

## 05. General Model diagnosis

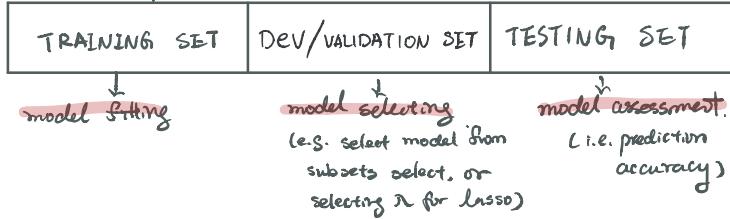
For classic models like linear regression models and GLM, the goodness of fit and predictability can be assessed by:

- various hypothesis test (e.g. Wald test, Score test, and likelihood ratio test); and.
- model selection criteria (e.g. AIC, BIC and various IC)  $\hookrightarrow$  make adjustment to training error to estimate testing error indirectly.

General model assessment methods:

- resampling: CV  $\rightarrow$  estimate testing error directly (more in 03); bootstraps (statistical inference like CI)
- Gains & ROC  $\rightarrow$  better if applied to predictions made on a validation set. (to help model selection).
- Actual versus Expected prediction (fitted) plots
- Partial dependence plots

First we need to split the data into :

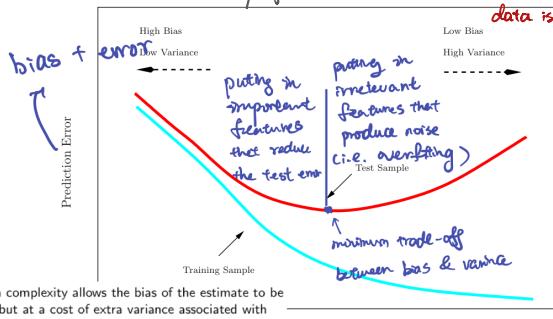


\* usually split into 2:1:1 for small dataset (<10,000)

\* usually split into 98:1:1 for large dataset (1,000,000)

Why do we use test data to assess model performance:

- TEST error is the average error that results from using a model to predict the response on new observations
- TRAN error is the residuals that result from a model fitting on observations
- If we use TRAN error, we tend to over-fit the data, and include unnecessary terms in the model;  
TRAN error doesn't tell us how well the model performs on new data.  $\Rightarrow$  a model that fits well on training data is not indicative of true performance.



The extra complexity allows the bias of the estimate to be reduced, but at a cost of extra variance associated with estimating parameters.

\* model complexity ↑, Te error ↓  
Te error ↑ (due to overfitting)

Model Complexity  $\rightarrow$  e.g. # feature from low to high, degree of polynomial from low to high.

3 / 44

### (2) Gains curve & ROC curve.

#### 2.1 Gains Curve

method: low predictions should correspond to low response observations, and high predictions should correspond to high response observations

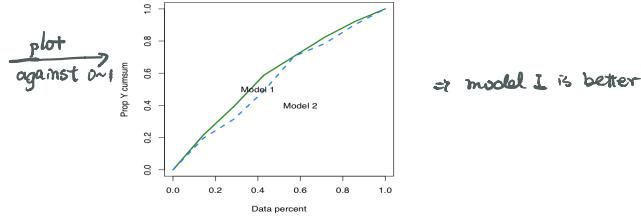
$\Rightarrow$  order the observations based on predicted response descendingly, and see how well the actual responses are sorted descendingly.

$\Rightarrow$  then include the cumulative sums of observations, as well as the % cumsum/total sum.

$\Rightarrow$  better model will have larger cumsum earlier since the larger predicted responses appear earlier.

steps: best way to describe is to use an example: suppose we have 2 candidates  $f_1, f_2$  to predict  $Y$ :

observ. model1 model2			$f_1(X_i)$		$f_2(X_i)$		$f_1(X_i)$		$f_2(X_i)$		$f_1(X_i)$		$f_2(X_i)$		
$Y_i$	$f_1(X_i)$	$f_2(X_i)$	$f_1(X_i)$	$Y_i$	$f_2(X_i)$	$Y_i$	$f_1(X_i)$	$Y_i$	$f_2(X_i)$	$Y_i$	$f_1(X_i)$	$Y_i$	$f_2(X_i)$	$Y_i$	
2.0	2.0	3.2	6.0	5.5	4.5	5.0	6.0	5.5	4.5	5.0	6.0	5.5	4.5	5.0	0.20
2.5	2.5	2.5	5.0	4.5	4.0	3.0	5.0	4.5	4.0	3.0	5.0	4.5	4.0	3.0	0.31
3.0	3.0	3.0	4.5	5.0	3.5	4.5	4.5	5.0	15.0	15.5	0.59	3.5	4.5	12.5	0.49
3.0	4.0	4.0	4.0	3.0	3.4	5.5	4.0	3.0	18.0	0.71	3.4	5.5	18.0	0.71	
4.5	5.0	3.5	3.0	3.0	3.2	2.0	3.0	3.0	21.0	0.82	3.2	2.0	20.0	0.78	
5.0	4.5	4.5	2.5	2.5	3.0	3.0	2.5	2.5	23.5	0.92	3.0	3.0	23.0	0.90	
5.5	6.0	3.4	2.0	2.0	2.5	2.5	2.0	2.0	25.5	1.00	2.5	2.5	25.5	1.00	



mathematic representations:

Gains curve plots  $\frac{\sum_{i=0}^j Y_{k_i}}{\sum_{k=1}^n Y_{k_i}}$  versus  $\frac{j}{n}$ ,  $j = 0, 1, \dots, n$

where  $k_1, k_2, \dots, k_n$  are ordered index of  $f(x_{k_1}) \geq f(x_{k_2}) \dots \geq f(x_{k_n}) \Rightarrow f(x_{k_n}) = f(x_{(n)})$

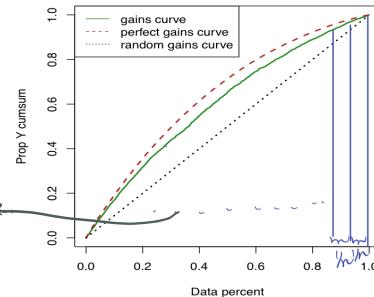
Thus, the points can be represented as:  $(G_x(t), G_y(t))$  where

$$G_x(t) = \frac{\sum_{i=0}^n I(f(x_i) > t)}{n}, \quad G_y(t) = \frac{\sum_{i=0}^n Y_i I(f(x_i) > t)}{\sum_{k=1}^n Y_{k_i}}$$

Compare with 'perfect curve' and 'random curve':

- perfect curve: prediction is observation itself

- random curve: randomly prediction i.e. 45° line



AUC (Area under curve):

$$AUC = (2n \sum_{k=1}^n Y_{k_i})^{-1} \sum_{j=0}^{n-1} \left[ \sum_{i=0}^j Y_{k_i} + \sum_{i=j+1}^{n-1} Y_{k_i} \right]$$

problem: only compares models in terms of how well the ordering of observations is concordant with that of their prediction  $\Rightarrow$  not able to assess goodness of fit.

ROC (receiver operating characteristic) (similar to Gains curve but  $Y$  is binary)

method: First construct a confusion matrix; then plot sensitivity rate against false positive rate

also try different cut-off value  $t$ ; we prefer models with larger tp and tn, and smaller fp and fn.

steps:

(1)

construct:

		prediction		Actual	total
		class 0	class 1		
class 0	tn	fp	tn + fp		
	fn	tp	fn + tp		
Total	tn + fn	fp + tp	n		

not useful

$$\begin{aligned} \text{sensitivity rate} &= \frac{tp}{tp + fn} && \begin{aligned} &\text{if cut-off value = 0} \\ &(\% \text{ of real is predicted as 1}) \end{aligned} \\ \text{specificity rate} &= \frac{tn}{tn + fp} && \begin{aligned} &\text{if cut-off value = 1} \\ &(\% \text{ of real is predicted as 0}) \end{aligned} \\ \text{false positive rate} &= \frac{fp}{tn + fp} && (\% \text{ of real is predicted as 1}) \\ \text{false negative rate} &= \frac{fn}{fn + tp} && (\% \text{ of real is predicted as 0}) \end{aligned}$$

depends on cut-off value, or  $f(x)$ .

(2) order observations based on predictions descendingly, i.e. find  $k_1, k_2, \dots, k_n$  st.  $f(x_{k_1}) \geq f(x_{k_2}) \dots \geq f(x_{k_n})$  where  $f(x)$  represents the prob. of observation being classified as 1. e.g.  $f(x_i)$

$$(3) \text{ compute sensitivity rate } R_s(j) = \frac{tp_j}{tp_j + fn_j} = \frac{\sum_{i=1}^j Y_{k_i}}{\sum_{k=1}^n Y_{k_i}} \text{ and}$$

$$\text{false positive rate } R_{fp}(j) = \frac{fp_j}{fp_j + tn_j} = \frac{\sum_{i=1}^j (1 - Y_{k_i})}{\sum_{k=1}^n (1 - Y_{k_i})}, \quad j = 1, \dots, n.$$

(4) plot  $R_s(j)$  against  $R_{fp}(j)$ ,  $j = 0, \dots, n$

(5) let  $t$  be the cut-off value, the point on ROC ( $R_s(t)$ ,  $R_{fp}(t)$ ) is:

$$R_s(t) = \frac{\sum_{i=0}^n Y_i I(f(x_i) > t)}{\sum_{k=1}^n Y_{k_i}}, \quad R_{fp}(t) = \frac{\sum_{i=0}^n (1 - Y_i) I(f(x_i) > t)}{\sum_{k=1}^n (1 - Y_{k_i})}$$

$$(6) AUC = \frac{1}{2} \sum_{j=0}^{n-1} (R_s(j) + R_s(j+1))(R_{fp}(j+1) - R_{fp}(j)).$$

$f(x_i)$	$Y_i$	$R_s(t=5)$	$R_{fp}(t=5)$
0.99	1	1/50	45/50
0.98	1	7/50	43/50
0.97	1	14/50	36/50
0.96	1	21/50	29/50
0.95	1	28/50	22/50
0.94	0	35/50	15/50
0.93	0	36/50	14/50
0.92	0	37/50	13/50
0.91	0	38/50	12/50
0.90	0	39/50	11/50
0.89	0	40/50	10/50
0.88	0	41/50	9/50
0.87	0	42/50	8/50
0.86	0	43/50	7/50
0.85	0	44/50	6/50
0.84	0	45/50	5/50
0.83	0	46/50	4/50
0.82	0	47/50	3/50
0.81	0	48/50	2/50
0.80	0	49/50	1/50
0.79	0	50/50	0/50

$n=100, X_p=50, X_n=50$

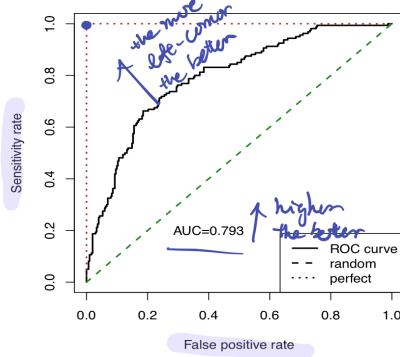
Compare with 'perfect curve' and 'random curve'

- perfect curve: prediction is observation itself  
i.e.  $f_p = 0, t_p = 1$
- random curve: randomly prediction i.e.  $45^\circ$  line

relationship with Gains curve:

- Gains curve is a 'flattened' version of ROC.
- if one model dominates on one, it dominates on the other as well.

**problem:** ignore the actual probability but focus on the ordering.



### (3) AVE

steps: (1) categorise  $X_{ij}^*$ ,  $j: X_{ij}^* \in k^{\text{th}}$  category,  $i=1, \dots, n$

$$(2) \text{ average}_k = \frac{\sum_{i=1}^n Y_i I(X_{ij}^* = k)}{\sum_{i=1}^n I(X_{ij}^* = k)}, \text{ expected}_k = \frac{\sum_{i=1}^n f(x_i) I(X_{ij}^* = k)}{\sum_{i=1}^n I(X_{ij}^* = k)}$$

(3) plot  $\text{average}_k$  against  $\text{expected}_k$ ,  $k=1, \dots, k$ .

(4) replace  $X_{ij}^*$  with  $f(x_i)$  to plot AVE from  $f(x)$

### (4) Partial dependence plots

methods: suppose we have  $p$  predictors, split them into  $X_s$  and  $X_c$ . So,  $f(x) = f(X_s, X_c)$

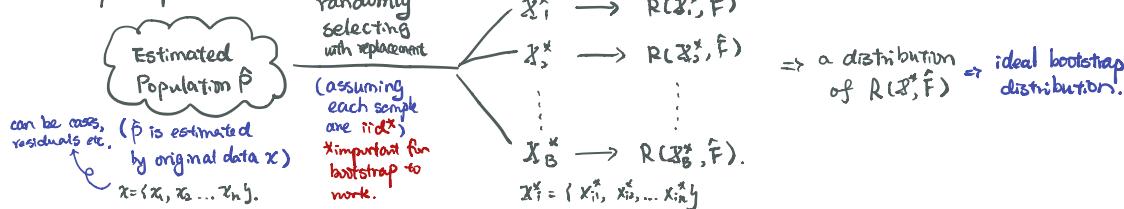
Then the partial dependence of  $Y$  on  $X_s$  is defined as  $f_s(X_s) = E_{X_c} f(X_s, X_c)$  with empirical

$$\text{approximation } \hat{f}_s(X_s) = \frac{1}{B} \sum_{i=1}^B f(X_s, \tilde{X}_{-i})$$

Then we can plot  $\hat{f}_s$  against  $X_s$ .

### (1) Bootstrap. (a special case of MC)

principles:



note:

- bootstrap cannot estimate prediction error because about  $2/3$  of original data appear in each bootstrap sample.  
 $\Rightarrow$  bootstrap will seriously underestimate prediction error as the model is trained with  $2/3$  of data.
- there total  $n^n$  possible sample, in which  $\binom{n-1}{n}$  are different.  $\Leftrightarrow B \ll \binom{n-1}{n}$ .

Theoretical Background:

- True population  $\hat{x} \sim F$ , we can use the ecdf of the original sample  $\hat{F}$  to estimate  $F$ .
- Thus  $R(\hat{x}^*, \hat{F})$  can be an estimator of  $R(\hat{x}, F)$ , which is a plug-in estimator, e.g.  $\hat{\mu} = \bar{x}_i$  is a plug-in estimator of population mean,  $\hat{\sigma}^2 = (\bar{x}_i - \hat{\mu})^2$  is a plug-in estimator of population variance.
- The bootstrap principle says  $R(\hat{x}, F) \approx R(\hat{x}^*, \hat{F})$ .

Methods:

#### (1) Nonparametric & parametric:

**Nonparametric:** (1) Bootstrap samples are drawn from  $x = \{x_1, \dots, x_n\}$  directly  $\Rightarrow \hat{x}^* = \{\hat{x}_1^*, \dots, \hat{x}_n^*\}$  for  $i=1, \dots, B$   
(2) Then ecdf of  $\{R(\hat{x}_i^*, \hat{F})\}$  is used to estimate ecdf of  $R(\hat{x}^*, \hat{F})$  and thus  $R(\hat{x}, \hat{F})$ .

**Parametric:** (1) Assume  $x_1, \dots, x_n \stackrel{\text{iid}}{\sim} F(x|\theta)$ , which can be estimated by  $\hat{\theta}$  instead of being estimated by ecdf  $\hat{F}$ .

Then draw bootstrap samples from  $F(x|\hat{\theta})$  assuming each sample is iid.  $\Rightarrow \hat{x}^* = \{\hat{x}_1^*, \dots, \hat{x}_n^*\} \stackrel{\text{iid}}{\sim} F(x|\hat{\theta})$

(2) Then ecdf of  $\{R(\hat{x}_i^*, \hat{F}(x|\hat{\theta}))\}$  is used to estimate ecdf of  $R(\hat{x}^*, \hat{F}(x|\hat{\theta}))$  and thus  $R(x, \hat{F}(x|\hat{\theta}))$

Comparison between nonparametric & parametric:

Approx. of	Nonparametric BS	Parametric BS
# samples	B samples randomly drawn from $\mathbb{X}$	
true dist.	$\mathbb{X} = \{x_1, \dots, x_n\} \stackrel{iid}{\sim} F(x \theta) \approx F(x \hat{\theta})$	$\mathbb{X} = \{x_1, \dots, x_n\} \stackrel{iid}{\sim} F(x \theta) \approx F(x \hat{\theta})$
BS dist.	$\mathbb{X}_i^* = \{x_{i1}^*, x_{i2}^*, \dots, x_{in}^*\} \stackrel{iid}{\sim} \hat{F}$ for $i=1, \dots, B$	$\mathbb{X}_i^* = \{x_1^*, x_2^*, \dots, x_n^*\} \stackrel{iid}{\sim} F(x \hat{\theta})$ for $i=1, \dots, B$
ideal BS dist.	$R(\mathbb{X}^*, \hat{F}) \approx ECDF$ of $\{R(\mathbb{X}_i^*, \hat{F}), i=1, \dots, B\} \stackrel{iid}{\sim} \hat{F}$	$R(\mathbb{X}^*, F(x \hat{\theta})) \approx ECDF$ of $\{R(\mathbb{X}_i^*, F(x \hat{\theta})), i=1, \dots, B\} \stackrel{iid}{\sim} \hat{F}$
quantity dist.	$R(X, F) \approx \hat{F}^*$	$R(X, F(x \hat{\theta})) \approx \hat{F}^*$

## (2) Bootstrap in regression:

context: consider a regression model:  $Y_i = x_i^\top \beta + \epsilon_i$  for  $i=1, \dots, n$  and  $\epsilon_i \stackrel{iid}{\sim} F(0, \sigma^2)$ .

observed data are  $\{(x_1, y_1), \dots, (x_n, y_n)\}$

two ways to bootstrap sample from observed data are: bootstrap the residuals and bootstrap the case.

### 2.1 bootstrap the residuals (nonparametric/parametric).

steps:

- all these must be done in one process
- (1) fit the model to the observed data. Obtain  $\hat{y}_i = x_i^\top \hat{\beta}$  and  $\hat{\epsilon}_i = y_i - \hat{y}_i$
  - (2) nonparametric: bootstrap the residuals  $\{\hat{\epsilon}_1, \dots, \hat{\epsilon}_n\}$  to obtain  $\{\hat{\epsilon}_1^*, \dots, \hat{\epsilon}_n^*\}$  (note: important to check if  $\{\hat{\epsilon}_1, \dots, \hat{\epsilon}_n\}$  are iid.)
  - (3) parametric: simulate  $\{\hat{\epsilon}_1^*, \dots, \hat{\epsilon}_n^*\}$  from some distribution that we assume
  - (4) create a bootstrap response:  $\hat{Y}_i^* = \hat{y}_i + \hat{\epsilon}_i^*$  for  $i=1, \dots, n$
  - (5) fit the model to  $\{(x_i, \hat{Y}_i^*), i=1, \dots, B\}$  to get bootstrap estimates of  $\hat{\beta}^*$
  - (6) repeat this for  $B$  times to get  $\{\hat{\beta}_1^*, \dots, \hat{\beta}_B^*\}$ , from which an ecdf can be built for statistical inference.

### 2.2 bootstrap the cases (nonparametric).

steps:

- all these must be done in one process
- (1) bootstrap the sample  $\{\mathbb{Z}_1, \dots, \mathbb{Z}_n\}$  to obtain  $\{\mathbb{Z}_1^*, \dots, \mathbb{Z}_n^*\}$  (note: here we assume  $\{\mathbb{Z}_1, \dots, \mathbb{Z}_n\}$  are iid).
  - (2) fit the model to  $\{\mathbb{Z}_1^*, \dots, \mathbb{Z}_n^*\}$  to get bootstrap estimates of  $\hat{\beta}^*$ .
  - (3) repeat this for  $B$  times to get  $\{\hat{\beta}_1^*, \dots, \hat{\beta}_B^*\}$ , from which an ecdf can be built for statistical inference.

\* Bootstrapping the cases is less sensitive to violations in the regression model assumptions (i.e. adequacy of the model and constancy of  $\sigma^2$ ) than bootstrapping the residuals.

### 2.3 bootstrap estimation of $bias(\hat{\theta})$ and $se(\hat{\theta})$ .

$$bias(\hat{\theta}) = \frac{1}{B} \sum_{i=1}^B \hat{\theta}_i^* - \hat{\theta} = \bar{\theta}^* - \hat{\theta} \quad se(\hat{\theta}) = \sqrt{\frac{1}{B-1} \sum_{i=1}^B (\hat{\theta}_i^* - \bar{\theta}^*)^2}$$

## (3) Bootstrap inference contents: (Bootstrap CI).

### 3.1 Normal

$$[2\hat{\theta} - \bar{\theta}^* - Z_{1-\alpha/2} \cdot se(\hat{\theta}), 2\hat{\theta} - \bar{\theta}^* + Z_{1-\alpha/2} \cdot se(\hat{\theta})]$$

$$2\hat{\theta} - \bar{\theta}^* = \hat{\theta} - bias(\hat{\theta}) = \hat{\theta} - (\bar{\theta}^* - \hat{\theta})$$

$$= 2\hat{\theta} - \bar{\theta}^* = \text{unbiased estimate of } \theta.$$

✓ follow the 'generic template'  
✗ may be useful for large sample size or normal population.

### 3.2 Percentile (or basic percentile)

$$[\hat{\theta}_{d/2(B+1)}^*, \hat{\theta}_{1-d/2(B+1)}^*]$$

✗ highly asymmetric

✗ vulnerable to bias & inaccurate coverage prob.

✓ may work for location parameters.

### 3.3 Basic (or residuals)

$$[2\hat{\theta} - \bar{\theta}_{1-d/2(B+1)}^*, 2\hat{\theta} - \bar{\theta}_{d/2(B+1)}^*]$$

✓ adjusted for bias.

### 3.4 BCa (Bias corrected and accelerated).

$$[\hat{\theta}_{(B+1)P_1}^*, \hat{\theta}_{(B+1)P_0}^*]$$

where  $P_1 = \phi(\frac{C_0 - \hat{\theta}}{1 - d(C_0 + \hat{\theta})} + C_1)$ ,  $P_0 = \phi(\frac{C_0 + \hat{\theta}}{1 - d(C_0 + \hat{\theta})} + C_0)$ ,

$C_0$  is determined by the relative position of  $\hat{\theta}$  among  $\{\hat{\theta}_1^*, \dots, \hat{\theta}_B^*\}$

$d$  is determined by skewness of  $\{\hat{\theta}_1^*, \dots, \hat{\theta}_B^*\}$

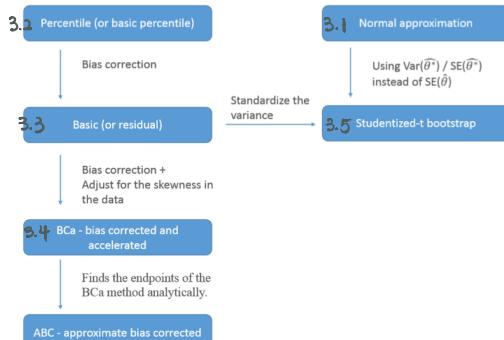
✓ adjusted for bias & skewness  
✓ reasonable computation time.  
✗ need EIF for parametric BS.

### 3.5 Studentized

$$[\hat{\theta} - \xi_{1-\alpha/2}(\hat{G}^*)\sqrt{V(\hat{\theta})}, \hat{\theta} - \xi_{\alpha/2}(\hat{G}^*)\sqrt{V(\hat{\theta})}]$$

where  $\hat{G}^*$  is the distribution of  $R(Z^*, \hat{\theta}) = \frac{\hat{\theta}^* - \hat{\theta}}{\sqrt{V(\hat{\theta}^*)}}$ ,  
 $\xi_\alpha(\hat{G}^*)$  is the  $\alpha$  quantile of  $\hat{G}^*$ .

- ✓ symmetrical CI by using  $\frac{\hat{\theta}^* - \hat{\theta}}{\sqrt{V(\hat{\theta}^*)}}$ .
- ✓ reasonable for location and variance-stabilized estimators.
- ✗ sensitive to presence of outliers.
- ✗ performs badly for heavy tailed dist.
- ✗ computationally expensive - requires double bootstrap to calculate  $V(\hat{\theta}^*)$ .



### (4) Bootstrap hypothesis testing.

$H_0: \theta = \theta_0$  vs.  $H_1: \theta \neq \theta_0$

can use statistics  $\frac{\hat{\theta}^* - \hat{\theta}}{\hat{\sigma}^*}$  or can use permutation test.

reducing MC error:  $(\frac{1}{B} \sum_{j=1}^B (\bar{X}_j^* - \bar{X}) \neq 0 \text{ due to MC variation})$ .

(1) balanced bootstrap: (all observations occur with the same frequency as in the original sample)

The simplest way to get  $B$  balanced bootstrap samples is to concatenate  $B$  copies of the observed sample of size  $n$ , randomly permute this series, and then read off  $B$  blocks of size  $n$  sequentially. The  $j$ th block becomes the  $j$ th bootstrap sample  $X_j^*$ . For the permutation involved, balanced bootstrap is also called **permutation bootstrap**.

### (2) antithetic bootstrap:

- ▶ For a sample of univariate data,  $x_1, \dots, x_n$ , denote the ordered data as  $x_{(1)}, \dots, x_{(n)}$ . Let  $\pi(i) = n - i + 1$  be a permutation operator that reverses the order statistics.
- ▶ Then for each bootstrap sample  $\mathcal{X}^* = \{X_1^*, \dots, X_n^*\}$ , let  $\mathcal{X}^{**} = \{X_1^{**}, \dots, X_n^{**}\}$  denote the sample obtained by substituting  $X_{(\pi(i))}$  for every instance of  $X_{(i)}$  in  $\mathcal{X}^*$ . Thus, for example, if  $\mathcal{X}^*$  has an unrepresentative predominance of the larger observed data values, the smaller observed values will predominate  $\mathcal{X}^{**}$ .
- ▶ Using this strategy, each bootstrap draw provides two estimators:  $R(\mathcal{X}^*, \hat{F})$  and  $R(\mathcal{X}^{**}, \hat{F})$ . The two estimators are often **negatively correlated**.

▶ Let  $R_a(\mathcal{X}^*, \hat{F}) = \frac{1}{2}[R(\mathcal{X}^*, \hat{F}) + R(\mathcal{X}^{**}, \hat{F})]$ . Then  $R_a$  has the following desirable property

$$\begin{aligned} \text{Var}[R_a(\mathcal{X}^*, \hat{F})] &= \frac{1}{4} ([\text{Var}[R(\mathcal{X}^*, \hat{F})]] + [\text{Var}[R(\mathcal{X}^{**}, \hat{F})]]) \\ &\quad + 2\text{Cov}[R(\mathcal{X}^*, \hat{F}), R(\mathcal{X}^{**}, \hat{F})]) \\ &\leq \text{Var}[R(\mathcal{X}^*, \hat{F})] \end{aligned}$$

if the covariance is negative.

▶ The above strategy of reducing Monte Carlo error in bootstrap is referred to as **antithetic bootstrap**.

▶ It is also possible to establish ordering of multivariate data to permit an antithetic bootstrap strategy.

e.g.  $B=3, n=3$

$Z^*: 112 \ 223 \ 133 \rightarrow R(Z^*)$

$Z^{**}: 332 \ 221 \ 311 \rightarrow R(Z^{**})$

permutation BS samples:  $R(Z^*)$

antithetic BS samples:  $R_a(Z^*) = \frac{1}{2}[R(Z^*) + R(Z^{**})]$