

LECTURE 3, PART I: ERROR AND VALIDATION

Text references: Chapter 3 in Shalizi

Review/Introductory Remark

Suppose X represents different covariates and Y is our response variable, and suppose that we are interested in the relationships between X and Y . Here are three *different* but related concepts:

1. association

2. causation

3. prediction

Association does not imply causation. Do you remember when you would be able to make causal statements?

Furthermore, a strong association between X and Y does not necessarily mean that X is an important predictor in the model. Why? Give an example.

In Lecture 2, Parts III and IV, we discussed causal inference, and the pitfalls of interpreting relationships or the output of a multiple regression program as casual. In this lecture, we are going to take a purely *predictive viewpoint* of regression. We start by describing the different sources of errors in prediction and how to estimate error in practice.

The Predictive Viewpoint. Set-Up.

Given data $(X_1, Y_1), \dots, (X_n, Y_n)$ we have two goals:

estimation: Find an estimate $\hat{r}(x)$ of the regression function $r(x)$.

prediction: Given a *new* X , predict Y ; we use $\hat{Y} = \hat{r}(X)$ as the prediction.

Note the following (what is random? what is fixed?):

1 _____

2 _____

3 _____

4 _____

5 _____

6 _____

7 _____

8 _____

9 _____

10 **Q:** How well does the the function \hat{r} , constructed on the training set/sample,
11 predict our new test point (X, Y) ?

12 To measure this, we define the **expected test error** (which is often referred
13 to as **prediction error** or **prediction risk**):

14 _____

15 _____

16 _____

17 _____

18 where the expectation is over all that is random (training set, and test
19 point).

Q: Why would we want to measure this? Here are two major reasons:

1. *Model assessment*: sometimes we would just like to know how well we can predict a future observation, in absolute terms
2. *Model selection*: often we will want to choose between different fitting functions; this could mean choosing between two different model classes entirely (e.g., linear regression versus some other fixed method) or choosing an underlying tuning parameter for a single method (e.g., choosing k in k -nearest neighbors). How to do this? A common way is to compare prediction errors (which are in practice *estimated* from data).

Note that there are two estimation problems:

The Bias-Variance Decomposition

First some intuition

The expected test error or prediction risk R above has an important property: it decomposes into informative parts. But, before we go through the mathematics of the decomposition, let's think about a few points:

1 1. Can we ever predict Y from X with zero prediction error? Likely, no.

2 Even if our function \hat{r} happened to capture the ideal relation underlying
3 X, Y , i.e., the true regression function $r(X) = \mathbb{E}(Y|X)$, we would
4 still incur error, due to noise. We call this *irreducible error*; it can represent
5

6
7
8
9
10 2. What happens if our fitted function \hat{r} belongs to a model class that is
11 far from the true regression function r ? E.g., we may choose to fit a
12 linear model in a setting where the true relationship is far from linear?
13 We'll often refer to this as a *misspecified model*. As a result, we encounter
14 error, what we call *estimation bias*; it represents
15
16
17
18
19
20
21

22 3. What happens if our fitted (random) function \hat{r} is itself quite variable?
23 In other words, for different training samples of size n , we end up constructing
24 substantially different functions \hat{r} ? This is another source of

error, that we'll call *estimation variance*; it represents

Point Estimation Revisited. The Bias-Variance Decomposition of MSE

Before we return to regression, first recall the problem of **point estimation**, i.e. the problem of providing a single “best guess” of some fixed quantity of interest. More formally let Y_1, \dots, Y_n be n IID data points from some distribution. A point estimator $\hat{\theta}_n$ of a parameter θ is some function of Y_1, \dots, Y_n (what is fixed? what is random?):

The bias of a point estimator is defined by:

1 The distribution of $\hat{\theta}_n$ is called the **sampling distribution**. Think about this
2 figure:

3 The quality of a point estimate is sometimes assessed by the **mean squared**
4 **error**, or MSE, defined by:

5 _____
6 _____

7 The MSE can be written as:

8 _____
9 _____

10 Now compare **point estimation** with regression and **curve estimation**.

11 Think about this figure and compare with the one above!

1 The Bias-Variance Decomposition of the MSE in Regression

2 If X and Y are random variables, recall the rule of iterated expectations

$$3 \quad \mathbb{E}[g(X, Y)] = \mathbb{E}[\mathbb{E}[g(x, Y)|X = x]],$$

4 where the inner expectation is taken with respect to $Y|X$ and the outer one
5 is taken with respect to the marginal distribution of X . Throughout the
6 following section, we use this rule, **conditioning on X** to obtain the risk
7 function $R(x)$ at $X = x$.

8 Hence, for our expected test error or prediction error, we have

13 Now look at $R(x)$, the *prediction error conditional on $X = x$* for some fixed
14 value of x :

1 _____
2 _____
3 _____
4 _____
5 _____
6 _____
7 The first term is just the *irreducible error*. The second term can be further
8 decomposed just as the MSE of a point estimator (NOTE: $\hat{r}(x)$ is a *random*
9 *variable* itself, calculated at a fixed value of x) into a bias component and a
10 variance component at x . Therefore, altogether,

11 _____
12 _____
13 which is called the **bias-variance decomposition**.

14 Typical trend: *underfitting* means high bias and low variance, *overfitting*
15 means low bias but high variance. E.g., think about k in k -nearest-neighbors
16 regression: relatively speaking, how do the bias and variance behave for
17 small k , and for large k ? [See R Demo 1 and HW Set 1, Problem 1]

1
2
3
4
5 Finally, we have for the prediction error R :

6
7
8 and R_{av} is called the **average prediction risk**.

9 To summarize: We wish to know R , the prediction risk. R_{av} provides an
10 excellent approximation, but R_{av} is not a quantity that we can readily cal-
11 culate empirically because we do not know $R(X_i)$. Let us next explore why
12 it is challenging to calculate R .

13 **The Optimism of the Training Error**

14 Define the **training error**

15
16
17 What wrong with the training error or expected training error? (As before,
18 we will condition on x_1, \dots, x_n in our discussion. If needed, one can always
19 derive unconditional expectations by the “law of total expectations”, and
20 unconditional variances by the “law of total variance”.)

We might guess that $\hat{R}_{\text{training}}$ estimates the prediction error (R) well but this is not true. The reason is that we used the observed pairs (x_i, Y_i) to obtain $\hat{Y}_i = \hat{r}(x_i)$. As a consequence Y_i and \hat{Y}_i are correlated. Typically \hat{Y}_i “predicts” Y_i better than it predicts a new Y at the same x_i . Let us explore this formally. Let $\bar{r}_i = \mathbb{E}(\hat{r}(x_i))$ and compute

$$\begin{aligned}\mathbb{E}(Y_i - \hat{Y}_i)^2 &= \mathbb{E}(Y_i - r(x_i) + r(x_i) - \bar{r}(x_i) + \bar{r}(x_i) - \hat{Y}_i)^2 \\ &= \sigma^2 + \mathbb{E}(r(x_i) - \bar{r}(x_i))^2 + \mathbb{V}(\hat{r}(x_i)) - 2\text{Cov}(\hat{Y}_i, Y_i).\end{aligned}$$

Note: this time the cross-product involving the 1st and 3rd terms is not 0 because $\text{Cov}(\hat{Y}_i, Y_i) \neq 0$. This is because Y_i is a particular observation from which we calculated \hat{Y}_i , hence the two terms are correlated. This introduces a bias into the estimate of risk:

Typically, $\text{Cov}(\hat{Y}_i, Y_i) > 0$ and so $\hat{R}_{\text{training}}$ underestimates the risk. Later, we shall see how to estimate the prediction risk.

Lecture 3, Part I: The Basics of Error and Validation

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20