

Topic 6: Bias-Variance Problem

6.1	Bias-Variance Tradeoff for a Predictive Model	141
6.2	Bias-Variance Tradeoff for Data Samples	142
6.3	Training vs Testing Errors	144
6.4	Resampling Methods	145
6.4.1	Cross-Validation	145
6.4.2	Bootstrap Approach	147

According to Section 7.5 of <https://bookdown.org/jefftemplewebb/IS-6489/>, **overfitting** is a major hazard in predictive analytics, especially when using machine learning algorithms like random forest (Topic 5) which, without proper tuning, can learn sample data almost perfectly, essentially fitting noise. When such a model is used to predict new data, with different noises, model performance can be shockingly bad. Overfitting error can cause very funny result as pointed out by Dr Kilian Q. Weinberger (<http://kilian.cs.cornell.edu/>) in his license plate recognition example(?) Cornell CS4780 SP17 lecture (<https://www.youtube.com/watch?v=zj-5nkNKAow>).

In this topic, we explore the difference between a predictive model and the true model using the notion of bias-variance tradeoff and the *generalisation error*, which measures how accurately a predictive model is able to predict outcome values for previously unseen data. The generalisation error is related to overfitting issue. Since it is impossible to know the bias-variance and the generalisation errors except for very rare special cases, we need resampling methods to estimate them. Two of the most commonly used resampling methods are the *cross-validation* (Section 6.4.1) and the *bootstrap*.

6.1 Bias-Variance Tradeoff for a Predictive Model

If we have a predictive model $h(x)$ and an true model $f(x) + \epsilon$ (with $\mathbb{E}_{X,\epsilon}[\epsilon] = 0$), the mean square error (MSE) is usually used to measure their difference for the regression problem:

$$MSE = \mathbb{E}_{X,\epsilon}[(h(X) - (f(X) + \epsilon))^2]. \quad (6.1)$$

Note that the true model $f(x) + \epsilon$ have a range and we don't just measure the difference between the predictive model and the true model but we measure the average differences over the range of x .

Using probability theory, (6.1) can be expanded to

$$\begin{aligned} MSE &= \mathbb{E}_{X,\epsilon}[(h(X) - f(X))^2 + 2(h(X) - f(X))\epsilon + \epsilon^2] \\ &= \mathbb{E}_X[(h(X) - f(X))^2] + 2\mathbb{E}_X(h(X) - f(X))\mathbb{E}_{X,\epsilon}[\epsilon] + \mathbb{E}_{X,\epsilon}[\epsilon^2] \\ &= \underbrace{\mathbb{E}_X[(h(X) - f(X))^2]}_{\text{bias error}} + \underbrace{\mathbb{E}_{X,\epsilon}[\epsilon^2]}_{\text{error variance}}. \end{aligned}$$

The bias-variance for a predictive model can be applied to physics or chemistry experiments where ϵ is the measurement or physical noise and $h(x)$ and $f(x)$ may be the same. It can also be applied to study the expected test error of a trained model h_D .

6.2 Bias-Variance Tradeoff for Data Samples

In machine learning, the predictive model $h(D, x)$ is usually a function of the training data D and input features x , therefore, the bias-variance analysis in the previous section is insufficient.

According to Wikipedia, the bias-variance tradeoff of a predictive model is a property that the variance of the parameter estimated across samples can be reduced by increasing the bias in the estimated parameters. The bias-variance dilemma is the conflict in trying to simultaneously minimise these two sources of error that prevent predictive models from generalising beyond their training set:

- the **bias error** is an error from erroneous assumptions in the predictive model. High bias can cause the predictive model to miss the relevant relations between features and target outputs (underfitting).
- The **variance** is an error from sensitivity to small fluctuations in the training set. High variance may result from the predictive model which models the random noise in the training data (overfitting).

In comparison to the MSE in the previous section, the MSE for predictive models with sampling can be decomposed into three terms — the bias, variance, and irreducible error.

Statisticians assume that the training set D , consisting of n inputs, are drawn i.i.d. from some distribution P , which is denoted as $D \sim P^n$. The machine learning algorithm \mathcal{A} learns on the data set D a predictive model $h(D, x) = \mathcal{A}(D)$. The MSE can be extended to the *Expected Test Error*:

$$ETE = \mathbb{E}_{\substack{(\mathbf{x}, y) \sim P \\ D \sim P^n}} [(h(D, \mathbf{x}) - y)^2] = \int_D \int_{\mathbf{x}} \int_y (h(D, \mathbf{x}) - y)^2 \mathbb{P}(\mathbf{x}, y) \mathbb{P}(D) d\mathbf{x} dy dD. \quad (6.2)$$

Decomposition of Equation 6.2:

$$\begin{aligned} ETE &= \mathbb{E}_{\mathbf{x}, y, D} \left[\left[(h(D, \mathbf{x}) - \bar{h}(\mathbf{x})) + (\bar{h}(\mathbf{x}) - y) \right]^2 \right] \\ &= \mathbb{E}_{\mathbf{x}, D} \left[(h(D, \mathbf{x}) - \bar{h}(\mathbf{x}))^2 \right] + 2 \mathbb{E}_{\mathbf{x}, y, D} \left[(h(D, \mathbf{x}) - \bar{h}(\mathbf{x})) (\bar{h}(\mathbf{x}) - y) \right] + \mathbb{E}_{\mathbf{x}, y} \left[(\bar{h}(\mathbf{x}) - y)^2 \right]. \end{aligned} \quad (6.3)$$

The middle term of the above equation is 0 as we show below

$$\begin{aligned} \mathbb{E}_{\mathbf{x}, y, D} \left[(h(D, \mathbf{x}) - \bar{h}(\mathbf{x})) (\bar{h}(\mathbf{x}) - y) \right] &= \mathbb{E}_{\mathbf{x}, y} \left[\mathbb{E}_D \left[h(D, \mathbf{x}) - \bar{h}(\mathbf{x}) \right] (\bar{h}(\mathbf{x}) - y) \right] \\ &= \mathbb{E}_{\mathbf{x}, y} \left[(\mathbb{E}_D [h(D, \mathbf{x})] - \bar{h}(\mathbf{x})) (\bar{h}(\mathbf{x}) - y) \right] \\ &= \mathbb{E}_{\mathbf{x}, y} \left[(\bar{h}(\mathbf{x}) - \bar{h}(\mathbf{x})) (\bar{h}(\mathbf{x}) - y) \right] = \mathbb{E}_{\mathbf{x}, y} [0] = 0 \end{aligned}$$

Note that \bar{h} is a weighted average of trained functions on various data D .

$$\bar{h}(\mathbf{x}) = \mathbb{E}_{D \sim P^n} [h(D, \mathbf{x})] = \int_D h(D, \mathbf{x}) \mathbb{P}(D) dD$$

where $\mathbb{P}(D)$ is the probability of drawing D from P^n .

Returning to the earlier expression (6.3), we're left with the variance and another term

$$\mathbb{E}_{\mathbf{x}, y, D} \left[(h(D, \mathbf{x}) - y)^2 \right] = \underbrace{\mathbb{E}_{\mathbf{x}, D} \left[(h(D, \mathbf{x}) - \bar{h}(\mathbf{x}))^2 \right]}_{\text{Variance}} + \mathbb{E}_{\mathbf{x}, y} \left[(\bar{h}(\mathbf{x}) - y)^2 \right] \quad (6.4)$$

Let the expected label $\bar{y}(\mathbf{x}) = \mathbb{E}_{y|\mathbf{x}}[y(\mathbf{x})] = \int y \mathbb{P}(y|\mathbf{x})dy$. We can break down the second term in the equation (6.4) as follows:

$$\mathbb{E}_{\mathbf{x},y} \left[(\bar{h}(\mathbf{x}) - y)^2 \right] = \mathbb{E}_{\mathbf{x},y} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) + (\bar{y}(\mathbf{x}) - y)^2 \right] \quad (6.5)$$

$$= \underbrace{\mathbb{E}_{\mathbf{x},y} \left[(\bar{y}(\mathbf{x}) - y)^2 \right]}_{\text{Noise}} + \underbrace{\mathbb{E}_{\mathbf{x}} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))^2 \right]}_{\text{Bias}^2} + 2 \mathbb{E}_{\mathbf{x},y} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) (\bar{y}(\mathbf{x}) - y) \right] \quad (6.6)$$

The third term in the equation above is 0, as we show below

$$\begin{aligned} \mathbb{E}_{\mathbf{x},y} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) (\bar{y}(\mathbf{x}) - y) \right] &= \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{y|\mathbf{x}} [\bar{y}(\mathbf{x}) - y] (\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{y|\mathbf{x}} [\bar{y}(\mathbf{x}) - y] (\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[(\bar{y}(\mathbf{x}) - \mathbb{E}_{y|\mathbf{x}} [y]) (\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) \right] \\ &= \mathbb{E}_{\mathbf{x}} \left[(\bar{y}(\mathbf{x}) - \bar{y}(\mathbf{x})) (\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x})) \right] = \mathbb{E}_{\mathbf{x}} [0] = 0. \end{aligned}$$

This gives us the decomposition of expected test error (6.2) as follows

$$\underbrace{\mathbb{E}_{\mathbf{x},y,D} \left[(h(D, \mathbf{x}) - y)^2 \right]}_{\text{ExpectedTestError}} = \underbrace{\mathbb{E}_{\mathbf{x},D} \left[(h(D, \mathbf{x}) - \bar{h}(\mathbf{x}))^2 \right]}_{\text{Variance}} + \underbrace{\mathbb{E}_{\mathbf{x},y} \left[(\bar{y}(\mathbf{x}) - y)^2 \right]}_{\text{Noise}} + \underbrace{\mathbb{E}_{\mathbf{x}} \left[(\bar{h}(\mathbf{x}) - \bar{y}(\mathbf{x}))^2 \right]}_{\text{Bias}^2}. \quad (6.7)$$

The “Variance” captures how much the classifier changes if we train on a different training set. How “over-specialised” is the classifier to a particular training set (i.e. is it overfitting)? If we have the best possible model for our training data, how far off are we from the average classifier?

The data-intrinsic “noise” measures the ambiguity due to our data distribution and feature representation. We can never beat this, it is an aspect of the data.

The inherent error that we obtain from the classifier even with infinite training data is captured by the “bias”, which is due to the classifier being “biased” to a particular kind of solution (e.g. linear classifier). In other words, “bias” is inherent to our model.

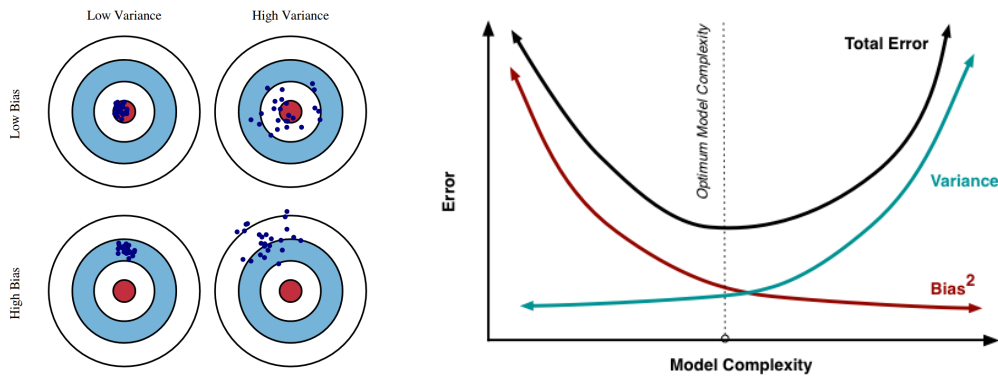


Figure 6.1: (a) Graphical illustration of bias and variance with the output y being a 2D vector. (b) The variation of Bias and Variance with the model complexity. This is similar to the concept of overfitting and underfitting. More complex models overfit while the simplest models underfit. Source: <http://scott.fortmann-roe.com/docs/BiasVariance.html>

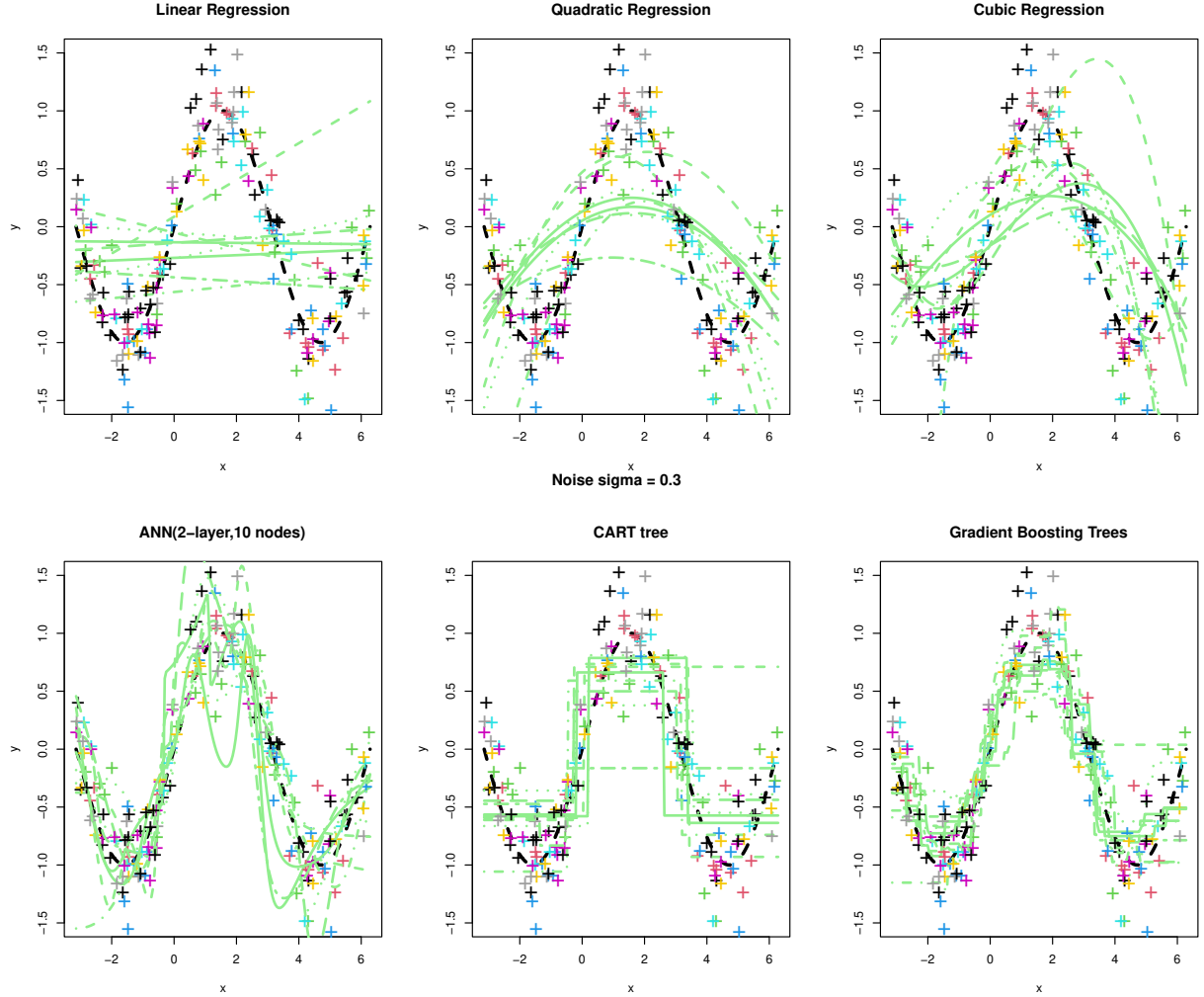
Remark 6.2.1. The reason to limit ourselves to regression problems is due to the many shortcomings of bias-variance decomposition theory for 0-1 loss function according to http://rasbt.github.io/mlxtend/user_guide/evaluate/bias_variance_decomp/.

Example 6.2.2. Consider a similar problem to Example 1.7.6 below:

$$y = \sin(x) + R, \quad -\pi < x < 2\pi, \quad R \sim \text{Normal}(0, 0.3^2).$$

By using nine data with 20 samples D_i , $i = 1, \dots, 9$, we train the predictive models linear regression to gradient boosting trees and we can see the less flexible models like linear regression, quadratic regression, etc. have large bias while the flexible models like ANN, CART tree and gradient boosting tree can have lower bias but have large variance.

Noise sigma = 0.3



□

6.3 Training vs Testing Errors

In the bias-variance tradeoff analysis, the data set D is assumed to be infinite while in Example 6.2.2, we only consider nine data set with only 20 samples. However, in real-world situation, we usually need to split the data to training and testing data. If a classifier is under-performing (e.g. if either of the test or training error is too high), there are several ways to improve the performance. To find out which of these techniques is the right one for the situation, the first step is to determine the root of the problem.

The graph in Figure plots the training error and the test error and can be divided into two overarching regimes. In the first regime (on the left side of the graph), training error is below the desired error threshold (denoted by ϵ), but test error is significantly higher. In the second regime (on the right side of the graph), test error is remarkably close to training error, but both are above the desired tolerance of ϵ .

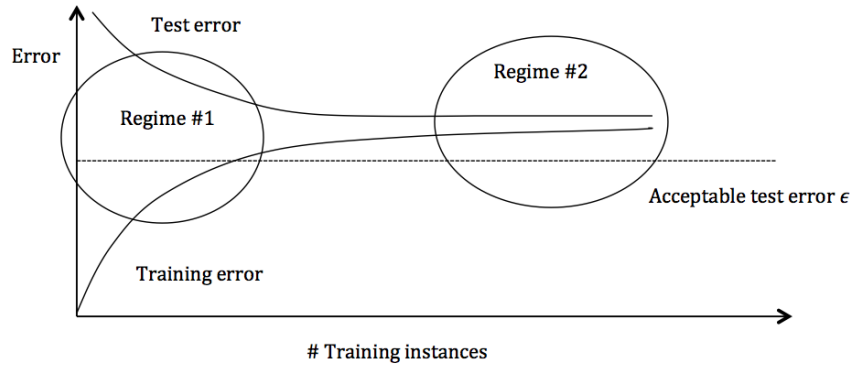


Figure 6.2: Test and training error as the number of training instances increases.

Regime 1 (High Variance)

In the first regime, the cause of the poor performance is “high variance”.

Symptoms	Remedies
Training error is much lower than test error	Add more training data
Training error is lower than ϵ	Reduce model complexity — complex models are prone to high variance
Test error is above ϵ	Bagging

Regime 2 (High Bias)

Unlike the first regime, the second regime indicates high bias: the model being used is not robust enough to produce an accurate prediction.

Symptoms	Remedies
Training error is higher than ϵ	<ul style="list-style-type: none"> • Use more complex model (e.g. kernelize, use non-linear models) • Add features • Boosting

6.4 Resampling Methods

Resampling methods are statistical techniques for estimating statistics from samples. In particular, they can be used to estimate the generalisation error for predictive models. Two most popular resampling methods are the cross-validation approach and the bootstrap approach.

6.4.1 Cross-Validation

Suppose the data D is used to train a supervised learning model \hat{h}_D and suppose that we have a new test data

$$(\mathbf{x}_1^{new}, y_1^{new}), \dots, (\mathbf{x}_L^{new}, y_L^{new}) \sim P. \quad (6.8)$$

The general concept associated with “test error” is the *loss function* of individual test data $(\mathbf{x}_i^{new}, y_i^{new})$, i.e.

$$\frac{1}{L} \sum_{i=1}^L \ell(y_i^{new}, \hat{h}_D(\mathbf{x}_i^{new})) \rightarrow \mathbb{E}[\ell(Y^{new}, \hat{h}_D(X^{new}))] \quad (6.9)$$

where the right is the *theoretical test set error* for a predictive model training on a particular training set h_D . The *theoretical generalisation error* is defined as the average of theoretical test set error over all possible training sets, i.e.

$$\mathbb{E}_D \left[\mathbb{E}_{(X^{new}, Y^{new})} [\ell(Y^{new}, \hat{h}_D(X^{new}))] \right]. \quad (6.10)$$

For classification problems, the usual form of test error (6.9) and ℓ is

$$\frac{1}{L} \sum_{i=1}^L I(y_i^{new} \neq \hat{h}_D(\mathbf{x}_i^{new})).$$

For regression problems, the usual form of test error (6.9) and ℓ is

$$\frac{1}{L} \sum_{i=1}^L (y_i^{new} - \hat{h}_D(\mathbf{x}_i^{new}))^2.$$

A true test data (6.8) is not available at the time we want to run and tune our supervised learning algorithm. Therefore, we need artificial ways to construct training and testing datasets:

1. Holdout method/Split-test validation (Section 1.8.5): shuffle the data and split into two sets (train-test) or three sets (train-validate-test for hyperparameter selection). The performance measure obtained is dependent on the sampling/shuffling;
2. Random subsampling: average out multiple holdout method;
3. K-fold cross-validation (CV) and Leave-one-out cross-validation (LOOCV) (Section 2.9): the shuffle data is splitted into k equal partitions (when $k = n$, the CV is called LOOCV) and use each partition as test data and the remainder as training data. From this, we can get k performance metrics to estimate theoretical generalisation error which would usually be less bias compare to the holdout method.

In holdout method, the estimation of the generalisation error (6.10) is

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i^{new}, \hat{h}_D(\mathbf{x}_i^{new})).$$

To estimate a confidence interval for accuracy in holdout method, we apply the binomial experiment and the normal distribution approximation. The confidence interval of accuracy, $acc = X/n$ with mean p and variance $p(1-p)/n$ follows:

$$P\left(-Z_{\alpha/2} \leq \frac{acc - p}{\sqrt{p(1-p)/n}} \leq Z_{1-\alpha/2}\right) = 1 - \alpha$$

where $-Z_{\alpha/2}$ and $Z_{1-\alpha/2}$ are the lower and upper bounds are related a standard normal distribution at a confidence level $1 - \alpha$.

Consider a model with an accuracy of 80% when evaluated on 100 test records. For a confidence level of $\alpha = 0.05$, $Z_{\alpha/2} = Z_{1-\alpha/2} = 1.96$ and the confidence interval is (may contain error in calculation, need checking):

$$\begin{aligned} & \frac{2 \times n \times acc + Z_{\alpha/2}^2 \pm Z_{\alpha/2} \sqrt{Z_{\alpha/2}^2 + 4n(acc - acc^2)}}{2(n + Z_{\alpha/2}^2)} \\ &= \frac{2 \times 100 \times 0.8 + 1.96^2 \pm 1.96 \sqrt{1.96^2 + 4 \times 100 \times (0.8 - 0.8^2)}}{2(100 + 1.96^2)} = [0.7246, 0.8830] \end{aligned}$$

To compare the performance of two models M_1 and M_2 , we suppose they are evaluated on two independent test sets, D_1 and D_2 . Suppose the error rate for M_1 on D_1 is e_1 and the error rate for M_2 on D_2 is e_2 . If $|D_1|$ and $|D_2|$ are sufficiently large, e_1, e_2 follows normal distributions. Then $d = e_1 - e_2$ also follows normal distribution and the $(1 - \alpha)$ confidence

interval is

$$|e_1 - e_2| \pm Z_{\alpha/2} \sqrt{\frac{e_1(1-e_1)}{|D_1|} + \frac{e_2(1-e_2)}{|D_2|}}.$$

Random subsampling (or Monte Carlo crossvalidation according to R.R. Picard and R.D. Cook, J. Am. Stat. Assoc., 79 (1984) 575–583) repeat the holdout method several times (with different random samples) to improve the estimation of the classifier’s performance. In contrast to a full CV method, random subsampling has been shown to be asymptotically consistent (J. Shao, J. Am. Stat. Assoc., 88 (1993) 486–494.) resulting in more pessimistic predictions of the test data compared with CV.

In k-fold CV, suppose that the data D are randomly partition their indices $\{1, \dots, n\}$ into $\{B_k\}$ such that $\bigcup_{k=1}^K B_k = \{1, \dots, n\}$ and $B_j \cap B_k = \emptyset$ ($j \neq k$). If a partition with B_k removed is used to produce a training model $\hat{h}_{n-|B_k|}^{(-B_k)}$. The estimation of the generalisation error (6.10) is

$$\mathbb{E}_D \left[\mathbb{E}_{(X^{new}, Y^{new})} [\ell(Y^{new}, \hat{h}_{n-|B_k|}^{(-B_k)}(X^{new}))] \right] \approx \frac{1}{K} \sum_{i=1}^K \frac{1}{|B_k|} \sum_{i \in B_k} \ell(y_i^{new}, \hat{h}_{n-|B_k|}^{(-B_k)}(\mathbf{x}_i^{new})).$$

In LOOCV, we used the i th sample point as test data and the remaining $n-1$ sample points as training data to train a predictive model $\hat{h}_{n-1}^{(-i)}$. The estimation of the generalisation error (6.10) is

$$\mathbb{E}_D \left[\mathbb{E}_{(X^{new}, Y^{new})} [\ell(Y^{new}, \hat{h}_{n-1}^{(-i)}(X^{new}))] \right] \approx \frac{1}{n} \sum_{i=1}^n \ell(y_i^{new}, \hat{h}_{n-1}^{(-i)}(\mathbf{x}_i^{new})).$$

Since training sample D is close to n , the bias is small. The variance is usually high because n training sets are similar to each other due to strong correlation:

$$\text{Var} \left(\frac{1}{n} \sum_{i=1}^n \ell(y_i^{new}, \hat{h}_{n-1}^{(-i)}(\mathbf{x}_i^{new})) \right) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \underbrace{\text{Cov} \left(\ell(y_i^{new}, \hat{h}_{n-1}^{(-i)}(\mathbf{x}_i^{new})), \ell(y_j^{new}, \hat{h}_{n-1}^{(-j)}(\mathbf{x}_j^{new})) \right)}_{\text{may be large due to correlation}}.$$

Note that this is not true in general as P. Burman (1989)’s paper “A comparative study of ordinary cross-validation, v -fold cross-validation and the repeated learning-testing methods” proved that for least squares linear regression, LOOCV has the smallest asymptotic bias and variance.

6.4.2 Bootstrap Approach

According to [https://en.wikipedia.org/wiki/Bootstrapping_\(statistics\)](https://en.wikipedia.org/wiki/Bootstrapping_(statistics)), *bootstrapping* is any test or metric that uses **random sampling with replacement**. The bootstrap is a method for estimating the variance of an estimator and for finding approximate confidence intervals for parameters.

Let $X_1, \dots, X_n \sim P$. Recall that the empirical distribution P_n is defined by

$$P_n(A) = \frac{1}{n} \sum_{i=1}^n I(X_i \in A).$$

In other words, P_n puts mass $1/n$ at each X_i . A parameter of the form $\theta = T(P)$ is called a *statistical function* and that the plug-in estimator is $\hat{\theta}_n = T(P_n)$.

An iid sample of size n drawn from P_n is called a *bootstrap* sample, denoted by

$$X_1^*, \dots, X_n^* \sim P$$

Bootstrap samples play an important role in what follows. Note that drawing an iid sample X_1^*, \dots, X_n^* from P_n is equivalent to drawing n observations, with replacement, from the original data $\{X_1, \dots, X_n\}$. Thus, bootstrap sampling is often described as “resampling the data”.

Let $\hat{\theta}_n = g(X_1, \dots, X_n)$ denote some estimator. We would like to find the variance of $\hat{\theta}_n$. Let

$$\text{Var}_P(\hat{\theta}_n) = \text{Var}_P(g(X_1, \dots, X_n)) =: S_n(P).$$

Let s^2 be the sample variance of $\hat{\theta}_n^{(1)}, \dots, \hat{\theta}_n^{(B)}$. By law of large numbers,

$$s^2 = \frac{1}{B} \sum_{j=1}^n (\hat{\theta}_n^{(j)})^2 - \left(\frac{1}{B} \sum_{j=1}^n \hat{\theta}_n^{(j)} \right)^2 \xrightarrow{P} \mathbb{E}[\hat{\theta}_n^2] - \mathbb{E}[\hat{\theta}_n]^2 = \text{Var}_P(\hat{\theta}_n).$$

Since we can take B as large as we want, we have that $s^2 \approx \text{Var}_P(\hat{\theta}_n)$. In other words, we can **approximate $S_n(P)$ by repeatedly simulating n observations from P .**

But we don't know P . So we estimate $S_n(P)$ with $S_n(P_n)$ where P_n is the empirical distribution. Since P_n is a *consistent estimator*, we expect that $S_n(P_n) \approx S_n(P)$. In other words, bootstrap approximation of the variance,

$$\widehat{\text{Var}_P(\hat{\theta}_n)} = \text{Var}_{P_n}(\hat{\theta}_n).$$

The steps for the Bootstrap Variance Estimator algorithm are

1. Draw a bootstrap sample $X_1^*, \dots, X_n^* \sim P_n$. Compute $\hat{\theta}_n^* = g(X_1^*, \dots, X_n^*)$.
2. Repeat the previous step, B times, yielding estimators $\hat{\theta}_{n,1}^*, \dots, \hat{\theta}_{n,B}^*$.
3. Compute (using population variance rather than sample variance as in James et al. [2013, SS5.2]):

$$\hat{s} = \sqrt{\frac{1}{B} \sum_{j=1}^B (\hat{\theta}_{n,j}^* - \bar{\theta})^2}$$

$$\text{where } \bar{\theta} = \frac{1}{B} \sum_{j=1}^B \hat{\theta}_{n,j}^*.$$

4. Output \hat{s} .

A very simple illustration is given in Practical 2 to estimate the standard deviation of the population.