Cross-validation (k-fold):

Whole data set $D$:

| $C_1$ | $C_2$ | $C_3$ | $\cdots$ | $C_K$ |

$n_1 \quad n_2 \quad n_3 \qquad\qquad n_K$

$C_i$ : labels# of our data points corresponding to the $i$th subset

$|C_i|$ = "cardinality of set $C_i$" = # of elements in $C_i$ ; $n_i$

$D = \bigcup\limits_{i=1}^{K} C_i$ ; $C_i \cap C_j = \emptyset$ ( subsets are disjoint ).

$T_i = D \setminus C_i$ $\qquad\qquad = \bigcup\limits_{\substack{j=1 \\ j \neq i}}^{K} C_j$ ; $T_i$ is $i$th training set

$\underset{\text{"set difference"}}{\uparrow}$ $\qquad\qquad$ $V_i = C_i$ : the $i$th validation set

① For $i = 1, \ldots, K$

   (a) "train"/ fit our model on $T_i$ , validate it on $V_i$.

   (b) get a "discrepancy" $\text{Discr}(\hat{\underset{\sim}{Y}}^{(i)}, \underset{\sim}{Y}^{(i)})$ between predicted values $\hat{\underset{\sim}{Y}}^{(i)}$ (based on $T_i$) for $\underset{\sim}{Y}^{(i)}$ (from $V_i$).

   "Discr": MSE or misclassification rate.

② Aggregate the discrepancies from $i = 1, \ldots, K$ into a single measure.

How to calibrate/estimate tuning parameters using K-fold CV?

Examples:

1) polynomial regression: tuning par. is degree of polynomial.

2) GAM, ridge regression, lasso: need to choose the penalty parameter $\lambda$ (for "regularization").

3) model selection: need to pick a subset of covariates.

In practice, the "aggregated discrepancy" $(AD)$ depends on the tuning parameters $(\tau)$.

$\Rightarrow$ optimize/minimize AD with respect to $\tau$.

# Cross-Validation for Classification Problems

- We divide the data into $K$ roughly equal-sized parts $C_1, C_2, \ldots C_K$. $C_k$ denotes the indices of the observations in part $k$. There are $n_k$ observations in part $k$: if $n$ is a multiple of $K$, then $n_k = n/K$.  $\circledast$  1. assume

  $n_k = const$, so

  $n_k/n = 1/K$

- Compute

$$\text{CV}_K = \sum_{k=1}^{K} \frac{n_k}{n} \underset{z_k}{\boxed{\text{Err}_k}}$$

$$CV_k = \bar{z} = \frac{1}{k} \sum_{j=1}^{K} z_j$$

  where $\text{Err}_k = \sum_{i \in C_k} I(y_i \neq \hat{y}_i)/n_k$.

$\circledast$
- The estimated standard deviation of $\text{CV}_K$ is

$$\underset{Err_k}{\widehat{\text{SE}}(\text{CV}_K)} = \sqrt{\sum_{k=1}^{K} (\text{Err}_k - \overline{\text{Err}_k})^2/(K-1)}$$

- This is a useful estimate, but strictly speaking, not quite valid.

1. To estimate $\text{Var}(z_j)$, we can use sample variance.

2. To estimate, note $\text{Var}(\bar{Z})$

$$\text{Var}(\bar{Z}) = \text{Var}\left(\frac{1}{k}\sum_{j=1}^{K} z_j\right) = \frac{1}{K^2}\text{Var}\left(\sum_{j=1}^{k} z_j\right)$$

$$= \frac{1}{k^2}\sum_{i=1}^{K}\sum_{j=1}^{K}\text{Cov}(z_i, z_j)$$

$$\underbrace{\sum_{j=1}^{k}\text{Var}(z_j)} + \underbrace{\sum_{i\neq j}\sum \text{Cov}(z_i, z_j)}$$
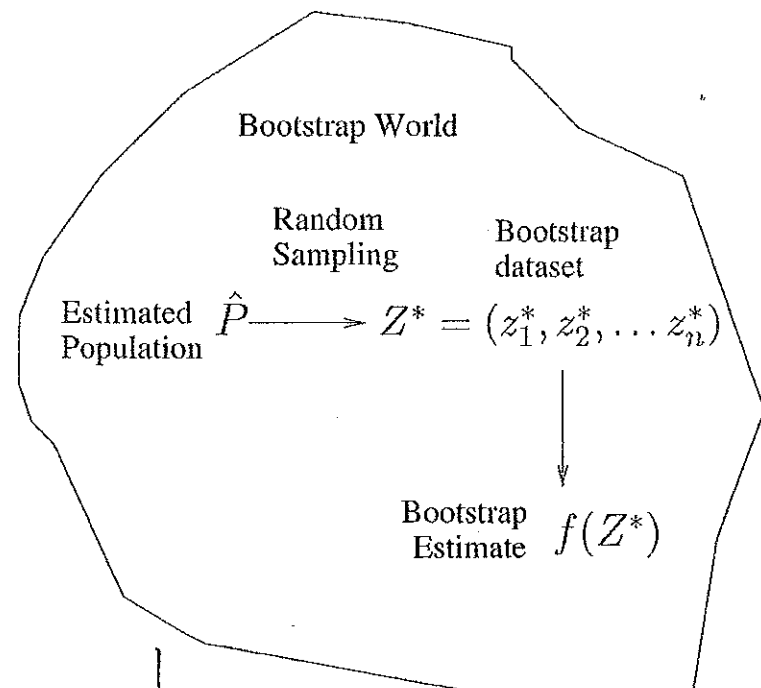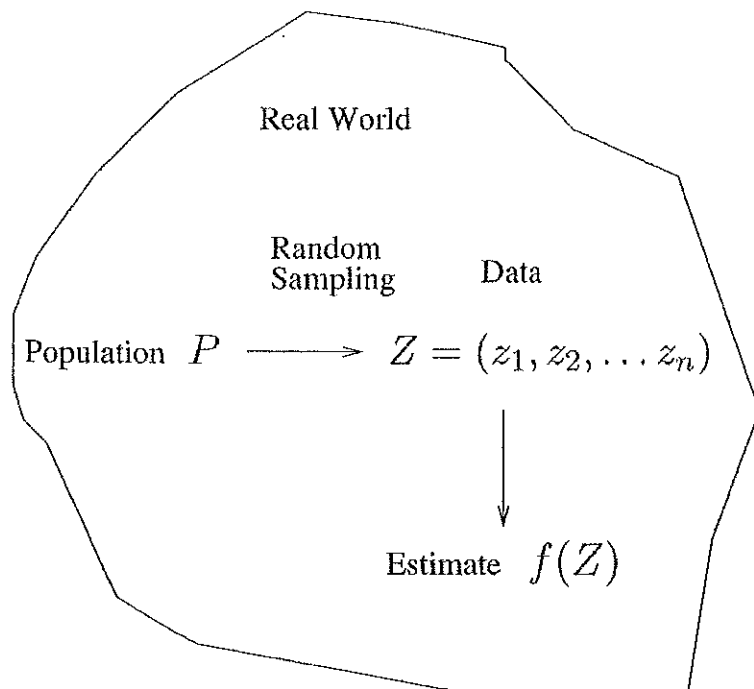
cannot be neglected for $SE(CV_k)$

# A general picture for the bootstrap

Case 1: finite populations.
$$P = \{ a_1, a_2, \ldots, a_N \}.$$

$$\hat{P} = \{ z_1, \ldots, z_n \} : \text{our dataset}$$

**Real World**

Random
Sampling    Data

Population $P \longrightarrow Z = (z_1, z_2, \ldots z_n)$

$\downarrow$

Estimate $f(Z)$

**Bootstrap World**

Random        Bootstrap
Sampling      dataset

Estimated $\hat{P} \longrightarrow Z^* = (z_1^*, z_2^*, \ldots z_n^*)$
Population

$\downarrow$

Bootstrap
Estimate $f(Z^*)$

$X$: a rvs for sampling with replacement.
$$Pr(X = a_i) = 1/N$$
$$Pr(X \leq x) = \sum_{a_i \leq x} Pr(X = a_i) = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(a_i \leq x)$$

If $n$ is large, $\hat{P}$ is "representative" of $P$.
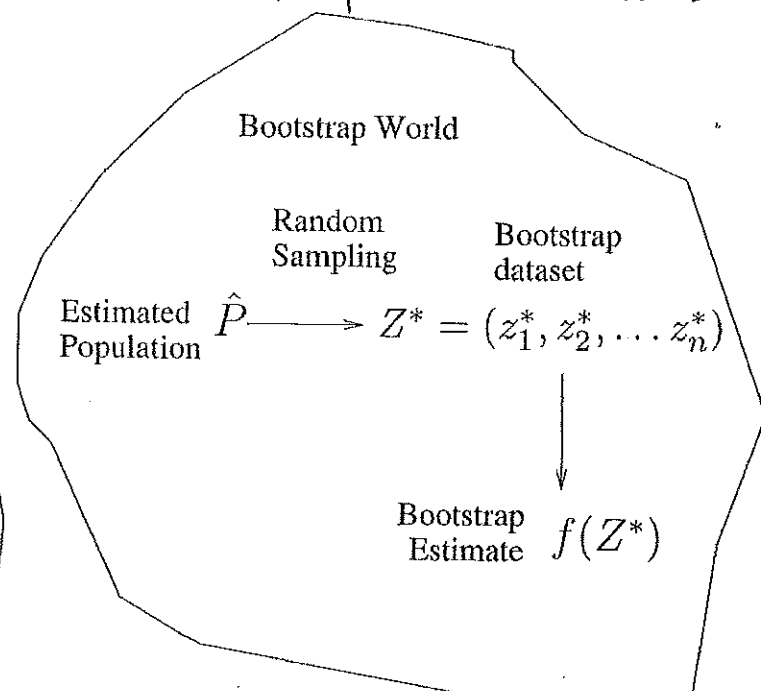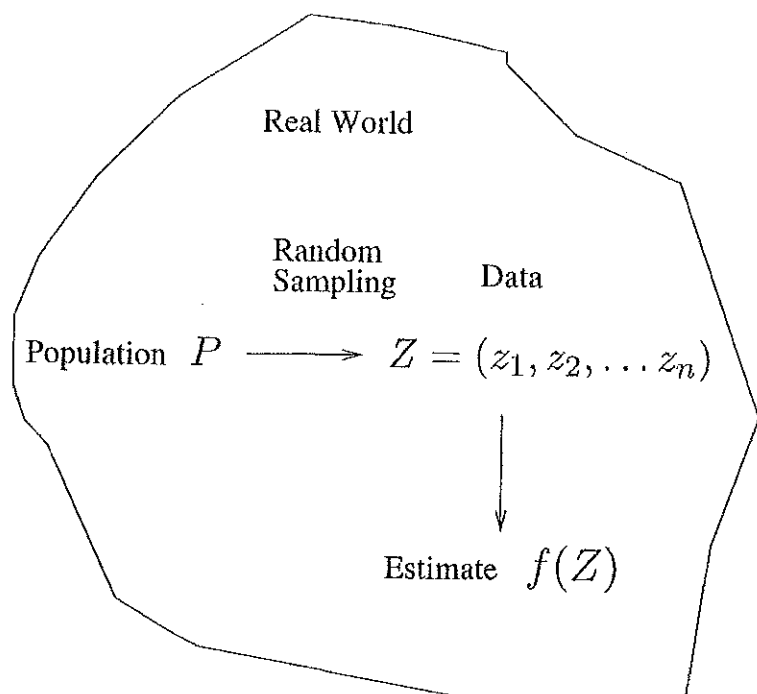
$X^*$: a rv for sampling with replacement
$$Pr(X^* = z_i) = 1/n$$
$$Pr(X^* \leq x) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{I}(z_i \leq x)$$

# A general picture for the bootstrap

Case 2: infinite populations: Think in terms of distribution functions

"non parametric Bootstrap".

**Real World**

Population $P$ $\longrightarrow$ $Z = (z_1, z_2, \ldots z_n)$

Random Sampling    Data

$\downarrow$

Estimate $f(Z)$

**Bootstrap World**

Estimated Population $\hat{P} \longrightarrow Z^* = (z_1^*, z_2^*, \ldots z_n^*)$

Random Sampling    Bootstrap dataset

$\downarrow$

Bootstrap Estimate $f(Z^*)$

$P$ is specified by a distribution function: $F(x) = \Pr(X \leq x)$

$\uparrow$

true cdf of the $X$

<span style="color:red">F(x) = Pr(X <= x)</span>

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

$\uparrow$ empirical distribution function.

$$\hat{F}_n(x) \to F(x) \text{ as } n \to \infty,$$

In R, ecdf function.

# Parametric Bootstrap:

Suppose we knew that the true cdf $F$ is Normal $(\mu, \sigma^2)$; $\mu, \sigma^2$: both unknown.

How can we estimate $F(x)$ or $F$ (the whole function)?

Plug-in estimator: estimate $\mu$ by $\bar{X}_n$, sample mean

$\sigma^2$ by $S^2$, sample var.

then Normal $(\mu = \bar{X}, \sigma^2 = S^2)$ is our estimator for $F$.