

1. a) The equations is  $dist = \beta_0 + \beta_1 speed + \beta_2 speed^2 + \varepsilon$ .
- b) The hypotheses are:
- $H_0 : \beta_0 = 0$  versus  $H_1 : \beta_0 \neq 0$
- $H_0 : \beta_1 = 0$  versus  $H_1 : \beta_1 \neq 0$
- $H_0 : \beta_2 = 0$  versus  $H_1 : \beta_2 \neq 0$
- $H_0 : \beta_1 = \beta_2 = 0$  versus  $H_1 : \beta_1$  and  $\beta_2$  not both zero.
- c) The statement is false. The problem is that  $speed^2$  cannot be held fixed as speed varies, which is necessary for the interpretation given.
- d) The coefficient  $p$ -values say that if  $speed$  is in the model then  $speed^2$  need not be, and conversely, if  $speed^2$  is in the model then  $speed$  need not be. The overall F-statistic  $p$ -value says that at least one of  $speed$  and  $speed^2$  is required in the model. There is no conflict.
2. a) Heteroscedasticity. The plot is designed to detect a dependence of variability on mean value. This problem can often be fixed by transforming the response.
- b) Regress the response on all the except one of the independent variables and then regress that variable on the other independent variables. Plot the residuals from the first of these fits against those from the second.
- A partial residual plot would look different from the corresponding plot of residuals against independent variable if their were significant correlations amount the independent variables.
- c) A plot showing the  $i$ -th residual plotted against the  $i - 1$ -st residual is used as a check for serial correlation. If a linear trend where present in such a plot it would be an indication of serial correlation.
- d) The errors in regression are assumed to be normally distributed. This can be checked with a qq-plot, histogram or a formal hypothesis test.
3. a) A point has high leverage when the fitted value at that point is very sensitive to the observed value at that point. The hat matrix diagonal values give a direct measure of leverage.
- b) A large covariance ratio indicates that standard errors and/or correlations are changed a lot by the omission of that point. covariance ratios should be about 1.
- c) Leaving the observation out of the analysis would produce a large increase in the fitted value at that point.

4. a) The value of  $R^2$  increases whenever a variable is added to the regression equation. Using it to guide variable selection would lead to all variables being in the equation. Mallows's  $C_p$  is more useful.
- b) Stepwise regression generally leads to too many variables being in the regression equation (so that they can be pruned out by hand). Not all the variables need show up as significant.