### Simple Linear Regression

2016年10月3日 0:01

#### Assumption

- 1. Uncorrelated/independent
- 2. Common variance  $\sigma^2$ ,  $Var(Y_i) = \sigma^2$
- 3. Linear,  $\mathbb{E}(Y_i|x_i) = \beta_0 + \beta_1 x_i$

## s. e. $(\widehat{\beta_1}) = \sqrt{\operatorname{Var}(\widehat{\beta_1})} = \sqrt{\frac{\widehat{\sigma}^2}{S_{XX}}}$ s. e. $(\widehat{\beta_0}) = \sqrt{\operatorname{Var}(\widehat{\beta_0})} = \sqrt{\widehat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{XX}}\right)}$

#### **Least Squares Estimation**

Least squares estimation
$$Q = \sum_{i=1}^{n} \{y_i - E(Y_i | x_i)\}^2 = \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial Q}{\partial \beta_0} = 0 \implies \widehat{\beta_0} = \overline{y} - \widehat{\beta_1} \overline{x}$$

$$\frac{\partial Q}{\partial \beta_1} = 0 \implies \widehat{\beta_1} = \frac{S_{XY}}{S_{XX}}$$

$$S_{XX} = \sum_{i=1}^{n} (x_i - \overline{x})^2$$

$$S_{XY} = \sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y}) = \sum_{i=1}^{n} (x_i - \overline{x})\overline{y}$$

#### Properties

 $\widehat{\beta_0}$  and  $\widehat{\beta_1}$  are unbiased and consistent estimators of  $\beta_0$  and  $\beta_1$ , respectively

estimato.

$$\mathbb{E}(\overline{Y_1}) = \beta_0 + \beta_1 \overline{x}$$

$$\mathbb{E}(\widehat{\beta_1}) = \beta_1$$

$$\mathbb{E}(\widehat{\beta_0}) = \beta_0 \stackrel{\checkmark}{=}$$

$$\mathbb{E}(\widehat{\beta_0}) = \sigma^2$$

$$\operatorname{Var}(\widehat{\beta_1}) = \frac{\sigma^2}{S_{XX}}$$

$$\operatorname{Var}(\widehat{\beta_0}) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{XX}}\right)$$

$$\operatorname{var}(\beta_0) = \sigma^2 \left( \frac{1}{n} + \frac{1}{S_{XX}} \right)$$

$$\operatorname{cov}(\widehat{\beta_0}, \widehat{\beta_1}) = \operatorname{cov}\left( \sum_i \left( \frac{1}{n} - \frac{\bar{x}(x_i - \bar{x})}{S_{XX}} \right) Y_i, \sum_i \frac{x_i - \bar{x}}{S_{XX}} Y_i \right)$$

$$= \sum_i \left( \frac{1}{n} - \frac{\bar{x}(x_i - \bar{x})}{S_{XX}} \right) \left( \frac{x_i - \bar{x}}{S_{XX}} \right) \operatorname{Var}(Y_i) = -\frac{\bar{x}}{S_{XX}} \sigma^2$$

$$\operatorname{cov}(\overline{Y}, \widehat{\beta_1}) = \operatorname{cov}\left(\sum \frac{1}{n} Y_i, \sum \frac{x_i - \overline{x}}{S_{XX}} Y_i\right) = 0$$

 $\widehat{eta_0}$  and  $\widehat{eta_1}$  are expressed as a function of the random variables  $Y_i$  instead of the obseved data  $y_i$ 

#### **Alternative Formulation**

 $S_{YY} = \sum_{i} (y_i - \bar{y})^2$ 

$$\mathbb{E}(Y_i|x_i) = \gamma + \beta_1(x_i - \bar{x})$$

So that  $\gamma$  is the expected response at  $x=\bar{x}$ , rather than at x=0. Least Squares Estimation,  $\widehat{\gamma}=\bar{y}$ ,  $\widehat{\beta_1}=\frac{S_{XY}}{S_{YY}}$ 

$$Var(\hat{\gamma}) = \frac{\sigma^2}{n}, cov(\hat{\gamma}, \widehat{\beta_1}) = 0$$

$$\operatorname{Var}(\widehat{\beta_0} + \widehat{\beta_1} x_*) = \operatorname{Var}(\widehat{\gamma} + \widehat{\beta_1} (x_* - \bar{x})) = \sigma^2 \left( \frac{1}{n} + \frac{(x_* - \bar{x})^2}{S_{XX}} \right)$$

#### Inference: Confidence Interval (Expected Response)

$$\mathbb{E}(Y|x_*) = \widehat{\beta_0} + \widehat{\beta_1}x_*$$

$$\widehat{\mathbb{E}}(Y|x_*) = \widehat{\beta_0} + \widehat{\beta_1}x_* = \overline{Y} + \widehat{\beta_1}(x_* - \overline{x})$$

Using ( ) or properties of variance,

$$\operatorname{Var}\left(\widehat{\mathbb{E}}(Y|X_*)\right) = \sigma^2 \left(\frac{1}{n} + \frac{(x_* - \bar{x})^2}{S_{XX}}\right)$$

#### $T \sim t_{n-1}$

Inference: Prediction Interval (Future Response)

$$Y_* = \beta_0 + \beta_1 x_* + \varepsilon$$
, where  $\varepsilon \sim N(0, \sigma^2)$ 

$$Var(Y_*) = \sigma^2 \left( 1 + \frac{1}{n} + \frac{(x_* - \bar{x})^2}{S_{XX}} \right)$$

$$T \sim t_{n-2}$$

#### Inference (When Normal Distributed)

# $\frac{\widehat{\beta_1} - \beta_1}{\sqrt{\widehat{\sigma}^2/S_{XX}}} \sim t_{n-2}$

#### **Residual and Regression Sums of Squares**

$$e_i = y_i - \widehat{y}_i = (y_i - \overline{y}) - \widehat{\beta}_1(x_i - \overline{x})$$

RSS = 
$$\sum e_i^2 = S_{YY} - \frac{S_{XY}^2}{S_{XX}} = S_{YY} - \text{regressionSS} = \text{totalSS} - \text{regressionSS}$$

$$\mathbb{E}(RSS) = (n-2)\sigma^2$$

$$\hat{\sigma}^2 = s^2 = \frac{1}{n-2} RSS = \frac{RSS}{DF} = residual Mean Square$$

$$Q = \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_i)^2 = S_{XX} \left( \beta_1 - \frac{S_{XY}}{S_{XX}} \right)^2 + n(\beta_0 - \bar{y} - \beta_1 \bar{x})^2 + S_{YY} - \frac{S_{XY}^2}{S_{XX}}$$

Coefficient of determination (Multiple R-Squared)

$$R^{2} = \frac{\text{regressionSS}}{\text{totalSS}} = \frac{S_{XY}^{2}/S_{XX}}{S_{YY}} = \frac{S_{XY}^{2}}{S_{XX}S_{YY}}$$

#### ΔΝΟVΔ

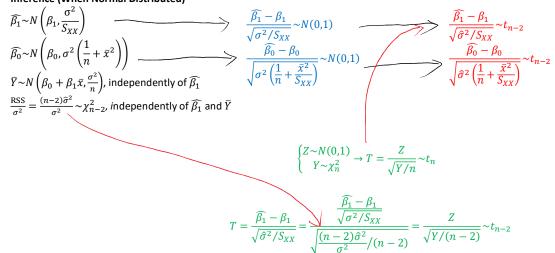
ANOVA					
SOURCE	DF	SS	MS	F	Р
Regression	1	$\frac{S_{XY}^2}{S_{XX}}$	$\frac{S_{XY}^2}{S_{XX}}$	$\frac{S_{XY}^2}{S_{XX}\hat{\sigma}^2}$	
Residual	n – 2	$S_{YY} - \frac{S_{XY}^2}{S_{XX}}$	$\hat{\sigma}^2$		
Total	n – 1	S <sub>YY</sub>			

$$MS = \frac{SS}{DF}$$

$$F = \frac{\text{regressionMS}}{\text{residualMS}} = \frac{\text{regressionSS}}{1} / \frac{\text{residualSS}}{n-2} \sim F_{1,n-2}$$

F can be used to test whether the expected response depends on the explanatory variable.





$$\widehat{y_i} = \widehat{\beta_0} + \widehat{\beta_1} x_i = \overline{y} + \widehat{\beta_1} (x_i - \overline{x}) = \overline{y} + (x_i - \overline{x}) \sum_{j=1}^n \frac{(x_j - \overline{x}) y_j}{S_{XX}}$$

$$= \sum_{j=1}^n \left( \frac{1}{n} + \frac{(x_i - \overline{x}) (x_j - \overline{x})}{S_{XX}} \right) y_j = \sum_{j=1}^n a_{ij} y_j$$

Influence Matrix:  $A = \{a_{ij}\}$ 

Leverage is  $A_{ii}$ , a measure of how much that particular data point influences its own fitted value. (potential)

Large Residual + High Leverage = Strong Influence

Small Residual + High Leverage = not too great influence (while inconsistent)

Large Residual + Low Leverage = not much effect on the fit

Cook's distance

$$d_i = \frac{1}{3\widehat{\sigma}^2} \sum_{i=1}^n (\widehat{y_{j(i)}} - \widehat{y_j})^2$$

 $\widehat{y_{i(i)}}$  is the i'th fitted value obtained when the j'th point is excluded when computing fthe regression model.

3 = p + 1, where p is the number of unknown parameters.

Cook's distance measures the influence that each data point is in fact exerting on the fit. Large Cook's Distance = potential outliers

LSM Page 2

variable

#### **Analysis of Residuals**

$$\varepsilon_i = Y_i - \mathbb{E}(Y_i | x_i) = Y_i - \beta_0 - \beta_1 x_i$$

Assumptions (derive from original assumptions):

- 1. Uncorrelated/independent,  $cov(\varepsilon_i, \varepsilon_i | x) = 0$
- 2. Common variance  $\sigma^2$ ,  $Var(\varepsilon_i|x) = \sigma^2$
- 3.  $\mathbb{E}(\varepsilon_i|x)=0$
- 4. If  $Y_i$  is Normal,  $\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$

Estimate  $\varepsilon_i$ ,

$$e_i = y_i - \widehat{\beta_0} - \widehat{\beta_1} x_i = y_i - \overline{y} - \widehat{\beta_1} (x_i - \overline{x})$$

$$p_{ij} = \frac{1}{n} + \frac{(x_i - \overline{x})(x_j - \overline{x})}{S_{XX}}$$

Then, (corresponding random variable)

$$cov(E_j, E_j|x) = -pij\sigma^2$$

$$Var(E_i|x) = (1 - p_{ij})\sigma^2 = \left(1 - \frac{1}{n} - \frac{(x_i - \bar{x})^2}{S_{XX}}\right)\sigma^2$$

$$\mathbb{E}(E_i|x)=0$$

If  $Y_i$  are Normal, so are the  $E_i$ .

If n is large and none of the quantities  $|x_i - \bar{x}|$  is large then the  $p_{ij}$  are small.

#### **Standardized Residuals**

$$r_i = \frac{e_i}{\text{estimated standard error of } E_i} = \frac{e_i}{\sqrt{\hat{\sigma}^2 \left(1 - \frac{1}{n} - \frac{(x_i - \bar{x})^2}{S_{XX}}\right)}}$$

- i. Plot  $r_i$  against  $x_i$ , look for sign that  $\mathbb{E}(R_i|x_i)$ ,  $\text{Var}(R_i|x_i)$  depends on  $x_i$
- ii. See if any unusually large values,  $|r_i| > 2.5 > 3$
- iii. Plot Normal probability of  $r_i$

LSM Revision Part 2 Author: s1680642

### Multiple Regression

1.  $\mathbb{E}(\mathbf{Y}|\mathbf{X}) = \mathbf{X}\beta$ ,  $Var(\mathbf{Y}|\mathbf{X}) = \sigma^2 \mathbf{I}_n$ Specially, in simple linear regression,  $y_i = \beta_0 + \beta_1 x_i$ ,

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \ \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}, \ \mathbb{E}(\mathbf{Y}|\mathbf{X}) = \mathbf{1}_n \beta_0 + \mathbf{x}\beta_1, \ \operatorname{Var}(\mathbf{Y}|\mathbf{X}) = \sigma^2 \mathbf{I}_n$$

2. Least squares estimation:  $Q = \sum_{i=1}^{n} \{y_i - \mathbb{E}(Y_i|\mathbf{X})\}^2 = \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\beta + \beta^T \mathbf{X}^T \mathbf{X}\beta$ 

Least squares unbiased estimator:  $\hat{\boldsymbol{\beta}} = (\mathbf{X}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$  $\mathbb{E}(\hat{\boldsymbol{\beta}}|\mathbf{X}) = \boldsymbol{\beta}, \operatorname{Var}(\hat{\boldsymbol{\beta}}|\mathbf{X}) = \sigma^2(\mathbf{X}^T\mathbf{X})^{-1}, \operatorname{Var}(\mathbf{c}^T\hat{\boldsymbol{\beta}}|\mathbf{X}) = \sigma^2\mathbf{c}^T(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{c}$ 

- 3. Vector of residuals:  $\mathbf{e} = \mathbf{y} \mathbf{X}\hat{\boldsymbol{\beta}} = (\mathbf{I}_n \mathbf{P}_{\mathbf{X}})\mathbf{y}$ , where  $\mathbf{P}_{\mathbf{X}} = \mathbf{X}(\mathbf{X}\mathbf{X})^{-1}\mathbf{X}^T$  is  $n \times n$ , symmetric, idempotent, rank p,  $(\mathbf{I}_n \mathbf{P}_{\mathbf{X}})\mathbf{X} = \mathbf{0}$ ,  $(\mathbf{I}_n \mathbf{P}_{\mathbf{X}})\mathbf{P}_{\mathbf{X}} = \mathbf{0}$   $\mathbb{E}(\mathbf{E}|\mathbf{X}) = (\mathbf{I}_n \mathbf{P}_{\mathbf{X}})\mathbb{E}(\mathbf{Y}|\mathbf{X}) = (\mathbf{I}_n \mathbf{P}_{\mathbf{X}})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$   $\operatorname{Var}(\mathbf{E}|\mathbf{X}) = (\mathbf{I}_n \mathbf{P}_{\mathbf{X}})\sigma^2\mathbf{I}_n(\mathbf{I}_n \mathbf{P}_{\mathbf{X}}) = \sigma^2(\mathbf{I}_n \mathbf{P}_{\mathbf{X}})$   $\operatorname{RSS} = \mathbf{e}^T\mathbf{e} = \mathbf{y}^T\mathbf{y} \hat{\boldsymbol{\beta}}^T\mathbf{X}^T\mathbf{y}, \ \hat{\sigma}^2 = \frac{\operatorname{RSS}}{n-p} = \frac{\mathbf{y}^T\mathbf{y} \hat{\boldsymbol{\beta}}^T\mathbf{X}^T\mathbf{y}}{n-p}$
- 4. Alternative formulation (for models with an intercept)  $\mathbb{E}(Y_{i}|\mathbf{X}) = \gamma + \beta_{1}(x_{i1} \bar{x}_{1}) + \beta_{2}(x_{i2} \bar{x}_{2}) + \dots + \beta_{q}(x_{iq} \bar{x}_{q})$   $\mathbb{E}(\mathbf{Y}|\mathbf{X}) = \gamma \mathbf{1}_{n} + \dot{\mathbf{X}}\dot{\boldsymbol{\beta}}, \text{ where } \dot{\mathbf{X}}_{ij} = x_{ij} \bar{x}_{j}, \, \dot{\boldsymbol{\beta}} = (\beta_{1} \cdots \beta_{q})^{T}, \, \gamma = \beta_{0} + \beta_{1}\bar{x}_{1} + \dots + \beta_{q}\bar{x}_{q}$ Least squares unbiased estimators:  $\hat{\gamma} = \bar{y}, \, \dot{\boldsymbol{\beta}} = (\dot{\mathbf{X}}_{T}\dot{\mathbf{X}})^{-1} \, \dot{\mathbf{X}}^{T}\mathbf{y}$   $\operatorname{Var}(\hat{\boldsymbol{\beta}}|\mathbf{X}) = \sigma^{2}(\dot{\mathbf{X}}^{T}\dot{\mathbf{X}})^{-1}, \, \operatorname{Var}(\hat{\gamma}|\mathbf{X}) = n^{-1}\sigma^{2}, \, \operatorname{cov}(\hat{\boldsymbol{\beta}}, \hat{\gamma}|\mathbf{X}) = \mathbf{0}$
- 5. Distributional results:
  - $\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^T\mathbf{X})^{-1})$
  - regression (model) SS  $\mathbf{Y}^T \mathbf{P}_{\mathbf{X}} \mathbf{Y} \sim \sigma^2 \chi^2(q, \sigma^{-2} \dot{\boldsymbol{\beta}}^T \dot{\mathbf{X}}^T \dot{\mathbf{X}} \dot{\boldsymbol{\beta}})$
  - RSS  $\mathbf{Y}^T(\mathbf{I}_n \mathbf{P}_{\mathbf{X}})\mathbf{Y} \sim \sigma^2 \chi^2(n-q-1,0)$
  - $\bullet$  RSS and regression SS are independent
  - $\frac{\mathbf{c}^T \hat{\boldsymbol{\beta}} \mathbf{c}^T \boldsymbol{\beta}}{\sigma \sqrt{\mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c}}} \sim N(0, 1) \text{ and }$  $\frac{\mathbf{c}^T \hat{\boldsymbol{\beta}} \mathbf{c}^T \boldsymbol{\beta}}{\hat{\sigma} \sqrt{\mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c}}} \sim t(n p) \text{ (test hypotheses about linear funcs of the parameters)}$
- 6. 95% Confidence interval:  $\mathbf{c}^T \hat{\boldsymbol{\beta}} \pm t_{0.025} \hat{\sigma} \sqrt{\mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c}}$ CI for future response:  $\mathbf{x}_*^T \hat{\boldsymbol{\beta}} \pm t_{0.025} \hat{\sigma} \sqrt{\mathbf{x}_*^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*}$ Prediction interval:  $\mathbf{x}_*^T \hat{\boldsymbol{\beta}} \pm t_{0.025} \hat{\sigma} \sqrt{1 + \mathbf{x}_*^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*}$
- 7. To test  $\beta = 0$ ,  $F = \frac{\text{regression SS}}{\text{RSS}} \sim F(p, n p)$  (chk simple linear regression, but with different SS)
- 8. To test a more general linear hypothesis about the coefficients of the model,  $\mathbf{C}\boldsymbol{\beta} = \mathbf{d}$ , Extra  $SS = \left(\mathbf{C}\hat{\boldsymbol{\beta}} \mathbf{d}\right)^T \left(\mathbf{C}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{C}^T\right)^{-1} \left(\mathbf{C}\hat{\boldsymbol{\beta}} \mathbf{d}\right)$   $F = \frac{(ESS \text{ for } H_0)/c}{RMS} = \frac{(RSS \text{ under } H_0 RSS \text{ under full model})/c}{RMS} \sim F(k, n p)$