

## Multiple regression

Situation: more than one independent variable; want to predict  $Y$  from  $x_1, x_2, \dots, x_k$ .

- ex:**
- IRS predicts the amount of money to be recovered in an audit using (among other variables) amt. of deduction for charitable gifts, amt. of real estate losses, etc.
  - House appraiser predicts sale price of a house based on sq. ft., # bedrooms, ave. sale price in neighborhood, etc.

Idea: mean of  $Y$  is taken to be a linear function of the predictor variables:

$$E(Y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \cdots + \beta_k x_k$$

With just two predictor variables (not functionally dependent), this equation is a **plane**.

**Model:**

$$Y \sim \text{normal}(\beta_0 + \beta_1 x_1 + \beta_2 x_2 \cdots + \beta_k x_k, \sigma^2)$$

**Different types of predictor variables:**

- ordinary quantitative variables
- indicator variables (AOV is a regression!)

3 treatments; means  $\mu_1, \mu_2, \mu_3$

$$x_1 = \begin{cases} 1 & \text{if observation is from trt 1} \\ 0 & \text{otherwise} \end{cases}$$

$$x_2 = \begin{cases} 1 & \text{if observation is from trt 2} \\ 0 & \text{otherwise} \end{cases}$$

$$E(Y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$\mu_1 = \beta_0 + \beta_1$$

$$\mu_2 = \beta_0 + \beta_2$$

$$\mu_3 = \beta_0$$

- nonlinear terms, e.g.

$$E(Y) = \beta_0 + \beta_1 x + \beta_2 x^2$$

(Note: this is still a **linear statistical model** because the parameters appear linearly)

- interactions, e.g.

$$E(Y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$$

**Estimates** of unknown parameters are conveniently represented with **matrix notation**

(your task: learn to multiply matrices; learn what an inverse matrix is)

Data:  $(y_1, x_{11}, x_{12}, \dots, x_{1k})$   
 $(y_2, x_{21}, x_{22}, \dots, x_{2k})$   
 $\vdots$   
 $(y_n, x_{n1}, x_{n2}, \dots, x_{nk})$

Matrices:

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix}_{(n \times (k+1))}$$

**(design matrix)**

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{(n \times 1)} \quad \hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{bmatrix}_{((k+1) \times 1)} \quad \text{ML, LS}$$

Matrix representations:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} \quad \text{predicted values}$$

$$\mathbf{e} = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \quad \text{residuals}$$

$$\text{SS}(\text{residuals}) = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})' (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

sum of squared errors (minimized at  $\hat{\boldsymbol{\beta}}$ )

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} \\ 1 & x_{21} & x_{22} \\ 1 & x_{31} & x_{32} \\ 1 & x_{41} & x_{42} \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

$$\mathbf{X}\boldsymbol{\beta} = \begin{bmatrix} \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} \\ \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} \\ \beta_0 + \beta_1 x_{31} + \beta_2 x_{32} \\ \beta_0 + \beta_1 x_{41} + \beta_2 x_{42} \end{bmatrix}$$

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ x_{11} & x_{21} & x_{31} & x_{41} \\ x_{12} & x_{22} & x_{32} & x_{42} \end{bmatrix} \begin{bmatrix} 1 & x_{11} & x_{12} \\ 1 & x_{21} & x_{22} \\ 1 & x_{31} & x_{32} \\ 1 & x_{41} & x_{42} \end{bmatrix}$$

$$= \begin{bmatrix} n & \sum x_{i1} & \sum x_{i2} \\ \sum x_{i1} & \sum x_{i1}^2 & \sum x_{i1} x_{i2} \\ \sum x_{i2} & \sum x_{i1} x_{i2} & \sum x_{i2}^2 \end{bmatrix}$$

$$(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{X}) = \mathbf{I} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$\mathbf{X}'\mathbf{X}$  matrix of sums of squares & crossproducts  
of the predictors (symmetric,  $(k + 1) \times (k + 1)$  )

$$(\mathbf{X}'\mathbf{X})\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y} \quad \text{“normal equations”}$$

(minimizing the sum of squared errors results in  
a system of  $k + 1$  linear equations in  $k + 1$  unknowns)

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

ML & LS estimates of the parameters  $\beta_0, \beta_1, \dots, \beta_k$

$$\hat{\sigma}^2 = \frac{1}{(n - k - 1)} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})' (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

unbiased (adjusted ML) estimate of  $\sigma^2$

Note: different ways of writing the sum of squared  
residuals result from matrix-algebraic rearrangements

$$\begin{aligned} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})' (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) &= \mathbf{y}'\mathbf{y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{y} \\ &= \mathbf{y}'[\mathbf{I} - \mathbf{H}]\mathbf{y} \end{aligned}$$

where  $\mathbf{I}$  = identity matrix,  $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  “hat matrix”

## Inferences for multiple regression

### The “AOV table”

source of variation	SS	df	MS
regression	SS(regression)	$k$	$\frac{SS(\text{regression})}{k}$
error	SS(residual)	$n - k - 1$	$\frac{SS(\text{residual})}{n - k - 1}$
total	SS(total)	$n - 1$	

$$SS(\text{regression}) = \left( \mathbf{X}\hat{\boldsymbol{\beta}} - \frac{1}{n}\mathbf{J}\mathbf{y} \right)' \left( \mathbf{X}\hat{\boldsymbol{\beta}} - \frac{1}{n}\mathbf{J}\mathbf{y} \right)$$

$$SS(\text{residual}) = \left( \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \right)' \left( \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \right)$$

$$SS(\text{total}) = \left( \mathbf{y} - \frac{1}{n}\mathbf{J}\mathbf{y} \right)' \left( \mathbf{y} - \frac{1}{n}\mathbf{J}\mathbf{y} \right)$$

where  $\mathbf{J}$  is an  $n \times n$  matrix of 1's

## Test of the model vs the mean

$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$   
(just the overall mean is used for predicting  $Y$ )

$H_a: \neq$   
(all predictor variables are included in the model)

Test statistic (from LR):  $f = \frac{SS(\text{regression})/k}{SS(\text{residual})/(n-k-1)}$

Rejection region: reject  $H_0$  if  $f \geq f_\alpha$  where the F distribution has  $k$  and  $n - k - 1$  df.

## Coefficient of determination

$$r^2 = \frac{SS(\text{regression})}{SS(\text{total})} = 1 - \frac{SS(\text{residual})}{SS(\text{total})}$$

(proportional reduction in the prediction error attained by using the multiple regression model instead of the overall mean)