

---

# The role of tail dependence in estimating posterior expectations

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Many tasks in modern probabilistic machine learning and statistics require estimating  
2 expectations over posterior distributions. While many algorithms have  
3 been developed to approximate these expectations, reliably assessing their performance  
4 in practice, in absence of ground truth, remains a significant challenge.  
5 In this work, we observe that the well-known  $k$ -hat diagnostic for importance sampling  
6 (IS) [1] can be unreliable, as it fails to account for the fact that the common  
7 self-normalized IS (SNIS) estimator is a ratio. First, we demonstrate that examining  
8 separate  $k$ -hat statistics for the numerator and denominator can be insufficient.  
9 Then, we propose a new statistic that accounts for the dependence between the  
10 estimators in the ratio. In particular, we find that the concept of tail dependence  
11 between numerator and denominator weights contains essential information for  
12 determining effective performance of the SNIS estimator.

## 13 1 Introduction and background

14 Algorithms for Bayesian computation continue to be used for increasingly complex probabilistic  
15 models, remaining an active research field [2]. Yet, in the absence of ground truth, it remains  
16 challenging in practice to determine how and in which sense an approximate inference algorithm  
17 has found a “good” solution, as studied by several recent works, for Markov Chain Monte Carlo  
18 (MCMC) [3–5], variational inference (VI) [6–8], and importance sampling [1, 9–11] (the latter two  
19 being closely connected). In this work, we focus on diagnostics that apply to IS and VI algorithms.

20 **Problem statement.** Let  $\theta \in \Theta$  (commonly,  $\mathbb{R}^{d_\theta}$ ) be the parameter of a Bayesian statistical model  
21  $\{p(y|\theta)\}_\theta$  for data  $y \in \mathcal{Y}$  with posterior PDF  $\pi(\theta|\mathcal{D}) \stackrel{\text{def}}{=} Z_\pi^{-1} \cdot \tilde{\pi}(\theta|\mathcal{D}) = Z_\pi^{-1} \cdot \prod_n p(y_n|\theta) \cdot \pi(\theta)$   
22 with  $\mathcal{D} \stackrel{\text{def}}{=} \{y_n\}_{n=1}^N$ ,  $Z_\pi$  the normalizer and prior PDF  $\pi(\theta)$ . Formally, we aim at constructing Monte  
23 Carlo estimates of a posterior expectation  $I \in \mathbb{R}_{>0}$ , defined as

$$I \stackrel{\text{def}}{=} \mathbb{E}_{\pi(\theta|\mathcal{D})}[f(\theta)] = \int f(\theta)\pi(\theta|\mathcal{D})d\theta, \quad (1)$$

24 where  $f : \Theta \rightarrow \mathbb{R}_{\geq 0}$  is a suitably integrable test function. In particular, we are interested in  
25 obtaining diagnostics to determine the quality of an estimator  $\hat{I}$ . As a concrete example, when we  
26 set  $f(\theta) = p(y^{(n+1)}|\theta)$  for a test point  $y^{(n+1)}$ ,  $I$  is often written as  $p(y^{(n+1)}|\mathcal{D})$ , i.e., the evaluation  
27 of the posterior predictive PDF  $p(y|\mathcal{D})$  at point  $y^{(n+1)}$ .<sup>1</sup>

28 **Self-normalized IS, combination with VI.** Approximating integrals like in Eq. (1) accurately is  
29 challenging. MCMC is a natural solution, but there are notable cases where it is not appropriate. For

---

<sup>1</sup>Such integrals can be used for estimating the predictive performance of a posterior [12] or the influence of a particular observation.

30 example, when even exact i.i.d. sampling from  $\pi(\theta|\mathcal{D})$  is inefficient, or when it is too expensive. In  
 31 these cases one usually resorts to IS [13], where we obtain samples from a chosen proposal PDF  $q$ ,  
 32 as  $\theta^{(s)} \stackrel{\text{i.i.d.}}{\sim} q(\theta)$ , and construct estimators for  $I$  as

$$\hat{I}_{\text{SNIS}} = \sum_{s=1}^S \bar{w}^{(s)} f(\theta^{(s)}) \quad , \quad \bar{w}^{(s)} \stackrel{\text{def}}{=} \frac{w^{(s)}}{\sum_{s'=1}^S w^{(s')}} \quad , \quad w^{(s)} = w(\theta^{(s)}) = \frac{\tilde{\pi}(\theta^{(s)}|\mathcal{D})}{q(\theta^{(s)})}. \quad (2)$$

33 Many theoretical properties of this estimator are known (see, e.g., [14] for a review). When the  
 34 normalizing constant  $Z_\pi$  is unknown (i.e., almost always), the normalization of the weights in Eq. (2)  
 35 is not optional. In practice, it is difficult to find a good proposal, i.e., leading to estimates that are  
 36 close to  $I$  in suitable ways. It is natural to use proposals that are the result of a VI algorithm [6],  
 37 which is done implicitly or explicitly in the VI literature. See [6, 15–24] as examples for the many  
 38 connections between VI and IS. A consequence of using a bad proposal is that the distribution of the  
 39 weights  $w_s$  tends to have a few very large values.

40 **Pareto-smoothed IS.** Exploiting this observation, [1] proposed Pareto-smoothed IS (PSIS), which  
 41 replaces the largest  $M$  unnormalized weights<sup>2</sup> to get SNIS estimators with better behaviour. They  
 42 fit a generalized Pareto distribution (GPD) to the weights  $\{w^{(s)}\}_{s=1}^S$ , whose PDF we denote as  
 43  $p(w)$ .<sup>3</sup> The new weights introduce bias but reduce variance.. The GPD has three parameters, the  
 44 most important of which is the shape parameter  $k$ . [1] propose to use an estimate of  $k$ , i.e.,  $\hat{k}$ , as a  
 45 diagnostic for IS.

46 **The  $\hat{k}$  diagnostic.** [1] use the estimated value of  $k$ , i.e.,  $\hat{k}$ , as a diagnostic for deciding whether  
 47 the SNIS estimates with PSIS-corrected weights are reliable. The GPD has  $1/k$  finite fractional  
 48 moments when the true  $k > 0$ , which suggests finite variance as soon as  $k < 0.5$ . Note that  
 49 this guarantees finite variance only for the normalizing constant estimator  $\hat{Z}_\pi = 1/S \sum_{s=1}^S w^{(s)}$ ,  
 50 which is implicit in the denominator of SNIS [25]. [1] find empirically that when using  $S > 2000$ ,  
 51 estimation with PSIS-corrected weights is reliable for  $\hat{k} < 0.7$ , a threshold less stringent than 0.5.  
 52 An advantage of  $\hat{k}$  is that it is not an IS estimate itself, unlike the effective sample size (ESS) [10],  
 53 attempting to address the issues with variance-based diagnostics [9].

## 54 2 Methodology

55 Several works [25–27] have shown theoretically and empirically that accurately estimating posterior  
 56 expectations such as  $I$  in Eq. (1) involves more than simply finding a proposal  $q(\theta)$  that is close to  
 57 the posterior  $\pi(\theta|\mathcal{D})$ . This is because the SNIS estimator is a ratio estimator, as  $I$  itself is the ratio  
 58 of two integrals,

$$I = \frac{\int f(\theta) \tilde{\pi}(\theta|\mathcal{D}) d\theta}{\int \tilde{\pi}(\theta|\mathcal{D}) d\theta} \stackrel{\text{def}}{=} \frac{I_{\text{num}}}{Z_\pi} \stackrel{\text{def}}{=} \frac{I_{\text{num}}}{I_{\text{den}}}, \quad (3)$$

59 where we relabelled the normalizing constant  $I_{\text{den}}$ . Therefore, we can write the SNIS estimator as

$$\hat{I}_{\text{SNIS}} = \frac{\frac{1}{S} \sum_{s=1}^S w^{(s)} f(\theta^{(s)})}{\frac{1}{S} \sum_{s=1}^S w^{(s)}} = \frac{\hat{I}_{\text{num}}}{\hat{I}_{\text{den}}}, \quad \theta^{(s)} \stackrel{\text{i.i.d.}}{\sim} q(\theta), \quad (4)$$

60 where the two estimators  $\hat{I}_{\text{num}}$  and  $\hat{I}_{\text{den}}$  are unbiased, but  $\hat{I}_{\text{SNIS}}$  is not. As elaborated in [25], the  
 61 asymptotic variance of the SNIS estimator is driven by the variance of the numerator estimator, the  
 62 variance of the denominator, and the covariance between them. For convenience, we define two  
 63 unnormalized importance weight functions, the one used in the numerator for  $\hat{I}_{\text{num}}$  and the one used  
 64 in  $\hat{I}_{\text{den}}$ , as

$$w_{\text{num}}(\theta) = \frac{f(\theta) \tilde{\pi}(\theta|\mathcal{D})}{q(\theta)}, \quad w_{\text{den}}(\theta) = \frac{\tilde{\pi}(\theta|\mathcal{D})}{q(\theta)}. \quad (5)$$

65 We can then write the SNIS estimator as a ratio of two unbiased IS estimators,

$$\hat{I}_{\text{SNIS}} = \frac{\frac{1}{S} \sum_{s=1}^S w_{\text{num}}(\theta^{(s)})}{\frac{1}{S} \sum_{s=1}^S w_{\text{den}}(\theta^{(s)})}, \quad \theta^{(s)} \stackrel{\text{i.i.d.}}{\sim} q(\theta). \quad (6)$$

<sup>2</sup>See [1] for the choice of  $M$ .

<sup>3</sup>We use the common abuse of notation of using lowercase letters for both the random variable itself and its realised values.

66 Given that there are two IS weights,  $w_{\text{num}}(\theta^{(s)})$ ,  $w_{\text{den}}(\theta^{(s)})$  in the above, it is natural to consider  
 67 that one may track reliability  $\hat{I}_{\text{SNIS}}$  by computing two diagnostics  $\hat{k}_{\text{num}}$ ,  $\hat{k}_{\text{den}}$  separately for weights  
 68  $\{w_{\text{num}}^{(s)}\}_{s=1}^S$  and  $\{w_{\text{den}}^{(s)}\}_{s=1}^S$ . [1] explored this option empirically, reporting that in their experiments  
 69 it was sufficient to take  $\max(\hat{k}_{\text{num}}, \hat{k}_{\text{den}})$  to determine reliability of the ratio. In this work, we will  
 70 argue that this heuristic misses useful information and propose a new diagnostic.

## 71 2.1 Capturing error cancellation with tail dependence

72 The diagnostics  $\hat{k}_{\text{num}}$  and  $\hat{k}_{\text{den}}$  describe how well  $\hat{I}_{\text{num}}$  and  $\hat{I}_{\text{den}}$  respectively approximate  $I_{\text{num}}$  and  
 73  $I_{\text{den}}$ , serving as an (improved) substitute for estimates of variance (like the ESS). Yet, the variance  
 74 of the SNIS estimator  $\hat{I}_{\text{SNIS}}$  is not only affected by the variance of the numerator of Eq. (6), the  
 75 variance of the denominator. It is also affected by the covariance between them,  $\text{Cov}_q[\hat{I}_{\text{num}}, \hat{I}_{\text{den}}]$   
 76 [25]. Note that  $\text{Cov}_q[\hat{I}_{\text{num}}, \hat{I}_{\text{den}}]$  is proportional to  $\text{Cov}_q[w_{\text{num}}, w_{\text{den}}]$ , the random variables we will  
 77 focus on from now.

78 A straightforward idea to capture this missing piece of information from  $\hat{k}_{\text{num}}$  and  $\hat{k}_{\text{den}}$  is to construct  
 79 an estimate of  $\text{Cov}_q[\hat{I}_{\text{num}}, \hat{I}_{\text{den}}]$ , using the same samples from  $q$  used to estimate  $I$ . Yet, doing so  
 80 would suffer the same drawbacks of variance-based diagnostics, which was a motivation for  $\hat{k}$  [1].  
 81 Thus, we will develop a diagnostic that is not a direct estimate of  $\text{Cov}_q[w_{\text{num}}, w_{\text{den}}]$ . Like [1], we  
 82 also exploit the fact that the distribution of  $w_{\text{num}}$  and  $w_{\text{den}}$  can be well approximated with a power-  
 83 law distribution in the tails. Specifically, we will look at a suitable notion of dependence between  
 84 the tails of  $w_{\text{num}}$  and  $w_{\text{den}}$ . This notion will replace the covariance  $\text{Cov}_q[w_{\text{num}}, w_{\text{den}}]$  as our target  
 85 estimate. In fact, covariance, up to normalization, is equivalent to Pearson's correlation  $\rho$ , which is  
 86 only a very specific form of dependence, with many known limitations [28].

87 **Dependence and error cancellation.** An intuition for why higher covariance between the estima-  
 88 tors  $\text{Cov}_q[w_{\text{num}}, w_{\text{den}}]$ , or other dependence metrics, can lead to lower error is that, in a ratio, error  
 89 cancellation can happen. Error cancellation in ratios has been exploited to derive better convergence  
 90 rates for other numerical integration methods [29]. In IS, it is known that large IS weights lead to  
 91 high errors. Therefore, error cancellation in the ratio of Eq. (6) could happen when a large weight in  
 92 the numerator is offset by another similarly large weight in the denominator. We now formalize this  
 93 using the notion of tail dependence.

94 **Definition 1 (Upper tail dependence coefficient and tail dependence)** *Let  $W_1, W_2$  be two real-*  
 95 *valued random variables. Let their (continuous) marginal CDFs be  $F_1, F_2$ . Then,*

$$\lim_{q \rightarrow 1^-} \mathbb{P}[W_2 > F_2^{-1}(q) | W_1 > F_1^{-1}(q)] = \lambda_U, \quad (7)$$

96 *provided the limit exists, is known as upper tail dependence coefficient  $\lambda_U \in [0, 1]$ . If  $\lambda_U > 0$ ,*  
 97 *we say that  $W_1, W_2$  are asymptotically tail dependent, with the magnitude of  $\lambda_U$  determining the*  
 98 *strength of dependence.*

99 Next, we discuss how to relate the above concept to the estimation of  $I$ .

## 100 2.2 Proposed reliability checks

101 We propose to diagnose whether the estimate in Eq. (6) is reliable by examining three quantities:  
 102  $\hat{k}_{\text{num}}$ ,  $\hat{k}_{\text{den}}$  and a new diagnostic that is constructed as an approximation of the tail dependence  
 103 coefficient  $\lambda_U$  between  $w_{\text{num}}, w_{\text{den}}$ . Our aim is to study how these quantities relate to the effective  
 104 performance of  $\hat{I}_{\text{SNIS}}$  as an estimator of  $I$ , which we define as follows.

105 **Definition 2 (Effective performance)** *We define the effective performance of an estimator  $\hat{I}$  of  $I$*   
 106 *as ensuring that the value of  $(\hat{I}/I)$  is close to 1 with high probability. This takes into account the*  
 107 *possibility of  $I$  being very small, e.g.,  $10^{-7}$  following the recommendation of [9]. In log-space, it is*  
 108 *equivalent to look at how  $\log I - \log \hat{I}$  is close to zero (recall  $I > 0$ ).*

109 **Semi-parametric estimation of tail dependence** In mathematical finance, various estimators  
 110 of tail dependence have been developed [30–32]. We begin by studying semi-parametric estima-  
 111 tors, following the assumption used by [1] and common in heavy-tailed distribution inference [33].

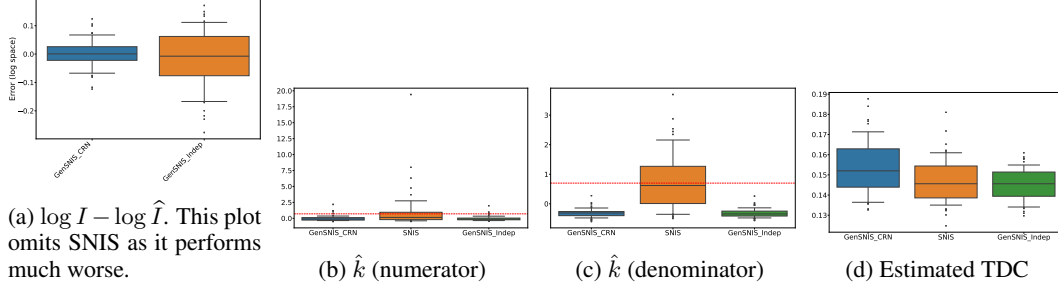


Figure 1: Results ( $d_\theta = 3$ ) over 50 replications. We compare SNIS, GenSNIS (see Section 3) with a common random number (CRN) and GenSNIS with independent marginals. From Fig. 1a, we see that GenSNIS with CRN performs best; this cannot be captured by  $\hat{k}$  values, but by the higher TDC.

112 Specifically, we assume the distribution of  $w_{\text{num}}, w_{\text{den}}$  is well approximated by a GPD in the tails.  
 113 Similarly, to estimate tail dependence, we assume the *copula* of their joint distribution is well ap-  
 114 proximated by an extreme value copula [34], also only in the tails.<sup>4</sup> We hypothesize that tail depen-  
 115 dence between  $w_{\text{num}}$  and  $w_{\text{den}}$  improves  $\hat{I}_{\text{SNIS}}$  performance, similar to the effect of  $\text{Cov}_q[w_{\text{num}}, w_{\text{den}}]$ ,  
 116 but easier to estimate and more reliable. To model this, we fit a Student-t copula, a simple parametric  
 117 choice that also serves as an extreme value copula [31].

### 118 3 Preliminary results on Bayesian linear regression and conclusions

119 We look at the distribution of  $\log I - \log \hat{I}$  over different replications. We consider estimating the  
 120 posterior predictive of a Bayesian linear regression (BLR) model where we can compute the exact  
 121 value of  $I$ . That is, from Eq. (1), we set  $f(\theta) = p(y^{(n+1)}|\theta)$  for a test point  $y^{(n+1)}$ , and  $\pi(\theta|\mathcal{D})$  is a  
 122 Gaussian with known mean and covariance (BLR posterior).<sup>5</sup>

123 To validate our hypothesis that tail dependence contains useful information, we check the behaviour  
 124 of the diagnostics  $\hat{k}_{\text{num}}, \hat{k}_{\text{den}}$  our tail dependence diagnostic  $\hat{\lambda}_U$  estimated from a Student-t copula  
 125  $C(u_1, u_2; \rho, \nu)$ , which is given by  $\hat{\lambda}_U = 2t_{\hat{\nu}+1}(-\sqrt{\hat{\nu}+1}\sqrt{1-\hat{\rho}}/\sqrt{1+\hat{\rho}})$  where  $t$  is the Student-t  
 126 CDF.<sup>6</sup> We find that, when  $k$ -diagnostics between competitors are similar for numerator and de-  
 127 nominator, a higher tail dependence coefficient (TDC) explains the better performance. To explain  
 128 our results, we need to introduce a recent generalization of the SNIS estimator proposed in [25],  
 129 i.e., sampling from an extended space  $\mathbb{R}^{d_\theta} \times \mathbb{R}^{d_\theta}$ , as  $\hat{I}_{\text{GenSNIS}} = \frac{\frac{1}{S} \sum_{s=1}^S w_{\text{num}}(\theta_1^{(s)})}{\frac{1}{S} \sum_{s=1}^S w_{\text{den}}(\theta_2^{(s)})}, [\theta_1^{(s)}, \theta_2^{(s)}] \stackrel{\text{i.i.d.}}{\sim}$   
 130  $q_{1,2}(\theta_1, \theta_2)$ . SNIS is a special case where the joint is a degenerate joint with  $\theta_1 = \theta_2$ . Another  
 131 special case is taking  $q_{1,2}(\theta_1, \theta_2) = q_1(\theta_1)q_2(\theta_2)$ , which is done in previous works including no-  
 132 tably target-aware Bayesian inference [26]. Finally, for these experiments we consider the choice  
 133 of  $q_{1,2}(\theta_1, \theta_2)$  that uses a common random number (CRN) for numerator and denominator, but has  
 134 different marginals. Concretely we used Gaussian proposals  $\mathcal{N}(\theta_1; \mu_1, \Sigma_1)$  and  $\mathcal{N}(\theta_2; \mu_2, \Sigma_2)$  for  
 135 numerator and denominator, respectively. The parameters are set to the optimal ones (given by the  
 136 BLR true posteriors for numerator and denominator) perturbed by an error term. The SNIS esti-  
 137 mator uses only one distribution  $q(\theta)$ , so we take the midpoint between the two optimal IS means  
 138 and covariances for its parameters. Fig. 1 shows the results. We indeed find in other settings (for  
 139  $d_\theta$ , noise variance, and covariate distributions) that when  $\hat{k}$  values are similar for numerator and  
 140 denominator, tail dependence explains the remaining performance if a difference exists. We plan to  
 141 test further TDC metrics and Bayesian models.

## 142 References

- 143 [1] Aki Vehtari, Daniel Simpson, Andrew Gelman, Yuling Yao, and Jonah Gabry. Pareto smoothed  
 144 importance sampling. *Journal of Machine Learning Research*, 25(72):1–58, 2024.

<sup>4</sup>A copula of a bivariate joint distribution is the distribution on  $[0, 1]^2$  after transforming the marginals to the uniform distribution. Many parametric copula families exist [35].

<sup>5</sup>See [36] for expressions about BLR including closed form posterior predictives.

<sup>6</sup>We use the estimate of  $\hat{\rho}$  from the Python statsmodels package, while setting  $\nu$  manually.

- 145 [2] Steven Winter, Trevor Campbell, Lizhen Lin, Sanvesh Srivastava, and David B. Dunson.  
146 Emerging Directions in Bayesian Computation. *Statistical Science*, 39(1):62 – 89, 2024.
- 147 [3] Vivekananda Roy. Convergence diagnostics for markov chain monte carlo. *Annual Review of*  
148 *Statistics and Its Application*, 7(1):387–412, 2020.
- 149 [4] Dootika Vats, James M. Flegal, and Galin L. Jones. *Monte Carlo Simulation: Are We There*  
150 *Yet?*, pages 1–15. John Wiley & Sons, Ltd, 2021.
- 151 [5] Galin L Jones and Qian Qin. Markov chain monte carlo in practice. *Annual Review of Statistics*  
152 *and Its Application*, 9(1):557–578, 2022.
- 153 [6] Yuling Yao, Aki Vehtari, Daniel Simpson, and Andrew Gelman. Yes, but did it work?: Eval-  
154 uating variational inference. In *International Conference on Machine Learning*, pages 5581–  
155 5590. PMLR, 2018.
- 156 [7] Yu Wang, Mikolaj Kasprzak, and Jonathan H Huggins. A targeted accuracy diagnostic for vari-  
157 ational approximations. In *International Conference on Artificial Intelligence and Statistics*,  
158 pages 8351–8372. PMLR, 2023.
- 159 [8] Manushi Welandawe, Michael Riis Andersen, Aki Vehtari, and Jonathan H. Huggins. A frame-  
160 work for improving the reliability of black-box variational inference. *Journal of Machine*  
161 *Learning Research*, 25(219):1–71, 2024.
- 162 [9] Sourav Chatterjee and Persi Diaconis. The sample size required in importance sampling. *The*  
163 *Annals of Applied Probability*, 28(2):1099–1135, 2018.
- 164 [10] Víctor Elvira, Luca Martino, and Christian P Robert. Rethinking the effective sample size.  
165 *International Statistical Review*, 90(3):525–550, 2022.
- 166 [11] Medha Agarwal, Dootika Vats, and Víctor Elvira. A principled stopping rule for importance  
167 sampling. *Electronic Journal of Statistics*, 16(2):5570–5590, 2022.
- 168 [12] Aki Vehtari, Andrew Gelman, and Jonah Gabry. Practical bayesian model evaluation using  
169 leave-one-out cross-validation and waic. *Statistics and computing*, 27:1413–1432, 2017.
- 170 [13] Art B. Owen. *Monte Carlo theory, methods and examples*. [https://artowen.su.domains/  
171 mc/](https://artowen.su.domains/mc/), 2013.
- 172 [14] Nicolas Chopin, Omiros Papaspiliopoulos, et al. *An introduction to sequential Monte Carlo*,  
173 volume 4. Springer, 2020.
- 174 [15] Andriy Mnih and Danilo Rezende. Variational inference for monte carlo objectives. In *Inter-  
175 national Conference on Machine Learning*, pages 2188–2196. PMLR, 2016.
- 176 [16] Joseph Sakaya and Arto Klami. Importance sampled stochastic optimization for variational  
177 inference. In *33rd Conference on Uncertainty in Artificial Intelligence*. AUAI Press, 2017.
- 178 [17] Justin Domke and Daniel R Sheldon. Importance weighting and variational inference. In  
179 *Advances in Neural Information Processing Systems (NeurIPS)*, pages 4470–4479, 2018.
- 180 [18] Axel Finke and Alexandre H Thiery. On importance-weighted autoencoders.  
181 <https://arxiv.org/abs/1509.00519>, 2019.
- 182 [19] Akash Kumar Dhaka, Alejandro Catalina, Manushi Welandawe, Michael R Andersen,  
183 Jonathan Huggins, and Aki Vehtari. Challenges and opportunities in high dimensional vari-  
184 ational inference. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages  
185 7787–7798, 2021.
- 186 [20] Lu Zhang, Bob Carpenter, Andrew Gelman, and Aki Vehtari. Pathfinder: Parallel quasi-  
187 Newton variational inference. *Journal of Machine Learning Research*, 23(1):13802–13850,  
188 2022.

- 189 [21] Pierre-Alexandre Mattei and Jes Frelsen. Uphill roads to variational tightness: Monotonicity  
190 and Monte Carlo objectives. <https://arxiv.org/abs/2201.10989>, 2022.
- 191 [22] Oskar Kviman, Harald Melin, Hazal Koptagel, Victor Elvira, and Jens Lagergren. Multiple im-  
192 portance sampling ELBO and deep ensembles of variational approximations. In *International*  
193 *Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 10687–10702, 2022.
- 194 [23] Arnaud Doucet, Eric Moulines, and Achille Thin. Differentiable samplers for deep latent  
195 variable models. *Philosophical Transactions of the Royal Society A*, 381(2247):20220147,  
196 2023.
- 197 [24] Thomas Guilmeau, Nicola Branchini, Emilie Chouzenoux, and Víctor Elvira. Adaptive im-  
198 portance sampling for heavy-tailed distributions via alpha-divergence minimization. In *Inter-*  
199 *national Conference on Artificial Intelligence and Statistics*, pages 3871–3879. PMLR, 2024.
- 200 [25] Nicola Branchini and Víctor Elvira. Generalizing self-normalized importance sampling with  
201 couplings. *arXiv preprint arXiv:2406.19974*, 2024.
- 202 [26] Tom Rainforth et al. Target-aware bayesian inference: how to beat optimal conventional  
203 estimators. *Journal of Machine Learning Research*, 2020.
- 204 [27] Topi Paananen, Juho Piironen, Paul-Christian Bürkner, and Aki Vehtari. Implicitly adaptive  
205 importance sampling. *Statistics and Computing*, 31(2):16, 2021.
- 206 [28] Dag Tjøstheim, Håkon Otneim, and Bård Støve. Statistical dependence: Beyond pearson’s  $\rho$ .  
207 *Statistical science*, 37(1):90–109, 2022.
- 208 [29] John F Monahan. *Numerical methods of statistics*. Cambridge University Press, 2011.
- 209 [30] Gabriel Frahm, Markus Junker, and Rafael Schmidt. Estimating the tail-dependence coeffi-  
210 cient: properties and pitfalls. *Insurance: mathematics and Economics*, 37(1):80–100, 2005.
- 211 [31] Stefano Demarta and Alexander J McNeil. The t copula and related copulas. *International*  
212 *statistical review*, 73(1):111–129, 2005.
- 213 [32] Rafael Schmidt and Ulrich Stadtmüller. Non-parametric estimation of tail dependence. *Scan-*  
214 *dinavian journal of statistics*, 33(2):307–335, 2006.
- 215 [33] Jayakrishnan Nair, Adam Wierman, and Bert Zwart. *The fundamentals of heavy tails: Proper-*  
216 *ties, emergence, and estimation*, volume 53. Cambridge University Press, 2022.
- 217 [34] Gordon Gudendorf and Johan Segers. Extreme-value copulas. In *Copula Theory and Its*  
218 *Applications: Proceedings of the Workshop Held in Warsaw, 25-26 September 2009*, pages  
219 127–145. Springer, 2010.
- 220 [35] Claudia Czado. Analyzing dependent data with vine copulas. *Lecture Notes in Statistics*,  
221 Springer, 222, 2019.
- 222 [36] Kevin P Murphy. *Probabilistic machine learning: Advanced topics*. MIT press, 2023.