

PART 2:
LINEAR REGRESSION

TMA4267 L7
07.02.2017

Model definition [F3.1.0]

Y = variable of primary interest (response, dependent variable)

X_1, X_2, \dots, X_n = regressors, explanatory variables
independent variables, covariates

Assumptions:

$$Y = \underbrace{f(x_1, x_2, \dots, x_n)}_{\text{systematic component}} + \varepsilon$$

↑
error term

- 1) Systematic component is a linear combination of the covariates.

$$f(x_1, x_2, \dots, x_n) = \underbrace{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n}_{\substack{\text{multiple} \\ \text{simple}}}$$

$$x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} \quad f(x) = x^T \beta$$

$p \times 1$
 $p = k+1$

- 2) Additive errors $Y = x^T \beta + \varepsilon$

Restrictive? Maybe \rightarrow transformations?

Data and design matrix

We collect independent data (Y_i, X_i) for $i=1, \dots, n$

response

$$\downarrow \\ Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

↑ pair
↓
 $(1, x_{i1}, x_{i2}, \dots, x_{in})^\top$

$$\underset{n \times p}{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ \vdots & & & & \\ 1 & x_{n1} & x_{n2} & & x_{nk} \end{bmatrix} \leftarrow \begin{array}{l} \text{individual/observational} \\ \text{unit} \end{array}$$

↑ covariates

Design matrix . We will assume that $n \gg p$.

$$\text{and } \text{rank}(X) = p$$

n = number of observations

p = number of covariates + 1 (intercept)

Q: What can make $\text{rank}(X) < p$?

Ex: Munich rent index: $n = 3082$

$$Y = \text{rent or rent pr sq.m} \rightarrow Y = \begin{bmatrix} 5.26 \\ 0.41 \\ \vdots \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} 1 & 20 & 1970 & 1 & 0 & 0 & 0 & 91b \\ 1 & 41 & 1941 & 1 & 0 & 1 & 1 & 206 \\ \vdots & \vdots \\ 1 & \vdots & \vdots & 3 & \vdots & 1 & 1 & \vdots \\ \uparrow & \uparrow & \uparrow & \text{location} & \uparrow & \text{bath} & \uparrow & \text{district} \\ \text{area} & \text{year of} & \text{ } & \text{C. house} \end{bmatrix}$$

\mathbf{X} is chosen so that we have a good model.

The classical linear model

$$Y_i = x_i^T \beta + \varepsilon_i \quad \text{one person}$$

$$Y = \begin{matrix} \mathbf{X} \beta \\ n \times 1 \end{matrix} + \begin{matrix} \varepsilon \\ n \times p \\ p \times 1 \\ n \times 1 \end{matrix} \quad \text{unobserved random vector}$$

$$1) E(\varepsilon) = 0$$

$$2) \text{Cor}(\varepsilon) = E(\varepsilon \varepsilon^T) - E(\varepsilon)E(\varepsilon)^T = \sigma^2 I$$

$$\text{Var}(\varepsilon_i) = \sigma^2 \text{ for all } i \leftarrow \text{homoscedastic errors}$$

$$\text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \text{ for } i \neq j \leftarrow \text{uncorrelated errors}$$

$$3) \text{The design matrix has rank } \text{rank}(\mathbf{X}) = h+1 = p$$

$$4) \text{If we in addition assume that } \varepsilon \sim N_n(0, \sigma^2 I) \text{ then we have a normal linear regression model.}$$

What does this imply for the distribution of Y ?

$Y \sim N_n$ since $Y = \underbrace{\mathbb{X}\beta}_{\text{constant}} + \underbrace{\varepsilon}_{N_n}$, and

$$E(Y) = E(\mathbb{X}\beta + \varepsilon) = \mathbb{X}\beta + E(\varepsilon) = \mathbb{X}\beta$$

$$\text{Cov}(Y) = \text{Cov}(\mathbb{X}\beta + \varepsilon) = 0 + \text{Cov}(\varepsilon) = \sigma^2 I$$

$$Y \sim N_n(\mathbb{X}\beta, \sigma^2 I)$$

The covariates \mathbb{X} may be regarded as random variables, and then the assumptions (1)+(2) are made conditional on $\mathbb{X} = x$, so $E(\varepsilon | \mathbb{X} = x) = 0$ and $\text{Cov}(\varepsilon | \mathbb{X} = x) = \sigma^2 I$

If we instead assume that

$$\begin{bmatrix} Y \\ \mathbb{X}_1 \\ \mathbb{X}_2 \\ \vdots \\ \mathbb{X}_n \end{bmatrix} \sim N_{n+1} \Rightarrow E(Y | \mathbb{X} = x) = \text{linear in } x$$

$$\text{Var}(Y | \mathbb{X} = x) = \text{not dependent on } x$$

Model parameters, estimates and residuals [F 3.1.1]

$$Y = \mathbf{X}\beta + \varepsilon, E(\varepsilon) = 0, \text{Cov}(\varepsilon) = \sigma^2 I$$

The model parameters are β, σ^2
 the unknown $p \times 1$

We will develop estimators:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Y \quad \begin{matrix} (\text{LS}) \\ \text{by least squares and} \\ \text{maximum likelihood.} \end{matrix}$$

$$\hat{\sigma}^2 = \frac{1}{n-p} (Y - \mathbf{X}\hat{\beta})^T (Y - \mathbf{X}\hat{\beta}) \quad \begin{matrix} (\text{RL}) \\ \text{by restricted} \\ \text{maximum likelihood} \\ (\text{REML}) \end{matrix}$$

Further: Y is a random vector with mean

$\mathbf{X}\beta$, and estimator for $E(Y) = \mathbf{X}\beta$ is $\hat{Y} = \mathbf{X}\hat{\beta}$.

The error ε is a random vector with $E(\varepsilon) = 0$ and $\text{Cov}(\varepsilon) = \sigma^2 I$, but ε is not observed.

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

$\begin{matrix} \downarrow & \downarrow & \downarrow \\ \text{observed} & \text{unknown} & \text{unobserved} \\ \uparrow & \uparrow & \uparrow \\ \text{observed} & \text{observed} & \text{unobserved} \end{matrix}$

Our best guess for the error is the residual vector $(\hat{\varepsilon}, e)$

$$\hat{\varepsilon} = Y - \hat{Y} = Y - \mathbf{X}\hat{\beta}$$

So, the residuals can be calculated, and we may think of the residuals as predictions of the errors

Be aware: don't mix errors ε (unobserved) with residuals $\hat{\varepsilon}$ ("observed").

The residuals will be used to assess model assumptions as proxies for the errors.

Model assumption [F 3.1.2]

1) Linearity of covariates $Y = \sum \beta_i x_i + \varepsilon$

$$y_i = \beta_0 + \beta_1 x_{i1} + \varepsilon_i \quad \text{ok}$$

$$\beta_0 + \beta_1 z_{i1}^2 + \varepsilon_i \quad \text{ok}$$

$$\beta_0 + \beta_1 \log(z_{i1}) + \varepsilon_i \quad \text{ok}$$

$$\beta_0 + \beta_1 \sin(\beta_2 z_{i1}) + \varepsilon_i \quad \text{not ok}$$

If the relationship between Y and x_1 is nonlinear \Rightarrow ok to use polynomial (or similar) in x_1 . More advanced: nonparametric regression

2) Homoscedastic error variance: $\text{Cov}(\varepsilon) = \sigma^2 I$

Need to check that $\text{Var}(\varepsilon)$ does not vary systematically across observations.

Look at covariate vs residuals \rightarrow trend?
fan out, fan in.

Solution if problem: If we knew $\text{Cov}(\varepsilon) = \sum \delta$

We may use a so-called general linear model,
with weighted least squares (lecEx3.P4),
but Σ is in general unknown.

Remaining : 3) uncorrelated errors
4) additive error .

4