

## 26 Julio - Statistical Learning

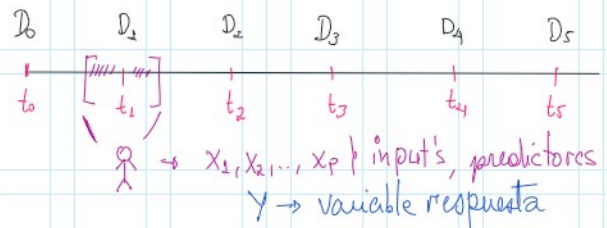
lunes, 26 de julio de 2021 10:50

**Statistical learning** refers to a vast set of tools for **understanding data**. These tools can be classified as **supervised** or **unsupervised**. Broadly speaking, supervised statistical learning involves **building a statistical model for predicting**, or estimating, an **output** based on one or more **inputs**. Problems of this nature occur in fields as diverse as business, medicine, astrophysics, and

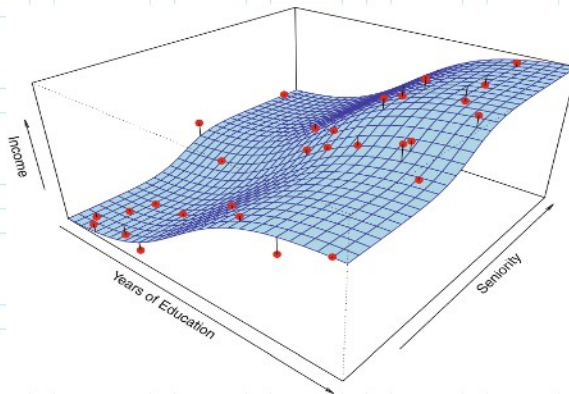
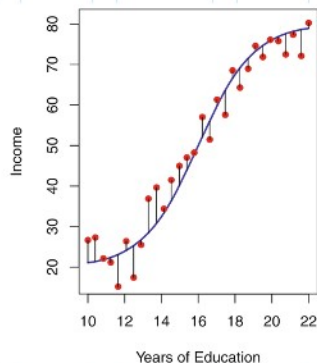
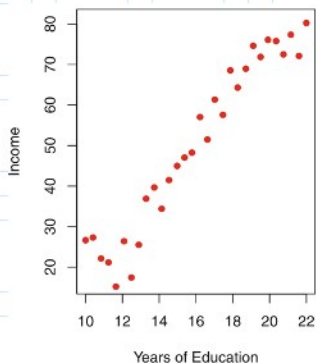
More generally, suppose that we observe a quantitative response  $Y$  and  $p$  different predictors,  $X_1, X_2, \dots, X_p$ . We assume that there is some relationship between  $Y$  and  $X = (X_1, X_2, \dots, X_p)$ , which can be written in the very general form

$$Y = f(X) + \epsilon \quad (2.1)$$

$f(X_1, X_2, \dots, X_p)$        $\epsilon \rightarrow$  aleatorio



Here  $f$  is some fixed but unknown function of  $X_1, \dots, X_p$ , and  $\epsilon$  is a **random error term**, which is independent of  $X$  and has mean zero. In this formulation,  $f$  represents the **systematic information** that  $X$  provides about  $Y$ .



**ESENCIA:** Aprendizaje estadístico es compuesto por un conjunto de propuestas o enfoques (*approaches*) para estimar  $f$

**Por que estimar  $f$ :** Basicamente para predicción e inferencia

### Prediction

In many situations, a set of inputs  $X$  are readily available, but the output  $Y$  cannot be easily obtained. In this setting, since the error term averages to zero, we can predict  $Y$  using

$$\hat{Y} = \hat{f}(X), \quad (2.2)$$

where  $\hat{f}$  represents our estimate for  $f$ , and  $\hat{Y}$  represents the resulting prediction for  $Y$ . In this setting,  $\hat{f}$  is often treated as a *black box*, in the sense that one is not typically concerned with the exact form of  $\hat{f}$ , provided that it yields accurate predictions for  $Y$ .

Consider a given estimate  $\hat{f}$  and a set of predictors  $X$ , which yields the prediction  $\hat{Y} = \hat{f}(X)$ . Assume for a moment that both  $\hat{f}$  and  $X$  are fixed. Then, it is easy to show that

$$\begin{aligned} E(Y - \hat{Y})^2 &= E[f(X) + \epsilon - \hat{f}(X)]^2 \\ &= \underbrace{[f(X) - \hat{f}(X)]^2}_{\text{Reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible}}, \end{aligned} \quad (2.3)$$

where  $E(Y - \hat{Y})^2$  represents the average, or *expected value*, of the squared difference between the predicted and actual value of  $Y$ , and  $\text{Var}(\epsilon)$  represents the *variance* associated with the error term  $\epsilon$ .

### Inference

We are often interested in **understanding the way that  $Y$  is affected as  $X_1, \dots, X_p$  change**. In this situation we wish to estimate  $f$ , but our goal is **not necessarily to make predictions for  $Y$** . We instead want to understand the relationship between  $X$  and  $Y$ , or more specifically, to understand how  $Y$  changes as a function of  $X_1, \dots, X_p$ . Now  $\hat{f}$  cannot be treated as a black box, because we need to know its exact form. In this setting, one may be interested in answering the following questions:

- **Which predictors are associated with the response?** It is often the case that only a small fraction of the available predictors are substantially

$$Y = f(X) + \epsilon$$

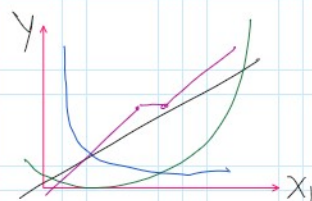
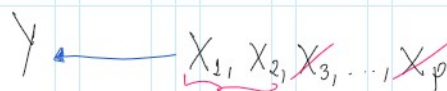
$$\hat{Y} = \hat{f}(X); \hat{f} \rightarrow \text{estimar}$$

$$\hat{f}(6.6) = \hat{Y}_{6.6} \rightarrow \text{predicción}$$



interested in answering the following questions:

- Which predictors are associated with the response? It is often the case that only a small fraction of the available predictors are substantially associated with  $Y$ . Identifying the few important predictors among a large set of possible variables can be extremely useful, depending on the application.
- What is the relationship between the response and each predictor? Some predictors may have a positive relationship with  $Y$ , in the sense that increasing the predictor is associated with increasing values of  $Y$ . Other predictors may have the opposite relationship. Depending on the complexity of  $f$ , the relationship between the response and a given predictor may also depend on the values of the other predictors.
- Can the relationship between  $Y$  and each predictor be adequately summarized using a linear equation, or is the relationship more complicated? Historically, most methods for estimating  $f$  have taken a linear form. In some situations, such an assumption is reasonable or even desirable. But often the true relationship is more complicated, in which case a linear model may not provide an accurate representation of the relationship between the input and output variables.



$f$  es desconocido  
 $\hookrightarrow \hat{f}$  estimado  
 Datos de entrenamiento

## How Do We Estimate $f$ ?

These observations are called the **training data** because we will use these observations to train, or teach, our method how to estimate  $f$ . Let  $x_{ij}$  represent the value of the  $j$ th predictor, or input, for observation  $i$ , where  $i = 1, 2, \dots, n$  and  $j = 1, 2, \dots, p$ . Correspondingly, let  $y_i$  represent the response variable for the  $i$ th observation. Then our training data consist of  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$  where  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$ .

### Parametric Methods

Parametric methods involve a two-step model-based approach.

1. First, we make an **assumption about the functional form**, or shape, of  $f$ . For example, one very simple assumption is that  $f$  is linear in  $X$ :

$$f(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p. \quad (2.4)$$

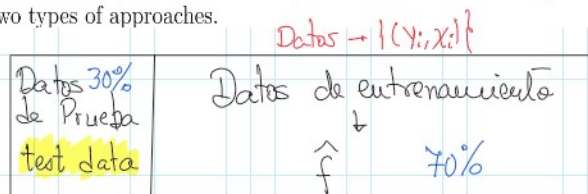
2. After a model has been selected, we need a procedure that uses the training data to **fit or train** the model. In the case of the linear model (2.4), we need to estimate the **parameters**  $\beta_0, \beta_1, \dots, \beta_p$ . That is, we want to find values of these parameters such that

$$Y \approx \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p.$$

### Non-parametric Methods

Non-parametric methods do not make explicit assumptions about the functional form of  $f$ . Instead they seek an estimate of  $f$  that gets as close to the data points as possible without being too rough or wiggly. Such approaches can have a major advantage over parametric approaches: by avoiding the assumption of a particular functional form for  $f$ , they have the potential to accurately fit a wider range of possible shapes for  $f$ . Any parametric

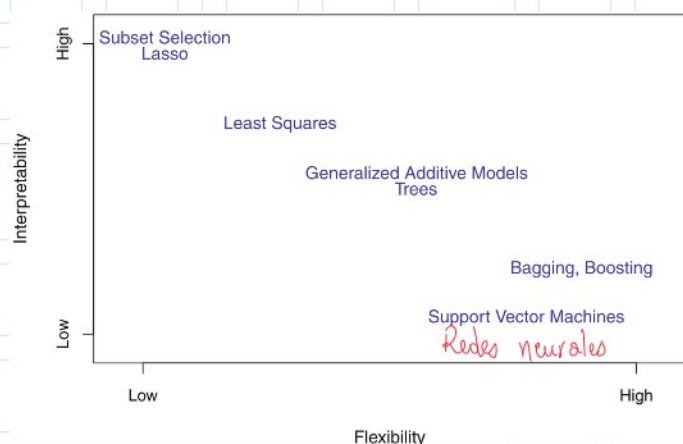
Our goal is to apply a statistical learning method to the training data in order to estimate the unknown function  $f$ . In other words, we want to find a function  $\hat{f}$  such that  $Y \approx \hat{f}(X)$  for any observation  $(X, Y)$ . Broadly speaking, most statistical learning methods for this task can be characterized as either *parametric* or *non-parametric*. We now briefly discuss these two types of approaches.



Ventaja: Es flexible en relación al enfoque paramétrico  
 $\hookrightarrow$  No asume una forma para " $f$ ".

Desventaja: El número de parámetros es bastante grande

## The Trade-Off Between Prediction Accuracy and Model Interpretability



Low

High

Flexibility

## Measuring the Quality of Fit $\rightarrow \hat{f}(x) = \hat{y}$

In order to evaluate the performance of a statistical learning method on a given data set, we need some way to measure how well its predictions actually match the observed data. That is, we need to quantify the extent to which the predicted response value for a given observation is close to the true response value for that observation. In the regression setting, the most commonly-used measure is the **mean squared error (MSE)**, given by

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2, \quad (2.5)$$

where  $\hat{f}(x_i)$  is the prediction that  $\hat{f}$  gives for the  $i$ th observation. The MSE will be small if the predicted responses are very close to the true responses, and will be large if for some of the observations, the predicted and true responses differ substantially.

To state it more mathematically, suppose that we fit our statistical learning method on our training observations  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ , and we obtain the estimate  $\hat{f}$ . We can then compute  $\hat{f}(x_1), \hat{f}(x_2), \dots, \hat{f}(x_n)$ . If these are approximately equal to  $y_1, y_2, \dots, y_n$ , then the **training MSE** given by (2.5) is small. However, we are really not interested in whether  $\hat{f}(x_i) \approx y_i$ ; instead, we want to know whether  $\hat{f}(x_0)$  is approximately equal to  $y_0$ , where  $(x_0, y_0)$  is a *previously unseen test observation not used to train the statistical learning method*. We want to choose the method that gives the lowest **test MSE**, as opposed to the lowest training MSE. In other words, if we had a large number of test observations, we could compute

$$\text{Ave}(y_0 - \hat{f}(x_0))^2, \quad (2.6)$$

the average squared prediction error for these test observations  $(x_0, y_0)$ . We'd like to select the model for which the average of this quantity—the test MSE—is as small as possible.

## Linear Regression

*Simple linear regression* lives up to its name: it is a very straightforward approach for predicting a quantitative response  $Y$  on the basis of a single predictor variable  $X$ . It assumes that there is approximately a linear relationship between  $X$  and  $Y$ . Mathematically, we can write this linear relationship as

$$Y \approx \beta_0 + \beta_1 X. \quad (3.1)$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x,$$

where  $\hat{y}$  indicates a prediction of  $Y$  on the basis of  $X = x$ . Here we use a *hat* symbol,  $\hat{\cdot}$ , to denote the estimated value for an unknown parameter or coefficient, or to denote the predicted value of the response.

## Estimating the Coefficients

In practice,  $\beta_0$  and  $\beta_1$  are unknown. So before we can use (3.1) to make predictions, we must use data to estimate the coefficients. Let

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

represent  $n$  observation pairs, each of which consists of a measurement of  $X$  and a measurement of  $Y$ . In the **Advertising** example, this data

Let  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  be the prediction for  $Y$  based on the  $i$ th value of  $X$ . Then  $e_i = y_i - \hat{y}_i$  represents the  $i$ th *residual*—this is the difference between the  $i$ th observed response value and the  $i$ th response value that is predicted by our linear model. We define the *residual sum of squares (RSS)* as

$$RSS = e_1^2 + e_2^2 + \dots + e_n^2,$$

The MSE in (2.5) is computed using the training data that was used to fit the model, and so should more accurately be referred to as the **training MSE**. But in general, we do not really care how well the method works on the training data. Rather, *we are interested in the accuracy of the predictions that we obtain when we apply our method to previously unseen test data.*



$$\text{RSS} = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + (y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^2 + \dots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2. \quad (3.3)$$

The least squares approach chooses  $\hat{\beta}_0$  and  $\hat{\beta}_1$  to minimize the RSS. Using some calculus, one can show that the minimizers are

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \\ \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x}, \end{aligned} \quad (3.4)$$

where  $\bar{y} \equiv \frac{1}{n} \sum_{i=1}^n y_i$  and  $\bar{x} \equiv \frac{1}{n} \sum_{i=1}^n x_i$  are the sample means. In other words, (3.4) defines the *least squares coefficient estimates* for simple linear regression.