

## Chapter 867

# Confidence Intervals for the Interaction Odds Ratio in Logistic Regression with Two Binary X's

## Introduction

Logistic regression expresses the relationship between a binary response variable and one or more independent variables called *covariates*. This procedure calculates sample size for the case when there are two binary covariates ( $X$  and  $Z$ ) and their interaction ( $XZ$ ) in the logistic regression model and a Wald statistic is used to calculate a confidence interval for the interaction odds ratio ( $OR_{int}$ ). Often,  $Y$  is called the *response* variable, the first binary covariate,  $X$ , is referred to as the *exposure* variable and the second binary covariate,  $Z$ , is referred to as the *confounder* variable. For example,  $Y$  might refer to the presence or absence of cancer and  $X$  might indicate whether the subject smoked or not, and  $Z$  is the presence or absence of a certain gene.

## Sample Size Calculations

Using the *logistic model*, the probability of a binary event is

$$X \Pr(Y = 1|X, Z) = \frac{\exp(\beta_0 + \beta_1 X + \beta_2 Z + \beta_3 XZ)}{1 + \exp(\beta_0 + \beta_1 X + \beta_2 Z + \beta_3 XZ)}$$

This formula can be rearranged so that it is linear in  $X$  as follows

$$\log\left(\frac{\Pr(Y = 1|X, Z)}{1 - \Pr(Y = 1|X, Z)}\right) = \beta_0 + \beta_1 X + \beta_2 Z + \beta_3 XZ$$

In the logistic regression model, the magnitude of the association of  $XZ$  and  $Y$  is represented by the slope  $\beta_3$ . Since  $XZ$  is binary, only two cases need be considered:  $XZ = 0$  and  $XZ = 1$ .

## Confidence Intervals for the Interaction Odds Ratio in Logistic Regression with Two Binary X's

The logistic regression model defines the baseline probability as

$$P_0 = \Pr(Y = 1|X = 0, Z = 0) = \frac{\exp(\beta_0)}{1 + \exp(\beta_0)}$$

The odds ratio between Y and XZ is defined as

$$OR_{int} = \exp(\beta_3)$$

It well known that the distribution of the maximum likelihood estimate of  $\beta_3$  is asymptotically normal. A confidence interval for this slope is commonly formed from the Wald statistic

$$z = \frac{\hat{\beta}_3}{s_{\hat{\beta}_3}}$$

A  $(1 - \alpha)\%$  two-sided confidence interval for  $\beta_3$  is

$$\hat{\beta}_3 \pm z_{1-\frac{\alpha}{2}} s_{\hat{\beta}_3}$$

By transforming this interval into the odds ratio scale by exponentiating both limits, a  $(1 - \alpha)\%$  two-sided confidence interval for  $OR_{int}$  is

$$(OR_{LL}, OR_{UL}) = \exp\left(\hat{\beta}_3 \pm z_{1-\frac{\alpha}{2}} s_{\hat{\beta}_3}\right)$$

Note that this interval is not symmetric about  $OR_{int}$ .

Often, the goal during this part of the planning process is to find the sample size that reduces the width of the interval to a certain value  $D = OR_{UL} - OR_{LL}$ . A suitable  $D$  is found using a simple search of possible values of  $N$ .

Usually, the value of  $s_{\hat{\beta}_3}$  is not known before the study so this quantity must be estimated. Demidenko (2008) gives a method for calculating an estimate of this variance from various quantities that can be set at the planning stage. Let  $p_x$  be the probability that  $X = 1$  in the sample. Similarly, let  $p_z$  be the probability that  $Z = 1$  in the sample.

Define the relationship between X and Z as a logistic regression as follows

$$\Pr(X = 1|Z) = \frac{\exp(\gamma_0 + \gamma_1 Z)}{1 + \exp(\gamma_0 + \gamma_1 Z)}$$

The value of  $\gamma_0$  is found from

$$\exp(\gamma_0) = \frac{Q + \sqrt{Q^2 + 4p_x(1 - p_x)\exp(\gamma_1)}}{2(1 - p_x)\exp(\gamma_1)}$$

$$Q = p_x(1 + \exp(\gamma_1)) + p_z(1 - \exp(\gamma_1)) - 1$$

The information matrix for this model is

$$I = \begin{bmatrix} L + F + J + R & F + R & J + R & R \\ F + R & F + R & R & R \\ J + R & R & J + R & R \\ R & R & R & R \end{bmatrix}$$

where

$$L = \frac{(1 - p_z)\exp(\beta_0)}{(1 + \exp(\gamma_0))(1 + \exp(\beta_0))^2}$$

$$R = \frac{p_z\exp(\beta_0 + \beta_1 + \beta_2 + \beta_3 + \gamma_0 + \gamma_1)}{(1 + \exp(\gamma_0 + \gamma_1))(1 + \exp(\beta_0 + \beta_1 + \beta_2 + \beta_3))^2}$$

## Confidence Intervals for the Interaction Odds Ratio in Logistic Regression with Two Binary X's

$$F = \frac{(1 - p_z)\exp(\beta_0 + \beta_1 + \gamma_0)}{(1 + \exp(\gamma_0))(1 + \exp(\beta_0 + \beta_1))^2}$$

$$J = \frac{p_z \exp(\beta_0 + \beta_2)}{(1 + \exp(\gamma_0 + \gamma_1))(1 + \exp(\beta_0 + \beta_2))^2}$$

The value of  $\sqrt{N}S_{\hat{\beta}_3}$  is the (4, 4) element of the inverse of  $I$ .

The values of the regression coefficients are input as  $P_0$  and the following odds ratio as follows

$$OR_{int} = \exp(\beta_3)$$

$$OR_{yx} = \exp(\beta_1)$$

$$OR_{yz} = \exp(\beta_2)$$

$$OR_{xz} = \exp(\gamma_1)$$

The value of  $\beta_0$  is calculated from  $P_0$  using

$$\beta_0 = \log\left(\frac{P_0}{1 - P_0}\right)$$

Thus, the confidence interval can be specified in terms of several odds ratios and  $P_0$ . Of course, these results are only approximate. The width of the final confidence interval depends on the actual data values.

## Procedure Options

This section describes the options that are specific to this procedure. These are located on the Design tab. For more information about the options of other tabs, go to the Procedure Window chapter.

## Design Tab

The Design tab contains most of the parameters and options that you will be concerned with.

### Solve For

#### Solve For

This option specifies the parameter to be solved for from the other parameters. The parameters that may be selected are *Precision (Confidence Interval Width)*, *Confidence Level*, or *Sample Size*.

### One-Sided or Two-Sided Interval

#### Interval Type

Specify whether the confidence interval will be two-sided, one-sided with an upper limit, or one-sided with a lower limit.

## Confidence Intervals for the Interaction Odds Ratio in Logistic Regression with Two Binary X's

---

### Confidence

#### Confidence Level (1 – Alpha)

This option specifies one or more values of the proportion of confidence intervals (constructed with this same confidence level, sample size, etc.) that would have the same width.

The range of possible values is between 0 and 1. However, the range is usually between 0.5 and 1. Common choices are 0.9, 0.95, and 0.99. You should select a value that expresses the needs of this study.

You can enter a single value such as *0.7* or a series of values such as *0.7 0.8 0.9* or *0.7 to 0.95 by 0.05*.

---

### Sample Size

#### N (Sample Size)

This option specifies the total number of observations in the sample. You may enter a single value or a list of values.

---

### Precision

#### Distance from ORint to Limit

In a one-sided confidence interval (sometimes called a confidence bound), this is the distance between the upper or lower confidence limit of ORint and the value of ORint. As the sample size increases, this value decreases and thus the interval becomes more precise.

Since an odds ratio is typically between 0.2 and 10, it is reasonable that the value of this distance is also between 0.2 and 10. By definition, only positive values are possible.

You can enter a single value such as *1* or a series of values such as *0.5 1 1.5* or *0.5 to 1.5 by 0.2*.

#### Width of ORint Confidence Interval

In a two-sided confidence interval, this is the difference between the upper and lower confidence limits of ORint. As the sample size increases, this width decreases and thus the interval becomes more precise.

Since an odds ratio is typically between 0.2 and 10, it is reasonable that the value of this width is also between 0.2 and 10. By definition, only positive values are possible.

You can enter a single value such as *1* or a series of values such as *0.5 1 1.5* or *0.5 to 1.5 by 0.2*.

---

### Baseline Probability

#### P0 [Pr(Y = 1 | X = 0, Z = 0)]

This gives the value of the baseline probability of a response,  $P_0$ , when neither the exposure nor confounder are present.

$P_0$  is a probability, so it must be between zero and one.

---

### Odds Ratios

#### ORint (X,Z Interaction Odds Ratio)

Specify one or more values of the XZ-interaction Odds Ratio. This is the value that you expect to calculate from the data.

You can enter a single value such as *1.5* or a series of values such as *1.5 2 2.5* or *0.5 to 0.9 by 0.1*.

The range of this parameter is  $0 < \text{ORint} < \infty$  (typically,  $0.1 < \text{ORint} < 10$ ).

**Confidence Intervals for the Interaction Odds Ratio in Logistic Regression with Two Binary X's****OR<sub>yx</sub> (Y,X Odds Ratio)**

Specify one or more values of the Odds Ratio of Y and X, a measure of the effect size (event rate) that is to be detected by the study. This is the ratio of the odds of the outcome Y given that the exposure X = 1 to the odds of Y = 1 given X = 0.

You can enter a single value such as *1.5* or a series of values such as *1.5 2 2.5* or *0.5 to 0.9 by 0.1*.

The range of this parameter is  $0 < \text{OR}_{yx} < \infty$  (typically,  $0.1 < \text{OR}_{yx} < 10$ ).

**OR<sub>yz</sub> (Y,Z Odds Ratio)**

Specify one or more values of the Odds Ratio of Y and Z, a measure of the relationship between Y and Z. This is the ratio of the odds of the outcome Y given that the exposure Z = 1 to the odds of Y = 1 given Z = 0.

You can enter a single value such as *1.5* or a series of values such as *1.5 2 2.5* or *0.5 to 0.9 by 0.1*.

The range of this parameter is  $0 < \text{OR}_{yz} < \infty$  (typically,  $0.1 < \text{OR}_{yz} < 10$ ).

**OR<sub>xz</sub> (X,Z Odds Ratio)**

Specify one or more values of the Odds Ratio of X and Z, a measure of the relationship between X and Z. This is the ratio of the odds of the exposure X given that the confounder Z = 1 to the odds that X = 1 given Z = 0.

You can enter a single value such as *1.5* or a series of values such as *1.5 2 2.5* or *0.5 to 0.9 by 0.1*.

The range of this parameter is  $0 < \text{OR}_{xz} < \infty$  (typically,  $0.1 < \text{OR}_{xz} < 10$ ).

---

**Prevalences****Percent with X = 1**

This is the percentage of the sample in which X = 1. It is often called the prevalence of X.

You can enter a single value or a range of values. The permissible range is 1 to 99.

**Percent with Z = 1**

This is the percentage of the sample in which Z = 1. It is often called the prevalence of Z.

You can enter a single value or a range of values. The permissible range is 1 to 99.

## Example 1 – Find Sample size

A study is to be undertaken to study the association between the occurrence of a certain type of cancer (response variable) and the presence of a certain food in the diet. A second variable, the presence or absence of a certain gene, is also thought to impact the result. The researchers want a sample size large enough to guarantee that the width of a confidence interval about ORint is less than 0.90.

The baseline cancer event rate is 5%. They estimate ORint will be 0.5 and they want a confidence level of 0.95. They want to look at the sensitivity of the analysis to the specification of the other odds ratios, so they want to obtain the results ORyz = 1, 1.5, 2 and ORxz = 1, 1.5, 2. They estimate ORyx at 1.5. The researchers estimate that about 40% of the sample eat the food being studied and about 25% will have the gene of interest.

### Setup

This section presents the values of each of the parameters needed to run this example. First, from the PASS Home window, load the **Confidence Intervals for the Interaction Odds Ratio in Logistic Regression with Two Binary X's** procedure. You may then make the appropriate entries as listed below, or open **Example 1** by going to the **File** menu and choosing **Open Example Template**.

<u>Option</u>	<u>Value</u>
<b>Design Tab</b>	
Solve For .....	<b>Sample Size</b>
Interval Type .....	<b>Two-Sided</b>
Confidence Level .....	<b>0.95</b>
Width of ORyx Confidence Interval .....	<b>0.9</b>
P0 [Pr(Y=1 X=0, Z=0)] .....	<b>0.05</b>
ORint (X, Z Interaction Odds Ratio).....	<b>0.5</b>
ORyx (Y, X Odds Ratio) .....	<b>1.5</b>
ORyz (Y, Z Odds Ratio).....	<b>1 1.5 2</b>
ORxz (X, Z Odds Ratio).....	<b>1 1.5 2</b>
Percent with X = 1 .....	<b>40</b>
Percent with Z = 1 .....	<b>25</b>

### Annotated Output

Click the Calculate button to perform the calculations and generate the following output.

#### Numeric Results

Numeric Results for Two-Sided Confidence Interval of ORint											
Conf Level	N	C.I. Width	ORint	Lower	Upper	ORyx	ORyz	ORxz	P0	Pct X=1	Pct Z=1
				C.L.	C.L.						
0.950	2995	0.8999	0.500	0.223	1.123	1.500	1.000	1.000	0.050	40.0	25.0
0.950	2868	0.8999	0.500	0.223	1.123	1.500	1.000	1.500	0.050	40.0	25.0
0.950	2845	0.8999	0.500	0.223	1.123	1.500	1.000	2.000	0.050	40.0	25.0
0.950	2253	0.9000	0.500	0.223	1.123	1.500	1.500	1.000	0.050	40.0	25.0
0.950	2169	0.8999	0.500	0.223	1.123	1.500	1.500	1.500	0.050	40.0	25.0
0.950	2156	0.9000	0.500	0.223	1.123	1.500	1.500	2.000	0.050	40.0	25.0
0.950	1884	0.8998	0.500	0.223	1.123	1.500	2.000	1.000	0.050	40.0	25.0
0.950	1821	0.8998	0.500	0.223	1.122	1.500	2.000	1.500	0.050	40.0	25.0
0.950	1813	0.8999	0.500	0.223	1.123	1.500	2.000	2.000	0.050	40.0	25.0

## Confidence Intervals for the Interaction Odds Ratio in Logistic Regression with Two Binary X's

### References

- Demidenko, Eugene. 2007. 'Sample size determination for logistic regression revisited', *Statistics in Medicine*, Volume 26, pages 3385-3397.
- Demidenko, Eugene. 2008. 'Sample size and optimal design for logistic regression with binary interaction', *Statistics in Medicine*, Volume 27, pages 36-46.
- Rochon, James. 1989. 'The Application of the GSK Method to the Determination of Minimum Sample Sizes', *Biometrics*, Volume 45, pages 193-205.

### Report Definitions

Logistic regression equation:  $\text{Log}(P/(1-P)) = \beta_0 + \beta_1 \times X + \beta_2 \times Z + \beta_3 \times XZ$ , where  $P = \text{Pr}(Y = 1|X, Z)$  and  $X$  and  $Z$  are binary.

Confidence Level is the proportions of studies with the same settings that produce a confidence interval that includes the true ORint.

$N$  is the sample size.

C.I. Width is the distance between the two boundaries of the confidence interval.

ORint is the expected sample value of the interaction odds ratio. It is the value of  $\exp(\beta_3)$ .

ORyx is the expected sample value of the odds ratio. It is the value of  $\exp(\beta_1)$ .

ORyz =  $\exp(\beta_2)$  is the odds ratio of  $Y$  versus  $Z$ .

ORxz is the odds ratio of  $X$  versus  $Z$  in a logistic regression of  $X$  on  $Z$ .

C.I. of ORint Lower Limit is the lower limit of the confidence interval of ORint.

C.I. of ORint Upper Limit is the upper limit of the confidence interval of ORint.

$P_0$  is the response probability at  $X = 0$ . That is,  $P_0 = \text{Pr}(Y = 1|X = 0, Z = 0)$ .

Percent  $X=1$  is the percent of the sample in which the exposure is 1 (present).

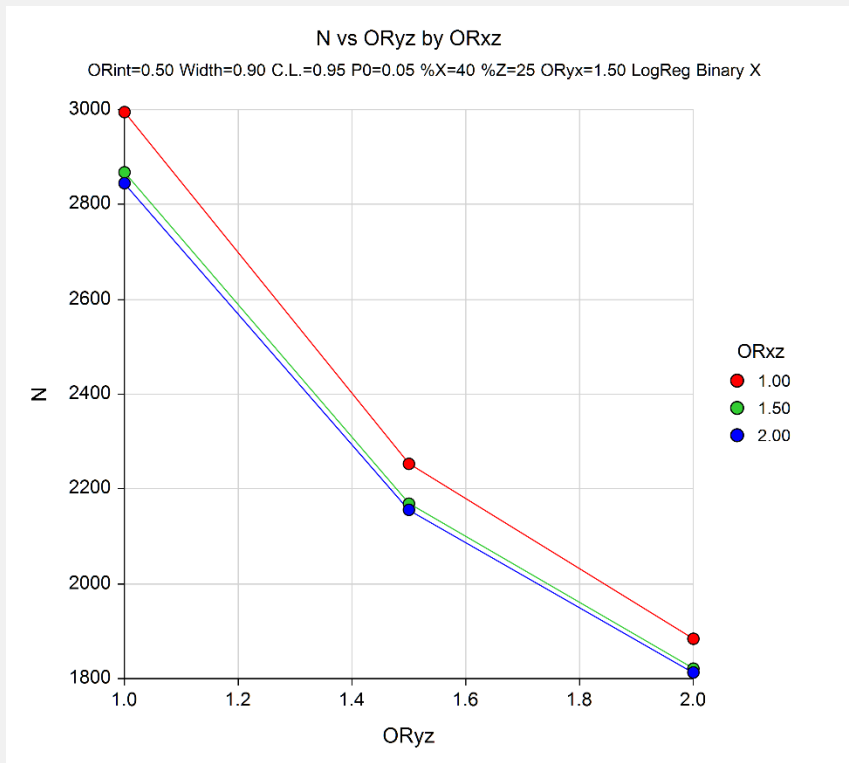
Percent  $Z=1$  is the percent of the sample in which the confounder is 1.

### Summary Statements

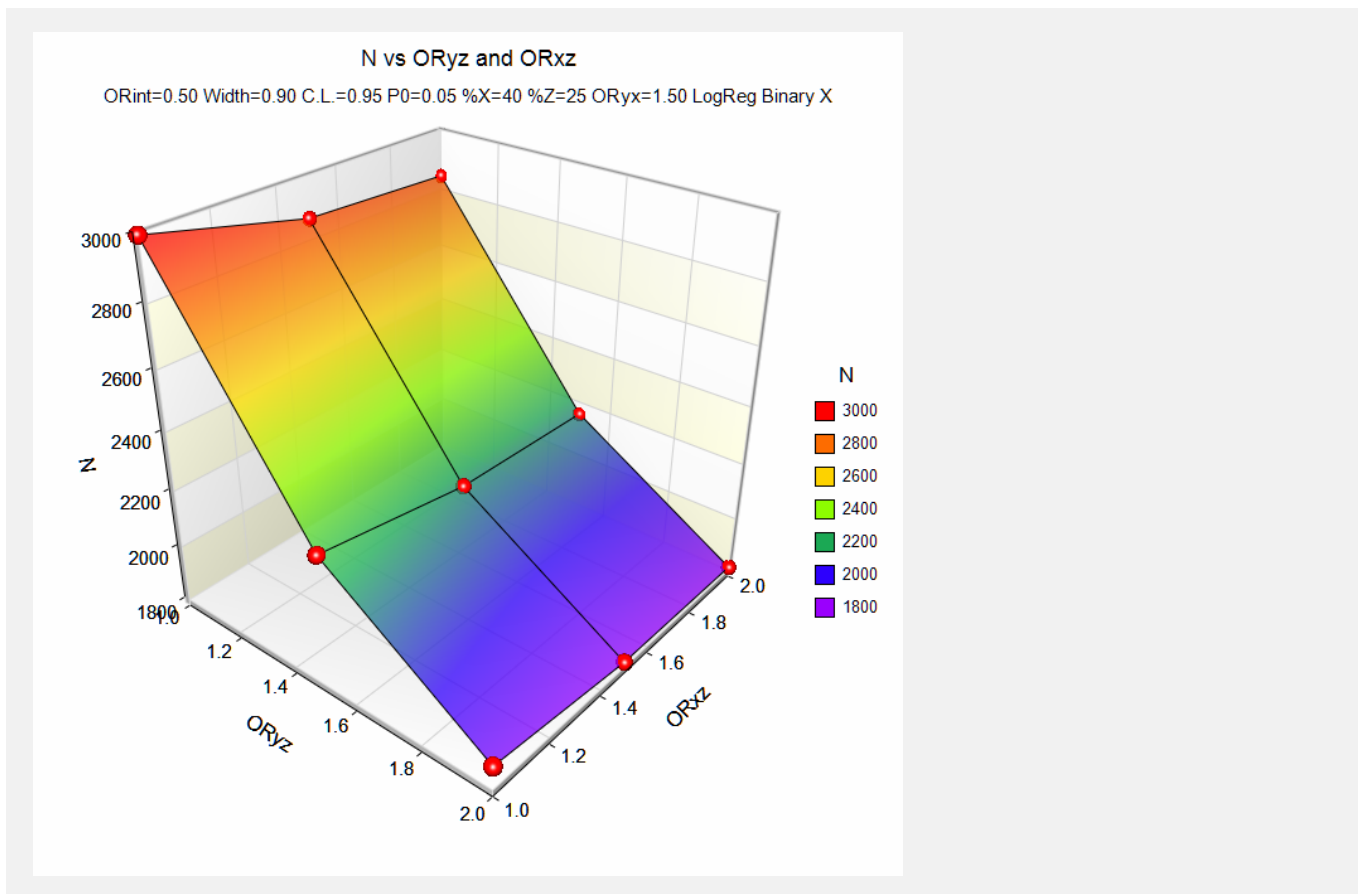
A logistic regression of a binary response variable ( $Y$ ) on two binary independent variables ( $X$  and  $Z$ ) and their interaction with a sample size of 2995 observations at a 0.950 confidence level produces a two-sided confidence interval for the odds ratio of the interaction with a width of 0.8999. The sample odds ratio of the interaction is assumed to be 0.500. Other settings are  $OR_{yx} = 1.500$ ,  $OR_{yz} = 1$ ,  $OR_{xz} = 1.000$ , and  $P_0$  (prevalence of  $Y$  given  $X = 0$  and  $Z = 0$ ) = 0.050. The prevalence of  $X$  is 40.0% and the prevalence of  $Z$  is 25.0%. A Wald statistic is used to construct the confidence interval.

This report shows the sample size for each of the scenarios.

### Plot Section



### Confidence Intervals for the Interaction Odds Ratio in Logistic Regression with Two Binary X's



These plots show the sample size for the various values of the other parameters.



## Example 2 – Validation for the Interaction

We could not find a direct validation result in the literature, so we will create one by hand. This is easy to do in this case because we can create a dataset, analyze it with a statistical program such as NCSS, and then compare these results to those obtained with the above formulas in PASS.

Here is a summary of the data that was used to generate this example. The numeric values are counts of the number of items in the corresponding cell.

Group	Y=1	Y=0	Total
X=1, Z=1	5	10	15
X=1, Z=0	3	21	24
X=0, Z=1	17	3	20
X=0, Z=0	9	7	16
<b>Total</b>	<b>34</b>	<b>41</b>	<b>75</b>

Here is a printout from NCSS showing the estimated odds ratio (0.79412) and confidence interval (0.08304 to 7.59380). The width is 7.51076.

Odds Ratios				
Independent Variable	Regression Coefficient	Odds Ratio	Lower 95% Confidence Limit	Upper 95% Confidence Limit
X	b(i)	Exp(b(i))		
Intercept	0.43853	1.55042	0.64509	3.72629
(X=1)	-2.19722	0.11111	0.02331	0.52968
(Z=1)	1.48329	4.40741	0.91195	21.30077
(X=1)*(Z=1)	-0.23052	0.79412	0.08304	7.59380

Note that the value of  $P_0$  is  $9 / 16 = 0.5625$ . The value of *Percent with X = 1* is  $100 \times 39 / 75 = 52\%$ . The value of *Percent with Z = 1* is  $100 \times 35 / 75 = 46.67\%$ . Also, a logistic regression of X on Z produced an  $OR_{xz}$  of 0.50.

## Setup

This section presents the values of each of the parameters needed to run this example. First, from the PASS Home window, load the **Confidence Intervals for the Interaction Odds Ratio in Logistic Regression with Two Binary X's** procedure. You may then make the appropriate entries as listed below, or open **Example 2** by going to the **File** menu and choosing **Open Example Template**.

Option	Value
<b>Design Tab</b>	
Solve For .....	Precision (C.I. Width)
Interval Type .....	Two-Sided
Confidence Level .....	0.95
N (Sample Size).....	75
$P_0$ [Pr(Y=1 X=0, Z=0)] .....	0.5625
ORint (Interaction Odds Ratio) .....	0.79412
ORyx (Y, X Odds Ratio) .....	0.1111
ORyz (Y, Z Odds Ratio).....	4.40741
ORxz (X, Z Odds Ratio).....	0.50
Percent with X = 1 .....	52
Percent with Z = 1.....	46.666667

## Output

Click the Calculate button to perform the calculations and generate the following output.

### Numeric Results

#### Numeric Results for Two-Sided Confidence Interval of ORint

Conf Level	N	C.I. Width	ORint	Lower C.L. ORint	Upper C.L. ORint	ORyx	ORyz	ORxz	P0	Pct X=1	Pct Z=1
0.950	75	7.51103	0.79412	0.08304	7.59407	0.11110	4.40741	0.50000	0.56250	52.0	46.7

Using the above settings, **PASS** calculates the confidence interval to be (0.08304, 7.59407) which leads to a C. I. Width of 7.51103. This matches the results obtained from the data to within rounding.