

community project

encouraging academics to share statistics support resources

All stcp resources are released under a Creative Commons licence

stcp-rothwell-med

Medical Statistics definitions

Medical statistics is a branch of statistics which focuses on medical applications. It introduces new methods for analysing proportions of events, which can be defined as $p = \frac{a}{n} = \frac{\text{number of events occurring}}{\text{number in group}}$ where a is a subset of n , in terms of risk. For example, the proportion of people who ate toast for breakfast would have a as the number of people who ate toast from the sample questioned and n as the total number of people in the sample.

This sheet will briefly explain various terms which arise in medical statistics regularly.

Risk/ Prevalence (P): The prevalence (or risk) of a disease is calculated as

$$P_{\text{disease}} = \frac{\text{number of people with disease}}{\text{total number of people within group}}$$

Risks are not always negative e.g. risk of surviving or probability of winning the lottery are calculated in the same way.

We are often interested in comparing risks from different groups and there are several ways of doing this. The following table will be used to demonstrate the formulae involved.

	Event occurs (Develop disease/ Died)	Event does not occur (Does not develop disease/ Survived)	Total	Risk of disease (death) by group
Exposed/ Treated	a	b	a+b	$p_{\text{exp}} = \frac{a}{(a+b)}$
Not Exposed/ Not treated	c	d	c+d	$p_{\text{unexp}} = \frac{c}{(c+d)}$

Relative risk: This measures how much more likely the event is to occur in one group compared to another.

The risk of developing the disease for the exposed population is $P_{\text{exp}} = \frac{a}{(a+b)}$.

The risk of developing the disease for the unexposed population is $P_{\text{unexp}} = \frac{c}{(c+d)}$.

$$\text{Relative risk} = RR = \frac{\text{risk of developing disease in exposed group}}{\text{risk of developing disease in unexposed group}} = \frac{P_{\text{exp}}}{P_{\text{unexp}}} = \frac{\frac{a}{(a+b)}}{\frac{c}{(c+d)}} = \frac{a(c+d)}{c(a+b)}$$

This is sometimes called a Risk Ratio. If $RR > 1$ then the risk of disease for the exposed group is larger than the risk of disease for the unexposed group.

Example: A randomised controlled trial investigated mortality rates within a year for 300 patients with lung cancer. The first group received a new chemotherapy treatment for lung cancer (New treatment) and the other received the standard chemotherapy treatment (Control treatment).

	Died	Survived	Total	Risk of dying
Control treatment	50	150	200	50/200 = 0.25
New treatment	10	90	100	10/100 = 0.1

The relative risk is $\frac{\text{risk of dying in control group}}{\text{risk of dying in new treatment group}} = \frac{0.25}{0.1} = 2.5$. This means that those in the control group were 2.5 times more likely to die than those in the treatment group. When calculating relative risks, it is easier to use the group with the highest risk in the numerator.

Sometimes a confidence interval is reported with the Relative Risk calculated from a sample. A confidence interval gives a range of likely values for the population relative risk.

95% Confidence Interval (CI) for a RR: For large samples this can be calculated using the natural logarithm (ln) because the confidence interval is not symmetrical.

First the variance of $\ln(RR)$ needs to be calculated:

$$Var(\ln RR) = \frac{1}{a} - \frac{1}{a+c} + \frac{1}{b} - \frac{1}{b+d} = \frac{1}{50} + \frac{1}{50+10} + \frac{1}{150} + \frac{1}{150+90} = 0.0475$$

95% Confidence interval for the $\ln(RR)$

$$\ln(RR) \pm 1.96 \times \sqrt{Var(\ln RR)} = \ln(2.5) \pm 1.96 \times \sqrt{0.0475} = 0.916 \pm 0.427 = (0.489, 1.343)$$

To get the confidence interval for the actual relative risk, take the exponential of the upper and lower value, so the 95% confidence interval will be

$$(e^{\ln(RR)-1.96 \times \sqrt{Var(\ln RR)}}, e^{\ln(RR)+1.96 \times \sqrt{Var(\ln RR)}}) = (e^{0.489}, e^{1.343}) = (1.63, 3.83).$$

The relative risk for the whole population is likely to be between 1.63 and 3.83. If the confidence interval includes 1, the risk in one group is not significantly higher than the risk in the second group. Here, both values are above 1 so the risk in the control group is significantly higher (RR 2.5, 95% CI: 1.63 to 3.83).

The **Relative Risk Difference (RRD)** is given by $RRD = 1 - RR$ when ($RR < 1$).

Put the smallest risk (treatment group) on top to get a RR under 1: $\frac{0.1}{0.25} = 0.4$. Therefore the RRD is $1 - 0.4 = 0.6$. The risk of dying is reduced by 60% in the treatment group.

The **Absolute Risk Difference (ARD)** is given by

$$ARD = |P_{exp} - P_{unexp}| = |0.25 - 0.1| = 0.15. \text{ The absolute risk has decreased by 15\%.}$$

The **Number Needed to Treat (NNT)** is the additional number of people you would need to give a new treatment to in order to cure one extra person compared to the

old treatment and is given by $NNT = \frac{1}{ARD} = \frac{1}{0.15} = 6.67$. So 7 people would need to receive the new treatment for one extra person to survive compared to the old treatment.

Another common measure used in medical statistics is the **Odds Ratio (OR)**. First, odds are calculated by $odds = \frac{p}{1-p}$, where p is the probability of an event occurring.

Therefore, the odds of disease in the exposed group would be $odds_{exp} = \frac{P_{exp}}{1-P_{exp}}$ and similarly the odds of disease in the unexposed group would be $odds_{unexp} = \frac{P_{unexp}}{1-P_{unexp}}$.

Then the odds ratio is $OR = \frac{odds_{exp}}{odds_{unexp}} = \frac{ad}{bc}$, which would be the odds of disease in the exposed group compared to the unexposed group. If the $OR > 1$ then the odds of disease occurring in the exposed group are larger than the odds of disease in the unexposed group, so exposure to the factor has increased the risk of contracting the disease. For our example the odds ratio is $OR = \frac{ad}{bc} = \frac{50 \times 90}{10 \times 150} = 3$, The odds ratio looks at the odds of being in a particular treatment group given that you had the disease. So, in this example, those who improved were 3 times more likely to have received the new treatment than the control treatment.

Note: The odds ratio and relative risk are similar if the total sample is large and the disease is rare.

95% Confidence Interval (CI) for an OR: This can be calculated for large samples and must be carried out using the natural logarithm (ln) because the confidence interval is not symmetrical.

First the variance of $\ln(OR)$ needs to be calculated.

$$Var(\ln OR) = \frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d} = \frac{1}{50} + \frac{1}{150} + \frac{1}{10} + \frac{1}{90} = 0.138$$

95% CI: $\ln(OR) \pm 1.96 \times \sqrt{Var(\ln OR)} = \ln(3) \pm 1.96 \times \sqrt{0.138} = (0.362, 1.818)$ then take the exponential of the upper and lower value, so the 95% confidence interval will be $(e^{\ln(OR)-1.96 \times \sqrt{Var(\ln OR)}}, e^{\ln(OR)+1.96 \times \sqrt{Var(\ln OR)}}) = (e^{0.362}, e^{1.818}) = (1.43, 6.16)$

The odds ratio comparing death rates after the standard treatment to the new treatment was 3 (95% CI: 1.43, 6.16) with those on the standard treatment being more likely to die.

Diagnostic Tests

		True Diagnosis		
		Disease +ve	Disease -ve	Total
Test Results	+ve	a	b	a+b
	-ve	c	d	c+d
		a+c	b+d	N

Sometimes there is a need to establish how good a diagnostic test is in detecting disease. One would have a table similar to the one above. A number of different measures can be gained from this information.

Sensitivity: $sens = \frac{a}{(a+c)}$ This is the probability of getting a positive test result *given that the person has the disease*. $P(+ve|D)$

Specificity: $spec = \frac{d}{(b+d)}$ This is the probability of getting a negative test result *given that the person does not have the disease*. $P(-ve|ND)$

Positive Predictive Value*: $PPV = \frac{a}{(a+b)}$ This is the probability of the person having the disease *given they get a positive test result*. $P(D|+ve)$

Negative Predictive Value*: $NPV = \frac{d}{(c+d)}$ This is the probability of the person not having the disease *given they get a negative test result*. $P(ND|-ve)$

Positive Likelihood Ratio: $LR^+ = \frac{sens}{1-spec}$ This gives a ratio of the test being positive for patients with disease compared with those without disease. Aim to be much greater than 1 for a good test.

Negative Likelihood Ratio: $LR^- = \frac{1-sens}{spec}$ This gives a ratio of the test being negative for patients with disease compared with those without disease. Aim to be much less than 1 for a good test.

General rule – A screening test needs high sensitivity, a diagnostic test needs high specificity.

*these tests must have a **random** sample of the whole population; they depend on the prevalence of the disease which cannot be calculated if the sample is not random.

For tests with a continuous outcome, such as a blood biomarker measurement, one can determine a good cut-off point for the test using an ROC curve. This plots sensitivity against (1-specificity). A good diagnostic test will be the point closest to the top left corner of the plot.

There will be a trade-off between high sensitivity and high specificity. Ideally both of these values should be high.

