

## Journal of the American Statistical Association

Publication details, including instructions for authors and subscription information:

<http://www.tandfonline.com/loi/uasa20>

### Two-Stage Importance Sampling With Mixture Proposals

Wentao Li <sup>a</sup>, Zhiqiang Tan <sup>b</sup> & Rong Chen <sup>b</sup>

<sup>a</sup> Department of Mathematics and Statistics, Fylde College, Lancaster University, Bailrigg, Lancaster LA1 4YF, United Kingdom

<sup>b</sup> Department of Statistics, Rutgers University, Piscataway, NJ, 08854

Accepted author version posted online: 24 Aug 2013. Published online: 19 Dec 2013.



To cite this article: Wentao Li, Zhiqiang Tan & Rong Chen (2013) Two-Stage Importance Sampling With Mixture Proposals, Journal of the American Statistical Association, 108:504, 1350-1365, DOI: [10.1080/01621459.2013.831980](https://doi.org/10.1080/01621459.2013.831980)

To link to this article: <http://dx.doi.org/10.1080/01621459.2013.831980>

PLEASE SCROLL DOWN FOR ARTICLE

Taylor & Francis makes every effort to ensure the accuracy of all the information (the "Content") contained in the publications on our platform. However, Taylor & Francis, our agents, and our licensors make no representations or warranties whatsoever as to the accuracy, completeness, or suitability for any purpose of the Content. Any opinions and views expressed in this publication are the opinions and views of the authors, and are not the views of or endorsed by Taylor & Francis. The accuracy of the Content should not be relied upon and should be independently verified with primary sources of information. Taylor and Francis shall not be liable for any losses, actions, claims, proceedings, demands, costs, expenses, damages, and other liabilities whatsoever or howsoever caused arising directly or indirectly in connection with, in relation to or arising out of the use of the Content.

This article may be used for research, teaching, and private study purposes. Any substantial or systematic reproduction, redistribution, reselling, loan, sub-licensing, systematic supply, or distribution in any form to anyone is expressly forbidden. Terms & Conditions of access and use can be found at <http://www.tandfonline.com/page/terms-and-conditions>

# Two-Stage Importance Sampling With Mixture Proposals

Wentao LI, Zhiqiang TAN, and Rong CHEN

For importance sampling (IS), multiple proposals can be combined to address different aspects of a target distribution. There are various methods for IS with multiple proposals, including Hesterberg's stratified IS estimator, Owen and Zhou's regression estimator, and Tan's maximum likelihood estimator. For the problem of efficiently allocating samples to different proposals, it is natural to use a pilot sample to select the mixture proportions before the actual sampling and estimation. However, most current discussions are in an empirical sense for such a two-stage procedure. In this article, we establish a theoretical framework of applying the two-stage procedure for various methods, including the asymptotic properties and the choice of the pilot sample size. By our simulation studies, these two-stage estimators can outperform estimators with naive choices of mixture proportions. Furthermore, while Owen and Zhou's and Tan's estimators are designed for estimating normalizing constants, we extend their usage and the two-stage procedure to estimating expectations and show that the improvement is still preserved in this extension.

KEY WORDS: Control variates; Normalizing constant; Pilot samples.

## 1. INTRODUCTION

Importance sampling (IS) is a very useful Monte Carlo method for approximating analytically intractable integrals. Normally IS technique is used to approximate two types of integrals:

1.  $Z = \int \pi(x)dx$ , where  $\pi(x)$  is a nonnegative integrable function.
2.  $\mu = \int h(x)\pi^*(x)dx$ , where  $\pi^*(x)$  is a probability density and  $h(x)$  is a real function.

The integral  $Z$  can be treated as the normalizing constant of some probability density  $\pi^*(x)$ . This type of integration arises in many areas, including missing data analysis, marginal likelihood calculation, estimation of free energies in physics (Gelman and Meng 1998), and communication system (Smith, Shafi, and Gao 1997). The second type of integral can be treated as the expectation of  $h(x)$  under  $\pi^*(x)$ . Applications of IS to this type of integral are also popular in many areas including rare event simulation (Denny 2001), reliability (Hesterberg 1995), computational finance (Owen and Zhou 1999), and computer graphics (Veach and Guibas 1995).

The idea of IS for approximating  $Z$  is based on the identity

$$\int \pi(x)dx = \int \frac{\pi(x)}{q(x)}q(x)dx,$$

where  $q(x)$  is a probability density, called the proposal or trial density. Then  $Z$  can be treated as the expectation under density  $q(x)$ . With a sample  $(x_1, \dots, x_n)$  from  $q(x)$ ,  $Z$  can be approximated by the sample average  $n^{-1} \sum \pi(x_i)/q(x_i)$ . For approximating  $\mu$ , similar ideas hold. In this case, even when it is not difficult to generate observations from  $\pi^*(x)$ , IS estimator can

have higher efficiency than using observations from  $\pi^*(x)$ , by appropriately choosing  $q(x)$  so that it takes the shape of both  $h(x)$  and  $\pi^*(x)$  into account. In some applications, the improvement can be of several orders of magnitude. One example is given in Hesterberg (1995). In rare event simulation,  $h(x)$  is an indicator function with support in the tail of  $\pi^*(x)$ . A trial distribution  $q(x)$  that focuses on the "important" part of the integrand  $h(x)\pi^*(x)$  will be more efficient than generating samples directly from  $\pi^*(x)$ .

However, in practice it is usually challenging to implement IS efficiently. One reason is that the complexity of the integrand makes it difficult to find  $q(x)$  that covers all its important parts. For example, when the integrand is multimodal, a unimodal  $q(x)$  will not be efficient. Some of such multimodal integrand can be found in Owen and Zhou (1999). Another well-known reason is that when  $q(x)$  has a "lighter" tail than the integrand, that is,  $\pi(x)/q(x)$  or  $h(x)\pi^*(x)/q(x)$  are unbounded, IS estimator can have infinite variance. When there is a lack of knowledge of the integrand in some regions, unexpected large values of integrand may result in inaccurate results. See Ford and Gregory (2007) for an example.

For both problems, a general remedy is to consider multiple proposal distributions to address different aspects of the integrand. For multimodal integrands, Oh and Berger (1993) used a family of Student's  $t$  distributions and Owen and Zhou (1999) used a family of beta distributions to model each mode of integrand individually. West (1993) and Givens and Raftery (1996) used a kernel estimate of the integrand as the proposal, which is a mixture of normal or  $t$  distributions. Even for a unimodal target distribution, one can construct a mixture of two proposals where one mimics the center of target and the other dominates the tail. Such a construction was used in Giordani and Kohn (2010), although in a different scenario. The requirement that the tail of integrand needs to be dominated by the proposal distribution can be met by including some heavy-tailed distributions in the mixture as "protection." For example, Hesterberg

Wentao Li is a Senior Research Associate, Department of Mathematics and Statistics, Fylde College, Lancaster University, Bailrigg, Lancaster LA1 4YF, United Kingdom (E-mail: [w.li@lancaster.ac.uk](mailto:w.li@lancaster.ac.uk)). Zhiqiang Tan is Associate Professor, Department of Statistics, Rutgers University, Piscataway, NJ 08854 (E-mail: [ztan@stat.rutgers.edu](mailto:ztan@stat.rutgers.edu)). Rong Chen is Professor, Department of Statistics, Rutgers University, Piscataway, NJ 08854 (E-mail: [rongchen@stat.rutgers.edu](mailto:rongchen@stat.rutgers.edu)). Chen's research is sponsored in part by NSF grants DMS 0800183, DMS 0905763, and DMS0915139. Tan's research is sponsored in part by NSF grant DMS 0749418. The authors thank the two editors, an associate editor, and two anonymous referees for their helpful comments.

(1995) included the target distribution itself as one of the components to provide an upper bound for the estimation variance, and Owen and Zhou (2000) used uniform distribution to bound the sample weights in a bounded domain case. Liang, Liu, and Carroll (2007) divided the state domain into subregions and used the mixture of truncated target distributions in all subregions as the proposal, which leads to bounded importance weights. In Bayesian analysis, the prior density of the parameters can serve as the heavy-tailed component, as used in Ford and Gregory (2007).

Given multiple potentially useful proposals, a straightforward combination method is to use their mixture as the new proposal. This method has two issues. One is that the mixture proposal may contaminate the good components in the mixture. Owen and Zhou (2000) showed that a mixture can lose efficiency by several orders of magnitude if the original proposal is nearly perfect. Another problem is that the mixture proportions need to be determined. Proper mixture proportions can increase the efficiency by an order of magnitude, as shown in Emond, Raftery, and Steele (2001).

Owen and Zhou (2000) suggested a regression method to combine multiple proposals using control variates to deal with the contamination problem. Control variate is a useful technique for variance reduction (Rubinstein and Kroese 2008). Owen and Zhou's method has the property that it will not perform worse than using the best of component proposals individually, if the sample size assigned to the best component is the same as in the mixture case. Such a lower bound lessens the contamination problem. Tan (2004) proposed to use nonparametric maximum likelihood estimation (MLE) in place of regression and showed that the MLE method is the most efficient among several classes of estimators including those in Owen and Zhou (2000), Hesterberg (1995), and Veatch and Guibas (1995).

To determine appropriate proportions, Fan et al. (2006) and Hesterberg (1995) followed some heuristic rules derived from experience or interpretation of proposals. A more sophisticated approach is to use a pilot study to determine the optimal proportions via minimizing some criterion. The estimated proportions are then used to generate the sample and construct the estimators. The criterion was selected to be the asymptotic variance of IS estimator with mixture proposal in Raghavan and Cox (1998), and the variation coefficient of pilot sample in Oh and Berger (1993). However, few theoretical properties have been investigated.

In this article, we propose a similar two-stage procedure and investigate its theoretical properties. In the first stage, pilot sample is drawn from a mixture proposal with predetermined proportions. The optimal mixture proportions are then estimated by minimizing the estimated asymptotic variance of Owen's regression estimator or Tan's MLE based on control variates. In the second stage, the sample is drawn from the mixture proposal with the estimated proportions. Integral estimators are constructed using all observations, including those from the pilot stage. Then we establish a theoretical framework of such a two-stage procedure. It is shown that under very weak conditions, the integral estimators constructed by the two-stage procedure are consistent and asymptotic normal with minimum asymptotic variance over all mixture proportions. Therefore, the two-stage procedure is adaptive toward using the optimal

mixture proportions. The optimal sample size used for the pilot stage is also calculated in the sense of minimizing an approximated mean square error (MSE) in higher order. Furthermore, we extend Owen's regression estimator and Tan's MLE to ratio estimators (Rubinstein and Kroese 2008). When estimating  $\mu$ , if one can evaluate  $\pi^*(x)$  only up to a normalizing constant, a ratio estimator is used, with the numerator being the estimated unnormalized integral and the denominator being the estimated normalizing constant. We show that the two-stage procedure for this extension also has the desirable asymptotic properties.

The remainder of this article is organized as follows. Section 2 reviews in detail some existing techniques related to IS. Section 3 proposes the new two-stage procedure and establishes its theoretical framework. Section 4 provides the extension of Owen's regression estimator and Tan's MLE to ratio estimators and their two-stage procedures. In Section 5, we demonstrate the two-stage approach with several numerical examples including estimation of a rare event probability and Value at Risk (VaR) under a Bayesian GARCH model.

## 2. REVIEW OF IS TECHNIQUES

In Sections 2 and 3, we assume that all functions in integrals (1) and (2) can be evaluated exactly. Since

$$\mu = \int h(x)^+ \pi^*(x) dx - \int h(x)^- \pi^*(x) dx,$$

where  $h(x)^+ = h(x)1_{\{h(x) \geq 0\}}$  and  $h(x)^- = -h(x)1_{\{h(x) < 0\}}$ , the estimation of  $\mu$  can be achieved by estimation of these two integrals, both with positive integrands. With this approach,  $\mu$  becomes a special case of  $Z$  and therefore we only consider estimating  $Z$  in these two sections.

### 2.1 Mixture Importance Sampling

Assume observations  $\{x_1, \dots, x_n\}$  are taken iid from a proposal distribution  $q(x)$ . The integral  $Z = \int \pi(x) dx$  can be estimated by

$$\hat{Z}_{\text{IS}} = \frac{1}{n} \sum_{i=1}^n \frac{\pi(x_i)}{q(x_i)}. \quad (1)$$

Under mild conditions, the asymptotic variance is  $\text{var}_q[\pi(x)/q(x)]$ , where  $\text{var}_q$  is the variance under distribution  $q(x)$  (Robert and Casella 2004). The optimal proposal is  $\pi(x)/Z$ , suggesting that the proposal  $q(x)$  should be chosen to mimic the shape of  $\pi(x)$  so that the high- and low-density regions of  $q(x)$  coincide with those of  $\pi(x)$ . With such a proposal, the majority of Monte Carlo sample from  $q(x)$  fall in the high-density region of  $\pi(x)$ , the importance region. In some scenarios, more than one  $q(x)$  may be needed. For example, for a multimodal  $\pi(x)$ , it is helpful to use several proposal distributions, each targeted at one importance region. Suppose  $q_1(x), \dots, q_p(x)$  are  $p$  probability densities serving as proposals. Given a mixture proportion vector  $\alpha = (\alpha_1, \dots, \alpha_p)$  satisfying  $\sum_{k=1}^p \alpha_k = 1$ , we can use the mixture distribution as the proposal and estimate  $Z$  by

$$\hat{Z}_{\text{MIS}} = \frac{1}{n} \sum_{i=1}^n \frac{\pi(x_i)}{q_\alpha(x_i)}, \quad (2)$$

where  $q_\alpha = \sum_{k=1}^p \alpha_k q_k(x)$  and  $\{x_1, \dots, x_n\}$  are generated from  $q_\alpha$ . In addition, the variance of  $\hat{Z}_{\text{IS}}$  also demands that the ratio  $\pi(X)/q(X)$  has a finite variance. A mixture distribution certainly makes it easy to satisfy such a condition as one can simply include a proposal distribution  $q_1(X)$  having  $\text{var}[\pi(X)/q_1(X)] < \infty$ , such as a uniform distribution if the domain is bounded. Such a proposal distribution sets an upper bound to the estimating variance and therefore plays the role of “safeguard” in IS, which is the key idea of defensive IS (Hesterberg 1988).

## 2.2 Stratified Sampling

Instead of generating samples directly from the mixture distribution as that in (2), stratified samples  $\{x_{k1}, \dots, x_{kn_k}\}$  can be taken with deterministic size  $n_k = \alpha_k n$  from the  $k$ th proposal  $q_k$ , which leads to the estimator in Hesterberg (1988)

$$\hat{Z}_{\text{SIS}}(\alpha) = \frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{\pi(x_{ki})}{q_\alpha(x_{ki})}. \quad (3)$$

Veatch and Guibas (1995) considered the estimator

$$\sum_{k=1}^p \frac{1}{n_k} \sum_{i=1}^{n_k} \omega_k(x_{ki}) \frac{\pi(x_{ki})}{q_k(x_{ki})}, \quad (4)$$

where  $\{\omega_k(x)\}_{k=1}^p$  is a group of coefficient functions for the sample weights and satisfies  $\sum_{k=1}^p \omega_k(x) = 1$ . They showed that  $\hat{Z}_{\text{SIS}}$  is a suboptimal choice in this large class. Raghavan and Cox (1998) proposed a two-stage algorithm to construct  $\hat{Z}_{\text{SIS}}$  with estimated optimal mixture proportions in the sense of minimizing the asymptotic variance.

## 2.3 Importance Sampling With Control Variates

One problem of using a mixture proposal distribution is the possible loss of efficiency due to mixing of good proposal distributions with poor ones (Owen and Zhou 2000). It is a premium to pay for the insurance of valid IS, but can be reduced by combining importance sampling and control variates. Given an unbiased estimator  $X_n$  of  $Z$ , improvement can be gained by constructing a proper control variate vector  $Y$  and using  $X_n - \beta^T(Y - E[Y])$  to estimate  $Z$ . The optimal  $\beta$  can be estimated using a regression approach to minimize asymptotic variance (Cochran 1977). In Owen and Zhou (2000), combining  $\hat{Z}_{\text{SIS}}$  and control variates  $\mathbf{g}(x) = (q_2(x) - q_1(x), \dots, q_p(x) - q_1(x))^T$  results in the estimator

$$\hat{Z}_{\text{Reg}}(\alpha) = \frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{\pi(x_{ki}) - \hat{\beta}_\alpha^T \mathbf{g}(x_{ki})}{q_\alpha(x_{ki})}, \quad (5)$$

where

$$\hat{\beta}_\alpha = \widetilde{\text{var}} \left[ \frac{\mathbf{g}(X)}{q_\alpha(X)} \right]^{-1} \widetilde{\text{cov}}^T \left[ \frac{\pi(X)}{q_\alpha(X)}, \frac{\mathbf{g}(X)}{q_\alpha(X)} \right],$$

and  $\widetilde{\text{var}}$  and  $\widetilde{\text{cov}}$  denote the pooled-sample variance and covariance. There are two appealing properties of  $\hat{Z}_{\text{Reg}}$ . First, its asymptotic variance is zero when  $\pi(x)$  is a linear combination of the proposals. Second,  $\hat{Z}_{\text{Reg}}$  has smaller asymptotic variance than every IS estimator constructed solely with  $q_k$  with  $n_k$  samples,  $k = 1, \dots, p$ . That is,  $\hat{Z}_{\text{Reg}}$  is always at least as good as the best one among the individual proposals.

## 2.4 Likelihood Approach

All previous integration methods directly approximate the target integrals. On the other hand, in Kong et al. (2003), Monte Carlo integration is treated as a statistical inference problem where the Monte Carlo sample serves as observations, the underlying measure in target integral, usually Lebesgue measure or counting measure, is treated as an unknown nonnegative measure, and the Monte Carlo sample is modeled using a semi-parametric model. Then by nonparametric maximum likelihood, the unknown measure is estimated by a discrete measure with the Monte Carlo sample as support, and the target integral is estimated by the integration over the discrete measure. As an example, with  $\{x_1, \dots, x_n\}$  generated identically and independently from  $q_1$  under Lebesgue measure, the model assumes that  $x_i$  is distributed as  $q_1(x)d\nu / \int q_1(x)d\nu$ , where  $\nu$  is an unknown nonnegative measure. The nonparametric maximum likelihood estimator of  $\nu$  is

$$\hat{\nu} \propto \frac{\hat{P}(\{x\})}{q_1(x)},$$

where  $\hat{P}$  has the support on  $\{x_1, \dots, x_n\}$  with mass  $n^{-1}$  at each point. Then  $Z = \int q(x)d\nu$  can be estimated by

$$\frac{\int q_1(x)d\hat{\nu}}{\int q(x)d\hat{\nu}} = \frac{1}{n} \sum_{i=1}^n \frac{q(x_i)}{q_1(x_i)}.$$

This is the same as the IS estimator with proposal distribution  $q_1(x)$ .

Given multiple proposals  $q_1, \dots, q_p$  and control variates  $\mathbf{g}(x)$ , Tan (2004) proposed to restrict the measure  $\nu$  in the set  $\{\nu : \int q_k(x)d\nu = \int q_1(x)d\nu, k = 1, \dots, p\}$ . The nonparametric MLE of  $\nu$  under such a restriction is

$$\hat{\nu} \propto \frac{\hat{P}(\{x\})}{q_\alpha(x) + \hat{\zeta}^T \mathbf{g}(x)},$$

where

$$\hat{\zeta} = \arg \max_{\zeta} \sum_{k=1}^p \sum_{i=1}^{n_k} \log [q_\alpha(x_{ki}) + \zeta^T \mathbf{g}(x_{ki})],$$

and the integral estimator is given by

$$\hat{Z}_{\text{MLE}}(\alpha) = \frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{\pi(x_{ki})}{q_\alpha(x_{ki}) + \hat{\zeta}^T \mathbf{g}(x_{ki})}. \quad (6)$$

It is shown that  $\hat{Z}_{\text{Reg}}$  is a first-order approximation of  $\hat{Z}_{\text{MLE}}$  and hence has the same asymptotic efficiency. The estimator  $\hat{Z}_{\text{MLE}}$  also achieves the highest asymptotic efficiency among the class of estimators in the form of

$$\sum_{k=1}^p \frac{1}{n_k} \sum_{i=1}^{n_k} \omega_k(x_{ki}) \frac{\pi(x_{ki}) - \beta_k^T(x_{ki}) \mathbf{g}(x_{ki})}{q_k(x_{ki})}, \quad (7)$$

where  $\omega_1(x), \dots, \omega_p(x)$  and  $\beta_1(x), \dots, \beta_p(x)$  satisfy that  $\omega_k(x) = 0$  when  $q_k(x) = 0$ ,  $\sum_{i=1}^p \omega_k(x) = 1$  and  $\sum_{i=1}^p \omega_k(x) \beta_k(x) = \mathbf{b}$  for some constant vector  $\mathbf{b}$ , and therefore dominates the class of estimators in (4).

Because  $\hat{Z}_{\text{Reg}}$  and  $\hat{Z}_{\text{MLE}}$  dominate the other estimators, we will only discuss the two-stage procedure for these two estimators. Furthermore, there is another important benefit of using  $\hat{Z}_{\text{Reg}}$  and  $\hat{Z}_{\text{MLE}}$  in that their asymptotic variance is a



convex function of  $\alpha$  and hence can be easily minimized. See Remark 1.

### 3. TWO-STAGE PROCEDURE

#### 3.1 The Algorithm

Suppose  $p$  proposal distributions  $q_1, \dots, q_p$  are given and the sample size is budgeted at  $n$ . Let  $\Theta = [\delta, 1 - \delta]^p$ , where  $\delta$  is some constant close to 0. The following algorithm is proposed to select mixture proportions  $\alpha$  and construct estimators:

1. First stage: Given a  $p$ -dimensional vector  $\gamma$  satisfying  $\sum_{k=1}^p \gamma_k = 1$ , generate  $n_0$  independent stratified observations  $\{x_i\}_{i=1}^{n_0}$  from  $q_\gamma(x) = \sum_{k=1}^p \gamma_k q_k(x)$ , that is,  $n_0 \gamma_k$  observations from  $q_k(x)$ ,  $k = 1, \dots, p$ . Obtain  $\hat{\alpha}$  by minimizing

$$\hat{\sigma}_Z^2(\alpha) = \frac{1}{n_0} \sum_{i=1}^{n_0} \frac{[\pi(x_i) - \hat{\beta}_\alpha^T g(x_i)]^2}{q_\alpha(x_i) q_\gamma(x_i)}, \quad (8)$$

where

$$\hat{\beta}_\alpha = \left( \frac{1}{n_0} \sum_{i=1}^{n_0} \frac{g(x_i) g(x_i)^T}{q_\alpha(x_i) q_\gamma(x_i)} \right)^{-1} \left( \frac{1}{n_0} \sum_{i=1}^{n_0} \frac{\pi(x_i) g(x_i)}{q_\alpha(x_i) q_\gamma(x_i)} \right)$$

and  $g(x) = (q_2(x) - q_1(x), \dots, q_p(x) - q_1(x))^T$ , with respect to  $\alpha$  over  $\Theta$ .

2. Second stage: Generate  $n - n_0$  independent stratified observations  $\{x_i\}_{i=n_0+1}^n$  from  $q_{\tilde{\alpha}}(x) = \sum_{k=1}^p \tilde{\alpha}_k q_k(x)$ . Estimate integral  $Z$  by  $\hat{Z}(\tilde{\alpha})$  with all  $n$  observations, where

$$\tilde{\alpha} = \frac{n_0}{n} \gamma + \frac{n - n_0}{n} \hat{\alpha} \quad (9)$$

and  $\hat{Z}(\tilde{\alpha})$  can be either  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$  or  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$ .

Some rationale and implementation remarks are as follows:

- (i) Criterion of selecting  $\alpha$ : In the first stage, the optimal  $\alpha$  is estimated using the  $n_0$  samples and it is desirable to select  $\alpha$  that gives the smallest asymptotic variance of the final estimator. Let  $\text{var}_\alpha$  denotes the variance taken with respect to  $q_\alpha(x)$ . We set the following conditions:

(C1) The union of supports of  $q_k(x)$  contains the support of  $\pi(x)$ .

(C2)  $\alpha_i > 0$  for  $i = 1, \dots, p$ .

(C3)  $\text{var}_\alpha [\pi(X)/q_\alpha(X)] < \infty$  for some  $\alpha \in \Theta$ .

Owen and Zhou (2000) and Tan (2004) showed that, under the above conditions,  $\hat{Z}_{\text{Reg}}(\alpha)$  and  $\hat{Z}_{\text{MLE}}(\alpha)$  are asymptotic normal and have the same asymptotic variance

$$\begin{aligned} \sigma_Z^2(\alpha) &= \text{var}_\alpha \left[ \frac{\pi(X) - \beta_\alpha^T g(X)}{q_\alpha(X)} \right] \\ &= \int \frac{(\pi(x) - \beta_\alpha^T g(x))^2}{q_\alpha(x)} dx - Z^2, \end{aligned} \quad (10)$$

where

$$\begin{aligned} \beta_\alpha &= \text{var}_\alpha \left[ \frac{g(X)}{q_\alpha(X)} \right]^{-1} \text{cov}_\alpha^T \left[ \frac{\pi(X)}{q_\alpha(X)}, \frac{g(X)}{q_\alpha(X)} \right] \\ &= \left( \int \frac{g(x) g(x)^T}{q_\alpha(x)} dx \right)^{-1} \left( \int \frac{\pi(x) g(x)}{q_\alpha(x)} dx \right). \end{aligned}$$

Conditions (C1) to (C3) are satisfied when we have at least one proposal component dominating the tail of  $\pi(x)$ . With the sample  $\{x_i\}_{i=1}^{n_0}$  from the pilot stage,  $\sigma_Z^2(\alpha) + Z^2$  is estimated by the IS estimator  $\hat{\sigma}_Z^2(\alpha)$  in (8) and the optimal  $\alpha$  is obtained by minimizing  $\hat{\sigma}_Z^2(\alpha)$ .

- (ii) Optimization range for  $\alpha$ : The purpose of restricting  $\alpha$  in  $[\delta, 1 - \delta]^p$  for some small  $\delta$  is to avoid unreliable estimators of  $\sigma_Z^2(\alpha)$  or  $\beta_\alpha$ . When  $\alpha_i = 0$  for some  $i$ ,  $\int \pi(x)^2 / q_\alpha(x) dx$  can be infinite if  $q_i$  is the only proposal that dominates certain part of  $\pi(x)$ 's tail, or  $\int g(x) g(x)^T / q_\alpha(x) dx$  and  $\int \pi(x) g(x) / q_\alpha(x) dx$  can be infinite if  $q_i$  is the only proposal that dominates some other proposals. In this case, if  $\alpha_i$  is too close to 0, the estimator  $\hat{\sigma}_Z^2(\alpha)$  or  $\hat{\beta}_\alpha$  is unreliable. Experience shows that  $\delta = 0.001$  is a reasonable choice.
- (iii) Choice of the initial proportions  $\gamma$ :  $\gamma$  is preferred to be close to the optimal proportion vector  $\alpha^*$ . If there is no any prior knowledge about  $\alpha^*$ , it is recommended to use  $\gamma$  with equal components in the first stage so that pilot sample is generated from each proposal equally.
- (iv)  $\hat{Z}$  in second stage: Instead of using  $n - n_0$  observations to construct the estimator  $\hat{Z}(\tilde{\alpha})$ , we use all  $n$  observations to construct the estimator  $\hat{Z}(\tilde{\alpha})$  where the mixture proportions  $\tilde{\alpha}$  account for the proportions of the combined sample.

#### 3.2 Theoretical Properties

Let  $\alpha^*$  be the minimizer of  $\sigma_Z^2(\alpha)$  under restriction  $\alpha \in \Theta$ . We assume the following additional conditions:

(C4)  $n_0 = o(n)$  and  $n_0 \rightarrow \infty$  as  $n \rightarrow \infty$ .

(C5)  $\pi(x)$  is not a linear combination of  $q_1(x), \dots, q_p(x)$ .

(C6)  $\alpha^*$  is in the interior of  $\Theta$ , that is,  $\alpha^* \in (\delta, 1 - \delta)^p$ .

Condition (C4) ensures  $\tilde{\alpha}$  converges to  $\alpha^*$ . Condition (C5) is necessary since if  $\pi(x)$  is a linear combination of  $q_1(x), \dots, q_p(x)$ ,  $\sigma^2(\alpha)$  will be 0 for all  $\alpha_1$ . Some discussions of condition (C6) are given in Remark 7.

##### 3.2.1 First-Order Properties

*Theorem 1.* Under conditions (C1) to (C5),  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$  and  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$  are consistent and

$$\begin{aligned} \sqrt{n} (\hat{Z}_{\text{MLE}}(\tilde{\alpha}) - Z) &\xrightarrow{L} N(0, \sigma_Z^2(\alpha^*)) \\ \text{and } \sqrt{n} (\hat{Z}_{\text{Reg}}(\tilde{\alpha}) - Z) &\xrightarrow{L} N(0, \sigma_Z^2(\alpha^*)). \end{aligned}$$

Therefore, the two-stage procedure achieves the minimum asymptotic variance that Owen and Zhou's and Tan's estimators can achieve among all possible mixture proportions. Furthermore, since  $\hat{Z}_{\text{Reg}}(\alpha)$  and  $\hat{Z}_{\text{MLE}}(\alpha)$  are better than the stratified sampling estimator  $\hat{Z}_{\text{SIS}}(\alpha)$ , the two-stage procedure

**outperforms** all estimators introduced in Section 2 in asymptotic variance. The proof is given in Appendix A.

*Remark 1.* It is important to point out that  $\sigma_Z^2(\alpha)$  and its estimator  $\hat{\sigma}_Z^2(\alpha)$  are **strictly convex** by Lemma 1 in Appendix A. This guarantees a unique solution and applicability of convex optimization algorithms in the pilot stage. This property, or equivalently the strict convexity of the function  $\sigma^2(\alpha, \beta) = \text{var}_\alpha[(\pi(X) - \beta^T \mathbf{g}(X))/q_\alpha(X)]$ , also ensures the consistency and asymptotic normality with convergence rate  $\sqrt{n_0}$  of random proportion vector  $\hat{\alpha}$  under mild conditions, by asymptotic theory for  $M$ -estimation with a convex criterion function (Haberman 1989). Therefore, larger  $n_0$  gives more reliable  $\hat{\alpha}$ .

*Remark 2.* For  $\hat{Z}_{\text{SIS}}(\alpha)$ , the optimal mixture proportions are the ones that make the mixture proposal  $q_\alpha$  the closest to the target distribution  $\pi$ . Therefore, knowledge about the target density surface can help to find an approximate choice of  $\alpha$ . However, the optimal mixture proportions  $\alpha^*$  for  $\sigma_Z^2(\alpha)$  sometimes can be counterintuitive. For instance, in Example 1(B2) of Section 5, the target distribution is a mixture of a normal distribution and a  $t$  distribution, with mixing probability 0.8 and 0.2, respectively. When the same normal distribution is used as one of the proposal distributions, its optimal mixture proportion is only 0.1%. This is because, for  $\hat{Z}_{\text{Reg}}(\alpha)$  and  $\hat{Z}_{\text{MLE}}(\alpha)$ , the numerator of  $\sigma_Z^2(\alpha)$  involves  $\beta_\alpha$ , a function of  $\alpha$ , which complicates the determination of the optimal proportions. Hence, an automatic selection for mixture proportions becomes necessary for  $\hat{Z}_{\text{Reg}}(\alpha)$  and  $\hat{Z}_{\text{MLE}}(\alpha)$ .

*Remark 3.* If  $\tilde{\alpha}$  in (9) can be replaced by some other random proportion vector, as long as it is consistent to  $\alpha^*$  as  $n \rightarrow \infty$ , the same asymptotic results hold. For example, one can choose the mixture proportions of the second stage so that the combined sample (of both the pilot stage and second stage) is as close to the estimated optimal proportion vector  $\hat{\alpha}$  as possible. For example, if  $n_0 \gamma_k < n \hat{\alpha}_k$  for all  $k = 1, \dots, p$ , one can use  $(n \hat{\alpha} - n_0 \gamma)/(n - n_0)$  in the second stage, which results in the combined sample having the exact estimated optimal proportion  $\hat{\alpha}$ . In this case, actually one should use  $n_0$  as large as possible until it violates the above condition.

*Remark 4.* Similar asymptotic properties for  $\hat{Z}_{\text{MIS}}(\tilde{\alpha})$  and  $\hat{Z}_{\text{SIS}}(\tilde{\alpha})$  are presented in Lemma 3 in Appendix A. They are always inferior to the control variate based estimators and hence of less interest.

**3.2.2 High-Order Properties.** Theorem 1 shows that the selection of the pilot sample size  $n_0$  does not affect the first-order property of  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$  and  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$  as long as  $n_0 = o(n)$  and  $n_0 \rightarrow \infty$ . Therefore, an optimal choice of  $n_0$  needs to be determined by higher-order properties of  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$  and  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$ . Consider the convergence rate of  $\tilde{\alpha} - \alpha^*$ , a weighted average of  $\gamma - \alpha^*$  and  $\hat{\alpha} - \alpha^*$  with weights  $n_0/n$  and  $1 - n_0/n$ . Since  $\gamma - \alpha^*$  is biased, one would want to have a smaller  $n_0$ . However, a large  $n_0$  makes  $\hat{\alpha} - \alpha^*$  closer to 0, at the rate  $O(1/\sqrt{n_0})$ . Therefore, the optimal  $n_0$  is chosen to balance the effects of these two rates. The following proposition gives the higher-order asymptotic expansions of  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$  and  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$ .

*Proposition 1.* Under conditions (C1)–(C6),  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$  and  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$  can be expanded as  $\hat{Z}^* + o(n_0/(n\sqrt{n})) + o(1/(n_0\sqrt{n}))$

and  $\hat{Z}^* = Z + g_1(\tilde{\alpha}) + g_2(\tilde{\alpha})$ , where

$$g_1(\tilde{\alpha}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(x_i) - \beta_{\alpha^*} g(x_i)}{q_{\alpha^*}(x_i)} - \int \frac{\pi(x) - \beta_{\alpha^*} g(x)}{q_{\alpha^*}(x)} q_{\tilde{\alpha}}(x) dx,$$

and

$$g_2(\tilde{\alpha}) = O\left(\frac{n_0}{n\sqrt{n}}\right) + O\left(\frac{1}{n_0\sqrt{n}}\right).$$

The explicit forms of  $g_2(\tilde{\alpha})$  are tedious and therefore presented in the Appendices. The selection of optimal  $n_0$  is based on minimizing the MSE of  $\hat{Z}^*$ , which is an approximation of the MSE of  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$  and  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$ . Such an approximation of moments, as the criterion of second-order optimality, has been widely used in higher-order asymptotic theory, for example, Rothenberg (1984).

*Theorem 2.* Under conditions (C1)–(C6) and (C7)  $\int \pi(x)^4/q_\alpha(x)^4 dx < \infty$  for some  $\alpha \in \Theta$ ,

it holds that

$$E[\hat{Z}^* - Z] = O\left(\frac{1}{n}\right) \quad \text{and} \quad \text{var}[\hat{Z}^* - Z] = \frac{1}{n} \sigma_Z^2(\alpha^*) + O\left(\frac{n_0}{n^2}\right) + O\left(\frac{1}{nn_0}\right).$$

Therefore,  $\text{MSE}[\hat{Z}^*] - n^{-1} \sigma_Z^2(\alpha^*) = O(\frac{n_0}{n^2}) + O(\frac{1}{nn_0})$ .

The above result gives the approximate mean squared error with higher-order terms beyond the usual asymptotic variance  $n^{-1} \sigma_Z^2(\alpha^*)$ . The order can be attributed to three sources of variability. See the Appendices for details. One source of variability is due to using the pilot sample with mixture proportions  $\gamma \neq \alpha^*$ , which leads to terms of order  $O(n_0/n^2)$ . The second source is the variability of estimator  $\hat{\alpha}$ , which is of the order  $O(1/(nn_0))$ . The third source is the variability of estimating  $\beta_{\alpha^*}$ , which is the optimal coefficient of control variates, in  $\sigma_Z^2(\alpha^*)$ . In  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$ , the estimator of  $\beta_{\alpha^*}$  is  $\hat{\beta}_{\tilde{\alpha}}$ . In  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$ , a similar estimator is used, as can be seen from the proof of Theorem 1. This variability is of the order  $O(1/(n\sqrt{n}))$ , which is also  $O(n_0/n^2) + O(1/(nn_0))$  because  $2/\sqrt{n} \leq n_0/n + 1/n_0$  by the inequality  $2ab \leq a^2 + b^2$ .

*Remark 5.* By minimizing the order of difference, the optimal  $n_0$  is  $O(\sqrt{n})$  and hence  $\text{MSE}[\hat{Z}^*] - n^{-1} \sigma_Z^2(\alpha^*)$  is of order  $O(1/(n\sqrt{n}))$ . The asymptotic rate shows that how  $n_0$  should change with the total sample size  $n$ . In practice, another consideration of selecting  $n_0$  is the coverage of the target distribution with pilot samples. A poor coverage can lead to poorly estimated asymptotic variance and result in inaccurate  $\hat{\alpha}$ . Our experience shows one should choose  $n_0$  at least  $\sqrt{n}$  and possibly larger according to the complexity of problem and the quality of proposal distributions. On the other hand, one can assess  $\hat{\alpha}$  by estimating its standard error after the pilot stage. If the standard error is larger than some criterion, such as 10% of  $\hat{\alpha}$ , one can add additional pilot samples. The standard error formula is given in Appendix B.

*Remark 6.* One essential fact leading to Theorem 2 is  $\tilde{\alpha} - \alpha^* = O(1/\sqrt{n_0}) + O(n_0/n)$ . Therefore, when  $\tilde{\alpha}$  is replaced by

some other construction that is consistent but with different rate (e.g., Remark 3), the orders in Theorem 2 may change.

*Remark 7.* When some coordinates of  $\alpha^*$  are on the boundary of  $[\delta, 1 - \delta]$ , the exact second-order property is complicated. However, it is still reasonable to use the same  $n_0$  as indicated in Theorem 2. For example, when  $\alpha_1^*$  is on the boundary,  $\hat{\alpha}$ , as an  $M$ -estimator, will converge to  $\alpha^*$  with a rate faster than or equal to  $O(1/\sqrt{n_0})$  (Geyer 1994). In the proof of Theorem 2, when the convergence rate of  $\hat{\alpha}_1$  changes from  $O(1/\sqrt{n_0})$  to  $O(1/n_0^\varepsilon)$  with  $\varepsilon \geq \frac{1}{2}$ , the second order  $O(n_0/n^2) + O(1/(n_0n))$  changes to  $O(n_0/n^2) + O(1/(n_0^\varepsilon n))$ . Then by choosing  $n_0 = O(\sqrt{n})$ ,  $\text{MSE}[\hat{Z}^*] - n^{-1}\sigma_Z^2(\alpha^*)$  is still  $O(1/(n\sqrt{n}))$  and the accuracy of  $\hat{Z}_{\text{Reg}}(\hat{\alpha})$  and  $\hat{Z}_{\text{MLE}}(\hat{\alpha})$  remains the same.

## 4. EXTENSION TO RATIO ESTIMATORS

### 4.1 Extension of IS Techniques to Ratio Estimators

As mentioned in Section 1, the integral (2) can be estimated by the ratio estimator

$$\hat{\mu}_{\text{IS}} = \frac{\frac{1}{n} \sum_{i=1}^n h(x_i) \pi(x_i) / q(x_i)}{\frac{1}{n} \sum_{i=1}^n \pi(x_i) / q(x_i)}, \quad (11)$$

(Liu 2008; Rubinstein and Kroese 2008). By the delta method, it is easy to show that the asymptotic variance of  $\hat{\mu}_{\text{IS}}$  is

$$\text{var}_q \left( \frac{h(x)\pi(x) - \mu\pi(x)}{q(x)} \right). \quad (12)$$

In the sense of minimizing (12), the optimal choice of  $q(x)$  is the probability density proportional to  $|h(x)\pi(x) - \mu\pi(x)|$ . Therefore, it is preferred to choose  $q(x)$  that mimics the shape of  $|h(x)\pi(x) - \mu\pi(x)|$ . Similar to estimating the normalizing constant, multiple proposals may be needed and the techniques in Section 2 may be beneficial.

Given  $p$  proposal distributions  $q_1(x), \dots, q_p(x)$  and mixture proportions  $\{\alpha_k\}_{k=1}^p$  satisfying  $\sum_{k=1}^p \alpha_k = 1$ . Observations  $\{x_{k1}, \dots, x_{kn_k}\}$  are generated from proposal  $q_k$  with size  $n_k = \alpha_k n$  for each  $k$ . In Hesterberg (1995), the mixture IS and stratified sampling were applied to  $\hat{\mu}_{\text{IS}}$  by using the mixture proposal  $q_\alpha$  in numerator and denominator separately as follows:

$$\hat{\mu}_{\text{SIS}} = \frac{\frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} h(x_i) \pi(x_{ki}) / q_\alpha(x_{ki})}{\frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \pi(x_{ki}) / q_\alpha(x_{ki})}.$$

Control variates and likelihood approach can also be applied to  $\hat{\mu}_{\text{IS}}$ . With the same control variates  $\mathbf{g}(x)$  as in (8),  $\mu$  can be estimated by the following:

$$\begin{aligned} \hat{\mu}_{\text{Reg}} &= \frac{\frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{h(x_{ki})\pi(x_{ki}) - \hat{\beta}_1^T \mathbf{g}(x_{ki})}{q_\alpha(x_{ki})}}{\frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{\pi(x_{ki}) - \hat{\beta}_2^T \mathbf{g}(x_{ki})}{q_\alpha(x_{ki})}}, \\ \hat{\mu}_{\text{MLE}} &= \frac{\frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{h(x_{ki})\pi(x_{ki})}{q_\alpha(x_{ki}) + \tilde{\zeta}^T \mathbf{g}(x_{ki})}}{\frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{\pi(x_{ki})}{q_\alpha(x_{ki}) + \tilde{\zeta}^T \mathbf{g}(x_{ki})}}, \end{aligned}$$

where

$$\begin{aligned} \hat{\beta}_1 &= \widetilde{\text{var}} \left( \frac{\mathbf{g}(X)}{q_\alpha(X)} \right)^{-1} \widetilde{\text{cov}}^T \left( \frac{h(X)\pi(X)}{q_\alpha(X)}, \frac{\mathbf{g}(X)}{q_\alpha(X)} \right) \quad \text{and} \\ \hat{\beta}_2 &= \widetilde{\text{var}} \left( \frac{\mathbf{g}(X)}{q_\alpha(X)} \right)^{-1} \widetilde{\text{cov}}^T \left( \frac{\pi(X)}{q_\alpha(X)}, \frac{\mathbf{g}(X)}{q_\alpha(X)} \right) \\ \tilde{\zeta} &= \underset{\zeta}{\text{argmax}} \sum_{k=1}^p \sum_{i=1}^{n_k} \log [q_\alpha(x_{ki}) + \zeta^T \mathbf{g}(x_{ki})]. \end{aligned}$$

*Remark 8.* The optimality of the above estimators can be seen by extending the optimality results of  $\hat{Z}_{\text{Reg}}$  in Owen and Zhou (2000) and  $\hat{Z}_{\text{MLE}}$  in Tan (2004) from scalar case to vector case. Specifically, under conditions (C1)–(C3) for  $\pi(x)$  and  $h(x)\pi(x)$ , the two estimators

$$\begin{aligned} &\left( \frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{h(x_{ki})\pi(x_{ki}) - \hat{\beta}_1^T \mathbf{g}(x_{ki})}{q_\alpha(x_{ki})} \right) \quad \text{and} \\ &\left( \frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{\pi(x_{ki}) - \hat{\beta}_2^T \mathbf{g}(x_{ki})}{q_\alpha(x_{ki})} \right) \\ &\left( \frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{h(x_{ki})\pi(x_{ki})}{q_\alpha(x_{ki}) + \tilde{\zeta}^T \mathbf{g}(x_{ki})} \right) \\ &\left( \frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{\pi(x_{ki})}{q_\alpha(x_{ki}) + \tilde{\zeta}^T \mathbf{g}(x_{ki})} \right) \end{aligned}$$

can be shown to be consistent and asymptotic normal with the minimum covariance matrix among all estimators in the form of

$$\begin{aligned} &\left( \frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{h(x_{ki})\pi(x_{ki})}{q_\alpha(x_{ki})} \right) - \left( \frac{\beta_1^T}{\beta_2^T} \right) \frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{\mathbf{g}(x_{ki})}{q_\alpha(x_{ki})} \\ &\left( \frac{1}{n} \sum_{k=1}^p \sum_{i=1}^{n_k} \frac{\pi(x_{ki})}{q_\alpha(x_{ki})} \right) \end{aligned}$$

for arbitrary real vectors  $\beta_1$  and  $\beta_2$ . Here  $A \geq B$  means  $A - B$  is nonnegative definite for two square matrices  $A$  and  $B$ . Then by the delta method, it is straightforward to show the optimality of  $\hat{\mu}_{\text{Reg}}$  and  $\hat{\mu}_{\text{MLE}}$ . Their asymptotic variances are identical and equal to

$$\sigma_\mu^2(\alpha) = \frac{1}{Z^2} \text{var}_\alpha \left( \frac{h(X)\pi(X) - \mu\pi(X) - \beta_\alpha^T \mathbf{g}(X)}{q_\alpha(X)} \right), \quad (13)$$

where

$$\beta_\alpha = \text{var} \left( \frac{\mathbf{g}(X)}{q_\alpha(X)} \right)^{-1} \text{cov}^T \left( \frac{h(X)\pi(X) - \mu\pi(X)}{q_\alpha(X)}, \frac{\mathbf{g}(X)}{q_\alpha(X)} \right).$$

### 4.2 Two-Stage Procedure For Ratio Estimators

Take  $\hat{\mu}_{\text{Reg}}$  and  $\hat{\mu}_{\text{MLE}}$  as functions of  $\alpha$  and denote by  $\hat{\mu}_{\text{Reg}}(\alpha)$  and  $\hat{\mu}_{\text{MLE}}(\alpha)$ . The two-stage procedure in Section 3 can be applied here:

1. First stage: Given initial proportion  $\gamma = (\gamma_1, \dots, \gamma_p)$  satisfying  $\sum_{k=1}^p \gamma_k = 1$ , generate  $n_0$  independent stratified

sample  $\{x_i\}_{i=1}^{n_0}$  from  $q_Y(x)$ . Obtain  $\hat{\alpha}$  by minimizing

$$\hat{\tau}^2(\alpha) = \frac{1}{n_0} \sum_{i=1}^{n_0} \frac{[h(x_i)\pi(x_i) - \hat{\mu}\pi(x_i) - \hat{\beta}_\alpha g(x_i)]^2}{q_\alpha(x_i)q_Y(x_i)}, \quad (14)$$

where

$$\hat{\mu} = \frac{1}{n_0} \sum_{i=1}^{n_0} \frac{h(x_i)\pi(x_i)}{q_Y(x_i)} \bigg/ \frac{1}{n_0} \sum_{i=1}^{n_0} \frac{\pi(x_i)}{q_Y(x_i)},$$

and

$$\hat{\beta}_\alpha = \left( \frac{1}{n_0} \sum_{i=1}^{n_0} \frac{g(x_i)g(x_i)^T}{q_\alpha(x_i)q_Y(x_i)} \right)^{-1} \times \left[ \frac{1}{n_0} \sum_{i=1}^{n_0} \frac{(h(x_i)\pi(x_i) - \hat{\mu}\pi(x_i))g(x_i)}{q_\alpha(x_i)q_Y(x_i)} \right],$$

with respect to  $\alpha$  over  $\Theta$ .

2. Second stage: Generate  $n - n_0$  independent stratified observations  $\{x_i\}_{i=n_0+1}^n$  from  $q_{\hat{\alpha}}(x)$ . Estimate integral  $\mu$  by  $\hat{\mu}_{\text{Reg}}(\tilde{\alpha})$  or  $\hat{\mu}_{\text{MLE}}(\tilde{\alpha})$  with all  $n$  observations, where  $\tilde{\alpha} = n_0/n \cdot \gamma + (n - n_0)/n \cdot \hat{\alpha}$ .

In the first stage,  $\hat{\tau}^2(\alpha)$  is the Monte Carlo estimate of  $Z^2\sigma_\mu^2(\alpha)$ . Similar to the results in Section 3.3.1,  $\hat{\mu}_{\text{Reg}}(\tilde{\alpha})$  and  $\hat{\mu}_{\text{MLE}}(\tilde{\alpha})$  for  $\mu$  have proper asymptotic results and the case for two proposal distributions is stated below.

*Theorem 3.* Under conditions (C1)–(C5) with  $\pi(x)$  replaced by  $h(x)\pi(x) - \mu\pi(x)$ ,  $\hat{\mu}_{\text{Reg}}(\tilde{\alpha})$  and  $\hat{\mu}_{\text{MLE}}(\tilde{\alpha})$  are consistent and

$$\sqrt{n}(\hat{\mu}_{\text{Reg}}(\tilde{\alpha}) - \mu) \xrightarrow{L} N(0, \sigma_\mu^2(\alpha^*))$$

and

$$\sqrt{n}(\hat{\mu}_{\text{MLE}}(\tilde{\alpha}) - \mu) \xrightarrow{L} N(0, \sigma_\mu^2(\alpha^*)),$$

where  $\alpha^*$  is the minimizer of  $\sigma_\mu^2(\alpha)$ .

### 4.3 Consideration of Selecting Component Proposal Distributions

In this article, we focus on finding the optimal mixture weights to construct a mixture proposal distribution for IS, assuming that the set of component proposals to be included in the mixture has been preselected. Since the proposed mixture proportion determination automatically discriminates the high-quality proposals from the poor ones, our procedure in a way alleviates the difficulty of selecting the set of proposal distributions. It also allows a larger set of proposals to be considered as the procedure serves as a selection tool. Nevertheless, preselection of the proposals is extremely important as it provides the basis for efficient inference of optimal mixture weights. This is an area of active research. Here we provide some remarks and practical guidance.

Consider the asymptotic variances of  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$ ,  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$ ,  $\hat{\mu}_{\text{Reg}}(\tilde{\alpha})$ , and  $\hat{\mu}_{\text{MLE}}(\tilde{\alpha})$  in (10) and (13). Owen and Zhou (2000) and Tan (2004) showed that when  $\pi(x)$  is a linear combination of the component proposals,  $\sigma_Z^2(\alpha) = 0$  for any  $\alpha$ . Therefore for estimating  $Z$ , it is preferred that the component proposals have a linear combination close to the shape of  $\pi(x)$ . This can be achieved by using proposals that separately approximate the modes and tails of  $\pi(x)$ . Alternatively, one can decompose  $\pi(x)$

into a linear combination

$$\pi(x) = \sum_{k=1}^r c_k \pi_k(x). \quad (15)$$

Then the component proposals can be obtained by approximating each  $\pi_k(x)$ . Owen and Zhou (2000) gave some illustrations of this strategy.

For the ratio estimator, it can be shown similarly that when  $h(x)\pi(x) - \mu\pi(x)$  is a linear combination of the component proposals,  $\sigma_\mu^2(\alpha) = 0$  for any  $\alpha$ . Therefore, the strategy used for estimating  $Z$  can be used here as well. In particular, we can find a decomposition

$$h(x)\pi(x) - \mu\pi(x) = \sum_{k=1}^r c_k h(x)\pi_k(x) - \sum_{k=1}^r \mu c_k^* \pi_k^*(x)$$

and find component proposals to approximate the individual terms. If  $h(x)$  takes negative values, additional terms corresponding to  $h(x) = h^+(x) - h^-(x)$  will be needed. Example 3 in Section 5 provides an illustration of this approach.

Another consideration is the tail requirement. For estimating  $Z$ ,  $q_{\alpha^*}(x)$  needs to have heavier tail than  $\pi(x)$ ; and for estimating  $\mu$ ,  $q_{\alpha^*}(x)$  needs to have heavier tail than  $h(x)\pi(x) - \mu\pi(x)$ . In cases where  $\pi(x)$ 's tail decreases exponentially, the requirements can be satisfied by including some Student's  $t$  distributions or other heavy-tailed distributions in the set of component proposals (Geweke 1989).

Oh and Berger (1993) and West (1993) proposed adaptive procedures to find better proposal distributions. Liang, Liu, and Carroll (2007) proposed a stochastic approximation procedure to partition the sample domain and used truncations of the target distribution in the subregions as component proposal distributions. The normalizing constant of each component is estimated in a pilot stage. These procedures can be used here for finding the component proposals in our setting. In fact, the pilot stage of our proposed procedure can also be used as well. The estimated optimal mixture weights from the pilot stage may provide hints on potentially useful proposals to be considered. For example, a large weight for a component proposal that mainly covers the tail in one direction may suggest to use additional proposals to cover the more extreme part of the tail in that direction. However, caution should be exercised when considering the removal of a proposal distribution because of its small weight, since it may be used to serve as a defensive proposal that guarantees finite variance of the IS estimator.

## 5. EMPIRICAL STUDIES

Here, we present several examples to illustrate the performance of the proposed procedure. In all examples, the standard restricted optimization algorithm BFGS (Broyden-Hletcher-Goldfarb-Shanno) (Battiti and Masulli 1990) is used in the pilot stage to find  $\hat{\alpha}$ .

*Example 1.* Let  $\phi(x; \sigma)$  be the normal density with mean 0 and standard error  $\sigma$ , and  $\psi_k(x)$  be the density of  $t$  distribution with degree of freedom  $k$ . In this example, we consider two target distributions and two sets of proposal distributions. The combination is listed in Table 1. The case (A1) represents the



Table 1. Parameter settings of four cases in Example 1

Target distribution	Proposal distributions			
	$q_1 = \prod_{i=1}^{10} \psi_k(x_i)$ and $q_2 = \prod_{i=1}^{10} \phi(x_i; \sigma)$			
	$k = 1$	$k = 1$	$k = 1$	$k = 2$
	$\sigma = 1.1$	$\sigma = 0.4$	$\sigma = 1$	$\sigma = 1$
$\prod_{i=1}^{10} \phi(x_i; 1)$	(A1)	(A2)		
$0.2 \prod_{i=1}^{10} \psi_4(x_i) + 0.8 \prod_{i=1}^{10} \phi(x_i; 1)$			(B1)	(B2)

situation that one of the proposal distribution,  $q_2(x)$ , is a good approximation to  $\pi^*(x)$  by itself, and  $q_1(x)$ , being a product of Cauchy distributions, is a relatively poor proposal. We expect that the two-stage procedure will be helpful to decrease the contamination of  $q_2(x)$ . The case (A2) represents the situation that both proposals are not good approximation to the target and an appropriate proportion is not immediately clear. Both (B1) and (B2) represent the situation that one of the proposals,  $q_2(x)$ , is a good approximation to the center of the target, but with a lighter tail, and the other proposal,  $q_1(x)$ , has a heavier tail, for protection. The case (B1) uses a more conservative protection (Cauchy) and (B2) is more aggressive ( $t_2$ ).

We compare five methods. The first three methods generate independent and stratified observations  $\{x_i\}_{i=1}^n$  from  $q_{\alpha_0}(x) = \alpha_0 q_1(x) + (1 - \alpha_0)q_2(x)$ , where  $\alpha_0 = (0.5, 1 - 0.5)$ . The last two methods generate independent and stratified observations  $\{x_i\}_{i=1}^{n_0}$  from  $q_{\alpha_0}(x)$  and  $\{x_i\}_{i=n_0+1}^n$  from  $q_{\tilde{\alpha}}(x)$ , where  $\tilde{\alpha}_1 = \alpha_0 n_0 / n + \hat{\alpha}_1(n - n_0) / n$  and  $\hat{\alpha}_1$  is obtained by the corresponding method. Since the simulation results of regression method are nearly identical to the likelihood approach, we only list MLE and 2MLE here. Specifically, the methods are as follows. For simplicity, only formulas for estimating  $Z$  are listed.

UIS (Unprotected Importance Sampling): This is estimator (1) with  $q(x) = q_2(x)$ .

SIS (Stratified Importance Sampling): This is estimator (3) with  $\alpha = \alpha_0$ .

MLE (MLE method): This is estimator (6) with  $\alpha = \alpha_0$ .

2SIS (Two-Stage Stratified Importance Sampling):

$$\frac{1}{n} \sum_{i=1}^n \frac{\pi(x_i)}{q_{\tilde{\alpha}}(x_i)}, \text{ where}$$

$$\hat{\alpha}_1 = \underset{\alpha}{\operatorname{argmin}} \left( \alpha_1 \widetilde{\operatorname{var}}_1 \left[ \frac{\pi(x)}{q_{\alpha_0}(x)} \right] + (1 - \alpha_1) \widetilde{\operatorname{var}}_2 \left[ \frac{\pi(x)}{q_{\alpha_0}(x)} \right] \right)$$

and  $\widetilde{\operatorname{var}}_k$  denotes the sample variance with the subset of  $\{x_i\}_{i=1}^{n_0}$ , which comes from  $q_k(x)$ . This is the method used in Raghavan and Cox (1998).

2MLE (Two-Stage MLE): This is our proposed method.

The results are shown in Table 2 for estimating  $Z$  and  $\mu$ . Simulation is replicated for 1000 times independently with  $n = 4000$  and  $n_0 = 400$  in each simulation. We report the means of  $\hat{Z}$  or  $\hat{\mu}$ , the means of  $\hat{\alpha}$  and the MSE

$$n\hat{V} = \frac{n}{1000} \sum_{i=1}^{1000} (\hat{Z}_i - Z)^2 \text{ or } \frac{n}{1000} \sum_{i=1}^{1000} (\hat{\mu}_i - \mu)^2,$$

where  $Z$  and  $\mu$  are theoretical values.

It is seen that, in (A1), where  $q_2$  is a good proposal by itself, 2SIS and 2MLE choose  $\alpha_1$  close or equal to the smallest allowed value (0.001) for  $q_1$ , which minimizes its contamination, and achieves the same efficiency as UIS (using the good proposal only). They are more efficient than SIS and MLE, which use equal proportions for both proposal distribution.

In (A2), both 2SIS and 2MLE choose  $\hat{\alpha}_1 = 0.98$ , giving much higher proportion to the heavy-tail  $t$  proposal. It is seen that the normal proposal has a much lighter tail ( $\sigma = 0.4$ ) than that the target ( $\sigma = 1$ ). In this case,  $q_1(x)$  is the better proposal. UIS, which uses  $q_2(x)$  exclusively, does not have finite variance. Comparing to one stage MLE, the two-stage procedure reduces MSE by about 43% and 34% for estimating  $Z$  and  $\mu$ , respectively.

In (B1) and (B2), UIS has the largest variance as expected. By using control variates, 2MLE and MLE perform much better than SIS and 2SIS. With the estimated mixture proportions, 2MLE reduces MSE by 10% and 30% for estimating  $Z$  in (B1) and (B2), respectively, comparing the one-stage MLE. Note that 2MLE obtains a larger estimated optimal proportion for  $q_1(x)$  in (B2) than in (B1). Intuitively this is because  $q_1(x)$  in (B2) is “closer” to the target integrand. In estimating  $\mu$ , 2MLE and

Table 2. Comparison of methods for Example 1, with each column for one setting

		$Z$				$\mu$			
Method		(A1)	(A2)	(B1)	(B2)	(A1)	(A2)	(B1)	(B2)
$\hat{\alpha}_1$	UIS	0	0	0	0	0	0	0	0
	SIS	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50
	2SIS	0.001	0.98	0.21	0.13	0.001	0.93	0.40	0.37
	MLE	0.50	0.50	0.50	0.50	0.50	0.50	0.50	0.50
	2MLE	0.004	0.98	0.72	0.999	0.001	0.91	0.42	0.30
$n\text{MSE}$	UIS	0.16	$9.4 \times 10^3$	3.0	0.47	0.19	$1.2 \times 10^3$	66	68
	SIS	0.45	28	0.15	0.16	0.34	3.2	0.38	0.27
	2SIS	0.15	16	0.087	0.028	0.20	2.1	0.37	0.26
	MLE	0.27	28	0.041	0.0094	0.34	3.2	0.37	0.16
	2MLE	0.15	16	0.037	0.0066	0.20	2.1	0.35	0.15

NOTE:  $\hat{\alpha}_1$  is the mean of 1000 estimated mixture proportions and MSE is mean square error of integral estimators.

Table 3. Comparison between finite sample and asymptotic results

	Z				$\mu$			
	(A1)	(A2)	(B1)	(B2)	(A1)	(A2)	(B1)	(B2)
$\hat{\alpha}_1$	0.004	0.98	0.72	0.999	0.001	0.91	0.42	0.30
$\alpha_1^*$	0.001	0.98	0.77	0.999	0.001	0.999	0.42	0.30
$n\hat{V}$	0.150	15.5	0.037	0.0066	0.20	2.06	0.35	0.15
$\sigma^2(\alpha^*)$	0.155	15.9	0.035	0.0061	0.19	1.97	0.36	0.15

NOTE:  $\alpha_1^*$  is the mixture proportion giving the minimum asymptotic variance,  $\hat{V}$  is the sample variance of integral estimators, and  $\sigma^2(\alpha^*)$  is the minimum asymptotic variance.

MLE perform better than SIS and 2SIS, but the two-stage 2MLE and one-stage MLE are similar, since the estimated optimal proportions are close to 0.5.

To check the convergence properties of 2MLE, in all four cases we report, in Table 3, a comparisons between the theoretical minimum asymptotic variances and the sample variance of 2MLE, as well as a comparison between the optimal proportions and the average estimated proportions. It is seen that both of them are quite close to the optimal values.

*Example 2.* Consider a rare event problem in Hesterberg (1995). Let  $X$  be a three-dimensional random variable with independent components  $(X_1, X_2, X_3)$  and

$$X = (X_1, X_2, X_3) = \max(0, Y_1 + 10d - Z_1 - Z_2 - \max(500, 3000 - Y_2 - 40d)),$$

where  $Y_1 \sim N((1600, 1650, 1600), 100^2 I_3)$ ,  $Y_2 \sim N((1600, 1700, 1600), 100^2 I_3)$ ,  $Z_1 \sim \Gamma(100\mathbf{1}_3, (5, 6, 7))$  with  $\Gamma(\text{scale}, \text{shape})$  denoting the gamma distribution,  $Z_2$  has density proportional to  $e^{x/100} I_{x \in (0, 300)}$ , and  $d = \max(0, 60 - t)$ , where  $t \sim N((54, 52, 55), 5^2 I_3)$ . Denote the density of  $X$  to be  $f(x) = \prod_{j=1}^3 f_j(x_j)$ . The targets of interest are

$$P = P \left[ \sum_{i=1}^3 X_i > 1200 \right]$$

and

$$\mu = E \left[ 80 \cdot \max \left( \sum_{i=1}^3 X_i - 1200, 0 \right) \right].$$

The true value of  $P$  is about 0.003 and therefore the probability measures the area in the tail of  $f(x)$ . Hesterberg (1995) used  $\hat{Z}_{\text{SIS}}$  to estimate  $P$  and  $\mu$  and constructed the proposal distributions by exponential tilting, using  $q(x) = c(\beta) \exp(\sum_{j=1}^3 \beta_j x_j) f(x)$  with parameters  $\beta = (\beta_1, \beta_2, \beta_3)$ . Seven proposals are constructed by setting  $\beta = c \cdot (I_1, I_2, I_3)$ , where  $I_j$  is binary and  $c$  is set so that  $E[\sum_{i=1}^3 X_i]$  is equal to some predetermined value, where  $(X'_1, X'_2, X'_3)$  follows  $q(x)$ . Including  $f(x)$  as another proposal component, there are eight proposal components. Hesterberg (1995) provided preset mixture proportions for these proposals, listed in Table 4.

Here, we compare the proposed two-stage procedure with the estimator used in Hesterberg (1995) for estimating both  $P$  and  $\mu$ . The results are listed in Table 5. Again, simulation is replicated 1000 times independently with  $n = 4000$  and  $n_0 = 400$  in each simulation. We report the sample means and variances of  $\hat{P}$  and  $\hat{\mu}$ , and the means of the mixture proportion  $\hat{\alpha}_i$  for  $i = 1, \dots, 8$ .

Table 4. Parameters setting of the mixture proposal

Proposal	$(I_1, I_2, I_3)$	$E[\sum_{j=1}^3 X'_j]$	$\alpha_i$
$q_1(x) = f(x)$			0.5
$q_2(x)$	(1, 0, 0)	1416	0.0035
$q_3(x)$	(0, 1, 0)	1266	0.028
$q_4(x)$	(0, 0, 1)	1616	0.0005
$q_5(x)$	(1, 1, 0)	1482	0.236
$q_6(x)$	(1, 0, 1)	1832	0.018
$q_7(x)$	(0, 1, 1)	1682	0.0635
$q_8(x)$	(1, 1, 1)	1898	0.151

NOTE: Each  $q_i(x)$  is proportional to  $\exp(\sum_{j=1}^3 \beta_j x_j) f(x)$ , where  $\beta = c \cdot (I_1, I_2, I_3)$  and  $c$  is selected such that the expectation is equal to the corresponding expectation, for example, 1416.  $\alpha_i$  is the mixture proportion for  $q_i(x)$ . Here  $(X'_1, X'_2, X'_3)$  has density  $q_i(x)$ .

Comparing to SIS, it is seen that while the means are the same, 2MLE reduces the variance by 47% for estimating  $P$  and 44% for estimating  $\mu$ . When comparing the proportion set selected by 2MLE and the predetermined proportion set used by SIS, it is seen that some of the proposal considered to be important for SIS is also determined important by 2MLE, such as  $q_5(x)$  and  $q_8(x)$ . The major difference is that SIS puts too much proportion on  $q_1(x)$ , while 2MLE only selects a very small proportion for it, indicating that only a small proportion is needed for  $q_1(x)$  to guarantee the bounded estimating variance.

*Example 3.* In this example, we examine the performance of 2MLE on estimating VaR using a Bayesian GARCH(1,1) model for S&P500 index series. Given a probability  $p$  and a time horizon  $d$ , VaR is the value that a portfolio would encounter a loss greater than or equal to, with probability  $p$  over the horizon.

Suppose at time  $T$ , we have historical log returns  $y = \{y_1, \dots, y_T\}$ . Let  $R(y_d) = \sum_{k=1}^d y_{T+k}$  be the cumulative return in the next  $d$  periods, where  $y_d = (y_{T+1}, \dots, y_{T+d})$  and denote  $F_{y_d}$  as the cumulative distribution function (CDF) of  $R$ . Then the  $d$  days ahead VaR is defined as

$$\text{VaR}_p = \inf \{x \in \mathbb{R} | F_{y_d}(x) \leq p\}.$$

VaR is a widely used measure of market risk (Duffie and Pan 1997; Jorion 1997). To obtain the CDF  $F_{y_d}$ , we model the return series using GARCH model (Engle 1982; Bollerslev 1986), a commonly used model for return series and modeling volatility dynamics. Specifically, we use a Bayesian GARCH(1,1) model with normal innovations (Geweke 1994; Bauwens and Lubrano 2008),

$$y_t = \varepsilon_t h_t^{1/2}, \quad \varepsilon_t \stackrel{\text{iid}}{\sim} N(0, 1), \quad h_t = \phi_0 + \phi_1 y_{t-1}^2 + \beta h_{t-1},$$

where  $\phi_0 \geq 0$ ,  $\phi_1 \geq 0$ , and  $\beta \geq 0$  and  $\phi_1 + \beta < 1$  to ensure stationarity. Following Geweke (1994), the prior distributions of  $\log \phi_0$  and  $(\phi_1, \beta)$  are selected to be  $N(a_0, \sigma_a^2)$  and  $U(\phi_1 \geq 0, \beta \geq 0, \phi_1 + \beta < 1)$ . Here  $\phi_0$  is transformed to have the real line as domain, and  $(\phi_1, \beta)$  follows a uniform distribution in the stationary domain. The hyperparameters  $a_0$  and  $\sigma_a^2$  are set to be 1 and 2, respectively. We also use the sample variance for  $h_0^2$  for simplicity.

The Bayesian approach has the advantage of taking into account of parameter estimation variability in the estimation of VaR. Due to the complexity, Monte Carlo method is used. Since

Table 5. Comparison between two methods of Example 2

				Mixture proportions							
Method		Mean	Var	$\hat{\alpha}_1$	$\hat{\alpha}_2$	$\hat{\alpha}_3$	$\hat{\alpha}_4$	$\hat{\alpha}_5$	$\hat{\alpha}_6$	$\hat{\alpha}_7$	$\hat{\alpha}_8$
$P$	SIS	$3.4 \times 10^{-3}$	$3.0 \times 10^{-8}$	0.500	0.0035	0.028	0.0005	0.236	0.018	0.064	0.151
	2MLE	$3.4 \times 10^{-3}$	$1.6 \times 10^{-8}$	0.001	0.0040	0.038	0.0009	0.420	0.051	0.170	0.310
$\mu$	SIS	41	1.8	0.500	0.0035	0.028	0.0005	0.236	0.018	0.064	0.151
	2MLE	41	1.0	0.001	0.002	0.021	0.0003	0.380	0.040	0.150	0.410

NOTES: SIS is the method of Hesterberg (1995) and 2MLE is our method.  $\bar{\mu}$  and  $\bar{P}$  are the means of 1000 point estimators,  $\hat{\alpha}_i$  are the average mixture proportions, and  $\hat{V}$  is the sample variance of 1000 estimators.

VaR is largely a tail property, an appropriate implementation of IS may significantly improve the efficiency. Although VaR is not in the form of integral, it can be estimated easily by empirical quantiles from the Monte Carlo samples. Note that, CDF and probability are in the form of integral. Related literatures about estimating VaR using IS can be found in Hoogerheide and Van Dijk (2010), Glasserman, Heidelberg, and Shahabuddin (2000), and Dunkel and Weber (2007).

The two-stage algorithm is tested to estimate VaR with  $p = 0.05$  and  $0.01$  and horizons 1, 2, and 5 days, corresponding to 4-, 5-, and 8-dimensional problem, as there are three parameters in the GARCH(1,1) model. Denote  $\theta = (\log \phi_0, \phi_1, \beta)$ . For each VaR, following the strategy discussed in Section 4.3, we construct the proposal distributions based on the asymptotic variance of the empirical posterior CDF at VaR

$$\sigma_p^2(\alpha) = \int \frac{[1_{\{R(y_d) \leq \text{VaR}\}}(y_d) - p] \pi(y_d, \theta) - \beta_\alpha^T g(y_d, \theta)]^2}{q_\alpha(y_d, \theta)} \times dy_d d\theta, \quad (16)$$

where  $\pi(y_d, \theta) = \prod_{k=1}^{T+d} p(y_k | y_{k-1}, \theta) p(\theta)$ ,  $p(y_k | y_{k-1}, \theta)$  is the innovation density and  $p(\theta)$  is the prior density of  $\theta$ .

Expression (16) is not the variance of the VaR estimator. However, Hoogerheide and Van Dijk (2010) showed that the asymptotic variance of  $\text{VaR}_p$  can be approximated by  $\sigma_p^2(\alpha)$  times a constant, which does not depend on the proposal density. Since it is difficult to sample from  $\pi(y_d, \theta)$  directly, we approximate it by the mixture of

$$q_1(y_d, \theta) = \prod_{k=1}^{T+d} p(y_k | y_{1:k-1}, \theta) q_N(\theta)$$

and

$$q_2(y_d, \theta) = q_*(y_{T+d} | y_{T+d-1}, \theta) \prod_{k=1}^{T+d-1} p(y_k | y_{1:k-1}, \theta) q_N(\theta),$$

where  $q_N(\theta)$  is the normal distribution with the mean vector being the MLE  $\hat{\theta}$  and the covariance matrix  $\Sigma_N$  being the negative inverse Hessian matrix of  $\pi(y, \theta)$  at  $\hat{\theta}$ , inflated by a constant to allow a wider coverage. We use  $q_*(y_{T+d} | y_{1:T+d-1}, \theta) \sim N(-h_{T+d}^{1/2}, h_{T+d})$  for the proposal  $q_2(y_d, \theta)$ . It tries to cover the tail (large loss on the last day of the horizon). Similar proposals can be constructed by considering other potential situations of large loss, but only this one is included in the current example.

With the approximation of  $\pi(y_d, \theta)$ , the heavier tail components can be constructed by modifying the tails of  $q_1$  and  $q_2$ . Then the following two proposals are included as the heavier

tail components:

$$q_3(y_d, \theta) = \prod_{k=1}^{T+d} p(y_k | y_{1:k-1}, \theta) q_t(\theta)$$

and

$$q_4(y_d, \theta) = q_*(y_{T+d} | y_{1:T+d-1}, \theta) \prod_{k=1}^{T+d-1} p(y_k | y_{1:k-1}, \theta) q_t(\theta),$$

where  $q_t(\theta)$  is the product of three location-scale generalization of  $t_1$  densities with the means being  $\hat{\theta}$  and the squared scale parameters being the diagonal elements of  $\Sigma_N$ , and  $q_*$  is the same as in the construction of  $q_2$ . Since  $\pi(y_d, \theta)/q_3(y_d, \theta) = p(\theta)/q_t(\theta)$  and  $p(\theta)$  is the prior distribution with exponentially decreasing tail,  $q_3(y_d, \theta)$  has heavier tail than  $\pi(y_d, \theta)$ . Similarly,  $q_3$  and  $q_4$  have heavier tail than  $q_1, q_2$ , and proposals below, and therefore only mixture proportions for  $q_3$  and  $q_4$  need to be restricted.

To incorporate the integrand as discussed in Section 4.3, we further extend  $q_1(y_d, \theta)$  and  $q_2(y_d, \theta)$  to include

$$q_5(y_d, \theta) \propto 1_{\{y_{T+d} \leq \text{VaR}_{0.05} - \sum_{k=1}^{d-1} y_{T+k}\}}(y_{T+d}) q_1(y_d, \theta)$$

and

$$q_6(y_d, \theta) \propto 1_{\{y_{T+d} \leq \text{VaR}_{0.05} - \sum_{k=1}^{d-1} y_{T+k}\}}(y_{T+d}) q_2(y_d, \theta),$$

Here the truncation is done only on  $y_{T+d}$ , instead of the more accurate but computationally expensive truncation of  $\sum_{k=1}^d y_{T+k} \leq \text{VaR}_{0.05}$  under joint normal distribution.

The estimation of  $\text{VaR}_{0.01}$  can be done simultaneously by including the following component proposals:

$$q_7(y_d, \theta) \propto 1_{\{y_{T+d} \leq \text{VaR}_{0.01} - \sum_{k=1}^{d-1} y_{T+k}\}}(y_{T+d}) q_1(y_d, \theta)$$

and

$$q_8(y_d, \theta) \propto 1_{\{y_{T+d} \leq \text{VaR}_{0.01} - \sum_{k=1}^{d-1} y_{T+k}\}}(y_{T+d}) q_2(y_d, \theta).$$

Overall,  $q_1(y_d, \theta)$  to  $q_8(y_d, \theta)$  are used as component proposal distributions.

Since our objective is to estimate  $\text{VaR}_{0.05}$  and  $\text{VaR}_{0.01}$  simultaneously, in the pilot stage we estimate the optimal mixture proportions by minimizing the sum of variances of the two estimators. Since  $q_5, \dots, q_8$  involve the unknown  $\text{VaR}_{0.05}$  and  $\text{VaR}_{0.01}$ , the first stage sampling is modified as follows.

1. Generate pilot samples from  $q_1$  to  $q_4$  with sample size  $n_0/8$  each.
2. Estimate  $\text{VaR}_{0.05}$  using the pilot samples from step 1. Replace  $\text{VaR}_{0.05}$  in  $q_5$  and  $q_6$  with the estimate and generate pilot samples from them, with sample size  $n_0/8$  each.

Table 6. Comparison between MLE and 2MLE in Example 3

Horizon	Method	$p = 0.05$		$p = 0.01$	
		$\widehat{\text{VaR}}$	$\widehat{V}$	$\widehat{\text{VaR}}$	$\widehat{V}$
1 day	MLE	-1.332	14e-5	-1.894	20e-5
	2MLE	-1.333	3.8e-5	-1.895	4.6e-5
2 days	MLE	-1.886	5.1e-4	-2.773	12e-4
	2MLE	-1.886	1.5e-4	-2.771	3.5e-4
5 days	MLE	-2.997	17e-4	-4.432	5.9e-3
	2MLE	-2.996	5.4e-4	-4.424	1.8e-3

NOTE:  $\widehat{\text{VaR}}$  is the average of 300 point estimators and  $\widehat{V}$  is the sample variance of 300 estimators.

3. Estimate  $\text{VaR}_{0.01}$  using the pilot samples from steps 1 and 2. Replace  $\text{VaR}_{0.01}$  in  $q_7$  and  $q_8$  with the estimate and generate pilot samples from them, with sample size  $n_0/8$  each.
4. Obtain  $\widehat{\alpha}$  by minimizing  $\widehat{\tau}_{0.05}^2(\alpha) + \widehat{\tau}_{0.01}^2(\alpha)$ , where  $\widehat{\tau}_p^2(\alpha)$  is the estimator for  $\sigma_p^2(\alpha)$  using all samples in the first three steps.

Here, we compare the two-stage procedure 2MLE with the one-stage MLE with equal mixture proportions. The log returns of S&P500 index from September 28, 2010 to July 13, 2011 are used, with total 200 observations. The simulation is replicated for 300 times independently with  $n = 4 \times 10^6$  and  $n_0 = 8 \times 10^4$  in each simulation.  $\delta$  is selected to be 0.001.

The summary of estimation results and the estimated mixture proportions  $\widehat{\alpha}$  are listed in Tables 6 and 7. From Table 6, it is seen that 2MLE's Monte Carlo variance is about 23%–32% of the variance of MLE, while there is almost no difference in the mean. Table 7 shows that the two-stage algorithm assigns most of the mixture proportions to  $q_3$  and  $q_4$ , which indicates that these two heavy-tail component proposals are more important than the others. This is probably because  $q_1$  and  $q_2$  do not

Table 7. Summary of mixture proportions estimated from stage 1

VaR	$\widehat{\alpha}_1$	$\widehat{\alpha}_2$	$\widehat{\alpha}_3$	$\widehat{\alpha}_4$	$\widehat{\alpha}_5$	$\widehat{\alpha}_6$	$\widehat{\alpha}_7$	$\widehat{\alpha}_8$
1 day	1e-3	8e-4	3.4e-1	6.2e-1	9e-3	1.7e-2	8e-3	2e-3
2 days	1e-3	1e-3	3.5e-1	6.0e-1	2.4e-2	3e-3	9e-3	2e-3
5 days	1e-3	6e-4	4.3e-1	5.4e-1	2.0e-2	8e-4	4e-3	5e-4

NOTE: The average over 300 simulations are reported.  $\widehat{\alpha}_1$  to  $\widehat{\alpha}_8$  correspond to the mixture proportions assigned to  $q_1$  to  $q_8$ .

cover the high-density area of target distribution sufficiently, resulted in the preference to  $q_3$  and  $q_4$ . Compared with MLE, the optimization in the pilot stage of 2MLE requires additional computing time, which is about 20% more in practice.

Finally, we report some interesting insights on the comparison between MLE and 2MLE. By multiplying a scaling constant  $c^2$  to the covariance matrix  $\Sigma_N$  used in the proposal  $q_1$ , all the related component proposals are made either more dispersed for  $c > 1$  or more concentrated for  $c < 1$ . Since  $\sigma_p^2(\alpha)$  is proportional to the estimation variance of VaR and the proportion does not depend on the proposal density, the trajectories of estimated  $\sigma_p^2(\alpha^*)$  and  $\sigma_p^2(\alpha_0)$  as function of  $c$  are given in Figure 1 to illustrate how the quality of proposal distribution affect the performance of 2MLE and MLE.

It is seen that 2MLE is always better than MLE. Most interestingly, it shows that the performance of both methods depends on the quality of the proposal distributions, but 2MLE is much less sensitive to the proposal distributions and has more robust performance than MLE. This is due to 2MLE's ability to automatically adjust mixture proportion for the most efficient estimation. The simulation results (not shown here) show that, when  $c$  is small, 2MLE tends to assign most of the mixture proportions to the heavy tail  $q_3$  and  $q_4$ . This insight reinforces the notion that the two-stage approach not only improves upon the one-stage approach, but also alleviates to some extent the difficulty of selecting proposal distributions.

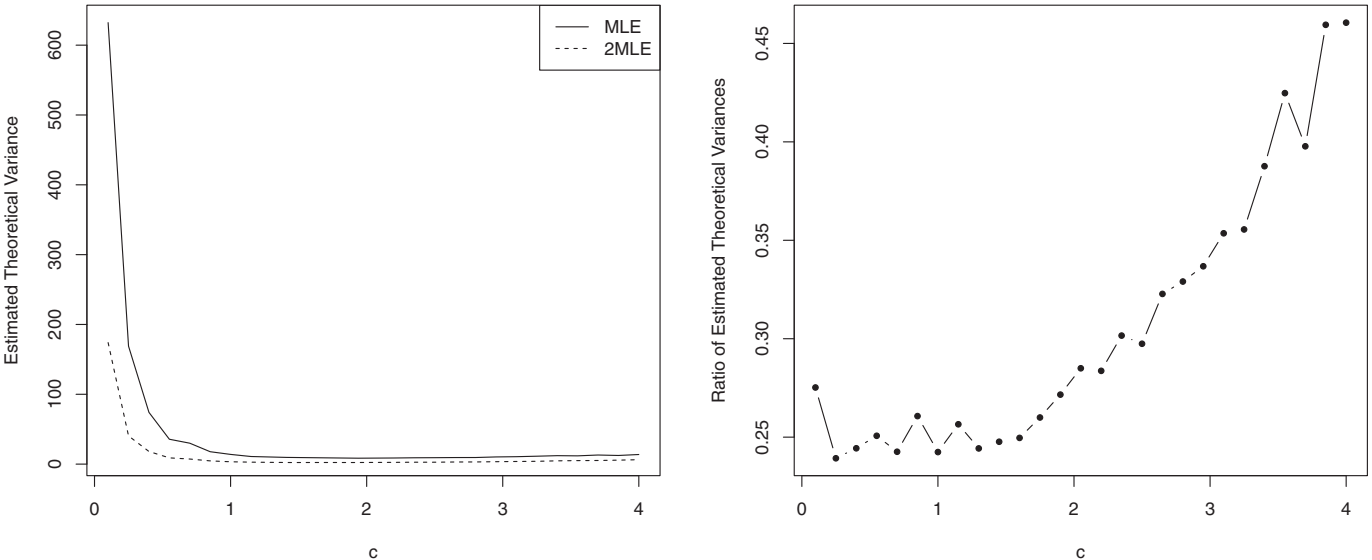


Figure 1. The left figure gives trajectories of estimated  $\sigma_p^2(\alpha^*)$  and  $\sigma_p^2(\alpha_0)$ , corresponding to MLE and 2MLE methods, respectively, with respect to the scaling constant  $c$ .  $c$  ranges from 0.1 to 4. For each  $c$ , the theoretical variances are estimated using one Monte Carlo sample, and the average over 10 replicates is reported. The right figure gives the trajectory of ratio of estimated  $\sigma_p^2(\alpha^*)$  over estimated  $\sigma_p^2(\alpha_0)$ .



## 6. SUMMARY

In this article, we proposed a two-stage procedure to select the optimal mixture proportions for the regression estimator in Owen and Zhou (2000) and MLE estimator in Tan (2004), and established the corresponding theoretical framework. The two-stage procedure significantly improved the existing methods in four aspects. First, the proposed estimator is asymptotically the best among all the estimators proposed in Owen and Zhou (2000), Tan (2004), and Raghavan and Cox (1998). Second, the criterion function of our pilot stage optimization is convex in its arguments, and therefore it is guaranteed that the optimization converges to the global minimum. Third, since there is no simple intuition in selecting the proportions for Owen and Zhou's (2000) regression estimator and Tan's (2004) MLE estimator, the proposed automatic procedure makes it much easier and safer to use mixture distributions for IS. Finally, the automatic determination of the mixture proportion alleviates the difficulty of choosing the set of proposal distributions to be considered in the mixture, as it serves as a selection and discrimination tool and hence allows users to include more potential proposal distributions for consideration.

## APPENDIX A. PROOF OF RESULTS

For simplicity, we only consider two proposal distributions. Then  $\alpha = (\alpha_1, 1 - \alpha_1)$  and  $\gamma = (\gamma_1, 1 - \gamma_1)$ . The proofs can be extended to the case of more than two proposals. To begin with, we establish the consistency of  $\hat{\alpha} = (\hat{\alpha}_1, 1 - \hat{\alpha}_1)$ . Note that  $\hat{\alpha}_1$  is equivalently a component of the bivariate  $M$ -estimate  $(\hat{\alpha}_1, \hat{\beta}) = \operatorname{argmin}_{\alpha, \beta} n_0^{-1} \sum_{i=1}^{n_0} m(x_i; \alpha, \beta)$ , where  $m(x; \alpha_1, \beta) = [\pi(x) - \beta g(x)]^2 / [q_\alpha(x) q_\gamma(x)]$ . Let  $M(\alpha_1, \beta) = \int m(x; \alpha_1, \beta) q_\gamma(x) dx$  and  $(\alpha_1^*, \beta^*) = \operatorname{argmin}_{\alpha_1, \beta} M(\alpha_1, \beta)$ . Meanwhile,  $M(\alpha_1, \beta)$  and  $\sigma_Z^2(\alpha)$  are strictly convex functions.

*Lemma 1.* It holds that

$$\begin{aligned} (\hat{\alpha}_1, \hat{\beta}) &\xrightarrow{P} (\alpha_1^*, \beta^*), \\ (\hat{\alpha}_1, \hat{\beta}) &= (\alpha_1^*, \beta^*) - \frac{1}{2\sqrt{n_0}} V^{-1} \hat{U} + o_p\left(\frac{1}{\sqrt{n_0}}\right), \end{aligned}$$

where  $V$  and  $\hat{U}$  are given in Lemma 2. Then  $\tilde{\alpha} = (\tilde{\alpha}_1, 1 - \tilde{\alpha}_1) \xrightarrow{P} (\alpha_1^*, 1 - \alpha_1^*)$ . Meanwhile,  $M(\alpha_1, \beta)$ ,  $\sigma_Z^2(\alpha)$ , and  $\hat{\sigma}^2(\alpha)$  are strictly convex functions.

*Proof.* Note that  $m(x; \alpha, \beta)$  is convex since its Hessian matrix

$$D^2 m(x; \alpha_1, \beta) = \frac{2g(x)^2}{q_\alpha(x)q_\gamma(x)} \begin{pmatrix} \frac{(\pi(x) - \beta g(x))^2}{q_\alpha(x)^2} & \frac{\pi(x) - \beta g(x)}{q_\alpha(x)} \\ \frac{\pi(x) - \beta g(x)}{q_\alpha(x)} & 1 \end{pmatrix}$$

is a positive semidefinite matrix. Then the consistency of  $(\hat{\alpha}_1, \hat{\beta})$  can be proved by verifying conditions 1–3 in Haberman (1989) for  $M$ -estimators by convex minimization. First, the parameter set  $\Theta = [\delta, 1 - \delta] \times \mathbb{R}$  of  $(\alpha_1, \beta)$  is convex and closed. Second,  $(\alpha_1^*, \beta^*)$  is unique. By Durrett (1996, Appendix 9), the differentiation and integration in  $M(\alpha, \beta)$  can be exchanged so that

$$\begin{aligned} D^2 M(\alpha_1, \beta) &= \begin{pmatrix} 2 \int \frac{[\pi(x) - \beta g(x)]^2 g(x)^2}{q_\alpha(x)^3} dx & 2 \int \frac{[\pi(x) - \beta g(x)] g(x)^2}{q_\alpha(x)^2} dx \\ 2 \int \frac{[\pi(x) - \beta g(x)] g(x)^2}{q_\alpha(x)^2} dx & 2 \int \frac{g(x)^2}{q_\alpha(x)} dx \end{pmatrix}. \end{aligned}$$

For any bivariate vector  $v$ ,  $v^T \{D^2 M(\alpha_1, \beta)\} v \geq 0$  and the equality holds only when  $\pi(x) \equiv c_1 q_1(x) + c_2 q_2(x)$  for some  $c_1$  and  $c_2$ . By condition (C5),  $D^2 M(\alpha_1, \beta)$  is positive definite. Therefore,  $M(\alpha_1, \beta)$  is strictly convex and  $(\alpha_1^*, \beta^*)$  is unique. Third, let  $W = (\delta, 1 - \delta) \times \mathbb{R}$ . By condition (C3),  $M(\alpha_1, \beta) < \infty$  for any  $(\alpha_1, \beta) \in W$ .

The expansion of  $(\hat{\alpha}_1, \hat{\beta})$  can be found in the proof of Haberman (1989, Theorem 6.1) by verifying his conditions 7 and 10. First,  $D^2 M(\alpha_1^*, \beta^*)$  is positive definite as mentioned above. Second, the gradient of  $m(x; \alpha_1, \beta)$  satisfies  $E|Dm(x; \alpha_1, \beta)|^2 < \infty$ . Therefore, the convergence of  $(\hat{\alpha}_1, 1 - \hat{\alpha}_1)$  holds because  $\tilde{\alpha} = \gamma n_0/n + \hat{\alpha}(n - n_0)/n$  and  $n_0 = o(n)$ .

Finally, with the strict convexity of  $M(\alpha_1, \beta)$ , which is stated above, the strict convexity of  $\sigma_Z^2(\alpha)$  can be seen by the facts that  $\sigma_Z^2(\alpha) = \min_{\beta} M(\alpha_1, \beta)$  and

$$\begin{aligned} \min_{\beta} M(\lambda \alpha_1 + (1 - \lambda) \alpha_2, \beta) \\ = \min_{\beta_1} \min_{\beta_2} M(\lambda \alpha_1 + (1 - \lambda) \alpha_2, \lambda \beta_1 + (1 - \lambda) \beta_2) \end{aligned}$$

for any  $\alpha_1, \alpha_2$  and  $\lambda \in [0, 1]$ . The strict convexity of  $\hat{\sigma}^2(\alpha)$  can be proved similarly.  $\square$

The following expansion of  $(\hat{\alpha}_1, \hat{\beta})$  will be used in the higher-order calculation of  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$  and  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$ .

*Lemma 2.* It holds that

$$\begin{aligned} (\hat{\alpha}_1, \hat{\beta}) &= (\alpha_1^*, \beta^*) + \frac{1}{2\sqrt{n_0}} V^{-1} \hat{U} - \frac{1}{2n_0} V^{-1} \\ &\quad \times \{(\hat{V} - \sqrt{n_0} V) V^{-1} \hat{U} + \hat{W}\} + o_p\left(\frac{1}{n_0}\right), \end{aligned}$$

where  $\hat{W}$  is a random variable of order  $O_p(1)$ ,

$$\begin{aligned} \hat{U} &= \begin{pmatrix} \frac{1}{\sqrt{n_0}} \sum_{i=1}^{n_0} \frac{(\pi(x_i) - \beta^* g(x_i))^2 g(x_i)}{q_{\alpha^*}(x_i)^2 q_\gamma(x_i)} \\ \frac{1}{\sqrt{n_0}} \sum_{i=1}^{n_0} \frac{2(\pi(x_i) - \beta^* g(x_i)) g(x_i)}{q_{\alpha^*}(x_i) q_\gamma(x_i)} \end{pmatrix}, \\ V &= \begin{pmatrix} \int \frac{[\pi(x) - \beta^* g(x)]^2 g(x)^2}{q_{\alpha^*}(x)^3} dx & \int \frac{[\pi(x) - \beta^* g(x)] g(x)^2}{q_{\alpha^*}(x)^2} dx \\ \int \frac{[\pi(x) - \beta^* g(x)] g(x)^2}{q_{\alpha^*}(x)^2} dx & \int \frac{g(x)^2}{q_{\alpha^*}(x)} dx \end{pmatrix}, \end{aligned}$$

and

$$\hat{V} = \begin{pmatrix} \frac{1}{\sqrt{n_0}} \sum_{i=1}^{n_0} \frac{[\pi(x_i) - \beta^* g(x_i)]^2 g(x_i)^2}{q_{\alpha^*}(x_i)^3 q_\gamma(x_i)} & \frac{1}{\sqrt{n_0}} \sum_{i=1}^{n_0} \frac{[\pi(x_i) - \beta^* g(x_i)] g(x_i)^2}{q_{\alpha^*}(x_i)^2 q_\gamma(x_i)} \\ \frac{1}{\sqrt{n_0}} \sum_{i=1}^{n_0} \frac{[\pi(x_i) - \beta^* g(x_i)] g(x_i)^2}{q_{\alpha^*}(x_i)^2 q_\gamma(x_i)} & \frac{1}{\sqrt{n_0}} \sum_{i=1}^{n_0} \frac{g(x_i)^2}{q_{\alpha^*}(x_i) q_\gamma(x_i)} \end{pmatrix}.$$

*Proof.* Note that  $\hat{U} = n_0^{-1/2} \sum_{i=1}^{n_0} Dm(x_i; \alpha_1^*, \beta^*)$ ,  $\hat{V} = n_0^{-1/2} \sum_{i=1}^{n_0} D^2 m(x_i; \alpha_1^*, \beta^*)$ , and  $V = \int D^2 m(x; \alpha_1^*, \beta^*) dx$ . Then by Taylor expansion around  $(\alpha_1^*, \beta^*)$  on  $n_0^{-1} \sum_{i=1}^{n_0} Dm(x_i; \hat{\alpha}_1, \hat{\beta}) = 0$  and the convergence of  $(\hat{\alpha}_1, \hat{\beta})$ , we have

$$0 = -\frac{1}{\sqrt{n_0}} \hat{U} + \frac{2}{\sqrt{n_0}} \hat{V} \begin{pmatrix} \hat{\alpha}_1 - \alpha_1^* \\ \hat{\beta} - \beta^* \end{pmatrix} + \frac{1}{n_0} \hat{W} + o_p\left(\frac{1}{n_0}\right),$$

then

$$\begin{pmatrix} \hat{\alpha}_1 - \alpha_1^* \\ \hat{\beta} - \beta^* \end{pmatrix} = \frac{1}{2\sqrt{n_0}} V^{-1} \hat{U} - \frac{1}{n_0} V^{-1} \left\{ (\hat{V} - \sqrt{n_0} V) \cdot \sqrt{n_0} \begin{pmatrix} \hat{\alpha}_1 - \alpha_1^* \\ \hat{\beta} - \beta^* \end{pmatrix} + \hat{W} \right\} + o_p \left( \frac{1}{n_0} \right).$$

The expansion of  $(\hat{\alpha}_1, \hat{\beta})$  follows by substituting  $(\hat{\alpha}_1 - \alpha_1^*, \hat{\beta} - \beta^*)$  to the right-hand side (RHS) of the above equation.  $\square$

The combined sample  $\{x_1, \dots, x_n\}$  can be split into four parts by distributions  $q_1$  or  $q_2$  and first or second stages. Denote  $I_{jk}$  to be the index set of observations from the  $j$ th stage and  $q_k$ , that is,  $I_{11} = \{1, \dots, n_0\gamma_1\}$ ,  $I_{12} = \{n_0\gamma_1 + 1, \dots, n_0\}$ ,  $I_{21} = \{n_0 + 1, \dots, n_0 + \lceil (n - n_0)\hat{\alpha}_1 \rceil\}$ ,  $I_{22} = \{n_0 + \lceil (n - n_0)\hat{\alpha}_1 \rceil + 1, \dots, n\}$ , where  $\lceil x \rceil$  means the largest integer smaller than  $x$ , and  $n_{jk}$  to be the size of  $I_{jk}$ . Here, we can use for the index, where  $\lfloor x \rfloor$  is the largest integer smaller than  $x$ , to define  $I_{jk}$ . But for investigating the asymptotic behavior, the difference can be ignored. We will use the decomposition

$$\begin{aligned} G_n \tau(x) &= \sqrt{\frac{n_0}{n}} \left\{ \sum_{k=1}^2 \sqrt{\gamma_k} \cdot \sqrt{n_{1k}} \left( \frac{1}{n_{1k}} \sum_{i \in I_{1k}} \tau(x_i) - \int \tau(x) q_k(x) dx \right) \right\} \\ &\quad + \sqrt{\frac{n - n_0}{n}} \left\{ \sum_{k=1}^2 \sqrt{\hat{\alpha}_k} \cdot \sqrt{n_{2k}} \left( \frac{1}{n_{2k}} \sum_{i \in I_{2k}} \tau(x_i) - \int \tau(x) q_k(x) dx \right) \right\} \\ &\equiv \sqrt{\frac{n_0}{n}} \{ \sqrt{\gamma_1} G_{11} \tau(x) + \sqrt{\gamma_2} G_{12} \tau(x) \} \\ &\quad + \sqrt{\frac{n - n_0}{n}} \{ \sqrt{\hat{\alpha}_1} G_{21} \tau(x) + \sqrt{\hat{\alpha}_2} G_{22} \tau(x) \}. \end{aligned} \quad (A.1)$$

The following lemma shows the convergence of  $\hat{Z}_{\text{SIS}}$  with  $\tilde{\alpha}$  as mixture proportion.

**Lemma 3.** For any integrable function  $h(x)$  satisfying  $\text{var}_{\alpha}[h(X)/q_{\alpha}(X)] < \infty$  for every  $\alpha_1 \in [\delta, 1 - \delta]$ , it holds that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \frac{h(x_i)}{q_{\tilde{\alpha}}(x_i)} &\xrightarrow{P} \int h(x) dx, \\ \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n \frac{h(x_i)}{q_{\tilde{\alpha}}(x_i)} - \int h(x) dx \right) &\xrightarrow{L} N \left( 0, \sum_{k=1}^2 \alpha_k^* \text{var}_k \left[ \frac{h(X)}{q_{\alpha^*}(X)} \right] \right), \end{aligned}$$

where  $\text{var}_k$  denote the variance under distribution density  $q_k(x)$ .

*Proof.* We only need to prove asymptotic normality since it implies the consistency. Using decomposition (A.1), we have

$$\begin{aligned} \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n \frac{h(x_i)}{q_{\tilde{\alpha}}(x_i)} - \int h(x) dx \right) &= \sqrt{\frac{n_0}{n}} \left\{ \sqrt{\gamma_1} G_{11} \left( \frac{h(x)}{q_{\tilde{\alpha}}(x)} \right) + \sqrt{\gamma_2} G_{12} \left( \frac{h(x)}{q_{\tilde{\alpha}}(x)} \right) \right\} \\ &\quad + \sqrt{\frac{n - n_0}{n}} \left\{ \sqrt{\hat{\alpha}_1} G_{21} \left( \frac{h(x)}{q_{\tilde{\alpha}}(x)} \right) + \sqrt{\hat{\alpha}_2} G_{22} \left( \frac{h(x)}{q_{\tilde{\alpha}}(x)} \right) \right\}. \end{aligned} \quad (A.2)$$

The asymptotic normality will be implied by showing that the first two terms in (A.2) are of order  $o_p(1)$  and the remaining is asymptotic normal.

For the first two terms, we prove that the collection of functions  $\{h(x)/q_{\alpha}(x)\}_{\alpha_1 \in [\delta, 1 - \delta]}$  is a Donsker class under either probability measures  $q_1$  or  $q_2$  by verifying the three conditions in van der Vaart (2000, Example 19.7). In fact, the parameter  $\alpha$  is in a bounded set;  $|h(x)/q_{\alpha_1}(x) - h(x)/q_{\alpha_2}(x)| \leq |m(x, \alpha_1, \alpha_2)| \cdot |\alpha_1 - \alpha_2|$  for every  $\alpha_1, \alpha_2$ , where  $m(x, \alpha_1, \alpha_2) = h(x)g(x)/(q_{\alpha_1}(x)q_{\alpha_2}(x))$ , and

$\int |m(x, \alpha_1, \alpha_2)|^2 q_k(x) dx < \infty$ . By van der Vaart (2000, Lemma 19.24) and Lemma 1, we have  $G_{1k}(h/q_{\tilde{\alpha}}) = G_{1k}(h/q_{\alpha^*}) + o_p(1)$ ,  $k = 1, 2$ . Then by Central Limit Theorem and  $n_0 = o(n)$ , the first two terms in (A.2) are of order  $o_p(1)$ .

For the last two terms, similarly, we argue that  $G_{2k}(h/q_{\tilde{\alpha}}) = G_{2k}(h/q_{\alpha^*}) + o_p(1)$  by a modification of van der Vaart (2000, Lemma 19.24) to handle random sample size. In fact, the key condition for his results, namely, weak convergence of  $G_{2k}(h/q_{\alpha^*})$ , is guaranteed by van Der Vaart and Wellner (1996, Theorem 3.5.1). Then by the independence between  $\hat{\alpha}_1$  and observations in  $\{x_i\}_{i=n_0+1}^n$  and an extension of Chow and Teicher (2003, sec. 9.4), we have

$$\begin{pmatrix} G_{21}(h/q_{\tilde{\alpha}}) \\ G_{22}(h/q_{\tilde{\alpha}}) \end{pmatrix} \xrightarrow{L} N \left( 0, \begin{pmatrix} \text{var}_1 \left[ \frac{h(X)}{q_{\alpha^*}(X)} \right] & 0 \\ 0 & \text{var}_2 \left[ \frac{h(X)}{q_{\alpha^*}(X)} \right] \end{pmatrix} \right)$$

and by Slutsky's theorem,  $\sqrt{\hat{\alpha}_1} G_{21}(\frac{h(x)}{q_{\tilde{\alpha}}(x)}) + \sqrt{\hat{\alpha}_2} G_{22}(\frac{h(x)}{q_{\tilde{\alpha}}(x)}) \xrightarrow{L} N(0, \sum_{k=1}^2 \alpha_k^* \text{var}_k[\frac{h(X)}{q_{\alpha^*}(X)}])$ . Therefore the lemma holds.  $\square$

In the above proof, only the consistency of  $\tilde{\alpha}$  is used. If  $\tilde{\alpha}$  is replaced by other consistent mixture proportion, the convergence properties still hold.

**Corollary 1.** For any  $\alpha$  satisfying  $\alpha \xrightarrow{P} \alpha^*$ , it holds that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \frac{h(x_i)}{q_{\alpha}(x_i)} &\xrightarrow{P} \int h(x) dx, \\ \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n \frac{h(x_i)}{q_{\alpha}(x_i)} - \int h(x) dx \right) &\xrightarrow{L} N \left( 0, \sum_{k=1}^2 \alpha_k^* \text{var}_k \left[ \frac{h(X)}{q_{\alpha^*}(X)} \right] \right). \end{aligned}$$

We also need the convergence of  $\tilde{\zeta}$  in  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$ , where  $\tilde{\zeta} = \text{argmin}_{\zeta} \sum_{i=1}^n \log[q_{\tilde{\alpha}}(x_i) + \zeta^T g(x_i)]$ .

**Lemma 4.** The following convergence properties for  $\tilde{\zeta}$  hold:

$$\tilde{\zeta} \xrightarrow{P} 0 \text{ and } \sqrt{n} \tilde{\zeta} \xrightarrow{L} N \left( 0, \frac{\sum_{k=1}^2 \alpha_k^* \text{Var}_k[g(X)/q_{\alpha^*}(X)]}{(\int g(x)^2/q_{\alpha^*}(x) dx)^2} \right).$$

*Proof.* The random variable  $\sqrt{n} \tilde{\zeta}$  is the minimizer of convex function  $\psi(s) = \sum_{i=1}^n \log[q_{\tilde{\alpha}}(x_i) + g(x_i)s/\sqrt{n}]$ . By verifying the condition of Hjort and Pollard (1994, basic corollary), we have the expansion

$$\sqrt{n} \tilde{\zeta} = \frac{n^{-1/2} \sum_{i=1}^n g(x_i)/q_{\tilde{\alpha}}(x_i)}{\int g(x)^2/q_{\alpha^*}(x) dx} + o_p(1)$$

and then the convergence of  $\tilde{\zeta}$  follows by Lemma 3. By Taylor expansion around 0, we have

$$\begin{aligned} \psi(s) &= \sum_{i=1}^n \log[q_{\tilde{\alpha}}(x_i)] + \left\{ \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{g(x_i)}{q_{\tilde{\alpha}}(x_i)} \right\} s \\ &\quad - \left\{ \frac{1}{2n} \sum_{i=1}^n \frac{g(x_i)^2}{q_{\tilde{\alpha}}(x_i)^2} \right\} s^2 + R_n(s), \end{aligned}$$

where

$$R_n(s) = \left\{ \frac{1}{3n\sqrt{n}} \sum_{i=1}^n \frac{g(x_i)^3}{q_{\tilde{\alpha}}(x_i)^3} \right\} s^3$$

and  $\xi_1$  is between  $\tilde{\alpha}_1$  and  $\tilde{\alpha}_1 + s/\sqrt{n}$ .

For fixed  $s$ ,  $\xi \xrightarrow{P} \alpha^*$  by  $\tilde{\alpha} \xrightarrow{P} \alpha^*$ . Then  $R_n(s) \xrightarrow{P} 0$  by Corollary 1, and condition of Hjort and Pollard (1994, basic corollary) holds.  $\square$

*Proof of Theorem 1.* By Taylor expansion of  $n^{-1} \sum_{i=1}^n g(x_i)/(q_{\tilde{\alpha}}(x_i) + \tilde{\zeta}g(x_i)) = 0, \tilde{\zeta}$  can be expanded as

$$\begin{aligned}\tilde{\zeta} &= \frac{\frac{1}{n} \sum_{i=1}^n g(x_i)/q_{\tilde{\alpha}}(x_i)}{\frac{1}{n} \sum_{i=1}^n g(x_i)^2/q_{\tilde{\alpha}}(x_i)^2 + \tilde{\zeta} \cdot \frac{1}{n} \sum_{i=1}^n g(x_i)^3/(q_{\tilde{\alpha}}(x_i) + \tilde{\zeta}g(x_i))^2} \\ &\equiv S_g/S_{gg}, \text{ where } \dot{\zeta} \text{ is between 0 and } \tilde{\zeta}.\end{aligned}$$

By Taylor expansion, we have

$$\begin{aligned}\hat{Z}_{\text{MLE}}(\tilde{\alpha}) &= \frac{1}{n} \sum_{i=1}^n \frac{\pi(x_i)}{(\tilde{\alpha}_1 + \tilde{\zeta})q_1(x_i) + (\tilde{\alpha}_2 - \tilde{\zeta})q_2(x_i)} \\ &= S_{\pi} - (S_{gg}^{-1}S_{\pi g} - \beta^*)S_g,\end{aligned}\quad (\text{A.3})$$

where

$$S_{\pi} = \frac{1}{n} \sum_{i=1}^n \frac{\pi(x_i) - \beta^*g(x_i)}{q_{\tilde{\alpha}}(x_i)}, \quad S_{\pi g} = \frac{1}{n} \sum_{i=1}^n \frac{\pi(x_i)g(x_i)}{(q_{\tilde{\alpha}}(x_i) + \dot{\zeta}g(x_i))^2},$$

and  $\dot{\zeta}$  is between 0 and  $\tilde{\zeta}$ .

Since  $\tilde{\alpha} \xrightarrow{P} \alpha^*$ ,  $\tilde{\zeta} \xrightarrow{P} 0$ , and  $q_{\tilde{\alpha}}(x_i) + \dot{\zeta}g(x_i) = q_{\tilde{\alpha}+(\dot{\zeta}, -\dot{\zeta})}(x_i)$ , we have

$$S_{\pi} \xrightarrow{P} Z, \quad \sqrt{n}(S_{\pi} - Z) \xrightarrow{L} N\left(0, \text{var}_{\alpha^*}\left[\frac{\pi(X) - \beta^*g(X)}{q_{\alpha^*}(X)}\right]\right),$$

$$\sqrt{n}S_g \xrightarrow{L} N\left(0, \text{var}_{\alpha^*}\left[\frac{g(X)}{q_{\alpha^*}(X)}\right]\right),$$

$$S_{\pi g} \xrightarrow{P} \text{cov}_{\alpha^*}\left[\frac{\pi(X)}{q_{\alpha^*}(X)}, \frac{g(X)}{q_{\alpha^*}(X)}\right]$$

and

$$S_{gg} \xrightarrow{P} \text{var}_{\alpha^*}\left[\frac{g(X)}{q_{\alpha^*}(X)}\right].$$

by Lemma 3 and Corollary 1. Then plugging the above results in (A.3), Slutsky's theorem gives that

$$\begin{aligned}\hat{Z}_{\text{MLE}}(\tilde{\alpha}) &\xrightarrow{P} Z \text{ and } \sqrt{n}(\hat{Z}_{\text{MLE}}(\tilde{\alpha}) - Z) \\ &\xrightarrow{L} N\left(0, \text{var}_{\alpha^*}\left[\frac{\pi(X) - \beta^*g(X)}{q_{\alpha^*}(X)}\right]\right).\end{aligned}$$

Similarly, the consistency and asymptotic normality of  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$  hold by the decomposition

$$\begin{aligned}\hat{Z}_{\text{Reg}}(\tilde{\alpha}) &= \frac{1}{n} \sum_{i=1}^n \frac{\pi(x_i) - \hat{\beta}_{\tilde{\alpha}}g(x_i)}{q_{\tilde{\alpha}}(x_i)} \\ &= S_{\pi} - \left[\widehat{\text{var}}\left(\frac{g(X)}{q_{\tilde{\alpha}}(X)}\right)^{-1} \widehat{\text{cov}}\left(\frac{\pi(X)}{q_{\tilde{\alpha}}(X)}, \frac{g(X)}{q_{\tilde{\alpha}}(X)}\right) - \beta^*\right] \cdot S_g.\end{aligned}$$

□

*Proof of Proposition 1.* Denote  $G_n\tau(x) = \sqrt{n}[n^{-1} \sum_{i=1}^n \tau(x_i) - \int \tau(x)q_{\tilde{\alpha}}(x)dx]$ . By (A.3) and Taylor expansion around  $\alpha_1^*$ , we have

$$\begin{aligned}\hat{Z}_{\text{MLE}}(\tilde{\alpha}) - Z &= \frac{1}{\sqrt{n}} \left\{ G_n \frac{\pi(x) - \beta^*g(x)}{q_{\tilde{\alpha}}(x)} \right\} + \frac{1}{\sqrt{n}} \left\{ G_n \frac{g(x)}{q_{\tilde{\alpha}}(x)} \right\} (S_{gg}^{-1}S_{\pi g} - \beta^*) \\ &= \frac{1}{\sqrt{n}} \left\{ G_n \frac{\pi(x) - \beta^*g(x)}{q_{\alpha^*}(x)} \right\} + \frac{1}{\sqrt{n}} \left\{ G_n \frac{(\pi(x) - \beta^*g(x))g(x)}{q_{\alpha^*}(x)^2} \right\} \\ &\quad \times (\tilde{\alpha}_1 - \alpha_1^*) \\ &\quad + \frac{1}{\sqrt{n}} \left\{ G_n \frac{(\pi(x) - \beta^*g(x))g(x)^2}{q_{\alpha^*}(x)^3} \right\} (\tilde{\alpha}_1 - \alpha_1^*)^2 \\ &\quad + o\left(\frac{1}{\sqrt{n}} \left(\frac{n_0}{n} + \frac{1}{\sqrt{n_0}}\right)^2\right) \\ &\quad + \frac{1}{\sqrt{n}} \left\{ G_n \frac{g(x)}{q_{\alpha^*}(x)} \right\} (\tilde{\beta} - \beta^*) + o\left(\frac{1}{n}\right),\end{aligned}$$

where

$$\tilde{\beta} = \frac{\frac{1}{n} \sum_{i=1}^n \pi(x_i)g(x_i)/q_{\tilde{\alpha}}(x_i)^2}{\int g(x)^2/q_{\alpha^*}(x)dx} \left\{ 2 - \frac{\frac{1}{n} \sum_{i=1}^n g(x_i)^2/q_{\tilde{\alpha}}(x_i)^2}{\int g(x)^2/q_{\alpha^*}(x)dx} \right\}.$$

Note that  $S_{gg}^{-1}S_{\pi g} - \beta^* = \tilde{\beta} - \beta^* + o(1/\sqrt{n})$  by Taylor expansion. The expansion of  $\tilde{\alpha}_1$  in Lemma 2 can be plugged into the above equation. After some algebra and note that  $1/\sqrt{n} \leq n_0/n + 1/n_0$  by the inequality  $2ab \leq a^2 + b^2$ , we obtain the expansion of  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$  as follows:

$$\begin{aligned}Z &+ \frac{1}{\sqrt{n}} G_n \frac{\pi(x) - \beta^*g(x)}{q_{\alpha^*}(x)} \\ &+ \frac{n_0}{n\sqrt{n}} \left\{ G_n \frac{(\pi(x) - \beta^*g(x))g(x)}{q_{\alpha^*}(x)^2} \right\} (\gamma_1 - \alpha_1^*) \\ &+ \frac{1}{\sqrt{n_0}\sqrt{n}} \cdot \left\{ G_n \frac{(\pi(x) - \beta^*g(x))g(x)}{q_{\alpha^*}(x)^2} \right\} A_{n_0} \\ &+ \frac{1}{n_0\sqrt{n}} \left\{ G_n \frac{(\pi(x) - \beta^*g(x))g(x)^2}{q_{\alpha^*}(x)^3} \right\} \cdot A_{n_0}^2 \\ &- G_n \frac{(\pi(x) - \beta^*g(x))g(x)}{q_{\alpha^*}(x)^2} \cdot B_{n_0} \left\{ \right. \\ &\quad \left. + \frac{1}{n} \left\{ G_n \frac{g(x)}{q_{\alpha^*}(x)} \right\} \left\{ \sqrt{n}(\tilde{\beta} - \beta^*) \right\} + o\left(\frac{n_0}{n\sqrt{n}}\right) + o\left(\frac{1}{n_0\sqrt{n}}\right) \right\} \\ &\equiv Z + g_1(\tilde{\alpha}) + g_2(\tilde{\alpha}) + o\left(\frac{n_0}{n\sqrt{n}}\right) + o\left(\frac{1}{n_0\sqrt{n}}\right),\end{aligned}\quad (\text{A.4})$$

where

$$\begin{aligned}A_{n_0} &= (1, 0) \cdot V^{-1}\hat{U}/2, \\ B_{n_0} &= (1, 0) \cdot V^{-1}((\hat{V} - \sqrt{n_0}V)V^{-1}\hat{U}/2 + \hat{W}), \\ g_1(\tilde{\alpha}) &= \frac{1}{\sqrt{n}} G_n \frac{\pi(x) - \beta^*g(x)}{q_{\alpha^*}(x)}\end{aligned}$$

and

$$\begin{aligned}g_2(\tilde{\alpha}) &= \frac{n_0}{n\sqrt{n}} \left\{ G_n \frac{(\pi(x) - \beta^*g(x))g(x)}{q_{\alpha^*}(x)^2} \right\} (\gamma_1 - \alpha_1^*) \\ &+ \frac{1}{\sqrt{n_0}\sqrt{n}} \cdot \left\{ G_n \frac{(\pi(x) - \beta^*g(x))g(x)}{q_{\alpha^*}(x)^2} \right\} A_{n_0} \\ &+ \frac{1}{n_0\sqrt{n}} \left\{ G_n \frac{(\pi(x) - \beta^*g(x))g(x)^2}{q_{\alpha^*}(x)^3} \right\} \cdot A_{n_0}^2 \\ &- G_n \frac{(\pi(x) - \beta^*g(x))g(x)}{q_{\alpha^*}(x)^2} \cdot B_{n_0} \left\{ \right. \\ &\quad \left. + \frac{1}{n} \left\{ G_n \frac{g(x)}{q_{\alpha^*}(x)} \right\} \left\{ \sqrt{n}(\tilde{\beta} - \beta^*) \right\} \right\}.\end{aligned}$$

The expansion of  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$  follows similarly, except the definition of  $\tilde{\beta}$  is changed to

$$\tilde{\beta} = \frac{\widehat{\text{cov}}[\pi(X)/q_{\tilde{\alpha}}(X), g(X)/q_{\tilde{\alpha}}(X)]}{\int g(x)^2/q_{\alpha^*}(x)dx} \left\{ 2 - \frac{\widehat{\text{var}}[g(X)/q_{\tilde{\alpha}}(X)]}{\int g(x)^2/q_{\alpha^*}(x)dx} \right\}.$$

□

*Proof of Theorem 2.* The calculation of moments of  $\hat{Z}^*$  involves calculating the moments of (A.4) including

$$\begin{aligned}E &\left[ \left( \frac{1}{n} \sum_{i=1}^n \frac{h_1(x_i)}{q_{\alpha^*}(x_i)} - \int \frac{h_1(x)}{q_{\alpha^*}(x)} q_{\tilde{\alpha}}(x)dx \right)^{k_1} \right. \\ &\quad \left. \left( \frac{1}{n_0} \sum_{i=1}^{n_0} \frac{h_2(x_i)}{q_{\alpha^*}(x_i)} - \int \frac{h_2(x)}{q_{\alpha^*}(x)} q_Y(x)dx \right)^{k_2} (\tilde{\beta} - \beta^*)^{k_3} \right]\end{aligned}$$

for functions  $h_1(x)$  and  $h_2(x)$ ,  $k_1 = 1, 2$ ,  $k_2 = 0, 1, 2$ , and  $k_3 = 0, 1, 2$ .

From (A.4), note that the calculation of  $E[\hat{Z}^* - Z]$  involves calculating the cases of  $k_1 = 1$ . For  $(k_1, k_2, k_3) = (1, 1, 0)$  and  $(1, 2, 0)$ , by

plugging in the decomposition

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \frac{h_1(x_i)}{q_{\alpha^*}(x_i)} - \int \frac{h_1(x)}{q_{\alpha^*}(x)} q_{\tilde{\alpha}}(x) dx \\ &= \frac{n_0}{n} \left( \frac{1}{n_0} \sum_{i=1}^{n_0} \frac{h_1(x_i)}{q_{\alpha^*}(x_i)} - \int \frac{h_1(x)}{q_{\alpha^*}(x)} q_{\gamma}(x) dx \right) \\ &+ \frac{n-n_0}{n} \left( \frac{1}{n-n_0} \sum_{i=1}^{n-n_0} \frac{h_1(x_i)}{q_{\alpha^*}(x_i)} - \int \frac{h_1(x)}{q_{\alpha^*}(x)} q_{\tilde{\alpha}}(x) dx \right), \quad (\text{A.5}) \end{aligned}$$

the expectations are of order  $O(1/n)$  and  $O(1/(n_0n))$ , respectively, by the law of iterated expectations conditioning on  $\{x_i\}_{i=1}^{n_0}$ . For  $(k_1, k_2, k_3) = (1, 0, 1)$ , the expectation is of order  $O(1/n)$  by the inequality  $2ab \leq a^2 + b^2$  and the fact  $E[\sqrt{n}(\tilde{\beta} - \beta^*)]^2 < \infty$ . For  $(k_1, k_2, k_3) = (1, 0, 0)$ , the expectation is 0. Therefore,  $E[\tilde{Z}^* - Z] = O(1/n)$ .

From (A.4), note that the calculation of  $\text{var}[\tilde{Z}^* - Z]$  involves calculating the cases of  $k_1 = 2$ . For  $(k_1, k_2, k_3) = (2, 0, 0)$ ,  $(2, 1, 0)$ , and  $(2, 2, 0)$ , the expectations are of order  $O(1/n)$ ,  $O(n_0/n^2)$ , and  $O(1/(n_0n))$ , respectively, by (A.5) and the law of total variance conditioning on  $\{x_i\}_{i=1}^{n_0}$ . For  $(k_1, k_2, k_3) = (2, 0, 1)$  and  $(2, 0, 2)$ , the expectations are of order  $O(1/(n\sqrt{n}))$  and  $O(1/n^2)$ , respectively, by the inequality  $2ab \leq a^2 + b^2$  and the fact  $E[\sqrt{n}(\tilde{\beta} - \beta^*)]^4 < \infty$ . The other terms are dominated by  $O(n_0/n^2) + O(1/(n_0n))$ . Therefore by noting that  $1/\sqrt{n} \leq n_0/n + 1/n_0$ ,

$$\text{var}[\tilde{Z}^* - Z] = \frac{1}{n} \text{var} \left[ G_n \frac{\pi(x) - \beta^* g(x)}{q_{\alpha^*}(x)} \right] + O\left(\frac{1}{nn_0}\right) + O\left(\frac{n_0}{n^2}\right).$$

Again by (A.5) and the law of total variance conditioning on  $\{x_i\}_{i=1}^{n_0}$ , some algebra gives that

$$\begin{aligned} & \frac{1}{n} \text{var} \left[ G_n \frac{\pi(x) - \beta^* g(x)}{q_{\alpha^*}(x)} \right] \\ &= \frac{1}{n} \sigma_Z^2(\alpha^*) + \frac{n_0}{n^2} \left\{ \text{var}_{\gamma} \left( \frac{\pi(x) - \beta^* g(x)}{q_{\alpha^*}(x)} \right) - \sigma_Z^2(\alpha^*) \right\} \\ &+ \frac{1}{n} \left( 1 - \frac{n_0}{n} \right) \left\{ \text{var}_1 \left( \frac{\pi(x) - \beta^* g(x)}{q_{\alpha^*}(x)} \right) - \text{var}_2 \left( \frac{\pi(x) - \beta^* g(x)}{q_{\alpha^*}(x)} \right) \right\} \\ &\times E(\tilde{\alpha}_1 - \alpha_1^*) \quad (\text{A.6}) \\ &\approx \frac{1}{n} \sigma_Z^2(\alpha^*) + O\left(\frac{n_0}{n^2}\right) + O\left(\frac{1}{nn_0}\right). \end{aligned}$$

Therefore,  $\text{var}[\tilde{Z}^* - Z] = \frac{1}{n} \sigma_Z^2(\alpha^*) + O(n_0/n^2) + O(1/(n_0n))$ .  $\square$

The proofs of Proposition 1 and Theorem 2 reveal the sources of the higher orders  $O(n_0/n^2)$  and  $O(1/(n_0n))$  in  $\text{MSE}[\tilde{Z}^*] - n^{-1} \sigma_Z^2(\alpha^*)$ . These two orders come from three sources, which can be seen by investigating each term in (A.4) and (A.6). One source is due to using pilot samples, which leads to terms

$$\frac{n_0}{n\sqrt{n}} \left\{ G_n \frac{(\pi(x) - \beta^* g(x))g(x)}{q_{\alpha^*}(x)^2} \right\} (\gamma_1 - \alpha_1^*) \text{ in (A.4)}$$

and

$$\frac{n_0}{n^2} \left\{ \text{var}_{\gamma} \left( \frac{\pi(x) - \beta^* g(x)}{q_{\alpha^*}(x)} \right) - \sigma_Z^2(\alpha^*) \right\} \text{ in (A.6),}$$

and results in the order  $O(n_0/n^2)$ . When  $\gamma = \alpha^*$ , these two terms are equal to 0 and thus they are derived from the difference between  $\gamma$  and  $\alpha^*$ . Another one is the variability of random coefficient of control variates, which leads to the term

$$\frac{1}{n} \left\{ G_n \frac{g(x)}{q_{\alpha^*}(x)} \right\} \{ \sqrt{n}(\tilde{\beta} - \beta^*) \} \text{ in (A.4).}$$

This variability results in the order  $O(1/(n\sqrt{n}))$ , which is also  $O(n_0/n^2) + O(1/(n_0n))$  when  $n_0 = \sqrt{n}$ , because  $2/\sqrt{n} \leq n_0/n +$

$1/n_0$ . The other source is the variability of estimated mixture proportion  $\tilde{\alpha}$ , which leads to all other terms in (A.4) and (A.6) except the previous three terms and  $n^{-1} \sigma_Z^2(\alpha^*)$ . This variability results in the order  $O(1/(n_0n))$ .

For the asymptotic properties of  $\hat{\mu}_{\text{Reg}}(\tilde{\alpha})$  and  $\hat{\mu}_{\text{MLE}}(\tilde{\alpha})$ , the proof differs in two aspects with that of  $\hat{Z}_{\text{Reg}}(\tilde{\alpha})$  and  $\hat{Z}_{\text{MLE}}(\tilde{\alpha})$ . One is the  $M$ -estimator  $\hat{\alpha}$  contains an estimated parameter  $\hat{\mu}$  in the criterion function. The other one is the ratio form of  $\hat{\mu}_{\text{Reg}}(\tilde{\alpha})$  and  $\hat{\mu}_{\text{MLE}}(\tilde{\alpha})$ .

*Lemma 5.*  $\hat{\alpha}_1 \xrightarrow{P} \alpha_1^*$  as  $n \rightarrow \infty$  for  $\hat{\alpha}_1$  defined in (14).

*Proof.*  $\hat{\alpha}_1$  can be equivalently obtained as a component of the bivariate estimator  $(\hat{\alpha}_1, \hat{\beta}) = \arg\min_{\alpha_1, \beta \in \Theta} n_0^{-1} \sum_{i=1}^{n_0} \rho(x; \alpha_1, \beta, \hat{\mu})$ , where  $\rho(x; \alpha_1, \beta, \mu) = \frac{[h(x)\pi(x) - \mu\pi(x) - \beta g(x)]^2}{q_{\alpha}(x)q_{\gamma}(x)}$  and  $\Theta = [\delta, 1 - \delta] \times \mathbb{R}$ . The proof of consistency of  $(\hat{\alpha}_1, \hat{\beta})$  contains two steps.

First, although the domain of  $\beta$  is unbounded,  $\hat{\beta}$  stays in a compact set almost surely when  $n \rightarrow \infty$ , because

$$|\hat{\beta}| = \left| \frac{\frac{1}{n_0} \sum_{i=1}^{n_0} \frac{(h(x_i)\pi(x_i) - \hat{\mu}\pi(x_i))g(x_i)}{q_{\hat{\alpha}}(x_i)q_{\gamma}(x_i)}}{\frac{1}{n_0} \sum_{i=1}^{n_0} \frac{g(x_i)^2}{q_{\hat{\alpha}}(x_i)q_{\gamma}(x_i)}} \right| \leq \frac{\frac{1}{n_0} \sum_{i=1}^{n_0} \left( \frac{|h(x_i)\pi(x_i)|}{q_{\gamma}(x_i)} + \frac{|\hat{\mu}\pi(x_i)|}{q_{\gamma}(x_i)} \right) \frac{2}{\delta}}{\frac{1}{n_0} \sum_{i=1}^{n_0} \frac{g(x_i)^2}{(q_1(x_i) + q_2(x_i))q_{\gamma}(x_i)}}$$

and the RHS converges to a constant almost surely since  $\hat{\mu} \rightarrow \mu$  almost surely. Then the consistency of  $(\hat{\alpha}_1, \hat{\beta})$  and the minimizer  $(\hat{\alpha}'_1, \hat{\beta}')$  restricted in some compact set  $C \subset \Theta$ , that is,  $\arg\min_{\alpha_1, \beta \in C} n_0^{-1} \sum_{i=1}^{n_0} \rho(x; \alpha_1, \beta, \hat{\mu})$ , are equivalent since  $P((\hat{\alpha}_1, \hat{\beta}) \in C) \rightarrow 1$ .

Second, the consistency of  $(\hat{\alpha}'_1, \hat{\beta}')$  and the minimizer with  $\hat{\mu}$  replaced by  $\mu$ , that is,  $\arg\min_{\alpha_1, \beta \in C} n_0^{-1} \sum_{i=1}^{n_0} \rho(x; \alpha_1, \beta, \mu)$ , are equivalent because

$$\begin{aligned} & \sup_{\alpha_1, \beta \in C} \left| \frac{1}{n_0} \sum_{i=1}^{n_0} \rho(x; \alpha_1, \beta, \hat{\mu}) - \frac{1}{n_0} \sum_{i=1}^{n_0} \rho(x; \alpha_1, \beta, \mu) \right| \\ &\leq (\hat{\mu}^2 - \mu^2) \max_{\alpha_1, \beta \in C} \frac{1}{n_0} \sum_{i=1}^{n_0} \frac{\pi(x_i)^2}{q_{\alpha_1}(x_i)q_{\gamma}(x_i)} \\ &+ (\hat{\mu} - \mu) \max_{\alpha_1, \beta \in C} \frac{2}{n_0} \sum_{i=1}^{n_0} \frac{(h(x_i)\pi(x_i) - \beta g(x_i))\pi(x_i)}{q_{\alpha_1}(x_i)q_{\gamma}(x_i)} \\ &\rightarrow 0 \text{ almost surely,} \end{aligned}$$

and the argument similar to van der Vaart (2000, Theorem 5.7). Then since the consistency of  $\arg\min_{\alpha_1, \beta \in C} n_0^{-1} \sum_{i=1}^{n_0} \rho(x; \alpha_1, \beta, \mu)$  holds by replacing  $\pi(x)$  in Lemma 1 by  $h(x)\pi(x) - \mu\pi(x)$ , the consistency of  $(\hat{\alpha}_1, \hat{\beta})$  follows.  $\square$

*Proof of Theorem 3.* The consistency and asymptotic normality of  $\hat{\mu}_{\text{MLE}}(\tilde{\alpha})$  follow the extension of proof of Theorem 1 to random vector

$$\sqrt{n} \left\{ \frac{1}{n} \sum_{i=1}^n \left( \frac{h(x_i)\pi(x_i)/q_{\tilde{\alpha}+\tilde{\gamma}}}{\pi(x_i)/q_{\tilde{\alpha}+\tilde{\gamma}}} \right) - \left( \frac{\int h(x)\pi(x)/q_{\alpha^*}(x)dx}{\int \pi(x)/q_{\alpha^*}(x)dx} \right) \right\}$$

and the delta method. The proof for  $\hat{\mu}_{\text{Reg}}(\tilde{\alpha})$  is similar.  $\square$

## APPENDIX B. VARIANCE MATRICES FOR $\hat{\alpha}$

Denote

$$I_{jkl} = \int \frac{(\pi(x) - \beta^{*T} g(x))^j g(x) g(x)^T}{q_{\alpha^*}(x)^k q_{\gamma}(x)^l} dx,$$

$$A = I_{230} - I_{120} I_{010}^{-1} I_{120}, \quad B = I_{010} - I_{120} I_{230}^{-1} I_{120},$$

$$C = I_{441} - 2I_{331} I_{010}^{-1} I_{120}, \quad \text{and } D = I_{331} - 2I_{221} I_{010}^{-1} I_{120}.$$

When estimating  $Z$ ,  $(\hat{\alpha}_2, \dots, \hat{\alpha}_p)$  has the asymptotic variance matrix

$$\frac{1}{n_0} (A^{-1} C A^{-1} - 2I_{230}^{-1} I_{120} B^{-1} D A^{-1}).$$



When estimating  $\mu$ , similar expression can be obtained by replacing  $\pi(x)$  in  $I_{jkl}$  by  $(h(x) - \mu)\pi(x)$ .

[Received March 2012. Revised June 2013]

## REFERENCES

- Battiti, R., and Masulli, F. (1990), "BFGS Optimization for Faster and Automated Supervised Learning," in *Proceedings of the International Neural Network Conference (INNC 90)-Paris-France*, pp. 757–760. [1356]
- Bauwens, L., and Lubrano, M. (2008), "Bayesian Inference on GARCH Models Using the Gibbs Sampler," *The Econometrics Journal*, 1, 23–46. [1358]
- Bollerslev, T. (1986), "Generalized Autoregressive Conditional Heteroskedasticity," *Journal of Econometrics*, 31, 307–327. [1358]
- Chow, Y., and Teicher, H. (2003), *Probability Theory: Independence, Interchangeability, Martingales*, New York: Springer Verlag. [1362]
- Cochran, W. (1977), *Sampling Techniques*, New York: Wiley. [1352]
- Denny, M. (2001), "Introduction to Importance Sampling in Rare-Event Simulations," *European Journal of Physics*, 22, 403. [1350]
- Duffie, D., and Pan, J. (1997), "An Overview of Value at Risk," *The Journal of Derivatives*, 4, 7–49. [1358]
- Dunkel, J., and Weber, S. (2007), "Efficient Monte Carlo Methods for Convex Risk Measures in Portfolio Credit Risk Models," in *Proceedings of the 2007 Winter Simulation Conference*, pp. 958–966. [1359]
- Durrett, R. (1996), *Probability: Theory and Examples*, Pacific Grove, CA: Duxbury Press. [1361]
- Emond, M., Raftery, A., and Steele, R. (2001), "Easy Computation of Bayes Factors and Normalizing Constants for Mixture Models via Mixture Importance Sampling," Technical Report No. 398, Department of Statistics, Washington University Seattle. [1351]
- Engle, R. (1982), "Autoregressive Conditional Heteroscedasticity With Estimates of the Variance of United Kingdom Inflation," *Econometrica: Journal of the Econometric Society*, 50, 987–1007. [1358]
- Fan, S., Chenney, S., Hu, B., Tsui, K., and Lai, Y. (2006), "Optimizing Control Variate Estimators for Rendering," in *Computer Graphics Forum* (Vol. 25), eds. E. Gröller and L. Szirmay-Kalos, Goslar, Germany: Eurographics Association, pp. 351–357. [1351]
- Ford, E., and Gregory, P. (2007), "Bayesian Model Selection and Extrasolar Planet Detection," in *Statistical Challenges in Modern Astronomy IV* (Vol. 371), eds. G. J. Babu and E. D. Feigelson, San Francisco, CA: Astron. Soc. Pacific, p. 189. [1350]
- Gelman, A., and Meng, X. (1998), "Simulating Normalizing Constants: From Importance Sampling to Bridge Sampling to Path Sampling," *Statistical Science*, 13, 163–185. [1350]
- Geweke, J. (1989), "Bayesian Inference in Econometric Models Using Monte Carlo Integration," *Econometrica: Journal of the Econometric Society*, 57, 1317–1339. [1356]
- (1994), "Bayesian Comparison of Econometric Models," Working Paper No. 532, Federal Reserve Bank of Minneapolis. [1358]
- Geyer, C. (1994), "On the Asymptotics of Constrained M-Estimation," *The Annals of Statistics*, 22, 1993–2010. [1355]
- Giordani, P., and Kohn, R. (2010), "Adaptive Independent Metropolis-Hastings by Fast Estimation of Mixtures of Normals," *Journal of Computational and Graphical Statistics*, 19, 243–259. [1350]
- Givens, G., and Raftery, A. (1996), "Local Adaptive Importance Sampling for Multivariate Densities With Strong Nonlinear Relationships," *Journal of the American Statistical Association*, 91, 132–141. [1350]
- Glasserman, P., Heidelberger, P., and Shahabuddin, P. (2000), "Variance Reduction Techniques for Estimating Value-at-Risk," *Management Science*, 46, 1349–1364. [1359]
- Haberman, S. (1989), "Concavity and Estimation," *The Annals of Statistics*, 17, 1631–1661. [1354,1361]
- Hesterberg, T. (1988), "Advances in Importance Sampling," Ph.D. dissertation, Department of Statistics, Stanford University. [1352]
- (1995), "Weighted Average Importance Sampling and Defensive Mixture Distributions," *Technometrics*, 37, 185–194. [1350,1351,1355,1358]
- Hjort, N., and Pollard, D. (1994), "Asymptotics for Minimisers of Convex Processes," Statistical Research Report, Department of Mathematics, University of Oslo. [1362]
- Hoogerheide, L., and Van Dijk, H. (2010), "Bayesian Forecasting of Value at Risk and Expected Shortfall Using Adaptive Importance Sampling," *International Journal of Forecasting*, 26, 231–247. [1359]
- Jorion, P. (1997), *Value At Risk: The New Benchmark for Controlling Market Risk*, Chicago: McGraw-Hill. [1358]
- Kong, A., McCullagh, P., Meng, X., Nicolae, D., and Tan, Z. (2003), "A Theory of Statistical Models for Monte Carlo Integration," *Journal of the Royal Statistical Society, Series B*, 65, 585–604. [1352]
- Liang, F., Liu, C., and Carroll, R. J. (2007), "Stochastic Approximation in Monte Carlo Computation," *Journal of the American Statistical Association*, 102, 305–320. [1351,1356]
- Liu, J. (2008), *Monte Carlo Strategies in Scientific Computing*, New York: Springer Verlag. [1355]
- Oh, M., and Berger, J. (1993), "Integration of Multimodal Functions by Monte Carlo Importance Sampling," *Journal of the American Statistical Association*, 88, 450–456. [1350,1351,1356]
- Owen, A., and Zhou, Y. (2000), "Safe and Effective Importance Sampling," *Journal of the American Statistical Association*, 95, 135–143. [1351,1352,1353,1355,1356,1361]
- (1999), "Adaptive Importance Sampling by Mixtures of Products of Beta Distributions," Technical Report No. 1999-1, Department of Statistics, Stanford University. [1350]
- Raghavan, N., and Cox, D. (1998), "Adaptive Mixture Importance Sampling," *Journal of Statistical Computation and Simulation*, 60, 237–260. [1351,1352,1357,1361]
- Robert, C., and Casella, G. (2004), *Monte Carlo Statistical Methods*, New York: Springer Verlag. [1351]
- Rothenberg, T. (1984), "Approximating the Distributions of Econometric Estimators and Test Statistics," in *Handbook of Econometrics* (Vol. 2), eds. Z. Griliches and M. D. Intriligator, Amsterdam: North Holland, pp. 881–935. [1354]
- Rubinstein, R., and Kroese, D. (2008), *Simulation and the Monte Carlo Method*, Hoboken, NJ: Wiley. [1351,1355]
- Smith, P., Shafi, M., and Gao, H. (1997), "Quick Simulation: A Review of Importance Sampling Techniques in Communications Systems," *IEEE Journal on Selected Areas in Communications*, 15, 597–613. [1350]
- Tan, Z. (2004), "On a Likelihood Approach for Monte Carlo Integration," *Journal of the American Statistical Association*, 99, 1027–1036. [1351,1352,1353,1355,1356,1361]
- van der Vaart, A. (2000), *Asymptotic Statistics*, Cambridge: Cambridge University Press. [1362,1364]
- van der Vaart, A., and Wellner, J. (1996), *Weak Convergence and Empirical Processes*, New York: Springer Verlag. [1362]
- Veach, E., and Guibas, L. (1995), "Optimally Combining Sampling Techniques for Monte Carlo Rendering," in *Proceedings of the 22nd Annual Conference on Computer Graphics and Interactive Techniques*, pp. 419–428. [1350,1351,1352]
- West, M. (1993), "Approximating Posterior Distributions by Mixture," *Journal of the Royal Statistical Society, Series B*, 55, 409–422. [1350,1356]