# 11.8 Correlation Model

- In the regression context, we estimate $E(Y)$ and predict $Y$ based on fixed values of $X$.

- In some cases, we view both $X$ and $Y$ as random variables and seek only to study the association between $X$ and $Y$.

- Recall from STAT 511 that the correlation coefficient $\rho$ measures the linear association between $X$ and $Y$.

- If $(X_1, Y_1), ..., (X_n, Y_n)$ are iid from a bivariate normal distribution, then the MLE of $\rho$ is the sample correlation coefficient:

Since the SLR slope estimate
$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$,

- If $(X, Y) \sim$ bivariate normal and if the linear regression model holds,

- Tests of $H_0: \rho = 0$ (which under bivariate normality correspond to testing whether $X$ and $Y$ are independent) are equivalent to our $t$-tests for $H_0: \beta_1 = 0$ in SLR.
- The test statistic is

which is algebraically equivalent to the SLR test statistic

- This test statistic has a _____ with _____ under $H_0$.
- Therefore the usual t-test RR can be used to test $H_0 : \rho = 0$ against whichever appropriate alternative.
- To test $H_0 : \rho = \rho_0$ for some nonzero $\rho_0$, we must use <u>Fisher's z-transformation.</u>
- For large samples,



So we use the test statistic

- The RR for an $\alpha$-level test is:

- A $100(1-\alpha)\%$ large-sample CI $[L, U]$ for

                             can be obtained via:

- The $100(1-\alpha)\%$ CI for $\rho$ can then be found by back-transforming the endpoints $L$ and $U$:

Example 1: The weights (X, in thousands of pounds) and gas mileages (Y, in mpg) for 32 randomly selected cars were collected. Assume bivariate normality for $(X, Y)$. Summary calculations yield $S_{xy} = -158.617$, $S_{xx} = 29.679$, $S_{yy} = 1126.047$. Estimate $\rho$ and test (at $\alpha = .05$) whether X and Y are independent.

Example 1(a): Test (at $\alpha = .05$) whether the correlation coefficient between X and Y is less than $-0.7$. Find a 95% CI for $\rho$.

- The cor.test function in R can perform these calculations easily.

# The Coefficient of Determination

- Note that $S_{yy} = \sum (y_i - \bar{y})^2$ is a measure of the sample variation in the y-values.

- Recall that $SSE = \sum (y_i - \hat{y}_i)^2$ measures the amount of variation in the y's <u>unexplained</u> by the linear regression model.

- In the context of SLR,

is called the <u>coefficient of determination</u> and measures the <u>proportion</u> of variation in the y-values <u>explained</u> by the linear model.

<u>Example 1</u>:  r =

So

<u>Note</u>:  In MLR, a <u>coefficient of</u>
<u>multiple</u> <u>determination</u> is defined
similarly:

and has a similar interpretation.