



Training and Test Datasets

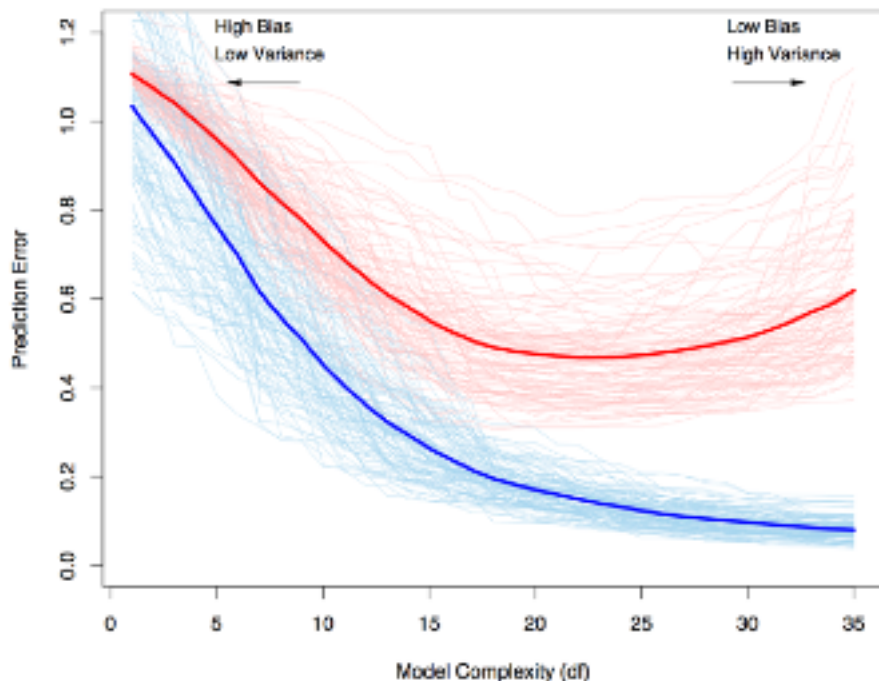


FIGURE 7.1. Behavior of test sample and training sample error as the model complexity is varied. The light blue curves show the training error $\overline{\text{err}}$, while the light red curves show the conditional test error $\text{Err}_{\mathcal{T}}$ for 100 training sets of size 50 each, as the model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

Source: Elements of
Statistical Learning

Option 1: Separate Validation Set

Training Set vs.
Validation Set vs.
Test Set

- If you have “enough” data and plan to do some model tuning, you should really partition your data into three parts – Training, Validation and Test sets.
- There is **no general rule** for how you should partition the data and it will depend on how strong the signal in your data is, but an example could be: 50% Train, 25% Validation and 25% Test.



Validation is for Model
Tuning

- The validation set is used **strictly for model tuning** (via validation of models with different parameters) and the test set is used to make a final estimate of the generalization error.

Training and Test Data Sets

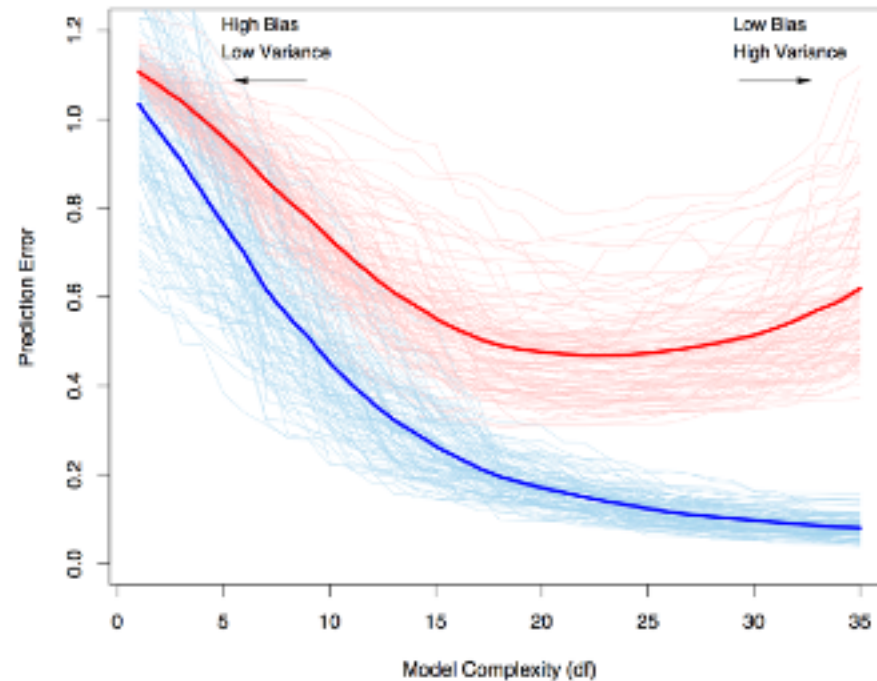


FIGURE 7.1. Behavior of test sample and training sample error as the model complexity is varied. The light blue curves show the training error \bar{err} , while the light red curves show the conditional test error Err_T for 100 training sets of size 50 each, as the model complexity is increased. The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

Source: Elements of
Statistical Learning