

1 A statistical test of the equality of latent orders

2 Michael L. Kalish^{a,*}, John C. Dunn^b, Oleg P. Burdakov^c, Oleg Sysoev^c

3 ^a*Syracuse University, USA*

4 ^b*University of Adelaide, Australia*

5 ^c*Linköping University, Sweden*

6 Abstract

It is sometimes the case that a theory proposes that the population means on two variables should have the same rank order across a set of experimental conditions. This paper presents a test of this hypothesis. The test statistic is based on the coupled monotonic regression algorithm developed by the authors. The significance of the test statistic is determined by comparison to an empirical distribution specific to each case, obtained via non-parametric or semi-parametric bootstrap. We present an analysis of the power and Type I error control of the test based on numerical simulation. Partial order constraints placed on the variables may sometimes be theoretically justified. These constraints are easily incorporated into the computation of the test statistic and are shown to have substantial effects on power. The test can be applied to any form of data, as long as an appropriate statistical model can be specified.

7 *Keywords:* state-trace analysis, monotonic regression, hypothesis test

*Corresponding address: Prof. M. L. Kalish, Department of Psychology, Syracuse University, Syracuse, NY, 13244

Email addresses: mlkalish@syr.edu (Michael L. Kalish), john.c.dunn@adelaide.edu.au (John C. Dunn), oleg.burdakov@liu.se (Oleg P. Burdakov), oleg.sysoev@liu.se (Oleg Sysoev)

1 Introduction

2 Consider an experiment in which data are obtained on two different variables
3 across k different conditions. We would like to know if these data are drawn
4 from populations whose means on the two variables have different orders. That
5 is, we ask if the variables have unequal *latent orders*. This question arises in the
6 theory of *state trace analysis* (STA) where inferences concerning the number of
7 latent variables underlying changes in two or more dependent variables depend
8 on the ordinal arrangements of their respective population means (Bamber,
9 1979; Prince et al., 2012a). STA contrasts a *one-dimensional model*, in which
10 changes in the dependent variables are mediated by one latent variable, and a
11 *two-dimensional model*, in which changes are mediated by more than one latent
12 variable (Loftus et al., 2004; Newell & Dunn, 2008). Under the assumption of the
13 one-dimensional model that each dependent variable is a (distinct) monotonic
14 function of the single latent variable, this model predicts that the latent orders
15 of the two variables are equal. It follows that if the variables have different
16 latent orders across a set of experimental conditions then the effects must be
17 mediated by more than one latent variable.

18 Implementation of STA requires a statistical procedure to test whether two
19 sets of population means have the same order across a set of conditions. To
20 our knowledge, at least three previous approaches to this problem have been
21 proposed in the psychological literature. The first of these, described by Loftus
22 et al. (2004), relies on reducing sampling error to near zero thereby using the
23 observed sample means as a proxy for the population means. Clearly, this ap-
24 proach cannot be applied in situations with non-negligible sampling error and
25 it lacks a means of quantifying when the sampling error is small enough to be
26 ignored. The second approach, described by Pratte & Rouder (2012), quanti-
27 fies the effects of sampling error but is limited to particular theory-dependent
28 dependent variables and to a fixed two-by-two factorial design. The third ap-
29 proach, described by Prince et al. (2012a), uses Bayesian model selection to

1 test whether two sets of population means have the same or different orders.
2 While the approach is in principle quite general, the particular implementation
3 described by Prince et al. (2012a) applies only to binomial data and to a rela-
4 tively constrained factorial design. We discuss this approach in greater detail
5 below and compare it to the test that we develop.

6 The test we present here is a null hypothesis statistical test (NHST), based
7 on the computation of an empirical p -value of the data given the null hypothe-
8 sis. Despite the well known problems with p -values (Wagenmakers, 2007), the
9 evidence provided by them remains useful; e.g., it predicts future replicability
10 (Open Science Collaboration, 2015).

11 The outline of the paper is as follows. First, we describe more fully the logic
12 of our statistical test, based on an extension of monotonic regression (Burdakov
13 et al., 2012). In so doing, we introduce the concept of partial order constraints
14 and foreshadow how they may be used to increase statistical power. Second,
15 we describe a null hypothesis significance test of the equality of latent orders
16 based on a bootstrap resampling procedure for estimating the empirical sam-
17 pling distribution of the test statistic. Third, we examine the statistical power
18 of our procedure for a fully randomized design with and without partial order
19 constraints. Finally, we extend the procedure to binomial data and compare it
20 to the Bayesian model selection approach developed by Prince et al. (2012a).

21 *The orders of sample and population means*

22 Consider two different dependent variables, x and y , observed across k
23 different experimental conditions. Let $x_1, \dots, x_k, y_1, \dots, y_k$, be the k pop-
24 ulation means of each variable and let $X_1, \dots, X_k, Y_1, \dots, Y_k$, be the corre-
25 sponding sample means. We define the (latent) order of x as a permutation,
26 $O(x) = (i_1, i_2, \dots, i_k)$, such that, $x_{i_1} \leq x_{i_2} \leq \dots \leq x_{i_k}$. We wish to test the
27 hypothesis that $O(x) = O(y)$, given the data. A desirable feature of such a
28 test is that it should be sensitive to both the number and magnitude of dif-

1 ferences in the two orders. Intuitively, given equal latent orders, numerically
 2 small violations of equality of the orders of the observed means are more likely
 3 than numerically large violations. This property is a feature of monotonic (or
 4 isotonic) regression (Robertson et al., 1988). Our test is based on this method.

5 *Monotonic Regression*

6 Monotonic regression addresses the problem of finding the best approxi-
 7 mation, \hat{X} , to a set of observed values, X , under the constraint that $O(\hat{X})$ is
 8 known, either completely or partially. Let K be the set of integers, $\{1, 2, \dots, k\}$.
 9 We represent a partial (or total) order on K by means of a subset of or-
 10 dered pairs $(i, j) \in E \subseteq K \times K^1$. An order, $O(\hat{X})$, is *consistent* with E if
 11 $\hat{X}_i \leq \hat{X}_j, \forall (i, j) \in E$. Formally, let X be a set of k values, let v be a set of
 12 corresponding weights, and let E be a partial order. Then monotonic regres-
 13 sion finds a set of values, \hat{X} , consistent with E , that best approximates X in a
 14 weighted least-squares sense. That is, \hat{X} solves the monotonic regression (MR)
 15 problem,

$$\min \sum_{i=1}^k v_i (X_i - \hat{X}_i)^2, \text{ subject to } \hat{X}_i \leq \hat{X}_j, \text{ for all } (i, j) \in E \quad (1)$$

16 The choice of weights is critical for obtaining a meaningful 'best' \hat{X} . In
 17 this respect, we are guided by the property that the solution of Equation (1)
 18 is the maximum likelihood estimate if the observations in each condition are
 19 independent and normally distributed with weights given by the precision of
 20 the data weighted by the number of observations in each condition (Robertson
 21 et al., 1988). That is,

$$\begin{aligned} v_i &= \frac{n_{x_i}}{S_{X_i}^2} \\ w_i &= \frac{n_{y_i}}{S_{Y_i}^2} \end{aligned} \quad (2)$$

¹Unless otherwise stated, a partial order, E , is assumed to be transitively closed.

1 where $S_{X_i}^2$ is the sample variance of variable x in condition i and $S_{Y_i}^2$ is the
2 sample variance of variable y in condition i .

3 In many situations the observations in each condition are not independent,
4 as when conditions are manipulated within participants rather than between.
5 In this case the maximum likelihood estimate depends on the entire covariance
6 matrix and the sets of weights, v_i and w_i , are replaced by appropriate matrices.
7 For this reason, we generalize Equation (2) in the following way. Suppose there
8 are g groups of participants of size n_i , $i = 1, \dots, g$, each measured under m
9 different conditions on variable x . The total number of conditions is thus $k =$
10 gm . Let \mathbf{S}_i be the $m \times m$ covariance matrix for group i . Then the corresponding
11 weight matrix is given by the following block-diagonal matrix,

$$V = \begin{bmatrix} n_1 \mathbf{S}_1^{-1} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & n_g \mathbf{S}_g^{-1} \end{bmatrix} \quad (3)$$

12 The weight matrix, W , for variable y is similarly defined². \mathbf{S}_i^{-1} approximates
13 the inverse of the population covariance matrix, Σ_i^{-1} . A better estimate of
14 Σ_i^{-1} can be obtained by first ‘shrinking’ \mathbf{S}_i , which reduces the unreliable off-
15 diagonal elements but does not necessarily set all of them to zero (Ledoit &
16 Wolf, 2004). We use Ledoit-Wolf method to adjust the weight matrices in our
17 current approach.

18 Let X be a vector of k sample means and let \hat{X} be a vector of values. Then,
19 with the weight matrix V defined by Equation (3), the MR problem is given by,

$$\min \left(X - \hat{X} \right)^T V \left(X - \hat{X} \right), \text{ subject to } \hat{X}_i \leq \hat{X}_j, \text{ for all } (i, j) \in E \quad (4)$$

20 We write the problem corresponding to Equation (4) as $\text{MR}(X, V, E)$ and

²We assume that observations on x and y are themselves independent.

1 the minimum value as $\omega(X, V, E)$, or, in shorthand form, as ω_X . Finding the
 2 solution to the MR problem is not trivial, but fast algorithms have been de-
 3 veloped. If E is a total order then the MR problem can be solved using the
 4 *pool-adjacent-violators algorithm* (PAVA), a version of which was used in the
 5 original development of non-metric multidimensional scaling (Kruskal, 1964).
 6 Otherwise, the problem as posed in Equation (4) can be solved using quadratic
 7 programming algorithms (de Leeuw et al., 2009). The functions *lsqlin* (equiv-
 8 alently, *quadprog*) and *lsei* implement this algorithm in MATLAB[®] and R (R
 9 Core Team, 2013) respectively. In addition, a rapid approximate solution may
 10 also be obtained using the *generalized pool-adjacent-violators* (GPAV) algorithm
 11 developed by Burdakov et al. (2006).

12 *Coupled monotonic regression*

13 Monotonic regression can be extended to incorporate the additional con-
 14 straint that the fitted values on two variables are themselves monotonically
 15 ordered. That is, $O(\hat{X}) = O(\hat{Y})$. This defines the following *coupled monotonic*
 16 *regression* (CMR) problem: Given two sets of values X and Y , corresponding
 17 weight matrices, V and W , and a partial order³, E , we wish to find \hat{X} and \hat{Y} that
 18 are solutions to $\text{MR}(X, V, E)$ and $\text{MR}(Y, W, E)$, respectively, while satisfying
 19 the additional coupled monotonicity constraint,

$$\begin{aligned}
 \hat{X}_i < \hat{X}_j &\Rightarrow \hat{Y}_i \leq \hat{Y}_j \\
 \hat{Y}_i < \hat{Y}_j &\Rightarrow \hat{X}_i \leq \hat{X}_j
 \end{aligned}
 \tag{5}$$

20 This constraint can also be expressed succinctly as follows. If Equation (5) holds
 21 then there is no (i, j) such that,

$$(\hat{X}_i - \hat{X}_j)(\hat{Y}_i - \hat{Y}_j) < 0.
 \tag{6}$$

22 If there is an (i, j) such that Equation (6) is true then the corresponding pair of
 23 points is called *infeasible* and the sets, \hat{X} and \hat{Y} , are called infeasible solutions.

³Note that in the CMR problem, but not in the MR problem, E can be empty.

1 To formalize the CMR problem, we note that for a given partial order, E ,
2 there is a set of all total orders, $\mathcal{L}(E) \supset E$, called the *linear extensions* of E .
3 The CMR problem can then be stated as the problem of finding sets, \hat{X} , \hat{Y} , and
4 \hat{E} , that solve,

$$\begin{aligned} & \min \left[\left(X - \hat{X} \right)^T V \left(X - \hat{X} \right) + \left(Y - \hat{Y} \right)^T W \left(Y - \hat{Y} \right) \right] \\ & \text{subject to, } \hat{X}_i \leq \hat{X}_j, \hat{Y}_i \leq \hat{Y}_j \text{ for all } (i, j) \in \hat{E}, \hat{E} \in \mathcal{L}(E) \end{aligned} \quad (7)$$

5 We write the problem corresponding to Equation (7) as $\text{CMR}(X, Y, V, W, E)$,
6 shorthand $\text{CMR}(E)$, and the minimum value as $\omega(X, Y, V, W, E)$, shorthand
7 ω_{XY} .

8 One way of solving the CMR problem defined by Equation (7) is by direct
9 search. While this is guaranteed to find a global minimum, it can be exception-
10 ally slow, as it requires evaluation of a potentially very large number of total
11 orders. For example, for $k = 10$ and $E = \emptyset$, there are $k! = 3,628,800$ orders to
12 search. To circumvent this problem, Burdakov et al. (2012) recently devised the
13 *CMR algorithm* that finds a solution in approximately exponential rather than
14 factorial time. We briefly describe that algorithm here and provide pseudo-code
15 in the Appendix.

16 The CMR algorithm is a branch-and-bound algorithm that can be viewed
17 as an intelligent search through the set of linear extensions of a specified partial
18 order, E . Given E , which may be empty, it progressively adds additional order
19 constraints until an optimal solution is reached.

20 On each iteration, a new extension, $E' \supset E$, is considered. For this E' ,
21 if the corresponding MR solutions, X' and Y' , are feasible, i.e. they satisfy
22 monotonicity constraint (5), then the fit of these values provides an upper bound
23 on ω_{XY} (improved, if possible, on each iteration). If the sets X' and Y' are
24 infeasible, however, the corresponding fit provides a lower bound on ω_{XY} for
25 any extension $E'' \supset E'$. The algorithm then chooses an infeasible $(i, j) \notin E'$,

1 and branches by generating two new extensions, $E' \cup \{(i, j)\}$ and $E' \cup \{(j, i)\}$.
 2 These extensions inherit the lower bound associated with E' and are added to
 3 the set of those to be further considered (tested). This set forms a queue because
 4 its elements, all extensions of E , are sorted in increasing value of their inherited
 5 lower bounds and the solution, \hat{E} , is guaranteed to be an extension of at least
 6 one member of the queue. In addition, on each iteration, the algorithm generates
 7 a feasible solution based on extending E' in several ways and choosing the one
 8 with the best fit. This fit is used for possible improvement of the currently
 9 available upper bound on ω_{XY} .

10 If the obtained fit for any E' is greater than the current upper bound then
 11 E' , as well as all its extensions, can be eliminated from the search. This leads
 12 to the improvement in performance over direct search. The algorithm continues
 13 branching and eliminating until the queue is empty or if the inherited upper
 14 bound of the first member in the queue is greater than the current best upper
 15 bound. The final upper bound is the fit of the optimal least-squares solution,
 16 ω_{XY} .

17 In a worst case scenario involving uncorrelated variables and $E = \emptyset$, simu-
 18 lations confirm that the CMR algorithm converges to the optimal solution as
 19 a function of $\exp(k)$ rather than $k!$. Even in this case, the relative speed up is
 20 substantial. For example, for $k = 10$, the CMR algorithm evaluates on aver-
 21 age about 25 sub-problems in contrast to a direct search of over three million
 22 sub-problems. In addition, to the extent that the variables are correlated over
 23 conditions and order constraints are specified in E , the algorithm will converge
 24 at an even faster rate.

25 Insert Figure 1 here

Example application of the CMR algorithm

Figure 1 shows a state-trace plot based on results found by Nosofsky et al. (2005) in their Experiment 1. The axes correspond to performance on two different categorization tasks (called “RB” and “II”, respectively). The experimental conditions consisted of a sequence of eight blocks of training trials followed by two blocks of re-training trials that differed between the two groups: a button-switch group who exchanged the position of the response buttons between training and re-training, and a control group who did not. The plot shown in Figure 1 was first generated by Dunn et al. (2012) who used it to discuss whether these data constituted evidence for the existence of more than one latent variable. The first step in answering this question is to solve the CMR problem and determine the fit of the best-fitting monotonically-related set of points. Dunn et al. were unable to solve this problem previously for two reasons. First, they only had direct search method available to them which was unable to find a solution in a practical period of time⁴. Second, the relevant data is a mixture of conditions, one of which was varied within-participants (trial block), the other between-participants (response switch vs. no switch). This requires use of the corresponding weight matrices defined by Equation (3).

Figure 1 also shows the optimal CMR solution, connected by dashed lines to aid visibility. The actual fit value, ω_{XY} , corresponding to the solution of Equation (7), was 3.514. This value depends upon the sample means, X and Y , the weight matrices, V and W , computed according to Equation (3), and the pre-defined partial order, E (empty in this case).

The partial order, E , may be used to specify prior knowledge concerning an expected order of the population means over a sub-set of conditions. In the present case, each group participated in 10 blocks of learning trials with the

⁴On a standard desktop, finding the CMR solution for the current problem by direct search would take approximately 10 hours. In contrast, the CMR algorithm produced the solution in approximately 0.1 seconds.

1 first eight corresponding to successive blocks of training on the same task. It
 2 is reasonable to assume that the population means should not decrease over
 3 these blocks. It may be similarly argued that the population means should not
 4 decrease across the two post-switch blocks, 9 and 10, in each group. Based on
 5 these considerations, it is possible to impose a partial order constraint on the
 6 solution to the CMR problem. Note that within this partial order, although the
 7 first eight blocks and the last two blocks are ordered for each group and task,
 8 there is no constraint on the order of blocks 8 and 9. Indeed, the possibility
 9 of different orders between these conditions on the RB and II variables in the
 10 button-switch group was the main theoretical question posed by Nosofsky et al.

11 If no partial order is specified, the fit value is 3.514 (as noted above). If the
 12 partial order constraint is specified then the fit value cannot decrease, and may
 13 increase. In the present case, the model fit increases slightly to 3.774 suggesting
 14 that the observed means, X and Y , conform closely to the assumed partial
 15 order. One reason for imposing a partial order constraint on the solution is that
 16 it may lead to a more powerful test of the hypothesis of equal orders. In this
 17 case, the test statistic is the difference in fit between a model that assumes only
 18 the partial order constraint and a model that assumes both the partial order
 19 constraint and coupled monotonicity. We discuss this in the next section.

20 *Hypothesis test*

21 While the CMR algorithm allows us to find a value for ω_{XY} , a substantial
 22 problem remains in determining whether this value is large enough to reject the
 23 null hypothesis that the population means have the same order. To do this, we
 24 first define two models of interest. The one-dimensional model (conditional on
 25 E) is defined as follows:

$$M_1 : O(x) = O(y) \quad \& \quad O(x), O(y) \in \mathcal{L}(E)$$

26 This states that the order of the population means on x is the same as the order
 27 on y and that this order is a linear extension of the specified partial order, E .

1 This model is nested within a two-dimensional model (conditional on E) defined
2 as follows:

$$M_2 : O(x), O(y) \in \mathcal{L}(E)$$

3 This states only that the orders on x and y are both (potentially different)
4 linear extensions of the specified partial order, E . Fitting M_2 does not require
5 the CMR algorithm as it consists of two standard MR problems, one in X and
6 one in Y . Further, if $E = \emptyset$, the fit of M_2 is necessarily equal to zero.

7 At present there is no statistical test of the loss in fit from M_2 to M_1 . In
8 the simpler case of (ordinary) monotonic regression, some work has been done
9 on developing a test of the hypothesis, $O(x) \in \mathcal{L}(E)$, against an unconstrained
10 alternative based on the sampling distribution of ω_X . It is known that under
11 this hypothesis, the test statistic follows a $\bar{\chi}^2$ (chi-bar squared) distribution
12 (Robertson et al., 1988). This is a mixture of χ^2 distributions of different
13 degrees of freedom with mixture weights, called level probabilities, which depend
14 in complex ways on the number of conditions, the number of participants, and
15 the partial order, E . As a result, $\bar{\chi}^2$ distributions have been calculated for only
16 a few relatively simple cases. While it may be possible to extend this approach
17 to coupled monotonic regression, we have not attempted this, as it seems likely
18 that calculation of the theoretical distribution would encounter even greater
19 difficulties.

20 Our test of the fit of M_1 against the fit of M_2 is constructed by empirically
21 estimating the sampling distribution of the difference in respective fits under
22 the assumption that M_1 is the true model. The method is adapted from the
23 bootstrap re-sampling procedure described by Wagenmakers et al. (2004). As
24 these authors point out, their procedure cannot be directly applied when the
25 models to be compared are nested. Since M_1 is nested in M_2 , M_2 always fits
26 better than M_1 . For this reason, the fit of M_1 can only be compared against
27 the fit of M_2 . The steps in this procedure are as follows:

1 Let \mathbf{X} and \mathbf{Y} and be two data sets, let X and Y be vectors of the corresponding
 2 sample means, and let V and W be the corresponding weight matrices. Let E
 3 be a specified partial order.

- 4 1. Using the CMR algorithm, find the observed fit of M_1 , $\omega_{XY} =$
 5 $\omega(X, Y, V, W, E)$. Using any suitable MR algorithm, find $\omega_X = \omega(X, V, E)$
 6 and $\omega_Y = \omega(Y, W, E)$, and calculate the observed fit of M_2 , $\omega_{X+Y} =$
 7 $\omega_X + \omega_Y$. If $E = \emptyset$ then $\omega_X = \omega_Y = 0$. Calculate the observed difference
 8 in fits, $\delta = \omega_{XY} - \omega_{X+Y}$.
- 9 2. Generate two non-parametric bootstrap samples, \mathbf{X}' and \mathbf{Y}' , from the
 10 corresponding data sets. This step is undertaken in order to incorpo-
 11 rate sampling error in parameter estimation. Calculate the corresponding
 12 sample means, X' and Y' , and weight matrices, V' and W' .
- 13 3. Solve the CMR problem for the bootstrap samples and, using X', Y', V'
 14 and W' , find the best-fitting values, \hat{X}' and \hat{Y}' .
- 15 4. Transform the original data so that the means are now equal \hat{X}' and \hat{Y}' .
 16 That is, form new samples, $\mathbf{X}_T = \mathbf{X} - X + \hat{X}'$ and $\mathbf{Y}_T = \mathbf{Y} - Y + \hat{Y}'$,
 17 and, from these, draw a second set of non-parametric bootstrap samples,
 18 \mathbf{X}'_T and \mathbf{Y}'_T . Calculate the corresponding sample means, X'_T and Y'_T ,
 19 and weight matrices, V'_T and W'_T , respectively.
- 20 5. Using the CMR algorithm, find the observed fit of M_1 ,
 21 $\omega'_{XY} = \omega(X'_T, Y'_T, V'_T, W'_T, E)$. Using any suitable MR algorithm,
 22 find $\omega'_X = \omega(X'_T, V'_T, E)$ and $\omega'_Y = \omega(Y'_T, W'_T, E)$, and calculate the
 23 observed fit of M_2 , $\omega'_{X+Y} = \omega'_X + \omega'_Y$. Calculate and store the sample
 24 difference in fits (for current iteration i), $\delta'_i = \omega'_{XY} - \omega'_{X+Y}$.
- 25 6. Repeat Steps 2-5 N times where N is a sufficiently large number (e.g.,
 26 10,000).
- 27 7. Calculate, p , the proportion of values of δ'_i that are greater than or equal
 28 to δ . If $p < \alpha$ then reject the null hypothesis.

29 The above procedure can also be adapted to test the fit of M_2 for $E \neq \emptyset$.
 30 In this case, the procedure is modified by replacing M_1 by M_2 and replacing

1 M_2 by the unconstrained model, the fit of which is necessarily zero. The steps
 2 of this procedure are as follows:

3 Let \mathbf{X} and \mathbf{Y} and be two data sets, let X and Y be vectors of the corresponding
 4 sample means, and let V and W be the corresponding weight matrices. Let E
 5 be a specified partial order.

- 6 1. Using any suitable MR algorithm, find $\omega_X = \omega(X, V, E)$ and $\omega_Y =$
 7 $\omega(Y, W, E)$, and calculate the observed fit of M_2 , $\omega_{X+Y} = \omega_X + \omega_Y$.
 8 Calculate the observed difference in fits⁵, $\delta = \omega_{X+Y} - 0$.
- 9 2. Generate two non-parametric bootstrap samples, \mathbf{X}' and \mathbf{Y}' , from the
 10 corresponding data sets. Calculate the corresponding sample means, X'
 11 and Y' , and weight matrices, V' and W' .
- 12 3. Solve the MR problems for each of the bootstrap samples and, using
 13 X', Y', V' and W' , find the best-fitting values, \hat{X}' and \hat{Y}' .
- 14 4. Form new samples, $\mathbf{X}_T = \mathbf{X} - X + \hat{X}'$ and $\mathbf{Y}_T = \mathbf{Y} - Y + \hat{Y}'$, and,
 15 from these, draw a second set of non-parametric bootstrap samples, \mathbf{X}'_T
 16 and \mathbf{Y}'_T . Calculate the corresponding sample means, X'_T and Y'_T , and
 17 associated weight matrices, V'_T and W'_T , respectively.
- 18 5. Using any MR algorithm, find $\omega'_X = \omega(X'_T, V'_T, E)$ and
 19 $\omega'_Y = \omega(Y'_T, W'_T, E)$, and calculate the fit of M_2 , $\omega'_{X+Y} = \omega'_X + \omega'_Y$.
 20 Calculate and store the sample difference in fits, $\delta'_i = \omega'_{X+Y} - 0$.
- 21 6. Repeat Steps 2-5 N times where N is a sufficiently large number (e.g.,
 22 10,000).
- 23 7. Calculate, p , the proportion of values of δ'_i that are greater than or equal
 24 to δ . If $p < \alpha$ then reject the null hypothesis.

25 Each of the hypothesis tests outlined above rely on two principal elements,
 26 the CMR algorithm and bootstrap re-sampling. Because both of these are quite

⁵We include the notional subtraction of zero, the fit of the unconstrained model, to highlight the parallels between the two procedures.

1 general, the procedure can be applied to a wide variety of research designs.
 2 The experimental conditions can be fully randomized across participants, ap-
 3 plied entirely within-participants, or any combination of between- and within-
 4 participant treatments. The procedure may also be adapted for discrete data
 5 although, in this case, the model-consistent data, \mathbf{X}_T and \mathbf{Y}_T , are derived from
 6 a parametric bootstrap of the observed data (in step 3 above). However, this is
 7 not a substantial concern as this relevant distribution is entirely specified by the
 8 data so parametric and non-parametric re-sampling are equivalent. We discuss
 9 the application of the method to binomial data in a later section.

10 Insert Figure 2 here

11 *Example application*

12 To illustrate the application of the hypothesis testing procedure, we return
 13 to the state-trace plot shown in Figure 1. Figure 2 shows two empirical dis-
 14 tributions of δ' , each based on 10,000 iterations, and two observed values of δ .
 15 The dashed line and unfilled triangle are based on the assumption of no partial
 16 order, $E = \emptyset$. In this case, the two-dimensional model fits perfectly (as it is un-
 17 constrained) and δ is equal to the observed fit of the one-dimensional model and
 18 has the value of 3.514 (as noted earlier). The corresponding empirical p -value is
 19 0.77 from which it is concluded (for $\alpha = .05$) that the hypothesis $O(x) = O(y)$
 20 cannot be rejected.

21 If the partial order, E , described earlier in relation to the data shown in
 22 Figure 1, is implemented then the testing procedure differs. The first step is to
 23 test the fit of M_2 which has an observed fit of 0.929. The corresponding empirical
 24 p -value is 0.72 from which it is concluded that the hypothesis $O(x), O(y) \in \mathcal{L}(E)$
 25 cannot be rejected. Following this, the next step is to test the difference in fit
 26 between M_1 and M_2 . The solid line in Figure 2 shows the estimated empirical
 27 distribution of δ' and the filled triangle shows the observed value of δ . As

1 stated earlier, the observed fit of the one-dimensional model (M_1) is fractionally
2 increased to 3.774. However, the value of δ is now $3.774 - 0.929 = 2.845$, and the
3 corresponding empirical p -value is 0.57. We again conclude that the hypothesis,
4 $O(x) = O(y)$, given $O(x), O(y) \in \mathcal{L}(E)$, cannot be rejected.

5 Although in this case, both analyses (with and without assuming a prior
6 partial order) lead to the same conclusion, inspection of Figure 2 illustrates the
7 increase in statistical power that may result from the addition of an appropriate
8 partial order constraint. Although δ has decreased from the no-partial-order to
9 the partial-order case, this difference is relatively small compared to the differ-
10 ence in the shapes of the corresponding empirical distributions. Specifically, the
11 distribution of δ' , when the partial order is specified (filled curve), is contracted
12 leftwards compared to the sampling distribution of δ' , when no partial order is
13 specified (dashed curve). As a result, relatively less mass falls to the right of
14 the observed value of δ leading to a lower p -value and an associated increase in
15 statistical analysis. The reason for this is that, if the population means satisfy
16 the partial order constraint, E , then the fit of M_2 will be close to zero. However,
17 many of the bootstrap samples of M_1 may violate the partial order in which
18 case the fit of M_2 will be substantially greater than zero, thereby contracting
19 the distribution of δ' .

20 Analyzing power

21 It is desirable that our proposed test have sufficient power to reject the null
22 hypothesis of equal latent orders when it is false. We address this issue in the
23 present section where our goals are; (1) to define a measure of effect size in
24 post-hoc power analyses, (2) to show how power can be estimated for any given
25 effect size, (3) to discuss the problem of estimating effect size for proactive power
26 analyses, and finally (4) to demonstrate the effect on power of imposing partial
27 order constraints.

1 We consider in detail a measure of effect size for a fully-randomized between-
 2 participant experiment with n participants in each of k conditions. In this case
 3 the true effect size is the fit, ω_{xy} , of the solution to the following CMR problem:

$$\begin{aligned} \omega_{xy} = \min & \left[(x - \hat{x})^T \Upsilon (x - \hat{x}) + (y - \hat{y})^T \Psi (y - \hat{y}) \right] \\ \text{subject to, } & (\hat{x}_i - \hat{x}_j) (\hat{y}_i - \hat{y}_j) \geq 0, \text{ for all } (i, j) \end{aligned} \quad (8)$$

4 where, $\Upsilon = \text{diag}(\sigma_{x_1}^2, \dots, \sigma_{x_k}^2)^{-1}$, $\Psi = \text{diag}(\sigma_{y_1}^2, \dots, \sigma_{y_k}^2)^{-1}$.

5 For convenience we set $\sigma_{x_i}^2 = \sigma_{y_i}^2 = \sigma^2$ for all i . Both the number of
 6 violations of monotonicity and the size (relative to the population precision) of
 7 each violation determine the value of ω_{xy} , so in order to explore the power of
 8 the CMR test we varied both of these over a wide range. We adopted the case
 9 where $k = 8$, and set $x = \{1, \dots, 8\}$. We manipulated the number of violations
 10 of monotonicity from 1 to 28 by choosing y as a permutation of $\{1, \dots, 8\}$ to
 11 produce the desired number of violations. For each number of violations, we then
 12 varied σ^2 in order to generate effect sizes ranging from 0.1 to 10. This process
 13 resulted in a set of 398 combinations of means, variances, and associated effect
 14 sizes which were used to estimate power for various sample sizes.

15  Insert Figure 3 here

16 For each of these 398 sets of population parameters we drew a sample data
 17 set consisting of k independent, normally distributed, samples, each of size n ,
 18 for $n = \{10, 20, 30, 40, 50\}$ for each variable, x and y . For each data set, we
 19 followed the 7-step procedure presented earlier to determine whether M_1 could
 20 be rejected for two levels controlling the Type I error rate, $\alpha = .05$ and $\alpha = .01$.
 21 Because each data set was drawn from a population in which the monotonic
 22 component of M_1 is false, the observed proportion of correct rejections is an
 23 estimate of the power, $(1 - \beta)$, of the test. The results of these simulations are
 24 shown in Figure 3. Each power curve corresponds to the best fitting logistic
 25 function of the effect size, ω_{xy} , for each value of n .

1 The curves shown in Figure 3 can be used to estimate the number of partic-
 2 ipants that an experimenter may need in order to achieve a given level of power
 3 for the fully randomized design considered above. To do so, it is necessary to
 4 estimate ω_{xy} . For the present equal- n design, an obvious estimates of that ω_{xy} is
 5 given by ω_{XY}/n . For designs with unequal n between groups, the corresponding
 6 estimate is ω_{XY}/\bar{n} , where \bar{n} is the mean n over groups⁶ These curves allow a
 7 researcher to make a rough claim about the scale of the observed effect size. In
 8 the case of Cohen's (1988) d , the scale relates to power as follows: a small effect
 9 has a power of .1 with $n = 20$, medium has a power of .2, and large about .4.
 10 For the CMR test, this corresponds to δ of 0.1, 0.2, 0.4 as power is nearly linear
 11 at that level with $n = 20$. A very large effect (power .8) would be $\delta = 0.50$.

12 *Power under partial order constraints*

13 In this section, we re-examine the potential increase in power due to the
 14 addition of a partial order constraint. Because there are a very large number
 15 of possible partial order constraints, we focus on one that naturally arises in a
 16 factorial design. Consider an experiment with two between-participant factors,
 17 A and B , such that A has two levels and B has 4 levels (i.e., $k = 8$). A prior
 18 belief may exist concerning the orders of the dependent variables on each factor.
 19 We suppose that for each level of B , level 1 of A will produce smaller values on
 20 both dependent variables (e.g., less accurate responding, lower response times)
 21 than will level 2. We further suppose that for each level of A , the levels of B will
 22 conform to a particular total order. By way of an example, an experiment may
 23 examine the effect on recognition memory of a change in the format of visually
 24 presented words and study duration. In this case, factor A is presentation format
 25 (two levels: same format at study and test, different formats) and B is study
 26 duration (4 levels: say, 0.25 sec, 0.5 sec, 1 sec, 2 sec). Based on prior knowledge,
 27 it is plausible to assume that memory for words presented in the same format

⁶Although an obvious approach, it is likely that reliance on ω_{XY} may underestimate ω_{xy} .
 Further research on this question is required.

1 will be no worse than memory for words presented in different formats, while,
2 for each format, memory will not decrease with increasing study duration. In
3 order to illustrate the effect of this prior partial order on statistical power, we
4 simulated the case of $n = 10$ for each group, using the procedure described
5 previously.

6 Insert Figure 4 here

7 Figure 4 reveals the gain in power that results from imposing the proposed
8 partial order. The addition of this constraint leads to a nearly five-fold increase
9 in the rate of increase of the power curve compared to the no-partial order
10 case. The relevant measure of effect size when there is a partial order is the
11 difference between M_1 and M_2 , δ_{xy} . In order to achieve power equal to 0.8
12 at $\alpha = .05$, we found that the observed effect size in the partial order case was
13 $\delta_{xy} = 0.13$, a value substantially less than the observed effect size in the no-
14 partial order case, $\delta_{xy} = 0.78$. The corresponding population variances were
15 0.18 and 0.03, respectively. In order to give some sense of how this may appear
16 in the data, we drew a random data set from populations with each of these
17 variances and summarized these in the state-trace plots shown in Figure 5.
18 The larger variance in the partial-order case is striking. In our experience,
19 measurements with variability of this magnitude are not difficult to find in
20 psychological experiments.

21 Insert Figure 5 here

22 As noted earlier, the imposition of a partial order reduces the variance of
23 the distribution of δ , the difference in fit between M_1 and M_2 , as long as the
24 population means conform to the partial order. On the other hand, if the
25 population means do not conform to the partial order then both M_1 and M_2 are
26 false. Because power is necessarily limited, Type II errors are always possible.

1 The test of the partial order model, M_2 , is at best a check that the experiment
2 has been correctly designed. Furthermore, a partial order should not be adopted
3 merely to facilitate rejection of M_1 . In order to be logically coherent, any partial
4 order should be defined prior to conducting the experiment and be based on a
5 compelling and universally accepted motivation.

6 The power analysis presented above is useful for post-hoc analyses, where
7 the effect size can be estimated from data. However, its use in prospective
8 power estimation is limited because the estimate of the effect size depends on
9 the particular design. For example, in the previous simulations, we assumed
10 a uniform spacing of x and y which may be unlikely to occur in practice. In
11 the context of state-trace analysis, the optimal design is one which maximizes
12 δ_{xy} given a particular two-dimensional manifold of possible latent means in the
13 state space. This, in turn, will depend upon the configuration of latent means
14 selected from the manifold through selection of the experimental factors and the
15 number and nature of their levels. Similarly, repeated measures will affect power
16 in ways that are dependent on the particulars of the variance-covariance matrix.
17 A prospective power analysis will thus require the experimenter to essentially
18 replicate a sub-set of our procedure for the design under consideration.

19 **Control of Type I error**

20 Our method is based on bootstrap resampling. An advantage of this ap-
21 proach is that no assumption is required concerning the nature of the distribu-
22 tion of observations⁷. However, bootstrap samples may underestimate variance
23 for small n (Chernick, 2007) which can lead to a corresponding inflation of the
24 Type I error rate. For this reason we conducted a series of simulations in which
25 we replaced the bootstrap samples with samples from the known distribution

⁷Although, of course, if the data are not normally distributed the obtained values of ω and δ will not be maximum likelihood estimates.

1 from which the data were drawn (in this case, a normal distribution). In each
 2 simulation, the population means were monotonically related; they were, for
 3 each variable, simply the integers 1 to 8, and no partial order was assumed.
 4 We manipulated the variance of each distribution and the sample size, both of
 5 which were assumed to be constant over conditions and variables. On each sim-
 6 ulation, for a given variance and sample size, a sample data set was drawn and
 7 the CMR procedure applied to generate an empirical distribution of fits (based
 8 on 10,000 samples). The procedure was applied both in its bootstrap form (as
 9 described earlier) and in a form in which the bootstrap step was replaced by
 10 re-sampling from the normal distributions used to generate the data. We then
 11 used the latter, parametric, empirical distribution to identify cut-offs for dif-
 12 ferent percentiles including the 95th and 99th percentiles corresponding to $\alpha =$
 13 0.05 and $\alpha = 0.01$, respectively. We then calculated the proportion of cases that
 14 exceeded these cut-offs in the empirical distribution derived from the bootstrap
 15 method. So long as resampling did not produce degenerate cases (which did not
 16 occur with $n > 8$ in our simulations) the percent of the cases that exceeded the
 17 cut-off deviated very little from the expected proportions.

18 **Extension of the CMR procedure to binomial data**

19 In this section, we describe how the CMR procedure can be extended to bi-
 20 nomial data structures. We also take the opportunity to compare this procedure
 21 to the Bayesian model selection approach developed by Prince et al. (2012a),
 22 highlighting their similarities and differences.

23 Some notations are introduced first. Let n_x be a (column) k -vector of the
 24 number of Bernoulli trials for variable x on each of k conditions. Let a_x be the
 25 (column) k -vector of the number of successes in each condition and let b_x be
 26 the corresponding vector of the number of failures, where $n_x = a_x + b_x$. Let
 27 X be the vector of the observed mean proportion of successes for variable x
 28 across k conditions, i.e. $X = a_x/n_x$, where the division is understood to be

1 element-wise. The same kind of notation can be introduced for variable y . We
2 seek to solve the CMR problem given by Equation (7).

With $V = \text{diag}(n_x)$ and $W = \text{diag}(n_y)$, the least-squares solution to the problem given by Equation (7) is also the maximum likelihood solution. This follows from Theorem 12 of Robertson et al. (1988, p. 32) which states that the solution, \hat{X} , to the least-squares monotonic regression on X with weights, n_x , is also the maximum likelihood solution. Because the solution to Equation (7) is the sum of two monotonic regression problems for some \hat{E} , it follows that it is also the maximum likelihood solution. The only difference in applying it to binomial data is that evaluation of sub-problems in the CMR algorithm is based on the actual likelihood function rather than evaluation of Equation (7). Equivalently, it can be based on the following negative log-likelihood function:

$$f(\hat{X}, \hat{Y}) = -(a_x^T \ln(\hat{X}) + b_x^T \ln(1 - \hat{X}) + a_y^T \ln(\hat{Y}) + b_y^T \ln(1 - \hat{Y}))$$

Because the value of this function is non-zero when the fit is perfect, it is convenient to subtract the corresponding value of the perfect fit, $f(X, Y)$. This leads to an equivalent formulation in terms of the G^2 -statistic:

$$G^2 = 2[f(\hat{X}, \hat{Y}) - f(X, Y)]$$

3 *Application to binomial data*

4 Prince et al. (2012b) analyzed a set of binomial data using the Bayesian
5 model selection procedure described by Prince et al. (2012a). These data were
6 obtained from a two-alternative forced-choice recognition memory experiment
7 that investigated the face-inversion effect, based on a similar study by Loftus
8 et al. (2004). The stimuli were pictures of faces or houses which defined the
9 dependent variables of interest (i.e., memory accuracy for faces and memory ac-
10 curacy for houses). Performance was tested under the orthogonal combination
11 of two factors; stimulus orientation (upright vs. inverted), and study duration

1 (short, medium, and long). All experimental factors (stimulus type, orienta-
2 tion, and duration) were manipulated within-participants. The data for each
3 participant (as well as data aggregated over participants) consists of counts of
4 successes (i.e., selecting the correct test item) and counts of failures (i.e., se-
5 lecting the incorrect test item) for each stimulus type under each of the six
6 experimental conditions⁸.

7 The three different study durations imply a partial order on performance.
8 Namely, the proportion of successes should not decrease from short to medium
9 and from medium to long durations for both upright and inverted presentation
10 formats for both face recognition and house recognition. For consistency with
11 Prince et al. we did not place a partial order on the upright and inverted
12 conditions, although this could readily be included.

13 Insert Figure 6 here

14 Figure 6 shows the state-trace plot based on the mean proportion of successes
15 averaged over all participants. The dashed line shows the best fitting monotonic
16 curve. It is clear that for each dependent variable, the effect of study duration
17 is consistent with the assumed partial order. These data may be analyzed in
18 three different ways using CMR. First, the mean scores of proportion correct
19 (corresponding to the points plotted in Figure 6) can be analyzed using the
20 original CMR procedure described earlier, assuming a normal distribution of
21 means across participants. In this case, the empirical p -value based on 10,000
22 iterations is 0.044, which implies rejection of the monotonic model, M_1 , at
23 $\alpha = 0.05$. Second, the counts of successes and failures can be aggregated over
24 participants and these data analyzed using the binomial CMR procedure. In this
25 case, the empirical p -value of δ based on 10,000 iterations is 0.017, also implying
26 rejection of M_1 . However, as Prince et al. have pointed out, aggregation over

⁸The authors are grateful to Melissa Prince and colleagues for making these data available.

1 participants has the potential to distort the underlying pattern of the data.
 2 For this reason, they analyzed each participant separately, which leads to the
 3 third way in which the data can be analyzed using binomial CMR. In this case,
 4 consistent with the analysis of the aggregated data, none of the p -values for
 5 M_2 were significant (minimum $p = 0.079$). On the other hand, none of the
 6 p -values for M_1 against M_2 reached significance (minimum $p = .062$). This
 7 is to be expected given the low power associated with the smaller number of
 8 observations for each participant. Given this, it is desirable to combine this
 9 evidence in a manner that does not lead to distortions due to averaging (Davis-
 10 Stober et al., In Press). This can be done by conducting a test of the *sum* of the
 11 individual fits. Such a test is equivalent to using the binomial CMR procedure
 12 to fit M_1 and M_2 to a concatenated set of kn conditions with a partial order
 13 constraint and a monotonicity constraint applied to each set of k conditions for
 14 each of the n participants. In practice, the relevant statistics can be obtained
 15 from the individual analyses already conducted - the sum of the model fits across
 16 participant is compared against the distribution of the sum of random samples
 17 drawn from the individual empirical distributions obtained from the bootstrap
 18 procedure. Consistent with the aggregated data which exactly conform to the
 19 partial order constraint, the combined p -value for the test of M_2 is not significant
 20 ($p = 0.817$). However, the combined p -value for the test of M_1 against M_2 fell
 21 short of significance ($p = 0.084$)⁹.

22 Insert Figure 7 here

23 *Comparison to Bayesian model selection approach*

24 In order to compare the results of the binomial CMR procedure with the
 25 Bayesian model selection developed by Prince, et al. (2012), is necessary to

⁹Based on 100,000 combined samples each corresponding to the sum of 18 individual random samples from the individual empirical distributions.

1 explain their approach in some detail and to identify the points of similarity
 2 and difference with the CMR approach. Figure 7 summarizes the main features
 3 of the two approaches. The left hand side of Figure 7 shows a binary tree
 4 generated by the sequential addition of order constraints. The top-most model
 5 is the unconstrained model (called the *encompassing model* by Prince et al.),
 6 which, by definition, fits the observed data perfectly. The second level contrasts
 7 two models defined by the addition of the partial order constraint, $O(x), O(y) \in$
 8 $\mathcal{L}(E)$, where $\mathcal{L}(E)$ is the set of linear extensions of the specified partial order, E .
 9 The model for which this constraint is true is called the *trace model* by Prince
 10 et al., and the model for which it is false is called the *non-trace model*. The
 11 Bayesian procedure directly compares these models and selects the one with
 12 the greater posterior model probability. In contrast, the CMR procedure tests
 13 if the addition of the partial order constraint leads to a statistically significant
 14 decrease in goodness of fit. Following the Bayesian procedure, if the trace model
 15 is selected¹⁰ then two additional models are contrasted at the third level, defined
 16 by the addition of the monotonicity constraint, $O(x) = O(y)$. The model for
 17 which this constraint is true is called the *monotonic model* by Prince et al.,
 18 and the model for which it is false is called the *multidimensional model*. Again,
 19 the Bayesian procedure directly compares these two models while the CMR
 20 procedure tests the loss of fit caused by the additional monotonicity constraint.

21 Finally, Prince et al. proposed a binary contrast at a fourth level, between
 22 two complementary models called the *overlap* and *non-overlap models*. In the
 23 experimental design used by Prince et al., non-overlap means that the effect
 24 of stimulus orientation (upright vs. inverted) is sufficiently large that there is
 25 no overlap between the sets of data points corresponding to the three stimulus
 26 durations. If this occurs, the resulting state-trace is trivially monotonic and
 27 Prince et al. advised that the experiment should be re-designed. Let $\mathcal{L}'(E)$ and

¹⁰Prince et al. describe both a sequential and simultaneous model evaluation procedure. We describe the sequential approach for expository purposes.

1 $\mathcal{L}''(E)$ be a partition of $\mathcal{L}(E)$ such that $\mathcal{L}'(E)$ is the set of linear extensions of
2 E consistent with overlap and $\mathcal{L}''(E)$ is the set of extensions inconsistent with
3 overlap. The final constraint is therefore that, $O(x), O(y) \in \mathcal{L}'(E)$.

4 An apparent advantage of the Bayesian procedure is that it allows the weight
5 of evidence for pairs of disjoint models at each level of constraint to be directly
6 compared. In contrast, a null hypothesis statistical test, which forms the heart of
7 our procedure, tests whether the addition of a constraint leads to a statistically
8 significant loss of fit. Offsetting this advantage is the necessity of assuming a
9 prior distribution over the set of all possible orders of conditions. Depending on
10 the context, different priors are possible and each choice will lead to a different
11 outcome in model selection. Prince et al. assumed that this prior is uniform.

12 Analogous to the combined p -value, Prince et al. calculated a group poste-
13 rior model probability based on combined Bayes factors, essentially the product
14 of individual Bayes factors, and found the probability of the trace model com-
15 pared to the non-trace model was greater than 0.95. This is analogous to our
16 test of M_2 (against the unconstrained model) which had a combined p -value of
17 0.85. Consistent with this, the rank order of individual participants' posterior
18 probabilities of the non-trace model is similar (but not identical) to the rank
19 order of the individual fits of M_2 , Kendall's tau = 0.73, $p < 0.0001$. Prince
20 et al. also found that the group posterior model probability of the monotonic
21 model compared to the multidimensional model was less than 0.05. In contrast,
22 our analogous test of M_1 against M_2 had a combined p -value of 0.070 which
23 fell short of significance ($\alpha = 0.05$). However, the rank order of participants'
24 posterior probabilities for the multidimensional model is similar (but not iden-
25 tical) to the rank order of the difference in fit between M_1 and M_2 , Kendall's
26 tau = 0.42, $p = 0.007$. Thus, while the two methods are based on different
27 theoretical orientations and procedures, and technically test different models,
28 their commonalities are such that they may well lead to similar conclusions.

29 Unlike Prince et al., we do not incorporate a test of overlap into our pro-

1 cedure. We have not pursued this option for three reasons. First, it is not
 2 essential to the principal question of testing the model of equal orders. Second,
 3 the concept appears to be most relevant to the kind of factorial design investi-
 4 gated by Prince et al. It is not clear how it might be relevant to other designs,
 5 such as that used by Nosofsky et al. (2005). Finally, it is not clear that the
 6 concept of non-overlap is sufficiently inclusive. Given a set of populations that
 7 have different orders (i.e., where M_1 is false), there are many configurations of
 8 sample means that will be trivially monotonically ordered¹¹. Non-overlap is but
 9 one example. In our view, the failure to reject M_1 requires further analyses of
 10 the data to determine whether this is due to the configuration of sample means.
 11 Such follow-up analysis is analogous to inspection of the scatterplot to aid in-
 12 terpretation of a correlation coefficient. If the data are trivially monotonic, the
 13 pattern of points will suggest possible changes to the levels of the experimental
 14 factors to increase the chance of rejecting M_1 (assuming it is false). Prince et al.
 15 made similar recommendations and suggested that, in attempting to maximize
 16 power, it may be useful to adopt non-standard factorial designs.

17 We endorse consideration of non-standard factorial designs. In such designs,
 18 the levels of one factor may differ across levels of the other factor. For example,
 19 in the face-inversion study conducted by Prince et al., stimulus durations for the
 20 more difficult inverted condition may be longer than corresponding durations
 21 for the easier upright condition. Such choices maximize the chance that some
 22 pairs of points in the state-trace plot will violate monotonicity. It must be
 23 remembered that even if the underlying state-trace is two dimensional (with
 24 unequal latent orders), this will only be revealed in the observed data if the
 25 configuration of points contains violations of monotonicity. This, in turn, will
 26 depend in complex ways on the levels of the factors that have been manipulated.
 27 Depending on these levels, violations may or may not be observed.

¹¹For the design used by Prince et al., other examples include the lack of an effect of either
 or both experimental factors, or a ‘staircase’ arrangement of points in the state space which
 suggest two-dimensionality but fail to produce any violations of monotonicity.

1 Conclusion

2 We have presented a comprehensive procedure for testing for the equality
3 of latent orders. The procedure consists of two main parts: (1) The CMR
4 algorithm that finds the best single order on two dependent variables over k
5 conditions and returns a measure of the lack-of-fit of that order to the data; (2)
6 a significance test for this lack-of-fit, based on bootstrap resampling. Consis-
7 tent with experience of the bootstrap (Chernick, 2007), we showed that this test
8 controls Type I error rate for sample sizes greater than eight. We also showed
9 that the power of the test was a function of effect size and sample size for a
10 fully randomized, equal n , design and that it obtained reasonably high levels of
11 power (> 0.80) for data that could plausibly occur in typical psychology exper-
12 iments. We also demonstrated the role of partial orders, or pre-experimental
13 order constraints on conditions, in substantially increasing power in the case
14 where the partial order is true.

15 Although we presented the CMR procedure principally in relation to con-
16 tinuous data, we showed how it can be readily extended to discrete data and
17 discussed the binomial case in some detail. A feature of the procedure for con-
18 tinuous data is that it permits a non-parametric bootstrap. Thus, it is not nec-
19 essary to make any distributional assumptions. Nor is it necessary to assume
20 equal variances or equal n , at least in a fully randomized design, as unequal
21 precisions are explicitly built into the monotonic regression weights.

22 No discussion of hypothesis testing should ignore the crucial differences be-
23 tween Bayesian and frequentist approaches. Our bootstrap method provides a
24 frequentist estimate of the variability of the CMR fit estimate. It should be pos-
25 sible to construct an alternative Bayesian approach to examining latent orders
26 using CMR, and Bayesian hypothesis tests for state-trace applications of latent
27 order testing without CMR already exist (Davis-Stober et al., In Press; Prince
28 et al., 2012b). One critical feature that divides the Bayesian and frequentist

1 approaches is the treatment of model complexity. The equal-order model is less
 2 complex than the alternative, where each variable follows its own (partial) order.
 3 Our frequentist approach does not penalize the separate-order model, because
 4 its complexity is unknown. Because the common-order model is nested within
 5 the separate-order model, the latter will always fit better than the former. We
 6 recommend rejection of the common-order model when the probability of the
 7 fit being as bad as is observed is small. The Bayesian approach does penalize
 8 for complexity, by specifying priors for both models. The separate-order model
 9 will have a more diffuse prior than the common-order model, making it possible
 10 to compare the models to each other and accept either one. This bi-directional
 11 decision is enabled only by making specific assumptions about what the appro-
 12 priate prior should be for both models. Such priors equate to theories about
 13 the data generating processes. On the one hand, such theories are critical to
 14 advancing our understanding of the process that give rise to observed data. On
 15 the other hand, disagreement about what theories are reasonable will necessarily
 16 extend to the results of Bayesian hypothesis testing. We have argued that there
 17 is a role for a procedure that makes minimal assumptions about the distribution
 18 of latent orders, and we believe that our NHST approach is informative within
 19 that context.

20 We motivated the development of the CMR procedure by reference to its
 21 relevance to state-trace analysis where the presence of different latent orders
 22 implies that the dependent variables are functions of more than one latent vari-
 23 able. For this reason, we discussed the application of the CMR procedure to
 24 two dependent variables, as commonly used in STA. However, the procedure
 25 can also be readily generalized to test the equality of latent orders over any
 26 number of dependent variables.

27 A further, intriguing, challenge is to consider the more complex case in which
 28 the latent orders of N dependent variables conform to a linear space of $d < N$
 29 dimensions (Dunn & James, 2003). For $N = 2$ dependent variables, equal latent

1 orders implies that $d = 1$. For $N > 2$ and $d > 1$, different constraints will apply
2 to generate sets of permitted N -tuples of orders. While this problem poses a
3 number of significant difficulties, its solution would lead to a general test of
4 latent orders beyond simple equality.

5 Bamber, D. (1979). State-trace analysis: A method of testing simple theories
6 of causation. *Journal of Mathematical Psychology*, 19, 137–181.

7 Burdakov, O. P., Dunn, J. C., & Kalish, M. L. (2012). An approach to solving
8 decomposable optimization problems with coupling constraints.

9 Burdakov, O. P., Sysoev, O., Grimvall, A., & Hussian, M. (2006). An $O(n^2)$
10 algorithm for isotonic regression. In G. di Pillo, & M. Roma (Eds.), *Large-*
11 *Scale Nonlinear Optimization* (pp. 25–33). New York: Springer volume 83 of
12 *Nonconvex Optimization and Its Applications*.

13 Chernick, M. R. (2007). *Bootstrap methods: A guide for practitioners and*
14 *researchers*. New York: Wiley.

15 Cohen, J. (1988). *Statistical Power Analysis for the Behavioral Sciences (2nd*
16 *Edition)*. Routledge.

17 Davis-Stober, C., Morey, R. D., Gretton, M., & Heathcote, A. (In Press). Bayes
18 factors for state-trace analysis. *Journal of Mathematical Psychology*, .

19 Dunn, J. C., & James, R. N. (2003). Signed difference analysis: Theory and
20 application. *Journal of Mathematical Psychology*, 47, 389–416.

21 Dunn, J. C., Newell, B. R., & Kalish, M. L. (2012). The effect of feedback
22 delay and feedback type on perceptual category learning: The limits of mul-
23 tiple systems. *Journal of Experimental Psychology: Learning, Memory and*
24 *Cognition*, 38, 840–859.

25 Kruskal, J. (1964). Nonmetric multidimensional scaling: A numerical method.
26 *Psychometrika*, 29, 115–129.

- 1 Ledoit, O., & Wolf, M. (2004). Honey, i shrunk the sample covariance matrix.
2 *The Journal of Portfolio Management*, 30, 110–119.
- 3 de Leeuw, J., Hornik, K., & Mair, P. (2009). Isotone optimization in r: Pool-
4 adjacent-violators algorithm (pava) and active set methods. *Journal of Sta-*
5 *tistical Software*, 32, 1–24.
- 6 Loftus, G. R., & Masson, M. E. J. (1994). Using confidence intervals in within-
7 subject designs. *Psychonomic Bulletin & Review*, 1, 476–490.
- 8 Loftus, G. R., Oberg, M. A., & Dillon, A. M. (2004). Linear theory, dimensional
9 theory, and the face-inversion effect. *Psychological Review*, 111, 835–863.
- 10 Newell, B. R., & Dunn, J. C. (2008). Dimensions in data: Testing psychological
11 models using state-trace analysis. *Trends in Cognitive Sciences*, 12, 285–290.
- 12 Nosofsky, R. M., Stanton, R. D., & Zaki, S. R. (2005). Procedural interference in
13 perceptual classification: Implicit learning or cognitive complexity? *Memory*
14 *& Cognition*, 33, 1256–1271.
- 15 Open Science Collaboration (2015). Estimating the reproducibility of psy-
16 chological science. *Science*, 349. URL: [http://www.sciencemag.org/](http://www.sciencemag.org/content/349/6251/aac4716.abstract)
17 [content/349/6251/aac4716.abstract](http://www.sciencemag.org/content/349/6251/aac4716.abstract). doi:10.1126/science.aac4716.
18 [arXiv:http://www.sciencemag.org/content/349/6251/aac4716.full.pdf](http://www.sciencemag.org/content/349/6251/aac4716.full.pdf).
- 19 Pratte, M. S., & Rouder, J. N. (2012). Assessing the dissociability of recollection
20 and familiarity in recognition memory. *Journal of Experimental Psychology:*
21 *Learning, Memory and Cognition*, 38, 1591–1607.
- 22 Prince, M., Brown, S., & Heathcote, A. (2012a). The design and analysis of
23 state-trace experiments. *Psychological Methods*, 17.
- 24 Prince, M., Hawkins, G., Love, J., & Heathcote, A. (2012b). An r package for
25 state-trace analysis. *Behavior Research Methods*, 44, 644–655.

- 1 R Core Team (2013). *R: A Language and Environment for Statistical Com-*
2 *puting*. R Foundation for Statistical Computing Vienna, Austria. URL:
3 <http://www.R-project.org/>.
- 4 Robertson, T., Wright, F. T., & Dykstra, R. L. (1988). *Order restricted statis-*
5 *tical inference*. Chichester, UK: John Wiley & Sons.
- 6 Wagenmakers, E.-J. (2007). A practical solution to the pervasive problems of p
7 values. *Psychonomic Bulletin & Review*, 14, 779–804.
- 8 Wagenmakers, E.-J., Ratcliff, R., Gomez, P., & Iverson, G. J. (2004). Assessing
9 model mimicry using the parametric bootstrap. *Journal of Mathematical*
10 *Psychology*, 48, 28–50.

1 Appendix

2 *CMR algorithm*

3 The following pseudo-code describes the CMR algorithm (Burdakov et al., 2012).
 4 Here, X and Y vectors of means, V and W are corresponding weight matrices,
 5 E is a specified partial order partial order and $F(\hat{X}, \hat{Y})$ is the objective function
 6 value in (3) computed for the vectors \hat{X} and \hat{Y} . L is a list of pairs of the form
 7 (e, f) where e is a partial order and f is the value of the corresponding inherited
 8 lower bound.

9 Input: X, Y, V, W, E . Output: $\hat{X}, \hat{Y}, F(\hat{X}, \hat{Y})$.
 10 $L = \{(E, -\infty)\}$, $F_U = \infty$, $F_L = -\infty$
 11 **while** $(|L| > 0)$ & $(F_L < F_U)$ **do**
 12 $(E', F_L) \leftarrow L(1)$
 13 **if** $F_L < F_U$ **then**
 14 find X' that solves $\text{MR}(X, v, E')$ and Y' that solves $\text{MR}(Y, w, E')$
 15 compute $F(X', Y')$
 16 **if** $F(X', Y') < F_U$ **then**
 17 **if** (X', Y') is feasible **then**
 18 $F_U \leftarrow F(X', Y')$, $(\hat{X}, \hat{Y}) \leftarrow (X', Y')$
 19 **else**
 20 generate feasible solution (X'', Y'') and compute $F(X'', Y'')$
 21 **if** $F(X'', Y'') < F_U$ **then**
 22 $F_U \leftarrow F(X'', Y'')$, $(\hat{X}, \hat{Y}) \leftarrow (X'', Y'')$
 23 **end**
 24 find (i, j) such that $(X'_i - X'_j)(Y'_i - Y'_j) < 0$
 25 $E'_{ij} \leftarrow E' \cup \{(i, j)\}$, $E'_{ji} \leftarrow E' \cup \{(j, i)\}$
 26 append $(E'_{ij}, F(X', Y'))$ and $(E'_{ji}, F(X', Y'))$ to L
 27 reorder $L = \{\dots, (e, f), \dots\}$ in increasing values of f
 28 **end**
 29 **end**


```
1     end
2 end
```

3 **Acknowledgements**

4 The authors wish to thank Ben Newell, EJ Wagenmakers, Andrew Heath-
1 cote, John Kruschke, Laura Anderson, and Don Bamber for their helpful dis-
2 cussions on many related topics. The authors gratefully acknowledge the con-
3 tinued support of the Australian Research Council (Discovery Grants: 0877510,
4 0878630, 110100751, and 130101535), the National Science Foundation (Award
5 1256959) to Kalish, and two visiting fellowships from Linköping University to
6 Dunn.

7 Figure Captions

- 8 1. Data from Nosofsky, Stanton, and Zaki (2005, Experiment 1). State-trace
1 plot of mean proportion correct on RB and II category structures for each
2 block of trials in the learning or pre-switch phase (Blocks 1-8 only) and in
3 the post-switch or transfer phase (final two blocks for each group). In the
4 control condition, the same response assignment was maintained across the
5 two phases. In the button switch condition, the response assignment was
6 switched between learning and transfer phases. Error bars indicate stan-
7 dard errors. Filled symbols correspond to performance in the pre-switch
8 phase. Unfilled symbols correspond to performance in the post-switch
9 phase. Dashed line and crosses indicate the best-fitting monotonic model.
10 Adapted from Figure 1b in *The effect of feedback delay and feedback type*
11 *on perceptual category learning: The limits of multiple systems*, by J. C.
12 Dunn, B. R. Newell, & M. L. Kalish, 2012, *Journal of Experimental Psy-*
13 *chology: Learning, Memory, & Cognition*, 38(4), pp. 840-859. Copyright
14 2012 by the American Psychological Association.
- 15 2. Empirical distributions of statistic, δ , based on analysis of data from
16 Nosofsky, Stanton, and Zaki (2005, Experiment 1). In the partial order
17 condition, a non-decreasing order is assumed over blocks 1 to 8 and over
18 blocks 9 to 10 in both the control and button-shift groups. Also shown are
19 the observed fit statistics for the data with and without the above partial
20 order, filled and unfilled triangles, respectively.
- 21 3. Power plots for the CMR effect size statistic, ω_{xy} , with no partial order
22 constraints and $k = 8$ conditions. (a) Power, $(1 - \beta)$, as a function of
23 effect size, ω_{xy} , and sample size, n_i , for $\alpha = 0.05$. (b) Power, $(1 - \beta)$, as
24 a function of effect size, ω_{xy} , and sample size, n_i , for $\alpha = 0.01$. Note the
25 different scales on the ordinates.
- 26 4. Power plots for the CMR effect size statistic, ω_{xy} , with a partial order con-
27 straint on $k = 8$ conditions (see text for constraint) compared to without

- 28 a partial order constraints. (a) Power, $(1 - \beta)$, as a function of effect size,
29 ω_{xy} , and sample size, n_i , for $\alpha = 0.05$. (b) Power, $(1 - \beta)$, as a function
1 of effect size, ω_{xy} , and sample size, n_i , for $\alpha = 0.01$.
- 2 5. State-trace plots of 4 x 2 factorial design corresponding to power of 0.80.
3 (a) Sample means and standard errors under no partial order. (b) Sample
4 means and standard errors under partial order defined on both factors.
- 5 6. State-trace plot of mean proportion correct (averaged over participants)
6 from Prince, Hawkins, Love and Heathcote (2012). The dashed line in-
7 dicates the best-fitting monotonic curve based on the CMR procedure.
8 Error bars indicate within-participant standard errors calculated accord-
9 ing to the Loftus-Masson procedure (Loftus & Masson, 1994).
- 10 7. Model structure tested by the CMR and Bayesian procedures. The left
11 hand side shows the model tree proposed by Prince, Brown and Heathcote
12 (2012) and tested by their Bayesian model selection procedure. The two
13 models at each level are the complements of each other and the Bayesian
14 procedure selects which of each pair is more strongly supported by the
15 data. The right hand side shows the constraints that added at each level
16 of the tree. The CMR procedure tests if the addition of each constraint
17 leads to a significant decrease in model fit. See text for a definition of each
18 term.













