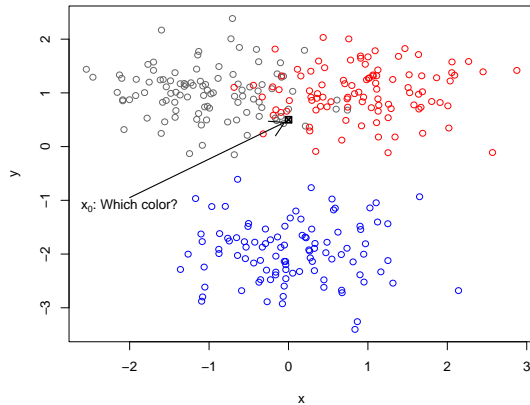


Review: Classification problem

Situation: Have observations $(x_1, y_1), \dots, (x_n, y_n)$, where $y_i \in \{0, 1, \dots, J-1\}$ gives a class. Have new observation x_0 , want to **predict the corresponding class y_0** .



Review

We have also discussed:

- **k-nearest neighbour algorithm** with tuning parameter k .
- Evaluation of classification rules

▶ **Misclassification rate:**

$$P(y_0 \neq \hat{y}(x_0; (x_1, y_1), \dots, (x_n, y_n)))$$

▶ **Apparent error rate:** (–too optimistic)

$$\frac{1}{n} \sum_{i=1}^n 1(y_i \neq \hat{y}(x_i; (x_1, y_1), \dots, (x_n, y_n)))$$

▶ Split training data into a training set and a test set.
(–too pessimistic)

Review: Model

We have: $p_j = P(Y = j)$, $f(x|y = j) = f_j(x)$

•

$$\pi_j(x_0) = P(Y_0 = j|x_0) = \frac{p_j f_j(x_0)}{\sum_{i=0}^{J-1} p_i f_i(x_0)}$$

•

$$ECM(i) = E(c(i|Y)|x_0) = \frac{\sum_{j=0}^{J-1} c(i|j) p_j f_j(x_0)}{\sum_{i=0}^{J-1} p_i f_i(x_0)}$$

•

$$\hat{y}_0 = \operatorname{argmin}_j ECM(i)$$

$$\stackrel{0/1\text{-loss}}{=} \operatorname{argmax}_j \{p_i f_i(x_0)\}$$

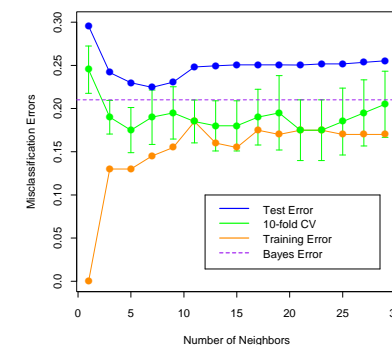
⇒

$$x|y = j \sim \mathcal{N}(\mu_j, \Sigma_j) \begin{cases} \text{LDA} & \Sigma_0 = \dots = \Sigma_{J-1} \\ \text{QDA} & \text{different } \Sigma_j \end{cases}$$

Review: Cross-validation (CV)

- **Leave-one-out cross validation**
- **K-fold cross validation**, where $K = 5$ or $K = 10$ is often used.

Determining k in knn-algorithm using CV:



Standard errors for cross validation

See blackboard

Bootstrap



http://tradingconsequences.blogs.edina.ac.uk/files/2013/10/Dr_Martens_black_old.jpg

Bootstrap Bill Turner



“Bootstrap” Bill Turner from Pirates of the Caribbean.

... Barbossa tied Bill to a cannon by his bootstraps and sent him to the bottom of the sea.

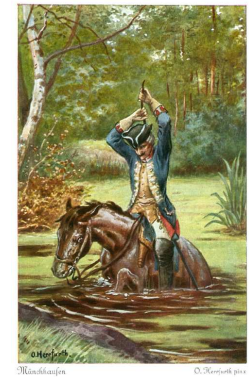
<http://kidstvmovies.about.com/od/piratesofthecaribbean3/ig/Pirates-At-World-s-End/-Bootstrap-Bill.htm>

... 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

Bootstrapping in statistics

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

See blackboard for rough idea

Show animation in R: `boot.iid` in `animation` package.

Bootstrap principle

Let θ 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 θ 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 .

Bootstrap principle

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 :

$$\hat{F}(x) = \frac{1}{n} \mathbf{1}(x_i \leq x)$$

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

Examples

Thus

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

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

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

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

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

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

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 θ . 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 $\mathcal{X}^* = \{X_1^*, X_2^*, X_3^*\}$.

Bootstrap estimate for standard error

- Parameter of interest: $\theta = T(F)$
- Our estimator for θ : $\hat{\theta} = s(x)$
- **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 $\hat{\theta}$.

Note: Except for very small n , $SD_{\hat{F}}(\hat{\theta}^*)$ cannot be computed. (Number of possible bootstrap sample: n^n .)

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

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

How large do we need B ?

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

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 $SE_F(\hat{\theta})$.

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

The variability in \widehat{SE}_B depends on n and B .

The parametric bootstrap

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

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

another estimate of F may be employed.

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

Bootstrapping regression

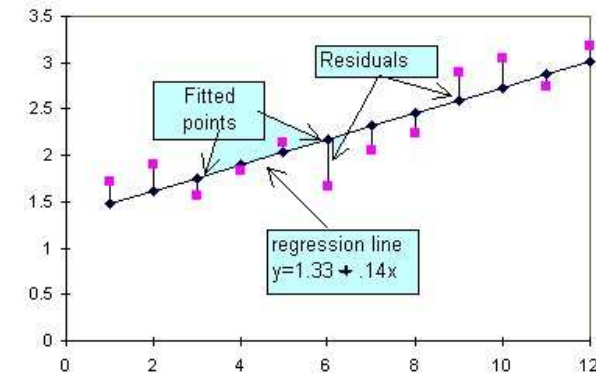
Consider the ordinary multiple regression model

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

where ϵ_i are iid mean zero random variables with constant variance.

- Naive: Bootstrapping by resampling from response variables to get distribution of $\hat{\boldsymbol{\beta}}^*$. However $Y_i|\mathbf{x}_i$ are not iid.
- Correct: Bootstrap the residuals.

Review: Residuals



<http://fsweb.bainbridge.edu/dbyrd/statistics/regression.htm>

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 $\hat{\boldsymbol{\beta}}^*$.

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

Bootstrapping residuals: Remarks

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

Bootstrapping the residuals is reliant on

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

Otherwise, a different scheme is recommended.

Paired bootstrap

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

⇒ Data pairs $z_i = (x_i, y_i)$ can be regarded as iid realisation from $Z_i = (X_i, Y_i)$ drawn from a **joint response-predictor distribution**.

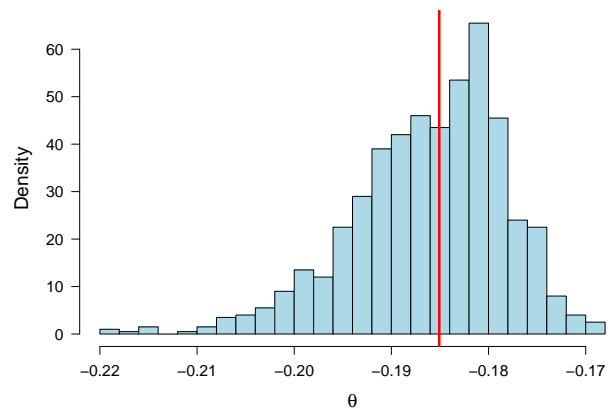
Bootstrap:

- Sample Z_1^*, \dots, Z_n^* completely at random with replacement from z_1, \dots, 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.

Histogram of 10000 bootstrap estimates



Show R-code demo-pairedBootstrap.R

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 $\hat{\theta} = \hat{\beta}_1/\hat{\beta}_0 = -0.185$.

Bootstrap bias correction

The mean value of

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

among the pseudo datasets is -0.00125 .

The **bias-corrected bootstrap estimate** of β_1/β_0 is $-0.18507 - (-0.00125) = -0.184$.

Bootstrapping dependent data

Critical requirement: Bootstrapped quantities are iid.

Consider a **first-order stationary autoregressive process**, the AR(1) model:

$$X_t = \alpha X_{t-1} + \epsilon_t$$

where $|\alpha| < 1$ and ϵ_t are iid with mean zero and constant variance.

Here, a method akin to bootstrapping the residuals for linear regression can be applied.

AR(1) model: A model based approach

1. Use a standard method to estimate α
2. Define the estimated innovations $\hat{\epsilon}_t = X_t - \hat{\alpha}X_{t-1}$ for $t = 2, \dots, n$ and let $\bar{\epsilon}$ be the mean of these.
3. Recenter $\hat{\epsilon}_t$ to have mean zero by defining $\hat{\epsilon}_t = \hat{\epsilon}_t - \bar{\epsilon}$.
4. Resample $n + 1$ values from the set $\{\hat{\epsilon}_2, \dots, \hat{\epsilon}_n\}$ with replacement to yield pseudo innovations $\{\epsilon_0^*, \dots, \epsilon_n^*\}$.
5. Generate pseudo data as $X_0^* = \epsilon_0^*$ and $X_t^* = \hat{\alpha}X_{t-1}^* + \epsilon_t^*$ for $t = 1, \dots, n$.

AR(1) model: A model based approach

Issue: Pseudo-data series is not stationary.

Remedy: Sample larger number of pseudo innovations and generate data series earlier, i.e. X_k^* for k much less than zero. The first portion of the data can be discarded as burn-in.