

Old and New Matrix Algebra Useful for Statistics

Thomas P. Minka

December 28, 2000

Contents

1	Derivatives	1
2	Kronecker product and vec	6
3	Vec-transpose	7
4	Multilinear forms	8
5	Hadamard product and diag	10
6	Inverting partitioned matrices	12
7	Polar decomposition	14
8	Hessians	15

Warning: This paper contains a large number of matrix identities which cannot be absorbed by mere reading. The reader is encouraged to take time and check each equation by hand and work out the examples. This is advanced material; see Searle (1982) for basic results.

1 Derivatives

Maximum-likelihood problems almost always require derivatives. There are six kinds of derivatives that can be expressed as matrices:

	Scalar	Vector	Matrix
Scalar	$\frac{dy}{dx}$	$\frac{d\mathbf{y}}{dx} = \left[\frac{\partial y_i}{\partial x} \right]$	$\frac{d\mathbf{Y}}{dx} = \left[\frac{\partial y_{ij}}{\partial x} \right]$
Vector	$\frac{dy}{d\mathbf{x}} = \left[\frac{\partial y}{\partial x_j} \right]$	$\frac{d\mathbf{y}}{d\mathbf{x}} = \left[\frac{\partial y_i}{\partial x_j} \right]$	
Matrix	$\frac{dy}{d\mathbf{X}} = \left[\frac{\partial y}{\partial x_{ji}} \right]$		

The partials with respect to the numerator are laid out according to the shape of \mathbf{Y} while the partials with respect to the denominator are laid out according to the transpose of \mathbf{X} . For example, $d\mathbf{y}/dx$ is a column vector while $dy/d\mathbf{x}$ is a row vector (assuming \mathbf{x} and \mathbf{y} are column vectors—otherwise it is flipped). Each of these derivatives can be tediously computed via partials, but this section shows how they instead can be computed with matrix manipulations. The material is based on Magnus and Neudecker (1988).

Define the differential $dy(x)$ to be that part of $y(x+dx) - y(x)$ which is linear in dx . Unlike the classical definition in terms of limits, this definition applies even when x or y are not scalars.

For example, this equation:

$$\mathbf{y}(\mathbf{x} + d\mathbf{x}) = \mathbf{y}(\mathbf{x}) + \mathbf{A}d\mathbf{x} + (\text{higher order terms}) \quad (1)$$

is well-defined for any \mathbf{y} satisfying certain continuity properties. The matrix \mathbf{A} is the derivative, as you can check by setting all but one component of $d\mathbf{x}$ to zero and making it small. The matrix \mathbf{A} is also called the Jacobian matrix $\mathbf{J}_{x \rightarrow y}$. Its transpose is the gradient of \mathbf{y} , denoted $\nabla \mathbf{y}$. The Jacobian is useful in calculus while the gradient is useful in optimization.

Therefore, the derivative of any expression involving matrices can be computed in two steps:

1. compute the differential
2. massage the result into canonical form

after which the derivative is immediately read off as the coefficient of dx , $d\mathbf{x}$, or $d\mathbf{X}$.

The differential of an expression can be computed by iteratively applying the following rules:

$$d\mathbf{A} = 0 \quad (\text{for constant } \mathbf{A}) \quad (2)$$

$$d(\alpha \mathbf{X}) = \alpha d\mathbf{X} \quad (3)$$

$$d(\mathbf{X} + \mathbf{Y}) = d\mathbf{X} + d\mathbf{Y} \quad (4)$$

$$d(\text{tr}(\mathbf{X})) = \text{tr}(d\mathbf{X}) \quad (5)$$

$$d(\mathbf{X}\mathbf{Y}) = (d\mathbf{X})\mathbf{Y} + \mathbf{X}d\mathbf{Y} \quad (6)$$

$$d(\mathbf{X} \otimes \mathbf{Y}) = (d\mathbf{X}) \otimes \mathbf{Y} + \mathbf{X} \otimes d\mathbf{Y} \quad (\text{see section 2}) \quad (7)$$

$$d(\mathbf{X} \circ \mathbf{Y}) = (d\mathbf{X}) \circ \mathbf{Y} + \mathbf{X} \circ d\mathbf{Y} \quad (\text{see section 5}) \quad (8)$$

$$d\mathbf{X}^{-1} = -\mathbf{X}^{-1}(d\mathbf{X})\mathbf{X}^{-1} \quad (9)$$

$$d|\mathbf{X}| = |\mathbf{X}| \text{tr}(\mathbf{X}^{-1}d\mathbf{X}) \quad (10)$$

$$d \log |\mathbf{X}| = \text{tr}(\mathbf{X}^{-1}d\mathbf{X}) \quad (11)$$

$$d\mathbf{X}^* = (d\mathbf{X})^* \quad (12)$$

where $*$ is any operator that rearranges elements, e.g. transpose, vec, and vec-transpose (section 3). The rules can be iteratively applied because of the chain rule, e.g. $d(\mathbf{A}\mathbf{X} + \mathbf{Y}) = d(\mathbf{A}\mathbf{X}) + d\mathbf{Y} = \mathbf{A}d\mathbf{X} + (d\mathbf{A})\mathbf{X} + d\mathbf{Y} = \mathbf{A}d\mathbf{X} + d\mathbf{Y}$. Most of these rules can be derived by subtracting $\mathbf{F}(\mathbf{X} + d\mathbf{X}) - \mathbf{F}(\mathbf{X})$ and taking the linear part. For example,

$$(\mathbf{X} + d\mathbf{X})(\mathbf{Y} + d\mathbf{Y}) = \mathbf{X}\mathbf{Y} + (d\mathbf{X})\mathbf{Y} + \mathbf{X}d\mathbf{Y} + (d\mathbf{X})(d\mathbf{Y})$$

from which (6) follows.

To derive $d\mathbf{X}^{-1}$, note that

$$0 = d\mathbf{I} = d\mathbf{X}^{-1}\mathbf{X} = (d\mathbf{X}^{-1})\mathbf{X} + \mathbf{X}^{-1}d\mathbf{X}$$

from which (9) follows.

The next step is to massage the differential into one of the six canonical forms (assuming \mathbf{x} and \mathbf{y} are column vectors):

$dy = a dx$	$d\mathbf{y} = \mathbf{a} d\mathbf{x}$	$d\mathbf{Y} = \mathbf{A} d\mathbf{x}$
$dy = \mathbf{a} d\mathbf{x}$	$d\mathbf{y} = \mathbf{A} d\mathbf{x}$	
$dy = \text{tr}(\mathbf{A} d\mathbf{X})$		

This is where the operators and identities developed in the following sections are useful. For example, since the derivative of \mathbf{Y} with respect to \mathbf{X} cannot be represented by a matrix, it is customary to use $d\text{vec}(\mathbf{Y})/d\text{vec}(\mathbf{X})$ instead (vec is defined in section 2). If the purpose of differentiation is to equate the derivative to zero, then this transformation doesn't affect the result. So after expanding the differential, just take vec of both sides and use the identities in sections 2 and 3 to get it into canonical form.

One particularly helpful identity is:

$$\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA}) \quad (13)$$

Examples:

$$\begin{aligned} \frac{d}{d\mathbf{X}} \text{tr}(\mathbf{AXB}) &= \mathbf{BA} \\ \text{because } d\text{tr}(\mathbf{AXB}) &= \text{tr}(\mathbf{A}(d\mathbf{X})\mathbf{B}) = \text{tr}(\mathbf{BA}d\mathbf{X}) \end{aligned} \quad (14)$$

$$\begin{aligned} \frac{d}{d\mathbf{X}} \text{tr}(\mathbf{AX'BXC}) &= \mathbf{CAX'B} + \mathbf{A'C'X'B'} \\ \text{because } d\text{tr}(\mathbf{AX'BXC}) &= \text{tr}(\mathbf{AX'B}(d\mathbf{X})\mathbf{C}) + \text{tr}(\mathbf{A}(d\mathbf{X})'\mathbf{BXC}) \\ &= \text{tr}((\mathbf{CAX'B} + \mathbf{A'C'X'B'})d\mathbf{X}) \end{aligned} \quad (15)$$

$$\begin{aligned} \frac{d}{d\mathbf{X}} \text{tr}(\mathbf{AX}^{-1}\mathbf{B}) &= -\mathbf{X}^{-1}\mathbf{BAX}^{-1} \\ \text{because } d\text{tr}(\mathbf{AX}^{-1}\mathbf{B}) &= -\text{tr}(\mathbf{AX}^{-1}(d\mathbf{X})\mathbf{X}^{-1}\mathbf{B}) \\ &= -\text{tr}(\mathbf{X}^{-1}\mathbf{BAX}^{-1}d\mathbf{X}) \end{aligned} \quad (16)$$

$$\begin{aligned} \frac{d}{d\mathbf{X}} \text{tr}(\mathbf{A}(\mathbf{X}\Sigma\mathbf{X}')^{-1}\mathbf{B}) &= -\Sigma\mathbf{X}'(\mathbf{X}\Sigma\mathbf{X}')^{-1}(\mathbf{BA} + \mathbf{A'B}')(\mathbf{X}\Sigma\mathbf{X}')^{-1} \\ &\quad (\text{for symmetric } \Sigma) \\ \text{because } d\text{tr}(\mathbf{A}(\mathbf{X}\Sigma\mathbf{X}')^{-1}\mathbf{B}) &= -\text{tr}(\mathbf{A}(\mathbf{X}\Sigma\mathbf{X}')^{-1}((d\mathbf{X})\Sigma\mathbf{X}' + \mathbf{X}\Sigma(d\mathbf{X})')(\mathbf{X}\Sigma\mathbf{X}')^{-1}\mathbf{B}) \\ &= -\text{tr}(\Sigma\mathbf{X}'(\mathbf{X}\Sigma\mathbf{X}')^{-1}(\mathbf{BA} + \mathbf{A'B}')(\mathbf{X}\Sigma\mathbf{X}')^{-1}d\mathbf{X}) \end{aligned} \quad (17)$$

$$\frac{d}{d\mathbf{X}} |\mathbf{X}| = |\mathbf{X}| \mathbf{X}^{-1} \quad (18)$$

$$\begin{aligned}
\frac{d}{d\mathbf{X}} |\mathbf{X}'\mathbf{X}| &= 2 |\mathbf{X}'\mathbf{X}| (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \\
\text{because } d |\mathbf{X}'\mathbf{X}| &= |\mathbf{X}'\mathbf{X}| \text{tr}((\mathbf{X}'\mathbf{X})^{-1} d(\mathbf{X}'\mathbf{X})) \\
&= |\mathbf{X}'\mathbf{X}| \text{tr}((\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}' d\mathbf{X} + (d\mathbf{X})' \mathbf{X})) \\
&= 2 |\mathbf{X}'\mathbf{X}| \text{tr}((\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' d\mathbf{X})
\end{aligned} \tag{19}$$

$$\begin{aligned}
\frac{d}{d\mathbf{X}} f(\mathbf{X}\mathbf{z}) &= \mathbf{z} \left. \frac{d}{d\mathbf{x}} f(\mathbf{x}) \right|_{\mathbf{x}=\mathbf{X}\mathbf{z}} \\
\text{because } df(\mathbf{x}) &= \left(\frac{d}{d\mathbf{x}} f(\mathbf{x}) \right) d\mathbf{x} \quad (\text{by definition}) \\
df(\mathbf{X}\mathbf{z}) &= \left. \frac{d}{d\mathbf{x}} f(\mathbf{x}) \right|_{\mathbf{x}=\mathbf{X}\mathbf{z}} (d\mathbf{X})\mathbf{z} \\
&= \text{tr}(\mathbf{z} \left. \frac{d}{d\mathbf{x}} f(\mathbf{x}) \right|_{\mathbf{x}=\mathbf{X}\mathbf{z}} d\mathbf{X})
\end{aligned} \tag{20}$$

Constraints Sometimes we want to take the derivative of a function whose argument must be symmetric. In this case, $d\mathbf{X}$ must be symmetric, so we get

$$dy(\mathbf{X}) = \text{tr}(\mathbf{A}d\mathbf{X}) \Rightarrow \frac{dy(\mathbf{X})}{d\mathbf{X}} = (\mathbf{A} + \mathbf{A}') - (\mathbf{A} \circ \mathbf{I}) \tag{21}$$

where $\mathbf{A} \circ \mathbf{I}$ is simply \mathbf{A} with off-diagonal elements set to zero. The reader can check this by expanding $\text{tr}(\mathbf{A}d\mathbf{X})$ and merging identical elements of $d\mathbf{X}$. An example of this rule is:

$$\frac{d}{d\Sigma} \log |\Sigma| = 2\Sigma^{-1} - (\Sigma^{-1} \circ \mathbf{I}) \tag{22}$$

when Σ must be symmetric. This is usually easier than taking an unconstrained derivative and then using Lagrange multipliers to enforce symmetry.

Similarly, if \mathbf{X} must be diagonal, then so must $d\mathbf{X}$, and we get

$$dy(\mathbf{X}) = \text{tr}(\mathbf{A}d\mathbf{X}) \Rightarrow \frac{dy(\mathbf{X})}{d\mathbf{X}} = (\mathbf{A} \circ \mathbf{I}) \tag{23}$$

Example: Principal Component Analysis Suppose we want to represent the zero-mean random vector \mathbf{x} as one random variable a times a constant unit vector \mathbf{v} . This is useful for compression or noise removal. Once we choose \mathbf{v} , the optimal choice for a is $\mathbf{v}'\mathbf{x}$, but what is the best \mathbf{v} ? In other words, what \mathbf{v} minimizes $E[(\mathbf{x} - a\mathbf{v})'(\mathbf{x} - a\mathbf{v})]$, when a is chosen optimally for each \mathbf{x} ?

Let $\Sigma = E[\mathbf{x}\mathbf{x}']$. We want to maximize

$$f(\mathbf{v}) = \mathbf{v}'\Sigma\mathbf{v} - \lambda(\mathbf{v}'\mathbf{v} - 1)$$

where λ is a Lagrange multiplier. Taking derivatives gives

$$\nabla f(\mathbf{v}) = 2\Sigma\mathbf{v} - 2\lambda\mathbf{v}$$

so the gradient is zero at any eigenvector of Σ . (Recall that the gradient is the transpose of the derivative.) If \mathbf{v} is an eigenvector then $f(\mathbf{v}) = \lambda$ so the maximum is attained when \mathbf{v} has the largest eigenvalue.

Example: Blind source separation Suppose we have k microphones listening to k overlapped sound sources. Can we recover the individual sources? More generally, suppose we've observed data \mathbf{x} generated by the function $\mathbf{x} = \mathbf{A}^{-1}\mathbf{u}$ where \mathbf{u} is a set of independent hidden causes and \mathbf{A} is an unknown square mixing matrix. Assume each u_i is distributed according to some density $f_i(w_i)$ with unknown parameter w_i . We want to find the mixing matrix which maximizes the likelihood of the data, so that we can then recover the hidden causes.

$$\mathbf{p}(\mathbf{u}|\mathbf{w}) = \prod_i f_i(u_i|w_i) \quad (24)$$

$$\mathbf{p}(\mathbf{x}|\mathbf{A}, \mathbf{w}) = |\mathbf{A}| \prod_i f_i(u_i|w_i) \quad (25)$$

$$\log \mathbf{p}(\mathbf{x}|\mathbf{A}, \mathbf{w}) = \log |\mathbf{A}| + \sum_i \log f_i(u_i|w_i) \quad (26)$$

$$\nabla_{\mathbf{A}} \log \mathbf{p}(\mathbf{x}|\mathbf{A}, \mathbf{w}) = \mathbf{A}^{-T} + \left[\frac{\nabla_{u_i} f_i(u_i|w_i)}{f_i(u_i|w_i)} \right] \mathbf{x}' \quad (27)$$

$$\nabla_{w_i} \log \mathbf{p}(\mathbf{x}|\mathbf{A}, \mathbf{w}) = \frac{\nabla_{w_i} f_i(u_i|w_i)}{f_i(u_i|w_i)} \quad (28)$$

These equations can be used in a gradient-based optimization to find \mathbf{A} and \mathbf{w} . This approach comes from Pearlmutter and Parra (1996).

Example: Gaussian covariance Suppose we've observed vectors \mathbf{x}_i independently sampled from a zero-mean Gaussian distribution, i.e.

$$\mathbf{p}(\mathbf{x}|\Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} \mathbf{x}' \Sigma^{-1} \mathbf{x}\right)$$

We want to determine the most likely covariance matrix Σ , keeping in mind that the solution must be symmetric. Maximizing the log-likelihood gives:

$$\sum_i \log \mathbf{p}(\mathbf{x}_i|\Sigma) = \sum_i \left(-(d/2) \log(2\pi) - (1/2) \log |\Sigma| - \frac{1}{2} \mathbf{x}_i' \Sigma^{-1} \mathbf{x}_i \right) \quad (29)$$

$$\frac{d}{d\Sigma} = \sum_i \left(-\Sigma^{-1} + \Sigma^{-1} \mathbf{x}_i \mathbf{x}_i' \Sigma^{-1} - \frac{1}{2} ((-\Sigma^{-1} + \Sigma^{-1} \mathbf{x}_i \mathbf{x}_i' \Sigma^{-1}) \circ \mathbf{I}) \right) = 0 \quad (30)$$

$$\Sigma = \left(\sum_i \mathbf{x}_i \mathbf{x}_i' \right) / \left(\sum_i 1 \right) \quad (31)$$

2 Kronecker product and vec

The Kronecker product (Lancaster and Tismenetsky, 1985) (Horn and Johnson, 1991) is

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} \end{bmatrix} \quad (32)$$

which, like ordinary matrix product, is associative and distributive but not commutative.

$$(\mathbf{A} \otimes \mathbf{B})' = \mathbf{A}' \otimes \mathbf{B}' \quad (33)$$

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{AC} \otimes \mathbf{BD} \quad (34)$$

which implies $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$.

If \mathbf{A} and \mathbf{B} are square, then the eigenvectors and eigenvalues of $(\mathbf{A} \otimes \mathbf{B})$ are given by

$$\mathbf{A} \otimes \mathbf{B} = \mathbf{V}_A \Lambda_A \mathbf{V}_A^{-1} \otimes \mathbf{V}_B \Lambda_B \mathbf{V}_B^{-1} = (\mathbf{V}_A \otimes \mathbf{V}_B)(\Lambda_A \otimes \Lambda_B)(\mathbf{V}_A \otimes \mathbf{V}_B)^{-1} \quad (35)$$

which implies

$$\text{rank}(\mathbf{A} \otimes \mathbf{B}) = \text{rank}(\mathbf{A})\text{rank}(\mathbf{B}) \quad (36)$$

$$\text{tr}(\mathbf{A} \otimes \mathbf{B}) = \text{tr}(\Lambda_A \otimes \Lambda_B) = \text{tr}(\mathbf{A})\text{tr}(\mathbf{B}) \quad (37)$$

$$|\mathbf{A} \otimes \mathbf{B}| = |\mathbf{A}|^{\text{rank}(\mathbf{B})} |\mathbf{B}|^{\text{rank}(\mathbf{A})} \quad (38)$$

Define $\text{vec}(\mathbf{A})$ to be the stacked columns of \mathbf{A} :

$$\text{vec}\left(\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}\right) = \begin{bmatrix} a_{11} \\ a_{21} \\ a_{12} \\ a_{22} \end{bmatrix} \quad (39)$$

Then the main result is

$$\text{vec}(\mathbf{ABC}) = (\mathbf{C}' \otimes \mathbf{A})\text{vec}(\mathbf{B}) \quad (40)$$

Example The Lyapunov equation is

$$\mathbf{AX} + \mathbf{XB} = \mathbf{C} \quad (41)$$

$$(\mathbf{I} \otimes \mathbf{A} + \mathbf{B}' \otimes \mathbf{I})\text{vec}(\mathbf{X}) = \text{vec}(\mathbf{C}) \quad (42)$$

which can be solved for $\text{vec}(\mathbf{X})$.

The other properties of vec will be presented in the context of vec -transpose.

3 Vec-transpose

Vec-transpose is a new operator that generalizes vec and transpose. It is essential for expressing derivatives of Kronecker products and is also useful for expressing multilinear forms. It was called “vector transposition” by Marimont and Wandell (1992).

Define

$$\begin{aligned} \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \\ a_{41} & a_{42} \\ a_{51} & a_{52} \\ a_{61} & a_{62} \end{bmatrix}^{(2)} &= \begin{bmatrix} a_{11} & a_{31} & a_{51} \\ a_{21} & a_{41} & a_{61} \\ a_{12} & a_{32} & a_{52} \\ a_{22} & a_{42} & a_{62} \end{bmatrix} \\ \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \\ a_{41} & a_{42} \\ a_{51} & a_{52} \\ a_{61} & a_{62} \end{bmatrix}^{(3)} &= \begin{bmatrix} a_{11} & a_{41} \\ a_{21} & a_{51} \\ a_{31} & a_{61} \\ a_{12} & a_{42} \\ a_{22} & a_{52} \\ a_{32} & a_{62} \end{bmatrix} \end{aligned}$$

and similarly $\mathbf{A}^{(p)}$ for any integer p dividing $\text{rows}(\mathbf{A})$.

The basic properties are:

$$\mathbf{A}^{(1)} = \mathbf{A}' \tag{43}$$

$$\mathbf{A}^{(\text{rows}(\mathbf{A}))} = \text{vec}(\mathbf{A}) \tag{44}$$

$$\text{vec}(\mathbf{A})^{(r)} = \text{reshape}(\mathbf{A}, r, c) \text{ (in Matlab notation)} \tag{45}$$

$$\mathbf{A}^{(p)(p)} = \mathbf{A} \tag{46}$$

$$(\alpha\mathbf{A} + \beta\mathbf{B})^{(p)} = \alpha\mathbf{A}^{(p)} + \beta\mathbf{B}^{(p)} \tag{47}$$

We can freely apply vec-transpose inside of a trace expression:

$$\text{tr}(\mathbf{A}'\mathbf{B}) = \text{tr}((\mathbf{A}^{(p)})'\mathbf{B}^{(q)}) \tag{48}$$

assuming conformability. This generalizes $\text{tr}(\mathbf{A}'\mathbf{B}) = \text{vec}(\mathbf{A})'\text{vec}(\mathbf{B})$ as well as $\text{tr}(\mathbf{A}'\mathbf{B}) = \text{tr}(\mathbf{AB}')$. In fact,

$$\text{tr}(\mathbf{A}'\mathbf{B}) = \text{tr}((\mathbf{A}^*)'\mathbf{B}^*) \tag{49}$$

for any operator $*$ that rearranges elements.

We can generalize equations 33 and 40:

$$(\mathbf{A} \otimes \mathbf{B})^{(p)} = \mathbf{A}' \otimes \mathbf{B}^{(p)} \quad (50)$$

$$\begin{aligned} ((\mathbf{D}' \otimes \mathbf{A})\mathbf{B}\mathbf{C})^{(p)} &= (\mathbf{C}' \otimes \mathbf{A})\mathbf{B}^{(p)}\mathbf{D} \\ \text{where } p &= \text{cols}(\mathbf{A}) \end{aligned} \quad (51)$$

We now have the tools to express the derivative of a Kronecker product:

$$\begin{aligned} \frac{d}{d\mathbf{A}} \text{tr}(\mathbf{X}'(\mathbf{A} \otimes \mathbf{B})\mathbf{Y}) &= \frac{d}{d\mathbf{A}} \text{tr}((\mathbf{X}^{(p)})'((\mathbf{A} \otimes \mathbf{B})\mathbf{Y})^{(p)}) \\ &= \frac{d}{d\mathbf{A}} \text{tr}((\mathbf{X}^{(p)})'(\mathbf{I} \otimes \mathbf{B})\mathbf{Y}^{(p)}\mathbf{A}') \\ &= (\mathbf{Y}^{(p)})'(\mathbf{I} \otimes \mathbf{B}')\mathbf{X}^{(p)} \end{aligned} \quad (52)$$

where p is uniquely defined by conformability to be $\text{cols}(\mathbf{B})$.

Equation 51 gives us the following rule for pulling a matrix out of nested vec-transposes:

$$((\mathbf{A}\mathbf{B})^{(p)}\mathbf{C})^{(p)} = (\mathbf{C}' \otimes \mathbf{I})\mathbf{A}\mathbf{B} = (\mathbf{A}^{(p)}\mathbf{C})^{(p)}\mathbf{B} \quad (53)$$

This formula is useful in fitting multilinear forms (see the next section). Unlike regular transpose, it is not true in general that $(\mathbf{B}^{(p)}\mathbf{C})^{(p)} = \mathbf{C}^{(p)}\mathbf{B}$, as we can see by setting $\mathbf{A} = \mathbf{I}$ in (53).

4 Multilinear forms

Multilinear statistical models are more expressive than linear models yet still easy to use. A multilinear form $f(\mathbf{x}, \mathbf{y}, \dots, \mathbf{z})$ is linear in each component separately, i.e.

$$f(\dots, \alpha\mathbf{y}_1 + \beta\mathbf{y}_2, \dots) = \alpha f(\dots, \mathbf{y}_1, \dots) + \beta f(\dots, \mathbf{y}_2, \dots)$$

For example, face images can be modeled as linear in identity and linear in lighting conditions, yielding a bilinear model (Tenenbaum and Freeman, 1997). Another example is colored objects under colored light.

Just as every bilinear form can be written as $\mathbf{x}'\mathbf{G}\mathbf{y} = (\mathbf{y}' \otimes \mathbf{x}')\text{vec}(\mathbf{G})$, every multilinear form can be written as $(\mathbf{z}' \otimes \dots \otimes \mathbf{y}' \otimes \dots \otimes \mathbf{x}')\text{vec}(\mathbf{G})$ (Magnus and Neudecker, 1988) (Prasolov, 1991) (Dodson and Poston, 1991). \mathbf{G} is the tensor defining the form.

For example, the trilinear form

$$\sum_{ijk} G_{ijk} x_i y_j z_k = (\mathbf{z}' \otimes \mathbf{y}' \otimes \mathbf{x}')\text{vec}(\mathbf{G}) = (\mathbf{y}' \otimes \mathbf{x}')\mathbf{G}\mathbf{z} = \mathbf{x}'(\mathbf{G}\mathbf{z})^{(p)}\mathbf{y}$$

Thus we can express a multilinear form either with tensor products or with vec-transpose. Matlab users have often used this kind of reshaping to manipulate higher-dimensional objects.

Just as $\mathbf{Ax} = [\mathbf{c}_1 \ \cdots \ \mathbf{c}_n] \mathbf{x} = \sum_i \mathbf{c}_i x_i$, a linear combination of vectors, we can write

$$(\mathbf{Ax})^{(p)} = ([\mathbf{C}_1^{(p)} \ \cdots \ \mathbf{C}_n^{(p)}] \mathbf{x})^{(p)} = (\sum_i \mathbf{C}_i^{(p)} x_i)^{(p)} = \sum_i \mathbf{C}_i x_i$$

which is a linear combination of matrices.

Therefore, we can think of \mathbf{G} in the trilinear form as a three-dimensional stack of matrices. The formula $\mathbf{x}'(\mathbf{G}\mathbf{z})^{(p)}\mathbf{y}$ says to first combine the stack according to \mathbf{z} , then combine the columns according to \mathbf{y} , and finally to combine the elements according to \mathbf{x} .

The vector-valued bilinear form is:

$$(\mathbf{y}' \otimes \mathbf{x}' \otimes \mathbf{I})\text{vec}(\mathbf{G}) = (\mathbf{G}\mathbf{y})^{(p)}\mathbf{x}$$

which is the same as the scalar-valued trilinear form, except that the three-dimensional tensor \mathbf{G} is only being summed in two dimensions. This form was used in Marimont and Wandell (1992) and subsequently by Tenenbaum and Freeman (1997). To fit this model to data, note that by (53) we have

$$((\mathbf{G}\mathbf{y})^{(p)}\mathbf{x})^{(p)} = (\mathbf{G}^{(p)}\mathbf{x})^{(p)}\mathbf{y}$$

so given an observation and the value of \mathbf{y} we can solve for \mathbf{x} and vice-versa, by applying vec-transpose to the observation. Therefore we can iterate from an initial guess until we reach a fixed point. This method generalizes to any multilinear form. Compare this to the complex algorithm without vec-transpose given in Magnus and Neudecker (1988).

What if \mathbf{G} must be learned as well as \mathbf{x} and \mathbf{y} ? In this case, we need an entire observation matrix $\mathbf{D} = (\mathbf{G}\mathbf{Y})^{(p)}\mathbf{X}$. Without loss of generality, we can assume that \mathbf{X} and \mathbf{Y} are orthogonal matrices, since \mathbf{G} can always be chosen to make this so. (This can be proven with a polar decomposition of \mathbf{X} and \mathbf{Y} .) Therefore, if we know \mathbf{Y} , we can solve for \mathbf{X} by singular-value decomposition of \mathbf{D} , and vice-versa with $\mathbf{D}^{(p)}$. Once \mathbf{X} and \mathbf{Y} have settled, it is easy to solve for $\text{vec}(\mathbf{G})$. This algorithm comes from Marimont and Wandell (1992).

5 Hadamard product and diag

The Hadamard product is simply the product of corresponding elements:

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \circ \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} = \begin{bmatrix} a_{11}b_{11} & a_{12}b_{12} \\ a_{21}b_{21} & a_{22}b_{22} \end{bmatrix} \quad (54)$$

Schur's product theorem (Horn and Johnson, 1991) says

$$\mathbf{A} \geq 0, \mathbf{B} \geq 0 \Rightarrow \mathbf{A} \circ \mathbf{B} \geq 0 \quad (55)$$

This is also true for Kronecker product (by (35)), but not for regular matrix product. To prove it for Hadamard product, define a random vector $\mathbf{z} = \mathbf{x} \circ \mathbf{y}$, where \mathbf{x} and \mathbf{y} are independent random vectors with covariance Σ_x and Σ_y . Then the covariance of \mathbf{z} can be shown to be $\Sigma_x \circ \Sigma_y$. Since every covariance matrix is nonnegative definite, the theorem follows.

Define

$$\text{diag}\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1 & 0 \\ 0 & x_2 \end{bmatrix} \quad (56)$$

$$\text{diag}^{-1}\left(\begin{bmatrix} x_1 & a \\ b & x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (57)$$

diag^{-1} is a kind of pseudoinverse because

$$\text{diag}^{-1}(\text{diag}(\mathbf{x})) = \mathbf{x} \quad (58)$$

but

$$\text{diag}(\text{diag}^{-1}(\mathbf{D})) = \mathbf{D} \quad (59)$$

only for diagonal \mathbf{D} .

The basic properties are:

$$\text{diag}(\alpha\mathbf{x} + \beta\mathbf{y}) = \alpha\text{diag}(\mathbf{x}) + \beta\text{diag}(\mathbf{y}) \quad (60)$$

$$\text{diag}(\mathbf{x} \circ \mathbf{y}) = \text{diag}(\mathbf{x}) \circ \text{diag}(\mathbf{y}) \quad (61)$$

$$\text{diag}(\mathbf{x} \otimes \mathbf{y}) = \text{diag}(\mathbf{x}) \otimes \text{diag}(\mathbf{y}) \quad (62)$$

$$\text{diag}^{-1}(\mathbf{A} \circ \mathbf{B}) = \text{diag}^{-1}(\mathbf{A}) \circ \text{diag}^{-1}(\mathbf{B}) \quad (63)$$

$$\text{diag}^{-1}(\mathbf{A} \otimes \mathbf{B}) = \text{diag}^{-1}(\mathbf{A}) \otimes \text{diag}^{-1}(\mathbf{B}) \quad (64)$$

$$\text{diag}^{-1}((\mathbf{A} \circ \mathbf{B})\mathbf{C}') = \text{diag}^{-1}(\mathbf{A}(\mathbf{B} \circ \mathbf{C})') = \text{diag}^{-1}(\mathbf{B}(\mathbf{A} \circ \mathbf{C})') \quad (65)$$

$$\text{vec}(\mathbf{A} \circ \mathbf{B}) = \text{diag}(\text{vec}(\mathbf{A}))\text{vec}(\mathbf{B}) \quad (66)$$

Equation 66 can be used to remove all Hadamard products from an expression.

Another way to remove Hadamard products is with the diag^{-1} matrix, which is the unique matrix \mathbf{R}_n satisfying

$$\mathbf{R}_n \text{vec}(\mathbf{A}) = \text{diag}^{-1}(\mathbf{A}) \quad (67)$$

where \mathbf{A} is n by n . \mathbf{R}_n is an $n \times n^2$ matrix with orthogonal rows (each row picks out one element of \mathbf{A}). Some useful properties are:

$$\mathbf{R}^+ = \mathbf{R}' \quad (\text{pseudo-inverse is the transpose}) \quad (68)$$

$$\mathbf{R}'\mathbf{x} = \text{vec}(\text{diag}(\mathbf{x})) \quad (69)$$

$$\mathbf{A} \circ \mathbf{B} = \mathbf{R}(\mathbf{A} \otimes \mathbf{B})\mathbf{R}' \quad (70)$$

These properties cause Hadamard product and diag to have a kind of duality with Kronecker product and vec , as seen in the following table:

$$\text{vec}(\mathbf{ABC}) = (\mathbf{C}' \otimes \mathbf{A})\text{vec}(\mathbf{B}) \quad \text{diag}^{-1}(\mathbf{A}\text{diag}(\mathbf{x})\mathbf{C}) = (\mathbf{C}' \circ \mathbf{A})\mathbf{x} \quad (71)$$

$$\text{diag}(\mathbf{x} \circ \mathbf{y} \circ \mathbf{z}) = (\mathbf{z}' \otimes \mathbf{x}) \circ \text{diag}(\mathbf{y}) \quad (72)$$

A useful special case of (71) is

$$\text{diag}^{-1}(\mathbf{AB}') = (\mathbf{A} \circ \mathbf{B})\mathbf{1} \quad (73)$$

If we factor $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^{-1}$ then by (71),

$$\text{diag}^{-1}(\mathbf{A}) = (\mathbf{V} \circ \mathbf{V}^{-T})\text{diag}^{-1}(\mathbf{\Lambda}) \quad (74)$$

which relates the diagonal of a matrix to its eigenvalues. Many facts about the matrix $\mathbf{V} \circ \mathbf{V}^{-T}$ can be found in Horn and Johnson (1991).

Many identities for $\text{diag}^{-1}(\mathbf{A})$ also apply to $\text{tr}(\mathbf{A})$, because

$$\text{tr}(\mathbf{A}) = \mathbf{1}'\text{diag}^{-1}(\mathbf{A}) \quad (75)$$

where $\mathbf{1}$ is a column vector of ones. For example, $\text{tr}(\mathbf{A}\text{diag}(\mathbf{x})\mathbf{C}) = \mathbf{1}'(\mathbf{C}' \circ \mathbf{A})\mathbf{x} = \text{diag}^{-1}(\mathbf{CA})\mathbf{x}$.

Similarly to (48) we have

$$\mathbf{x}'\mathbf{y} = \text{tr}(\text{diag}(\mathbf{x})'\text{diag}(\mathbf{y})) \quad (76)$$

which allows us to compute

$$\begin{aligned} \frac{d}{d\mathbf{A}}\mathbf{x}'(\mathbf{A} \circ \mathbf{B})\mathbf{y} &= \frac{d}{d\mathbf{A}}\text{tr}(\text{diag}(\mathbf{x})'\mathbf{B}\text{diag}(\mathbf{y})\mathbf{A}') = \\ &= \text{diag}(\mathbf{y})\mathbf{B}'\text{diag}(\mathbf{x}) = \mathbf{B}' \circ \mathbf{y}\mathbf{x}' \end{aligned} \quad (77)$$

cf (52).

6 Inverting partitioned matrices

If we partition \mathbf{P} into $\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}$ then the Schur complement of \mathbf{A} in \mathbf{P} (Prasolov, 1991) is

$$(\mathbf{P}|\mathbf{A}) = \mathbf{D} - \mathbf{CA}^{-1}\mathbf{B} \quad (78)$$

Similarly,

$$(\mathbf{P}|\mathbf{D}) = \mathbf{A} - \mathbf{BD}^{-1}\mathbf{C} \quad (79)$$

$$-(\mathbf{P}|\mathbf{B}) = \mathbf{C} - \mathbf{DB}^{-1}\mathbf{A} \quad (\text{if } \mathbf{B}^{-1} \text{ exists}) \quad (80)$$

$$-(\mathbf{P}|\mathbf{C}) = \mathbf{B} - \mathbf{AC}^{-1}\mathbf{D} \quad (\text{if } \mathbf{C}^{-1} \text{ exists}) \quad (81)$$

Define

$$(\mathbf{A}|\mathbf{P}) = (\mathbf{P}|\mathbf{A})^{-1} \quad (82)$$

Then the main result is

$$\mathbf{P}^{-1} = \begin{bmatrix} (\mathbf{D}|\mathbf{P}) & -(\mathbf{B}|\mathbf{P}) \\ -(\mathbf{C}|\mathbf{P}) & (\mathbf{A}|\mathbf{P}) \end{bmatrix} \quad (83)$$

This formula still holds if all inverses are replaced by pseudo-inverses. The reader may want to check this formula when \mathbf{P} is a 2×2 matrix.

Since $\mathbf{PP}^{-1} = \mathbf{I}$ and $\mathbf{P}^{-1}\mathbf{P} = \mathbf{I}$, we know

$$\mathbf{A}(\mathbf{B}|\mathbf{P}) = \mathbf{B}(\mathbf{A}|\mathbf{P}) \quad (84)$$

$$\mathbf{C}(\mathbf{D}|\mathbf{P}) = \mathbf{D}(\mathbf{C}|\mathbf{P}) \quad (85)$$

$$(\mathbf{A}|\mathbf{P})\mathbf{C} = (\mathbf{C}|\mathbf{P})\mathbf{A} \quad (86)$$

$$(\mathbf{B}|\mathbf{P})\mathbf{D} = (\mathbf{D}|\mathbf{P})\mathbf{B} \quad (87)$$

These identities define $(\mathbf{B}|\mathbf{P})$ and $(\mathbf{C}|\mathbf{P})$ when \mathbf{B} and \mathbf{C} are singular. Since $(\mathbf{A}|\mathbf{P})(\mathbf{D} - \mathbf{CA}^{-1}\mathbf{B}) = \mathbf{I}$, we know

$$\begin{aligned} (\mathbf{A}|\mathbf{P}) &= \mathbf{D}^{-1} + (\mathbf{A}|\mathbf{P})\mathbf{CA}^{-1}\mathbf{BD}^{-1} \\ &= \mathbf{D}^{-1} + (\mathbf{C}|\mathbf{P})\mathbf{BD}^{-1} \\ &= \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}(\mathbf{D}|\mathbf{P})\mathbf{BD}^{-1} \end{aligned} \quad (88)$$

$$\text{similarly } (\mathbf{D}|\mathbf{P}) = \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{B}(\mathbf{A}|\mathbf{P})\mathbf{CA}^{-1} \quad (89)$$

which are handy for rewriting (83) solely in terms of $(\mathbf{D}|\mathbf{P})$ or $(\mathbf{A}|\mathbf{P})$.

The Schur complement has the flavor of a division. Equation 84, for example, tells us that $(\mathbf{B}|\mathbf{P})(\mathbf{P}|\mathbf{A}) = \mathbf{A}^{-1}\mathbf{B}$ which is a kind of cancellation of \mathbf{P} . The clearest example of this is the formula for the determinant of \mathbf{P} :

$$|\mathbf{P}| = |(\mathbf{P}|\mathbf{A})| |\mathbf{A}| = |(\mathbf{P}|\mathbf{D})| |\mathbf{D}| \quad (90)$$

When $\mathbf{A} = \mathbf{I}$ and $\mathbf{D} = \mathbf{I}$ (not necessarily the same size) in (90) we get the handy formula

$$|\mathbf{I} + \mathbf{CB}| = |\mathbf{I} + \mathbf{BC}| \quad (91)$$

for any matrices \mathbf{B} and \mathbf{C} .

Example Magnus and Neudecker give a formula for $\begin{bmatrix} \mathbf{A} & \mathbf{B} & \mathbf{C} \\ \mathbf{B}' & \mathbf{D} & 0 \\ \mathbf{C}' & 0 & \mathbf{E} \end{bmatrix}^{-1}$ but leave the proof to the reader. This is easily handled using Schur complements. Define $\mathbf{X} = [\mathbf{B} \ \mathbf{C}]$ and $\mathbf{Y} = \begin{bmatrix} \mathbf{D} & 0 \\ 0 & \mathbf{E} \end{bmatrix}$. Then

$$\mathbf{P}^{-1} = \begin{bmatrix} \mathbf{A} & \mathbf{X} \\ \mathbf{X}' & \mathbf{Y} \end{bmatrix}^{-1} = \begin{bmatrix} (\mathbf{Y}|\mathbf{P}) & -(\mathbf{X}|\mathbf{P}) \\ -(\mathbf{X}|\mathbf{P})' & (\mathbf{A}|\mathbf{P}) \end{bmatrix}$$

where

$$(\mathbf{Y}|\mathbf{P}) = (\mathbf{A} - \mathbf{BD}^{-1}\mathbf{B}' - \mathbf{CE}^{-1}\mathbf{C}')^{-1} \quad (92)$$

$$(\mathbf{X}|\mathbf{P}) = (\mathbf{Y}|\mathbf{P})\mathbf{XY}^{-1} \quad (93)$$

$$= [(\mathbf{Y}|\mathbf{P})\mathbf{BD}^{-1} \quad (\mathbf{Y}|\mathbf{P})\mathbf{CE}^{-1}] \quad (94)$$

$$(\mathbf{A}|\mathbf{P}) = \mathbf{Y}^{-1} + \mathbf{Y}^{-1}\mathbf{X}'(\mathbf{Y}|\mathbf{P})^{-1}\mathbf{XY}^{-1} \quad (95)$$

$$= \begin{bmatrix} \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{B}'(\mathbf{Y}|\mathbf{P})\mathbf{BD}^{-1} & \mathbf{D}^{-1}\mathbf{B}'(\mathbf{Y}|\mathbf{P})\mathbf{CE}^{-1} \\ \mathbf{E}^{-1}\mathbf{C}'(\mathbf{Y}|\mathbf{P})\mathbf{BD}^{-1} & \mathbf{E}^{-1} + \mathbf{E}^{-1}\mathbf{C}'(\mathbf{Y}|\mathbf{P})\mathbf{CE}^{-1} \end{bmatrix} \quad (96)$$

Example: Conditioning a Gaussian density Suppose we partition a zero-mean Gaussian random vector as $\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}$ with partitioned covariance matrix $\mathbf{K} = \begin{bmatrix} \mathbf{K}_{xx} & \mathbf{K}_{xy} \\ \mathbf{K}_{yx} & \mathbf{K}_{yy} \end{bmatrix}$. The distribution of \mathbf{x} conditioned on \mathbf{y} is $\mathbf{p}(\mathbf{x}|\mathbf{y}) = \mathbf{p}(\mathbf{x}, \mathbf{y})/\mathbf{p}(\mathbf{y})$ which is proportional to the joint distribution. Knowing that this conditional is also Gaussian, we can immediately derive its mean and variance by dropping the terms in the joint distribution that depend on \mathbf{y} :

$$\begin{aligned} \mathbf{p}(\mathbf{x}|\mathbf{y}) &\propto \exp\left(-\frac{1}{2} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}' \begin{bmatrix} \mathbf{K}_{xx} & \mathbf{K}_{xy} \\ \mathbf{K}_{yx} & \mathbf{K}_{yy} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}\right) \\ &\propto \exp\left(-\frac{1}{2} \mathbf{x}'(\mathbf{K}_{yy}|\mathbf{K})\mathbf{x} + \mathbf{x}'(\mathbf{K}_{xy}|\mathbf{K})\mathbf{y}\right) \\ &= \exp\left(-\frac{1}{2} \mathbf{x}'(\mathbf{K}_{yy}|\mathbf{K})\mathbf{x} + \mathbf{x}'(\mathbf{K}_{xy}|\mathbf{K})\mathbf{m}\right) \quad (\text{defining } \mathbf{m}) \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m})'(\mathbf{K}_{yy}|\mathbf{K})(\mathbf{x} - \mathbf{m})\right) \end{aligned}$$

where the conditional mean $\mathbf{m} = (\mathbf{K}|\mathbf{K}_{yy})(\mathbf{K}_{xy}|\mathbf{K})\mathbf{y} = \mathbf{K}_{xy}\mathbf{K}_{yy}^{-1}\mathbf{y}$. The conditional variance is therefore $(\mathbf{K}|\mathbf{K}_{yy}) = \mathbf{K}_{xx} - \mathbf{K}_{xy}\mathbf{K}_{yy}^{-1}\mathbf{K}_{yx}$.

7 Polar decomposition

Suppose a set of points \mathbf{B} has been subjected to an unknown rotation and then jittered by white Gaussian noise to give a new set of points \mathbf{A} . What is the most likely rotation? More generally, what unitary matrix minimizes $f(\mathbf{U}) = \text{tr}((\mathbf{A} - \mathbf{UB})'(\mathbf{A} - \mathbf{UB}))$?

Expanding $f(\mathbf{U})$ gives

$$f(\mathbf{U}) = \text{tr}(\mathbf{A}'\mathbf{A}) - 2\text{tr}(\mathbf{A}'\mathbf{UB}) + \text{tr}(\mathbf{B}'\mathbf{B})$$

so the problem reduces to maximizing $\text{tr}(\mathbf{A}'\mathbf{UB}) = \text{tr}(\mathbf{UBA}')$. Define unitary \mathbf{V} and \mathbf{W} and positive diagonal \mathbf{S} so that $\mathbf{BA}' = \mathbf{VSW}'$. This is the singular value decomposition of \mathbf{BA}' . Then

$$\text{tr}(\mathbf{UBA}') = \text{tr}(\mathbf{UVSW}') = \text{tr}(\mathbf{W}'\mathbf{UVS}) \stackrel{\text{def}}{=} \text{tr}(\mathbf{XS})$$

where \mathbf{X} is also unitary. Since \mathbf{S} is diagonal,

$$\text{tr}(\mathbf{XS}) = \sum_i s_{ii}x_{ii}$$

which is maximum when $x_{ii} = 1$; that is, $\mathbf{X} = \mathbf{I}$. Therefore $\mathbf{U} = \mathbf{WV}'$ is the desired solution.

If $\mathbf{B} = \mathbf{I}$, this solution minimizes $\|\mathbf{A} - \mathbf{U}\|$, i.e. it is the closest unitary matrix to an arbitrary matrix \mathbf{A} . This \mathbf{U} has the property that there exists a positive definite \mathbf{P} such that $\mathbf{A} = \mathbf{PU}$. These two matrices are called the polar decomposition of \mathbf{A} : \mathbf{U} is the rotation and \mathbf{P} is the magnitude, exactly analogous to the decomposition of a complex number.

Scott and Longuet-Higgins (1991) used this technique to match columns in \mathbf{A} with those in \mathbf{B} . Matching assumes that \mathbf{U} is a permutation matrix, but finding \mathbf{U} in this case is difficult. So they first found the optimal unitary matrix and obtained a permutation from it.

8 Hessians

The Hessian matrix is a matrix of second derivatives. The Hessian of a scalar function with respect to a vector argument is

$$\frac{dy}{d\mathbf{x}d\mathbf{x}'} = \left[\frac{\partial y}{\partial x_i \partial x_j} \right]$$

This section is based on Magnus and Neudecker (1988).

The Hessian is the derivative of the first derivative. The first derivative $\mathbf{a}(\mathbf{x}) = \left[\frac{\partial y}{\partial x_j} \right]$ is a row vector function of \mathbf{x} , and the derivative of this function with respect to \mathbf{x}' is a matrix

$$\frac{d\mathbf{a}}{d\mathbf{x}'} = \left[\frac{\partial a_j}{\partial x_i} \right] = \left[\frac{\partial y}{\partial x_i \partial x_j} \right]$$

The Hessian can also be defined by the Taylor expansion of y :

$$y(\mathbf{x} + d\mathbf{x}) = y(\mathbf{x}) + \mathbf{a}'d\mathbf{x} + \frac{1}{2}d\mathbf{x}'\mathbf{H}d\mathbf{x} + (\text{higher order terms}) \quad (97)$$

The matrix \mathbf{H} is the Hessian, as you can check by setting all but two components of $d\mathbf{x}$ to zero.

The Hessian can be computed in three steps:

1. Compute the first differential
2. Compute the differential of the first differential
3. Massage the result into canonical form

The only new differential rule we need is:

$$d(d\mathbf{x}) = 0 \quad (98)$$

because $d\mathbf{x}$ is not a function of \mathbf{x} .

The second differential has three canonical forms:

$d^2y = h(dx)^2$
$d^2y = d\mathbf{x}'\mathbf{H}d\mathbf{x}$
$d^2y = d\text{vec}(\mathbf{X})'\mathbf{H}d\text{vec}(\mathbf{X})$

where \mathbf{H} must be symmetric.

Some of these forms require rewriting the differential in terms of $d\text{vec}(\mathbf{X})$, which can be tricky. Equation 48 is particularly helpful for introducing vec into an expression. To eliminate terms

like $\text{vec}(\mathbf{X}')$, Magnus and Neudecker (1988) define the *commutation matrix* \mathbf{K}_{nm} , which is the permutation matrix satisfying

$$\mathbf{K}_{nm} \text{vec}(\mathbf{X}') = \text{vec}(\mathbf{X}) \quad (99)$$

where \mathbf{X} is n by m . Also,

$$\mathbf{K}'_{nm} = \mathbf{K}_{nm}^{-1} = \mathbf{K}_{mn} \quad (100)$$

Examples:

$$\begin{aligned} \frac{d^2}{d\text{vec}(\mathbf{X})d\text{vec}(\mathbf{X})'} \text{tr}(\mathbf{X}'\mathbf{X}) &= 2\mathbf{I}_{n^2} \\ \text{because } d\text{tr}(\mathbf{X}'\mathbf{X}) &= 2\text{tr}(\mathbf{X}'d\mathbf{X}) \\ d^2\text{tr}(\mathbf{X}'\mathbf{X}) &= 2\text{tr}(d\mathbf{X}'d\mathbf{X}) = 2\text{vec}(d\mathbf{X})'\text{vec}(d\mathbf{X}) \end{aligned} \quad (101)$$

$$\frac{d^2}{d\text{vec}(\mathbf{X})d\text{vec}(\mathbf{X})'} \text{tr}(\mathbf{X}'\mathbf{A}\mathbf{X}) = 2\mathbf{I}_n \otimes \mathbf{A} \quad (102)$$

$$\begin{aligned} \text{because } d\text{tr}(\mathbf{X}'\mathbf{A}\mathbf{X}) &= 2\text{tr}(\mathbf{X}'\mathbf{A}d\mathbf{X}) \\ d^2\text{tr}(\mathbf{X}'\mathbf{A}\mathbf{X}) &= 2\text{tr}(d\mathbf{X}'\mathbf{A}d\mathbf{X}) \end{aligned} \quad (103)$$

$$= 2\text{vec}(d\mathbf{X})'\text{vec}(\mathbf{A}d\mathbf{X}) \quad (104)$$

$$= 2\text{vec}(d\mathbf{X})'(\mathbf{I} \otimes \mathbf{A})\text{vec}(d\mathbf{X}) \quad (105)$$

$$\frac{d^2}{d\text{vec}(\mathbf{X})d\text{vec}(\mathbf{X})'} \text{tr}(\mathbf{X}^2) = 2\mathbf{K}_{nn} \quad (106)$$

$$\begin{aligned} \text{because } d\text{tr}(\mathbf{X}^2) &= 2\text{tr}(\mathbf{X}d\mathbf{X}) \\ d^2\text{tr}(\mathbf{X}^2) &= 2\text{tr}((d\mathbf{X})^2) = 2\text{vec}(d\mathbf{X})'\text{vec}(d\mathbf{X}) \\ &= 2\text{vec}(d\mathbf{X})'\mathbf{K}_{nn}\text{vec}(d\mathbf{X}) \end{aligned}$$

$$\frac{d^2}{d\text{vec}(\mathbf{X})d\text{vec}(\mathbf{X})'} \log |\mathbf{X}| = -\mathbf{K}_{nn}(\mathbf{X}^{-T} \otimes \mathbf{X}^{-1}) \quad (107)$$

$$\begin{aligned} \text{because } d \log |\mathbf{X}| &= \text{tr}(\mathbf{X}^{-1}d\mathbf{X}) \\ d^2 \log |\mathbf{X}| &= -\text{tr}(\mathbf{X}^{-1}(d\mathbf{X})\mathbf{X}^{-1}d\mathbf{X}) \\ &= -\text{vec}(d\mathbf{X})'\text{vec}(\mathbf{X}^{-1}(d\mathbf{X})\mathbf{X}^{-1}) \\ &= -\text{vec}(d\mathbf{X})'(\mathbf{X}^{-T} \otimes \mathbf{X}^{-1})\text{vec}(d\mathbf{X}) \end{aligned}$$

$$\frac{d^2}{d\text{vec}(\mathbf{X})d\text{vec}(\mathbf{X})'} \text{tr}(\mathbf{X}^{-1}) = \mathbf{K}_{nn}(\mathbf{X}^{-2T} \otimes \mathbf{X}^{-1} + \mathbf{X}^{-2} \otimes \mathbf{X}^{-T}) \quad (108)$$

$$\begin{aligned} \text{because } d\text{tr}(\mathbf{X}^{-1}) &= -\text{tr}(\mathbf{X}^{-1}(d\mathbf{X})\mathbf{X}^{-1}) \\ d^2\text{tr}(\mathbf{X}^{-1}) &= 2\text{tr}((d\mathbf{X})\mathbf{X}^{-1}(d\mathbf{X})\mathbf{X}^{-2}) \end{aligned}$$

$$\begin{aligned}
&= 2\text{vec}(\text{d}\mathbf{X}')'\text{vec}(\mathbf{X}^{-1}(\text{d}\mathbf{X})\mathbf{X}^{-2}) \\
&= 2\text{vec}(\text{d}\mathbf{X}')'(\mathbf{X}^{-2\text{T}} \otimes \mathbf{X}^{-1})\text{vec}(\text{d}\mathbf{X}) \\
&= \text{vec}(\text{d}\mathbf{X}')'(\mathbf{X}^{-2\text{T}} \otimes \mathbf{X}^{-1} + \mathbf{X}^{-2} \otimes \mathbf{X}^{-\text{T}})\text{vec}(\text{d}\mathbf{X}) \text{ (for symmetry)}
\end{aligned}$$

Constraints Sometimes we want to compute the Hessian of a function whose argument must be symmetric. In this case, the conversion from canonical form is

$$\text{d}^2y = \text{dvec}(\mathbf{X})'\mathbf{H}\text{dvec}(\mathbf{X}) \Rightarrow \quad (109)$$

$$\frac{\text{d}^2y}{\text{dvec}(\mathbf{X})\text{dvec}(\mathbf{X})'} = (\mathbf{I}_{n^2} + \mathbf{K}_{nn} - \text{diag}(\text{vec}(\mathbf{I}_n)))\mathbf{H}(\mathbf{I}_{n^2} + \mathbf{K}_{nn} - \text{diag}(\text{vec}(\mathbf{I}_n))) \quad (110)$$

$$= \mathbf{D}_n\mathbf{D}_n'\mathbf{H}\mathbf{D}_n\mathbf{D}_n' \quad (111)$$

where \mathbf{K}_{nn} is the commutation matrix discussed earlier and \mathbf{D}_n is defined below. Since ∂x_{ij} and ∂x_{ji} are identical, this formula adds together $\frac{\partial y}{\partial x_{ij}^2}$, $\frac{\partial y}{\partial x_{ij}\partial x_{ji}}$, $\frac{\partial y}{\partial x_{ji}\partial x_{ij}}$, and $\frac{\partial y}{\partial x_{ji}^2}$, when $i \neq j$. To derive it, we make $\text{d}\mathbf{X}$ symmetric by substituting $\text{d}\mathbf{X} + \text{d}\mathbf{X}' - (\text{d}\mathbf{X} \circ \mathbf{I})$ (cf (21)) and get

$$\text{vec}(\text{d}\mathbf{X}) \Rightarrow \text{vec}(\text{d}\mathbf{X} + \text{d}\mathbf{X}' - (\mathbf{I} \circ \text{d}\mathbf{X})) \quad (112)$$

$$= (\mathbf{I}_{n^2} + \mathbf{K}_{nn} - \text{diag}(\text{vec}(\mathbf{I}_n)))\text{vec}(\text{d}\mathbf{X}) \quad (113)$$

by (99) and (66).

However, we may not want the full Hessian, but only the submatrix corresponding to unique elements of \mathbf{X} . That is, we want $\frac{\text{d}^2y}{\text{dvech}(\mathbf{X})\text{dvech}(\mathbf{X})'}$, where $\text{vech}(\mathbf{X})$ (Searle, 1982) is $\text{vec}(\mathbf{X})$ with elements above the diagonal deleted. For example,

$$\text{vec}\left(\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}\right) = \begin{bmatrix} a_{11} \\ a_{21} \\ a_{12} \\ a_{22} \end{bmatrix} \quad (114)$$

$$\text{vech}\left(\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}\right) = \begin{bmatrix} a_{11} \\ a_{21} \\ a_{22} \end{bmatrix} \quad (115)$$

To convert between $\text{vec}(\mathbf{X})$ and $\text{vech}(\mathbf{X})$, Magnus and Neudecker (1988) define the *duplication matrix* \mathbf{D}_n , which is the permutation matrix satisfying

$$\mathbf{D}_n\text{vech}(\mathbf{X}) = \text{vec}(\mathbf{X}) \quad (116)$$

where \mathbf{X} is n by n . This leads to the rule

$$\text{d}^2y = \text{dvec}(\mathbf{X})'\mathbf{H}\text{dvec}(\mathbf{X}) \quad (117)$$

$$= \text{dvech}(\mathbf{X})'\mathbf{D}_n'\mathbf{H}\mathbf{D}_n\text{dvech}(\mathbf{X}) \quad (118)$$

$$\Rightarrow \frac{\text{d}^2y}{\text{dvech}(\mathbf{X})\text{dvech}(\mathbf{X})'} = \mathbf{D}_n'\mathbf{H}\mathbf{D}_n \quad (119)$$

Furthermore, it can be shown that the matrix $\mathbf{I}_{n^2} + \mathbf{K}_{nn} - \text{diag}(\text{vec}(\mathbf{I}_n))$ above is equal to $\mathbf{D}_n\mathbf{D}_n'$.

Example: Hessian of a Gaussian likelihood Let $l(\mathbf{m}, \mathbf{V})$ be the logarithm of a Gaussian likelihood at \mathbf{x} :

$$l(\mathbf{m}, \mathbf{V}) = -\frac{1}{2} \log |2\pi \mathbf{V}| - \frac{1}{2} (\mathbf{x} - \mathbf{m})' \mathbf{V}^{-1} (\mathbf{x} - \mathbf{m}) \quad (120)$$

By (107), the first term has Hessian

$$\frac{1}{2} \mathbf{D}_n \mathbf{D}_n' (\mathbf{V}^{-1} \otimes \mathbf{V}^{-1}) \mathbf{D}_n \mathbf{D}_n' \quad (121)$$

where the symmetry of \mathbf{V} has been invoked.

The differential of the second term is

$$(\mathbf{x} - \mathbf{m})' \mathbf{V}^{-1} d\mathbf{m} + \frac{1}{2} \text{tr}(\mathbf{V}^{-1} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})' \mathbf{V}^{-1} d\mathbf{V}) \quad (122)$$

and the second differential is

$$\begin{aligned} d^2 &= -d\mathbf{m}' \mathbf{V}^{-1} d\mathbf{m} - \text{tr}((d\mathbf{V}) \mathbf{V}^{-1} (d\mathbf{m})(\mathbf{x} - \mathbf{m})' \mathbf{V}^{-1}) \\ &\quad - \text{tr}((d\mathbf{V}) \mathbf{V}^{-1} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})' \mathbf{V}^{-1} (d\mathbf{V}) \mathbf{V}^{-1}) \end{aligned} \quad (123)$$

$$\begin{aligned} &= -d\mathbf{m}' \mathbf{V}^{-1} d\mathbf{m} - \text{vec}(d\mathbf{V})' (\mathbf{V}^{-1} (\mathbf{x} - \mathbf{m}) \otimes \mathbf{V}^{-1}) \text{vec}(d\mathbf{m}) \\ &\quad - \text{vec}(d\mathbf{V})' (\mathbf{V}^{-1} \otimes \mathbf{V}^{-1} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})' \mathbf{V}^{-1}) \text{vec}(d\mathbf{V}) \end{aligned} \quad (124)$$

So the full Hessian is

$$\frac{d^2 l(\mathbf{m}, \mathbf{V})}{d\mathbf{m} d\mathbf{m}'} = -\mathbf{V}^{-1} \quad (125)$$

$$\frac{d^2 l(\mathbf{m}, \mathbf{V})}{d\mathbf{m} d\text{vec}(\mathbf{V})'} = -\mathbf{D}_n \mathbf{D}_n' (\mathbf{V}^{-1} (\mathbf{x} - \mathbf{m}) \otimes \mathbf{V}^{-1}) \quad (126)$$

$$\frac{d^2 l(\mathbf{m}, \mathbf{V})}{d\text{vec}(\mathbf{V}) d\text{vec}(\mathbf{V})'} = \mathbf{D}_n \mathbf{D}_n' (\mathbf{V}^{-1} \otimes (\frac{1}{2} \mathbf{V}^{-1} - \mathbf{V}^{-1} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})' \mathbf{V}^{-1})) \mathbf{D}_n \mathbf{D}_n' \quad (127)$$

And the reduced Hessian involving unique elements of \mathbf{V} is

$$\frac{d^2 l(\mathbf{m}, \mathbf{V})}{d\mathbf{m} d\text{vech}(\mathbf{V})'} = -\mathbf{D}_n' (\mathbf{V}^{-1} (\mathbf{x} - \mathbf{m}) \otimes \mathbf{V}^{-1}) \quad (128)$$

$$\frac{d^2 l(\mathbf{m}, \mathbf{V})}{d\text{vech}(\mathbf{V}) d\text{vech}(\mathbf{V})'} = \mathbf{D}_n' (\mathbf{V}^{-1} \otimes (\frac{1}{2} \mathbf{V}^{-1} - \mathbf{V}^{-1} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})' \mathbf{V}^{-1})) \mathbf{D}_n \quad (129)$$

Acknowledgements

Tony Jebara and Martin Szummer helped clarify the presentation.

References

- [1] C. T. J. Dodson and T. Poston. *Tensor geometry: the geometric viewpoint and its uses*. Springer-Verlag, 1991.
- [2] Roger A. Horn and Charles R. Johnson. *Topics in Matrix Analysis*. Cambridge University Press, 1991.
- [3] Harold Jeffreys. *Theory of Probability*. Clarendon Press, Oxford, third edition, 1961.
- [4] Peter Lancaster and Miron Tismenetsky. *The theory of matrices*. Academic Press, 1985.
- [5] Jan R. Magnus and Heinz Neudecker. *Matrix differential calculus with applications in statistics and econometrics*. John Wiley & Sons, 1988.
- [6] David H. Marimont and Brian A. Wandell. Linear models of surface and illuminant spectra. *Journal of the Optical Society of America*, 9(11):1905–1913, November 1992.
- [7] Barak A. Pearlmutter and Lucas C. Parra. Maximum likelihood blind source separation: A context-sensitive generalization of ICA. In Michael C. Mozer, Michael Jordan, and Thomas Petsche, editors, *Advances in Neural Information Processing Systems 9*, Cambridge, MA, 1997. MIT Press. <http://www.cs.unm.edu/~bap/publications.html>.
- [8] Viktor V. Prasolov. *Problems and Theorems in Linear Algebra*. American Mathematical Society, 1991.
- [9] G. Scott and H. Longuet-Higgins. An algorithm for associating the features of two images. *Proc. Royal Society of London*, B(244):21–26, 1991.
- [10] Shayle R. Searle. *Matrix Algebra Useful for Statistics*. John Wiley & Sons, New York, NY, 1982.
- [11] Joshua B. Tenenbaum and William T. Freeman. Separating style and content. In Michael C. Mozer, Michael Jordan, and Thomas Petsche, editors, *Advances in Neural Information Processing Systems 9*, Cambridge, MA, 1997. MIT Press.