

# Regression (Chapter 11)

*as opposed to qualitative or categorical*

**Setup:** Have data on two quantitative variables —  $X$  and  $Y$  — on a sample of  $n$  subjects.

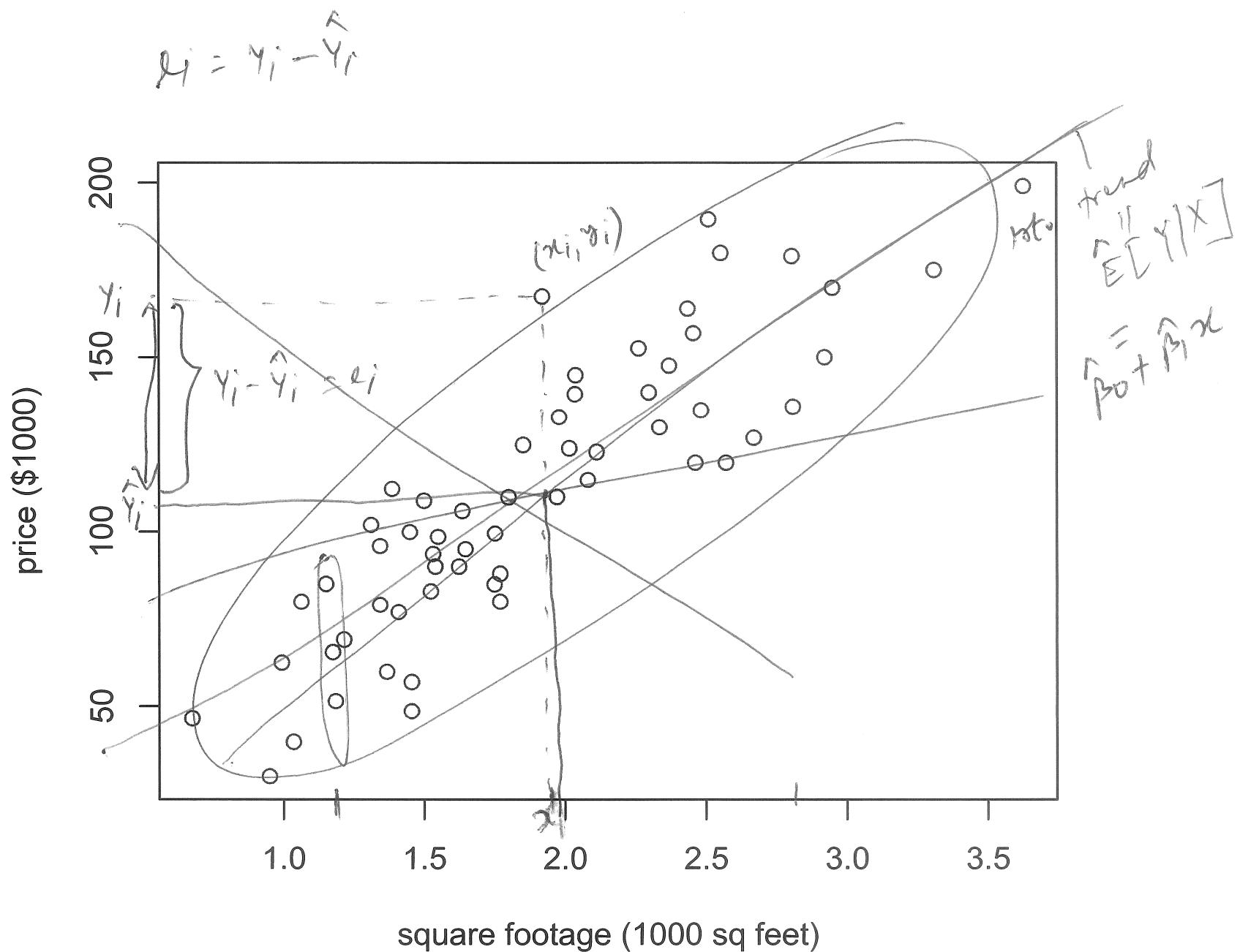
**Q:** Is there any association between  $X$  and  $Y$ ? What kind?

**Scatterplot:**

Data:  $(y_i, x_i), i = 1, 2, \dots, n$ .

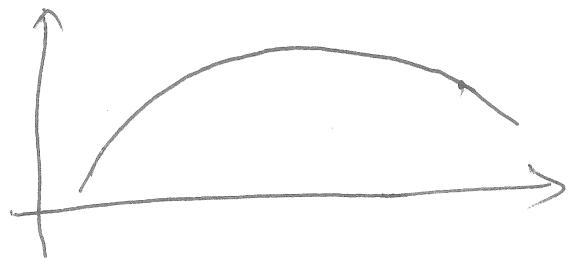
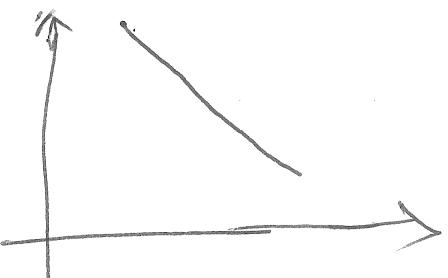
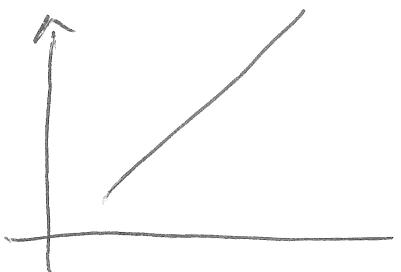
T coming from  $i$ th subject.  $\{E[Y|X]\}$

- Plot  $y$  against  $x$
- Look for the trend in the plot — a smooth curve that shows how the average value of  $Y$  changes with  $x$
- Trend may be linear or non-linear
- If there is a trend, then the two variables are associated. In this case,  $x$  may be used to predict  $y$
- Trend may be strong or weak. It is strong if the points are tightly clustered around the trend (small scatter)
- No trend: No association — i.e., the variables are independent, and  $x$  is not helpful for predicting  $y$ .
- Predicted response :  $\hat{Y} = \hat{E}[Y|X]$



## Overall pattern in a scatterplot

Form:



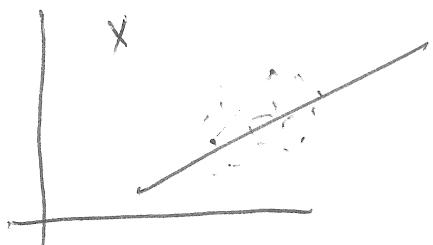
Direction: (only for linear trend)

- +ve or -ve

Strength - assess the scatter of the points: around the trend:

Small scatter  $\Rightarrow$  relationship b/w Y and X is strong  
Large scatter  $\Rightarrow$  " " " " " is weak.

Outliers - observations that don't follow overall pattern:



Focus on  
linear trend

# Correlation Coefficient

**Population correlation:** A measure of linear relationship between  $X$  and  $Y$ . It is defined as,

$$\rho = \frac{\text{cov}(X, Y)}{\text{sd}(X)\text{sd}(Y)},$$

where  $\text{cov}(X, Y) = E\{(X - E(X))(Y - E(Y))\}$  is covariance between  $X$  and  $Y$ .

**Sample correlation:** Estimator of  $\rho$ . It is defined as

$$r = \frac{S_{xy}}{S_x S_y},$$

where  $S_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$  — sample covariance

$$S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

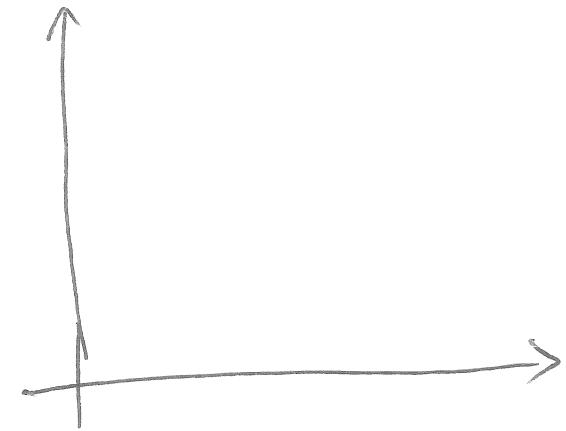
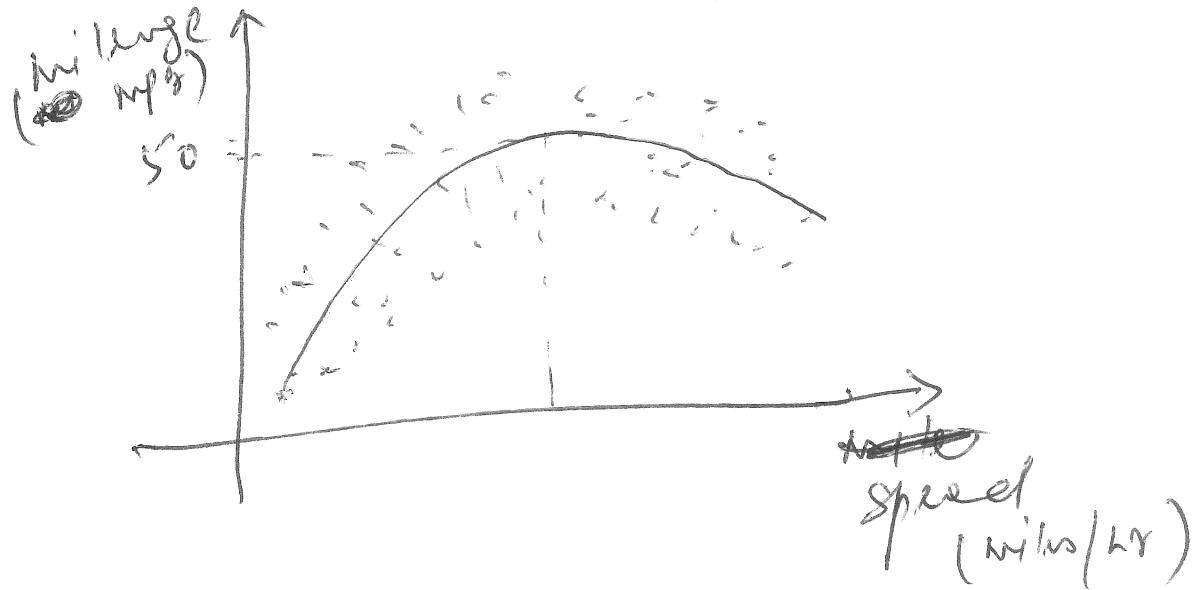
$$S_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$



## Properties of $\rho$ and $r$ :

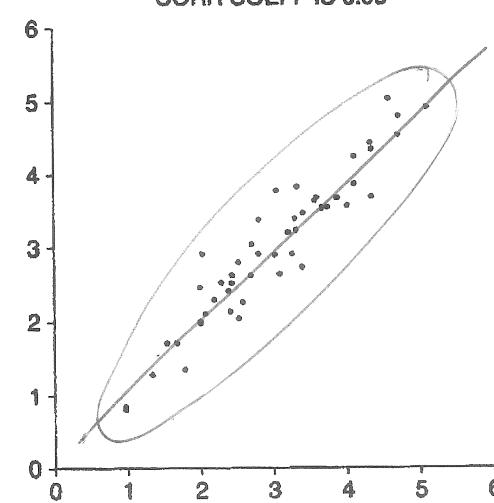
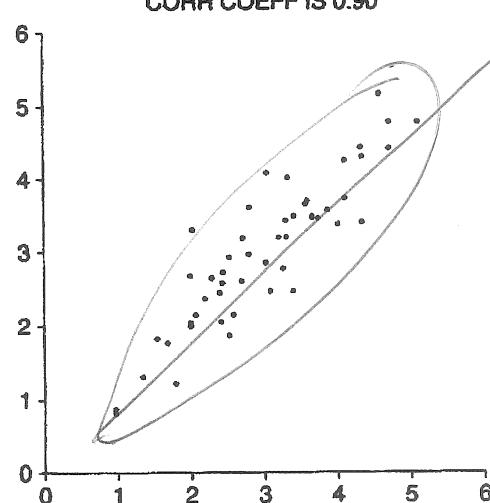
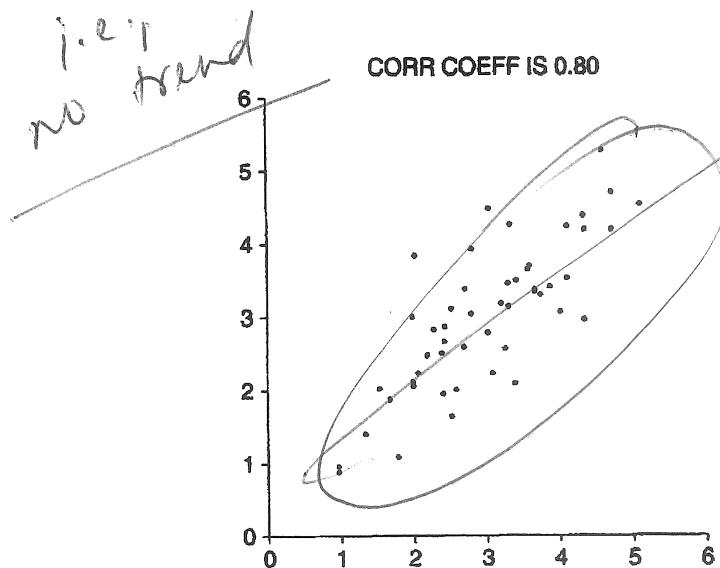
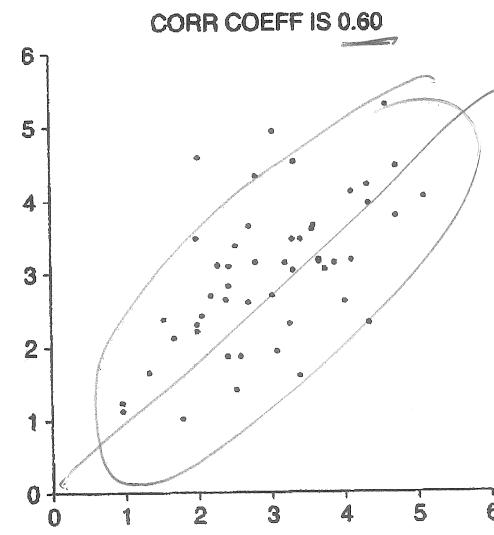
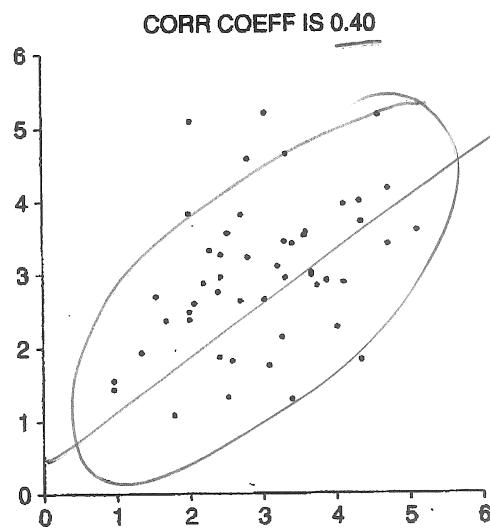
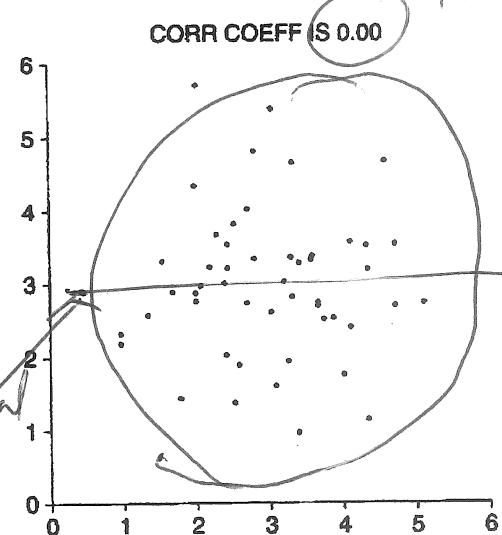
- Range between  $-1$  to  $1$
- Sign tells us direction of the linear relationship
- Absolute value tells us strength of the linear relationship
- Perfect correlation:  $|\rho|=1 \approx |r|=1 \Rightarrow Y = a+bX$
- Unit free
- No change if  $X$  and  $Y$  are interchanged or if  $X$  is replaced by  $aX + b$  and/or  $Y$  is replaced by  $cY + d$ , where  $a$  and  $c$  have the same sign. The sign will reverse if  $a$  and  $c$  have different signs.
- Zero correlation: No linear relationship. But there may be non-linear relationship.
- Independent  $X$  and  $Y$ : zero correlation, but the converse may not be true

Example non-linear relationship:

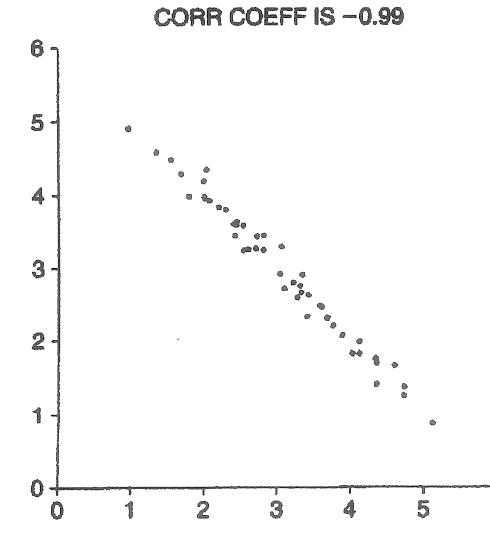
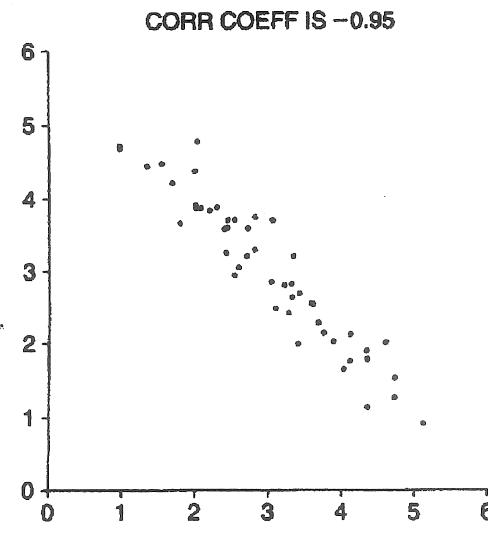
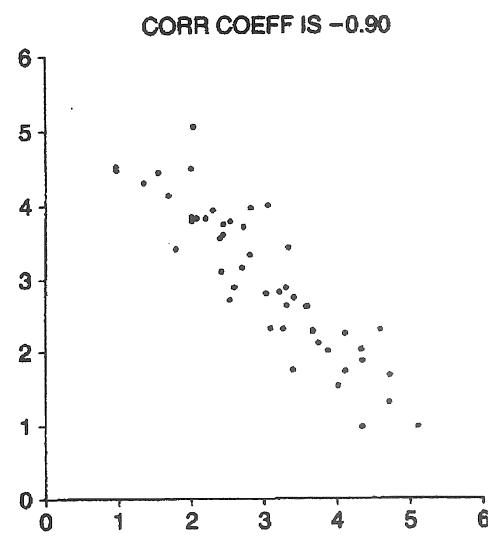
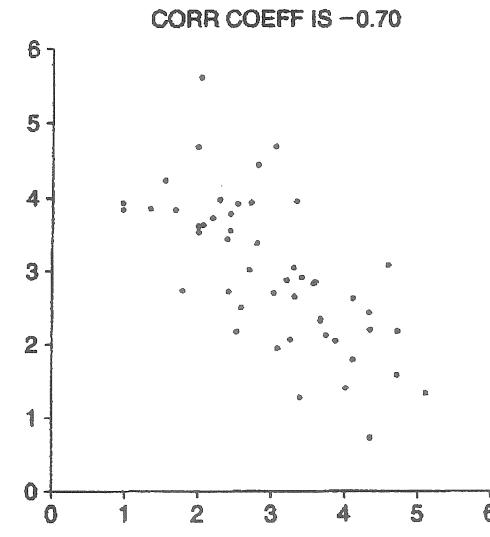
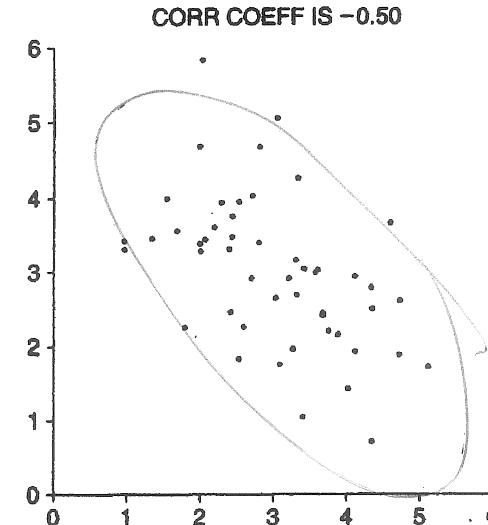
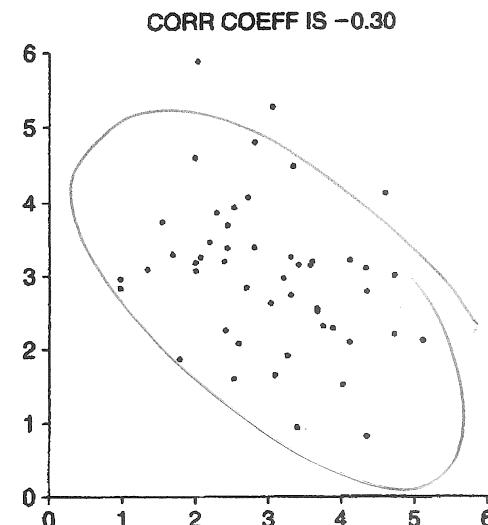


# Examples of Positive Correlation

$\rightarrow$  independent



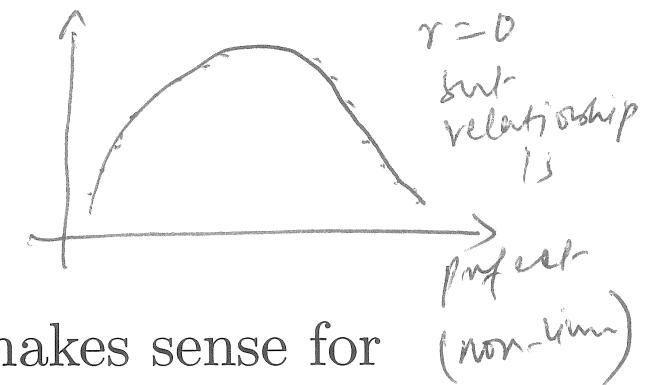
# Examples of Negative Correlation



# Caution: Non-linear relationships

Ex: Scatter plot of speed and mileage (miles per gallon) of an automobile.

More non-linear patterns: Any ~~thing~~ trend that is not linear is non-linear



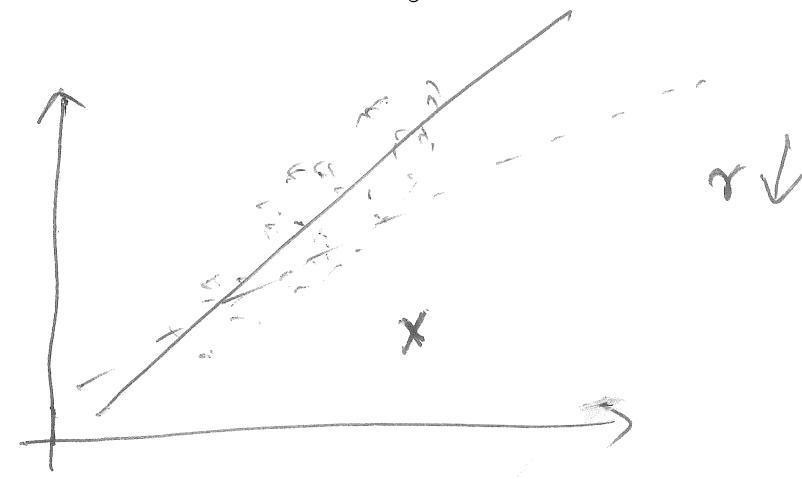
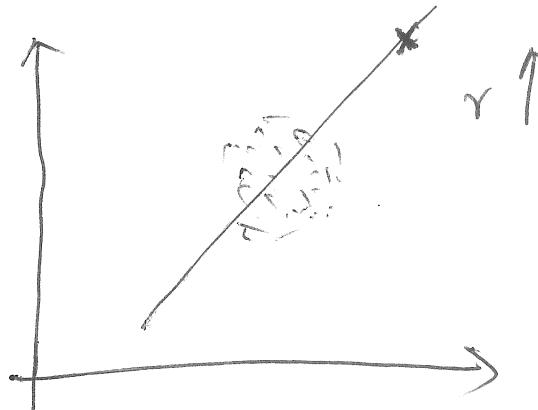
You can always compute  $r$  — but it doesn't make sense for curve patterns.

**Lesson:** Correlation only measures the strength of **linear** association — i.e., how close the points are to a straight line.

(case of friend)  
Linear trend

## Caution: Outliers

The position of an outlier relative to the rest of the (“cloud”) of points determines how it affects  $r$ . Outliers may decrease or increase the value of  $r$ .



**Note:** Just knowing the value of  $r$  will give no information about whether outliers are present. That's why it is important to look at scatterplots.

for the linear trend

# Simple Linear Regression

**Setup:** Have data  $(X_i, Y_i), i = 1, \dots, n$ , on two quantitative variables  $X$  &  $Y$ . Their scatterplot shows a linear relationship. Need an equation that would allow us to predict  $Y$  from  $X$ .

**Response variable ( $Y$ ):** variable to be predicted (or modeled), aka, *dependent variable*.

**Predictor ( $X$ ):** variable used to predict  $Y$ , aka, *independent or explanatory variable or covariate*.

**Regression model:** A function that models mean response —  $E(Y|X = x)$  — as a function of  $x$

**Simple linear regression:**  $E(Y|X = x) = \beta_0 + \beta_1 x$

- Assumes mean response changes *linearly* with  $x$
- $\beta_0$  : intercept —  $E(Y|X = 0)$
- $\beta_1$  : slope — rate of change of mean response. It represents the change in mean when  $x$  increases by 1 unit.

- The regression coefficients are estimated from data.
- Let  $(\hat{\beta}_0, \hat{\beta}_1)$  = estimator of  $(\beta_0, \beta_1)$ .

**Observed response:**  $Y_i$  when  $X = x_i, i = 1, \dots, n$ .

**Fitted response:**  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x$   $\hat{Y} = \hat{E}[Y|x=x] = \hat{\beta}_0 + \hat{\beta}_1 x$

$\uparrow$  Estimated trend.  
 $\uparrow$  Estimated II  
 $\uparrow$  Estimated  
 $\uparrow$  (Fitted)  
 $\uparrow$  regression  
 $\uparrow$  line.

- Estimated mean response when  $X = x$
- Response predicted by the regression line
- $(\hat{Y}, x)$  falls on the regression line
- Fitted values:  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i, i = 1, \dots, n$ .

**Residuals:**  $e_i = Y_i - \hat{Y}_i, i = 1, \dots, n$ .

- Vertical distance between observed and predicted  $Y$ 's
- Error in prediction
- Large residuals: observed and fitted  $Y$ 's are too far

**Least squares method for estimating coefficients:** Find  $(\hat{\beta}_0, \hat{\beta}_1)$  that minimize the *sum of squares of residuals*

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2.$$

measures  
 how close the  
 line is to  
 data.

our focus.

- Results in the **line of best fit** — the line is such that the fitted  $Y$ s are “closest” to the observed  $Y$ s
- Fitted regression line:  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x$
- Other criteria possible, e.g., minimizing  $\sum_{i=1}^n |e_i|$ , but the resulting estimates don’t have simple expressions

(but note that  $\sum_{i=1}^n e_i$  is not a valid)

To minimize  $\sum_{i=1}^n e_i^2$  wrt  $(\beta_0, \beta_1)$ , solve the normal equations

*RSS* → *residual sum of squares*

$$\frac{\partial \sum_{i=1}^n e_i^2}{\partial \beta_0} = 0, \quad \frac{\partial \sum_{i=1}^n e_i^2}{\partial \beta_1} = 0,$$

resulting in the **least squares estimates**

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}, \quad \hat{\beta}_1 = r S_y / S_x,$$

where  $r$  is **sample correlation**, and  $S_x$  and  $S_y$  are **standard deviations** of  $x$  and  $y$  samples, respectively.

Recall that:

$$S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad S_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2,$$

$$r = \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{S_x S_y}$$

# The fitted regression line

Fitted regression line:  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x$ . Plugging-in  $\hat{\beta}_0$  and  $\hat{\beta}_1$ ,

$$\begin{aligned}\hat{Y} &= (\bar{y} - \hat{\beta}_1 \bar{x}) + \hat{\beta}_1 x = \bar{y} + \hat{\beta}_1(x - \bar{x}) \\ &= \bar{y} + r \frac{s_y}{s_x} (x - \bar{x}) \Rightarrow \frac{(\hat{Y} - \bar{Y})}{s_y} = r\end{aligned}$$

$$\frac{(x - \bar{x})}{s_x}$$

"z-score of  $x$ "

implying that

$$\frac{\hat{Y} - \bar{Y}}{s_y} = r \frac{(x - \bar{x})}{s_x}.$$

"z-score of  $\hat{Y}$ "

$$|r| \leq 1$$

- If  $x$  is 1 SD away from its mean  $\bar{x}$ ,  $\hat{Y}$  is  $r$  SD away from its mean  $\bar{Y}$ . Since  $|r| \leq 1$ , this means  $\hat{Y}$  is **closer** to  $\bar{Y}$  (in units of SD) than  $x$  is to  $\bar{x}$  — **regression toward mean**.
- The fitted line passes through the points  $(\bar{x}, \bar{y})$ .
- The sign of slope  $\hat{\beta}_1$  is same as the sign of  $r$ .
- The sum of residuals,  $\sum_{i=1}^n e_i =$
- The average of fitted values,  $(1/n) \sum_{i=1}^n \hat{Y}_i =$