1. Introduction

Welcome to the series of E-learning module on Fisher's Z transformation and its uses. In this module we are going cover the basic concept of correlation, Fisher's Z transformation, uses and applications of Z transformation, and role of Fisher's Z transformation in testing of hypothesis.

By the end of this session, you will be able to:

- Describe correlation
- Explain the role and importance of Fisher's Z transformation
- Explain the application of Z-transformation
- Explain Significant tests of correlation using Z transformation

In statistics, hypotheses about the value of the population correlation coefficient Rho between variables X and Y can be tested using the Fisher's transformation applied to the sample correlation coefficient, represented by small letter 'R'

What is correlation?

Correlation is one of the most common and most useful statistics. A correlation is a single number that describes the degree of relationship between two variables.

Correlation Example

Let's assume that we want to look at the relationship between two variables, height (in inches) and self esteem. Perhaps we have a hypothesis that how tall you are, effects your self esteem (incidentally, we don't think we have to worry about the direction of causality here-it's not likely that self esteem causes your height!).

Let's say we collect some information on (say) twenty individuals (all male, we know that the average height differs for males and females so, to keep this example simple we'll just use males). Height is measured in inches. Self esteem is measured based on the average one to five rating items (where higher scores mean higher self esteem) and then we can study the associated relationship between heights of males and their self esteem by computing a correlation coefficient for the collected data.

2. Properties and Types of Correlation and Testing the Significance of a Correlation

Properties of correlation coefficient

• r only measures the strength of a linear relationship. There are other kinds of relationships besides linear

• r is always between minus one and one inclusive. Minus one means perfect negative linear correlation and plus one means perfect positive linear correlation

- r has the same sign as the slope of the regression (best fit) line
- r does not change if the independent (x) and dependent (y) variables are interchanged
- r does not change if the scale on either variable is changed
- r has a Student's t distribution

Causation

If there is a significant linear correlation between two variables, then one of five situations can be true.

- There is a direct cause and effect relationship
- There is a reverse cause and effect relationship
- The relationship may be caused by a third variable
- The relationship may be caused by complex interactions of several variables
- The relationship may be coincidental

Types of correlation:

A positive correlation indicates that if one variable increases or decreases, the other variable also increases or decreases. That is, the two variables move in the same direction.

A negative correlation coefficient indicates that if one variable increases, the other decreases and vice versa. That is, the two variables move in the opposite direction.

Zero correlation indicates no relationship between the two variables.

One or minus one indicates a linear relationship, such that if one variable is known, the second can be accurately predicted.

Testing the Significance of a Correlation

Once we have computed a correlation, one can determine the probability that the observed correlation occurred by chance. That is, we can conduct a significance test. Most often we are interested in determining the probability that the correlation is a real one and not a chance occurrence. In this case, we are testing the mutually exclusive hypothesis H naught: 'r' is equal to zero against the alternative hypothesis H one: r is not equal to zero or 'r' is less than zero or 'r' is greater than zero.

The easiest way to test this hypothesis is to find a statistics book that has a table of critical values of 'r'. Most introductory statistics texts would have a table like this.

As in all hypotheses testing, you need to first determine the significance level. The generally used most common significance level of alpha is equal to zero point zero five.

This means that we are conducting a test where the odds that the correlation is a chance occurrence is no more than five out of one hundred .

Before we look up the critical value in a table we have to compute the degrees of freedom or 'df'.

The df is simply equal to 'n' minus two (in case of 't' test). Finally, one has to decide whether you are doing a one tailed or two tailed test.

In the examples, where we have no strong prior theory to suggest whether the relationship between two variables (such as height and self esteem as in the above mentioned example) would be positive or negative, we will opt for the two-tailed test.

With these three pieces of information : the significance level (alpha is equal to zero point zero five, degrees of freedom and type of test, we can now test the significance of the correlation which we found.

When we look up this value in the handy little table at the back of any statistics book suppose we find that the critical value as zero point four four three eight, it this means that if our correlation is greater than zero point four four three eight or less than minus zero point four four three eight (remember, this is a two-tailed test) we can conclude that the odds are less than five out of one hundred that this is a chance occurrence.

Suppose our correlation is zero point seven three or higher, we conclude that it is not a chance finding and that the correlation is "statistically significant" (given the parameters of the test). We can reject the null hypothesis and accept the alternative.

3. Fisher's Transformation and Its Applications

The Fisher transformation is an approximate variance stabilizing transformation for 'r' when X and Y follow a bivariate normal distribution. This means that the variance of 'Z' is approximately constant for all values of the population correlation coefficient Rho.

Without the Fisher transformation, the variance of 'R' grows smaller as modulus of Rho gets closer to one.

Since the Fisher transformation is approximately the identity function when modulus of 'R' is less than one by two, it is sometimes useful to remember that the variance of 'R' is well approximated by one by N as long as modulus of Rho is not too large and N is not too small. This is related to the fact that the asymptotic variance of R is one for bivariate normal data.

The behaviour of this transform has been extensively studied since Fisher introduced it in Nineteen fifteen. Fisher himself found the exact distribution of Z' for data from a bivariate normal distribution in Nineteen twenty one.

In Nineteen fifty one Gayen determined the exact distribution of 'Z' for data from a bivariate Type.

<u>Hotelling</u>, in Nineteen fifty three calculated the Taylor series expressions for the moments of Z and several related statistics and in Nineteen eighty nine Hawkins discovered the asymptotic distribution of Z for virtually any data.

Other Uses

While the Fisher transformation is mainly associated with the Pearson product-moment correlation coefficient for bivariate normal observations, it can also be applied to Spearman's rank correlation coefficient in more general cases.

A similar result for the <u>asymptotic distribution</u> applies, but with a minor adjustment factor.

Applications of Fisher's *z* Transformation

Fisher describes the following practical applications of the transformation:

- Testing whether a population correlation is equal to a given value
- Testing for equality of two population correlations
- Combining correlation estimates from different samples

Correlation between two populations

To test if a population correlation Rho one from a sample of 'n' one observations with sample correlation 'R' one is equal to a given Rho naught, first apply the Z transformation to 'R' and Rho naught:

Z one is equal to one by two into Ln of (one plus 'R' one by one minus 'R' one) which is equal to tan H inverse of 'R' one and zhi naught is equal to one by two into Ln of (one plus Rho naught by one minus Rho one) which is equal to tan H inverse of Rho naught where Ln is the natural logarithm.

The p-value is then computed by treating Z one minus zhi one is equal to Rho naught by 2 into ('n' one minus one) as a normal random variable with mean zero and variance one by (n one minus three)

Assume that sample correlations 'R' one and 'R' two are computed from two independent samples of 'n' one and 'n' two observations, respectively.

To test whether the two corresponding population correlations, Rho one and Rho to, are equal, first apply the Z transformation to the two sample correlations:

Z one is equal to one by two into Ln of (one plus 'R' one by one minus 'R' one) and Z two is equal to one by two into Ln of (one plus r two by one minus 'R' two)

The P value is derived under the null hypothesis of equal correlation. That is, the difference Z one minus Z two is distributed as a normal random variable with mean zero and variance one by ('n' one minus three) plus ('n' two minus three)

Assuming further that the two samples are from populations with identical correlation, a combined correlation estimate can be computed. The weighted average of the corresponding Z values is

Z bar is equal to ('n' one minus three) into Z one plus ('n' two minus three) into Z two divided by ('n' one plus 'n' two minus six), where the weights are inversely proportional to their variances.

4. Fisher's Transformation Contd

Note that this approach can be extended to include more than two samples.

The best known technique for transforming correlation coefficient (r) values into weighted additive quantities is the 'R' to 'Z' transformation by Fisher.

Fisher's 'R' to 'Z' transformation is an elementary transcendental function called the inverse hyperbolic tangent function.

The reverse, a Z to 'R' transformation, is therefore a hyperbolic tangent function.

These transformations are needed to compute a weighted mean correlation coefficient and for hypothesis testing. Note that averaged correlation coefficients are not computable directly from raw 'R' values.

Indeed, it is not possible to add, subtract, average or take standard deviations out of raw r values. The sampling distribution of <u>Pearson's r</u> is not <u>normally distributed</u>.

Fisher developed a transformation now called "Fisher's Z transformation" that converts Pearson's R's to the normally distributed variable Z. It is not important to understand how Fisher came up with the above mentioned formula.

What is important are two attributes of the distribution of the z' statistic:

• It is normal and

• It has a known standard error of: sigma Z is equal to one by square root of 'n' minus three

Fisher's Z is used for computing confidence intervals on Pearson's correlation and for confidence intervals on the difference between correlations. You can use the R to Z table, to convert from R to Z and back.

Using the Fisher R to Z transformation, the calculated value of Z can be applied to assess the significance of the difference between two correlation coefficients, R one and R two, found in two independent samples. If R one is greater than R two, the resulting value of Z will have a positive sign.

If R one is smaller than R two, the sign of Z will be negative.

We have already discussed about Pearson correlation as a measure of the STRENGTH of a relationship between two variables.

But any relationship should be assessed for its SIGNIFICANCE as well as its strength.

A general discussion of significance tests for relationships between two continuous variables.

• Identify the factors in relationships between two variables

The strength of the relationship is indicated by the correlation coefficient: 'R'

However, it is actually measured by the coefficient of determination: 'R' square.

The significance of the relationship is expressed in probability levels 'P' (e.g., significant at p is equal to zero point zero five)

This tells how unlikely a given correlation coefficient, 'R' will occur given no relationship in the

population.

NOTE: The smaller the P level, the more significant the relationship. BUT! The larger the correlation, the stronger the relationship.

Consider the classical model for testing significance. It assumes that you have a sample of cases from a population.

The question is whether your observed statistic for the sample is *likely* to be observed given some assumption of the corresponding population parameter.

If your observed statistic does not exactly match the population parameter, perhaps the difference is due to sampling error.

The fundamental question: Is the difference between what you observe and what you expect given the assumption of the population large enough to be significant to reject the assumption?

The greater the difference the more the sample statistic deviates from the population parameter. That is, the less likely (small probability values) that the population assumption is true.

The classical model makes some assumptions about the population parameter:

R is equal to correlation between two variables in the sample.

Rho is equal to correlation between the two same variables in the population

A common assumption is that there is NO relationship between X and Y in the population: Rho is equal to zero

Under this common null hypothesis in correlational analysis R is equal to zero.

Testing for the significance of the correlation coefficient, R

When the test is against the null hypothesis: R XY is equal to zero.

We will be interested to know, What is the likelihood of drawing a sample with RXY is equal to zero?

The sampling distribution of R is approximately normal (but bounded at minus one and plus one) when N is large and distributes T when N is small.

Points to be noted:

- Note that a relationship can be strong and yet not significant
- Conversely, a relationship can be weak but significant
- The key factor is the size of the sample

• For small samples, it is easy to produce a strong correlation by chance and one must pay attention to significance, to keep from jumping to conclusions: that is, rejecting a true null hypothesis, which means making a type one error

• For large samples, it is easy to achieve significance, and one must pay attention to the strength of the correlation to determine if the relationship explains very much

Testing the significance of R when R is NOT assumed to be zero: The test requires first transforming the sample R to a new value, Z. Which is seldom used.

5. Commonly Encountered Correlation Tests

Most Commonly Encountered Correlation Tests

We often speak to researchers wanting to compare the significance of two correlations. The two scenarios most commonly encountered are:

1) comparing dependent correlations; and

2) comparing independent correlations

1. Two dependent correlations

This is the scenario when you have three variables, X, Y, and Z, and you want to compare the XY correlation with the XZ correlation.

It comes up often when you want to know which of two variables are more related to a third variable. In this sense it is often related to approaches that attempt to assess variable importance using multiple regression.

2. Two independent correlations.

This is the scenario when two correlations are obtained from different samples and you want to test whether they are significantly different. An example is where a researcher wants to know whether intelligence test scores and performance are correlated the same in different social groups.

In most cases, this is very similar to testing for a group by IV (independent variable) interaction effect.

Thus, moderator regression is often a more appropriate means of testing the relationship. However, you can also run a specific test of statistical significance on the difference between the two correlations.

Other scenarios include:

3) Testing whether a correlation is significantly different from some target value, usually zero, but possibly another value

4) A significance test on a correlation matrix

5) Structural Equation Modelling software can also be used to test more general hypotheses about patterns in correlation matrices.

Hence to test the significance of an observed sample correlation coefficient from an uncorrelated Bivariate Normal population, T test is used.

However in random sample of size 'n' from a Bivariate Normal population in which a population correlation coefficient is not equal to zero, Prof. R.A Fisher proved that the distribution of a sample correlation coefficient is by no means Normal and in the neighbourhood of Rho is equal to plus or minus one, its probability curve is extremely skewed even for large 'n'.

Hence if Rho is not equal to zero, Fisher suggested the transformation Z is equal to one by

two into Ln of (one plus R by one minus R) which is equal to tan H inverse of R and proved that even for small samples the distribution of Z is approximately Normal with mean Zhi is equal to one by two into Ln of (one plus Rho by one minus Rho) which is equal to tan H inverse of Rho and variance one by (N minus three) and for large N say (greater than thirty), the approximation is fairly good.

Here's a summary of our learning in this session where we have understood:

- Basic concept of correlation
- Concept and Role of Fisher's Z transformation
- Uses and applications of Fisher's Z transformation
- Tests of significance using Z transformation