

Introduction to Statistics for Research

Significance Tests

Martin Bland

Emeritus Professor of Health Statistics
University of York

<http://www-users.york.ac.uk/~mb55/>

An Example: the Sign Test

Consider a two treatment cross-over trial of pronethalol vs. placebo for the treatment of angina (Pritchard *et al.*, 1963).

Patients received placebo for two periods of two weeks and pronethalol for two periods of two weeks, in random order.

Completed diaries of attacks of angina.

Pritchard BNC, Dickinson CJ, Alleyne GAO, Hurst P, Hill ID, Rosenheim ML, Laurence DR. Report of a clinical trial from Medical Unit and MRC Statistical Unit, University College Hospital Medical School, London. *BMJ* 1963; 2: 1226-7.

Results of a trial of pronethalol for the treatment of angina pectoris (Pritchard *et al.*, 1963)

Patient	Placebo	Pronethalol	Placebo - Pronethalol
1	71	29	42
2	323	348	-25
3	8	1	7
4	14	7	7
5	23	16	7
6	34	25	9
7	79	65	14
8	60	41	19
9	2	0	2
10	3	0	3
11	17	15	2
12	7	2	5

These 12 patients are a sample from the population of all patients.

Would the other members of this population experience fewer attacks while using Pronethalol?

In a significance test, we ask whether the difference observed was small enough to have occurred by chance if there were really no difference in the population.

If it were so, then the evidence in favour of there being a difference between the treatment periods would be weak.

On the other hand, if the difference were much larger than we would expect due to chance if there were no real population difference, then the evidence in favour of a real difference would be strong.

Results of a trial of pronethalol for the treatment of angina pectoris (Pritchard *et al.*, 1963)

Patient	Placebo	Pronethalol	Placebo – Pronethalol
1	71	29	42
2	323	348	-25
3	8	1	7
4	14	7	7
5	23	16	7
6	34	25	9
7	79	65	14
8	60	41	19
9	2	0	2
10	3	0	3
11	17	15	2
12	7	2	5

Is there good evidence that Pronethalol reduces the number of attacks?
Most patients experience fewer attacks on Pronethalol.

To carry out the test of significance we suppose that, in the population, there is no difference between the two treatment periods.

The hypothesis of 'no difference' or 'no effect' in the population is called the **null hypothesis**.

We compare this with the **alternative hypothesis** of a difference between the treatments, in either direction.

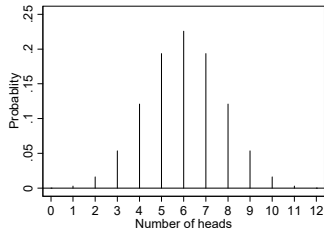
We find the probability of getting data as extreme as those observed if the null hypothesis were true.

If this probability is large the data are consistent with the null hypothesis; if it is small the data are unlikely to have arisen if the null hypothesis were true and the evidence is in favour of the alternative hypothesis.

The sign test

The number of negatives would behave in exactly the same way as the number of heads if we toss a coin 12 times.

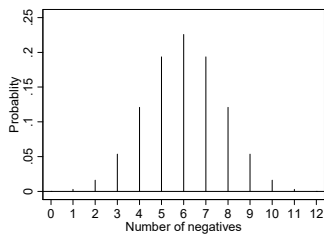
This is quite easy to investigate mathematically. We call it the Binomial Distribution with $n = 12$ and $p = 0.5$.



The sign test

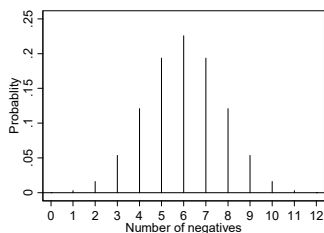
The number of negatives would behave in exactly the same way as the number of heads if we toss a coin 12 times.

This is quite easy to investigate mathematically. We call it the Binomial Distribution with $n = 12$ and $p = 0.5$.



The sign test

If there were any subjects who had the same number of attacks on both regimes we would omit them, as they provide no information about the direction of any difference between the treatments. In this test, n is the number of subjects for whom there is a difference, one way or the other.



Distribution of number of negatives if null hypothesis were true.

The sign test

The expected number of negatives under the null hypothesis is 6. The number of negative differences is 1. What is the probability of getting a value as far from this as is that observed?

-ves	Probability	-ves	Probability
0	0.00024	7	0.19336
1	0.00293	8	0.12085
2	0.01611	9	0.05371
3	0.05371	10	0.01611
4	0.12085	11	0.00293
5	0.19336	12	0.00024
6	0.22559		

The sign test

The expected number of negatives under the null hypothesis is 6. The number of negative differences is 1. What is the probability of getting a value as far from this as is that observed?

-ves	Probability	-ves	Probability
0	0.00024	7	0.19336
1	0.00293	8	0.12085
2	0.01611	9	0.05371
3	0.05371	10	0.01611
4	0.12085	11	0.00293
5	0.19336	12	0.00024
6	0.22559		

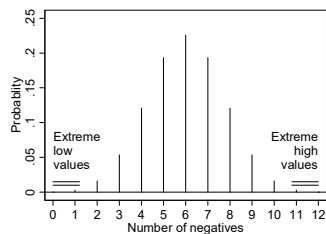
The sign test

The expected number of negatives under the null hypothesis is 6. The number of negative differences is 1. What is the probability of getting a value as far from this as is that observed?

-ves	Probability	-ves	Probability
0	0.00024	7	0.19336
1	0.00293	8	0.12085
2	0.01611	9	0.05371
3	0.05371	10	0.01611
4	0.12085	11	0.00293
5	0.19336	12	0.00024
6	0.22559		

The sign test

The expected number of negatives under the null hypothesis is 6. The number of negative differences is 1. What is the probability of getting a value as far from this as is that observed?



The sign test

The expected number of negatives under the null hypothesis is 6. The number of negative differences is 1. What is the probability of getting a value as far from this as is that observed?

-ves	Probability
0	0.00024
1	0.00293
11	0.00293
12	0.00024

Total	0.00634

The sign test

The probability of getting as extreme a value as that observed, in either direction, is 0.00634.

If the null hypothesis were true we would have a sample which is so extreme that the probability of it arising by chance is 0.006, less than one in a hundred.

Thus, we would have observed a very unlikely event if the null hypothesis were true.

The data are not consistent with null hypothesis, so we can conclude that there is strong evidence in favour of a difference between the treatment periods.

(Since this was a double blind randomized trial, it seems reasonable to suppose that this was caused by the activity of the drug.)

The sign test

The sign test is an example of a test of significance.

The number of negative changes is called the **test statistic**, something calculated from the data which can be used to test the null hypothesis.

Principles of significance tests

The general procedure for a significance test is as follows:

1. Set up the null hypothesis and its alternative.
2. Check any assumptions of the test.
3. Find the value of the test statistic.
4. Refer the test statistic to a known distribution which it would follow if the null hypothesis were true.
5. Find the probability of a value of the test statistic arising which is as or more extreme than that observed, if the null hypothesis were true.
6. Conclude that the data are consistent or inconsistent with the null hypothesis.

Principles of significance tests

The general procedure for a significance test is as follows:

1. Set up the null hypothesis and its alternative.

Null hypothesis:

'No difference between treatments' OR 'Probability of a difference in number of attacks in one direction is equal to the probability of a difference in number of attacks in the other direction'.

Alternative hypothesis:

'A difference between treatments' OR 'Probability of a difference in number of attacks in one direction is not equal to the probability of a difference in number of attacks in the other direction'.

Principles of significance tests

The general procedure for a significance test is as follows:

1. Set up the null hypothesis and its alternative.
2. Check any assumptions of the test.

Assumption:

That the patients are independent.

Principles of significance tests

The general procedure for a significance test is as follows:

1. Set up the null hypothesis and its alternative.
2. Check any assumptions of the test.
3. Find the value of the test statistic.

Test statistic:

Number of negatives (= 1).

Principles of significance tests

The general procedure for a significance test is as follows:

1. Set up the null hypothesis and its alternative.
2. Check any assumptions of the test.
3. Find the value of the test statistic.
4. Refer the test statistic to a known distribution which it would follow if the null hypothesis were true.

Known distribution:

Binomial, $n = 12$, $p = 0.5$.

Principles of significance tests

The general procedure for a significance test is as follows:

1. Set up the null hypothesis and its alternative.
2. Check any assumptions of the test.
3. Find the value of the test statistic.
4. Refer the test statistic to a known distribution which it would follow if the null hypothesis were true.
5. Find the probability of a value of the test statistic arising which is as or more extreme than that observed, if the null hypothesis were true.

Probability:

$P = 0.006$

Principles of significance tests

The general procedure for a significance test is as follows:

1. Set up the null hypothesis and its alternative.
2. Check any assumptions of the test.
3. Find the value of the test statistic.
4. Refer the test statistic to a known distribution which it would follow if the null hypothesis were true.
5. Find the probability of a value of the test statistic arising which is as or more extreme than that observed, if the null hypothesis were true.
6. Conclude that the data are consistent or inconsistent with the null hypothesis.

Conclusion: inconsistent.

Principles of significance tests

There are many different significance tests, all of which follow this pattern.

Statistical significance

If the data are not consistent with the null hypothesis, the difference is said to be **statistically significant**.

If the data are consistent with the null hypothesis, the difference is said to be **not statistically significant**.

We can think of the significance test probability as an index of the strength of evidence against the null hypothesis.

The probability of such an extreme value of the test statistic occurring if the null hypothesis were true is often called the **P value**.

It is **not** the probability that the null hypothesis is true. The null hypothesis is either true or it is not; it is not random and has no probability.

Significance levels and types of error

How small is small? A probability of 0.006, as in the example above, is clearly small and we have a quite unlikely event. But what about 0.06, or 0.1?

Suppose we take a probability of 0.01 or less as constituting reasonable evidence against the null hypothesis. If the null hypothesis is true, we shall make a wrong decision one in a hundred times.

Deciding against a true null hypothesis is called an **error of the first kind, type I error, or α (alpha) error**.

We get an **error of the second kind, type II error, or β (beta) error** if we decide in favour of a null hypothesis which is in fact false.

Significance levels and types of error

The smaller we demand the probability be before we decide against the null hypothesis, the larger the observed difference must be, and so the more likely we are to miss real differences.

By reducing the risk of an error of the first kind we increase the risk of an error of the second kind.

	Null hypothesis true	Alternative hypothesis true
Test not significant	No error	Type II error, beta error
Test significant	Type I error, alpha error.	No error

Significance levels and types of error

The smaller we demand the probability be before we decide against the null hypothesis, the larger the observed difference must be, and so the more likely we are to miss real differences.

By reducing the risk of an error of the first kind we increase the risk of an error of the second kind.

The conventional compromise is to say that differences are significant if the probability is less than 0.05.

This is a reasonable guideline, but should not be taken as some kind of absolute demarcation.

If we decide that the difference is significant, the probability is sometimes referred to as the **significance level**.

Interpreting the P value

As a rough and ready guide, we can think of P values as indicating the strength of evidence like this:

P value	Evidence for a difference or relationship
Greater than 0.1:	Little or no evidence
Between 0.05 and 0.1:	Weak evidence
Between 0.01 and 0.05:	Evidence
Less than 0.01:	Strong evidence
Less than 0.001:	Very strong evidence

Significant, real and important

If a difference is not statistically significant, it could still be real.

We may simply have too small a sample to show that a difference exists.

Furthermore, the difference may still be important.

'Not significant' does not imply that there is no effect.

It means that we have not demonstrated the existence of one.

Presenting P values

Computers print out the exact P values for most test statistics.

These should be given, rather than change them to 'not significant', 'ns' or $P > 0.05$.

Similarly, if we have $P = 0.0072$, we are wasting information if we report this as $P < 0.01$.

This method of presentation arises from the pre-computer era, when calculations were done by hand and P values had to be found from tables.

Personally, I would quote this to one significant figure, as $P = 0.007$, as figures after the first do not add much, but the first figure can be quite informative.

Presenting P values

Sometimes the computer prints 0.0000. This may be correct, in that the probability is less than 0.00005 and so equal to 0.0000 to four decimal places.

The probability can rarely be **exactly** zero, so we usually quote this as $P < 0.0001$.

Significance tests and confidence intervals

Often involve similar calculations.

If CI does not include the null hypothesis value, the difference is significant.

E.g. for a difference between two proportions, null hypothesis value = 0.

If 95% CI contains zero, difference is not significant.

If 95% CI does not contain zero, difference is significant.

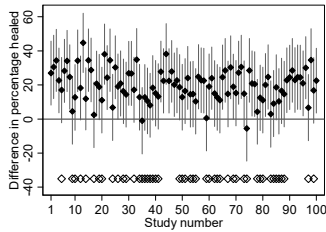
E.g. ulcer healing 63% (31/49) vs. 50% (26/52).

95% CI for difference: -7 to +33 percentage points.

Difference could be zero. Not significant.

Significance tests and confidence intervals

Ulcer healing simulation:



Open symbols denote no significant differences.

Significance tests and confidence intervals

The null hypothesis may contain information about the standard error.

E.g. comparison of two proportions, the standard error for the difference depends on the proportions themselves.

If the null hypothesis is true we need only one estimate of the proportion.

This alters the standard error for the difference.

Confidence interval: $SE = 0.0977$

Significance test: $SE = 0.0987$

95% CI and 5% significance test sometimes give different answers near the cut-off point.

Multiple significance tests

If we test a null hypothesis which is in fact true, using 0.05 as the critical significance level, we have a probability of 0.95 of coming to a 'not significant' (i.e. correct) conclusion.

If we test two independent true null hypotheses, the probability that neither test will be significant is $0.95 \times 0.95 = 0.90$.

If we test twenty such hypotheses the probability that none will be significant is $0.95^{20} = 0.36$.

This gives a probability of $1 - 0.36 = 0.64$ of getting at least one significant result.

We are more likely to get one than not.

The expected number of spurious significant results is $20 \times 0.05 = 1$.

Multiple significance tests

Many medical research studies are published with large numbers of significance tests.

These are not usually independent, being carried out on the same set of subjects, so the above calculations do not apply exactly.

If we go on testing long enough we will find something which is 'significant'.

We must beware of attaching too much importance to a lone significant result among a mass of non-significant ones.

It may be the one in twenty which we should get by chance alone.

Multiple significance tests

- ❖ Many subgroups.
- ❖ Many outcome variables.

Many subgroups

Williams *et al.* (1992) randomly allocated elderly patients discharged from hospital to two groups: timetabled visits by health visitor assistants versus no visit unless there was perceived need.

Patients assessed for physical health, disability, and mental state using questionnaire scales.

No significant differences overall between the intervention and control groups.

Williams, E.I., Greenwell, J., and Groom, L.M. (1992) The care of people over 75 years old after discharge from hospital: an evaluation of timetabled visiting by Health Visitor Assistants. *Journal of Public Health Medicine* 14, 138-44.

Many subgroups

Williams *et al.* (1992)

Among women aged 75-79 living alone the control group showed significantly greater deterioration in physical score than did the intervention group (P=0.04), and among men over 80 years the control group showed significantly greater deterioration in disability score than did the intervention group (P=0.03).

The authors stated that 'Two small sub-groups of patients were possibly shown to have benefited from the intervention. . . . These benefits, however, have to be treated with caution, and may be due to chance factors.'

Many subgroups: Bonferroni correction

Multiply the P values by the number of tests.

If any is then significant, the test of the overall composite null hypothesis is significant.

E.g. Williams *et al.* (1992).

Subjects were cross-classified by age groups, whether living alone, and sex, so there were at least eight subgroups, if not more.

Even if we consider the three scales separately, the true P values are $8 \times 0.04 = 0.32$ and $8 \times 0.03 = 0.24$.

Composite null hypothesis: there is a difference between the treatments in at least one group of subjects.

Many subgroups: Bonferroni correction

Composite null hypothesis: there is a difference between the treatments in at least one group of subjects.

This is **not** the same as: the difference between the treatments varies between different group of subjects.

This needs a test of interaction (regression lecture).

Multiple outcome measurements

E.g. Newnham *et al.* (1993) randomized pregnant women to receive a series of Doppler ultrasound blood flow measurements or to control.

They found a significantly higher proportion of birthweights below the 10th and 3rd centiles ($P=0.006$ and $P=0.02$).

These were only two of many comparisons. At least 35 were reported in the paper, though only these two were reported in the abstract.

Birthweight was not the intended outcome variable for the trial.

Newnham, J.P., Evans, S.F., Con, A.M., Stanley, F.J., Landau, L.I. (1993) Effects of frequent ultrasound during pregnancy: a randomized controlled trial. *Lancet* 342, 887-91.

Multiple outcome measurements

These tests are not independent, because they are all on the same subjects, using variables which may not be independent.

The proportions of birthweights below the 10th and 3rd centiles are clearly not independent, for example.

We can apply the Bonferroni correction.

For the example, the P values could be adjusted by $35 \times 0.006 = 0.21$ and $35 \times 0.02 = 0.70$.

Because the tests are not independent, the adjusted P value is too big.

Test is **conservative**.

One- and two-sided tests of significance

In the pronethalol example, the alternative hypothesis was that there was a difference in one or other direction.

This is called a **two sided** or **two tailed** test, because we used the probabilities of extreme values in both directions.

One sided or **one tailed** test:

Alternative hypothesis: in the population, the number of attacks on the placebo is greater than the number on pronethalol.

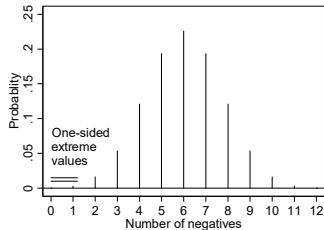
Null hypothesis: in the population, the number of attacks on the placebo is less than or equal to the number on pronethalol.

$P = 0.003$, and of course, a higher significance level than the two sided test.

One- and two-sided tests of significance

One sided null hypothesis: the number of attacks on the placebo is less than or equal to the number on pronethalol.

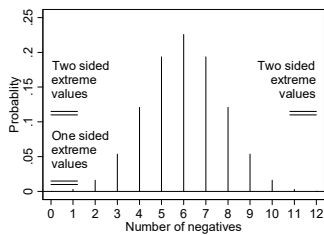
One sided alternative hypothesis: the number of attacks on the placebo is greater than the number on pronethalol.



One- and two-sided tests of significance

Two sided null hypothesis: the number of attacks on the placebo is equal to the number on pronethalol.

Two sided alternative hypothesis: the number of attacks on the placebo is not equal to the number on pronethalol.



One- and two-sided tests of significance

One sided or one tailed test:

Alternative hypothesis: in the population, the number of attacks on the placebo is greater than the number on pronethalol.

Null hypothesis: in the population, the number of attacks on the placebo is less than or equal to the number on pronethalol.

This implies that a decrease in the placebo direction would have the same interpretation as no change.

Seldom true in health research.

Tests should be two sided unless there is a good reason not to do this.

Pitfalls of significance tests

You should never, ever, conclude that there is no difference or relationship because it is not significant.

You should not rely on significance tests alone if you can give confidence intervals. Particularly useful when the test is not significant.

You should give exact P values where possible, not $P < 0.05$ or $P = NS$, though only one significant figure is necessary.

You should avoid multiple testing. Be clear what the main hypothesis and outcome variable are.
