

Regression

Sofia Olhede



November 2, 2020

- 1 Linear Regression
 - Least squares regression
 - Residuals
 - Confidence intervals for coefficients and variance
 - Confidence intervals for coefficients and variance
 - Regression Diagnostics and Distribution Plots

Set-up

- Consider a set of measurements given by the response variable Y_i and with a corresponding set of predictor variables x_{i1}, \dots, x_{ip} . Hence the data set is

$$\{y_i, x_{i1}, \dots, x_{ip}\}_{i=1}^n.$$

- Definition: A linear model is

$$\mathbb{E}\{Y\} = X\beta,$$

$n \times 1$ $n \times p \cdot p \times 1$

where $Y = (Y_1 \dots Y_n)^T$, is the vector of observations, X is the known $n \times p$ design matrix and $\beta = (\beta_1 \dots \beta_p)^T$ is the $p \times 1$ parameter vector.

- We are trying to quantify the systemic variation in Y due to $X\beta$.

Linear Regression

- Example: polynomial regression. This can be written as

$$\mathbb{E}\{Y_i\} = \beta_0 + \beta_1 x_i + \cdots + \beta_p x_i^p,$$

where x_i is the i th predictor variable corresponding to Y_i .

- For example we might fit a linear model of the form

$$\mathbb{E}\{Y_i\} = \beta_1 + \beta_2 x_{i1} + \beta_3 x_{i2} + \beta_4 x_{i1} x_{i2} + \beta_5 x_{i3}^2,$$

where x_{ki} is the value of the k th predictor for observation i .

- Note that

$$E(Y_i) = \beta_1 + \beta_2 x^{\beta_3},$$

is not a linear model.

- We will assume $p < n$ (full rank).

A matrix is said to have full rank if its rank equals the largest possible for a matrix of the same dimensions, which is the lesser of the number of rows and columns. The maximum number of linearly independent rows (columns) in a matrix A is called the row (column) rank of A .

- The rank of the matrix X is the dimension of the space spanned by the columns of X . Assume $\text{rank}(X)=p$.

Linear Regression

- We can also add further assumptions

Second-order assumptions (SOA) $\text{var}(Y) = \sigma^2 I_n$ where σ^2 is unknown. Thus $\text{var}(Y_i) = \sigma^2$ for all i and the Y_i s are uncorrelated.

Normal theory assumptions (NTA) The Y_i s are independently and normally distributed with common unknown variance σ^2 so


$$Y \sim N(X\beta, \sigma^2 I_n).$$

- NTA implies SOA but for now we will only assume the weaker SOA.

Linear Regression

- The linear model can be rewritten as

$$Y = X\beta + \epsilon$$

 noise

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}.$$

where $E(\epsilon) = 0$ and $\text{var}(\epsilon) = \sigma^2 I_n$.

- Minimise the difference between the observed values and the model fit to it.

Linear Regression

the residual sum of squares (RSS), also known as the sum of squared residuals (SSR) or the sum of squared estimate of errors (SSE), is the sum of the squares of residuals (deviations predicted from actual empirical values of data). It is a measure of the discrepancy between the data and an estimation model. A small RSS indicates a tight fit of the model to the data. It is used as an optimality criterion in parameter selection and model selection

- Find $\hat{\beta}$ that minimise the residual sum of squares (RSS), i.e. find

$$\hat{\beta} = \arg \min_{\beta} (\epsilon^T \epsilon = \sum_{i=1}^n \epsilon_i^2).$$

from before:
 $\epsilon = Y - X\beta$

- Write $\theta = X\beta$. Then $\theta \in R(X) = \Theta$, (the vector space spanned by the columns of X).

least square estimate

- The lse is the $\hat{\theta}$ that minimises $\|Y - \theta\|^2$, the square of the length of $Y - \theta$. This is minimised when $Y - \hat{\theta}$ is perpendicular to Θ .

- v , is perpendicular to Θ if $X^T v = 0$. Thus

$$X^T(Y - \hat{\theta}) = 0 \quad \text{so} \quad \hat{\beta} = (X^T X)^{-1} X^T Y,$$

closed-form solution

if $X^T X$ is invertible.

Linear Regression

- Here, $\hat{\beta}$ is the **ordinary least squares estimate** of β and is **unique**.
- Or:

$$\begin{aligned}\epsilon^T \epsilon &= (Y - X\beta)^T (Y - X\beta) \\ &= Y^T Y - 2\beta^T X^T Y + \beta^T X^T X \beta,\end{aligned}$$

- $\beta^T X^T Y = Y^T X \beta$ (both are scalars).
 - Differentiating wrt β and setting to zero we see that
- closed form solution as we know it:*

$$-2X^T Y + 2X^T X \beta = 0$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y,$$

as

$$\frac{\partial}{\partial \beta} (a^T \beta) = a, \quad \frac{\partial}{\partial \beta} (\beta^T A \beta) = 2A\beta.$$

Linear Regression

- $\hat{\beta}$ is linear in Y , and $\hat{\beta}$ is unbiased for β :

$$\begin{aligned} E(\hat{\beta}) &= (X^T X)^{-1} X^T E(Y) \quad \leftarrow E(Y) = X\beta \\ &= \underbrace{(X^T X)^{-1} X^T (X\beta)}_{=I} = \beta, \end{aligned}$$

- Let $A = (X^T X)^{-1} X^T$:

$$\begin{aligned} \text{Var}(\hat{\beta}) &= \text{Var}(AY) \\ &= A \text{Var}(Y) A^T \quad \leftarrow \text{formula 375 Matrix Cookbook} \textcircled{f} \\ &= \sigma^2 A A^T \\ &= \sigma^2 (X^T X)^{-1} \underbrace{X^T X (X^T X)^{-1}}_{=I} \\ &= \sigma^2 (X^T X)^{-1}, \end{aligned}$$

as

$$\textcircled{1} \quad \text{Var}(AY) = A \text{Var}(Y) A^T.$$

Linear Regression

- **Gauss-Markov Theorem** Among all unbiased linear estimates of β for a full rank linear model satisfying SOA, any linear combination of the least squares estimator $\hat{\beta}$ has the smaller or equal variance to that of any other, e.g.
- $$\text{Var}\{a^T \hat{\beta}\} \leq \text{Var}\{a^T \tilde{\beta}\}$$

Proof Write another estimator $\tilde{\beta} = BY$ (linearity). We can calculate the expectation of this estimator to be

$$\begin{aligned} \mathbb{E}\{\tilde{\beta}\} &= B \mathbb{E}\{Y\} \\ &= BX\beta = \beta. \end{aligned} \quad (1)$$

This implies that $BX = I$. We define

$$C = B - (X^T X)^{-1} X^T \quad (2)$$

$$\tilde{\beta} = (C + (X^T X)^{-1} X^T) Y = \hat{\beta} + CY. \quad (3)$$

and $CX = 0$ to preserve unbiasedness.

Linear Regression

- For any constant vector \mathbf{a} we note

$$\begin{aligned}\mathbb{V}\text{ar}\{\mathbf{a}^T \tilde{\boldsymbol{\beta}}\} &= \mathbb{V}\text{ar}\{\mathbf{a}^T \{\hat{\boldsymbol{\beta}} + \mathbf{C}\mathbf{Y}\}\} \\ &= \mathbf{a}^T \mathbb{V}\text{ar}\{\hat{\boldsymbol{\beta}}\} \mathbf{a} + \mathbf{a}^T \mathbb{V}\text{ar}\{\mathbf{C}\mathbf{Y}\} \mathbf{a} + 2 \mathbb{C}\text{ov}\{\mathbf{a}^T \hat{\boldsymbol{\beta}}, \mathbf{a}^T \mathbf{C}\mathbf{Y}\}.\end{aligned}\tag{4}$$

We now only need to show that the covariance term is zero. As

$$\begin{aligned}\mathbb{C}\text{ov}\{\mathbf{a}^T \hat{\boldsymbol{\beta}}, \mathbf{a}^T \mathbf{C}\mathbf{Y}\} &= \mathbf{a}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{C}\text{ov}\{\mathbf{Y}, \mathbf{Y}\} \mathbf{C}^T \mathbf{a} \\ &= 0,\end{aligned}\tag{5}$$

and so the result follows. □.

Simple Linear Regression

- Let $(y = b + mx + \text{noise})$

$$Y_i = \beta_1 + \beta_2 x_i + \epsilon_i, \quad i = 1, \dots, n.$$
- $Y^T = (Y_1, \dots, Y_n)$, $\beta^T = (\beta_1, \beta_2)$ and

$$X = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}.$$

- Assume SOA and NO x_i s are equal

$$X^T X = \begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n 1 & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{bmatrix}$$

$$X^T X = \begin{pmatrix} n & n\bar{x} \\ n\bar{x} & \sum x_i^2 \end{pmatrix}$$

$$(X^T X)^{-1} = \frac{1}{n \sum x_i^2 - n^2 \bar{x}^2} \begin{pmatrix} \sum x_i^2 & -n\bar{x} \\ -n\bar{x} & n \end{pmatrix}$$

$$X^T Y = \begin{pmatrix} n\bar{y} \\ \sum x_i Y_i \end{pmatrix}.$$

inverse of 2x2 has a closed form solution:
 if $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$
 then $A^{-1} = \frac{1}{\det A} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$

Simple Linear Regression

Now we can find $\hat{\beta} = (X^T X)^{-1} X^T Y$, hence

$$\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \frac{1}{\sum x_i^2 - n\bar{x}^2} \times \begin{pmatrix} \bar{Y} \sum x_i^2 - \bar{x} \sum x_i Y_i \\ \sum x_i Y_i - n\bar{x}\bar{Y} \end{pmatrix}.$$

$$\hat{\beta}_2 = \frac{S_{xy}}{S_{xx}} = \frac{\sum (x_i - \bar{x})(Y_i - \bar{Y})}{\sum (x_i - \bar{x})^2}$$

$$\hat{\beta}_1 = \bar{Y} - \hat{\beta}_2 \bar{x}.$$

$$\text{Var}(\hat{\beta}) = \frac{\sigma^2}{nS_{xx}} \begin{pmatrix} \sum x_i^2 & -n\bar{x} \\ -n\bar{x} & n \end{pmatrix}.$$

Simple Linear Regression

- If $\bar{x} = 0$ everything becomes easy: the covariance matrix is diagonal and $\hat{\beta}_1 = \bar{Y}$.
- To get a diagonal covariance we adopt the alternative linear model

$$Y_i = \beta_1 + \beta_2(x_i - \bar{x}) + \epsilon_i, \quad i = 1, \dots, n.$$

Then we find that $\hat{\beta}_1 = \bar{Y}$, $\hat{\beta}_2 = S_{xy}/S_{xx}$ and

$$\text{var}(\hat{\beta}) = \begin{pmatrix} n^{-1} & 0 \\ 0 & S_{xx}^{-1} \end{pmatrix}.$$

This idea could be generalised to orthogonal polynomials.

↪?


Linear Regression

- Let $\hat{Y} = X\hat{\beta}$. We found $\hat{\beta}$ by minimising the RSS (Residual Sum of Squares),

$$\begin{aligned}
 e^T e &= \min_{\beta} \epsilon^T \epsilon \\
 &= (Y - X\hat{\beta})^T (Y - X\hat{\beta}) \\
 &= Y^T Y - 2\hat{\beta}^T X^T Y + \hat{\beta}^T X^T X \hat{\beta} \quad \text{from before} \\
 &= Y^T Y - \hat{\beta}^T X^T Y \\
 &\quad + \hat{\beta}^T (X^T X \hat{\beta} - X^T Y) \\
 &= (Y^T - \hat{\beta}^T X^T) Y \\
 &= Y^T (Y - X\hat{\beta}) \\
 &= Y^T Y - \hat{\beta}^T X^T X \hat{\beta}.
 \end{aligned}$$

Linear Regression


- Also the RSS is given by


$$RSS = e^T e = Y^T Y - \hat{Y}^T \hat{Y},$$

the difference between the squares of the observed and fitted Y values.

- The **residuals** of the model are given by the difference between the observed and fitted values so that

$$\begin{aligned} e &= Y - \hat{Y} \\ &= Y - X\hat{\beta} \\ &= \{I_n - X(X^T X)^{-1} X^T\} Y \\ &= (I_n - P)Y, \end{aligned}$$



$P = X(X^T X)^{-1} X^T$ is known as the “hat” matrix and relates the fitted and observed responses as $\hat{Y} = PY$.

Linear Regression

- The hat matrix has a number of known properties:
 1. P is a symmetric $n \times n$ matrix
 2. P is idempotent so that $P^2 = P$
 3. The rank of P is the same as rank X (i.e. both of rank p). From this note $\text{rank}(I_n - P) = n - \text{rank}(P) = n - p$ and that $(I_n - P)$ is also idempotent as

$$(I_n - P)^2 = I_n^2 - 2P + P^2 = I_n - P,$$

as $P^2 = P$.



- Firstly we find the $E(e) = 0$ as

$$E(e) = (I_n - P)E(Y) = (I_n - P)X\beta = 0,$$

$\hookrightarrow e = (I_n - P)Y$ and $E(Y) = \beta X$
 from prev. slide

as

$$\begin{aligned} PX &= X \overbrace{(X^T X)^{-1} X^T}^{= I} X \\ &= X \end{aligned}$$

Linear Regression

- More is known about the residuals:

Theorem The residual sum of squares is an unbiased estimator of $(n - p)\sigma^2$.

- Thus we know that

$$\begin{aligned}
 \hat{\sigma}^2 &= \frac{RSS}{n - p} \\
 &= \frac{(Y - X\hat{\beta})^T (Y - X\hat{\beta})}{n - p} \\
 &= \frac{Y^T Y - \hat{Y}^T \hat{Y}}{n - p},
 \end{aligned}$$

slide 16

is an unbiased estimator of σ^2 .

Linear Regression

- Note that

total sum of squares fitted sum of squares

↑ ↗

$$\begin{aligned}
 \mathbb{E}\{RSS\} &= \mathbb{E}\{Y^T Y - \hat{Y}^T \hat{Y}\} \\
 &= \mathbb{E}\{((I - P)Y)^T ((I - P)Y)\} \\
 &= \mathbb{E}\{\text{trace}\{(I - P)Y\} \{(I - P)Y\}^T\} \\
 &= \mathbb{E}\{\text{trace}\{(I - P)YY^T \{(I - P)\}^T\}\} \\
 &= \sigma^2 \text{trace}(I - P) \\
 &= \sigma^2 \{n - p\}.
 \end{aligned}$$

In linear algebra, the trace of a square matrix A , is defined to be the sum of elements on the main diagonal (from the upper left to the lower right) of A

The result thus follows.

Maximum likelihood approach

- Let $Y \sim N(X\beta, \sigma^2 I_n)$, i.e. NTA.
- The log-likelihood of the data is

Normal Theory Assumption

$$L(\beta, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (Y - X\beta)^T (Y - X\beta).$$

- maximising L with respect to β is equivalent to minimising $(Y - X\beta)^T (Y - X\beta)$
- The maximum likelihood estimate to σ^2 is RSS/n .

i.e. same estimator as Least Square

Maximum likelihood approach

- With NTA:

$$\hat{\beta} \sim N(\beta, \overbrace{\sigma^2(X^T X)^{-1}}^{\text{Covariance}})$$

$$V = \frac{(n-p)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p}^2 \quad \left. \begin{array}{l} \text{independent} \\ \hookrightarrow \text{chi-square dist} \end{array} \right\}$$

- $\hat{\beta}$ and $\hat{\sigma}^2$ are independent.

Theorem 15 If $A = \{a_{ij}\} = (X^T X)^{-1}$ (so $\text{var}(\hat{\beta}) = \sigma^2 A$), then under NTA, the following are $100(1 - \alpha)\%$ confidence intervals for the β_j s and σ^2 :

- $(\hat{\beta}_j - t_{1-\alpha/2} \hat{\sigma} \sqrt{a_{jj}}, \hat{\beta}_j + t_{1-\alpha/2} \hat{\sigma} \sqrt{a_{jj}})$
 - $\left(\frac{(n-p)\hat{\sigma}^2}{\chi_{1-\alpha/2}^2}, \frac{(n-p)\hat{\sigma}^2}{\chi_{\alpha/2}^2} \right)$
- ?

Maximum likelihood approach

repeated slide?

- With NTA:

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

$$V = \frac{(n-p)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p}^2$$

- $\hat{\beta}$ and $\hat{\sigma}^2$ are independent.

Theorem 15 If $A = \{a_{ij}\} = (X^T X)^{-1}$ (so $\text{var}(\hat{\beta}) = \sigma^2 A$), then under NTA, the following are $100(1 - \alpha)\%$ confidence intervals for the β_j s and σ^2 :

1. $(\hat{\beta}_j - t_{1-\alpha/2} \hat{\sigma} \sqrt{a_{jj}}, \hat{\beta}_j + t_{1-\alpha/2} \hat{\sigma} \sqrt{a_{jj}})$
2. $\left(\frac{(n-p)\hat{\sigma}^2}{\chi_{1-\alpha/2}^2}, \frac{(n-p)\hat{\sigma}^2}{\chi_{\alpha/2}^2} \right)$

Residuals

Let

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

but that the analyst incorrectly assumes that

$$Y_i = \beta_0 + \epsilon_i$$

Then

$$\begin{aligned} E\{e_i\} &= E\{Y_i - \hat{\beta}_0\} \\ &= E\left\{Y_i - \frac{1}{n} \sum Y_i\right\} \\ &= \frac{n-1}{n}(\beta_1 x_i) + \frac{1}{n} \sum_{j \neq i} (\beta_1 x_j) \\ &= \beta_1(x_i - \bar{x}) \end{aligned} \tag{6}$$

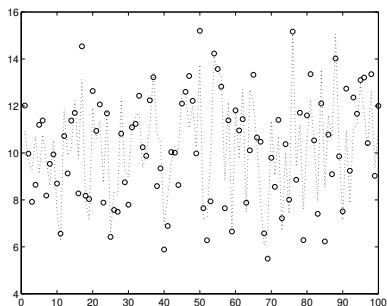


Figure:

Here $Y_i = 10 + 2x_i + 3\epsilon_i$. This is not apparent from the plot, of Y_i (dots) and $E_{Y|\beta, \sigma^2}(Y_i)$ (dotted line).

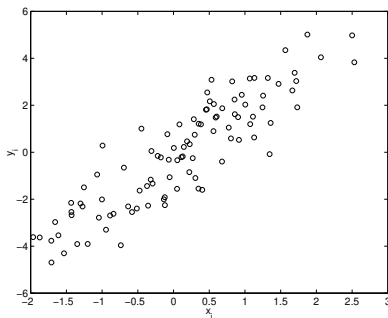


Figure:

Looking at a plot of the residuals against the explanatory variable gives a different opinion.

last slide missing