# Chapter 7: Model Assessment and Selection

DD3364

April 20, 2012

# Introduction

- Have target variable $Y$ to estimate from a vector of inputs $X$.

- A prediction model $\hat{f}(X)$ has been estimated from training data $\mathcal{T} = \{(x_1, y_1), \ldots, (x_n, y_n)\}$.

- The loss function $L(Y, \hat{f}(X))$ measures the errors between $Y$ and $\hat{f}(X)$.

- Common loss functions are

$$L(Y, \hat{f}(X)) = \begin{cases} (Y - \hat{f}(X))^2 & \textbf{squared error}, \\ |Y - \hat{f}(X)| & \textbf{absolute error} \end{cases}$$

- Have target variable $Y$ to estimate from a vector of inputs $X$.

- A prediction model $\hat{f}(X)$ has been estimated from training data $\mathcal{T} = \{(x_1, y_1), \ldots, (x_n, y_n)\}$.

- The loss function $L(Y, \hat{f}(X))$ measures the errors between $Y$ and $\hat{f}(X)$.

- Common loss functions are

$$L(Y, \hat{f}(X)) = \begin{cases} (Y - \hat{f}(X))^2 & \text{squared error,} \\ |Y - \hat{f}(X)| & \text{absolute error} \end{cases}$$

- Have target variable $Y$ to estimate from a vector of inputs $X$.

- A prediction model $\hat{f}(X)$ has been estimated from training data $\mathcal{T} = \{(x_1, y_1), \ldots, (x_n, y_n)\}$.

- The loss function $L(Y, \hat{f}(X))$ measures the errors between $Y$ and $\hat{f}(X)$.

- Common loss functions are

$$L(Y, \hat{f}(X)) = \begin{cases} (Y - \hat{f}(X))^2 & \textbf{squared error}, \\ |Y - \hat{f}(X)| & \textbf{absolute error} \end{cases}$$

**Test Error** a.k.a. **generalization error**

$$\mathsf{Err}_{\mathcal{T}} = E[\, L(Y, \hat{f}(X)) \,|\, \mathcal{T} \,]$$

- Prediction error over an independent test sample.

- $X$ and $Y$ are drawn randomly from $p(X, Y)$.

- The training set $\mathcal{T}$ if fixed.

**Expected Prediction Error** (expected test error)

$$\mathsf{Err} = E[\,L(Y, \hat{f}(X))\,] = E[\,\mathsf{Err}_{\mathcal{T}}\,]$$

- In this case take expectation over all the random quantities including the training set.

Which quantities interest us

- Would like to estimate $\mathsf{Err}_{\mathcal{T}}$.

- But in most cases it is easier to estimate Err. Why??

**Expected Prediction Error** (expected test error)

$$\mathsf{Err} = E[\, L(Y, \hat{f}(X))\,] = E[\,\mathsf{Err}_{\mathcal{T}}\,]$$

- In this case take expectation over all the random quantities including the training set.

**Which quantities interest us**

- Would like to estimate $\mathsf{Err}_{\mathcal{T}}$.

- But in most cases it is easier to estimate Err. Why??

**Training error**

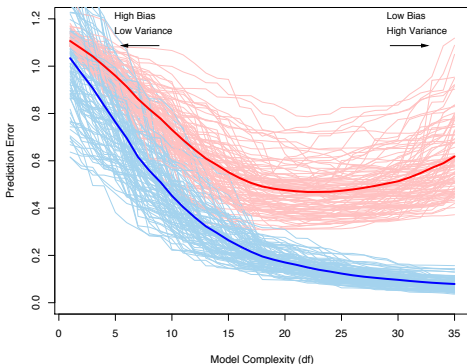$$\overline{\text{err}} = \frac{1}{n} \sum_{i=1}^{n} L(y_i, \hat{f}(x_i))$$

Already know as complexity of $\hat{f}$ increases

- then $\overline{\text{err}} \to 0$,

- but there is a tendency to overfit and $\text{Err}_{\mathcal{T}}$ increases

- $\overline{\text{err}}$ is not a good estimate of $\text{Err}_{\mathcal{T}}$ or Err.

We will be revisiting the **Bias**-**Variance** trade-off.

Test and training error as model complexity increases.



1. Light blue curve: training error $\overline{\text{err}}$.
2. Solid blue curve: expected training error $E[\overline{\text{err}}]$.
3. Light red curve: conditional test error $\text{Err}_{\mathcal{T}}$.
4. Solid red curve: expected test error Err.

- Have target categorical variable $G \in \{1, \dots, K\}$ to estimate from a vector of inputs $X$.

- Typically model $p_k(X) = P(G = k|X)$ and define

$$\hat{G}(X) = \arg \max_k p_k(X)$$

- Common loss functions are

  ① **0-1 loss**

  $$L(G, \hat{G}(X)) = \mathsf{Ind}(G \neq \hat{G}(X))$$

  ② **log-likelihood** a.ka. deviance

  $$L(G, \hat{p}(X)) = -2 \sum_{k=1}^{K} \mathsf{Ind}(G = k) \, \hat{p}_k(X) = -2 \log \hat{p}_G(X)$$

- Have target categorical variable $G \in \{1, \ldots, K\}$ to estimate from a vector of inputs $X$.

- Typically model $p_k(X) = P(G = k|X)$ and define

$$\hat{G}(X) = \arg\max_k p_k(X)$$

- Common loss functions are

  1. 0-1 loss

$$L(G, \hat{G}(X)) = \text{Ind}(G \neq \hat{G}(X))$$

  2. log-likelihood a.ka. deviance

$$L(G, \hat{p}(X)) = -2 \sum_{k=1}^{K} \text{Ind}(G = k)\, \hat{p}_k(X) = -2 \log \hat{p}_G(X)$$

- Have target categorical variable $G \in \{1, \ldots, K\}$ to estimate from a vector of inputs $X$.

- Typically model $p_k(X) = P(G = k | X)$ and define

$$\hat{G}(X) = \arg \max_k p_k(X)$$

- Common loss functions are

  **❶ 0-1 loss**

  $$L(G, \hat{G}(X)) = \mathsf{Ind}(G \neq \hat{G}(X))$$

  **❷ log-likelihood** a.ka. deviance

  $$L(G, \hat{p}(X)) = -2 \sum_{k=1}^{K} \mathsf{Ind}(G = k) \, \hat{p}_k(X) = -2 \log \hat{p}_G(X)$$

- **Test Error**

$$\mathrm{Err}_{\mathcal{T}} = E[\, L(G, \hat{G}(X)) \,|\, \mathcal{T} \,]$$

- **Training Error** one common definition

$$\overline{\mathrm{err}} = -\frac{2}{n} \sum_{i=1}^{n} \log \hat{p}_{g_i}(x_i)$$

- $\hat{f}_\alpha(x)$ typically has a tunable parameter $\alpha$ controlling its complexity.

- Want to find the value of $\alpha$ s.t.

$$\hat{\alpha} = \arg\min_\alpha E[\, L(Y, \hat{f}_\alpha(X)) \,]$$

- Estimate $E[\, L(Y, \hat{f}_\alpha(X)) \,]$ for different values of $\alpha$.

- This chapter presents methods how to do this.

- Choose the $\alpha$ with minimum estimate.

**Model selection**

Estimate the performance of different models in order to choose the best one.

**Model Assessment**

Having chosen a final model, estimate its prediction error on new data.

Randomly divide the dataset into 3 parts

| Train | Validation | Test |
|:---:|:---:|:---:|

Common split ratio 50%, 25%, 25%.

## Model Selection

- Use **training set** to fit each model.

- Use **validation set** to estimate $\text{Err}_{\mathcal{T}}$ for each model.

- Choose model with lowest $\text{Err}_{\mathcal{T}}$ estimate.

## Model Assessment of the chosen model

- Use the **test set** - unseen until this stage - to estimate $\text{Err}_{\mathcal{T}}$.

Approximate the **validation step** either

- analytically with approaches such as
    1. Akaike Information Criterion
    2. Baysian Information Criterion
    3. Minimum Description Length
    4. Structural Risk Minimization

  or

- with efficient sample re-use
    1. cross-validation
    2. the bootstrap

Each method also provides estimates of Err or $Err_{\mathcal{T}}$ of the final chosen model.

Approximate the **validation step** either

- analytically with approaches such as

  1. Akaike Information Criterion
  2. Baysian Information Criterion
  3. Minimum Description Length
  4. Structural Risk Minimization

  or

- with efficient sample re-use

  1. cross-validation
  2. the bootstrap

Each method also provides estimates of Err or $\text{Err}_{\mathcal{T}}$ of the final chosen model.

# The Bias-Variance Decomposition

- Will assume an additive model

$$Y = f(X) + \epsilon$$

where $\mathrm{E}[\epsilon] = 0$ and $\mathrm{Var}[\epsilon] = \sigma_\epsilon^2$.

- Then the **expected prediction error** of $\hat{f}(X)$ at $X = x_0$

$$\mathrm{Err}(x_0) = \mathrm{E}[(Y - \hat{f}(x_0))^2 | X = x_0]$$

can be expressed as

$$\mathrm{Err}(x_0) = \text{Irreducible Error} + \text{Bias}^2 + \text{Variance}$$

**Irreducible error**:    $\sigma_\epsilon^2$,

**Bias**:    $\mathrm{E}[\hat{f}(x_0) - f(x_0)]$,

**Variance**:    $\mathrm{Var}[\hat{f}(x_0)]$

- Will assume an additive model

$$Y = f(X) + \epsilon$$

where $\mathrm{E}[\epsilon] = 0$ and $\mathrm{Var}[\epsilon] = \sigma_\epsilon^2$.

- Then the **expected prediction error** of $\hat{f}(X)$ at $X = x_0$

$$\mathrm{Err}(x_0) = \mathrm{E}[(Y - \hat{f}(x_0))^2 | X = x_0]$$

can be expressed as

$$\mathrm{Err}(x_0) = \text{Irreducible Error} + \text{Bias}^2 + \text{Variance}$$

| | |
|---|---|
| **Irreducible error**: | $\sigma_\epsilon^2$, |
| **Bias**: | $\mathrm{E}[\hat{f}(x_0) - f(x_0)]$, |
| **Variance**: | $\mathrm{Var}[\hat{f}(x_0)]$ |

$$\mathsf{Err}(x_0) = \sigma_\epsilon^2 + \left[ f(x_0) - \frac{1}{k} \sum_{l=1}^{k} f(x_{(l)}) \right]^2 + \frac{\sigma_\epsilon^2}{k}$$

squared bias        variance

- Complexity of model is inversely related $k$.

- As $k$ increases the variance decreases.

- As $k$ increases the squared bias increases.

The above expression was computed by assuming the $x_i$'s are fixed.

Have a linear model

$$\hat{f}_p(x) = x^t \hat{\beta}$$

where $\hat{\beta}$ is $p$-dimensional and fit by least squares, then

$$\text{Err}(x_0) = \sigma_\epsilon^2 + \left[ f(x_0) - E[\hat{f}_p(x_0)] \right]^2 + \|h(x_0)\|^2 \sigma_\epsilon^2$$

with $h(x_0) = \mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}x_0$ and $\hat{f}_p(x_0) = x_0^t(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t y$.

Have a linear model

$$\hat{f}_{p,\alpha}(x) = x^t \hat{\beta}_\alpha$$

where $\hat{\beta}_\alpha$ is $p$-dimensional and fit via ridge regression, then

$$\text{Err}(x_0) = \sigma_\epsilon^2 + \left[ f(x_0) - E[\hat{f}_{p,\alpha}(x_0)] \right]^2 + \|h_\alpha(x_0)\|^2 \, \sigma_\epsilon^2$$

with

$$h_\alpha(x_0) = \mathbf{X}(\mathbf{X}^t\mathbf{X} + \alpha I)^{-1} x_0$$

$$\hat{f}_{p,\alpha}(x_0) = x_0^t (\mathbf{X}^t\mathbf{X} + \alpha I)^{-1} \mathbf{X}^t \, y$$

Therefore this regression fit model has a different bias and variance to the least square fit.

Have a linear model

$$\hat{f}_{p,\alpha}(x) = x^t \hat{\beta}_\alpha$$

where $\hat{\beta}_\alpha$ is $p$-dimensional and fit via ridge regression, then

$$\text{Err}(x_0) = \sigma_\epsilon^2 + \left[ f(x_0) - E[\hat{f}_{p,\alpha}(x_0)] \right]^2 + \|h_\alpha(x_0)\|^2 \sigma_\epsilon^2$$

with

$$h_\alpha(x_0) = \mathbf{X}(\mathbf{X}^t\mathbf{X} + \alpha I)^{-1} x_0$$

$$\hat{f}_{p,\alpha}(x_0) = x_0^t(\mathbf{X}^t\mathbf{X} + \alpha I)^{-1}\mathbf{X}^t y$$

Therefore this regression fit model has a different bias and variance to the least square fit.

Let $\beta_*$ denote the parameters of the best-fitting linear approx to $f$:

$$\beta_* = \arg\min_{\beta} E[\,(f(X) - X^t\beta)^2\,]$$

Can write the **averaged squared bias**

$$E_{x_0}\left[(f(x_0) - E[\hat{f}_\alpha(x_0)])^2\right]$$

as

$$E_{x_0}[\,(f(x_0) - x_0^t\beta_*)^2\,] + E_{x_0}[\,(x_0^t\beta_* - E[x_0^t\hat{\beta}_\alpha])^2\,]$$

Ave[Model Bias]$^2$       Ave[Estimation Bias]$^2$

- Estimation bias is zero for ordinary least sq. estimate.
- Estimation bias is positive for ridge regression estimate.

Let $\beta_*$ denote the parameters of the best-fitting linear approx to $f$:

$$\beta_* = \arg \min_{\beta} E[\,(f(X) - X^t\beta)^2\,]$$

Can write the **averaged squared bias**
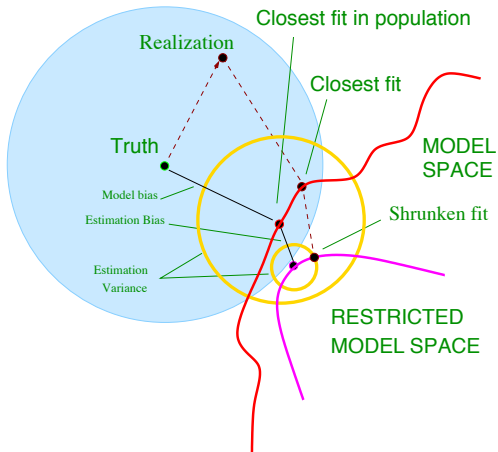
$$E_{x_0}\left[(f(x_0) - E[\,\hat{f}_\alpha(x_0)\,])^2\right]$$

as

$$E_{x_0}[\,(f(x_0) - x_0^t\beta_*)^2\,] + E_{x_0}[\,(x_0^t\beta_* - E[x_0^t\hat{\beta}_\alpha])^2\,]$$

$\text{Ave[Model Bias]}^2 \qquad\qquad \text{Ave[Estimation Bias]}^2$

- Estimation bias is zero for ordinary least sq. estimate.
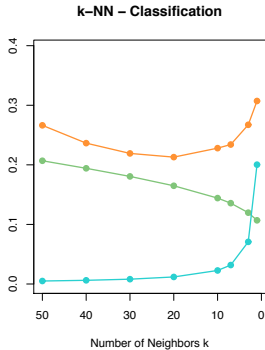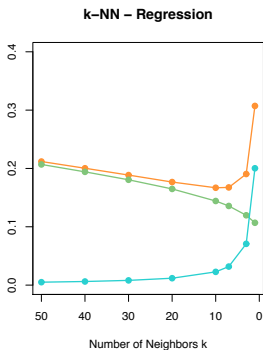- Estimation bias is positive for ridge regression estimate.

**The Set-up**

- Have $n = 80$ observations and $p = 20$ predictors.

- $X$ is uniformly distributed in $[0, 1]^{20}$ and

$$Y = \begin{cases} 0 & \text{if } X_1 \leq .5 \\ 1 & \text{if } X_1 > .5 \end{cases}$$

- Apply $k$-nn to perform both the classification and regression tasks.

- Use **squared error loss** to measure **Err** for the **regression task**.

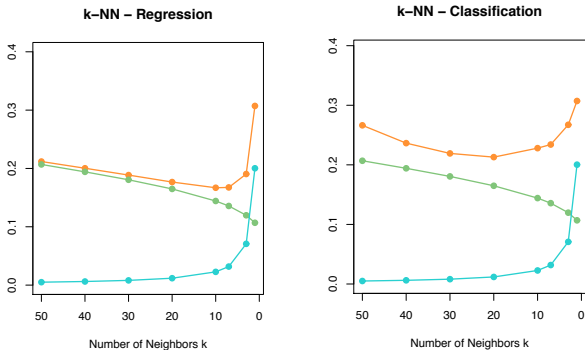- Use **0-1 loss** to measure **Err** for the **classification task**.

**Expected prediction error** as $k$ varies



- Orange curve: **expected prediction error**
- Green curve: **squared bias**
- Blue curve: **variance**

Note prediction error curves are not the same as the loss functions differ.

**Expected prediction error** as $k$ varies



- Orange curve: **expected prediction error**
- Green curve: **squared bias**
- Blue curve: **variance**

Note prediction error curves are not the same as the loss functions differ.
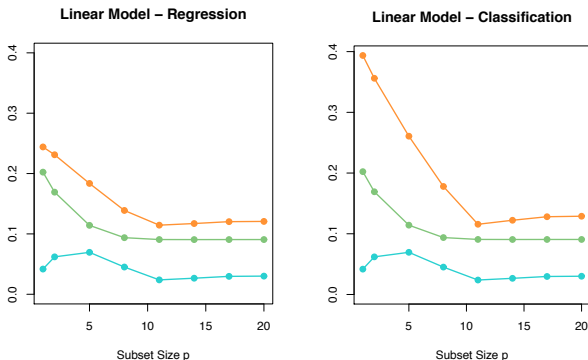
**The Set-up**

- Have $n = 80$ observations and $p = 20$ predictors.

- $X$ is uniformly distributed in $[0, 1]^{20}$ and

$$Y = \begin{cases} 1 & \text{if } \sum_{j=1}^{10} X_j > 5 \\ 0 & \text{otherwise} \end{cases}$$

- Use best subset linear regression of size $p$ for classification and regression tasks.

- Use **squared error loss** to measure **Err** for the **regression task**.

- Use **0-1 loss** to measure **Err** for the **classification task**.

**Expected prediction error** as $p$ varies



- Orange curve: **expected prediction error**
- Green curve: **squared bias**
- Blue curve: **variance**

Note prediction error curves are not the same as the loss functions differ.

# Optimism of the Training Error Rate

- Training error $\overline{\text{err}} \ll \text{Err}_{\mathcal{T}}$ as it uses $\mathcal{T}$ for both fitting and assessment.

- **One factor is**

  The training and test input vectors

  - for $\overline{\text{err}}$ are the same.

  - while for $\text{Err}_{\mathcal{T}}$ they differ.

- Can begin to understand the optimism of $\overline{\text{err}}$ if we focus on **in-sample error**

  $$\text{Err}_{\text{in}} = \frac{1}{n} \sum_{i=1}^{n} E_{Y'}[\, L(y'_i, \hat{f}(x_i))|\mathcal{T}\,]$$

  where expectation is over new responses $y'_i$ at each training point $x_i$.

- Training error $\overline{\text{err}} \ll \text{Err}_{\mathcal{T}}$ as it uses $\mathcal{T}$ for both fitting and assessment.

- **One factor is**

  The training and test input vectors

  - for $\overline{\text{err}}$ are the same.

  - while for $\text{Err}_{\mathcal{T}}$ they differ.

- Can begin to understand the optimism of $\overline{\text{err}}$ if we focus on **in-sample error**

$$\text{Err}_{\text{in}} = \frac{1}{n} \sum_{i=1}^{n} E_{Y'}[\, L(y_i', \hat{f}(x_i))|\mathcal{T}\,]$$

  where expectation is over new responses $y_i'$ at each training point $x_i$.

- Training error $\overline{\text{err}} \ll \text{Err}_{\mathcal{T}}$ as it uses $\mathcal{T}$ for both fitting and assessment.

- **One factor is**

  The training and test input vectors
  - for $\overline{\text{err}}$ are the same.
  - while for $\text{Err}_{\mathcal{T}}$ they differ.

- Can begin to understand the optimism of $\overline{\text{err}}$ if we focus on **in-sample error**

$$\text{Err}_{\text{in}} = \frac{1}{n} \sum_{i=1}^{n} E_{Y'}[\, L(y_i', \hat{f}(x_i)) | \mathcal{T} \,]$$

  where expectation is over new responses $y_i'$ at each training point $x_i$.

- Define the **optimism** as

$$\text{op} = \text{Err}_{\text{in}} - \overline{\text{err}}$$

- The **average optimism** is

$$\omega = \text{E}_{\mathbf{y}}[\text{op}]$$

where
  - the training input vectors are held fixed,
  - the expectation is over the training output values.

- For many loss functions

$$\omega = \frac{1}{n} \sum_{i=1}^{n} \text{Cov}(\hat{y}_i, y_i)$$

$$\omega = \frac{1}{n} \sum_{i=1}^{n} \text{Cov}(\hat{y}_i, y_i)$$

- The more strongly $y_i$ affects its prediction $\hat{y}_i$ the larger $\omega$.

- The larger $\omega$ the greater the optimism of $\overline{\text{err}}$.

- In summary get the important relation

$$\mathsf{E}_{\mathbf{y}}[\mathsf{Err}_{\mathsf{in}}] = \mathsf{E}_{\mathbf{y}}[\overline{\text{err}}] + \frac{1}{n} \sum_{i=1}^{n} \text{Cov}(\hat{y}_i, y_i)$$

**Option 1**

- Estimate the optimism and add it to $\overline{\text{err}}$

- The methods $C_p, \text{AIC}, \text{BIC}$ work in this way for a special class estimates.

- Can use **in-sample** error for model selection but not a good estimate of Err.

**Option 1**

- Estimate the optimism and add it to $\overline{\text{err}}$

- The methods $C_p, \text{AIC}, \text{BIC}$ work in this way for a special class estimates.

- Can use **in-sample** error for model selection but not a good estimate of Err.

**Option 2**

- Use cross-validation and bootstrap as direct estimates of the **extra-sample** Err.

# Estimates of In-Sample Prediction Error

- If $\hat{y}_i$ is obtained by a linear fit with $d$ inputs then

$$\sum_{i=1}^{n} \mathsf{Cov}(\hat{y}_i, y_i) = d\,\sigma_\epsilon^2$$

for the additive error model $Y = f(X) + \epsilon$.

- And so

$$\mathsf{E}_\mathbf{y}[\mathsf{Err}_{in}] = \mathsf{E}_\mathbf{y}[\overline{\mathsf{err}}] + 2\,\frac{d}{n}\sigma_\epsilon^2$$

- Adapting this expression leads to the $C_p$ statistic

$$C_p = \overline{\mathsf{err}} + 2\,\frac{d}{n}\hat{\sigma}_\epsilon^2$$

where $\hat{\sigma}_\epsilon^2$ is an estimate of the noise variance.

- If $\hat{y}_i$ is obtained by a linear fit with $d$ inputs then

$$\sum_{i=1}^{n} \text{Cov}(\hat{y}_i, y_i) = d\,\sigma_\epsilon^2$$

  for the additive error model $Y = f(X) + \epsilon$.

- And so

$$\boxed{\mathsf{E}_\mathbf{y}[\text{Err}_{\text{in}}] = \mathsf{E}_\mathbf{y}[\overline{\text{err}}] + 2\,\frac{d}{n}\sigma_\epsilon^2}$$

- Adapting this expression leads to the $C_p$ statistic

$$C_p = \overline{\text{err}} + 2\,\frac{d}{n}\hat{\sigma}_\epsilon^2$$

  where $\hat{\sigma}_\epsilon^2$ is an estimate of the noise variance.

- If $\hat{y}_i$ is obtained by a linear fit with $d$ inputs then

$$\sum_{i=1}^{n} \text{Cov}(\hat{y}_i, y_i) = d\,\sigma_\epsilon^2$$

  for the additive error model $Y = f(X) + \epsilon$.

- And so

$$\mathsf{E}_{\mathbf{y}}[\text{Err}_{\text{in}}] = \mathsf{E}_{\mathbf{y}}[\overline{\text{err}}] + 2\,\frac{d}{n}\sigma_\epsilon^2$$

- Adapting this expression leads to the $C_p$ statistic

$$C_p = \overline{\text{err}} + 2\,\frac{d}{n}\hat{\sigma}_\epsilon^2$$

  where $\hat{\sigma}_\epsilon^2$ is an estimate of the noise variance.

rewards the fit between the model and the data

$$\text{AIC} = -\frac{2}{n}\text{loglik} + 2\frac{d}{n}$$

penalty for including extra predictors in the model

where

$$\text{loglik} = \sum_{i=1}^{n} \log P_{\hat{\theta}}(y_i)$$

and $\hat{\theta}$ is the MLE of $\theta$.

**Note**: AIC can be seen as an estimate of $\text{Err}_{\text{in}}$ in this case with a log-likelihood loss.

- Have a set of models $f_\alpha(x)$ indexed by $\alpha$,

- $\overline{\mathrm{err}}$ is the training error,

- $d(\alpha)$ the # of parameters for each model.

then

$$\mathsf{AIC}(\alpha) = \overline{\mathrm{err}}(\alpha) + 2\frac{d(\alpha)}{n}\hat{\sigma}_\epsilon^2$$
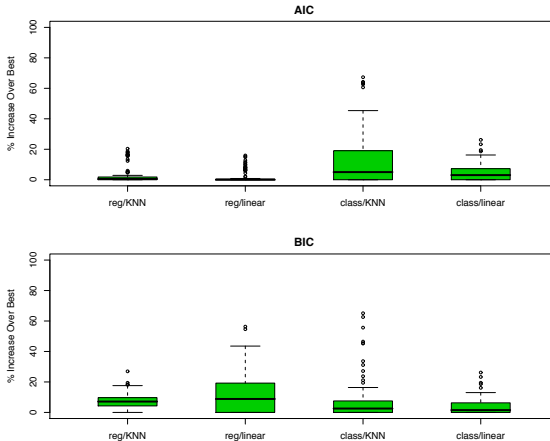
**Note**: AIC can be seen as an estimate of $\mathrm{Err}_{\mathrm{in}}$ in this case with a squared-error loss.

- Classifier is a logistic regression function with an expansion of $M$ spline basis functions.

- AIC is used to estimate $\mathrm{Err}_{\mathsf{in}}$ with a log-likelihood loss,

- AIC does well except when $M = 256$ is large and $n = 1000$.

# How well AIC and BIC perform wrt model selection



- Boxplots show $100 \dfrac{\mathrm{Err}(\hat{\alpha}) - \min_{\alpha} \mathrm{Err}(\alpha)}{\max_{\alpha} \mathrm{Err}(\alpha) - \min_{\alpha} \mathrm{Err}(\alpha)}$ where $\hat{\alpha}$ is the best parameter found via the selection method under investigation.
- 100 training sets were used.

# The Effective Number of Parameters

- For **regularized fitting** need to generalize the concept of *number of parameters*.

- Consider regularized linear fitting - ridge regression, cubic smoothing splines

$$\hat{y} = \mathbf{S}\, y$$

where

**1** $y = (y_1, y_2, \ldots, y_n)^t$ is the vector of training outputs,

**2** $\hat{y} = (\hat{y}_1, \ldots, \hat{y}_n)$ is the vector of predictions,

**3** $\mathbf{S}$ is an $n \times n$ matrix - depends on $x_1, \ldots, x_n$ but not $y_1, \ldots, y_n$.

- Then the **effective number of parameters** is defined as

$$\text{df}(\mathbf{S}) = \text{trace}(\mathbf{S})$$

- For **regularized fitting** need to generalize the concept of *number of parameters*.

- Consider regularized linear fitting - ridge regression, cubic smoothing splines

$$\hat{y} = \mathbf{S}\, y$$

where

1. $y = (y_1, y_2, \ldots, y_n)^t$ is the vector of training outputs,

2. $\hat{y} = (\hat{y}_1, \ldots, \hat{y}_n)$ is the vector of predictions,

3. $\mathbf{S}$ is an $n \times n$ matrix - depends on $x_1, \ldots, x_n$ but not $y_1, \ldots, y_n$.

- Then the **effective number of parameters** is defined as

$$\mathrm{df}(\mathbf{S}) = \mathrm{trace}(\mathbf{S})$$

- If $y$ arises from an additive-error model

$$Y = f(X) + \epsilon$$

with $\mathsf{Var}(\epsilon) = \sigma_\epsilon^2$ then

$$\sum_{i=1}^{n} \mathsf{Cov}(\hat{y}_i, y_i) = \mathsf{trace}(\mathbf{S})\sigma_\epsilon^2$$

- The more general definition of **effective dof** is then

$$\mathsf{df}(\hat{y}) = \frac{\sum_{i=1}^{n} \mathsf{Cov}(\hat{y}_i, y_i)}{\sigma_\epsilon^2}$$

- If $y$ arises from an additive-error model

$$Y = f(X) + \epsilon$$

with $\mathsf{Var}(\epsilon) = \sigma_\epsilon^2$ then

$$\sum_{i=1}^{n} \mathsf{Cov}(\hat{y}_i, y_i) = \mathsf{trace}(\mathbf{S})\sigma_\epsilon^2$$

- The more general definition of **effective dof** is then

$$\mathsf{df}(\hat{y}) = \frac{\sum_{i=1}^{n} \mathsf{Cov}(\hat{y}_i, y_i)}{\sigma_\epsilon^2}$$

# The Bayesian Approach and BIC

**Bayesian Information Criterion**

$$\boxed{\text{BIC} = -2\,\text{loglik} + \log(n)\,d}$$

- Assuming **Gaussian model** and known variance $\sigma_\epsilon^2$ then

$$-2\,\text{loglik} = \frac{1}{\sigma_\epsilon^2} \sum_{i=1}^{n} (y_i - \hat{h}(x_i))^2 = \frac{n\,\overline{\text{err}}}{\sigma_\epsilon^2}$$

and

$$\text{BIC} = \frac{n}{\sigma_\epsilon^2} \left( \overline{\text{err}} + \log(n)\,\frac{d}{n}\sigma_\epsilon^2 \right)$$

- Note BIC $\propto$ AIC, but BIC penalizes complex model more heavily than AIC.

**Bayesian Information Criterion**

$$\boxed{\text{BIC} = -2\,\text{loglik} + \log(n)\,d}$$

- Assuming **Gaussian model** and known variance $\sigma_\epsilon^2$ then

$$-2\,\text{loglik} = \frac{1}{\sigma_\epsilon^2} \sum_{i=1}^{n} (y_i - \hat{h}(x_i))^2 = \frac{n\,\overline{\text{err}}}{\sigma_\epsilon^2}$$

and

$$\boxed{\text{BIC} = \frac{n}{\sigma_\epsilon^2} \left( \overline{\text{err}} + \log(n)\,\frac{d}{n}\sigma_\epsilon^2 \right)}$$

- Note BIC $\propto$ AIC, but BIC penalizes complex model more heavily than AIC.

- **Starting point:**
  Have $\{\mathcal{M}_1, \ldots, \mathcal{M}_M\}$ a set of candidate models and their corresponding parameters $\theta_1, \ldots, \theta_m$.

- **Goal:**
  Choose the best model $\mathcal{M}_i$.

- **How:**

  - Have training data $\mathbf{Z} = \{(x_1, y_1), \ldots, (x_n, y_n)\}$

  - Have priors $p(\theta_m | \mathcal{M}_m)$.

  - The posterior of model $\mathcal{M}_m$ is

    $$P(\mathcal{M}_m | \mathbf{Z}) \propto P(\mathcal{M}_m)\, p(\mathbf{Z} | \mathcal{M}_m)$$
    $$\propto P(\mathcal{M}_m) \int p(\mathbf{Z} | \theta_m, \mathcal{M}_m)\, p(\theta_m | \mathcal{M}_m)\, d\theta_m$$

- **Starting point:**

  Have $\{\mathcal{M}_1, \ldots, \mathcal{M}_M\}$ a set of candidate models and their corresponding parameters $\theta_1, \ldots, \theta_m$.

- **Goal:**

  Choose the best model $\mathcal{M}_i$.

- **How:**

  - Have training data $\mathbf{Z} = \{(x_1, y_1), \ldots, (x_n, y_n)\}$

  - Have priors $p(\theta_m | \mathcal{M}_m)$.

  - The posterior of model $\mathcal{M}_m$ is

    $$P(\mathcal{M}_m | \mathbf{Z}) \propto P(\mathcal{M}_m) \, p(\mathbf{Z} | \mathcal{M}_m)$$
    $$\propto P(\mathcal{M}_m) \int p(\mathbf{Z} | \theta_m, \mathcal{M}_m) \, p(\theta_m | \mathcal{M}_m) \, d\theta_m$$

- **Starting point:**

  Have $\{\mathcal{M}_1, \ldots, \mathcal{M}_M\}$ a set of candidate models and their corresponding parameters $\theta_1, \ldots, \theta_m$.

- **Goal:**

  Choose the best model $\mathcal{M}_i$.

- **How:**

  - Have training data $\mathbf{Z} = \{(x_1, y_1), \ldots, (x_n, y_n)\}$

  - Have priors $p(\theta_m | \mathcal{M}_m)$.

  - The posterior of model $\mathcal{M}_m$ is

  $$P(\mathcal{M}_m | \mathbf{Z}) \propto P(\mathcal{M}_m) \, p(\mathbf{Z} | \mathcal{M}_m)$$
  $$\propto P(\mathcal{M}_m) \int p(\mathbf{Z} | \theta_m, \mathcal{M}_m) \, p(\theta_m | \mathcal{M}_m) \, d\theta_m$$

- The posterior of model $\mathcal{M}_m$ is

$$P(\mathcal{M}_m|\theta_m) \propto P(\mathcal{M}_m) \int p(\mathbf{Z}|\theta_m, \mathcal{M}_m)\, p(\theta_m|\mathcal{M}_m)\, d\theta_m$$

- Usually assume uniform prior: $P(\mathcal{M}_m) = 1/M$.

- Approximate the above integral by simplification and Laplace approximation to get

$$\log P(\mathbf{Z}|\mathcal{M}_m) = \log P(\mathbf{Z}|\hat{\theta}_m, \mathcal{M}_m) - \frac{d_m}{2} \log n + O(1)$$

where $\hat{\theta}_m$ is a MLE and $d_m$ is # free parameters in $\mathcal{M}_m$.

- Then BIC $\propto -2 \log P(\mathcal{M}_m|\mathbf{Z})$

- The posterior of model $\mathcal{M}_m$ is

$$P(\mathcal{M}_m|\theta_m) \propto P(\mathcal{M}_m) \int p(\mathbf{Z}|\theta_m, \mathcal{M}_m) \, p(\theta_m|\mathcal{M}_m) \, d\theta_m$$

- Usually assume uniform prior: $P(\mathcal{M}_m) = 1/M$.

- Approximate the above integral by simplification and Laplace approximation to get

$$\log P(\mathbf{Z}|\mathcal{M}_m) = \log P(\mathbf{Z}|\hat{\theta}_m, \mathcal{M}_m) - \frac{d_m}{2} \log n + O(1)$$

where $\hat{\theta}_m$ is a MLE and $d_m$ is # free parameters in $\mathcal{M}_m$.

- Then BIC $\propto -2 \log P(\mathcal{M}_m|\mathbf{Z})$

- The posterior of model $\mathcal{M}_m$ is

$$P(\mathcal{M}_m|\theta_m) \propto P(\mathcal{M}_m) \int p(\mathbf{Z}|\theta_m, \mathcal{M}_m)\, p(\theta_m|\mathcal{M}_m)\, d\theta_m$$

- Usually assume uniform prior: $P(\mathcal{M}_m) = 1/M$.

- Approximate the above integral by simplification and Laplace approximation to get

$$\log P(\mathbf{Z}|\mathcal{M}_m) = \log P(\mathbf{Z}|\hat{\theta}_m, \mathcal{M}_m) - \frac{d_m}{2} \log n + O(1)$$

where $\hat{\theta}_m$ is a MLE and $d_m$ is $\#$ free parameters in $\mathcal{M}_m$.

- Then BIC $\propto -2 \log P(\mathcal{M}_m|\mathbf{Z})$

- The posterior of model $\mathcal{M}_m$ is

$$P(\mathcal{M}_m|\theta_m) \propto P(\mathcal{M}_m) \int p(\mathbf{Z}|\theta_m, \mathcal{M}_m)\, p(\theta_m|\mathcal{M}_m)\, d\theta_m$$

- Usually assume uniform prior: $P(\mathcal{M}_m) = 1/M$.

- Approximate the above integral by simplification and Laplace approximation to get

$$\log P(\mathbf{Z}|\mathcal{M}_m) = \log P(\mathbf{Z}|\hat{\theta}_m, \mathcal{M}_m) - \frac{d_m}{2} \log n + O(1)$$

  where $\hat{\theta}_m$ is a MLE and $d_m$ is # free parameters in $\mathcal{M}_m$.

- Then BIC $\propto -2 \log P(\mathcal{M}_m|\mathbf{Z})$

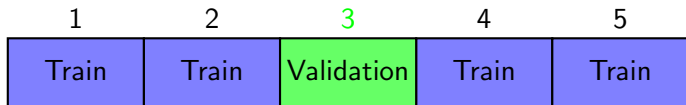- If $\mathcal{M}_{\text{true}} \in \{\mathcal{M}_1, \ldots, \mathcal{M}_M\}$ then as $n \to \infty$

  - BIC will select $\mathcal{M}_{\text{true}}$. ✓

  - AIC will not. It tends to choose too complex models as $n \to \infty$. ✗

- However, when $n$ is small

  - BIC often chooses models which are too simple. ✗

- If $\mathcal{M}_{\mathsf{true}} \in \{\mathcal{M}_1, \ldots, \mathcal{M}_M\}$ then as $n \to \infty$

  - BIC will select $\mathcal{M}_{\mathsf{true}}$. ✓

  - AIC will not. It tends to choose too complex models as $n \to \infty$. ✗

- However, when $n$ is small

  - BIC often chooses models which are too simple. ✗

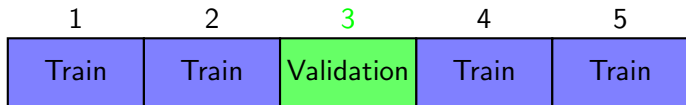# Cross-Validation

# $K$-Fold Cross-Validation

**General Approach**

- Split the data into $K$ roughly equal-size parts.

| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| Train | Train | Validation | Train | Train |

- For the $k$th part calculate the prediction error of the model fit using the other $K-1$ parts.

- Do this for $k = 1, 2, \ldots, K$ and combine the $K$ estimates of the prediction error.

**General Approach**

- Split the data into $K$ roughly equal-size parts.

| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| Train | Train | Validation | Train | Train |

- For the $k$th part calculate the prediction error of the model fit using the other $K - 1$ parts.

- Do this for $k = 1, 2, \ldots, K$ and combine the $K$ estimates of the prediction error.

**When and why**

- It is applied when labelled training data is relatively sparse.

- This method directly estimates $\mathsf{Err} = \mathsf{E}[L(Y, \hat{f}(X))]$.

- The mapping $\kappa : \{1, \ldots, n\} \to \{1, \ldots, K\}$ indicates observation $i$ belongs to partition $\kappa(i)$.

- $\hat{f}^{-k}(x)$ is the function fitted with the $k$th part of the data removed.

- **Cross-validation** estimate of the prediction error is

$$\mathrm{CV}(\hat{f}) = \frac{1}{n} \sum_{i=1}^{n} L(y_i, \hat{f}^{-\kappa(i)}(x_i))$$

- Typical choices for $K$ are 5 or 10.

- The case $K = n$ is known as leave-one-out cross-validation.

- The mapping $\kappa : \{1, \ldots, n\} \to \{1, \ldots, K\}$ indicates observation $i$ belongs to partition $\kappa(i)$.

- $\hat{f}^{-k}(x)$ is the function fitted with the $k$th part of the data removed.

- **Cross-validation** estimate of the prediction error is

$$\mathrm{CV}(\hat{f}) = \frac{1}{n} \sum_{i=1}^{n} L(y_i, \hat{f}^{-\kappa(i)}(x_i))$$

- Typical choices for $K$ are 5 or 10.

- The case $K = n$ is known as leave-one-out cross-validation.

- The mapping $\kappa : \{1, \ldots, n\} \to \{1, \ldots, K\}$ indicates observation $i$ belongs to partition $\kappa(i)$.

- $\hat{f}^{-k}(x)$ is the function fitted with the $k$th part of the data removed.

- **Cross-validation** estimate of the prediction error is

$$\mathrm{CV}(\hat{f}) = \frac{1}{n} \sum_{i=1}^{n} L(y_i, \hat{f}^{-\kappa(i)}(x_i))$$

- Typical choices for $K$ are 5 or 10.

- The case $K = n$ is known as leave-one-out cross-validation.

- Have models $f(x, \alpha)$ indexed by a parameter $\alpha$.

- $\hat{f}^{-k}(x, \alpha)$ is $\alpha$th model fit with $k$th part of the data removed.

- Then define

$$\mathrm{CV}(\hat{f}, \alpha) = \frac{1}{n} \sum_{i=1}^{n} L(y_i, \hat{f}^{-\kappa(i)}(x_i, \alpha))$$

- Choose the model

$$\hat{\alpha} = \arg\min_{\alpha} \mathrm{CV}(\hat{f}, \alpha)$$

- Have models $f(x, \alpha)$ indexed by a parameter $\alpha$.

- $\hat{f}^{-k}(x, \alpha)$ is $\alpha$th model fit with $k$th part of the data removed.

- Then define

$$\mathrm{CV}(\hat{f}, \alpha) = \frac{1}{n} \sum_{i=1}^{n} L(y_i, \hat{f}^{-\kappa(i)}(x_i, \alpha))$$

- Choose the model

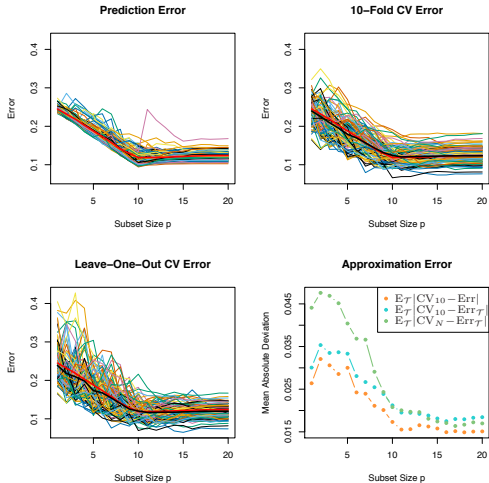$$\hat{\alpha} = \arg\min_{\alpha} \mathrm{CV}(\hat{f}, \alpha)$$

**Intuition says**

- When $K = 5$ or 10 then $CV(\hat{f}) \approx \mathrm{Err}$ as training sets for each fold are fairly different.

- When $K = n$ then $CV(\hat{f}) \approx \mathrm{Err}_{\mathcal{T}}$ as training sets for each fold are almost identical.

**Intuition says**

- When $K = 5$ or 10 then $CV(\hat{f}) \approx \mathrm{Err}$ as training sets for each fold are fairly different.

- When $K = n$ then $CV(\hat{f}) \approx \mathrm{Err}_{\mathcal{T}}$ as training sets for each fold are almost identical.

**Book's simulation experiments say**

- Cross-validation, really only effectively estimates Err.

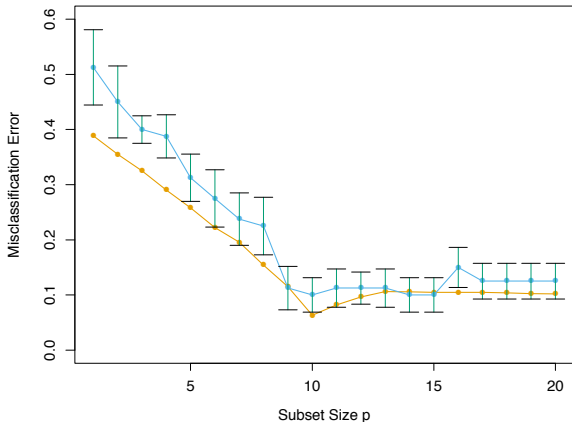# What quantity does $K$-fold validation estimate?



- **Thick red curve: Err**
- **Thick black curve:** $\mathbf{E}_{\mathcal{T}}[\mathbf{CV}_K]$

- When $K = n$
  - $CV(\hat{f})$ is approx an unbiased estimate of Err. ✓

  - $CV(\hat{f})$ has high variance as the $n$ *training sets* are similar. ✗

  - Computational burden is high. ✗ (except for a few exceptions)

- When $K = 5$ (is lowish)
  - $CV(\hat{f})$ has low variance. ✓

  - $CV(\hat{f})$ is potentially an upward biased estimate of Err. ✗

    Only occurs if at each fold there is not enough training data to fit a good model.

- Orange curve: $\mathrm{Err}_{\mathcal{T}}$

- Blue curve: $\mathrm{CV}_{10}(\hat{f})$

**Right & Wrong way to do Cross-validation**

What's wrong with this strategy?

1. **Screen the predictors** Find a subset of *good* predictors that are correlated with the class labels.

2. Build a classifier based on the subset of good predictors.

3. Perform cross-validation to estimate the unknown tuning parameters and to estimate $\mathrm{Err}$ of the final model.

What's wrong with this strategy?

1. **Screen the predictors** Find a subset of *good* predictors that are correlated with the class labels.

2. **Build a classifier** based on the subset of good predictors.

3. Perform cross-validation to estimate the unknown tuning parameters and to estimate Err of the final model.

**What's wrong with this strategy?**

1. **Screen the predictors** Find a subset of *good* predictors that are correlated with the class labels.

2. **Build a classifier** based on the subset of good predictors.

3. **Perform cross-validation** to estimate the unknown tuning parameters and to estimate $\mathrm{Err}$ of the final model.

What's wrong with this strategy?

1. **Screen the predictors** Find a subset of *good* predictors that are correlated with the class labels.

2. **Build a classifier** based on the subset of good predictors.

3. **Perform cross-validation** to estimate the unknown tuning parameters and to estimate $\mathrm{Err}$ of the final model.

The good predictors were chosen after seeing all the data.

1. Divide the samples into $K$ groups randomly.

2. For each fold $k = 1, \ldots, K$

   - Find a subset of *good* predictors using all the samples minus the $k$th fold.

   - Build a classifier using all the samples minus the $k$th fold.

   - Use the classifier to predict the labels for the samples in the $k$th fold.
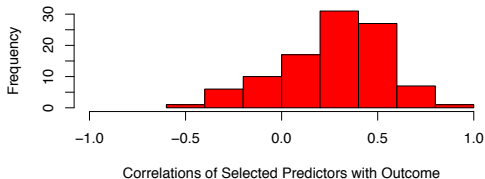
**Set-up**

- Have a binary classification problem.

- $n = 50$ with an equal number of points from each class.

- Have $p = 5000$ quantitative predictors that are independent of the class labels.

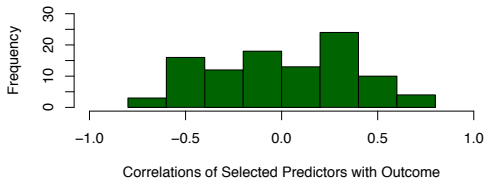- The true error rate of any classifier is 50%.

✗ If one performs pre-selection of 100 predictors and then builds a 1-nn classifier the average CV error rate was 3% over 50 simulations ! ✗

**Wrong way**

Frequency vs Correlations of Selected Predictors with Outcome

**Right way**

Frequency vs Correlations of Selected Predictors with Outcome

- Cross-validation must be applied to the entire sequence of modelling steps.

- Samples must be **left out** before any selection or filtering is applied which uses the labels.

- One exception: An unsupervised screening step can use all the samples.

- Cross-validation must be applied to the entire sequence of modelling steps.

- Samples must be **left out** before any selection or filtering is applied which uses the labels.

- One exception: An unsupervised screening step can use all the samples.

# Bootstrap Method

- Have a training set $\mathbf{Z} = (z_1, z_2, \ldots, z_n)$ with each $z_i = (x_i, y_i)$.

- The **bootstrap** idea is

  for $b = 1, 2, \ldots, B$

  1. Randomly draw $n$ samples with replacement from $\mathbf{Z}$ to get $\mathbf{Z}^{*b}$ that is

     $$\mathbf{Z}^{*b} = (z_{b_1}, z_{b_2}, \ldots, z_{b_n}) \quad \text{with } b_i \in \{1, \ldots, n\}$$

  2. Refit the model using $\mathbf{Z}^{*b}$ to get $S(\mathbf{Z}^{*b})$

  Examine the behaviour of the $B$ fits

  $$S(\mathbf{Z}^{*1}), S(\mathbf{Z}^{*2}), \ldots, S(\mathbf{Z}^{*B}).$$

- Have a training set $\mathbf{Z} = (z_1, z_2, \ldots, z_n)$ with each $z_i = (x_i, y_i)$.

- The **bootstrap** idea is

for $b = 1, 2, \ldots, B$

   **1** Randomly draw $n$ samples with replacement from $\mathbf{Z}$ to get $\mathbf{Z}^{*b}$ that is

$$\mathbf{Z}^{*b} = (z_{b_1}, z_{b_2}, \ldots, z_{b_n}) \quad \text{with } b_i \in \{1, \ldots, n\}$$

   **2** Refit the model using $\mathbf{Z}^{*b}$ to get $S(\mathbf{Z}^{*b})$

Examine the behaviour of the $B$ fits

$$S(\mathbf{Z}^{*1}), S(\mathbf{Z}^{*2}), \ldots, S(\mathbf{Z}^{*B}).$$

- For example its variance

$$\widehat{\mathrm{Var}}[S(\mathbf{Z})] = \frac{1}{B-1} \sum_{b=1}^{B} (S(Z^{*b}) - \bar{S}^*)^2$$

where

$$\bar{S}^* = \frac{1}{B} \sum_{b=1}^{B} S(Z^{*b})$$

**Attempt 1**

$$\widehat{\mathrm{Err}}_{\mathrm{boot}} = \frac{1}{B}\frac{1}{n}\sum_{b=1}^{B}\sum_{i=1}^{n} L(y_i, \hat{f}^{*b}(x_i))$$

where $\hat{f}^{*b}(x_i)$ is the predicted value at $x_i$ using the model computed from $\mathbf{Z}^{*b}$.

- Why is this not a good estimate??
  - Overlap between training and test sets

- How could we do better?
  - Mimic cross-validation

**Attempt 1**

$$\widehat{\mathrm{Err}}_{\mathrm{boot}} = \frac{1}{B}\frac{1}{n}\sum_{b=1}^{B}\sum_{i=1}^{n}L(y_i, \hat{f}^{*b}(x_i))$$

where $\hat{f}^{*b}(x_i)$ is the predicted value at $x_i$ using the model computed from $\mathbf{Z}^{*b}$.

- Why is this not a good estimate??
  - Overlap between training and test sets

- How could we do better?
  - Mimic cross-validation

**Attempt 1**

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

where $\hat{f}^{*b}(x_i)$ is the predicted value at $x_i$ using the model computed from $\mathbf{Z}^{*b}$.

- Why is this not a good estimate??
  - Overlap between training and test sets

- How could we do better?
  - Mimic cross-validation

**Attempt 1**

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

where $\hat{f}^{*b}(x_i)$ is the predicted value at $x_i$ using the model computed from $\mathbf{Z}^{*b}$.

- Why is this not a good estimate??
  - Overlap between training and test sets

- How could we do better?
  - Mimic cross-validation

**Attempt 2**: Leave-one-out bootstrap

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

where $C^{-i}$ is the set of bootstrap samples $b$ not containing observation $i$.

- Either make

- Make $B$ large enough so $|C^{-i}| > 0$ for all $i$ **or**

- Omit observation $i$ from testing if $|C^{-i}| = 0$.

**Attempt 2**: Leave-one-out bootstrap

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

where $C^{-i}$ is the set of bootstrap samples $b$ not containing observation $i$.

- Either make

- Make $B$ large enough so $|C^{-i}| > 0$ for all $i$ **or**

- Omit observation $i$ from testing if $|C^{-i}| = 0$.

- **Pros**:
  1. avoids the overfitting problem of $\widehat{\mathrm{Err}}_{\mathrm{boot}}$

- **Cons**:
  1. Has the training-set-size bias of cross-validation
  2. The $P(\text{observation } i \in \mathbf{Z}^{*b})$ is

  $$1 - \left(1 - \frac{1}{n}\right)^n \approx 1 - e^{-1} = .632$$

     Therefore the average number of distinct observations in $\mathbf{Z}^{*b}$ is $.632\, n$.
  3. $\widehat{\mathrm{Err}}^{(1)}$'s bias is thus similar to twofold cross-validation.

**Attempt 3**: The .632 estimator

$$\widehat{\mathrm{Err}}^{(.632)} = .368\,\overline{\mathrm{err}} + .632\,\widehat{\mathrm{Err}}^{(1)}$$

- Compromise between the **training error** $\overline{\mathrm{err}}$ and the **leave-one-out bootstrap** estimate.

- Its derivation is not easy.

- Obviously the constant .632 relates to $P(\text{observation } i \in \mathbf{Z}^{*b})$.

The .632 estimator **does not do well** if predictor **overfits**.

**Attempt 3**: The .632 estimator

$$\widehat{\mathrm{Err}}^{(.632)} = .368\,\overline{\mathrm{err}} + .632\,\widehat{\mathrm{Err}}^{(1)}$$

- Compromise between the **training error** $\overline{\mathrm{err}}$ and the **leave-one-out bootstrap** estimate.

- Its derivation is not easy.

- Obviously the constant .632 relates to $P(\text{observation } i \in \mathbf{Z}^{*b})$.

The .632 estimator **does not do well** if predictor **overfits**.

**No-information error rate**:

$$\hat{\gamma} = \frac{1}{n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} L(y_i, \hat{f}(x_j))$$

- Estimate of the error rate of $\hat{f}$ if inputs and outputs were independent.

- Note the prediction rule, $\hat{f}$, is evaluated on all possible combinations of targets $y_i$ and predictors $x_j$.

**No-information error rate**:

$$\hat{\gamma} = \frac{1}{n^2} \sum_{i=1}^{n} \sum_{j=1}^{n} L(y_i, \hat{f}(x_j))$$

- Estimate of the error rate of $\hat{f}$ if inputs and outputs were independent.

- Note the prediction rule, $\hat{f}$, is evaluated on all possible combinations of targets $y_i$ and predictors $x_j$.

**Relative overfitting rate**:

$$\hat{R} = \frac{\widehat{\mathrm{Err}}^{(1)} - \overline{\mathrm{err}}}{\hat{\gamma} - \overline{\mathrm{err}}}$$

- $0 \leq \hat{R} \leq 1$

- $\hat{R} = 0 \implies$ no overfitting.

- $\hat{R} = 1 \implies$ overfitting equals no-information value $\hat{\gamma} - \overline{\mathrm{err}}$.

**Attempt 4**: The .632+ estimator

$$\boxed{\widehat{\mathrm{Err}}^{(.632+)} = (1 - \hat{w})\,\overline{\mathrm{err}} + \hat{w}\,\widehat{\mathrm{Err}}^{(1)}}$$

with

$$\hat{w} = \frac{.632}{1 - .368\hat{R}}$$

- $.632 \leq \hat{w} \leq 1$ as $\hat{R}$ ranges from 0 to 1.

- $\widehat{\mathrm{Err}}^{(.632+)}$ ranges from $\widehat{\mathrm{Err}}^{(.632)}$ to $\widehat{\mathrm{Err}}^{(1)}$.

- $\widehat{\mathrm{Err}}^{(.632+)}$ is a compromise between $\widehat{\mathrm{Err}}^{(.632)}$ and $\overline{\mathrm{err}}$ that depends on the amount of overfitting.

- Derivation of the above eqn is non-trivial.
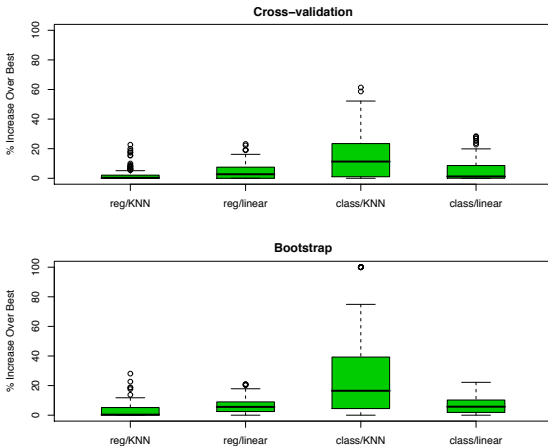
**Attempt 4**: The .632+ estimator

$$\widehat{\mathrm{Err}}^{(.632+)} = (1 - \hat{w})\,\overline{\mathrm{err}} + \hat{w}\,\widehat{\mathrm{Err}}^{(1)}$$

with

$$\hat{w} = \frac{.632}{1 - .368\hat{R}}$$

- $.632 \leq \hat{w} \leq 1$ as $\hat{R}$ ranges from 0 to 1.

- $\widehat{\mathrm{Err}}^{(.632+)}$ ranges from $\widehat{\mathrm{Err}}^{(.632)}$ to $\widehat{\mathrm{Err}}^{(1)}$.

- $\widehat{\mathrm{Err}}^{(.632+)}$ is a compromise between $\widehat{\mathrm{Err}}^{(.632)}$ and $\overline{\mathrm{err}}$ that depends on the amount of overfitting.

- Derivation of the above eqn is non-trivial.

- Boxplots show $100 \dfrac{\text{Err}(\hat{\alpha}) - \min_\alpha \text{Err}(\alpha)}{\max_\alpha \text{Err}(\alpha) - \min_\alpha \text{Err}(\alpha)}$ where $\hat{\alpha}$ is the best parameter found via the selection method under investigation.
- 100 training sets were used.