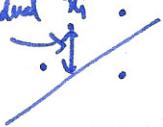


## 2.5 Warnings about correlation and regression

residuals: residual  $\hat{y}_i$ :  residual plot: 

• correlation does not imply causation.

### 10.1 Inference for linear regression

step 0 : plot data, check it looks linear.

model for simple linear regression:  $(x_1, y_1), \dots, (x_n, y_n)$   $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$

$\epsilon_i$  independent dist as  $N(0, \sigma^2)$  parameters for the model are  $\beta_0, \beta_1, \sigma$ .

recall, best fit line has slope  $b_1 = r \frac{s_y}{s_x}$  and intercept  $b_0 = \bar{y} - b_1 \bar{x}$ .

(there are no estimates for  $\beta_1$  and  $\beta_0$ ). what about  $\sigma$ ? look at residuals.

$e_i = \text{observation-predicted} = y_i - \hat{y}_i = y_i - b_0 - b_1 x_i$

Fact estimate for  $\sigma^2$  is  $s^2 = \frac{\sum e_i^2}{n-2} = \frac{\sum (y_i - \hat{y}_i)^2}{n-2}$   $s = \sqrt{s^2}$  is model standard deviation.

#### Model assumptions

- sample is an SRS
- linear relation in data (plot it)
- common standard deviation of residuals:  
- which is normal: use qqplot.  

Confidence intervals for slopes claim level C:  $b_1 \pm t^* \text{SE}_{b_1}$

$t^*$  critical C-value from  $t_{n-2}$ ,  $\text{SE}_{b_1} = \frac{s}{\sqrt{\sum (x_i - \bar{x})^2}}$   $\text{SE}_{b_0} = s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2}}$

#### Hypothesis tests for slopes

$H_0: \beta_1 = 0$  compute test statistic  $t = \frac{b_1}{\text{SE}_{b_1}}$  dist as  $t_{n-2}$ .

Estimated mean value at  $x^*$ :  mean  $\hat{y} = b_0 + b_1 x^*$

Level C confidence interval for the mean response:  $\hat{y} \pm t^* \text{SE}_{\hat{y}} = s \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$

Estimated value at  $x^*$ :  $\hat{y} = b_0 + b_1 x^*$  confidence interval  $\hat{y} \pm t^* \text{SE}_{\hat{y}} = s \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$

## §12.1 One way analysis of variance (ANOVA)

recall §7: comparing two means. aim: compare means of many groups.

<u>Example 12.2</u>	store	$\bar{x}_i$	$s_i$	$n_i$	pp	
A	38	8	50		$\mu_1$	$\sigma$
B	48	13	50		$\mu_2$	$\sigma$
C	42	11	50		$\mu_3$	$\sigma$
D	28	7	50		$\mu_4$	$\sigma$
E	35	10	50		$\mu_5$	$\sigma$

data points.

$x_{ij}$   $\uparrow \downarrow$   $i$ -th pop  $j$ -th person in sample.

$H_0$ : the means are all equal

$H_a$ : the means are not all equal.

model: each population has distribution  $N(\mu_i, \sigma)$  (here  $\mu_1, \dots, \mu_5$ )  $\sim \mu_A, \dots, \mu_E$ .

note: s.d. assumed to be same in each pop!

estimates for parameters:  $\bar{x}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij}$ . estimate for  $\mu_i$

standard deviation of each sample:  $s_i = \sqrt{\frac{\sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2}{n_i - 1}}$ .

note: ANOVA not particularly sensitive to differing  $\sigma_i$ , rough guide: if  $s_i$  vary by less than a factor of 2 usually OK.

pooled estimate for  $\sigma$ :  $s_p^2 = \frac{(n_1-1)s_1^2 + \dots + (n_I-1)s_I^2}{(n_1-1) + (n_2-1) + \dots + (n_I-1)}$   $s_p = \sqrt{s_p^2}$

note:  $s_p$  not average of  $s_i$ , it's weighted by size of sample.

stores.aov <- aov(age ~ store, data = stores)

numerical variable	categorical group	summary of fit				
		Df	Sum Sq	Mean Sq	F value	Pr(>F)
store	4	11154	2788.5	27.39	$10^{-16}$	
residuals	245	24974	101.8			

summary(stores.aov):

if  $H_0$  is true all mean are the same,  $\Rightarrow F = \frac{MSB}{MSE} \approx 1$  and dist as F-dist w/ df = I-1

$MSE = s_p^2$  estimate for within group variance

$MSB =$  estimate for variance between groups.

use stores.csv data for this