

INFERENCE FOR MULTIPLE LINEAR REGRESSION

Terminology: Similar to terminology for simple linear regression

- $\hat{y}_i = \hat{\eta}^T \underline{u}_i$ (i^{th} fitted value or i^{th} fit)
- $\hat{e}_i = y_i - \hat{y}_i$ (i^{th} residual)
- $RSS = RSS(\hat{\eta}) = \sum (y_i - \hat{y}_i)^2 = \sum \hat{e}_i^2$ (residual sum of squares)

Results similar to those in simple linear regression:

- $\hat{\eta}_j$ is an unbiased estimator of η_j .
- $\hat{\sigma}^2 = \frac{1}{n-k} RSS$ is an unbiased estimator of σ^2 .
- $\hat{\sigma}^2$ is a multiple of a χ^2 distribution with $n-k$ degrees of freedom -- so we say $\hat{\sigma}^2$ and RSS have $df = n-k$.

Note: In simple regression, $k = 2$.

Example: Haystacks

Additional Assumptions Needed for Inference:

- (3) $Y|\underline{x}$ is normally distributed
(Recall that this will be the case if \underline{X}, Y are multivariate normal.)
- (4) The y_i 's are independent observations from the $Y|\underline{x}_i$'s.

Consequences of Assumptions (1) - (4) for Inference for Coefficients:

- $Y|\underline{x} \sim N(\underline{\eta}^T \underline{u}, \sigma^2)$
- There is a formula for $s.e.(\hat{\eta}_j)$. (We'll use software to calculate it.)
- $\frac{\hat{\eta}_j - \eta_j}{s.e.(\hat{\eta}_j)} \sim t(n-k)$ for each j .

Example: Haystacks

Inference for Means:

In simple regression, we saw

$$\text{Var}(\hat{E}(Y|x)) = \text{Var}(\hat{E}(Y|x) | x_1, \dots, x_n) = \sigma^2 \left(\frac{1}{n} + \frac{(x - \bar{x})^2}{SXX} \right).$$

So

$$\text{s.e.}(\hat{E}(Y|x)) = \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{SXX}} = \hat{\sigma} \text{ times a function of the } x_i\text{'s (but not the } y_i\text{'s)}$$

An analogous computation (best done by matrices -- see Section 7.9) in the multiple regression model gives

$$\text{Var}(\hat{E}(Y|\underline{x})) = \text{Var}(\hat{E}(Y|\underline{x}) | \underline{x}_1, \dots, \underline{x}_n) = h\sigma^2,$$

where $h = h(\underline{u})$ ($= h(\underline{x})$ by abuse of notation) is a function of $\underline{u}_1, \underline{u}_2, \dots, \underline{u}_n$, called the *leverage*. (The name will be explained later.)

In simple regression,

$$h(x) = \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{SXX}}$$

Note that $(x - \bar{x})^2$ is a measure of the distance from x to \bar{x} . Similarly, in multiple regression, $h(\underline{x})$ is a type of measure of the distance from \underline{u} to the *centroid*

$$\bar{\underline{u}} = \begin{bmatrix} 1 \\ \bar{u}_1 \\ \text{M} \\ \bar{u}_{k-1} \end{bmatrix},$$

(that is, it is a monotone function of $\sum (u_j - \bar{u}_j)^2$.) In particular:

The further \underline{u} is from $\bar{\underline{u}}$, the larger $\text{Var}(\hat{E}(Y|x))$ is, so the less precisely we can estimate $E(Y|x)$ or y . (Thus an outlier could give a large h , and hence make inference less precise.)

Example: 1 predictor

Define:

$$\text{s.e.}(\hat{E}(Y|\underline{x})) = \hat{\sigma} \sqrt{h(\underline{u})}$$

Summarizing:

- The larger the leverage, the larger s.e. ($\hat{E}(Y|x)$) is, so the less precisely we can estimate $E(Y|x)$.
- The leverage depends just on the \underline{x}_i 's, not on the y_i 's.

Similarly to simple regression:

$$\frac{\hat{E}(Y | \underline{x}) - E(Y | \underline{x})}{\text{s.e.}(\hat{E}(Y | \underline{x}))} \sim t(n-k).$$

Thus we can do hypothesis tests and find confidence intervals for the conditional mean response $E(Y|\underline{x})$

Prediction: Results are similar to simple regression:

- Prediction error = $Y|\underline{x} - \hat{E}(Y|\underline{x})$
- $\text{Var}(Y|\underline{x} - \hat{E}(Y|\underline{x})) = \sigma^2(1 + h(\underline{x})) = \sigma^2 + \text{Var}(E(Y|\underline{x}))$
- Define s.e. ($Y_{\text{pred}}|\underline{x}$) = $\hat{\sigma}\sqrt{1+h}$
- $\frac{Y|\underline{x} - \hat{E}(Y|\underline{x})}{\text{se}(y_{\text{pred}}|\underline{x})} \sim t(n-k)$, so we can form prediction intervals.

Example: Haystacks