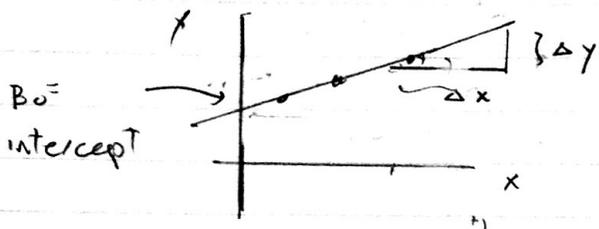


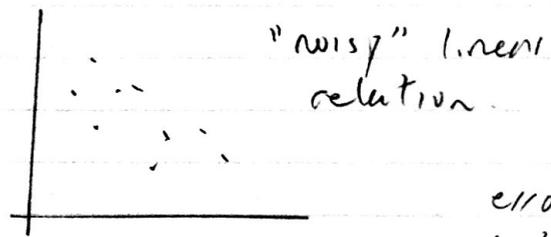
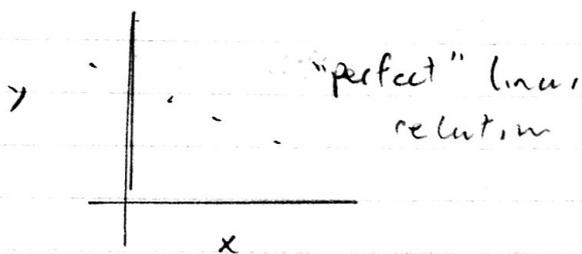
# Simple Linear Regression

$$|\text{cor}(\underline{x}, \underline{y})| = 1 \iff \underline{y} = \beta_0 \underline{1} + \beta_1 \underline{x} \quad \text{some } \beta_0, \beta_1$$



$$\begin{aligned} \frac{\Delta y}{\Delta x} &= \frac{(\beta_0 + \beta_1(x + \Delta x)) - (\beta_0 + \beta_1 x)}{\Delta x} \\ &= \beta_1 \Delta x / \Delta x = \beta_1 = \text{"slope"} \end{aligned}$$

More typically,  $|\text{cor}(\underline{x}, \underline{y})| < 1$



error,  
disturbance,  
deviation,

$$\begin{aligned} \underline{y} &= \beta_0 \underline{1} + \beta_1 \underline{x} \\ y_i &= \beta_0 + \beta_1 x_i \end{aligned}$$

$$\begin{aligned} \underline{y} &= \beta_0 \underline{1} + \beta_1 \underline{x} + \underline{e} \\ y_i &= \beta_0 + \beta_1 x_i + e_i \end{aligned}$$

what values of  $\beta_0, \beta_1$  provide the "best fit"?

## OLS Regression Line

$$\begin{aligned} \text{RSS}(\beta_0, \beta_1) &= \sum (y_i - (\beta_0 + \beta_1 x_i))^2 = \|\underline{y} - (\beta_0 \underline{1} + \beta_1 \underline{x})\|^2 \\ &= \sum e_i^2 = \|\underline{e}\|^2 \end{aligned}$$

$e_i = y_i - (\beta_0 + \beta_1 x_i)$  = residual for the particular values of  $\beta_0, \beta_1$ .

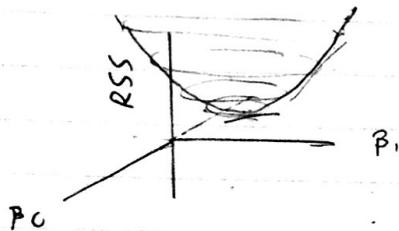
Defn: The OLS values of  $(\beta_0, \beta_1)$  are the values  $(\hat{\beta}_0, \hat{\beta}_1)$  that minimize  $\text{RSS}(\beta_0, \beta_1)$

$$(\hat{\beta}_0, \hat{\beta}_1) = \arg \min_{(\beta_0, \beta_1)} \text{RSS}(\beta_0, \beta_1)$$

## Optimization:

Method 1: geometry (will do later)

Method 2: calculus



\*  $RSS(\beta_0, \beta_1)$  is a quadratic (convex) function of  $(\beta_0, \beta_1)$

\* global minimum occurs where derivative (gradient) is zero.

$$\frac{\partial RSS}{\partial \beta_0} = -2 \sum (y_i - (\beta_0 + \beta_1 x_i)) = 0 \Leftrightarrow \sum (y_i - (\beta_0 + \beta_1 x_i)) = 0$$

$$\frac{\partial RSS}{\partial \beta_1} = -2 \sum x_i (y_i - (\beta_0 + \beta_1 x_i)) = 0 \Leftrightarrow \sum x_i (y_i - (\beta_0 + \beta_1 x_i)) = 0$$

The OLS values  $(\hat{\beta}_0, \hat{\beta}_1)$  therefore satisfy...

$$\textcircled{1} \quad \sum (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)) = 0$$

$$\textcircled{2} \quad \sum x_i (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)) = 0$$

These are called the normal equations for simple linear regression.

Why "normal equations"? at  $\hat{e}_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) =$  "resid at OLS val." or just "residual"

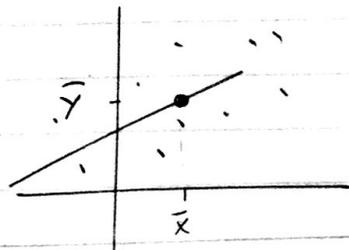
$$\text{then } \textcircled{1} \Rightarrow \sum \hat{e}_i = 0 \quad \hat{\underline{e}} \cdot \underline{1} = 0$$

$$\textcircled{2} \Rightarrow \sum x_i \hat{e}_i = 0 \quad \hat{\underline{e}} \cdot \underline{x} = 0$$

$\Rightarrow \hat{\underline{e}}$  is normal (orthogonal) to the vectors  $\underline{1}, \underline{x}$ .

$$\text{Also, note from } \textcircled{1} \text{ that } \sum y_i = \sum (\hat{\beta}_0 + \hat{\beta}_1 x_i)$$

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



regression line goes through the point  $(\bar{x}, \bar{y})$

## Solving the normal equations

$$\begin{aligned} \textcircled{1} \quad \sum (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)) &= 0 \\ \bar{y} &= \hat{\beta}_0 + \hat{\beta}_1 \bar{x} \\ \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \end{aligned}$$

Plug this into  $\textcircled{2}$ :

$$\begin{aligned} \sum x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ \sum x_i (y_i - \bar{y} + \hat{\beta}_1 \bar{x} - \hat{\beta}_1 x_i) &= 0 \end{aligned}$$

$$\sum x_i (y_i - \bar{y}) = \hat{\beta}_1 \sum x_i (x_i - \bar{x}) \quad \#$$

Now note the following trick:  $\sum (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i (y_i - \bar{y}) + \sum \bar{x} (y_i - \bar{y})$   
 $= \sum x_i (y_i - \bar{y}) + \bar{x} \sum (y_i - \bar{y})$   
Similarly,  $\sum x_i (x_i - \bar{x}) = \sum (x_i - \bar{x})(x_i - \bar{x}) = \sum x_i (x_i - \bar{x}) = \sum x_i (y_i - \bar{y})$   
 $= \sum (x_i - \bar{x})^2$

$$\begin{aligned} \text{let } S_{XX} &= \sum (x_i - \bar{x})^2 \\ S_{XY} &= \sum (x_i - \bar{x})(y_i - \bar{y}) \end{aligned}$$

Then  $\#$  says  $S_{XY} = \hat{\beta}_1 S_{XX}$ ,  $\hat{\beta}_1 = \underline{\underline{S_{XY} / S_{XX}}}$

Finally, the OLS values are:  $\begin{bmatrix} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 = S_{XY} / S_{XX} \end{bmatrix}$

## Relation to correlation

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

← looks like covariance or correlation  
← looks like variance of  $x$

$$= \frac{S_{XY}}{S_{XX}} = \left( \frac{S_{XY}}{S_{XX}} \right)^{1/2} \frac{S_{XY}}{S_{XY}^{1/2} S_{XX}^{1/2}}$$

$$= \underline{\underline{\left( \frac{S_y}{S_x} \right) \times \text{Cor}(x, y)}}$$

What you need to find  $\hat{\beta}_0, \hat{\beta}_1$ :  $\begin{bmatrix} \bar{x}, \bar{y} \\ s_x^2, s_y^2 \\ \text{Cor}(x, y) \end{bmatrix}$

### Regression Line

$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{e}_i$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x = \text{"predicted" or "fitted" value, at } x.$$

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

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

$$= \bar{y} + \frac{s_y}{s_x} \times \text{Cor}(x, y) \times (x - \bar{x})$$

### Questions

1) How does the line change if we change the location of the  $x$ 's?

2) How does the line change if we change the scale of the  $x$ 's?

3) What aspects of the regression line are all invariant to

- a) location changes?
- b) scale changes?
- c) loc. + scale changes?