

# BIOS 6312: Modern Biostatistics Methodology II

**Andrew J. Spieker, Ph.D.**

Associate Professor of Biostatistics  
Vanderbilt University

Set 10: Predictive capacity of regression models

Version: 04/26/2025

# TABLE OF CONTENTS

- 1 Framing prediction problems
- 2 The bias-variance trade-off
- 3 Training and test error
- 4 Prediction for binary outcomes

## Ideas:

- Much of our focus has been on interpretation and estimations of coefficients. In this unit, our goals are about prediction:
  - ▶  $E[Y|\mathbf{X} = \mathbf{x}] = \mathbf{x}^T \boldsymbol{\beta}$  (linear model).
  - ▶ Predicted mean:  $\hat{Y}(\mathbf{x}) = \mathbf{x}^T \hat{\boldsymbol{\beta}}$ .
  - ▶ We even discussed prediction *intervals* for continuous outcomes. Our current discussion of prediction is a bit different in the sense that we want to learn about how close  $\hat{Y}(\mathbf{x})$  is to  $E[Y|\mathbf{X} = \mathbf{x}]$ .
- Key questions:
  - ▶ How do we know how well a method has done?
  - ▶ How do we compare methods?
- Questions for another course (e.g., statistical learning):
  - ▶ How do we improve our models without overfitting the data?

## Defining prediction error: Metrics

- Example:  $E[Y|\mathbf{X} = \mathbf{x}] = \mathbf{x}^T \boldsymbol{\beta}$  (linear model).
  - ▶ Predicted mean:  $\hat{Y}(\mathbf{x}) = \mathbf{x}^T \hat{\boldsymbol{\beta}}$ .
- Many ways we could think to evaluate predictive ability.
  - ▶ Mean squared error:  $E[(Y - \hat{Y})^2]$ .
  - ▶ Mean absolute error:  $E[|Y - \hat{Y}|]$ .
- In the two above examples, lower values indicate better predictive ability (difficult to compare methods across different prediction error metrics).
- For now, let's focus on mean squared error,  $E[(Y - \hat{Y})^2]$ .

## **Improving prediction error:** Naive method

- Possible method: include more covariates.
- Example: MRI data
  - ▶  $Y$ : DSST
  - ▶  $X$ : one covariate? A whole bunch of covariates?
- Compare root mean squared error (RMSE) between models of differential complexity.

## Improving prediction error: Naive method

```
1 ## Read in data set
2 dat <- read.csv("mri.csv")
3
4 ## Fit model (output not relevant for the current goal)
5 model <- regress("mean", dsst ~ male, data = dat)
6
7 ## RMSE (relevant for the current goal)
8 > model$sigma
9 [1] 12.64
```

## Improving prediction error: Naive method

```
1 ## Characterize factor variables as such
2 dat$race <- factor(dat$race)
3 dat$chf <- factor(dat$chf)
4 dat$genhlth <- factor(dat$genhlth)
5 dat$numinf <- factor(dat$numinf)
6
7 ## Fit model (output not relevant for the current goal)
8 model <- regress("mean", dsst ~ male + age + race + weight +
9                 height + packyrs + yrsquit + alcoh +
10                physact + chf + diabetes + genhlth +
11                ldl + alb + crt + sbp + aai + fev +
12                atrophy + numinf, data = dat)
13
14 ## RMSE (relevant for the current goal)
15 > model$sigma
16 [1] 10.98
```

## Improving prediction error: Naive method

- This method suffers from the fact that a reduced RMSE could be explained by the fact that your model...
  - ▶ ... better predicts the outcome in the population.
  - ▶ ... merely serves as a better fit to *your data*.
  - ▶ ... or a combination of these two!
- Not interesting to know well your model predicts the observations used to fit that model. Of greater interest: knowing how well model predicts observations *not* in your sample.

# TABLE OF CONTENTS

- 1 Framing prediction problems
- 2 The bias-variance trade-off**
- 3 Training and test error
- 4 Prediction for binary outcomes

# THE BIAS-VARIANCE TRADE-OFF

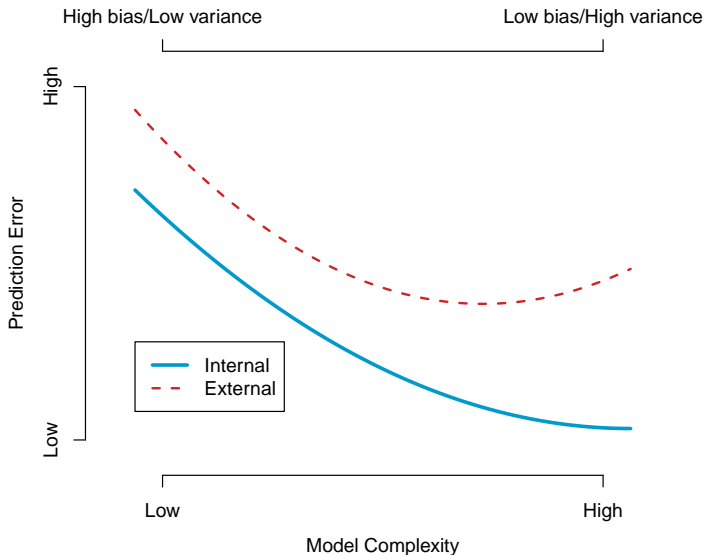
## Prediction error: Bias-variance decomposition

- Consider general case:  $E[Y|\mathbf{X} = \mathbf{x}] = f(\mathbf{x})$ .
- Estimate  $f$  on a *training* set to obtain  $\hat{f}$ .
- Consider random out-of-sample observation,  $(\mathbf{X}_0, Y_0)$ .
- Mean squared error at  $\mathbf{X}_0$  is:

$$\begin{aligned} E[(Y_0 - \hat{f}(\mathbf{X}_0))^2] &= E[(Y_0 - E[Y_0])^2] + E[(E[Y_0] - E[\hat{f}(\mathbf{X}_0)])^2] \\ &\quad + E[(\hat{f}(\mathbf{X}_0) - E[\hat{f}(\mathbf{X}_0)])^2] \\ &= \sigma_{Y_0}^2 + \text{Bias}^2(\hat{f}(\mathbf{X}_0)) + \text{Var}(\hat{f}(\mathbf{X}_0)). \end{aligned}$$

- Prediction error: trade-off between bias and variance.
  - ▶ **Bias**: Model *correctly* captures  $\mathbf{x}/E[Y|\mathbf{X} = \mathbf{x}]$  relationship?
  - ▶ **Variance**: Model *precisely* captures relationship?
- *Training* error (internal) does not estimate the *test* error (external) well because it does not properly account for model complexity.

# THE BIAS-VARIANCE TRADE-OFF



# TABLE OF CONTENTS

- 1 Framing prediction problems
- 2 The bias-variance trade-off
- 3 Training and test error**
- 4 Prediction for binary outcomes

## Estimating test error: Simple method

- 1 Split your sample (say, 80/20); label the larger portion as *training* set and the other as a *test* set.
- 2 Fit your model on the training set.
- 3 Use fitted model to predict the values on the test set and estimate the test prediction error.

**Idea:** Mimic the process of getting two independent sets of data.

## R: Sample splitting

```
1 ## Sample size
2 n <- dim(dat)[1]
3
4 ## Set see for reproducibility
5 set.seed(6312)
6
7 ## Random split
8 samp <- sample(c(rep(0, 588), rep(1,147)), replace = FALSE)
9 dat0 <- dat[samp == 0,]
10 dat1 <- dat[samp == 1,]
```

## R: Training and test error for simpler model

```
1 ## Fit simpler model
2 model <- regress("mean", dsst ~ male, data = dat0)
3
4 ## Training error
5 > model$sigma
6 [1] 12.58
7
8 ## Test error
9 MSE.test <- mean((dat1$dsst - predict(model, newdata = dat1)
10                [,1])^2, na.rm = TRUE)
11
12 > sigma.test
13 [1] 12.91
```

# TRAINING AND TEST ERROR

## R: Training and test error for more complex model

```
1 ## Fit more complex model
2 model <- regress("mean", dsst ~ male + age + race + weight +
3               height + packyrs + yrsquit + alcohol +
4               physact + chf + diabetes + genhlth +
5               ldl + alb + crt + sbp + aai + fev +
6               atrophy + numinf, data = dat0)
7
8 ## Training error
9 > model$sigma
10 [1] 10.93
11
12 ## Test error
13 MSE.test <- mean((dat1$dsst - predict(model, newdata = dat1)
14               [,1])^2, na.rm = TRUE)
15
16 > sigma.test
17 [1] 11.44
```

**Example:** Training/test error for DSST (MRI data)

- Results (RMSE).

	<b>Training</b>	<b>Test</b>
<b>Single predictor</b>	12.6	12.9
<b>Many predictors</b>	10.9	11.4

- What do you notice?

## Ideas: Model complexity in linear models

- Model complexity can be characterized by the number of degrees of freedom used by the model,  $df_M$ .
- Linear regression:

$$df_M = \frac{1}{\sigma^2} \sum_{i=1}^N \text{Cov}(\hat{Y}_i, Y_i).$$

- With  $K$  predictors plus an intercept,  $df_M = K + 1$ .
  - ▶ Absent an intercept, there are  $K$  degrees of freedom.
- Highly complex models tend to have lower bias as compared to less complex models, and will provide better prediction error on the *training* data used to fit the model.

# TABLE OF CONTENTS

- 1 Framing prediction problems
- 2 The bias-variance trade-off
- 3 Training and test error
- 4 Prediction for binary outcomes

## Binary outcomes:

- Logistic regression can be used to characterize how well a set of predictors can correctly classify an outcome.
- Model:

$$\text{logit}(P(Y = 1|\mathbf{X} = \mathbf{x})) = \mathbf{x}^T \boldsymbol{\beta}.$$

- Predicted values:

$$\hat{y}_i = \hat{Y}(\mathbf{x}_i) = \hat{P}(Y = 1|\mathbf{X} = \mathbf{x}_i) = \text{expit}(\mathbf{x}_i^T \boldsymbol{\beta}).$$

- General idea: Those with a true value of  $Y = 1$  should have higher predicted values than those with a true value of  $Y = 0$ .

## Binary outcomes:

- We can use the receiver operating characteristic (ROC) curve to characterize the classification ability.
- Consider a cutoff point,  $c$  (with  $0 \leq c \leq 1$ ), above which all with  $\hat{Y} > c$  are classified as  $Y_c = 1$  on the basis of the logistic regression model.
- Can characterize two metrics of predictive ability:
  - ▶  $P(Y_c = 0|Y = 0) = P(\hat{Y} \leq c|Y = 0)$ : proportion classified as  $Y = 0$  among those with a true value of  $Y = 0$ .
  - ▶  $P(Y_c = 1|Y = 1) = P(\hat{Y} > c|Y = 1)$ : proportion classified as  $Y = 1$  among those with a true value of  $Y = 1$ .
- These correspond to the specificity and sensitivity.
- ROC curve: plot of  $(1 - \text{specificity})$  and sensitivity across different values of  $c$ .

## Binary outcomes:

- Area under the ROC curve (AUC) summarizes the predictive ability in a single quantity.
- An AUC of 1.00 corresponds to perfect classification; an AUC of 0.50 signifies no classification ability.
- Note: AUC corresponds to a scaled variant of Wilcoxon Rank Sum test statistic, and estimates  $P(\hat{y}^1 > \hat{y}^0)$ , where, e.g.,  $\hat{y}^1$  denotes the predicted probability (of  $Y = 1$ ) for a randomly sampled individual with  $Y = 1$ .

## **Example:** Diabetes in MRI data

- $Y$ : Diabetes
- Predictors: age, sex, smoking history (pack years), physical activity, weight, height.
- Model:

$$\text{logit}(P(Y = 1|\mathbf{X} = \mathbf{x})) = \mathbf{x}^T \boldsymbol{\beta}.$$

- Naturally, we want to divide the data into a training and test set in order to characterize predictive ability.

## R: Training and test AUC

```
1 ## Fit model
2 model <- regress("odds", diabetes ~ age + male + packyrs +
3                 physact + weight + height, data = dat0)
4
5 ## TRAINING ERROR
6 ## Predicted values
7 y0hat <- predict(model, type = "response")
8
9 ## Clever way to get the AUC
10 y0hat.1 <- y0hat[dat0$diabetes == 1]
11 y0hat.0 <- y0hat[dat0$diabetes == 0]
12 n1 <- length(y0hat.1)
13 n0 <- length(y0hat.0)
14 AUC0 <- (wilcox.test(y0hat.1, y0hat.0)$statistic) / (n0 * n1)
15 > as.numeric(AUC0)
16 [1] 0.6881
17
18 ## TEST ERROR (Totally analogous)
19 y1hat <- predict(model, newdata = dat1)
20 y1hat.1 <- y1hat[dat1$diabetes == 1]
21 y1hat.0 <- y1hat[dat1$diabetes == 0]
22 n1 <- length(y1hat.1)
23 n0 <- length(y1hat.0)
24 AUC1 <- (wilcox.test(y1hat.1, y1hat.0)$statistic) / (n0 * n1)
25 > as.numeric(AUC1)
26 [1] 0.6136
```

## **Example:** Diabetes in MRI data

- Unsurprisingly, training AUC higher relative to the test AUC.
- There are many other metrics of predictive ability that can be applied to binary outcomes.

## **This unit:**

- The bias-variance trade-off.
- Training error does not estimate test error.
- Measures of predictive ability (continuous/binary)
- A lot let unsaid.
  - ▶ Forward- and backward-selection.
  - ▶ Model optimism.
  - ▶ Penalized regression.
  - ▶ Regression trees.
  - ▶ Support vector machines.
  - ▶ Neural nets.
- This unit is just to orient you to the basics.

## So far:

- Review.
- Simple linear regression.
- Multiple linear regression (foundations).
- Multiple linear regression (interactions and strata).
- Transformations and basis expansions.
- Regression with binary outcomes.
- Regression with nominal, ordinal, and count outcomes.
- Introduction to clustered data.
- Methods for time-to-event outcomes.
- Predictive capacity of regression models.

## The end!

- Comprehensive examinations for biostatistics students.
- Summer research.
- For those taking advanced regression: Review your linear algebra!
  - ▶ Eigenvalues and eigenvectors.
  - ▶ Projection matrices.
  - ▶ Positive definite matrices.
  - ▶ Singular value decomposition.
- Don't stop learning!