

Regression II

Sofia Olhede



November 10, 2020

1 Distributional checks

- Leverage
- Weighted Least Squares

Last slide (missing) *had last week ↴*

- Assume that Y_1, \dots, Y_n are i.i.d. $N(\mu, \sigma^2)$.

$$\begin{aligned} P(Y_i \leq y) &= P\left(\frac{Y_i - \mu}{\sigma} \leq \frac{y - \mu}{\sigma}\right) \\ &= \Phi\left(\frac{y - \mu}{\sigma}\right) \end{aligned}$$

Then as $n \rightarrow \infty$

$$F_n(y_i) \rightarrow \Phi\left(\frac{y_i - \mu}{\sigma}\right)$$

(recall from prev chapter that $F_n(y)$ is the Empirical CDF and so

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

Set-up

- We can generalize this to

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

and so as $n \rightarrow \infty$

$$F_n(y) \rightarrow F(y)$$

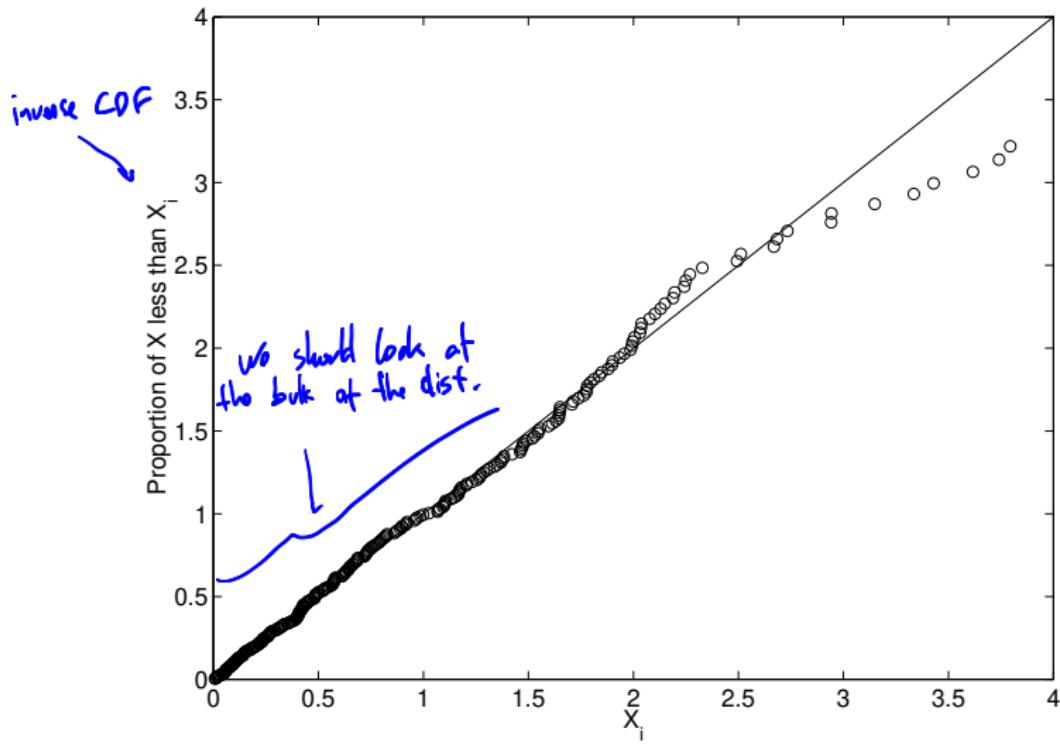
$$\begin{aligned} F^{-1}(F_n(y_i)) &\approx y_i \\ F^{-1}(\text{prop. of obs. } \leq y_i) &\approx y_i. \end{aligned}$$

*↓
proportion of observations*

"Gormorov-Smirnov (?) test is useless because it is not powerful enough, because of the asymmetri because NULL and the other hypothesis. Also residuals are not a good idea because they're not iid. Q-Q plots is much better, do Q-Q plot" (next slide)

QQ-plot

A Q-Q plot is a scatterplot created by plotting two sets of quantiles against one another. If both sets of quantiles came from the same distribution, we should see the points forming a line that's roughly straight.



Set-up

- Note that

$$\text{Var}\{\mathbf{e}\} = \sigma^2(\mathbf{I}_n - \mathbf{P})$$

hat matrix
↑

?

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

$$p_{ii} > \frac{2p}{n}.$$

0

Weighted Least Squares

- Consider the linear model

$$\mathbb{E}\{Y_i\} = x_i \beta$$

and

$$\text{Var}(Y_i) = \frac{\sigma^2}{w_i},$$

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.

In statistics, a vector of random variables is heteroscedastic if the variability of the random disturbance is different across elements of the vector. Here, variability could be quantified by the variance or any other measure of statistical dispersion.

Heteroscedasticity is the absence of homoscedasticity. A typical example is the set of observations of income in different cities.

The existence of heteroscedasticity is a major concern in regression analysis and the analysis of variance, as it invalidates statistical tests of significance that assume that the modelling errors all have the same variance. While the ordinary least squares estimator is still unbiased in the presence of heteroscedasticity, it is inefficient and generalized least squares should be used instead.^{[5][6]}

In this case: $w_i \rightarrow \infty \Rightarrow \text{Var is small}$
 $w_i \rightarrow 0 \Rightarrow \text{Var is huge}$ 

Weighted Least Squares

- Find an estimate for β by minimising the weighted sum of squares:

$$\begin{aligned} S(\beta) &= \sum_{i=1}^n w_i (Y_i - x_i^T \beta)^2 \\ &= \sigma^2 \sum_{i=1}^n \frac{(Y_i - E\{Y_i\})^2}{\text{Var}\{Y_i\}} \end{aligned}$$

Weighted least squares (WLS), also known as weighted linear regression, is a generalization of ordinary least squares and linear regression in which the errors covariance matrix is allowed to be different from an identity matrix.

Weighted Least Squares

- In vector form we then have:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{D}\boldsymbol{\epsilon}$$

new matrix D (diagonal)

$$\text{Var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_n$$

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

and

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

Weighted Least Squares

- We then multiply through the linear equation by

$$\begin{aligned}
 D^{-1}Y &= D^{-1}X\beta + \epsilon \\
 \tilde{Y} &= \tilde{X}\beta + \epsilon \\
 \text{Var}\{\epsilon\} &= \sigma^2 I_n \\
 \mathbb{E}\{\epsilon\} &= 0
 \end{aligned}
 \tag{1}$$

$D^{-1}Y = \sigma^2(X\beta + D\epsilon)$
 $\Rightarrow D^{-1}X\beta + \epsilon$
 replace $\tilde{Y} = D^{-1}Y$ and $\tilde{X} = D^{-1}X$
 } from before

This is recognizable as a liner model. The β estimate is given by

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

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

Testing in the Least Squares Set-Up

- Assume that $\mathbf{Y} \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$.
- Definition: If $\mathbf{Z} \sim N(\boldsymbol{\mu}, \mathbf{I}_n)$ where $\boldsymbol{\mu} \neq 0$ then $\mathbf{Z}^T \mathbf{Z}$ is said to have a non-central χ^2 distribution on n d. o. f. and non-centrality parameter $\delta > 0$ given by $\delta^2 = \boldsymbol{\mu}^T \boldsymbol{\mu}$. degrees of freedom follows from def. of χ^2 dist.
- otherwise: if $\boldsymbol{\mu} = 0 \Rightarrow \chi_n^2$ as the non-centrality parameter is then zero.
- We normally for the general distribution write $U = \mathbf{Z}^T \mathbf{Z} \sim \chi_n^2(\delta)$.
- The distribution of $\mathbf{Z}^T \mathbf{Z}$ depends on $\boldsymbol{\mu}$ only via δ .
- $\mathbb{E}(U) = n + \delta^2$. sum of squares
- $\text{Var}(U) = 2n + 4\delta^2$
- If $U_i \sim \chi_{n_i}^2(\delta_i)$ for $i = 1, \dots, k$ and if the $\{U_i\}$ are all independent then

i.e. if all U_i are chi-squares
then sum is also chi-square with
a degree of freedom equal to sum of all
d.o.f. and a new non-centrality parameters
where $n = \sum_{i=1}^k n_i$ and $\delta^2 = \sum_{i=1}^k \delta_i^2$.

Testing in the Least Squares Set–Up

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

In linear algebra, an idempotent matrix is a matrix which, when multiplied by itself, yields itself.
i.e. $A = A^2$

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

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

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

Testing in the Least Squares Set-Up

- Furthermore it follows that

$$\begin{aligned}
 Z^T AZ &= (PV)^T A(PV) \\
 &= V^T P^T A P V \quad \xrightarrow{\text{P}^T A P = D \text{ from prev slide}} \\
 &= V^T D V \quad \xrightarrow{\text{D is diagonal and only zeros or ones}} \\
 &= V^T D^T D V \quad \hookrightarrow D = D^T D \text{ e.g. } \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \\
 &= (DV)^T (DV) \\
 &= \text{sum of squares of } r \text{ components} \\
 &= \chi_r^2(d),
 \end{aligned}$$

for some d . In fact

→ $d^2 = \mathbb{E}(DV)^T \mathbb{E}(DV) = (DP^T \mu)^T (DP^T \mu) = \mu^T P D P^T \mu = \mu^T A \mu.$

Testing in the Least Squares Set–Up

- Lemma: If $Z \sim N(\mu, I_n)$ 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

$$\text{Cov}\{A_1 Z, A_2 Z\} = A_1 \text{Cov}(Z) A_2^T \quad (3)$$

$$= A_1 I A_2 = 0. \quad (4)$$

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, \dots, 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.

for all i in 1,...,k

Testing in the Least Squares Set-Up

Lemma: If A_1, \dots, A_k are symmetric $n \times n$ matrices such that $\sum A_i = I$ and such that $\text{rank}(A_i) = r_i$ 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 $i = 1, \dots, k$.

ie if one property holds,
the other two will hold too!

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

$$Z^T A_i Z$$

*→ this will be important later.
We will find many relation in this form*

are independent *distributions of the type:*

$$\chi_{r_i}^2(\delta_i)$$

*↓
e.g. if we take quadratic form of Z
i.e. $Z^T Z$ and we multiply on idempotent
matrix, we end up with many χ^2 dist.*

$$\text{where } \delta_i^2 = \mu^T A_i \mu.$$

Testing in the Least Squares Set-Up

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

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

Because they are mutually orthogonal by assumption this implies that

$$Z^T A_i Z$$

are independent.

Testing in the Least Squares Set–Up

- Assume we want to test the hypothesis

$$H_0 : A\beta = 0$$

versus

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

where $\text{rank}(A) = s = p - p_0$.

- Under H_0 we get the simpler linear model

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

where β_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

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

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

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

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

→ NULL model
→ alternative model

- Consider a monotonic increasing function of t : → take to power $\frac{2}{n}$ and subtract 1:

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

?

or

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

?

Testing in the Least Squares Set-Up

- Use the Fisher-Cochran theorem:

$$\mathbf{I}_n = (\mathbf{I}_n - \mathbf{P}) + (\mathbf{P} - \mathbf{P}_0) + \mathbf{P}_0$$

7
o

with ranks

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

Let

$$\mathbf{A}_1 = (\mathbf{I}_n - \mathbf{P})$$

$$\mathbf{A}_2 = (\mathbf{P} - \mathbf{P}_0)$$

$$\mathbf{A}_3 = \mathbf{P}_0$$

These are are symmetric and idempotent.

Testing in the Least Squares Set-Up

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

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

\mathcal{I}_0

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

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

and

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

independently (the non-centrality parameters vanish.)

Testing in the Least Squares Set-Up

- We then have

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

- 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 β_0 . This can be summarized in an ANOVA (ANalysis Of VAriance) table.

Testing in the Least Squares Set-Up

ANOVA table:

P_0 is the new hat matrix, slide 17 (?)

Source	d.o.f.	Sum of Squares	Mean squares	F statistics
Red reduction	$p - s$	$\underline{y}^T P_0 \underline{y}$		
H_0	s	$\underline{y}^T (P - P_0) \underline{y}$	$M_1 = \frac{\underline{y}^T (P - P_0) \underline{y}}{s}$	$\frac{M_1}{M_2}$
Residual	$n - p$	RSS: $\underline{y}^T (I - P) \underline{y}$	$\textcolor{blue}{M_2} = \frac{\underline{y}^T (I - P) \underline{y}}{n - p}$	
Σ total	n	$\underline{y}^T \underline{y}$		

(1)

• $M_2 = \frac{RSS}{n-p}$ is an *unbiased estimate of σ^2* .

• Reject the null hypothesis at level α if

sum of squares (final)
= unexplained/residual sum of sq.
on the full model

+ sum of sq for testing $\underline{y}^T (P - P_0) \underline{y}$

+ sum of reduction due to P_0 $\underline{y}^T P_0 \underline{y}$
(also sum of $\underline{y}_1^T \underline{y}_1$)

$y^T y$

$= \underline{y}^T (P_0 + P - P_0 + I - P) \underline{y} = \underline{y}^T y$

$$F > f_\alpha$$

a quantile for the F-distribution

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

where

NOTE: this is one of the most used methods for linear models in the world. Here we assume y is Gaussian. Later, we'll learn generalized linear models that allow y not to be Gaussian!

Assumptions in the Least Squares Set-Up

↓
not weighted



- Four basic assumptions inherent in the Gaussian linear regression model:
- Linearity: $\mathbb{E}\{Y\}$ is linear in X .
- Homoskedasticity: $\text{Var}\{\epsilon_j\} = \sigma^2$ for all j . (before we mentioned heteroscedasticity)
- Gaussian Distribution: errors are normally distributed. →
- Uncorrelated Errors: ϵ_i uncorrelated with ϵ_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.

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

- 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 is the standardized residual
 $= \text{residual divided by estimation of variance}$

$$r_i = \frac{e_i}{\sqrt{s^2(1 - p_{ii})}} \rightarrow \chi^2 \text{ dist.}$$

where

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

χ^2 dist.

Residuals we see before (Gaussian)

variance of residuals (?)

p_{ii} is taken from hat matrix

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 β_{-j} be the LSE when model is fitted to data without case j and let $\hat{Y}_{-j} = X\beta_{-j}$ be the fitted value.



Outliers (more)

- Define Cook's distance

$$\rightarrow C_j = \frac{1}{ps^2} \left\{ \hat{Y} - \hat{Y}_{-j} \right\}^T \left\{ \hat{Y} - \hat{Y}_{-j} \right\},$$

number of model terms \downarrow
 variance

- This measures the scaled distance between the predictions, and recall that:

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

variance

- It is possible to show that

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

hat matrix element

and thus it can be seen that a large C_j implies either 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.

very nice / useful slide

Diagnostics

How do we diagnose if a model appears to be fitting well?



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

(last slide missing)