

# STATS305A - Lecture 12

John Duchi

Scribed by Michael Howes

10/28/21

## Contents

|          |                                       |          |
|----------|---------------------------------------|----------|
| <b>1</b> | <b>Announcements</b>                  | <b>1</b> |
| <b>2</b> | <b>Model Selection and prediction</b> | <b>1</b> |
| 2.1      | Motivation                            | 1        |
| 2.2      | In sample and out of sample risk      | 1        |
| 2.3      | Bias/variance decomposition           | 3        |
| 2.4      | Comparing models                      | 4        |
| 2.5      | Estimating risk                       | 4        |

## 1 Announcements

- Etude 2 due today 5pm.
- No class next Tuesday.

## 2 Model Selection and prediction

### 2.1 Motivation

Up to this point we've treated the model  $Y = X\beta + \varepsilon$  as "god-given". This is a bit inaccurate. In real life we will typically have data and no model and have to figure it out and select a model. When selecting a model we have two desiderata:

- Identify important features that are relating  $x$  to our response  $y$ .
- Pure predictive accuracy: how well can we predict  $y$  from  $x$ ?

These two are intertwined. We don't always have to choose one over the other.

### 2.2 In sample and out of sample risk

Suppose we are in a setting where  $y = f(x) + \varepsilon$  and  $\mathbb{E}[\varepsilon|x] = 0$ . This is equivalent to having  $f(x) = \mathbb{E}[Y|X = x]$  since if  $\varepsilon = y - f(x)$ , then

$$\mathbb{E}[\varepsilon|x] = \mathbb{E}[y|x] - f(x).$$

Thus  $\mathbb{E}[\varepsilon|x] = 0$  if and only if  $f(x) = \mathbb{E}[y|x]$ . Define  $\sigma^2(x) = \mathbb{E}[\varepsilon^2|x]$  which is the conditional variance of  $\varepsilon$ .

Our goal is to fit a predictor  $\hat{f}$  using a sample  $\{(x_i, y_i)\}_{i=1}^n$ . Note that if we think of the sample of  $\{(x_i, y_i)\}_{i=1}^n$  as random, then the predictor  $\hat{f}$  is random (like how  $\hat{\beta}$  is random in the linear model). Thus we can take the expectation of quantities involving  $\hat{f}$  over all samples  $\{(x_i, y_i)\}_{i=1}^n$ . This idea will be used many times over the course of this lecture.

**Definition 1.** If we have a predictor  $\hat{f}$  of a model  $y = f(x) + \varepsilon$ , then we define the *in-sample (MSE) risk* of  $\hat{f}$  to be

$$R_{in}(\hat{f}) = \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f(x_i))^2 \right],$$

where the above expectation is taken over all samples  $\{(x_i, y_i)\}_{i=1}^n$  with  $x_i$  fixed. (That is we fix  $x$  and calculate  $\hat{f}$  using different samples  $(x, y)$ , we then calculate the quantity  $\frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f(x_i))^2$  and take the expectation over all samples  $(x, y)$ .)

**Aside 1.** Sometimes the in-sample risk is called the  $L^2(P_n)$  risk. This is because  $R_{in}$  is the expectation of the squared  $L^2$  norm error of  $\hat{f} - f$  with respect to the distribution

$$P_n = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i}.$$

**Definition 2.** Sometimes the insample risk is defined with respect to a fresh sample  $\{Y_i^*\}_{i=1}^n$  where

$$Y_i^* = \text{a new sample of } Y_i = f(x_i) + \varepsilon_i^*,$$

where  $\varepsilon_i^*$  is an independent copy of  $\varepsilon_i$ . We then define

$$R_{in}^*(\hat{f}) = \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (Y_i^* - \hat{f}(x_i))^2 \right],$$

where here the expectation is over both  $Y_1, \dots, Y_n$  (used to calculate  $\hat{f}$ ) and over  $Y_1^*, \dots, Y_n^*$  (used to calculate  $(Y_i^* - \hat{f}(x_i))^2$ ).

Note that

$$\begin{aligned} R_{in}^*(\hat{f}) &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (Y_i^* - \hat{f}(x_i))^2 \right] \\ &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (Y_i^* - f(x_i))^2 \right] + \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f(x_i))^2 \right] \\ &= \frac{1}{n} \sum_{i=1}^n \sigma^2(x_i) + R_{in}(\hat{f}). \end{aligned}$$

We call  $\frac{1}{n} \sum_{i=1}^n \sigma^2(x_i)$  the irreducible error.

Now suppose that we have a function  $g : \mathcal{X} \rightarrow \mathbb{R}$  where  $\mathcal{X}$  is the space  $X$  lives in. Note that  $g$  is different to  $\hat{f}$ . The predictor  $\hat{f}$  is something that depends on the sample  $(x, y)$  used to fit  $\hat{f}$ . The function  $g$  is simply a function. It is a way of taking an  $X$  and producing a number. With this in mind we have our next definition.

**Definition 3.** Given a function  $g : \mathcal{X} \rightarrow \mathbb{R}$ , the *out of sample (MSE) risk* of  $g$  is

$$R_{out}(g) = \mathbb{E}[(Y - g(X))^2] = \int_{\mathcal{X}} \mathbb{E}[(Y - g(x))^2 | X = x] p(x) dx.$$

Here the expectation is over both  $Y$  and  $X$  (hence out of sample - we are allowing  $X$  to change).

Note that

$$\begin{aligned} R_{out}(g) &= \mathbb{E}[(Y - f(X) + f(X) + g(X))^2] \\ &= \mathbb{E}[(Y - f(X))^2] + \mathbb{E}[(f(X) - g(X))^2] + 2\mathbb{E}[(Y - f(X))(f(X) - g(X))] \\ &= \mathbb{E}[\sigma^2(X)] + \mathbb{E}[(f(X) - g(X))^2]. \end{aligned}$$

We again call  $\mathbb{E}[\sigma^2(X)]$  the irreducible error and we could call  $\mathbb{E}[(f(X) - g(X))^2]$  the “error in mean prediction” (this last term is just a term John used - he said that there isn’t really a term in literature for  $\mathbb{E}[(f(X) - g(X))^2]$ ).

In the out of sample risk we average over all the  $X$ ’s we could possibly draw. In the in sample we fix the value  $x_i$  and average over all possible  $y_i$ . Note that if our data is i.i.d., then

$$R_{out}(g) = \mathbb{E}[(g(X_{n+1}) - Y_{n+1})^2].$$

The quantity  $R_{out}(\hat{f})$  is a number but it is random since it depends on the sample  $(x_i, y_i)_{i=1}^n$  used to fit  $\hat{f}$  and the sample  $(x_i, y_i)_{i=1}^n$  is random. Thus  $\mathbb{E}[R_{out}(\hat{f})]$  is our expected out of sample risk over all samples used to fit  $\hat{f}$ .

### 2.3 Bias/variance decomposition

Let  $\hat{f}$  be any estimator of  $f$ . For any  $x \in \mathcal{X}$  we can write

$$\begin{aligned} \mathbb{E}[(\hat{f}(x) - f(x))^2] &= \mathbb{E}\left[\left(\hat{f}(x) - \mathbb{E}[\hat{f}(x)] + \mathbb{E}[\hat{f}(x)] - f(x)\right)^2\right] \\ &= \mathbb{E}\left[\left(\hat{f}(x) - \mathbb{E}[\hat{f}(x)]\right)^2\right] + \left(\mathbb{E}[\hat{f}(x)] - f(x)\right)^2 \\ &= \text{Var}(\hat{f}(x)) + \text{Bias}(\hat{f}(x))^2, \end{aligned}$$

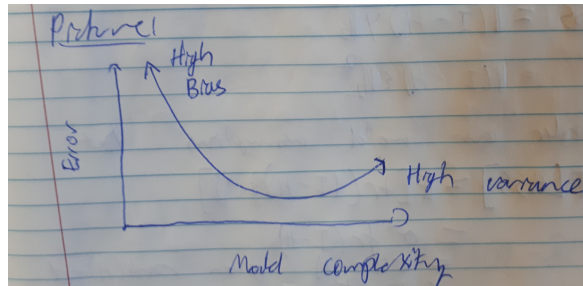
where, as in  $R_{in}(\hat{f})$ , the above expectation is taken with respect to the samples used to fit  $\hat{f}$ . Thus we have

$$R_{in}(\hat{f}) = \frac{1}{n} \sum_{i=1}^n \left( \text{Bias}(\hat{f}(x_i))^2 + \text{Var}(\hat{f}(x_i)) \right),$$

and

$$\begin{aligned} \mathbb{E}[R_{out}(\hat{f})] &= \mathbb{E}[\sigma^2(X)] + \mathbb{E}[\text{Bias}(\hat{f}(X))^2] + \mathbb{E}[\text{Var}(\hat{f}(X))] \\ &= \text{irreducible error} + \text{expected bias squared} + \text{expected variance}. \end{aligned}$$

This decomposition means that in practice we often see the following curve when we plot error against model complexity. The error initially goes down as complexity increases but then increases. This is because initially the model has high bias when the model is simple but once the model is very complicated it has high variance.



Although one word of warning: We sometimes don't know how to properly measure model complexity. This means that when we use a proxy for model complexity we might not see the above picture. This is sometimes the case with machine learning algorithms where the error continues to go down even as the number of parameters is getting very very high. The complexity of the model may not be as high as the number of parameters suggests.

## 2.4 Comparing models

The above analysis gives us the idea that we should choose a model that trades optimally between bias and variance. Unfortunately we can't do this exactly since we can only approximate the bias and variance of a model. We thus will have to estimate the prediction error. We can then use this estimate to choose a model. There is a challenge though, the natural quantity

$$\frac{1}{n}RSS = \frac{1}{n} \sum_{i=1}^n (Y_i \hat{Y}_i)^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{f}(x_i))^2,$$

is biased downwards since we chose  $\hat{f}$  to minimize RSS. Thus  $\frac{1}{n}RSS$  will tend to be smaller than  $R_{in}^*$  since in  $R_{in}^*$  we have a new independent sample of  $Y$ 's.

**Example 1.** Consider for concreteness the linear model  $Y = X\beta + \varepsilon$  where  $\varepsilon \sim (0, \sigma^2 I_n)$  and  $X$  is full rank with  $\text{rank}(X) = d$ . Let  $H = X(X^T X)^{-1} X^T$ , we know that

$$RSS = \|(I - H)Y\|_2^2 = \|(I - H)\varepsilon\|_2^2.$$

Thus we have

$$\begin{aligned} \mathbb{E}[RSS] &= \mathbb{E}[\text{tr}((I - H)\varepsilon\varepsilon^T)] \\ &= \sigma^2 \text{tr}(I - H) \\ &= \sigma^2(n - d) \\ &< \sigma^2 n. \end{aligned}$$

Thus  $\mathbb{E}[\frac{1}{n}RSS] = \frac{n-d}{n}\sigma^2 < \sigma^2$  which is a problem since the irreducible error is  $\sigma^2$  so we will always have  $R_{in}^* \geq \sigma^2$ .

There are two approaches to overcome this issue:

- (a) Penalized risk estimates.
- (b) Validation.

We will discuss the first approach today.

## 2.5 Estimating risk

How do we estimate the risk of a model? Let's think about the sum of square errors. Note that

$$\begin{aligned} \mathbb{E}[(Y_i - \hat{Y}_i)^2] &= \mathbb{E}[(Y_i - f(x_i) + f(x_i) - \hat{Y}_i)^2] \\ &= \sigma_i^2 + 2\mathbb{E}[(Y_i - f(x_i))(f(x_i) - \hat{Y}_i)] + \mathbb{E}[(f(x_i) - \hat{f}(x_i))^2] \\ &= \sigma_i^2 + \mathbb{E}[(f(x_i) - \hat{f}(x_i))^2] - 2\text{Cov}(Y_i, \hat{Y}_i) \end{aligned}$$

The last equality holds because  $f(x_i) = \mathbb{E}[Y_i]$  and so for any constant  $c_i$

$$\mathbb{E} \left[ (Y_i - f(x_i))(\hat{Y}_i - c_i) \right] = \mathbb{E} \left[ (Y_i - f(x_i))(\hat{Y}_i - \mathbb{E}[\hat{Y}_i]) \right] = \text{Cov}(Y_i, \hat{Y}_i).$$

Thus we have

$$\begin{aligned} \mathbb{E}[RSS] &= \mathbb{E} \left[ \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \right] \\ &= \sum_{i=1}^n \sigma_i^2 + \sum_{i=1}^n \mathbb{E} \left[ (f(x_i) - \hat{f}(x_i))^2 \right] - 2 \text{Cov}(Y, \hat{Y}) \\ &= \sum_{i=1}^n \sigma_i^2 + nR_{in}(\hat{f}) - 2 \text{Cov}(Y, \hat{Y}). \end{aligned}$$

Equivalently,

$$\begin{aligned} R_{in}^*(\hat{f}) &= R_{in}(\hat{f}) + \frac{1}{n} \sum_{i=1}^n \sigma_i^2 \\ &= \frac{1}{n} \mathbb{E}[RSS] + \frac{2}{n} \text{Cov}(Y, \hat{Y}). \end{aligned}$$

The quantity  $\text{Cov}(Y, \hat{Y})$  is called the “effective degrees of freedom”. The upshot of these calculations is that if we can estimate  $\frac{2}{n} \text{Cov}(Y, \hat{Y})$ , then we can estimate the in-sample risk of  $\hat{f}$ .

**Example 2.** Let’s calculate the  $\text{Cov}(Y, \hat{Y})$  in some special cases. Suppose we have a “true” linear model  $Y = X\beta + \varepsilon$  with  $\varepsilon \sim (0, \sigma^2 I_n)$  and we fit  $\hat{Y} = HY$ . Then

$$\begin{aligned} \text{Cov}(Y, \hat{Y}) &= \text{Cov}(\varepsilon, H\varepsilon) \\ &= \mathbb{E}[\varepsilon^T H \varepsilon] \\ &= \sigma^2 \text{tr}(H) \\ &= \sigma^2 d, \end{aligned}$$

if  $H$  has rank  $d$ . Suppose now that we don’t know that  $Y$  comes from a linear model but we have the more general model  $Y_i = f(x_i) + \varepsilon_i$ . Suppose we still fit  $\hat{f}(x) = x^T \hat{\beta}$  where  $\hat{\beta} = (X^T X)^{-1} X^T Y$ . We then still have  $\hat{Y} = HY$  and so

$$\begin{aligned} \text{Cov}(Y, \hat{Y}) &= \text{Cov}(\varepsilon, HY) \\ &= \text{Cov}(\varepsilon, H[f(x_i) + \varepsilon_i]_{i=1}^n) \\ &= \text{Cov}(\varepsilon, H\varepsilon) \\ &= \sum_{i=1}^n \sigma_i^2 H_{ii}, \end{aligned}$$

where  $\sigma_i^2 = \text{Var}(Y_i) = \text{Var}(\varepsilon_i)$ . If we have homoscedastic noise, then  $\sigma_i^2 = \sigma$  for all  $i$  and

$$\text{Cov}(Y, \hat{Y}) = \sigma^2 d,$$

where again  $d$  is the rank of  $H$ . These two examples give motivation for calling  $\text{Cov}(Y, \hat{Y})$  the effective degrees of freedom.

**Definition 4.** *Mallow's  $C_p$  statistic* is the statistic.

$$C_p := \frac{1}{n}RSS + 2\frac{\sigma^2}{n}\text{rank}(X) = \frac{1}{n}RSS + 2\frac{\sigma^2}{n}p,$$

where  $p = \text{rank}(X)$ .

**Definition 5.** *Mallow's method* is the following.

- Given a number of possible models, calculate  $C_p$  for each model.
- Choose the model that minimizes  $C_p$

Unfortunately we do not know  $\sigma^2$  and so to use Mallow's method we have to estimate  $\sigma^2$ . Some options are

- (a) If we have a pair  $(Y_i^{(1)}, Y_i^{(2)})$  for each  $x_i$ , then  $\mathbb{E} \left[ \left( Y_i^{(1)} - Y_i^{(2)} \right)^2 \right] = 2\sigma$  and so we can use

$$\hat{\sigma}^2 = \frac{1}{2n} \sum_{i=1}^n \left( Y_i^{(1)} - Y_i^{(2)} \right)^2.$$

- (b) If we do not have paired data, then we can find  $x$ 's that are close to each other. We can then pair these  $x$ 's and use the above estimator.
- (c) We can also use the most complicated model to estimate  $\sigma^2$ . Suppose that we have  $Z \in \mathbb{R}^{n \times q}$  where  $q$  is very big. Then we can define

$$\hat{\gamma} = (Z^T Z)^{-1} Z^T Y \in \mathbb{R}^q,$$

and use the estimator

$$\hat{\sigma}^2 = \frac{1}{n-q} \sum_{i=1}^n (Y_i - Z_i^T \hat{\gamma})^2.$$