



## Recovering a Basic Space from Issue Scales in R

**Keith T. Poole**  
University of Georgia

**Jeffrey B. Lewis**  
University of California,  
Los Angeles

**Howard Rosenthal**  
New York University

**James Lo**  
University of Southern California

**Royce Carroll**  
Rice University

---

### Abstract

**basicspace** is an R package that conducts Aldrich-McKelvey and Blackbox scaling to recover estimates of the underlying latent dimensions of issue scale data. We illustrate several applications of the package to survey data commonly used in the social sciences. Monte Carlo tests demonstrate that the procedure can recover latent dimensions and reproduce the matrix of responses at moderate levels of error and missing data.

*Keywords:* multivariate, R, scaling.

---

## 1. Introduction

Scaling techniques are used by political scientists in empirical models of voting (also known as ideal point models) to infer legislator locations in an abstract policy or ideological space from their legislative votes. The recovered scores have had wide applicability to the study of Congress (Poole and Rosenthal 1997; McCarty, Poole, and Rosenthal 2006), elections (Herron and Lewis 2007), courts (Martin and Quinn 2002), and in non-legislative voting bodies such as the United Nations (Voeten 2001).<sup>1</sup> However, the use of such models may not always be possible because roll call data is often not available or recorded and surveys may provide the only information with which preferences may be measured. In these instances, the basic space

---

<sup>1</sup>For a more extensive review of applications of spatial modeling in the social sciences, see Poole (2005). The most prominent ideal point model in the political science literature, W-NOMINATE (Poole and Rosenthal 1997), estimates the policy preferences of legislators using observed roll call votes as the primary source of data. The **wnominate** package (Poole, Lewis, Lo, and Carroll 2011) contains software used to estimate NOMINATE scores.

model shown here presents an attractive alternative estimator.<sup>2</sup>

The **basicspace** package is available from the Comprehensive R Archive Network at <https://CRAN.R-project.org/package=basicspace> and enables the spatial<sup>3</sup> analysis of self-placement and/or perceptual survey data in R (R Core Team 2015). Issue scales, where respondents place themselves and/or stimuli on a numeric scale, are a common form of data gathered and analyzed by survey researchers and social scientists. For example, since 1968 the American National Election Studies have gathered seven-point scale data on a variety of issues. Respondents are shown scales with labeled endpoints such as “liberal” and “conservative” and are then asked to place themselves and political figures on the scales. As with other forms of response data, researchers are often interested in understanding the extent to which a set of issue scale placements are driven by an underlying latent dimension. This package contains software designed to recover the latent dimensions – i.e., a basic space – from issue scale data such as surveys. The functions contained in **basicspace** will recover spatial estimates of respondent positions and scale them and the stimuli into a common space. The package implements the Blackbox method (Poole 1998) which generalizes the Aldrich-McKelvey scaling procedure (Aldrich and McKelvey 1977) to multiple dimensions and missing data. These methods have been used in a number of previous social science studies including Palfrey and Poole (1987), Poole (1998), and Saiegh (2009).

The motivation for recovering spatial information from issue scales, such as those in political surveys, is to detect the underlying dimensions behind the reported attitudes of survey respondents that explain the basic relationships among the respondents and stimuli. An early technique for analyzing such data, Aldrich and McKelvey’s method (Aldrich and McKelvey 1977), made use of respondent information regarding the positions of stimuli (e.g., politicians or parties) to estimate the perceptual bias of each respondent and obtain estimated locations for both stimuli and respondents along a single issue scale dimension (e.g., liberal-conservative). Poole (1998) developed the Blackbox scaling procedure as a generalization of Aldrich-McKelvey’s method, which is implemented in this package as the **blackbox** and **blackbox\_transpose** functions. The Basic Space method applies to a wide range of issue scale problems because it incorporates information from multiple issue scales to scale in multiple dimensions and because it allows for missing data (e.g., survey non-response). The **blackbox** function recovers  $N$ -dimensional ideal points (i.e., spatial coordinates) for respondents based on their own preference data across any number of issue scales. The **blackbox\_transpose** function, meanwhile, recovers the spatial location of stimuli based on respondent estimates.

This paper proceeds in four steps. First, we begin with a description of the mathematics that underlie the Blackbox estimator, which performs a singular value decomposition of a rectangular matrix containing missing elements. We then provide three examples using the **basicspace** package to implement this method in R. First, we show a Monte Carlo analysis that suggests the estimator produces an accurate decomposition of our simulated data matrices, even with 30 per cent of the data missing. Secondly, we show how the procedure can be applied to self-placement survey data from the 1980 National Election Survey. Next, we proceed with an application of the model to perceptual data from the 1980 National Election Study, where various political candidates are ranked along a 7 point liberal-conservative scale.

---

<sup>2</sup>For a more comprehensive review of the advantages and disadvantages of different data sources for spatial models, see Saiegh (2009).

<sup>3</sup>“Spatial” in this case specifically refers not to geography, but to the spatial model of voting popularized in the work of Downs (1957).

Finally, we describe the earlier estimator developed by [Aldrich and McKelvey \(1977\)](#) which is also included as a function in this package.

## 2. Model

The exposition of the model presented here closely follows [Poole \(1998\)](#). Consider a matrix of survey data  $X_0$  with  $N$  respondents and  $M$  issue scales, with individuals on the rows and issues on the columns. Some cells of the matrix  $X_0$  are missing, and we let  $X$  denote the version of  $X_0$  that has no missing data. In each cell  $x_{ij}$ , respondent  $i$  ( $i = 1, \dots, N$ ) reports their position on issue scale  $j$  ( $j = 1, \dots, M$ ), with some responses missing.<sup>4</sup> Now let  $\psi_{ik}$  be the  $i$ th individual's position on the  $k$ th basic dimension ( $k = 1, \dots, s$ ),  $W$  be an  $M$  by  $s$  matrix of weights that map individual positions from the basic space to the issue dimensions,  $c$  be a vector of issue dimension intercept terms of length  $M$ ,  $J_N$  by an  $N$  length vector of ones, and  $E_0$  be error terms in the data matrix. The model that we seek to estimate is:

$$X_0 = [\Psi W^\top + J_N c^\top]_0 + E_0 \quad (1)$$

Without loss of generality, we also assume that  $E_0$  is drawn from a symmetric distribution with mean 0 and the centroid of the basic space coordinates is at the origin (i.e.,  $J_N^\top \Psi = 0$ ). Substituting into the model equation, this implies that  $J_N^\top [X - J_N c^\top] = 0_M^\top$ , where  $0_M$  is an  $M$  length vector of zeroes. Then in the situation where  $X_0$  has no missing data, the parameters of interest can all be recovered using singular value decomposition. To see why this is true, recall that for an  $N$  by  $M$  matrix of real elements with  $N \geq M$ , there exists an  $N$  by  $M$  orthogonal matrix  $U$ , an  $M$  by  $M$  orthogonal matrix  $V$ , and an  $M$  by  $M$  matrix  $\Lambda$  such that:

$$X = U \Lambda V^\top \quad (2)$$

where  $\Lambda$  is a diagonal matrix of singular values.<sup>5</sup> To solve Equation 1, set  $c$  equal to the column means of  $X$ , or  $c_j = N^{-1} \sum_{i=1}^N x_{ij} = \bar{x}_j$ . Then using Equation 2, the singular value decomposition of  $X - J_N c^\top$  can be expressed as:

$$X - J_N c^\top = U \Lambda V^\top = \Psi W^\top$$

This implies that in the absence of missing data, one solution for  $\Psi$  and  $W$  is:

$$\begin{aligned} \Psi &= U \Lambda^{0.5} \\ W &= V \Lambda^{0.5} \end{aligned}$$

with  $\Lambda^{0.5}$  being a diagonal matrix where diagonal elements are the square roots of  $\Lambda$ . While other solutions to this problem exist, [Eckart and Young \(1936\)](#) have shown that the least squares approximation in  $s$  dimensions of a matrix  $X$  can be found by using only the first  $s$  singular values of  $X$  along the diagonal of  $\Lambda$  and re-multiplying  $U \Lambda V^\top$ .

In the presence of missing data in data matrix  $X_0$ , the use of singular value decomposition to solve for  $W$  and  $\Psi$  is no longer possible, and we instead estimate  $\hat{W}$  and  $\hat{\Psi}$  using an alternating

<sup>4</sup>The '0' subscript indicates that some elements are missing from the matrix.

<sup>5</sup>A more general form of this equation can be written in which  $\Lambda$  is instead an  $N$  by  $M$  matrix and  $U$  is  $N$  by  $N$ .

least squares (ALS) technique that is similar to the procedures used in [Carroll and Chang \(1970\)](#) and [Takane, Young, and Leeuw \(1977\)](#). The objective function to be minimized is the sum of the squared deviations across all cells in  $X$  after the columns have been adjusted for column means, or:

$$\xi = \sum_{i=1}^N \sum_{j=1}^{m_i} \left\{ \left[ \sum_{k=1}^s \psi_{ik} w_{jk} \right] + c_j - x_{ij} \right\}^2$$

where  $m_i$  is the number of non-missing entries on row  $i$ . In minimizing this objective function, two constraints from the earlier analysis with no missing data are applied. First, we exploit the fact that  $\Psi$  and  $W$  are orthogonal matrices, which implies that  $\Psi^\top \Psi = W^\top W$ .<sup>6</sup> Secondly, following our earlier restriction that  $J_N^\top [X - J_N c^\top] = 0_M^\top$ ,  $J_N^\top U = J_N^\top \Psi = 0_M^\top$  as well (where  $0$  is a null vector). These restrictions produce the Lagrangian multiplier problem:

$$\mu = \xi + 2\gamma^\top [\Psi^\top J_N] + \text{tr}[\Phi(\Psi^\top \Psi - W^\top W)]$$

where  $\Phi$  is a symmetric  $s$  by  $s$  matrix of Lagrangian multipliers and  $\gamma$  is an  $s$  length vector of Lagrangian multipliers. Since all Lagrangian multipliers are zero,<sup>7</sup> the partial derivatives of  $\xi$  are:

$$\frac{\partial \mu}{\partial \psi_{ik}} = 2 \sum_{j=1}^{m_i} \left[ \left( \sum_{l=1}^s w_{jl} \psi_{il} \right) + c_j - x_{ij} \right] w_{jk} \quad (3)$$

$$\frac{\partial \mu}{\partial w_{jk}} = 2 \sum_{i=1}^{n_j} \left[ \left( \sum_{l=1}^s w_{jl} \psi_{il} \right) + c_j - x_{ij} \right] \psi_{jl} \quad (4)$$

$$\frac{\partial \mu}{\partial c_j} = 2 \sum_{i=1}^{n_j} \left[ \left( \sum_{l=1}^s w_{jl} \psi_{il} \right) + c_j - x_{ij} \right] \quad (5)$$

Let  $W^*$  be an  $m_i$  by  $s$  matrix with appropriate rows corresponding to missing entries in  $X_0$  removed,  $x_{0i}$  be the length  $m_i$  row of  $X_0$ , and  $c_0$  be the length  $m_i$  vector of constants corresponding to the elements of  $x_{0i}$ . Then if  $W^{*\top} W^*$  exists, the  $i$ th row of  $\Psi$  can be estimated by setting Equation 3 to zero, collecting the  $s$  partial derivatives of the  $i$ th row of  $\Psi$  into a vector and solving for  $\psi_j$  as:

$$\hat{\psi}_i = (W^{*\top} W^*)^{-1} W^{*\top} [x_{0i} - c_0] \quad (6)$$

which can of course be estimated using ordinary least squares. Similarly, let  $\Psi_j^* = [\Psi_0 | J_0]$  be an  $n_j$  by  $s + 1$  matrix with the appropriate rows corresponding to missing data removed and bordered by ones,  $w_j$  be the  $s$  length vector of row  $j$  in  $W$ ,  $c_j$  be the  $j$ th element of  $c$ , and  $x_{0j}$  be the  $j$ th column of  $X_0$ . Then if  $\Psi_j^{*\top} \Psi_j^*$  exists,  $w_j$  and  $c_j$  can be jointly estimated by combining Equations 4 and 5 as:

$$\frac{\hat{w}_j}{\hat{c}_j} = (\Psi_j^{*\top} \Psi_j^*)^{-1} \Psi_j^{*\top} x_{0j} \quad (7)$$

Equations 6 and 7 represent the core set of equations that are used to solve for  $\hat{W}$ ,  $\hat{c}$ , and  $\hat{\Psi}$ . Once a set of starting values has been generated, Equations 6 and 7 are iterated until

<sup>6</sup>More specifically,  $\Psi^\top \Psi = \Lambda^{0.5} U^\top U \Lambda^{0.5} = \Lambda^{0.5} I_M \Lambda^{0.5} = \Lambda = W^\top W$ .

<sup>7</sup>See Appendix A in [Poole \(1998\)](#) for a full proof that all Lagrangian multipliers are zero.

convergence. Generation of appropriate start values is conducted one dimension at a time, and a more detailed justification of the procedure can be found in [Poole \(1998\)](#). On the first dimension, starting values are generated using the following three equations:

$$\hat{c}_j = \frac{\sum_{i=1}^{n_j} x_{ij}}{n_j} = \bar{x}_j \quad (8)$$

$$w_{j1} = \text{diag}(\Gamma) \quad (9)$$

where  $\Gamma$  is an  $M$  by  $M$  diagonal matrix with diagonal elements either set to 1 or -1 that maximizes the number of positive elements in the  $M$  by  $M$  covariance matrix  $\Gamma[X_0 - J_N c^\top]^\top [X_0 - J_N c^\top] \Gamma$ .  $\Gamma$  is found by a simple iterative process similar to that used to speed eigenvector/eigenvalue decomposition ([Poole 1998](#)). Given Equations 8 and 9, starting values for  $\psi$  are:

$$\hat{\psi}_{i1} = \frac{\sum_{j=1}^{m_i} \hat{w}_{j1} (x_{ij} - \hat{c}_j)}{m_i} \quad (10)$$

If more than one dimension is to be estimated ( $s > 1$ ), start values for other dimensions can be generated simply by replacing the data matrix  $X_0$  with the matrix of residuals  $E_{0s}$  in Equations 9 and 10. However, no further estimation of start values for  $\hat{c}$  is required. The matrix of residuals to be used for generating start values on dimension  $s$  is:

$$E_{0s} = X_0 - \sum_{j=1}^s \hat{\Psi}_s \hat{w}_s^\top - J_N \hat{c}^\top$$

This residual matrix allows the generation of higher-dimension start values by iterating  $\Gamma$  to maximize the positive elements in  $E_{0s}$ . The starting values are now:

$$\hat{\psi}_{is} = \frac{\sum_{j=1}^{m_i} \hat{w}_{js} e_{(s-1)ij}}{\sum_{j=1}^{m_j} \hat{w}_{js}^2} \quad (11)$$

where the initial  $\hat{w}_{js}$  values of +1s and -1s are used to obtain  $\hat{\psi}_{is}$  starting values. The starting values of  $\hat{w}_{js}$  are now:

$$\hat{w}_{js} = \frac{\sum_{i=1}^{n_j} \hat{\psi}_{is} e_{(s-1)ij}}{\sum_{i=1}^{n_j} \hat{\psi}_{is}^2} \quad (12)$$

Summarizing the preceding discussion in full, the basic space technique decomposes an  $N$  by  $M$  matrix  $X_0$  with  $N \geq M$  following Equation 1. Estimation of Equation 1 proceeds in three steps. In the first stage, starting values on the first dimension are generated for  $\hat{c}_j$ ,  $\hat{w}_{j1}$ , and  $\hat{\psi}_{i1}$  by iterating Equations 8–10 until convergence. In the second stage, if  $s > 1$ , higher dimensional starting values for  $\hat{\psi}_{is}$  and  $\hat{w}_{is}$  are generated dimension by dimension using Equations 11–12. Finally, the starting values generated in the preceding two stages are improved by iterating Equations 6–7 until convergence.

### 3. Monte Carlo test

In this section, we present the first of four motivating examples. We begin with a Monte Carlo example that tests the basic space technique against simulated data. Four key variables should be set in each simulation: the number of respondents  $N$  (set here to  $N = 1000$ ), the number of issue scales (also referred to as stimuli, and set here to  $M = 20$ ), the number of explanatory dimensions (set here to  $s = 2$ ), the fraction of observations that are missing (set here as 0.3), and the distribution of error terms (set here as random uniform draws from  $-0.5$  to  $0.5$ ). These variables can be changed for other simulations, but the restriction that  $N \geq M$  must be hold true. In cases where  $M \geq N$  please refer to the second example that uses `blackbox_transpose` and `aldmck`.

```
R> set.seed(1231)
R> library("basicspace")
R> N <- 1000
R> M <- 20
R> s <- 2
R> fraction.missing <- 0.3
R> E <- matrix(runif(N * M, min = -0.5, max = 0.5), nrow = N, ncol = M)
```

To generate the  $X$  matrix (i.e., the matrix in Equation 1 before missing values are introduced), separately generate the matrices that produce the singular value decomposition of  $X$  following Equation 2. Also generate the  $J_n$  and  $c$  vectors from Equation 1. While  $X$  can be generated directly in one step, creating the components separately enjoys two significant advantages. First, recovery of the true values of  $\Psi$  and  $W$  is simplified. Secondly, the creation of  $\Lambda$  separately allows us to more easily tune the dimensionality of the matrix as desired.

```
R> U <- matrix(runif(N * s), nrow = N, ncol = s)
R> D <- diag(seq(from = 2.1, by = -0.2, length.out = s))
R> V.prime <- matrix(runif(s * M), nrow = s, ncol = M)
R> c <- rnorm(M)
R> Jn <- rep(1, N)
```

With the intermediate matrices just generated, we can produce our  $X$  matrix by using Equation 1 and the true  $\Psi$  and  $W$  matrices using:  $\Psi = U\Lambda^{0.5}$  and  $W = V\Lambda^{0.5}$ .

```
R> X.true <- U %*% D %*% V.prime + Jn %o% c
R> X.0 <- X.true + E
R> Psi.true <- U %*% sqrt(D)
R> W.true <- t(V.prime) %*% sqrt(D)
```

$X_0$  is simply the  $X$  matrix with missing data values included completely at random, so we insert our missing data code (999 in the example) into the appropriate fraction of values as follows:

```
R> missing <- sample(1:(N * M), round(fraction.missing * N * M))
R> X.0[missing] <- 999
```

The final step before estimation is to assign row and column names to the data set prior to input. In most applications these names are generally pulled from a survey, but they can also be generated manually:

```
R> rownames(X.0) <- paste("Legis", 1:N, sep = "")
R> colnames(X.0) <- paste("V", 1:M, sep = "")
```

Estimation of the Monte Carlo data after formatting is trivial. The function that applies the basic space decomposition described in this paper is `blackbox`. It takes four arguments: the matrix to be decomposed, a vector of missing data values, a Boolean flag indicating whether verbose output is desired, the number of dimensions to estimate, and the minimum number of issue scales that an individual needs to provide responses to if they are to be included in the estimation.

```
R> result <- blackbox(X.0, missing = 999, verbose = TRUE, dims = 3,
+   minscale = 8)
```

```
Beginning Blackbox Scaling...20 stimuli have been provided.
```

```
Blackbox estimation completed successfully.
```

```
R> names(result)[1:4]
```

```
[1] "stimuli"      "individuals"  "fits"        "Nrow"
```

The output object contains multiple data frames summarizing the results of the estimation. The key data frames are `stimuli`, which contain estimates of  $\hat{W}$  and  $\hat{c}$ , as well as `individuals`, which contain estimates of  $\hat{\Psi}$ . The other quantities are fit statistics described in greater detail in the standard documentation for the function.

With the estimates complete, we are now able to test the recovery of our parameters of interest. In general, scaling problems are not fully identified. Stated differently, given  $X = \Psi W^\top$ ,  $\Psi$  and  $W^\top$  are not unique solutions because  $X = \Psi K K^{-1} W^\top$  for any conformable and invertible matrix  $K$ , so  $X$  can always be decomposed instead as  $X = \Psi^* W^{*\top}$  where  $\Psi^* = \Psi K$  and  $W^{*\top} = K^{-1} W^\top$ . When evaluating parameter fit, we are therefore largely concerned with finding monotonic relationships between the true and estimated parameters of interest. Figure 1 compares the true vs. estimated values of  $\Psi$  across two dimensions, and the results suggest a reasonable model fit. For this comparison,  $\Psi$  and  $\Psi^*$  are mean centered and rotated.

```
R> Psi.hat <- cbind(result$individuals[[2]]$c1, result$individuals[[2]]$c2)
R> c.hat <- result$stimuli[[2]]$c
R> xrow <- sapply(1:N, function(i) length(rep(1, s)[!is.na(Psi.hat[i, ])]))
R> Psi.hat <- Psi.hat[!(xrow < 2), ]
R> Psi.true <- Psi.true[!(xrow < 2), ]
R> Psi.hat[,1] <- Psi.hat[, 1] - mean(Psi.hat[, 1])
R> Psi.hat[,2] <- Psi.hat[, 2] - mean(Psi.hat[, 2])
R> Psi.true[,1] <- Psi.true[, 1] - mean(Psi.true[, 1])
```

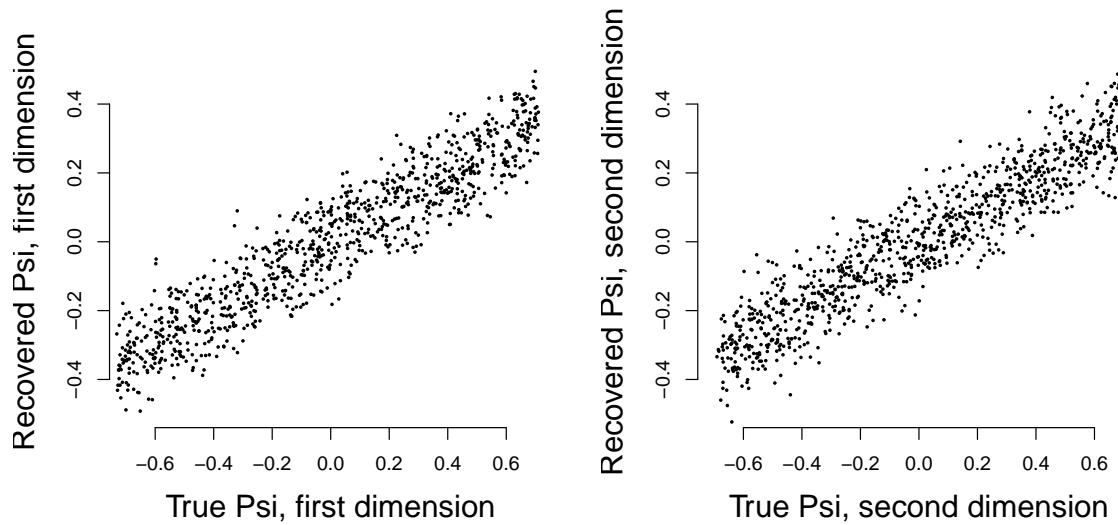


Figure 1: Plots of true vs. estimated  $\Psi$  scores, first and second dimension.

```
R> Psi.true[,2] <- Psi.true[, 2] - mean(Psi.true[, 2])
R> C <- t(Psi.true) %*% Psi.hat
R> svddecomp <- svd(C)
R> U.rotate <- svddecomp$u
R> V.rotate <- svddecomp$v
R> T <- V.rotate %*% t(U.rotate)
R> Psi.hatrotate <- Psi.hat %*% T
R> par(mfrow = c(1, 2))
R> plot(Psi.true[, 1], Psi.hatrotate[, 1], xlim = c(-0.7, 0.7),
+       ylim = c(-0.5, 0.5), pch = 20, cex = 0.4, cex.lab = 1.6, bty = "n",
+       xlab = "True Psi, first dimension",
+       ylab = "Recovered Psi, first dimension")
R> plot(Psi.true[, 2], Psi.hatrotate[, 2], xlim = c(-0.7, 0.7),
+       ylim = c(-0.5, 0.5), pch = 20, cex = 0.4, cex.lab = 1.6, bty = "n",
+       xlab = "True Psi, second dimension",
+       ylab = "Recovered Psi, second dimension")
```

Figure 2 shows the results for the same procedure applied to  $W$ . In Figure 3 we repeat this analysis for  $c$ , which is a column mean that is only estimated in one dimension. In both cases the estimates for  $\hat{W}$  and  $\hat{c}$  are a monotonic transformation of the true parameters as expected.<sup>8</sup>

```
R> W.hat <- cbind(result$stimuli[[2]]$w1, result$stimuli[[2]]$w2)
R> W.hatrotate <- W.hat %*% T
R> par(mfrow = c(1,2))
```

<sup>8</sup>In other estimates, the relationship may only be affine because  $X = \Psi W^T$  implies  $X = -(\Psi) - (W^T)$  as well.



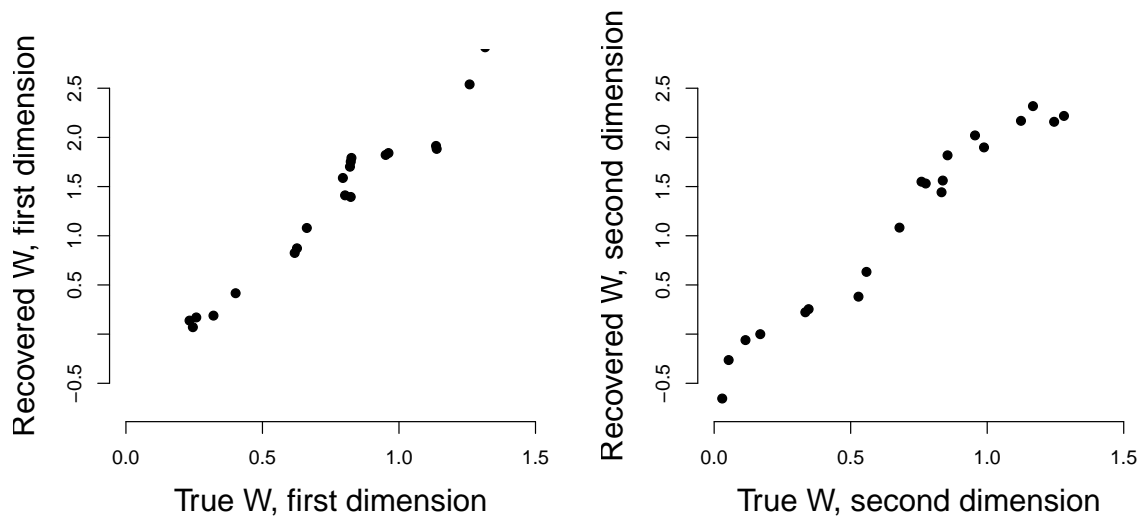
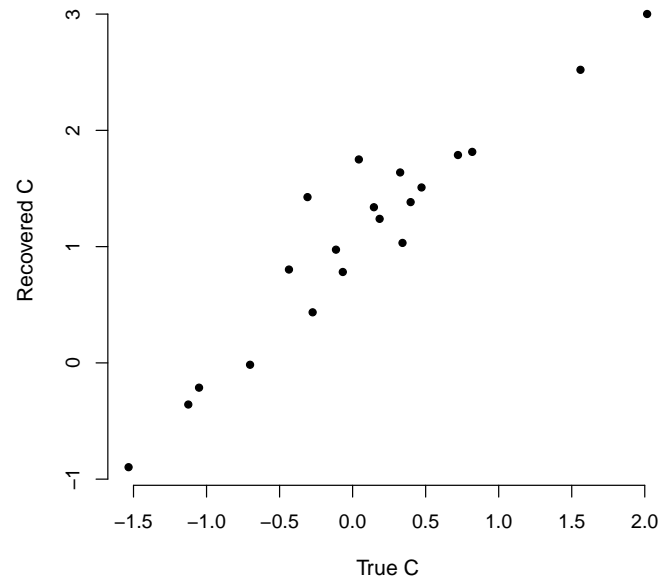
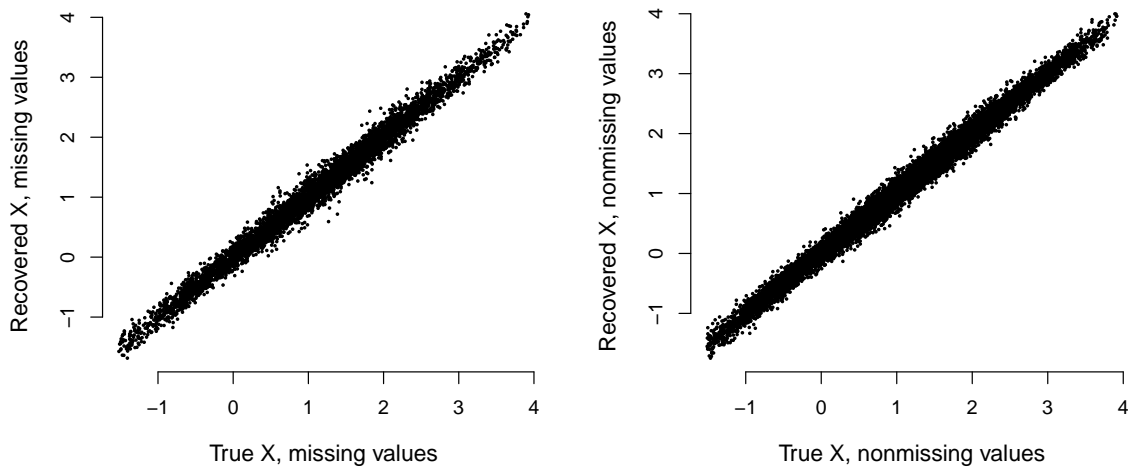


Figure 2: Plots of true vs. estimated  $W$  scores, first and second dimension.

```
R> plot(W.true[, 1], W.hatrotate[, 1], xlim = c(0.00, 1.50),
+       ylim = c(-0.75, 2.75), pch = 20, cex = 1.5, cex.lab = 1.6, bty = "n",
+       xlab = "True W, first dimension",
+       ylab = "Recovered W, first dimension")
R> plot(W.true[, 2], W.hatrotate[, 2], xlim = c(0.00, 1.50),
+       ylim = c(-0.75, 2.75), pch = 20, cex = 1.5, cex.lab = 1.6, bty = "n",
+       xlab = "True W, second dimension",
+       ylab = "Recovered W, second dimension")
R> par(mfrow = c(1, 1))
R> plot(c, c.hat, pch = 20, cex = 1.2, cex.lab = 1.1, bty = "n",
+       xlab = "True C", ylab = "Recovered C")
```

Finally, we pool our estimates of  $\hat{W}$ ,  $\hat{\Psi}$ , and  $\hat{c}$  together to estimate the full matrix  $\hat{X}$  following Equation 1. While social scientists are principally concerned with estimation of  $\hat{W}$  and  $\hat{\Psi}$ , others seeking to conduct singular value decomposition of matrices with missing data may find  $\hat{X}$  to be of value. One obvious application of  $\hat{X}$  is its potential use as an imputation tool for missing data.<sup>9</sup> To test the viability of this idea, we separately plot the true values of  $X$  against the estimated values of  $\hat{X}$  separately for the cells retained in the estimation, and compare those results to estimates of  $\hat{X}$  in cells that were discarded prior to estimation to simulate the missing data mechanism. Figure 4 presents our results for retained vs. imputed  $X$ . What is particularly notable about this result is the close similarity between these plots – the imputed values not only appear reasonable (i.e., line up with the true values along a 45° line), but imputed values do not appear to have significantly higher mean squared error than the values that were retained (i.e., variance along the 45° line is similar in both plots). These results suggest that the use of the techniques demonstrated here may have greater applicability beyond survey research. Further discussion of imputation can be found in the

<sup>9</sup>The simulation presented here simulates missing data under the missing completely at random (MCAR) assumption – nevertheless, this should also work under conditions where data are instead missing at random (MAR).

Figure 3: Plot of true vs. estimated  $c$  scores.Figure 4: Plots of true vs. estimated  $X$  scores for missing vs. non-missing values.

unpublished appendix to [Poole \(1998\)](#).

```
R> W.hat <- cbind(result$stimuli[[2]]$w1, result$stimuli[[2]]$w2)
R> Psi.hat <- cbind(result$individuals[[2]]$c1,
+   result$individuals[[2]]$c2)
R> X.hat <- Psi.hat %*% t(W.hat) + Jn %o% result$stimuli[[2]]$c
R> par(mfrow = c(1, 2))
R> plot(X.true[missing], X.hat[missing], pch = 20, cex = 0.4, cex.lab = 1.2,
+   bty = "n", xlab = "True X, missing values",
+   ylab = "Recovered X, missing values")
```

```
R> plot(X.true[!(1:(N * M) %in% missing)], X.hat[!(1:(N * M) %in% missing)],
+      pch = 20, cex = 0.4, cex.lab = 1.2, bty = "n",
+      xlab = "True X, nonmissing values",
+      ylab = "Recovered X, nonmissing values")
```

#### 4. Example 1: 1980 NES issue scales

In this section we present an application of the basic space model to a set of issue scales from the 1980 National Election Study. This survey contains  $N = 1,614$  respondents who were asked to place themselves on scales about desired levels of defense spending, inflation, tax cuts, abortion, liberal-conservative scales, the role of women, the role of government in providing jobs, busing, and other similar issues. We assume that each respondent has a location in a common ideological space and attempt to recover estimates of those locations, which is represented as  $\Psi$  in Equation 1. The data is simply stored in a standard matrix or data frame with respondents on the rows and survey questions (i.e., stimuli) on the columns as follows:

```
R> data("Issues1980")
R> Issues1980[1:10, 1:4]
```

	libcon1	defense	govserv	inflation
1	0	7	5	4
2	4	4	6	7
3	6	3	0	0
4	5	6	2	8
5	3	4	2	4
6	5	5	4	0
7	8	2	6	5
8	2	7	7	6
9	6	7	2	2
10	5	4	2	5

Virtually all surveys contain missing data, and for the two survey questions about abortion, ‘7’ is used as a missing data code. However, many of the other scales in this data set use 7 point scales, so we need to recode the missing data for those questions. For all questions, 0, 8, and 9 are missing data codes.<sup>10</sup>

```
R> Issues1980[Issues1980[, "abortion1"] == 7, "abortion1"] <- 8
R> Issues1980[Issues1980[, "abortion2"] == 7, "abortion2"] <- 8
```

Estimation of the scores is now trivial using the `blackbox` function, which takes the same arguments already described in the Monte Carlo example:

```
R> Issues1980_bb <- blackbox(Issues1980, missing = c(0, 8, 9),
+   verbose = FALSE, dims = 3, minscale = 8)
```

<sup>10</sup>Data used for this estimator are typically opinion surveys where significant amounts of missing data are commonplace – thus, recoding of this sort will typically be necessary for most applications.

Objects of class `blackbox` can also be summarized using the `summary` function, although the summaries largely provide only summaries of the stimuli. For each dimension estimated, the summary provides the intercept ( $c$ ) and stretch ( $w_1 \dots w_3$ ) parameters for each question, as well as the number of respondents and various fit statistics.

```
R> summary(Issues1980_bb)
```

SUMMARY OF BLACKBOX OBJECT

```
-----
```

	N	c	w1	R2
libcon1	875	4.280	-3.028	0.414
defense	1163	5.210	-1.754	0.123
govserv	1119	4.323	4.302	0.450
inflation	816	4.106	2.015	0.159
abortion1	1238	2.856	0.627	0.031
taxcut	836	2.839	-1.074	0.055
libcon2	949	4.369	-2.755	0.414
govhelpmin	1160	4.542	-3.400	0.412
russia	1152	3.891	-3.034	0.231
womenrole	1223	2.845	-2.866	0.204
govjobs	1131	4.377	-4.488	0.518
equalrights	1144	2.663	-3.297	0.381
busing	1219	6.051	-2.699	0.255
abortion2	1246	2.675	0.724	0.047

	N	c	w1	w2	R2
libcon1	875	4.300	-2.966	0.954	0.424
defense	1163	5.214	-1.779	0.899	0.147
govserv	1119	4.368	4.331	3.042	0.617
inflation	816	4.152	2.088	2.940	0.393
abortion1	1238	2.856	0.512	-2.211	0.290
taxcut	836	2.818	-1.103	-0.667	0.071
libcon2	949	4.377	-2.758	0.459	0.423
govhelpmin	1160	4.535	-3.456	-0.119	0.424
russia	1152	3.887	-3.140	0.241	0.247
womenrole	1223	2.872	-2.466	6.007	0.771
govjobs	1131	4.350	-4.595	-2.417	0.635
equalrights	1144	2.673	-3.148	2.438	0.491
busing	1219	6.049	-2.741	0.059	0.263
abortion2	1246	2.676	0.629	-2.112	0.318

	N	c	w1	w2	w3	R2
libcon1	875	4.294	-2.976	0.708	-1.180	0.448
defense	1163	5.200	-1.806	1.586	2.562	0.315
govserv	1119	4.410	4.295	3.707	2.929	0.778
inflation	816	4.169	1.998	3.286	1.111	0.451
abortion1	1238	2.856	0.497	-2.004	1.174	0.312

taxcut	836	2.813	-1.049	-0.902	-0.891	0.091
libcon2	949	4.367	-2.785	0.265	-0.557	0.437
govhelpmin	1160	4.534	-3.457	0.140	0.961	0.440
russia	1152	3.831	-3.255	1.558	5.590	0.695
womenrole	1223	2.891	-2.372	5.602	-2.868	0.805
govjobs	1131	4.341	-4.632	-2.176	1.392	0.648
equalrights	1144	2.680	-3.159	1.860	-2.372	0.563
busing	1219	6.042	-2.819	0.329	1.282	0.306
abortion2	1246	2.675	0.587	-1.980	0.906	0.329

```

Dimensions Estimated: 3
Number of Rows: 1270
Number of Columns: 14
Total Number of Data Entries: 15271
Number of Missing Entries: 2509
Percent Missing Data: 14.11%
Sum of Squares (Grand Mean): 52705.13

```

When using `blackbox` for applied research, the researcher's principal goal is the recovery of the individual parameters stored as the `individuals` data frame. These typically represent our estimate of the individual's ideological location in the basic space. Due to the model identification issue discussed above, these measures are defined only up to an affine transformation of the true space. In particular, the rotation of the estimate is not specified, so if the ideological location is to be substantively measured as a liberalism/conservatism score, its rotation should be validated so that it can be transformed if necessary. Here we conduct such a check by correlating our recovered scores with self-reported liberal-conservative scores, where higher scores indicate higher levels of conservatism. The correlation is negative, suggesting that as the recovered scores increase, the respondents become more liberal. Since the norm in political science research is to orient liberal-conservative scores to increase as conservatism increases, the researcher may wish to rotate the scores (i.e., by multiplying them by -1) before using them for auxiliary analyses.

```

R> cor(Issues1980_bb$individuals[[1]]$c1, Issues1980[, "libcon1"],
+      use = "pairwise")

```

```
[1] -0.2310037
```

## 5. Example 2: 1980 NES liberal-conservative scale

In our previous example applying the basic space model to analyze respondent self-placement on issue scales, we considered an example where the bias and stretch parameters  $c$  and  $w$  were estimated for the column parameters. However, we may instead wish to estimate a version of the model where  $c$  and  $w$  are estimated for the row parameters (i.e., the survey respondents) instead. This is simply a transposed version of the basic space model, where  $M \geq N$  instead of  $N \leq M$ . In this example we analyze perceptual data from the 1980 National Election Study. A total of  $N = 888$  respondents were asked to place six stimuli (Carter, Reagan,

Kennedy, Anderson, the Republicans, and the Democrats) on a 7 point liberal-conservative scale. Our objective is to estimate the locations of the six stimuli in the basic space, which each respondent perceives with some bias and stretch parameter. The data is input in a manner identical to before, with survey respondents on the rows and stimuli on the columns. One very important difference between `blackbox` and `blackbox_transpose` is that in most survey data sets, the number of respondents is very large relative to the number of stimuli. This typically means that `blackbox_transpose` takes much longer to estimate because it estimates both a bias  $c$  and stretch  $W$  parameter for each respondent. To estimate the 1980 liberal-conservative placements using `blackbox_transpose`, we simply load the data and call the function as follows:

```
R> data("LC1980")
R> LCdat <- LC1980[, -1]
R> LCdat[1:10, ]
```

	Carter	Reagan	Kennedy	Anderson	Republicans	Democrats
1	2	6	1	7	5	5
8	4	6	4	7	6	4
9	3	6	3	3	6	2
10	6	4	3	3	5	4
11	7	2	5	5	7	5
13	6	6	2	5	7	4
14	3	6	2	5	6	3
16	3	7	4	2	7	3
17	5	3	5	2	8	8
19	3	6	4	5	6	2

To avoid the long runtime of

```
R> LC1980_bbt <- blackbox_transpose(LCdat, missing = c(0, 8, 9),
+   dims = 3, minscale = 5, verbose = TRUE)
```

the precomputed results can be loaded

```
R> data("LC1980_bbt")
```

In an effort to simplify interpretation of results from `blackbox_tranpose`, we include two plot functions. These functions plot the location of the stimuli against a probability and cumulative distribution plot of locations of the population weights (see Figure 5).

```
R> par(mfrow = c(1, 2))
R> plot(LC1980_bbt)
R> plotcdf.blackbt(LC1980_bbt)
```

We can also produce summary reports of the stimuli as follows:

```
R> summary(LC1980_bbt)
```

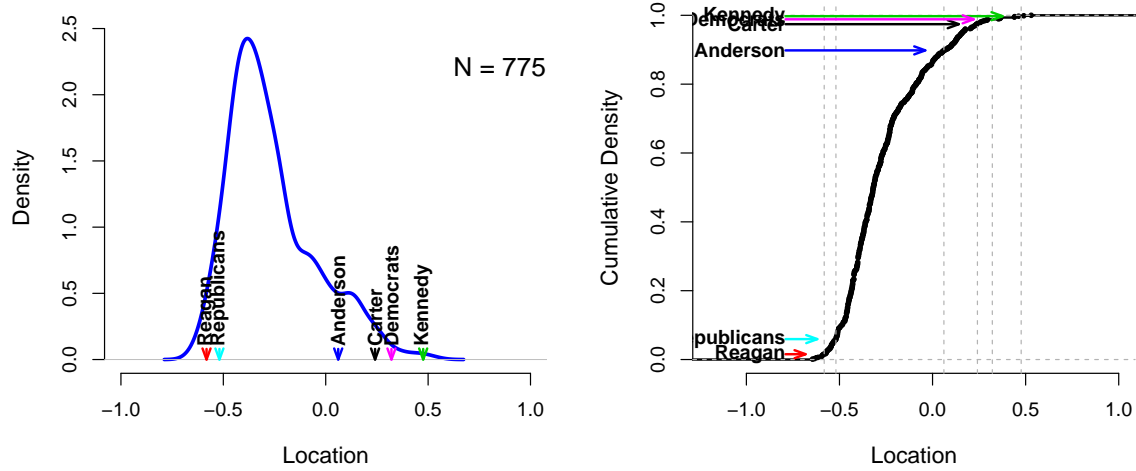


Figure 5: Blackbox transpose PDF and CDF plots.

SUMMARY OF BLACKBOX TRANSPOSE OBJECT

	N	coord1D	R2
Carter	768	0.241	0.563
Reagan	765	-0.582	0.822
Kennedy	754	0.476	0.648
Anderson	689	0.061	0.230
Republicans	771	-0.519	0.757
Democrats	774	0.321	0.651

	N	coord1D	coord2D	R2
Carter	768	0.238	-0.407	0.720
Reagan	765	-0.580	-0.101	0.839
Kennedy	754	0.481	0.013	0.680
Anderson	689	0.059	0.864	0.946
Republicans	771	-0.518	-0.117	0.767
Democrats	774	0.321	-0.252	0.718

	N	coord1D	coord2D	coord3D	R2
Carter	768	0.191	-0.261	-0.663	0.918
Reagan	765	0.216	0.556	0.141	0.856
Kennedy	754	0.162	-0.510	0.697	0.981
Anderson	689	-0.911	0.053	-0.002	1.000
Republicans	771	0.210	0.498	0.055	0.780
Democrats	774	0.131	-0.335	-0.228	0.765

Dimensions Estimated: 3  
 Number of Rows: 6  
 Number of Columns: 775

Total Number of Data Entries: 4521  
 Number of Missing Entries: 129  
 Percent Missing Data: 2.77%  
 Sum of Squares (Grand Mean): 12683.93

The second dimension is picking up John Anderson, a Representative from Illinois who ran as a third party candidate in 1980. Respondents clearly had trouble placing Anderson on the liberal-conservative scale. The second dimension is picking up this ambiguity of position.

## 6. Example 3: Aldrich and McKelvey's estimator

The transposed basic space model is a generalization of a model developed by Aldrich and McKelvey (1977), which was restricted to analyzing matrices with no missing values in only one dimension. For historical purposes, we include the original Aldrich-McKelvey estimator with this package. The Aldrich-McKelvey model is:

$$Y_{ij} = Z_j + \epsilon_{ij}$$

where  $Z_j$  is the true location of  $j$  and  $\epsilon_{ij}$  is a random variable with mean 0, positive variance that is independent of  $i$  and  $j$  (homoskedastic), and zero covariance across the  $i$ 's and  $j$ 's. Aldrich and McKelvey then introduce two distortion parameters,  $c_i$  and  $w_i$ , that transform the perceived candidate position into a reported candidate position  $R_{ij}$ , according to:

$$R_{ij} = \frac{1}{w_i}(Y_{ij} - c_i)$$

A least-squares minimization procedure is then used to obtain estimates of  $\{Z_j\}_{j=1}^J$  and  $\{w_i, c_i\}_{i=1}^I$ .

We begin by reestimating the earlier results using the 1980 liberal-conservative scale with the Aldrich-McKelvey estimator. While the `aldmck` function accepts nearly identical arguments the `blackbox_transpose`, one notable difference appears by default. `aldmck` also accepts a column in the data matrix, specified by the `respondent` argument, that specifies the respondent's self placement on the issue scale. The reported respondent rating is then transformed into an ideology score by applying the respondent's personal stretch and bias parameters to that score, with the results shown in Figure 6. Note that the results largely correspond to those shown earlier with `blackbox_transpose`.

```
R> data("LC1980")
R> result <- aldmck(data = LC1980, polarity = 2, respondent = 1,
+   missing = c(0, 8, 9), verbose = TRUE)
```

Beginning Aldrich-McKelvey Scaling...

```
Column 'Self' is set as the self placement.
Column 'Carter' is set as the left-leaning stimulus.
646 of 888 observations are complete.
6 stimuli have been provided.
```

Aldrich-McKelvey estimation completed successfully.



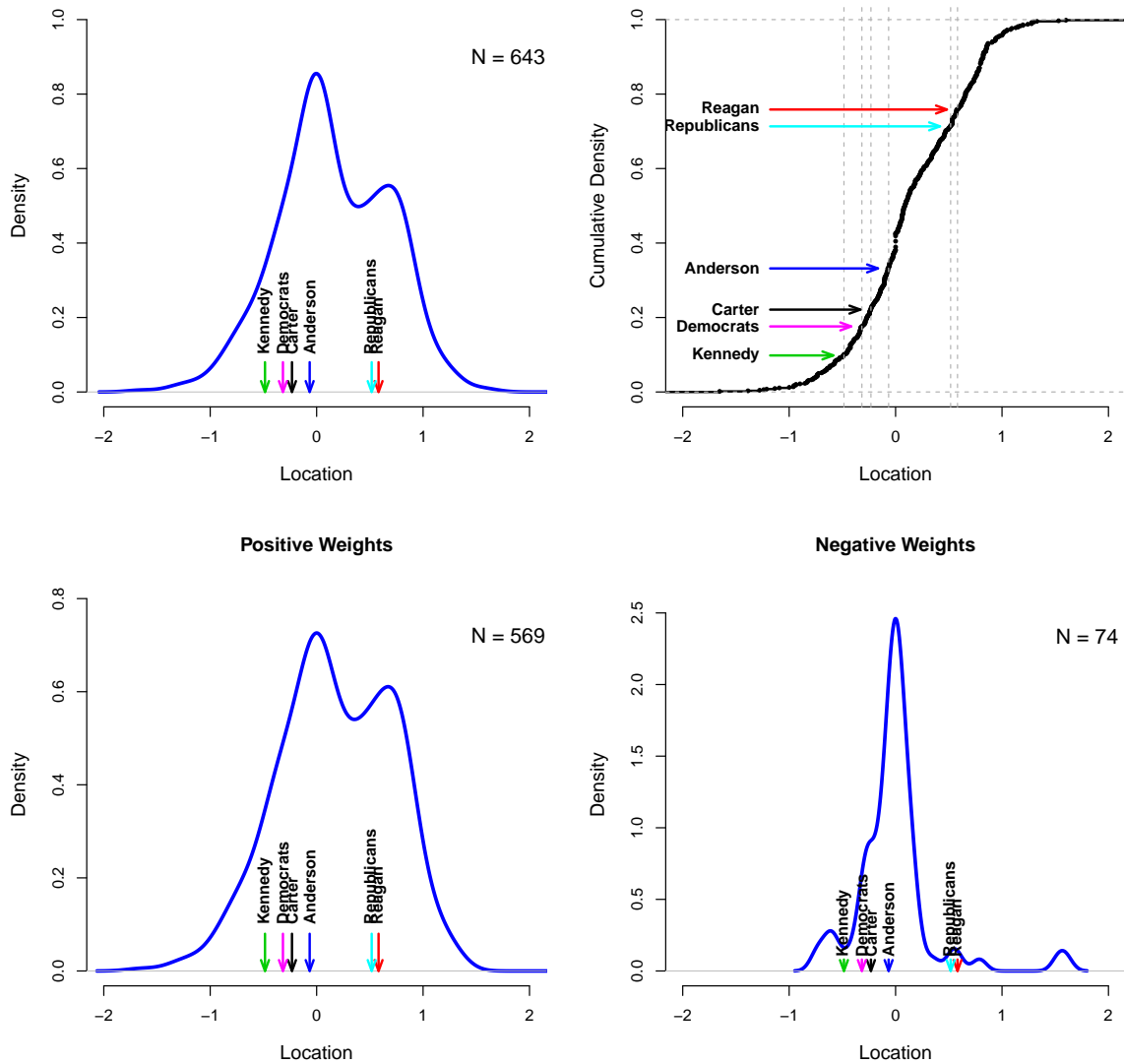


Figure 6: Aldrich-McKelvey plots.

```
R> summary(result)
```

SUMMARY OF ALDRICH-MCKELVEY OBJECT

-----

```
Number of Stimuli: 6
Number of Respondents Scaled: 643
Number of Respondents (Positive Weights): 569
Number of Respondents (Negative Weights): 74
Reduction of normalized variance of perceptions: 0.14
```

```
Location
Kennedy    -0.485
Democrats  -0.317
```

```

Carter      -0.232
Anderson    -0.065
Republicans  0.517
Reagan      0.582

```

Note from this summary that the Aldrich-McKelvey function identifies a number of individuals with negative weights. These represent the set of individuals who see the space “backwards” (i.e., they see Reagan and the Republicans to the left of Carter and the Democrats). In general, the inclusion of a smaller number of these individuals does not significantly affect our estimates, although it can sometimes be helpful to identify them.

```
R> plot.aldmck(result)
```

Estimation of uncertainty for estimates using Aldrich-McKelvey can be obtained via the non-parametric bootstrap (Efron and Tibshirani 1993). To simulate 100 samples from the 1980 liberal-conservative scales and estimate the standard error of the stimuli, we do the following:

```

R> result <- boot_aldmck(data = LC1980, polarity = 2, respondent = 1,
+   missing = c(0, 8, 9), iter = 100)
R> apply(result, 2, sd)

```

```

      Carter      Reagan      Kennedy      Anderson Republicans      Democrats
0.012458959 0.006188856 0.011004480 0.016906648 0.006787753 0.008501512

```

The Aldrich-McKelvey function is primarily intended for scaling perceptual data from surveys, though it can also be used to replicate previously published Monte Carlo results from Palfrey and Poole (1987). Palfrey and Poole find that the Aldrich-McKelvey algorithm is robust in the presence of heteroskedasticity, and test this by replacing the assumed homoskedastic error term  $\epsilon_{ij}$  with a respondent-specific  $\epsilon_i$ . In this example we replicate their result in a single trial, and show that the recovered stimuli  $\hat{Z}_j$  is almost perfectly correlated with the true  $Z_j$  (note that the correlation can be negative because polarity is set randomly).

```

R> Nstimuli <- 6
R> Nresp <- 500
R> Z_j <- rnorm(6)
R> Z_j <- (Z_j - mean(Z_j))/sd(Z_j)
R> respondent.sd <- runif(Nresp, min = 0.3, max = 0.9)
R> error_heteroskedastic <- matrix(NA, Nresp, Nstimuli)
R> for(i in 1:Nresp) error_heteroskedastic <- rnorm(Nstimuli,
+   sd = respondent.sd)
R> w_i <- runif(Nresp, min=0, max=1)
R> c_i <- rnorm(Nresp)
R> Y_ij <- rep(1,500) %o% Z_j
R> Y_ij <- Y_ij + error_heteroskedastic
R> R_ij <- 1/w_i %o% rep(1, Nstimuli) * (Y_ij -
+   c_i %o% rep(1, Nstimuli))
R> result <- aldmck(R_ij, polarity = 6, missing = 999)

```

```
R> cor(Z_j, result$stimuli)
```

```
[1] -0.9995579
```

Although we only show one trial in this paper, the result shown here is reproducible over multiple simulations.<sup>11</sup>

## 7. Conclusion

The `basicspace` package includes a number of functions that enable the estimation of a basic space using self-placement and/or perceptual survey data in R. These include the following functions:

- *Estimation functions*: `aldmck`, `blackbox`, `blackbox_transpose`.
- *Convenience extraction functions*: `individuals`, `fit`, `stimuli`.
- *Generic functions*: `predict`, `plot`, `summary`.
- *Bootstrap functions*: `boot_aldmck`, `boot_blackbt`, and `plot` functions for these objects.

In addition to the functions listed above, three example data sets have also been included. These include the 1980 NES issue scales (`Issues1980`), the 1980 NES liberal-conservative scales (`LC1980`), and the 2004 PELA liberal-conservative scales (`columbia`).

Social scientists often wish to infer the locations of survey respondents – such as voters or legislators – in an abstract policy or ideological space. Perceptual data-oriented estimators such as the basic space technique described here have broad applicability. Given the abundance of perceptual data questions found in most social science surveys, there will continue to be numerous potential applications of the estimators included with this package. An R package that facilitates the analysis of perceptual data in a popular statistics environment will hopefully enable broader use of this method.

## Acknowledgments

This research was supported by a grant from the National Science Foundation (NSF-SBS-0611974). James Lo also acknowledges support from SFB 884, “Political Economy of Reforms”.

## References

Aldrich JH, McKelvey RD (1977). “A Method of Scaling with Applications to the 1968 and 1972 Presidential Elections.” *American Political Science Review*, **71**(1), 111–130. doi: [10.2307/1956957](https://doi.org/10.2307/1956957).

---

<sup>11</sup>Replication of the Monte Carlo test with separate groups of informed and uninformed individuals, from Palfrey and Poole (1987, pg. 515), can be conducted by simply changing the error deviations in the code above to have 250 respondents with  $\sigma_i = 0.3$  and 250 respondents with  $\sigma_i = 0.9$ .

- Carroll JD, Chang JJ (1970). “Analysis of Individual Differences in Multidimensional Scaling via an  $N$ -Way Generalization of Eckart-Young Decomposition.” *Psychometrika*, **35**(3), 283–319. doi:10.1007/bf02310791.
- Downs A (1957). *An Economic Theory of Democracy*. Harper.
- Eckart C, Young G (1936). “The Approximation of One Matrix by Another of Lower Rank.” *Psychometrika*, **1**(3), 211–218. doi:10.1007/bf02288367.
- Efron B, Tibshirani R (1993). *An Introduction to the Bootstrap*, volume 57. Chapman & Hall/CRC.
- Herron M, Lewis J (2007). “Did Ralph Nader Spoil a Gore Presidency? A Ballot-Level Study of Green and Reform Party Voters in the 2000 Presidential Election.” *Quarterly Journal of Political Science*, **2**(3), 205–226. doi:10.1561/100.00005039.
- Martin AD, Quinn KM (2002). “Dynamic Ideal Point Estimation via Markov Chain Monte Carlo for the U.S. Supreme Court, 1953–1999.” *Political Analysis*, **10**(2), 134. doi:10.1093/pan/10.2.134.
- McCarty NM, Poole KT, Rosenthal H (2006). *Polarized America: The Dance of Ideology and Unequal Riches*. MIT Press, Cambridge.
- Palfrey TR, Poole KT (1987). “The Relationship between Information, Ideology, and Voting Behavior.” *American Journal of Political Science*, **31**(3), 511–530. doi:10.2307/2111281.
- Poole K, Lewis J, Lo J, Carroll R (2011). “Scaling Roll Call Votes with `wnominate` in R.” *Journal of Statistical Software*, **42**(14), 1–21. doi:10.18637/jss.v042.i14.
- Poole KT (1998). “Recovering a Basic Space from a Set of Issue Scales.” *American Journal of Political Science*, **42**(3), 954–993. doi:10.2307/2991737.
- Poole KT (2005). *Spatial Models of Parliamentary Voting*. Cambridge University Press.
- Poole KT, Rosenthal H (1997). *Congress: A Political-Economic History of Roll Call Voting*. Oxford University Press, New York.
- R Core Team (2015). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- Saiegh SM (2009). “Recovering a Basic Space from Elite Surveys: Evidence from Latin America.” *Legislative Studies Quarterly*, **34**(1), 117–145. doi:10.3162/036298009787500349.
- Takane Y, Young FW, Leeuw JD (1977). “Nonmetric Individual Differences Multidimensional Scaling: An Alternating Least Squares Method with Optimal Scaling Features.” *Psychometrika*, **42**(1), 7–67. doi:10.1007/bf02293745.
- Voeten E (2001). “Outside Options and the Logic of Security Council Action.” *American Political Science Review*, **95**(4), 845–858. doi:10.1017/s000305540101005x.

**Affiliation:**

Keith T. Poole  
University of Georgia  
Department of Political Science  
Baldwin Hall  
Athens, GA 30602, United States of America  
E-mail: [kpoole@uga.edu](mailto:kpoole@uga.edu)  
URL: <http://www.voteview.com/>

Jeffrey B. Lewis  
University of California, Los Angeles  
Political Science Department, Bunche Hall  
Los Angeles, CA 90095, United States of America  
E-mail: [jblewis@ucla.edu](mailto:jblewis@ucla.edu)  
URL: <http://www.polisci.ucla.edu/faculty/lewis/>

Howard Rosenthal  
New York University  
Department of Politics  
19 W. 4th Street, New York, 10012  
E-mail: [howardrosenthal@nyu.edu](mailto:howardrosenthal@nyu.edu)  
URL: <http://politics.as.nyu.edu/object/HowardRosenthal>

James Lo  
University of Southern California  
Department of Political Science  
3518 Trousdale Parkway, VKC 327  
Los Angeles, CA 90089, United States of America  
E-mail: [jameslo@princeton.edu](mailto:jameslo@princeton.edu)

Royce Carroll  
Rice University  
Department of Political Science, MS 24  
PO Box 1892  
Houston, Texas 77251-1892, United States of America  
E-mail: [rcarroll@rice.edu](mailto:rcarroll@rice.edu)  
URL: <http://rcarroll.web.rice.edu/>