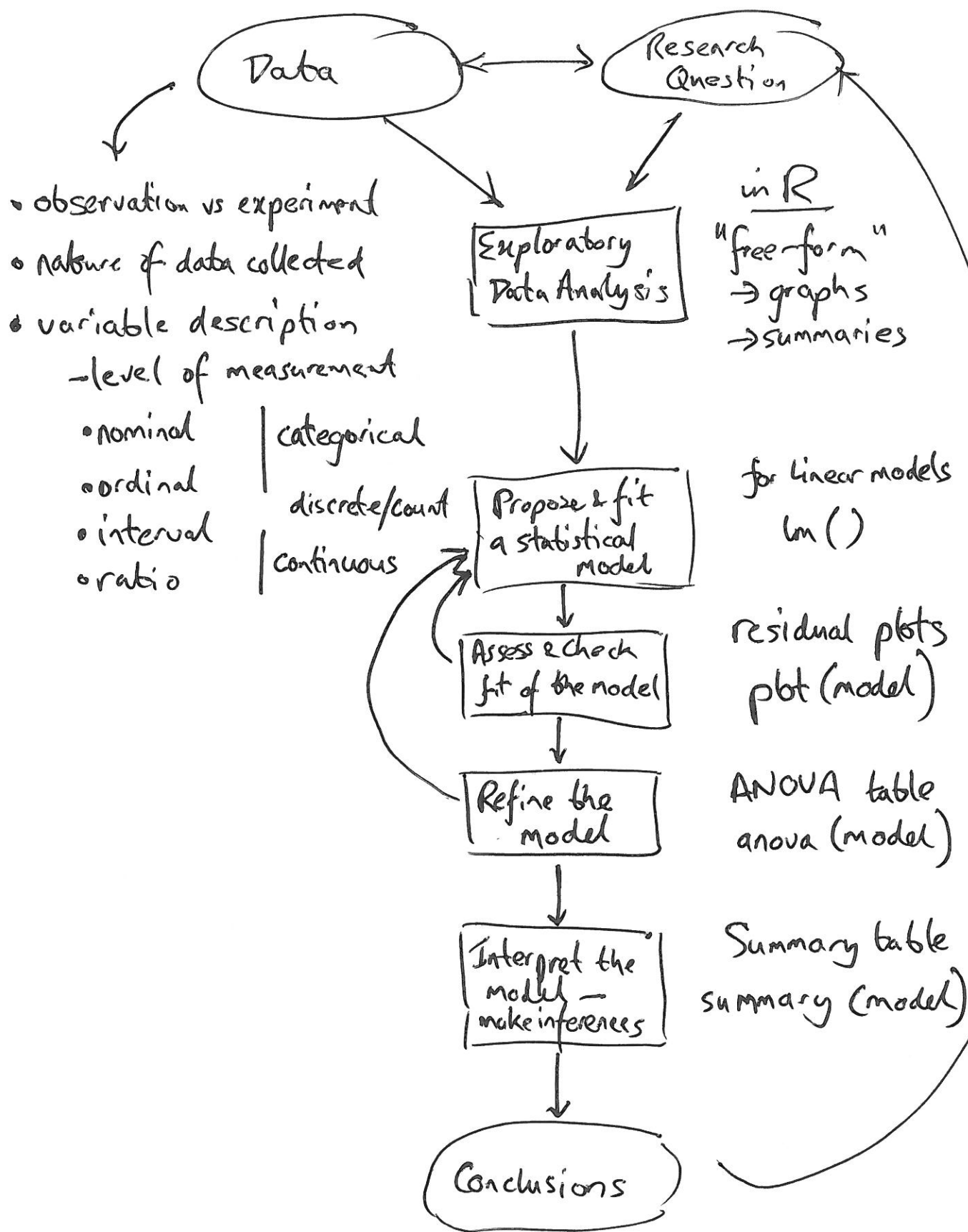
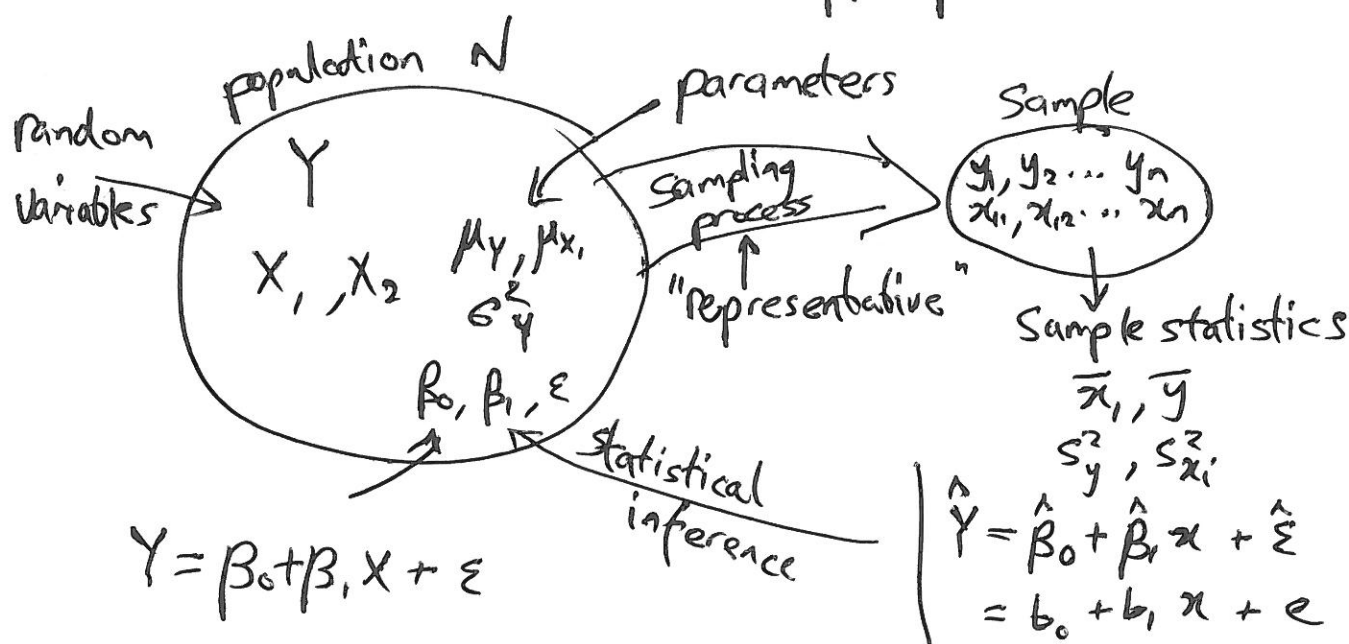


# The modelling process (as part of a research process)



# General Assumptions (underlying all statistical models)

- Sample of data is representative of the population of interest (as defined by the research question)
- that all the error is in the direction of the response or unknown variable  $Y$ , i.e. that the  $X$ 's are "known" (or measured with negligible error in comparison to  $Y$ )
- chosen model is appropriate



Notation capital letters  $Y, X_1, X_2$  denote random variables

Small letters  $y, x_1, \dots$  denote sample observations

Greek letters  $\mu, \pi, \sigma$  denote population parameters

Roman equivalents  $m, p, s$  denote sample statistics which are typically point estimates of the equiv.

population parameters — though we can use  $\hat{\phantom{x}}$  for this purpose

$$\hat{\mu}_x = \underset{\substack{\uparrow \\ \text{not used}}}{m_x} = \bar{x}, \quad \hat{\pi} = p, \quad \hat{\beta}_0 = b_0, \quad \hat{\beta}_1 = b_1$$

Assumptions underlying Regression Models:

Multiple Linear Regression Model or General Linear Model

$$Y = \underbrace{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k}_{\text{deterministic part}} + \underbrace{\varepsilon}_{\text{stochastic part}}$$

$E[Y | X_1, X_2 \dots X_k]$       [probability model]  
 mean model      variance model

Model-specific assumptions are about the variance model

$$\varepsilon's \stackrel{iid}{\sim} N(0, \sigma^2)$$

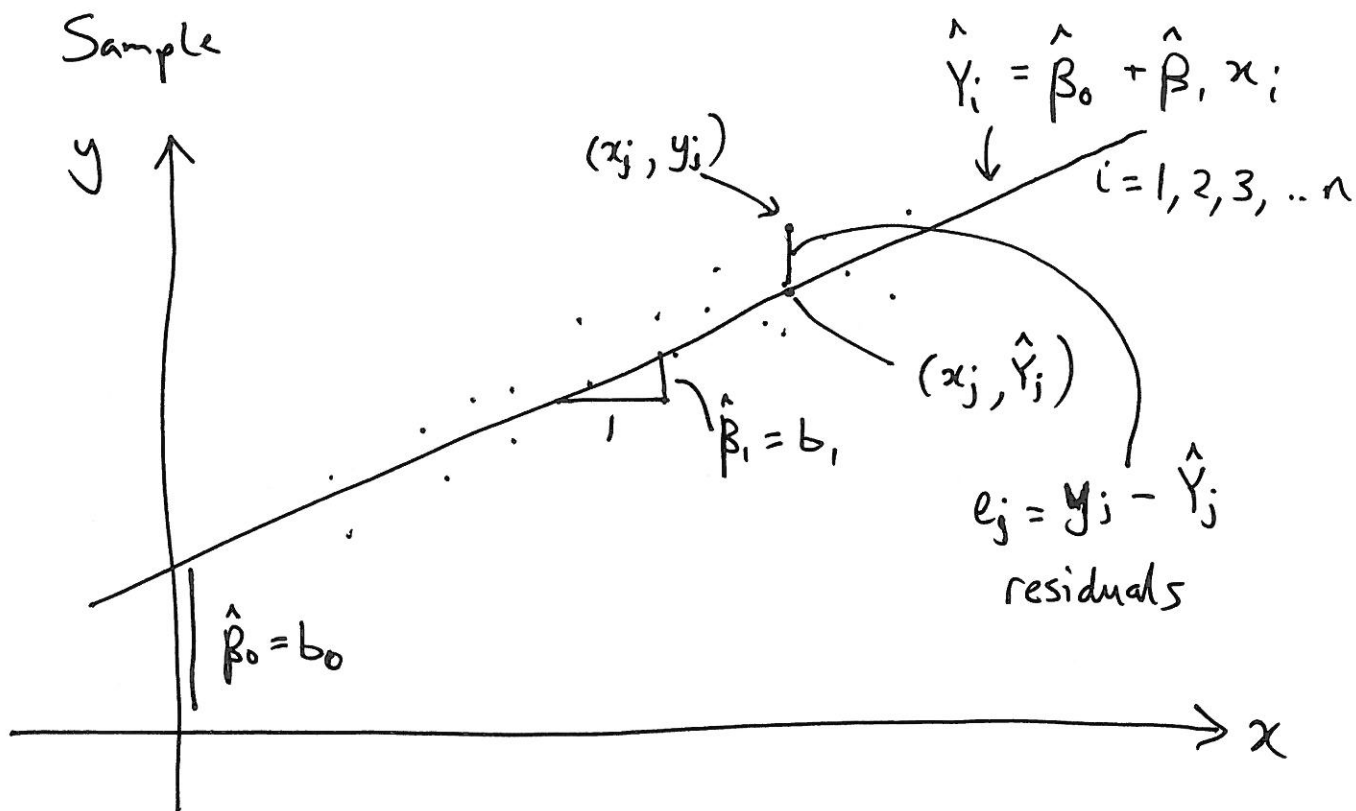
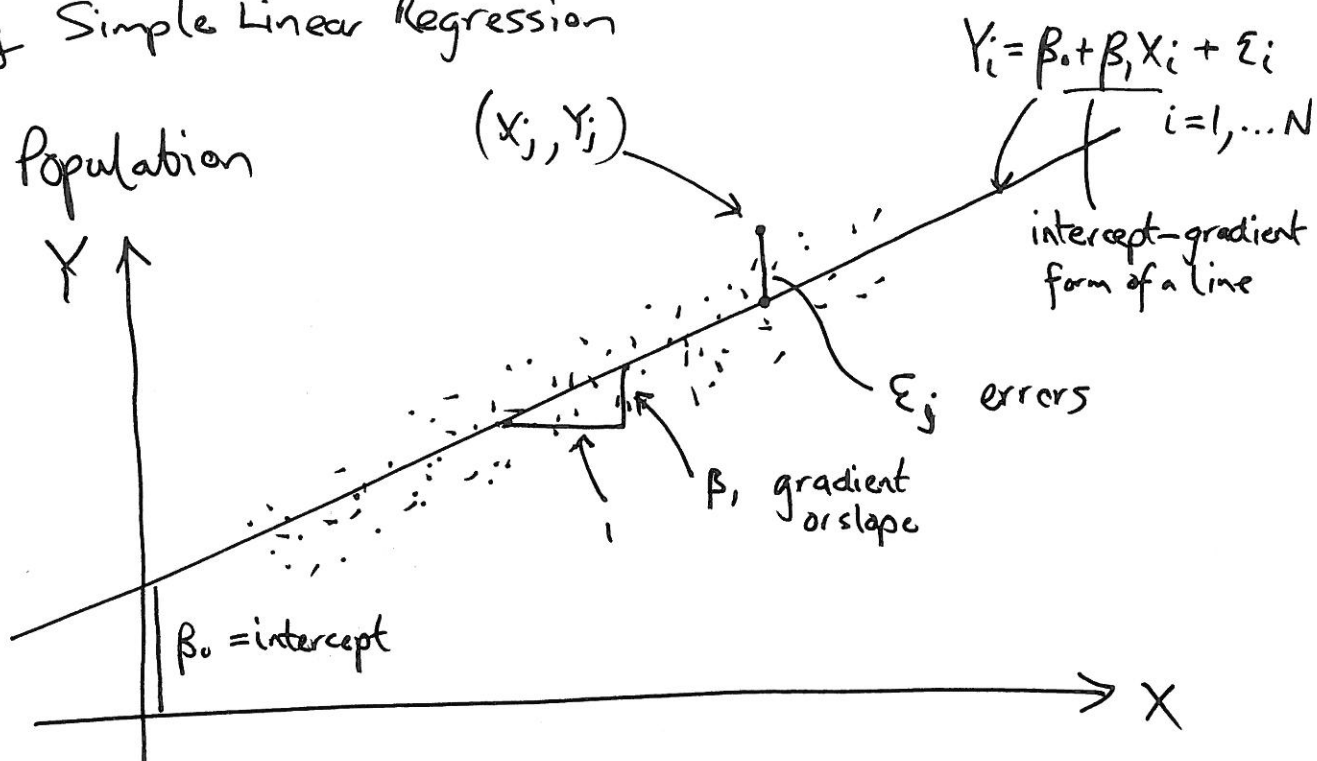
The errors are independent & identically (normally) distributed with constant variance

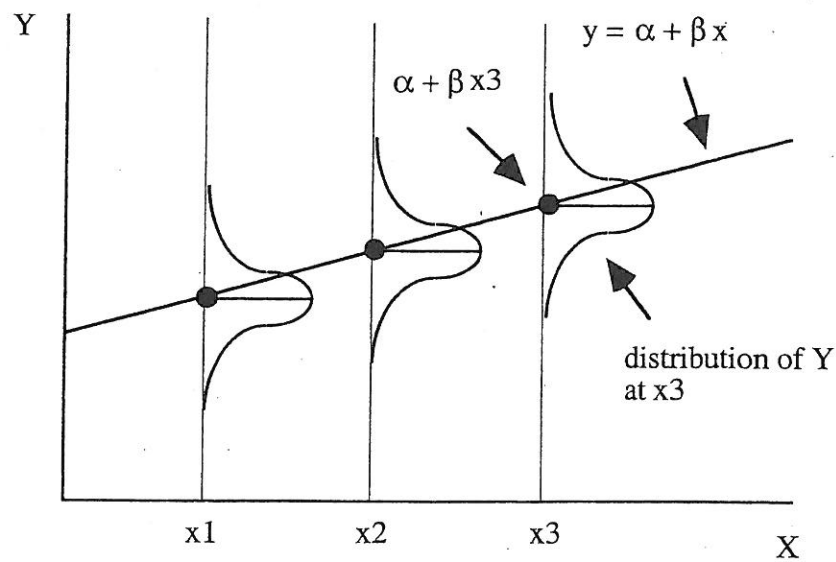
→ these assumptions are all about the errors, but we assess them using the residuals, which are the observed or estimated errors

Key assumptions are:

1. independence (no systematic patterns in the residuals)
2. constant variance (no systematic changes in the variance of the residuals)
3. normality (residuals are normally distributed → needs a qq plot)

eg Simple Linear Regression





Assumption of constant variance  
of the errors (i.e. same distribution with same  $\sigma^2$   
for all values of  $X$ )

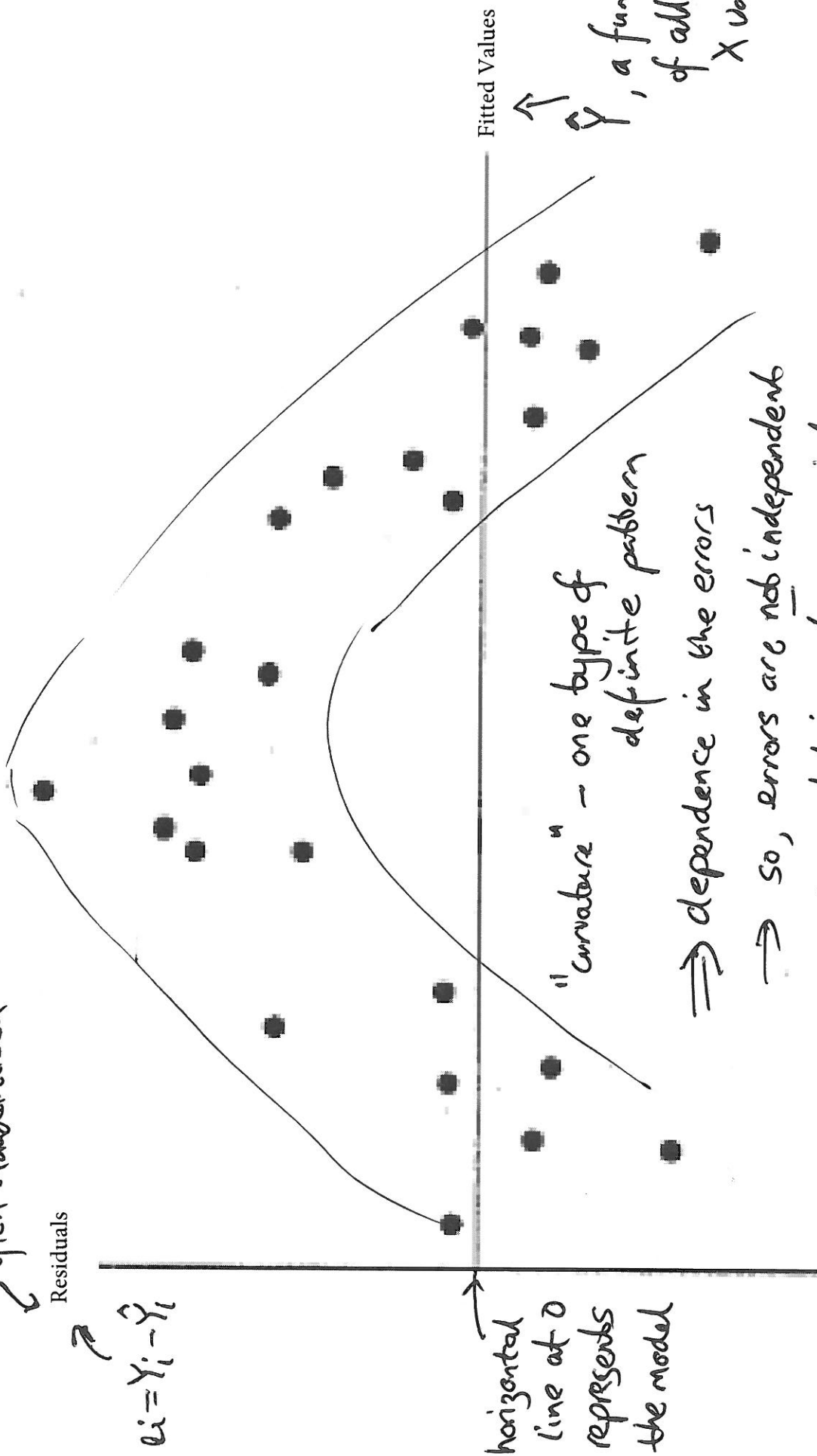
## Main residual plot

Plot I

often standardised

Residuals

$$e_i = Y_i - \hat{Y}_i$$



horizontal line at 0 represents the model

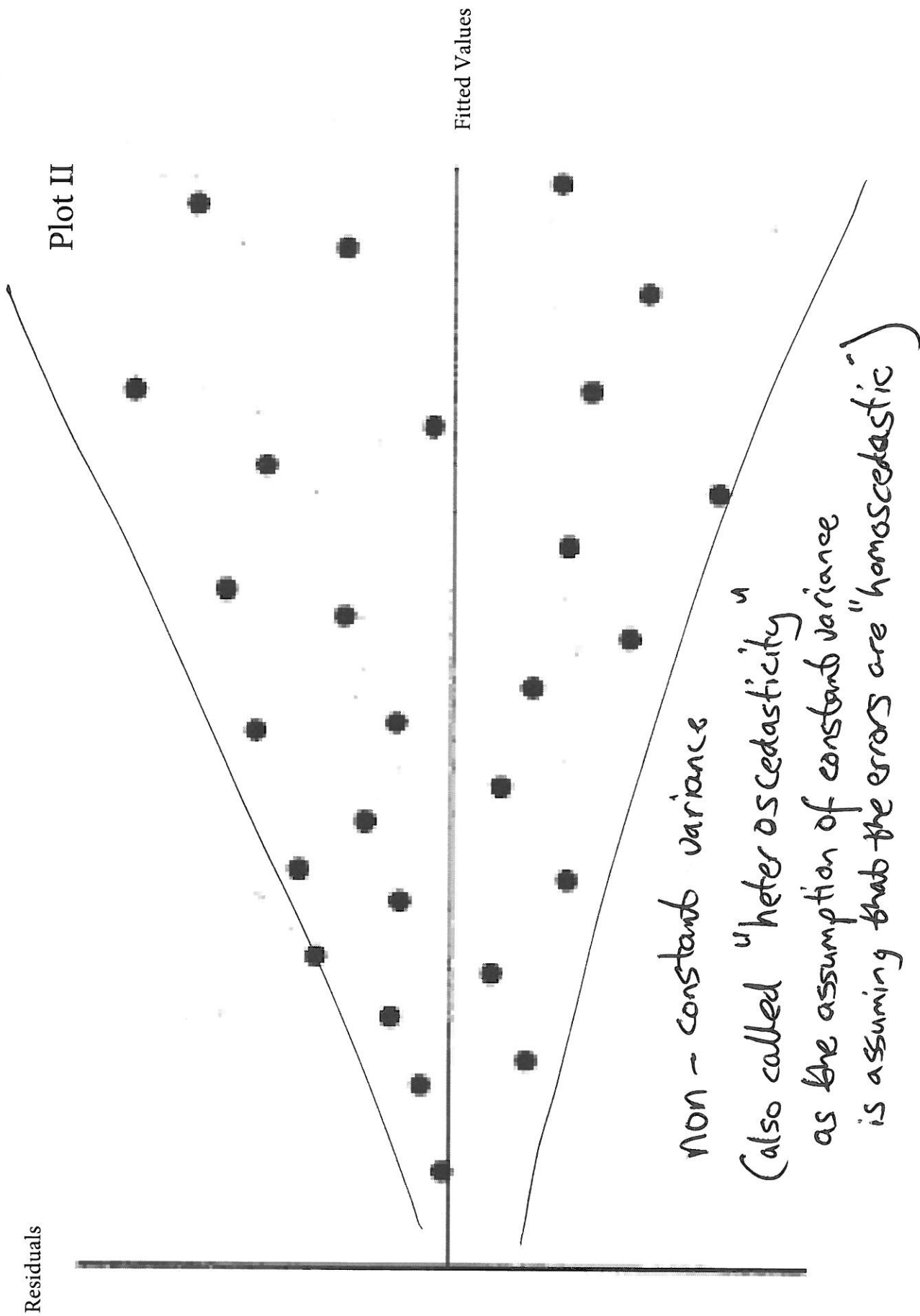
"curvature" - one type of definite pattern

⇒ dependence in the errors

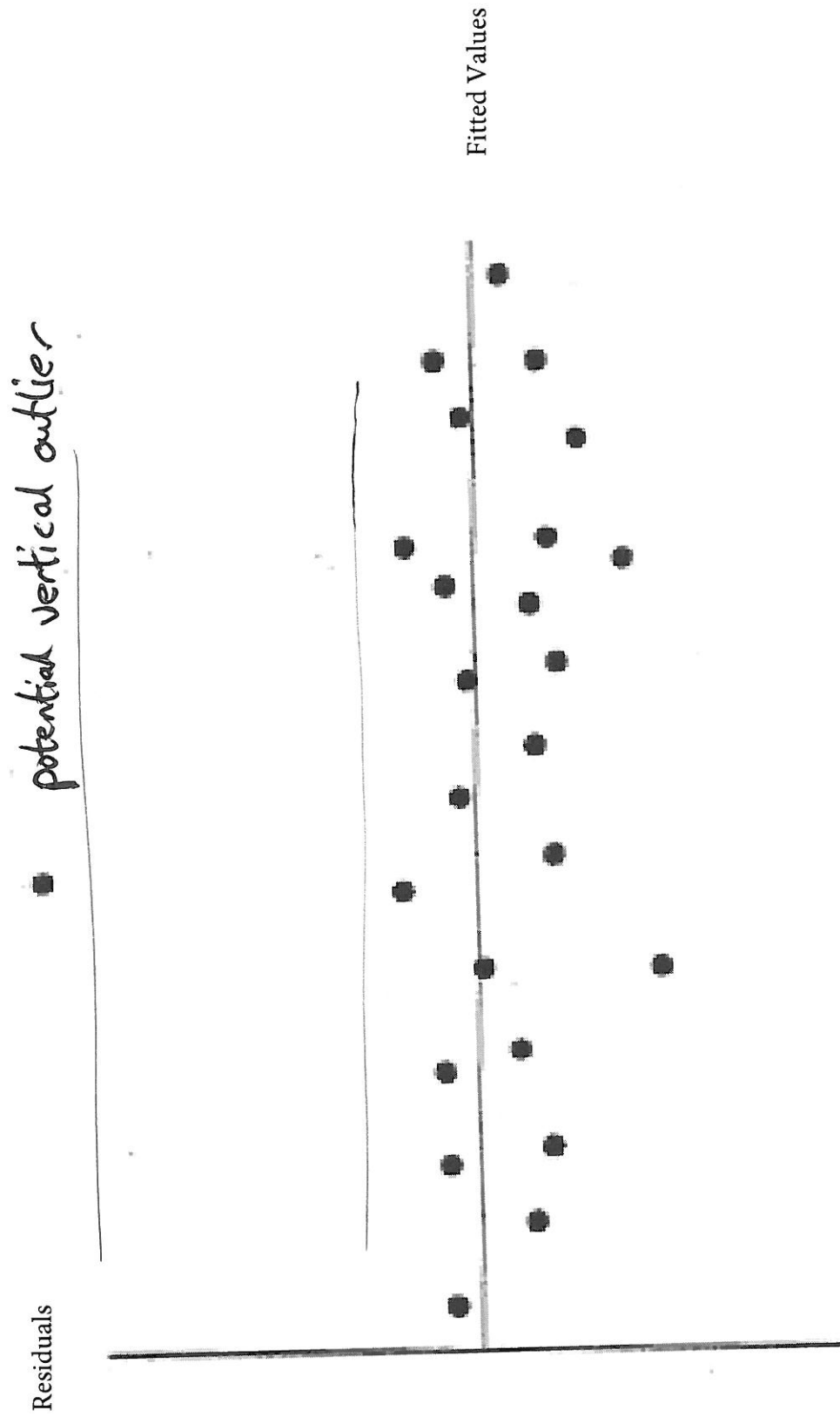
⇒ so, errors are not independent

& model is not appropriate

$\hat{Y}$ , a function of all the  $X$  variables

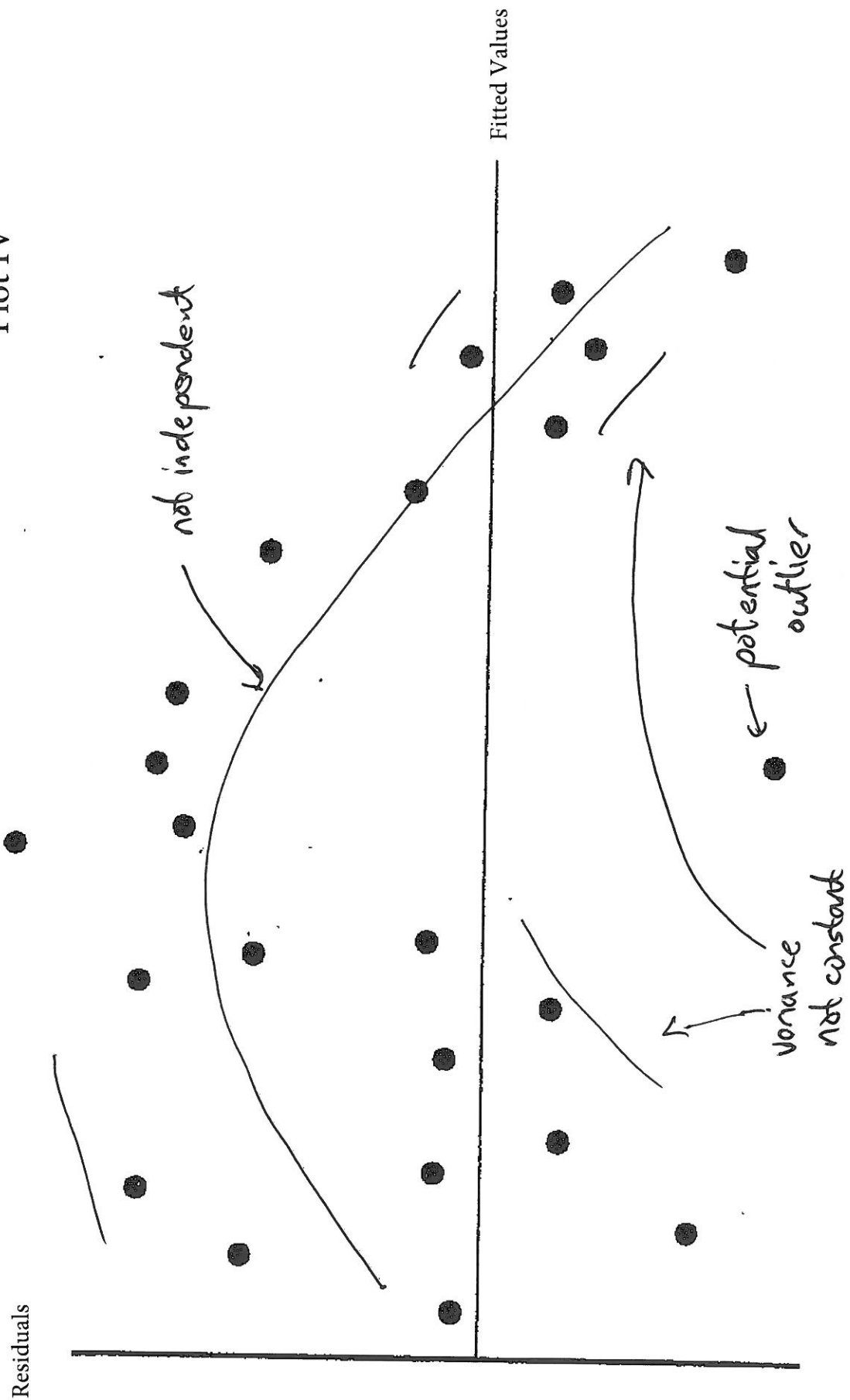


Plot III

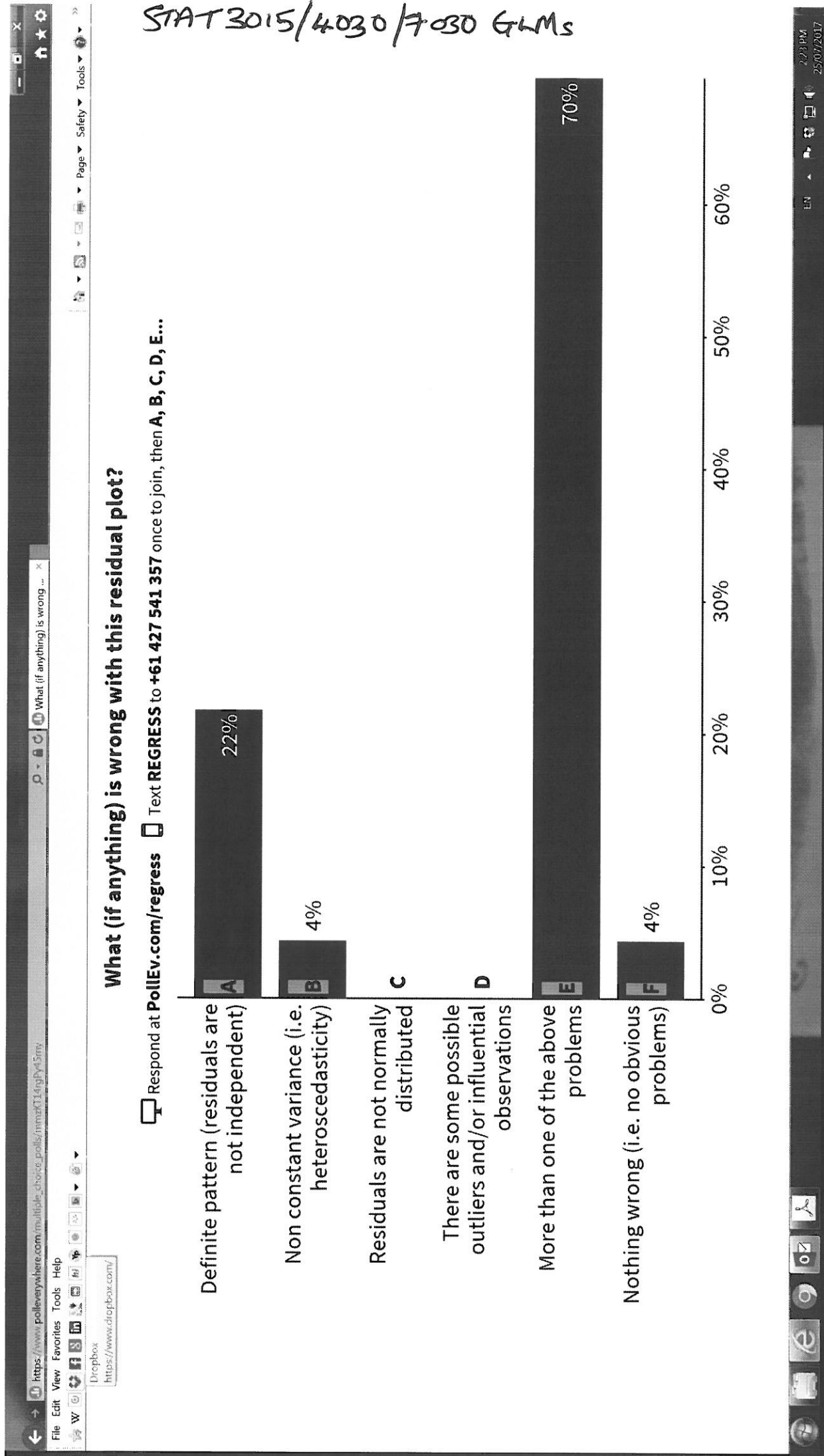




Plot IV

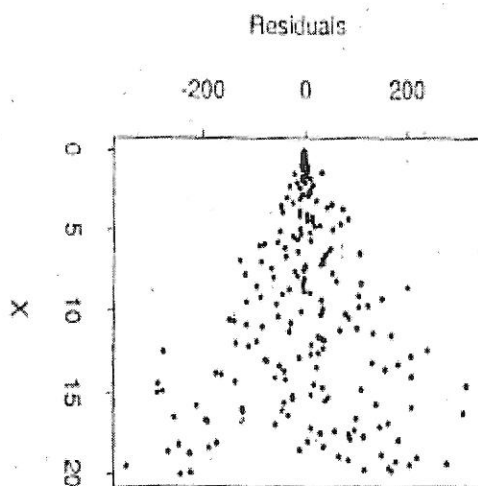


# Results of Poll Everywhere poll for Plot IV

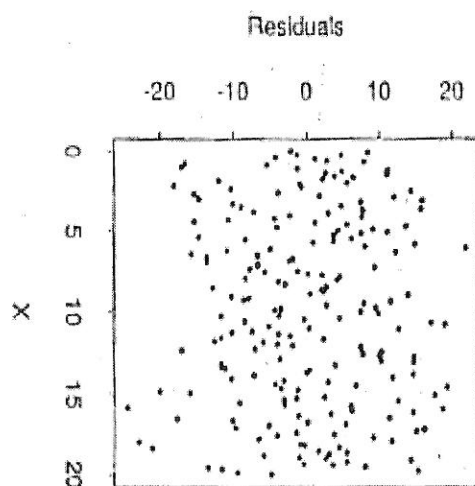


25/7  
(10)

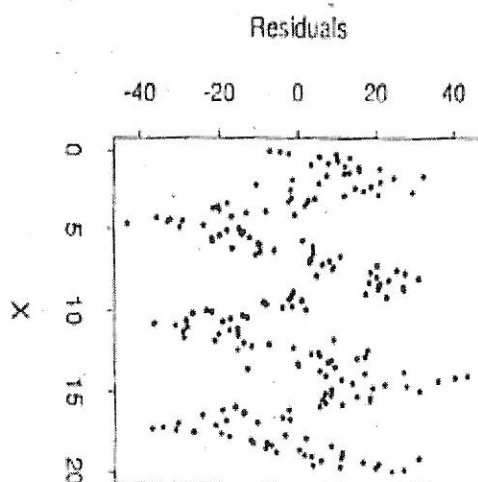
More examples for you to try:



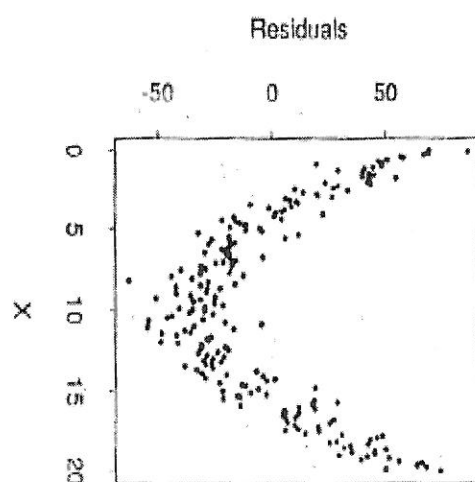
Plot VII



Plot V

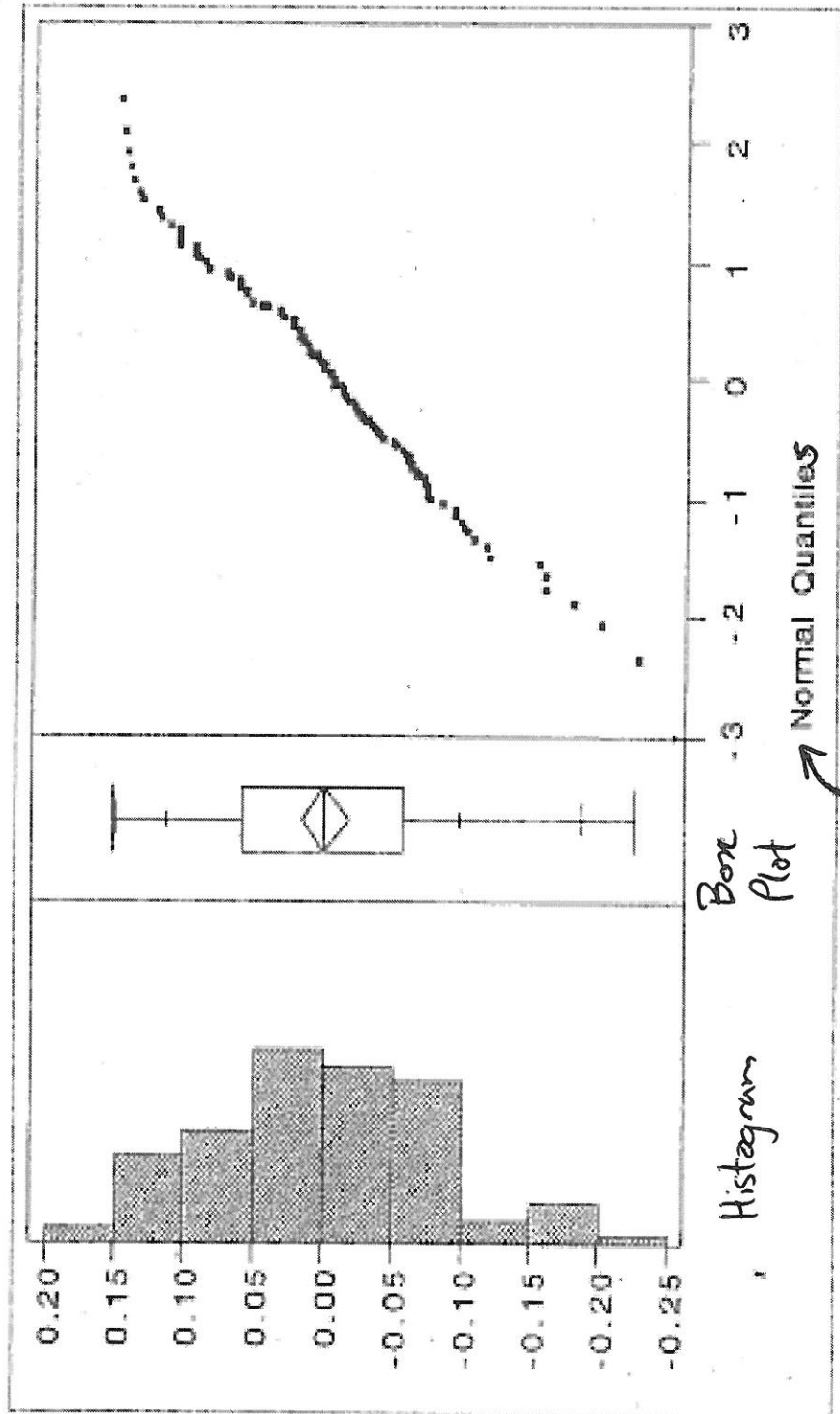


Plot VIII

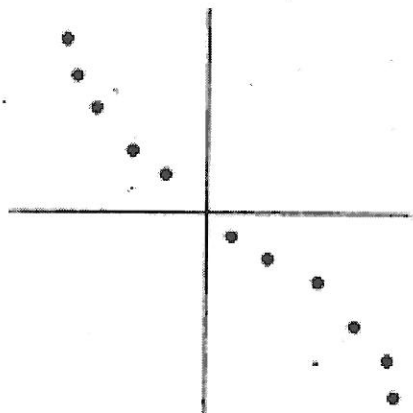


Plot VI

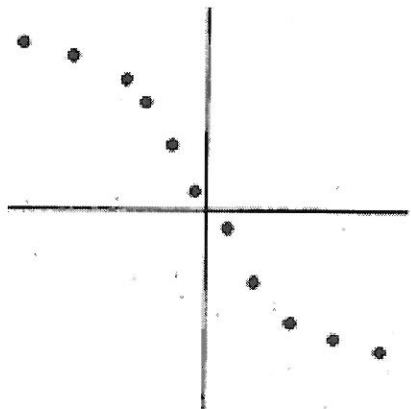
Assessing normality (only assumption not tested by main residual plot)



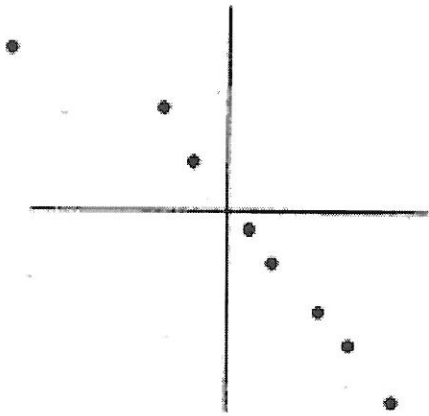
Normal quantile (qq) plot using equivalent quantiles from a standard normal distribution



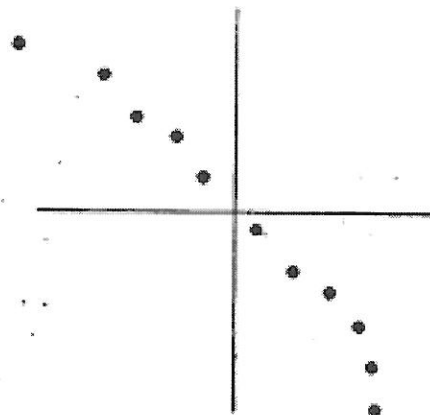
Shorter Tails



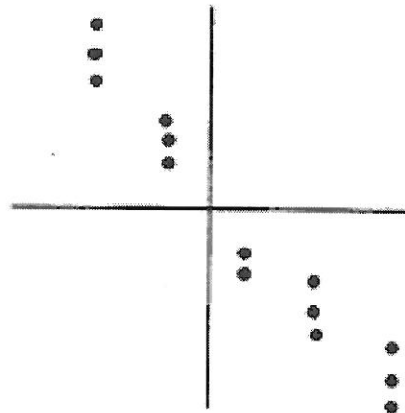
Longer Tails



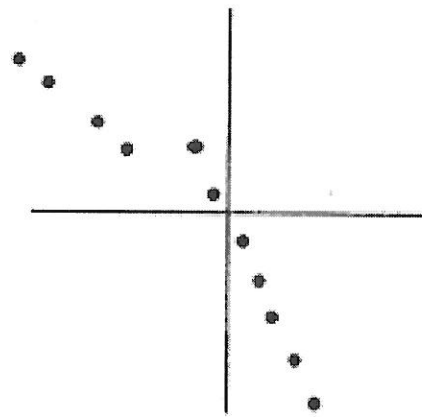
An outlier



Asymmetry

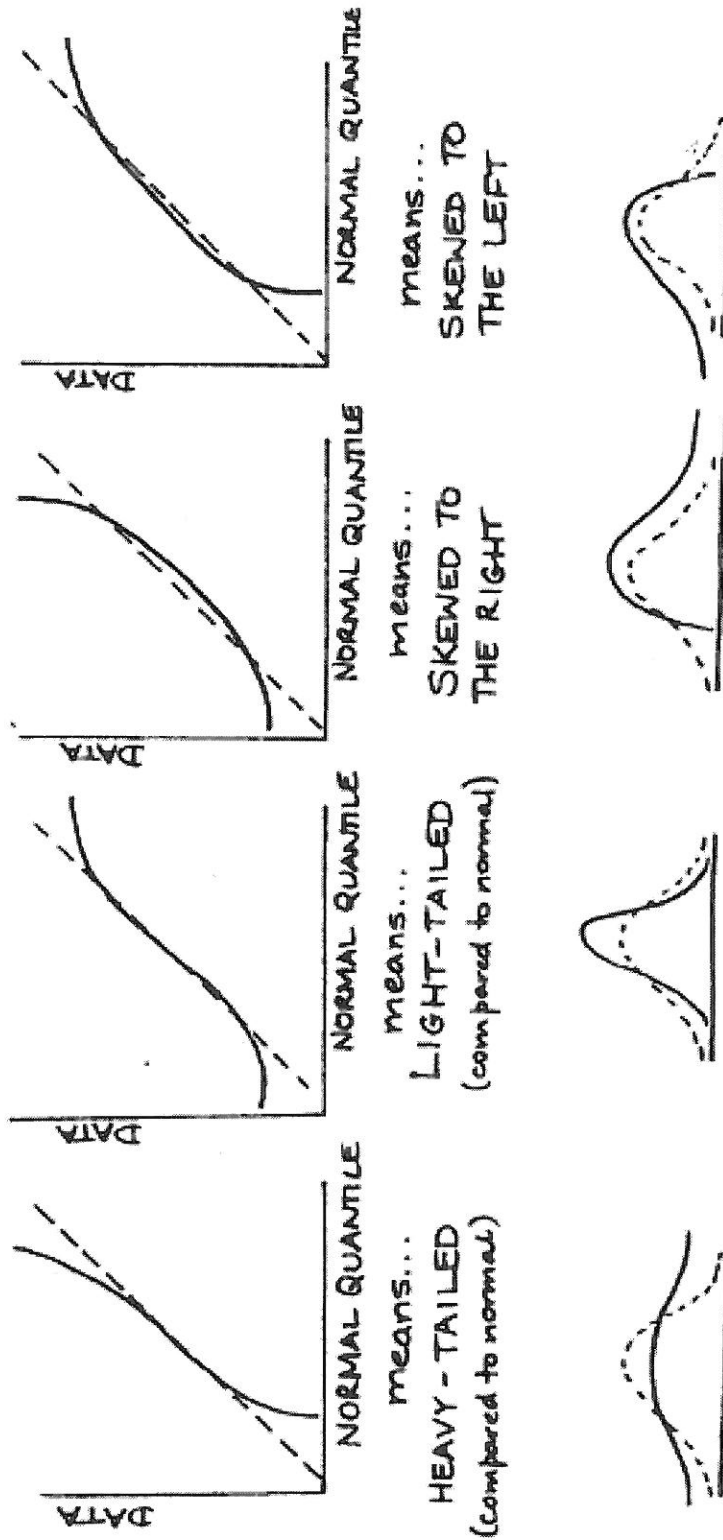


Rounding



Clustering

Features you might see in a normal quantile plot



Departures from normality