

# Big Data + Bayesian Nonparametrics

**Matt Taddy**

Microsoft Research and University of Chicago

<http://taddylab.com>

# Big Data

The sample sizes are enormous.

- ▶ We'll see 21 and 200 million today.
- ▶ Data can't fit in memory, or even storage, on a single machine.

The data are super weird.

- ▶ Internet transaction data distributions have a big spike at zero and spikes at other discrete values (e.g., 1 or \$99).
- ▶ Big tails (e.g., \$12 mil/month eBay user spend) that matter.
- ▶ The potential feature space is unmanageably large.
- ▶ We cannot write down or measure believable models.

Both 'Big' and 'Strange' beg for nonparametrics.

## Bayesian nonparametrics

In usual BNP you *model* a complex generative process with flexible priors, then apply that model directly in prediction and inference.

$$\text{e.g., } y = f(\mathbf{x}) + \varepsilon \quad \text{or} \quad f(y|\mathbf{x}) = f(\mathbf{x}, y)/f(\mathbf{x})$$

'Flexibility' comes from a huge number of variables,

$$\text{e.g., } f(\mathbf{z}) = \sum_{k=1}^{\infty} \omega_k N(\mathbf{z}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k),$$

and learning requires integration over uncertainty on nuisance parameters, e.g.,  $k_i$  component membership for observation  $\mathbf{z}_i$ .

For years this integration was done via MCMC. But these algorithms did not scale for big data. [Can we do scalable BNP?](#)

## Classical distribution-free nonparametrics

Frequentists are great at finding simple procedures (e.g.  $[\mathbf{X}'\mathbf{X}]^{-1}\mathbf{X}'\mathbf{y}$ ) and describing how they work under a variety of possible true DGP.

( DGP = Data Generating Process )

This is 'distribution free' nonparametrics.

- 1: Find some statistic that is useful regardless of DGP.
- 2: Derive the distribution for this stat under minimal assumptions.

Practitioners apply the simple stat and feel happy that it will work.

## Distribution-free Bayesian nonparametrics

Find some *statistic of the DGP* that you care about:

- ▶ derive from first principles, e.g. moment conditions
- ▶ *an algorithm* that we know works, e.g. CART
- ▶ think about geometric projections, e.g. OLS

Call this statistic  $\theta(g)$  where  $g(\mathbf{z})$  is the DGP (e.g., for  $\mathbf{z} = [\mathbf{x}, y]$ ).

Then you write down a flexible model for the DGP  $g$ , and study properties of the posterior on  $\theta(g)$  induced by the posterior over  $g$ .

## A flexible model for the DGP

Say  $\mathbf{z} = [\mathbf{x}, y]$  is a single *independent* data point.

Each data point assumes one of a *finite* number of possible values,  $[\zeta_1 \dots \zeta_L]$ , with probabilities proportional to  $[\theta_1 \dots \theta_L]$ .

$$g(\mathbf{z}) = \frac{1}{|\boldsymbol{\theta}|} \sum_{l=1}^L \theta_l \mathbb{1}_{[\mathbf{z}=\zeta_l]}$$

We complete specification with a conjugate prior on the weights:

$$\frac{\boldsymbol{\theta}}{|\boldsymbol{\theta}|} \sim \text{Dir}(\mathbf{a}) \propto \frac{1}{|\boldsymbol{\theta}|^{L(a-1)}} \prod_l \theta_l^{a_l-1} \quad \text{where } a, \theta_l > 0.$$

This is the Dirichlet-multinomial sampling model (Ferguson 1973).

Now you've observed some data, say  $\mathbf{Z} = \{\mathbf{z}_1 \dots \mathbf{z}_n\}$ .  
(say every  $\mathbf{z}_i = [\mathbf{x}_i, y_i]$  is unique).

The posterior over weights has  $\theta_l \stackrel{ind}{\sim} \text{Exp}(a + \mathbb{1}_{[\zeta_l \in \mathbf{Z}]})$ .

## A convenient limiting case

$a \rightarrow 0$  leads to  $p(\theta_l = 0) = 1$  for  $\zeta_l \notin \mathbf{Z}$ .

In this case, we can focus on only the *observed support* and write the posterior for our DGP

$$g(\mathbf{z}) = \frac{1}{|\theta|} \sum_{i=1}^n \theta_i \mathbb{1}[\mathbf{z} = \mathbf{z}_i], \quad \theta_i \stackrel{iid}{\sim} \text{Exp}(1).$$

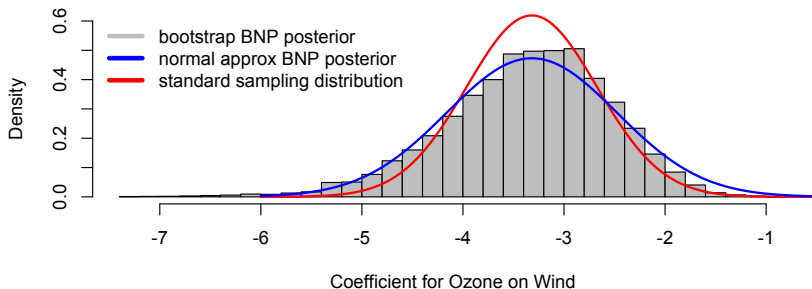
This is just the Bayesian bootstrap. (Rubin 1981)

## Example: Ordinary Least Squares

*Population* OLS is a posterior functional

$$\beta = (\mathbf{X}'\Theta\mathbf{X})^{-1}\mathbf{X}'\Theta\mathbf{y}$$

where  $\Theta = \text{diag}(\theta)$ . This is a random variable. (sample via BB)





## What is the blue line?

BB sampling is great, but analytic approximations are also useful.

Consider a first-order Taylor series approximation,

$$\tilde{\beta} = \hat{\beta} + \nabla\beta|_{\theta=1}(\theta - 1)$$

where  $\hat{\beta} = [\mathbf{X}'\mathbf{X}]^{-1}\mathbf{X}'\mathbf{y} = \beta|_{\theta=1}$  is the sample OLS line.

We can derive *exact* posterior moments for  $\tilde{\beta}$ . e.g.,

$$\text{var}(\beta) \approx \text{var}(\tilde{\beta}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\text{diag}(\mathbf{e}^2)\mathbf{X}'(\mathbf{X}'\mathbf{X})^{-1}$$

where  $\mathbf{e} = \mathbf{y} - \mathbf{X}\hat{\beta}$  are the sample OLS residuals.

This is the economist's Huber-White 'Sandwich' variance formula.

(Lancaster 2003 or Poirier 2011.)

## Example: Heterogeneous Treatment Effects

eBay runs lots of experiments: they make changes to the marketplace (website) for random samples of users.

Every experiment has response  $y$  and treatment  $d$  [0/1].

In our illustrative example,  $d_i = 1$  for bigger pictures in my eBay.

$i \in \text{control} : d_i = 0$

$i \in \text{treatment} : d_i = 1$



This is a typical 'A/B experiment'.

What is 'heterogeneity in treatment effects'? (HTE)

Different units [people, devices] respond differently to some treatment you apply [change to website, marketing, policy].

I imagine it exists.

We know  $\mathbf{x}_i$  about user  $i$ . About 400 features in our example.

- ▶ Their previous spend, items bought, items sold...
- ▶ Page view counts, items watched, searches, ...
- ▶ All of the above, broken out by product, fixed v. auction, ...

Can we accurately measure heterogeneity: index it on  $\mathbf{x}$ ?

## Treatment effect regression

**Potential outcomes:**  $v_i(d)$  is the response for user  $i$  if  $d_i = d$ .

We'd love to model  $\mathbb{E}[v_i(1) - v_i(0)|\mathbf{x}_i]$ , but we can't directly.

Since  $d \perp\!\!\!\perp [\mathbf{x}, \mathbf{v}]$ , a good alternative is to consider

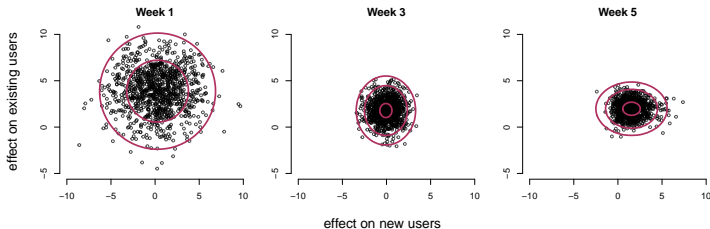
$$\mathbb{E}[v_i(1)|d_i = 1, \mathbf{x}_i] - \mathbb{E}[v_i(0)|d_i = 0, \mathbf{x}_i].$$

Even simpler, just diff the OLS projections for each group

$$\beta_t - \beta_c = (\mathbf{X}'_t \Theta_t \mathbf{X}_t)^{-1} \mathbf{X}'_t \Theta_t \mathbf{y}_t - (\mathbf{X}'_c \Theta_c \mathbf{X}_c)^{-1} \mathbf{X}'_c \Theta_c \mathbf{y}_c$$

e.g., coefficient on *[made purchase this year]* vs an intercept:

*million users:*                      7.45                      10.97                      13.22



Sample is from posterior, contours are normal approximation.

The most basic application of an A/B experiment is estimating the average treatment effect,

$$\text{ATE} = \mathbb{E}[v(t)] - \mathbb{E}[v(c)]$$

The usual estimator is  $\bar{y}_t - \bar{y}_c$ . But many authors advocate instead the 'regression adjusted'  $\bar{\mathbf{x}}'(\hat{\beta}_t - \hat{\beta}_c)$  (Lin, Berk+, Deng+, 2013).

Both are [frequentist] unbiased.

But quantifying uncertainty about the regression adjustment, and any variance reduction it provides, is tough without assumptions.

From our BNP perspective, both the usual and adjusted ATE estimators are random functions of the DGP.

$$\begin{aligned} \text{ATE}_g^{\text{obs}} &= \frac{1}{|\boldsymbol{\theta}_t|} \boldsymbol{\theta}_t' \mathbf{y}_t - \frac{1}{|\boldsymbol{\theta}_c|} \boldsymbol{\theta}_c' \mathbf{y}_c \\ \text{ATE}_g^{\text{lin}} &= \frac{1}{|\boldsymbol{\theta}|} \boldsymbol{\theta}' \mathbf{X} \left( (\mathbf{X}_t' \boldsymbol{\Theta}_t \mathbf{X}_t)^{-1} \mathbf{X}_t' \boldsymbol{\Theta}_t \mathbf{y}_t - (\mathbf{X}_c' \boldsymbol{\Theta}_c \mathbf{X}_c)^{-1} \mathbf{X}_c' \boldsymbol{\Theta}_c \mathbf{y}_c \right). \end{aligned}$$

Using our usual tricks, we can resolve any *controversy*:

$$\text{var} \left( \text{ATE}_g^{\text{lin}} \right) \approx \text{var}(\text{ATE}_g^{\text{obs}}) - \left( \frac{R_t^2 s_{y_t}^2}{n_t^2} + \frac{R_c^2 s_{y_c}^2}{n_c^2} \right),$$

so that big variance reduction comes only for small sample  $n_d$  or large  $R_d^2 = \text{cor}(y_d, \hat{y}_d)$ . In digital experiments we have big samples and low signal, so we shouldn't expect much variance reduction.

## Example: Decision Trees

Trees are great: nonlinearity, deep interactions, heteroskedasticity.



The 'optimal' decision tree is a statistic we care about (s.w.c.a).



## CART: greedy growing with optimal splits

Given node  $\{\mathbf{x}_i, y_i\}_{i=1}^n$  and DGP weights  $\theta$ , find  $x$  to minimize

$$\begin{aligned} |\theta| \sigma^2(x, \theta) = & \sum_{k \in \text{left}(x)} \theta_k (y_k - \mu_{\text{left}(x)})^2 \\ & + \sum_{k \in \text{right}(x)} \theta_k (y_k - \mu_{\text{right}(x)})^2 \end{aligned}$$

for a regression tree. Classification impurity can be Gini, etc.

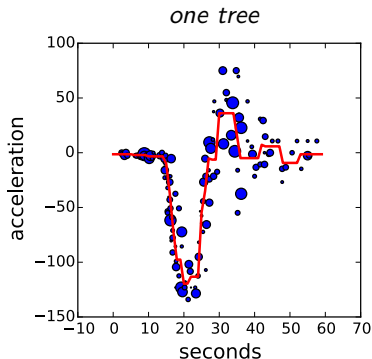
Population-CART might be a statistic we care about.

Or, in settings where greedy CART would do poorly (big  $p$ ), a randomized splitting algorithm might be a better s.w.c.a.

## Bayesian Forests: a posterior for CART trees

For  $b = 1 \dots B$ :

- draw  $\theta^b \stackrel{iid}{\sim} \text{Exp}(\mathbf{1})$
- run weighted-sample CART to get  $\mathcal{T}_b = \mathcal{T}(\theta^b)$



Random Forest  $\approx$  Bayesian forest  $\approx$  posterior over CART fits.

## Theoretical **trunk** stability

Given forests as a posterior, we can start talking about *variance*.

Consider the first-order approximation

$$\begin{aligned}\sigma^2(x, \boldsymbol{\theta}) &\approx \sigma^2(x, \mathbf{1}) + \nabla \sigma^2|_{\boldsymbol{\theta}=\mathbf{1}}(\boldsymbol{\theta} - \mathbf{1}) \\ &= \frac{1}{n} \sum_i \theta_i [y_i - \bar{y}_i(x)]^2\end{aligned}$$

with  $\bar{y}_i(x)$  the sample mean in  $i$ 's node when splitting on  $x$ .

Based on this approx, we can say that for data at a given node,

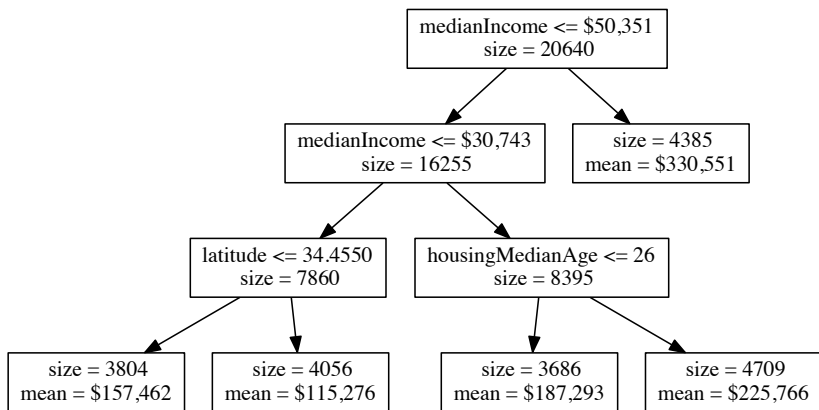
$$p(\text{optimal split matches sample CART}) \gtrsim 1 - \frac{p}{\sqrt{n}} e^{-n},$$

with  $p$  split locations and  $n$  observations.

(Taddy, Chen, Yu, Wyle 2015)

## California Housing Data

20k observations on median home prices in zip codes.



Above is the trunk you get setting min-leaf-size of 3500.



- ▶ sample tree occurs 62% of the time.
- ▶ 90% of trees split on income twice, and then latitude.
- ▶ 100% of trees have 1st 2 splits on median income.

Empirically and theoretically: trees are stable, at the trunk.

## Random Forest **bottlenecks**

RFs (or BFs) are awesome, but when the data are too big to fit in memory or on a single machine they get extremely expensive. (e.g., Google PLANET, Panda 2009).

A common solution is a 'sub-sampling forest': instead of drawing with-replacement, or re-weighting, draw  $m \ll n$  sub-samples.

This defeats the whole purpose of trees: these are rules that are designed to grow in complexity with the amount of data available. (that's why we bother storing so much data!)

If you starve the individual trees of data you lose.

## Empirical Bayesian Forests (EBF)

RFs are expensive. Sub-sampling hurts bad.

Instead:

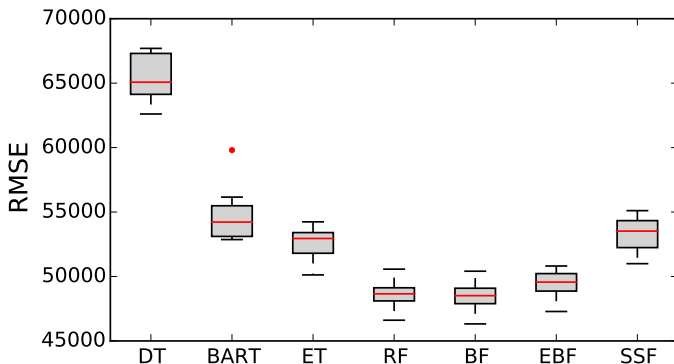
- ▶ fit a single tree to a shallow **trunk**.
- ▶ Map data to each **branch**.
- ▶ Fit a full forest on the smaller branch datasets.

Since the trunks are all similar for each tree in a full forest, our EBF looks nearly the same at a fraction of computational cost.

It's an updated version of classical **Empirical Bayes**:

use plug-in estimates at high levels in a hierarchical model,  
focus effort at full Bayesian learning for the the hard bits.

## OOS predictive performance on California Housing

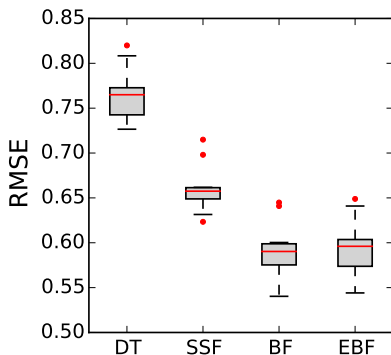


Here EBF and BF give nearly the same results. *SSF does not.*

EBFs crunch more data faster without hurting performance.



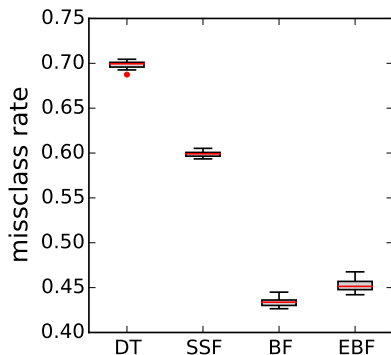
## EBFs work all over the place



RMSE	% WTB	
0.5905	0.0	BF
0.5953	0.8	EBF
0.6607	11.9	SSF
0.7648	29.5	DT

Predicting wine rating from chemical profile

## EBFs work all over the place



MCR	% WTB	
0.4341	0.0	BF
0.4531	4.4	EBF
0.5989	38.0	SSF
0.6979	60.8	DT

or beer choice from demographics

## Choosing the trunk depth

Distributed computing perspective: **fix only as deep as you must!**

How big is each machine? Make that your branch size.

	CA housing			Wine			Beer		
<i>Min Leaf Size in <math>10^3</math></i>	6	3	1.5	2	1	0.5	20	10	5
<i>% Worse Than Best</i>	1.6	2.4	4.3	0.3	0.8	2.2	1.0	4.4	7.6

Still, open questions: e.g., more trees vs shallower trunk?

## Catching **Bad Buyer Experiences** at eBay

BBE: 'not as described', delays, etc.

$p(\text{BBE})$  is a key input to search rankings, and elsewhere.

Big Data axiom: best way to improve prediction is more data.

(a.k.a. 'data beats model' )

On 200 million transactions, EBF with 32 branches yields a 1.5% drop in misclassification over the SSF alternatives.

EBFs via Spark  $\Rightarrow$  more data in less time.

Putting it into production requires some careful engineering, but this really is a very simple algorithm. **Big gain, little pain.**

## Back to HTE: Treatment Effect Trees

Recall our earlier example: A/B experiment with user covariates  $\mathbf{x}_i$ .

Athey+Imbens propose indexing user HTE by fitting CART to

$$y_i^* = y_i \frac{d_i - q}{q(1 - q)} = \begin{cases} y_i / (1 - q) & \text{if } d_i = 0 \\ y_i / q & \text{if } d_i = 1 \end{cases}$$

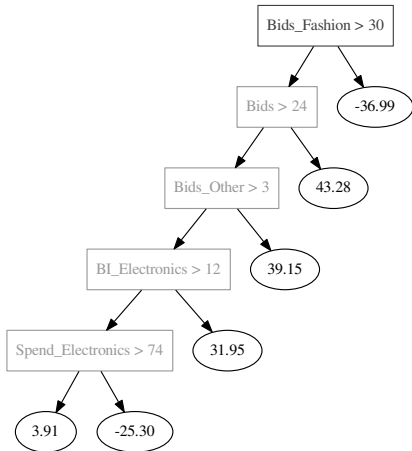
where  $q$  is the probability of treatment (2/3 in our example).

This works because

$$\mathbb{E}[y_i^* | \mathbf{v}_i] = v_i(1) - v_i(0).$$

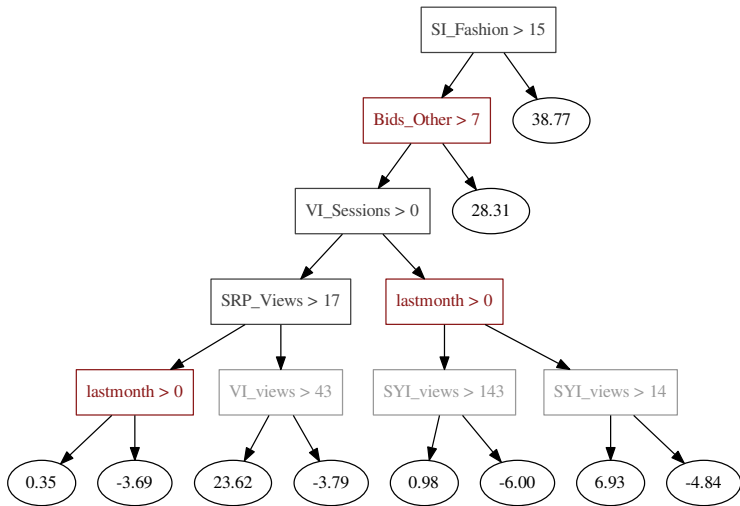
How sure are we of the fitted tree? [Bayes forest is a posterior.](#)

e.g., Pruned CART with after one week:



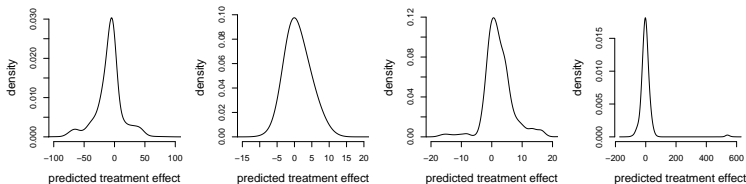
prob variable is node in tree:  $p < \frac{1}{3}$ ,  $p \in [\frac{1}{3}, \frac{1}{2})$ , and  $p \geq \frac{1}{2}$ .

After 5 weeks the tree is much more stable.

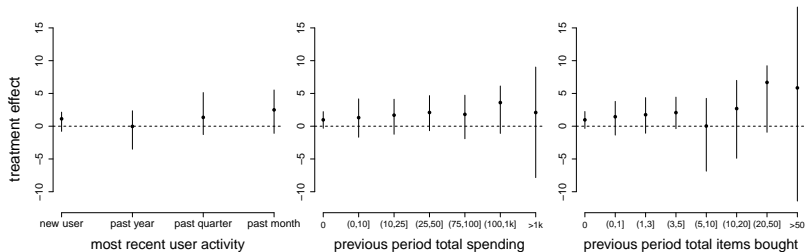


But you should still use a forest below a certain point.

*It is wise to look to the full forest instead:*



Treatment effect posteriors for sampled users (5-week data).

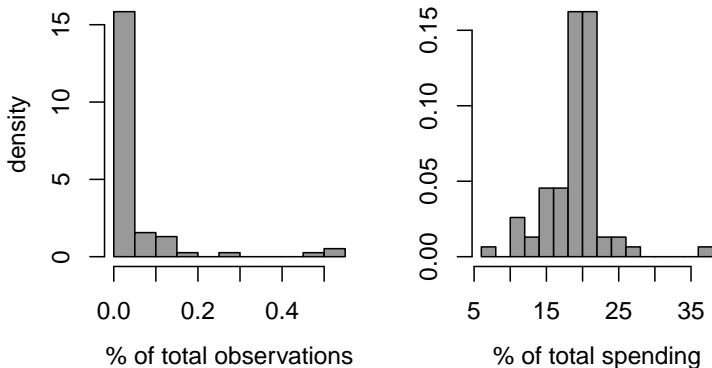


Variable average effects for three important user features.



## Scalable semiparametrics for heavy-tailed distributions

In some settings we have really heavy tailed data.



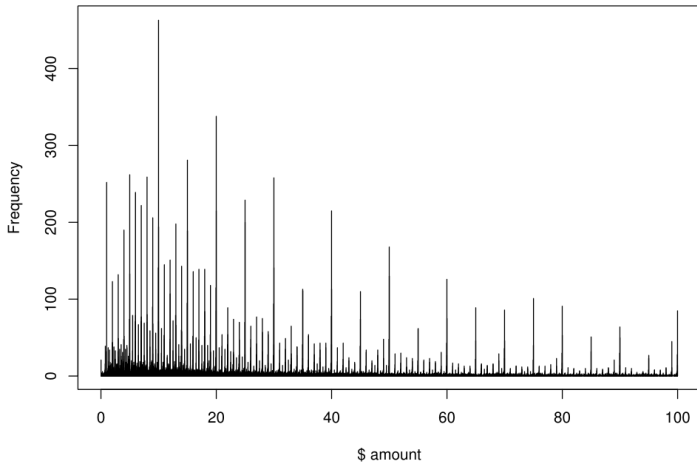
The standard bootstrap fails for infinite variance distributions.

We should suspect our Bayesian bootstrap will have similarly bad frequentist properties. **e.g., the posterior will not be consistent.**

## Lots of spikes!

14,521 unique \$ amounts from 50,000 users.

78.18% of the \$ amounts are less than \$100.



So the distribution below the tail is also screwy.

## Semi-parametrics: only model where you really need it.

We can replace our DGP with a model that includes a parametric tail above some threshold  $u$ :

$$p(z) = \frac{1}{|\boldsymbol{\theta}|} \sum_{l=1}^L \theta_l \mathbb{1}[z = \zeta_l] + \frac{\theta_{L+1}}{|\boldsymbol{\theta}|} \text{GPD}(z - u; \xi, \sigma) \mathbb{1}[z \geq u]$$

where  $\text{GPD}(\cdot; \xi, \sigma)$  is the generalized Pareto density function

$$p(y) = \frac{1}{\sigma} \left(1 + \xi \frac{y}{\sigma}\right)^{-\left(\frac{1}{\xi} + 1\right)}$$

with  $\xi > 0$  the tail index and  $\sigma > 0$  the scale.

Say  $\mu = \mathbb{E}[z]$  and  $\lambda = \mathbb{E}[z - u | z \geq u]$ .

$\mathbb{E}[\mu]$  and  $\text{var}(\mu)$  are available in closed form up-to  $\mathbb{E}[\lambda]$  and  $\text{var}[\lambda]$ .

## Inference for tail parameters

Everything for the weights is as before (exponential posterior).

For the GPD tail parameters, we use the prior

$$\pi(\sigma, \xi) = \frac{1}{\sigma} \xi^{a-1} (1 - \xi)^{b-1} \mathbb{1}_{\xi \in (0,1)},$$

For noninformative inference, use  $a = b = 1$ . But it is cooler to use a bank of related 'background' tail transactions to estimate  $\hat{\alpha}, \hat{\beta}$ .

We have a super simple independence MH posterior sampler that is based upon re-sampling from a parametric bootstrap.

Or, you can use a laplace approximation.

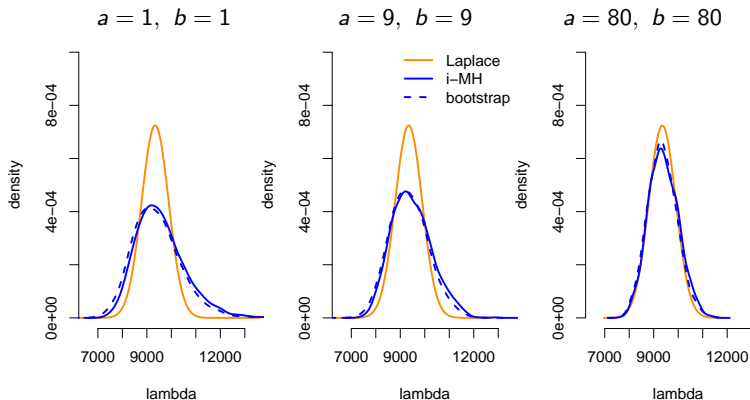
## Bootstrap-based posterior sampling

- ▶ 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  $n$  from the MAP GPD model.
  - ▶ Obtain new MAP estimates  $[\hat{\xi}_b, \hat{\sigma}_b]$  for this simulated sample.
- ▶ Get a kernel estimate, say  $r(\xi, \sigma)$ , for the bootstrap density.
- ▶ Replace each  $[\hat{\xi}_b, \hat{\sigma}_b]$  with  $[\hat{\xi}_{b-1}, \hat{\sigma}_{b-1}]$  with probability

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

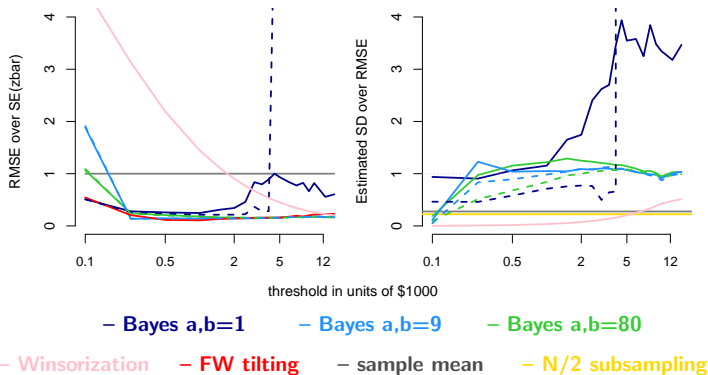
Super fast to run and code. And it relates Bayes to Frequentism!

Inference for the mean exceedance  $\lambda = \sigma/(1 - \xi)$ ,  
beyond  $u = \$9000$ , in one of our treatment groups.



These results give all we need for inference about the overall mean.

Performance over 100 resamples of  $N = 50,000$  from  $10^7$  total.



## Consistency and thresholds

The *semiparametric bootstrap* is consistent if  $u_N = O(N^{\xi/(1+2\delta\xi)})$ .  
Which is cool, because the nonparametric bootstrap is not.

In this limit,  $\sigma_N = \xi u_N$ . So, one can increase  $u$  until this holds.  
Indeed, it holds roughly in our examples for good performing  $u$ .

In practice, **results are robust to a range of tail thresholds**.

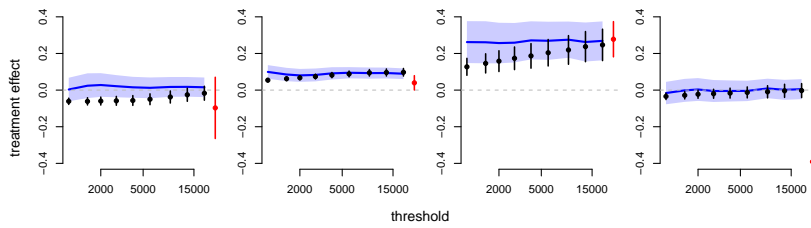
This also suggests why informative priors on  $\sigma$  are not useful:  
it is already informed by both the threshold and the tail index.



# Semi-parametrics in treatment effect estimation

Semiparametric inference makes a big difference in A/B analysis. Looking at ATEs,  $\mu_1 - \mu_2$ , in a few example experiments:

sp-bayes, capped, naive



Such care in uncertainty quantification is essential whenever you need to understand the full posterior (or sampling) distribution. e.g., bandit learning with heavy-tailed rewards is another application.

## Big Data and distribution-free BNP

I think about BNP as a way to analyze (and improve) algorithms.  
Decouple action/prediction from the full generative process model.

This type of analysis is key for taking these algorithms from ML to fields like Economics that really care about uncertainty.

## Efficient Big Data analysis

To cut computation without hurting performance, we need to think about what portions of the 'model' are **hard** or **easy** to learn.

Once we figure this out, we can use a little bit of the data to learn the easy stuff and direct our full data at the hard stuff.

This is *the* future for Big Data.

**thanks!**