

Elements of statistical learning: Chapter 3

July 29, 2020

1 Linear models

- Sampling properties of $\hat{\beta}$
- Gauss-Markov Theorem

2 Multiple regression

- From simple univariate to multiple regressions

LS estimator

Let \mathbf{X} be an $N \times (p + 1)$ matrix of explanatory variables and \mathbf{y} an $N \times 1$ vector of outputs. Then we know the LS estimator $\hat{\beta}$

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y},$$

[see lecture slides "ESL1" for recap and proof].

The "hat" matrix

As such for the fitted linear model

$$\begin{aligned}\hat{\mathbf{y}} &= \mathbf{X}\hat{\beta} \\ &= \underbrace{\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'}_H \mathbf{y} \\ &= H\mathbf{y}\end{aligned}$$

where H is commonly referred to as the hat matrix.

H the projection matrix

Let us denote the column vectors of \mathbf{X} by $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_p$ with $\mathbf{x}_0 \equiv 1$.

- These vectors span a subspace of \mathbb{R}^N , also referred to as a column vector of \mathbf{X} .
- We minimize $RSS(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2$ by choosing $\hat{\beta}$, so that the residual vector $\mathbf{y} - \hat{\mathbf{y}}$ is orthogonal to this subspace.
- the hat matrix H computes the orthogonal projection, and hence it is also known as the projection matrix.

Variance-covariance matrix

Assumptions

- 1 Observations y_i are uncorrelated have constant variance σ^2
- 2 x_i are fixed (i.e. non-stochastic)

$$\begin{aligned}\text{var}(\hat{\beta}) &= \text{var} [(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}] \\ &= \text{var} [(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\beta + \epsilon)] \\ &= \text{var} [(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{X})\beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\epsilon] \\ &= \text{var} [(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\epsilon] \\ &= \mathbb{E} \{ (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\epsilon[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\epsilon]'\} \\ &= \mathbb{E} \{ (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\epsilon\epsilon'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \} \\ &= \mathbb{E} \{ (\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{X})\epsilon\epsilon'(\mathbf{X}'\mathbf{X})^{-1} \}\end{aligned}$$

Note that ϵ is the error term and has zero mean and also remember that \mathbf{X} is fixed, and thus

$$\mathbb{E}[aZ] = a\mathbb{E}[Z]$$

where Z is a random variable and a is a constant. Therefore,

$$\begin{aligned} \text{var}(\hat{\beta}) &= \mathbb{E} \{ \epsilon \epsilon' (\mathbf{X}' \mathbf{X})^{-1} \} \\ &= (\mathbf{X}' \mathbf{X})^{-1} E \{ \epsilon \epsilon' \} \\ &= (\mathbf{X}' \mathbf{X})^{-1} \sigma^2 \end{aligned}$$

where σ^2 can be calculated by

$$\sigma^2 = \frac{1}{N - p - 1} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

thus, assuming the errors are further Gaussian

$$\hat{\beta} \sim N(\beta, (\mathbf{X}' \mathbf{X})^{-1} \sigma^2)$$

Gauss-Markov Theorem

Least squares estimator of parameter β has the smallest variance among all linear unbiased estimators. Why is the LS estimator unbiased?

Proof.

$$\begin{aligned}\hat{\beta} &= \mathbb{E}[\hat{\beta}] \\ &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}] \\ &= \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\beta + \epsilon)] \\ &= \mathbb{E}[\beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\epsilon] \\ &= \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbb{E}[\epsilon] \\ &= \beta\end{aligned}$$



From simple univariate to multiple regressions

Suppose first we have a univariate model with no intercept

$$Y_i = X_i\beta + \varepsilon_i$$

The least squares estimates and residuals are

$$\hat{\beta} = \frac{\sum_{i=1}^N x_i y_i}{\sum_{i=1}^N x_i^2}$$

with residuals

$$r_i = y_i - x_i \hat{\beta}$$

which in vector notation can be expressed as

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^N x_i y_i = \mathbf{x}' \mathbf{y}$$

which is the inner product between \mathbf{x} and \mathbf{y} .

Thus, the OLS estimator $\hat{\beta}$ can be expressed as follows

$$\hat{\beta} = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle},$$

Suppose now that we have p inputs $\mathbf{x}_1, \dots, \mathbf{x}_p$, which are the columns of the matrix \mathbf{X} and are orthogonal, such that $\langle \mathbf{x}_i, \mathbf{x}_k \rangle = 0$ for all $i \neq j$. When the inputs are orthogonal, the multiple least squares estimates $\hat{\beta}_j$ are equal to the univariate estimates - i.e.

$$\hat{\beta}_j = \frac{\langle \mathbf{x}_j, \mathbf{y} \rangle}{\langle \mathbf{x}_j, \mathbf{x}_j \rangle}$$

In other words, the inputs are orthogonal and have no impact on each other's parameters estimates in the model.

Consider the case of an intercept and a single input \mathbf{x} , then the least squares coefficient of \mathbf{x} has the form

$$\hat{\beta}_1 = \frac{\langle \mathbf{x} - \bar{x}1, \mathbf{y} \rangle}{\langle \mathbf{x} - \bar{x}1, \mathbf{x} - \bar{x}1 \rangle}$$

The steps of the algorithm can be seen as follows

- 1 Regress \mathbf{x} on 1 to obtain $\bar{x}1$
- 2 Obtain the residuals $\mathbf{z} = \mathbf{x} - \bar{x}1$
- 3 Regress \mathbf{y} on \mathbf{z} to obtain the coefficient $\hat{\beta}_1$

$$\hat{\beta}_1 = \frac{\langle \mathbf{z}, \mathbf{y} \rangle}{\langle \mathbf{z}, \mathbf{z} \rangle}$$

Step 1 orthogonalizes x with respect to $x_0 = 1$.