

## Example 1

The data below show the sugar content (as a % of weight) of several national brands of children's and adults' cereals.

$\bar{X} = \text{sugar content in a typical children's cereal}$  —  $E[\bar{X}]$   
 $\bar{Y} = \text{sugar content in a typical adult's cereal}$  —  $E[\bar{Y}]$

**Children's cereals:** 40.3, 55, 45.7, 43.3, 50.3, 45.9, 53.5, 43, 44.2, 44, 47.4, 44, 33.6, 55.1, 48.8, 50.4, 37.8, 60.3, 46.5  
**Adults' cereals:** 20, 30.2, 2.2, 7.5, 4.4, 22.2, 16.6, 14.5, 21.4, 3.3, 6.6, 7.8, 10.6, 16.2, 14.5, 4.1, 15.8, 4.1, 2.4, 3.5, 8.5, 10, 1, 4.4, 1.3, 8.1, 4.7, 18.4

- (a) Is it reasonable to assume that each sample comes from a normal distribution? No, may be ok for children's cereals but not for adults' cereals
- (b) Can the variances of the two distributions be assumed to be equal? Justify your answer.

- (c) Compute an appropriate 95% confidence interval for difference in mean sugar contents of the two cereal types. What assumptions did you make, if any, to construct the CI?

$$95\% \text{ CI for } \mu_{x-y} = [32.5, 40.8]$$

Plausible values

Clearly: ① Since the value of  $\bar{x}$  is more, we can conclude that  $\mu_x > \mu_y$ .

② plausible values for  $\mu_x - \mu_y$ :  
 $32.5 \text{ to } 40.8$

Assumption: normality which may not justified, esp.  
for the second sample. But given that the diff. for the two dist. is so close, the conclusions are still valid.

- (d) What do you conclude on the basis of your answer in (c)?
- Can we say that children's cereals have more sugar on average than adult cereals? If yes, by how much? Justify your answers.

Yes,  
see previous prob.  
the proportion of  
sugar content in  
children's cereals is  
given by  
Mr. Vade  
& we  
will  
find out  
the  
proportion  
of sugar  
in  
adult  
cereals.

## Example 3

Consider the dataset stored in the file bp.txt. This dataset contains one measurement of systolic blood pressure (in mmHg) made by each of two methods — a finger method and an arm method — from the same 200 patients.

- (a) Perform an exploratory analysis of the data by examining the distributions of the measurements from the two methods using boxplots. Comment on what you see. Do the two distributions seem similar? Justify your answer.

Now in class  
/

see & understand  
the 2 methods  
for contrasts

- (b) Use histograms and QQ plots to examine the shapes of the two distributions. Comment on what you see. Does the assumption of normality seem reasonable? Justify your answer.

I discussed in my

- (c) Construct an appropriate 95% confidence interval for the difference in the means of the two methods. Interpret your results. Can we conclude that the two methods have identical means? Justify your answer. What assumptions, if any, did you make to construct the interval? Do the assumptions seem to hold?

$$95\% \text{ CI for } \mu_x - \mu_y: [-6.3, -2.3]$$

arm-finger

$\Rightarrow$  Since the entire interval is below zero, we can conclude that  $\mu_x - \mu_y < 0 \Rightarrow \mu_x < \mu_y$  and also plausible values for  $\mu_x - \mu_y$  is below zero.

Since the entire interval is below zero, we can conclude that  $\mu_x - \mu_y < 0 \Rightarrow \mu_x < \mu_y$  and also plausible values for  $\mu_x - \mu_y$  is below zero.

Normality for normality seemed sufficient since no sample size is not needed.

# Bootstrap

$$F(x) = P[X \leq x]$$

**Set up:** Data  $X_1, \dots, X_n$  — i.i.d. as  $X$  with population cdf  $F$  (which is not completely known). — from  $m$  sample problems.

**Parameter of interest:**  $\theta$ , estimated by  $\hat{\theta}$

**Examples:** Mean, variance, median, quantiles, etc.

**Issue:** Need to get the *sampling distribution* of  $\hat{\theta}$  so that we can compute, e.g., standard error of  $\hat{\theta}$ , or confidence interval for  $\theta$ ?

**Q:** Why not use the methods that we have learnt?

$$\frac{\hat{\theta}}{\bar{X}}$$

### Sampling distribution

- $N\left[\mu, \frac{\sigma^2}{n}\right]$  if pop. is normal
- $\approx N\left[\mu, \frac{\sigma^2}{n}\right]$  if  $n$  is large
- $\approx N\left[\theta, \frac{1}{n}\right]$  if  $n$  is large.

$$\hat{\theta}_{MLE}$$

- Sample median }
  - sample variance }
    - sample quantiles }
  - complicated.

Basics

**Bootstrap:** A simulation based technique that allows us to approximate the sampling distribution of  $\hat{\theta}$ . Assumes large  $n$ , but its value needed for validity of bootstrap is typically less than that for the usual large-sample procedure.

Original sample:  $X_1, \dots, X_n \sim$  i.i.d. with cdf  $F$

Bootstrap (re)sample:  $X_1^*, \dots, X_n^* \sim$  i.i.d. with cdf  $\hat{F}$ , where  
 $\hat{F} = \text{estimated cdf (which is completely known)}$

## Parametric bootstrap:

- Functional form of  $F$  is known (e.g., normal), but  $F$  may depend on unknown parameter  $\theta$ .
  - $\hat{F}$  is same as  $F$  but with  $\theta$  replaced by its MLE,  $\hat{\theta}$ . In other words,  $\hat{F}$  is the cdf of the fitted model.
  - Ex:  $F = N(\mu, \sigma^2)$ ,  $\hat{F} = N[\bar{X}, \frac{s^2}{n}]$
  - Often easy to simulate i.i.d. draws  $X_1^*, \dots, X_n^*$  from  $\hat{F}$ .

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

## Nonparametric bootstrap:

- Functional form of  $F$  is unknown.

- $\hat{F}$  = empirical cdf, where

Suppose:  $x_1, x_2, \dots, x_n$  — obs. data  
 $\hat{F}(x)$  is a step function changing at  $x_1, x_2, \dots, x_n$  with step size  $= \frac{1}{n}$

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x)$$

- Think of  $\hat{F}$  as a discrete distribution that assigns  $1/n$  probability to each of the sample observations,  $X_1, \dots, X_n$ .

- Get  $X_1^*, \dots, X_n^*$  by sampling  $n$  times with replacement from  $X_1, \dots, X_n$ .

✓ consider a rv  $X^*$  that takes  $x_1, \dots, x_n$  as possible values, look with equal prbs.  $\frac{1}{n}$  — CDF of  $X^*$  is  $\hat{F}$ .  
 Need to get  $x_1^*, x_2^*, \dots, x_n^*$  as a random sample  
 from  $F$  — use sample fn. in R.

$$I(A) = \begin{cases} 1 & \text{if } A \text{ occurs} \\ 0 & \text{otherwise} \end{cases}$$

## Bootstrap distribution of estimator $\hat{\theta}$ of $\theta$ :

Original sample:  $X_1, \dots, X_n$  — gives  $\hat{\theta}$

- Simulate a large number  $b$  of *bootstrap resamples*, and compute  $\hat{\theta}^*$  from each resample. *exactly the way we compute  $\hat{\theta}$  from the original sample*

Sample	Est.	Resample #	$\hat{\theta}^*$
$X_1, X_2, \dots, X_n$	$\hat{\theta}$		
			$\hat{\theta}_1^*, \hat{\theta}_2^*, \dots, \hat{\theta}_n^*$
			$\hat{\theta}_1^*, \hat{\theta}_2^*, \dots, \hat{\theta}_n^*$
			$\hat{\theta}_1^*, \hat{\theta}_2^*, \dots, \hat{\theta}_n^*$

b draws.

- This process gives a large number of draws,  $\hat{\theta}_1^*, \dots, \hat{\theta}_b^*$ .
- These draws are coming from the *bootstrap distribution* of  $\hat{\theta}$ .
- How to see this distribution?  
Histogram, box plot etc.
- It approximates the *sampling distribution* of  $\hat{\theta}$ .
- Use the draws  $\hat{\theta}_1^*, \dots, \hat{\theta}_b^*$  to estimate features of sampling distribution of  $\hat{\theta}$  that may be of interest. — *especially, writing a Monte Carlo method to approximate features of a dist.*

## Estimating a feature $\eta$ of distribution of $\hat{\theta}$ :

- Get a large number  $b$  of draws,  $\hat{\theta}_1^*, \dots, \hat{\theta}_b^*$ .

- $\hat{\eta}^*$  = same feature computed from these draws.

This is Monte Carlo method

$$\text{Ex 1: } \eta = E(\hat{\theta}). \quad \hat{\eta}^* = \left( \frac{1}{b} \sum_{k=1}^b \hat{\theta}_k^* \right) = \bar{\hat{\theta}}^*$$

$$\text{Ex 2: } \eta = \text{var}(\hat{\theta}). \quad \hat{\eta}^* = \frac{1}{b-1} \sum_{k=1}^b (\hat{\theta}_k^* - \bar{\hat{\theta}}^*)^2 \Rightarrow \text{SE}(\hat{\theta}) = SD \text{ of the draws}$$

$$\text{Ex 3: } \eta = \text{bias of } \hat{\theta} = E(\hat{\theta}) - \theta. \quad \hat{\eta}^* = \bar{\hat{\theta}}^* - \theta$$

Ex 4:  $\eta = \alpha\text{-th quantile of } \hat{\theta}$ .  $\hat{\eta}^* = \alpha\text{-th sample quantile of the draws } \hat{\theta}_1^*, \dots, \hat{\theta}_b^*$

$$\text{Ex 5: } \eta = \alpha\text{-th quantile of } (\hat{\theta} - \theta). \quad \hat{\eta}^* = \hat{\theta}_{(b\alpha)} \rightarrow \text{will use: } \hat{\theta}_{((b+1)\alpha)}$$

$$\Rightarrow \hat{\eta}^* = \hat{\theta}_{(b\alpha)} - \theta$$

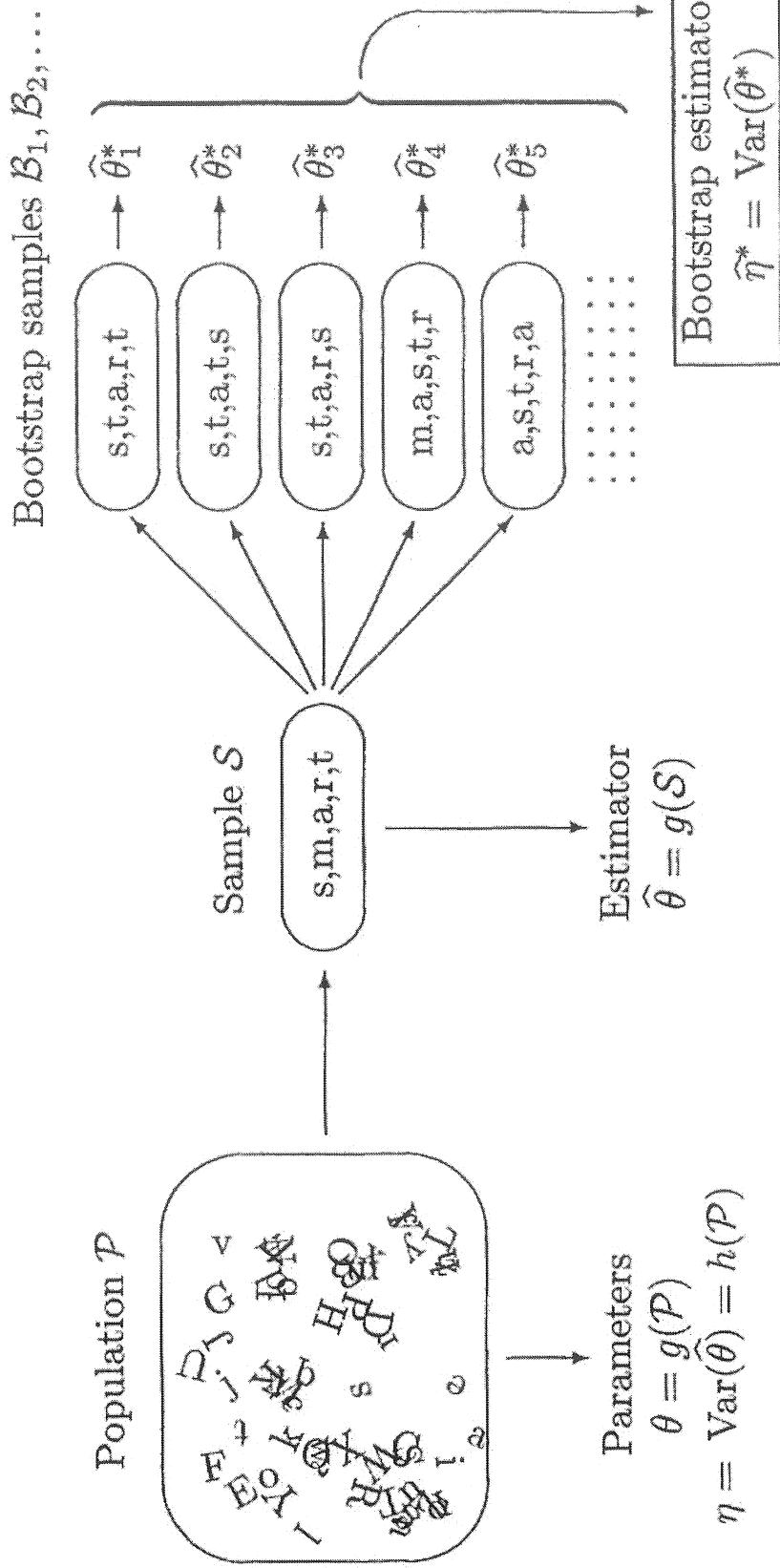


FIGURE 10.6: Bootstrap procedure estimates  $\eta = \text{Var}(\hat{\theta})$  by the variance of  $\hat{\theta}_i^*$ 's, obtained from bootstrap samples.