

Decompositions of the Prediction Error

Stephen Vardeman
Analytics Iowa LLC
ISU Statistics and IMSE

Quantifying the performance of a predictor

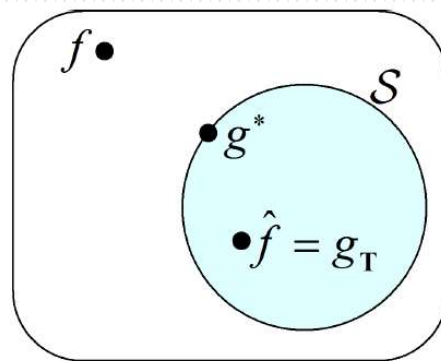
- We suppose that training cases (\mathbf{x}_i, y_i) for $i = 1, 2, \dots, N$ and test case (\mathbf{x}, y) are iid from distribution P and that predictor \hat{f} is built using the training set \mathbf{T}
- Then expected prediction loss (alternatively called the “prediction error,” “test error,” and “generalization error”) suffered using \hat{f} is

$$\text{Err} \equiv \mathbf{E}^T \mathbf{E}^{(x,y)} \left[L \left(\hat{f}(\mathbf{x}), y \right) \right]$$

- Two decompositions of this help make clear what must be controlled in order to make prediction error small

General decomposition of prediction error

- Use the notation
 - f for the (theoretically) optimal predictor
 - \mathcal{S} for a class of functions g from which a predictor can be chosen
 - $\hat{f} = g_{\mathbf{T}}$ for the training-set-dependent element of \mathcal{S} used for prediction
 - g^* for a (fixed) element of \mathcal{S} with minimum expected prediction loss
- The situation can then be pictured as below in terms of the optimal, restricted optimal, and fitted predictors



General decomposition of prediction error

- The optimal f is potentially (typically) outside \mathcal{S}
- The “closest” (in terms of prediction error) one can get to it inside \mathcal{S} is g^*
- For any fixed training set $\hat{f} = g_T$ can be no better than g^*
- As the training set varies randomly, how much better g^* is than \hat{f} varies randomly

General decomposition of prediction error

- So since here

$$\text{Err} = \mathbf{E}^T \mathbf{E}^{(x,y)} L(\hat{f}(\mathbf{x}), y) = \mathbf{E}^T \mathbf{E}^{(x,y)} L(g_T(\mathbf{x}), y)$$

we have

$$\begin{aligned} \text{Err} = & \mathbf{E}^{(x,y)} L(f(\mathbf{x}), y) + \left(\mathbf{E}^{(x,y)} L(g^*(\mathbf{x}), y) - \mathbf{E}^{(x,y)} L(f(\mathbf{x}), y) \right) \\ & + \left(\mathbf{E}^T \mathbf{E}^{(x,y)} L(g_T(\mathbf{x}), y) - \mathbf{E}^{(x,y)} L(g^*(\mathbf{x}), y) \right) \end{aligned}$$

- This can be thought of as

$$\begin{aligned} \text{Err} = & \text{minimum expected loss possible} + \text{modeling penalty} \\ & + \text{fitting penalty} \end{aligned}$$

SEL decomposition of prediction error

- A more detailed and illuminating decomposition of Err is possible for SEL
- Begin with a measure of prediction performance at input vector \mathbf{x}

$$\text{Err}(\mathbf{x}) \equiv \mathbf{E}^T \mathbf{E} \left[\left(\hat{f}(\mathbf{x}) - y \right)^2 \mid \mathbf{x} \right]$$

(noting that $\text{Err} = \mathbf{E}^X \text{Err}(\mathbf{x})$)

- Then

$$\begin{aligned} \text{Err}(\mathbf{x}) &= \mathbf{E}^T \left\{ \left(\hat{f}(\mathbf{x}) - \mathbf{E}[y \mid \mathbf{x}] \right)^2 + \mathbf{E} \left[(y - \mathbf{E}[y \mid \mathbf{x}])^2 \mid \mathbf{x} \right] \right\} \\ &= \mathbf{E}^T \left\{ \left(\hat{f}(\mathbf{x}) - \mathbf{E}^T \hat{f}(\mathbf{x}) \right)^2 + \left(\mathbf{E}^T \hat{f}(\mathbf{x}) - \mathbf{E}[y \mid \mathbf{x}] \right)^2 \right\} + \text{Var}[y \mid \mathbf{x}] \\ &= \text{Var}^T \left(\hat{f}(\mathbf{x}) \right) + \left(\mathbf{E}^T \hat{f}(\mathbf{x}) - \mathbf{E}[y \mid \mathbf{x}] \right)^2 + \text{Var}[y \mid \mathbf{x}] \end{aligned}$$

SEL decomposition of prediction error

- $\text{Var}^{\mathbf{T}}(\hat{f}(\mathbf{x}))$ is the variance of prediction at \mathbf{x}
- $\left(\mathbb{E}^{\mathbf{T}} \hat{f}(\mathbf{x}) - \mathbb{E}[y | \mathbf{x}]\right)^2$ is a kind of squared bias of prediction at \mathbf{x}
- $\text{Var}[y | \mathbf{x}]$ is an unavoidable variance in outputs at \mathbf{x}

- Clearly then

$$\text{Err} = \mathbb{E}^{\mathbf{x}} \text{Var}^{\mathbf{T}}(\hat{f}(\mathbf{x})) + \mathbb{E}^{\mathbf{x}} \left(\mathbb{E}^{\mathbf{T}} \hat{f}(\mathbf{x}) - \mathbb{E}[y | \mathbf{x}]\right)^2 + \mathbb{E}^{\mathbf{x}} \text{Var}[y | \mathbf{x}]$$

and SEL prediction error is a sum of averages (against the marginal of \mathbf{x}) of the three quantities at fixed inputs

SEL decomposition of prediction error

- Further analysis of $E^{\mathbf{x}} \left(E^{\mathbf{T}} \hat{f}(x) - E[y|x] \right)^2$ provides additional insight
- Suppose that \mathbf{T} is used to select a function $g_{\mathbf{T}}$ from some linear subspace \mathcal{S} of the space of functions h with $E^{\mathbf{x}} (h(\mathbf{x}))^2 < \infty$ and $\hat{f} = g_{\mathbf{T}}$.
Further, let

$$g^*(\mathbf{x}) = \arg \min_{g \in \mathcal{S}} E^{\mathbf{x}} (g(\mathbf{x}) - E[y|\mathbf{x}])^2$$

(the best approximation to the optimal predictor in \mathcal{S} and projection of $E[y|\mathbf{x}]$ onto \mathcal{S}).

SEL decomposition of prediction error

- It's then possible to argue that

$$\mathbb{E}^{\mathbf{x}} \left(\mathbb{E}^T \hat{f}(\mathbf{x}) - \mathbb{E}[y|\mathbf{x}] \right)^2 = \mathbb{E}^{\mathbf{x}} \left(\mathbb{E}^T \hat{f}(\mathbf{x}) - g^*(\mathbf{x}) \right)^2 + \mathbb{E}^{\mathbf{x}} \left(\mathbb{E}[y|\mathbf{x}] - g^*(\mathbf{x}) \right)^2$$

- The first term on the right is an average (across inputs) squared fitting bias
- The second term is an average (across inputs) squared model bias
- So ultimately for SEL

$$\begin{aligned} \text{Err} &= \mathbb{E}^{\mathbf{x}} \text{Var}[y|\mathbf{x}] + \mathbb{E}^{\mathbf{x}} \left(\mathbb{E}[y|\mathbf{x}] - g^*(\mathbf{x}) \right)^2 \\ &\quad + \mathbb{E}^{\mathbf{x}} \left(\mathbb{E}^T \hat{f}(\mathbf{x}) - g^*(\mathbf{x}) \right)^2 + \mathbb{E}^{\mathbf{x}} \text{Var}^T \left(\hat{f}(\mathbf{x}) \right) \end{aligned}$$

SEL decomposition of prediction error

- The SEL decomposition is related to the general one in that

$$\begin{aligned} \text{minimum expected loss possible} &= \text{expected (across } \mathbf{x} \text{) response variance} \\ &= \mathbf{E}^{\mathbf{x}} \text{Var} [y|\mathbf{x}], \end{aligned}$$

$$\begin{aligned} \text{modeling penalty} &= \text{expected (across } \mathbf{x} \text{) squared model bias} \\ &= \mathbf{E}^{\mathbf{x}} (\mathbf{E} [y|\mathbf{x}] - g^*(\mathbf{x}))^2, \end{aligned}$$

$$\begin{aligned} \text{fitting penalty} &= \left(\begin{array}{c} \text{expected (across } \mathbf{x} \text{)} \\ \text{squared fitting bias} \end{array} \right) + \left(\begin{array}{c} \text{expected (across } \mathbf{x} \text{)} \\ \text{prediction variance} \end{array} \right) \\ &= \mathbf{E}^{\mathbf{x}} \left(\mathbf{E}^T \hat{f}(\mathbf{x}) - g^*(\mathbf{x}) \right)^2 + \mathbf{E}^{\mathbf{x}} \text{Var}^T \left(\hat{f}(\mathbf{x}) \right) \end{aligned}$$

SEL decomposition of prediction error

- The facts that
 - the modeling and fitting penalties have elements of both bias and variance
 - complex predictors tend to have low bias and high variance in comparison to simple ones

lead to the necessity of balancing these elements in predictor development and the so-called **variance-bias trade-off**

- Once more, in qualitative terms, it is the matching of predictor complexity to real information content of a training set that is at issue here