



## DERIVATIVES OF THE PROCRUSTUS TRANSFORMATION WITH APPLICATIONS

JAN DE LEEUW AND KEKONA KAI SORENSON

ABSTRACT. We derive expressions for the derivatives of the Procrustus Transformation. As byproducts we find the derivatives of singular values and vectors and of the direct and inverse symmetric square root. The results are applied to minimization algorithms over Stiefel manifolds, and in particular to Polynomial Component and Factor Analysis.

### 1. PROCRUSTUS TRANSFORMATION

Suppose  $X$  is an  $n \times m$  matrix, with  $n \geq m$  and  $\mathbf{rank}(X) = m$ . The *Procrustus Transformation*  $\Pi(X)$  is the  $n \times m$  orthonormal matrix that minimizes the sum-of-squares  $\mathbf{SSQ}(Y - X)$ , or, equivalently, the  $n \times m$  orthonormal matrix that maximizes  $\mathbf{tr} Y'X$ . See Cliff [1966].

The stationary equation for this optimization problem are

$$(1a) \quad X = YM,$$

$$(1b) \quad Y'Y = I,$$

with  $M$  a symmetric matrix of Lagrange Multipliers. It follows that  $M^2 = X'X$ , i.e.  $M = (X'X)^{\frac{1}{2}}$ , the unique symmetric square root of  $X'X$ . And thus  $\Pi(X) = X(X'X)^{-\frac{1}{2}}$ .

---

*Date:* Saturday 14<sup>th</sup> April, 2012 — 15h 9min — Typeset in LUCIDA BRIGHT.

*Key words and phrases.* Matrix derivatives, Stiefel manifolds, majorization algorithms, polynomial factor analysis.

It is common to define the Procrustus Transformation using the *Singular Value Decomposition*. We write the SVD as  $X = K\Lambda L'$ , where the  $n \times m$  matrix  $K$  satisfies  $K'K = I$ , the  $m \times m$  matrix  $L$  satisfies  $L'L = LL' = I$ , and the diagonal matrix  $\Lambda$  has positive elements in non-increasing order along the diagonal. Then  $X'X = \Lambda^2 L'$ , and  $(X'X)^{\frac{1}{2}} = \Lambda L'$ , and  $(X'X)^{-\frac{1}{2}} = L\Lambda'$ . Thus  $\Pi(X) = KL'$ .

## 2. USING THE SVD

The singular value decomposition is defined as the solution to the equations

$$(2a) \quad XL = K\Lambda,$$

$$(2b) \quad X'K = L\Lambda,$$

$$(2c) \quad K'K = I,$$

$$(2d) \quad L'L = I.$$

Suppose the elements  $x_{ij}$  of  $X$  are differentiable functions of a single real variable, say  $\theta$ . We use notation such as  $\mathcal{D}X$  for matrices of derivatives with respect to  $\theta$ . Thus  $\{\mathcal{D}X\}_{ij} = \mathcal{D}x_{ij}$  is an  $n \times m$  matrix of functions, just like  $X$ .

Differentiating (2) gives

$$(3a) \quad (\mathcal{D}X)L + X(\mathcal{D}L) = K(\mathcal{D}\Lambda) + (\mathcal{D}K)\Lambda,$$

$$(3b) \quad (\mathcal{D}X')K + X'(\mathcal{D}K) = L(\mathcal{D}\Lambda) + (\mathcal{D}L)\Lambda,$$

$$(3c) \quad (\mathcal{D}K)'K + K'(\mathcal{D}K) = 0,$$

$$(3d) \quad (\mathcal{D}L)'L + L'(\mathcal{D}L) = 0.$$

Let  $A \triangleq K'(\mathcal{D}X)L$ . It follows immediately that  $\mathcal{D}\Lambda = \mathbf{diag}(A)$ . Also define  $\bar{A} \triangleq A - \mathbf{diag}(A)$ . Then

$$\bar{A} = K'(\mathcal{D}K)\Lambda - \Lambda L'(\mathcal{D}L),$$

$$\bar{A}' = L'(\mathcal{D}L)\Lambda - \Lambda K'(\mathcal{D}K).$$

Thus

$$\Lambda \bar{A} = \Lambda K'(\mathcal{D}K)\Lambda - \Lambda^2 L'(\mathcal{D}L),$$

$$\bar{A}' \Lambda = L'(\mathcal{D}L)\Lambda^2 - \Lambda K'(\mathcal{D}K)\Lambda,$$

and thus

$$(4a) \quad \Lambda \bar{A} + \bar{A}' \Lambda = L'(\mathcal{D}L)\Lambda^2 - \Lambda^2 L'(\mathcal{D}L).$$

Similarly

$$\bar{A}\Lambda = K'(\mathcal{D}K)\Lambda^2 - \Lambda L'(\mathcal{D}L)\Lambda,$$

$$\Lambda \bar{A}' = \Lambda L'(\mathcal{D}L)\Lambda - \Lambda^2 K'(\mathcal{D}K),$$

and thus

$$(4b) \quad \bar{A}\Lambda + \Lambda \bar{A}' = K'(\mathcal{D}K)\Lambda^2 - \Lambda^2 K'(\mathcal{D}K).$$

Suppose all singular values are different. Define the  $m \times m$  anti-symmetric matrices  $E$  and  $F$  with zero diagonal elements and off-diagonal elements

$$(5a) \quad e_{j\ell} \triangleq \frac{\lambda_j a_{j\ell} + \lambda_\ell a_{\ell j}}{\lambda_\ell^2 - \lambda_j^2},$$

$$(5b) \quad f_{j\ell} \triangleq \frac{\lambda_\ell a_{j\ell} + \lambda_j a_{\ell j}}{\lambda_\ell^2 - \lambda_j^2}.$$

From (4a)  $L'(\mathcal{D}L) = E$  and thus  $\mathcal{D}L = LE$ . From (4b)  $K'(\mathcal{D}K) = F$ . Choose an  $n \times (n - m)$  matrix  $K_0$  satisfying  $K'K_0 = 0$  and  $K'_0K_0 = I$ . Also define  $A_0 \triangleq K'_0(\mathcal{D}X)L$ . From (3a) we have  $A_0 = K'_0(\mathcal{D}K)\Lambda$ , which means  $\mathcal{D}K = KF + K_0F_0$ , where  $F_0 \triangleq A_0\Lambda^{-1}$ . To write our result more compactly we define  $\vec{K} \triangleq [K \mid K_0]$  and  $\vec{F} \triangleq \begin{bmatrix} F \\ F_0 \end{bmatrix}$ .

**Lemma 2.1.** *For the singular value decomposition  $X = K\Lambda L'$  of an  $n \times m$  matrix  $X$ , with  $n \geq m$ , we find*

$$(6a) \quad \mathcal{D}\Lambda = \mathbf{diag}(A),$$

$$(6b) \quad \mathcal{D}L = LE,$$

$$(6c) \quad \mathcal{D}K = \vec{K}\vec{F}.$$

**Theorem 2.2.**

$$(7) \quad \mathcal{D}\Pi = K(F - E)L' + K_0F_0L' = \vec{K}\vec{F}L' - KEL'$$

*Proof.* Use  $\mathcal{D}\Pi = K(\mathcal{D}L') + (\mathcal{D}K)L'$  and apply the lemma.  $\square$

Note that for  $i \neq j$

$$(8) \quad f_{j\ell} - e_{j\ell} = \frac{a_{j\ell} - a_{\ell j}}{\lambda_j + \lambda_\ell},$$

which is again anti-symmetric.

### 3. USING THE SYMMETRIC SQUARE ROOT

For archival purposes we give an alternative derivation that uses the definition  $\Pi(X) = X(X'X)^{-\frac{1}{2}}$ . Differentiating gives

$$(9) \quad \mathcal{D}\Pi = (\mathcal{D}X)(X'X)^{-\frac{1}{2}} + X\mathcal{D}((X'X)^{-\frac{1}{2}}).$$

Now

$$(10) \quad \mathcal{D}((X'X)^{-\frac{1}{2}}) = -(X'X)^{-\frac{1}{2}}\mathcal{D}((X'X)^{\frac{1}{2}})(X'X)^{-\frac{1}{2}}.$$

We also differentiate

$$(11) \quad (X'X)^{\frac{1}{2}}(X'X)^{\frac{1}{2}} = X'X,$$

and find

$$(12) \quad \mathcal{D}((X'X)^{\frac{1}{2}})(X'X)^{\frac{1}{2}} + (X'X)^{\frac{1}{2}}\mathcal{D}((X'X)^{\frac{1}{2}}) = (\mathcal{D}X)'X + X'(\mathcal{D}X).$$

Now  $(X'X)^{\frac{1}{2}} = L\Lambda L'$ . Define  $H \triangleq L'\mathcal{D}((X'X)^{\frac{1}{2}})L$ . Then

$$(13) \quad H\Lambda + \Lambda H = A'\Lambda + \Lambda A,$$

and thus

$$(14) \quad h_{j\ell} = \frac{\lambda_j a_{j\ell} + \lambda_\ell a_{\ell j}}{\lambda_j + \lambda_\ell}.$$

We see that  $\mathcal{D}((X'X)^{\frac{1}{2}}) = LHL'$  and  $\mathcal{D}((X'X)^{-\frac{1}{2}}) = -L\Lambda^{-1}H\Lambda^{-1}L'$ .

**Lemma 3.1.**

$$(15a) \quad \mathcal{D}((X'X)^{\frac{1}{2}}) = LHL',$$

$$(15b) \quad \mathcal{D}((X'X)^{-\frac{1}{2}}) = -L\Lambda^{-1}H\Lambda^{-1}L',$$

where  $H$  is given by (14).

**Theorem 3.2.**

$$(16) \quad \mathcal{D}\Pi = (\mathcal{D}X)(X'X)^{-\frac{1}{2}} - KHA^{-1}L'.$$

*Proof.*

□

For both (7) and (16) we see that  $K'_0(\mathcal{D}\Pi)L = K'_0(\mathcal{D}X)L\Lambda^{-1}$  and  $K'(\mathcal{D}\Pi)L = F - E = (A - H)\Lambda^{-1}$ . Thus the two expressions for  $\mathcal{D}\Pi$  indeed give identical results.

#### 4. PARTIAL DERIVATIVES

$$\frac{\partial \{\Pi\}_{is}}{\partial x_{kt}} = e'_i \{e_k e'_t (X'X)^{-\frac{1}{2}} - KHA^{-1}L'\} e_s,$$

where

$$h_{st} = \frac{\lambda_s a_{st} + \lambda_t a_{ts}}{\lambda_s + \lambda_t}.$$

#### 5. MINIMIZATION OVER THE ORTHONORMALS

Consider the constrained minimization problem of minimizing  $f(X)$  over the  $n \times m$  matrices  $X$  satisfying  $X'X = I$ . This is equivalent to the unconstrained problem of minimizing  $g(X) \triangleq f(\Pi(X))$  over all  $n \times m$  matrices  $X$ . By the chain rule the derivative is

$$\mathcal{D}g(X) = \mathcal{D}f(\Pi(X))\mathcal{D}\Pi(X),$$

which shows the relevance of our previous calculations.

Consider the Procrustus example of maximizing  $f(X) = \mathbf{tr} Z'X$  over  $X'X = I$ . Instead we maximize  $g(X) = \mathbf{tr} Z'\Pi(X)$  over all  $X$ .

Here is some simple **R** code using the `optimx` package [Nash and Varadhan, 2011].

```

1 require ("optimx")
2
3 set.seed (12345)
4 z <- matrix (rnorm (36), 9, 4)
5 x <- matrix (sample (1 : 36), 9, 4)
6
7 crussy <- function (x, z) {
8   x <- matrix (x, nrow (z), ncol (z))
9   s <- svd (x)
10  p <- tcrossprod (s $ u, s $ v)
11  return (sum (p * z))
12 }
13
14 res <- optimx (par = as.vector (x), fn = crussy, method = "
    BFGS", control = list (trace = 1, maximize = TRUE), z =
    z)

```

The code uses numerical differentiation, so it does not need a function to compute derivatives. We can add a routine to compute the derivatives.

```

1 grussy <- function (x, z) {
2   n <- nrow (z)
3   m <- ncol (z)
4   x <- matrix (x, n, m)
5   g <- matrix (0, n, m)
6   s <- svd (x)
7   k <- s $ u
8   l <- s $ v
9   d <- s $ d
10  q <- l %*% (t (1) / d)
11  for (i in 1 : n) {
12    for (j in 1 : m) {
13      o <- outer (k [i, ], l [j, ])
14      h <- (d * o + t (d * o)) / outer (d, d, "+")
15      u <- k %*% h %*% (t (1) / d)

```

```

16         v <- outer (ifelse (i == (1 : n), 1, 0), q[j,
17                   ]) - u
18         g [i, j] <- sum (z * v)
19     }
20     }
21     return (as.vector (g))
22 }
23 resg <- optimx (par = as.vector (x), fn = crussy, gr =
24               grussy, method = c("BFGS"),
25               control = list (trace = 1, maximize = TRUE), z = z)

```

## 6. APPLICATION TO POLYNOMIAL COMPONENT ANALYSIS

$$f(X, B) = \text{SSQ}(Y - \text{pol}(\Pi(X))B).$$

$$g(X) = \min_B f(X, B) = f(X, B(X)).$$

$$\mathcal{D}g(X) = \mathcal{D}_1 f(X, B(X))$$

Minimize  $g$  over all  $X$ , without constraints.

## APPENDIX A. CODE

```

1  library(polynom)
2
3  procdev<-function(n, m, y = 0, p = 4, eps = 1e-6){
4    set.seed(12345)
5    nm <- n * m
6    f <- rep (list(1 : p), n * m)
7    for (i in 1 : nm) {
8      f[[i]] <- polynomial (rnorm (p))
9    }
10   x <- matrix (0, n, m)
11   g <- matrix (0, n, m)
12   for (i in 1 : nm) {
13     x[i] <- predict (f[[i]], y)
14     g[i] <- predict (deriv (f[[i]]), y)
15   }
16   ss <- svd(x)
17   sk <- ss$u
18   s1 <- ss$v
19   sd <- ss$d
20   px <- tcrossprod (sk, s1)
21   sq <- tcrossprod (s1, (s1 %%% diag (sd)))
22   si <- tcrossprod (s1, (s1 %%% diag (1/sd)))
23   a <- crossprod (sk, g %%% s1)
24   la <- sd * a
25   e <- -(la + t(la)) / outer (sd ^ 2, sd ^ 2, "-")
26   diag (e) <- 0
27   la <- a %%% diag (sd)
28   f <- -(la + t(la)) / outer (sd ^ 2, sd ^ 2, "-")
29   diag (f) <- 0
30   d1 <- s1 %%% e
31   dk <- sk %%% f + (diag (n) - tcrossprod (sk)) %%% g %
      %%% s1 %%% diag (1/sd)
32   dr <- tcrossprod(sk, d1) + tcrossprod (dk, s1)
33   return (list (x=x, px=px, dr=dr))
34 }

```



REFERENCES

N. Cliff. Orthogonal Rotation to Congruence. *Psychometrika*, 31: 33-42, 1966.

J.C. Nash and R. Varadhan. Unifying Optimization Algorithms to Aid Software System Users: optimx for R. *Journal of Statistical Software*, 43(9):1-13, 2011. URL <http://www.jstatsoft.org/v43/i09/paper>.

DEPARTMENT OF STATISTICS, UNIVERSITY OF CALIFORNIA, LOS ANGELES, CA 90095-1554

*E-mail address*, Jan de Leeuw: [deleeuw@stat.ucla.edu](mailto:deleeuw@stat.ucla.edu)

*E-mail address*, Kekona Kai Sorenson: [kekona@gmail.com](mailto:kekona@gmail.com)

*URL*, Jan de Leeuw: <http://gifi.stat.ucla.edu>

*URL*, Kekona Kai Sorenson: <http://www.kekonakai.com/Site/Welcome.html>