

Introduction to Probability and Statistics  
for Engineers and Scientists  
*Sheldon M Ross*  
Notes and Exercises

Anirudh Krishnan

September 13, 2021

# Contents

## 1 Regression

1.1	Introduction . . . . .	
1.2	Least squares estimators of regression parameters . . . . .	
1.3	Distribution of the estimators . . . . .	
1.4	Statistical Inferences about regression parameters . . . . .	
1.5	Coefficient of Determination and Sample correlation coefficient . . . . .	
1.6	Analysis of residuals:assessing the model . . . . .	
1.7	Transforming to linearity . . . . .	
1.8	Weighted Least Squares . . . . .	
1.9	Polynomial Regression . . . . .	
1.10	Multiple Linear Regression . . . . .	

# Chapter 1

## Regression

*“Are you sure you’ve gotten rid of any multicollinearity in the inputs?”*

### 1.1 Introduction

The problem of determining the relationship between a set of inputs  $\{x_i\}$  and the resulting output  $Y$  is a frequent problem in engineering. This is further complicated by the lack of prior knowledge about the nature of this dependence.

A simplistic model of response depending on a set of inputs uses a linear combination of these inputs along with some noise  $e$ , which is assumed to have zero mean.

$$Y = \beta_0 + \beta_1 x_1 + \cdots + \beta_r x_r + e \quad (1.1)$$

$$\mathbb{E}[Y|\mathbf{x}] = \beta_0 + \beta_1 x_1 + \cdots + \beta_r x_r \quad (1.2)$$

The set of inputs  $\{x_i\}$  are called *independent variables* and the response  $Y$  which is some function of the inputs is called the *dependent variable*.

The set of coefficients  $\{\beta_i\}$  are called the *regression coefficients*, and are to be determined based on an observed data-set. The special case of  $r = 1$  is called a *simple* regression, while  $r > 1$  is the much more complicated *multiple* regression problem.

$$Y = \alpha + \beta x + e \quad (1.3)$$

The choice of a simple linear regression model is appropriate when the data appears to follow a straight line relationship subject to random error when visualized as a scatter plot.

## 1.2 Least squares estimators of regression parameters

In a simple regression problem, let the estimators of  $\alpha, \beta$  be  $A, B$ . For a given set of inputs  $\{x_i\}$  and responses  $\{Y_i\}$ , the estimated response is

$$\hat{Y}_i = A + Bx_i \quad (1.4)$$

Using an expression for the squared difference between the observed and estimated responses, it is possible to find an estimator that minimizes the sum of these squares.

$$SS = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - A - Bx_i)^2 \quad (1.5)$$

The usual method of setting partial derivatives to zero in order to find the minimum, yields a system of 2 linear equations in  $A, B$  called the *normal equations*

$$\sum Y_i = A n + B \sum x_i \quad (1.6)$$

$$\sum x_i Y_i = A \sum x_i + B \sum x_i^2 \quad (1.7)$$

Using the shorthand notation  $\bar{Y}, \bar{x}$  for their sample means, and solving for the least squares estimators gives the estimated regression line as  $y = A + Bx$ , where

$$A = \bar{Y} - B\bar{x} \quad B = \frac{\sum (x_i - \bar{x})Y_i}{\sum x_i^2 - n\bar{x}^2} \quad (1.8)$$

### 1.3 Distribution of the estimators

In order to determine the distribution of the estimators  $A, B$ , an additional assumption about the random errors  $e$  is used. For some value  $\sigma^2$  which is a constant independent of the input,

$$e \sim \mathcal{N}(0, \sigma^2) \quad Y_i \sim \mathcal{N}(\alpha + \beta x_i, \sigma^2) \quad (1.9)$$

From the above expression for the estimator  $B$ , it is some linear combination of normal RVs  $\{Y_i\}$ . Substituting  $B$  into  $A$  shows that  $A$  is also a normal RV,

$$\mathbb{E}[B] = \beta \quad \text{Var}(B) = \frac{\sigma^2}{\sum x_i^2 - n\bar{x}^2} \quad (1.10)$$

$$\mathbb{E}[A] = \alpha \quad \text{Var}(A) = \frac{\sigma^2}{\sum x_i^2 - n\bar{x}^2} \frac{\sum x_i^2}{n} \quad (1.11)$$

From the above expectation values,  $A, B$  are both unbiased estimators of  $\alpha, \beta$ .

*Residuals* : The difference between each observed value and its estimate is called the residual. To determine the variance of the error, the sum of squares of these residuals can be rearranged into a known RV.

$$r_i = Y_i - \hat{Y}_i = Y_i - A - Bx_i \quad (1.12)$$

$$SS_R = \sum r_i^2 = \sum (Y_i - A - Bx_i)^2 \quad (1.13)$$

Rearranging the  $SS_R$  into a chi-square distribution gives an unbiased estimator of  $\sigma^2$

$$\frac{SS_R}{\sigma^2} \sim \chi_{n-2}^2 \quad \mathbb{E} \left[ \frac{SS_R}{n-2} \right] = \sigma^2 \quad (1.14)$$

Similar to the sample mean and sample variance from a normal population being independent, the variance of the noise  $\sigma^2$  is independent of the estimators  $A, B$ .

The MLE of  $\alpha, \beta$  also happen to be the least squares estimators  $A, B$ .

*Shorthand notation for sums of squares :* For convenience, some shorthand expressions for the sums of squares are listed here.

$$S_{xY} = \sum (x_i - \bar{x}) (Y_i - \bar{Y}) = \sum x_i Y_i - n\bar{x}\bar{Y} \quad (1.15)$$

$$S_{xx} = \sum (x_i - \bar{x})^2 = \sum x_i^2 - n\bar{x}^2 \quad (1.16)$$

$$S_{YY} = \sum (Y_i - \bar{Y})^2 = \sum Y_i^2 - n\bar{Y}^2 \quad (1.17)$$

Using the above shorthand, the estimators  $A, B$  can be described as

$$B = \frac{S_{xY}}{S_{xx}} \sim \mathcal{N}\left(\beta, \frac{\sigma^2}{S_{xx}}\right) \quad (1.18)$$

$$A = \bar{Y} - B\bar{x} \sim \mathcal{N}\left(\alpha, \frac{\sigma^2}{S_{xx}} \frac{\sum x_i^2}{n}\right) \quad (1.19)$$

$$SS_R = \frac{S_{xx} S_{YY} - S_{xY}^2}{S_{xx}} \quad (1.20)$$

The above relation for the  $SS_R$  can only be established using brute-force computations, so a theoretical justification is not outlined here.

## 1.4 Statistical Inferences about regression parameters

The problem of constructing hypothesis tests is straightforward given the transformation into known RVs outlined above.

*Inferences on  $\beta$*  : A hypothesis that the average response is independent of the input requires rearranging the estimator  $B$  into a t-RV and then defining the hypothesis test as

$$\begin{aligned} \frac{B - \beta}{\sigma/\sqrt{S_{xx}}} &\sim Z \quad \text{and} \quad \frac{SS_R}{\sigma^2} \sim \chi_{n-2}^2 \\ \sqrt{\frac{(n-2)S_{xx}}{SS_R}} (B - \beta) &\sim t_{n-2} \\ Y &= \alpha + \beta x + e \end{aligned} \tag{1.21}$$

Define a hypothesis test of significance level  $\gamma$  as

$$\begin{aligned} H_0 : \beta &= 0 \quad \text{vs.} \quad H_1 : \beta \neq 0 \\ \text{reject } H_0 &\text{ if } \sqrt{\frac{(n-2)S_{xx}}{SS_R}} |B| > t_{\gamma/2, n-2} \\ \text{accept } H_0 &\quad \text{otherwise} \end{aligned} \tag{1.22}$$

The corresponding  $100(1 - \gamma)\%$  confidence interval for  $\beta$  is give by

$$\beta \in \left[ B \pm \sqrt{\frac{SS_R}{(n-2)S_{xx}}} t_{\gamma/2, n-2} \right] \tag{1.23}$$

*Regression to the mean* : A linear simple regression with the parameter  $\beta \in [0, 1]$ , displays the following property,

$$\begin{aligned} \mathbb{E}[Y] &= \alpha + \beta x \\ \mathbb{E}[Y] &< x \quad \forall x > \frac{\alpha}{1 - \beta} \\ \mathbb{E}[Y] &> x \quad \text{otherwise} \end{aligned} \tag{1.24}$$

This trend is very common in real-world datasets and was first observed in the comparison of height at a given age between successive generations of a family. Even though the variables are positively correlated,  $\beta \in [0, 1]$  causes the extreme values to regress towards the  $y = x$  line.

A test of the hypothesis  $\beta \geq 1$  can be rejected in such real-world datasets as confirmation of regression to the mean.

*Regression fallacy* : The false attribution to some outside influence on the observed phenomenon of regression to the mean, when it might be happening simply because of chance.

*Inferences on  $\alpha$*  : Similar to the previous results for  $B$ , the confidence interval for  $\alpha$  is given by,

$$\sqrt{\frac{n(n-2)S_{xx}}{SS_R \sum x_i^2}} (A - \alpha) \sim t_{n-2} \quad (1.25)$$

$$\alpha \in \left[ A \pm \sqrt{\frac{SS_R \sum x_i^2}{n(n-2)S_{xx}}} t_{\gamma/2, n-2} \right] \quad (1.26)$$

Similar to the hypothesis tests on  $\beta$ , the tests on  $\alpha$  are constructed as follows.

$$\begin{array}{ll} H_0 : \alpha = 0 & \text{vs.} \quad H_1 : \alpha \neq 0 \\ \text{reject } H_0 & \text{if } \sqrt{\frac{n(n-2)S_{xx}}{SS_R \sum x_i^2}} |A| > t_{\gamma/2, n-2} \\ \text{accept } H_0 & \text{otherwise} \end{array} \quad (1.27)$$

*Inferences on the mean response* : A point estimator for  $Y(x_0)$  in a simple linear regression model is clearly  $\hat{Y} = A + Bx_0$ . This is additionally an unbiased estimator since  $\mathbb{E}[A] = \alpha$ ,  $\mathbb{E}[B] = \beta$ .

Since  $A, B$  are themselves normal RVs with pre-calculated mean and variance,



$$\mathbb{E}[A + Bx_0] = \alpha + \beta x_0 \quad (1.28)$$

$$\text{Var}(A + Bx_0) = \sigma^2 \left[ \frac{1}{n} + \frac{(\bar{x} - x_0)^2}{S_{xx}} \right] \quad (1.29)$$

$$A + Bx_0 \sim N \left( \alpha + \beta x_0, \sigma^2 \left[ \frac{1}{n} + \frac{(\bar{x} - x_0)^2}{S_{xx}} \right] \right) \quad (1.30)$$

Using the fact that the mean response is independent of the random error, which can be rearranged to form a chi-square RV, a t-RV can be formulated and then used to define a confidence interval.

$$\alpha + \beta x_0 \in A + Bx_0 \pm t_{\gamma/2, n-2} \sqrt{\left( \frac{SS_R}{n-2} \right) \left( \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)} \quad (1.31)$$

Note that the lack of direct knowledge about  $\sigma^2$ , forced the use of a t-RV instead of the naive choice of a normal-RV when defining the confidence interval.

*Prediction interval of a future response* : Unlike the mean response, a point estimate of the response to a new input  $x_0$  is defined as  $Y(x_0) = \alpha + \beta x_0 + e$ . A point prediction of  $Y(x_0)$  is simply  $A + Bx_0$ , since  $A, B$  are the unbiased point estimators of  $\alpha, \beta$  respectively.

Given that the mean, median and mode of a normal RV are all equal, the question of choosing one over the other does not arise here. Using the fact that both the next response  $Y$  and the predicted value  $A + Bx_0$  are normal RVs, gives

$$Y - A - Bx_0 \sim N \left( 0, \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(\bar{x} - x_0)^2}{S_{xx}} \right] \right) \quad (1.32)$$

Once again rearranging the  $Y - A - Bx_0$  and the known chi-square RV  $SS_R/\sigma^2$  into a t-RV gives the confidence interval for the prediction as

$$Y \in A + Bx_0 \pm t_{\gamma/2, n-2} \sqrt{\left( \frac{SS_R}{n-2} \right) \left( 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)} \quad (1.33)$$

Note that the results of a linear simple regression cannot be used to make predictions about the response to an input very different from the data used initially to perform the regression.

## 1.5 Coefficient of Determination and Sample correlation coefficient

The variation in successive values of the response  $\{Y_i\}$  can result from the variation in the corresponding inputs  $\{x_i\}$  as well as the variance of the random noise  $e$  which can cause even the same input to produce different responses.

The coefficient of determination  $R^2$  is defined as the proportion of the variation that is explained by the change in input values.

$$R^2 = \frac{S_{YY} - SS_R}{S_{YY}} \quad (1.34)$$

$R^2 \in [0, 1]$  with values close to 1 implying that the variation in response is almost fully explained by the change in input. A large value of  $R$  indicates that the linear regression is a good fit.

*Relation to sample correlation coefficient* :  $r$  has been defined previously as

$$r^2 = \frac{S_{xY}^2}{S_{xx} S_{YY}} = \frac{S_{xx} S_{YY} - SS_R S_{xx}}{S_{xx} S_{YY}} = R^2 \quad (1.35)$$

$$|r| = \sqrt{R} \quad (1.36)$$

While  $r \in [-1, 1]$  denotes the correspondence between increase in values of one variable and another, its square represents the extent to which a simple linear regression can explain the set of two data points.

## 1.6 Analysis of residuals: assessing the model

Even though visual inspection of the scatter plot is a good method to determine how appropriate the choice of linear simple regression model is, any further doubts can be taken care of using the residuals  $Y_i - (A + Bx_i)$ .

*Standardized residuals* are obtained by transforming the residuals into standard normal RVs as follows,

$$Z \sim \frac{Y_i - (A + Bx_i)}{\sqrt{SS_R/(n-2)}} \quad \forall i \in \{1, \dots, n\} \quad (1.37)$$

Any observable pattern in a plot of the residuals that indicates deviation from  $Z$  is immediate evidence against choosing a simple linear regression model.

## 1.7 Transforming to linearity

In cases where the response is not a linear function of the input, it is useful to transform the relation using logarithms, exponentiation or other tools into a linear equation. This extends the power of the linear regression method to systems governed by non-linear relations.

For example, consider an exponential decay relationship between the number of items  $N$  and time  $t$ ,

$$\begin{aligned} N(t) &= u \exp(-vt) \\ \log(N) &= \log(u) - v \log(t) \\ \text{substituting } Y = \log(N) \quad \alpha = \log(u) \quad \beta = -v \quad x = \log(t) \\ Y &= \alpha + \beta x + e \end{aligned}$$

## 1.8 Weighted Least Squares

When the assumption that the variance in the response is a constant independent of the input is no longer reasonable, the above procedure has to be modified to incorporate a weighting

term in the variance,

$$\text{Var}(Y_i) = \frac{\sigma^2}{w_i} \quad (1.38)$$

The estimators  $A, B$  must now minimize the sum of squares with these weights  $\{w_i\}$ . This is once again solved by setting the partial derivatives to zero and obtaining the *normal equations*.

$$SS_w = \frac{1}{\sigma^2} \sum_{i=1}^n w_i (Y_i - A - Bx_i)^2 \quad (1.39)$$

$$\sum w_i Y_i = A \sum w_i + B \sum w_i x_i \quad (1.40)$$

$$\sum w_i x_i Y_i = A \sum w_i x_i + B \sum w_i x_i^2 \quad (1.41)$$

The least squares estimator, even if unbiased, need not be the best estimator of the mean of a normal RV.

For a sample from a normal RV, the weighted least squares estimators of  $\alpha, \beta$  happen to also be the MLE. Alternatively, consider the transformation to the error  $e$ , which removes its dependence on the input.

$$\begin{aligned} Y &= \alpha + \beta x + e \\ Y\sqrt{w} &= \alpha\sqrt{w} + \beta x\sqrt{w} + e\sqrt{w} \end{aligned} \quad (1.42)$$

The new error term  $e\sqrt{w}$  is now independent of the input. It has mean 0 and constant variance. Now, the least squares estimators of  $\alpha, \beta$  would be the ones that minimize

$$\begin{aligned} & \text{minimize } \sum_{i=1}^n (Y_i \sqrt{w_i} - \alpha \sqrt{w_i} - \beta x_i \sqrt{w_i})^2 \\ & \text{minimize } \sum_{i=1}^n w_i (Y_i - \alpha - \beta x_i)^2 \end{aligned}$$

The weighted least squares method has the added advantage of giving greater emphasis to data points with the lower variance (and thus greater  $w_i$ ).

For the special case when  $Y$  is a Poisson RV, a common transformation is to model  $\sqrt{Y}$  as a linear simple regression, with the added approximation that  $\text{Var}(\sqrt{Y}) \approx 0.25$  regardless of the parameter  $\lambda$ .

## 1.9 Polynomial Regression

The extension of the linear regression model to include terms corresponding to higher powers of the single input is straightforward, despite the heavy use of matrix operations.

Polynomial regression follows the same outline as simple regression with the minimization of  $SS_R$  leading to a set of *normal equations*.

For convenience, define  $\sum_1^n (x_i)^j = K_j$  and for the RHS,  $\sum_i^n (x_i)^j Y_i = M_j$

$$\begin{bmatrix} K_0 & K_1 & K_2 & \dots & K_r \\ K_1 & K_2 & K_3 & \dots & K_{r+1} \\ K_2 & K_3 & K_4 & \dots & K_{r+2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ K_r & K_{r+1} & K_{r+2} & \dots & K_{2r} \end{bmatrix} \begin{bmatrix} B_0 \\ B_1 \\ B_2 \\ \vdots \\ B_r \end{bmatrix} = \begin{bmatrix} M_0 \\ M_1 \\ M_2 \\ \vdots \\ M_r \end{bmatrix} \quad (1.43)$$

Since the sets  $\{K_i\}$ ,  $\{M_i\}$  are both fixed for a given dataset, the above system of linear equations is uniquely solved using matrix inversion. The value of  $r$  chosen should be as small as possible in order to avoid overfitting. This choice is usually made through visual inspection of the scatter diagram.

An unnecessarily large choice of degree  $r$  makes the regression model significantly worse at predicting the response to inputs far from the set of data points used to calculate it. In

real-world problems, this also leads to incorrect inferences about the underlying physical mechanisms.

## 1.10 Multiple Linear Regression

Most real-world systems of interest are not governed by a single input. This requires generalizing the simple linear regression problem to deal with multiple inputs  $\{x_i\}$ .

$$\begin{aligned} Y &= \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k + e \\ e &\sim \mathcal{N}(0, \sigma^2) \\ \mathbb{E}[Y_i] &= \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik} \end{aligned} \tag{1.44}$$

Once again, minimizing the  $SS_R$  by setting the derivative to zero yields a system of normal equations.

$$\text{minimizing } \sum_{i=1}^n (Y - B_0 - B_1 x_{i1} - B_2 x_{i2} - \cdots - B_k x_{ik})^2 \tag{1.45}$$

Defining the matrices for shorthand notation as follows,

$$\mathbf{e} = \begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ \vdots \\ e_n \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} B_0 \\ B_1 \\ B_2 \\ \vdots \\ B_k \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \\ \vdots \\ Y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ 1 & x_{31} & x_{32} & \cdots & x_{3k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix} \tag{1.46}$$

The regression model and thus the normal equations can now be written as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \quad (1.47)$$

$$\mathbf{X}^\top \mathbf{X} \mathbf{B} = \mathbf{X}^\top \mathbf{Y} \quad (1.48)$$

Since the existence of the inverse  $(\mathbf{X}^\top \mathbf{X})^{-1}$  is not an issue in most real world datasets, simple matrix operations can be used to obtain the estimator matrix  $\mathbf{B}$ .

The set of least squares estimators  $\mathbf{B}$  also happen to be unbiased estimators of  $\boldsymbol{\beta}$ . The variance of these estimators requires further analysis of the matrix  $\mathbf{C} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$

$$\mathbf{B} = \mathbf{C} \mathbf{Y}$$

$$B_{i-1} = \sum_{m=1}^n C_{im} Y_m \quad (1.49)$$

$$\text{Cov}(B_{i-1}, B_{j-1}) = \sigma^2 (\mathbf{C} \mathbf{C}^\top)_{ij} \quad (1.50)$$

$$\text{Cov}(\mathbf{B}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \quad (1.51)$$

Since  $\text{Cov}(B_i, B_i) = \text{Var}(B_i)$ , the diagonal elements of the matrix of covariances  $\text{Cov}(\mathbf{B})$ , give the variances of the least squares estimators (scaled by  $\sigma^2$ ).

The value of  $\sigma^2$  can be estimated by rearranging the sum of squares of residuals  $SS_R$  into a chi-squared RV.

$$SS_R = \sum_{i=1}^n (Y - B_0 - B_1 x_{i1} - B_2 x_{i2} - \dots - B_k x_{ik})^2$$

$$\frac{SS_R}{\sigma^2} \sim \chi_{n-(k+1)}^2 \quad (1.52)$$

$$\mathbb{E} \left[ \frac{SS_R}{n - (k + 1)} \right] = \sigma^2 \quad (1.53)$$

The above expression is an unbiased estimator of  $\sigma^2$  and is also independent of the set of estimators  $\mathbf{B}$ .

Defining the residual as  $r_i$  along with the column matrix of the set of residuals  $\mathbf{r}$  yields a useful computational formula for  $SS_R$

$$\mathbf{r} = \mathbf{Y} - \mathbf{XB} \quad (1.54)$$

$$SS_R = \mathbf{r}^\top \mathbf{r} \quad (1.55)$$

$$= \mathbf{Y}^\top \mathbf{Y} - \mathbf{B}^\top \mathbf{X}^\top \mathbf{Y} \quad (1.56)$$

*Coefficient of multiple determination* : The quantity  $R^2$  is defined as follows. It measures the decrease in the sum of squares of residuals  $SS_r$  when using a multiple regression model.

$$R^2 = 1 - \frac{SS_R}{\sum (Y_i - \bar{Y})^2} \quad (1.57)$$

$$Y = \beta_0 + \beta_1 x_1 + \cdots + \beta_n x_n + e \quad \text{vs} \quad Y = \beta_0 + e \quad (1.58)$$

*Predicting future responses* : A point estimate of the mean response is simply

$$\begin{aligned} \mathbb{E}[Y|\mathbf{x}] &= \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k \\ \widehat{\mathbb{E}[Y|\mathbf{x}]} &= \sum_{i=0}^k B_i x_i \end{aligned} \quad (1.59)$$

Here the first term is special because  $x_o \equiv 1$ . In order to find an interval estimate, consider the distribution of the above point estimator, which is a linear combination of normal RVs  $\{Y_i\}$ .



$$\mathbb{E} \left[ \sum_{i=0}^k B_i x_i \right] = \sum_{i=0}^k \beta_i x_i \quad (1.60)$$

$$\text{Var} \left( \sum_{i=0}^k B_i x_i \right) = \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x} \sigma^2 \quad (1.61)$$

$$(1.62)$$

Here,  $\mathbf{x}$  is the column matrix composed of the set of input variables  $\{x_0, \dots, x_k\}$ . Additionally, replacing  $\sigma$  with its estimator  $\sqrt{SS_R/(n-k-1)}$ , yields a t-RV with  $n-k-1$  DOF.

$$\frac{\sum_{i=0}^k B_i x_i - \sum_{i=0}^k \beta_i x_i}{\sqrt{\mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x} \left( \frac{SS_R}{n-k-1} \right)}} \sim t_{n-k-1} \quad (1.63)$$

The above distribution can be used to construct confidence intervals for  $\mathbb{E}[Y|\mathbf{x}]$ .

As opposed to predicting the mean value of the response, when the experiment is to be performed only once in order to generate a single data point, it is more relevant to predict the response  $Y(\mathbf{x})$  itself.

$$Y(\mathbf{x}) = \sum_{i=0}^k \beta_i x_i + e \quad \text{where} \quad x_0 \equiv 1 \quad (1.64)$$

The point estimator is clearly just  $\sum B_i x_i$  based on the previous  $n$  data points. Since  $\sum B_i x_i$  are thus independent of  $Y(\mathbf{x})$ ,

$$\text{Var} \left[ Y(\mathbf{x}) - \sum_{i=0}^k B_i x_i \right] = \sigma^2 + \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x} \sigma^2 \quad (1.65)$$

$$\frac{Y(\mathbf{x}) - \sum_{i=0}^k B_i x_i}{\sqrt{(1 + \mathbf{x}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}) \left( \frac{SS_R}{n - k - 1} \right)}} \sim t_{n-k-1} \quad (1.66)$$

*Dummy variables for categorical data* : A binary categorical variable can be represented as a dummy variable in the regression that takes values of  $\{0, 1\}$  depending on the data point belonging to the category or not. This is only necessary when the dataset is too small to be fragmented into multiple datasets each analyzed separately.

For large enough datasets, it is better to avoid the use of dummy variables and simply perform many separate regression analyses on the different data-subsets.

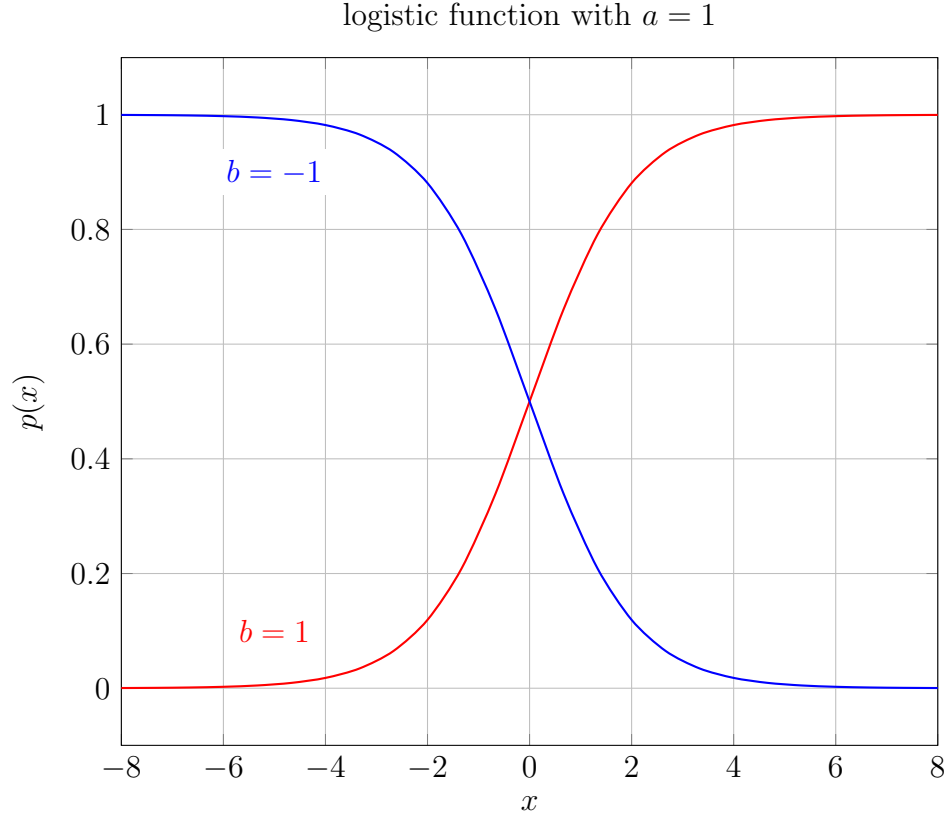
*Logistic regression models for binary response data* : Consider an experiment which only gives a binary response, with the probability of success  $p(x)$  governed by the distribution. Additionally define the odds of success  $o(x)$  as the ratio of the win to loss probabilities.

$$p(x) = \frac{\exp(a + bx)}{1 + \exp(a + bx)} \quad (1.67)$$

$$o(x) = \frac{p(x)}{1 - p(x)} = \exp(a + bx) \quad (1.68)$$

$$\log(o(x)) = a + bx \quad (1.69)$$

From the plot below, the asymptotic values of this function at  $-\infty, \infty$  are 0, 1 depending on the sign of  $b$ .



In order to find the maximum likelihood estimators, consider the joint PDF of a set of binary responses  $\{Y_1, \dots, Y_k\}$ ,

$$\log(P\{Y_i = y_i; i = 1, \dots, k\}) = \sum_{i=1}^k y_i(a + bx_i) - \sum_{i=1}^k \log(1 + \exp(a + bx_i)) \quad (1.70)$$

Even though an analytical minimization of the above expression is not possible, there are several iterative computational approaches possible.

# Chapter 2

## Analysis of Variance

*“Ch10 Quote Here.”*

### 2.1 Introduction