rho \otimes \bar{\eta}_2^{(1)}|$, where ρ is a vector of n Rademacher random variables generated independently.
- (4) $\tilde{\xi}_i = \xi_i I_{\xi_i \neq 0} + t_{2m/10} I_{\xi_i = 0}$ with $\mathcal{S}_4 = |\tilde{\xi} - \rho \otimes \bar{\eta}_2^{(2)}|$, where ρ is a vector of n independent Bernoulli random variables with probability of success 0.75.

We vary m from 10 to 200 to investigate the impact of noise. The MSEs are obtained by averaging over $N = 500$ replications. The results for scenarios S1 and S2 are summarized in Table 1 and in Figures 3 and 4 wherein ASUS.j and Aux-Scr.j correspond to versions of ASUS and Aux-Scr that rely on the side information in the auxiliary sequence $\mathcal{S}_j, j = 1, \dots, 4$.

From the left panels of Figures 3 and 4, we see that ASUS exhibits the best performance when compared against EBT, EJS, and SureShrink estimators. In particular, ASUS.1, ASUS.2 outperform their counterparts ASUS.3, ASUS.4. This reveals how the usefulness of the latent sequence ξ would affect the performance of ASUS. Nonetheless, ASUS.3 and ASUS.4 still provide improvements over, and, crucially, are never worse than the SureShrink estimator. This reveals the impact of the accuracy of the auxiliary sequence \mathcal{S} (in capturing the information in ξ) on the performance of ASUS. The right panels of Figures 3 and 4 present the risk comparison between ASUS and Aux-Scr using the auxiliary sequences $\mathcal{S}_1, \dots, \mathcal{S}_4$. Not surprisingly, ASUS and Aux-Scr have almost identical risk performance using the auxiliary sequences $\mathcal{S}_1, \mathcal{S}_2$, and \mathcal{S}_3 for large m . As m increases, the accuracy of these auxiliary sequences increase but the negative Bernoulli perturbations in \mathcal{S}_4 interferes with its magnitude so that a smaller $|S_{i4}|$ may correspond to a signal coordinate. The Aux-Scr procedure which discards observations based on the magnitude of the auxiliary sequence may miss important signal coordinates while relying on \mathcal{S}_4 . ASUS, however, does not discard any observations and continues to exploit the available information in the noisy auxiliary sequences.

In Table 1, we report risk estimates and estimates of \mathcal{T} for ASUS when $m = 200$. The estimates of the hyper-parameters of Aux-Scr are provided in Table 2 of the supplementary materials and we only report its risk estimates here in Table 1. We can

Table 1. One-sample estimation with side information: risk estimates and estimates of \mathcal{T} for ASUS at $m = 200$.

		One-sample estimation with side information	
		Scenario S1	Scenario S2
OR	τ^*	2	1.003
	t_1^*, t_2^*	4.114, 0.138	4.073, 0.133
	n_1^*, n_2^*	4750, 250	4008, 992
	Risk	0.095	0.224
ASUS.1	τ	1.342	0.979
	t_1, t_2	4.114, 0.107	4.073, 0.156
	n_1, n_2	4748, 252	4008, 992
	Risk	0.097	0.243
ASUS.2	τ	11.229	5.82
	t_1, t_2	4.115, 0.106	4.073, 0.137
	n_1, n_2	4748, 252	4008, 992
	Risk	0.095	0.228
ASUS.3	τ	1.777	1.778
	t_1, t_2	4.089, 0.662	3.422, 0.441
	n_1, n_2	4271, 729	3606, 1394
	Risk	0.146	0.357
ASUS.4	τ	7.785	8.524
	t_1, t_2	1.360, 3.653	0.745, 3.864
	n_1, n_2	1775, 3225	2249, 2751
	Risk	0.165	0.356
Aux-Scr.1	Risk	0.097	0.243
Aux-Scr.2	Risk	0.095	0.232
Aux-Scr.3	Risk	0.147	0.360
Aux-Scr.4	Risk	0.186	0.414
SureShrink	Risk	0.191	0.429
EBT	Risk	0.253	0.692
EJS	Risk	0.408	0.652

NOTE: Here, $n_k^* = |\mathcal{I}_k^*|$ and $n_k = |\widehat{\mathcal{I}}_k^*|$ for $k = 1, 2$.

see that ASUS.1 and ASUS.2 choose similar thresholding hyper-parameters (t_1, t_2) as those of the oracle estimator. Moreover, ASUS.4 demonstrates a lower estimation risk than Aux-Scr.4 using the same auxiliary sequence \mathcal{S}_4 .

4.2.2. Two-Sample Estimation With Side Information

We consider the problem of estimating the difference of two Gaussian mean vectors. An auxiliary sequence can be constructed from data by following Example 3 in Section 2.2. We first simulate

$$\xi_{1i} \sim (1 - p_1)\delta_0 + p_1 \text{Unif}(3, 7), \quad \xi_{2i} \sim (1 - p_2)\delta_0 + p_2 \delta_{\{4\}},$$

where $\delta_{\{4\}}$ is the Dirac delta at 4 and then generate $\mu_{i,1} = \xi_{1i} + \eta_{1i}$ and $\mu_{i,2} = \xi_{2i} + \eta_{2i}$ with $\eta_{1i}, \eta_{2i} \stackrel{iid}{\sim} N(0, 0.01)$. The parameter of interest is $\theta = \mu_1 - \mu_2$ and the associated latent side information vector is $\xi = \xi_1 - \xi_2$. The observations based on the simulated mean vectors are generated as $U_i \sim N(\mu_{i,1}, \sigma_{i,1}^2), V_i \sim N(\mu_{i,2}, \sigma_{i,2}^2)$. Finally, the primary and auxiliary statistics are obtained as $Y_i = U_i - V_i, S_i = |U_i + \kappa_i V_i|$. We fix $p_1 = n^{-0.6}, p_2 = n^{-0.3}, \kappa_i = \sigma_{i,1}/\sigma_{i,2}$ and consider two scenarios where $\sigma_{i,1} = \sigma_{i,2} = 1$ under scenario S1 and $(\sigma_{1,i}^2, \sigma_{2,i}^2) \stackrel{iid}{\sim} \text{Unif}(0.1, 1)$ under scenario S2. The estimates of risks are obtained by averaging over $N = 1000$ replications. We vary n from 500 to 5000 to investigate the impact of the strength of side information. The simulation results are reported in Table 2 and Figure 5.

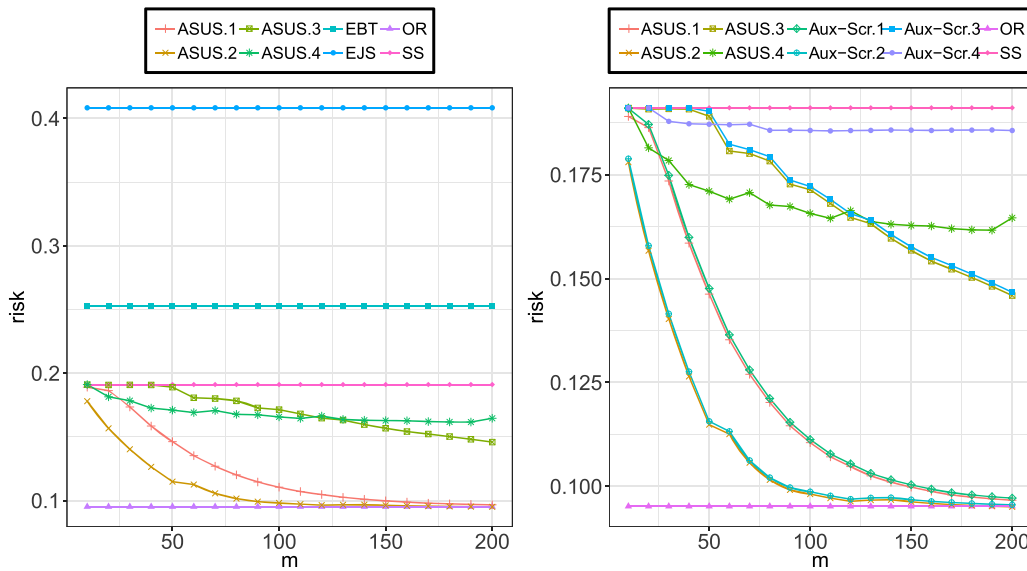


Figure 3. One-sample estimation with side information for scenario S1: estimated risks of different estimators. Left: ASUS versus EBT and EJS. Right: ASUS versus Aux-Scr.

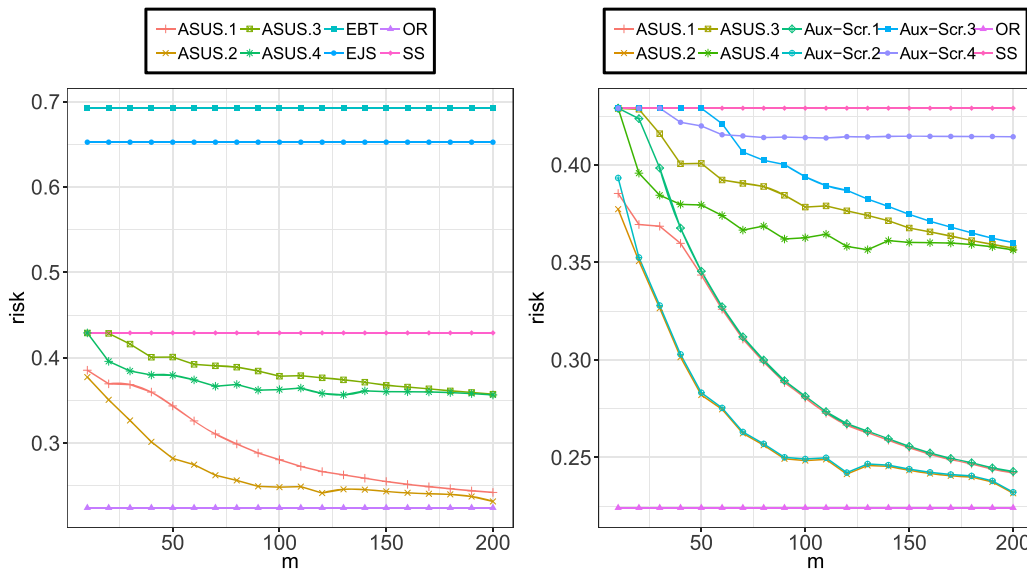


Figure 4. One-sample estimation with side information for scenario S2: estimated risks of different estimators. Left: ASUS versus EBT and EJS. Right: ASUS versus Aux-Scr.

We see that ASUS uses the side information in \mathcal{S} and exhibits the best performance across both scenarios. In scenario S2, the variances of Y_i are smaller, which leads to an improved risk performance of ASUS over scenario S1. Similar to the previous simulation study, the risk of ASUS would not exceed the risk of the SureShrink estimator across both the scenarios. Different magnitudes of the thresholding hyper-parameters (t_1, t_2) in Table 2 further corroborates the importance of the auxiliary statistics S_i in constructing groups with disparate sparsity levels and thereby improving the overall estimation accuracy. This is particularly true in the case of scenario S2 where EBT and SureShrink are competitive but ASUS is far more efficient because it has constructed two groups where one group holds majority of the signals and ASUS uses the smaller threshold t_2 to retain the signals. The other group holds majority of the noise wherein ASUS uses the larger threshold t_1 to shrink them to zero. Moreover, we notice that ASUS provides a better risk performance than Aux-Scr across both the scenarios. Using the

side information in \mathcal{S} , Aux-Scr discards observations that have $|S_i| \leq \tau$ thereby eliminating some potentially information rich signal coordinates and thus returns a higher risk than ASUS.

4.3. Analysis of RNA Sequence Data

We compare the performance of ASUS against the SureShrink (SS) estimator for analysis of the RNA sequence data described in the introduction. The goal is to estimate the true expression levels θ of the n genes that are infected with VZV strain. Through previous studies conducted in the lab, expression levels corresponding to other four experimental conditions, including uninfected cells (C1, 3 replicates), a fibrosarcoma cell line (C2, 3 replicates) and cells treated with interferons gamma (C3, 2 replicates), alpha (C4, 3 replicates), were also collected. Let X_i be the mean expression level of gene i across the four experimental conditions. Set $S_i = |X_i|$ with $K = 2$. Let $\hat{\theta}_i^S(t)$ denote the

Table 2. Two-sample estimation with side information: risk estimates and estimates of \mathcal{T} for ASUS at $n = 5000$.

		Two-sample estimation with side information	
		Scenario S1	Scenario S2
OR	τ^*	1.947	1.363
	t_1^*, t_2^*	4.106, 0.137	4.106, 0.424
	n_1^*, n_2^*	4584, 416	4583, 417
	Risk	0.185	0.132
	<hr/>		
ASUS	τ	3.167	2.504
	t_1, t_2	1.223, 0.253	3.058, 0.323
	n_1, n_2	4570, 430	4195, 805
	Risk	0.610	0.239
	<hr/>		
Aux-Scr	τ	14.385	2.768
	t_1, t_2	0.955, 0.002	5.708, 0.498
	n_1, n_2	4991, 9	3681, 1319
	Risk	0.688	0.258
	<hr/>		
SureShrink	Risk	0.688	0.318
EBT	Risk	0.761	0.311
EJS	Risk	0.891	0.600

NOTE: Here, $n_k^* = |\mathcal{I}_k^{\tau^*}|$ and $n_k = |\widehat{\mathcal{I}}_k^\tau|$ for $k = 1, 2$.

SureShrink estimator of θ_i based on Y_i , the mean expression level of gene i under the VZV condition. The standard deviation σ_i for the mean expression level pertaining to gene i across the 3 replicates of the VZV strain is derived from the study conducted in Sen et al. (2018).

On the right panel of Figure 6, the dotted line represents the minimum of the SURE risk of $\hat{\theta}^S(t)$, which is minimized at $t = 0.61$. The solid line represents the minimum of the SURE risk of a class of two-group estimators over a grid of τ values. ASUS chooses τ that minimizes the SURE risk (the red dot in Figure 6). The resulting risk is 1.99% at $\hat{\tau} = (1.25, 1.16, 0)$, a significant reduction compared to the risk estimate of 3.69% for $\hat{\theta}^S(t)$. To evaluate the results in a predictive framework, we next use only two replicates of the VZV strain for calibrating the hyper-parameters and calculate the prediction errors based on the hold out third replicate. The risk reduction by ASUS over SureShrink is about 30%.

In this example, a reduction in risk is possible because ASUS has efficiently exploited the sparsity information about

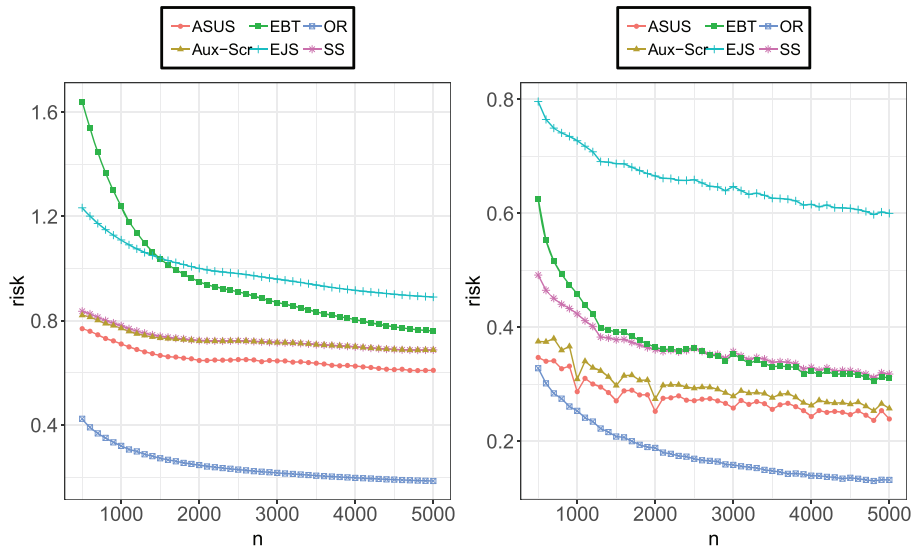


Figure 5. Two-sample estimation with side information: average risks of different estimators. Left: Scenario S1 and Right: Scenario S2.

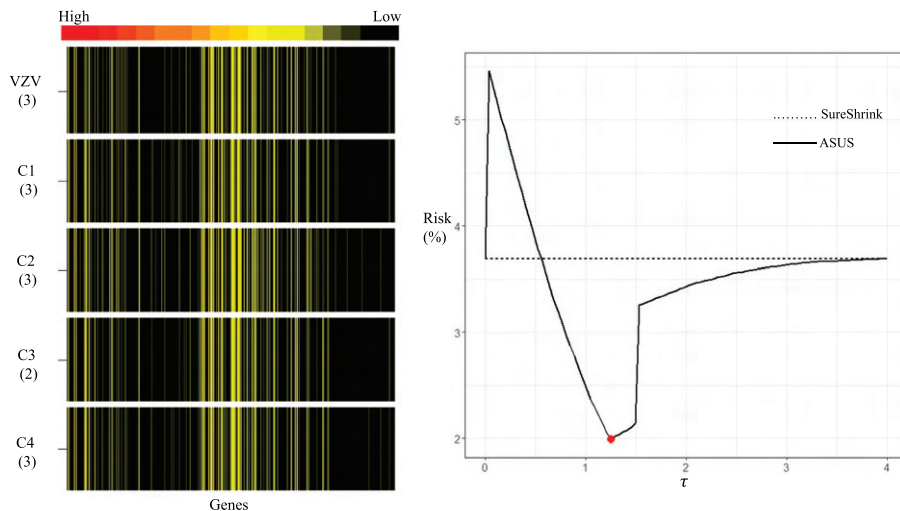


Figure 6. Left: Heat map showing the following from top to bottom: average expression levels of VZV, C1, C2, C3, and C4 across their respective replicates (in parenthesis). Right: SURE estimate of the risk of $\hat{\theta}_i^S(t)$ at $t = 0.61$ versus an unbiased estimate of the risk of ASUS for different values of τ .

Table 3. Summary of SureShrink and ASUS methods (RNA-seq data).

		RNA Seq
n		53,216
SureShrink	t	0.61
	SURE estimate	3.69
ASUS	τ	1.25
	t_1	1.16
	t_2	0
	n_1	39,535
	n_2	13,681
	SURE estimate	1.99

NOTE: $n_k = |\widehat{\mathcal{I}}_k^t|$ for $k = 1, 2$.

θ encoded by S . This can be seen, for example, from (i) the stark contrast between the magnitudes of thresholding hyperparameters t_1, t_2 for the two groups in Table 3 and (ii) the heat maps in Figure 6 where the genes expressions under the four experimental conditions follow the expression pattern of VZV. Moreover, the risk of Aux-Scr for this example was seen to be no better than the SureShrink estimator and thus has been excluded from the results reported in Table 3. Figure 7(a) presents the distribution of gene expression for genes that belong to groups $\widehat{\mathcal{I}}_1^t$ and $\widehat{\mathcal{I}}_2^t$. ASUS exploits the side information in S to partition the estimation units into two groups with very different sparsity levels and therefore returns a much smaller risk.

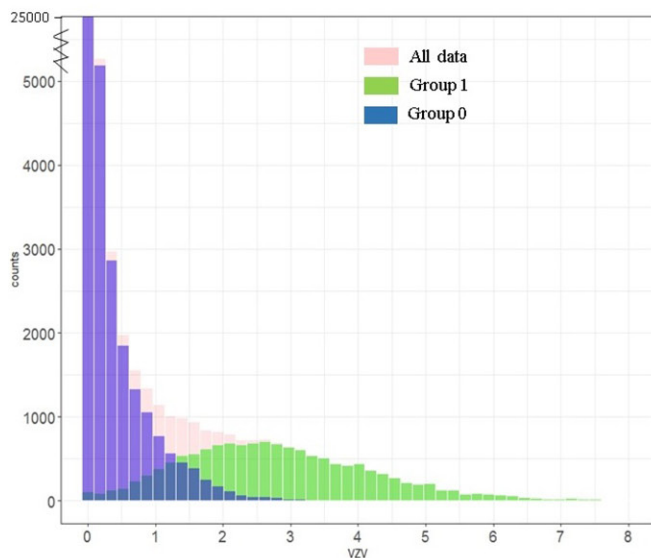
The ASUS estimator $\hat{\theta}^{SI}(\widehat{T})$ results in the discovery of 114 new genes than those discovered by using $\hat{\theta}^{SI}(t)$. Figure 7(b) shows the network of protein-protein interactions of 20 such genes. The interaction network is generated using NetworkAnalyst (Xia, Gill, and Hancock 2015) that maps the chosen genes to a comprehensive high-quality protein-protein interaction

(PPI) database based on InnateDB. A search algorithm is then performed to identify first-order neighbors (genes that directly interact with a given gene) for each of these mapped genes. The resulting nodes and their interaction partners are returned to build the network. In case of the RNA-seq data, the interaction network of the 20 new genes indicates that ASUS may help reveal important biological synergies between genes that have a high estimated expression level for VZV and other genes in the human genome.

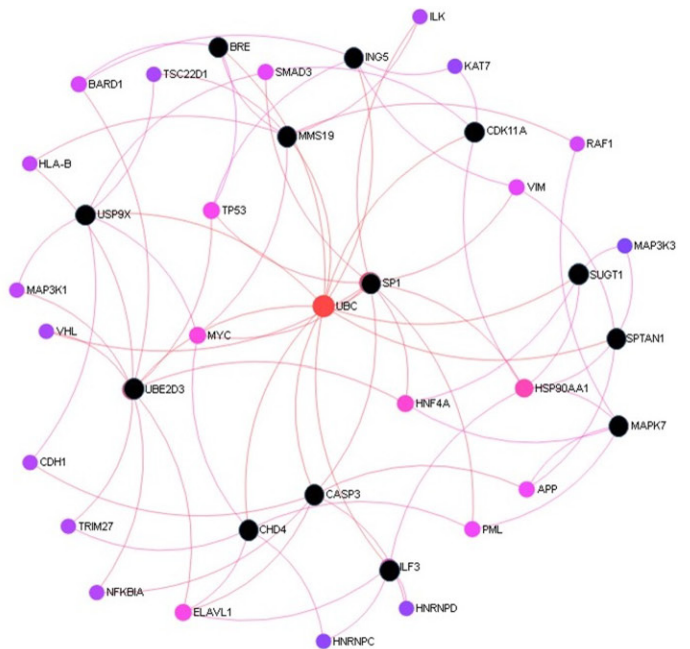
5. Discussion

In high-dimensional estimation and testing problems, the sparsity structure can be encoded in various ways; we have considered three basic settings where the structural information on sparsity may be extracted from (i) prior or domain-specific knowledge, (ii) covariate sequence based on the same data, or (iii) summary statistics based on secondary data sources. This article develops a general integrative framework for sparse estimation that is capable of handling all three scenarios. We use higher order minimax optimality tools to establish the adaptivity and robustness of ASUS. Numerical studies using both simulated and real data corroborate the improvement of ASUS over existing methods.

We conclude the article with a discussion of several open issues. Firstly, in large-scale compound estimation problems, various data structures such as sparsity, heteroscedasticity, dependency, and hierarchy are often available alongside the primary summary statistics. ASUS can only handle the sparsity structure; and it is desirable to develop a unified framework that can effectively incorporate other types of structures into



(a)



(b)

Figure 7. (a) Histogram of gene expressions for VZV. Group 1 is $\widehat{\mathcal{I}}_2^t$ and Group 0 is $\widehat{\mathcal{I}}_1^t$. (b) A network of 20 new genes highlighted in black with their interaction partners.

inference. New theoretical frameworks will be needed to characterize the usefulness of various types of side information and to establish precise conditions under which the new integrative method is asymptotically optimal. Secondly, in situations where there are multiple auxiliary sequences, it is unclear how to modify the ASUS framework to construct groups using an auxiliary matrix. The computation involved in the search for the optimal group-wise thresholds, which requires the evaluation of the SURE function for every possible combination of group-wise thresholds, quickly becomes prohibitively expensive as the number of columns increases. Finally, the higher dimension would affect the stability of an integrative procedure adversely. A promising idea for handling multiple auxiliary sequences is to construct a new auxiliary sequence that represents the “optimal use” of all available side information. However, the search for this optimal direction of projection is quite challenging. It would be of great interest to explore these directions in future research.

Supplementary Materials

This supplement contains a detailed description of the Auxiliary Screening procedure (Aux-Scr), proofs of the results in Section 2 and 3 of the main paper, additional simulation experiments, a real data analysis and an example that demonstrates a data driven procedure for choosing K .

Acknowledgments

We thank Ann Arvin and Nandini Sen for helpful discussions on the virology application. We thank the AE and two referees for the constructive suggestions that have greatly helped to improve the presentation of the article. In particular, we are grateful to an excellent comment from a referee that leads to the Bayesian interpretation of ASUS in [Section 2.4](#).

Funding

The research of WS was supported in part by NSF grants DMS-CAREER 1255406 and DMS-1712983. TB and GM were partially supported by NSF DMS-1811866 and by the Zumberge individual award from the University of Southern California’s James H. Zumberge Faculty Research and Innovation Fund.

References

- Abramovich, F., Benjamini, Y., Donoho, D. L., and Johnstone, I. M. (2006), “Adapting to Unknown Sparsity by Controlling the False Discovery Rate,” *The Annals of Statistics*, 34, 584–653. [2053,2055,2059]
- Abramovich, F., Grinshtein, V., and Pensky, M. (2007), “On Optimality of Bayesian Testimation in the Normal Means Problem,” *The Annals of Statistics*, 35, 2261–2286. [2055]
- Bickel, P. (1983), “Minimax Estimation of a Normal Mean Subject to Doing Well at a Point,” in *Recent Advances in Statistics*, eds. M. H. Rizvi, J. S. Rustagi, and D. Siegmund, New York: Academic Press, pp. 511–528. [2060]
- Brown, L. D. (2008), “In-Season Prediction of Batting Averages: A Field Test of Empirical Bayes and Bayes Methodologies,” *The Annals of Applied Statistics*, 2, 113–152. [2061]
- Brown, L. D., and Greenshtein, E. (2009), “Nonparametric Empirical Bayes and Compound Decision Approaches to Estimation of a High-Dimensional Vector of Normal Means,” *The Annals of Statistics*, 37, 1685–1704. [2055,2058]
- Brown, L. D., Mukherjee, G., and Weinstein, A. (2018), “Empirical Bayes Estimates for a 2-Way Cross-Classified Additive Model,” *The Annals of Statistics* 46, 1693–1720. [2058]

- Cai, T. T., Low, M., and Ma, Z. (2014), “Adaptive Confidence Bands for Nonparametric Regression Functions,” *Journal of the American Statistical Association*, 109, 1054–1070. [2055]
- Cai, T. T., and Sun, W. (2017), “Optimal Screening and Discovery of Sparse Signals With Applications to Multistage High Throughput Studies,” *Journal of the Royal Statistical Society, Series B*, 79, 197–223. [2059]
- Cai, T. T., Sun, W., and Wang, W. (2019), “CARS: Covariate Assisted Ranking and Screening for Large-Scale Two-Sample Inference,” *Journal of the Royal Statistical Society, Series B*, 81, 187–234. [2056]
- Calvano, S. E., Xiao, W., Richards, D. R., Felciano, R. M., Baker, H. V., Cho, R. J., Chen, R. O., Brownstein, B. H., Cobb, J. P., Tschoeke, S. K., and Miller-Graziano, C. (2005), “A Network-Based Analysis of Systemic Inflammation in Humans,” *Nature*, 437, 1032–1037. [2053]
- Collier, O., Comminges, L., and Tsybakov, A. B. (2017), “Minimax Estimation of Linear and Quadratic Functionals on Sparsity Classes,” *The Annals of Statistics*, 45, 923–958. [2055]
- Cover, T. M., and Thomas, J. A. (2012), *Elements of Information Theory*, New York: Wiley. [2055]
- Donoho, D., and Jin, J. (2004), “Higher Criticism for Detecting Sparse Heterogeneous Mixtures,” *The Annals of Statistics*, 32, 962–994. [2053]
- Donoho, D. L., and Johnstone, I. M. (1995), “Adapting to Unknown Smoothness via Wavelet Shrinkage,” *Journal of the American Statistical Association*, 90, 1200–1224. [2057,2061]
- (1998), “Minimax Estimation via Wavelet Shrinkage,” *The Annals of Statistics*, 26, 879–921. [2059]
- Efron, B. (2011), “Tweedie’s Formula and Selection Bias,” *Journal of the American Statistical Association*, 106, 1602–1614. [2058]
- Erickson, S., and Sabatti, C. (2005), “Empirical Bayes Estimation of a Sparse Vector of Gene Expression Changes,” *Statistical Applications in Genetics and Molecular Biology*, 4, 1132. [2053]
- Holland, D., Wang, Y., Thompson, W. K., Schork, A., Chen, C.-H., Lo, M.-T., Witoelar, A., Werge, T., O’Donovan, M., Andreassen, O. A., and Dale, A. M. (2016), “Estimating Effect Sizes and Expected Replication Probabilities From GWAS Summary Statistics,” *Frontiers in Genetics*, 7, 15. [2053]
- Johnstone, I. M. (1994), “On Minimax Estimation of a Sparse Normal Mean Vector,” *The Annals of Statistics*, 22, 271–289. [2055,2060]
- (2015), “Gaussian Estimation: Sequence and Wavelet Models,” Draft Version. [2054,2055,2059]
- Johnstone, I. M., and Silverman, B. W. (1997), “Wavelet Threshold Estimators for Data With Correlated Noise,” *Journal of the Royal Statistical Society, Series B*, 59, 319–351. [2059]
- (2004), “Needles and Straw in Haystacks: Empirical Bayes Estimates of Possibly Sparse Sequences,” *The Annals of Statistics*, 32, 1594–1649. [2053,2058,2061]
- Ke, T., Jin, J., and Fan, J. (2014), “Covariance Assisted Screening and Estimation,” *Annals of Statistics*, 42, 2202–2242. [2055]
- Kou, S., and Yang, J. J. (2015), “Optimal Shrinkage Estimation in Heteroscedastic Hierarchical Linear Models,” arXiv no. 1503.06262. [2055]
- Li, C., Li, M., Lange, E. M., and Watanabe, R. M. (2008), “Prioritized Subset Analysis: Improving Power in Genome-Wide Association Studies,” *Human Heredity*, 65, 129–141. [2056]
- Mallat, S. (2008), *A Wavelet Tour of Signal Processing: The Sparse Way* (3rd ed.), Boston, MA: Academic Press. [2054]
- Matsui, S. (2013), “Genomic Biomarkers for Personalized Medicine: Development and Validation in Clinical Studies,” *Computational and Mathematical Methods in Medicine*, 2013, 865980. [2053]
- Mukherjee, G., and Johnstone, I. M. (2015), “Exact Minimax Estimation of the Predictive Density in Sparse Gaussian Models,” *Annals of Statistics*, 43, 937. [2059]
- Sen, N., Sung, P., Panda, A., and Arvin, A. M. (2018), “Distinctive Roles for Type I and Type II Interferons and Interferon Regulatory Factors in the Host Cell Defense Against Varicella-Zoster Virus,” *Journal of Virology*, 92, e01151-18. [2053,2064]
- Sun, W., and Wei, Z. (2011), “Multiple Testing for Pattern Identification, With Applications to Microarray Time-Course Experiments,” *Journal of the American Statistical Association*, 106, 73–88. [2053]

- Tan, Z. (2015), "Improved Minimax Estimation of a Multivariate Normal Mean Under Heteroscedasticity," *Bernoulli*, 21, 574–603. [2055]
- Tibshirani, R. J. (2014), "Adaptive Piecewise Polynomial Estimation via Trend Filtering," *The Annals of Statistics*, 42, 285–323. [2055]
- Watanabe, S., Kuzuoka, S., and Tan, V. Y. (2015), "Nonasymptotic and Second-Order Achievability Bounds for Coding With Side-Information," *IEEE Transactions on Information Theory*, 61, 1574–1605. [2055]
- Weinstein, A., Ma, Z., Brown, L. D., and Zhang, C.-H. (2018), "Group-Linear Empirical Bayes Estimates for a Heteroscedastic Normal Mean," *Journal of the American Statistical Association*, 113, 698–710. [2055]
- Wyner, A. (1975), "On Source Coding With Side Information at the Decoder," *IEEE Transactions on Information Theory*, 21, 294–300. [2055]
- Xia, J., Gill, E. E., and Hancock, R. E. (2015), "Networkanalyst for Statistical, Visual and Network-Based Meta-Analysis of Gene Expression Data," *Nature Protocols*, 10, 823–844. [2065]
- Xie, X., Kou, S., and Brown, L. D. (2012), "SURE Estimates for a Heteroscedastic Hierarchical Model," *Journal of the American Statistical Association*, 107, 1465–1479. [2055,2058]
- Zerboni, L., Sen, N., Oliver, S. L., and Arvin, A. M. (2014), "Molecular Mechanisms of Varicella Zoster Virus Pathogenesis," *Nature Reviews Microbiology*, 12, 197–210. [2053]

Supplementary Material for “Adaptive Sparse Estimation with Side Information”

This supplement contains a detailed description of the Auxiliary Screening procedure (Aux-Scr) (Section A), proofs of the results in Section 2 and 3 of the main paper (Sections B and C respectively), additional simulation experiments (Section D), a real data analysis (Section E) and an example that demonstrates a data driven procedure for choosing K (Section F).

A The Auxiliary Screening approach

We consider a potential competitor of ASUS, called Aux-Scr, which uses the auxiliary sequence \mathbf{S} to conduct a preliminary screening of the primary data thereby discarding data instances that contain little information and retains the potentially information rich primary data for estimation. Using the notation described in the main paper, we define Aux-Scr for $K = 2$ groups as follows:

Let \mathbf{Y} and \mathbf{S} denote the primary statistics and auxiliary sequence obeying the models (1)-(3) described in Section 2.1 of the main paper. Let $\eta_t(\cdot)$ be a soft-thresholding operator such that

$$\eta_t(Y_i) = \begin{cases} -Y_i\sigma_i^{-1}, & \text{if } |Y_i\sigma_i^{-1}| \leq t; \\ -t \operatorname{sign}(Y_i\sigma_i^{-1}), & \text{otherwise.} \end{cases}$$

The Aux-Scr estimator operates in two steps: first it constructs $K = 2$ groups using the magnitude of \mathbf{S} , where in group 1 $|S_i| \leq \tau$ and in group 2 $|S_i| > \tau$. Then it conducts soft-thresholding estimation using the primary statistics \mathbf{Y} in group 2 and estimates $\hat{\theta}_i = 0$ for all coordinates that belong to group 1. The tuning parameters for both grouping and shrinkage are determined using the SURE criterion.

Procedure 1. For $k = 1, 2$, denote $\widehat{\mathcal{I}}_1^T = \{i : |S_i| \leq \tau\}$ and $\widehat{\mathcal{I}}_2^T = \{i : |S_i| > \tau\}$. Consider the following class of shrinkage estimators:

$$\hat{\theta}_i^{SI}(\mathcal{T}) := Y_i + \sigma_i \eta_{t_k}(Y_i) \text{ if } i \in \widehat{\mathcal{I}}_k^T,$$

where, $\mathcal{T} = \{\tau, t_2\}$, t_2 varies in $[0, t_n]$ with $t_n = (2 \log n)^{1/2}$ and $t_1 = \max\{|Y_i \sigma_i^{-1}| : i \in \widehat{\mathcal{I}}_1^\tau\}$. Thus, the set of all possible hyper-parameter \mathcal{T} values is $\mathcal{H}_n = \mathbf{R}_+ \times [0, t_n]$. Define the SURE function

$$S(\mathcal{T}, \mathbf{Y}, \mathbf{S}) = n^{-1} \left[\sum_{i=1}^n \sigma_i^2 + \sum_{k=1}^K \sum_{i \in \widehat{\mathcal{I}}_k^\tau} \{ \sigma_i^2 (|Y_i \sigma_i^{-1}| \wedge t_k)^2 - 2 \sigma_i^2 I(|Y_i \sigma_i^{-1}| \leq t_k) \} \right]. \quad (1)$$

Let $\widehat{\mathcal{T}} = \arg \min_{\mathcal{T} \in \mathcal{H}_n} S(\mathcal{T}, \mathbf{Y}, \mathbf{S})$. Then, the Aux-Scr estimator is given by $\widehat{\theta}_i^{SI}(\widehat{\mathcal{T}})$ with $t_1 = \max\{|Y_i \sigma_i^{-1}| : i \in \widehat{\mathcal{I}}_1^\tau\}$.

Following the arguments in the proof of Proposition 1, it can be shown that equation (1) is an unbiased estimate of the true risk. Moreover, unlike ASUS, the thresholding hyper-parameter t_1 is fixed and the SURE criteria is used to select the grouping hyper-parameter τ and the thresholding hyper-parameter t_2 .

When compared to Aux-Scr, ASUS has three distinct advantages: optimality, robustness and adaptivity. First, the screening strategy does not address the important issue on how to set an optimal group-wise cutoff in the screening stage; this issue has been resolved by the SURE criterion in ASUS. Second, the ‘‘divide-and-threshold’’ strategy adopted by ASUS is clearly more effective than the ‘‘screening’’ strategy that directly throws away a lot of data. When \mathbf{S} is imperfect in capturing the sparsity structure, the screening step would inevitably miss important signal coordinates. By contrast, ASUS is more robust to noisy side information as it only utilizes \mathbf{S} to divide \mathbf{Y} into groups; no coordinates are discarded directly. Finally, Aux-Scr uses the same threshold for all coordinates that pass the preliminary screening stage. By contrast, ASUS is more adaptive to the unknown sparsity as it sets varied group-wise thresholds to reflect the possibly varied sparsity levels across groups.

B Proofs of the results in Section 2

Proof of Proposition 1 - Recall that $r_n(\mathcal{T}; \boldsymbol{\theta}) = n^{-1} \mathbb{E} \|\widehat{\boldsymbol{\theta}}^{SI}(\mathcal{T}) - \boldsymbol{\theta}\|^2$ where the expectation is taken with respect to the joint distribution of (Y_i, S_i) given ξ_i for $i = 1, 2, \dots, n$. Now, expanding $\|\widehat{\boldsymbol{\theta}}^{SI}(\mathcal{T}) - \boldsymbol{\theta}\|^2$ as $\|\mathbf{Y} - \boldsymbol{\theta}\|^2 + \|\widehat{\boldsymbol{\theta}}^{SI}(\mathcal{T}) - \mathbf{Y}\|^2 + 2\langle \mathbf{Y} - \boldsymbol{\theta}, \widehat{\boldsymbol{\theta}}^{SI}(\mathcal{T}) - \mathbf{Y} \rangle$ and taking expectation, we have,

$$n r_n(\mathcal{T}; \boldsymbol{\theta}) = \sum_{i=1}^n \sigma_i^2 + \sum_{i=1}^n \sum_{k=1}^K \sigma_i^2 \mathbb{E} \left\{ \eta_{t_k}^2(Y_i) I(i \in \mathcal{I}_k^\tau) \right\} + 2 \sum_{i=1}^n \sum_{k=1}^K \sigma_i \mathbb{E} \left\{ (Y_i - \theta_i) \eta_{t_k}(Y_i) I(i \in \mathcal{I}_k^\tau) \right\}$$

Observe that from Models (1) to (3), the pairs $\left\{ (Y_i - \theta_i) \eta_{t_k}(Y_i), I(i \in \widehat{\mathcal{I}}_k^\tau) \right\}$ are uncorrelated for each i and for all $k = 1, \dots, K$. Further, note that by Lemma 1 of [Stein \(1981\)](#)

$$\mathbb{E} \left\{ \eta'_{t_k}(Y_i) \right\} = \sigma_i^{-1} \int_{\mathbf{R}} \eta'_{t_k}(u) \phi\left(\frac{u - \theta_i}{\sigma_i}\right) du = \sigma_i^{-2} \mathbb{E} \left\{ \eta_{t_k}(Y_i) (Y_i - \theta_i) \right\}.$$

Thus, $-\sigma_i^2 \mathbb{E}\{I(|Y_i \sigma_i^{-1}| \leq t_k)\} = \sigma_i \mathbb{E}\{\eta_{t_k}(Y_i)(Y_i - \theta_i)\}$ which completes the proof.

Proof of Theorem 1, statement (a) - First note that we can decompose $S(\mathcal{T}, \mathbf{Y}, \mathbf{S})$ into K components:

$S(\mathcal{T}, \mathbf{Y}, \mathbf{S}) = \sum_{k=1}^K S_k(\mathcal{T}_k, \mathbf{Y}, \mathbf{S})$ where

$$S_k(\mathcal{T}_k, \mathbf{Y}, \mathbf{S}) = n^{-1} \sum_{i=1}^n \left\{ \sigma_i^2 - 2\sigma_i^2 I(|Y_i \sigma_i^{-1}| \leq t_k) + \sigma_i^2 (|Y_i \sigma_i^{-1}| \wedge t_k)^2 \right\} I\{S_i \in (\tau_{k-1}, \tau_k]\}$$

and $\mathcal{T}_k = \{\tau_{k-1}, \tau_k, t_k\}$. Let

$$S_i^k(\mathcal{T}_k, Y_i, S_i) = \sigma_i^2 - 2\sigma_i^2 I(|Y_i \sigma_i^{-1}| \leq t_k) + \sigma_i^2 (|Y_i \sigma_i^{-1}| \wedge t_k)^2 I\{S_i \in (\tau_{k-1}, \tau_k]\}$$

and notice that $S_i^k(\mathcal{T}_k, Y_i, S_i)$ is bounded above by $\sigma_i^2(1 + t_k^2)$ for each i . Also, we can decompose the risk $r_n(\mathcal{T}, \boldsymbol{\theta}) = \sum_{k=1}^K r_n^k(\mathcal{T}_k, \boldsymbol{\theta})$ where

$$r_n^k(\mathcal{T}_k, \boldsymbol{\theta}) = n^{-1} \sum_{i=1}^n \mathbb{E} \left[\left(Y_i + \sigma_i \eta_{t_k}(Y_i) - \theta_i \right) I\{X_i \in (\tau_{k-1}, \tau_k]\} \right]^2 = n^{-1} \sum_{i=1}^n r_i^k(\mathcal{T}_k, \theta_i)$$

and noting that $r_i^k(\mathcal{T}_k, \theta_i) \leq 2\sigma_i^2(1 + t_k^2)$. The last inequality follows from the upper bound on the risk of soft thresholding estimator at threshold t_k . Now, by triangle inequality, it is enough to show

$$c_n \mathbb{E} \left\{ \sup_{\mathcal{T}_k \in \mathbb{R}^2 \times [0, t_n]} \left| S_k(\mathcal{T}_k, \mathbf{Y}, \mathbf{S}) - r_n^k(\mathcal{T}_k; \boldsymbol{\theta}) \right| \right\} < \infty \text{ for all } i \text{ and for all large } n \quad (2)$$

Based on the form of $S_k(\mathcal{T}_k, \mathbf{Y}, \mathbf{S})$, we consider a re-parametrization of the problem with respect to $0 \leq \tilde{\tau}_{k-1} < \tilde{\tau}_k \leq 1$ where $\tilde{\tau}_k = \max_{i \in \mathcal{I}_k^\tau} F_{S_i}(\tau_k)$, $\tilde{\tau}_{k-1} = \min_{i \in \mathcal{I}_k^\tau} F_{S_i}(\tau_{k-1})$ and F_{S_i} is the distribution function of S_i . The only τ_{k-1}, τ_k dependent quantity in the expression of $S_k(\mathcal{T}_k, \mathbf{Y}, \mathbf{S})$ is $\widehat{\mathcal{I}}_k^\tau = \{i : \tau_{k-1} < S_i \leq \tau_k\}$ which is re-parametrized to $\widehat{\mathcal{I}}_k^{\tilde{\tau}} = \{i : \tilde{\tau}_{k-1} < F_{S_i}(s_i) \leq \tilde{\tau}_k\}$. This facilitates the analysis since now the supremum with respect to $\tilde{\mathcal{T}}_k = \{\tilde{\tau}_{k-1}, \tilde{\tau}_k, t_k\}$ is actually over a compact set.

We will mimic the proof of Proposition 1 of [Donoho and Johnstone \(1995\)](#), hereafter referred to as DJ95P1, to prove equation (2). For the other terms similar arguments will continue to hold.

Let $S_k(\tilde{\mathcal{T}}_k, \mathbf{Y}, \mathbf{S}) - r_n(\tilde{\mathcal{T}}_k, \boldsymbol{\theta}) = n^{-1} \sum_{i=1}^n U_i(\tilde{\mathcal{T}}_k) = V_n(\tilde{\mathcal{T}}_k)$ where $\mathbb{E} U_i(\tilde{\mathcal{T}}_k) = 0$ and from the upper bounds on $S_i^k(\mathcal{T}_k, Y_i, S_i)$ and $r_i^k(\mathcal{T}_k, \theta_i)$, $|U_i(\tilde{\mathcal{T}}_k)| \leq 3(1 + t_n^2)\sigma_i^2$. Now replace $Z_n(t)$ in DJ95P1 with $V_n(\tilde{\mathcal{T}}_k)$ and notice that Hoeffding's inequality gives, for a fixed $\tilde{\mathcal{T}}_k$ and (for now) arbitrary $r_n > 1$,

$$\mathbb{P} \left\{ \left| V_n(\tilde{\mathcal{T}}_k) \right| > \frac{r_n}{\sqrt{n}} \right\} \leq 2 \exp \left\{ - \frac{nr_n}{18(1 + t_n^2)^2 \sum_{i=1}^n \sigma_i^4} \right\}$$

Next, for a perturbation $\tilde{\mathcal{T}}_k' = \{\tilde{\tau}'_{k-1}, \tilde{\tau}'_k, t'_k\}$ of $\tilde{\mathcal{T}}_k$ where $\tilde{\tau}'_k > \tilde{\tau}_k$, $\tilde{\tau}'_{k-1} > \tilde{\tau}_{k-1}$ and $t'_k > t_k$, we wish

to bound the increments $\left|V_n(\tilde{\mathcal{T}}_k) - V_n(\tilde{\mathcal{T}}'_k)\right|$. To that effect, define $\tilde{\mathcal{T}}_k^{(r,s)}$ to be $\tilde{\mathcal{T}}_k$ but with components (r, s) replaced by the components (r, s) of $\tilde{\mathcal{T}}'_k$, $r < s = 1, 2, 3$. Then we can write, dropping the subscript k from $\tilde{\mathcal{T}}_k$ for brevity, that $|S_k(\tilde{\mathcal{T}}, \mathbf{Y}, \mathbf{S}) - S_k(\tilde{\mathcal{T}}', \mathbf{Y}, \mathbf{S})|$ is bounded above by the sum of three terms: $|S_k(\tilde{\mathcal{T}}, \mathbf{Y}, \mathbf{S}) - S_k(\tilde{\mathcal{T}}^{(3)}, \mathbf{Y}, \mathbf{S})|$, $|S_k(\tilde{\mathcal{T}}^{(3)}, \mathbf{Y}, \mathbf{S}) - S_k(\tilde{\mathcal{T}}^{(2,3)}, \mathbf{Y}, \mathbf{S})|$ and $|S_k(\tilde{\mathcal{T}}^{(2,3)}, \mathbf{Y}, \mathbf{S}) - S_k(\tilde{\mathcal{T}}', \mathbf{Y}, \mathbf{S})|$. The first term is bounded by $n^{-1}(2 + t'^2 - t^2)N_n(t, t')$ which follows directly from the proof of DJ95P1 with $N_n(t, t') = \sum_{i=1}^n \sigma_i^2 I(t < |Y_i \sigma_i^{-1}| \leq t')$. The second term is bounded by $n^{-1}(3 + t'^2)M_n(\tilde{\tau}_k, \tilde{\tau}'_k)$ where $M_n(\tilde{\tau}, \tilde{\tau}') = \sum_{i=1}^n \sigma_i^2 I(\tilde{\tau} < F_{S_i}(s_i) < \tilde{\tau}')$ and similarly the third term is bounded by $n^{-1}(3 + t'^2)M_n(\tilde{\tau}_{k-1}, \tilde{\tau}'_{k-1})$.

For the risk $r_n^k(\tilde{\mathcal{T}}, \boldsymbol{\theta})$, we follow the same decomposition and upper bound $|r_n^k(\tilde{\mathcal{T}}, \boldsymbol{\theta}) - r_n^k(\tilde{\mathcal{T}}', \boldsymbol{\theta})|$ by

$$\left|r_n^k(\tilde{\mathcal{T}}, \boldsymbol{\theta}) - r_n^k(\tilde{\mathcal{T}}^{(3)}, \boldsymbol{\theta})\right| + \left|r_n^k(\tilde{\mathcal{T}}^{(3)}, \boldsymbol{\theta}) - r_n^k(\tilde{\mathcal{T}}^{(2,3)}, \boldsymbol{\theta})\right| + \left|r_n^k(\tilde{\mathcal{T}}^{(2,3)}, \boldsymbol{\theta}) - r_n^k(\tilde{\mathcal{T}}', \boldsymbol{\theta})\right|$$

From the proof of DJ95P1, we upper bound the first term above by $5n^{-1}\delta_{0n}t_n \sum_{i=1}^n \sigma_i^2$ as long as $|t - t'| < \delta_{0n}$ for some $\delta_{0n} > 0$. The second and the third terms are upper-bounded by $2n^{-1}\delta_{1n}(1 + t_n^2) \sum_{i=1}^n \sigma_i^2$ and $2n^{-1}\delta_{2n}(1 + t_n^2) \sum_{i=1}^n \sigma_i^2$ respectively as long as $|\tilde{\tau}_k - \tilde{\tau}'_k| < \delta_{1n}$ and $|\tilde{\tau}_{k-1} - \tilde{\tau}'_{k-1}| < \delta_{2n}$ for some $\delta_{1n}, \delta_{2n} > 0$.

Hence, we can bound $n\left|V_n(\tilde{\mathcal{T}}_k) - V_n(\tilde{\mathcal{T}}'_k)\right|$ by

$$\begin{aligned} & \left(2 + t'^2 - t^2\right)N_n(t, t') + \left(3 + t'^2\right)\left\{M_n(\tilde{\tau}_k, \tilde{\tau}'_k) + M_n(\tilde{\tau}_{k-1}, \tilde{\tau}'_{k-1})\right\} + \\ & 5\delta_{0n}t_n \sum_{i=1}^n \sigma_i^2 + 2\left(\delta_{1n} + \delta_{2n}\right)\left(1 + t_n^2\right) \sum_{i=1}^n \sigma_i^2 \end{aligned} \quad (3)$$

Following the proof of DJ95P1, we choose $\delta_{0n}, \delta_{1n}, \delta_{2n}$ such that $\delta_{0n}t_n, \delta_{1n}t_n^2$ and $\delta_{2n}t_n^2$ are all $o(n^{-1/2})$ and for large n we use $\mathbb{E} N_n(t, t') + (3 + t_n^2)\{\mathbb{E} M_n(\tilde{\tau}_k, \tilde{\tau}'_k) + \mathbb{E} M_n(\tilde{\tau}_{k-1}, \tilde{\tau}'_{k-1})\} \leq c_0n\delta_{0n} + c_1n\delta_{1n} + c_2n\delta_{2n} = O(r_n n^{1/2})$ for some absolute constants c_0, c_1, c_2 . This and the bound in equation (3) establish $r_n/\sqrt{n} = O(c_n^{-1})$ directly from the proof of DJ95P1 which proves the desired ℓ_1 convergence of equation (2).

Proof of Theorem 1, statement (b) - Due to the result of theorem 1 part (a), proving the result in part (b) essentially reduces to showing $c_n \mathbb{E} \left[\sup_{\mathcal{T} \in \mathcal{H}_n} \left| l_n\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\mathcal{T})\} - r_n(\mathcal{T}, \boldsymbol{\theta}) \right| \right] < \infty$. Note that the loss $l_n\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\mathcal{T})\}$ decomposes as the sum of K losses: $l_n\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\mathcal{T})\} = \sum_{k=1}^K l_n^k\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\mathcal{T}_k)\}$ where $l_n^k\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\mathcal{T}_k)\} = n^{-1} \sum_{i=1}^n \left\{ Y_i + \sigma_i \eta_{t_k}(Y_i) - \theta_i \right\}^2 I\left\{ S_i \in (\tau_{k-1}, \tau_k] \right\}$ and $\mathcal{T}_k = \{\tau_{k-1}, \tau_k, t_k\}$. As the risk is just the expectation of the loss, by triangle inequality, it is enough to show

$$c_n \mathbb{E} \left[\sup_{\mathcal{T}_k \in \mathbf{R}^2 \times [0, t_n]} \left| l_n^k\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\mathcal{T}_k)\} - \mathbb{E} l_n^k(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\mathcal{T}_k)) \right| \right] < \infty \text{ for all } i \text{ and for all large } n$$

Note that each of the losses l_n^k again decomposes into two parts:

$$\begin{aligned} A_n &= n^{-1} \sum_{i=1}^n \left\{ \sigma_i Z_i - \sigma_i t_k \text{sign}(\theta_i + \sigma_i Z_i) \right\}^2 I \left\{ S_i \in (\tau_{k-1}, \tau_k] \right\} I \left\{ |\theta_i + \sigma_i Z_i| > \sigma_i t_k \right\} \\ B_n &= n^{-1} \sum_{i=1}^n \theta_i^2 I \left\{ S_i \in (\tau_{k-1}, \tau_k] \right\} I \left\{ |\theta_i + \sigma_i Z_i| \leq \sigma_i t_k \right\} \end{aligned}$$

where Z_i 's are i.i.d $N(0, 1)$ random variables.

We next prove that for some ϵ_0 , (i) there exist functions \mathfrak{g}_n and \mathfrak{h}_n such that $\mathbb{P} \left\{ c_n \sup_{\mathcal{T}_k} |A_n - \mathbb{E} A_n| > \epsilon \right\} \leq \mathfrak{g}_n(\epsilon)$ and $\mathbb{P} \left\{ c_n \sup_{\mathcal{T}_k} |B_n - \mathbb{E} B_n| > \epsilon \right\} \leq \mathfrak{h}_n(\epsilon)$ for all $\epsilon > \epsilon_0$ and for all large n , and (ii) both $\overline{\lim} \int_{\epsilon_0}^{\infty} \mathfrak{g}_n(\epsilon) d\epsilon < \infty$, $\overline{\lim} \int_{\epsilon_0}^{\infty} \mathfrak{h}_n(\epsilon) d\epsilon < \infty$. This establishes the desired result.

We deal with B_n first and, without loss of generality, establish the bound for

$$B_n = n^{-1} \sum_{i=1}^n \theta_i^2 I \left\{ S_i \in (\tau_{k-1}, \tau_k] \right\} I \left\{ \theta_i + \sigma_i Z_i \leq \sigma_i t_k \right\}. \text{ Now}$$

$$\mathbb{P} \left\{ c_n \sup_{\mathcal{T}_k} |B_n - \mathbb{E} B_n| > \epsilon \right\} = \mathbb{P} \left\{ c_n \sup_{\mathcal{T}_k} |B_n - \mathbb{E} B_n| > \epsilon \text{ and } \mathcal{F}_n \right\} + \mathbb{P} \left(\mathcal{F}_n^c \right) \quad (4)$$

where the set $\mathcal{F}_n = \{ \max_{i=1, \dots, n} |Z_i| \leq (1 + \epsilon) \sqrt{2 \log n} \}$, and $\mathbb{P} \left(\mathcal{F}_n^c \right) \leq \phi(0) n^{-\epsilon}$ for all large n . We bound the first term on the right side of equation (4) by using the Glivenko-Cantelli theorem for weighted empirical measures (Singh, 1975). As $t_n = \sqrt{2 \log n}$ and $t_k \in [0, t_n]$, on \mathcal{F}_n the weights in B_n can be positive only when $\theta_i^2 \leq (2 + \epsilon)^2 \sigma_i^2 2 \log n$. We next use the inequality in equation (6) of Singh (1975) with a in that equation equaling $\epsilon c_n^{-1} \sum_{i=1}^n \sigma_i^4$. Further, note that for all large n , $\epsilon c_n^{-1} \sum_{i=1}^n \sigma_i^4 \geq \sqrt{(2 + \epsilon)^4 (2 \log n)^2 \sum_{i=1}^n \sigma_i^4}$ which is the maximum possible ℓ_2 norm of the weights θ_i^2 . This, along with assumption A1 and the fact that $\sum_{i=1}^n \theta_i^2 \leq \sqrt{n} (\sum_{i=1}^n \theta_i^4)^{1/2}$ gives

$$\mathbb{P} \left\{ \sup_{\mathcal{T}_k} c_n |B_n - \mathbb{E} B_n| > \epsilon, \mathcal{F}_n \right\} < 4 \frac{n \epsilon \log^\delta n}{\sqrt{\sum_{i=1}^n \theta_i^4}} \exp \left\{ - \frac{\epsilon^2 \log^{2(\delta-1)} n}{2(2 + \epsilon)^4} \right\}$$

Now if $\sum_{i=1}^n \theta_i^4 = o(c_n^{-1})$, then the desired bound on $\mathbb{E} \sup_{\mathcal{T}_k} c_n |B_n - \mathbb{E} B_n|$ is obvious; else the above probability is bounded above by $\mathfrak{h}_n(\epsilon) = 4n^2 \epsilon \exp \left\{ - \epsilon^2 \log^{2(\delta-1)} n / 2(2 + \epsilon)^4 \right\}$ which satisfies the aforementioned integrability condition. Thus, the proof of the result for B_n is complete.

We now turn our attention to A_n . Again, without loss of generality, we prove the bound for $A_n = n^{-1} \sum_{i=1}^n (Z_i - t_k)^2 I \left\{ S_i \in (\tau_{k-1}, \tau_k] \right\} I \left\{ \theta_i + Z_i > t_k \right\}$. As we saw in the case of B_n , the variances σ_i^2 appear only through $n^{-1} \sum_{i=1}^n \sigma_i^2$ which is finite by assumption A1. Thus, we take $\sigma_i = 1$ for all i and decompose A_n as sum of three parts by expanding $(Z_i - t_k)^2 = Z_i^2 - 2Z_i t_k + t_k^2$. The bound on the third term follows directly by the traditional Glivenko-Cantelli theorem and by noting that $t_k^2 \leq 2 \log n$. Here we establish the ℓ_1 convergence result for the first term. The proof for the second term is very similar.

We further reduce the problem. Without loss of generality, we assume $\theta_i = 0$ and prove the ℓ_1 convergence result for $A_n = n^{-1} \sum_{i=1}^n Z_i^2 I\{S_i \in (\tau_{k-1}, \tau_k]\} I\{Z_i > t_k\}$. We again apply the same technique as with B_n and control the probability $\mathbb{P}\left\{\sup_{\mathcal{T}_k} c_n |A_n - \mathbb{E} A_n| > \epsilon, \mathcal{F}_n\right\}$. Similarly as with B_n , but now conditioned on $\{Z_i : i = 1, \dots, n\}$, the above probability is easily controlled at the desired rate by applying equation (6) of [Singh \(1975\)](#), i.e., $\mathbb{P}\left\{\sup_{\mathcal{T}_k} c_n |A_n - \mathbb{E} A_n| > \epsilon, \mathcal{F}_n \mid Z_1, \dots, Z_n\right\} \leq \mathfrak{g}_n(\epsilon)$ where \mathfrak{g}_n does not depend on Z_i and for some $\epsilon_0 > 0$, $\int_{\epsilon_0}^{\infty} \mathfrak{g}_n(\epsilon) d\epsilon < \infty$ for all large n with $\sum_{i=1}^n Z_i^2 c_n \rightarrow \infty$ as $n \rightarrow \infty$.

This establishes the desired ℓ_1 result for A_n and completes the proof.

Proof of Corollary 1 - Both statements of this corollary follow from result (b) of Theorem 1. For statement (a), note that for any $\epsilon > 0$, the probability $\mathbb{P}\left[l_n\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\hat{\mathcal{T}})\} \geq l_n\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\mathcal{T}^{OL})\} + c_n^{-1}\epsilon\right]$ is bounded above by $\mathbb{P}\left[l_n\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\hat{\mathcal{T}})\} - S(\hat{\mathcal{T}}, \mathbf{Y}, \mathbf{S}) \geq l_n\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\mathcal{T}^{OL})\} - S(\mathcal{T}^{OL}, \mathbf{Y}, \mathbf{S}) + c_n^{-1}\epsilon\right]$, which converges to 0 by Theorem 1 (b).

Statement (b) of this corollary follows as the difference $l_n\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\hat{\mathcal{T}})\} - l_n\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\mathcal{T}^{OL})\}$ can be decomposed as sum of $l_n\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\hat{\mathcal{T}})\} - S(\hat{\mathcal{T}}, \mathbf{Y}, \mathbf{S})$, $S(\mathcal{T}^{OL}, \mathbf{Y}, \mathbf{S}) - l_n\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\mathcal{T}^{OL})\}$ and $S(\hat{\mathcal{T}}, \mathbf{Y}, \mathbf{S}) - S(\mathcal{T}^{OL}, \mathbf{Y}, \mathbf{S})$. By definition, the last term is not positive. Thus, the sum is bounded above by $2 \sup_{\mathcal{T} \in \mathcal{H}_n} |S(\mathcal{T}, \mathbf{Y}, \mathbf{S}) - l_n\{\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{SI}(\mathcal{T})\}|$ which converges to 0 at the prescribed rate by Theorem 1 (b).

C Detailed proofs of the results of section 3 of the main paper

We use Theorem 2 of [Johnstone \(1994\)](#), which provides an explicit higher order evaluation of the maximal risk of the soft threshold estimator with the best possible choice of threshold. We restate the theorem abet for the symmetric case, which slightly increases the maximal risk presented in equation (17) of the aforementioned theorem.

Result 1. Consider the class of univariate soft-threshold estimators $\hat{\theta}_\lambda^S(x) = \text{sign}(x)(x - \lambda)_+$ for $\lambda \geq 0$. If the parameter $\theta \in \mathbf{R}$ is such that $P(\theta = 0) \leq 1 - \eta$. Then as $\eta \rightarrow 0$, the best choice of threshold is $f(t) + O(t^{-3} \log t)$ and the minimal possible risk is $H(t) = \eta(h(t) + 36 t^{-2} \log t + O(t^{-2}))$ where $t = \sqrt{2 \log \eta^{-1}}$ and

$$f(t) = \sqrt{t^2 - 6 \log t + 2 \log \phi(0)} \text{ and } h(t) = f^2(t) + 5. \quad (5)$$

Proof of Theorem 2 - Directly applying the above result we have,

$$\mathcal{R}_n^{NS}(\alpha, \beta, \rho_n) = (\pi_{1,n} n^{-\alpha} + \pi_{2,n} n^{-\beta}) H(t_n) \bar{\sigma}_n^2, \text{ where, } \bar{\sigma}_n^2 = n^{-1} \sum_{i=1}^n \sigma_i^2,$$

and $t_n^2 = 2 \log(\pi_{1,n} n^{-\alpha} + \pi_{2,n} n^{-\beta})^{-1}$ as the density level is at most $\pi_{1,n} n^{-\alpha} + \pi_{2,n} n^{-\beta}$ in $\Theta_n(\alpha, \beta, \rho_n)$.

Now, if we completely know the latent side information then again applying equation (17) of [Johnstone \(1994\)](#) separately to the two groups: $\{i : \xi_i \leq \tau_n^*\}$ and $\{i : \xi_i > \tau_n^*\}$ we have:

$$\mathcal{R}_n^{OS}(\alpha, \beta, \rho_n) = \{\pi_{1,n} n^{-\alpha} H(t_{1,n}) + \pi_{2,n} n^{-\beta} H(t_{2,n})\} \bar{\sigma}_n^2 \text{ where } t_{1,n}^2 = 2\alpha k_n, t_{2,n}^2 = 2\beta k_n.$$

Also, $t_n^2 = t_{1,n}^2 - 2 \log \pi_{1,n} + O(\pi_{2,n} \pi_{1,n}^{-1} n^{\alpha-\beta})$. By Assumption (A2.1) there exists $\epsilon > 0$ such that $\pi_{2,n} \pi_{1,n}^{-1} n^{\alpha-\beta} < n^{-\epsilon}$ for all large n . Thus, as $n \rightarrow \infty$ with $c_0 = 5 + 2 \log \phi(0)$, and $\tilde{k}_n = k_n / \log k_n$,

$$\begin{aligned} \mathcal{R}_n^{OS} &= \pi_{1,n} p_{1,n} \bar{\sigma}_n^2 \{2\alpha k_n - 3 \log(2\alpha k_n) + c_0 + O(\tilde{k}_n^{-1})\} \\ \mathcal{R}_n^{NS} &= \pi_{1,n} p_{1,n} \bar{\sigma}_n^2 [2\alpha k_n + 2 \log \pi_{1,n}^{-1} - 3 \log(2\alpha k_n) - 3 \log \{1 + (\alpha k_n)^{-1} \log \pi_{1,n}^{-1}\} + c_0 + O(\tilde{k}_n^{-1})], \end{aligned}$$

from which the lemma follows.

To understand the phenomenon here in a simplifier lens, consider the first order approximations:

$$\mathcal{R}_n^{NS} \sim \bar{\sigma}_n^2 \pi_{1,n} p_{1,n} f(t_n), \mathcal{R}_n^{OS} \sim \bar{\sigma}_n^2 \pi_{1,n} p_{1,n} f(t_{1,n}), f(t_n) - f(t_{1,n}) \sim 2 \log \pi_{1,n}^{-1} \text{ and}$$

$$\mathcal{R}_n^{NS} - \mathcal{R}_n^{OS} \sim (2 \log \pi_{1,n}^{-1}) \pi_{1,n} p_{1,n} \bar{\sigma}_n^2$$

Thus, the gain due to incorporation of side information is essentially due to the fact that we can use a lower threshold for the subgroup with smaller sparsity than that used in the agglomerative case with no side information and this is exactly the phenomenon depicted in Figure 2 of the main paper.

Proofs of Theorems 3, 4 and Lemma 2 - Note that in our asymptotic set-up, there exists $\epsilon > 0$ such that

$$p_{1,n} \pi_{1,n} (p_{2,n} \pi_{2,n})^{-1} \geq n^\epsilon \text{ for all large } n. \quad (6)$$

We will be using this property to simplify our calculations by restricting ourselves to dominant terms.

As such we will be ignoring terms which are $o(p_{1,n} \pi_{1,n} \bar{\sigma}_n^2 k_n^{-1})$. Without loss of generality we assume that $S|\xi$ has monotone likelihood ratio in S and consider $q_{i,n}^{jk}(\tau) := \mathbb{P}_n(\hat{I}_i^j | I_i^k)$ for $j, k \in \{1, 2\}$, $i \in \{1, \dots, n\}$, where, $\hat{I}_i^1 = \{S_i \leq \tau\}$, $I_i^1 = \{\xi_i \leq \tau_n^*\}$, $\hat{I}_i^2 = \mathbf{R} \setminus \hat{I}_i^1$, $I_i^2 = \mathbf{R} \setminus I_i^1$ and $q_n^{jk}(\tau) = n \sum_{i=1}^n q_{i,n}^{jk}(\tau) \sigma_i^2 / \bar{\sigma}_n^2$.

Proof of Lemma 2 - Dividing the difference in Theorem 2 by the expression of \mathcal{R}_n^{OS} from the display

above it we get

$$\mathcal{R}_n^{NS}/\mathcal{R}_n^{OS} = 1 + [2 \log \pi_{1,n}^{-1} - 3 \log \{1 + (\alpha k_n)^{-1} \log \pi_{1,n}^{-1}\}] / \{2\alpha k_n - 3 \log(2\alpha k_n) + c_0\} + O(k_n^{-\nu}), \quad (7)$$

where $c_0 = 5 + 2 \log \phi(0)$ and $\nu < 2$. The first result of the lemma now follows by noting $\log \pi_{1,n}^{-1} \leq \gamma_0 < \beta - \alpha$ which is due to Assumption (A2.1).

Note that, iff $c_n \rightarrow 0$ then $\mathcal{R}_n^{NS}/\mathcal{R}_n^{OS} \rightarrow 1$. From the above display it follows that $k_n^{1+\delta}(\mathcal{R}_n^{NS}/\mathcal{R}_n^{OS} - 1) \geq k_n^\delta \log \pi_{1,n}^{-1} \{1 - 1.5(\alpha k_n)^{-1}\} + O(k_n^{-\nu+1+\delta})$ where $\nu < 2$. Thus, we have the first part of the third result. Its second part follows directly from the proof of Theorem 3, which is present after the proof of this lemma.

Next, we establish the upper bound on the maximal risk of ASUS given in the second statement of lemma 2. Let \mathcal{R}_n^{KS} denotes the maximal risk of ASUS when we can set any possible thresholds in ASUS including those depending on the density levels $p_{1,n}, p_{2,n}$ as well as the mixing probabilities $\pi_{1,n}$ and $\pi_{2,n}$. However, we do not know the latent variable or its subsequent oracle optimal groups $\{i : \xi_i \leq \tau_n^*\}$ and $\{i : \xi_i > \tau_n^*\}$. Thus, by definition $\mathcal{R}_n^{OS} \leq \mathcal{R}_n^{KS} \leq \mathcal{R}_n^{NS}$. Now, ASUS always chooses the thresholds and the segmentation hyper-parameter in a data-dependent fashion minimizing the SURE criterion. We next apply theorem 1 which tells us that the maximal risk of ASUS \mathcal{R}_n^{AS} can not be much bigger than \mathcal{R}_n^{KS} . As such, theorem 1 compounded with theorem 4a of DJ95 implies $\mathcal{R}_n^{AS} - \mathcal{R}_n^{KS} \leq \mathcal{R}_n^F I\{\mu_n^2 \leq 3d_n\} + o(\pi_{1,n} p_{1,n} \sigma_n^2 k_n^{-1})$ where \mathcal{R}_n^F is the risk of ASUS with fixed threshold of $\sqrt{2k_n}$ and $d_n = n^{-1/2} \log^{3/2} n$ and $\mu_n^2 = n^{-1} \sum_{i=1}^n \theta_i^2 \wedge (2k_n)$. By Lemma 8.3 of [Johnstone \(2015\)](#), $\mathcal{R}_n^F \leq n^{-1} + n^{-1} \sum_{i=1}^n \{\theta_i^2 \wedge (1 + 2k_n)\} \sigma_i^2$. Thus, $\mathcal{R}_n^{AS} \leq \mathcal{R}_n^{NS} + o(\pi_{1,n} p_{1,n} \sigma_n^2 k_n^{-1})$ and the result follows from (7).

Proof of Theorem 3 - First consider the situation where the sparsity levels $p_{1,n}$ and $p_{2,n}$ are known. Due to the product structure of our ASUS estimator, we first concentrate on its maximal risk for each of the i th coordinate. This reduces to an univariate risk analysis. If noise variance equals 1, univariate soft-threshold estimators with threshold λ has:

- (a) the risk at the origin: $g_1(\lambda)(1 + O(\lambda^{-2}))$ for large λ where $g_1(\lambda) = 4\phi(\lambda)/\lambda^3$
- (b) the maximal risk at the non-origin points: $g_2(\lambda) = 1 + \lambda^2$ and the maximum is attained when the parametric value is $\pm\infty$.

Now, if the probability of the parameter θ_i being non-zero is p then the maximal risk of the soft-threshold estimator with threshold λ is $g(p, \lambda) = (1 - p)g_1(\lambda) + pg_2(\lambda)$.

As θ_i is generated from the two group model of equation (8) of the main paper with density levels

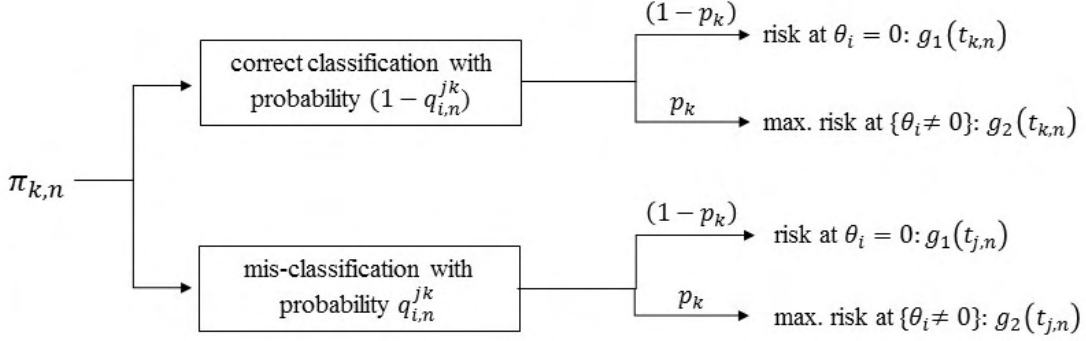


Figure 1: Pictorial representation of coordinate-wise decomposition of maximal risk \mathcal{R}^{AS} . Here $j, k = 1, 2$ and $j \neq k$.

$p_{k,n}$, $q_{i,n}^{12}$ and $q_{i,n}^{21}$ are the probabilities of mis-classifying group 2 and group 1 respectively and thresholds $t_{k,n}$ were used for those detected in group $k = 1, 2$. Note, that without mis-classification the maximal risk at each coordinate i is weighted by the group probabilities $\pi_{k,n}$ and the optimal threshold choices are $f(\sqrt{2\alpha k_n})$ and $f(\sqrt{2\beta k_n})$ where, f is defined in (5). However, under mis-classification these thresholds will change and the optimal thresholds will be $f(m_{1,n}^{\text{opt}}[i])$ and $f(m_{2,n}^{\text{opt}}[i])$ where,

$$m_{k,n}^{\text{opt}}[i] = \left\{ -2 \log \left(\frac{\pi_{k,n} q_{i,n}^{kk} p_{k,n} + \pi_{j,n} q_{i,n}^{kj} p_{j,n}}{\pi_{k,n} q_{i,n}^{kk} + \pi_{j,n} q_{i,n}^{kj}} \right) \right\}^{1/2} \quad \text{for } j \neq k. \quad (8)$$

These can not be used as $q_{i,n}^{jk}$ are not known while constructing the estimator. We are interesting in deriving upper bounds on the maximal risk of ASUS and so, unlike the optimal thresholds which depend on i , here we consider thresholds $t_{1,n}$ and $t_{2,n}$ which are uniform over the groups. With mis-classification, the probabilities $q_{i,n}^{jk}$ will be also involved into the expression for maximal risk as now θ_i coming from group 1 (say) might be treated with either threshold $t_{1,n}$ (when correctly classified) or $t_{2,n}$ (when incorrectly classified). Figure 1 provides a pictorial representation of how the probabilities $q_{i,n}^{jk}$ enter this decomposition. The maximal risk for coordinate i is given by

$$\sum_{j=1}^2 \sum_{k=1}^2 \pi_{k,n} q_{i,n}^{jk} g(p_{k,n}, t_{j,n}) \sigma_i^2 \{1 + o(1)\}. \quad (9)$$

We fix a threshold of $t_{1,n} = f(\sqrt{2\alpha k_n})$ and $t_{2,n} = f(\sqrt{2\gamma k_n})$ where γ is allowed to vary in $[\alpha, \beta]$. These thresholds when substituted in (9) will produce an upper bound on the maximal risk. Doing so and

using (6), we see that the maximal risk for the i th coordinate is upper bounded by

$$\sigma_i^2 \left[\pi_{1,n} p_{1,n} \left\{ q_{i,n}^{11} g_2(t_{1,n}) + q_{i,n}^{21} g_2(t_{2,n}) \right\} + \sum_{j=1}^2 \sum_{k=1}^2 \pi_{k,n} q_{i,n}^{jk} g_1(t_{j,n}) + O\left(\frac{\pi_{1,n} p_{1,n}}{k_n}\right) \right] \quad (10)$$

Now, consider the second term in equation (10). As $t_{2,n} \geq t_{1,n}$, we lower bound and upper bound it by $A_n + \pi_{1,n} q_{i,n}^{11} g_1(t_{1,n})$ and $A_n + \pi_{1,n} g_1(t_{1,n})$ where $A_n = \pi_{2,n} q_{i,n}^{22} g_1(t_{2,n}) + \pi_{2,n} q_{i,n}^{12} g_1(t_{1,n})$.

Define $\tilde{k}_n = k_n / \log k_n$. Note that $g_1(t_{1,n}) = 4p_{1,n} \{1 + O(\tilde{k}_n^{-1})\}$ and $g_1(t_{2,n}) = 4\phi(0)t_{2,n}^{-3}\phi(t_{2,n})(1 + O(k_n^{-1}))$. Thus, with $\kappa_n(\gamma) = n^\alpha \phi(0)\phi(t_{2,n})t_{2,n}^{-3}$, $t_{2,n} = f(\sqrt{2\gamma k_n})$ the second term in (10) is bounded above by

$$4\pi_{2,n} p_{2,n} \{q_{i,n}^{22} \rho_n \kappa_n(\gamma) + \rho_n q_{i,n}^{12} + 1\} \{1 + O(\tilde{k}_n^{-1})\} \quad (11)$$

as $\rho_n = \pi_{2,n} / \pi_{1,n}$. Now, consider the first term in equation (10). We have,

$$q_{i,n}^{11} g_2(t_{1,n}) + q_{i,n}^{21} g_2(t_{2,n}) = g_2(t_{1,n}) + q_{i,n}^{21} \{g_2(t_{2,n}) - g_2(t_{1,n})\}$$

and $g_2(t_{2,n}) - g_2(t_{1,n}) = 2(\gamma - \alpha)k_n - 3 \log \log(\gamma/\alpha) := \delta_n(\gamma)$. Thus, the first term in equation (10),

$$\pi_{1,n} p_{1,n} \left\{ q_{i,n}^{11} g_2(t_{1,n}) + q_{i,n}^{21} g_2(t_{2,n}) \right\} = \pi_{1,n} p_{1,n} \left\{ g_2(t_{1,n}) + q_{i,n}^{21} \delta_n(\gamma) \right\} \quad (12)$$

Now, $g_2(t_{1,n}) = 1 + t_{1,n}^2 = h(\sqrt{2\alpha k_n}) - 4$, and so, from equations (11) and (12), maximal risk of ASUS for coordinate i is upper bounded by

$$\pi_{1,n} p_{1,n} \sigma_i^2 \left[h(\sqrt{2\alpha k_n}) + \{q_{i,n}^{21} \delta_n(\gamma) + 4\rho_n q_{i,n}^{12} + 4q_{i,n}^{22} \rho_n \kappa_n(\gamma)\} \{1 + O(\tilde{k}_n^{-1})\} + O\left(\tilde{k}_n^{-1}\right) \right]$$

and therefore the maximal risk over the n coordinates of the ASUS estimator when thresholds can be directly chosen depending on the density levels $p_{1,n}$ and $p_{2,n}$ is

$$\mathcal{R}_n^{KS} \leq \pi_{1,n} p_{1,n} \sigma_n^2 \left[h(\sqrt{2\alpha k_n}) + \{q_n^{21} \delta_n(\gamma) + 4\rho_n q_n^{12} + 4q_n^{22} \rho_n \kappa_n(\gamma)\} \{1 + O(\tilde{k}_n^{-1})\} + O\left(\tilde{k}_n^{-1}\right) \right]$$

where $q_n^{jk}(\tau) = \sum_{i=1}^n q_{i,n}^{jk}(\tau) \sigma_i^2 / \sum_{i=1}^n \sigma_i^2$ for $j, k \in \{1, 2\}$. Recall, \mathcal{R}_n^{KS} denotes the maximal risk of ASUS when we can set any possible thresholds in ASUS including those depending on the density levels $p_{1,n}, p_{2,n}$ as well as the mixing probabilities $\pi_{1,n}$ and $\pi_{2,n}$. Now, consider the general case when those are unknown. ASUS always chooses the thresholds $t_{1,n}$ and $t_{2,n}$ and the segmentation hyper-parameter τ_n in a data-dependent fashion minimizing the SURE criterion. We next apply theorem 1 similarly as in

the proof of lemma 2 which provides us with $\mathcal{R}_n^{AS} - \mathcal{R}_n^{KS} \leq o(\pi_{1,n} p_{1,n} \bar{\sigma}_n^2 k_n^{-1})$ Also, from calculation in the previous subsection for any $\nu < 1$, we know $\mathcal{R}_n^{OS} \geq \pi_{1,n} p_{1,n} \bar{\sigma}_n^2 \{h(\sqrt{2\alpha k_n}) + o(k_n^{-\nu})\}$ and so

$$\mathcal{R}_n^{AS} - \mathcal{R}_n^{OS} \leq \pi_{1,n} p_{1,n} \bar{\sigma}_n^2 \left[\{q_n^{21} \delta_n(\gamma) + 4\rho_n q_n^{12} + 4q_n^{22} \rho_n \kappa_n(\gamma)\} \{1 + O(\tilde{k}_n^{-1})\} + O(\tilde{k}_n^{-1}) \right]. \quad (13)$$

Again, in our asymptotic set-up there exists $\epsilon > 0$ such that $\rho_n \leq n^{\gamma_0 - \epsilon}$ for all large n . Choosing $\gamma = \alpha + \gamma_0$, we have $\delta_n(\gamma) = 2\gamma_0 k_n - 3 \log \log(1 + \gamma_0/\alpha)$ and $\rho_n \kappa_n = o(n^{-\epsilon'})$ for some $0 < \epsilon' < \epsilon$. Thus, with this choice of γ , based on (13) the controls on q_n^{12} and q_n^{21} stated in the theorem implies $\mathcal{R}_n^{AS} - \mathcal{R}_n^{OS} \leq O(\tilde{k}_n^{-1})$. This, along with theorem 2 provide us with the desired result for the theorem as well as the third result of lemma 2.

Proof of Theorem 4, statement (a) - Consider case (ii) first. Note that $\lim_{n \rightarrow \infty} n \rho_n q_n^{21}(\tau_n) = 0$ implies $\rho_n q_{i,n}^{21}(\tau_n) = o(1)$ as $n \rightarrow \infty$ for all i . Hence for each coordinate i , the optimal threshold for group 1 considering two groups in the data is $f(m_{1,n}^{\text{opt}}[i])$ (see equation (8)). As $\rho_n q_{i,n}^{21}(\tau_n) = o(1)$, we have $m_{1,n}^{\text{opt}} = \sqrt{2\alpha k_n} \{1 + o(k_n^{1-\epsilon})\}$ for any $\epsilon > 0$. Thus, the threshold here asymptotically equals $t_{1,n}$ used in the proof of theorem 3: part b, before. Concentrating on only the $j = 1, k = 1$ and $j = 2, k = 1$ terms in equation (9) by the previously conducted analysis we have the maximal risk of the i th coordinate bounded below by $\pi_{1,n} p_{1,n} \sigma_i^2 \left[h(\sqrt{2\alpha k_n}) + 4\rho_n q_{i,n}^{12} \{1 + O(\tilde{k}_n^{-1})\} \right]$. Thus, if ASUS considers two groups then $\mathcal{R}_n^{KS} \geq \pi_{1,n} p_{1,n} \bar{\sigma}_n^2 [h(\sqrt{2\alpha k_n}) + 4\rho_n q_{i,n}^{12} \{1 + O(\tilde{k}_n^{-1})\}]$ and the ratio $(\mathcal{R}_n^{NS} - \mathcal{R}_n^{OS})/(\mathcal{R}_n^{KS} - \mathcal{R}_n^{OS})$ diverges to ∞ unless $\overline{\lim} \rho_n q_n^{12} k_n^{-1} < \infty$ as $\log \pi_{1,n}^{-1} = \gamma_0 k_n$. In that case we use a uniform choice of threshold $t_{1,n} = t_{2,n} = f(m_n^{\text{opt}})$ where $m_n^{\text{opt}} = \{-2 \log(\pi_{1,n} p_{1,n} + \pi_{2,n} p_{2,n})\}^{1/2}$. This along with part (a) completes the proof for case (ii).

For case (i), note that the $m_{1,n}^{\text{opt}}$ for the i th coordinate defined in (8) equals

$$m_{1,n}^{\text{opt}}[i] = \{2 \log p_{1,n}^{-1} + 2 \log(1 + \rho_n q_{i,n}^{21}/q_{i,n}^{11})\}^{1/2} \{1 + o(1)\} \text{ as } n \rightarrow \infty.$$

Now, considering the risk at the origin for $j = 1, k = 1$ term in equation (9), we see that it will contain at least an extra additive component of $O(\pi_{1,n} p_{1,n} \rho_n q_{i,n}^{21}/q_{i,n}^{11} \sigma_i^2)$ over $h(\sqrt{2\alpha k_n})$. Thus, the average maximal risk over the n coordinates is bounded below by

$$\pi_{1,n} p_{1,n} \left\{ h(\sqrt{2\alpha k_n}) \bar{\sigma}_n^2 + O\left(\rho_n n^{-1} \sum_{i=1}^n q_{i,n}^{21}/q_{i,n}^{11} \sigma_i^2 \right) \right\}$$

As $q_{i,n}^{11} \leq 1$ for all i , the second term on right side above is bounded below by $O(\pi_{1,n} p_{1,n} \rho_n q_n^{21} \bar{\sigma}_n^2)$ which provides us with the desired result.

Proof of Theorem 4, statement (b) - By definition, $\mathcal{R}_n^{KS} \leq \mathcal{R}_n^{NS}$ and by application of theorem 1 (as

in the proof of lemma 2) we have $\mathcal{R}_n^{AS} - \mathcal{R}_n^{KS} \leq o(\pi_{1,n} p_{1,n} \bar{\sigma}_n^2 k_n^{-1})$. Also, by assumption A2.2, there exists some $\nu < 1$ such that $\lim k_n^{\nu/2} (1 - \pi_1) = \infty$. This implies $k_n^\nu \log \pi_{1,n}^{-2}$, which coupled with theorem 2 gives us with the desired result.

Proof of Lemma 1 - Without loss of generality, assume that the marginal distribution of the auxiliary sequence S given the latent parameter ξ has monotone likelihood ratio in the statistic S . Also, let $\underline{\mu}_{2,n} = \inf\{\mu_i : i \in I_{2,n}^*\} > \sup\{\mu_i : i \in I_{1,n}^*\} = \bar{\mu}_{1,n}$. Then, $\underline{\mu}_{2,n} = \bar{\mu}_{1,n} + d_n$. Under this reduction $q_{i,n}^{jk}(\tau) := \mathbb{P}_n(\hat{I}_i^j | I_i^k)$ for $j, k \in \{1, 2\}$, $i \in \{1, \dots, n\}$, where, the sets $\hat{I}_i^1 = \{S_i \leq \tau\}$, $I_i^1 = \{\xi_i \leq \tau_n^*\}$, $\hat{I}_i^2 = \mathbf{R} \setminus \hat{I}_i^1$, $I_i^2 = \mathbf{R} \setminus I_i^1$. Let $\bar{b}_n = \sup_i \max(2\nu_i^2, b_i)$. Now, set $\tau_n = \bar{\mu}_{1,n} + \bar{b}_n^\gamma (2k_n + \log k_n)^\gamma$. Thus, for all large n the condition imposed on d_n in the lemma implies $\tau_n \leq \underline{\mu}_{2,n} - \bar{b}_n^\gamma (\log \rho_n + \log k_n)^\gamma$. Next, note that, $q_n^{21}(\tau_n) \leq \bar{\sigma}_n^2 \sup_i q_{i,n}^{21}(\tau_n) \leq \bar{\sigma}_n^2 \sup_i P(S_i \leq \tau_n | \bar{\mu}_{1,n})$ and thereafter, using the tail bounds of the sub-Exponential/Gaussian distributions, we have for all i ,

$$P(S_i \leq \tau_n | \bar{\mu}_{1,n}) \leq \exp\{-(\tau_n - \bar{\mu}_{1,n})^{1+I(b_i=0)} / \max(2\nu_i^2, b_i)\} \leq k_n^{-2} / \log k_n$$

for all large n . Similarly, it follows that $q_n^{12}(\tau_n) \leq \rho_n^{-1} / \log k_n$, which completes the proof.

Statement and proof of Lemma 3 -

Lemma 3 - *If our parametric space $\Theta(\alpha, \beta, \rho_n)$ is such that $\overline{\lim}_{n \rightarrow \infty} n^{1/2} \pi_{1,n} p_{1,n} < \infty$, ASUS convergences in probability to the SureShrink procedure with the fixed threshold choice of $\sqrt{2 \log n}$.*

Proof. Define, $\mu_n^2 = n^{-1} \sum_{i=1}^n \theta_i^2 \wedge (5k_n)$ for some prefixed $\epsilon > 0$. Note that $\mu_n^2 \leq 5\pi_{1,n} p_{1,n} k_n (1 + o(1))$ and so, by the condition of the lemma $\mu_n^2 / d_n \rightarrow 0$ where $d_n = n^{-1/2} \log^{3/2} n$. Define $s_n^2 = n^{-1} \sum_{i=1}^n y_i^2 \wedge (2k_n) - 1$ where $Y_i = \theta_i + Z_i$ and Z_i are i.i.d. $N(0, 1)$. Let $F_n = \{\max_i z_i^2 < 3k_n\}$. Note that $P(s_n^2 \leq d_n) \leq P(F_n^c) + P(s_n^2 \leq d_n \text{ and } F_n)$. The first term converge to 0 and the second term on the right side above is bounded by $P(n^{-1} \sum_{i=1}^n Y_i^2 - 1 \leq d_n, n^{-1} \sum_{i=1}^n \theta_i^2 \leq \mu_n^2)$ which converges to 0 as $\mu_n^2 / d_n \rightarrow 0$ (see proof of theorem 4 (b) of DJ95). This completes the proof.

D Additional simulation experiments

In this section, we present a number of simulation experiments demonstrating the asymptotic performance of ASUS as n increases. We fix $m = 50$ and allow n to vary from 500 to 5000 in increments of 100. To simulate the parameter vector θ , we continue to use the setup of the one sample problem discussed in Section 4.2.1 of the main paper and simulate η_{1i} as before but vary the sparsity levels in θ under scenarios S1 and S2 as follows:

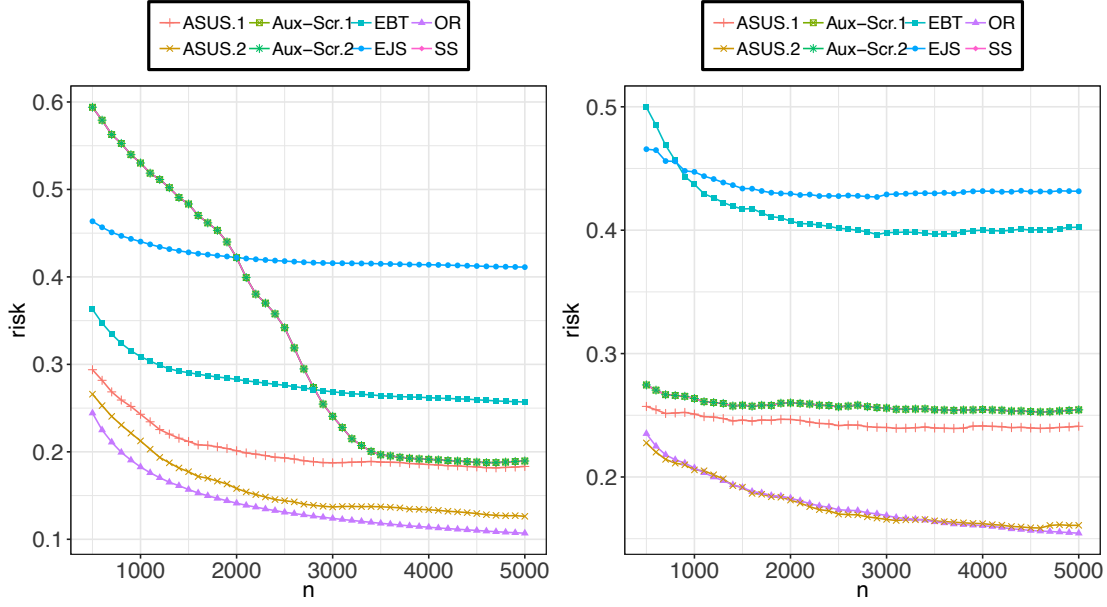


Figure 2: Asymptotic performance of ASUS: Average risks of different estimators. Dashed line represents the risk of the oracle estimator $\hat{\theta}_i^{SI}(\mathcal{T}^{OR})$. Left: Scenario S1 and Right: Scenario S2.

$$(S1) \quad \xi \sim \left(\underbrace{\text{Unif}(6, 7)}_{1\% \text{ of } n}, \underbrace{\text{Unif}(2, 3)}_{4\% \text{ of } n}, \underbrace{0, \dots, 0}_{n - 5\% \text{ of } n} \right)$$

$$(S2) \quad \xi \sim \left(\underbrace{\text{Unif}(4, 8)}_{4\% \text{ of } n}, \underbrace{\text{Unif}(1, 3)}_{16\% \text{ of } n}, \underbrace{0, \dots, 0}_{n - 20\% \text{ of } n} \right)$$

with $\theta = \xi + \eta_1$, $Y_i \sim N(\theta_i, \sigma_i^2)$ where $\sigma_i^2 = 1$ for all i under S1 and $\sigma_i^2 \stackrel{i.i.d}{\sim} \text{Unif}(0.1, 1)$ under S2. For scenario S1, we consider two side information \mathcal{S}_1 and \mathcal{S}_2 as follows: (ASUS.1) $S_{i1}|i \in \mathcal{I}_1^* = |N(\mu_0, \sigma_i^2) + \bar{\eta}_{2i}|$, $S_{i1}|i \in \mathcal{I}_2^* = |N(\mu_1, \sigma_i^2) + \bar{\eta}_{2i}|$ with $\mu_0 = \sqrt{\log k_n}$, $\mu_1 = 0$, and (ASUS.2) $S_{i2}|i \in \mathcal{I}_1^* = |N(\mu_0, \sigma_i^2) + \bar{\eta}_{2i}|$, $S_{i2}|i \in \mathcal{I}_2^* = |N(\mu_1, \sigma_i^2) + \bar{\eta}_{2i}|$ with $\mu_0 = \sqrt{k_n}$, $\mu_1 = 0$. Here $\sigma_i^2 \stackrel{i.i.d}{\sim} \text{Unif}(0.1, 1)$ and $\bar{\eta}_{2i}$ is the average over m samples of $N(0, 0.01)$. Similarly for scenario S2, we consider two side information \mathcal{S}_1 and \mathcal{S}_2 as follows: (ASUS.1) $S_{i1}|i \in \mathcal{I}_1^* = |\chi_{1+\sqrt{\log k_n}}^2 + \bar{\eta}_{2i}|$, $S_{i1}|i \in \mathcal{I}_2^* = |\chi_1^2 + \bar{\eta}_{2i}|$, and (ASUS.2) $S_{i2}|i \in \mathcal{I}_1^* = |\chi_{1+k_n}^2 + \bar{\eta}_{2i}|$, $S_{i2}|i \in \mathcal{I}_2^* = |\chi_1^2 + \bar{\eta}_{2i}|$. Thus \mathcal{S}_1 and \mathcal{S}_2 differ in the separation of the means of their conditional distributions with \mathcal{S}_2 in scenario S2 having a near optimal separation in the means as prescribed by Lemma 1 in Section 3.3.

We repeat this sampling scheme for $N = 1000$ repetitions and report the results in table 1 and figure 2. As observed in the one sample estimation problem, ASUS continues to exhibit the best performance amongst all the competing estimators. However, the efficiency of ASUS in exploiting side information clearly depends on the magnitude of separation of the conditional means of \mathcal{S} under the two groups. For example ASUS.2, which has a bigger separation between the conditional means of side information \mathcal{S}_2 , exhibits the best performance across all scenarios. In fact under scenario S2, \mathcal{S}_2 is sub-exponential

Table 1: Asymptotic performance of ASUS: risk estimates and estimates of \mathcal{T} for ASUS at $n = 5000$. Here $n_k^* = |\mathcal{I}_k^{\tau^*}|$ and $n_k = |\widehat{\mathcal{I}}_k^\tau|$ for $k = 1, 2$.

Asymptotic performance of ASUS			
		Scenario S1	Scenario S2
OR	τ^*	2.00	0.980
	t_1^*, t_2^*	4.115, 0.130	4.073, 0.062
	n_1^*, n_2^*	4750, 250	4000, 1000
	risk	0.107	0.154
ASUS.1	τ	1.936	1.904
	t_1, t_2	1.253, 3.520	0.740, 1.281
	n_1, n_2	3460, 1540	2857, 2143
	risk	0.183	0.241
ASUS.2	τ	1.60	2.918
	t_1, t_2	0.420, 4.104	0.139, 1.8
	n_1, n_2	446, 4554	1024, 3976
	risk	0.126	0.161
Aux-Scr.1	τ	5.920	34.218
	t_1, t_2	0.424, -	0,-
	n_1, n_2	5000, 0	5000,0
	risk	0.189	0.254
Aux-Scr.2	τ	7.375	50.958
	t_1, t_2	0.424, -	0,-
	n_1, n_2	5000, 0	5000, 0
	risk	0.189	0.254
SureShrink	risk	0.189	0.254
EBT	risk	0.257	0.402
EJS	risk	0.411	0.431

and the $\log n$ separation between the conditional means brings the risk of ASUS.2 closer to the risk of the oracle estimator. As far as ASUS.1 is concerned, the relatively smaller separation in the conditional means of \mathcal{S}_1 does not allow ASUS to optimally partition the n coordinates into heterogeneous groups in terms of sparsity and hence it performs only marginally better than the SureShrink estimator. Moreover, we observe that across both the scenarios, the risk of the auxiliary screening procedure, Aux-Scr, is indistinguishable from the risk of the SureShrink estimator, thus demonstrating that discarding observations based on the magnitude of the side information may lead to missing important signal coordinates especially when the side information is corrupted with noise. ASUS, however, does not discard any observations and continues to exploit the available information in the noisy auxiliary sequence. In table 1, we report risk estimates and estimates of \mathcal{T} for ASUS and Aux-Scr at $n = 5000$.

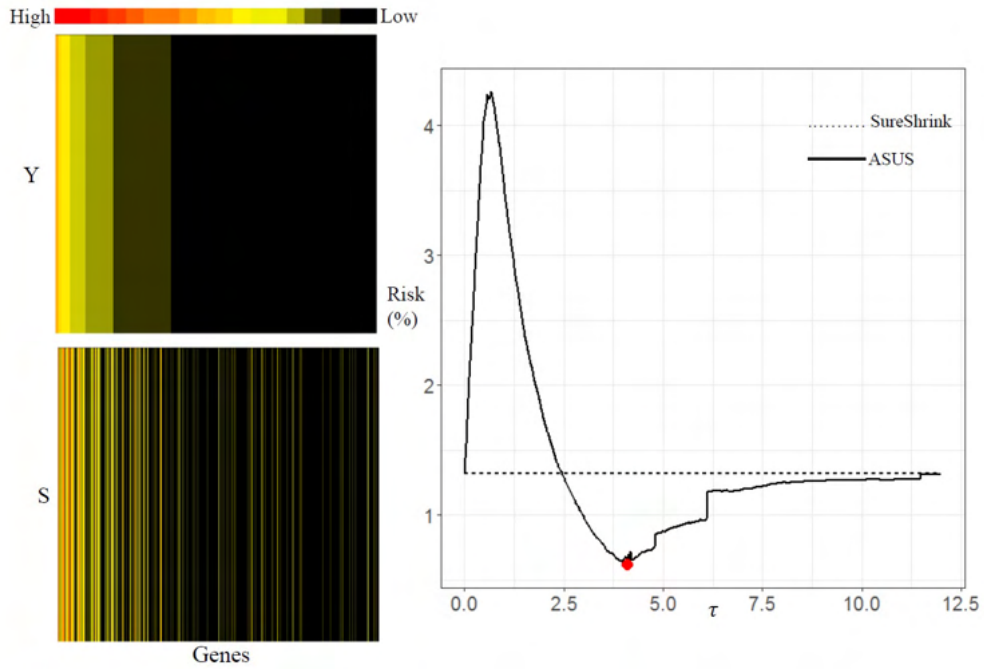
In table 2, we report the estimates of the hyper-parameters of Aux-Scr under the setting of the simulation experiment described in section 4.2.1 of the main paper.

Table 2: Risk estimates and estimates of \mathcal{T} for Aux-Scr at $n = 5000$ under the setting of the simulation experiment described in section 4.2.1 of the main paper. Here $n_k^* = |\mathcal{I}_k^{\tau^*}|$ and $n_k = |\mathcal{I}_k^{\tau}|$ for $k = 1, 2$.

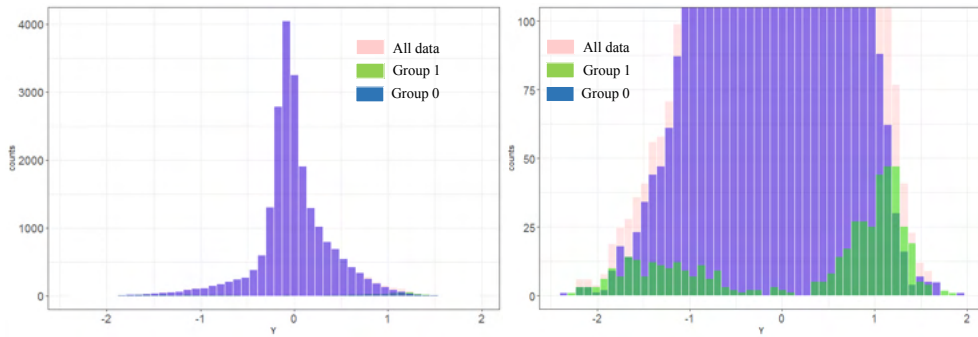
Asymptotic performance of Aux-Scr			
		Scenario S1	Scenario S2
Aux-Scr.1	τ	1.326	0.981
	t_1, t_2	4.612, 0.109	4.618, 0.159
	n_1, n_2	4747, 253	4002, 998
	risk	0.097	0.243
Aux-Scr.2	τ	11.231	10.752
	t_1, t_2	4.601, 0.106	4.580, 0.160
	n_1, n_2	4748, 252	3976, 1024
	risk	0.095	0.232
Aux-Scr.3	τ	1.774	1.460
	t_1, t_2	4.668, 0.665	4.760, 0.544
	n_1, n_2	4261, 739	3290, 1710
	risk	0.147	0.360
Aux-Scr.4	τ	5.072	2.004
	t_1, t_2	3.036, 1.043	3.500, 0.851
	n_1, n_2	2023, 2977	984, 4016
	risk	0.186	0.414

E Microarray Time Course (MTC) Data

Our second real data application is an MTC dataset collected by [Calvano et al. \(2005\)](#) for studying systemic inflammation in humans and is an example of a setting where ASUS can be used for 2 sample estimation problems. This dataset contains eight study subjects which are randomly assigned to case and control groups and then administered with endotoxin and placebo, respectively. The expression levels of $n = 22,283$ genes in human leukocytes are measured before infusion (0 hour) and at 2, 4, 6, 9, and 24 hours afterwards. One of the goals of this experiment is to identify, in the case group, early to middle response genes that are differentially expressed within 4 hours and thus reveal meaningful early activation gene sequence that governs the immune responses. As discussed in [Sun and Wei \(2011\)](#) the early activation sequence quickly activates many secreted pro-inflammatory factors in response to exterior intrusion. These activated factors subsequently trigger the expression of several transcription factors to initiate the immune response. In the late period, the expression levels of a number of transcription factors limiting the immune response are increased. Finally the whole system concludes with full recovery and a normal phenotype. To identify the genes that regulate this sequence, we take time point 0 as the baseline and time points 4 and 24 as the interval over which differential gene expression will be estimated. We follow the data preprocessing steps outlined in [Sun and Wei \(2011\)](#) and denote $Y_{i,j}$ as the arcsinh transformed average gene expression value for gene i at time point j . Let $Y_i = \tilde{Y}_{i,4} - \tilde{Y}_{i,24}$ where $\tilde{Y}_{i,j} = Y_{i,j} - Y_{i,0}$ denotes the baseline adjusted expression level of gene i at time point j . The side information that we use in this setting is $X_i = \tilde{Y}_{i,4} + \kappa_i \tilde{Y}_{i,24}$ with $S_i = |X_i|$, $K = 2$, $\kappa_i = \tilde{\sigma}_{i,4}/\tilde{\sigma}_{i,24}$



(a)



(b)

Figure 3: (a) Left: Heatmap of gene expressions \mathbf{Y} and side information \mathbf{S} . Right: SURE estimate of the risk of $\hat{\theta}_i^S(t)$ at $t = 1.13$ versus an unbiased estimate of the risk of ASUS for different values of τ . (b) Left: Histogram of gene expressions \mathbf{Y} . Group 1 is $\hat{\mathcal{I}}_2^\tau$ and Group 0 is $\hat{\mathcal{I}}_1^\tau$. Right: A magnified plot to show $\hat{\mathcal{I}}_2^\tau$.

where $\tilde{\sigma}_{i,j}$ is the observed standard deviation of $\tilde{Y}_{i,j}$ across the 4 replicates. In figure 3a right, the dotted line represents the SURE estimate of risk of $\hat{\theta}^S(t)$ at $t = 1.13$. ASUS uses information in (Y_i, S_i) and returns an estimate of risk (the red dot) that is significantly smaller than the risk estimate returned by $\hat{\theta}^S(t)$. In order to evaluate the results in a predictive framework, we next used 3 out of the 4 replicates for calibrating the hyper-parameters and calculated the prediction errors of the ASUS and SureShrink procedures based on the held out fourth replicate. Here, the risk reduction by ASUS compared to SureShrink is almost 9%

Figure 3a left presents heatmap of expression level Y_i in the top panel and heatmap of the associated side information S_i in the bottom panel for gene i . The expression levels for the genes are ordered in terms of their magnitude. Notice how the magnitude of side information in S_i follows the pattern of the expression levels Y_i , largely indicating that Y_i is small whenever S_i is small. ASUS exploits this extra information in S_i and thus performs better than the SureShrink estimator that only relies on the information in Y_i . Figure 3b left presents the distribution of gene expression for genes that belong to the groups $\hat{\mathcal{I}}_1^\tau$ and $\hat{\mathcal{I}}_2^\tau$. In this example, group $\hat{\mathcal{I}}_2^\tau$ holds only about 2% of the n genes and is therefore inconspicuous in this plot. We present a magnified version of this plot in right that demonstrates in green the distribution of gene expression for genes that belong to $\hat{\mathcal{I}}_2^\tau$ and summarize the results in table 3. In

Table 3: Summary of the performance of SureShrink and ASUS on MTC data. Here $n_k = |\hat{\mathcal{I}}_k^\tau|$ for $k = 1, 2$.

		MTC
		n
		22,283
SureShrink	t	1.13
	SURE estimate	1.32
		τ
		4.11
		t_1
		1.3
		t_2
		0.04
ASUS	n_1	21,791
	n_2	492
	SURE estimate	0.62

this real data example, a reduction in risk is possible because ASUS has efficiently exploited the sparsity information encoded in \mathcal{S} . This can be seen, for example, from the stark contrast between the magnitudes of thresholding hyper-parameters t_1 and t_2 in table 3. Moreover, the risk of Aux-Scr for this example was seen to be no better than the SureShrink estimator and thus has been excluded from the results reported in table 3.

F Choice of K

We consider the toy example discussed in section 2.3 of the main paper and let $K \geq 2$. For each candidate value of K we plot the SURE estimate of risk of ASUS in figure 4 left. An estimate of K may be taken to be the one that appears at the elbow of this plot and that implies $\hat{K} = 2$. Often a large value of K , say $K = 5$ or 6 , may continue to provide a marginal reduction in overall risk as opposed to $K = 2$ as seen in this example but such a reduction in risk comes at a cost of increased computational burden of conducting a search over $O(m_n^{K-1})$ points for large n . A cross validation based approach for selecting K in such scenarios is often useful. On a related note, in figure 4 right we demonstrate how ASUS reaps the benefits

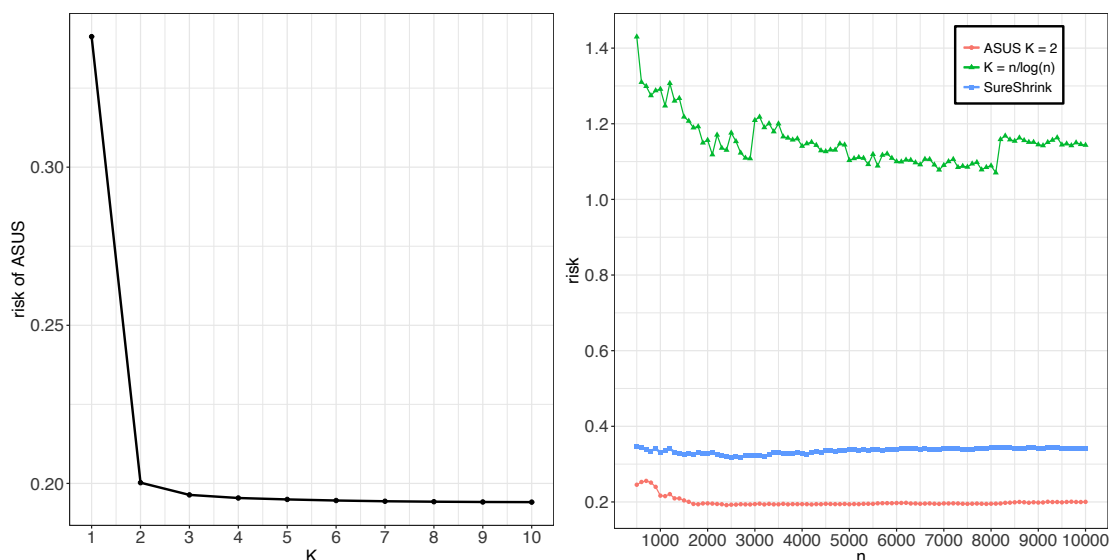


Figure 4: Toy example of section 2.3. Left: SURE estimate of risk of ASUS as K varies. Right: Estimate of risks for ASUS with $K = 2$, SureShrink and ASUS with no side information but with $K = n/\log n$.

of adapting to the informativeness of side information. We continue with the toy example of section 2.3 and estimate the hyper-parameter \mathcal{T} for $K = 2$. In contrast to this scheme, we also construct a version of ASUS where the number of groups K is automatically set to $n/\log n$ without the aid of any side information and only the thresholding hyper-parameters t_k are determined using the data driven hybrid scheme of [Donoho and Johnstone \(1995\)](#). The risk of this estimator is denoted by the green dotted line in figure 4 which clearly indicates that ASUS provides a better risk performance when the segmentation hyper-parameter τ is chosen in a data driven adaptive fashion.

References

- Calvano, S. E., W. Xiao, D. R. Richards, R. M. Felciano, H. V. Baker, R. J. Cho, R. O. Chen, B. H. Brownstein, J. P. Cobb, S. K. Tschoeke, et al. (2005). A network-based analysis of systemic inflammation in humans. *Nature* 437(7061), 1032–1037.
- Donoho, D. L. and I. M. Johnstone (1995). Adapting to unknown smoothness via wavelet shrinkage. *Journal of the American Statistical Association* 90(432), 1200–1224.
- Johnstone, I. M. (1994). On minimax estimation of a sparse normal mean vector. *The Annals of Statistics*, 271–289.
- Johnstone, I. M. (2015). Gaussian estimation: sequence and wavelet models. *Draft version*.
- Singh, R. S. (1975). On the glivenko-cantelli theorem for weighted empiricals based on independent random variables. *The Annals of Probability*, 371–374.
- Stein, C. M. (1981). Estimation of the mean of a multivariate normal distribution. *The Annals of Statistics*, 1135–1151.
- Sun, W. and Z. Wei (2011). Multiple testing for pattern identification, with applications to microarray time-course experiments. *Journal of the American Statistical Association* 106(493), 73–88.