

## What have we learned until now: Part 1

- Direct simulation from probability densities
- Several different methods
  - ▶ inversion sampling
  - ▶ methods based on relationship between RV
  - ▶ rejection sampling (\*)
  - ▶ importance sampling (\*)
- We get samples that are independent of each other

## What have we learned until now: Part 1 and 2

### Bayesian paradigm

- Likelihood  $\pi(y|x)$
- Prior  $\pi(x)$
- Posterior  $\pi(x|y) \propto \pi(y|x)\pi(x)$

## What have we learned until now: Part 2

Hierarchical models are an extremely useful tool in Bayesian model building.

### Three parts:

- **Observation model  $y|x$ :** Encodes information about observed data.
- **The latent model  $x|\theta$ :** The unobserved process.
- **Hyperpriors for  $\theta$ :** Models for all of the parameters in the observation and latent processes.

Note: here we indicate the observed data by  $y$  while  $x$  and  $\theta$  are parameters

## What have we learned until now: Part 2

### MCMC algorithm:

- **Problem:** Sample from  $\pi(x)$ ,  $x \in S$ .
- **MCMC idea:**
  - ▶ Construct **Markov chain with  $\pi(x)$  as limiting distribution.**
  - ▶ Simulate the Markov chain for a long time so that it has time to converge.
  - ▶ **Most MCMC samplers are based on reversible Markov chains**  
 $\Rightarrow$  Their convergence is proved by checking the detailed balance equation.
- Can be applied to virtually any bayesian model
- Convergence and slow mixing can be a big issue

## What have we learned until now: Part 2

Integrated nested Laplace approximation:

- Can be applied only on a (large) class of models: Latent Gaussian models
- No sampling involved, based on numerical approximation
- The focus is on posterior marginals

TMA4300 - Part 3

Last part of this course

TMA4300 - Part 3

⇒ Not closely related to the two first parts

- ▶ no more MCMC
- ▶ mostly non-Bayesian perspective

⇒ Two topics (not closely related to each other):

- ▶ Bootstrapping
- ▶ Expectation-Maximization algorithm

\*

---

\*Slides are partially based on lecture notes kindly provided by Håkon Tjelmeland, Andrea Riebler, and Sara Martino.

## Bootstrap

## Bootstrap



[http://tradingconsequences.blogs.edina.ac.uk/files/2013/10/Dr\\_Martens\\_black\\_old.jpg](http://tradingconsequences.blogs.edina.ac.uk/files/2013/10/Dr_Martens_black_old.jpg)

### An example for introduction

| Group       | Survival Time                      | Sample size | Mean  | Estimated SE |
|-------------|------------------------------------|-------------|-------|--------------|
| Treatment   | 94,197,16,38<br>99,141,23          | 7           | 86.86 |              |
| Control     | 52,104,146,10,51,46<br>30,40,27,46 | 9           | 56.22 |              |
| Differance: |                                    |             | 30.63 |              |

- Is the difference in mean significant?
- 

### An example for introduction

| Group       | Survival Time                      | Sample size | Mean  | Estimated SE |
|-------------|------------------------------------|-------------|-------|--------------|
| Treatment   | 94,197,16,38<br>99,141,23          | 7           | 86.86 | 25.24        |
| Control     | 52,104,146,10,51,46<br>30,40,27,46 | 9           | 56.22 | 14.14        |
| Differance: |                                    |             | 30.63 | 30.93        |

- Is the difference in mean significant?
-

## An example for introduction

| Group       | Survival Time                      | Sample size | Mean  | Estimated SE |
|-------------|------------------------------------|-------------|-------|--------------|
| Treatment   | 94,197,16,38<br>99,141,23          | 7           | 86.86 | 25.24        |
| Control     | 52,104,146,10,51,46<br>30,40,27,46 | 9           | 56.22 | 14.14        |
| Differance: |                                    |             | 30.63 | 30.93        |

- Is the difference in mean significant?
  - What if we want to compare the medians instead?
- Show code Bootstrap\_intro.R

## The bootstrap

- Bootstrap is a computer-based technique for doing statistical inference (usually with a minimum of assumptions)
- It is not Bayesian

## ... pull oneself up by one's bootstraps

*To begin an enterprise or recover from a setback without any outside help; to succeed only on one's own effort or abilities.*

Wiktionary

The term is sometimes attributed to Rudolf Erich Raspe's story "The Surprising Adventures of Baron Munchausen", where the main character pulls himself (and his horse) out of a swamp by his hair



[http://redstateeclectic.typepad.com/redstate\\_commentary/2010/11/sustainability-isnt-sustainable.html](http://redstateeclectic.typepad.com/redstate_commentary/2010/11/sustainability-isnt-sustainable.html)

## Important concepts

- empirical distribution function
- plug in estimator
- bootstrap sample

## Today's lecture

- Bootstrap
  - ▶ Non-parametric
  - ▶ Parametric
- Bootstrap estimate of SD
- Bootstrap estimate of bias

## Plug in estimator

Let  $\theta$  be an interesting feature of  $F$ ,  $\theta = t(F)$ .

For example:

$$\theta = E(X) = \int xf(x)dx$$

$$\theta = \text{Var}(X) = \int (x - E(X))^2 f(x)dx$$

The **plug-in estimator** for  $\theta$  is defined by:

$$\hat{\theta} = t(\hat{F})$$

The plug-in principle is quite good, if the only information about  $F$ , comes from the sample  $x$ .

## Empirical distribution function

Assume we have **iid** observations from an (unknown) distribution  $F$ :

$$F \rightarrow (x_1, \dots, x_n)$$

The **empirical distribution function**  $\hat{F}$  is the CDF that puts mass  $1/n$  at each data point  $x_i$ :

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

where  $1(\cdot)$  denotes the indicator function.

For iid samples  $\hat{F}$  is a sufficient estimator for  $F$ .

## Examples

Thus

$$\theta = E(X) \Rightarrow \hat{\theta} = E_{\hat{F}}(X) = \sum_{i=1}^n x_i \frac{1}{n} = \bar{x}$$

## Examples

Thus

$$\theta = E(X) \Rightarrow \hat{\theta} = E_{\hat{F}}(X) = \sum_{i=1}^n x_i \frac{1}{n} = \bar{x}$$

$$\begin{aligned}\theta = \text{Var}(X) &\Rightarrow \hat{\theta} = \text{Var}_{\hat{F}}(X) = E_{\hat{F}}[(X - E_{\hat{F}}(X))^2] \\ &= \sum_{i=1}^n (x_i - E_{\hat{F}}(X))^2 \frac{1}{n} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\end{aligned}$$

## Examples

Thus

$$\theta = E(X) \Rightarrow \hat{\theta} = E_{\hat{F}}(X) = \sum_{i=1}^n x_i \frac{1}{n} = \bar{x}$$

$$\begin{aligned}\theta = \text{Var}(X) &\Rightarrow \hat{\theta} = \text{Var}_{\hat{F}}(X) = E_{\hat{F}}[(X - E_{\hat{F}}(X))^2] \\ &= \sum_{i=1}^n (x_i - E_{\hat{F}}(X))^2 \frac{1}{n} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\end{aligned}$$

$$\begin{aligned}\theta = \text{SD}(X) &\Rightarrow \hat{\theta} = \text{SD}_{\hat{F}}(X) = \sqrt{\text{Var}_{\hat{F}}(X)} \\ &= \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}\end{aligned}$$

## Example

$$E_F \left[ \left( \frac{X - \mu}{\sigma} \right)^3 \right] = \frac{E[X^3] - 3\mu\sigma^2 - \mu^3}{\sigma^3}$$

See notes

## Setting

Assume we have :

$$F \rightarrow (x_1, \dots, x_n)$$

Thus  $\hat{F}$  gives mass  $\frac{1}{n}$  to each observed value.

A **bootstrap sample** is defined to be a random sample of size  $n$  from  $\hat{F}$ , say  $\mathbf{x}^* = (x_1^*, \dots, x_n^*)$

$$\hat{F} \rightarrow (x_1^*, \dots, x_n^*)$$

## Simple illustration

Suppose  $n = 3$  univariate data points, namely

$$\{x_1, x_2, x_3\} = \{1, 2, 6\}$$

are observed as an iid sample from  $F$  that has mean  $\theta$ . At each observed data value,  $\hat{F}$  places mass  $1/3$ . Suppose the estimator to be bootstrapped is the sample mean  $\hat{\theta}$ .

There are  $3^3 = 27$  possible outcomes for  $\mathcal{X}^* = \{X_1^*, X_2^*, X_3^*\}$ .

## Simple illustration (II)

| $\mathcal{X}^*$ | $\hat{\theta}^*$ | $P^*(\hat{\theta}^*)$ | Observed frequency |
|-----------------|------------------|-----------------------|--------------------|
| 1 1 1           | 3/3              | 1/27                  | 36/1000            |
| 1 1 2           | 4/3              | 3/27                  | 101/1000           |
| 1 2 2           | 5/3              | 3/27                  | 123/1000           |
| 2 2 2           | 6/3              | 1/27                  | 25/1000            |
| 1 1 6           | 8/3              | 3/27                  | 104/1000           |
| 1 2 6           | 9/3              | 6/27                  | 227/1000           |
| 2 2 6           | 10/3             | 3/27                  | 131/1000           |
| 1 6 6           | 13/3             | 3/27                  | 111/1000           |
| 2 6 6           | 14/3             | 3/27                  | 102/1000           |
| 6 6 6           | 18/3             | 1/27                  | 40/1000            |

## Bootstrap estimate for standard error

- Parameter of interest:  $\theta = t(F)$
- Our estimator for  $\theta$ :  $\hat{\theta} = s(x)$
- Want (to estimate)  $SD_F(\hat{\theta})$ .

A bootstrap replication of  $\hat{\theta}$  is

$$\hat{\theta}^* = s(x^*)$$

Use plug-in principle to estimate  $SD_F(\hat{\theta})$ .

The bootstrap estimate of the standard error of  $\hat{\theta} = s(x)$  is  $SD_{\hat{F}}(\hat{\theta}^*)$ .

This is called the ideal bootstrap estimate of standard error of  $\hat{\theta}$ .

## Ideal bootstrap estimate of standard error

- For the sample mean it can be computed analytically
- For (very) small sample sizes it can be computed using all the possible bootstrap replicates. (Number of possible bootstrap sample:  $n^n$ .)
- In other cases it can be approximated via Monte Carlo techniques

## Computational way of obtaining a good estimate

We can estimate  $SD_{\hat{F}}(\hat{\theta}^*)$  by simulation:

1. Generate **B bootstrap samples**  $x^{1*}, \dots, x^{B*}$ .
2. Evaluate the corresponding parameter estimates

$$\hat{\theta}^*(b) = s(x^{b*}), \quad b = 1, 2, \dots, B$$

3. Estimate  $SD_{\hat{F}}(\hat{\theta}^*)$  by

$$\widehat{SE}_B = \sqrt{\frac{\sum_{b=1}^B (\hat{\theta}^*(b) - \hat{\theta}^*(\cdot))^2}{B-1}}$$

where

$$\hat{\theta}^*(\cdot) = \frac{1}{B} \sum_{b=1}^B \hat{\theta}^*(b)$$

Note

$$\lim_{B \rightarrow \infty} \widehat{SE}_B = \widehat{SE}_{\infty} = \widehat{SD}_{\hat{F}}(\hat{\theta}^*)$$

## Example

Setting

$$\theta = E(X)$$

$$\hat{\theta} = s(x) = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$

$$\hat{\theta}^* = s(x^*) = \frac{1}{n} \sum_{i=1}^n x_i^* = \bar{x}^*$$

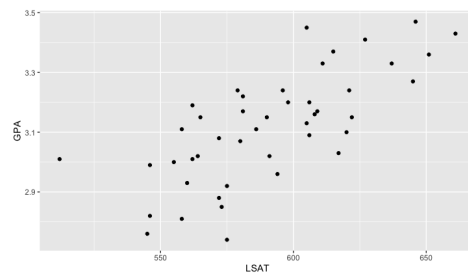
Here, the ideal bootstrap estimate exists

see blackboard

## Example: The correlation coefficient

Scores for 15 law schools in the USA

$$y_i = (LSAT_i, GPA_i), \quad i = 1, \dots, 15$$



The correlation between the two scores is estimated to be 0.78, but what is its standard error?

## Example: The correlation coefficient

- 1000 bootstrap replicates

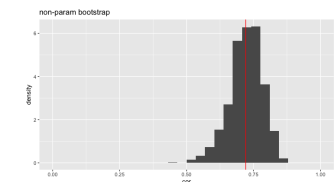
$$y^{1*}, \dots, y^{1000*}$$

- For each replicates compute

$$\hat{\theta}^{i*} = s(y^{i*})$$

- Estimate bootstrap SE

$$\widehat{SD}_{\hat{F}}(\theta) = 0.121$$





## How large do we need $B$ ?

Intuitively we understand that the  $\widehat{SE}_B$  has larger standard deviation than  $\widehat{SE}_\infty$ .

Theory, not to be discussed here, gives the following rules of thumb:

1. Even a small  $B$  is informative, say  $B = 25$  or  $B = 50$  is often enough to get a good estimate of  $SE_F(\hat{\theta})$ .
2. Very seldomly more than  $B = 200$  is necessary to estimate  $SE_F(\hat{\theta})$ .

## Again ...

... we can/must estimate  $SD_{\hat{F}_{\text{par}}}(\hat{\theta}^*)$  by simulation:

1. Generate  $B$  bootstrap samples  $x^{1*}, \dots, x^{B*}$ , where

$$x^{b*} = (x_1^{b*}, \dots, x_n^{b*})$$

with  $x_1^{b*}, \dots, x_n^{b*} \stackrel{\text{iid}}{\sim} \hat{F}_{\text{par}}$ .

2. Evaluate the corresponding parameter estimates

$$\hat{\theta}^*(b) = s(x^{b*}), \quad b = 1, 2, \dots, B$$

3. Estimate  $SD_{\hat{F}_{\text{par}}}(\hat{\theta}^*)$  by

$$\widehat{SE}_B = \sqrt{\frac{\sum_{b=1}^B (\hat{\theta}^*(b) - \hat{\theta}^*(\cdot))^2}{B-1}}$$

where

$$\hat{\theta}^*(\cdot) = \frac{1}{B} \sum_{b=1}^B \hat{\theta}^*(b)$$

## The parametric bootstrap

When data are modeled to originate from a parametric distribution, so

$$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} F(x, \xi),$$

another estimate of  $F$  may be employed.

Suppose that the observed data are used to estimate  $\xi$  by  $\hat{\xi}$ . Then each **parametric bootstrap** pseudo-dataset  $\mathcal{X}^*$  can be generated by drawing  $X_1^*, \dots, X_n^* \stackrel{\text{iid}}{\sim} F(x, \hat{\xi}) = \hat{F}_{\text{par}}$ .

## Example: Correlation coefficients

We assume now that

$$y_i = (LSAT_i, GPA_i) \sim \mathcal{N}(\mu, \Sigma), \text{ i.i.d}$$

where  $\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix}$  Estimate  $\mu$  and  $\Sigma$  and obtain:

$$\hat{F}_{(\hat{\mu}, \hat{\Sigma})}$$

## Example: The correlation coefficient

- 1000 bootstrap replicates

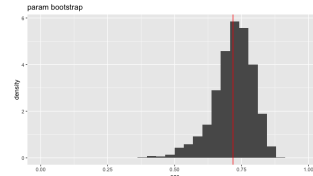
$$y^{1*}, \dots, y^{1000*} \sim \hat{F}_{(\hat{\mu}, \hat{\Sigma})}$$

- For each replicates compute

$$\hat{\theta}^{i*} = s(y^{i*})$$

- Estimate bootstrap SE

$$\hat{SD}_{\hat{F}}(\theta)$$



## Bootstrapping regression

Consider the ordinary multiple regression model

$$Y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \epsilon_i, \quad \text{for } i = 1, \dots, n,$$

where  $\epsilon_i$  are iid mean zero random variables with constant variance.

- Parameters of interest  $\boldsymbol{\beta}$
- Want to estimate  $SD(\hat{\boldsymbol{\beta}})$

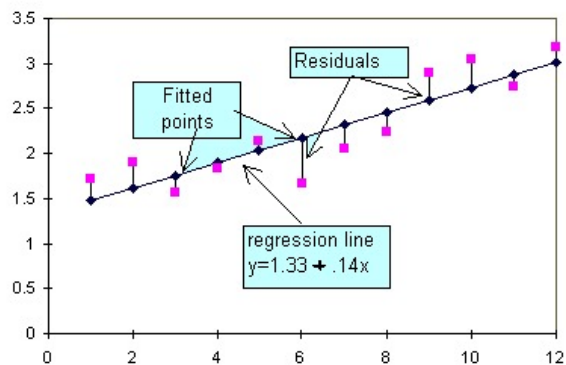
## Review: Linear Regression

- Least square estimate of  $\boldsymbol{\beta}$

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin}\left\{\sum (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2\right\} \Rightarrow \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

- Residuals

$$e_i = Y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}$$



## Bootstrap regression

Alternative 1: Bootstrap the residuals  $e_i = Y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}$

Alternative 2: Bootstrap the pairs  $\mathbf{Z}_i = (\mathbf{X}_i, Y_i)$

## Bootstrap the residuals

1. Fit the regression model to the observed data and obtain the fitted responses  $\hat{y}_i$  and residuals  $\hat{\epsilon}_i$ .
2. Sample a bootstrap set of residuals  $\hat{\epsilon}_1^*, \dots, \hat{\epsilon}_n^*$  from the set of fitted residuals completely at random and with replacement.
3. Generate a bootstrap set of pseudo responses

$$Y_i^* = \hat{y}_i + \hat{\epsilon}_i^*, \quad \text{for } i = 1, \dots, n.$$

4. Regress  $Y^*$  on  $\mathbf{x}$  to obtain a bootstrap estimate  $\hat{\beta}^*$ .

Repeat this process to get an empirical distribution of  $\hat{\beta}^*$ .

## Bootstrap the pair $\mathbf{Z}_i = (\mathbf{X}_i, Y_i)$

Suppose response and predictors are measured from a collection of individuals selected at random

⇒ Data pairs  $\mathbf{z}_i = (\mathbf{x}_i, y_i)$  can be regarded as iid realisation from  $\mathbf{Z}_i = (\mathbf{X}_i, Y_i)$  drawn from a joint response-predictor distribution.

Bootstrap:

- Sample  $\mathbf{Z}_1^*, \dots, \mathbf{Z}_n^*$  completely at random with replacement from  $\mathbf{z}_1, \dots, \mathbf{z}_n$ .
- Apply regression model on pseudo dataset to get  $\hat{\beta}^*$ .

Repeat this approach many times.

Note: Paired bootstrap is less sensitive to violation of assumptions, e.g. adequacy of regression model, than bootstrapping the residuals.

## Bootstrapping residuals: Remarks

This approach is also used for autoregressive models, for example.

Note: Bootstrapping the residuals is reliant on

- The model provides an appropriate fit
- The residuals have a constant variance

Otherwise, a different scheme is recommended.

Comment: No need to bootstrap for linear regression model with least squares estimation, as analytical results are then available.

## Copper-nickel alloy

Data: 13 measurements of corrosion loss ( $y_i$ ) in copper-nickel alloys, each with a specific iron content ( $x_i$ ).

Question: Change in corrosion loss in the alloys as the iron content increases, relative to corrosion loss where there is no iron, i.e.

$$\theta = \beta_1 / \beta_0.$$

|       |       |       |       |       |       |       |       |
|-------|-------|-------|-------|-------|-------|-------|-------|
| $x_i$ | 0.01  | 0.48  | 0.71  | 0.95  | 1.19  | 0.01  | 0.48  |
| $y_i$ | 127.6 | 124.0 | 110.8 | 103.9 | 101.5 | 130.1 | 122.0 |
| $x_i$ | 1.44  | 0.71  | 1.96  | 0.01  | 1.44  | 1.96  |       |
| $y_i$ | 92.3  | 113.1 | 83.7  | 128.0 | 91.4  | 86.2  |       |

The observed data yield  $\hat{\theta} = \hat{\beta}_1 / \hat{\beta}_0 = -0.185$ .

## Bias of an estimator

- We observe  $X_1, X_2, \dots, X_n \sim F$  iid
- Parameter of interest  $\theta = t(F)$
- Estimator  $\hat{\theta} = s(X)$   
(may or may not be based on the plug-in principle)
- Bias definition

$$\text{bias}_F(\hat{\theta}, \theta) = E_F[\hat{\theta}] - \theta = E_F[s(\mathbf{x})] - t(F)$$

## Bootstrap estimate of bias

We want to estimate

$$\text{bias}_F(\hat{\theta}, \theta) = E_F[s(\mathbf{x})] - t(F)$$

Idea: Apply the plug-in principle and define the bootstrap estimate of bias as:

$$\text{bias}_{\hat{F}} = E_{\hat{F}}[s(\mathbf{x}^*)] - t(\hat{F})$$

where  $\hat{F}$  is an estimate of  $F$  (for example the empirical distribution)

## Bias estimate of the bias

1. Generate  **$B$  bootstrap samples**  $x^{1*}, \dots, x^{B*}$ .
2. Evaluate the corresponding parameter estimates

$$\hat{\theta}^*(b) = s(x^{b*}), \quad b = 1, 2, \dots, B$$

3. Approximate the bootstrap expectation  $E_{\hat{F}}[s(\mathbf{x}^*)]$  as:

$$\hat{\theta}^*(\cdot) = \frac{1}{B} \sum_{b=1}^B \hat{\theta}^*(b)$$

4. Approximate the ideal bootstrap estimate for bias as

$$\widehat{\text{bias}}_B = \hat{\theta}^*(\cdot) - t(\hat{F})$$

## Bias corrected estimate

Once we have estimated the bias we can compute the bias-corrected estimator

$$\hat{\theta}_c = \hat{\theta} - \widehat{\text{bias}}_B = \hat{\theta} - [\hat{\theta}^*(\cdot) - t(\hat{F})]$$

## Bias corrected estimate

Once we have estimated the bias we can compute the bias-corrected estimator

$$\hat{\theta}_c = \hat{\theta} - \widehat{\text{bias}}_B = \hat{\theta} - [\hat{\theta}^*(\cdot) - t(\hat{F})]$$

Note: Bias correction will not always give an improved estimator.

We have that  $\text{Var}(\hat{\theta}_c) \geq \text{Var}(\hat{\theta})$  so if the bias is small is better not to do bias correction.

## Bootstrap bias correction

### Copper-nickel alloy example

The mean value of

$$\hat{\theta}^* - \hat{\theta}$$

among the pseudo datasets is about  $-0.00125$ .

The **bias-corrected bootstrap estimate** of  $\beta_1/\beta_0$  is  $-0.18507 - (-0.00125) = -0.184$ .

## Confidence intervals (percentile method)

A “simple-minded” two-sided confidence interval with coverage  $(1 - \alpha)$  for a parameter  $\alpha$  is given by

$$[q_{\alpha/2}^*, q_{1-\alpha/2}^*]$$

where  $q_{\alpha}^*$  is the  $\alpha$ -bootstrap quantile in the distribution of  $\hat{\theta}^*$ .

**Experience:** Often good, but often **too low coverage**, i.e the true  $\alpha$  for the interval is lower than the specified value.

**Note:** Better bootstrap confidence intervals exist and often have better coverage accuracy — at the price of being somewhat more difficult to implement

Show R-code `bootstrap_regression.R`