

Semi-parametric inference for the means of heavy-tailed distributions

Matt Taddy
Hedibert Lopes
David Goldberg
Matt Gardner

Semi-parametric inference for the means of heavy-tailed distributions

Matt Taddy, Hedibert Lopes, David Goldberg and Matt Gardner

Abstract

Heavy tailed distributions present a tough setting for inference. They are also common in industrial applications, particularly with Internet transaction datasets, and machine learners often analyze such data without considering the biases and risks associated with the misuse of standard tools. This article outlines a procedure for inference about the (possibly conditional) mean of a heavy tailed distribution that combines nonparametric inference for the bulk of the support with parametric inference – motivated from extreme value theory – for the heavy tail. We are able to derive analytic posterior conditional means and variances for the expected value of a heavy tailed distribution. We also introduce a simple and novel independence Metropolis Hastings algorithm that samples from the distribution for tail parameters via small adjustments to a parametric bootstrap, and through this algorithm are able to provide comparisons between our framework and frequentist semiparametric inference. We also provide a modeling extension that shrinks tails across distributions to an overall background tail. We illustrate on two examples: treatment effect estimation on a set of 72 A/B experiments, and the fitting of regression trees for prediction of user spending. Both use data from tens of millions of users of eBay.com.

1. Introduction

We refer to a data generating process (DGP) as *heavy tailed* when the distribution on exceedances beyond extreme thresholds cannot be bounded by an exponential distribution. Heavy tails are quite common in measures of user (e.g., a registered website user or a recognized device) activity on the internet (??). For example, Figure 1 illustrates spending, in US\$ spent on bought merchandise,

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

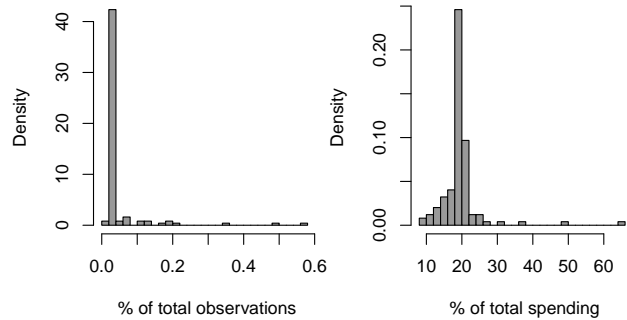


Figure 1. The proportion of observations (left) and of total spending (right) that is due to users spending greater than \$2000 on eBay.com in each of 144 A/B experiment treatment groups.

across samples of users on eBay.com. Each sample,¹ ranging in size from 10^5 to 10^8 users, corresponds to a treatment group in one of the 72 A/B experiments (randomized controlled trials) studied in Section 5. In our modal treatment group, there is less than 0.1% of users who spend more than \$2000; however, these users account for 20% of the total spending.

Such heavy tails imply that observations in high percentiles are both high variance and will have a large influence on sample means. In some cases (including, it appears, 3/4 of the groups studied in Section 5) the tail variance will be infinite. Even when the variance is merely near-infinite, these heavy tails have important consequences for our inference.

- The learning rate for mean inference is slower than \sqrt{n} , such that the usual standard error estimators underestimate uncertainty (e.g., ?).
- Common Gaussian error assumptions are invalid both in finite samples and asymptotically (?).
- Nonparametric bootstrap estimators of sampling uncertainty about the mean will fail: they are inconsistent for the true sampling distribution (?).

There are real practical implications for these issues. For example, the over-sized influence of large observations on the sample mean is well recognized by practitioners who

¹These contain targeted user subsets and are not from current traffic. They are not representative of eBay’s aggregated revenue.

measure on-line transactions. A common ad-hoc solution is to use Winsorization (?), wherein values above a threshold are replaced by that threshold. However, estimation is then sensitive to the Winsorization threshold and, due to the inconsistency of the nonparametric bootstrap, there are no obvious tools available for its optimal selection. Transaction distributions also include density spikes at zero and at other discrete points (e.g., 0.99, 100), making fully parametric modeling impractical. As another example, nonparametric learners such as random forests will tend to over-fit (and generally perform poorly) in the presence of extreme values (?).

This article resolves these issues by combining nonparametric inference for the bulk of a distribution with parametric inference for tail data above a fixed threshold. Related approaches have been proposed in the literature (see below), but we make a number of unique contributions.

- We are able to derive exact posterior moments for the distribution’s mean. These expressions require only the posterior mean and variance for the expected tail exceedance, which we make available via both analytical and efficient computational approximation.
- For inference about the tail parameters, we present a novel independence Metropolis Hastings algorithm that samples from the posterior through adjustment of the results from a parametric bootstrap. The algorithm is trivial to code, and provides information about the distance between Bayesian and frequentist inference.
- Our procedure is closely related to a frequentist algorithm that combines nonparametric and parametric bootstraps. We are able to show that this algorithm is consistent for the true sampling distribution, and the theory provides guidance on the choice of threshold.
- We describe how to shrink individual tails towards an overall background distribution. This is used both for improved estimation for average treatment effects, but also in a novel algorithm for construction of regression trees on heavy tailed data.

Finally, we test and illustrate our work in two real applications: the analysis of A/B experiments and in building robust random forests. The resulting framework is very powerful, and performs better than the alternatives, but is also very simple. We hope that it will find users amongst the large community of analysts dealing with heavy tailed data who currently rely upon sensitive ad-hoc techniques.

1.1. Related Literature

A related Bayesian approach to extreme value analysis is proposed in ? : they combine a Dirichlet process mixture

model below a threshold with a GPD above. All parameters, including the value of the threshold itself, are sampled from their joint posterior in an MCMC algorithm. Our approach is more simple and scalable: we allow for analytic expression of many of the relevant posterior statistics of interest and require only a simple bootstrap-based sampler for the tail.

? describes estimation for the mean of a heavy tailed distribution that combines the sample mean below a threshold with the mean of a maximum likelihood estimated model above that threshold. The point estimates from this approach will be similar to ours, and will converge with enough data, but our approaches to uncertainty quantification are distinct. Johansson’s asymptotic variance formulas depend upon unknown model parameters.

? provide a completely different approach to the problem, based upon sub-sampling. While the nonparametric bootstrap fails for heavy tails, Romano & Wolf show that algorithms based upon without-replacement sub-sampling can provide consistent approximations to the sampling distribution. We discuss and compare to their approach in our applications.

Finally, ? estimate the tail distribution for small samples through exponential tilting of models fit on larger samples. While their approach is totally different from our Bayesian hierarchical shrinkage technique, both works share the strategy of using background datasets to inform difficult estimation for individual tails.

2. A Semiparametric model for heavy tailed data generating processes

Our inference strategy is built around the use of Dirichlet-multinomial sampling as a flexible representation for an arbitrary data generating process (DGP). In its standard application, this model treats the observed sample as a draw from a multinomial distribution over a large but finite set of support points. A Dirichlet prior is placed on the probabilities in this multinomial, and the posterior distribution over possible DGPs is induced by the posterior on these probabilities. The approach has a long history. It was introduced by Ferguson (?), it serves as the foundation for the Bayesian bootstrap (?), and it has been studied by numerous authors (?????).

Our work presents an extension of the standard Dirichlet-multinomial scheme. Consider a univariate random variable, say z . We assume the usual fully-nonparametric model below a certain fixed *threshold*, say u . That is, the prior DGP for $z < u$ is a Multinomial draw, with Dirichlet distributed probability, from a large-but-finite number of support points. At the same time, with some probability our realized z is instead drawn as $u + v$ where $v > 0$ is

165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219

a random *exceedance* from some distribution. Hence, the full DGP model is a mixture of Multinomial sampling on arbitrary support and a parametric tail distribution.

We model our tail exceedances as realizations from a generalized Pareto distribution (GPD)², with $p(V < v) = 1 - (1 + \xi v/\sigma)^{-1/\xi}$ and density function on $v > 0$

$$\text{GPD}(e; \xi, \sigma) = \frac{1}{\sigma} \left(1 + \xi \frac{v}{\sigma}\right)^{-(\frac{1}{\xi}+1)} \quad (1)$$

for tail index $\xi > 0$ and scale $\sigma > 0$. The generalized Pareto is a commonly applied tail model (????) with theoretical justification as the limiting distribution for exceedance beyond large u for a wide family of processes (???). As $\xi \rightarrow 0$ the GPD converges to an exponential distribution, and for $\xi > 0$ the tails are heavier-than-exponential. For $\xi \geq 1/2$ the variance of v is infinite, and for $\xi \geq 1$ the mean is infinite. Thus our analysis will focus on GPD models with $\xi \in (0, 1)$, so that the tail is heavy enough to cause problems when bootstrapping³ but not so heavy that the mean does not exist.

Combining the GPD and Dirichlet-multinomial sampling yields our semi-parametric prior for heavy tailed DGPs,

$$g(z; \boldsymbol{\theta}) = \frac{1}{|\boldsymbol{\theta}|} \sum_{i=1}^L \theta_i \mathbb{1}_{[z=\zeta_i]} + \frac{\theta_{L+1}}{|\boldsymbol{\theta}|} \text{GPD}(z-u; \xi, \sigma) \mathbb{1}_{[z \geq u]} \quad (2)$$

where $\mathcal{Z} = \{\zeta_1 \dots \zeta_L\}$, all elements less than u , is the support for the bulk of the DGP $g(z)$ ⁴ and $\boldsymbol{\theta} = [\theta_1 \dots \theta_{L+1}]'$ is a vector of random weights with $\theta_l \geq 0 \forall l$.⁵

Observations are assumed drawn independently from (2) by first sampling l_i with probability θ_{l_i} and then assigning $z_i = \zeta_{l_i}$ for $l_i \leq L$ and otherwise drawing $z_i - u \sim \text{GPD}$. A posterior over g is induced by the posterior over the model parameters: $\boldsymbol{\theta}$, ξ , and σ . Functionals of g , such as $\mathbb{E}_g f(z)$ for arbitrary function f and where \mathbb{E}_g implies expectation over $z \sim g$, are thus random variables.

2.1. Inference on the sampling weights

A conjugate prior for the weights, $\boldsymbol{\theta}$, places independent exponential distributions on each element: $\theta_l \sim \text{Exp}(a_l)$ for $l = 1, \dots, L+1$, where $\mathbb{E}[\theta_l] = a_l$, $\text{var}(\theta_l) = a_l^2$,

²However, our development does not depend upon this specific tail model; you can replace the GPD with your preferred distribution while making use of the ideas in this paper.

³If you estimate $\hat{\xi} \approx 0$, then standard Bayesian bootstrap methods (e.g., ?) should apply and there is no need to model a parametric tail.

⁴We will often suppress $\boldsymbol{\theta}$ and write $g(\cdot)$ for $g(\cdot; \boldsymbol{\theta})$ unless the weights need to be made explicit.

⁵ $|\boldsymbol{\theta}|$ denotes $\sum_i |\theta_i|$, the L_1 norm.

and we call $a_l > 0$ the prior ‘rate’.⁶ We use a single⁷ rate parameter a , such that $\mathbf{a} = [a \dots a]'$. After observing a sample $\mathbf{Z} = [z_1 \dots z_N]'$, each weight remains independent in the posterior with exponential distribution $\theta_l | \mathbf{Z} \sim \text{Exp}\left(a + \sum_{i=1}^N \mathbb{1}_{[l_i=l]}\right)$.

We focus on the limiting prior that arises as $a \rightarrow 0$; see Chamberlain and Imbens (?) and Taddy et al. (??) for discussion. This ‘non-informative’ limit yields a massive computational convenience: as $a \rightarrow 0$ the weights for unobserved support points converge to a degenerate random variable at zero: $p(\theta_l = 0 | \mathbf{Z}) = 1$ if $l \neq l_i$ for any i . Our posterior for the DGP is then a multinomial sampling model with random positive weights on only the *observed data points* and on the tail model ($l_i = L+1$).

To simplify notation, say $z_i < u$ for $i \leq m$ and $z_i \geq u$ for $i = m+1, \dots, m+n$ with $N = m+n$. We then overload notation and re-write $\boldsymbol{\theta} = [\theta_1, \dots, \theta_m, \theta_{m+1}]'$ as the posterior vector of weights on observations z_1, \dots, z_m (all less than u) and on the tail. We can then write a posterior DGP realization, conditional upon GPD parameters, as

$$g(z) | \mathbf{z}, \xi, \sigma \quad (3)$$

$$= \frac{1}{|\boldsymbol{\theta}|} \sum_{i=1}^m \theta_i \mathbb{1}_{[z=z_i]} + \frac{\theta_{m+1}}{|\boldsymbol{\theta}|} \text{GPD}(z-u; \xi, \sigma) \mathbb{1}_{[z \geq u]},$$

$$\theta_i \stackrel{iid}{\sim} \text{Exp}(1) \forall i \leq m \text{ and } \theta_{m+1} \sim \text{Exp}(n).$$

This defines our posterior distribution over DGPs, with details on the GPD tail posterior deferred until Section 3.

2.2. Inference for the DGP mean

The mean of $g(z)$ is a random variable that can be written

$$\mu = \mathbb{E}z = \sum_{i=1}^m \frac{\theta_i z_i}{|\boldsymbol{\theta}|} + \frac{\theta_{m+1}}{|\boldsymbol{\theta}|} \left(u + \frac{\sigma}{1-\xi}\right). \quad (4)$$

Uncertainty about $\mathbb{E}z$ is induced by the posterior on weights $\boldsymbol{\theta}$ and on the mean exceedance $\lambda = \sigma/(1-\xi)$. Because u is fixed, we have $\boldsymbol{\theta} \perp \lambda$.

It is easy to see that $\mathbb{E}\mu = \frac{1}{m+n} \sum_{i=1}^m z_i + \frac{n}{m+n}(u + \mathbb{E}\lambda)$, while the law of total variation yields posterior variance $\text{var}\mu = \mathbb{E}[\text{var}(\mu | \lambda)] + \text{var}(\mathbb{E}[\mu | \lambda])$. Given the properties

⁶This is equivalent to the more common specification of a Dirichlet distribution on the normalized weights, written $\text{Dir}(\boldsymbol{\theta}/|\boldsymbol{\theta}|; \mathbf{a}) \propto \prod_{l=1}^{L+1} (\theta_l/|\boldsymbol{\theta}|)^{a_l-1}$.

⁷One could also use multiple rate parameters; e.g., a larger value a_l on θ_l associated with $\zeta_l = 0$ if zeros are common, or a larger a_{L+1} on θ_{L+1} corresponding to the tail data.

275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329

of the Dirichlet posterior on $\theta/|\theta|$, the first term is

$$\begin{aligned} \mathbb{E}[\text{var}(\mu|\lambda)] &= \mathbb{E} \left[\frac{\sum_{i=1}^m (z_i - \mu_\lambda)^2 + n(u + \lambda - \mu_\lambda)^2}{(m+n)(m+n+1)} \right] \\ &= \frac{\sum_{i=1}^m (z_i - \mathbb{E}\mu_\lambda)^2 + n(u + \mathbb{E}\lambda - \mathbb{E}\mu)^2}{(m+n)(m+n+1)} \\ &\quad + \frac{n^2(m+n-2)\text{var}(\lambda)}{(m+n)^2(m+n+1)} \end{aligned} \quad (5)$$

where $\mu_\lambda = [\sum_{i=1}^m z_i + n(u + \lambda)] / (m+n)$ and $\mathbb{E}\mu$ is the posterior expectation from above. The second term is

$$\text{var}(\mathbb{E}[\mu|\lambda]) = \frac{n^2}{(m+n)^2} \text{var}(\lambda) \quad (6)$$

and thus the full expression is

$$\begin{aligned} \text{var}\mu &= \frac{\sum_{i=1}^m (z_i - \mathbb{E}\mu)^2 + n(u + \mathbb{E}\lambda - \mathbb{E}\mu)^2}{(m+n)(m+n+1)} \\ &\quad + \frac{2n^2(m+n-0.5)}{(m+n)^2(m+n+1)} \text{var}(\lambda). \end{aligned} \quad (7)$$

As detailed below, $\mathbb{E}\lambda$ and $\text{var}(\lambda)$ are available through either Laplace approximation or Markov chain Monte Carlo.

3. Inference for tail parameters

In this section we describe Bayesian modeling and posterior inference for the GPD parameters, ξ and σ , conditional upon the sample of exceedances $\{v_i = z_{m+i} - u\}_1^n$.

As discussed above, we are focusing on heavy tails with finite mean exceedances that correspond to $\xi \in (0, 1)$. On this range, σ can take any positive value. A simple independent prior setup would then be

$$\begin{aligned} \pi(\sigma, \xi) &= \text{Be}(\xi; a, b) \text{Ga}(\sigma; c, d) \\ &\propto \xi^{a-1} (1-\xi)^{b-1} \sigma^{c-1} e^{-d\sigma} \end{aligned} \quad (8)$$

where $\text{Ga}(\cdot; a, b)$ denotes a beta density with mean $a/(a+b)$ and $\text{Ga}(\cdot; c, d)$ a gamma density with mean c/d , with $a, b, c, d > 0$. We will work with versions of this prior throughout, however our sampling algorithm is trivially adjusted to work for alternative specifications.

A useful version of (8) sets $a = b = 1$ and takes the non-informative limit $c, d \rightarrow 0$ to obtain

$$\pi(\sigma, \xi) = \frac{1}{\sigma} \mathbb{1}_{\xi \in (0,1)}, \quad (9)$$

the combination of a unit uniform on ξ and an improper uniform prior on $\log \sigma$. Following ? and ?, the posterior for GPD parameters will be proper under the prior in (9) given a minimum of three observations.⁸ This noninformative default is used the absence of any background information.

⁸This holds under any prior specification that combines $\pi(\sigma) \propto \sigma^{-1}$ with an independent proper prior on ξ .

Finally, this specification combines with the GPD likelihood to yield a log posterior proportional to

$$\begin{aligned} l(\sigma, \xi) &= -n \log \sigma - \left(\frac{1}{\xi} + 1 \right) \sum_i \log \left(1 + \xi \frac{v_i}{\sigma} \right) \\ &\quad + (a-1) \log \xi + (b-1) \log(1-\xi) + (c-1) \log \sigma - d\sigma, \end{aligned} \quad (10)$$

which simplifies considerably under the model in (9).

Maximization of (10) leads to MAP estimates of the parameters, say $[\hat{\xi}, \hat{\sigma}]$. The related problem of MLE estimation for GPDs is well studied by ? and his algorithm is easily adapted for fast MAP estimation within our domain $[\xi, \sigma] \in (0, 1) \times \mathbb{R}^+$.

3.1. Laplace posterior approximation

The main object of interest is actually the posterior for the GPD mean, $\sigma/(1-\xi)$. We make the transformation

$$\lambda = \frac{\sigma}{1-\xi} \Leftrightarrow \sigma = \lambda(1-\xi), \quad (11)$$

with inverse Jacobian $|J| = 1 - \xi$, to obtain the posterior

$$\begin{aligned} p(\lambda, \xi | \mathbf{v}) &\propto \\ &\frac{\xi^{a-1} e^{-d\lambda(1-\xi)}}{\lambda^{n-c+1} (1-\xi)^{n-b-c+1}} \prod_i \left(1 + \frac{\xi}{1-\xi} \frac{v_i}{\lambda} \right)^{-\left(\frac{1}{\xi}+1\right)}. \end{aligned} \quad (12)$$

Note that the MAP estimate for λ is just $\hat{\lambda} = \hat{\sigma}/(1-\hat{\xi})$. The Laplace approximation (?) to the *marginal* posterior distribution on λ is available as

$$\hat{p}(\lambda | \mathbf{x}) = \text{N} \left(\hat{\lambda}, -\nabla_{\lambda\lambda}^{-1} \Big|_{[\hat{\lambda}, \hat{\xi}]} \right), \quad (13)$$

where $\nabla_{\lambda\lambda}$ is the curvature of the log posterior⁹ with respect to λ . Thus

$$\begin{aligned} \text{var}(\lambda | \mathbf{x}) &= -\nabla_{\lambda\lambda}^{-1} \Big|_{[\hat{\lambda}, \hat{\xi}]} \\ &= -\hat{\lambda}^2 \left[n - c + 1 + \left(\frac{1}{\hat{\xi}} + 1 \right) \sum_i (\hat{q}_i^2 - 2\hat{q}_i) \right]^{-1}, \end{aligned} \quad (14)$$

with $\hat{q}_i = \hat{\xi} v_i / [(1-\hat{\xi})\hat{\lambda} + \hat{\xi} v_i]$, is the approximate posterior variance for λ .

3.2. Independence-MH via the parametric bootstrap

For small tail samples the Laplace approximation will dramatically underestimate posterior uncertainty. Instead, we propose a novel independence Metropolis Hastings (i-MH) algorithm (see, e.g., ?) that uses a parametric bootstrap of

⁹Via $\partial \log p(\lambda, \xi | \mathbf{v}) / \partial \lambda = [(1/\xi + 1) \sum_i q_i - n + c - 1] / \lambda - d(1-\xi)$ where $q_i = \xi v_i / [(1-\xi)\lambda + \xi v_i]$.

the MAP estimates as a proposal distribution in Markov Chain Monte Carlo. This approach is similar in spirit to the bootstrap reweighting of (?), but unlike that work it does not require an analytic expression for the sampling distribution of the statistics of interest. The algorithm proceeds as follows.

- Fit the MAP parameter estimates $[\hat{\xi}, \hat{\sigma}]$.
- Obtain B draws $[\hat{\xi}_b, \hat{\sigma}_b]$ from the parametric bootstrap:
 - Generate a sample of size n by simulating from the MAP estimated GPD model.
 - Obtain new MAP estimates $[\hat{\xi}_b, \hat{\sigma}_b]$ conditional upon this simulated sample.
- Estimate the bivariate bootstrap distribution, say $r(\xi, \sigma)$, via kernel smoothing on $\{\hat{\xi}_b, \hat{\sigma}_b\}_{b=1}^B$.
- For $b = 2 \dots B$, replace $[\hat{\xi}_b, \hat{\sigma}_b]$ with $[\hat{\xi}_{b-1}, \hat{\sigma}_{b-1}]$ with probability

$$1 - \min \left\{ \frac{r(\hat{\xi}_{b-1}, \hat{\sigma}_{b-1}) \exp[l(\hat{\xi}_b, \hat{\sigma}_b)]}{r(\hat{\xi}_b, \hat{\sigma}_b) \exp[l(\hat{\xi}_{b-1}, \hat{\sigma}_{b-1})]}, 1 \right\}$$

where l is the log posterior objective in (10).¹⁰

In addition to being fast and simple, this algorithm offers a bridge between Bayesian and frequentist inference: if the acceptance probabilities are high, then there is little difference between the sampling distribution and the posterior.

4. A semiparametric heavy tailed bootstrap

This section studies a bootstrap algorithm that is closely related to our semiparametric Bayesian procedure. Consider frequentist inference about $Q_N = \sqrt{N}(\hat{\mu}_N - \mu)$ for a sample of N observations drawn from true distribution function $F(z)$, with $\int_0^\infty z dF(z) = \mu < \infty$ and where $\hat{\mu}_N$ denotes the MLE for μ based upon a size- N sample from F .

A frequentist bootstrap replaces $F \approx \hat{F}_N$ and uses this to obtain $b = 1, \dots, B$ draws of $Q_N^b = \sqrt{N}(\hat{\mu}_N^b - \hat{\mu}_N)$ where $\hat{\mu}_N^b$ is the MLE based upon a size- N sample from \hat{F}_N . The targeted sampling distribution, $G_N(q) = p(Q_N < q)$, is then estimated as $\hat{G}_N(q) = B^{-1} \sum_{b=1}^B \mathbb{1}_{\{Q_N^b < q\}}$.

Standard results on bootstrap consistency (??) require that \hat{G}_N converges in distribution (i.e., weakly) to G_∞ *uniformly* across all \hat{F}_N in a neighborhood, say \mathcal{F} , containing F and also \hat{F}_N for N big enough.¹¹ Convergence in probability for $\hat{F}_N(z)$ to $F(z) \forall z$ then implies consistency of

¹⁰This ‘rejection’ probability – the probability of *not* moving states – is equal to one minus the familiar acceptance probability.

¹¹In addition, the mapping $F \mapsto G_\infty$ must be continuous.

\hat{G}_N in that, as $N \rightarrow \infty$, $p(|\hat{G}_N(q) - G_N(q)| < \epsilon) \rightarrow 0$ for all $q, \epsilon > 0$.

? shows that the nonparametric bootstrap – using the empirical distribution function (EDF) as \hat{F}_N – is inconsistent for the distribution of the sample mean for data that has infinite variance. As explained by ?, in this setting \hat{G}_N based upon samples from \hat{F}_N does not converge uniformly to G_∞ because sums of the largest re-sampled observations, $\sum_{i=N-r}^N z_{(i)}^b$ for $r \geq 1$, can be dominated by repeats of the largest sample observation, $z_{(N)}$.

Instead, define a semiparametric bootstrap that takes the MLE tail parameters, $[\hat{\xi}_n, \hat{\sigma}_n]$, and for $b = 1, \dots, B$

- draw $m_b \sim \text{Bin}(m/N, N)$ and set $n_b = N - m_b$;
- sample with replacement m_b observations z_i where $i \leq m$ (i.e., $z_i < u$), say $\{z_1^b, \dots, z_{m_b}^b\}$;
- generate $v_1^b \dots v_{n_b}^b$ from $\text{GPD}(\hat{\xi}_n, \hat{\sigma}_n)$ and obtain the corresponding MLE, $\hat{\lambda}_{n_b}^b = \hat{\sigma}_{n_b}^b / (1 - \hat{\xi}_{n_b}^b)$; and
- set $\hat{\mu}_N^b = \left(\sum_{i=1}^{m_b} z_{i_b} + n_b(u + \hat{\lambda}_{n_b}^b) \right) / N$.

The results can then be applied in estimation of the sampling distribution, e.g., for $\sqrt{N}(\hat{\mu}_N - \mu)$ as approximately equal to the sampled distribution on $\sqrt{N}(\hat{\mu}_N^b - \hat{\mu}_N)$.

The distribution for $\hat{\mu}_N$ implied by our semiparametric bootstrap is the combination of three bootstrap estimators, for distributions on $\frac{1}{m} \sum_{i=1}^N z_i \mathbb{1}_{[z_i < u]}$, on m/N , and on $\hat{\lambda}_n$. Consistency of the nonparametric bootstrap for the first two statistics¹² can be established through standard arguments (?). To show consistency for our semiparametric bootstrap, we need to confirm that the parametric bootstrap using $\hat{F}_N(z - u | z \geq u) = \text{GPD}(\hat{\xi}_n, \hat{\sigma}_n)$ converges to the correct distribution for $\hat{\lambda}_n$.

? considers DGPs with distribution functions $F(z) = 1 - cz^{-1/\zeta}(1+z^{-\delta}L(z))$, where $c, \delta > 0$ and $L(tz)/L(z) \rightarrow 1$ with $z \rightarrow \infty$ for $t > 0$. This defines a wide class of heavy tailed distributions, and for u_N large enough the distribution $F(z - u_N | z \geq u_N)$ approaches a GPD(ξ, σ_N) where $\sigma_N = u_N \xi$. Following the same steps as Johansson, which apply results from ? on the asymptotic distribution for MLEs $[\hat{\xi}_n, \hat{\sigma}_n]$, you can show that for $F(z)$ with $\xi \in (0, 1)$ and $z^{-\delta}L(z)$ non-increasing, if $u_N = O(N^{\xi/(1+2\delta\xi)})$ then

$$\sqrt{n}(\hat{\lambda}_n - \mathbb{E}_F[z - u_N | z \geq u_N]) \rightarrow_d N(0, q_n) \quad (15)$$

where $q_n = \hat{\sigma}_N(1 + \xi)(1 - \xi + 2\xi^2)/(1 - \xi)^4$. Thus our bootstrap sample generator, $\text{GPD}(\hat{\xi}_n, \hat{\sigma}_n)$, converges

¹²The first is the mean from sampling a bounded domain, the second is the proportion of successes in a binomial trial of size N .

495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549

to $F(z - u_N | z \geq u_N)$ along a sequence of distributions with means $\hat{\lambda}_n$ that are asymptotically normal around the target of interest, $\mathbb{E}_F[z - u_N | z \geq u_N]$. From ?, this is enough to establish consistency of this tail bootstrap, and hence of our full semiparametric bootstrap.

Intuitively, the parametric tail bootstrap succeeds here because our MLEs allow us to converge quickly to the ‘true’ GPD model; inference is then based upon *new* samples from this distribution and, unlike resamples from the EDF, these are not overly influenced by large order statistics in the original sample. From (15), this convergence holds so long as u_N is growing at the right rate. Theoretically, the true GPD tail that we converge to with u_N has $\sigma_N = \xi u_N$. One can compare this σ_N to $\hat{\sigma}_n$, the sample MLE for the same parameter, to see if they roughly agree. If not, u_N might need to increase. In practice, however, we also advocate repeating inference over a range of thresholds u and using results from the region where they stabilize.

We now return to a Bayesian framework for the remainder of this article, but the reader who is interested in frequentist inference may feel free to use the semiparametric bootstrap of this section instead. The two algorithms are very closely related: there is little difference between the Bayesian bootstrap and the nonparametric bootstrap on the bulk of the distribution (?) and the i-MH sampler of Section 3.2 explicitly agrees with the frequentist parametric tail bootstrap whenever acceptance probabilities are close to one. Given this connection, the discussion above indicates that our semiparametric Bayesian inference will also have good frequentist properties in large samples.

5. A/B Experiments

Our motivating application for these ideas involves A/B experiments: two independent heavy tailed samples are obtained, one from a group receiving a treatment and another from a group receiving a control. The object of interest is then

$$\gamma = \mu_1 - \mu_0 \quad (16)$$

where μ_1 is the mean of the treatment group and μ_0 the mean of the control group.

Our main inferential target is the average treatment effect in an A/B experiment: the difference in *average* response, say z , between a group that received some treatment (e.g., a change to the website, some marketing, a new search algorithm) and a control group that experiences the status quo. That is, for a binary treatment indicator d , where $d_i = 1$ for treated and $d_i = 0$ for untreated, we seek to estimate the treatment group means $\mu_d = \mathbb{E}[z|d]$ and the average treatment effect

$$\gamma = \mu_1 - \mu_0. \quad (17)$$

We wish to estimate and test for $\gamma \neq 0$ in a setting where

each treatment group has $\text{var}(z|d) = \infty$.

A standard analysis estimates γ with $\bar{z}_1 - \bar{z}_0$, the difference of sample means. A naive standard error for this estimate is $\sqrt{\text{sse}(z_0)/N_0^2 + \text{sse}(z_1)/N_1^2}$, where $N_d = m_d + n_d$ in the notation from above: m_d and n_d are the numbers of observations below and at-or-above the tail cutoff, respectively. The results for such an analysis are shown in Figure 2; out of our 72 experiments there are 8 significant p -values at a 10% False Discovery Rate. However, both point and standard-error estimates here will be sensitive to extreme values in the sample. Moreover, the standard error formula is incorrect for distributions with infinite variance, where the learning rate is less than $\sqrt{N_d}$.

there is a prior belief that the experiments do not make much of a difference in the tail. That is, we suspect that in most cases the users spending a lot of money will act the same regardless of what website experience they are met with. Unfortunately, for infinite variance data it is common that a few very large observations can have an outsized influence on the mean estimates (and those of γ). Thus we wish to accurately quantify uncertainty in the presence of heavy tails and avoid having them overly influence our inference.

5.1. Tail analysis results

Using a cutoff of $u = 2000$, we analyzed the exceedances across our 72 experiments. Four example experiments are in Figure 3; each analysis used 10,000 iterations of the I-MH sampler.

- Compared to results for the MCMC (I-MH) posterior samples, the Laplace approximation variance is far too small.
- The parametric bootstrap proposal distribution is very similar to the I-MH sampler. Indeed, 85% of our experiments the average acceptance probability was greater than 0.9. Thus our tail analysis is very similar to that of a parametric bootstrap.
- Posterior predictive fits look good. As predicted by extreme value theory, these improve even more for larger u (e.g., $u = 3000$) at the expense of smaller tail-sample-sizes.
- All tail indices ξ have 100% of posterior weight between 0 and 1. This implies that the exceedences have a finite mean but infinite variance.

5.2. Semiparametric Bayesian results

Combining the parametric tail analysis with nonparametric inference below the cutoff, we obtain a new set of effect and uncertainty estimates. The plot in Figure 4 shows inference

550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604

605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659

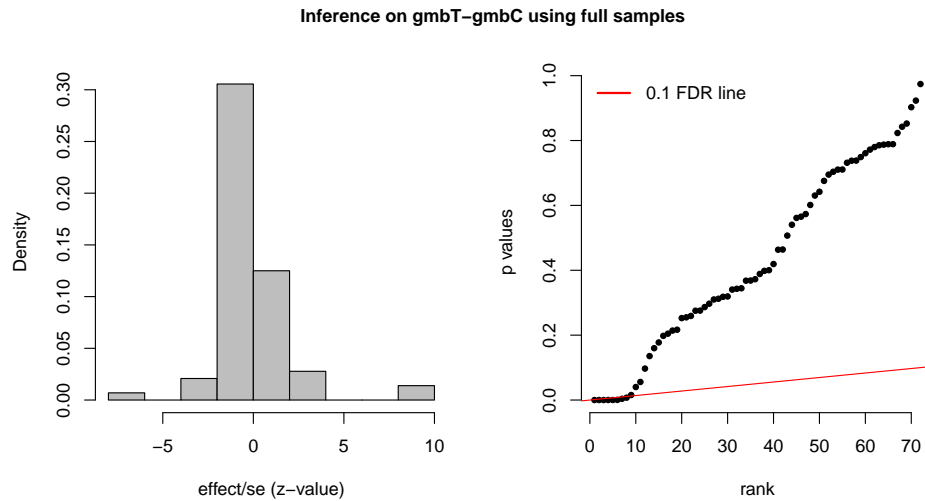


Figure 2. Inference using naive sample mean differences and standard errors.

analogous to our earlier results in Figure 2 (we are abusing concepts and reporting standardized effects and p-values based on the Bayesian posterior moments). The significant effect estimates are roughly similar: we'd make nearly the same decisions at 10% FDR. However, there are bigger differences than it first appears. In the next section we see, for example, that the naive confidence intervals can be far different from the correct posterior intervals.

5.3. Capping

We now investigate the effect of *capping* on estimation. In this procedure, one replaces response values above a certain 'cap' with that value. It is viewed as a variance reduction technique. Results for four of our experiments are shown Figure 5; in each the naive sample mean differences and standard errors are applied to the capped values, as well as to the raw uncapped values (red line), and we compare to the semiparametric Bayesian analysis (for $u = 2000$).

- The point estimates from higher capped values are near to the posterior mean than those from either low caps or the uncapped values.
- The capped value standard errors are always much smaller than the true posterior standard deviation. Note that these standard errors are correct for the mean of capped values; it is for the true overall mean – our object of interest – that they underestimate uncertainty.
- A safe-ish procedure appears to be to use a high cap to get the point estimate, but use the standard error calculated on the raw values to assess uncertainty. Note

that, since we're learning at a rate $< \sqrt{N}$, these values still likely tend to underestimate uncertainty.

5.4. Sensitivity to cutoff

We consider how the semiparametric Bayesian analysis responds to changes in the tail cutoff, u . Figure 6 shows that results are quite stable to this choice; e.g., compare to figure 3 to see that the capping procedure is much more sensitive even over this limited range.

6. Extensions

There are some nice things we could do with this.

6.1. Hierarchical modeling

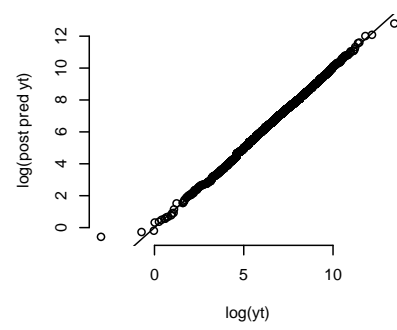
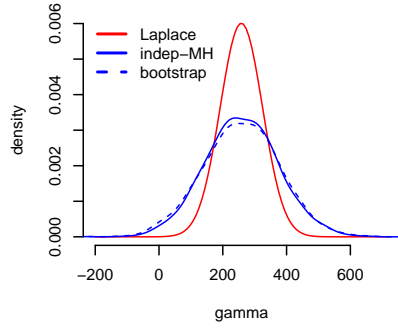
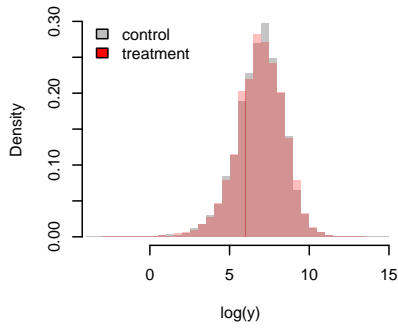
We (or rather eBay analysts) have a strong prior that the treatments do not make much of a difference in the tail. We could build this into the analysis by shrinking all of the tails to an overall mean. An efficient Empirical Bayes procedure would first calculate a GPD model for the aggregate tail, then use this to create a 'prior' model for the individual tail GPD parameters. You could also have things evolve in time, etc.

6.2. Heavy tailed regression

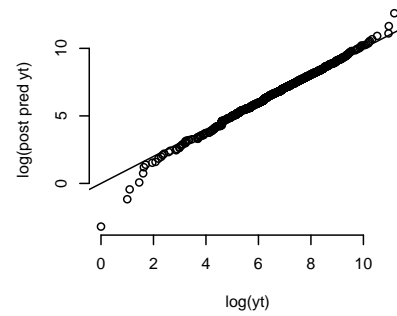
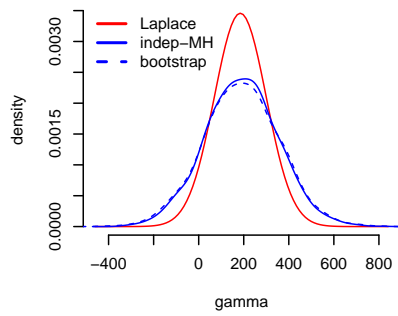
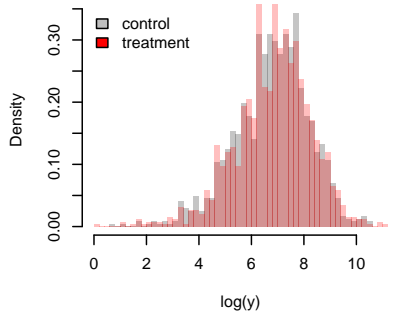
These same considerations apply to any regression model that we might want to fit where the response is heavy tailed (e.g., if the response is GMB). It would be good to look at both ordinary least-squares (OLS) and trees (CART, random forests) to see what we can improve.

770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824

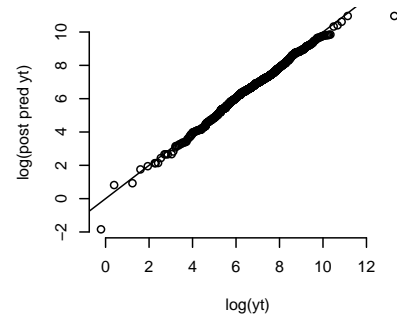
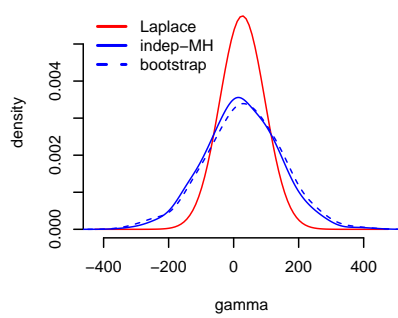
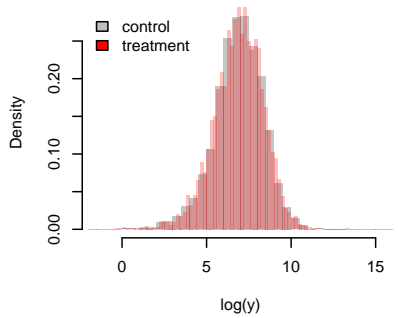
treatment 33203



treatment 32350



treatment 33444



treatment 32099

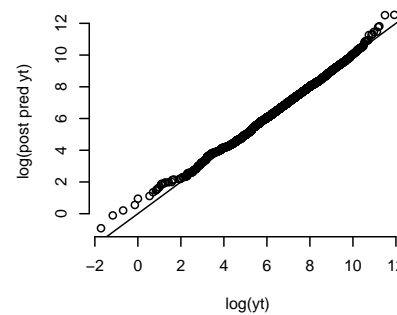
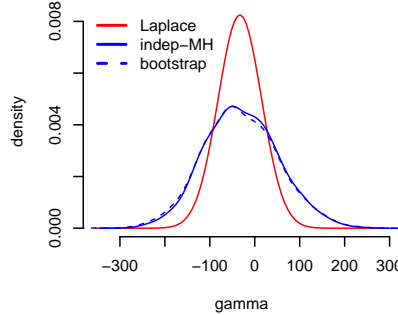
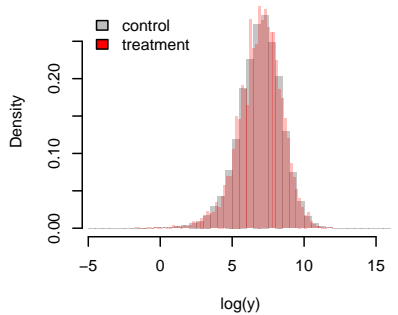


Figure 3. Tail analysis for exceedances beyond $z = 2000$. The left panel shows the histograms of exceedances by treatment group, the center panel shows both Laplace approximate and Independence MH posterior distributions for the treatment effect in the tail ($\lambda_1 - \lambda_0$), and the right plot compares quantiles for the observed tail data in the treatment group to those of a posterior predictive sample of the same size. The dashed lines in the center panel are the parametric bootstrap proposal distribution.

825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879

Posterior inference on gmbT-gmbC

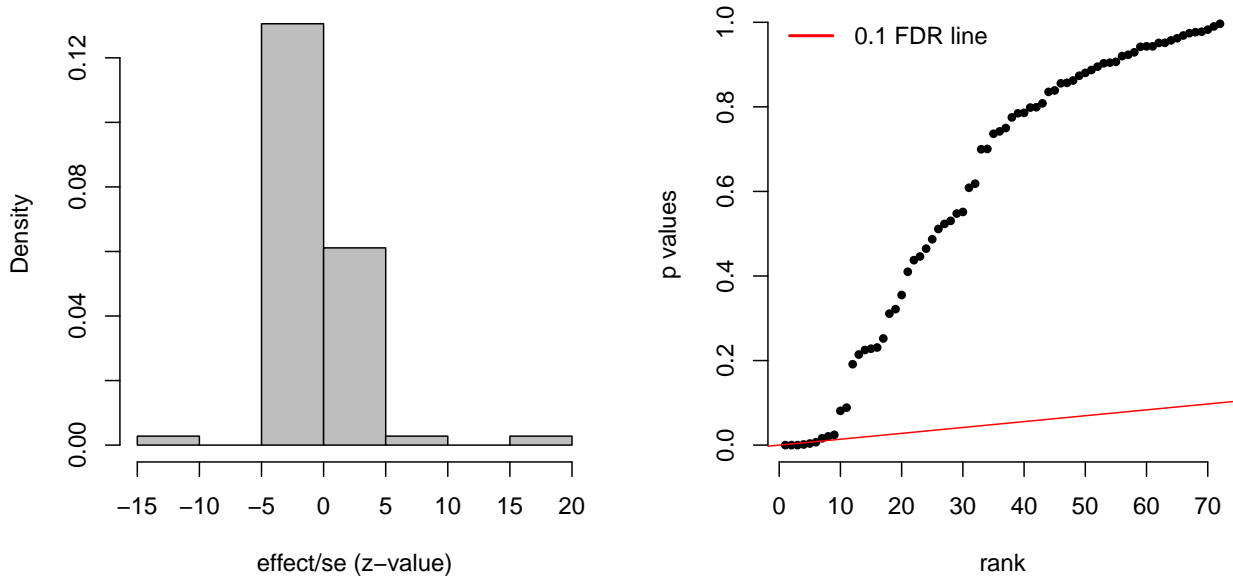


Figure 4. Posterior inference results when combining the parametric tail analysis with nonparametric inference below the cutoff at $u = 2000$. The ‘effects’ here are posterior means, while the ‘se’ are posterior standard deviations. For comparison with Figure 2, we report results on the the same scale of standardized effects and Normal distribution ‘p-values’.

7. Conclusion

This is great Big Data analysis: avoiding modeling and difficulty for the easy bits (the middle of the distribution), and applying good solid statistical modeling on the hard bits (the tail). Since the posterior distribution and the parametric bootstrap for the tail models are virtually indistinguishable, you don’t need to be Bayesian to buy the conclusions.

A. MAP estimation under the default prior

Maximization of the log posterior in (10) is unaffected by a re-parametrization¹³, so we make the convenient translation $\alpha = 1/\xi$ and $\tau = \xi/\sigma$ to obtain the maximization objective

$$l(\alpha, \tau) = (n + 1)(\log \alpha + \log \tau) - (\alpha + 1) \sum_i \log(1 + \tau v_i). \quad (18)$$

A similar replacement is advocated for maximum likelihood estimation in (?)¹⁴. The gradient with respect to α

¹³Note that we are just solving for a maximum, not translating the distribution to a different parameter space.

¹⁴Where $k = -\xi$.

is

$$\nabla_\alpha = \frac{n + 1}{\alpha} - \sum_i \log(1 + \tau v_i). \quad (19)$$

This is solvable as a function of τ , so that at the maximum of (18) we can make the replacement

$$\hat{\alpha}(\tau) = \frac{n + 1}{\sum_i \log(1 + \tau v_i)}. \quad (20)$$

By substitution into (18), the resulting *profile* objective function is then proportional to

$$l(\tau) = \log[\tau \hat{\alpha}(\tau)] - \frac{1}{\hat{\alpha}(\tau)} \quad (21)$$

with gradient

$$\nabla_\tau = \frac{n + 1}{\tau} - (\hat{\alpha}(\tau) + 1) \sum_i \frac{v_i}{1 + \tau v_i}. \quad (22)$$

Obtaining roots for ∇_τ over the domain $\tau > 0$ is straightforward. We just note that MAP $\hat{\tau}$ will be very small for many fat tailed datasets and thus you need to set an appropriately small convergence tolerance (10^{-8} is fine for our examples). Finally, the roots for (22) yield MAP estimates

$$\left[\hat{\xi}, \hat{\sigma} \right] = \left[\frac{1}{\hat{\alpha}(\hat{\tau})}, \frac{1}{\hat{\tau} \hat{\alpha}(\hat{\tau})} \right]. \quad (23)$$

880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934

935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989

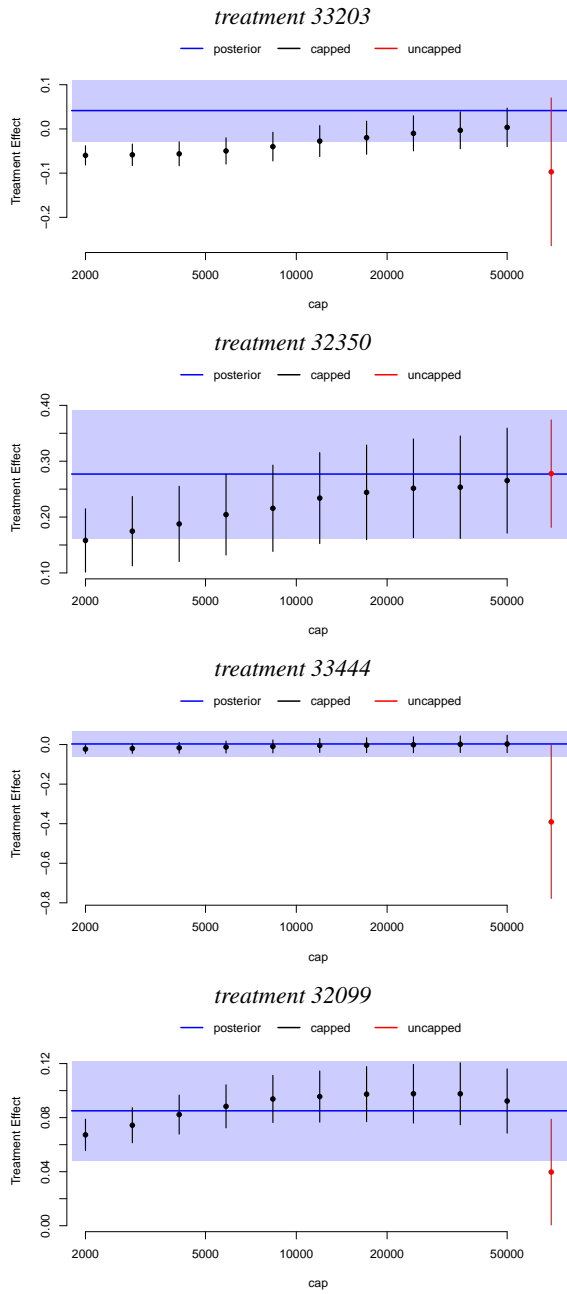


Figure 5. A study of capping (and naive uncapped analysis) in comparison to the semiparametric Bayesian posterior. In each case, intervals plotted are $\pm 1sd$ or $\pm 1se$, as appropriate.

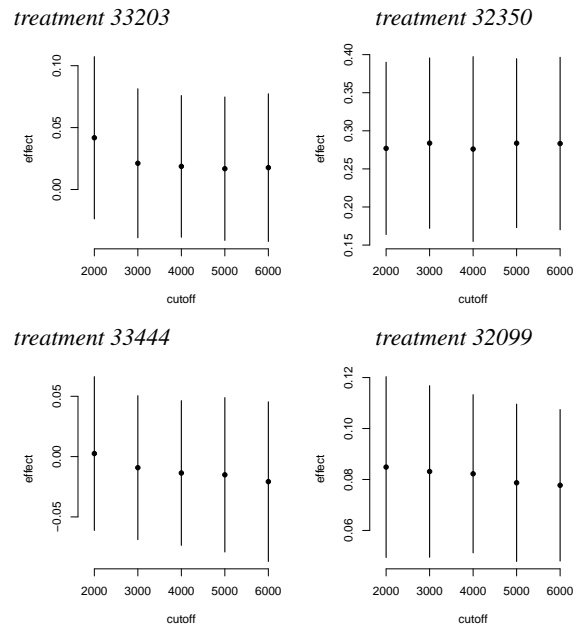


Figure 6. The posterior mean for γ , plus or minus 1 standard deviation, as a function of tail cutoffs $u = 2, 3, 4, 5$ thousand.

990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044

1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099