Diagnostics: residual plots

ST552 Lecture 16

Charlotte Wickham 2019-02-15

In rough order of importance:

Systematic form of the model, E(Y) = Xβ. If violated, the parameters in the model may be meaningless, estimates may be biased.

 $Y = X \beta + \varepsilon \quad \varepsilon \sim N(0, 0^{3})$ $\varepsilon : i d N(0, 0^{2})$

- Independence of errors, ε_i independent of ε_j for all i and j_i
 If violated, estimates are still unbiased, but standard errors are generally inappropriate.
 miduum miduum
- Constant variance, $Var(\epsilon_i) = \sigma^2$ for all *i*. If violated, variance in predictions may not be properly quantified. \rightarrow midum
- Normality, ε ~ N(). Can rely on CLT for large samples. If violated, prediction intervals are probably innappropriate.

- If our model is correct, $\epsilon \sim N(0, \sigma^2 I)$. But, we don't observe the errors.
- Usually, we use the residuals as our best guess for the errors, and examine them for problems with the assumptions.
- However, residuals by construction are not equal variance, or uncorrelated (you can try to standardize), but in practice the effects are small and ignored.
- We can't prove the assumptions are satisfied, but we can look for evidence of gross violations.

Graphical versus formal inferential methods

no lests

- I am a strong proponent of graphical methods over formal tests for assumption checking.
- Tests can only provide quantification of a deviation you are expecting, graphics reveal the unexpected.
- Tests tend to make you focus on statistical significance not practical significance.

For example, a large sample of data that is just a little non-Normal, will tend to give tiny p-values in a test of Normality, but for our purposes it isn't really a problem.

Residuals versus fitted values

- Residuals versus explanatories (both those included and those excluded from the model)
- Normal probability plot (Q-Q plot) of the residuals
 - Anything else you can think of that might reveal structure in the residuals. For example, if measurements are made over time or space, look for temporal or spatial patterns in the residuals.

Residuals versus fitted or explantories

An even width band vertically centered around zero as you move left to right (always put the residuals on the y-axis).





Handout (Charlotte will bring):

Part One Describe what you see in the residual plots that suggests a violation of assumptions.

Part Two The same five models are examined but in a random order, with a much smaller sample size. Can you match these diagnostics to those in Part Two??

Part Three Do you see any violations here? Tsick question Generated satisfying assumptions "Train your eyes"

Common problems and possible solutions

- Non-constant spread
 - transform response (background knowledge, trial & error, Box-Cox)
 - use more complicated models (glm, gee) gls
- Non-linearity
 - transform response
 - transform predictor
 - allow for curvature (add predictor², splines, gam)
 - use a non-linear model
- Non-normality
 - transform response
 - use more complicated models (glm) ← ST € 623
- Structure when examined against an excluded variable include it