

Start 135 Dec 19

Today Finish Regression (sec 14.2)

Thursday Bayesian Starts (sec 8.6)

Monday Review with Adam 9-11 AM here

Tuesday no class. GST extended OH,

Wednesday Final 9-12 here.

Simple Linear regression (sec 14.2)

Last time

95% Confidence Interval for  $\beta_1$  is

$$\hat{\beta}_1 \pm t_{n-2}(0.025) S_{\hat{\beta}_1}$$

where  $S_{\hat{\beta}_1} = \frac{s^2}{\text{var}(x)} = \frac{\text{RSS}}{(n-2) \text{var}(x)}$ .

SE of regression line

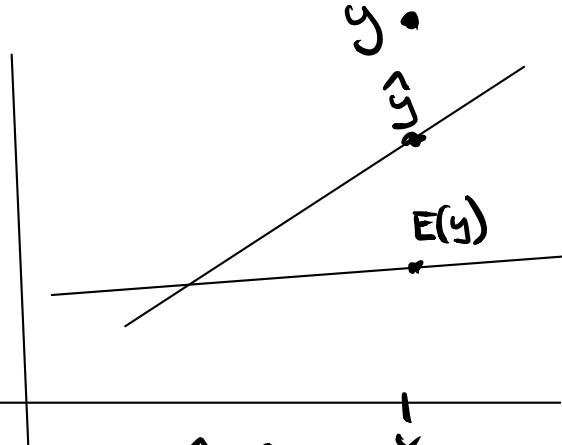
Let  $x$  be a pt on  $x$  axis.

If  $y$  is paired with  $x$ ,  $y = \beta_0 + \beta_1 x + e$

$\Rightarrow E(y) = \beta_0 + \beta_1 x$  so  $E(y)$  lies on True line,

Denote  $y$  on regression line by  $\hat{y}$

Picture



regression line

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

True line

$$E(y) = \beta_0 + \beta_1 x$$

$$E(\hat{y}) = E(\hat{\beta}_0 + \hat{\beta}_1 x) = \beta_0 + \beta_1 x \text{ so } \hat{y} \text{ is an unbiased est of } E(y).$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x = (\bar{y} - \hat{\beta}_1 \bar{x}) + \hat{\beta}_1 x$$

$$\hat{y} = \bar{y} + \hat{\beta}_1 (x - \bar{x})$$

$$\text{Var}(\hat{y}) = \text{Var}(\bar{y} + \hat{\beta}_1 (x - \bar{x}))$$

$$= \text{Var}(\bar{y}) + (x - \bar{x})^2 \text{Var}(\hat{\beta}_1)$$

$$+ 2(x - \bar{x}) \text{Cov}(\bar{y}, \hat{\beta}_1)$$

|| 0

Show  $\text{Cov}(\bar{y}, \hat{\beta}_1) = 0$

$$\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = -\frac{\sigma^2 n \bar{x}}{n^2 \text{Var}(x)}$$

||

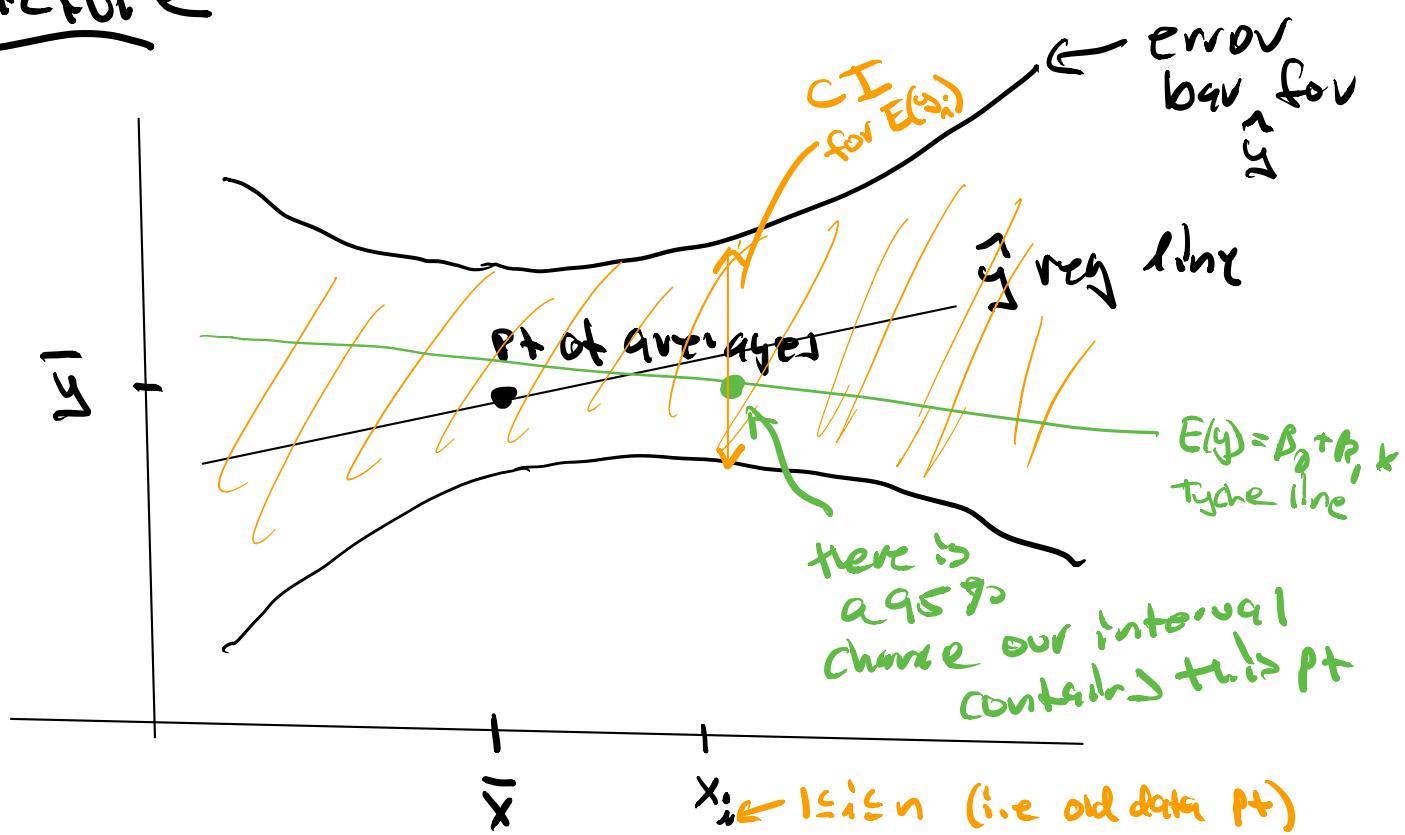
$$\text{Cov}(\bar{y} - \hat{\beta}_1 \bar{x}, \hat{\beta}_1) = \text{Cov}(\bar{y}, \hat{\beta}_1) - \bar{x} \text{Var}(\hat{\beta}_1)$$

+  $\frac{n \sigma^2}{n^2 \text{Var}(x)}$

$$\Rightarrow \text{Cov}(\bar{y}, \hat{\beta}_1) = 0 \quad \checkmark$$

$$\text{so } \text{Var}(\hat{y}) = \frac{\sigma^2}{n} + \underbrace{(x - \bar{x})^2 \sigma^2}_{\text{nvar}(x)}$$

# Picture



## Confidence Interval for True line $E(y)$

A 95% CI for  $E(y) = \beta_0 + \beta_1 x$  means at  $\hat{y}$

is  $\hat{y} = t_{n-2}(0.025) S_{\hat{y}}$  where

$$S_{\text{y}\bar{y}}^2 = S^2 \left( \frac{1}{n} + \frac{\overline{(x - \bar{x})^2}}{\text{var}(x)} \right)$$

Prediction Interval (PI) for a new  
new observation  $y_{n+1}$ . — not in  
book,

A confidence interval for a RV  
is called a prediction interval.

A new observation  $y_{n+1}$  is a RV.

You give  $x_{n+1}$  and you get  $y_{n+1}$ .

Find a 95% PI for  $y_{n+1}$ .

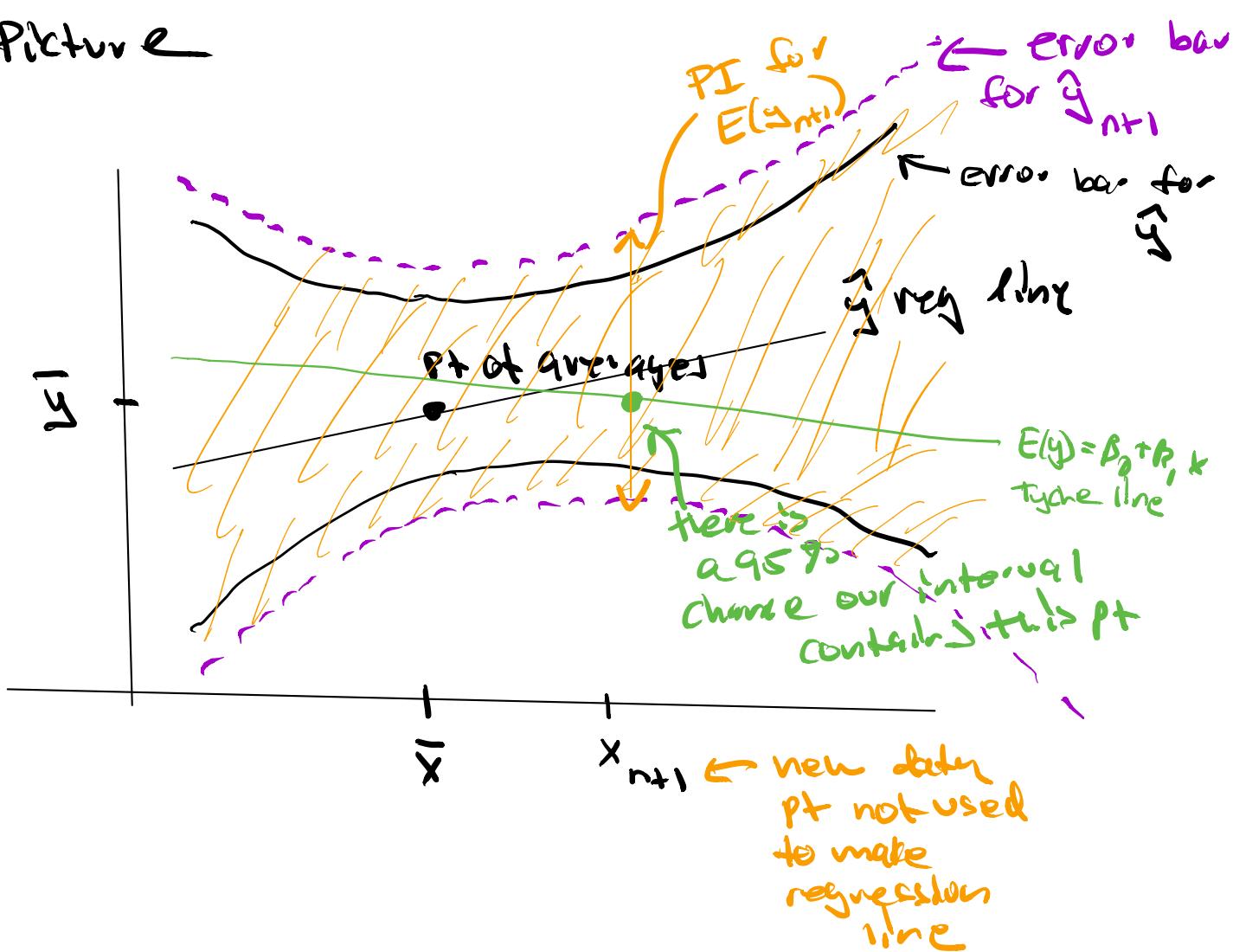
Then A 95% PI for a new observation  
 $y_{n+1}$  for predictor  $x_{n+1}$  is

$$\hat{y}_{n+1} \pm t_{n-2}(0.025) S \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\text{var}(x)}}$$

$\frac{\text{RSS}}{n-2}$

This is similar to 95% CI of  $E(y)$   
but has more uncertainty,

## Picture

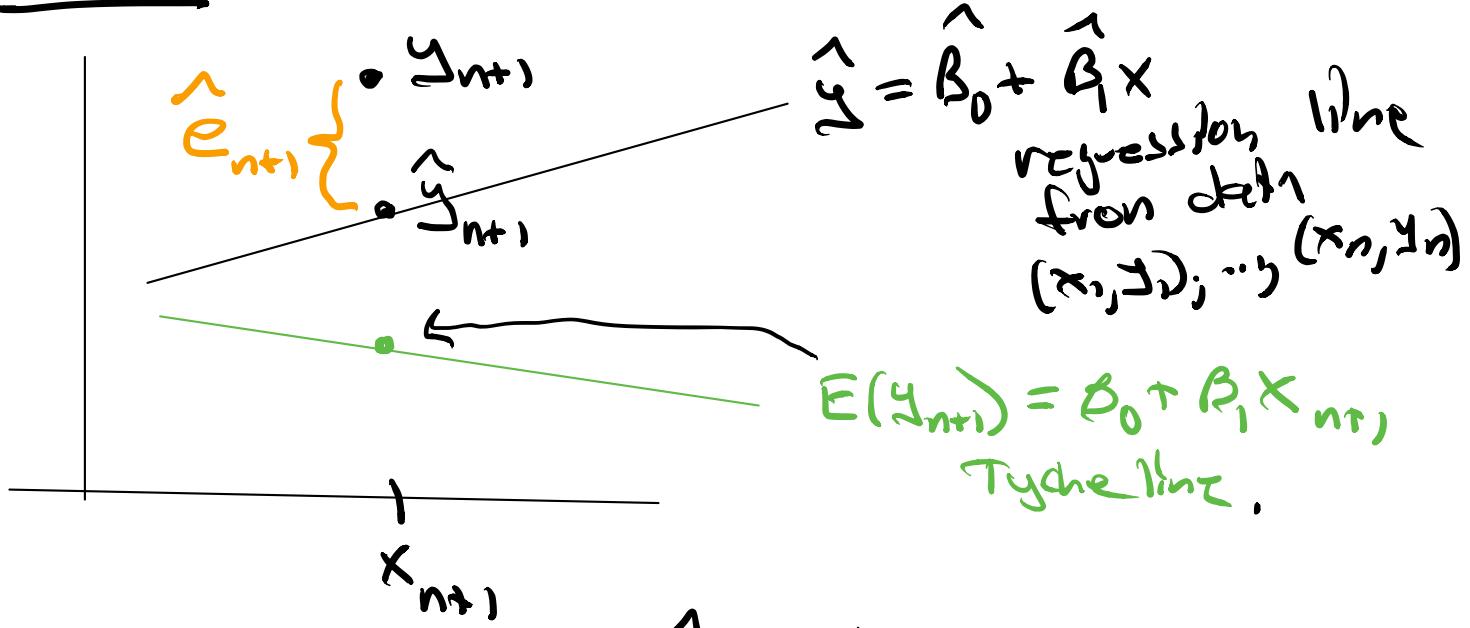


Notice that the error bar is wider for predicting  $\hat{y}_{n+1}$  given  $(x_1, y_1), \dots, (x_n, y_n)$  and  $x_{n+1}$ .

Proof :

When  $y_{n+1}$  is a new response, not used in the calculation of the regression line, it is independent of  $\hat{y}_{n+1}$ . Consequently the variance of the residual  $\hat{e}_{n+1}$  is larger than the variance of the residual  $\hat{e}_i$  for  $1 \leq i \leq n$ .

## Picture



The variance of  $\hat{e}_{n+1}$  is larger than the variance of  $\hat{y}_{n+1}$ .

We will see that

$$\begin{aligned}
 \text{var}(\hat{e}_{n+1}) &= \text{var}(y_{n+1}) + \text{var}(\hat{y}_{n+1}) \\
 &\quad \frac{\sigma^2}{n} + \frac{(x_{n+1} - \bar{x})^2 \sigma^2}{n \text{var}(x)} \\
 &= \sigma^2 \left( 1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{n \text{var}(x)} \right)
 \end{aligned}$$

so why is  $y_{n+1}$  and  $\hat{y}_{n+1}$  indep?

$\hat{y}_{n+1}$  is a pt on regression line that may calculated with  $(x_1, y_1), \dots, (x_n, y_n)$  not  $(x_{n+1}, y_{n+1})$ .

$$\text{we have } E(\hat{\epsilon}_{n+1}) = E(Y_{n+1} - \hat{Y}_{n+1}) \\ \approx E(Y_{n+1}) - E(\hat{Y}_{n+1}) = 0$$

$$E(\hat{\beta}_0 + \hat{\beta}_1 X_{n+1}) \\ \approx \beta_0 + \beta_1 X_{n+1}$$

$$\text{so } Y_{n+1} - \hat{Y}_{n+1} \sim N(0, \sigma^2 \left( 1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{n \text{var}(x)} \right))$$

$$\Rightarrow \frac{Y_{n+1} - \hat{Y}_{n+1}}{S \left( 1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{n \text{var}(x)} \right)} \sim t_{n-2}$$

$$\Rightarrow \hat{Y}_{n+1} \pm t_{n-2} (0.075) S \sqrt{\left( 1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{n \text{var}(x)} \right)}$$

is a 95% PI

□

----- break -----

## F statistic for regression — not in book,

For HW 6 you will show

$$\sum (y_i - \bar{y})^2 = \sum (y_i - \hat{y})^2 + \sum (\hat{y} - \bar{y})^2$$

$SS_{TOT}$

"

TSS

\  
total

$SS_{error}$

"

RSS

\  
residual

$SS_{model}$

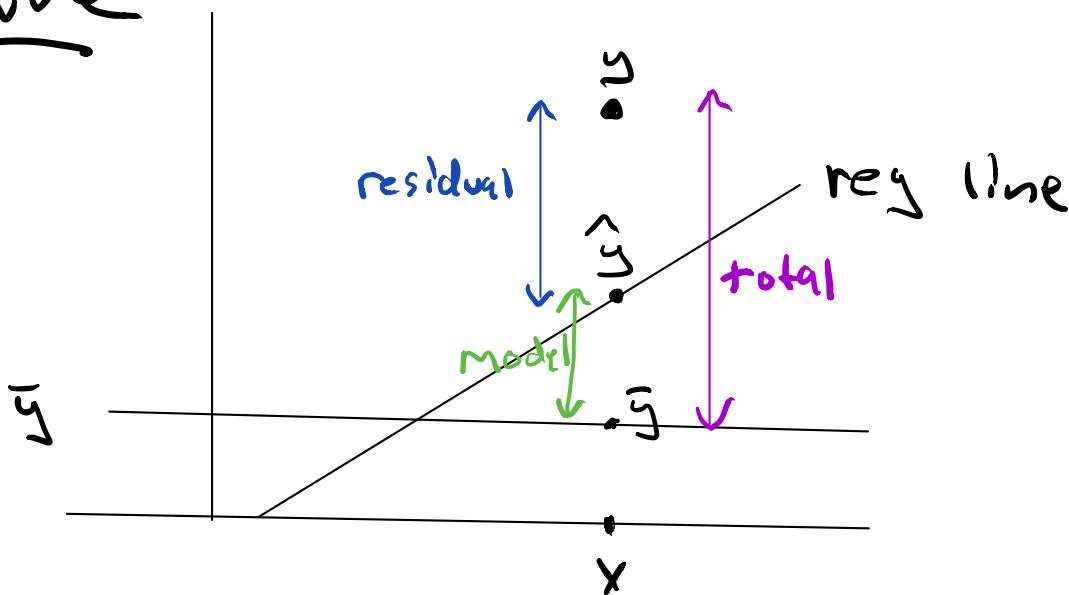
"

MSS

\  
model

i.e. we can divide the total variability of the dependent variable into two components > ; the variability of the regression line about the mean and the residual variation of observations about the regression line.

## Picture



We showed RSS has  $n-2$  d.f

TSS has  $n-1$  d.f

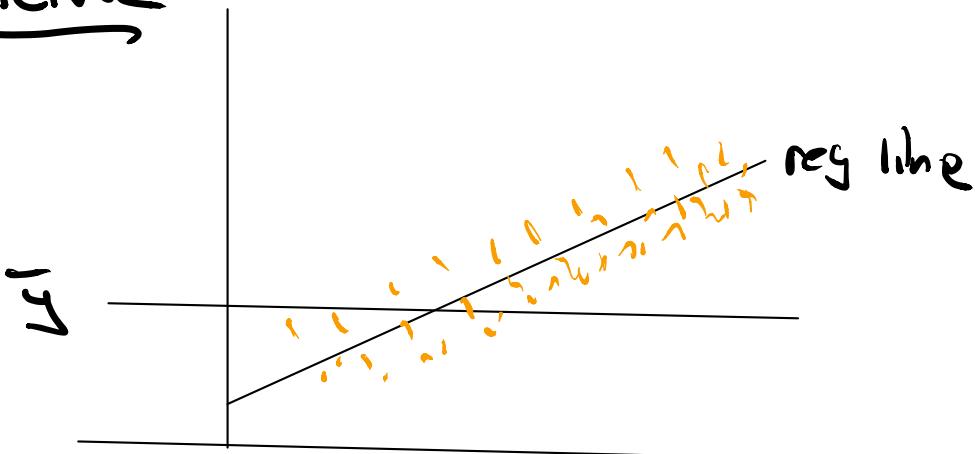
RSS and MSE are independent

$\chi^2$  so their d.f add.  $\Rightarrow$  d.f MSE is 1.

This allows us to do hypothesis testing,

null  $B_1 = 0$  (no relationship between)  
alt  $B_1 \neq 0$   $x$  and  $y$

Pictorial



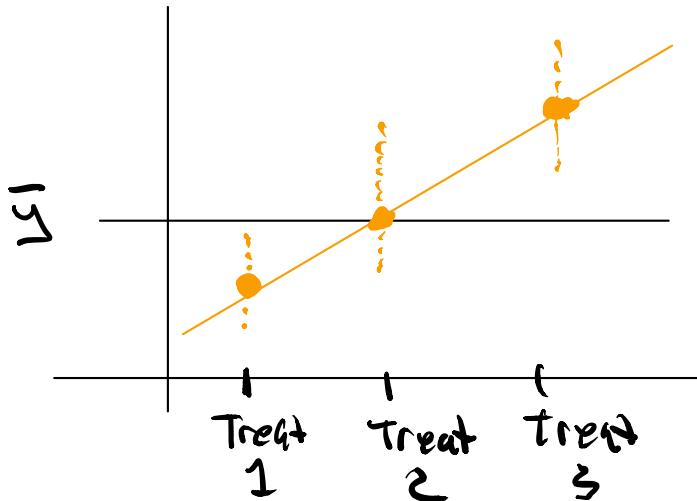
the spread of these  $\hat{y}$  values about  $\bar{y}$  is bigger  
the steeper  $B_1$  gets

If the null is true ( $B_1 = 0$ )  
then MSE is very small.

For a fixed TSS if MSS is large then RSS is small.

### An aside Anova Picture

$$SST = SSW + SSB$$



Informally, if you plot values of  $y$  for different treatment groups, the regression line goes through the pt of averages for each grp and you can see that  $RSS = SSW$ ,  $MSS = SSB$ . So ANOVA really is regression.

End of aside.

Recall (194 char 6) def" F dist :

Let  $U$  and  $V$  be indep  $\chi^2$  RV w/

d.f.  $m$  and  $n$  respectively. Then

$$\frac{U/m}{V/n} \sim F_{m,n}$$

with  $K$  predictors  
 $K = n - k - 1$

For vs,

$$\frac{\left(\frac{MSS}{\sigma^2}\right)_1}{\left(\frac{RSS}{\sigma^2}\right)_{n-2}} = \frac{MSS}{\left(\frac{RSS}{\sigma^2}\right)} \sim F_{1, n-2}$$

we reject the null when

$F_{1, n-2}$  is large.

Not surprisingly

$$F_{1, n-2} = t_{n-2}^2$$

For simple linear regression F test is equivalent to t test.

t-test  
for  $\beta_1 \neq 0$ .

$R^2$  - the coefficient of determination

In th 6 you will show  $\hat{\sigma}_{\epsilon}^2 = (1-r^2) \sigma_y^2$

we have  $\hat{\sigma}_{\epsilon}^2 = \frac{RSS}{n} = \frac{\sum (y_i - \hat{y}_i)^2}{n}$

and  $\sigma_y^2 = \frac{TSS}{n} = \frac{\sum (y_i - \bar{y})^2}{n}$

so  $RSS = (1-r^2)TSS$

$$\Rightarrow r^2 = 1 - \frac{RSS}{TSS} = \frac{TSS - RSS}{TSS} = \frac{MSS}{TSS}$$

called  $R^2$  in multiple regression,

## interpretation

$R^2$  is the fraction of the total variation of  $y$  "explained" by the model,

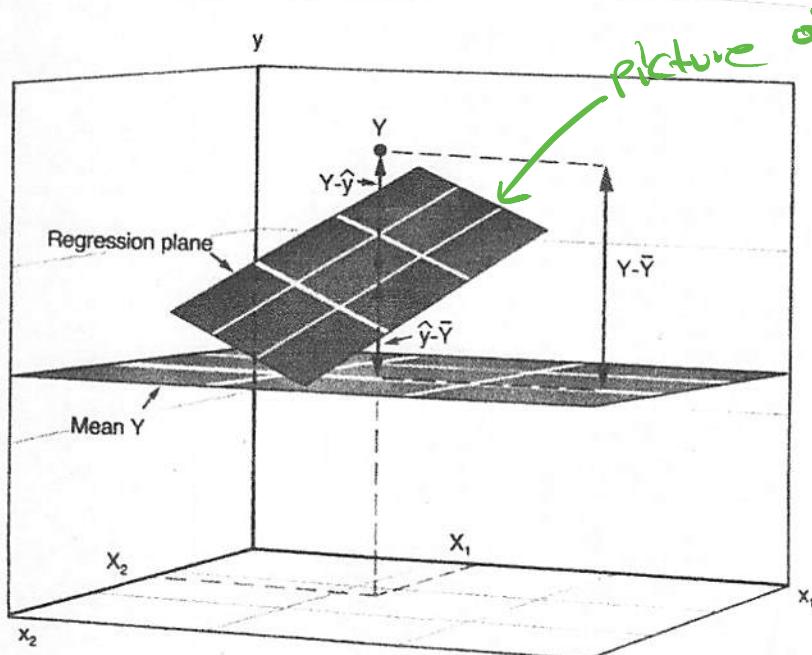
If  $MSE$  is very small the slope slope of the regression line is small, — i.e  $r^2$  close to zero.

If  $MSE = TS$  the data lies on a straight line, — i.e  $r^2 = 1$

----- End of Simple linear regression -----

Multiple linear regression is similar.

If we assume that our  $k$  predictors are independent we can do an F test whether  $\beta_1, \beta_2, \dots, \beta_k$  are all zero. For simple linear regression we only have  $\beta_1$  and F test is equivalent to t test as mentioned above.



**FIGURE 3-5** The deviation of the observed value of  $Y$  from the mean of all values of  $Y$ , ( $Y - \bar{Y}$ ), can be separated into two components: the deviation of the observed value of  $Y$  from the value on the regression plane ( $Y - \hat{Y}$ ) at the associated values of the independent variables  $X_1$  and  $X_2$ , and the deviation of the regression plane from the observed mean value of  $\bar{Y}$  ( $\hat{Y} - \bar{Y}$ ) (compare with Fig. 2-7).

We have

$$TSS = MSS + RSS$$

$\nwarrow n-1 \text{ d.f.} \quad \nwarrow K \text{ d.f.} \quad \nwarrow n-K-1 \text{ d.f.}$

Just like we had before,

$$\text{Let } R^2 = \frac{MSS}{TSS} \quad \text{coeff. of determination}$$

$$\text{Let } F = \frac{\frac{MSS}{K}}{\frac{RSS}{n-K-1}} = \frac{\frac{MSS}{TSS \cdot K}}{\frac{RSS}{TSS \cdot (n-K-1)}} = \frac{\frac{MSS}{TSS \cdot K}}{\frac{TSS - MSS}{TSS \cdot (n-K-1)}} = \frac{\frac{R^2}{K}}{\frac{1-R^2}{n-K-1}}$$

this is our  
T.S.

If we reject the null it is natural to ask which of the  $B_i$ ,  $i=1, 2, \dots, k$  are statistically different from zero? In other words which of the variables  $x_1, x_2, \dots, x_k$  have a relationship with  $y$ .

You will explore all of this hopefully in your next stats class !!