

LECTURE 16

HYPOTHESISTESTING FOR PEARSON'S CORRELATION

PSY2002

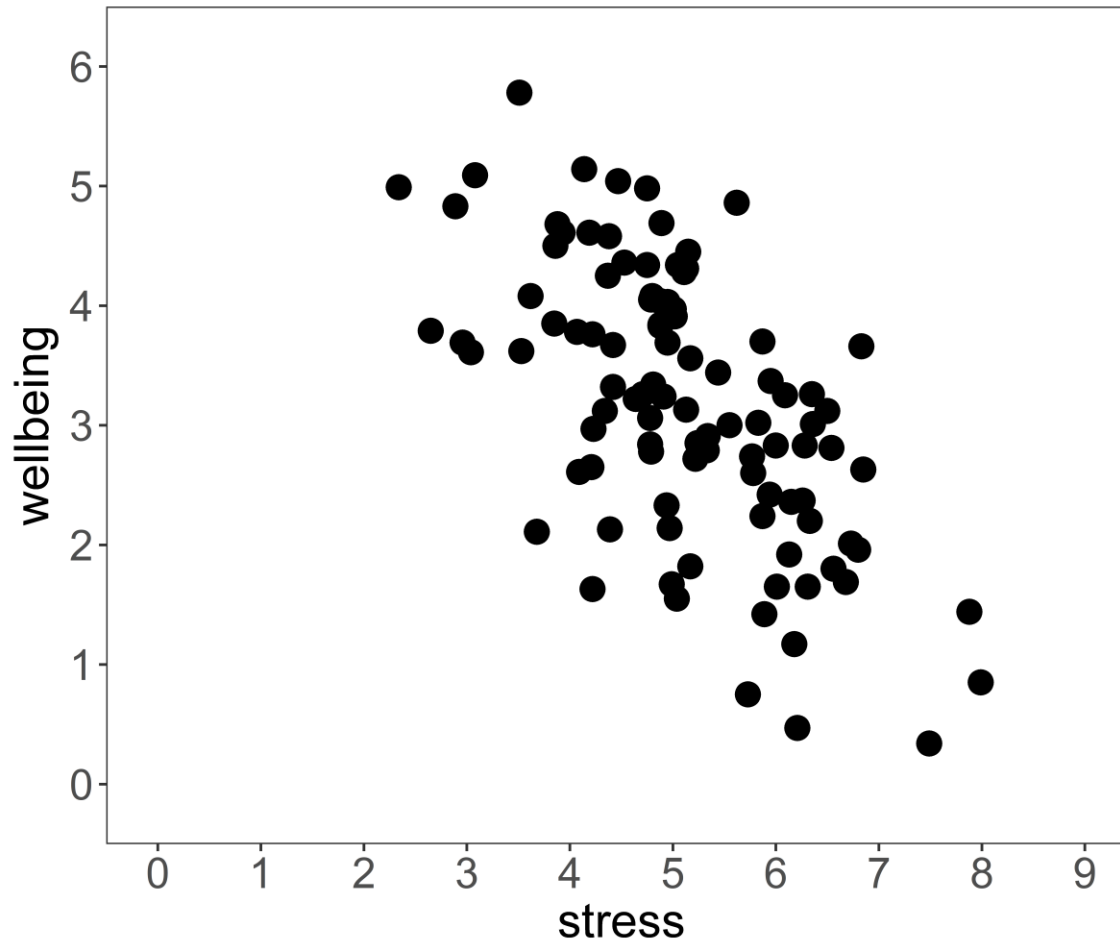
Hye Won Suk

A WORKING EXAMPLE

- A researcher wanted to examine the relationship between stress and wellbeing.
- She recruited 100 students at Sogang university and measured them on their stress level and wellbeing level.
- The wellbeing scores could take values between 0 and 6, and a higher score indicates a higher wellbeing level.
- The stress scores could take values between 0 and 9, and a higher score indicates a higher stress level.

SCATTER PLOT

- We can see a negative relationship between stress and wellbeing.



PEARSON'S CORRELATION COEFFICIENT

- The data set is given in the data file, *wellbeing_data.txt*.
- For the given data set, if you calculate the Pearson's correlation between the two variables (wellbeing and stress), you will obtain $r = -.620$.
- The following shows the result obtained using SAS.

Pearson Correlation Coefficients, N = 100 Prob > r under H0: Rho=0		
	stress	wellbeing
stress	1.00000	-0.61957 <.0001
wellbeing	-0.61957 <.0001	1.00000

INTERPRETATION OF CORRELATION

- How to interpret the obtained $r = -.620$?
 - Direction of the relationship
 - There is a negative relationship between stress and wellbeing.
 - Students with higher stress levels tend to show lower wellbeing levels.
 - Strength of the relationship
 - According to the Cohen's rule of thumb, we can say that there is a strong (linear) relationship between stress and wellbeing.

COEFFICIENT OF DETERMINATION

- Another way to interpret the correlation is to calculate the coefficient of determination and interpret its meaning.
- Coefficient of determination (결정 계수) is the squared correlation (r^2).
 - This quantity indicates the proportion of the variance in one variable that is accounted for by the other variable.
 - In the working example, $r^2 = (-.620)^2 = .384$.
 - This indicates that 38.4% of the variance of stress is accounted for by wellbeing.
 - At the same time, it also indicates that 38.4% of the variance of wellbeing is due to stress.

COEFFICIENT OF DETERMINATION

- The coefficient of determination (결정 계수), or r^2 is often used as an effect size measure.
- r^2 can vary between 0 and 1.
- A higher value of r^2 indicates a stronger relationship between the two variables.
 - $r^2 = 0$ indicates that 0% of the variance of one variable is associated with the other variable (no linear relationship).
 - $r^2 = 1$ indicates that 100% of the variance of one variable is accounted for by the other variable (perfect linear relationship).

HYPOTHESIS TESTING FOR A CORRELATION

- The sample Pearson's correlation coefficient describes the characteristic of the sample.
- However, it also serves as an estimator for the population correlation (ρ).
 - ρ indicates the population Pearson's correlation coefficient. It is read 'rho.'
- We can use a hypothesis testing to examine if the sample Pearson's correlation reflects an actual relationship in the population or appears just due to sampling.

HYPOTHESIS TESTING FOR A CORRELATION

- Again, follow the five steps of hypothesis testing.
 - Step 1: State the hypotheses
 - Step 2: Set the criteria for a decision
 - Step 3: Collect data and compute test statistics
 - A t -statistic will be calculated. (t -test)
 - Step 4: Make a decision
 - Step 5: State a conclusion

STEP I: STATE THE HYPOTHESES

- Null hypothesis (H_0)
 - H_0 : Stress and wellbeing are not linearly related.
 - H_0 : In the population, the Pearson's correlation between stress and wellbeing is 0.
 - $H_0: \rho = 0$

STEP I: STATE THE HYPOTHESES

- Alternative hypothesis (H_1)
 - H_1 : Stress and wellbeing are linearly related.
 - H_1 : In the population, the Pearson's correlation between stress and wellbeing is not 0.
 - $H_1: \rho \neq 0$

STEP 2: SET THE CRITERIA

- $\alpha = 0.05$
 - The alpha level (or level of significance) is a probability value that is used to define the concept of “very unlikely” in a hypothesis test.
 - By convention, we use $\alpha = 0.05$ unless otherwise specified. $\alpha = 0.05$ indicates that we will treat extreme 5% of the values as being unlikely to be observed under the null hypothesis.

STEP 3: COMPUTE TEST STATISTICS

- In step 3, we calculate a t -statistic as follows.

$$t = \frac{r - \rho}{\sqrt{\frac{1 - r^2}{n - 2}}}$$

- $\sqrt{\frac{1 - r^2}{n - 2}}$ indicates the standard error for r .
- n indicates the sample size.
- The t -statistic is known to follow a t distribution with $df = n - 2$ when the null hypothesis is true (and both variables are normally distributed).

STEP 3: COMPUTE TEST STATISTICS

- In the example, we can obtain the t -statistic and df as follows:

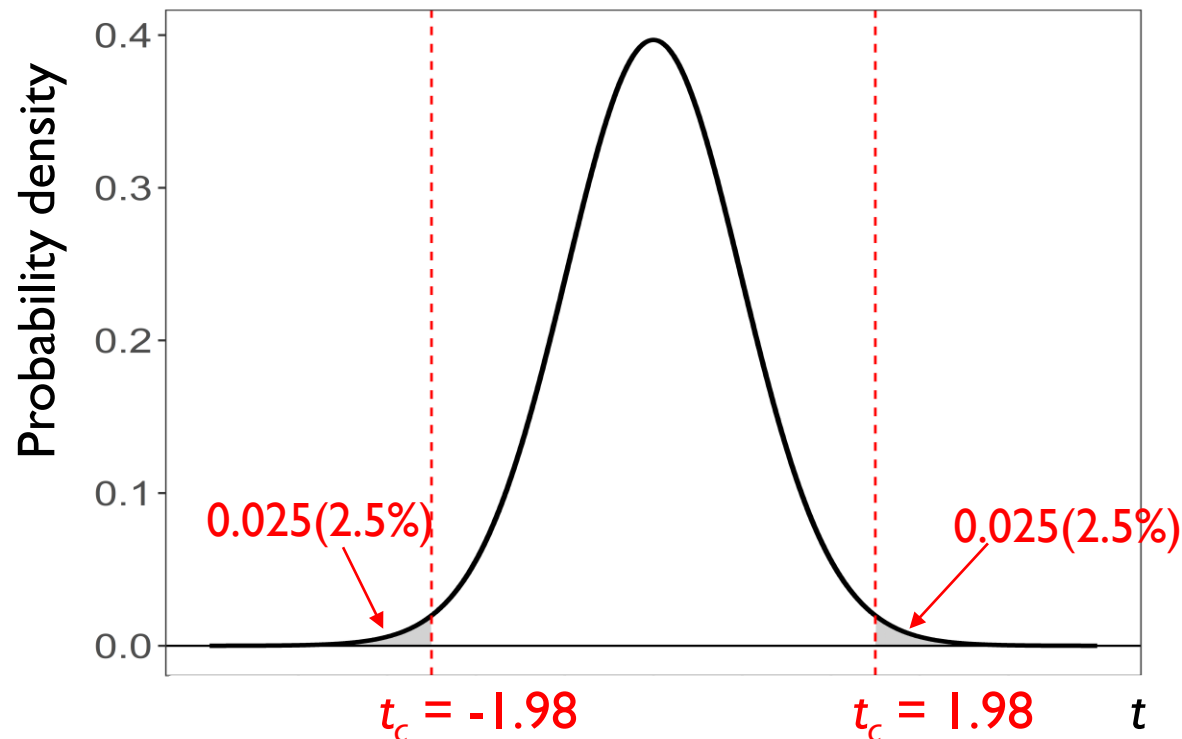
$$t = \frac{r - 0}{\sqrt{\frac{1 - r^2}{n - 2}}} = \frac{(-.620)}{\sqrt{\frac{1 - (-.620)^2}{100 - 2}}} = -7.82$$

$$df = n - 2 = 100 - 2 = 98$$

- The obtained t -statistic and df are reported as follows:
 - $t(98) = -7.82$

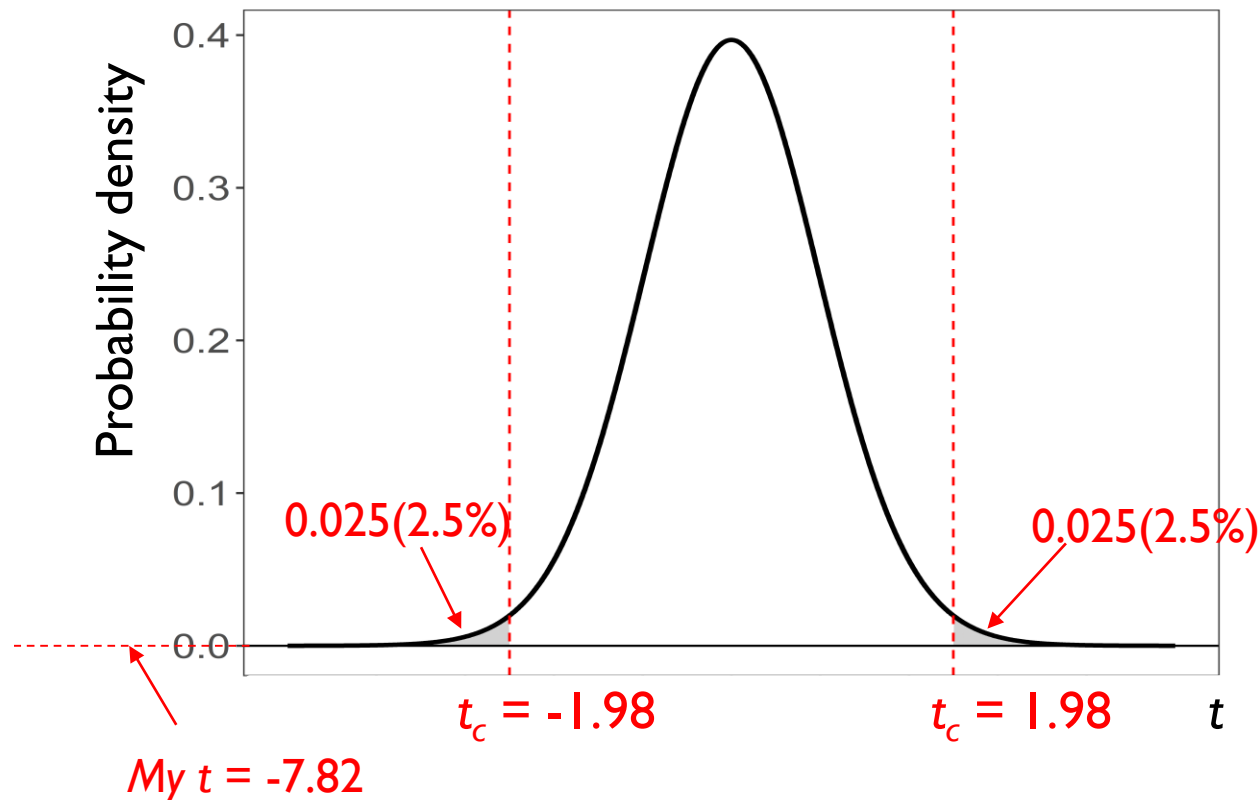
STEP 4: MAKE A DECISION

- Under the t -distribution with $df = 98$, the critical value for $\alpha = .05$ is 1.98. (<http://www.ttable.org/>)



STEP 4: MAKE A DECISION

- $|\text{My } t\text{-value}| > t_c; 7.82 > 1.98$
- My t -value is in the extreme zone (or in the critical region).
My t -value is a strong evidence against H_0 . \rightarrow Reject H_0 .



STEP 4: MAKE A DECISION

- Using SAS will provide the p-value.

Pearson Correlation Coefficients, N = 100 Prob > r under H0: Rho=0		
	stress	wellbeing
stress	1.00000	-0.61957 <.0001
wellbeing	-0.61957 <.0001	1.00000

p-value

- In this case, the p-value is very small ($<.0001$), and the exact p-value is not provided.
- However, we can see that the p-value is smaller than α , that is, $p < .05$, and thus we reject the null hypothesis.

STEP 5: STATE A CONCLUSION

- Stress and wellbeing showed a significant strong negative correlation ($r = -.620$, $p < .0001$, $r^2 = .384$). This indicates that students having higher stress levels tend to show lower wellbeing levels (or, students having higher wellbeing levels tend to show lower stress levels).

SAS OUTPUT

- The following tables provide the descriptive statistics for each variable, the estimated Pearson's correlation coefficient, and the p-value.

The CORR Procedure

2 Variables: stress wellbeing

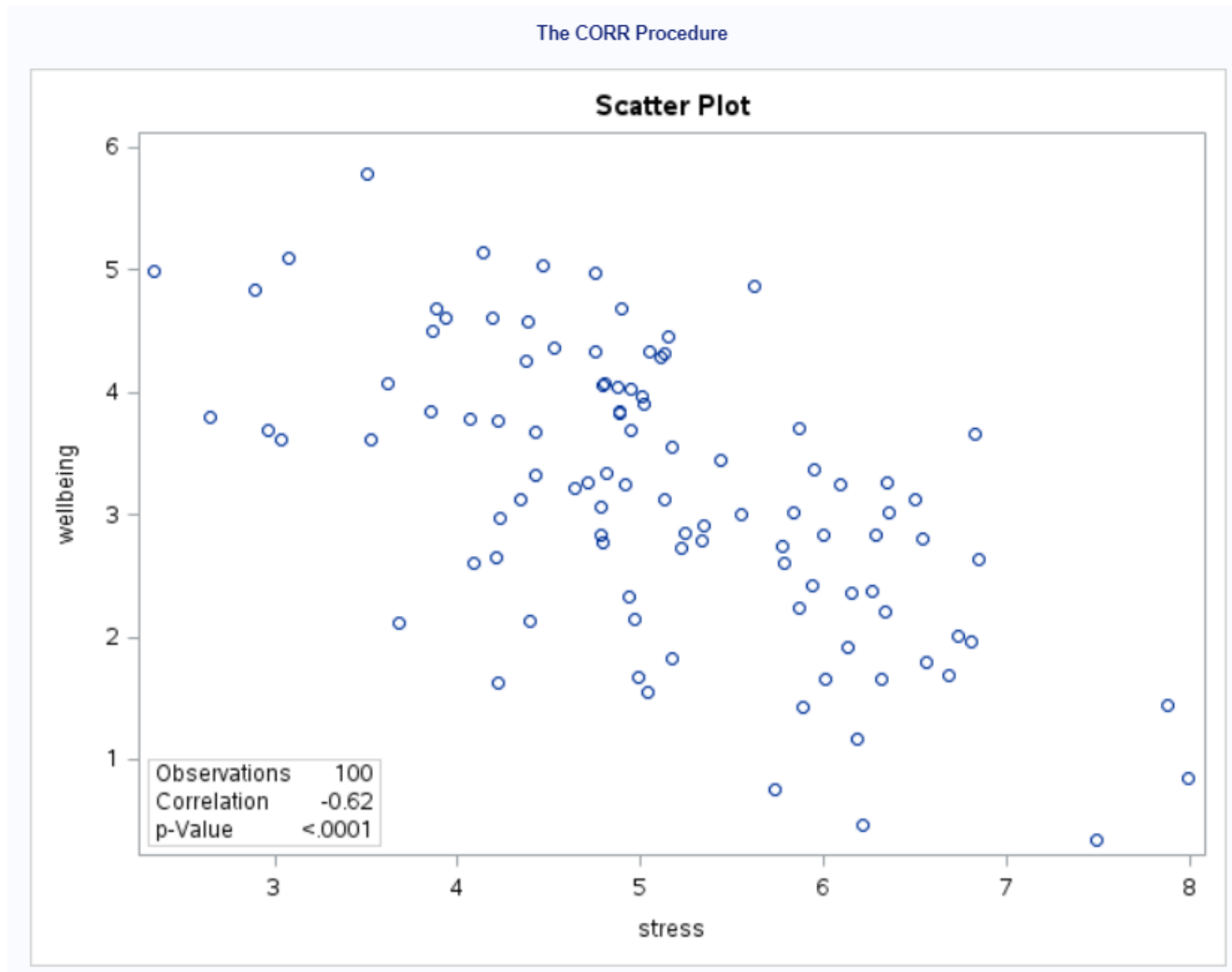
Simple Statistics

Variable	N	Mean	Std Dev	Sum	Minimum	Maximum
stress	100	5.11200	1.11336	511.20000	2.34000	7.99000
wellbeing	100	3.17730	1.14898	317.73000	0.34000	5.78000

Pearson Correlation Coefficients, N = 100 Prob > |r| under H0: Rho=0

	stress	wellbeing
stress	1.00000	-0.61957 <.0001
wellbeing	-0.61957 <.0001	1.00000

SAS OUTPUT



SUMMARY

- We can test the significance of the Pearson's correlation coefficient, in which a t -statistic is calculated.

$$t = \frac{r - \rho}{\sqrt{\frac{1 - r^2}{n - 2}}}$$

- The t -statistic is known to follow a t -distribution with $df = n - 2$ when the null hypothesis true (and both variables are normally distributed).
- We can use the typical 5 steps of hypothesis testing as usual.
- When interpreting the Pearson's correlation coefficient, we can use the coefficient of determination (r^2). It also serves as an effect size measure.