

MACHINE LEARNING

TRAINING VERSUS TEST ERROR

Sebastian Engelke

MASTER IN BUSINESS ANALYTICS



**UNIVERSITÉ
DE GENÈVE**

Training versus test error

Training Data

- ▶ Fit our model \hat{f} on the training data $\{(x_1, y_1), \dots, (x_n, y_n)\}$, the n measurements of (X, Y) .
- ▶ Compute the **training error** as

$$\begin{aligned}\text{MSE}_{\text{Tr}} &= \text{RSS}/n \\ &= \frac{1}{n} \sum_{i=1}^n \{y_i - \hat{f}(x_i)\}^2.\end{aligned}$$

- ▶ Since higher model complexity always reduces this error, training error is overly **optimistic**!

Training versus test error

Training Data

- Fit our model \hat{f} on the training data $\{(x_1, y_1), \dots, (x_n, y_n)\}$, the n measurements of (X, Y) .
- Compute the **training error** as

$$\begin{aligned}\text{MSE}_{\text{Tr}} &= \text{RSS}/n \\ &= \frac{1}{n} \sum_{i=1}^n \{y_i - \hat{f}(x_i)\}^2.\end{aligned}$$

- Since higher model complexity always reduces this error, training error is overly **optimistic**!

Test Data

- The **test data** are m new samples $\{(\tilde{x}_1, \tilde{y}_1), \dots, (\tilde{x}_m, \tilde{y}_m)\}$ of (X, Y) .
- We test our model \hat{f} fitted to the training data by predicting $\hat{y}_i = \hat{f}(\tilde{x}_i)$ at all \tilde{x}_i .
- The **test/generalization error** is

$$\text{MSE}_{\text{Te}} = \frac{1}{m} \sum_{i=1}^m \{\tilde{y}_i - \hat{f}(\tilde{x}_i)\}^2.$$

- The test error is an estimate of the **expected prediction error**

$$\text{Err}_{\hat{f}} = \mathbb{E}[\{Y - \hat{f}(X)\}^2],$$

Training versus test error

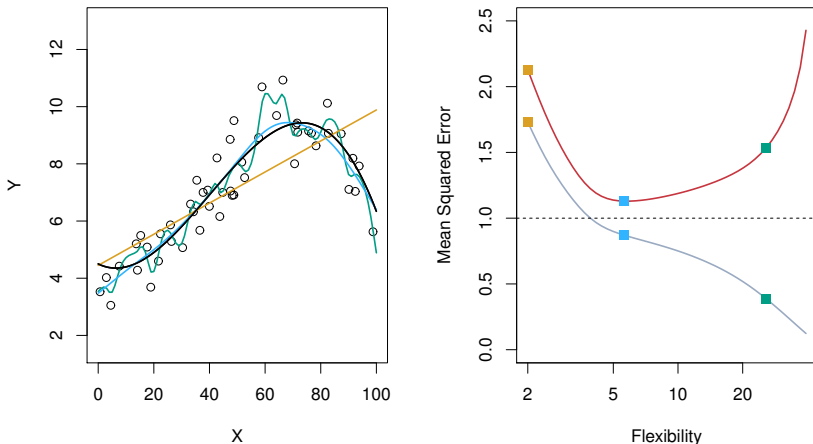


Figure: Left: Data simulated from f , the black line, and three estimates, the linear fit (orange) and higher order polynomials (blue and green). Right: Training error MSE_{Tr} (grey) and test error MSE_{Te} (red) for the three fits.

The Bias–Variance decomposition

- ▶ As usually, we assume $Y = f(X) + \epsilon$, with $E(\epsilon) = 0$ and $\text{Var}(\epsilon) = \sigma_\epsilon^2$.
- ▶ For a regression fit \hat{f} , we compute the expected prediction error at an input point $X = x_0$:

$$\begin{aligned}\text{Err}_{\hat{f}}(x_0) &= E[(Y - \hat{f}(X))^2 \mid X = x_0] \\ &\stackrel{\text{Exercise}}{=} \sigma_\epsilon^2 + [E\hat{f}(x_0) - f(x_0)]^2 + E[\hat{f}(x_0) - E\hat{f}(x_0)]^2 \\ &= \sigma_\epsilon^2 + \text{Bias}^2(\hat{f}(x_0)) + \text{Var}(\hat{f}(x_0)) \\ &= \text{Irreducible Error} + \text{Bias}^2 + \text{Variance}.\end{aligned}$$



The Bias–Variance decomposition

- ▶ As usually, we assume $Y = f(X) + \epsilon$, with $E(\epsilon) = 0$ and $\text{Var}(\epsilon) = \sigma_\epsilon^2$.
- ▶ For a regression fit \hat{f} , we compute the expected prediction error at an input point $X = x_0$:

$$\begin{aligned}\text{Err}_{\hat{f}}(x_0) &= E[(Y - \hat{f}(X))^2 \mid X = x_0] \\ &\stackrel{\text{Exercise}}{=} \sigma_\epsilon^2 + [E\hat{f}(x_0) - f(x_0)]^2 + E[\hat{f}(x_0) - E\hat{f}(x_0)]^2 \\ &= \sigma_\epsilon^2 + \text{Bias}^2(\hat{f}(x_0)) + \text{Var}(\hat{f}(x_0)) \\ &= \text{Irreducible Error} + \text{Bias}^2 + \text{Variance}.\end{aligned}$$

- ▶ The **Irreducible Error** is the error due to the noise variable ϵ and cannot be improved.
- ▶ The **Bias** is the difference between the average prediction $E\hat{f}(x_0)$ and the truth at x_0 .
- ▶ The **Variance** is the variability of the prediction $\hat{f}(x_0)$ when \hat{f} is fitted for different data sets.



The Bias–Variance decomposition

Example: The **k-Nearest-Neighbors** with tuning parameter $k \in \{1, 2, \dots\}$ gives

$$\hat{f}_k(x_0) = \frac{1}{k} \sum_{x_i \in N_k(x_0)} y_i.$$

The Bias–Variance decomposition

Example: The **k-Nearest-Neighbors** with tuning parameter $k \in \{1, 2, \dots\}$ gives

$$\hat{f}_k(x_0) = \frac{1}{k} \sum_{x_i \in N_k(x_0)} y_i.$$

We obtain from the above:

$$\begin{aligned} \text{Err}_{\hat{f}}(x_0) &= \mathbb{E}[(Y - \hat{f}_k(x_0))^2 \mid X_0 = x_0] \\ &= \sigma_\epsilon^2 + \left[f(x_0) - \frac{1}{k} \sum_{x_i \in N_k(x_0)} f(x_i) \right]^2 + \frac{\sigma_\epsilon^2}{k}. \end{aligned}$$

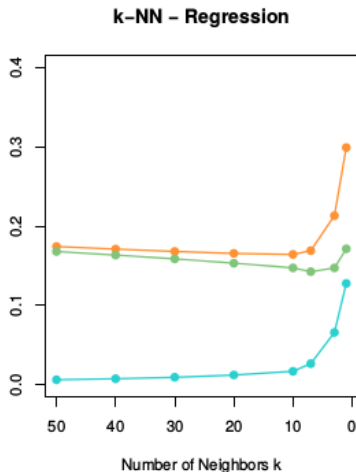


Figure: Expected prediction error (orange), squared bias (green) and variance (blue).

Bias–Variance trade off

Simple parametric methods such as **linear models**

- ▶ have low flexibility (**high bias**) but need less data to fit;
- ▶ are stable for different data sets (**low variance**);
- ▶ fit well if the data satisfies the model assumptions;
- ▶ are good for inference and **interpretation**.



Bias–Variance trade off

Simple parametric methods such as **linear models**

- ▶ have low flexibility (**high bias**) but need less data to fit;
- ▶ are stable for different data sets (**low variance**);
- ▶ fit well if the data satisfies the model assumptions;
- ▶ are good for inference and **interpretation**.

Complex (possibly non-parametric) methods such as **k-Nearest-Neighbors**

- ▶ can be very flexible (**low bias**) but need a lot of data;
- ▶ are prone to **overfitting** and vary strongly depending on the data set (**high variance**);
- ▶ are problematic in high dimensions (large p), so-called **curse of dimensionality**;
- ▶ are potentially better for **prediction**.



Training versus test error

