# Bootstrapping
(Givens/Hoeting Cht. 9 ; Gentle Ch 4

<u>Motivation</u>: The motivation for the bootstrap method is the same as that for the jackknife. We have an estimator T for some parameter $\theta$ of a distribution F and we would like to know characteristics (bias, variance, conf. intv) of the distribution of the estimator T.

   Ideally we would have many samples from F allowing us to generate multiple estimates $T_i$ & use these to obtain this information. However this is unrealistic and so again we turn to clever ways of utilizing a single sample.

   For the jackknife the basic idea was to remove portions of the sample & then recalc the est, eventually looking at the mean of the weighted differences of the samples. For the bootstrap the main idea is that of resampling.

In general, Resampling methods involve the use of many samples each drawn from a single sample from F. Then the conditional distribution of the new sample is used for statistical inference. Since the sample set is finite, it is often easy to compute statistical functions on the sample, and thus gain info even when very little is known about F.

For the bootstrap, the basic idea is that an observ. sample should contain all of the infor. mation about the underlying population & so the observed

Sample is considered to be the population. Hence, the distribution of the test statistic $T$ can be simulated by using random samples from the "population" (ie observed $\S$ original sample).

## More Formally. (using the text's notation).

Sec 9.1

- Let $\Theta = T(F)$ be our parameter of interest of a distribution $F$, expressed as a functional of $F$. Usually $T(F) = \int g(z)\, dF(z)$.
  For ex $T(F) = \int z\, dF(z)$ is the mean of the dist.

*$\int F \sim f$ means $f$ when $F$ is c.d.f.*

- Let $X_1, \ldots, X_n$ be ~~an~~ data observed as a real: of the r.v. $X_1, \ldots, X_n \sim$ i.i.d. $F$. ($F$ is c.d.f.) and $X = \{X_1, \ldots, X_n\}$ denote the entire data set.

*discrete uniform*

- Let $\hat{F}$ be the empirical dist of the observed data. (Basical weights each obs. $x_i$ w prob $1/n$).

*(as usual)*

- An est. of $\Theta$, $\hat{\Theta} = T(\hat{F})$.
  Ex if $\Theta$ is pop. mean $\hat{\Theta} = \int z\, d\hat{F}(z) = \sum_{i=1}^{n} x_i/n$.

Now, we want to know about the dist of our estimator $T(\hat{F})$. And so we ask ques. about $T(\hat{F})$ or some $R(X, F)$, a statistical function of the data $\&$ their unknown dist. $F$.

Ex/
$$R(X, F) = \frac{T(\hat{F}) - T(F)}{S(\hat{F})}$$
where $S(\hat{F})$ est. s.d. of $T(\hat{F})$

$R(X, F)$ could be the bias of $T(\hat{F})$, Var, etc. It is the function of interest. (More general than $T(\hat{F})$) and allows for easy know. of data set dependency).

As we mentioned before, the dist of $R(X, F)$ is <sub>prob.</sub> unknown & may be intractable. So we use the emp. dist of the obs. data (from F) to approx. the dist of $R(X, F)$.

- The bootstrap method:

sampled w/replacement.

Let $X^* = \{X_1^*, \ldots, X_n^*\}$ be i.i.d r.v. drawn from dist. $\hat{F}$. Then $X^*$ is called the bootstrap sample of pseudo data or pseudo dataset

The bootstrap strategy is to examine the dist. of $R(X^*, \hat{F})$ and use it to make inf. about $R(X, F)$. (Easy to resample $X^*$ from $X$).

Comment: in some special cases we can derive analytical results about the dist. of $R(X^*, \hat{F})$ however, usually it is done through sim.

9.1 pg 254

Ex/ Suppose $X = \{x_1, x_2, x_3\} = \{1, 2, 6\}$ is obs from F. and we wish to est. the mean of F, $\Theta$.

Clearly, $\hat{\Theta} = T(\hat{F})$ or $R(X, F)$ is the sample mean $\hat{\Theta} = 9/3$. we would like to have more info about the dist. of $\hat{\Theta}$.

The empirical dist of F is $\hat{F}$ which places mass $1/3$ at each observed value.

A bootstrap sample $X^* = \{X_1^*, X_2^*, X_3^*\}$ will be elements drawn i.i.d. from $\hat{F}$. There are $3^3 = 27$ possibilities for $X^*$, each with prob $1/27$.

  Ex ( $X^* = \{1, 1, 6\}$ )

• Let $\widehat{F}^*$ denote the empirical dist func. of $X^*$
ex ($\widehat{F}^*$ puts $2/3$ at 1 + $1/3$ at 6).

• and we have the new corresp. bootstrap est
$$\widehat{\theta}^* = T(\widehat{F}^*) \qquad ex\left(\widehat{\theta}^* = 2/3 + 6(1/3) = 8/3\right).$$

The bootstrap strategy is to examine
(analy. if poss. ow. via simulation the dist
of $\widehat{\theta}^*$).

For this case $\widehat{\theta}^*$ can take on 10 dist. values
since the order of the data does. matter.

| $X^*$ | $\widehat{\theta}^*$ | $P^*[\widehat{\theta}^*]$ |
|-------|--------------|----------------|
| 1, 1, 1 | $3/3$ | $1/27$ |
| 1, 1, 2 | $4/3$ | $3/27$ |
| 1, 2, 2 | $5/3$ | $3/27$ |
| 2, 2, 2 | $6/3$ | $1/27$ |
| 1, 1, 6 | $8/3$ | $3/27$ |
| 1, 2, 6 | $9/3$ | $6/27$ |
| 2, 2, 6 | $10/3$ | $3/27$ |
| 1, 6, 6 | $13/3$ | $3/27$ |
| 2, 6, 6 | $14/3$ | $3/27$ |
| 6, 6, 6 | $18/3$ | $1/27$ |

$\rightarrow$

Sum to 1.

• For each $X^*$ we calc. $\widehat{\theta}^*$
• For each $\widehat{\theta}^*$ we can write down its prob
$P^*[\widehat{\theta}^*]$ w.r.t. the bootstrap exp of
drawing $X^*_\wedge$ ~~from~~ $X$.   $(R(X^*, \widehat{F}))$
          conditional on

• The bootstrap principle is to equate the
dist of $R(X, F)$ w/ that of $R(X^*, \widehat{F})$.
and make our inference: (based on dist of $\widehat{\theta}^*$).

So

$$P^*\left[\hat{\theta}^* \leq \frac{6}{3}\right] = 8/27.$$

or a simple $25/27$ $(\sim 93\%)$ conf. interval for $\theta$ is $(4/3, 14/3)$.

## Sec 9.2. Basic Methods.
## Sec. 9.2.1 Nonparametric bootstrap.

For realistic sample sizes $(n)$. the number of potential bootstrap pseudo-data sets is very large $(n^n)$ and so complete enumer. is not possible. Instead, B ind. rand. bpd. are drawn from the empirical dist. $\hat{F}$. of the observed data.

- i.e. draw $X_i^* = \{X_{i,1}, \ldots, X_{i,n}\}$ for $i=1,\ldots,B$ and then use $R(X_i^*, \hat{F})$ $(i=1,\ldots,B)$ to approx. dist of $R(X, F)$.

Comments: — no parametric assumptions are required (don't need any info about unknown dist)
- can answer wider range of ques.
- often more accurate than standard. parametric theory.
- can make simulation error small by increasing B.

Ex/ Previous ex. non-parametric bootstrap.

$$\hat{F} = \begin{cases} 1 & 1/3 \\ 2 & 1/3 \\ 6 & 1/3 \end{cases}$$

Generate $X_i^*$ by sampling $X_{i,1}^*, X_{i,2}^*, X_{i,3}^*$ w/ replacement from $\{1,2,6\}$. Each $X^*$ yields a $\hat{\theta}^*$. Ex pg 254 shows freq of $\hat{\theta}^*$ when B=1000. These freq approx $P[\hat{\theta}^*]$.

**Sec 9.2.2.   Parametric Bootstrap.**

In ordinary bootstrap $X^*$ gener. by drawing $X_1^*, .. X_n^*$ ii.d. from $\hat{F}$. (The empirical dist w/ no other assumptions).

However, if it is believed that $F$ is a parametric distribution $F(x, \Theta)$ another method may be employed.

**Parametric Bootstrap.**

 - Draw $X_1, ..., X_n \sim$ ii.d $F(x, \Theta)$. and use $X$ to est. $\Theta$, $\hat{\Theta}$.
 - Each parametric bootstrap pseudo-data set $X^*$ is generated ~~from~~ by drawing $X_1^*, ..., X_n^* \sim$ iid. $F(x, \hat{\Theta})$.

**Ex/**

Assume $X_1, ..., X_n$ are $\text{Normal}(\Theta, 1)$.

compute $\hat{\Theta} = \Sigma X_i / n$

draw $X_1^*, ..., X_n^*$ $\text{Normal}(\hat{\Theta}, 1)$.

If the model is known or believed to be a good fit, this is a powerful tool, allowing inference & producing more accurate confidence intervals than standard asymptotic theory.

**Sec 9.2.4**
**Sec 4.1**    **Bootstrap Bias Correction**

We want to know the bias of $T(F) = \Theta$.
∴ we are interested in using bootstrap analysis on $R(X, F) = T(\hat{F}) - T(F) = \hat{\Theta} - \Theta$.

We want to know $\text{Exp}[\text{Bias}] = E[\hat{\theta} - \theta] = E[\hat{\theta}] - \theta$  exp w.r.t. $F$.

So using the bootstrap principle we calc.

$$E^*[\text{Bias } \hat{\theta}^* - \hat{\theta}] = E^*[\hat{\theta}^*] - \hat{\theta} = \overline{\theta}^* - \hat{\theta}$$

where $\quad \overline{\theta}^* = \sum_{j=1}^{B} \hat{\theta}_j^* / B$ .   each $\hat{\theta}_j^*$ based on $x_j^*$ drawn from $\hat{F}$.

And so

our new bias corrected estimator is

$$\theta_{BIAS} = \hat{\theta} - (\overline{\theta}^* - \hat{\theta}) = 2\hat{\theta} - \overline{\theta}^* .$$

and should have less bias than $\hat{\theta}$.

## Sec 4.2   Bootstrap Estimation of Variance

Suppose now we wish to estimate the variance of $T(\hat{F}) = \hat{\theta}$.

We create this estimate by calculating the sample variance of our bootstrap estimates $\hat{\theta}_j^*$.

$$\hat{V}(\hat{\theta}) = \hat{V}(\theta^*) = \frac{1}{B-1} \sum_{j=1}^{B} (\hat{\theta}_j^* - \overline{\theta}^*)^2$$

## Sec 9.3.1   Bootstrap Confidence Intervals — Percentile Method
## Sec 4.3

The simplest method for drawing inference about a univariate para. $\theta$ using bootstrap simulations is to construct a confidence interval using the percentile method. This method amounts to reading percentiles off the histogram of $\hat{\theta}^*$ values produced by bootstrapping.

We had ex. earlier, look pg. 258 Fig 9.1 we can find a bootstrap est of $1-\alpha$ conf. interval based on $((1-\alpha/2)100th)$ ! $(\alpha/2\ 100th)$ emp. percentiles

of the histogram.

Thus 95% conf. int fo $\Theta$ is $(-0.205, -.174)$.

In practice if we wish to use B bootstrap samples to est. our confidence interval of size $(1-\alpha)100th$ we obtain the interval

$$\left( t^*_{(\alpha/2)}, t^*_{(1-\alpha/2)} \right)$$

where $t^*_{(\pi)}$ is the $[\pi B]^{th}$ order statistic of our B sized bootstrap samples $\hat{\theta}^*$.

## Sec 9.3.1.1  Justification for the percentile Method.

Consider a strictly increasing transformation $\phi$ and a distribution function H that is cont. and symmetric $(H(z) = 1 - H(-z))$. with the property that

(*) $$P\left[ h_{\alpha/2} \leq \phi(\hat{\theta}) - \phi(\theta) \leq h_{1-\alpha/2} \right] = 1-\alpha$$

(we do not need to know H! $\neq \phi$). where $h_\alpha$ is the $\alpha$ quantile of H.

&/ if $\phi$ is the normalizing variance stabilizing transform the H is standard Normal. (doesn't matter).

Apply the bootstrap principal to eq (*).

isolate $\hat{\theta}^*$

$$1-\alpha \approx P^*\left[ h_{\alpha/2} \leq \phi(\hat{\theta}^*) - \phi(\hat{\theta}) \leq h_{1-\alpha/2} \right]$$

$$= P^*\left[ h_{\alpha/2} + \phi(\hat{\theta}) \leq \phi(\hat{\theta}^*) \leq h_{1-\alpha/2} + \phi(\hat{\theta}) \right]$$

$$= P^*\left[ \phi^{-1}\left( h_{\alpha/2} + \phi(\hat{\theta}) \right) \leq \hat{\theta}^* \leq \phi^{-1}\left( h_{1-\alpha/2} + \phi(\hat{\theta}) \right) \right]$$

And so we have rewritten the expression in terms of the confidence interval of $\hat{\theta}^{a}$ and the bootstrap dist.

Since the bootstrap dist. is observed by us, its percentiles are known quantities. i.e.

we know $\mathcal{E}_\alpha$ the $\alpha$ quantile of the empirical dist of $\hat{\theta}^{*}$.

Hence we know

$$P^*\left[\mathcal{E}_{\alpha/2} \leq \hat{\theta}^* \leq \mathcal{Z}_{1-\alpha/2}\right] = 1-\alpha$$

and we can approx.

$$\phi^{-1}\left(h_{\alpha/2} + \phi(\hat{\theta})\right) \approx \mathcal{E}_{\alpha/2}$$
$$\phi^{-1}\left(h_{1-\alpha/2} + \phi(\hat{\theta})\right) \approx \mathcal{E}_{1-\alpha/2}.$$

Now, going back to the original eq $(*)$ we have

$$1-\alpha = P\left[h_{\alpha/2} - \phi(\hat{\theta}) \leq -\phi(\theta) \leq h_{1-\alpha/2} - \phi(\hat{\theta})\right]$$

$$= P\left[-h_{1-\alpha/2} + \phi(\hat{\theta}) \leq \phi(\theta) \leq -h_{\alpha/2} + \phi(\hat{\theta})\right]$$

$$= P\left[\phi^{-1}\left[-h_{1-\alpha/2} + \phi(\hat{\theta})\right] \leq \theta \leq \phi^{-1}\left(-h_{\alpha/2} + \phi(\hat{\theta})\right)\right]$$

$$= P\left[\phi^{-1}\left[h_{\alpha/2} + \phi(\hat{\theta})\right] \leq \theta \leq \phi^{-1}\left(h_{1-\alpha/2} + \phi(\hat{\theta})\right)\right]$$

since by symm. $h_{\alpha/2} = -h_{1-\alpha/2}$ $(H(z) = 1-H(-z)).$

$$\approx P\left[\mathcal{E}_{\alpha/2} \leq \theta \leq \mathcal{E}_{1-\alpha/2}\right]$$

and so the conf. limits happily coincide w/ those for $\hat{\theta}^*$. we read off the quantiles for $\hat{\theta}^*$ from the bootstrap dist as our confidence limits for $\theta$.

# Sec 9.3.2 Pivoting

Although simple, the percentile method is prone to bias & inaccurate coverage prob. To ensure best performance, the bootstrapped statistic should be approx. pivital. (that is its dist. should not depend on the unknown parameter $\theta$).

Ex/ if $g$ is our variance stabalizing transformation. then the variance of $g(\hat{\theta})$ is ind. of $\hat{\theta}$, and is a good pivot.

Sec 9.3.2 discusses several pivoting tech.

# Sec 9.3.2.2   The bootstrap t.

An approx. pivot that is quite easy to implement is given by the bootstrap t method, also called the studentized bootstrap.

Suppose we wish to est $\theta = T(F)$ by $\hat{\theta} = T(\hat{F})$. We can est. the variance of $\hat{\theta}$ by $V(\hat{F})$. Doing so, it is reasonable to hope that

$$R(X, F) = \frac{T(\hat{F}) - T(F)}{\sqrt{V(\hat{F})}} = \frac{\hat{\theta} - \theta}{\sqrt{V(\hat{F})}}$$

is roughly pivital.

And so we bootstrap $R(X, F)$ obtaining a collection of $R(X^*, \hat{F})$.

Let $\hat{G}$ be the dist of $R(X, F)$ and $G^*$ be the dist of $R(X^*, \hat{F})$.

unknown
observed.

Then by def a $1-\alpha$ confidence interval for $\Theta$ can be obtained from

$$1-\alpha = P\left[\xi_{\alpha/2}(\hat{G}) \leq R(X,F) \leq \xi_{1-\alpha/2}(\hat{G})\right]$$

$$= P\left[\xi_{\alpha/2}(\hat{G}) \leq \frac{T(\hat{F})-T(F)}{\sqrt{V(\hat{F})}} \leq \xi_{1-\alpha/2}(\hat{G})\right]$$

$$= P\left[\hat{\Theta} - \sqrt{V(\hat{F})}\,\xi_{1-\alpha/2}(\hat{G}) \leq \Theta \leq \hat{\Theta} - \sqrt{V(\hat{F})}\,\xi_{\alpha/2}\right]$$

where $\xi_{\alpha}^{(\hat{G})}$ is the $\alpha$ quantile of $\hat{G}$. However $F$ (and hence $\hat{G}$) is unknown.

So we use the bootstrap principle to imply that $\hat{G}$ is roughly equal to $\hat{G}^*$, $\therefore \xi_{\alpha}(\hat{G}) \approx \xi_{\alpha}(\hat{G}^*) \; \forall \alpha$. And since we can obtain $\xi_{\alpha}(\hat{G}^*)$ from the histogram of bootstrap values $R(X^*, \hat{F}^*)$ we have the bootstrap confidence interval for $\Theta$:

$$\left(T(\hat{F}) - \sqrt{V(\hat{F})}\,\xi_{1-\alpha/2}(\hat{G}^*), \; T(\hat{F}) - \sqrt{V(\hat{F})}\,\xi_{\alpha/2}(\hat{G}^*)\right)$$

Comment: these are percentiles of the tails & thus.
   B should be very large (several thousand).
   - Can use bootstrap variance est. for $V(\hat{F})$.

Bootstrap-t usually provide confidence interval coverage rates that closely approx. the nominal conf. level. They are most reliable when $T(\hat{F})$ is approx a location statistic → a constant shift in all the data values will induce the same shift in $T(\hat{F})$. Also, can be sensitive to the presence of outliers & should be used with caution in these cases. Also, it performs poorly when the underlying dist. is heavy tailed.

Se. 9.3.3 Bootstrap for hypothesis testing.
Pg 268

Conducting a hyp. test is closely related to est. conf int. The simplest approach for bootstrap hyp. testing is to base the p-value on a bootstrap conf. int. Specifically, consider a null hyp based on a para whose est. can be bootstrapped.
- Obtain conf. int by percentile or other method
- if the $(1-\alpha)100\%$ b.c.i does not cover the null value, reject w/ p-value no greater than $\alpha$.
- This is a simple method & works better if the bootstrap sampling is done in a manner that reflects the null hyp.

To illus. what we mean, consider a null hyp about a univariate para $\theta$ with null value $\theta_0$. Let the test statistic be $R(X, F) = \hat{\theta} - \theta_0$; the null hyp would reject whenever $|\hat{\theta} - \theta_0|$ is large compared to a reference dist.

A tempting method would be to gen the ref. dist by resampling values $R(X^*, F) = \hat{\theta}^* - \theta_0$ via bootstrap. However if if the null is false, then this stat. does not have the correct ref. dist. If $\theta_0$ is far from true $\theta$, then $|\hat{\theta} - \theta_0|$ will not seem big compared to $|\hat{\theta}^* - \theta_0|$.

Instead use $R(X^*, \hat{F}) = \hat{\theta}^* - \hat{\theta}$ to generate bootstrap estimate. Then if $\theta_0$ is far from $\hat{\theta}$, the values $|\hat{\theta}^* - \hat{\theta}|$ will be small compared to $|\hat{\theta} - \theta_0|$. Thus comparing $\hat{\theta} - \theta_0$ to $\hat{\theta}^* - \hat{\theta}$ yields greater statistical power.

Sec 9.4
Gentle 4.5   Reducing MC Error. (Variance/bias reduction).

MC $\overset{BS}{\wedge}$ estimators have two sources of variation
due to ① the initial sampling
        ② the bootstrap sampling.

We will discuss 2 methods to reduce variance.

— Jackknife after Bootstrap

Gentle.   First we need to est. the variance of the
          bootstrap estimator. One way to do this
          is to use the jackknife. The brute force
          way to do this is to
              — do n sep. bootstraps on the original
                 sample w/ a diff obs removed each time.
              — Then use jackknife discussed earlier.
          This is not comp efficient.

                         draw m bs samples and
Another procedure $\wedge$ (Efron 1992) stores the indices
of the original sample included in each bootstrap sample
in an n×m matrix

m bootstrap
samples
$$\begin{pmatrix} x_1 & x_2 \\ x_4 & x_3 \\ x_7 & x_9 \\ \vdots & \vdots \\ \vdots & \vdots \\ \vdots & \vdots \end{pmatrix}$$

Then for each of the m
samples not containing
$x_i$, treat it as though it
came from the orig. sample
w/ $x_i$ omitted.

Bootstrap samples obt. either way have the same dist.
Then use these samples to gen. Jackknife.

— Can have problems if $x_j$ in all samples, but large
   m (rel to n) has a low prob. of this.

Sec 9.4.1    Balanced Bootstrap.

Consider a bootstrap bias correction of the sample mean.
The bias correction should be zero, because
$\overline{X}$ is unbiased for the true mean $\mu$.
Now
$$R(X, F) = \overline{X} - \mu$$
will have bootstrapped values

$$R(X_j^*, \hat{F}) = \overline{X}_j^* - \overline{X} \qquad j = 1, \dots, B.$$ whose mean is our bias est.

Although $\overline{X}$ is unbiased it is highly unlikely that ~~This yields bias est~~ a random select of
the pseudo-data sets will produce a dist
$R(X^*, \hat{F})$ whose mean is exactly $0$.
This is due to ordinary M.C. variation.

However, if each data value occurs in
$\{X_1^*, \dots X_B^*\}$ w/ the same rel freq as it
does in $X$ then the bootstrap bias est

$$\frac{1}{B} \sum_{j=1}^{B} R(X_j^*, \hat{F}) \quad \text{must be zero.}$$

This is called balancing the bootstrap.

The easiest way to do this is to concatenate
B copies of the data set $X$, permute their values
and then read off B blocks n values at a time.
The jth block becomes $X_j^*$.

Other comments:
Sec 9.5   Fairly common idea Bootstrap aggregating
(Bagging). Basically replace $\hat{\theta}$ w/ $\widetilde{\theta}$
$\widetilde{\theta}^* = \frac{1}{B} \sum_{j=1}^{B} \hat{\theta}_j^* \rightarrow$ model averaging. reduces var due to small chan in data.