

Empirical Bayes hierarchical models for regularizing maximum likelihood estimation in the matrix Gaussian Procrustes problem

Douglas L. Theobald* and Deborah S. Wuttke

Department of Chemistry and Biochemistry, UCB 215, University of Colorado, Boulder, CO 80309

Edited by Richard V. Kadison, University of Pennsylvania, Philadelphia, PA, and approved September 26, 2006 (received for review September 30, 2005)

Procrustes analysis involves finding the optimal superposition of two or more “forms” via rotations, translations, and scalings. Procrustes problems arise in a wide range of scientific disciplines, especially when the geometrical shapes of objects are compared, contrasted, and analyzed. Classically, the optimal transformations are found by minimizing the sum of the squared distances between corresponding points in the forms. Despite its widespread use, the ordinary unweighted least-squares (LS) criterion can give erroneous solutions when the errors have heterogeneous variances (heteroscedasticity) or the errors are correlated, both common occurrences with real data. In contrast, maximum likelihood (ML) estimation can provide accurate and consistent statistical estimates in the presence of both heteroscedasticity and correlation. Here we provide a complete solution to the nonisotropic ML Procrustes problem assuming a matrix Gaussian distribution with factored covariances. Our analysis generalizes, simplifies, and extends results from previous discussions of the ML Procrustes problem. An iterative algorithm is presented for the simultaneous, numerical determination of the ML solutions.

heteroscedasticity | morphometrics | Procrustes analysis | superpositions | least-squares

The goal of Procrustes analysis is to superpose non-identical shapes in an optimal manner via scalings, translations, and rotations (1–3). Classical Procrustes analysis uses the unweighted least-squares (LS) criterion for finding the optimal transformations. LS implicitly assumes that the landmarks describing the forms’ shapes are uncorrelated and that they have identical variances (i.e., that they are homoscedastic). However, in many practical applications these assumptions are known to be violated. For instance, when superpositioning macromolecular proteins, individual atoms (the landmarks) are connected via covalent chemical bonds, and thus the variance of a given atom can be correlated with the variance of atoms to which it is connected. Furthermore, certain atoms may have larger mobilities than others, or they may have spatial positions with relatively greater uncertainty due to experimental error. Similarly, when comparing the skulls from different members of a species, some homologous features may be highly variable relative to others. Hence, different landmarks can have widely different variances. Under these conditions, ordinary LS can give misleading results, even with large samples of data (4).

The method of maximum likelihood (ML) is a common alternative to LS and is widely considered to be fundamental for statistical modeling and parameter estimation (5). Given the proper model, ML can provide accurate and robust estimates of parameters in the presence of both heteroscedasticity and correlation. Incorporating heterogeneous variances and non-zero correlations into the ML Procrustes problem involves weighting by two covariance matrices (a landmark and a dimensional covariance matrix). However, estimation of these covariance matrices has been a major stumbling block prohibiting a viable ML Procrustes analysis (1, 3, 6–10). Without special assumptions, simultaneous estimation of the landmark covariance matrix and the translations is generally impossible. By parametrically constraining the landmark covariance

matrix using a hierarchical, empirical Bayes likelihood model, we provide a relatively simple solution that permits valid covariance estimation.

In addition to allowing simultaneous estimation of the landmark covariance matrix and the translations, we address several additional challenges that have thus far impeded successful ML inference in the Procrustes problem. (i) As usually described, Procrustes analysis gives degenerate estimates of the covariance matrices even when the translations are known (7). We suggest a method for uniquely identifying the covariance matrices, using the fact that the absolute reference frame in Procrustes analysis is arbitrary. (ii) The dimensional covariance matrix introduces further difficulties for optimal rotation estimation, and current methods involve a circuitous and computationally intensive numerical procedure. We provide a simplified, analytic solution to estimating the optimal rotations when using an arbitrary dimensional covariance matrix. (iii) We give a simplified ML estimator for the scaling factors and also present a hierarchical model for estimating the scalings that can improve estimation in common biological cases. (iv) The ML estimator of the mean form is known to be inconsistent (i.e., it converges to the wrong value as the number of forms increases). Based on the approximate χ^2 distribution of the Procrustes sum of squared distances, we derive an easily calculable bias-corrected ML estimator for the mean. Given that the errors are not large, this estimator is consistent for landmarks with uncorrelated, nonisotropic ellipsoidal variability. (v) Finally, we present a relatively simple and stable iterative algorithm for the simultaneous solution of these unknown parameters.

Formulation of the ML Procrustes Problem

In statistical shape theory, the shape of a geometrical configuration is commonly defined as the geometrical information that is unchanged after translation, rotation, and scaling of the object. Particular instances of an object, which are described by a set of points summarizing the object’s shape, we refer to as a “form.” We also define “size-and-shape” as those geometrical aspects of a form that are unchanged after translation and rotation. In its most general formulation, Procrustes analysis involves the optimal matching of two or more form matrices. Here we consider only sets of forms where each form matrix has the same dimensions.

Consider N different forms (X_i , $i = 1 \dots N$), each with K corresponding points or “landmarks” that depict the form’s shape. Each form is defined as a $K \times D$ matrix holding K rows

Author contributions: D.L.T. and D.S.W. designed research; D.L.T. performed research; D.L.T. contributed new reagents/analytic tools; D.L.T. and D.S.W. analyzed data; and D.L.T. and D.S.W. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS direct submission.

Freely available online through the PNAS open access option.

Abbreviations: LS, least squares; ML, maximum likelihood.

*To whom correspondence should be sent at the present address: Brandeis University, MS013, Waltham, MA 02454. E-mail: dtheobald@brandeis.edu.

© 2006 by The National Academy of Sciences of the USA

of landmarks, where each landmark is a D -vector in \mathbb{R}^D . For example, when superposing protein structures, the forms are independent protein models and the landmarks are the three-dimensional atomic coordinates of each atom in the protein structures. Furthermore, each form is observed in a different, arbitrary, and unknown coordinate system (1, 10, 11). The absolute reference frame for the forms is arbitrary; only relative rotational and translational transformations are of interest. We wish to estimate the optimal orthogonal rotations, translations, and scaling factors that, when applied to this set of form matrices, predict the observed data with the highest likelihood.

Gaussian Perturbation Model with Factored Covariances

To analyze the Procrustes problem from a likelihood perspective, one must choose a statistical model that is assumed to generate the observed data. We assume a perturbation model in which each form X_i is distributed according to a Gaussian (normal) probability density (7, 11). In this Gaussian perturbation model, each X_i can be considered to be an arbitrarily scaled, rotated, and translated Gaussian perturbation of a mean form M

$$X_i = \frac{1}{\beta_i} (M + E_i) R_i' - 1_K \mathbf{t}_i, \quad [1]$$

where R_i is a $D \times D$ rotation matrix, β_i is a scaling parameter, \mathbf{t}_i is a $1 \times D$ row vector for the translational offset, and 1_K denotes the $K \times 1$ column vector of ones. The matrix E_i is a zero-mean matrix of Gaussian random errors $E_i \approx N_{K,D}(0, \Sigma, \Xi)$ with the same dimensions as the mean form M , where Σ and Ξ are covariance matrices described in detail below. In other words, to generate each of the observed forms in D -dimensional Euclidean space, each landmark in an idealized mean form M is randomly perturbed according to a Gaussian distribution. The perturbations for each landmark may be unequal on average and may be correlated with each other, as specified by the covariance matrices. Each perturbed form is then arbitrarily scaled, rotated, and translated, the results of which constitute the observed data.

Covariance matrices describe the variability and correlations among the landmarks in the forms. In the above model, we use factored covariances by assuming that the overall $KD \times KD$ covariance matrix Ω can be decomposed adequately as the Kronecker product of Σ and Ξ (i.e., $\Omega = \Sigma \otimes \Xi$). Here, Σ is a $K \times K$ covariance matrix for the rows of the forms (the landmarks), and Ξ is a constrained $D \times D$ covariance matrix for the columns of the forms (the dimensions). In the protein structure example, the landmark covariance matrix Σ describes the variances and correlations among the atoms. Additionally, the dimensional covariance matrix Ξ describes the variances and correlations between the three coordinate axes, for example, if atomic positions are on average more variable in one dimension than another. We note that, in many cases, the factored covariance model is physically unrealistic, because it assumes that each of the landmarks have the same ellipsoidal principal axes of variability. However, we accept this assumption as a simplifying approximation, and note that it is more general than the assumptions underlying classical LS Procrustes analysis.

Although the product $\Sigma \otimes \Xi$ is uniquely identified, Σ and Ξ independently are not, because for all $c > 0$, $\Omega = \Sigma \otimes \Xi = c\Sigma \otimes 1/c \Xi$. Thus, to jointly identify both Σ and Ξ , we impose the constraint that $\text{tr} \Xi = D$. Furthermore, as a generality in Procrustes analysis, the absolute rotational frame of reference is arbitrary, and only the relative rotations among the forms are of interest. As a result, Ξ is only identifiable up to an arbitrary orthogonal similarity transformation (7). However, any given Ξ can be spectrally decomposed such that $\Xi = V\Delta V'$. Then, the diagonal matrix of eigenvalues Δ is equivalent to the Ξ that would result if the mean form M were rotated with the orthogonal matrix V' of left eigenvectors. Thus, for unique identifiability, we require that the axes of the mean

form's reference frame be aligned with the principle axes of the dimensional covariance matrix Ξ . This requirement also guarantees that Ξ is diagonal.

Procrustes Matrix Gaussian Likelihood Equation

A likelihood analysis requires an explicit statement of the likelihood function derived from the probability density function that describes the assumed statistical model. The full joint likelihood equation for the Gaussian Procrustes model given in Eq. 1 is obtained from a multivariate matrix normal distribution (12, 13). To simplify the appearance of the likelihood equation, we define the squared Frobenius Mahalanobis matrix norm as $\|U\|_{V,W}^2 = \text{tr}\{VU'WU\}$, which can be regarded as the sum of the squared elements of U with columns weighted by V and rows weighted by W . Letting $|U|$ be the determinant of matrix U , the full Procrustes log-likelihood $\ell(\beta, R, \mathbf{t}, M, \Xi, \Sigma; X) = \ell_P$ is given by

$$\ell_P = -\frac{1}{2} \sum_i^N \|\beta_i(X_i + 1_K \mathbf{t}_i) R_i - M\|_{\Xi^{-1}, \Sigma^{-1}}^2 - \frac{NKD}{2} \ln(2\pi) - \frac{NK}{2} \ln|\Xi| - \frac{ND}{2} \ln|\Sigma|. \quad [2]$$

The Procrustes log-likelihood equation above (Eq. 2) is the key conceptual contribution underlying the ML analysis that follows. In order to solve the ML Procrustes problem, ML estimators must be derived for each of the unknown parameters in the Procrustes log-likelihood equation. We provide the individual ML estimates for each of these unknowns: the Ξ , Σ , M , and the N β_i , \mathbf{t}_i , and R_i (for derivations see *Appendices*).

Regularization of the Covariance Matrices Allows Simultaneous Estimation

Estimation of the factored covariance matrices has been the primary impediment for a successful, practical implementation of a joint non-isotropic ML treatment of the Procrustes problem (1, 3, 6–10). We first present the unrestricted ML estimators of the covariance matrices and then discuss the difficulties that our approach remedies.

The unrestricted ML estimators of the covariance matrices, $\hat{\Xi}_U$ and $\hat{\Sigma}_U$, are given by straightforward variants of the usual covariance estimators for the matrix normal distribution (*Appendix I*) (4, 13).

$$\hat{\Sigma}_U = \frac{1}{ND} \sum_i^N (\beta_i \check{X}_i R_i - \hat{M}) \Xi^{-1} (\beta_i \check{X}_i R_i - \hat{M})' \quad [3]$$

$$\hat{\Xi}_U = \frac{1}{NK} \sum_i^N (\beta_i \check{X}_i R_i - \hat{M})' \Sigma^{-1} (\beta_i \check{X}_i R_i - \hat{M}), \quad [4]$$

where \check{X} is a row-weighted centered form.[†] Because one of the covariance matrices must be normalized for identifiability, we (arbitrarily) choose to normalize $\hat{\Xi}_U$ by

$$\tilde{\Xi} = \frac{D \hat{\Xi}_U}{\text{tr} \hat{\Xi}_U}, \quad [5]$$

and use $\tilde{\Xi}^{-1}$ in place of Ξ^{-1} in calculating $\hat{\Sigma}_U$ given in Eq. 3.

[†]Row-wise weighted centering can be applied to an uncentered structure X_i by $\check{X}_i = X_i + 1_K \mathbf{t}_i$, where \mathbf{t}_i is the ML estimate of the translation. The ML translation is independent of the rotations R_i , the scaling factors β_i , and the dimensional covariance matrix Ξ , and it is given by $-(1_K \Sigma^{-1} X_i / 1_K \Sigma^{-1} 1_K)$, where X_i is an uncentered structure (11).

Both of these covariance matrix estimators involve an inverse of a covariance matrix, as do several of the other ML solutions, a fact that presents three important and distinct difficulties for ML estimation in the Procrustes problem. First, in general, the ML estimators of the matrix normal covariance matrices only exist when the estimators are invertible, a requirement that obtains when $N \geq (K/D) + 1$ and $N \geq (D/K) + 1$ (13). In practice, these conditions are usually unsatisfied due to a paucity of available data. Second, Procrustes analysis requires centering the form matrices (e.g., Eq. 6), which imparts a common linear constraint on the columns of the forms. Even with ample data, then, the sample landmark covariance matrix is collinear and rank deficient, with at least one zero eigenvalue (3, 6), and consequently it is noninvertible. Third, and most fundamentally, the displacements due to Σ and the translations t_i are linearly entangled (see Eq. 1), and therefore unrestricted simultaneous estimation of both is impossible (6, 10). This latter problem is especially grave, because the unrestricted Procrustes estimators of the translations t_i and of the covariance matrix Σ are strongly interdependent; translations that closely superpose a given landmark decrease the estimated variance, and a small variance in turn preferentially weights the translations to center on that landmark. For example, it is always possible to find translation vectors that will perfectly superpose any selected landmark, thereby rendering its sample variance exactly zero, $\hat{\Sigma}$ singular, and the likelihood infinite.

To overcome the first two of these problems it is sufficient to guarantee the invertibility of the estimated landmark covariance matrix, i.e., all eigenvalues must be positive. In the third difficulty, the possibility of an infinite likelihood immediately points to a need for regularization, either of the translations (6) or of the covariance matrix, or of both.

In real data sets, the variances often are unable to take arbitrary values, but rather they describe coupled components of a larger, integrated system. For instance, the atoms in a macromolecule are linked by chemical bonds, and the bones in a skull are physically joined by ligaments and at sutures. In such cases, the variances for each landmark are similar in magnitude, and extremely small or large variances are improbable and physically unrealistic. Thus, the problem of simultaneous estimation of the covariance matrices and the translations can be circumvented by realistically restricting the variances of the landmarks.

An alternative, then, is to view each of the eigenvalues (λ_j) of the landmark covariance matrix as a draw from a positive distribution (where $\lambda_j > 0$). The eigenvalues of the covariance matrix can be considered as variances with all correlations eliminated, e.g., the eigenvalues of a diagonal covariance matrix are the variances themselves. The inverse gamma ($I\gamma$) distribution is used conventionally in Bayesian statistics as a natural proper prior for the variances in a multivariate Gaussian model (14). The flexibility of the $I\gamma$ distribution is convenient as it can adopt a wide variety of shapes approximating other distributions. Therefore, we regularize our matrix Gaussian Procrustes model using a hierarchical empirical Bayes approach (an extended likelihood) (5, 15) by assuming that the eigenvalues of the landmark covariance matrix belong to a common $I\gamma$ distribution. The full, regularized Procrustes log-likelihood (ℓ_r) is then the sum of the Procrustes log-likelihood from Eq. 2 and the log-likelihood of an $I\gamma$ distribution of the eigenvalues (see Eq. 15 in Appendix II): $\ell_r = \ell_p + \ell(\lambda)_{I\gamma}$.

Assuming an $I\gamma$ distribution for the eigenvalues of the landmark covariance matrix (with scale parameter α and shape parameter γ), the extended ML estimator $\hat{\Sigma}_{I\gamma}$ of Σ is a simple linear function of the unrestricted ML estimate $\hat{\Sigma}_U$ from Eq. 3 (see Appendix II)

$$\hat{\Sigma}_{I\gamma} = \frac{ND}{ND + 2(\gamma + 1)} \left(\frac{2\alpha}{ND} \mathbf{I} + \hat{\Sigma}_U \right). \quad [6]$$

In this hierarchical model, the $I\gamma$ parameters α and γ are point estimates determined by the data, unlike when using a bona fide Bayesian prior. Eq. 6 can be interpreted as a shrinkage estimator that contracts the eigenvalues of the covariance matrix towards the mode of the inverse gamma distribution. Given positive α and γ parameters, this regularized model constrains all eigenvalues (and variances) of the landmark covariance matrix to be positive, thereby guaranteeing that $\hat{\Sigma}_{I\gamma}$ is invertible. Furthermore, this treatment enables simultaneous estimation of the landmark covariance matrix Σ and the translations by decoupling the strong dependence of their estimators. Thus, the inverse gamma regularization is sufficient to overcome all three of the difficulties listed above. In our experience, the $\hat{\Sigma}_{I\gamma}$ estimator is also well conditioned. We note that our regularized hierarchical estimator of the multivariate Gaussian covariance matrix also has potential for wide applicability beyond Procrustes analysis.

Simplified ML Estimates of the Rotations

Maximizing the Procrustes likelihood with respect to the rotations is equivalent to maximizing the first term of Eq. 2, which requires weighting by the inverse of the dimensional covariance matrix Ξ . Unfortunately, there is apparently no general algebraic solution to this maximization problem using an arbitrary, symmetric, dimensional weight matrix (2, 11, 16). Koschat and Swayne have proposed a solution for an arbitrary diagonal $D \times D$ weight matrix using a computationally intensive, iterative numerical algorithm assuming that $\Sigma = \mathbf{I}$ (16). However, the Koschat–Swayne algorithm is reported to have poor convergence properties and is sensitive to the seed values (16).

In our ML Procrustes treatment, the dimensional weight matrix is not arbitrary; rather, we weight specifically by the inverse of the dimensional covariance matrix (i.e., Ξ_U^{-1}). In this case, there is a much simpler, analytic solution to maximizing the estimated ML target for the optimal rotation. Here the dimensionally weighted Procrustes problem reduces conveniently (see Appendix III) to maximizing

$$\sum_i^N \text{tr}(\Xi^{-1} \hat{M}' \hat{\Sigma}^{-1} \check{X}_i R_i). \quad [7]$$

The rotations that maximize Eq. 7 can be found with the usual Procrustes singular value decomposition (SVD) solution (11). Let the SVD of an arbitrary matrix D be $U\Lambda V'$. Then the maximizing rotations \hat{R}_i are given by

$$\begin{aligned} \Xi^{-1} \hat{M}' \hat{\Sigma}^{-1} \check{X}_i &= U\Lambda V' \\ \hat{R}_i &= VPU', \end{aligned} \quad [8]$$

where rotoinversions can be avoided by constraining the determinant of \hat{R}_i to be 1 by using $P = \mathbf{I}$ if $|V||U| = 1$ or $P = \text{diag}(1, \dots, 1, -1)$ if $|V||U| = -1$.

ML Estimation of the Scaling Factors

When superposing and comparing the shapes of different sized objects, scaling factors must be estimated to remove differences due to size alone. The unconstrained ML estimates of the scaling parameters β_i are given by

$$\beta_i = \frac{\text{tr}(\Xi^{-1} \hat{M}' \hat{\Sigma}^{-1} \check{X}_i R_i)}{\text{tr}(\Xi^{-1} R_i' \check{X}_i' \hat{\Sigma}^{-1} \check{X}_i R_i)}. \quad [9]$$

However, unlike the rotations which are generally arbitrary with real data, it may be reasonable to model the scaling factors themselves as random variables drawn from a relevant probability distribution. For instance, the log-normal distribution arises as the limiting distribution of many small, multiplicative

random effects, as may be expected for random deviations in size (17). Many biological entities subject to random variations in growth naturally fit a log-normal distribution. Assuming that $\beta_i \approx LN(\theta, \theta)$, where

$$P(\beta_i) = (2\pi\theta)^{-\frac{1}{2}}\beta_i^{-1}e^{-\frac{(\ln \beta_i - \theta)^2}{2\theta}}, \quad [10]$$

then the probabilities $P(\beta_i) = P(1/\beta_i)$, and the mode of the distribution is 1 as required. The ML estimate $\hat{\beta}_i$ for the lognormally distributed scaling factors is the positive root of (Appendix IV)

$$\beta_i^2 \text{tr}(\Xi^{-1}R_i\tilde{X}_i'\Sigma^{-1}\tilde{X}_iR_i) - \beta_i \text{tr}(\Xi^{-1}M'\Sigma^{-1}\tilde{X}_iR_i) + \frac{\ln \beta_i}{\theta} = 0. \quad [11]$$

The $(\ln \beta_i)/\theta$ term in Eq. 11 can be thought of as a “penalty term” that concentrates the $\hat{\beta}$ estimates around unity; otherwise it is identical to the nonparametric solution (Eq. 9).

ML Estimate of the Mean Form

Maximizing the likelihood with respect to the mean form M gives the familiar average (11)

$$\hat{M} = \bar{X} = \frac{1}{N} \sum_i \beta_i \tilde{X}_i R_i. \quad [12]$$

However, the average \bar{X} is an asymptotically biased estimate of the mean form M (it is too large) for constant K when superpositioning without scaling (i.e., setting $\beta_i = 1$ for all i) (1, 7, 18). For the case when Ξ and Σ are diagonal (i.e., no correlations), we offer the following easily computable, bias-corrected ML Procrustes estimator \tilde{M} of the mean (Appendix V)

$$\tilde{M} = \bar{X} \sqrt{1 - \frac{D^2 + D}{2\text{tr}(\Xi^{-1}\bar{X}'\Sigma^{-1}\bar{X})}}. \quad [13]$$

This bias-corrected estimate is consistent (for constant K as $N \rightarrow \infty$) provided that the errors are not very large, specifically when $\text{tr}(\Xi^{-1}\bar{X}'\Sigma^{-1}\bar{X}) \geq (D^2 + D)/2$. For three-dimensional forms, for example, this estimator is valid when the landmark variances are less than $K/6$ times the squared radius of gyration of the average form.

An Iterative Algorithm for ML Procrustes Analysis

In practice, the ML solutions given above must be solved simultaneously via numerical methods, because each of the unknown parameters is a function of the others. We have developed the following iterative algorithm (11, 13, 19) for determining solutions to each of the Procrustes ML estimators.

1. **Initialize.** Set $\hat{\Sigma} = I$, $\hat{\Xi} = I$, and $\beta_i = 1$ for all i ; choose one of the X_i to serve as \hat{M} .
2. **Translate.** Obtain each \tilde{X}_i by centering X_i according to footnote †.
3. **Rotate.** Calculate each \hat{R}_i according to Eq. 8, and set each $X_i = \tilde{X}_i \hat{R}_i$.
4. **Scale.** For lognormally distributed scale factors, first find the lognormal shape/location parameter θ with Eq. 19. Newton–Raphson can then be used to find the lognormal scale estimates with an initial guess equal to the arbitrary scale factor estimate Eq. 9 or to the scale factor from the previous iteration.
5. **Estimate the mean.** Find the new average \hat{M} according to Eq. 12 or 13. Return to step 3 and loop to convergence.
6. **Estimate the inverse gamma distributed eigenvalues.** Calculate $\hat{\Sigma}_U$ from Eq. 3. The $\hat{\Sigma}_U$ covariance matrix must be decomposed, and the positive eigenvalues fit iteratively to an

inverse gamma distribution. The rank of the sample covariance matrix is $\min(ND - D - 3, K - 3)$. At convergence, the rank is reduced by two due to the collinearity induced by the linear constraints from the scale and shape parameters of the inverse gamma regularization (see Eq. 6). The rank is further reduced (by one) due to additional collinearity from centering the form matrices (6). Thus, the maximally obtainable number of positive eigenvalues in our hierarchical treatment is $K - 3$; this rank deficiency is inherent in the Procrustes sample covariance matrix and cannot be eliminated by increasing the number of forms. Thus, we treat all zero eigenvalues as missing data using an Expectation–Maximization algorithm and include only the full-rank, positive eigenvalues of the sample covariance matrix when estimating the parameters of the inverse gamma distribution. If the covariance matrix is assumed to be diagonal, one must include only the $K - 3$ largest variances in the inverse gamma fit, as the smallest three are known *a priori* to be zero-valued eigenvalues. Estimation of the regularized eigenvalues thus cycles until convergence between two steps: (i) fit the current estimates of the full-rank eigenvalues to an inverse gamma distribution, and (ii) obtain updated estimates of the eigenvalues by modification of the sample eigenvalues analogous to Eq. 6. $\lambda_{I\gamma,j} = (ND\lambda_{U,j} + 2\alpha)/(2(1 + \gamma) + ND)$.

7. **Estimate the landmark covariance matrix.** Calculate $\hat{\Sigma}_{I\gamma}$ from Eq. 6 using the $\hat{\alpha}$ and $\hat{\gamma}$ inverse gamma parameters determined in the previous step.
8. **Estimate the dimensional covariance matrix.** Calculate $\hat{\Xi}$ from Eq. 5. Return to step 6 and loop to convergence.
9. **Align the superposition with $\hat{\Xi}$.** Let the spectral decomposition of $\hat{\Xi}$ be $V\Delta V'$. Rotate the entire superposition, including the average form \hat{M} , with the matrix V' of left eigenvectors (i.e. set $\tilde{M} = \hat{M}V'$) so that $\hat{\Xi}$ is equivalent to the diagonal matrix Δ .
10. **Loop.** Return to step 2 until convergence (e.g., stop if the relative difference between rotation matrices between iterations fall below a specified tolerance, such as 10^{-7}).

We have implemented this algorithm in THESEUS, a UNIX C program intended for likelihood analysis of biological macromolecules (25). Because of the complexity of the problem, the efficiency of the algorithm is difficult to determine generally. If Σ is assumed to be diagonal, then the algorithm is very efficient ($O(NK)$). In our experience, even with large three-dimensional problems, when assuming a diagonal landmark covariance matrix (e.g., $K = 200$ and $N = 500$), the algorithm converges rapidly after a few dozen iterations (from milliseconds to a few seconds on a 500 MHz G4 laptop). When using full, arbitrary covariance matrices on a low-dimensional problem, the algorithm is $O(K^3)$, as the limiting step is the inversion of the landmark covariance matrix $\hat{\Sigma}$; in this case, the calculation may take a few minutes.

Comparison with Previous Analyses

Our likelihood analysis offers several ML Procrustes estimators, including the covariance matrices ($\hat{\Xi}$ and $\hat{\Sigma}$), a bias-corrected average form (\tilde{M}), lognormally distributed scaling factors ($\hat{\beta}_i$), and simplified estimates of the rotations (\hat{R}_i) when weighting the dimensions. Several of our solutions differ significantly from corresponding solutions to similar problems proposed previously in the Procrustes literature (1, 7, 8, 10, 11, 16). Unlike the unconstrained covariance estimators reported in Goodall and recounted by others (1, 7, 8, 11), each of our covariance estimators is itself a function of the other covariance matrix. Equations 10.2 and 10.3 of Goodall (11) are only valid when either Ξ or Σ is constrained to be equal to the identity matrix. Additionally, to find the optimal rotations when weighting with a factored dimensional covariance matrix, Goodall proposes a method using the Koschat–Swayne algorithm (16). Goodall’s

The log-likelihood for the row-wise covariance matrix Σ is

$$\ell(\Sigma) = -\frac{1}{2} \sum_i^N \text{tr}(\mathbf{E}_i \Xi^{-1} \mathbf{E}_i' \Sigma^{-1}) - \frac{ND}{2} \ln |\Sigma| + C.$$

Taking the derivative with respect to Σ and setting it equal to zero gives

$$\begin{aligned} \frac{\partial \ell(\Sigma)}{\partial \Sigma} &= \sum_i^N \Sigma^{-1} \mathbf{E}_i \Xi^{-1} \mathbf{E}_i' \Sigma^{-1} - \frac{1}{2} \mathbf{I} \odot \left(\sum_i^N \Sigma^{-1} \mathbf{E}_i \Xi^{-1} \mathbf{E}_i' \Sigma^{-1} \right) \\ &\quad - ND \Sigma^{-1} + \frac{ND}{2} \mathbf{I} \odot \Sigma^{-1} = 0 \\ &= \sum_i^N \Sigma^{-1} \mathbf{E}_i \Xi^{-1} \mathbf{E}_i' \Sigma^{-1} - ND \Sigma^{-1} = 0, \end{aligned}$$

where $\mathbf{A} \odot \mathbf{B}$ denotes the Hadamard (element-wise) product of two matrices \mathbf{A} and \mathbf{B} . Thus, after simplification and multiplication by Σ twice, we have $\hat{\Sigma} = (1/ND) \sum_i^N \mathbf{E}_i \Xi^{-1} \mathbf{E}_i'$, which is equivalent to Eq. 3. Derivation of the ML estimator of Ξ proceeds similarly.

Appendix II: Inverse Gamma Distributed Variances. The probability distribution function for an inverse gamma distribution of the eigenvalues of the row-wise covariance matrix is

$$P(\lambda_j) = \frac{\alpha^\gamma}{\Gamma(\gamma)} \lambda_j^{-(1+\gamma)} e^{-\frac{\alpha}{\lambda_j}}.$$

The corresponding log-likelihood of the eigenvalues is

$$\begin{aligned} \ell(\lambda)_{I_\gamma} &= -(1 + \gamma) \sum_j^K \ln \lambda_j - \alpha \sum_j^K \frac{1}{\lambda_j} + K\gamma \ln \alpha - K \ln \Gamma(\gamma) \\ &= -(1 + \gamma) \ln |\Sigma| - \alpha \text{tr} \Sigma^{-1} + K\gamma \ln \alpha - K \ln \Gamma(\gamma). \end{aligned} \quad [14]$$

The unrestricted log-likelihood for Σ can be written as (compare Appendix I)

$$\ell(\Sigma) = -\frac{ND}{2} \sum_i^N \text{tr}(\hat{\Sigma}_U \Sigma^{-1}) - \frac{ND}{2} \ln |\Sigma| + C.$$

The Procrustes extended log-likelihood (5, 15) assuming inverse gamma distributed eigenvalues $\ell(\Sigma)_{I_\gamma}$ is simply the sum $\ell(\lambda)_{I_\gamma} + \ell(\Sigma)$. Taking the derivative of $\ell(\Sigma)_{I_\gamma}$ with respect to Σ and setting it equal to zero gives

$$\begin{aligned} \frac{\partial \ell(\Sigma)_{I_\gamma}}{\partial \Sigma} &= -(1 + \gamma) \Sigma^{-1} + \alpha \Sigma^{-2} \\ &\quad + \frac{ND}{2} \hat{\Sigma}^{-1} \hat{\Sigma}_U \Sigma^{-1} - \frac{ND}{2} \Sigma^{-1} = 0, \end{aligned}$$

which, after pre- and postmultiplication by Σ^{-1} and rearrangement, gives Eq. 6.

Appendix III: ML Rotations with Dimensional Weighting by Ξ^{-1} . If all structures have been centered, the likelihood equations simplify. For centered form matrices, the log-likelihood with respect to the rotations is

$$\ell(\mathbf{R}) = -\frac{1}{2} \sum_i^N \text{tr} \{ \Xi^{-1} (\beta_i \check{\mathbf{X}}_i \mathbf{R}_i - \mathbf{M})' \Sigma^{-1} (\beta_i \check{\mathbf{X}}_i \mathbf{R}_i - \mathbf{M}) \} + C. \quad [15]$$

The nuisance parameters Ξ , Σ , and \mathbf{M} are in general unknown, and as a consequence the likelihood given in Eq. 15 cannot be maximized directly. Therefore, as in Goodall (11), we use the method of estimated likelihood (5, 21, 22) to construct a modified target likelihood function by replacing these nuisance parameters with their ML estimates ($\hat{\Xi}$, $\hat{\Sigma}$, and $\hat{\mathbf{M}}$, respectively). To simplify the derivation here, we choose Σ to be the normalized member of the factored covariance matrices. After substitution with ML estimates, Eq. 15 can be expanded as the estimated likelihood $\ell_e(\mathbf{R}) = \ell(\mathbf{R}, \hat{\Xi}, \hat{\Sigma}, \hat{\mathbf{M}})$

$$\begin{aligned} \ell_e(\mathbf{R}) &= \beta_i \sum_i^N \text{tr}(\hat{\Xi}_U^{-1} \hat{\mathbf{M}}' \hat{\Sigma}^{-1} \check{\mathbf{X}}_i \mathbf{R}_i) \\ &\quad - N \text{tr}(\hat{\Xi}_U^{-1} \hat{\mathbf{M}}' \hat{\Sigma}^{-1} \hat{\mathbf{M}}) - \frac{1}{2} \text{tr} \\ &\quad \left\{ \hat{\Xi}_U^{-1} \sum_i^N (\beta_i^2 \mathbf{R}_i' \check{\mathbf{X}}_i' \hat{\Sigma}^{-1} \check{\mathbf{X}}_i \mathbf{R}_i - \hat{\mathbf{M}}' \hat{\Sigma}^{-1} \hat{\mathbf{M}}) \right\} + C. \end{aligned} \quad [16]$$

After substitution with Eq. 12 for the estimate of the mean form $\hat{\mathbf{M}}$, the ML estimate of the dimensional covariance matrix $\hat{\Xi}_U$ can be expressed as

$$\hat{\Xi}_U = \frac{1}{NK} \sum_i^N (\beta_i^2 \mathbf{R}_i' \check{\mathbf{X}}_i' \hat{\Sigma}^{-1} \check{\mathbf{X}}_i \mathbf{R}_i - \hat{\mathbf{M}}' \hat{\Sigma}^{-1} \hat{\mathbf{M}}). \quad [17]$$

Substitution of Eq. 17 into Eq. 16 gives

$$\begin{aligned} \ell_e(\mathbf{R}) &= \beta_i \sum_i^N \text{tr}(\hat{\Xi}_U^{-1} \hat{\mathbf{M}}' \hat{\Sigma}^{-1} \check{\mathbf{X}}_i \mathbf{R}_i) \\ &\quad - N \text{tr}(\hat{\Xi}_U^{-1} \hat{\mathbf{M}}' \hat{\Sigma}^{-1} \hat{\mathbf{M}}) - \frac{NK}{2} \text{tr}(\hat{\Xi}_U^{-1} \hat{\Xi}_U) + C. \end{aligned} \quad [18]$$

In the final term of Eq. 18 the two matrices cancel reducing to $NKD/2$. Hence, when maximizing the estimated likelihood $\ell_e(\mathbf{R})$ with respect to the rotations (holding other nuisance parameters fixed), both the last two terms of Eq. 18 are constant. Maximization of the first term alone (as given in Eq. 8) provides a ML estimate of the rotation \mathbf{R}_i .

Appendix IV: ML Procrustes Estimate of Lognormally Distributed Scale Factors. In the lognormal distribution given in Eq. 11, the shape and location parameter both equal θ to constrain the mode of the distribution to 1. Given a sample of N scale factors β_i , the ML estimate of θ is

$$\hat{\theta} = \frac{1}{2} \left(\sqrt{1 + \frac{4}{N} \sum_i^N (\ln \beta_i)^2} - 1 \right). \quad [19]$$

The log-likelihood for the scale factors β given a lognormal $\text{LN}(\theta, \theta)$ distribution is

$$\ell(\beta_i)_{LN} = -\frac{1}{2} \ln \theta - \frac{\theta}{2} - \frac{(\ln \beta_i)^2}{2\theta}.$$

The Procrustes extended log-likelihood $\ell_{P,LN}$ (5, 15) assuming lognormally distributed scale factors is then given by the sum of $\ell(\beta_i)_{LN}$ and ℓ_P (from Eq. 2). With respect to β_i the joint log-likelihood is

$$\begin{aligned} \ell(\beta_i)_{P,LN} = & -\ln \beta_i - \frac{1}{2} \ln \theta - \frac{\theta}{2} - \frac{(\ln \beta_i)^2}{2\theta} \\ & - \frac{1}{2} \text{tr}(\beta_i^2 \Xi^{-1} \mathbf{R}'_i \check{\mathbf{X}}'_i \Sigma^{-1} \check{\mathbf{X}}_i \mathbf{R}_i) \\ & + \frac{1}{2} \text{tr}(2\beta_i \Xi^{-1} \mathbf{M}' \Sigma^{-1} \check{\mathbf{X}}_i \mathbf{R}_i) + C. \end{aligned}$$

Taking the derivative of $\ell(\beta)_{P,LN}$ with respect to β_i , setting it equal to zero, and multiplying through by $-\beta_i$ yields Eq. 11.

Appendix V: Bias-Corrected ML Procrustes Estimate of the Mean Form.

In the following, we assume that Ξ and Σ are both diagonal, and that $\beta_i = 1$ for all i . The sum of squared distances from the mean for each landmark (G_j) in a set of optimally superimposed forms is given by

$$G_j = \sum_i^N (x_{i,j} \Xi^{-1} x'_{i,j}) - N \bar{x}_j \Xi^{-1} \bar{x}'_j, \quad [20]$$

where $x_{i,j}$ is the j th $1 \times D$ row vector in the optimally translated and rotated form matrix $\check{\mathbf{X}}_i \mathbf{R}_i$, such that $\check{\mathbf{X}}_i \mathbf{R}_i = [x'_{i,1} \dots x'_{i,j}]'$. It can be shown from standard multivariate analysis that (12)

$$\frac{1}{N} \sum_i^N x_{i,j} \Xi^{-1} x'_{i,j} \rightarrow \mu_j \Xi^{-1} \mu'_j + D \sigma_j^2 \text{ as } N \rightarrow \infty, \quad [21]$$

where similarly μ_j is the j th row vector in the mean form matrix \mathbf{M} . Assuming that G_j is approximately χ^2 distributed (1, 23, 24), then $G_j \approx \sigma_j^2 \chi_{F_j}^2$ with F_j degrees of freedom. For each of the K landmark row vectors in the N forms there are D degrees of freedom. However, for each of the N forms, D degrees of freedom are lost due to translational estimation, and $\frac{1}{2} D(D-1)$ degrees are lost due to rotational estimation (1, 24). Thus, by distributing the lost degrees of freedom over each of the K row vectors in proportion to the contribution of each landmark to the minimization criterion ($w_j = \mu_j \Xi^{-1} \mu'_j / [\sigma_j^2 \text{tr}(\Xi^{-1} \mathbf{M}' \Sigma^{-1} \mathbf{M})]$), we have

$$F_j = N \left(D - \frac{\mu_j \Xi^{-1} \mu'_j (D^2 + D)}{2 \sigma_j^2 \text{tr}(\Xi^{-1} \mathbf{M}' \Sigma^{-1} \mathbf{M})} \right).$$

It then follows trivially that (17)

$$\frac{G_j}{N} \rightarrow \sigma_j^2 \left(D - \frac{\mu_j \Xi^{-1} \mu'_j [D^2 + D]}{2 \sigma_j^2 \text{tr}(\Xi^{-1} \mathbf{M}' \Sigma^{-1} \mathbf{M})} \right) \text{ as } N \rightarrow \infty \quad [22]$$

Combining Eqs. 20–22 yields the asymptotic result as $N \rightarrow \infty$

$$\bar{x}_j \Xi^{-1} \bar{x}'_j \rightarrow \mu_j \Xi^{-1} \mu'_j \left(1 + \frac{D^2 + D}{2 \text{tr}(\Xi^{-1} \mathbf{M}' \Sigma^{-1} \mathbf{M})} \right). \quad [23]$$

The final factor on the right hand side of Eq. 23 is an estimate of the asymptotic bias of the average. Setting $\mu_j = B \bar{x}_j$ (implying $\mathbf{M} = B \bar{\mathbf{X}}$) and solving for B gives Eq. 13.

We thank Meredith Betterton, Ian Dryden, Colin Goodall, Subhash Lele, and Norm Pace for helpful comments and suggestions. This work was supported by the Arnold and Mable Beckman Foundation and National Institutes of Health Grant GM59414. D.L.T. was supported by Postdoctoral Fellowship Grant PF-04-118-01-GMC from the American Cancer Society.

1. Dryden IL, Mardia KV (1998) *Statistical Shape Analysis* (Wiley, New York).
2. Gower JC, Dijksterhuis GB (2004) *Procrustes Problems* (Oxford Univ Press, Oxford), Vol 30.
3. Lele S, Richtsmeier JT (2001) *An Invariant Approach to Statistical Analysis of Shapes* (Chapman and Hall/CRC, Boca Raton, FL).
4. Mardia K, Goodall C (1993) in *Multivariate Environmental Statistics*, eds Patil G, Rao C (North-Holland, New York), Vol 6, pp 347–385.
5. Pawitan Y (2001) *In All Likelihood: Statistical Modeling and Inference Using Likelihood*, Oxford Science Publications (Clarendon, Oxford).
6. Lele S, Richtsmeier JT (1990) *Syst Zool* 39:60–69.
7. Lele S (1993) *Math Geol* 25:573–602.
8. Glasbey C, Horgan G, Gibson G, Hitchcock D (1995) *Biometrical J* 37:481–495.
9. Goodall C (1991) *J Roy Stat Soc B Met* 53:334–339.
10. Goodall C (1995) in *Proceedings in Current Issues in Statistical Shape Analysis*, eds Mardia KV, Gill CA (Leeds Univ Press, Leeds, UK), pp 18–33.
11. Goodall C (1991) *J Roy Stat Soc B Met* 53:285–321.
12. Arnold SF (1981) *The Theory of Linear Models and Multivariate Analysis* (Wiley, New York).
13. Dutilleul P (1999) *J Stat Comput Sim* 64:105–123.
14. Broemeling L (1985) *Bayesian Analysis of Linear Models* (Dekker, New York).
15. Bjornstad JF (1996) *J Am Stat Assoc* 91:791–806.
16. Koschat MA, Swayne DF (1991) *Psychometrika* 56:229–239.
17. Evans M, Hastings N, Peacock JB (2000) *Statistical Distributions* (Wiley, New York), 3rd Ed.
18. Kent JT, Mardia KV (1997) *J Roy Stat Soc B Met* 59:281–290.
19. Dempster AP, Laird NM, Rubin DB (1977) *J Roy Stat Soc B Met* 39:1–38.
20. Dwyer PS (1967) *J Am Stat Assoc* 62:607–625.
21. Royall R (1991) *Statistical Evidence: A Likelihood Paradigm* (Chapman and Hall, New York), Vol 71.
22. Royall R, Tsou T-S (1995) *J Am Stat Assoc* 90:316–320.
23. Langron SP, Collins AJ (1985) *J Roy Stat Soc B Met* 47:277–284.
24. Sibson R (1979) *J Roy Stat Soc B Met* 41:217–229.
25. Theobald DL, Wuttke DS (2006) *Bioinformatics* 22:2171–2172.