

SL: LINEAR REGRESSION

In this chapter, we consider the problem of linear regression.

- Traditional approach: observations (X, Y) are assumed to arise from a parametric class of probability distributions, which is known up to some parameter θ :

$$(X, Y) \sim P_{\theta}, \text{ for } \theta \in \Theta$$

↑ parameter space

Specifically, we assume here that

$$Y = \beta_0 + \tilde{\beta}^t x + \varepsilon$$

where

- ▶ input variable x is conditioned on: fixed! Non random!
 $x \in \mathbb{R}^d$, $x = (x_1, \dots, x_d)^t$
 $d = \text{number of predictors/features/covariates}$
 $= \text{dimension of the input space } X = \mathbb{R}^d$
- ▶ $\beta_0 \in \mathbb{R}$, $\tilde{\beta} = (\beta_1, \dots, \beta_d)^t \in \mathbb{R}^d = \text{unknown coefficients}$
- ▶ $\varepsilon \in \mathbb{R} = \text{noise term} \equiv \text{what is left unexplained by the model}$

Typically, one assumes that $\varepsilon \sim \mathcal{N}(0, \sigma^2)$; where the noise variance is also unknown

- Consequences: - Parameter of interest is $\theta = (\beta_0, \tilde{\beta}, \sigma^2)$
 We want to estimate θ
- The conditional distribution of $Y|X=x$ is easily deduced from the distribution of ε :
 $Y|X=x \sim \mathcal{N}(\beta_0 + \tilde{\beta}^t x, \sigma^2)$

- Learning sample $\mathcal{L}_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ (size n)

$$\Rightarrow \begin{cases} Y_1 = \beta_0 + \tilde{\beta}^t x_1 + \varepsilon_1 \\ \vdots \\ Y_n = \beta_0 + \tilde{\beta}^t x_n + \varepsilon_n \end{cases}, \text{ where it is assumed that the } \varepsilon_1, \dots, \varepsilon_n \text{ are independent:}$$

Put $\underline{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^t \in \mathbb{R}^n$

$\Sigma_{\underline{\varepsilon}} = \text{covariance matrix of } \underline{\varepsilon}$
 $= \sigma^2 \mathbb{I}_n$

↑ $(n \times n)$ identity matrix

(also $E \underline{\varepsilon} = 0$)

($\underline{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbb{I}_n)$)

Matrix
Notation

$$\underline{Y} = (Y_1, \dots, Y_n)^t \in \mathbb{R}^n$$

$$\underline{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1d} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{nd} \end{pmatrix}$$

$n \times (d+1)$ matrix

$$\underline{Y} = \underline{X} \underline{\beta} + \underline{\varepsilon}$$

$(n \times 1)$ $(n \times (d+1))$ $(d+1 \times 1)$

$$\underline{\beta} = (\beta_0, \beta_1, \dots, \beta_d)^t \in \mathbb{R}^{d+1}$$

$$= (\beta_0, \tilde{\beta}^t)^t$$

- Challenges:
 - ▶ Parameter estimation & its properties
 - ▶ Model assessment
 - ▶ Prediction

I - MAXIMUM LIKELIHOOD ESTIMATOR & ITS PROPERTIES

The log-likelihood is:

$$l(\theta) = \log \left\{ \frac{1}{(2\pi\sigma^2)^{n/2}} \prod_{i=1}^n \exp\left(-\frac{1}{2\sigma^2} (y_i - \beta_0 - \tilde{\beta}^t x_i)^2\right) \right\}$$

$$l(\theta) = -\frac{n}{2} \log \sigma^2 - \frac{n}{2} \log 2\pi - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \tilde{\beta}^t x_i)^2 \quad (3)$$

⇒ Maximizing the log-likelihood with respect to $(\beta_0, \tilde{\beta})$ (we will deal with σ^2 later) is equivalent to minimize $\sum_{i=1}^n (y_i - \beta_0 - \tilde{\beta}^t x_i)^2$

$$\Leftrightarrow \min_{\beta_0, \tilde{\beta}} \sum_{i=1}^n l(y_i, \beta_0 + \tilde{\beta}^t x_i)$$

↑
square loss

$$\Leftrightarrow \min_{f \in \mathcal{F}} \sum_{i=1}^n l(y_i, f(x_i))$$

↑
 $\mathcal{F} = \{f \mid f(x) = \beta_0 + \tilde{\beta}^t x, \beta_0 \in \mathbb{R}, \tilde{\beta} \in \mathbb{R}^d\}$
= class of linear functions

$$\Leftrightarrow \min_{f \in \mathcal{F}} \hat{R}_n(f)$$

↑ empirical risk

MLE under the linear Gaussian model is equivalent to ERM with a square loss.

Remark: The quantity $\sum_{i=1}^n (y_i - \beta_0 - \tilde{\beta}^t x_i)^2$ is also known as the Residual Sum of Squares in the literature (RSS):

$$\begin{aligned} \text{RSS} &= \sum_{i=1}^n (y_i - \beta_0 - \tilde{\beta}^t x_i)^2 \\ &= (y - X\beta)^t (y - X\beta) \\ &= \|y - X\beta\|^2 \end{aligned}$$

"Least Squares"

↳ Optimization problem is $\hat{\beta} = \underset{\beta \in \mathbb{R}^{d+1}}{\text{argmin}} \|y - X\beta\|^2$

Assumption: $\text{rank } X = d+1$ (4)

↳ In other words, the columns of X are linearly independent

• Minimization of the RSS:

$$\begin{cases} \frac{\partial \text{RSS}}{\partial \beta} = -2y^t X + 2\hat{\beta}^t X^t X = 0 & \text{to find the minimum} \\ \frac{\partial^2 \text{RSS}}{\partial \beta^2} = 2X^t X & \text{to check if the value of } \beta \text{ obtained by} \\ & \text{setting the first derivative of the RSS} \\ & \text{to zero is indeed a minimum.} \end{cases}$$

↑
aka the HESSIAN

Toolbox:

$$\frac{\partial y^t A \beta}{\partial \beta} = y^t A$$

$$\frac{\partial \beta^t A \beta}{\partial \beta} = \beta^t (A + A^t)$$

$y, \beta = \text{vectors}$
 $A = \text{matrix of appropriate dim}$

And indeed,

$$\begin{aligned} \text{RSS} &= (y - X\beta)^t (y - X\beta) \\ &= y^t y - y^t X\beta - \beta^t X^t y + \beta^t X^t X \beta \\ &= y^t y - 2y^t X\beta + \beta^t X^t X \beta \quad (\text{so equal to its transpose}) \\ & \quad \text{symmetric matrix} \end{aligned}$$

$$\hookrightarrow \frac{\partial \text{RSS}}{\partial \beta} = -2y^t X + 2\beta^t X^t X$$

Solution is such that $\hat{\beta}^t X^t X = y^t X$ (5)
 $X^t X \hat{\beta} = X^t y$

Is $X^t X$ invertible?

Under the assumption that $\text{rank } X = d+1$, the answer is yes!

Indeed, suppose that there exists a $\beta \in \mathbb{R}^{d+1}$ such that $(X^t X) \beta = 0$ [i.e. suppose $X^t X$ not invertible]

Then $\beta^t (X^t X) \beta = 0$
 $\|X \beta\|^2$

which implies that $X \beta = 0$; which contradicts the assumption that $\text{rank } X = (d+1)$.

\Rightarrow The matrix $X^t X$ is positive definite.

\Rightarrow The Hessian is positive definite \Rightarrow the solution corresponds to a minimum indeed

\Rightarrow The least squares (LS) estimate is thus

$$\hat{\beta} = (X^t X)^{-1} X^t y$$

We get the estimate $\hat{y} = X \hat{\beta}$,

$$\hat{y} = X (X^t X)^{-1} X^t y =: H y$$

H = the 'hat' matrix: it puts a "hat" on y .

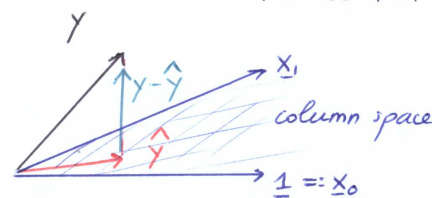
Claim: The least square solution corresponds to the orthogonal projection of $y = (y_1, \dots, y_n)^t$ onto the plane spanned by $(1, x_1, \dots, x_d)$ (6)

\hookrightarrow aka the column space of $X = \begin{bmatrix} 1 & x_1 & \dots & x_d \\ \vdots & \vdots & & \vdots \end{bmatrix}$, commonly denoted $\mathcal{C}(X)$.

$x_j \in \mathbb{R}^n = j$ -th covariate.

Indeed, let's show that $\langle \hat{y} - y, x_j \rangle = 0$

$\langle \cdot, \cdot \rangle =$ inner product in a Euclidean space.
 Recall that for $u, v \in \mathbb{R}^n$, $\langle u, v \rangle = u^t v$.



$$\begin{aligned} \langle \hat{y} - y, x_j \rangle &= (\hat{y} - y)^t x_j \\ &= (X(X^t X)^{-1} X^t y - y)^t x_j \\ &= y^t (X(X^t X)^{-1} X^t - I_n)^t x_j \end{aligned}$$

Next, $\underbrace{X(X^t X)^{-1} X^t}_{n \times (d+1)} \underbrace{x_j}_{(n \times 1)} = \underbrace{[X(X^t X)^{-1} X^t X]}_{n \times n \text{ matrix}} \underbrace{x_j}_{(n \times 1)} (= j\text{-th column})$
 $= [X]_j = x_j$

$\Rightarrow \langle \hat{y} - y, x_j \rangle = 0 \Leftrightarrow \hat{y} - y \perp x_j$
 i.e. The hat matrix H projects y onto the linear subspace spanned by the columns of X . $H =$ projection matrix.

Toolbox: Facts about projection matrices. (7)

- A square matrix P such that $P^2 = P$ [called IDEMPOTENT] is a projection.

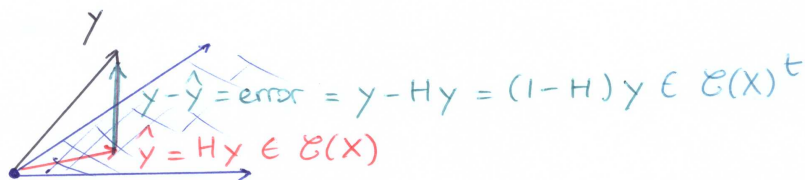
If in addition P is symmetric, then it corresponds to an orthogonal projection.

It is clear from its definition that $H = X(X^t X)^{-1} X^t$ is idempotent and symmetric.

Then $(I - P)$ is also idempotent and symmetric:

$$(I - P)^2 = I - 2P + P^2 = I - 2P + P = I - P$$

$\Rightarrow (I - P)$ is also an orthogonal projection.



$C(X)$ = column space of X ; of dimension $(d+1)$
 $C(X)^t$ = residual space; \perp to $C(X)$

\hookrightarrow We have the decomposition $y = y - \hat{y} + \hat{y}$
 $= \underbrace{(I - H)y}_{\in C(X)^t} + \underbrace{Hy}_{\in C(X)}$
 orthogonal vectors

Fact:

True in general: we write $\mathbb{R}^n = C(X) \oplus C(X)^t$
 $\begin{matrix} \mathbb{R}^n \\ \in \\ y \end{matrix}$
 ie. intersection is $\{0\}$, and vectors from $C(X)$ and $C(X)^t$ are orthogonal

- The rank of an idempotent matrix is equal to its trace: $\text{rank } P = \text{trace } P$. (8)

$$\text{rank } H = d+1 = \text{trace } H = \dim\{C(X)\}$$

$$\text{rank}(I - H) = \text{Tr}(I - H) = \text{Tr } I - \text{Tr } H = n - d - 1 = \dim\{C(X)^t\}$$

- The eigenvalues of a projection matrix are equal to 0 or 1:

Let λ = eigenvalue of P associated with x : $Px = \lambda x$.

Then $P^2 x = P(Px) = \lambda Px = \lambda^2 x$

" $Px = \lambda x$

$\Rightarrow \lambda^2 = \lambda \Rightarrow \lambda = 0$ or 1

end of Toolbox ■

Interpretation of the LS solution.

We have $\hat{z} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_d x_d \in \mathbb{R}^n$

Coefficient $\hat{\beta}_j$ = coordinate of \hat{z} in the base $(1, x_1, \dots, x_d)$ of the column space of X .

OCTOPUS XHO $\hat{\beta}_j$ are not the coordinates of the projection of y onto x_j . This is true only if the base $(1, x_1, \dots, x_d)$ is orthogonal. Look:

- Consider =

$$y = X\beta + \varepsilon, \quad X = \begin{pmatrix} x_0 & \dots & x_d \end{pmatrix}$$

$\begin{matrix} (nx1) & (nxd) & (dx1) & (nx1) \end{matrix}$

The LS fit is $\hat{\beta} = (X^t X)^{-1} X^t y$.

Then

$$\blacktriangleright X^t y = \begin{pmatrix} | & | & | \\ x_0 & \dots & x_d \\ | & | & | \end{pmatrix}^t y = \begin{pmatrix} x_0^t y \\ \vdots \\ x_d^t y \end{pmatrix} = \begin{pmatrix} \langle x_0, y \rangle \\ \vdots \\ \langle x_d, y \rangle \end{pmatrix}$$

$\begin{matrix} (d \times n) & (nx1) & (d \times 1) \end{matrix}$

$$\rightarrow X^t X = \begin{pmatrix} -x_0^t - \\ -x_d^t - \end{pmatrix} \begin{pmatrix} | & & | \\ x_0 & \dots & x_d \\ | & & | \end{pmatrix} = \begin{pmatrix} \langle x_0, x_0 \rangle & \dots & \langle x_0, x_d \rangle \\ \vdots & & \vdots \\ \langle x_d, x_0 \rangle & \dots & \langle x_d, x_d \rangle \end{pmatrix} \quad (9)$$

If the columns of X are orthogonal, then $\langle x_i, x_j \rangle = 0$ for $i \neq j$; and the matrix $X^t X$ is diagonal:

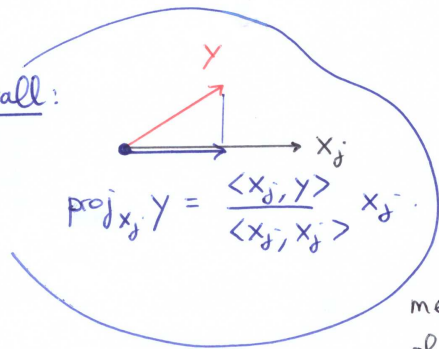
$$(X^t X)^{-1} = \begin{pmatrix} 1/\langle x_0, x_0 \rangle & & 0 \\ & \dots & \\ 0 & & 1/\langle x_d, x_d \rangle \end{pmatrix}$$

We obtain

$$\hat{\beta} = (X^t X)^{-1} X^t y = \begin{pmatrix} \langle x_j, y \rangle \\ \langle x_j, x_j \rangle \\ \vdots \end{pmatrix} \Rightarrow \boxed{\hat{\beta}_j = \frac{\langle x_j, y \rangle}{\langle x_j, x_j \rangle}}$$

$\hat{\beta}_j$ = coordinate of the projection of y onto x_j .
= contribution of input j to output y , independently of other inputs.

Recall:



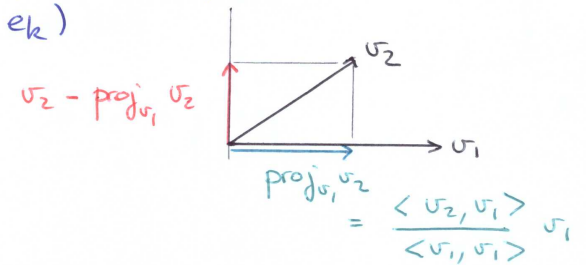
But then, what is the meaning of $\hat{\beta}_j$ when the columns of X are not orthogonal?

→ Use Gram-Schmidt orthogonalization procedure

Toolbox: Consider $\mathcal{M} = \text{span} \{v_1, \dots, v_k\}$, where v_1, \dots, v_k are not orthogonal

What does "span" mean? Well, it is the set of vectors that can be expressed as a linear combination of v_1, \dots, v_k .

The goal is to derive a new basis for \mathcal{M} , which is orthogonal, or, even better, orthonormal. We will call it (e_1, \dots, e_k) (10)



Then

$$\begin{array}{l|l} u_1 = v_1 & e_1 = u_1 / \|u_1\| \\ u_2 = v_2 - \text{proj}_{u_1} v_2 & e_2 = u_2 / \|u_2\| \\ u_3 = v_3 - \text{proj}_{u_1} v_3 - \text{proj}_{u_2} v_3 & e_3 = u_3 / \|u_3\| \\ \vdots & \vdots \\ u_k = v_k - \sum_{j=1}^{k-1} \text{proj}_{u_j} v_k & e_k = u_k / \|u_k\| \quad \blacksquare \end{array}$$

Back to our linear regression problem. To fix ideas, let's consider the case of a single predictor with intercept:

$$\underline{X} = \begin{pmatrix} 1 & x \\ \vdots & \vdots \\ 1 & x \end{pmatrix} \quad (n \times 2) \quad \rightarrow \quad X^t X = \begin{pmatrix} 1^t \\ x^t \end{pmatrix} (1 \ x) = \begin{pmatrix} n & \sum x_i \\ \sum x_i & \langle x, x \rangle \end{pmatrix}$$

$$\rightarrow (X^t X)^{-1} = \frac{1}{n \langle x, x \rangle - (\sum x_i)^2} \begin{pmatrix} \langle x, x \rangle & -\sum x_i \\ -\sum x_i & n \end{pmatrix}$$

$$\rightarrow X^t y = \begin{pmatrix} 1^t \\ x^t \end{pmatrix} y = \begin{pmatrix} \sum y_i \\ \langle x, y \rangle \end{pmatrix}$$

$$\Rightarrow \hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = (X^t X)^{-1} X^t y$$

We get

$$\hat{\beta}_1 = \frac{1}{n\langle x, x \rangle - (\sum x_i)^2} \left\{ -(\sum x_i)(\sum y_i) + n\langle x, y \rangle \right\} \quad (11)$$

$$= \frac{\langle x, y \rangle - \bar{x}(\sum y_i)}{\langle x, x \rangle - n^{-1}(\sum x_i)^2} \quad \rightarrow \quad \bar{x} = \frac{1}{n} \sum x_i$$

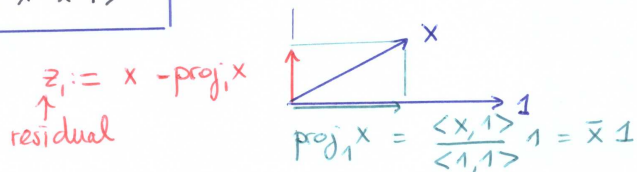
$$= \frac{\langle x, y \rangle - \langle \bar{x} \mathbf{1}, y \rangle}{\langle x, x \rangle - n(\bar{x})^2}$$

Next, since $\langle x - \bar{x}\mathbf{1}, x - \bar{x}\mathbf{1} \rangle = \langle x, x \rangle - \bar{x}\langle \mathbf{1}, x \rangle - \bar{x}\langle x, \mathbf{1} \rangle + (\bar{x})^2\langle \mathbf{1}, \mathbf{1} \rangle$
 $= \langle x, x \rangle - n(\bar{x})^2$,

we obtain

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

= Regress y on residual $z_1 = x - \bar{x}\mathbf{1}$ to get coefficient $\hat{\beta}_1$.



$\hat{\beta}_1$ = added contribution coming from x ("what is new") after x has been adjusted from $\mathbf{1}$.

More generally:

- GRAM-SCHMIDT
- (i) Initialize $z_0 := x_0 := \mathbf{1}$
 - (ii) For $j=1, \dots, d$
 - Regress x_j on z_0, \dots, z_{j-1} to produce coef: $\hat{\delta}_{lj} = \frac{\langle x_j, z_l \rangle}{\langle z_l, z_l \rangle}$, $l=0, \dots, j-1$
 - and $z_j = x_j - \sum_{l=0}^{j-1} \hat{\delta}_{lj} z_l =$ residual vector
 - (iii) Regress y on residual z_d to get $\hat{\beta}_d$

Remark: the z_j do not have unit length in general. (12)

In words = the regression coefficient $\hat{\beta}_j$ represents the additional contribution of x_j on y , after x_j has been adjusted for $x_0, \dots, x_{j-1}, x_{j+1}, \dots, x_d$ (what is left in x_j is non-informative since redundant with the information contained in the remaining covariates)

\Rightarrow Linear model is highly interpretable.

- The matrix form of the Gram-Schmidt orthogonalization procedure yields the QR decomposition of X :

$$X = Z \Gamma$$

$n \times (d+1)$ $n \times (d+1)$ $(d+1) \times (d+1)$

- Γ is upper triangular: rewrite step (ii) in the GS procedure page 11 as

$$x_j = z_j + \sum_{l=0}^{j-1} \hat{\delta}_{lj} z_l$$

- Z has orthogonal columns & the column space spanned by the first l columns of X and Z are the same.

$$\Gamma = \begin{pmatrix} 1 & \hat{\delta}_{01} & \hat{\delta}_{02} & \hat{\delta}_{03} & \dots \\ & 1 & \hat{\delta}_{12} & \hat{\delta}_{13} & \dots \\ & & 1 & \hat{\delta}_{23} & \dots \\ & & & 1 & \dots \\ & & & & \ddots \end{pmatrix}$$

$$Z = \begin{pmatrix} | & | & & | \\ z_0 & z_1 & \dots & z_d \\ | & | & & | \end{pmatrix}$$

- Put $D = \begin{pmatrix} \sqrt{\langle z_0, z_0 \rangle} & & & \\ & \ddots & & \\ & & \sqrt{\langle z_d, z_d \rangle} & \\ & & & \ddots \end{pmatrix} =$ diagonal matrix

Put $Q = Z D^{-1}$ $R = D \Gamma$ R is upper triangular

$n \times (d+1)$ $n \times (d+1)$ $(d+1) \times (d+1)$ $(d+1) \times (d+1)$

Why?

Because $Q^t Q = D^{-1} Z^t Z D^{-1} = D^{-1} \begin{pmatrix} \langle z_0, z_0 \rangle & & 0 \\ & \ddots & \\ 0 & & \langle z_d, z_d \rangle \end{pmatrix} D^{-1} = I_{d+1}$

$\Rightarrow Q$ has orthogonal columns that form a basis for the column space of X .

We write $X = QR$ "the QR decomposition of X ".

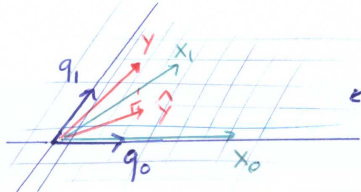
$n \times (d+1)$ $n \times (d+1)$ $(d+1) \times (d+1)$

This decomposition is very useful. Let $Q = (q_0, q_1, \dots, q_d)$. The LS solution is the projection of y onto the column space of X . The solution can be represented in an orthonormal basis (q_0, \dots, q_d) for $\mathcal{C}(X)$ since

$\rightarrow (X^t X)^{-1} X^t y = R^{-1} (R^t)^{-1} R^t Q^t y = R^{-1} Q^t y$

$X^t X = R^t \underbrace{Q^t Q}_I R = R^t R$

$\rightarrow \hat{y} = X (X^t X)^{-1} X^t y = (QR) (R^{-1} Q^t y) = Q Q^t y = \sum_{j=0}^d \langle y, q_j \rangle q_j$



$\leftarrow x_0$ and x_1 span the same plane as q_0 and q_1 , this should be clear.

Remark: Alternatively, you may use an SVD decomposition (14)

of $X = U \Lambda V^t$, rank $X = d+1$, where the columns of U provide an orthonormal basis for $\mathcal{C}(X)$. The LS solution can then be expressed as $\hat{y} = \sum \langle y, u_j \rangle u_j$, where $u_j = j$ -th column of U .

Note that the orthonormal basis provided by U and Q are not necessarily the same ones.

Further geometrical considerations.

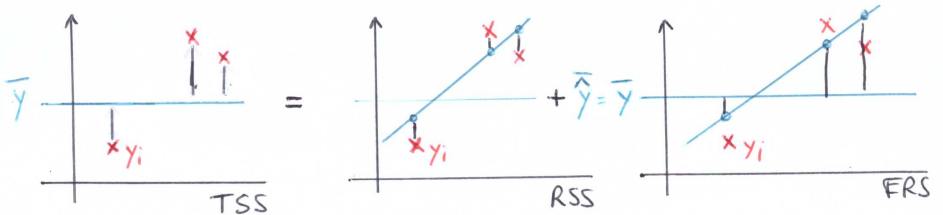
If the intercept is included in the model, we have the decomposition:

$\sum (y_i - \bar{y})^2 = \sum (y_i - \hat{y}_i)^2 + \sum (\hat{y}_i - \bar{y})^2$

Total Sum of Squares (TSS) = Residual Sum of Squares (RSS) + Explained Sum of Squares (ESS)

$\|y - \bar{y}1\|^2 = \|y - \hat{y}\|^2 + \|\hat{y} - \bar{y}1\|^2$

(where $\bar{y} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i$)



Note that if the intercept is in the model, then $\bar{y} = \bar{\hat{y}}$.

Indeed, making use of the first equation in $(X^t X) \hat{\beta} = X^t y$,

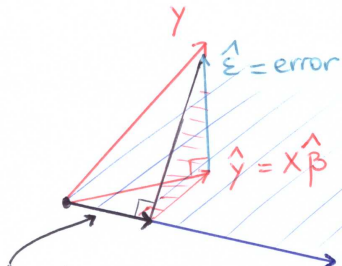
rewritten $X^t(y - X\hat{\beta}) = 0$, we obtain

$$\sum_{i=1}^n y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_d x_{id}) = 0$$

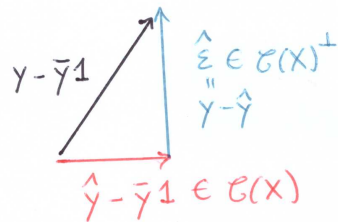
$$\Leftrightarrow \sum_{i=1}^n y_i - \hat{y}_i = 0$$

$$\Leftrightarrow \bar{y} = \bar{\hat{y}} \quad \blacksquare$$

The decomposition $TSS = ESS + RSS$ can be derived by direct computation, or can be visualized geometrically as follows:



zoom on the triangle with green/black/red sides:



Pythagoras theorem yields indeed

$$\|y - \bar{y}1\|^2 = \|y - \hat{y}\|^2 + \|\hat{y} - \bar{y}1\|^2$$

$$\text{proj}_1 \hat{y} = \frac{\langle \hat{y}, 1 \rangle}{\langle 1, 1 \rangle} 1 = \bar{\hat{y}} 1$$

$$\text{proj}_1 y = \frac{\langle y, 1 \rangle}{\langle 1, 1 \rangle} 1 = \bar{y} 1$$

Since $\bar{\hat{y}} = \bar{y}$, we conclude that $\text{proj}_1 \hat{y} = \text{proj}_1 y$.

\Rightarrow Introduce $R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}$

Interpretation: R^2 represents the proportion of variability in y that is explained using the linear model.

\approx measure of linear association between X and Y . No wonder that for simple linear regression ($d=1$), we have $R^2 = \hat{\rho}^2$,

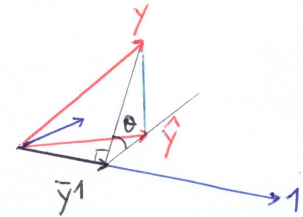
where $\hat{\rho}$ denotes the empirical correlation coefficient between the predictor and the response variable:

$$\hat{\rho} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

Geometrically, $R^2 = \cos^2 \theta = \frac{\|\hat{y} - \bar{y}1\|^2}{\|y - \bar{y}1\|^2} = \frac{ESS}{TSS}$;

$\theta =$ angle between y and \hat{y} , taken at $\bar{y}1$.

\hookrightarrow When $\theta = 0$, $R^2 = 1$, and all the variance is explained: $y \in C(X)$; so that y is exactly a linear combination of the covariates



Remarks: (i) R^2 does not necessarily lie between 0 and 1 if you do not include the constant in the model.

(ii) R^2 is not comparable across different response variables y s

Eg: In one model consider y , and in another $\ln y$. Since the scale is changed, the sum of squares as well, and R^2 is affected.

(iii) R^2 almost always increases when you add a new predictor (proof below)

Consequence: you might be tempted to add more and more variables into the model.

Alternative: Use the adjusted R^2 square,

$$R_a^2 = \frac{\|\hat{\epsilon}\|^2 / (n-d-1)}{\|y - \bar{y}1\|^2 / (n-1)} = 1 - \frac{n-1}{n-d-1} \frac{RSS}{TSS}$$

→ As the number of features increases, the RSS decreases, but the adjustment term $n^{-1}/(n-d-1)$ increases. (17)

• Proof that R^2 (almost) always increases when adding a new predictor.

Step I: Partitioned Regression.

Suppose that the regression model can be written

$$\begin{aligned} \rightarrow y &= X\beta + \varepsilon \\ &= X_1\beta_1 + X_2\beta_2 + \varepsilon \\ &= (X_1 \mid X_2) \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} + \varepsilon \end{aligned}$$

sets of variables ; $\beta_1/\beta_2 = \text{vectors}$

→ LS solution satisfies $(X^t X)\beta = X^t y$, with

$$X^t X = \begin{pmatrix} X_1^t X_1 & X_1^t X_2 \\ X_2^t X_1 & X_2^t X_2 \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} X_1^t X_1 & X_1^t X_2 \\ X_2^t X_1 & X_2^t X_2 \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} X_1^t y \\ X_2^t y \end{pmatrix} \quad \text{--- (A)}$$

→ Solve for $\hat{\beta}_1$: $(X_1^t X_1)\hat{\beta}_1 + (X_1^t X_2)\hat{\beta}_2 = X_1^t y$

$$\hat{\beta}_1 = (X_1^t X_1)^{-1} X_1^t y - (X_1^t X_1)^{-1} X_1^t X_2 \hat{\beta}_2$$

$$\hat{\beta}_1 = (X_1^t X_1)^{-1} X_1^t (y - X_2 \hat{\beta}_2) \quad \text{--- (B)}$$

Note that if $X_1 \perp X_2$, $X_1^t X_2 = 0$ and the estimated coefficient $\hat{\beta}_1$ remains the same, whether the variable X_2 is included in the model or not. This is not true if X_1 and X_2 are not orthogonal.

→ Substitute expression (B) for $\hat{\beta}_1$ back into the second equation in (A): (18)

$$\begin{aligned} (X_2^t X_1)\hat{\beta}_1 + (X_2^t X_2)\hat{\beta}_2 &= X_2^t y \\ (X_2^t X_1) \{ (X_1^t X_1)^{-1} X_1^t y - (X_1^t X_1)^{-1} X_1^t X_2 \hat{\beta}_2 \} + (X_2^t X_2)\hat{\beta}_2 &= X_2^t y \end{aligned}$$

$$\begin{aligned} \{ X_2^t X_2 - X_2^t X_1 (X_1^t X_1)^{-1} X_1^t X_2 \} \hat{\beta}_2 &= X_2^t y - (X_2^t X_1) (X_1^t X_1)^{-1} X_1^t y \\ &= X_2^t y - (X_2^t X_1) (X_1^t X_1)^{-1} X_1^t y \end{aligned}$$

$$X_2^t \{ I - \underbrace{X_1 (X_1^t X_1)^{-1} X_1^t}_{H_1} \} X_2 \hat{\beta}_2 = X_2^t \{ I - \underbrace{X_1 (X_1^t X_1)^{-1} X_1^t}_{H_1} \} y$$

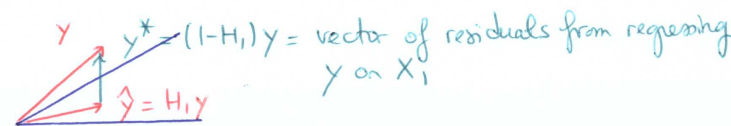
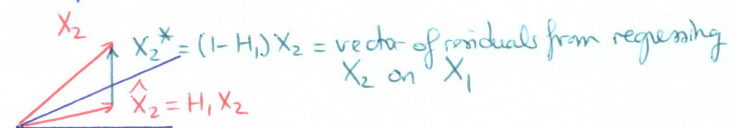
projection matrix onto $\mathcal{C}(X_1)$.

→ $\hat{\beta}_2 = (X_2^t (I - H_1) X_2)^{-1} (X_2^t (I - H_1) y)$
Since $(I - H_1)$ is idempotent, introducing variables

$$\begin{cases} X_2^* = (I - H_1) X_2 \\ y^* = (I - H_1) y \end{cases}$$

yields $\hat{\beta}_2 = (X_2^{*t} X_2^*)^{-1} X_2^{*t} y^* \quad \text{--- (*)}$

→ Interpretation:



⇒ To estimate $\hat{\beta}_2$, first regress X_2 and y on X_1 , and then regress residuals y^* on residuals X_2^* . (19)

Step II: R^2 and the addition of a new variable.

→ First, start with the model $y = X\beta + \varepsilon$, and the fit $y = X\hat{\beta} + \hat{\varepsilon}$
↑ vector of residuals $\hat{\varepsilon} = (I-H)y$
↑ Least square fit $(X^tX)^{-1}X^ty$

Next, add a single extra variable z , so that

$$y = X\hat{\beta}' + \hat{\gamma}z + \hat{\varepsilon}'$$

↑ LS coefficients
↑ new vector of residuals

Unless $z \perp X$, $\hat{\beta}$ and $\hat{\beta}'$ differ; see the remark bottom of page 17.

→ The RSS in the original model is $\|\hat{\varepsilon}\|^2$, and in the second model $\|\hat{\varepsilon}'\|^2 \Rightarrow$ the R^2 will increase if $\|\hat{\varepsilon}'\|^2 < \|\hat{\varepsilon}\|^2$.

→ LS solution $(\hat{\beta}', \hat{\gamma})$ satisfies

$$\begin{pmatrix} X^tX & X^tz \\ z^tX & z^tz \end{pmatrix} \begin{pmatrix} \hat{\beta}' \\ \hat{\gamma} \end{pmatrix} = \begin{pmatrix} X^ty \\ z^ty \end{pmatrix}$$

↳ Proceeding as in Step I, we see that the first line yield the estimate

$$\hat{\beta}' = (X^tX)^{-1}X^t(y - \hat{\gamma}z)$$

(compare with expression (B) page 17)

$$\hat{\beta}' = \hat{\beta} - \hat{\gamma}(X^tX)^{-1}X^tz \quad (C)$$

→ Moreover, $\hat{\varepsilon}' = y - X\hat{\beta}' - \hat{\gamma}z$
↑ Insert expression (C) here.

We get $\hat{\varepsilon}' = y - X\hat{\beta} + X(X^tX)^{-1}X^t(\hat{\gamma}z) - (\hat{\gamma}z)$
 $= y - X\hat{\beta} - (I - \underbrace{X(X^tX)^{-1}X^t}_{=H})(\hat{\gamma}z)$
 $= \hat{\varepsilon} - (I-H)(\hat{\gamma}z)$

Put $z^* = (I-H)z$
 $\hat{\varepsilon}' = \hat{\varepsilon} - \hat{\gamma}z^*$

→ From the last written expression, we conclude that
 $\|\hat{\varepsilon}'\|^2 = \|\hat{\varepsilon}\|^2 + \hat{\gamma}^2\|z^*\|^2 - 2\hat{\gamma}\hat{\varepsilon}^tz^* \quad (D)$

We turn our attention to this term

→ $\hat{\varepsilon} = (I-H)y = y^*$ = residual from regressing y on X
 $z^* = (I-H)z$ = residual from regressing z on X

⇒ Applying formula (*) gives

$$\hat{\gamma} = (z^{*t}z^*)^{-1}(z^{*t}y^*)$$

(with $X_1 = X$ and $X_2 = z$, in the notation of page 18)

⇒ It follows that $(z^{*t}z^*)\hat{\gamma} = z^{*t}y^*$
 $= z^{*t}\hat{\varepsilon}$
 $= \hat{\varepsilon}^tz^*$

↑
 Plug this expression of $\hat{\varepsilon}^tz^*$ back into (D):

$$\|\hat{\varepsilon}'\|^2 = \|\hat{\varepsilon}\|^2 + \hat{\gamma}^2\|z^*\|^2 - 2\hat{\gamma}^2\|z^*\|^2$$

$$\|\hat{\varepsilon}'\|^2 = \|\hat{\varepsilon}\|^2 - \hat{\gamma}^2\|z^*\|^2$$

⇒ $\|\hat{\varepsilon}'\|^2 < \|\hat{\varepsilon}\|^2$, unless $\hat{\gamma} = 0$

⇒ Unless the LS on the additional variable is zero, adding an extra variable will decrease R^2 .

Properties of the LS solution.

(21)

(1) $\hat{\beta}$ is unbiased for β .

Indeed, since $\hat{\beta} = (X^t X)^{-1} X^t y$, we have

$$\mathbb{E} \hat{\beta} = (X^t X)^{-1} X^t \mathbb{E} y = \beta, \text{ as required}$$

$= X\beta$ since $y = X\beta + \varepsilon$.

Remarks (i) Unbiasedness of $\hat{\beta}$ holds independently of the distribution of ε , provided $\mathbb{E}(\varepsilon | X) = 0$ and $\mathbb{E}(\varepsilon \varepsilon^t | X) = \sigma^2 I$.
 \Rightarrow Normal distribution is not needed.

(ii) In a random design (X random), $\hat{\beta}$ is still unbiased since

$$\mathbb{E} \hat{\beta} = \mathbb{E}_X \mathbb{E} \hat{\beta} | X = \mathbb{E}_X \beta = \beta.$$

(2) The covariance matrix of $\hat{\beta}$ is $\Sigma_{\hat{\beta}} = \mathbb{E}(\hat{\beta} - \beta)(\hat{\beta} - \beta)^t$, with

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

so that

$$\Sigma_{\hat{\beta}} = (X^t X)^{-1} X^t \underbrace{\mathbb{E}\{(y - X\beta)(y - X\beta)^t\}}_{= \sigma^2 I_n} X (X^t X)^{-1}$$

$$\Sigma_{\hat{\beta}} = \sigma^2 (X^t X)^{-1}$$

Remark: In a random design, use

$$\Sigma_{\hat{\beta}} = \mathbb{E}_X (\Sigma_{\hat{\beta} | X}) + \underbrace{\Sigma_{\mathbb{E}(\hat{\beta} | X)}}_{= 0} = \sigma^2 \mathbb{E} (X^t X)^{-1}.$$

(3) $\hat{\sigma}^2 := \frac{1}{n-d-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ is unbiased for σ^2 .

(22)

Preliminary calculations:

$$\begin{aligned} \sum (y_i - \hat{y}_i)^2 &= \|y - \hat{y}\|^2 \\ &= \|y - X\hat{\beta}\|^2 \\ &= \|\hat{\varepsilon}\|^2 \end{aligned}$$

where we defined $\hat{\varepsilon} = y - X\hat{\beta} = (I - H)y$
 $=$ vector of residuals

$\blacktriangleright \mathbb{E} \hat{\varepsilon} = \mathbb{E}(I - H)y = (I - H)\mathbb{E}y = 0$ since $\mathbb{E}y \in \mathcal{C}(X)$
 \blacktriangleright Since $y = X\beta + \varepsilon \Rightarrow (I - H)y = (I - H)\varepsilon$
 $\in \mathcal{C}(X)$ projects onto $\mathcal{C}(X)^\perp$.

Thus

$$\begin{aligned} \Sigma_{\hat{\varepsilon}} &= \mathbb{E}(\hat{\varepsilon} \hat{\varepsilon}^t) \text{ since } \mathbb{E} \hat{\varepsilon} = 0 \\ &= \mathbb{E}[(I - H)y][y^t (I - H)] \\ &= \mathbb{E}[(I - H)\varepsilon][\varepsilon^t (I - H)] \text{ since } (I - H)y = (I - H)\varepsilon \\ &= (I - H) \underbrace{\mathbb{E}(\varepsilon \varepsilon^t)}_{= \sigma^2 I_n} (I - H) \\ &= \sigma^2 (I - H) \text{ since } (I - H) \text{ is idempotent.} \end{aligned}$$

Back to $\mathbb{E} \sum (y_i - \hat{y}_i)^2 = \mathbb{E} \|\hat{\varepsilon}\|^2$ since a scalar is equal to its trace

Since \forall matrices

$$\text{Tr}(AA^t) = \text{Tr}(A^t A)$$

$$\begin{aligned} &= \mathbb{E} \text{Tr}(\hat{\varepsilon}^t \hat{\varepsilon}) \\ &= \mathbb{E} \text{Tr}(\hat{\varepsilon} \hat{\varepsilon}^t) \\ &= \text{Tr} \mathbb{E}(\hat{\varepsilon} \hat{\varepsilon}^t) \\ &= \text{Tr} \Sigma_{\hat{\varepsilon}} \\ &= \sigma^2 \text{Tr}(I - H) \\ &= \sigma^2 \text{Rk}(I - H) = \sigma^2 (n - d - 1) \end{aligned}$$

Remarks: (i) $\hat{\beta} = (X^t X)^{-1} X^t y$ is called a LINEAR estimator, in the sense that it is computed as a linear combination of the entries in y . Any estimator of the form By is called a linear estimator (of β). $\hat{\beta}$ satisfies an interesting property:

Theorem: (GAUSS-MARKOV THEOREM)
 $\hat{\beta}$ has minimum variance amongst all unbiased linear estimators of β

OK, $\hat{\beta}$ is a vector; so what do we mean by minimum variance? We mean that if $\tilde{\beta}$ has covariance matrix $\Sigma_{\tilde{\beta}}$ such that $E\tilde{\beta} = \beta$; with covariance matrix $\Sigma_{\tilde{\beta}}$,

Then $\Sigma_{\tilde{\beta}} \succeq \Sigma_{\hat{\beta}}$; in the sense that $\Sigma_{\tilde{\beta}} - \Sigma_{\hat{\beta}}$ is a positive semi-definite matrix.

(ii) Since $\hat{\beta}$ is a linear estimator and y (conditionally on $X=x$) is normally distributed, it follows that $\hat{\beta}$ is normally distributed and $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2(X^t X)^{-1})$

(iii) Unbiased estimator $\hat{\sigma}^2$ of σ^2 defined on page 15 does not correspond to the maximum likelihood estimator of σ^2 . \rightarrow biased estimator

MLE is $s^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$
 (log lik. given on page 3)

(iv) What about the distribution of $\hat{\sigma}^2$? (24)

To answer this question, we make use of COCHRAN theorem; which stipulates that the decomposition of a multivariate Gaussian vector into orthogonal subspaces produces independent random variables, whose laws are explicit.

Theorem = (COCHRAN)

Let $Z \sim \mathcal{N}(\mu, \sigma^2 I_n)$

- \mathcal{M} = a subspace of \mathbb{R}^n of dimension p
- P = matrix of orthogonal projection on \mathcal{M}
- P_{\perp} = matrix of orthogonal projection on \mathcal{M}^{\perp} , where \mathcal{M}^{\perp} = orthogonal complement of \mathcal{M}

Then

- (a) $PZ \sim \mathcal{N}(P\mu, \sigma^2 P)$ & $P_{\perp}Z \sim \mathcal{N}(P_{\perp}\mu, \sigma^2 P_{\perp})$
 (b) Vectors PZ and $P_{\perp}Z = (I-P)Z$ are indpt.
 (c) $\frac{1}{\sigma^2} \|P(Z-\mu)\|^2 \sim \chi_p^2$
 and $\frac{1}{\sigma^2} \|P_{\perp}(Z-\mu)\|^2 \sim \chi_{n-p}^2$

Consequences:

$\hat{\beta}$ and $\hat{\sigma}^2$ are independent since applying (b):

$$\begin{aligned} \hat{\beta} &= (X^t X)^{-1} X^t y \\ &= (X^t X)^{-1} X^t (X(X^t X)^{-1} X^t) y \\ &= (X^t X)^{-1} X^t H y \\ &= \text{function of } H y \end{aligned}$$

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{n-d-1} \|\hat{\epsilon}\|^2 = \frac{1}{n-d-1} \|(I-H)y\|^2 \\ &= \text{function of } (I-H)y \end{aligned}$$

For $\frac{(n-d-1)\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-d-1}$ since, in the notation of the theorem,

with $Z = \varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$, we have that

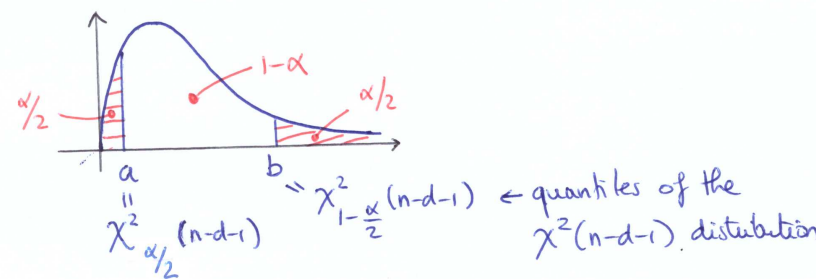
$$\hat{\sigma}^2 = \frac{\|\hat{\varepsilon}\|^2}{n-d-1} = \frac{1}{n-d-1} \|(I-H)\varepsilon\|^2 = \frac{1}{n-d-1} \|H_{\perp}(\varepsilon - E\varepsilon)\|^2 \quad \left. \vphantom{\hat{\sigma}^2} \right\} E\varepsilon=0$$

so that $\frac{(n-d-1)\hat{\sigma}^2}{\sigma^2} = \frac{1}{\sigma^2} \|H_{\perp}(\varepsilon - E\varepsilon)\|^2 \sim \chi^2_{n-d-1}$
 since $H_{\perp} = I-H$ projects onto a subspace of dimension $n-d-1$.

• Since the distributions of $\hat{\beta}$ and $\hat{\sigma}^2$ are known, we can make good use of them to construct confidence intervals:

For σ , we have that $P\left(a \leq \frac{(n-d-1)\hat{\sigma}^2}{\sigma^2} \leq b\right) = 1-\alpha$,

where



Re-arranging terms yields $\left[\frac{(n-d-1)\hat{\sigma}^2}{b}, \frac{(n-d-1)\hat{\sigma}^2}{a} \right]$
 a $100(1-\alpha)\%$ confidence interval for σ^2 .

For β , we have the choice to construct $(d+1)$ separate confidence intervals; each with nominal coverage $(1-\alpha)$; or construct a global confidence interval, so that all components are taken into account simultaneously (\equiv aka confidence regions)

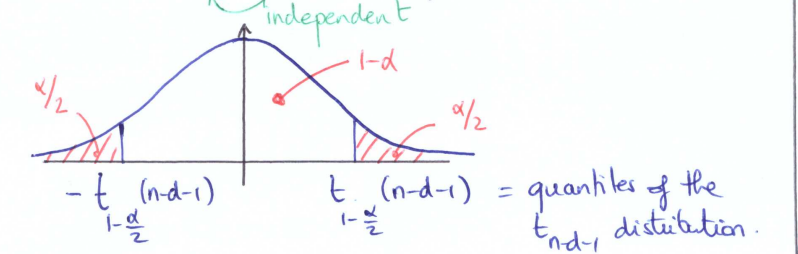
Put $v_j = [(X^T X)^{-1}]_{jj}$
 $= j$ -th diagonal element of the square matrix $(X^T X)^{-1} \rightarrow \text{Var } \hat{\beta}_j = \sigma^2 v_j$.

Then $\frac{\hat{\beta}_j - \beta_j}{\sqrt{\sigma^2 v_j}} \sim \mathcal{N}(0,1)$

$$\frac{\frac{\hat{\beta}_j - \beta_j}{\sqrt{\sigma^2 v_j}}}{\frac{\sqrt{\frac{(n-d-1)\hat{\sigma}^2}{\sigma^2}}}{(n-d-1)}} = \frac{\hat{\beta}_j - \beta_j}{\hat{\sigma} \sqrt{v_j}} \sim t_{n-d-1}$$

 Student-t distribution with $(n-d-1)$ degrees of freedom.

(This follows directly from the definition of the Student-t distribution: $T = \frac{Z}{\sqrt{S/n}} \sim t_n$, where $Z \sim \mathcal{N}(0,1)$ and $S \sim \chi^2_n$)



$\Rightarrow \left[\hat{\beta}_j - t_{\frac{1-\alpha}{2}}(n-d-1) \hat{\sigma} \sqrt{v_j}, \hat{\beta}_j + t_{\frac{1-\alpha}{2}}(n-d-1) \hat{\sigma} \sqrt{v_j} \right]$
 is a $100(1-\alpha)\%$ confidence interval for β_j .

Alternatively, you may want to test for $H_0: \beta_j = 0$ (27)
 ("absence of effect" = typically what you put inside a null hypothesis)

$$z_j := \hat{\beta}_j / \hat{\sigma} \sqrt{v_j} \quad \text{"z-score"}$$

\Rightarrow A large absolute value of z_j leads to the rejection of the null hypothesis; equivalently, a small p-value indicates the absence of statistical evidence that there is no association between the corresponding predictor X_j and the response.

(we never accept the null, we fail to reject the null)

OCTOPUS * H0 If we have a large number of predictors, say 100, then testing for the 100 coefficients individually not only is a waste of time, but also is inappropriate. Why? Because of statistical regularity: at the 95% nominal level for example, and assuming that all $\beta_j \neq 0$, 5% of the z-scores will appear in the distribution tails just by chance, leading to the rejection of the null.

\rightarrow Use a global test aka FISHER test.

In other words, test for several coefficients simultaneously:

Consider two candidate models, one being the reduced version of the other one, by removing some of the features of the larger model (we say that the models are nested)

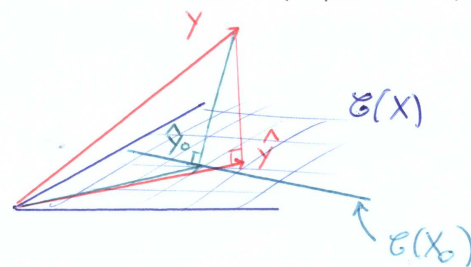
\Leftrightarrow a subset of the β_j are equal to zero.

- Full model: $(d_1 + 1)$ parameters (28)
- Reduced model: $(d_0 + 1)$ parameters, $d_0 \leq d_1$

\hookrightarrow Corresponds to $y = X_0 \beta_0 + \varepsilon$, $\text{rank } X_0 = (d_0 + 1)$
 $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$

where the columns of X_0 are a subset of the columns of X .

If the reduced model is correct, then \hat{y}_0 = orthogonal projection of y onto the column space of X_0 , is "close" to \hat{y} = \perp projection of y onto $\mathcal{C}(X)$.



\Rightarrow Check the size of $\|\hat{y}_0 - \hat{y}\|^2$ = square distance between \hat{y}_0 and \hat{y} .

It turns out that we are in a better position if we remove the scale effect on the output variable (square distance depends on the measurement unit, right?)

$$\Rightarrow \text{Consider } \frac{\|\hat{y}_0 - \hat{y}\|^2}{\|y - \hat{y}\|^2} \left(= \frac{\|(\hat{y}_0 - y) + (y - \hat{y})\|^2}{\|y - \hat{y}\|^2} \right) = \text{relative error (see p. 30)}$$

In fact, $(\hat{y}_0 - \hat{y})$ and $(y - \hat{y})$ live in subspaces of different dimension, and we need to readjust for the dimension to be able to derive the distribution of the relative error:

$$F := \frac{\|\hat{y}_0 - \hat{y}\|^2 / (d_1 - d_0)}{\|y - \hat{y}\|^2 / (n - d_1 - 1)}$$

Facts: the numerator and denominator are independent, (29)

and $\frac{1}{\sigma^2} \|y - \hat{y}\|^2 \sim \chi^2_{n-d_1-1}$ &

$\frac{1}{\sigma^2} \|\hat{y} - \hat{y}_0\|^2 \sim \chi^2_{d_1-d_0}$

under the null hypothesis that the reduced model is correct.

⇒ Under H_0 : $F \sim F_{d_1-d_0, n-d_1-1}$

"large" value, reject, bla bla bla...

F-distribution with parameters (d_1-d_0) and $(n-d_1-1)$.

Remark: Why are the num & den independent, intuitively? Because they live in orthogonal subspaces.

Remarks: (i) An alternative notation for F is widely used:

$$F = \frac{(RSS_0 - RSS_1) / (d_1 - d_0)}{RSS_1 / (n - d_1 - 1)}, \text{ where}$$

↳ RSS_1 = Residual Sum of Squares of the larger model with (d_1+1) parameters

↳ RSS_0 = Residual Sum of Squares of the reduced model with (d_0+1) parameters, $(d_0 < d_1)$

Indeed,

$$\|y - \hat{y}_0\|^2 = \|y - \hat{y} + \hat{y} - \hat{y}_0\|^2 = \|y - Hy + Hy - H_0 y\|^2$$

where H_0 = projection matrix onto $\mathcal{C}(X_0)$

$$\|y - \hat{y}_0\|^2 = \|(1-H)y + Hy - H_0 Hy\|^2$$

since $H_0 y = H_0 Hy$ as $\mathcal{C}(X_0) \subset \mathcal{C}(X)$.

$$= \|(1-H)y + (1-H_0)Hy\|^2 = \|\underbrace{(1-H)y}_{\in \mathcal{C}(X)^\perp} + \underbrace{(1-H_0)Hy}_{\in \mathcal{C}(X) \cap \mathcal{C}(X_0)^\perp}\|^2$$

perpendicular subspaces ⇒ Pythagora

$$= \|(1-H)y\|^2 + \|(1-H_0)Hy\|^2 = \|y - \hat{y}\|^2 + \|\hat{y} - \hat{y}_0\|^2$$

Thus

$$\|\hat{y} - \hat{y}_0\|^2 = \|y - \hat{y}_0\|^2 - \|y - \hat{y}\|^2 = RSS_0 - RSS_1$$

(ii) If the simpler model drops only one coefficient, the F-test is testing for the nullity of this coefficient. But this is exactly what the Student-t test was designed for! So, which test shall we be using? It turns out that the two tests are equivalent in this case, and it is possible to show that $F = z_j^2$ (the z-score associated with the corresponding coefficient)

(iii) What about confidence regions for β ?

Since $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (X^t X)^{-1})$, you may want to use $(\hat{\beta} - \beta)^t \Sigma_{\hat{\beta}}^{-1} (\hat{\beta} - \beta) \sim \chi^2_{d_H}$

but unfortunately $\Sigma_{\hat{\beta}} = \sigma^2 (X^t X)^{-1}$, with σ unknown. Big deal: replace σ^2 by $\hat{\sigma}^2$: (31)

→ Let R be an $r \times (d+1)$ matrix of rank r

Typically, we will take R of the form

$$R = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & 0 & \dots & 0 \end{pmatrix} \begin{matrix} \uparrow \\ \uparrow \\ \uparrow \\ \vdots \\ \downarrow \end{matrix}$$

← (d+1) →

Role of R = select a subset of the original coefficients $\beta_0, \beta_1, \dots, \beta_d$.

Consider $R(\hat{\beta} - \beta)$.
We know that $R\hat{\beta} \sim \mathcal{N}(R\beta, \sigma^2 R(X^t X)^{-1} R^t)$

$$\Rightarrow \frac{1}{\hat{\sigma}^2} [R(\hat{\beta} - \beta)]^t [R(X^t X)^{-1} R^t]^{-1} [R(\hat{\beta} - \beta)] \sim \chi_r^2$$

→ Also, $\frac{(n-d-1)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-d-1}^2$

→ Finally, $\hat{\sigma}^2$ and $\hat{\beta}$ are independent (consequence of Cochran theorem)

Thus:

$$\frac{\frac{1}{\hat{\sigma}^2} [R(\hat{\beta} - \beta)]^t [R(X^t X)^{-1} R^t]^{-1} [R(\hat{\beta} - \beta)] / r}{\frac{(n-d-1)\hat{\sigma}^2}{\sigma^2} / (n-d-1)} \sim F_{r, n-d-1}$$

Simplifying terms yields: (32)

$$\frac{1}{r\hat{\sigma}^2} [R(\hat{\beta} - \beta)]^t [R(X^t X)^{-1} R^t]^{-1} [R(\hat{\beta} - \beta)] \sim F_{r, n-d-1}$$

From this, we can construct a confidence region for simultaneously r parameters of the model, $\beta_{j_1}, \dots, \beta_{j_r}$:

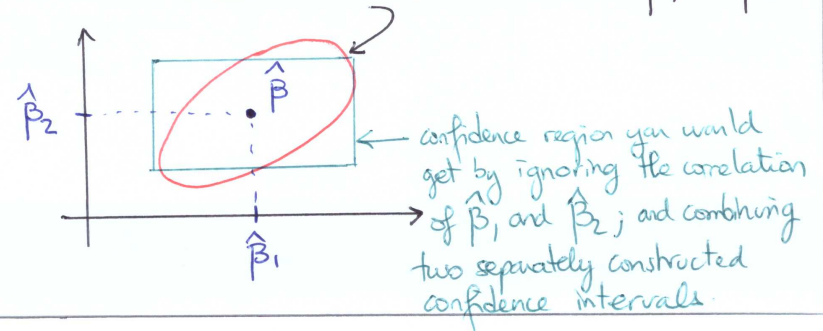
$$\left\{ (\beta_{j_1}, \dots, \beta_{j_r}) \mid \frac{1}{r\hat{\sigma}^2} [R(\hat{\beta} - \beta)]^t [R(X^t X)^{-1} R^t]^{-1} [R(\hat{\beta} - \beta)] \leq f_{1-\alpha}^*(r, n-d-1) \right\}$$

↑
Confidence region, with nominal coverage $(1-\alpha)$.
↑
 $f_{1-\alpha}^*(r, n-d-1)$ is the $(1-\alpha)$ quantile of the $F_{r, n-d-1}$ distribution

Example: taking $r=2$ and $R = \begin{bmatrix} 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \end{bmatrix}$,
 $R(\hat{\beta} - \beta) = \begin{pmatrix} \hat{\beta}_0 - \beta_0 \\ \hat{\beta}_1 - \beta_1 \end{pmatrix} \equiv$ keep one coefficient + intercept

$$\left\{ (\beta_0, \beta_1) \mid \frac{\sigma_{22}(\hat{\beta}_0 - \beta_0)^2 - 2\sigma_{12}(\hat{\beta}_0 - \beta_0)(\hat{\beta}_1 - \beta_1) + \sigma_{11}(\hat{\beta}_1 - \beta_1)^2}{2\hat{\sigma}^2(\sigma_{11}\sigma_{22} - \sigma_{12}^2)} \leq f_{1-\alpha}^*(2, n-d-1) \right\}$$

|| ellipsoid = confidence region which takes into account the correlation between $\hat{\beta}_1$ and $\hat{\beta}_2$.



II. MODEL ASSESSMENT

(33)

The linear model $y = X\beta + \varepsilon$ assumes that $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$, and $\text{rank } X = (d+1)$. Once we have an estimate of β , it is necessary to check these assumptions. We study:

- Residuals $\hat{\varepsilon} = y - \hat{y} = y - X\hat{\beta}$, and see how we can use them to check that $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ is legitimate.
- Projection Matrix $H = X(X^T X)^{-1} X^T$
- Matrix of Observation X and tools to detect colinearity.

II.1. Study of the residuals.

We derived page 22 that

$$\begin{aligned} E \hat{\varepsilon} &= 0 \\ \Sigma_{\hat{\varepsilon}} &= \sigma^2 (I - H) \end{aligned}$$

RESIDUALS

these should be compared with the model assumptions

$$\begin{aligned} E \varepsilon &= 0 \\ \Sigma_{\varepsilon} &= \sigma^2 I_n \end{aligned}$$

ERRORS

Residuals are correlated, while we assume uncorrelated errors. Even worse, $\text{Var } \hat{\varepsilon}_i = \sigma^2 (1 - h_{ii})$ ← $h_{ij} = [H]_{ij}$
= non-constant variance.

We get rid of this extra variability, and consider the normalized residuals $t_i := \frac{\hat{\varepsilon}_i}{\hat{\sigma}_i \sqrt{1 - h_{ii}}}$.

However, σ is unknown, and we shall make use of the standardized residuals instead: $t_i := \frac{\hat{\varepsilon}_i}{\hat{\sigma}_i \sqrt{1 - h_{ii}}}$. (34)

ratio of a normal RV with the square root of a χ^2 RV (up to a renormalization factor)
⇒ Expect t_i to have a Student distribution. Almost! Consider instead the studentized residuals

$$t_i^* = \frac{\hat{\varepsilon}_i}{\hat{\sigma}_{-i} \sqrt{1 - h_{ii}}}$$

where $\hat{\sigma}_{-i}$ = estimator of σ when observation i is removed from the learning sample

We are going to show that t_i^* has a Student distribution.

Derivation of the distribution of t_i^* .

Notations: $X_{-i} = \begin{pmatrix} \text{---} & x_1^t & \text{---} \\ \text{---} & x_{i-1}^t & \text{---} \\ \text{---} & x_{i+1}^t & \text{---} \\ \text{---} & x_n^t & \text{---} \end{pmatrix} = \text{matrix } X \text{ without the } i\text{-th row}$
($(n-1) \times (d+1)$)

$\hat{\beta}_{-i}$ = LS estimate of β obtained from the data set $\mathcal{D}_{-i} = \{(x_1, y_1), \dots, (x_{i-1}, y_{i-1}), (x_{i+1}, y_{i+1}), \dots, (x_n, y_n)\}$

$\tilde{y}_i = x_i^t \hat{\beta}_{-i}$ = prediction of y_i

A few facts: → $\hat{\beta}_{-i} \sim \mathcal{N}(\beta, \sigma^2 (X_{-i}^T X_{-i})^{-1})$

→ $y_i - \tilde{y}_i = x_i^t (\beta - \hat{\beta}_{-i}) + \varepsilon_i$

$$\rightarrow \text{Var}(y_i - \tilde{y}_i) = \text{Var}(x_i^t(\beta - \hat{\beta}_{-i}) + \varepsilon_i) \quad (35)$$

where $\hat{\beta}_{-i} = (X_{-i}^t X_{-i})^{-1} X_{-i}^t y$,
 which depends on $\varepsilon_1, \dots, \varepsilon_{i-1}, \varepsilon_{i+1}, \dots, \varepsilon_n$ only
 $\Rightarrow \hat{\beta}_{-i}$ is uncorrelated with ε_i .

$$\begin{aligned} &= \text{Var}\{x_i^t(\beta - \hat{\beta}_{-i})\} + \text{Var}\{\varepsilon_i\} \\ &= \text{Var}\{x_i^t \hat{\beta}_{-i}\} + \sigma^2 \\ &= \sigma^2(1 + x_i^t (X_{-i}^t X_{-i})^{-1} x_i) \end{aligned}$$

$\rightarrow y_i - \tilde{y}_i$ is normally distributed.

Thus

$$\frac{\frac{y_i - \tilde{y}_i}{\sigma \sqrt{1 + x_i^t (X_{-i}^t X_{-i})^{-1} x_i}}{\sqrt{\frac{(n-d-2)\hat{\sigma}_{-i}^2}{\sigma^2}} / (n-d-2)}}{\sim \mathcal{N}(0,1)} = \frac{y_i - \tilde{y}_i}{\hat{\sigma}_{-i} \sqrt{1 + x_i^t (X_{-i}^t X_{-i})^{-1} x_i}} \sim t_{n-d-2}$$

$\sim \chi_{n-d-2}^2$

(since we have $n-1$ observations and $d+1$ predictors)

We make use of the Sherman-Morrison-Woodbury theorem, which ensures that for a $(d \times d)$ nonsingular matrix A and $(d \times 1)$ column vectors u and v , we have

$$(A + uv^t)^{-1} = A^{-1} - \frac{A^{-1}uv^tA^{-1}}{1 + v^tA^{-1}u}$$

We derive an alternative expression for this term.

Apply SMW theorem to

$$X_{-i}^t X_{-i} = X^t X - x_i x_i^t \quad (\text{since } X^t X = \sum_{i=1}^n x_i x_i^t)$$

\downarrow

$$\begin{aligned} (X_{-i}^t X_{-i})^{-1} &= (X^t X - x_i x_i^t)^{-1} \\ &= (X^t X)^{-1} + \frac{(X^t X)^{-1} x_i x_i^t (X^t X)^{-1}}{1 - x_i^t (X^t X)^{-1} x_i} \end{aligned}$$

Thus

$$\begin{aligned} x_i^t (X_{-i}^t X_{-i})^{-1} x_i &= \underbrace{x_i^t (X^t X)^{-1} x_i}_{h_{ii}} + \frac{(x_i^t (X^t X)^{-1} x_i)(x_i^t (X^t X)^{-1} x_i)}{1 - x_i^t (X^t X)^{-1} x_i} \\ &= h_{ii} + \frac{h_{ii}^2}{1 - h_{ii}} = \frac{h_{ii}}{1 - h_{ii}} \end{aligned}$$

We obtain

$$1 + x_i^t (X_{-i}^t X_{-i})^{-1} x_i = 1 + \frac{h_{ii}}{1 - h_{ii}} = \frac{1}{1 - h_{ii}} \quad (1)$$

Moreover,

$$y_i - \hat{y}_i = (1 - h_{ii})(y_i - \tilde{y}_i) \quad (2)$$

\hookrightarrow Indeed, using $X_{-i}^t y_i = X^t y - x_i y_i$, we get

$$\begin{pmatrix} 1 & | & | & | & | \\ x_1 & x_{i-1} & x_{i+1} & x_n & \\ | & | & | & | & \\ \hline & & & & \end{pmatrix} \begin{pmatrix} y_i \\ y_{i-1} \\ y_{i+1} \\ y_n \end{pmatrix}$$

(d+1) x (n-1) (n-1) x 1

$$\begin{aligned} \hat{\beta}_{-i} &= (X_{-i}^t X_{-i})^{-1} X_{-i}^t y_i \\ &= (X_{-i}^t X_{-i})^{-1} (X^t y - x_i y_i) \end{aligned}$$

$$\begin{aligned} \hat{\beta}_{-i} &= \left[(X^t X)^{-1} + \frac{(X^t X)^{-1} x_i x_i^t (X^t X)^{-1}}{1 - x_i^t (X^t X)^{-1} x_i} \right] (X^t y - x_i y_i) \quad (37) \\ &= \hat{\beta} - (X^t X)^{-1} x_i y_i + \frac{(X^t X)^{-1} x_i x_i^t (X^t X)^{-1}}{1 - x_i^t (X^t X)^{-1} x_i} (X^t y - x_i y_i) \\ &= \hat{\beta} - \frac{(X^t X)^{-1} x_i y_i (1 - x_i^t (X^t X)^{-1} x_i) - (X^t X)^{-1} x_i x_i^t (X^t X)^{-1} (X^t y - x_i y_i)}{1 - x_i^t (X^t X)^{-1} x_i} \\ &= \hat{\beta} - \frac{(X^t X)^{-1} x_i y_i - (X^t X)^{-1} x_i y_i x_i^t (X^t X)^{-1} x_i - (X^t X)^{-1} x_i x_i^t (X^t X)^{-1} (X^t y - x_i y_i)}{1 - x_i^t (X^t X)^{-1} x_i} \\ &= \hat{\beta} - \frac{(X^t X)^{-1} x_i y_i - (X^t X)^{-1} x_i x_i^t (X^t X)^{-1} X^t y}{1 - x_i^t (X^t X)^{-1} x_i} = \hat{\beta} \\ &= \hat{\beta} - \frac{(X^t X)^{-1} x_i (y_i - x_i^t \hat{\beta})}{1 - x_i^t (X^t X)^{-1} x_i} = y_i - \hat{y}_i \\ &= \hat{\beta} - \frac{(X^t X)^{-1} x_i (y_i - x_i^t \hat{\beta})}{1 - h_{ii}} = h_{ii} \end{aligned}$$

$$\Rightarrow \hat{\beta}_{-i} = \hat{\beta} - \frac{(X^t X)^{-1} x_i (y_i - \hat{y}_i)}{1 - h_{ii}} \quad (\bullet)$$

We obtain

$$y_i - x_i^t \hat{\beta}_{-i} = y_i - x_i^t \hat{\beta} + \frac{x_i^t (X^t X)^{-1} x_i}{1 - h_{ii}} (y_i - \hat{y}_i)$$

$$(y_i - \tilde{y}_i) = (y_i - \hat{y}_i) + \frac{h_{ii}}{1 - h_{ii}} (y_i - \hat{y}_i)$$

$$= \frac{1}{1 - h_{ii}} (y_i - \hat{y}_i), \text{ as required}$$

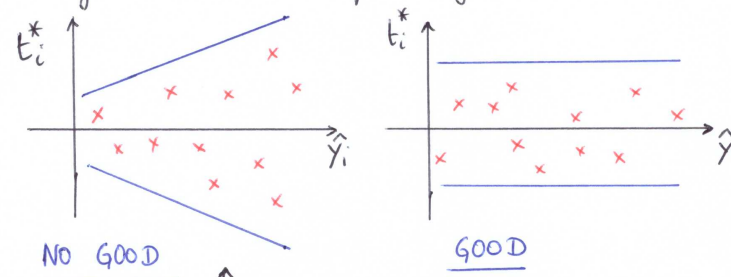
Putting (1) and (2) together,

$$\frac{\hat{\beta}_{-i} - \hat{\beta}}{\hat{\sigma}_{-i} \sqrt{1 + x_i^t (X_i^t X_i)^{-1} x_i}} = \frac{(y_i - \hat{y}_i) / (1 - h_{ii})}{\hat{\sigma}_{-i} \sqrt{1 / (1 - h_{ii})}} = \frac{y_i - \hat{y}_i}{\hat{\sigma}_{-i} \sqrt{1 - h_{ii}}} \sim t_{n-d-2} \quad \text{Good.}$$

The studentized residuals can be plotted in various ways (38) to detect possible anomalies, and check the model assumptions.

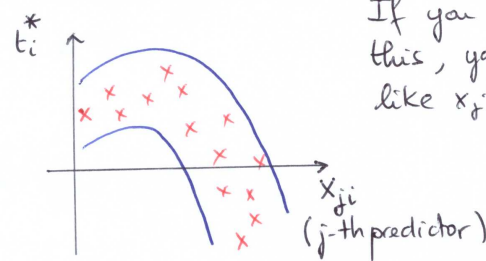
(a) Check constant variance.

By plotting t_i^* as a function of \hat{y}_i , you may visually check the assumption of constant variance:



You might consider transforming the response variable: \sqrt{y} , $\log y$ (if appropriate).

(b) Plot t_i^* against a predictor variable.



If you observe something like this, you should include terms like x_j^2 , x_j^3 , ... in the model.

(c) Check normality.

The residuals are not normally distributed & are correlated. However, $t_i^* \sim \text{Student}$, which is almost normal if $n \gg d$. You may use a QQ plot to check visually if the studentized residuals look normal indeed.

(d) Assumption of independence.

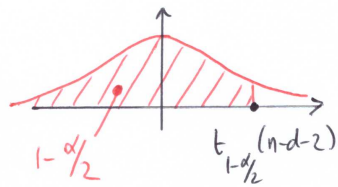
39

This assumption is technically hard to check. If the data is chronological (or spacial), a plot of the residuals versus time (or space) may reveal a systematic pattern.

(e) Use residuals to detect outliers.

↑
= points for which the response variable y_i is far from the predicted value.

For example, (x_i, y_i) = outlier if $|t_i^*| \gg t_{1-\alpha/2}(n-d-2)$



In practice, if $\alpha = 0.05$ and $n-d-2 \geq 30$, then $t_{1-\alpha/2}(n-d-2) \approx 2$, and (x_i, y_i) = outlier if $|t_i^*| \gg 2$

Outliers are BAD and should be omitted from the analysis, if this makes sense to do so. Outliers can have a dramatic effect on R^2 , by artificially increasing RSS, thus decreasing R^2 .

II.2. Analysis of the projection matrix.

Equation (2) page 36 states that

$$y_i - \hat{y}_i = (1 - h_{ii})(y_i - \tilde{y}_i)$$

error made when using the full dataset

error made when predicting y_i without information contained in x_i .

⇒ When h_{ii} is small (close to 0), the error we make,

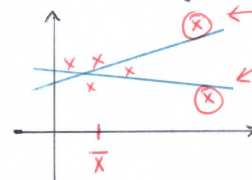
whether we include x_i in the dataset or not, is similar. In this case, x_i does not have a big impact on the prediction of y_i . The situation is reversed if h_{ii} is close to 1.

40

⇒ h_{ii} may be used to quantify which points in the learning sample influence the most the prediction.

Example: For simple linear regression ($d=1$) with intercept, we can show that

$$h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2} \in \left[\frac{1}{n}, 1 \right]$$



far away points influence the slope of the linear fit a lot. These points are away from the center of mass \bar{x} ⇒ They have a large value of h_{ii} , close to 1.

When $x_i = \bar{x}$, $h_{ii} = \frac{1}{n} = h_{min}$

x_i far from \bar{x} , then $h_{ii} \uparrow 1 = h_{max}$.

Lemma:

Let H = projection matrix onto $\mathcal{E}(X)$; $\dim \mathcal{E}(X) = d+1$.

Then

(i) $\text{Tr } H = d+1$

(ii) $\forall i \in \{1, \dots, n\}$, $0 \leq h_{ii} \leq 1$

(iii) If $h_{ii} = 0$ or 1 , then $h_{ij} = 0 \quad \forall j \neq i$

(iv) $\forall j \neq i$, $-\frac{1}{2} \leq h_{ij} \leq \frac{1}{2}$.

Since $\hat{y}_i = \sum_{j=1}^n h_{ij} y_j = h_{ii} y_i + \sum_{j \neq i} h_{ij} y_j$, we deduce from the lemma that

- ↳ If $h_{ii} = 1$, then $h_{ij} = 0$ for any $i \neq j$, and \hat{y}_i is completely determined by y_i since $\hat{y}_i = y_i$.
- ↳ If $h_{ii} = 0$, then, then $h_{ij} = 0 \quad \forall i \neq j$ and $\hat{y}_i = 0$; y_i has absolutely no influence on \hat{y}_i .
- ↳ Keeping all other values of y_j for $j \neq i$ fixed,
 (Change in \hat{y}_i) = $h_{ii} \times$ (Change in y_i)

A "high" leverage indicates that y_i provides a major contribution to the predicted value \hat{y}_i .

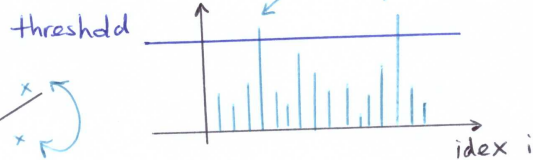
Note that $\text{Tr } H = d+1 = \sum h_{ii}$.

⇒ the average value of the h_{ii} is $\frac{d+1}{n}$.

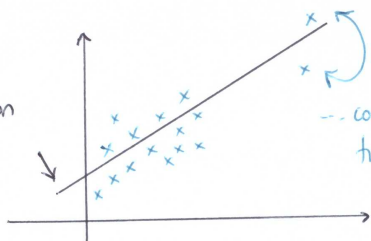
Use this to quantify what we mean by "high" leverage:

- Example:
- If $h_{ii} > \frac{2d}{n}$
 - If $h_{ii} > \frac{3d}{n}$

High leverage points



Regression line



--- correspond to these two observations.
 ↳ These are high leverage points; and not outliers! Indeed, the points are close to the regression line.

#Take Away Message

- Analysis of the (studentized) residuals allows us to detect observations (x_i, y_i) with unusual values of y_i (\equiv far from the regression line)
- Analysis of the projection matrix allows us to detect points (x_i, y_i) located away from the center of mass \bar{x} (\equiv points with a high influence on the position of the regression line)

↳ Cook (77) introduced a measure which combines information about the residuals (outliers) and about the projection matrix (high leverage points)

II.3. Cook's distance.

To evaluate the impact of observation (x_i, y_i) on the estimation of $\hat{\beta}$, we can naturally remove it from the learning sample, compute the least square estimate (denoted $\hat{\beta}_{-i}$), and evaluate the distance between $\hat{\beta}$ and $\hat{\beta}_{-i}$.

→ If the distance is "small", then (x_i, y_i) has little impact, while if the distance is "large", observation (x_i, y_i) influences the estimation of $\hat{\beta}$ significantly.

→ Since $\hat{\beta}$ and $\hat{\beta}_{-i} \in \mathbb{R}^{d+1}$, a Euclidean distance between $\hat{\beta}$ and $\hat{\beta}_{-i}$ based on an inner product can be written $d(\hat{\beta}, \hat{\beta}_{-i}) = \sqrt{(\hat{\beta} - \hat{\beta}_{-i})^T Q (\hat{\beta} - \hat{\beta}_{-i})}$, where Q is a positive definite and symmetrical matrix.

Many choices for Q are possible.

In view of the confidence region established on page 32, we may take $Q = \frac{1}{(d+1)\hat{\sigma}^2} (X^t X)$, with

$R = I_{d+1}$, so that

$$\left\{ \beta \in \mathbb{R}^{d+1} \mid (\hat{\beta} - \beta)^t \frac{(X^t X)}{(d+1)\hat{\sigma}^2} (\hat{\beta} - \beta) \leq f_{1-\alpha}(d+1, n-d-1) \right\}$$

is a $(1-\alpha)$ confidence region for β .

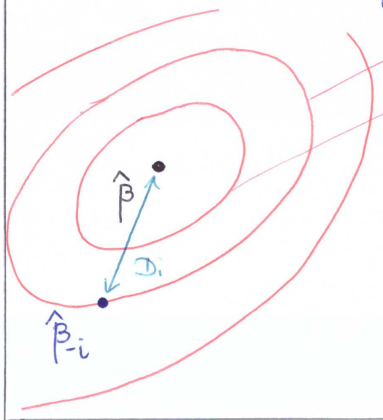
↑
Take this as a measure of the influence of each data point. Specifically, set

$$D_i = (\hat{\beta} - \hat{\beta}_{-i})^t \frac{(X^t X)}{(d+1)\hat{\sigma}^2} (\hat{\beta} - \hat{\beta}_{-i})$$

(COOK'S DISTANCE)

⇒ Evaluate Cook's distance for each observation.

If for example $D_i \approx f_{0.5}(d+1, n-d-1)$, then the removal of the i -th observation moves the LS estimate to the edge of the 50% confidence region. That's not good.



50% confidence region
10% confidence region

- ▶ Cook suggested that each $\hat{\beta}_{-i}$ to stay within a 10% confidence region to cause no concern.
- ▶ At first, it seems like computing $(n+1)$ regressions is costly. In fact,

We can show that Cook's distance admits a simple form:

$$D_i = \frac{1}{d+1} \frac{h_{ii}}{1-h_{ii}} t_i^2 \quad (*)$$

Two terms contribute

- to Cook's distance:
- t_i^2 = large value for outliers (p. 34)
 - $\frac{h_{ii}}{1-h_{ii}}$ = large for 'high' leverage points (corresponding to large values of h_{ii})

⇒ Cook's distance naturally combine these two elements in a unique coefficient.

Proof of (*)

Expression (1) page 37: $\hat{\beta}_{-i} - \hat{\beta} = -(y_i - \hat{y}_i) \frac{(X^t X)^{-1} x_i}{1-h_{ii}}$

Expression (2) page 36: $y_i - \hat{y}_i = (1-h_{ii})(y_i - \tilde{y}_i)$

$$\Rightarrow \hat{\beta}_{-i} - \hat{\beta} = -(y_i - \tilde{y}_i) (X^t X)^{-1} x_i$$

Thus

$$D_i = (y_i - \tilde{y}_i)^2 x_i^t (X^t X)^{-1} \frac{(X^t X)}{(d+1)\hat{\sigma}^2} (X^t X)^{-1} x_i$$

$$= (y_i - \tilde{y}_i)^2 \frac{1}{(d+1)\hat{\sigma}^2} x_i^t (X^t X)^{-1} x_i = h_{ii}$$

$$= \frac{h_{ii}}{(d+1)\hat{\sigma}^2} (y_i - \tilde{y}_i)^2$$

$$= \frac{h_{ii}}{(d+1)\hat{\sigma}^2} \frac{1}{(1-h_{ii})^2} (1-h_{ii})^2 (y_i - \tilde{y}_i)^2 = \frac{h_{ii}}{(d+1)\hat{\sigma}^2} \frac{1}{(1-h_{ii})^2} (y_i - \hat{y}_i)^2$$

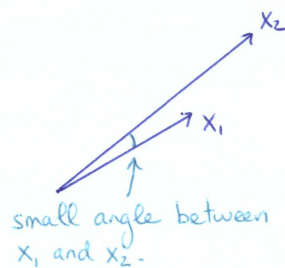
$$D_i = \frac{h_{ii}}{(1-h_{ii})^2} \frac{1}{(d+1)\hat{\sigma}^2} (y_i - \tilde{y}_i)^2 \quad (45)$$

$$= \frac{h_{ii}}{1-h_{ii}} \frac{1}{d+1} \left(\frac{y_i - \tilde{y}_i}{\hat{\sigma}\sqrt{1-h_{ii}}} \right)^2$$

$$= \frac{h_{ii}}{1-h_{ii}} \frac{t_i^2}{d+1}, \text{ as required} \quad \blacksquare$$

II. 4. Detecting colinearity.

The derivation of the LS estimate requires that $\text{rank } X = d+1$, otherwise $X^T X$ is not invertible. Even if this assumption is rarely violated in practice, it often happens that the columns of X are "nearly" linearly dependent. The consequences are rather bad.



If the angle between x_1 and x_2 is small, the plane spanned by x_1 and x_2 changes dramatically if one of the x_i or x_2 changes slightly.

Since $\hat{y} = \perp$ projection on $\mathcal{C}(X)$, the predicted value will also be greatly affected.
 \Rightarrow Prediction is unreliable and the estimation of β is unstable.

Mathematically, computing the inverse of $X^T X$ is numerically unstable: the matrix $X^T X$ is ill-conditioned; the determinant of $X^T X$ will be close to 0, and thus $(X^T X)^{-1}$ has large values.

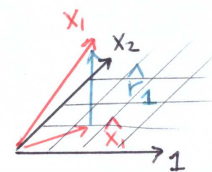
But remember that the covariance matrix of $\hat{\beta}$ is precisely $\sigma^2 (X^T X)^{-1}$. We see that colinearity amongst the predictors implies an increase in the variance of the coefficient estimates.

\hookrightarrow This impacts any hypothesis test of the form $H_0: \beta_i = 0$, since the z-score $z_i = \frac{\hat{\beta}_i}{\hat{\sigma}\sqrt{v_{ii}}}$ (page 21) is smaller than it should be (due to an artificial increase in v_{ii}), hence failing to reject the null when in fact you should. In other words, you conclude that there is no statistical evidence of association of a given predictor with the response variable, when in fact there is. The POWER of the test is reduced.
 $\hookrightarrow P(\text{reject } H_0 \mid H_1 \text{ is true})$

\Rightarrow We need ways to detect colinearity. We derive a simple coefficient (VIF) to assess colinearity. We proceed in two steps.

- Step I: Alternative expression for $\hat{\beta}_1$ (we treat the case $\hat{\beta}_1$ without loss of generality. All coefficients can be treated the same way).

Idea: Since we worry about colinearity amongst the predictors, let's regress one of them (here x_1) on the remaining ones ($1, x_2, \dots, x_d$); $x_1 = \hat{x}_1 + \hat{r}_1$



our output variable, predicted using $(1, x_2, \dots, x_d)$.

$$\Rightarrow \forall i=1, \dots, n \quad x_{i1} = \hat{x}_{i1} + \hat{r}_{i1}$$

- Next, recall that $X^t(y - X\hat{\beta}) = 0$; the second row of this equation is (47)

$$\sum_{i=1}^n x_{i1} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_d x_{id}) = 0$$

↓ replace with $\hat{x}_{i1} + \hat{r}_{i1}$

$$\diamond \quad \sum_{i=1}^n (\hat{x}_{i1} + \hat{r}_{i1})(y_i - \hat{\beta}_0 - \dots - \hat{\beta}_d x_{id}) = 0$$

- Likewise, rows 2 to (d+1) of $X^t(y - X\hat{\beta})$ yield $\sum_{i=1}^n x_{ik} (y_i - \hat{y}_i) = 0$

- Also, since \hat{x}_{i1} is a linear function of $1, x_{i2}, \dots, x_{id}$, we immediately see that $\sum_{i=1}^n (y_i - \hat{y}_i) \hat{x}_{i1} = 0$

$$(\langle y - \hat{y}, \hat{x}_1 \rangle = 0 \text{ since vectors } (y - \hat{y}) \text{ and } \hat{x}_1 \text{ are } \perp)$$

↓

We get from \diamond that

$$\sum_{i=1}^n (\hat{x}_{i1} + \hat{r}_{i1})(y_i - \hat{y}_i) = 0 \Rightarrow \sum_{i=1}^n \hat{r}_{i1} (y_i - \hat{y}_i) = 0$$

$$\square \quad \sum_{i=1}^n \hat{r}_{i1} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_d x_{id}) = 0$$

In addition, $\sum_{i=1}^n \hat{r}_{i1} = 0$ since the intercept is included in the model.

$$\sum_{i=1}^n \hat{r}_{i1} x_{ik} = \sum_{i=1}^n x_{ik} (x_{i1} - \hat{x}_{i1}) = 0$$

↑
k=2, ..., d
LS equation when regressing x_1 on $(1, x_2, \dots, x_d)$.

Equation \square becomes

$$\sum_{i=1}^n \hat{r}_{i1} (y_i - \hat{\beta}_1 x_{i1}) = 0$$

$$\Rightarrow \sum_{i=1}^n \hat{r}_{i1} (y_i - \hat{\beta}_1 \hat{x}_{i1} - \hat{\beta}_1 \hat{r}_{i1}) = 0$$
(48)

Next, notice that

$$\begin{aligned} \sum_{i=1}^n \hat{r}_{i1} \hat{x}_{i1} &= \sum_{i=1}^n \hat{r}_{i1} (\hat{\alpha}_0 + \hat{\alpha}_2 x_{i2} + \dots + \hat{\alpha}_d x_{id}) \\ &\quad \uparrow \quad \uparrow \quad \uparrow \\ &\quad \text{regression coefficients when} \\ &\quad \text{regressing } x_1 \text{ on } (1, x_2, \dots, x_d) \\ &= \hat{\alpha}_0 \underbrace{\sum_{i=1}^n \hat{r}_{i1}}_{=0} + \hat{\alpha}_2 \underbrace{\sum_{i=1}^n \hat{r}_{i1} x_{i2}}_{=0} + \dots + \hat{\alpha}_d \underbrace{\sum_{i=1}^n \hat{r}_{i1} x_{id}}_{=0} \\ &= 0 \end{aligned}$$

Thus $\sum_{i=1}^n \hat{r}_{i1} (y_i - \hat{\beta}_1 \hat{r}_{i1}) = 0$, which yields

$$\hat{\beta}_1 = \frac{\sum \hat{r}_{i1} y_i}{\sum \hat{r}_{i1}^2} \leftarrow > 0 \text{ as long as } x_1 \text{ is not an exact linear combination of } 1, x_2, \dots, x_d, \text{ which is ok if rank } X = (d+1).$$

We can go further since $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_d x_{id} + \varepsilon_i$

$$\sum \hat{r}_{i1} y_i = \beta_0 \underbrace{\sum \hat{r}_{i1}}_{=0} + \beta_1 \sum \hat{r}_{i1} x_{i1} + \dots + \beta_d \sum \hat{r}_{i1} x_{id} + \sum \hat{r}_{i1} \varepsilon_i = 0$$

$$\begin{aligned} &= \sum \hat{r}_{i1} (\hat{x}_{i1} + \hat{r}_{i1}) \\ &= \underbrace{\sum \hat{r}_{i1} \hat{x}_{i1}}_{=0} + \sum \hat{r}_{i1}^2 = \sum \hat{r}_{i1}^2 \\ &\quad \text{since } \hat{x}_1 \perp \hat{r}_1 \end{aligned}$$

$$\Rightarrow \sum \hat{r}_{i1} y_i = \beta_1 \sum \hat{r}_{i1}^2 + \sum \hat{r}_{i1} \varepsilon_i \Rightarrow \hat{\beta}_1 = \beta_1 + \frac{\sum \hat{r}_{i1} \varepsilon_i}{\sum \hat{r}_{i1}^2}$$

(*)

• Step I. Alternative expression for $\text{Var } \hat{\beta}_1$. (49)

Making use of (*) at bottom of page 18, we get

$$\text{Var } \hat{\beta}_1 = \frac{\sum \hat{r}_{i1}^2 \text{Var } \varepsilon_i}{(\sum \hat{r}_{i1}^2)^2} = \frac{\sigma^2}{\sum \hat{r}_{i1}^2}$$

Also, $R_1^2 = 1 - \frac{\text{RSS}_1}{\text{TSS}_1} \leftarrow \sum \hat{r}_{i1}^2$
 \uparrow
 R^2 coefficient obtained by regressing x_1 on all other predictors.

Thus $\sum \hat{r}_{i1}^2 = \text{TSS}_1 (1 - R_1^2)$, and $\text{Var } \hat{\beta}_1 = \frac{\sigma^2}{\text{TSS}_1} \frac{1}{1 - R_1^2}$.

Generalizing this expression for other coefficients,

$$\text{Var } \hat{\beta}_j = \frac{\sigma^2}{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2} \frac{1}{1 - R_j^2},$$

where $R_j^2 = R^2$ coefficient obtained by regressing x_j on all other predictors.

A few observations:

- A larger $\sigma^2 \Rightarrow$ a larger variance of the coefficient estimates
- More variability in $x_j \Rightarrow$ Better accuracy in the estimation of $\hat{\beta}_j$.
- Increasing n decreases $\text{Var } \hat{\beta}_j$.
- Smallest variance is obtained for $R_j^2 = 0$; i.e. when x_j has 0 sample correlation with every other predictor.
- Worst case scenario when $R_j^2 = 1$, i.e. when x_j is an exact linear combination of other predictors (colinearity!)

\Rightarrow We define the Variance Inflation Factor (VIF) (50) to be the ratio of the variance of $\hat{\beta}_j$ with its smallest possible value, given by $\sigma^2 / \sum (x_{ij} - \bar{x}_j)^2$.

$$\text{VIF}(\hat{\beta}_j) = \frac{1}{1 - R_j^2} = \frac{\text{Var } \hat{\beta}_j}{\min \text{var } \hat{\beta}_j}$$

A large VIF is a sign of colinearity of x_j with other predictors.

Rule of thumb: a $\text{VIF} > 5$ or 10 is problematic.

What shall we do when we detect colinearity?

- drop one or more variables
- combine them in some way (if meaningful to do so)
- use principal components instead.

III. MAKING PREDICTIONS.

Last step in the inference procedure: making predictions. There are two kinds of uncertainty associated with the prediction task:

↳ uncertainty about the coefficient estimates $\hat{\beta}_j$.

To quantify the source of uncertainty associated with the estimation of β_j , we construct CONFIDENCE INTERVALS using $\text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_d x_d)$

↑
new observation
 $x = (x_1, \dots, x_d)$

51
 ↪ Even if we knew the true parameters β_j , the response variable would not be predicted perfectly because of the additive error ε . Taking this source of uncertainty into account leads to a PREDICTION INTERVAL: the response variable is

$$Y_{nt} = \beta_0 + \beta_1 x_{nt,1} + \dots + \beta_d x_{nt,d} + \varepsilon_{nt}$$

New input point ↗

with
 - $E \varepsilon_{nt} = 0$
 - $\text{Var } \varepsilon_{nt} = \sigma^2$
 - $\text{Cov}(\varepsilon_j, \varepsilon_{nt}) = 0$
 $j=1, \dots, n$

→ Prediction is $\hat{y}_{nt} = \hat{\beta}_0 + \dots + \hat{\beta}_d x_{nt,d}$

→ Prediction error is

$$\hat{\varepsilon}_{nt} := Y_{nt} - \hat{y}_{nt} = x_{nt}^t (\beta - \hat{\beta}) + \varepsilon_{nt}$$

(where $x_{nt} = (1, x_{nt,1}, \dots, x_{nt,d})^t$.)

$$E \hat{\varepsilon}_{nt} = x_{nt}^t (\beta - E \hat{\beta}) + E \varepsilon_{nt} = 0$$

$$\text{Var } \hat{\varepsilon}_{nt} = \text{Var} \{ x_{nt}^t (\beta - \hat{\beta}) + \varepsilon_{nt} \}$$

(Since $\hat{\beta}$ is computed from the first n observations, it is uncorrelated with ε_{nt})

$$= \text{Var} \{ x_{nt}^t (\beta - \hat{\beta}) \} + \text{Var } \varepsilon_{nt}$$

$$= x_{nt}^t \sum_{\hat{\beta}} x_{nt} + \sigma^2$$

$$= \sigma^2 (1 + x_{nt}^t (X^t X)^{-1} x_{nt})$$

52
 Then, business as usual:

$$\frac{\hat{\varepsilon}_{nt}}{\sigma \sqrt{1 + x_{nt}^t (X^t X)^{-1} x_{nt}}} \sim \mathcal{N}(0, 1)$$

$$\frac{Y_{nt} - \hat{y}_{nt}}{\hat{\sigma} \sqrt{1 + x_{nt}^t (X^t X)^{-1} x_{nt}}} \sim t_{n-d-1}$$

$\sim \chi_{n-d-1}^2$

Denoting by $t_{\alpha}(n-d-1)$ the α -quantile of the t_{n-d-1} distribution, a $100(1-\alpha)\%$ prediction interval for Y_{nt} is

$$x_{nt}^t \hat{\beta} = \hat{y}_{nt} \pm t_{1-\alpha/2}(n-d-1) \hat{\sigma} \sqrt{1 + x_{nt}^t (X^t X)^{-1} x_{nt}}$$

PREDICTION INTERVAL (PI)

Remarks: (i) The confidence interval for $x_{nt}^t \beta$ is narrower than the PI, as $x_{nt}^t \hat{\beta} \sim \mathcal{N}(x_{nt}^t \beta, \sigma^2 x_{nt}^t (X^t X)^{-1} x_{nt})$, and is given by

$$x_{nt}^t \hat{\beta} = \hat{y}_{nt} \pm t_{1-\alpha/2}(n-d-1) \hat{\sigma} \sqrt{x_{nt}^t (X^t X)^{-1} x_{nt}}$$

CONFIDENCE INTERVAL (CI)

(ii) For simple linear regression ($d=1$), the PI simplifies to

$$\hat{\beta}_0 + \hat{\beta}_1 x_{nt} \pm t_{1-\alpha/2}(n-2) \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x_{nt} - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

⇒ Prediction error has minimum variance at the center of mass $x_{nt} = \bar{x}$. This can be shown to be true when $d > 1$ as well.

(iii) For known β_j s, CI has width 0, but not PI.