Multiple Linear Regression using Matrices - II I

- Geometrically, the vector **y** is a point in the *n* dimensional sample space.
- The vector β is a point in the *p* dimensional *parameter space.*, where p = k + 1.
- For each β, Xβ is a point in the sample space (i.e. the n×p matrix X is a mapping from the parameter space to the sample space).
- The values Xβ form a p dimensional linear subspace or plane within the sample space. It is also assumed that X is of rank p.
- This subspace is called the *expectation surface*, because $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$.

Multiple Linear Regression using Matrices - II II

• The sum of squares

$$S(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$$

is the squared length of the vector $\mathbf{y} - \mathbf{X}\boldsymbol{\beta}$, or the squared distance between the observed \mathbf{y} and the point $\mathbf{X}\boldsymbol{\beta}$.

- The method of least squares finds the value of β , called $\hat{\beta}$, for which the point $\mathbf{X}\hat{\beta}$ on the expectation surface is closest to the point \mathbf{y} .
- From geometric considerations, Xβ̂ must be the perpendicular projection of y on the expectation surface.

Hat matrix

- The 'hat' matrix, H = X(X^TX)⁻¹X^T, projects y onto the expectation surface at a point ŷ closest to y.
- So $\hat{\mathbf{y}} = H\mathbf{y} = \mathbf{X}[(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}] = \mathbf{X}\hat{\boldsymbol{\beta}}$, where $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$ is the least squares estimate of $\boldsymbol{\beta}$.
- Note that *HX* = *X*. This means, for example, that any vector *v* in the column space of *X*, *Hv* = *v*. As an example, *H*1 = 1.
- Note that **H** is symmetric, and that $H^2 = H$.

We assume that:

- $oldsymbol{y} = oldsymbol{X}oldsymbol{eta} + oldsymbol{\epsilon}$, where
 - ${\, \bullet \, }$ elements of ϵ are normally distributed with
 - **E**(*e*) = 0
 - $Cov(\epsilon) = \sigma^2 I$

equivalently

- elements of **y** are normally distributed with
- $E(y) = X\beta$ and
- Cov(y) = $\sigma^2 I$

Sampling distribution of $\hat{oldsymbol{eta}}$

$$\hat{\boldsymbol{eta}} = (\boldsymbol{X}^{\mathsf{T}} \boldsymbol{X})^{-1} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{y}$$

• elements of $\hat{\beta}$ are normally distributed, because linear combinations of normal random variables are normally distributed, with

•
$$E(\hat{\beta}) = (X^T X)^{-1} X^T E(y) = (X^T X)^{-1} X^T X \beta = \beta$$
 and

$$Cov(\hat{\boldsymbol{\beta}}) = Cov((\boldsymbol{X}^{T}\boldsymbol{X})^{-1}\boldsymbol{X}^{T}\boldsymbol{y})$$
$$= (\boldsymbol{X}^{T}\boldsymbol{X})^{-1}\boldsymbol{X}^{T}Cov(\boldsymbol{y})(\boldsymbol{X}^{T}\boldsymbol{X})^{-1}\boldsymbol{X}^{T})^{T}$$
$$= \boldsymbol{X}^{T}\boldsymbol{X})^{-1}\boldsymbol{X}^{T}\sigma^{2}\boldsymbol{I}\boldsymbol{X}(\boldsymbol{X}^{T}\boldsymbol{X})^{-1}$$
$$= \sigma^{2}(\boldsymbol{X}^{T}\boldsymbol{X})^{-1}$$

◆ロ ▶ ▲母 ▶ ▲目 ▶ ▲日 ▶ ▲日 ● ◆ ○ ○

Sampling distribution of \hat{y}

٠

$$\hat{y} = Hy$$

• elements of \hat{y} are normally distributed, because linear combinations of normal random variables are normally distributed, with

•
$$E(\hat{y}) = HE(y) = HX\beta = X(X^TX)^{-1}X^TX\beta = X\beta$$

$$Cov(\hat{\mathbf{y}}) = Cov(H\mathbf{y}) = HCov(\mathbf{y})H^{T} = H\sigma^{2}IH^{T}$$
$$= \sigma^{2}HH^{T} = \sigma^{2}H = \sigma^{2}X(X^{T}X)^{-1}X^{T}$$

Sampling distribution of residuals

$$\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{H})\mathbf{y}$$

 elements of *r* are normally distributed, because linear combinations of normal random variables are normally distributed, with

$$\boldsymbol{E}(\boldsymbol{r}) = (\boldsymbol{I} - \boldsymbol{H})\boldsymbol{E}(\boldsymbol{y}\boldsymbol{p}) = (\boldsymbol{I} - \boldsymbol{H})\boldsymbol{X}\boldsymbol{\beta} = \boldsymbol{0}$$

٩

۲

$$Cov(\mathbf{r}) = Cov((\mathbf{I} - \mathbf{H})\mathbf{y}) = (\mathbf{I} - \mathbf{H})\sigma^2 \mathbf{I}(\mathbf{I} - \mathbf{H})$$
$$= \sigma^2 (\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H})^T = \sigma^2 (\mathbf{I} - \mathbf{H})$$

• The variance of the i'th residual is $\sigma^2(1 - h_{ii})$ where h_{ii} is the i'th diagonal element of **H**. This means that the residuals will typically have different variances.

Covariance of residuals and fitted values

$$Cov(\mathbf{r}, \hat{\mathbf{y}}) = Cov((\mathbf{I} - \mathbf{H})\mathbf{y}, \mathbf{H}\mathbf{y})$$
$$= \sigma^{2}(\mathbf{I} - \mathbf{H})\mathbf{H} = \mathbf{0}$$

(Why?)

۵

- which means that residuals and fitted values are uncorrelated
- which means that residuals and fitted values are independent (as uncorrelated normal random variables are independent).

Estimation of σ^2

- The residuals are contained in the vector $\boldsymbol{e} = \boldsymbol{y} \boldsymbol{\hat{y}}$.
- The residual, or error, sum of squares is $SSE = \sum e_i^2 = e^T e$
- With k predictor variables, the error mean square is

$$MSE = \frac{SSE}{n - (k + 1)} = \frac{SSE}{n - p}$$

• $\hat{\sigma}^2 = MSE$ is an unbiased estimate of σ^2 because

$$\boldsymbol{E}(SSE) = \boldsymbol{E}(\boldsymbol{e}^{T}\boldsymbol{e}) = \boldsymbol{0}^{T}\boldsymbol{0} + \sigma^{2}tr(\boldsymbol{I}(\boldsymbol{I} - \boldsymbol{H})) = \sigma^{2}(tr\boldsymbol{I} - tr(\boldsymbol{H}))$$
$$= \sigma^{2}(n - p)$$

(Why?)

Confidence intervals I

- We will see that SSE is independent of $\hat{\beta}$, and
- where C_{jj} is the j + 1'st diagonal entry of $(\boldsymbol{X}^T \boldsymbol{X})^{-1}$, the standard error of $\hat{\beta}_j$ is $\hat{\sigma} \sqrt{C_{jj}}$, and
- $t = \frac{\hat{\beta}_j \beta_j}{\hat{\sigma}\sqrt{C_{jj}}}$ has a t distribution with n (k + 1) = n p degrees of freedom.
- It follows that a $100(1 \alpha)$ % confidence interval for β_j is given by

$$\hat{eta}_j \pm t_{lpha/2, n-p} \hat{\sigma} \sqrt{C_{jj}}$$

• where \mathbf{x}_0 is a vector of covariate values, a $100(1 - \alpha)\%$ confidence interval for the mean of y when $\mathbf{x} = \mathbf{x}_0$, $\mu_{y|\mathbf{x}_0} = E(y|\mathbf{x}_0)$, is given by

$$\hat{y}_0 \pm t_{\alpha/2,n-p} \hat{\sigma} \sqrt{\boldsymbol{x}_0^T (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{x}_0}$$

Prediction intervals (not responsible)

- In future, you plan to measure the value of y when the predictor variables equal x₀.
- A point estimate of the future y is given by $\hat{y}_0 = \mathbf{x}'_0 \hat{\beta}$.
- The goal is to find an interval in which the future y will lie, with probability 1α .
- A 100(1 α)% prediction interval for a future value of y when x = x₀ is given by

$$\hat{y}_0 \pm t_{\alpha/2,n-p} \hat{\sigma} \sqrt{1 + \boldsymbol{x}_0^T (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{x}_0}$$

Simultaneous confidence intervals for β_i 's

- When making more than one confidence interval, it is important to ensure that the joint coverage of all of the intervals constructed is $\geq 100(1 \alpha)$.
- If *m* intervals are constructed, with each interval having coverage probability $100(1 \alpha/m)$, it follows from Bonferroni's inequality that the probability that the *m* intervals will simulataneously cover their associated parameters is at least 1α .
- For example, the k + 1 intervals

$$\hat{\beta}_j \pm t_{lpha/(2(k+1)),n-1-k}\sqrt{MSE}\sqrt{C_{j,j}}, j=0,1,2,\ldots,k$$

have joint coverage of at least $100(1 - \alpha)$. In general, if constructing *m* simultaneous confidence intervals, just replace α by α/m , and use the usual procedure for the CI.

Bonferroni's inequality (Not responsible for this material)

Where E₁ and E₂ are events, and E^c is the complement of E, because the probability of the union is less than or equal to the sum of the probabilities, we know that

$$P(E_1^c \bigcup E_2^c) \le P(E_1^c) + P(E_2^c)$$

• then use the fact that $P(E) = 1 - P(E^c)$, similarly $P(E_1^c) = 1 - P(E_1)$, and let $E = E_1^c \bigcup E_2^c$, from which $E^c = E_1 \bigcap E^2$. Putting all of this together it follows that

$$P(E_1 \bigcap E_2) = 1 - P(E_1^c \bigcup E_2^c) \ge 1 - (P(E_1^c) + P(E_2^c))$$

• Letting $P(E_1) = 1 - \alpha/2$ and $P(E_2) = 1 - \alpha/2$, it follows that

$$P(E_1 \bigcap E_2) \ge 1 - (\alpha/2 + \alpha/2) = 1 - \alpha$$

Application to simultaneous inference - take this as a given.

- Let E_1 be the event that β_j is contained in the interval with endpoints $\hat{\beta}_j \pm t_{(\alpha/2)/2,n-p}\sqrt{MSE}\sqrt{C_{j+1,j+1}}$. The probability of this is $1 \alpha/2$.
- Let E_2 be the event that β_m is contained in the interval with endpoints $\hat{\beta}_m \pm t_{(\alpha/2)/2,n-p}\sqrt{MSE}\sqrt{C_{m+1,m+1}}$. The probability of this is $1 \alpha/2$.
- It follows from Bonferroni's inequality that the probability that EACH of β_j and β_m is simultaneously contained in their associated intervals is at least 1α .
- We say that the two intervals form a joint $100(1 \alpha)\%$ confidence region for (β_j, β_m) .
- What works for two the the β_j's, works for an arbitrary number, hence the form for the k + 1 simultanous intervals given above, where α is replaced by α/(k + 1).

A joint 100(1 - α)% confidence region for the vector β is given by those points β for which

$$(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \leq p \, \mathsf{MSE} \, F_{\alpha, p, n-p}$$

• This is a p dimensional ellipse centred at $\hat{\beta}$.