# Regression II

Sofia Olhede

**EPFL**

November 10, 2020

# Set–up

EPFL

- We can generalize this to

$$P(Y \leq y) = F(y)$$

and so as $n \to \infty$

$$F_n(y) \to F(y)$$

$$
\begin{aligned}
F^{-1}(F_n(y_i)) &\approx y_i \\
F^{-1}(\text{prop. of obs. } \leq y_i) &\approx y_i.
\end{aligned}
$$

# QQ–plot

# Set–up

**EPFL**

- Note that
$$\mathbb{V}\mathrm{ar}\{e\} = \sigma^2(I_n - P)$$

- If $p_{ii} \approx 1$ then the variance of the $i$th residual is very low.

- Totally determined by X, i.e. the design matrix is forcing the $i$th observation to have high impact.

- The $i$th observation has **high leverage**.

- $\sum_{i=1}^{n} p_{ii} = p$ so the "average" is $p/n$ and a rule of thumb is to take notice when
$$p_{ii} > \frac{2p}{n}.$$

# Weighted Least Squares

EPFL

- Consider the linear model

$$\mathbb{E}\{Y_i\} = \mathsf{x}_i \beta$$

and

$$\mathbb{V}\mathrm{ar}(Y_i) = \frac{\sigma^2}{w_i},$$

where $w_i$ are known weights.

- heteroscedastic variables.
- Using least squares is no longer optimal.
- Cases with small $w_i$ need to be downweighted with respect to the parameter estimation while those with $w_i$ large need to be given more weight.

# Weighted Least Squares

EPFL

- Find an estimate for $\beta$ by minimising the weighted sum of squares:

$$
\begin{aligned}
S(\beta) &= \sum_{i=1}^{n} w_i \left( Y_i - \mathsf{x}_i^T \beta \right)^2 \\
&= \sigma^2 \sum_{i=1}^{n} \frac{(Y_i - E\{Y_i\})^2}{\mathbb{V}\mathrm{ar}\{Y_i\}}
\end{aligned}
$$

# Weighted Least Squares

EPFL

- In vector form we then have:

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

$$\mathbb{V}\mathrm{ar}(\epsilon) = \sigma^2 I_n$$

$$\mathbb{E}\{\epsilon\} = 0$$

and

$$D = \begin{pmatrix} \frac{1}{\sqrt{w_1}} & 0 & \ldots & 0 \\ 0 & \frac{1}{\sqrt{w_2}} & \ldots & 0 \\ 0 & \ldots & \frac{1}{\sqrt{w_i}} & 0 \\ \ldots & \ldots & \ldots & \ldots \\ 0 & \ldots & 0 & \frac{1}{\sqrt{w_n}} \end{pmatrix}$$

# Weighted Least Squares

EPFL

- We then multiply through the linear equation by

$$
\begin{aligned}
\mathsf{D}^{-1}\mathsf{Y} &= \mathsf{D}^{-1}\mathsf{X}\beta + \epsilon \\
\tilde{\mathsf{Y}} &= \tilde{\mathsf{X}}\beta + \epsilon \\
\mathbb{V}\mathrm{ar}\{\epsilon\} &= \sigma^2\mathsf{I}_n \\
\mathbb{E}\{\epsilon\} &= 0
\end{aligned}
\tag{1}
$$

This is recognizable as a liner model. The $\beta$ estimate is given by

$$
\begin{aligned}
\hat{\beta} &= \left(\tilde{\mathsf{X}}^T\tilde{\mathsf{X}}\right)^{-1}\tilde{\mathsf{X}}^T\tilde{\mathsf{Y}} \\
&= \left(\left(\mathsf{D}^{-1}\mathsf{X}\right)^T\left(\mathsf{D}^{-1}\mathsf{X}\right)\right)^{-1} \\
&\quad \left(\mathsf{D}^{-1}\mathsf{X}\right)^T\mathsf{D}^{-1}\mathsf{Y} \\
&= \left(\mathsf{X}^T\mathsf{V}\mathsf{X}\right)^{-1}\mathsf{X}^T\mathsf{V}\mathsf{Y}
\end{aligned}
\tag{2}
$$

where $\mathsf{V} = \mathsf{D}^{-2}$.

**EPFL**

# Testing in the Least Squares Set–Up

- Assume that $\boldsymbol{Y} \sim N(X\beta, \sigma^2 I_n)$.
- Definition: If $Z \sim N(\boldsymbol{\mu}, I_n)$ where $\boldsymbol{\mu} \neq 0$ then $Z^T Z$ is said to have a non-central $\chi^2$ distribbution on $n$ d. o. f. and non-centrality parameter $\delta > 0$ given by $\delta^2 = \mu^T \mu$.
- $\mu = 0 \Rightarrow \chi_n^2$ as the non-centrality parameter is then zero.
- We normally for the general distribution write $U = Z^T Z \sim \chi_n^2(\delta)$.
- The distribution of $Z^T Z$ depends on $\mu$ only via $\delta$.
- $\mathbb{E}(U) = n + \delta^2$.
- $\mathbb{V}\text{ar}(U) = 2n + 4\delta^2$
- If $U_i \sim \chi_{n_i}^2(\delta_i)$ for $i = 1, \ldots k$ and if the $\{U_i\}$ are all independent then

$$\sum_{i=1}^{k} U_i \sim \chi_n^2(\delta),$$

where $n = \sum_{i=1}^{k} n_i$ and $\delta^2 = \sum_{i=1}^{k} \delta_i^2$.

**EPFL**

# Testing in the Least Squares Set–Up

- Lemma: If $Z \sim N(\boldsymbol{\mu}, I_n)$ and if A is a $n \times n$ symmetric and idempotent matrix of rank $r$ then

$$Z^T A Z \sim \chi_n^2(\delta),$$

where $\delta^2 = \mu^T A \mu$.

Proof: Let $A$ be a symmetric idempotent matrix ($A^2 = A$). Then $A$ has $r$ eigenvalues that are unity, and $n - r$ eigenvalues that are zero. Because $A$ is symmetric there is an orthogonal $P^T P = I_n$ matrix $P$ st

$$P^T A P = D,$$

where $D$ is diagonal with $r$ ones and $n - r$ zeros down the diagonal. Let $V = P^T Z$. Then

$$V \sim N(P^T \mu, I_n).$$

EPFL

## Testing in the Least Squares Set–Up

- Furthermore it follows that

$$
\begin{aligned}
Z^T A Z &= (PV)^T A (PV) \\
&= V^T P^T A P V \\
&= V^T D V \\
&= V^T D^T D V \\
&= (DV)^T (DV) \\
&= \text{sum of squares of } r \text{ components} \\
&= \chi_r^2(d),
\end{aligned}
$$

for some $d$. In fact
$$
d^2 = \mathbb{E}(DV)^T \, \mathbb{E}(DV) = (DP^T \mu)^T (DP^T \mu) = \mu^T PDP^T \mu = \mu^T A \mu.
$$

# Testing in the Least Squares Set–Up

- Lemma: If $Z \sim N(\boldsymbol{\mu}, I_n)$ and $A_1$ and $A_2$ are symmetric idempotent matrices such that $A_1 A_2 = 0$ then $Z^T A_1 Z$ and $Z^T A_2 Z$ are independent.
  Proof: $Z^T A_i Z = (A_i Z)^T (A_i Z)$ for $i = 1, 2$. Consider the two vectors $A_1 Z$ and $A_2 Z$ then

$$\mathbb{C}\text{ov}\{A_1 Z, A_2 Z\} = A_1 \, \mathbb{C}\text{ov}(Z) A_2^T \tag{3}$$
$$= A_1 I A_2 = 0. \tag{4}$$

  This means that every component of $A_1 Z$ is uncorrelated with every component of $A_2 Z$. By normality this means that the components are independent.

- Corollary: If $A_1, \ldots, A_k$ are symmetric and idempotent and if $A_i A_j = 0$ for $i \neq j$ then $\{Z^T A_i Z\}$ are mutually independent.

# Testing in the Least Squares Set–Up

- Lemma: If $A_1, \ldots A_k$ are symmetric $n \times n$ matrices such that $\sum A_i = I$ and such that $\mathrm{rank}(A_i) = r_i$ then the following are equivalent:

  (a) $\sum_i r_i = n$
  (b) $A_i A_j = 0 \ i \neq j$
  (c) $A_i$ are idempotent for $i = 1, \ldots, k$.

# Testing in the Least Squares Set–Up

- If

$$Z \sim N(\mu, I_n)$$

and

$$\sum_i A_i = I_n$$

where $A_i$ with ranks $r_i$ are symmetric $n \times n$ matrices such that at least one of

1. $\sum_i r_i = n$
2. $A_i A_j = 0$ for $i \neq j$
3. $A_i$ are idempotent.

holds (and therefore all of them, proof omitted) then

$$Z^T A_i Z$$

are independent

$$\chi^2_{r_i}(\delta_i)$$

where $\delta_i^2 = \mu^T A_i \mu$.

# Testing in the Least Squares Set–Up

- Proof: $A_i$ are assumed to be idempotent. By the lemma this means

$$Z^T A_i Z \sim \chi^2_{r_i}(\delta_i).$$

Because they are mutually orthogonal by assumption this implies that

$$Z^T A_i Z$$

are independent.

# Testing in the Least Squares Set–Up

EPFL

- Assume we want to test the hypothesis

$$H_0 : \ A\beta = 0$$

versus

$$H_1 : \ A\beta \neq 0$$

where $\mathrm{rank}A = s = p - p_0$.

- Under $H_0$ we get the simpler linear model

$$E\{Y\} = X_0\beta_0$$

where $\beta_0$ is $p_0 \times 1$. New hat matrix:

$$P_0 = X_0\left(X_0^T X_0\right)^{-1} X_0^T.$$

# Testing in the Least Squares Set–Up

EPFL

- $P_0$ has trace $p_0$.
- Consider the likelihood ratio:

$$t = \frac{\text{maximum liklihood under } H_1}{\text{maximum liklihood under } H_0}$$

- From MLE we get a biased estimate of $\sigma$ and the least squares estimates of $\widehat{\boldsymbol{\beta}}$.
- Plugging in:

$$t = \left(\frac{\widehat{\sigma}^2_{ML,0}}{\widehat{\sigma}^2_{ML}}\right)^{\frac{n}{2}}.$$

- Consider a monotonic increasing function of $t$:

$$f(t) = \frac{n-p}{p-p_0}\left(t^{2/n} - 1\right),$$

or

$$F = \frac{n-p}{p-p_0}\frac{RSS_0 - RSS}{RSS}.$$

# Testing in the Least Squares Set–Up

EPFL

- Use the Fisher-Cochran theorem:

$$I_n = (I_n - P) + (P - P_0) + P_0$$

with ranks

$$n = (n - p) + (p - p_0) + p_0.$$

Let

$$
\begin{aligned}
A_1 &= (I_n - P) \\
A_2 &= (P - P_0) \\
A_3 &= P_0
\end{aligned}
$$

These are are symmetric and idempotent.

# Testing in the Least Squares Set–Up

EPFL

- We may write $P_0 = XB$ for some $B$ of constants. Let

$$Z = \frac{1}{\sigma}Y$$

Note

$$RSS = Y^T A_1 Y = \sigma^2 Z^T A_1 Z,$$

$$RSS_0 - RSS = Y^T A_2 Y = \sigma^2 Z^T A_2 Z$$

and so with NTA by the Fisher-Cochran theorem

$$RSS/\sigma^2 \sim \chi^2_{n-p}$$

and

$$(RSS_0 - RSS)/\sigma^2 \sim \chi^2_{p-p_0}$$

independently (the non-centrality parameters vanish.)

# Testing in the Least Squares Set–Up

**EPFL**

- We then have

$$
\begin{aligned}
F &= \frac{\sigma^2}{\sigma^2} \frac{n-p}{p-p_0} \frac{RSS_0 - RSS}{RSS} \\
&\sim \frac{\chi^2_{p-p_0}/(p-p_0)}{\chi^2_{n-p}/(n-p)} \\
&\sim F_{p-p_0, n-p}
\end{aligned}
$$

- Decompose the **total sum of squares** by

$$
\begin{aligned}
Y^T Y &= Y^T (I_n - P) Y \\
&\quad + Y^T (P - P_0) Y + Y^T P_0 Y
\end{aligned}
$$

These are the **total sum of squares** (TSS), **residual sum of squares** (RSS), **sum of squares for testing** $H_0$ and the sum of squares due reduction due to $\beta_0$. This can be summarized in an ANOVA (ANalysis Of VAriance) table.

**EPFL**

# Testing in the Least Squares Set–Up

- Then

| Source | d.o.f. | Sum of Squares | Mean squares | F |
|--------|--------|----------------|--------------|---|
| Red | $p - s$ | $\underline{y}^T P_0 \underline{y}$ | | |
| $H_0$ | $s$ | $\underline{y}^T (P - P_0) \underline{y}$ | $M_1 = \dfrac{\underline{y}^T (P - P_0) \underline{y}}{s}$ | $\dfrac{M_1}{M_2}$ |
| Residual | $n - p$ | $\underline{y}^T (I - P) \underline{y}$ | $M_2 = \dfrac{\underline{y}^T (I - P) \underline{y}}{n - p}$ | |
| total | $n$ | $\underline{y}^T \underline{y}$ | | |

- $M_2 = \dfrac{RSS}{n-p}$ is an *unbiased* estimate of $\sigma^2$.
- Reject the null hypothesis at level $\alpha$ if

$$F > f_\alpha$$

where

$$P(F_{s, n-p} > f_\alpha) = \alpha.$$

EPFL

# Assumptions in the Least Squares Set–Up

- Four basic assumptions inherent in the Gaussian linear regression model:

- <u>Linearity</u>: $\mathbb{E}\{Y\}$ is linear in $X$.

- <u>Homoskedasticity</u>: $\mathbb{V}\text{ar}\{\epsilon_j\} = \sigma^2$ for all $j$.

- <u>Gaussian Distribution</u>: errors are normally distributed.

- <u>Uncorrelated Errors</u>: $\epsilon_i$ uncorrelated with $\epsilon_j$ for $i \neq j$.

- When one of these assumptions fails clearly, then Gaussian linear regression is inappropriate as a model for the data.

- Isolated problems, such as outliers and influential observations also deserve investigation. They may or may not decisively affect model validity.

EPFL

# Assumptions in the Least Squares Set–Up

- Scientific reasoning: impossible to validate model assumptions.
- Cannot prove that the assumptions hold. Can only provide evidence in favour (or against!) them.
- Strategy: Find implications of each assumption that we can check graphically (mostly concerning residuals).
- Construct appropriate plots and assess them (requires experience).
- 'Magical Thinking': Beware of overinterpreting plots!

# Outliers

**EPFL**

- An outlier is an observation that does not conform to the general pattern of the rest of the data.
- We *standardise* the residuals through:

$$r_i = \frac{e_i}{\sqrt{s^2(1 - p_{ii})}}$$

where

$$s^2 = \frac{RSS}{n - p}.$$

$s^2$ has $n - p$ degrees of freedom, and note that $r_i$ is *not* student $t$.

# Outliers (more)

- Outliers may be influential: they "stand out" in the "y-dimension".
- However an observation may also be influential because of unusual values in the "x–dimension".
- Such influential observations cannot be so easily detected through plots. But we may wish to automatically detect problems.
- How to find cases having strong effect on fitted model?
- Idea: see effect when case $j$ , i.e., $(x_j^T, Y_j)$ is not kept.
- Let $\boldsymbol{\beta}_{-j}$ be the LSE when model is fitted to data without case $j$ and let $\widehat{Y}_{-j} = X\boldsymbol{\beta}_{-j}$ be the fitted value.

# Outliers (more)

- Define <u>Cook's distance</u>

$$C_j = \frac{1}{ps^2} \left\{ \widehat{Y} - \widehat{Y}_{-j} \right\}^T \left\{ \widehat{Y} - \widehat{Y}_{-j} \right\}.$$

- This measures the scaled distance between the predictions and recall

$$s^2 = \frac{1}{n-p} \| Y - \widehat{Y} \|^2.$$

- It is possible to show that

$$C_j = \frac{r_j^2 \, p_{jj}}{p(1 - p_{jj})},$$

  and thus it can be seen that a large $C_j$ implies and/or large $r_j$ and/or large $p_{jj}$.

# Outliers (more)

- Cases with $C_j > 8/(n - 2p)$ are considered large.
- We therefore plot $C_j$ against $j$ and compare with this cut-off.

EPFL

## Diagnostics

- We plot $Y$ against columns of $X$ to check for linearity and outliers.
- We plot the standardized residuals r against the columns of $X$.
- We plot the standardized residuals r against covariates we left out.
- We plot r against $\hat{Y}$ to check homoscedasticity.
- We make qq plots to check distribution.
- We make the Cook distance plot to check for influential observations.