

GR 5243 Project 3: Overview of predictive modeling

Tian Zheng

March 10, 2017

Part I: Predictive modeling / Supervised Learning

(Supervised) Predictive Learning

- ▶ Y the outcome variable is given.
- ▶ Learning task I: using original or derived X to construct a model to predict Y .
- ▶ Learning task II: Identify subject-matter knowledge from the derived model.
- ▶ For categorical Y , it is referred to as a *Classification* problem.

(Supervised) Predictive Learning

- ▶ Y the outcome variable is given.
- ▶ Learning task I: using original or derived X to construct a model to predict Y .
- ▶ Learning task II: Identify subject-matter knowledge from the derived model.
- ▶ For categorical Y , it is referred to as a *Classification* problem.

(Supervised) Predictive Learning

- ▶ Y the outcome variable is given.
- ▶ Learning task I: using original or derived X to construct a model to predict Y .
- ▶ Learning task II: Identify subject-matter knowledge from the derived model.
- ▶ For categorical Y , it is referred to as a *Classification* problem.

(Supervised) Predictive Learning

- ▶ Y the outcome variable is given.
- ▶ Learning task I: using original or derived X to construct a model to predict Y .
- ▶ Learning task II: Identify subject-matter knowledge from the derived model.
- ▶ For categorical Y , it is referred to as a *Classification* problem.

Classification

Most classification methods can be discussed using

- ▶ Decision boundaries between classes;
 - ▶ Linear classifier versus non-linear classifier.
- ▶ Discriminant functions: for each class k , define $\hat{f}_k(x)$.
 - ▶ Prediction may then be computed as: $\underset{k}{\operatorname{argmax}} \hat{f}_k(x)$
 - ▶ Most current methods work in this domain.

Classification

Most classification methods can be discussed using

- ▶ Decision boundaries between classes;
 - ▶ **Linear classifier versus non-linear classifier.**
- ▶ Discriminant functions: for each class k , define $\hat{f}_k(x)$.
 - ▶ Prediction may then be computed as: $\underset{k}{\operatorname{argmax}} \hat{f}_k(x)$
 - ▶ Most current methods work in this domain.

Classification

Most classification methods can be discussed using

- ▶ Decision boundaries between classes;
 - ▶ Linear classifier versus non-linear classifier.
- ▶ Discriminant functions: for each class k , define $\hat{f}_k(x)$.
 - ▶ Prediction may then be computed as: $\underset{k}{\operatorname{argmax}} \hat{f}_k(x)$
 - ▶ Most current methods work in this domain.

Classification

Most classification methods can be discussed using

- ▶ Decision boundaries between classes;
 - ▶ Linear classifier versus non-linear classifier.
- ▶ Discriminant functions: for each class k , define $\hat{f}_k(x)$.
 - ▶ Prediction may then be computed as: $\underset{k}{\operatorname{argmax}} \hat{f}_k(x)$
 - ▶ Most current methods work in this domain.

Classification

Most classification methods can be discussed using

- ▶ Decision boundaries between classes;
 - ▶ Linear classifier versus non-linear classifier.
- ▶ Discriminant functions: for each class k , define $\hat{f}_k(x)$.
 - ▶ Prediction may then be computed as: $\underset{k}{\operatorname{argmax}} \hat{f}_k(x)$
 - ▶ Most current methods work in this domain.

- └ Predictive modeling: supervised Learning
- └ Decision theory for “model building” or “learning”

Part I.1: Decision theory for “model building” or “learning”

- └ Predictive modeling: supervised Learning
- └ Decision theory for “model building” or “learning”

Loss function

- ▶ Can be viewed as a distance metric between the observed data (Y) and the models ($E(Y|X) = f(X)$).
- ▶ For example, squared error loss (quadratic loss)

$$L(Y, f(X)) = (Y - f(X))^2$$

- ▶ Define *expected prediction error* (EPE)

$$\begin{aligned} \text{EPE}(f) &= E(Y - f(X))^2 \\ &= E_X E_{Y|X}((Y - f(X))^2 | X) \end{aligned}$$

- ▶ Under the above loss function, $f(x) = E(Y|X = x)$ minimizes $\text{EPE}(f)$, which is known as the *regression function*.

Loss function

- ▶ Can be viewed as a distance metric between the observed data (Y) and the models ($E(Y|X) = f(X)$).
- ▶ For example, squared error loss (quadratic loss)

$$L(Y, f(X)) = (Y - f(X))^2$$

- ▶ Define *expected prediction error* (EPE)

$$\begin{aligned} \text{EPE}(f) &= E(Y - f(X))^2 \\ &= E_X E_{Y|X}((Y - f(X))^2 | X) \end{aligned}$$

- ▶ Under the above loss function, $f(x) = E(Y|X = x)$ minimizes $\text{EPE}(f)$, which is known as the *regression function*.

Loss function

- ▶ Can be viewed as a distance metric between the observed data (Y) and the models ($E(Y|X) = f(X)$).
- ▶ For example, squared error loss (quadratic loss)

$$L(Y, f(X)) = (Y - f(X))^2$$

- ▶ Define *expected prediction error* (EPE)

$$\begin{aligned} \text{EPE}(f) &= E(Y - f(X))^2 \\ &= E_X E_{Y|X}((Y - f(X))^2 | X) \end{aligned}$$

- ▶ Under the above loss function, $f(x) = E(Y|X = x)$ minimizes $\text{EPE}(f)$, which is known as the *regression function*.

Loss function

- ▶ Can be viewed as a distance metric between the observed data (Y) and the models ($E(Y|X) = f(X)$).
- ▶ For example, squared error loss (quadratic loss)

$$L(Y, f(X)) = (Y - f(X))^2$$

- ▶ Define *expected prediction error* (EPE)

$$\begin{aligned} \text{EPE}(f) &= E(Y - f(X))^2 \\ &= E_X E_{Y|X}((Y - f(X))^2 | X) \end{aligned}$$

- ▶ Under the above loss function, $f(x) = E(Y|X = x)$ minimizes $\text{EPE}(f)$, which is known as the *regression function*.

- └ Predictive modeling: supervised Learning
- └ Decision theory for “model building” or “learning”

Different strategies

- ▶ Least squares method is trying to minimize the observed instance of $L(Y, f(X))$ (among the linear models) to estimate the f that minimizes $E(L(Y, f(X)))$.
- ▶ Nearest-neighbor method uses the average (mean) of $y_i | x_i \in N_k(x)$ to estimate $f(x) = E(Y | X = x)$. Here, *conditioning on x* was implemented in the form of estimation within a small neighborhood, $N_k(x)$, of x . $E(\cdot)$ is estimated using average.
More observations \Rightarrow more observations in the neighborhood of x , the tighter the $N_k(x)$, and the more stable $\hat{f}(x)$ is, given $E(Y | X)$ is continuous enough.
- ▶ Actually, it is not difficult to show that if $N, K \rightarrow \infty$ and $K/N \rightarrow 0$, then $\hat{f}(x) \rightarrow E(Y | X = x)$.

Different strategies

- ▶ Least squares method is trying to minimize the observed instance of $L(Y, f(X))$ (among the linear models) to estimate the f that minimizes $E(L(Y, f(X)))$.
- ▶ Nearest-neighbor method uses the average (mean) of $y_i | x_i \in N_k(x)$ to estimate $f(x) = E(Y | X = x)$. Here, *conditioning on x* was implemented in the form of estimation within a small neighborhood, $N_k(x)$, of x . $E(\cdot)$ is estimated using average.
More observations \Rightarrow more observations in the neighborhood of x , the tighter the $N_k(x)$, and the more stable $\hat{f}(x)$ is, given $E(Y | X)$ is continuous enough.
- ▶ Actually, it is not difficult to show that if $N, K \rightarrow \infty$ and $K/N \rightarrow 0$, then $\hat{f}(x) \rightarrow E(Y | X = x)$.

Different strategies

- ▶ Least squares method is trying to minimize the observed instance of $L(Y, f(X))$ (among the linear models) to estimate the f that minimizes $E(L(Y, f(X)))$.
- ▶ Nearest-neighbor method uses the average (mean) of $y_i | x_i \in N_k(x)$ to estimate $f(x) = E(Y | X = x)$. Here, *conditioning on x* was implemented in the form of estimation within a small neighborhood, $N_k(x)$, of x . $E(\cdot)$ is estimated using average.
More observations \Rightarrow more observations in the neighborhood of x , the tighter the $N_k(x)$, and the more stable $\hat{f}(x)$ is, given $E(Y | X)$ is continuous enough.
- ▶ Actually, it is not difficult to show that if $N, K \rightarrow \infty$ and $K/N \rightarrow 0$, then $\hat{f}(x) \rightarrow E(Y | X = x)$.

- └ Predictive modeling: supervised Learning
 - └ Decision theory for “model building” or “learning”

Another loss function

- ▶ Replace L_2 metric by L_1 metric.
- ▶ $L(Y, f(X)) = |Y - f(X)|$
- ▶ It can be shown that $f(x) = \text{median}(Y|X = x)$ minimizes EPE.
- ▶ Regression: least absolute deviations (LAD) regression or LAR (least absolute residuals).
- ▶ KNN: use median instead of average for each $N_k(x)$.

- └ Predictive modeling: supervised Learning
- └ Decision theory for “model building” or “learning”

Another loss function

- ▶ Replace L_2 metric by L_1 metric.
- ▶ $L(Y, f(X)) = |Y - f(X)|$
- ▶ It can be shown that $f(x) = \text{median}(Y|X = x)$ minimizes EPE.
- ▶ Regression: least absolute deviations (LAD) regression or LAR (least absolute residuals).
- ▶ KNN: use median instead of average for each $N_k(x)$.

- └ Predictive modeling: supervised Learning
- └ Decision theory for “model building” or “learning”

Another loss function

- ▶ Replace L_2 metric by L_1 metric.
- ▶ $L(Y, f(X)) = |Y - f(X)|$
- ▶ It can be shown that $f(x) = \text{median}(Y|X = x)$ minimizes EPE.
- ▶ Regression: least absolute deviations (LAD) regression or LAR (least absolute residuals).
- ▶ KNN: use median instead of average for each $N_k(x)$.

- └ Predictive modeling: supervised Learning
- └ Decision theory for “model building” or “learning”

Another loss function

- ▶ Replace L_2 metric by L_1 metric.
- ▶ $L(Y, f(X)) = |Y - f(X)|$
- ▶ It can be shown that $f(x) = \text{median}(Y|X = x)$ minimizes EPE.
- ▶ Regression: least absolute deviations (LAD) regression or LAR (least absolute residuals).
- ▶ KNN: use median instead of average for each $N_k(x)$.

Another loss function

- ▶ Replace L_2 metric by L_1 metric.
- ▶ $L(Y, f(X)) = |Y - f(X)|$
- ▶ It can be shown that $f(x) = \text{median}(Y|X = x)$ minimizes EPE.
- ▶ Regression: least absolute deviations (LAD) regression or LAR (least absolute residuals).
- ▶ KNN: use median instead of average for each $N_k(x)$.

- └ Predictive modeling: supervised Learning
- └ Decision theory for “model building” or “learning”

Loss function for categorical outcome G

- ▶ The loss function is defined based on a penalty matrix, $\mathbf{L} = [L(k, l)]_{K \times K}$ for G with K classes.
- ▶ Zero-one loss is corresponding to the number of misclassifications, and defined by $L(k, l) = 1$ if $k \neq l$; 0 otherwise.
- ▶ Define $\text{EPE} = \mathbb{E}_X \sum_{k=1}^K L(g_k, \hat{G}(X)) \Pr(g_k | X)$.
- ▶ With 0-1 loss, it can be shown that

$$\hat{G}(x) = g_k \text{ if } \Pr(g_k | X = x) = \max_{g \in G} \Pr(g | X = x)$$

minimizes EPE.

Loss function for categorical outcome G

- ▶ The loss function is defined based on a penalty matrix, $\mathbf{L} = [L(k, l)]_{K \times K}$ for G with K classes.
- ▶ Zero-one loss is corresponding to the number of misclassifications, and defined by $L(k, l) = 1$ if $k \neq l$; 0 otherwise.
- ▶ Define $\text{EPE} = \mathbb{E}_X \sum_{k=1}^K L(g_k, \hat{G}(X)) \Pr(g_k | X)$.
- ▶ With 0-1 loss, it can be shown that

$$\hat{G}(x) = g_k \text{ if } \Pr(g_k | X = x) = \max_{g \in G} \Pr(g | X = x)$$

minimizes EPE.

Loss function for categorical outcome G

- ▶ The loss function is defined based on a penalty matrix, $\mathbf{L} = [L(k, l)]_{K \times K}$ for G with K classes.
- ▶ Zero-one loss is corresponding to the number of misclassifications, and defined by $L(k, l) = 1$ if $k \neq l$; 0 otherwise.
- ▶ Define $\text{EPE} = \mathbb{E}_X \sum_{k=1}^K L(g_k, \hat{G}(X)) \Pr(g_k | X)$.
- ▶ With 0-1 loss, it can be shown that

$$\hat{G}(x) = g_k \text{ if } \Pr(g_k | X = x) = \max_{g \in G} \Pr(g | X = x)$$

minimizes EPE.

Loss function for categorical outcome G

- ▶ The loss function is defined based on a penalty matrix, $\mathbf{L} = [L(k, l)]_{K \times K}$ for G with K classes.
- ▶ Zero-one loss is corresponding to the number of misclassifications, and defined by $L(k, l) = 1$ if $k \neq l$; 0 otherwise.
- ▶ Define $\text{EPE} = \mathbb{E}_X \sum_{k=1}^K L(g_k, \hat{G}(X)) \Pr(g_k | X)$.
- ▶ With 0-1 loss, it can be shown that

$$\hat{G}(x) = g_k \text{ if } \Pr(g_k | X = x) = \max_{g \in G} \Pr(g | X = x)$$

minimizes EPE.

- └ Predictive modeling: supervised Learning
 - └ Decision theory for “model building” or “learning”

Bayes Classifier

- ▶ “Classify to the most probable class, using the conditional distribution $\Pr(G|X)$ ”
- ▶ Error rate of this classifier is called the *Bayes rate*.
- ▶ KNN attempts to estimate the Bayes classifier by estimating $\Pr(G|X = x)$ by $\hat{\Pr}(G|N_k(x))$.

Part I.2: Discriminant analysis methods

Linear discriminant analysis

- ▶ LDA assumes a Gaussian mixture with common variance-covariance matrix.
- ▶ $P(G|X) = \{\Pr(g_k|X=x) \propto \Pr(x|g_k)\pi_k\}$
- ▶ It is then sufficient to look at the log-ratio

$$\log \frac{\Pr(g_k|X=x)}{\Pr(g_l|X=x)} = \log \frac{\pi_k}{\pi_l} - \frac{1}{2}(\mu_k + \mu_l)^T \Sigma^{-1}(\mu_k - \mu_l) + x^T \Sigma^{-1}(\mu_k - \mu_l)$$

Linear discriminant analysis

- ▶ LDA assumes a Gaussian mixture with common variance-covariance matrix.
- ▶ $P(G|X) = \{\Pr(g_k|X = x) \propto \Pr(x|g_k)\pi_k\}$
- ▶ It is then sufficient to look at the log-ratio

$$\log \frac{\Pr(g_k|X = x)}{\Pr(g_l|X = x)} = \log \frac{\pi_k}{\pi_l} - \frac{1}{2}(\mu_k + \mu_l)^T \Sigma^{-1}(\mu_k - \mu_l) + x^T \Sigma^{-1}(\mu_k - \mu_l)$$

Linear discriminant analysis

- ▶ LDA assumes a Gaussian mixture with common variance-covariance matrix.
- ▶ $P(G|X) = \{\Pr(g_k|X = x) \propto \Pr(x|g_k)\pi_k\}$
- ▶ It is then sufficient to look at the log-ratio

$$\log \frac{\Pr(g_k|X = x)}{\Pr(g_l|X = x)} = \log \frac{\pi_k}{\pi_l} - \frac{1}{2}(\mu_k + \mu_l)^T \Sigma^{-1}(\mu_k - \mu_l) + x^T \Sigma^{-1}(\mu_k - \mu_l)$$

Linear discriminant analysis

- ▶ The *linear discriminant function* is then

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

and $\hat{G}(x) = \underset{k}{\operatorname{argmax}} \delta_k(x)$.

- ▶ Computational considerations

- ▶ $\hat{\pi}_k = N_k / N$

- ▶ $\hat{\mu}_k = \bar{X}_k$

- ▶ $\hat{\Sigma}$ is the pooled estimate of the variance-covariance matrix.

Linear discriminant analysis

- ▶ The *linear discriminant function* is then

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

and $\hat{G}(x) = \underset{k}{\operatorname{argmax}} \delta_k(x)$.

- ▶ **Computational considerations**

- ▶ $\hat{\pi}_k = N_k / N$

- ▶ $\hat{\mu}_k = \bar{X}_k$

- ▶ $\hat{\Sigma}$ is the pooled estimate of the variance-covariance matrix.

Linear discriminant analysis

- ▶ The *linear discriminant function* is then

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

and $\hat{G}(x) = \underset{k}{\operatorname{argmax}} \delta_k(x)$.

- ▶ Computational considerations

- ▶ $\hat{\pi}_k = N_k / N$

- ▶ $\hat{\mu}_k = \bar{X}_k$

- ▶ $\hat{\Sigma}$ is the pooled estimate of the variance-covariance matrix.

Linear discriminant analysis

- ▶ The *linear discriminant function* is then

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

and $\hat{G}(x) = \underset{k}{\operatorname{argmax}} \delta_k(x)$.

- ▶ Computational considerations

- ▶ $\hat{\pi}_k = N_k / N$

- ▶ $\hat{\mu}_k = \bar{X}_k$

- ▶ $\hat{\Sigma}$ is the pooled estimate of the variance-covariance matrix.

Linear discriminant analysis

- ▶ The *linear discriminant function* is then

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

and $\hat{G}(x) = \underset{k}{\operatorname{argmax}} \delta_k(x)$.

- ▶ Computational considerations

- ▶ $\hat{\pi}_k = N_k / N$

- ▶ $\hat{\mu}_k = \bar{X}_k$

- ▶ $\hat{\Sigma}$ is the pooled estimate of the variance-covariance matrix.

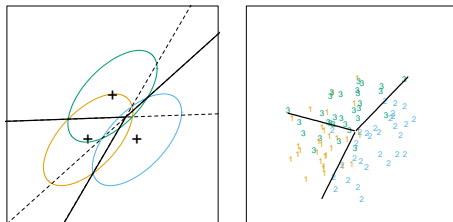


FIGURE 4.5. The left panel shows three Gaussian distributions, with the same covariance and different means. Included are the contours of constant density enclosing 95% of the probability in each case. The Bayes decision boundaries between each pair of classes are shown (broken straight lines), and the Bayes decision boundaries separating all three classes are the thicker solid lines (a subset of the former). On the right we see a sample of 30 drawn from each Gaussian distribution, and the fitted LDA decision boundaries.

Quadratic discriminant functions

- ▶ If we assume different variance-covariance matrices Σ_k for the classes $k = 1, \dots, K$, the discriminant function based on log-likelihood-ratio will be quadratic.
- ▶ That is,

$$\delta_k(x) = \log \pi_k - \frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)$$

- ▶ Thus, the discriminant boundary $(\{x : \delta_k(x) = \delta_l(x)\})$ is a quadratic function.

Quadratic discriminant functions

- ▶ If we assume different variance-covariance matrices Σ_k for the classes $k = 1, \dots, K$, the discriminant function based on log-likelihood-ratio will be quadratic.
- ▶ That is,

$$\delta_k(x) = \log \pi_k - \frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)$$

- ▶ Thus, the discriminant boundary $(\{x : \delta_k(x) = \delta_l(x)\})$ is a quadratic function.

Quadratic discriminant functions

- ▶ If we assume different variance-covariance matrices Σ_k for the classes $k = 1, \dots, K$, the discriminant function based on log-likelihood-ratio will be quadratic.
- ▶ That is,

$$\delta_k(x) = \log \pi_k - \frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)$$

- ▶ Thus, the discriminant boundary ($\{x : \delta_k(x) = \delta_l(x)\}$) is a quadratic function.

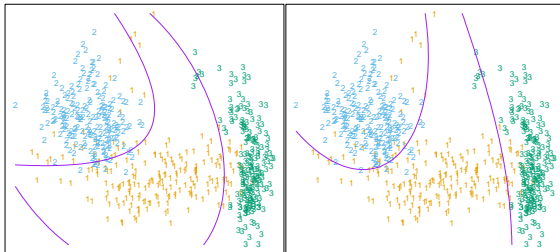


FIGURE 4.6. Two methods for fitting quadratic boundaries. The left plot shows the quadratic decision boundaries for the data in Figure 4.1 (obtained using LDA in the five-dimensional space $X_1, X_2, X_1X_2, X_1^2, X_2^2$). The right plot shows the quadratic decision boundaries found by QDA. The differences are small, as is usually the case.

The method of nearest neighbors

- ▶ Define $N_k(x)$ as the neighborhood of x , s.t. it contains k nearest points in the training set.
- ▶ Here, we implicitly implied a metric of distance is needed for KNN methods.
- ▶ The most popular choice is the *Euclidean distance*.
- ▶ The prediction is then $\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$.
- ▶ Obviously RSS minimizes at $k = 1$.
- ▶ Model fitted using KNN is less rigid than the linear methods such as LDA and generate non-linear prediction functions.

The method of nearest neighbors

- ▶ Define $N_k(x)$ as the neighborhood of x , s.t. it contains k nearest points in the training set.
- ▶ Here, we implicitly implied a metric of distance is needed for KNN methods.
- ▶ The most popular choice is the *Euclidean distance*.
- ▶ The prediction is then $\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$.
- ▶ Obviously RSS minimizes at $k = 1$.
- ▶ Model fitted using KNN is less rigid than the linear methods such as LDA and generate non-linear prediction functions.

The method of nearest neighbors

- ▶ Define $N_k(x)$ as the neighborhood of x , s.t. it contains k nearest points in the training set.
- ▶ Here, we implicitly implied a metric of distance is needed for KNN methods.
- ▶ The most popular choice is the *Euclidean distance*.
- ▶ The prediction is then $\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$.
- ▶ Obviously RSS minimizes at $k = 1$.
- ▶ Model fitted using KNN is less rigid than the linear methods such as LDA and generate non-linear prediction functions.

The method of nearest neighbors

- ▶ Define $N_k(x)$ as the neighborhood of x , s.t. it contains k nearest points in the training set.
- ▶ Here, we implicitly implied a metric of distance is needed for KNN methods.
- ▶ The most popular choice is the *Euclidean distance*.
- ▶ The prediction is then $\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$.
- ▶ Obviously RSS minimizes at $k = 1$.
- ▶ Model fitted using KNN is less rigid than the linear methods such as LDA and generate non-linear prediction functions.

The method of nearest neighbors

- ▶ Define $N_k(x)$ as the neighborhood of x , s.t. it contains k nearest points in the training set.
- ▶ Here, we implicitly implied a metric of distance is needed for KNN methods.
- ▶ The most popular choice is the *Euclidean distance*.
- ▶ The prediction is then $\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$.
- ▶ Obviously RSS minimizes at $k = 1$.
- ▶ Model fitted using KNN is less rigid than the linear methods such as LDA and generate non-linear prediction functions.

The method of nearest neighbors

- ▶ Define $N_k(x)$ as the neighborhood of x , s.t. it contains k nearest points in the training set.
- ▶ Here, we implicitly implied a metric of distance is needed for KNN methods.
- ▶ The most popular choice is the *Euclidean distance*.
- ▶ The prediction is then $\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$.
- ▶ Obviously RSS minimizes at $k = 1$.
- ▶ Model fitted using KNN is less rigid than the linear methods such as LDA and generate non-linear prediction functions.

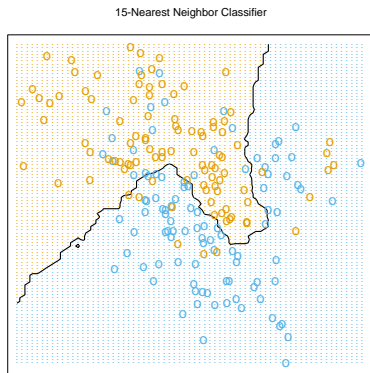


FIGURE 2.2. The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (BLUE = 0, ORANGE = 1) and then fit by 15-nearest-neighbor averaging as in (2.8). The predicted class is hence chosen by majority vote amongst the 15-nearest neighbors.

Properties

Linear methods: low variance, high bias

- ▶ Unbiased under the linear model. Biases if the data follow a nonlinear model.
- ▶ Stable, smooth. (i.e., individual observations are not very *influential*.)
- ▶ Nice analytical results

Nearest Neighbor methods: high variance, low bias

- ▶ No assumption on the model.
- ▶ High variance in the predictions since each \hat{Y} is calculated only by a few observations.

Properties

Linear methods: low variance, high bias

- ▶ Unbiased under the linear model. Biases if the data follow a nonlinear model.
- ▶ Stable, smooth. (i.e., individual observations are not very *influential*.)
- ▶ Nice analytical results

Nearest Neighbor methods: high variance, low bias

- ▶ No assumption on the model.
- ▶ High variance in the predictions since each \hat{Y} is calculated only by a few observations.

Properties

Linear methods: low variance, high bias

- ▶ Unbiased under the linear model. Biases if the data follow a nonlinear model.
- ▶ Stable, smooth. (i.e., individual observations are not very *influential*.)
- ▶ Nice analytical results

Nearest Neighbor methods: high variance, low bias

- ▶ No assumption on the model.
- ▶ High variance in the predictions since each \hat{Y} is calculated only by a few observations.

Properties

Linear methods: low variance, high bias

- ▶ Unbiased under the linear model. Biases if the data follow a nonlinear model.
- ▶ Stable, smooth. (i.e., individual observations are not very *influential*.)
- ▶ Nice analytical results

Nearest Neighbor methods: high variance, low bias

- ▶ No assumption on the model.
- ▶ High variance in the predictions since each \hat{Y} is calculated only by a few observations.

Properties

Linear methods: low variance, high bias

- ▶ Unbiased under the linear model. Biases if the data follow a nonlinear model.
- ▶ Stable, smooth. (i.e., individual observations are not very *influential*.)
- ▶ Nice analytical results

Nearest Neighbor methods: high variance, low bias

- ▶ No assumption on the model.
- ▶ High variance in the predictions since each \hat{Y} is calculated only by a few observations.

- └ Predictive modeling: supervised Learning
- └ Discriminant analysis methods

“Each method has its own situations for which it works best.”

- └ Predictive modeling: supervised Learning
- └ Discriminant analysis methods

Variants of these two strategies

- ▶ Kernel methods can be applied to distance metrics.
- ▶ Smoother kernel to replace the kNN “kernel”.
- ▶ Using local weights, or weights vary across dimensions to make linear methods less rigid.
- ▶ Linear models fit to a basis expansion of the original inputs: this expands the class of models considered.

- └ Predictive modeling: supervised Learning
- └ Discriminant analysis methods

Variants of these two strategies

- ▶ Kernel methods can be applied to distance metrics.
- ▶ Smoother kernel to replace the kNN “kernel”.
- ▶ Using local weights, or weights vary across dimensions to make linear methods less rigid.
- ▶ Linear models fit to a basis expansion of the original inputs: this expands the class of models considered.

- └ Predictive modeling: supervised Learning
- └ Discriminant analysis methods

Variants of these two strategies

- ▶ Kernel methods can be applied to distance metrics.
- ▶ Smoother kernel to replace the kNN “kernel”.
- ▶ Using local weights, or weights vary across dimensions to make linear methods less rigid.
- ▶ Linear models fit to a basis expansion of the original inputs: this expands the class of models considered.

Variants of these two strategies

- ▶ Kernel methods can be applied to distance metrics.
- ▶ Smoother kernel to replace the kNN “kernel”.
- ▶ Using local weights, or weights vary across dimensions to make linear methods less rigid.
- ▶ Linear models fit to a basis expansion of the original inputs: this expands the class of models considered.

Part II: Curse of dimensionality

Assumptions made in classification

- ▶ If you have infinite data points and unlimited computational power and storage, classification is easy.
- ▶ What *Bayes rate* tell us about theoretical limit of performance?
- ▶ For finite size training sample, you don't have *enough* observations everywhere you need to make prediction.
- ▶ Assumptions
 - ▶ Probability of class labels are continuous over feature values.
 - ▶ The distance metric or kernel function used are meaningful for the classification problem.
 - ▶ The test sample were drawn from the same distributions as the training sample.

Assumptions made in classification

- ▶ If you have infinite data points and unlimited computational power and storage, classification is easy.
- ▶ What *Bayes rate* tell us about theoretical limit of performance?
- ▶ For finite size training sample, you don't have *enough* observations everywhere you need to make prediction.
- ▶ Assumptions
 - ▶ Probability of class labels are continuous over feature values.
 - ▶ The distance metric or kernel function used are meaningful for the classification problem.
 - ▶ The test sample were drawn from the same distributions as the training sample.

Assumptions made in classification

- ▶ If you have infinite data points and unlimited computational power and storage, classification is easy.
- ▶ What *Bayes rate* tell us about theoretical limit of performance?
- ▶ For finite size training sample, you don't have *enough* observations everywhere you need to make prediction.
- ▶ Assumptions
 - ▶ Probability of class labels are continuous over feature values.
 - ▶ The distance metric or kernel function used are meaningful for the classification problem.
 - ▶ The test sample were drawn from the same distributions as the training sample.

Assumptions made in classification

- ▶ If you have infinite data points and unlimited computational power and storage, classification is easy.
- ▶ What *Bayes rate* tell us about theoretical limit of performance?
- ▶ For finite size training sample, you don't have *enough* observations everywhere you need to make prediction.
- ▶ **Assumptions**
 - ▶ Probability of class labels are continuous over feature values.
 - ▶ The distance metric or kernel function used are meaningful for the classification problem.
 - ▶ The test sample were drawn from the same distributions as the training sample.

Assumptions made in classification

- ▶ If you have infinite data points and unlimited computational power and storage, classification is easy.
- ▶ What *Bayes rate* tell us about theoretical limit of performance?
- ▶ For finite size training sample, you don't have *enough* observations everywhere you need to make prediction.
- ▶ Assumptions
 - ▶ Probability of class labels are continuous over feature values.
 - ▶ The distance metric or kernel function used are meaningful for the classification problem.
 - ▶ The test sample were drawn from the same distributions as the training sample.

Assumptions made in classification

- ▶ If you have infinite data points and unlimited computational power and storage, classification is easy.
- ▶ What *Bayes rate* tell us about theoretical limit of performance?
- ▶ For finite size training sample, you don't have *enough* observations everywhere you need to make prediction.
- ▶ Assumptions
 - ▶ Probability of class labels are continuous over feature values.
 - ▶ The distance metric or kernel function used are meaningful for the classification problem.
 - ▶ The test sample were drawn from the same distributions as the training sample.

Assumptions made in classification

- ▶ If you have infinite data points and unlimited computational power and storage, classification is easy.
- ▶ What *Bayes rate* tell us about theoretical limit of performance?
- ▶ For finite size training sample, you don't have *enough* observations everywhere you need to make prediction.
- ▶ Assumptions
 - ▶ Probability of class labels are continuous over feature values.
 - ▶ The distance metric or kernel function used are meaningful for the classification problem.
 - ▶ The test sample were drawn from the same distributions as the training sample.

High dimensions

- ▶ Let p be the dimension of the input space (usually, number of predictors.)
- ▶ Assume that the inputs uniformly distributed in a p -dimensional unit cube.
- ▶ Consider the K nearest neighborhood. Let $r = K/N$. The average size of the neighborhood can be approximated by a p -dimensional hyper-sphere with radius d , so that its volume πd^p equals r .
- ▶ Therefore, $d = \sqrt[p]{r/\pi}$.

High dimensions

- ▶ Let p be the dimension of the input space (usually, number of predictors.)
- ▶ Assume that the inputs uniformly distributed in a p -dimensional unit cube.
- ▶ Consider the K nearest neighborhood. Let $r = K/N$. The average size of the neighborhood can be approximated by a p -dimensional hyper-sphere with radius d , so that its volume πd^p equals r .
- ▶ Therefore, $d = \sqrt[p]{r/\pi}$.

High dimensions

- ▶ Let p be the dimension of the input space (usually, number of predictors.)
- ▶ Assume that the inputs uniformly distributed in a p -dimensional unit cube.
- ▶ Consider the K nearest neighborhood. Let $r = K/N$. The average size of the neighborhood can be approximated by a p -dimensional hyper-sphere with radius d , so that its volume πd^p equals r .
- ▶ Therefore, $d = \sqrt[p]{r/\pi}$.

High dimensions

- ▶ Let p be the dimension of the input space (usually, number of predictors.)
- ▶ Assume that the inputs uniformly distributed in a p -dimensional unit cube.
- ▶ Consider the K nearest neighborhood. Let $r = K/N$. The average size of the neighborhood can be approximated by a p -dimensional hyper-sphere with radius d , so that its volume πd^p equals r .
- ▶ Therefore, $d = \sqrt[p]{r/\pi}$.

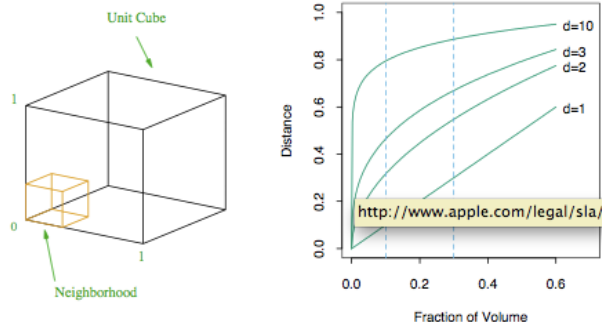


FIGURE 2.6. The curse of dimensionality is well illustrated by a subcubical neighborhood for uniform data in a unit cube. The figure on the right shows the side-length of the subcube needed to capture a fraction r of the volume of the data, for different dimensions p . In ten dimensions we need to cover 80% of the range of each coordinate to capture 10% of the data.

Curse of dimensionality

- ▶ Again, let p be the dimension of the input space (usually, number of predictors.) Assume that the inputs uniformly distributed in a p -dimensional unit cube.
- ▶ For a random point $\{x_1, \dots, x_p\}$, $x_i \sim \text{Unif}(0, 1)$, iid.
- ▶ It can be shown that $E \min(x_i) = \frac{1}{p+1}$.
- ▶ This implies that for any point, it is very close to at least one boundary. Inference at the boundary is usually difficult.

Curse of dimensionality

- ▶ Again, let p be the dimension of the input space (usually, number of predictors.) Assume that the inputs uniformly distributed in a p -dimensional unit cube.
- ▶ For a random point $\{x_1, \dots, x_p\}$, $x_i \sim \text{Unif}(0, 1)$, iid.
- ▶ It can be shown that $E \min(x_i) = \frac{1}{p+1}$.
- ▶ This implies that for any point, it is very close to at least one boundary. Inference at the boundary is usually difficult.

Curse of dimensionality

- ▶ Again, let p be the dimension of the input space (usually, number of predictors.) Assume that the inputs uniformly distributed in a p -dimensional unit cube.
- ▶ For a random point $\{x_1, \dots, x_p\}$, $x_i \sim \text{Unif}(0, 1)$, iid.
- ▶ It can be shown that $E \min(x_i) = \frac{1}{p+1}$.
- ▶ This implies that for any point, it is very close to at least one boundary. Inference at the boundary is usually difficult.

Curse of dimensionality

- ▶ Again, let p be the dimension of the input space (usually, number of predictors.) Assume that the inputs uniformly distributed in a p -dimensional unit cube.
- ▶ For a random point $\{x_1, \dots, x_p\}$, $x_i \sim \text{Unif}(0, 1)$, iid.
- ▶ It can be shown that $E \min(x_i) = \frac{1}{p+1}$.
- ▶ This implies that for any point, it is very close to at least one boundary. Inference at the boundary is usually difficult.

Variance-Bias decomposition

Take the expected prediction error (EPE) for L_2 loss

$$\begin{aligned} E_{y|x}(y - \hat{y})^2 &= E_{y|x}(y - E(y|x) + E(y|x) - E(\hat{y}|x) + E(\hat{y}|x) - \hat{y})^2 \\ &= \text{var}(y|x) + (E(y|x) - E(\hat{y}|x))^2 + \text{var}(\hat{y}|x) \\ &= \text{Irreducible error} + \text{Bias}^2 + \text{var}(\hat{y}|x) \end{aligned}$$

How dimensions affect estimation

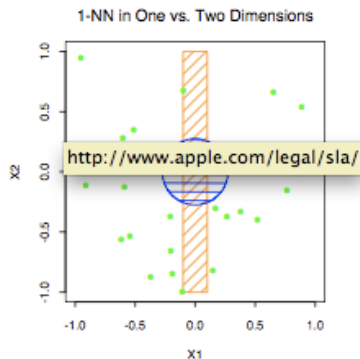
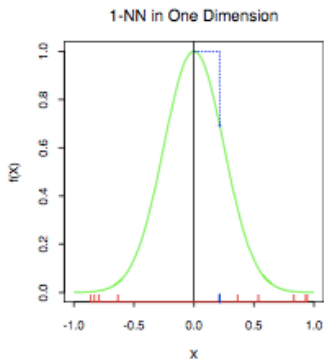
- ▶ Consider X in p dimensional space.
- ▶ $Y = e^{-8\|X\|^2}$.
- ▶ We are using nearest-neighbor method to estimate Y at $X = 0$.

How dimensions affect estimation

- ▶ Consider X in p dimensional space.
- ▶ $Y = e^{-8\|X\|^2}$.
- ▶ We are using nearest-neighbor method to estimate Y at $X = 0$.

How dimensions affect estimation

- ▶ Consider X in p dimensional space.
- ▶ $Y = e^{-8\|X\|^2}$.
- ▶ We are using nearest-neighbor method to estimate Y at $X = 0$.



How dimensions affect estimation

- ▶ As the number of dimensions increase, the distance of the nearest neighbor to $X = 0$ increases.
- ▶ The nearest neighbor estimate is therefore down-biased.
- ▶ The function $Y = f(X)$ is symmetric about different dimensions of X .
- ▶ Actually, it only depends on the distance between the nearest neighbor at $X = 0$.
- ▶ The variability of the estimate is then decided by the variability of the distance to NN.

How dimensions affect estimation

- ▶ As the number of dimensions increase, the distance of the nearest neighbor to $X = 0$ increases.
- ▶ The nearest neighbor estimate is therefore down-biased.
- ▶ The function $Y = f(X)$ is symmetric about different dimensions of X .
- ▶ Actually, it only depends on the distance between the nearest neighbor at $X = 0$.
- ▶ The variability of the estimate is then decided by the variability of the distance to NN.

How dimensions affect estimation

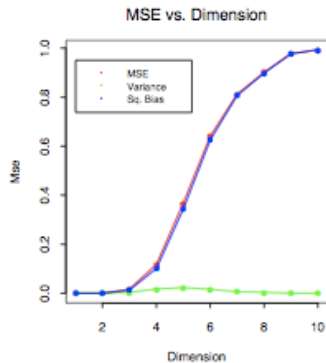
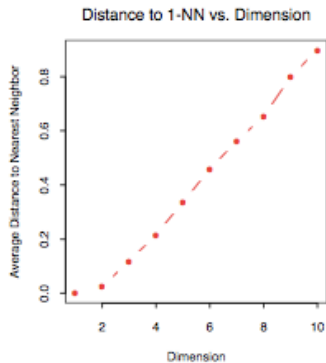
- ▶ As the number of dimensions increase, the distance of the nearest neighbor to $X = 0$ increases.
- ▶ The nearest neighbor estimate is therefore down-biased.
- ▶ The function $Y = f(X)$ is symmetric about different dimensions of X .
- ▶ Actually, it only depends on the distance between the nearest neighbor at $X = 0$.
- ▶ The variability of the estimate is then decided by the variability of the distance to NN.

How dimensions affect estimation

- ▶ As the number of dimensions increase, the distance of the nearest neighbor to $X = 0$ increases.
- ▶ The nearest neighbor estimate is therefore down-biased.
- ▶ The function $Y = f(X)$ is symmetric about different dimensions of X .
- ▶ Actually, it only depends on the distance between the nearest neighbor at $X = 0$.
- ▶ The variability of the estimate is then decided by the variability of the distance to NN.

How dimensions affect estimation

- ▶ As the number of dimensions increase, the distance of the nearest neighbor to $X = 0$ increases.
- ▶ The nearest neighbor estimate is therefore down-biased.
- ▶ The function $Y = f(X)$ is symmetric about different dimensions of X .
- ▶ Actually, it only depends on the distance between the nearest neighbor at $X = 0$.
- ▶ The variability of the estimate is then decided by the variability of the distance to NN.



How dimensions affect estimation

- ▶ Now consider another case where $Y = \frac{1}{2}(X_1 + 1)^3$.
- ▶ The function depends on only one dimension.
- ▶ i.e., the other dimensions are irrelevant for the learning of this function.
- ▶ This function does not peak at 0. The bias is then not as prominent.
- ▶ The variability of the estimate is decided on the distance to NN along X_1 , which increases as the number of irrelevant dimensions increases.

How dimensions affect estimation

- ▶ Now consider another case where $Y = \frac{1}{2}(X_1 + 1)^3$.
- ▶ The function depends on only one dimension.
- ▶ i.e., the other dimensions are irrelevant for the learning of this function.
- ▶ This function does not peak at 0. The bias is then not as prominent.
- ▶ The variability of the estimate is decided on the distance to NN along X_1 , which increases as the number of irrelevant dimensions increases.

How dimensions affect estimation

- ▶ Now consider another case where $Y = \frac{1}{2}(X_1 + 1)^3$.
- ▶ The function depends on only one dimension.
- ▶ i.e., the other dimensions are irrelevant for the learning of this function.
- ▶ This function does not peak at 0. The bias is then not as prominent.
- ▶ The variability of the estimate is decided on the distance to NN along X_1 , which increases as the number of irrelevant dimensions increases.

How dimensions affect estimation

- ▶ Now consider another case where $Y = \frac{1}{2}(X_1 + 1)^3$.
- ▶ The function depends on only one dimension.
- ▶ i.e., the other dimensions are irrelevant for the learning of this function.
- ▶ This function does not peak at 0. The bias is then not as prominent.
- ▶ The variability of the estimate is decided on the distance to NN along X_1 , which increases as the number of irrelevant dimensions increases.

How dimensions affect estimation

- ▶ Now consider another case where $Y = \frac{1}{2}(X_1 + 1)^3$.
- ▶ The function depends on only one dimension.
- ▶ i.e., the other dimensions are irrelevant for the learning of this function.
- ▶ This function does not peak at 0. The bias is then not as prominent.
- ▶ The variability of the estimate is decided on the distance to NN along X_1 , which increases as the number of irrelevant dimensions increases.

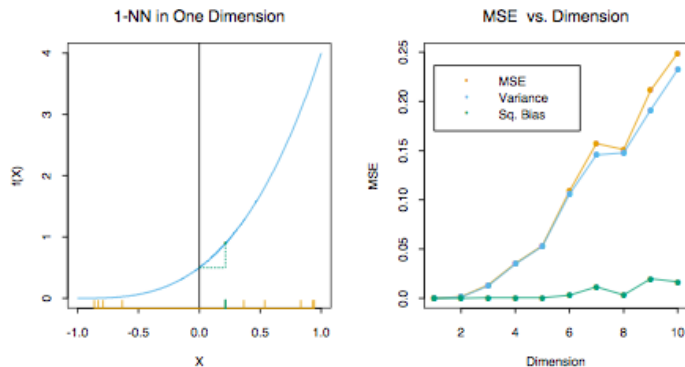


FIGURE 2.8. A simulation example with the same setup as in Figure 2.7. Here the function is constant in all but one dimension: $F(X) = \frac{1}{2}(X_1 + 1)^3$. The variance dominates.

Part III: Statistical models: an overview

The predictive relation

- ▶ The predictive relation between Y and X depends on the definition of “goodness of fit”, sometime in the form of a loss function.
- ▶ L_2 loss: $f(x) = E(Y|X = x)$.
- ▶ L_1 loss: $f(x) = \text{median}(Y|X = x)$.
- ▶ We have shown the nearest-neighbor methods can be viewed as local direct estimates of $f(x)$.
- ▶ We have also shown that the nearest-neighbor methods will run into trouble when the dimension of the input space becomes higher.
- ▶ It is also intuitive that if the relation between Y and X is known to be more structured, KNN methods are not optimal.

The predictive relation

- ▶ The predictive relation between Y and X depends on the definition of “goodness of fit”, sometime in the form of a loss function.
- ▶ L_2 loss: $f(x) = E(Y|X = x)$.
- ▶ L_1 loss: $f(x) = \text{median}(Y|X = x)$.
- ▶ We have shown the nearest-neighbor methods can be viewed as local direct estimates of $f(x)$.
- ▶ We have also shown that the nearest-neighbor methods will run into trouble when the dimension of the input space becomes higher.
- ▶ It is also intuitive that if the relation between Y and X is known to be more structured, KNN methods are not optimal.

The predictive relation

- ▶ The predictive relation between Y and X depends on the definition of “goodness of fit”, sometime in the form of a loss function.
- ▶ L_2 loss: $f(x) = E(Y|X = x)$.
- ▶ L_1 loss: $f(x) = \text{median}(Y|X = x)$.
- ▶ We have shown the nearest-neighbor methods can be viewed as local direct estimates of $f(x)$.
- ▶ We have also shown that the nearest-neighbor methods will run into trouble when the dimension of the input space becomes higher.
- ▶ It is also intuitive that if the relation between Y and X is known to be more structured, KNN methods are not optimal.

The predictive relation

- ▶ The predictive relation between Y and X depends on the definition of “goodness of fit”, sometime in the form of a loss function.
- ▶ L_2 loss: $f(x) = E(Y|X = x)$.
- ▶ L_1 loss: $f(x) = \text{median}(Y|X = x)$.
- ▶ We have shown the nearest-neighbor methods can be viewed as local direct estimates of $f(x)$.
- ▶ We have also shown that the nearest-neighbor methods will run into trouble when the dimension of the input space becomes higher.
- ▶ It is also intuitive that if the relation between Y and X is known to be more structured, KNN methods are not optimal.

The predictive relation

- ▶ The predictive relation between Y and X depends on the definition of “goodness of fit”, sometime in the form of a loss function.
- ▶ L_2 loss: $f(x) = E(Y|X = x)$.
- ▶ L_1 loss: $f(x) = \text{median}(Y|X = x)$.
- ▶ We have shown the nearest-neighbor methods can be viewed as local direct estimates of $f(x)$.
- ▶ We have also shown that the nearest-neighbor methods will run into trouble when the dimension of the input space becomes higher.
- ▶ It is also intuitive that if the relation between Y and X is known to be more structured, KNN methods are not optimal.

The predictive relation

- ▶ The predictive relation between Y and X depends on the definition of “goodness of fit”, sometime in the form of a loss function.
- ▶ L_2 loss: $f(x) = E(Y|X = x)$.
- ▶ L_1 loss: $f(x) = \text{median}(Y|X = x)$.
- ▶ We have shown the nearest-neighbor methods can be viewed as local direct estimates of $f(x)$.
- ▶ We have also shown that the nearest-neighbor methods will run into trouble when the dimension of the input space becomes higher.
- ▶ It is also intuitive that if the relation between Y and X is known to be more structured, KNN methods are not optimal.

The predictive relation

- ▶ For functions f and g that satisfy

$$f(x_i) = g(x_i), i = 1, \dots, n,$$

their fit to the observed data $(x_i, y_i), i = 1, \dots, n$ is the same.

- ▶ The above fact leads to the definition of some kinds of *equivalent models*.
- ▶ By constraining to one model family, the hope is the interception between the model class and those equivalent models is unique.
- ▶ When the interception is not unique, we will have the *identifiability issue*.

The predictive relation

- ▶ For functions f and g that satisfy

$$f(x_i) = g(x_i), i = 1, \dots, n,$$

their fit to the observed data $(x_i, y_i), i = 1, \dots, n$ is the same.

- ▶ The above fact leads to the definition of some kinds of *equivalent models*.
- ▶ By constraining to one model family, the hope is the interception between the model class and those equivalent models is unique.
- ▶ When the interception is not unique, we will have the *identifiability issue*.

The predictive relation

- ▶ For functions f and g that satisfy

$$f(x_i) = g(x_i), i = 1, \dots, n,$$

their fit to the observed data $(x_i, y_i), i = 1, \dots, n$ is the same.

- ▶ The above fact leads to the definition of some kinds of *equivalent models*.
- ▶ By constraining to one model family, the hope is the interception between the model class and those equivalent models is unique.
- ▶ When the interception is not unique, we will have the *identifiability issue*.

The predictive relation

- ▶ For functions f and g that satisfy

$$f(x_i) = g(x_i), i = 1, \dots, n,$$

their fit to the observed data $(x_i, y_i), i = 1, \dots, n$ is the same.

- ▶ The above fact leads to the definition of some kinds of *equivalent models*.
- ▶ By constraining to one model family, the hope is the interception between the model class and those equivalent models is unique.
- ▶ When the interception is not unique, we will have the *identifiability issue*.

The predictive relation

- ▶ The general belief is that the more complicated a model $f(x)$ is, the more likely for $f(x)$ to give prediction further away from $E(y|x)$ at x that is not close to x_i 's observed—the *principle of parsimony* or the *Occam's razor*.
- ▶ Therefore, constraints are usually *complexity* constraints.
- ▶ Usually, the modeling can be carried out using *structured model families*, basis expansions and *kernel/local regression*.
- ▶ The complexity of models are controlled accordingly.

The predictive relation

- ▶ The general belief is that the more complicated a model $f(x)$ is, the more likely for $f(x)$ to give prediction further away from $E(y|x)$ at x that is not close to x_i 's observed—the *principle of parsimony* or the *Occam's razor*.
- ▶ Therefore, constraints are usually *complexity constraints*.
- ▶ Usually, the modeling can be carried out using *structured model families*, basis expansions and *kernel/local regression*.
- ▶ The complexity of models are controlled accordingly.

The predictive relation

- ▶ The general belief is that the more complicated a model $f(x)$ is, the more likely for $f(x)$ to give prediction further away from $E(y|x)$ at x that is not close to x_i 's observed—the *principle of parsimony* or the *Occam's razor*.
- ▶ Therefore, constraints are usually *complexity* constraints.
- ▶ Usually, the modeling can be carried out using *structured model families*, *basis expansions* and *kernel/local regression*.
- ▶ The complexity of models are controlled accordingly.

The predictive relation

- ▶ The general belief is that the more complicated a model $f(x)$ is, the more likely for $f(x)$ to give prediction further away from $E(y|x)$ at x that is not close to x_i 's observed— the *principle of parsimony* or the *Occam's razor* .
- ▶ Therefore, constraints are usually *complexity* constraints.
- ▶ Usually, the modeling can be carried out using *structured model families*, basis expansions and *kernel/local regression*.
- ▶ The complexity of models are controlled accordingly.

Part IV: Model selection and selection

Model Assessment and Selection

- ▶ **Model Selection:** use training data to choose the (approximately) best model among a group of candidates.
- ▶ **Model Assessment:** use available data to evaluate the performance of the final model.
- ▶ Data division for model selection and assessment
 - ▶ In data-rich situation: training, validation, test.
 - ▶ No general rule on how to choose the number of observations for each partition.
 - ▶ Example given by the “statistical learning” book: 50%, 25%, 25%.

Model Assessment and Selection

- ▶ **Model Selection:** use training data to choose the (approximately) best model among a group of candidates.
- ▶ **Model Assessment:** use available data to evaluate the performance of the final model.
- ▶ Data division for model selection and assessment
 - ▶ In data-rich situation: training, validation, test.
 - ▶ No general rule on how to choose the number of observations for each partition.
 - ▶ Example given by the “statistical learning” book: 50%, 25%, 25%.

Model Assessment and Selection

- ▶ **Model Selection:** use training data to choose the (approximately) best model among a group of candidates.
- ▶ **Model Assessment:** use available data to evaluate the performance of the final model.
- ▶ **Data division for model selection and assessment**
 - ▶ In data-rich situation: training, validation, test.
 - ▶ No general rule on how to choose the number of observations for each partition.
 - ▶ Example given by the “statistical learning” book: 50%, 25%, 25%.

Model Assessment and Selection

- ▶ **Model Selection:** use training data to choose the (approximately) best model among a group of candidates.
- ▶ **Model Assessment:** use available data to evaluate the performance of the final model.
- ▶ Data division for model selection and assessment
 - ▶ In data-rich situation: training, validation, test.
 - ▶ No general rule on how to choose the number of observations for each partition.
 - ▶ Example given by the “statistical learning” book: 50%, 25%, 25%.

Model Assessment and Selection

- ▶ **Model Selection:** use training data to choose the (approximately) best model among a group of candidates.
- ▶ **Model Assessment:** use available data to evaluate the performance of the final model.
- ▶ Data division for model selection and assessment
 - ▶ In data-rich situation: training, validation, test.
 - ▶ No general rule on how to choose the number of observations for each partition.
 - ▶ Example given by the “statistical learning” book: 50%, 25%, 25%.

Model Assessment and Selection

- ▶ **Model Selection:** use training data to choose the (approximately) best model among a group of candidates.
- ▶ **Model Assessment:** use available data to evaluate the performance of the final model.
- ▶ Data division for model selection and assessment
 - ▶ In data-rich situation: training, validation, test.
 - ▶ No general rule on how to choose the number of observations for each partition.
 - ▶ Example given by the “statistical learning” book: 50%, 25%, 25%.

Test Error versus Training Error

- ▶ **Loss function** $L(Y, \hat{f}(x))$ defines prediction errors (e.g., squared error—quadratic loss; absolute error— L_1 loss);
- ▶ **Test error** is the expected prediction error over *an independent test sample*:

$$\text{Err} = E[L(Y, \hat{f}(X))].$$

- ▶ Note: this expectation takes into account the sampling variability in the training data that are used to fit \hat{f} (including effects from the sample size).

Test Error versus Training Error

- ▶ **Loss function** $L(Y, \hat{f}(x))$ defines prediction errors (e.g., squared error–quadratic loss; absolute error– L_1 loss);
- ▶ **Test error** is the expected prediction error over *an independent test sample*:

$$\text{Err} = E[L(Y, \hat{f}(X))].$$

- ▶ Note: this expectation takes into account the sampling variability in the training data that are used to fit \hat{f} (including effects from the sample size).

Test Error versus Training Error

- ▶ **Loss function** $L(Y, \hat{f}(x))$ defines prediction errors (e.g., squared error–quadratic loss; absolute error– L_1 loss;);
- ▶ **Test error** is the expected prediction error over *an independent test sample*:

$$\text{Err} = E[L(Y, \hat{f}(X))].$$

- ▶ **Note:** this expectation takes into account the sampling variability in the training data that are used to fit \hat{f} (including effects from the sample size).

Test Error versus Training Error

- ▶ The *test error* is also called *generalization error*.
- ▶ *Training error* is the average loss of the training set:

$$\overline{\text{err}} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}(x_i)).$$

- ▶ Training error is not a good estimate of the test error.

Test Error versus Training Error

- ▶ The *test error* is also called *generalization error*.
- ▶ *Training error* is the average loss of the training set:

$$\overline{\text{err}} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}(x_i)).$$

- ▶ Training error is not a good estimate of the test error.

Test Error versus Training Error

- ▶ The *test error* is also called *generalization error*.
- ▶ *Training error* is the average loss of the training set:

$$\overline{\text{err}} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}(x_i)).$$

- ▶ Training error is not a good estimate of the test error.

Test error and training error versus model complexity

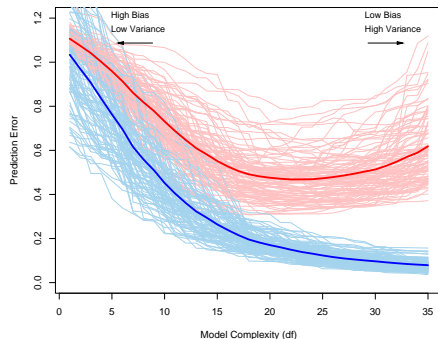
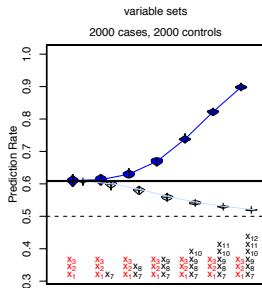
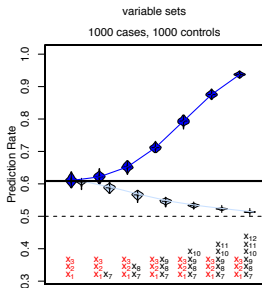
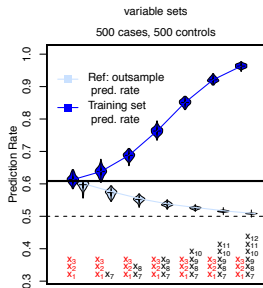


FIGURE 7.1. Behavior of test sample and training sample error as the model complexity is varied. The light blue curves show the training error $\overline{\text{err}}$, while the light red curves show the conditional test error Err_T for 100 training sets of size 50 each, as the model complexity is increased. The solid curves show the expected test error Err and the expected training error $\mathbb{E}[\overline{\text{err}}]$.

Test error and training error versus model complexity



variable sets

variable sets

variable sets

Cross-validation

- ▶ Cross-validation methods directly estimate the extra-sample prediction error (the generalization error) by using non-overlapping training data and test data.

- ▶ K-fold cross-validation Layout

1	2	...	K-1	K
Test	Train	...	Train	Train

⋮

1	2	...	K-1	K
Train	Train	...	Train	Test

K estimates of prediction errors are combined to estimate the prediction error of the training set.

Cross-validation

- ▶ Cross-validation methods directly estimate the extra-sample prediction error (the generalization error) by using non-overlapping training data and test data.

- ▶ K-fold cross-validation Layout

1	2	...	K-1	K
Test	Train	...	Train	Train

⋮

1	2	...	K-1	K
Train	Train	...	Train	Test

K estimates of prediction errors are combined to estimate the prediction error of the training set.

Cross-validation (cont'd)

- ▶ More specifically, prediction error is estimated by

$$CV = \frac{1}{N} \sum_{k=1}^K \sum_{i \in \text{block } k} L(y_i, \hat{f}_{-k}(x_i))$$

- ▶ $K = N$ is just leave-one-out validation: the training sets will be too similar.
- ▶ Small K , the training set is too small comparing to the total training set. Loss of efficiency.
- ▶ We are trying to estimate the prediction performance at the total training set size using K -fold cross-validation.
- ▶ K should be chosen (theoretically) so that the learning curve does not have a steep slope at $\frac{k-1}{K} N$.

Cross-validation (cont'd)

- ▶ More specifically, prediction error is estimated by

$$CV = \frac{1}{N} \sum_{k=1}^K \sum_{i \in \text{block } k} L(y_i, \hat{f}_{-k}(x_i))$$

- ▶ $K = N$ is just leave-one-out validation: the training sets will be too similar.
- ▶ Small K , the training set is too small comparing to the total training set. Loss of efficiency.
- ▶ We are trying to estimate the prediction performance at the total training set size using K -fold cross-validation.
- ▶ K should be chosen (theoretically) so that the learning curve does not have a steep slope at $\frac{k-1}{K} N$.

Cross-validation (cont'd)

- ▶ More specifically, prediction error is estimated by

$$CV = \frac{1}{N} \sum_{k=1}^K \sum_{i \in \text{block } k} L(y_i, \hat{f}_{-k}(x_i))$$

- ▶ $K = N$ is just leave-one-out validation: the training sets will be too similar.
- ▶ Small K , the training set is too small comparing to the total training set. Loss of efficiency.
- ▶ We are trying to estimate the prediction performance at the total training set size using K -fold cross-validation.
- ▶ K should be chosen (theoretically) so that the learning curve does not have a steep slope at $\frac{k-1}{K} N$.

Cross-validation (cont'd)

- ▶ More specifically, prediction error is estimated by

$$CV = \frac{1}{N} \sum_{k=1}^K \sum_{i \in \text{block } k} L(y_i, \hat{f}_{-k}(x_i))$$

- ▶ $K = N$ is just leave-one-out validation: the training sets will be too similar.
- ▶ Small K , the training set is too small comparing to the total training set. Loss of efficiency.
- ▶ We are trying to estimate the prediction performance at the total training set size using K -fold cross-validation.
- ▶ K should be chosen (theoretically) so that the learning curve does not have a steep slope at $\frac{k-1}{K} N$.

Cross-validation (cont'd)

- ▶ More specifically, prediction error is estimated by

$$CV = \frac{1}{N} \sum_{k=1}^K \sum_{i \in \text{block } k} L(y_i, \hat{f}_{-k}(x_i))$$

- ▶ $K = N$ is just leave-one-out validation: the training sets will be too similar.
- ▶ Small K , the training set is too small comparing to the total training set. Loss of efficiency.
- ▶ We are trying to estimate the prediction performance at the total training set size using K -fold cross-validation.
- ▶ K should be chosen (theoretically) so that the learning curve does not have a steep slope at $\frac{k-1}{K} N$.

Cross-validation (cont'd)

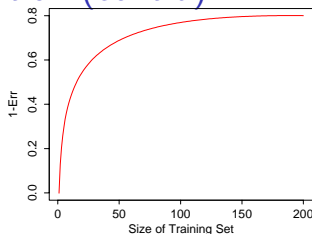


Figure 7.8: *Hypothetical learning curve for a classifier on a given task; a plot of $1 - \text{Err}$ versus the size of the training set N . With a dataset of 200 observations, fivefold cross-validation would use training sets of size 160, which would behave much like the full set. However, with a dataset of 50 observations fivefold cross-validation would use training sets of size 40, and this would result in a considerable overestimate of prediction error.*

Cross-validation (cont'd)

Model selection based on cross-validation

- ▶ The model with the best cross-validation prediction error—the “best” model.
- ▶ *One-standard-error rule*: since the K-fold cross-validation also allows one to estimate the standard error of prediction errors. We choose the most parsimonious (simplest) model whose error is no more than one standard error above the error of the “best” model

Cross-validation (cont'd)

Model selection based on cross-validation

- ▶ The model with the best cross-validation prediction error—the “best” model.
- ▶ *One-standard-error rule*: since the K-fold cross-validation also allows one to estimate the standard error of prediction errors. We choose **the most parsimonious (simplest) model whose error is no more than one standard error above the error of the “best” model**

Cross-validation (cont'd)

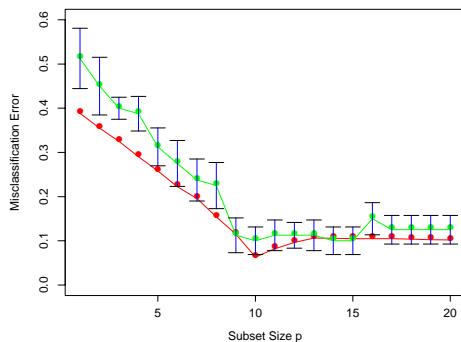


Figure 7.9: *Prediction error (red) and tenfold cross-validation curve (green) estimated from a single training set, from the scenario in the bottom right panel of Figure 7.3.*

Bootstrap methods

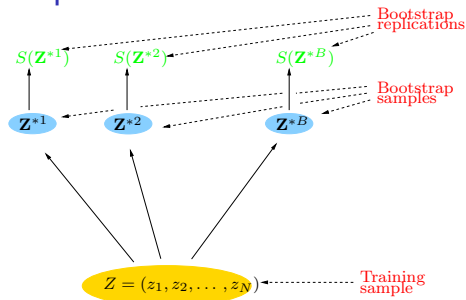


Figure 7.10: Schematic of the bootstrap process. We wish to assess the statistical accuracy of a quantity $S(\mathbf{Z})$ computed from our dataset. B training sets \mathbf{Z}^{*b} , $b = 1, \dots, B$ each of size N are drawn with replacement from the original dataset. The quantity of interest $S(\mathbf{Z})$ is computed from each bootstrap training set, and the values $S(\mathbf{Z}^{*1}), \dots, S(\mathbf{Z}^{*B})$ are used to assess the statistical accuracy of $S(\mathbf{Z})$.

Basic Bootstrap

- ▶ Nonparametric bootstrap: resample (w/o replacements) training data (x_i, y_i) .
- ▶ Parametric bootstrap: $y_i^* = \hat{f}(x_i) + \varepsilon_i^*$, where ε_i^* is generated from a parametric distribution.

Basic Bootstrap

- ▶ Nonparametric bootstrap: resample (w/o replacements) training data (x_i, y_i) .
- ▶ Parametric bootstrap: $y_i^* = \hat{f}(x_i) + \varepsilon_i^*$, where ε_i^* is generated from a parametric distribution.

Bootstrap for model assessment

- ▶ Prediction error estimated using Bootstrap

$$\widehat{\text{Err}}_{\text{boot}} = \frac{1}{B} \widehat{\text{Err}}^{*b}$$

where

$$\widehat{\text{Err}}^{*b} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}^{*b}(x_i))$$

- ▶ Here, the bootstrap sample are acting as the training samples, while the *original* training sample set is acting as the test sample.
- ▶ The problem is that a bootstrap sample and the original sample have a lot of observations in common.

Bootstrap for model assessment

- ▶ Prediction error estimated using Bootstrap

$$\widehat{\text{Err}}_{\text{boot}} = \frac{1}{B} \widehat{\text{Err}}^{*b}$$

where

$$\widehat{\text{Err}}^{*b} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}^{*b}(x_i))$$

- ▶ Here, the bootstrap sample are acting as the training samples, while the *original* training sample set is acting as the test sample.
- ▶ The problem is that a bootstrap sample and the original sample have a lot of observations in common.

Bootstrap for model assessment

- ▶ Prediction error estimated using Bootstrap

$$\widehat{\text{Err}}_{\text{boot}} = \frac{1}{B} \widehat{\text{Err}}^{*b}$$

where

$$\widehat{\text{Err}}^{*b} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}^{*b}(x_i))$$

- ▶ Here, the bootstrap sample are acting as the training samples, while the *original* training sample set is acting as the test sample.
- ▶ The problem is that a bootstrap sample and the original sample have a lot of observations in common.

Bootstrap for model validation

- ▶ For validation, for each observation (x_i, y_i) , we keep track of the bootstrap samples \mathbf{Z}^b that do not contain z_i . We call the set of such b 's C^{-i} .
- ▶ The probability for a bootstrap sample not containing z_i is $1 - (1 - \frac{1}{N})^N \approx 1 - e^{-1} = 0.632$.
- ▶ The leave-one-out bootstrap estimate of prediction error is then

$$\widehat{\text{Err}}^{(1)} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|C^{-i}|} \sum_{b \in C^{-i}} L(y_i, \hat{f}^{*b}(x_i))$$

Bootstrap for model validation

- ▶ For validation, for each observation (x_i, y_i) , we keep track of the bootstrap samples \mathbf{Z}^b that do not contain z_i . We call the set of such b 's C^{-i} .
- ▶ The probability for a bootstrap sample not containing z_i is $1 - (1 - \frac{1}{N})^N \approx 1 - e^{-1} = 0.632$.
- ▶ The leave-one-out bootstrap estimate of prediction error is then

$$\widehat{\text{Err}}^{(1)} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|C^{-i}|} \sum_{b \in C^{-i}} L(y_i, \hat{f}^{*b}(x_i))$$

Bootstrap for model validation

- ▶ For validation, for each observation (x_i, y_i) , we keep track of the bootstrap samples \mathbf{Z}^b that do not contain z_i . We call the set of such b 's C^{-i} .
- ▶ The probability for a bootstrap sample not containing z_i is $1 - (1 - \frac{1}{N})^N \approx 1 - e^{-1} = 0.632$.
- ▶ The leave-one-out bootstrap estimate of prediction error is then

$$\widehat{\text{Err}}^{(1)} = \frac{1}{N} \sum_{i=1}^N \frac{1}{|C^{-i}|} \sum_{b \in C^{-i}} L(y_i, \hat{f}^{*b}(x_i))$$

Part V: Model ensemble

Model averaging: bagging

- ▶ Based on bootstrap samples

$$\hat{f}_{\text{bag}} = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x)$$

- ▶ This is different from model built with bootstrap estimate of the model parameters, when the model is nonlinear.
- ▶ For classification problems, the \hat{f} usually takes the form estimated probabilities for classes. And the prediction is usually given as the most probable class.
- ▶ The bagged predictor for classification can take two forms:
 - ▶ bag the predictions (weighted votes)
 - ▶ bag the class probabilities

Model averaging: bagging

- ▶ Based on bootstrap samples

$$\hat{f}_{\text{bag}} = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x)$$

- ▶ This is different from model built with bootstrap estimate of the model parameters, when the model is nonlinear.
- ▶ For classification problems, the \hat{f} usually takes the form estimated probabilities for classes. And the prediction is usually given as the most probable class.
- ▶ The bagged predictor for classification can take two forms:
 - ▶ bag the predictions (weighted votes)
 - ▶ bag the class probabilities

Model averaging: bagging

- ▶ Based on bootstrap samples

$$\hat{f}_{\text{bag}} = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x)$$

- ▶ This is different from model built with bootstrap estimate of the model parameters, when the model is nonlinear.
- ▶ For classification problems, the \hat{f} usually takes the form estimated probabilities for classes. And the prediction is usually given as the most probable class.
- ▶ The bagged predictor for classification can take two forms:
 - ▶ bag the predictions (weighted votes)
 - ▶ bag the class probabilities

Model averaging: bagging

- ▶ Based on bootstrap samples

$$\hat{f}_{\text{bag}} = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x)$$

- ▶ This is different from model built with bootstrap estimate of the model parameters, when the model is nonlinear.
- ▶ For classification problems, the \hat{f} usually takes the form estimated probabilities for classes. And the prediction is usually given as the most probable class.
- ▶ The bagged predictor for classification can take two forms:
 - ▶ bag the predictions (weighted votes)
 - ▶ bag the class probabilities

Model averaging: bagging

- ▶ Based on bootstrap samples

$$\hat{f}_{\text{bag}} = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x)$$

- ▶ This is different from model built with bootstrap estimate of the model parameters, when the model is nonlinear.
- ▶ For classification problems, the \hat{f} usually takes the form estimated probabilities for classes. And the prediction is usually given as the most probable class.
- ▶ The bagged predictor for classification can take two forms:
 - ▶ bag the predictions (weighted votes)
 - ▶ bag the class probabilities

Model averaging: bagging

- ▶ Based on bootstrap samples

$$\hat{f}_{\text{bag}} = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x)$$

- ▶ This is different from model built with bootstrap estimate of the model parameters, when the model is nonlinear.
- ▶ For classification problems, the \hat{f} usually takes the form estimated probabilities for classes. And the prediction is usually given as the most probable class.
- ▶ The bagged predictor for classification can take two forms:
 - ▶ bag the predictions (weighted votes)
 - ▶ bag the class probabilities

Model averaging: weighted averages of candidate models

- ▶ *Committee methods*: the final model is an unweighted “average” of the fitted individual \hat{f}_m 's.

$$\hat{f}(x) = \frac{1}{M} \sum_{m=1}^M \hat{f}_m(x; \hat{\theta}_m)$$

- ▶ Models M_1, \dots, M_M can also be averaged based on $\Pr(M_m|\mathbf{Z})$.

$$\hat{f}(x) = \sum_{m=1}^M \hat{f}_m(x; \hat{\theta}_m) \Pr(M_m|\mathbf{Z})$$

- ▶ $\Pr(M_m|\mathbf{Z})$ can be well estimated by BIC. Full Bayesian computation can also derive these probabilities but is more computationally intensive, and the computational cost is not justified by performance.

Model averaging: weighted averages of candidate models

- ▶ *Committee methods*: the final model is an unweighted “average” of the fitted individual \hat{f}_m 's.

$$\hat{f}(x) = \frac{1}{M} \sum_{m=1}^M \hat{f}_m(x; \hat{\theta}_m)$$

- ▶ Models M_1, \dots, M_M can also be averaged based on $\Pr(M_m|\mathbf{Z})$.

$$\hat{f}(x) = \sum_{m=1}^M \hat{f}_m(x; \hat{\theta}_m) \Pr(M_m|\mathbf{Z})$$

- ▶ $\Pr(M_m|\mathbf{Z})$ can be well estimated by BIC. Full Bayesian computation can also derive these probabilities but is more computationally intensive, and the computational cost is not justified by performance.

Model averaging: weighted averages of candidate models

- ▶ *Committee methods*: the final model is an unweighted “average” of the fitted individual \hat{f}_m 's.

$$\hat{f}(x) = \frac{1}{M} \sum_{m=1}^M \hat{f}_m(x; \hat{\theta}_m)$$

- ▶ Models M_1, \dots, M_M can also be averaged based on $\Pr(M_m|\mathbf{Z})$.

$$\hat{f}(x) = \sum_{m=1}^M \hat{f}_m(x; \hat{\theta}_m) \Pr(M_m|\mathbf{Z})$$

- ▶ $\Pr(M_m|\mathbf{Z})$ can be well estimated by BIC. Full Bayesian computation can also derive these probabilities but is more computationally intensive, and the computational cost is not justified by performance.

Model averaging: weighted averages (cont'd)

- ▶ Models averages from previous slide can be viewed as *weighted average*.

$$\hat{f}(x) = \sum_{m=1}^M \omega_m \hat{f}_m(x)$$

- ▶ Searching for the best weight, one may consider

$$\hat{\omega} = \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^N L(y_i, \sum_{m=1}^M \omega_m \hat{f}_m(x_i))$$

- ▶ Will this lead to overfitting?

Model averaging: weighted averages (cont'd)

- ▶ Models averages from previous slide can be viewed as *weighted average*.

$$\hat{f}(x) = \sum_{m=1}^M \omega_m \hat{f}_m(x)$$

- ▶ Searching for the best weight, one may consider

$$\hat{\omega} = \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^N L(y_i, \sum_{m=1}^M \omega_m \hat{f}_m(x_i))$$

- ▶ Will this lead to overfitting?

Model averaging: weighted averages (cont'd)

- ▶ Models averages from previous slide can be viewed as *weighted average*.

$$\hat{f}(x) = \sum_{m=1}^M \omega_m \hat{f}_m(x)$$

- ▶ Searching for the best weight, one may consider

$$\hat{\omega} = \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^N L(y_i, \sum_{m=1}^M \omega_m \hat{f}_m(x_i))$$

- ▶ Will this lead to overfitting?

Model averaging: stacking

- Stacking or *stacked generalization*

$$\hat{\omega}^{st} = \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^N L \left(y_i, \sum_{m=1}^M \omega_m \hat{f}_m^{-i}(x_i) \right)$$

- If we search ω among vectors with elements 0 or 1, and with sum 1, stacking is equivalent to leave-one-out model selection.

Model averaging: stacking

- Stacking or *stacked generalization*

$$\hat{\omega}^{st} = \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^N L \left(y_i, \sum_{m=1}^M \omega_m \hat{f}_m^{-i}(x_i) \right)$$

- If we search ω among vectors with elements 0 or 1, and with sum 1, stacking is equivalent to leave-one-out model selection.