# Regression

Sofia Olhede

EPFL

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 $x_{i1}, \ldots, x_{ip}$. Hence the data set is

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

- Definition: A linear model is

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

where $Y = \begin{pmatrix} Y_1 & \ldots & Y_n \end{pmatrix}^T$, is the vector of observations, X is the known $n \times p$ design matrix and $\boldsymbol{\beta} = \begin{pmatrix} \beta_1 & \ldots & \beta_p \end{pmatrix}^T$ is the $p \times 1$ parameter vector.

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

# Linear Regression

EPFL

- 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 $p \leq n$ (full rank).
- The rank of the matrix X is the dimension of the space spanned by the columns of X. Assume rank(X)=$p$.

# Linear Regression

- We can also add further assumptions

  Second-order assumptions (SOA) $\text{var}(Y) = \sigma^2 I_n$ where $\sigma^2$ 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 $\sigma^2$ so
  $$Y \sim N(X\beta, \sigma^2 I_n).$$

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

# Linear Regression

- The linear model can be rewritten as

$$
\begin{aligned}
\mathsf{Y} &= \mathsf{X}\beta + \epsilon \\
\begin{pmatrix} Y_1 \\ \cdots \\ Y_n \end{pmatrix} &= \begin{pmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \\
&\quad + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}.
\end{aligned}
$$

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

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

# Linear Regression

EPFL

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

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

- Write $\theta = X\beta$. Then $\theta \in R(X) = \Theta$, (the vector space spanned by the columns of X).
- The lse is the $\widehat{\theta}$ that minimises $||Y - \theta||^2$, the square of the length of $Y - \theta$. This is minimised when $Y - \widehat{\theta}$ is perpendicular to $\Theta$.
- $v$, is perpendicular to $\Theta$ if $X^T v = 0$. Thus

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

if $X^T X$ is invertible.

# Linear Regression

**EPFL**

- Here, $\widehat{\beta}$ is the **ordinary least squares estimate** of $\beta$ 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 $\beta$ and setting to zero we see that

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

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

as

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

# Linear Regression

**EPFL**

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

$$\begin{aligned} E(\widehat{\beta}) &= (X^T X)^{-1} X^T E(Y) \\ &= (X^T X)^{-1} X^T (X\beta) = \beta, \end{aligned}$$

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

$$\begin{aligned} \mathbb{V}ar(\widehat{\beta}) &= \mathbb{V}ar(AY) \\ &= A \ \mathbb{V}ar(Y) \ A^T \\ &= \sigma^2 A A^T \\ &= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} \\ &= \sigma^2 (X^T X)^{-1}, \end{aligned}$$

as

$$\mathbb{V}ar(AY) = A \ \mathbb{V}ar(Y) \ A^T.$$

**EPFL**

# Linear Regression

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

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

$$\mathbb{E}\{\tilde{\beta}\} = B\,\mathbb{E}\{Y\}$$
$$= BX\boldsymbol{\beta} = \boldsymbol{\beta}. \tag{1}$$

This implies that $BX = I$. We define

$$C = B - (X^TX)^{-1}X^T \tag{2}$$
$$\tilde{\beta} = (C + (X^TX)^{-1}X^T)Y = \widehat{\beta} + CY. \tag{3}$$

and $CX = 0$ to preserve unbiasedness.

## Linear Regression

- For any constant vector a we note

$$
\begin{aligned}
\mathbb{Var}\{a^T\widetilde{\boldsymbol{\beta}}\} &= \mathbb{Var}\{a^T\{\widehat{\beta} + CY\}\} \\
&= a^T\,\mathbb{Var}\{\widehat{\beta}\}a + a^T\,\mathbb{Var}\{CY\}a + 2\,\mathbb{Cov}\{a^T\widehat{\beta}, a^T CY\}.
\end{aligned}
\tag{4}
$$

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

$$
\begin{aligned}
\mathbb{Cov}\{a^T\widehat{\beta}, a^T CY\} &= a^T(X^TX)^{-1}X^T\,\mathbb{Cov}\{Y, Y\}C^T a \\
&= 0,
\end{aligned}
\tag{5}
$$

and so the result follows.                                    $\square$.

# Simple Linear Regression

**EPFL**

- Let

$$Y_i = \beta_1 + \beta_2 x_i + \epsilon_i, \qquad i = 1, \ldots, n.$$

- $Y^T = (Y_1, \ldots, Y_n)$, $\beta^T = (\beta_1, \beta_2)$ and

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

- Assume SOA and NO $x_i$s are equal

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

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

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

EPFL

# Simple Linear Regression

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

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

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

$$
\widehat{\beta}_1 = \overline{Y} - \widehat{\beta}_2\overline{x}.
$$

$$
\mathbb{V}\mathrm{ar}(\widehat{\beta}) = \frac{\sigma^2}{nS_{xx}} \begin{pmatrix} \sum x_i^2 & -n\overline{x} \\ -n\overline{x} & n \end{pmatrix}.
$$

# Simple Linear Regression

- If $\bar{x} = 0$ everything becomes easy: the covariance matrix is diagonal and $\widehat{\beta_1} = \bar{Y}$.

- To get a diagonal covariance we adopting the alternative linear model

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

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

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

This idea could be generalised to orthogonal polynomials.

# Linear Regression

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

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

# Linear Regression

- Also the RSS is given by

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

  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

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

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

# Linear Regression

**EPFL**

- 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

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

  as

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

# 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

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

is an unbiased estimator of $\sigma^2$.

# Linear Regression

- Note that

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

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

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

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

# Maximum likelihood approach

- With NTA:

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

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

- $\hat{\beta}$ 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 $\beta_j$s and $\sigma^2$:

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. $\left( \frac{(n-p)\hat{\sigma}^2}{\chi^2_{1-\alpha/2}}, \frac{(n-p)\hat{\sigma}^2}{\chi^2_{\alpha/2}} \right)$

# Maximum likelihood approach

- With NTA:

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

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

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

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

  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. $\left( \frac{(n-p)\hat{\sigma}^2}{\chi^2_{1-\alpha/2}}, \frac{(n-p)\hat{\sigma}^2}{\chi^2_{\alpha/2}} \right)$

**EPFL**

# 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

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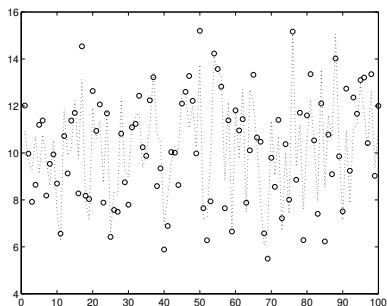&= \beta_1(x_i - \bar{x})
\end{aligned}
$$

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).
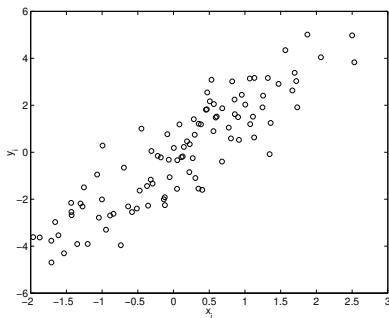
Figure:

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