

Recap: Violations of Assumptions

1) Bids

2) Non-Normality
- Outliers

$$y \sim N(Z\beta, \sigma^2 V)$$

3) Now: Heteroscedasticity
(non constant variance)

$$V \neq I$$

Heteroscedasticity

Suppose $y \sim N(Z\beta, \sigma^2 V)$

V is full rank, not the identity

- Example 1: $V = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$, as is the case
where different amounts of data goes into
each measurement

- Example: AR(1):

$$\begin{pmatrix} 1 & p & 0 \end{pmatrix}$$

$$V = \begin{pmatrix} p & i & p \\ 0 & p & i \end{pmatrix} \quad n=3$$

- If we know V , we can use generalized LS:

$$V = P^T \Lambda P \quad P^T P = P P^T = I$$

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$$

let: $D := \Lambda^{-\frac{1}{2}} P$

- we have: $\tilde{D} \tilde{y} = D(Z\beta + \varepsilon) = \tilde{D} Z \beta + \tilde{D} \varepsilon$

$$\boxed{\tilde{y} = \tilde{Z} \beta + \tilde{\varepsilon}}$$

Where: $\text{Var}(\tilde{\varepsilon}) = D \overset{V}{\text{Var}(\varepsilon)} D^T$
 $= D P^T \Lambda P D$
 $= \underbrace{\Lambda^{-\frac{1}{2}} P P^T}_{I} \underbrace{\Lambda}_{I} \underbrace{P P^T}_{I} \Lambda^{-\frac{1}{2}} = I$

- we take:

$$\hat{\beta}_{GLS} = (\tilde{Z}^T \tilde{Z})^{-1} \tilde{Z}^T \tilde{y}$$

so that $\hat{\beta}_{GLS}$ minimizes

$$\|Dy - DZ\beta\|^2$$

- If $V = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$

then we can use

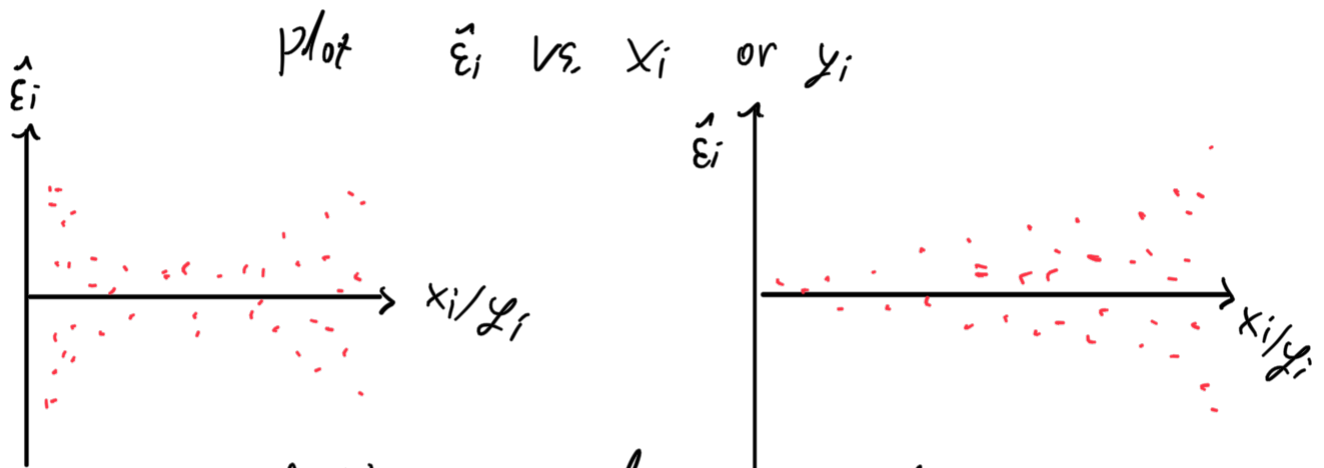
$$\hat{\sigma}_i^2 = \hat{u}_i^2, \quad \hat{u}_i = y_i - z_i^T \hat{\beta}_{OLS}$$

..... 1 + -1 +

Huber-White residuals

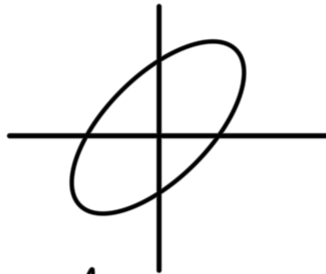
$$\hat{\beta}_{OLS} = (Z'Z)^{-1}Z'y$$

Detection



(things we don't want to see)

- Another approach is to plot $\hat{\epsilon}_i$ vs $\hat{\epsilon}_{i+1}$
a pattern like

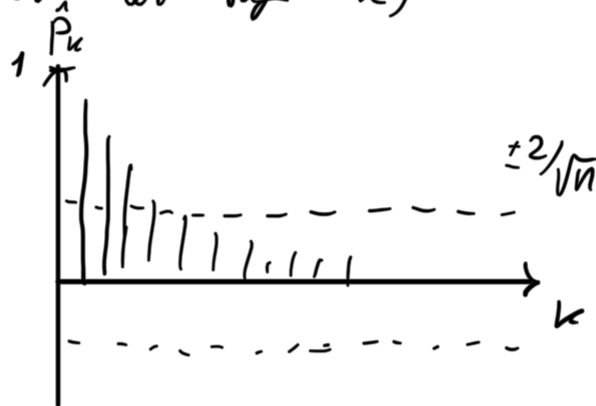


indicates dependence in errors (like in AR(1))
say

- This issue needs to be addressed using time-series methods
- While $\hat{\beta}$ is still unbiased, we get erroneous variance so we cannot do proper statistical tests.
- We can also compute the correlation between the residuals:

$$\hat{\rho}_k = \frac{1}{n} \sum_{i=k+1}^n \hat{\epsilon}_i \hat{\epsilon}_{i-k}$$

(auto correlation at lag k)



- under the null hypothesis when there is no correlation between the residuals,
 $P_k = 0, k \geq 1$
- Tests for non-zero correlation:
 Durbin-Watson & Ljung-Box ...
- See class on time-series analysis
 Stanford STATS 207

Course Summary

Additional Topics

- ANOVA with random effects:

$$Y_{ij} = \mu + \alpha_i + \epsilon_{ij}, \quad \epsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$i = 1, \dots, k$$

$$j = 1, \dots, n_i$$

$$\alpha_i \stackrel{iid}{\sim} N(0, \sigma_A^2)$$

- Two-Factor ANOVA:

$$Y_{ijk} = \mu_{ij} + \varepsilon_{ijk}, \quad k=1, \dots, N_{ij}, \quad i=1, \dots, I \\ j=1, \dots, J$$

$$\mu_{ij} = \alpha_i + \beta_j + \gamma_{ij} \quad 0 = \sum_i \alpha_i = \sum_j \beta_j = \sum_i \gamma_{ij} = \sum_j \gamma_{ij}$$

Also: Fixed x Fixed, Fixed x Random
Random x Random

- Causality

- Bootstrapped Regression:

- Pairs: resample n pairs from $(x_i, y_i)_{i=1}^n$ with repetition

- Residuals: resample n times from $\{\hat{\varepsilon}_i\}_{i=1}^n$
use $y_i^* = z_i^T \hat{\beta} + \hat{\varepsilon}_i$

Course Overview

- Applied Stats & Linear Model
- We have (x_i, y_i) pairs & we want to model how y depends on x
- For example: we want to predict y_{n+1} from x_{n+1}
- Also: how things generalize,

What can go wrong?

The Pipeline

- (1) Data. $\{(x_i, y_i)\}_{i=1}^n$, $y_i \in \mathbb{R}$, but x_i can be 1 group, k groups, $k \times n$ groups, or \mathbb{R}^d
- (2) Features. Z , e.g. $z_i^T = (1, x_{i1}, \dots, x_{id})$
 - Can also add-in non-linear features like x_{i3}^2 or $1(x_{i2} \leq 7)$
 - The inclusion/exclusion of features can be based on intuition, experience, science, or machine learning.
- (3) Model. $y \sim N(z\beta, \sigma^2 I)$
- (4) Fitting. Estimate $\hat{\beta}$ and $\hat{\sigma}^2$ (based on linear algebra)
- (5) Inference. Derive confidence intervals, p-values hypothesis testing, Power calculations (based on normal dist. theory)
- (6) Interpretation. Association between predictors and response. Concerns about causality and the meaning of a "true" β_j , because it depends on whether the k -th predictor is in the model or not.
- (7) Model Selection. Using methods like AIC, BIC, CV, lasso, ridge regression, optimized for prediction, but not necessarily for understanding of causal mechanisms.

(8) Problems & Fixes

- Non-normality: CLT fixes most cases
- Non-constant variance: Bootstrapping or Huber-White residuals in GLS
- = Bias. No good solutions, context and experience on topic are useful. Lack of fit sum of squares may help identify bias but can do little to solve it.
- Outliers: some methods exist but there are tough to overcome in high dimensions. If both y_i & z_i are outliers, it is a good indication that this sample should be removed.
- Correlations within residuals. Time series methods

Followups

- Generalized Linear Models (GLM):

$$E(y_i | z_i) = \psi(z_i^T \beta) \quad \psi: \mathbb{R} \rightarrow \mathbb{R}$$

- Look at $y_i \in \{0, 1\}$, $y_i \in \{0, 1, 2, \dots, k\}$

$$y_i \in \mathbb{R}^p$$

- More computationally sophisticated algorithms than LS:

$$E(y_i | z_i) = \psi_{\beta}(z_i^T \beta)$$

- Causal conclusions from observations
- Time series analysis (dependency between (x_i, y_i) over i)

GLM

$$\varphi: \mathbb{R} \rightarrow \mathbb{R}$$

$$E(Y_i | Z_i) = \varphi(Z_i^T \beta)$$

φ is called the "link function"

Examples:

$$\varphi(x) = x$$

$$\varphi(x) = \frac{1}{1 + e^{-x}}$$

$$\varphi(x) = \Phi(x)$$

$$\varphi(x) = e^x$$

logit model

probit model

Poisson Regression

Logist Regression: $\varphi(x) = \frac{1}{1 + e^{-x}}$

- we have $y_i \in \{0, 1\}$

- we continuous predictors x_1, \dots, x_n $x_i \in \mathbb{R}^p$

- we propose that the log-likelihood ratio between the classes is of the form:

$$\left(\log \left(\frac{\Pr(Y_i = 1 | x_i)}{\Pr(Y_i = 0 | x_i)} \right) \right) = x_i^T \beta \in \text{logit}$$

of $\Pr(Y_i=0|x_i)$ square

$$\text{So that } p = p(\beta, x_i) = E(Y_i|x_i) = \Pr(Y_i=1|x_i) \\ = \frac{1}{1 + e^{-x_i^T \beta}}$$

- The log-likelihood function:

$$l(\beta; \{y_i, x_i\}) = \sum_{i: y_i=1} \log(p(\beta; x_i)) + \sum_{i: y_i=0} \log(1-p(\beta; x_i))$$

- $l(\beta; \{y_i, x_i\})$ is convex, has a unique minimum