

SF2935: MODERN METHODS OF STATISTICAL LEARNING LECTURE 4 RESAMPLING METHODS: BOOTSTRAP AND CROSS-VALIDATION

Tatjana Pavlenko

7 November 2017



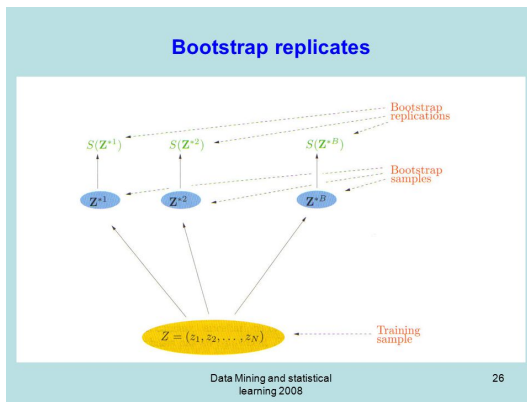
THE BOOTSTRAP: BASIC IDEA AND EXAMPLE

The bootstrap is a general resampling technique for assessing statistical accuracy. We first describe the technique in general.

- ▶ Suppose that we have a model fit to a set of training dataset $\mathbf{z} = (z_1, \dots, z_N)$ where $z_i = (x_i, y_i)$.
- ▶ The basic idea of the bootstrap sampling is to randomly draw datasets with replacement from the training data \mathbf{z} , each dataset of the same size N as the original dataset. This is done B times, generating B bootstrap datasets.
- ▶ We further refit the model to each of the bootstrap sample, and investigate the behavior of the fits over the B replications.



BOOTSTRAP (CONT.)



FIGUR: Schematic of the bootstrap technique. The goal is to evaluate the statistical accuracy of $S(\mathbf{Z})$ computed from the dataset $\mathbf{Z} = (Z_1, \dots, Z_N)$. B training samples, $\{\mathbf{Z}^{*b}\}_{b=1}^B$ each of size N are drawn **with replacement** from the original dataset. $S(\mathbf{Z})$ is then computed for each bootstrap training sample, and values $\{S(\mathbf{Z}^{*b})\}_{b=1}^B$ are used to evaluate the statistical accuracy of $S(\mathbf{Z})$. **Obs!** These notations are used locally to be consistent with the figure.

THE BOOTSTRAP: BASIC IDEA AND EXAMPLE

Let $S(\mathbf{Z})$ be any quantity (statistics) computed from \mathbf{Z} . Using bootstrap sampling we can estimate any aspect of the distribution of $S(\mathbf{Z})$! Assume that we seek to estimate $\text{Var}[S(\mathbf{Z})]$.

Bootstrap variance estimator

1. Draw a bootstrap sample $\mathbf{Z}^* = (Z_1^*, \dots, Z_N^*)$. Compute $S(\mathbf{Z}^*)$.
2. Repeat the previous step, B times, yielding samples $\{\mathbf{Z}^{*b}\}_{b=1}^B$ and estimators $S(\mathbf{Z}^{*b})$.
3. Compute

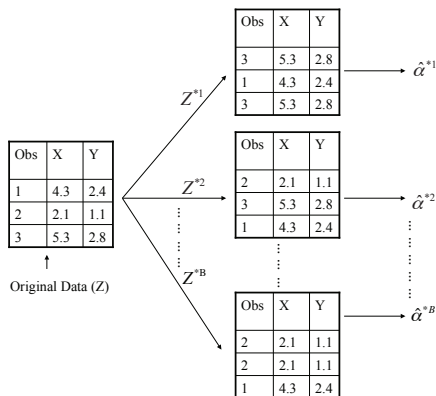
$$\widehat{\text{Var}}_{\text{boot}}[S(\mathbf{Z})] = \frac{1}{B-1} \sum_{b=1}^B \left(S(\mathbf{Z}^{*b}) - \bar{S}^* \right)^2,$$

where $\bar{S}^* = \sum_{b=1}^B S(\mathbf{Z}^{*b}) / B$.

$\widehat{\text{Var}}_{\text{boot}}[S(\mathbf{Z})]$ can be thought of as a Monte-Carlo estimate of the variance of $S(\mathbf{Z})$ under sampling from the *empirical distribution function* (EDF) for the data $\mathbf{z} = (z_1, \dots, z_N)$. More details on EDF will be presented later.



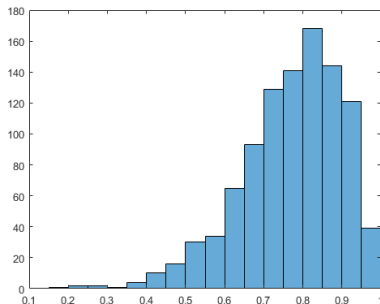
BOOTSTRAP METHODS (CONT.)



FIGUR: In the real world we only get to see the single value of $\hat{\alpha}$, but the bootstrap world is more generous: we can generate as many bootstrap replications $\hat{\alpha}^{*b}$ as we wish (or have capacity for), and directly estimate their variability using e.g. $\widehat{\text{Var}}_{\text{boot}}$ or $\widehat{\text{se}}_{\text{boot}}$.



BOOTSTRAP METHODS (CONT.)



FIGUR: Let α denote the correlation coefficient parameter, i.e $\alpha = \text{Corr}(X, Y)$. A histogram of the estimates of α obtained from $B = 1000$ bootstrap samples form the original data set. These notations are used locally to be consistent with the figure on slide 5.

BOOTSTRAP(CONT.)



Münchhausen

O. Herfurth pinx

FIGUR: The term **bootstrap** attributed to a story in *The Surprising Adventures of Baron Münchhausen* where the eponymous Baron pulls himself out of a swamp by his own bootstraps. Published in London in 1785 by Rudolf Erich Raspe.



WHY DOES THE BOOTSTRAP WORK? STATISTICAL ISSUES

There are two senses in which bootstrap resampling methods might *work*:

- ▶ Q1: Do these methods provide reliable results when being used with the data generated in applications?
- ▶ Q2: Under what theoretical (idealized) conditions will a bootstrap resampling technique provide results that are in some sense statistically correct?

An important role in answering these questions is played by the *empirical distribution function*, (EDF).



Assume the following.

- ▶ We have a sample of data, $x = (x_1, \dots, x_n)$ which are thought as the outcomes of iid r. v. X_1, \dots, X_n whose pdf and cdf are f and F , respectively.
- ▶ The sample is to be used to make inference about a population characteristic (for ex. parameter of F), denoted by θ ($\theta = \theta(F)$ which means that θ is some function of the true unknown distribution F) using a statistics $t = t(X)$ whose value at the sample is $t(x)$.
- ▶ The choice of t has been made, i.e it is an estimator of θ which we take to be a scalar. For instance, let θ be the mean of F , then

$$\theta(F) = \mathbb{E}_F(X) = \int xf(x)dx,$$

- ▶ We focus attention on the uncertainty of estimator $t(X)$. For example, what is its bias, its standard error $se(t(X))$ or its quantiles? How do we obtain confidence limits for θ using $t(X)$? All these questions are concerned with the **probability distribution of $t(X)$** .



THE BOOTSTRAP IN GENERAL (CONT.)

In the real world $\hat{\theta} = t(\mathbf{x})$ is obtained in two steps: first \mathbf{x} is generated by iid sampling from F and then $\hat{\theta}$ is calculated from \mathbf{x} according to $t(\cdot)$, i.e

$$F \xrightarrow{\text{iid}} \mathbf{x} \xrightarrow{t(\cdot)} \hat{\theta}.$$

Since we don't know F the natural way to proceed is

- ▶ replace F with its data-based approximation denoted by \hat{F} .
- ▶ use the Monte-Carlo sampling from the approximation \hat{F} to investigate the variation in $t(X)$.

The common strategy to obtain the approximation \hat{F} is to use the **empirical distribution function** (EDF).

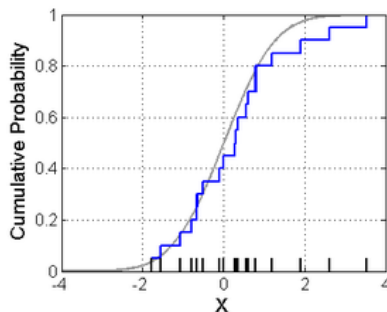
- ▶ The EDF associated with the dataset $\mathbf{x} = (x_1, \dots, x_n)$ is defined as a sample proportion

$$\hat{F}_n(u) = \frac{\#\{x_i \leq u\}}{n} = \frac{1}{n} \sum_{i=1}^n H(u - x_i),$$

$H(u)$ is the unit step function which jumps from 0 to 1 at $u = 0$



BOOTSTRAP METHODS (CONT.)



FIGUR: The blue line represents an empirical distribution function and the gray curve is the true cumulative distribution function. Notice that the values of EDF are fixed at $(0, \frac{1}{n}, \frac{2}{n}, \dots, \frac{n}{n})$, so the EDF is equivalent to its points of increase, the ordered values $(x_{(1)}, \dots, x_{(n)})$ of the data.

$\hat{F}_n(x)$ plays the role of fitted population model when no mathematical form is assumed for F . How to sample from $\hat{F}_n(x)$?

- ▶ Observe that $\hat{F}_n(u) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x_i \leq u\}}$ of $\mathbf{x} = (x_1, \dots, x_n)$ places equal weights of $1/n$ to each distinct observed value x_i .
- ▶ This implies that *a bootstrap sample \mathbf{x}^* is an iid sample drawn from \hat{F}_n* . Then each x^* independently has equal probability of $1/n$ of being any member of $\{x_1, \dots, x_n\}$
- ▶ This in turn implies that a sample \mathbf{x}^* of size n from the empirical distribution $\hat{F}_n(x)$ associated with the data $\mathbf{x} = (x_1, \dots, x_n)$ can be generated as follows:
 - ▶ sample independently n indices i_1, i_2, \dots, i_n from the uniform distribution over the set of integers $\{1, 2, \dots, n\}$;
 - ▶ let $\mathbf{x}^* = (x_{i_1}, x_{i_2}, \dots, x_{i_n})$.
- ▶ This sampling procedure draws n elements from the set $\{x_1, x_2, \dots, x_n\}$ with replacement! (Recall the bootstrap.)



SOME PROPERTIES OF \hat{F}_n .

- ▶ For a fixed u , $\mathbb{1}_{\{X_i \leq u\}}$ is $Be(p)$ random variable with $p = F(u)$. Hence $n\hat{F}_n(u) = \sum_{i=1}^n \mathbb{1}_{\{X_i \leq u\}}$ is $Bin(n, F(u))$. This implies that $\hat{F}_n(u)$ is an unbiased estimator of $F(u)$.
- ▶ By the *strong law of large numbers*, as $n \rightarrow \infty$

$$\hat{F}_n(u) \xrightarrow{a.s.} F(u),$$

i.e. $\hat{F}_n(u)$ converges to $F(u)$ almost surely, for every value of u .
Thus $\hat{F}_n(u)$ is a consistent estimator of the CDF $F(u)$.

- ▶ By CLT, point-wise, $\hat{F}_n(u)$ has asymptotically normal distribution as $n \rightarrow \infty$

$$\sqrt{n} \left(\hat{F}_n(u) - F(u) \right) \xrightarrow{D} N(0, F(u)(1 - F(u))),$$

with the standard \sqrt{n} rate of convergence.



.... OF \hat{F}_n .

- ▶ Now having observed data $\mathbf{x} = (x_1, \dots, x_n)$ and using the asymptotic properties of \hat{F}_n we may *replace* F with \hat{F}_n in such a way that any quantity involving F can now be approximated by plugging \hat{F}_n into the quantity instead as $\theta = \theta(F) \approx \hat{\theta} = \theta(\hat{F}_n)$.
For example if θ is the mean of F then

$$\theta = \mathbb{E}_F(X) \approx \hat{\theta} = \mathbb{E}_{\hat{F}_n}(X) = \frac{1}{n} \sum_{i=1}^n x_i.$$

- ▶ It can be shown that \hat{F}_n maximizes the probability of obtaining the observed sample \mathbf{x} under all possible choices of F in $\mathbf{x} = (x_1, x_2, \dots, x_n)$ given that we have observed $x_i \stackrel{iid}{\sim} F$, i.e it is the **nonparametric MLE** of F .



Q: How to investigate the uncertainty in estimator $t(\mathbf{X})$ of θ using bootstrap?

The uncertainty of $t(\mathbf{X})$ is expressed in terms of the *error distribution*, i.e. the distribution of $t(\mathbf{X}) - \theta$. We study the error distribution by considering its approximation $t(\mathbf{X}) - \hat{\theta}$ where the observed value of $\hat{\theta} = \hat{\theta}(\mathbf{x})$ replaces the unknown θ . (The error distribution is thus a shift of the $t(\mathbf{X})$ distribution.)

Bootstrap estimated error distribution

1. Construct \hat{F}_n using $\mathbf{x} = (x_1, \dots, x_n)$.
- 2*. Draw a sample $\mathbf{x}^* = (x_1^*, \dots, x_n^*)$ from \hat{F}_n and compute $t(\mathbf{x}^*) - \hat{\theta}$.
3. Repeat the previous step, B times, yielding samples $\{\mathbf{x}^{*b}\}_{b=1}^B$ and corresponding set of bootstrap estimates $\{t(\mathbf{x}^{*b}) - \hat{\theta}\}_{b=1}^B$.
4. The error distribution is estimated by considering the sample distribution of the $\{t(\mathbf{x}^{*b}) - \hat{\theta}\}_{b=1}^B$ (for ex by a histogram).

Obs! In the case of the empirical distribution, sampling from \hat{F}_n at step 2* is the same as drawing, with replacement from the dataset $\mathbf{x} = (x_1, \dots, x_n)$.



- ▶ Let F_ε be the error distribution associated with $\varepsilon = t(\mathbf{X}) - \theta$. A (equi-tailed) confidence set I_θ for θ with the coverage probability $1 - \alpha$ is then defined by the limits

$$I_\theta = \left(\hat{\theta} - F_\varepsilon^{-1}(1 - \alpha/2), \hat{\theta} - F_\varepsilon^{-1}(\alpha/2) \right).$$

A confidence set satisfies $P(\theta \in I_\theta) = 1 - \alpha$ exactly, but it not possible to specify such sets except in a few special parametric models.

- ▶ Bootstrap sampling algorithms can be applied to obtain confidence sets as follows:
- ▶ The sample of bootstrapped errors $\{\varepsilon_b^*\}_{b=1}^B = \{t(\mathbf{x}^{*b}) - \hat{\theta}\}_{b=1}^B$ is approximately distributed as F_ε and can be used to approximate the quantiles of F_ε .



Bootstrap confidence set

- ▶ The sample of bootstrapped errors $\{\varepsilon_b^*\}_{b=1}^B = \{t(\mathbf{x}^{*b}) - \hat{\theta}\}_{b=1}^B$ is approximately distributed as F_ε and can be used to obtain estimated quantiles of F_ε .
- ▶ The B simulated values $\{\varepsilon_b^*\}_{b=1}^B$ are ordered as $(\varepsilon_{(1)}^*, \dots, \varepsilon_{(B)}^*)$ and the α -quantile of F_ε is estimated by the $(B+1)\alpha$ th of these, i.e. for any $\alpha \in (0, 1)$, $F_\varepsilon^{-1}(p) \approx \varepsilon_{((B+1)\alpha)}^*$.
- ▶ The confidence limits of I_θ are then replaced by their bootstrap based counterparts which gives $1 - \alpha$ bootstrap confidence set for θ

$$I_\theta^B = \left(\hat{\theta} - \varepsilon_{((B+1)(1-\alpha/2))}^*, \hat{\theta} - \varepsilon_{((B+1)(\alpha/2))}^* \right).$$

THE NON-PARAMETRIC BOOTSTRAP

- ▶ The bootstrap estimation of the confidence set considered above is *non-parametric* in the sense that **no model is assumed for the distribution F underlying the data.**
- ▶ Sampling strategy was

$$\hat{F}_n \xrightarrow{\text{iid}} \{\mathbf{x}^{*b}\}_{b=1}^B \longrightarrow \{t(\mathbf{x}^{*b})\}_{b=1}^B,$$

(from the empirical distribution of the real data), i.e. drawings from the data with replacement.

- ▶ Sampling from the empirical distribution \hat{F}_n can be treated as sampling from a finite population, where the *population* is the observed data $(\mathbf{x}_1, \dots, \mathbf{x}_n)$. This process is therefore called for *resampling*. To be iid sampling has to be with replacement.
- ▶ The bootstrap methods work well when \hat{F}_n is close to the true unknown distribution F . This will usually be the case when the sample size n is large enough.



THE PARAMETRIC BOOTSTRAP

- ▶ In the non-parametric bootstrap, the distribution F underlying the data is assumed to be completely unknown (only iid data).
- ▶ In the parametric bootstrap, we assume that data $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ are generated as iid from a parametric family F_θ with $\theta \in \Theta$. Then to utilize the known structure of F we
 - ▶ first obtain an estimator of θ , $\hat{\theta} = \hat{\theta}(\mathbf{x})$ ($\hat{\theta}$ should be a good estimator of θ , for example MLE),
 - ▶ and draw new bootstrap samples $\{\mathbf{x}^{*b}\}_{b=1}^B$ from the (plug-in) distribution $\hat{F} = F(\hat{\theta})$, from which the usual bootstrap estimators $\{t(\mathbf{x}^{*b})\}_{b=1}^B$ are then constructed.
- ▶ The resampling strategy is

$$F(\hat{\theta}) \longrightarrow \{\mathbf{x}^{*b}\}_{b=1}^B \longrightarrow \hat{\theta}^*$$

and proceed as before to calculate e.g. $\hat{\text{se}}_{\text{boot}}$.



BOOTSTRAPPING IN REGRESSION

Bootstrap technique can be used for estimating prediction error in regression models. The specific version of bootstrap to be applied depends on the assumption on stochastic model generating the data. Specifically, whether random- \mathbf{X} or fixed- \mathbf{X} case is considered.

Assume the following:

- ▶ The learning data $\mathcal{D} = \{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$ are iid observations from the joint probability distribution of (\mathbf{X}, Y) (in $p + 1$ dimensional space) that are generated by the regression model,

$$Y = \beta_0 + \mathbf{X}'\beta + \varepsilon = \mu(\mathbf{X}) + \varepsilon,$$

$\mu(\mathbf{X}) = \mathbb{E}(Y|\mathbf{X}) = \beta_0 + \mathbf{X}'\beta$ is the true regression function, $\mathbb{E}(\varepsilon|\mathbf{X}) = 0$ and $\text{Var}(\varepsilon|\mathbf{X}) = \sigma^2$.

- ▶ Given \mathcal{D} we can estimate prediction error, PE by

$$\widehat{PE}(\hat{\mu}, \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\mu}(\mathbf{x}_i))^2 = \frac{RSS}{n}$$



BOOTSTRAPPING IN REGRESSION (CONT.)

- ▶ **Obs!** $\widehat{PE}(\hat{\mu}, \mathcal{D})$ is *apparent error rate* or resubstitution error rate for \mathcal{D} . $\widehat{PE}(\hat{\mu}, \mathcal{D})$ is computed by fitting OLS regression using the original data \mathcal{D} and then by applying that function to see how well it predicts those same members of \mathcal{D} .
- ▶ We expect that $\widehat{PE}(\hat{\mu}, \mathcal{D})$ will be too optimistic to estimate PE , i.e. $\widehat{PE}(\hat{\mu}, \mathcal{D}) < PE$.
- ▶ **Q:** How to improve the estimation of the prediction error?
- ▶ Computationally intensive correction strategy, i.e the bootstrap, that can be used to improve $\widehat{PE}(\hat{\mu}, \mathcal{D})$.
- ▶ **Obs!** With the joint distribution of (\mathbf{X}, Y) the resampling involves *sampling pairs* with replacement form $\{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$.



Bootstrap estimated prediction error

1. Draw a sample of (pairs) with replacement from \mathcal{D} , denote it by $\mathcal{D}^{*B} = \{(\mathbf{x}_i^{*b}, y_i^{*b}), i = 1, \dots, n\}_{b=1}^B$.
2. For each b regress y_i^{*b} on \mathbf{x}_i^{*b} to obtain OLS regression function $\hat{\mu}^{*b}(\mathbf{x})$.
3. Compute a *simple bootstrap estimator* of PE by applying $\hat{\mu}^{*b}$ to the *original* sample \mathcal{D} as

$$\widehat{PE}(\hat{\mu}^{*B}(\mathbf{x}), \mathcal{D}) = \frac{1}{B} \frac{1}{n} \sum_{b=1}^B \sum_{i=1}^n \left(y_i - \hat{\mu}^{*b}(\mathbf{x}_i) \right)^2.$$

4. Apply $\hat{\mu}^{*b}$ to \mathcal{D}^{*B} and compute

$$\widehat{PE}(\hat{\mu}^{*B}(\mathbf{x}), \mathcal{D}^{*B}) = \frac{1}{B} \frac{1}{n} \sum_{b=1}^B \sum_{i=1}^n \left(y_i^{*b} - \hat{\mu}^{*b}(\mathbf{x}_i) \right)^2.$$

5. Compute *the overall optimism*

$$opt^{*B} = \widehat{PE}(\hat{\mu}^{*B}(\mathbf{x}), \mathcal{D}) - \widehat{PE}(\hat{\mu}^{*B}(\mathbf{x}), \mathcal{D}^{*B})$$

6. Compute the bootstrap bias-corrected estimator of PE as

$$PR^* = \widehat{PE}(\hat{\mu}, \mathcal{D}) + opt^{*B}$$



- ▶ Observe that the simple bootstrap estimator $\widehat{PE}(\hat{\mu}^{*b}(\mathbf{x}), \mathcal{D})$ underestimates PE because there are observations *common to the bootstrap samples in \mathcal{D}^{*B} (operating as a learning data) and to the original data \mathcal{D} (operating as the test set)*.
- ▶ The probability that the observation i th (\mathbf{x}_i, y_i) from \mathcal{D} , is selected at least once to be in the b th bootstrap sample \mathcal{D}^{*b} is approximated when $n \rightarrow \infty$ as

$$\Pr\left((\mathbf{X}_i, Y_i) \in \mathcal{D}^{*b}\right) = 1 - \left(1 - \frac{1}{n}\right)^n \rightarrow 1 - e^{-1} \approx 0.632.$$

- ▶ On average, about 37% of observations in \mathcal{D} are left out of each bootstrap sample, which contains $0.632n$ distinct observations.
- ▶ One consequence: if p is close to n there could singularity or near singularity in $\mathcal{X}'\mathcal{X}$ ($\mathcal{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)'$ denotes the $n \times p$ design matrix).



BOOTSTRAPPING IN REGRESSION: FURTHER IMPROVEMENTS

- ▶ **Q:** How to improve PE^* ? Solution: The *leave-one-out bootstrap* strategy. Let $PE_{(1)}$ denote the expected bootstrap prediction error at those observations (\mathbf{x}_i, y_i) that are *not* included in the B bootstrap samples.
- ▶ In general, estimate $PE_{(i)}$ as $PE_{(i)}^* = \frac{1}{B_i} \sum_{b \in C_i} (y_i - \mu_b^*(\mathbf{x}_i))^2$ where C_i is the set of indices of the bootstrap samples that do not contain (\mathbf{x}_i, y_i) , and $B_i = |C_i|$ is the number of such bootstrap samples.
- ▶ The *0.632 bootstrap estimator of optimism* is given by

$$opt_{0.632}^B = 0.632(PE_{(1)}^* - \widehat{PE}(\hat{\mu}, \mathcal{D})).$$

- ▶ Replacing opt^{*B} with $opt_{0.632}^{*B}$ we obtain *0.632 bootstrap estimator of PE* as

$$PE_{0.632}^* = \widehat{PE}(\hat{\mu}, \mathcal{D}) + opt_{0.632}^{*B} = 0.368 \cdot \frac{RSS}{n} + 0.632 \cdot PE^*$$



- ▶ Non-parametric bootstrap
 - ▶ assumes that the data $\mathbf{x} = (x_1, \dots, x_n)$ are iid observations from an *unknown distribution*, i.e no assumption on F ;
 - ▶ is the resampling procedure where the simulated sample $\mathbf{x}^* = (x_1^*, \dots, x_n^*)$ is a random sample taken with replacement from the original data.
- ▶ Parametric bootstrap
 - ▶ assumes that the data $\mathbf{x} = (x_1, \dots, x_n)$ are iid observations from a specific parametric family F_θ .
 - ▶ resamples a *known distribution function*, whose parameters are estimated from the original data.

Some advantages and disadvantages of both approaches.

- ▶ In the non-parametric bootstrap, samples are drawn from a discrete set of n observations. For small n , might underestimate the amount of variation in the population we originally sampled. Generally, the samples of $n \leq 10$ are too small for obtaining reliable nonparametric bootstrap estimators.
- ▶ Parametric bootstrap usually requires less data to obtain reliable estimators due to stronger assumptions but is sensitive to the model choice. EDF (used in the non-parametric bootstrap) based on even quite small samples deviates little from that of its population one.
- ▶ ML estimators are commonly used in the parametric bootstrap. However their good properties nearly always are based upon their large sample behaviour.

Read more about bootstrap in e.g. *Bootstrap methods and their applications*, by A.C. Davison and D. V. Hinkley.



CROSS-VALIDATION: REVIEW OF THE LEARNING PROBLEM

Performance of a learning method naturally relates to its prediction accuracy on *independent* test data. We discuss *Cross-Validation*, CV, a simple and very efficient method for estimating prediction error. Initially derived by Stone (1974), but there are several versions.

- ▶ Let Y be a quantitative response variable, \mathbf{X} a vector of feature variables/predictors, and a prediction model $\hat{f}(\mathbf{X})$, which has been estimated using the training data $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$
- ▶ The common loss function for measuring errors (discrepancies) between Y and $\hat{f}(\mathbf{X})$ is the squared error

$$L(Y, \hat{f}(\mathbf{X})) = (Y - \hat{f}(\mathbf{X}))^2.$$

- ▶ We will be considering various types of estimation of the loss.



PERFORMANCE ASSESSMENT: TEST AND TRAINING ERRORS.

- ▶ *Test error*, also called for *generalization* error, is defined as

$$\text{Err}_{\mathcal{T}} = \mathbb{E} \left(L(Y, \hat{f}(\mathbf{X})) | \mathcal{T} \right).$$

- ▶ (\mathbf{X}, Y) is *new*, an independent (on \mathcal{T}) data drawn randomly from the joint population distribution of \mathbf{X} and Y , assumed to be the same as for \mathcal{T} .
 - ▶ The training set $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ is fixed.
- ▶ *Expected test error*, also known as expected prediction error, is

$$\text{Err} = \mathbb{E} \left[\mathbb{E} \left(L(Y, \hat{f}(\mathbf{X})) | \mathcal{T} \right) \right],$$

where the expectation averages over all the randomness, including that in the training set \mathcal{T} which produces the estimated model \hat{f} .



PERFORMANCE ASSESSMENT: TEST AND TRAINING ERRORS (CONT).

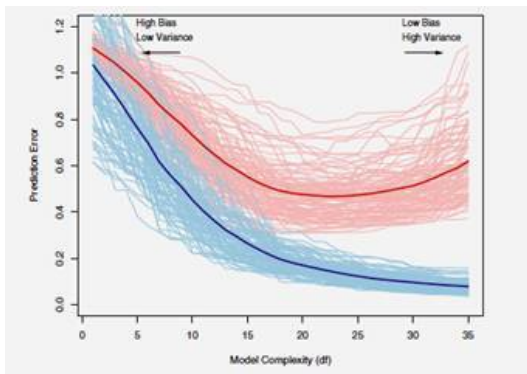
- ▶ *Training error* is also known as *apparent* or *re-substitution* error rate, is the average loss over the training sample

$$\overline{\text{err}} = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{f}(\mathbf{x}_i)).$$

- ▶ Is the most natural estimate of the error rate and shows how well we predict those *same* members of \mathcal{T} .
- ▶ For almost all learning procedures, $\overline{\text{err}}$ is readily available after training.
- ▶ Have misleadingly optimistic value since it estimates the predictive ability of the fitted model from the same $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ that was used to fit the model! But we are interested in the accuracy of the prediction which is obtained from applying the learned/fitted model to previously unseen test data.
- ▶ Generally, $\overline{\text{err}} \ll \text{Err}_{\mathcal{T}}$.
- ▶ This phenomena is presented on the figure below.



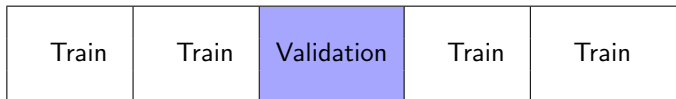
TRAINING- VS TEST-PERFORMANCE WITH VARYING MODEL COMPLEXITY



FIGUR: Behavior of the test and training error as a function of model complexity. Light blue curves represent a set of training errors. Solid blue curve represents the expected training error. Light red curves represent conditional test errors. Solid red curve represents the expected test error.

K-FOLD CROSS-VALIDATION: GENERIC IDEA

- ▶ Split randomly the set of available data into two parts: a *training set* and a *validation or hold-out set*
- ▶ For example, split into $K = 5$ roughly equal-size parts:



- ▶ Leave out part k (third above), fit model to the other $K - 1$ parts (combined together), specify the prediction error of the model fit when predicting the left-out k th part of the data.
- ▶ Perform this for $k = 1, 2, \dots, K$ and combine the K estimates (usually by averaging over splits) into validation-set error which provides an estimate of the *test error*.
- ▶ The beauty of CV is its simplicity and generality:
it allows for estimation of the test error rate by using (re-using) the training data!



K-FOLD CROSS-VALIDATION: MORE PRECISELY

- ▶ Let $\kappa : \{1, \dots, n\} \rightarrow \{1, \dots, K\}$ denote the indexing function that indicates the partition to which observation i is allocated by the randomization.
- ▶ Let $\hat{f}^{-k}(\mathbf{x})$ denote the function fitted with the k th part of the data removed.
- ▶ Then the *cross-validation* estimate of the prediction error is given by

$$\text{CV}_{(K)} = \frac{1}{K} \sum_{k=1}^K \sum_{(\mathbf{x}_i, y_i) \in \mathcal{T}_{\kappa(i)}} L(y_i, \hat{f}^{-\kappa(i)}(\mathbf{x}_i)).$$

- ▶ Usual choices for K are 5 or 10.
- ▶ The special case of $K = n$ is known as *leave-one-out* cross-validation, LOOCV. In this case, $\kappa(i) = i$, and for the i th observation the fit is computed using all the data except the i th.



K-FOLD CROSS-VALIDATION: MORE PRECISELY

- ▶ Given a set of models $f(x, \alpha)$ indexed by a tuning parameter α , denote by $\hat{f}^{-k}(\mathbf{x}, \alpha)$ the α th model fit with the k th part of the data removed.
- ▶ For this set of the models define

$$\text{CV}_{(K)}(\alpha) = \frac{1}{K} \sum_{k=1}^K \sum_{(\mathbf{x}_i, y_i) \in \mathcal{T}_{K(i)}} L(y_i, \hat{f}^{-k(i)}(\mathbf{x}_i, \alpha)).$$

- ▶ The function $\text{CV}(\alpha)$ provides an estimate of the test error curve; $\hat{\alpha}$, a minimizer of it, is the (CV) optimal value of the tuning parameter.
- ▶ The final chosen model is $f(x, \hat{\alpha})$ which we then fit to *all the data*.



K-FOLD CROSS-VALIDATION: MSE LOSS

- ▶ Let the K parts be C_1, \dots, C_k where C_k is the set of indices of the observations in part k , there are n_k obs. in part k .
- ▶ Assume that n is a multiple of K , so that $n_k = n/K$ and consider linear fitting under squared-error loss

$$L(\mathbf{y}, \hat{f}^{-\kappa(i)}(\mathbf{x})) = \text{MSE}^{-k} = \frac{1}{n_k} \sum_{i \in C_k} (y_i - \hat{y}_i)^2,$$

where \hat{y}_i is the fit for observation i obtained from the data with part k removed. Then

$$\text{CV}(\hat{f}) = \sum_{k=1}^K \frac{n_k}{n} \text{MSE}^{-k}$$

- ▶ Setting $K = n$ yields n -fold of leave-one out, LOOCV.



CV ON CLASSIFICATION PROBLEM

- ▶ Let Y be a qualitative response variable and \mathbf{X} a vector of feature variables/predictors. This corresponds to the classification model with $Y \in \mathcal{C} = \{1, \dots, c\}$ given classes. The training data set $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$. In the classification setting we then assign an observation \mathbf{x} to the class c for which

$$\hat{y}(\mathbf{x}) = \hat{y}_c = \arg \max_{c \in \mathcal{C}} \mathbb{P}(Y = c | \mathbf{X} = \mathbf{x}).$$

- ▶ 0-1 loss function, $L(Y, \hat{Y}(\mathbf{X})) = I(Y \neq \hat{Y}(\mathbf{X}))$, can be used to quantifying the accuracy of the classifier.
- ▶ Then the *training* error is $\overline{\text{err}} = \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$ with averaging over \mathcal{T} ,
- ▶ and the *test* error is $\text{Err}_{\mathcal{T}} = E(L(Y, \hat{Y}(\mathbf{X}) | \mathcal{T}))$.
- ▶ The test error associated with a *new* set of observations (independent on \mathcal{T} !) $\mathcal{T}_{\text{new}} = \{(\mathbf{x}_{\text{new}}, y_{\text{new}})\}$ is estimated by

$$\text{Err} = \text{Ave}(\text{Err}_i) \quad \text{where} \quad \text{Err}_i = I(y_i \neq \hat{y}_i) \text{ for } \mathcal{T}_{\text{new}}.$$



CV ON CLASSIFICATION PROBLEM (CONT.)

- ▶ **K -fold CV.** Split the data \mathcal{T} randomly into K parts, $\mathcal{T}_{C_1}, \dots, \mathcal{T}_{C_K}$ where C_k is the set of indices of the observations in part k and let $n_k = n/K$.
 - ▶ For each k compute $\text{Err}_k = \sum_{i \in \mathcal{T}_{C_k}} I(y_i \neq \hat{y}_i) / n_k$, n_k is the number of obs. in \mathcal{T}_{C_k} .
 - ▶ Compute $\text{CV}_{(K)} = \sum_{k=1}^K \frac{n_k}{n} \text{Err}_k$.
- ▶ Example from ISL: Logistic regression with a non-linear (quadratic) decision boundary.

$$\log \left(\frac{p}{1-p} \right) = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \beta_3 X_2 + \beta_2 X_2^2.$$

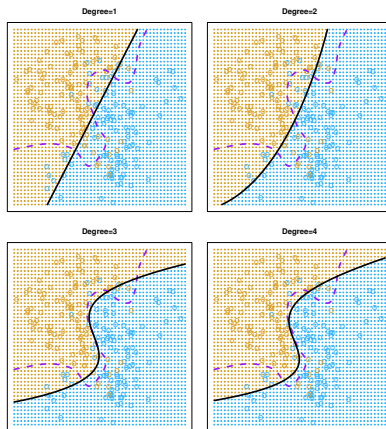
- ▶ For the simulated data, we can compute the *true* test error rate which is 0.201 and use the Bayes error rate, which is 0.133 as a benchmark. For real data, Bayes decision boundary and test error rates are unknown.

Q: How to choose between the four logistic models?

A: use CV to make decision.



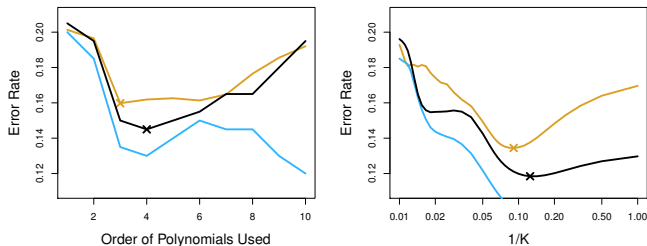
CV FOR CLASSIFICATION USING LOGISTIC REGRESSION



FIGUR: Simulated two-class bivariate data along with the logistic regression fit and Bayesian decision boundary (purple dashed line). Estimated classification boundaries from polynomial (degrees 1 to 4) logistic regression with the test errors estimated as 0.201, 0.197, 0.160 and 0.162, respectively for each polynomial function. Bayes (irreducible) error rate is 0.133.



CV FOR CLASSIFICATION USING LOGISTIC REGRESSION (CONT.)



FIGUR: Test error (brown), training error (blue) and 10-fold CV error (black). Left: polynomial logistic regression based classifier. Right: KNN-classifier (K here is the number of neighbors, not a number of folds).

RIGHT AND WRONG WAY TO PERFORM CROSS-VALIDATION

Example: Consider a two-class classification problem. Assume that the number of feature variables is $p = 5000$ and total sample size is $n = 50$, i.e. need to construct a classifier for a very high-dimensional case.

One possible strategy is to first select a subset of informative variables and then validate the performance of the resulting classifier, we call it for

Strategy I:

1. **Screen** the feature variables for their discriminative power. Select, say, 100 variables having largest correlation with the class label.
2. **Construct a classifier** based on the selected subset of highly discriminative features.
3. **Perform cross-validation** to estimate the unknown tuning parameters and to estimate Err of the final classifier.

Q: What is wrong with Strategy I ?



RIGHT AND WRONG WAY TO PERFORM CROSS-VALIDATION

Answer: The discriminative feature variables were chosen after seeing all the data!

Correct way is **Strategy II:**

1. **Split** the samples into K folds randomly.
2. For each fold $k = 1, \dots, K$
 - ▶ **Find** a subset of discriminative feature variables using all the samples except the k th.
 - ▶ **Construct the classifier** using all the samples except the k th. **Apply** the classifier to predict the class labels for the samples in k th fold.
 - ▶ **Select** the final model which minimizes CV error.

General approach to performing multi-step model fitting with CV:

- ▶ CV must be applied to the **entire sequence of model fitting steps**.
- ▶ Any supervised procedure which has seen the labels of the training data and used them, is a form of *training* and must be included in the validation procedure.

