## Correlation Coefficient

- $r = \dfrac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} = \dfrac{SSXY}{\sqrt{SSX \cdot SSY}}$

  - recall that $\hat{\beta}_1 = \dfrac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \dfrac{SSXY}{SSX}$

  - $\rightarrow \hat{\beta}_1 = r\,\dfrac{S_Y}{S_X}$

- Properties
  - $-1 \le r \le 1$
  - $r$ does not depend on units of measurement (dimensionless)
  - $r \approx 0 \rightarrow$ no *linear* relationship between $X$ and $Y$
  - larger $|r|$ means a stronger *linear* relationship
  - $r^2 = \dfrac{\sum (y_i - \bar{y})^2 - \sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} = \dfrac{SSY - SSE}{SSY}$
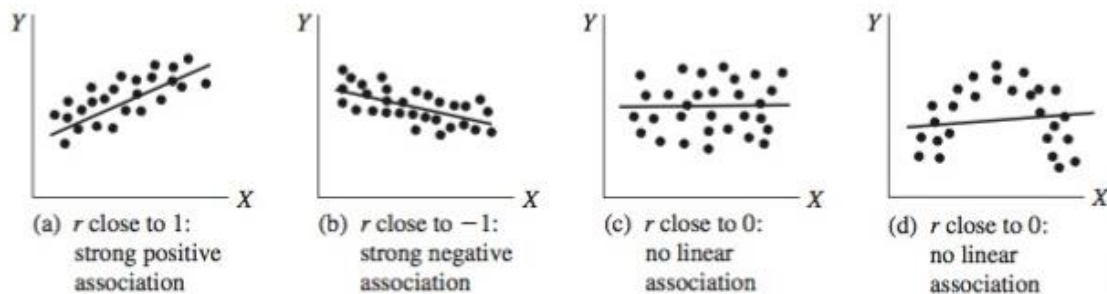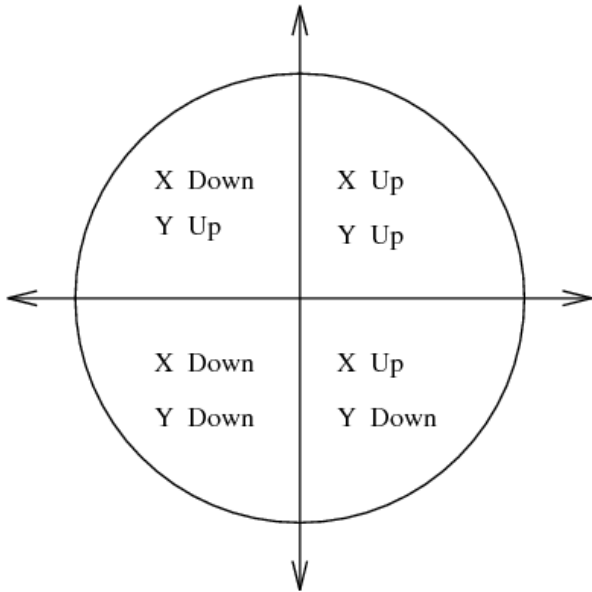


(a) *r* close to 1: strong positive association

(b) *r* close to −1: strong negative association

(c) *r* close to 0: no linear association

(d) *r* close to 0: no linear association

**FIGURE 6.1**   Correlation coefficient as a measure of association

# Transformations and Tukey's "Rule of the Bulge"

- Observe which way the curve bulges as suggested by a scatterplot of the data
- Transform **y** or **x** (or both) according to the signs of the corresponding quadrant:

  *up* → powers >1

  *down* → powers <1 (including logarithm)

|  | |
|---|---|
| X Down<br>Y Up | X Up<br>Y Up |
| X Down<br>Y Down | X Up<br>Y Down |

## Correlation Coefficient - Inference

$$r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} = \frac{SSXY}{\sqrt{SSX \bullet SSY}} \quad \rightarrow \quad \hat{\beta}_1 = r\frac{S_Y}{S_X}$$

- Testing Hypotheses about $\rho$, the *population* correlation coefficient
- $H_0$: $\rho = 0$
  - Sampling distribution of $r$ is symmetric & approx. normal
  - $t = \dfrac{r\sqrt{n-2}}{\sqrt{1-r^2}}$ with $n$-2 d.f. is equivalent to testing $H_0$: $\beta_1 = 0$ with $t = \dfrac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$

- $H_0$: $\rho = \rho_0$ where $\rho_0 \neq 0$
  - Sampling distribution of $r$ is *NOT* symmetric or approx. normal
  - *Fisher's "Z-transformation"* gives an approx. normally distributed statistic

  $$\frac{1}{2}\bullet\ln\left(\frac{1+r}{1-r}\right) \sim N\left(\frac{1}{2}\bullet\ln\left(\frac{1+\rho}{1-\rho}\right), \frac{1}{n-3}\right)$$

  - Testing $H_0$: $\rho = \rho_0$ $\quad Z = \dfrac{\frac{1}{2}\bullet\ln(1+r/1-r) - \frac{1}{2}\bullet\ln(1+\rho_0/1-\rho_0)}{1/\sqrt{n-3}} \sim N(0,1)$

  - 100(1-$\alpha$)% CI for $\dfrac{1}{2}\bullet\ln\left(\dfrac{1+\rho}{1-\rho}\right)$ : $\quad \dfrac{1}{2}\bullet\ln\left(\dfrac{1+r}{1-r}\right) \pm z_{\alpha/2}\bullet\dfrac{1}{\sqrt{n-3}} = (L_z, U_z)$

  - 100(1-$\alpha$)% CI for $\rho$: $(L_\rho, U_\rho) = $ (lower endpoint, upper endpoint)
    - Find $L_z$ & $U_z$ above and solve for $L_\rho$ & $U_\rho$ in $L_z = \dfrac{1}{2}\bullet\ln\left(\dfrac{1+L_\rho}{1-L_\rho}\right) \quad U_z = \dfrac{1}{2}\bullet\ln\left(\dfrac{1+U_\rho}{1-U_\rho}\right)$
    - $\rightarrow (L_\rho, U_\rho) = \left(\dfrac{e^{2L_z}-1}{e^{2L_z}+1}, \dfrac{e^{2U_z}-1}{e^{2U_z}+1}\right)$

- $H_0$: $\rho_1 = \rho_2$ (for 2 *independent* samples)
  - let $W_1 = \dfrac{1}{2}\bullet\ln\left(\dfrac{1+r_1}{1-r_1}\right)$ and $W_2 = \dfrac{1}{2}\bullet\ln\left(\dfrac{1+r_2}{1-r_2}\right)$
  - $Z = \dfrac{W_1 - W_2}{\sqrt{1/(n_1-3) + 1/(n_2-3)}} \sim N(0,1)$

## Bivariate Normal Distribution

$$f(x, y) = \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} e^{\frac{-1}{2(1-\rho^2)}\left[\left(\frac{x-\mu_x}{\sigma_x}\right)^2 - 2\rho\left(\frac{x-\mu_x}{\sigma_x}\right)\left(\frac{y-\mu_y}{\sigma_y}\right) + \left(\frac{y-\mu_y}{\sigma_y}\right)^2\right]}$$

$$E(Y \mid X = x) = \mu_y + \rho\frac{\sigma_y}{\sigma_x}(x - \mu_x)$$

$$\hat{\mu}_{Y\mid X=x} = \overline{y} + r\frac{s_y}{s_x}(x - \overline{x})$$