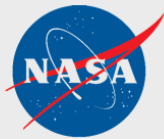


JPL-Caltech Virtual Summer School

Big Data Analytics

September 2 – 12, 2014

Amy Braverman
(Jet Propulsion Laboratory)
Subsampling



National Aeronautics and
Space Administration

Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

Outline

Whether Frequentist or Bayesian, we need to know how the sampling distribution of $\hat{\theta}$ depends on the (realized) value (of $\Theta \Rightarrow \theta$).

Can we get this information without making assumptions about the distribution*, $f_{\hat{\theta}|\Theta}(\hat{\theta}|\theta)$?

* Please forgive the abuse of notation: using $\hat{\theta}$ to represent both a random variable and its realized value.



National Aeronautics and
Space Administration

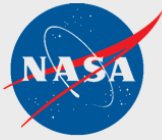
Jet Propulsion Laboratory
California Institute of Technology
Pasadena, California

Outline

Introduction to subsampling:

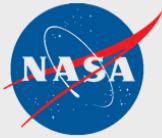
- ▶ Subsampling.
- ▶ Comments and cautions.

Material in this section based substantially on notes provided by Professors Anirban DasGupta of Purdue University and Charlie Geyer of the University of Minnesota, and on the Politis, Romano, and Wolf 1999 book, *Subsampling*.



- ▶ A subsample is a sample of size $M < N$ from Y_1, \dots, Y_N drawn *without* replacement. Denote the subsample by V_1, \dots, V_M .
- ▶ V_1, \dots, V_M is a sample from the *true* distribution, F .
- ▶ Compute $\hat{\theta}^* = g(V_1, \dots, V_M)$. Recall that $\hat{\theta} = g(Y_1, \dots, Y_N)$.
- ▶ Repeat the subsampling B times and obtain $\hat{\theta}_1^*, \dots, \hat{\theta}_B^*$.
- ▶ The CDF of $\hat{\theta}$ is estimated by $L_{Sub,M}(t)$:

$$L_{Sub,M}(t) = P_F(\sqrt{M}(\hat{\theta}^* - \hat{\theta}) \leq t) \approx \frac{1}{B} \sum_{b=1}^B 1 \left[\sqrt{M} (\hat{\theta}_b^* - \hat{\theta}) \leq t \right].$$



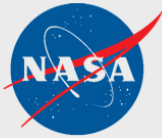
Only two things are required:

- ▶ We must assume that $N^\alpha (\hat{\theta} - \theta)$ converges (in probability) to J where J is some non-degenerate distribution, and α is some exponent. We don't have to know either of them- just that they exist.
- ▶ We have to choose M such that

$$M \longrightarrow \infty, \quad \frac{M}{N} \longrightarrow 0 \quad \text{as } N \longrightarrow \infty.$$

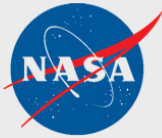
If so, then $L_{Sub,M}(t)$ converges (in probability) to $J(t)$ as $N \longrightarrow \infty$. *

* Stated imprecisely and without all technical conditions. See Politis, Romano, and Wolfe (1999) page 43 for details.

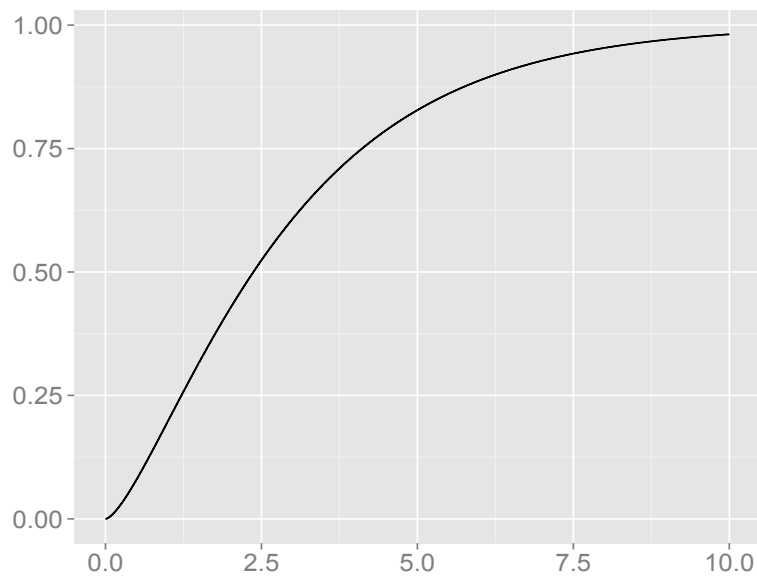


Remarks:

- ▶ Assumptions less restrictive than required for the Bootstrap.
- ▶ “Works” for situations where the Bootstrap won’t: stationary m -dependent sequences (time series), extremes, etc.
- ▶ The catch: must choose M to be big, but not too big. How?



F



Example of the concept: $M = 3$

$$\mathbf{y} = (2.82, 1.71, 2.07, 4.60, 3.39)^T, \quad \hat{\theta}(\mathbf{Y}) = 2.82,$$

$$\mathbf{v}_1 = (4.60, 1.71, 2.07)^T, \quad \hat{\theta}_1^*(\mathbf{Y}) = 2.07,$$

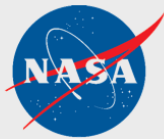
$$\mathbf{v}_2 = (2.07, 2.82, 1.71)^T, \quad \hat{\theta}_2^*(\mathbf{Y}) = 2.07,$$

$$\mathbf{v}_3 = (3.39, 1.71, 2.82)^T, \quad \hat{\theta}_3^*(\mathbf{Y}) = 2.82,$$

\vdots

$$\mathbf{v}_B = (1.71, 3.39, 2.82)^T, \quad \hat{\theta}_B^*(\mathbf{Y}) = 2.82.$$

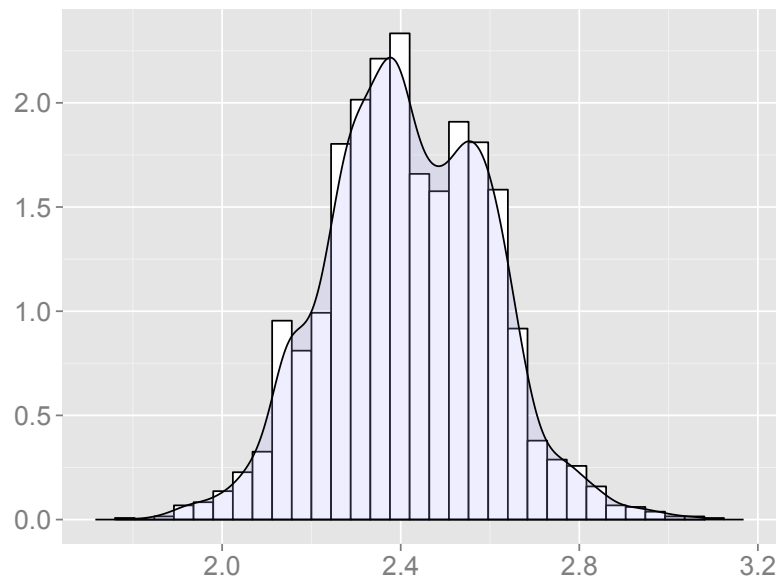
Yes, it's true that it all depends on the sample, but we are taking expectations over all possible samples.



Subsampling

Subsampling distribution of the sample median, $\hat{\theta}^*$

$$N = 3000, B = 3000, M = \lfloor N^{2/3} \rfloor$$

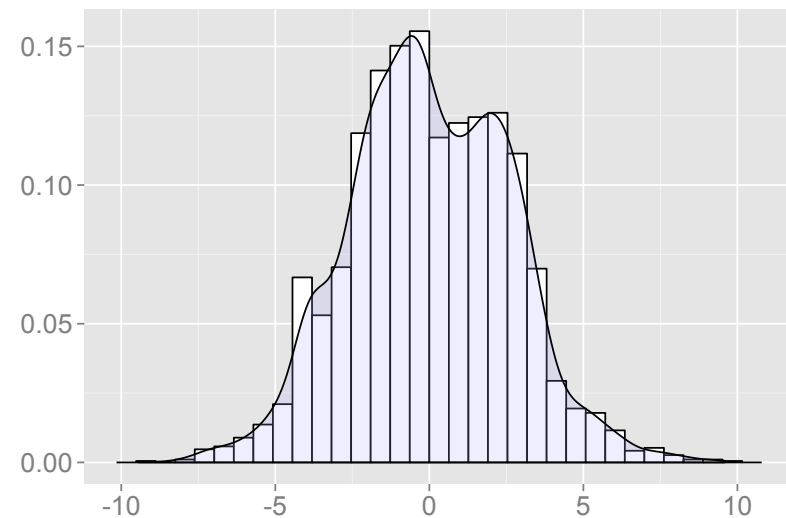


$$E(\hat{\theta}^*) = 2.418,$$
$$\text{var}(\hat{\theta}^*) = .033$$

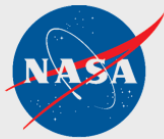
Subsampling distribution of the sample median,

$$\sqrt{M}(\hat{\theta}^* - \hat{\theta})$$

$$N = 3000, B = 3000, M = \lfloor N^{2/3} \rfloor$$

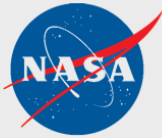


$$E(\sqrt{M}(\hat{\theta}^* - \hat{\theta})) = .015,$$
$$\text{var}(\sqrt{M}(\hat{\theta}^* - \hat{\theta})) = 6.798$$



	Expected value	Variance
True: $\sqrt{N}(\hat{\theta} - q_{.50})$	-.049	7.118
Bootstrap: $\sqrt{N}(\hat{\theta}^* - \hat{\theta})$.257	9.063
M/N bootstrap: $\sqrt{M}(\hat{\theta}^* - \hat{\theta})$.111	7.342
Subsampling: $\sqrt{M}(\hat{\theta}^* - \hat{\theta})$.015	6.798

So, which one to use?



Quoting/paraphrasing Charlie Geyer:

- ▶ The Right Thing is samples of the right size, N , from the right distribution, F .
- ▶ The Politis and Romano thing (subsampling) is samples of the wrong size, M , from the right distribution, F .
- ▶ The Efron thing (bootstrapping) is samples of the right size, N , from the wrong distribution, \hat{F}_N .
- ▶ Both Wrong Things are wrong. We would like to do the right thing, but we can't because we have only one sample.
- ▶ Which Wrong Thing do we want to do?