# The AB/BA Cross-over: How to perform the two-stage analysis if you can't be persuaded that you shouldn't.

Stephen Senn, Department of Statistical Science and Department of Epidemiology and Public Health, University College London

# The AB/BA Cross-over: How to perform the two-stage analysis if you can't be persuaded that you shouldn't.

**Stephen Senn, Department of Statistical Science and Department of Epidemiology and Public Health, University College London**

## 1. Introduction

It is both a pleasure and an honour to contribute an article to Roel van Strik's *Liber Amicorum*. Roel has never missed an International Society for Biostatistics (ISCB) meeting and, although I cannot quite match this record, I have tried to attend (and succeeded in attending) a good few myself over the years so that I can look back on many pleasant conversations with Roel on statistical and other matters, as well as some memorable social evenings at ISCB spent in his company over a beer or two. A recent example can be given. We set off together at ISCB17 in Budapest on a trip to view the Danube Bend, were rather surprised to find ourselves later in the evening witnessing (at alarmingly close quarters) a medieval tournament which convincingly blended point processes and survival analysis (the mosquitoes in the Camargue at ISCB11 in Nimes were the previous such demonstration I recall), and finished the evening in a nearby restaurant discussing, amongst other things, 'whatever happened to multivariate analysis?', 'is explained variation important?', 'which was the year that ISCB didn't meet?', 'and do we prefer red or white Hungarian wine?'.

Let me quote from the membership information for ISCB. '*The International Society for Clinical Biostatistics (ISCB) was founded in 1978 to stimulate research into the principles and methodology used in the design and analysis of clinical research and to increase the relevance of statistical theory to the real world of clinical medicine.*' This gives me the excuse for my contribution. Both Roel and I believe that although statistical theory is fascinating, medical statistics must also be judged by its relevance. (We should hardly have made ISCB such a favourite if we did not.) The AB/BA cross-over design is a powerful, useful and therefore relevant tool of medical research. It also raises some fascinating philosophical and statistical issues, some of which I shall address below.

## 2. A lightning tour of the AB/BA design

In an AB/BA cross-over, patients are allocated at random to receive in a first period, either treatment A followed (perhaps after a suitable wash-out) by treatment B in a second period, or treatment B followed by treatment A[1,2]. The cross-classification of treatments and sequences provides four cells for each of which (for continuous outcomes) we may calculate a mean. These statistics and their expectations (in terms of a possible model) are illustrated in table I below.

**Table I. Cell means and expectations for an AB/BA cross-over.**

| | | Period 1 | Period 2 |
|---|---|---|---|
| | **AB** | $\mu + \tau_A + \pi_1$ | $\mu + \tau_B + \pi_2 + \lambda_{AB}$ |
| | | $\overline{Y}_{11}$ | $\overline{Y}_{12}$ |
| **Sequence** | | | |
| | **BA** | $\mu + \tau_B + \pi_1$ | $\mu + \tau_A + \pi_2 + \lambda_{BA}$ |
| | | $\overline{Y}_{21}$ | $\overline{Y}_{22}$ |

The purpose of the trial is to estimate the treatment contrast $\tau = \tau_A - \tau_B$. The various other parameters, the period effects, $\pi_1$ and $\pi_2$, the two carry-over effects, $\lambda_{AB}$ and $\lambda_{BA}$, as well as the general level parameter, $\mu$, can be regarded as nuisance parameters we wish to eliminate, if possible, from any estimate. If $n_1$ patients are allocated to sequence AB and $n_2$ to sequence BA , the variance of individual measurements is $\sigma^2$ and the correlation between them is $\rho$,

then the variance-covariance matrix is as given in Table II below[2,3]. Cell means from different sequences are independent of each other, being calculated from measurements on different patients, so that their correlation is 0, whereas cell means calculated for the same sequence are calculated from the same patients and have correlation $\rho$.

**Table II. Variances and covariances for the cell means.**

|  | $\overline{Y}_{11}$ | $\overline{Y}_{12}$ | $\overline{Y}_{21}$ | $\overline{Y}_{22}$ |
|---|---|---|---|---|
| $\overline{Y}_{11}$ | $\sigma^2/n_1$ | | | |
| $\overline{Y}_{12}$ | $\rho\sigma^2/n_1$ | $\sigma^2/n_1$ | | |
| $\overline{Y}_{21}$ | 0 | 0 | $\sigma^2/n_2$ | |
| $\overline{Y}_{22}$ | 0 | 0 | $\rho\sigma^2/n_2$ | $\sigma^2/n_2$ |

These cell means may now be used to form contrasts which enable us to estimate parameters of interest. These are as given in Table III below[2, 4].

**Table III. Contrasts and expectations for an AB/BA design.**

| Name | Definition | Expectation |
|---|---|---|
| CROS | $[(\overline{Y}_{11} - \overline{Y}_{12}) + (\overline{Y}_{22} - \overline{Y}_{21})]/2$ | $(\tau_A - \tau_B) - (\lambda_{AB} - \lambda_{BA})/2 = \tau - \lambda/2$ |
| SEQ | $(\overline{Y}_{11} + \overline{Y}_{12}) - (\overline{Y}_{22} + \overline{Y}_{21})$ | $(\lambda_{AB} - \lambda_{BA}) = \lambda$ |
| PAR | $\overline{Y}_{11} - \overline{Y}_{21}$ | $(\tau_A - \tau_B) = \tau$ |

PAR is an unbiased estimate of the treatment effect and SEQ is an unbiased estimate of the carry-over effect. CROS is a potentially biased (in the presence of carry-over) estimate of the treatment effect. The reason that it is of interest is that it is a within-patient estimate based on all four cell means and hence potentially much more efficient than PAR. If $q = 1/n_1 + 1/n_2$, then Table IV gives the variance-covariance matrix for the three statistics[2]. It can be seen that CROS is more efficient than PAR and (since $\rho$ will almost certainly be positive and often much greater than zero) usually much more efficient.

**Table IV. Variances and covariances for three contrast of interest.**

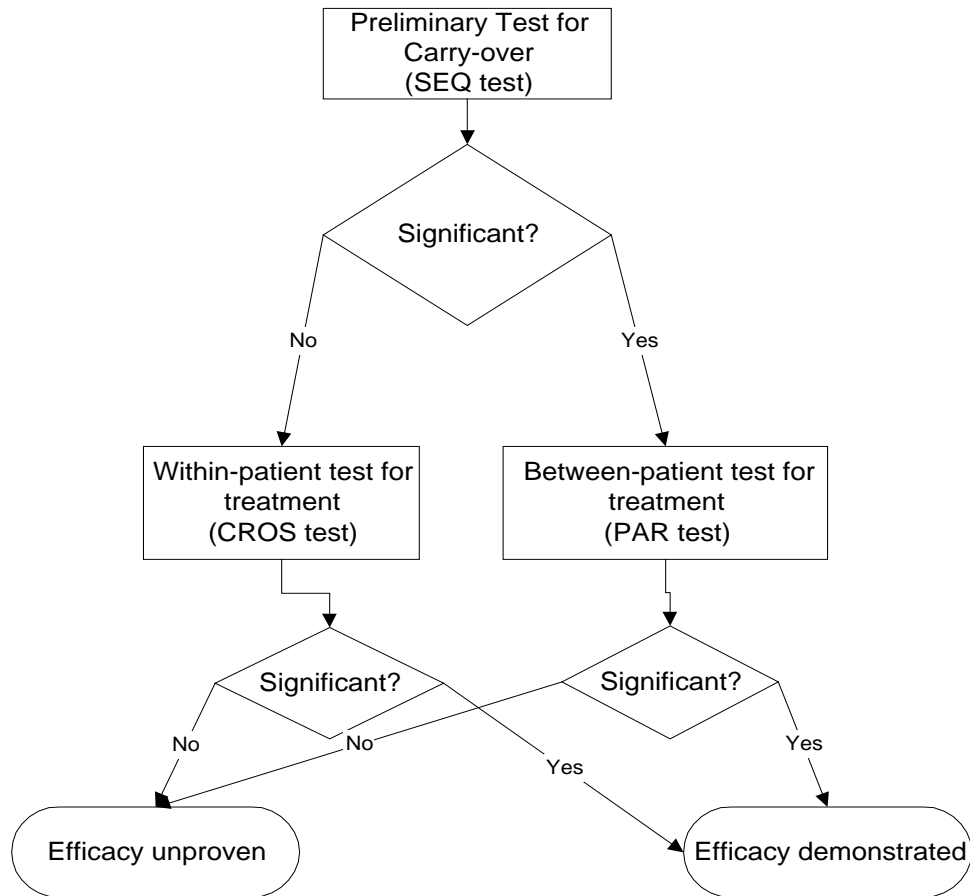|  | CROS ($L_c$) | PAR ($L_p$) | SEQ ($L_s$) |
|---|---|---|---|
| CROS ($L_c$) | $q(1 - \rho)\sigma^2/2$ | | |
| PAR ($L_p$) | $q(1 - \rho)\sigma^2$ | $q\sigma^2$ | |
| SEQ ($L_s$) | 0 | $q(1 + \rho)\sigma^2$ | $2q(1 + \rho)\sigma^2$ |

## 3. The two-stage procedure

For nearly quarter of a century, between its proposal in 1965 by Grizzle[3] and its critical examination in 1989 by Freeman[4], a procedure which was regarded as being acceptable by many statisticians was as given in Figure 1. This could be presented as an acceptable way of steering a course between the Scylla of always relying on the efficient but potentially biased CROS estimate and the Charybdis of always using the inefficient PAR estimate. The pilot in steering this delicate course was to be, SEQ, which forms the basis of a test of carry-over. (Grizzle proposed it should be carried out at the 10% level because of its low power[3].) What was not generally appreciated, however was that in trying to avoid the monsters, one had been created. What Freeman pointed[4] out was, that although CROS and SEQ were independent, PAR and SEQ were highly correlated (see Table IV), with a correlation which would, in practice, lie between $1/\sqrt{2}$ and 1. As a consequence, having observed a 'significant' carry-over effect, there was a very high probability that a significant value of PAR would be obtained and therefore the type I error rate of the two-stage procedure as a whole would exceed the nominal level claimed.

The implication of this result, that the two-stage procedure is misleading, has been generally accepted since, by statisticians researching into the design so that, for example, three books devoted to the cross-over trial suggest that it should not be used[1,5,6], and in particular those by myself[1] and Ratkowsky et al are quite categorical on this point[6]. More recently, however,

Jones and Lewis have investigated the AB/BA design[7] and have come to the conclusion, with which I agree and always have agreed, that the AB/BA design is an extremely useful one, and the further conclusion, with which I do not agree and never have agreed, that the two-stage analysis may be acceptable. (Or at the least they come to the conclusion that it is not much worse than CROS.) It is this latter point which I wish to explore below.

**Figure I. The two-stage procedure**



## 4. Power and size.
Jones and Lewis[7] perform a simulation to compare the power of PAR, CROS and the two-stage approach. In fact, provided we are happy to accept asymptotic results, numerical integration of the bivariate Normal distribution gives the necessary probabilities and I shall use this approach in due course. However, note first of all, that since the two-stage procedure has an inflated size, a comparison based on power alone is illegitimate. For example, the procedure of always declaring the treatment effect significant without looking at the data has a power of 100%. This is uniformly greater than any of the other three procedures, yet it is clearly inferior to all of them and no one would propose it. To compare such procedures it is either necessary to consider some linear combination of type I and type II errors or, if we are to be more strictly classical, to compare the power of procedures which have a type I error rate less than or equal to some target level, say 5%.

It turns out, however, that a simple argument may be used to correct the two-stage procedure. Suppose that we take the most pessimistic case and assume that the correlation with SEQ and PAR is 1 and assume further that there is no treatment effect and hence (which would reasonably follow) that there is no carry-over effect. If carried out at the conventional 10% level, SEQ will be 'significant' in 10% of all cases. In half of these significant cases it will be significant at the 5% level. Since PAR and SEQ are assumed to be perfectly correlated, it follows that wherever SEQ is significant at the 5% level, PAR will be significant also. Hence, *conditional* on obtaining a significant carry-over effect for SEQ, the two-stage procedure will give a significant treatment result in 50% of the cases. To reduce this to the desired level of

5% all that is necessary is to carry out the test of PAR at the nominal level of 0.5%. This adjustment will produce a corrected two-stage procedure which is conservative. The overall size of the test for treatment now lies between approximately 4.75% and 5%, rather than between 7% and 9.5% as in the form originally proposed by Grizzle[3].

Table V gives the values for the example considered by Jones and Lewis[7]. They supposed that a cross-over was performed with 22 patients per sequence for a clinically relevant difference of $\tau = 5$, and a total variance of 144 with a correlation $\rho$, between repeat measurements of 2/3. They then considered the power of the test of treatment for various values of $\lambda$, the difference in carry-overs, for CROS, PAR and the uncorrected two-stage procedure. The table gives asymptotic results for the power obtained by using Mathcad 6.0 Plus ® to perform double integration of the bivariate Normal whose parameters may be established by consideration of tables III and IV. As might be expected, the results of the integration are in very close agreement with the simulation of Jones and Lewis[7], if in general slightly higher (as befits asymptotic results). In addition the power for the *corrected* two stage analysis (2SC) proposed here has been given in the last column.

**Table V. Power for four approaches to analysing an AB/BA cross-over trial**

| $\lambda$ | CROS | PAR | 2SU | 2SC |
|---|---|---|---|---|
| 0 | 0.923 | 0.282 | 0.881 | 0.871 |
| 0.5 | 0.895 | 0.282 | 0.863 | 0.850 |
| 1 | 0.861 | 0.282 | 0.839 | 0.821 |
| 1.5 | 0.821 | 0.282 | 0.809 | 0.785 |
| 2 | 0.773 | 0.282 | 0.772 | 0.741 |
| 2.5 | 0.719 | 0.282 | 0.730 | 0.691 |
| 3 | 0.659 | 0.282 | 0.684 | 0.635 |
| 3.5 | 0.595 | 0.282 | 0.635 | 0.575 |
| 4 | 0.528 | 0.282 | 0.585 | 0.513 |
| 4.5 | 0.461 | 0.282 | 0.535 | 0.451 |
| 5 | 0.395 | 0.282 | 0.488 | 0.391 |

What the table shows, is that PAR is quite unacceptable on grounds of power. CROS is generally good but for values of $\lambda \geq 2$ the standard uncorrected two-stage procedure, 2SU has superior power. On this basis, Jones and Lewis were prepared to recommend it as a viable procedure[7]. What the table hides, however, is the size of the test. In the case of CROS and PAR it is 0.05 exactly. However, for 2SU it is 0.087, an inflation which is greater than that attendant upon a 'free look' in a sequential clinical trial (see for example, Pocock[8] page 148), a procedure which no regulator appears to regard as acceptable. To make a fair comparison of the two-stage procedure one needs to perform it in the corrected form. This has a size *very* slightly less than 0.05. (It is equal to 0.05 to at least 4 significant figures). The power of 2SC, however, shows it to be everywhere inferior to CROS and hence it cannot be recommended.
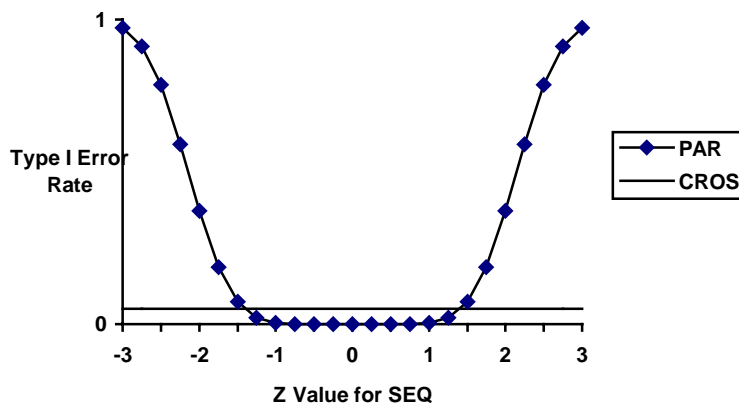
## 5. The philosophical issues
I promised some philosophical issues in the introduction and these will now be raised. They have to do with a) conditional power b) relevant subsets and c) the principle of total information.

Consider the possibility of further adjustments to the two-stage procedure. Quite apart from the issue as to whether the preliminary test should be carried out at the 10% level there is the issue of conditional size. For the calculations, I adjusted it so that the *conditional* type I error rate for PAR given that SEQ was significant was less than or equal to 5%. I could, however, have proposed dropping the level at which CROS was carried out and, by robbing Peter to pay Paul, used a higher conditional value for PAR. For example, using CROS at a level of 0.5% and PAR at a nominal level of 4.55%, which is a conditional level of 45.5% or less, leads to a procedure with overall size less than 0.9 x 0.005 + 0.1 x 0.455 = 0.05. Is this legitimate? The question is similar to that posed by Cox and Hinckley in the famous example of the two weighing machines[9]. The scientist will use the more or the less precise of two machines with a given probability. Is he required to control the same type I error rate for each machine separately or simply averaged over the two? The answer depends on approaches to inference.

Suppose we reject completely , as I consider we ought to, the two-stage analysis in any shape or form. Is it legitimate to perform an unconditional PAR analysis of an AB/BA design? Although this test may have low power, it does at least have correct size.  First, we may note that from one point of view, SEQ provides the means of dividing the sample space into what Fisher[10] (p35) called recognisable subsets. In fact, conditional on the value of SEQ, CROS and PAR have the same variance. They differ in terms of their bias. If we know that there is no carry-over, then large absolute values of SEQ are simply a means of identifying biased values of PAR. (This is one explanation of the poor performance of the two-stage analysis. It encourages us to use PAR where it is most biased[1,2].) That being so, should we not always interpret the result of testing PAR in the light of the observed value of SEQ? To do so simply brings us back to CROS. Or are we justified in ignoring the further information if, for whatever reason, we simply go ahead and use PAR? The importance of the question is illustrated by figure II, which shows a plot of the conditional type I error rate of the PAR and CROS tests (at the 5% level) as a function of the Z-value for the test of SEQ. The mean value for both curves is 0.05 but is this unconditional average (over the probability distribution of SEQ) relevant? Does not the Z-value for SEQ give us the means of recognising particular values of PAR?

**Figure II. Type I error rate for tests using PAR and CROS as a function of the standardised value of SEQ when $\rho = 2/3$**



Of course, all of this discussion is posited on the assumptions that neither the trial nor the carry-over effect are very large. These assumptions are reasonable in the context in which cross-over trials are carried out. If we look at the issue in mean square error terms we can see that for the usual sort of trial size and for plausible carry-over, the variance of PAR will dominate the bias of CROS. For *extremely* large trials, this might not be the case,  and there would be no reason not to use PAR but by the same token, there would be no purpose in carrying out a cross-over trial.

## 6. Is a Bayesian solution possible?
In principle the answer must be 'yes', but we haven't got it yet. Andy Grieve has done some most impressive work in this direction[11] but it still falls short of taking what, in my opinion, must be the crucial step: recognising that prior belief in carry-over must be dependent on belief in treatment. Until then, it seems to me, that a perfectly reasonable approach is to restrict the use of the design to cases where carry-over is strongly believed to be negligible. Where this is so, CROS clearly provides the key to a coherent Bayesian analysis.

### 7. Keep the design but dump the analysis.
Where does all this leave us? In my opinion, where any good ISCB member always should have been, which is considering the *relevance of statistical theory to the real world of clinical medicine.* It behoves every medical statistician who proposes a new approach to estimation to ask the following question: *am I seriously interested in finding out the effects of treatment*? To the extent that physicians between 1965 and 1989 followed our advice as a profession to carry-out the two-stage procedure and to the extent that it made any difference to the final result, we did them and their future patients a disservice. The solution to carry-over does not lie in yet more complicated analysis but in collaborating with our medical colleagues, striving to understanding pharmacology, and choosing designs appropriately. It is time that we restored the AB/BA design to them, encouraged them to use the simple CROS analysis and looked for more important things to discuss.

### 8. Finally, a word to Roel.
Whatever we find to discuss with our medical colleagues, Roel, I look forward to a good few more discussions together at future ISCB meetings, as well as some practical statistics. Let us agree to avoid the point processes and the survival analysis; even the cross-over trials can be dangerous (remember the Pont du Gard at Nimes!). Congratulations on your retirement. I look forward to seeing you at ISCB 18 in Boston for some judicious dose-finding (and I don't mean tea!).

### References
1.  Senn (1993) *Cross-over Trials in Clinical Research*, Chichester: John Wiley.
2.  S.J. Senn (1994) The AB/BA Cross-over Design, Past, Present and Future? *Statistical Methods in Medical Research,* **3**, 303-324
3.  Grizzle, J.E. (1965) The two-period change over design and its use in clinical trials, *Biometrics*, **21**, 467-480. Corrections Grizzle, 1974, *Biometrics*, **30**, 727 and Grieve, 1982, *Biometrics*, **38**, 517.
4.  Freeman, P.R. (1989) The performance of the two-stage analysis of two-treatment, two-period cross-over trials, *Statistics in Medicine*, **8**, 1421-1432.
5.  Jones, B.J. and Kenward, M.G. (1989) *Design and Analysis of Cross-Over Trials*, London: Chapman and Hall.
6.  Ratkowsky, D.A., Evans, M.A. and Alldredge, J.R.(1993) *Cross-over Experiments , Design, Analysis and Application*, New York: Marcel Dekker.
7.  Jones, B.J. and Lewis J. (1994) The case for cross-over trials in phase III. *Statistics in Medicine,* **14,** 1025-1038, 1995
8.  Pocock, S.J. (1983) *Clinical trials: a Practical Approach*, Chichester: Wiley.
9.  Cox, D.R. and Hinkley, D.V. (1974) *Theoretical Statistics*, London: Chapman and Hall.
10. Fisher, R.A. (1990) *Statistical Methods and Scientific Inference*, (3rd edition). Reprinted in *Statistical Methods, Experimental Design and Scientific Inference*, Oxford: Oxford Science Publications.
11. Grieve,A.P. (1985) A Bayesian analysis of the two-period cross-over trial, *Biometrics*, **41**, 979-990.