

WITH(OUT) A TRACE
MATRIX DERIVATIVES THE EASY WAY

Steven W. Nydick

University of Minnesota

May 16, 2012

OUTLINE

1 INTRODUCTION

- Notation
- History of Paper

2 TRACES

- Algebraic Trace Properties
- Calculus Trace Properties

3 TRACE DERIVATIVES

- Directional Derivatives
- Example 1: $\text{tr}(\mathbf{A}\mathbf{X})$
- Example 2: $\text{tr}(\mathbf{X}^T \mathbf{A}\mathbf{X}\mathbf{B})$
- Example 3: $\text{tr}(\mathbf{Y}^{-1})$
- Example 4: $|\mathbf{Y}|$

4 TRACE DERIVATIVE APPLICATIONS

- Application 1: Least Squares
- Application 2: Restricted Least Squares ($\mathbf{X} = \mathbf{X}^T$)
- Application 3: MLE Factor Analysis (LRC)

5 REFERENCES

NOTATION

- \mathbf{A} : A matrix
- \mathbf{A}_c : A matrix **held constant**

- \mathbf{x} : A vector
- y : A scalar (or a scalar function)

- \mathbf{x}^T or \mathbf{X}^T : The transpose of \mathbf{x} or \mathbf{X}
- x_{ij} : The element in the i th row and j th column of \mathbf{X}
- $(x^T)_{ij}$: The element in the i th row and j th column of \mathbf{X}^T

- $\frac{\partial \mathbf{Y}}{\partial \mathbf{x}}$: A matrix with elements $\frac{\partial y_{ij}}{\partial x}$
- $\frac{\partial y}{\partial \mathbf{X}}$: A matrix with elements $\frac{\partial y}{\partial x_{ij}}$

- $\langle \mathbf{x} \rangle_i$ or $\langle \mathbf{X} \rangle_{ij}$: The i th or ij th **place** of \mathbf{x} or \mathbf{X}

GRADIENT, JACOBIAN, HESSIAN

A Gradient is the derivative of a scalar with respect to a vector.

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \left(\left[\frac{\partial f(\mathbf{x})}{\partial x_1} \right] \quad \left[\frac{\partial f(\mathbf{x})}{\partial x_2} \right] \quad \cdots \quad \left[\frac{\partial f(\mathbf{x})}{\partial x_n} \right] \right)^T$$

If we have the function: $f(\mathbf{x}) = 2x_1x_2 + x_2^2 + x_1x_3^2$, then the Gradient is

$$\begin{aligned} \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} &= \left(\left[\frac{\partial f(\mathbf{x})}{\partial x_1} \right] \quad \left[\frac{\partial f(\mathbf{x})}{\partial x_2} \right] \quad \left[\frac{\partial f(\mathbf{x})}{\partial x_3} \right] \right)^T \\ &= [2x_2 + x_3^2 \quad 2x_1 + 2x_2 \quad 2x_1x_3]^T \end{aligned}$$

GRADIENT, JACOBIAN, HESSIAN

A Jacobian is the derivative of a vector with respect to a transposed vector.

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}^T} = \begin{pmatrix} \left[\frac{\partial f_1(\mathbf{x})}{\partial x_1} \right] & \cdots & \left[\frac{\partial f_1(\mathbf{x})}{\partial x_n} \right] \\ \vdots & \cdots & \vdots \\ \left[\frac{\partial f_k(\mathbf{x})}{\partial x_1} \right] & \cdots & \left[\frac{\partial f_k(\mathbf{x})}{\partial x_n} \right] \end{pmatrix}$$

If we have the function

$$\mathbf{f}(\mathbf{x}) = [3x_1^2 + x_2 \quad \ln(x_1) \quad \sin(x_2)]^T$$

Then the Jacobian is

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}^T} = \begin{pmatrix} 6x_1 & 1 \\ \frac{1}{x_1} & 0 \\ 0 & \cos(x_2) \end{pmatrix}$$

GRADIENT, JACOBIAN, HESSIAN

The Hessian is derivative of a Gradient with respect to a transposed vector.

$$\frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T} = \begin{pmatrix} \left[\frac{\partial f(\mathbf{x})}{\partial x_1^2} \right] & \cdots & \left[\frac{\partial f(\mathbf{x})}{\partial x_1 \partial x_n} \right] \\ \vdots & \ddots & \vdots \\ \left[\frac{\partial f(\mathbf{x})}{\partial x_n \partial x_1} \right] & \cdots & \left[\frac{\partial f(\mathbf{x})}{\partial x_n^2} \right] \end{pmatrix}$$

Because our above Gradient is

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = [2x_2 + x_3^2 \quad 2x_1 + 2x_2 \quad 2x_1x_3]^T$$

The Hessian would be

$$\frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T} = \begin{pmatrix} 0 & 2 & 2x_3 \\ 2 & 2 & 0 \\ 2x_3 & 0 & 2x_1 \end{pmatrix}$$

SIMPLIFYING CLASSES OF MATRIX DERIVATIVES

History of Schöneman's paper:

- 1 Wrote it while a post doc at UNC.
- 2 Originally submitted it to Psychometrika in 1965.
- 3 Editor mildly criticized paper.
 - 1 Compliment: reformulate certain problems (Lagrange multipliers) into interesting form (traces).
 - 2 Complaint: why would we want to do that?
- 4 Revised paper, resubmitted paper, but editorship changed hands, and took them almost a year to respond (asking for another revision).
 - 1 The new editor told him that a reviewer said: "nothing wrong with paper but not too important".

SIMPLIFYING CLASSES OF MATRIX DERIVATIVES

History of Schöneman's paper:

- ⑤ Later learned that the original delay was caused by a statistician with expertise in matrix derivatives who thought that the paper would be published eventually.
- ⑥ The paper **was** published eventually ... **20 years later** in MBR.
- ⑦ Wrote the article “Better Never than Late: Peer Review and the Preservation of Prejudice” in 2001.

SIMPLIFYING CLASSES OF MATRIX DERIVATIVES

There are two beneficial properties of Schöneman's paper:

- ① Derivatives are always in matrix form.
- ② No need for Dummy Matrices.

But, uses traces, and thus, uses trace properties.

So ... A Review of Traces/Trace Properties:

WHAT IS A TRACE?

Definition:

$$\operatorname{tr}(\mathbf{Y}) = \sum_i (y_{ii}), \quad \mathbf{Y} \text{ is square}$$

OK - that's simple, but what does that mean?

Well, take a square matrix and add up the diagonal elements

$$\mathbf{A} = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix}, \quad \operatorname{tr}(\mathbf{A}) = a_{11} + a_{22} + \cdots + a_{nn}$$

LINEARITY OF TRACES

The **MOST** important aspect of traces (for our later derivations):

$$\text{tr} : M(\mathbb{R})^n \rightarrow \mathbb{R}^1 \quad \text{is linear}$$

Thus:

$$\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B}) \quad (1)$$

and

$$\text{tr}(c\mathbf{A}) = c \text{tr}(\mathbf{A}) \quad (2)$$

TRANSPOSITION OF DEPENDENT VARIABLE

Traces have **SEVERAL OTHER** important properties.

Property 1: Transposition of Dependent Variable

We have:

$$\operatorname{tr}(\mathbf{Y}) = \operatorname{tr}(\mathbf{Y}^T) \quad (3)$$

Thus:

$$\frac{\partial \operatorname{tr}(\mathbf{Y})}{\partial \mathbf{X}} = \frac{\partial \operatorname{tr}(\mathbf{Y}^T)}{\partial \mathbf{X}}$$

CYCLIC PERMUTATION

Property 2: Cyclic Permutation

We have:

$$\operatorname{tr}(\mathbf{AB}) = \operatorname{tr}(\mathbf{BA}) \quad (4)$$

Why? Well, start from the left of Equation (4).

$$\begin{aligned} \operatorname{tr}(\mathbf{AB}) &= \operatorname{tr} \left[\begin{pmatrix} a_{11} & \cdots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nm} \end{pmatrix} \begin{pmatrix} b_{11} & \cdots & b_{1n} \\ \vdots & \ddots & \vdots \\ b_{m1} & \cdots & b_{mn} \end{pmatrix} \right] \\ &= a_{11}b_{11} + \cdots + a_{1m}b_{m1} + \sum_{i=1}^m (a_{2i}b_{i2}) + \cdots + \sum_{i=1}^m (a_{ni}b_{in}) \\ &= \sum_{j=1}^n \sum_{i=1}^m (a_{ji}b_{ij}) \end{aligned}$$

CYCLIC PERMUTATION

And also start from the right of Equation (4).

$$\begin{aligned} \operatorname{tr}(\mathbf{BA}) &= \operatorname{tr} \left[\begin{pmatrix} b_{11} & \cdots & b_{1n} \\ \vdots & \ddots & \vdots \\ b_{m1} & \cdots & b_{mn} \end{pmatrix} \begin{pmatrix} a_{11} & \cdots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nm} \end{pmatrix} \right] \\ &= \sum_{i=1}^m \sum_{j=1}^n (b_{ij} a_{ji}) = \sum_{j=1}^n \sum_{i=1}^m (b_{ij} a_{ji}) = \sum_{j=1}^n \sum_{i=1}^m (a_{ji} b_{ij}) \\ &= \operatorname{tr}(\mathbf{AB}) \end{aligned}$$

So:

$$\frac{\partial \operatorname{tr}(\mathbf{AB})}{\partial \mathbf{X}} = \frac{\partial \operatorname{tr}(\mathbf{BA})}{\partial \mathbf{X}}$$

Rotating the order **does not** change the trace of square matrices.

CYCLIC PERMUTATION

A consequence of the last derivation:

- Let \mathbf{U} and \mathbf{H} have the same dimensions.
- If you want to multiply paired entries (e.g. $u_{ij}h_{ij}$) and add all the multiplications:
 - Flip one of the matrices, multiply, and take the trace of that multiplication.

$$\sum_{j=1}^m \sum_{i=1}^n (u_{ij}h_{ij}) = \sum_{j=1}^m \sum_{i=1}^n \left((u^T)_{ji} h_{ij} \right) = \text{tr}(\mathbf{U}^T \mathbf{H}) \quad (5)$$

TRANSPPOSITION OF INDEPENDENT VARIABLE

Calculus Property 1: Transposition of Independent Variable

By definition:

$$\frac{\partial \text{tr}(\mathbf{Y})}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{Y})}{\partial x_{ij}}, \quad i = 1, \dots, n; j = 1, \dots, m$$

where $\frac{\partial \text{tr}(\mathbf{Y})}{\partial x_{ij}}$ is what we put in the ij th place in our derivative matrix.

Thus:

$$\frac{\partial \text{tr}(\mathbf{Y})}{\partial(\mathbf{X}^T)} = \frac{\partial \text{tr}(\mathbf{Y})}{\partial x_{ji}}, \quad (j = 1, \dots, m; i = 1, \dots, n) = \left(\frac{\partial \text{tr}(\mathbf{Y})}{\partial \mathbf{X}} \right)^T \quad (6)$$

because $\frac{\partial \text{tr}(\mathbf{Y})}{\partial x_{ji}}$ is what we put in the ij th place in our derivative matrix.

TRANSPOSITION OF INDEPENDENT VARIABLE

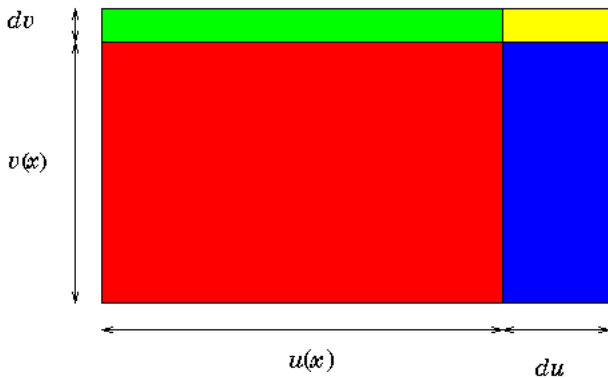
Deriving with respect to a transposed variable replaces each entry in the new matrix with the **derivative** of the **corresponding transposed component**.

Replacing every entry with the derivative of the transposed component
→ Transposing the entire matrix of partial derivatives.

PRODUCT RULE

Calculus Property 2: Product Rule

An illustration of the product rule:



PRODUCT RULE

Calculus Property 2: Product Rule

Based on the previous illustration:

$$\frac{d(uv)}{dx} = \left(\frac{du}{dx}\right)(v) + (u)\left(\frac{dv}{dx}\right)$$

In this case, u and v are scalar **functions** of x .

Now, we want to translate this to matrices and traces of matrices:

Pick any row \mathbf{u}_i and any column \mathbf{v}_j from \mathbf{U} and \mathbf{V} .

PRODUCT RULE

If we take the derivative of the matrix product with respect to a scalar:

$$\frac{\partial(\mathbf{UV})}{\partial x}$$

we find that the i,j th place in our new derivative matrix is

$$\begin{aligned}\frac{\partial(\mathbf{u}_i^T \mathbf{v}_j)}{\partial x} &= \frac{\partial(u_{i1}v_{1j} + \cdots + u_{in}v_{nj})}{\partial x} \\ &= \frac{\partial(u_{i1}v_{1j})}{\partial x} + \cdots + \frac{\partial(u_{in}v_{nj})}{\partial x} \\ &= \frac{\partial u_{i1}}{\partial x} v_{1j} + u_{i1} \frac{\partial v_{1j}}{\partial x} + \cdots + \frac{\partial u_{in}}{\partial x} v_{nj} + u_{in} \frac{\partial v_{nj}}{\partial x}\end{aligned}$$

PRODUCT RULE

So, now we want to collect terms:

$$\begin{aligned}
 \frac{\partial(\mathbf{u}_i^T \mathbf{v}_{\cdot j})}{\partial x} &= \frac{\partial u_{i1}}{\partial x} v_{1j} + u_{i1} \frac{\partial v_{1j}}{\partial x} + \cdots + \frac{\partial u_{in}}{\partial x} v_{nj} + u_{in} \frac{\partial v_{nj}}{\partial x} \\
 &= \left(\frac{\partial u_{i1}}{\partial x} v_{1j} + \cdots + \frac{\partial u_{in}}{\partial x} v_{nj} \right) + \left(u_{i1} \frac{\partial v_{1j}}{\partial x} + \cdots + u_{in} \frac{\partial v_{nj}}{\partial x} \right) \\
 &= \frac{\partial \mathbf{u}_i^T}{\partial x} \mathbf{v}_{\cdot j} + \mathbf{u}_i^T \frac{\partial \mathbf{v}_{\cdot j}}{\partial x}
 \end{aligned}$$

And because our element is **arbitrary**, we can generalize:

$$\begin{aligned}
 \frac{\partial(\mathbf{UV})}{\partial x} &= \frac{\partial \mathbf{U}}{\partial x} \mathbf{V} + \mathbf{U} \frac{\partial \mathbf{V}}{\partial x} \\
 &= \frac{\partial(\mathbf{UV}_c)}{\partial x} + \frac{\partial(\mathbf{U}_c \mathbf{V})}{\partial x}
 \end{aligned} \tag{7}$$

PRODUCT RULE

There are two notes on the product rule:

Note 1: For the product rule to make sense, both \mathbf{U} and \mathbf{V} should be functions of \mathbf{X} .

For the Univariate Case, let $u = x^2 + 2$ and $v = 2x + \sin(x)$. Then:

$$\begin{aligned}\frac{d(uv)}{dx} &= \left(\frac{du}{dx}\right)(v) + (u)\left(\frac{dv}{dx}\right) \\ &= (2x)(2x + \sin(x)) + (x^2 + 2)(2 + \cos(x)) \\ &= 4x^2 + 2x \sin(x) + 2x^2 + 4 + x^2 \cos(x) + 2 \cos(x) \\ &= x^2(6 + \cos(x)) + 2(x \sin(x) + \cos(x)) + 4\end{aligned}$$

We will discuss the multivariate case later.

PRODUCT RULE

There are two notes on the product rule:

Note 2: We can put \mathbf{V}_c and \mathbf{U}_c inside the derivative function, but they are now **constants** with respect to \mathbf{X} , even if they are functions of \mathbf{X} .

For the Univariate Case, let $u = x^3 + \ln(x)$ and $v = 3x^2$. Then:

$$\begin{aligned}\frac{d(uv_c)}{dx} &= \frac{d \left[\left(x^3 + \ln(x) \right) (3x^2)_c \right]}{dx} \\ &= \left(3x^2 + \frac{1}{x} \right) (3x^2) \\ &= 9x^4 + 3x\end{aligned}$$

We will discuss the multivariate case later.

MULTIDIMENSIONAL CHAIN RULE

Calculus Property 3: Chain Rule

Let:

$$z = 2x_1^2 + x_1 \cos(x_2)$$

Then, by definition:

$$\frac{\partial z}{\partial \mathbf{x}} = \left(\begin{array}{c} \left[\frac{\partial z}{\partial x_1} \right] \\ \left[\frac{\partial z}{\partial x_2} \right] \end{array} \right) = \left(\begin{array}{c} 4x_1 + \cos(x_2) \\ -x_1 \sin(x_2) \end{array} \right)$$

Our partial derivative with respect to x_1 is $4x_1 + \cos(x_2)$, and our partial derivative with respect to x_2 is $-x \sin(x_2)$. Furthermore, these go in the respective parts of our derivative matrix (replacing x_1 and x_2).

MULTIDIMENSIONAL CHAIN RULE

Now if:

$$x_1 = 3t \quad \text{and} \quad x_2 = t$$

Then:

$\frac{dz}{dt}$ is now the derivative with respect to t accounted for by x_1 **and** the derivative with respect to t accounted for by x_2 .

And we account for:

$$\begin{aligned} [4(3t) + \cos(t)] \frac{\partial x_1}{\partial t} & \quad \text{by } x_1 \\ -(3t) \sin(t) \frac{\partial x_2}{\partial t} & \quad \text{by } x_2 \end{aligned}$$

MULTIDIMENSIONAL CHAIN RULE

Thus:

$$\begin{aligned}\frac{dz}{dt} &= \left(\frac{\partial z}{\partial \mathbf{x}} \right)^T \frac{\partial \mathbf{x}}{\partial t} = \sum_{i=1}^2 \left(\frac{\partial z}{\partial x_i} \frac{\partial x_i}{\partial t} \right) \\ &= [12t + \cos(t)](3) + [-(3t) \sin(t)](1) \\ &= 36t + 3 \cos(t) - 3t \sin(t)\end{aligned}$$

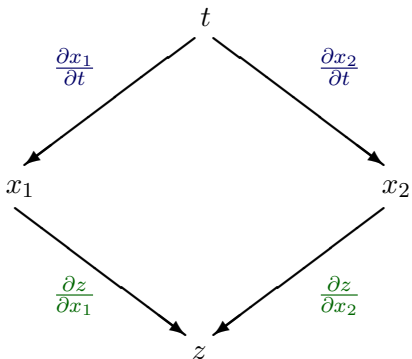
But: $z = 2x_1^2 + x_1 \cos(x_2) = 2(3t)^2 + (3t) \cos(t) = 18t^2 + 3t \cos(t)$

So, another way of getting the same result:

$$\begin{aligned}\frac{dz}{dt} &= \frac{d[18t^2 + 3t \cos(t)]}{dt} \\ &= 36t + 3t[-\sin(t)] + 3 \cos(t) = 36t + 3 \cos(t) - 3t \sin(t)\end{aligned}$$

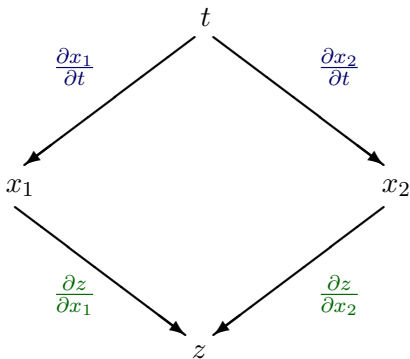
AN EASY WAY TO REMEMBER THE CHAIN RULE

Because effects are slopes and slopes are derivatives, writing out a path diagram from t to z would have the derivatives along the paths.



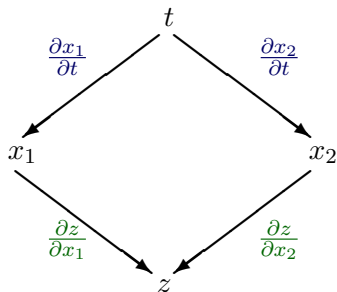
The **total effect** of t on z is found by multiplying the effects down each path and summing the total effects across paths.

AN EASY WAY TO REMEMBER THE CHAIN RULE



For example, let's find the total effect of a 1 unit change in t on z . Well, if t changes 1 unit, then x_1 changes $\frac{dx_1}{dt}$ units (because the derivative is the slope of t on x_1). Moreover, if x_1 changes 1 unit, then z changes $\frac{dz}{dx_1}$ units (because the derivative is the slope of x_1 on z).

AN EASY WAY TO REMEMBER THE CHAIN RULE



Therefore, if t changes 1 unit, then it's effect on z through x_1 would be the distance it travels in the x_1 direction:

$$x_1 \text{ distance} = \frac{dt}{dx_1} \times 1 = \frac{dt}{dx_1}$$

multiplied by how much a unit change in the x_1 direction changes z :

$$z \text{ distance through } x_1 = \frac{dz}{dx_1} \times (x_1 \text{ distance}) = \frac{dz}{dx_1} \frac{dt}{dx_1}$$

AN EASY WAY TO REMEMBER THE CHAIN RULE

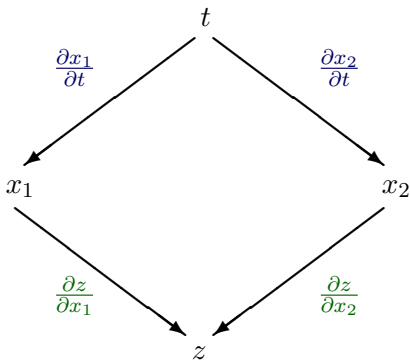
And if t changes 1 unit, then it's effect on z through x_2 would be:

$$z \text{ distance through } x_2 = \frac{dx_2}{dz} \times (x_2 \text{ distance}) = \frac{dx_2}{dz} \frac{dt}{dx_2}$$

Thus, if t moves 1 unit, it moves z : $\left(\frac{dx_1}{dz} \frac{dt}{dx_1}\right)$ through x_1 **and** it moves z : $\left(\frac{dx_2}{dz} \frac{dt}{dx_2}\right)$ through x_2 , so it **in total** moves z :

$$z \text{ total distance} = \frac{dx_1}{dz} \frac{dt}{dx_1} + \frac{dx_2}{dz} \frac{dt}{dx_2}$$

AN EASY WAY TO REMEMBER THE CHAIN RULE



Or, as written before, to find the **effect** of t on z :

$$\frac{dz}{dt} = \left(\frac{\partial z}{\partial x_1} \right) \left(\frac{\partial x_1}{\partial t} \right) + \left(\frac{\partial z}{\partial x_2} \right) \left(\frac{\partial x_2}{\partial t} \right) = \sum_{i=1}^2 \left(\frac{\partial z}{\partial x_i} \frac{\partial x_i}{\partial t} \right)$$

MATRICES CHAIN RULE

And because as our vector chain rule:

$$\frac{d[f(\mathbf{y})]}{dx} = \sum_i \left(\frac{d[f(\mathbf{y})]}{d[y_i(x)]} \frac{d[y_i(x)]}{dx} \right) \quad (8)$$

We can expand upon that to obtain a chain rule for matrices.

$$\frac{\partial[f(\mathbf{Y})]}{\partial x_{pq}} = \sum_i \sum_j \left(\frac{\partial[f(\mathbf{Y})]}{\partial[y_{ij}(x_{pq})]} \frac{\partial[y_{ij}(x_{pq})]}{\partial x_{pq}} \right) \quad (9)$$

In Equation (9) y_{ij} is a function of x_{pq} , and we have to take the derivative with respect to each of the elements in \mathbf{Y} .

DERIVATIVES

Here is the standard derivative definition:

$$Df(x) = \lim_{t \rightarrow 0} \frac{f(x+t) - f(x)}{t}$$

The equation is an infinitesimal form of $m = \frac{\Delta y}{\Delta x}$; it is finding the slope or **linear approximation** to this function as the distance between the points on the x -axis goes to 0.

If there is a large distance between points on the x -axis, and if the function is not linear, then the slope will not be a good representation of how the function is changing. However, as the distance between points on the x -axis goes to 0, the function becomes more linear.

DIRECTIONAL DERIVATIVES: VECTORS

In vector calculus, there is a similar equation.

$$D_{\mathbf{w}}f(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{w}) - f(\mathbf{x})}{t}$$

Now, our function is a surface (a scalar function of as many dimensional inputs as there are elements in \mathbf{x}), so the derivative will change at a multidimensional point \mathbf{x} based on the direction we travel from that point.

Think of what happens if you were to stand on a mountain and turn around in a circle: in some directions, the slope will be very steep (and you might fall off the mountain), but in other directions, there will barely be any slope at all.

DIRECTIONAL DERIVATIVES: VECTORS

Now \mathbf{w} tells us which direction we want to be facing when we calculate the derivative at a specific point \mathbf{x} .

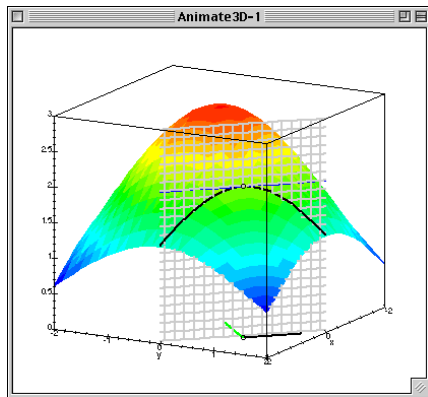
$$D_{\mathbf{w}}f(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{w}) - f(\mathbf{x})}{t}$$

The directional derivative is basically telling us what is the best **linear approximation** of this function at a particular point if we are facing up the mountain, down the mountain, at a 45 degree angle up the mountain, etc.

DIRECTIONAL DERIVATIVES: VECTORS

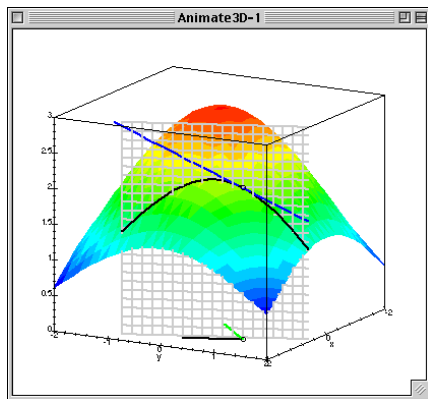
If our \mathbf{x} vector is **two** dimensional, then the function would form a mountain in **three** dimensional space.

One Direction (at a given point):



DIRECTIONAL DERIVATIVES: VECTORS

A Second Direction (at the same point):



Notice how the steepness of the slope changes at both points.

DIRECTIONAL DERIVATIVES: VECTORS

First -- pick an arbitrary unit length \mathbf{w} :

$$\mathbf{w}^T \mathbf{w} = 1$$

Second -- set up the standard, directional derivative definition:

$$D_{\mathbf{w}} f(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{w}) - f(\mathbf{x})}{t}$$

If $\mathbf{w}_{(i)} = (0, 0, 0, 0, 1, 0, \dots, 0)$ where 1 is in $\langle \mathbf{w} \rangle_i$, then

$$D_{\mathbf{w}_{(i)}} f(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{w}_{(i)}) - f(\mathbf{x})}{t}$$

will reduce to the **regular** partial derivative in the i th place.

DIRECTIONAL DERIVATIVES: VECTORS

Now, if we can find a \mathbf{u} , such that:

$$D_{\mathbf{w}}f(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{w}) - f(\mathbf{x})}{t} = \mathbf{w}^T \mathbf{u}$$

Then, for an arbitrary place i , in an arbitrary direction $\langle \mathbf{w} \rangle_i$:

$$D_{\mathbf{w}_{(i)}}f(\mathbf{x}) = \mathbf{w}_{(i)}^T \mathbf{u} \quad \text{reduces to} \quad D_{\mathbf{w}_{(i)}}f(\mathbf{x}) = u_i$$

where u_i is the partial derivative in the i th place.

Because the i th place is arbitrary:

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{u}$$

DIRECTIONAL DERIVATIVES: MATRICES

Now let $\mathbf{Y}_{(ij)}$ be a Matrix such that $y_{ij} = 1$ in $\langle \mathbf{Y} \rangle_{ij}$ and 0 elsewhere.

Then, extending our directional derivative definition to matrices:

$$D_{\mathbf{Y}} f(\mathbf{X}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{X} + t\mathbf{Y}) - f(\mathbf{X})}{t} \quad (10)$$

We can conclude that

$$D_{\mathbf{Y}_{(ij)}} f(\mathbf{X}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{X} + t\mathbf{Y}_{(ij)}) - f(\mathbf{X})}{t}$$

will “pick off” the partial derivative in the ij th place.

DIRECTIONAL DERIVATIVES: MATRICES

Now, if we can find a \mathbf{U} , such that

$$D_{\mathbf{Y}}f(\mathbf{X}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{X} + t\mathbf{Y}) - f(\mathbf{X})}{t} = \text{tr}(\mathbf{Y}^T \mathbf{U}) \quad (11)$$

Then, for an arbitrary place ij , in an arbitrary direction $\langle \mathbf{Y} \rangle_{ij}$

$$\begin{aligned} D_{\mathbf{Y}_{(ij)}}f(\mathbf{X}) &= \text{tr}(\mathbf{Y}_{(ij)}^T \mathbf{U}) = \sum_j \sum_i (y_{ij} u_{ij}) && \text{by (5)} \\ &= u_{ij} \end{aligned}$$

Because the ij th place is arbitrary:

$$\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} = \mathbf{U} \quad (12)$$

DIRECTIONAL DERIVATIVES: MATRICES

First -- put in the form of the definition:

$$D_{\mathbf{Y}}f(\mathbf{X}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{X} + t\mathbf{Y}) - f(\mathbf{X})}{t}$$

Second -- simplify until you can find the equality:

$$D_{\mathbf{Y}}f(\mathbf{X}) = \text{tr}(\mathbf{Y}^T \mathbf{U})$$

Third -- remove your \mathbf{U} , and note that:

$$\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} = \mathbf{U}$$

DEFINITION

Our 1st function:

$$f(\mathbf{X}) = \text{tr}(\mathbf{A}\mathbf{X})$$

Our objective is to find:

$$\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{A}\mathbf{X})}{\partial \mathbf{X}}$$

Only by simplifying the definition:

$$D_{\mathbf{Y}}f(\mathbf{X}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{X} + t\mathbf{Y}) - f(\mathbf{X})}{t}$$

CALCULATION

$$D_{\mathbf{Y}}f(\mathbf{X}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{X} + t\mathbf{Y}) - f(\mathbf{X})}{t} \quad \text{by (10)}$$

$$= \lim_{t \rightarrow 0} \frac{\text{tr}(\mathbf{A}[\mathbf{X} + t\mathbf{Y}]) - \text{tr}(\mathbf{A}\mathbf{X})}{t}$$

$$= \lim_{t \rightarrow 0} \frac{\text{tr}(\mathbf{A}\mathbf{X} + \mathbf{A}t\mathbf{Y}) - \text{tr}(\mathbf{A}\mathbf{X})}{t}$$

$$= \lim_{t \rightarrow 0} \frac{\text{tr}(t\mathbf{A}\mathbf{Y})}{t} \quad \text{by (1)}$$

$$= \lim_{t \rightarrow 0} \text{tr}(\mathbf{A}\mathbf{Y}) \quad \text{by (2)}$$

$$= \text{tr}(\mathbf{A}\mathbf{Y})$$

$$= \text{tr}([\mathbf{A}\mathbf{Y}]^T) \quad \text{by (3)}$$

$$= \text{tr}(\mathbf{Y}^T \mathbf{A}^T)$$

RESULT

So, we found that:

$$\begin{aligned} D_{\mathbf{Y}}f(\mathbf{X}) &= \lim_{t \rightarrow 0} \frac{f(\mathbf{X} + t\mathbf{Y}) - f(\mathbf{X})}{t} \\ &= \text{tr}(\mathbf{Y}^T \mathbf{A}^T) = \text{tr}(\mathbf{Y}^T \mathbf{U}) \end{aligned} \quad \text{by (11)}$$

And we can spot that in **this** case:

$$\mathbf{U} = \mathbf{A}^T$$

And thus, by Equation (12):

$$\begin{aligned} \mathbf{U} &= \frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} \\ &= \frac{\partial \text{tr}(\mathbf{A}\mathbf{X})}{\partial \mathbf{X}} = \mathbf{A}^T \end{aligned} \quad (13)$$

DEFINITION

Our 2nd function:

$$f(\mathbf{X}) = \text{tr}(\mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B})$$

Our objective is to find:

$$\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B})}{\partial \mathbf{X}}$$

Only by simplifying the definition:

$$D_{\mathbf{Y}} f(\mathbf{X}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{X} + t\mathbf{Y}) - f(\mathbf{X})}{t}$$

CALCULATION

$$D_{\mathbf{Y}} f(\mathbf{X}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{X} + t\mathbf{Y}) - f(\mathbf{X})}{t} \quad \text{by (10)}$$

$$= \lim_{t \rightarrow 0} \frac{\text{tr}([\mathbf{X} + t\mathbf{Y}]^T \mathbf{A} [\mathbf{X} + t\mathbf{Y}] \mathbf{B}) - \text{tr}(\mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B})}{t}$$

$$= \lim_{t \rightarrow 0} \frac{\text{tr}([\mathbf{X} + t\mathbf{Y}]^T \mathbf{A} [\mathbf{X} + t\mathbf{Y}] \mathbf{B} - \mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B})}{t} \quad \text{by (1)}$$

$$= \lim_{t \rightarrow 0} \frac{\text{tr}(\mathbf{X}^T \mathbf{A} t\mathbf{Y} \mathbf{B} + t\mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B} + t\mathbf{Y}^T \mathbf{A} t\mathbf{Y} \mathbf{B})}{t}$$

$$= \lim_{t \rightarrow 0} \frac{\text{tr}(t[\mathbf{X}^T \mathbf{A} \mathbf{Y} \mathbf{B} + \mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B} + t\mathbf{Y}^T \mathbf{A} \mathbf{Y} \mathbf{B}])}{t}$$

$$= \lim_{t \rightarrow 0} [\text{tr}(\mathbf{X}^T \mathbf{A} \mathbf{Y} \mathbf{B} + \mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B} + t\mathbf{Y}^T \mathbf{A} \mathbf{Y} \mathbf{B})] \quad \text{by (2)}$$

CALCULATION

Continuing:

$$\begin{aligned}
 D_{\mathbf{Y}} f(\mathbf{X}) &= \lim_{t \rightarrow 0} [\text{tr}(\mathbf{X}^T \mathbf{A} \mathbf{Y} \mathbf{B} + \mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B} + t \mathbf{Y}^T \mathbf{A} \mathbf{Y} \mathbf{B})] \\
 &= \lim_{t \rightarrow 0} [\text{tr}(\mathbf{X}^T \mathbf{A} \mathbf{Y} \mathbf{B} + \mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B})] + \lim_{t \rightarrow 0} [t \text{tr}(\mathbf{Y}^T \mathbf{A} \mathbf{Y} \mathbf{B})] \\
 & \qquad \qquad \qquad \text{by (1) \& (2)} \\
 &= \lim_{t \rightarrow 0} [\text{tr}(\mathbf{X}^T \mathbf{A} \mathbf{Y} \mathbf{B} + \mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B})] \\
 &= \text{tr}(\mathbf{X}^T \mathbf{A} \mathbf{Y} \mathbf{B} + \mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B}) \\
 &= \text{tr}(\mathbf{X}^T \mathbf{A} \mathbf{Y} \mathbf{B}) + \text{tr}(\mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B}) \qquad \qquad \text{by (1)} \\
 &= \text{tr}(\mathbf{B} \mathbf{X}^T \mathbf{A} \mathbf{Y}) + \text{tr}(\mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B}) \qquad \qquad \text{by (4)}
 \end{aligned}$$

CALCULATION

And Finally:

$$\begin{aligned} D_{\mathbf{Y}} f(\mathbf{X}) &= \text{tr}(\mathbf{B} \mathbf{X}^T \mathbf{A} \mathbf{Y}) + \text{tr}(\mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B}) \\ &= \text{tr}[(\mathbf{B} \mathbf{X}^T \mathbf{A} \mathbf{Y})^T] + \text{tr}(\mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B}) && \text{by (3)} \\ &= \text{tr}(\mathbf{Y}^T \mathbf{A}^T \mathbf{X} \mathbf{B}^T) + \text{tr}(\mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B}) \\ &= \text{tr}(\mathbf{Y}^T \mathbf{A}^T \mathbf{X} \mathbf{B}^T + \mathbf{Y}^T \mathbf{A} \mathbf{X} \mathbf{B}) && \text{by (1)} \\ &= \text{tr}(\mathbf{Y}^T [\mathbf{A}^T \mathbf{X} \mathbf{B}^T + \mathbf{A} \mathbf{X} \mathbf{B}]) \end{aligned}$$

RESULT

So, we found that:

$$\begin{aligned} D_{\mathbf{Y}} f(\mathbf{X}) &= \lim_{t \rightarrow 0} \frac{f(\mathbf{X} + t\mathbf{Y}) - f(\mathbf{X})}{t} \\ &= \text{tr}(\mathbf{Y}^T [\mathbf{A}^T \mathbf{X} \mathbf{B}^T + \mathbf{A} \mathbf{X} \mathbf{B}]) = \text{tr}(\mathbf{Y}^T \mathbf{U}) \end{aligned} \quad \text{by (11)}$$

And we can spot that in **this** case:

$$\mathbf{U} = \mathbf{A}^T \mathbf{X} \mathbf{B}^T + \mathbf{A} \mathbf{X} \mathbf{B}$$

And thus, by Equation (12):

$$\begin{aligned} \mathbf{U} &= \frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} \\ &= \frac{\partial \text{tr}(\mathbf{X}^T \mathbf{A} \mathbf{X} \mathbf{B})}{\partial \mathbf{X}} = \mathbf{A}^T \mathbf{X} \mathbf{B}^T + \mathbf{A} \mathbf{X} \mathbf{B} \end{aligned} \quad (14)$$

DEFINITION

Our 3rd function, assuming that \mathbf{Y} is non-singular and depends on \mathbf{X} :

$$f(\mathbf{X}) = \text{tr}(\mathbf{Y}^{-1})$$

Our objective is to find a better expression for

$$\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{Y}^{-1})}{\partial \mathbf{X}}$$

by working with previous trace derivative rules.

CALCULATION

We have:

$$\begin{aligned} \frac{\partial \text{tr}(\mathbf{Y}^{-1})}{\partial \mathbf{X}} &= \frac{\partial \text{tr}(\mathbf{Y}^{-2}\mathbf{Y})}{\partial \mathbf{X}} \\ &= \frac{\partial \text{tr}(\mathbf{Y}_c^{-2}\mathbf{Y})}{\partial \mathbf{X}} + \boxed{\frac{\partial \text{tr}(\mathbf{Y}^{-2}\mathbf{Y}_c)}{\partial \mathbf{X}}} \quad \text{by (7)} \end{aligned}$$

[Zoom In]

$$\begin{aligned} \boxed{\frac{\partial \text{tr}(\mathbf{Y}^{-1}\mathbf{Y}^{-1}\mathbf{Y}_c)}{\partial \mathbf{X}}} &= \frac{\partial \text{tr}(\mathbf{Y}_c\mathbf{Y}_c^{-1}\mathbf{Y}^{-1})}{\partial \mathbf{X}} + \frac{\partial \text{tr}(\mathbf{Y}^{-1}\mathbf{Y}_c^{-1}\mathbf{Y}_c)}{\partial \mathbf{X}} \\ &\quad \text{by (7) \& (4)} \\ &= \frac{\partial \text{tr}(\mathbf{Y}^{-1})}{\partial \mathbf{X}} + \frac{\partial \text{tr}(\mathbf{Y}^{-1})}{\partial \mathbf{X}} = \frac{2\partial \text{tr}(\mathbf{Y}^{-1})}{\partial \mathbf{X}} \quad (15) \end{aligned}$$

RESULT

Therefore:

$$\begin{aligned} \frac{\partial \text{tr}(\mathbf{Y}^{-1})}{\partial \mathbf{X}} &= \frac{\partial \text{tr}(\mathbf{Y}_c^{-2} \mathbf{Y})}{\partial \mathbf{X}} + \frac{\partial \text{tr}(\mathbf{Y}^{-2} \mathbf{Y}_c)}{\partial \mathbf{X}} \\ \frac{\partial \text{tr}(\mathbf{Y}^{-1})}{\partial \mathbf{X}} &= \frac{\partial \text{tr}(\mathbf{Y}_c^{-2} \mathbf{Y})}{\partial \mathbf{X}} + \frac{2 \partial \text{tr}(\mathbf{Y}^{-1})}{\partial \mathbf{X}} && \text{by (15)} \\ - \frac{\partial \text{tr}(\mathbf{Y}^{-1})}{\partial \mathbf{X}} &= \frac{\partial \text{tr}(\mathbf{Y}_c^{-2} \mathbf{Y})}{\partial \mathbf{X}} \end{aligned}$$

And, finally, after multiplying by (-1) on both sides:

$$\frac{\partial \text{tr}(\mathbf{Y}^{-1})}{\partial \mathbf{X}} = - \frac{\partial \text{tr}(\mathbf{Y}_c^{-2} \mathbf{Y})}{\partial \mathbf{X}} \quad (16)$$

We have turned a “derivative of the trace-inverse” problem into a standard “trace derivative” problem.

DEFINITION

Our 4th function, assuming that \mathbf{Y} is non-singular and depends on \mathbf{X} :

$$f(\mathbf{X}) = |\mathbf{Y}|$$

Our objective is to find a better expression for

$$\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} = \frac{\partial |\mathbf{Y}|}{\partial \mathbf{X}}$$

by working with previous trace derivative rules and determinant rules.

DETERMINANT REVIEW

We know from linear algebra:

$$\mathbf{Y}^{-1} = \frac{1}{|\mathbf{Y}|} \mathbf{Q} \quad \text{where } \mathbf{Q} \text{ is the adjoint matrix}$$

Process for calculating $(q^T)_{ij}$:

- 1 Cross out row i and column j of \mathbf{Y} .
- 2 Take determinant of the smaller $[(n-1) \times (n-1)]$ matrix.
- 3 If $i+j$ is odd, then negate the previous step.

Thus:

$$|\mathbf{Y}|\mathbf{I} = \mathbf{Q}\mathbf{Y} \quad (17)$$

ADJOINT REVIEW

- Note: $(q^T)_{ij} = q_{ji}$ does not depend on any of the elements in row i or column j of \mathbf{Y} .
- Thus, q_{ji} does not depend on y_{ij} .
- So: $\frac{\partial(q_{ji}y_{ij})}{\partial y_{ij}} = q_{ji}$

Based on the previous slide we have:

$$\begin{array}{c}
 \mathbf{Y}|\mathbf{I} = \mathbf{QY} \qquad \qquad \qquad \text{by (17)} \\
 \left(\begin{array}{cccc}
 |\mathbf{Y}| & 0 & \cdots & 0 \\
 0 & |\mathbf{Y}| & \vdots & \vdots \\
 \vdots & \cdots & \ddots & 0 \\
 0 & \cdots & 0 & |\mathbf{Y}|
 \end{array} \right) = \left(\begin{array}{cccc}
 \sum_p (q_{1p}y_{p1}) & & & \mathbf{0} \\
 & \ddots & & \\
 & & \ddots & \\
 \mathbf{0} & & & \sum_p (q_{np}y_{pn})
 \end{array} \right)
 \end{array}$$

DETERMINANT DERIVATIVE, PART 1

Thus, given any j such that $1 \leq j \leq n$:

$$|\mathbf{Y}| = \sum_p (q_{jp} y_{pj}) \quad (18)$$

And, for an arbitrary y_{ij} , pick the j th row of \mathbf{q} and column of \mathbf{y} :

$$\frac{\partial |\mathbf{Y}|}{\partial y_{ij}} = \frac{\partial \left(\sum_p (q_{jp} y_{pj}) \right)}{\partial y_{ij}} \quad \text{by (18)}$$

$$= \sum_p \left(q_{jp} \frac{\partial y_{pj}}{\partial y_{ij}} \right) = q_{ji} \frac{\partial y_{ij}}{\partial y_{ij}} = q_{ji} = (q^T)_{ij} \quad (19)$$

Because y_{ij} was arbitrary, for an entire matrix:

$$\frac{\partial |\mathbf{Y}|}{\partial \mathbf{Y}} = \mathbf{Q}^T \quad (20)$$

DETERMINANT DERIVATIVE, PART 2

Now, for an arbitrary pq th element of \mathbf{X} (where \mathbf{Y} depends on \mathbf{X}):

$$\frac{\partial |\mathbf{Y}|}{\partial x_{pq}} = \sum_i \sum_j \left(\frac{\partial |\mathbf{Y}|}{\partial y_{ij}} \frac{\partial y_{ij}}{\partial x_{pq}} \right) \quad \text{by (9)}$$

$$= \sum_i \sum_j \left(q_{ji} \frac{\partial y_{ij}}{\partial x_{pq}} \right) \quad \text{by (19)}$$

$$= \frac{\partial \left(\sum_i \sum_j (q_{cji} y_{ij}) \right)}{\partial x_{pq}}$$

$$= \frac{\partial \text{tr}(\mathbf{Q}_c \mathbf{Y})}{\partial x_{pq}} \quad \text{by (5)}$$

Because x_{pq} was arbitrary, for an entire matrix:

$$\frac{\partial |\mathbf{Y}|}{\partial \mathbf{X}} = \frac{\partial \text{tr}(\mathbf{Q}_c \mathbf{Y})}{\partial \mathbf{X}} \quad (21)$$

DEFINITION

Let's say that we have

$$\mathbf{A} = \mathbf{X} + \mathbf{E}$$

where \mathbf{A} is the observation matrix, \mathbf{E} is a matrix of stochastic fluctuations with a mean of 0, and \mathbf{X} is our approximation to \mathbf{A} .

In Least Squares, our objective is to minimize the sum of squared errors:

$$\begin{aligned} SSE &= e_{11}^2 + e_{12}^2 + \cdots + e_{1n}^2 + e_{21}^2 + \cdots + e_{2n}^2 + \cdots + e_{mn}^2 \\ &= \sum_i \sum_j e_{ij}^2 \\ &= \sum_i \sum_j (e_{ij}e_{ij}) \\ &= \text{tr}(\mathbf{E}^T \mathbf{E}) \end{aligned} \quad \text{by (5)}$$

DEFINITION

If we have no constraints on \mathbf{X} , then we are, equivalently, minimizing:

$$SSE = \text{tr}(\mathbf{E}^T \mathbf{E}) = \text{tr}[(\mathbf{A} - \mathbf{X})^T (\mathbf{A} - \mathbf{X})]$$

CALCULATION

A minimization:

$$\begin{aligned}
 \frac{\partial(SSE)}{\partial \mathbf{X}} &= \frac{\partial \operatorname{tr}(\mathbf{E}^T \mathbf{E})}{\partial \mathbf{X}} = \frac{\partial \operatorname{tr}[(\mathbf{A} - \mathbf{X})^T (\mathbf{A} - \mathbf{X})]}{\partial \mathbf{X}} \\
 &= \frac{\partial \operatorname{tr}(\mathbf{A}^T \mathbf{A} - \mathbf{A}^T \mathbf{X} - \mathbf{X}^T \mathbf{A} + \mathbf{X}^T \mathbf{X})}{\partial \mathbf{X}} \\
 &= \frac{\partial [\operatorname{tr}(\mathbf{A}^T \mathbf{A}) - \operatorname{tr}(\mathbf{A}^T \mathbf{X}) - \operatorname{tr}(\mathbf{X}^T \mathbf{A}) + \operatorname{tr}(\mathbf{X}^T \mathbf{X})]}{\partial \mathbf{X}} \quad \text{by (1)} \\
 &= \frac{\partial \operatorname{tr}(\mathbf{A}^T \mathbf{A})}{\partial \mathbf{X}} - \frac{\partial \operatorname{tr}(\mathbf{A}^T \mathbf{X})}{\partial \mathbf{X}} - \frac{\partial \operatorname{tr}(\mathbf{X}^T \mathbf{A})}{\partial \mathbf{X}} + \frac{\partial \operatorname{tr}(\mathbf{X}^T \mathbf{X})}{\partial \mathbf{X}} \\
 &= \mathbf{0} - \frac{\partial \operatorname{tr}(\mathbf{A}^T \mathbf{X})}{\partial \mathbf{X}} - \frac{\partial \operatorname{tr}(\mathbf{X}^T \mathbf{A})}{\partial \mathbf{X}} + \frac{\partial \operatorname{tr}(\mathbf{X}^T \mathbf{X})}{\partial \mathbf{X}}
 \end{aligned}$$

CALCULATION

Continuing:

$$\begin{aligned}
\frac{\partial(SSE)}{\partial \mathbf{X}} &= -\frac{\partial \operatorname{tr}(\mathbf{A}^T \mathbf{X})}{\partial \mathbf{X}} - \frac{\partial \operatorname{tr}(\mathbf{X}^T \mathbf{A})}{\partial \mathbf{X}} + \frac{\partial \operatorname{tr}(\mathbf{X}^T \mathbf{X})}{\partial \mathbf{X}} \\
&= -\mathbf{A} - \frac{\partial \operatorname{tr}(\mathbf{A}^T \mathbf{X})}{\partial \mathbf{X}} + \frac{\partial \operatorname{tr}(\mathbf{X}^T \mathbf{X})}{\partial \mathbf{X}} && \text{by (13) \& (3)} \\
&= -\mathbf{A} - \mathbf{A} + \frac{\partial \operatorname{tr}(\mathbf{X}^T \mathbf{X})}{\partial \mathbf{X}} && \text{by (13)} \\
&= -\mathbf{A} - \mathbf{A} + \frac{\partial \operatorname{tr}(\mathbf{X}^T \mathbf{I} \mathbf{X} \mathbf{I})}{\partial \mathbf{X}} \\
&= -\mathbf{A} - \mathbf{A} + [\mathbf{I}^T \mathbf{X} \mathbf{I}^T + \mathbf{I} \mathbf{X} \mathbf{I}] && \text{by (14)} \\
&= -\mathbf{A} - \mathbf{A} + [\mathbf{X} + \mathbf{X}] = -2\mathbf{A} + 2\mathbf{X} && (22)
\end{aligned}$$

RESULT

As in any **Least Squares** problem, we should set our derivative equal to 0 in order to find the minimum of the function.

$$\begin{aligned}\frac{\partial(SSE)}{\partial \mathbf{X}} &= -2\mathbf{A} + 2\mathbf{X} = 0 \\ 2\mathbf{X} &= 2\mathbf{A} \\ \hat{\mathbf{A}} = \mathbf{X} &= \mathbf{A}\end{aligned}$$

Surprisingly, without **any** constraints on \mathbf{X} , the best approximation of \mathbf{A} is \mathbf{A} itself.

Oh, the things you learn in calculus ☺!

LAGRANGE MULTIPLIERS

Pretend you have a function:

$$f(\mathbf{X})$$

To maximize or minimize less than mn restraints equivalent to

$$h(x_{11}, \dots, x_{mn})_{ij} = 0$$

use LaGrange Multipliers u_{ij} (one for **each** restraint), and set

$$g(\mathbf{X}) = f(\mathbf{X}) + \sum_i \sum_j (u_{ij} h_{ij}) \quad (23)$$

Finally, take the derivative with respect to \mathbf{X} , set equal to 0, and solve.

DEFINITION

We still want to find an \mathbf{X} that minimizes the *SSE* to best approximate \mathbf{A} ; however, we are now subject to the constraint that \mathbf{X} is a Symmetric Matrix.

\mathbf{X} Symmetric Means:

$$\begin{aligned}\mathbf{X} &= \mathbf{X}^T \\ \mathbf{X} - \mathbf{X}^T &= 0\end{aligned}$$

The Recipe:

- 1 We have our equation to minimize: $\text{tr}(\mathbf{E}^T \mathbf{E})$.
- 2 We have our constraint: $\mathbf{H} = \mathbf{X} - \mathbf{X}^T = 0$.
- 3 Put in LaGrange multiplier form.
- 4 Take the derivative.
- 5 Set the derivative equal to 0.
- 6 Solve for \mathbf{X} .

CALCULATION

First -- set up the problem:

$$g(\mathbf{X}) = f(\mathbf{X}) + \sum_i \sum_j (u_{ij} h_{ij}) \quad \text{by (23)}$$

$$= \text{tr}(\mathbf{E}^T \mathbf{E}) + \text{tr}(\mathbf{U}^T \mathbf{H}) \quad \text{by (5)}$$

$$= \text{tr}(\mathbf{E}^T \mathbf{E}) + \text{tr}[\mathbf{U}^T (\mathbf{X} - \mathbf{X}^T)]$$

$$= \text{tr}(\mathbf{E}^T \mathbf{E}) + \text{tr}(\mathbf{U}^T \mathbf{X}) - \text{tr}(\mathbf{U}^T \mathbf{X}^T) \quad \text{by (1)}$$

$$= \text{tr}(\mathbf{E}^T \mathbf{E}) + \text{tr}(\mathbf{U}^T \mathbf{X}) - \text{tr}(\mathbf{U} \mathbf{X}) \quad \text{by (3) \& (4)}$$

CALCULATION

Second -- take the derivative:

$$\begin{aligned}\frac{\partial g(\mathbf{X})}{\partial \mathbf{X}} &= \frac{\partial [\text{tr}(\mathbf{E}^T \mathbf{E}) + \text{tr}(\mathbf{U}^T \mathbf{X}) - \text{tr}(\mathbf{U} \mathbf{X})]}{\partial \mathbf{X}} \\ &= \frac{\partial \text{tr}(\mathbf{E}^T \mathbf{E})}{\partial \mathbf{X}} + \frac{\partial \text{tr}(\mathbf{U}^T \mathbf{X})}{\partial \mathbf{X}} - \frac{\partial \text{tr}(\mathbf{U} \mathbf{X})}{\partial \mathbf{X}} \\ &= -2\mathbf{A} + 2\mathbf{X} + \frac{\partial \text{tr}(\mathbf{U}^T \mathbf{X})}{\partial \mathbf{X}} - \frac{\partial \text{tr}(\mathbf{U} \mathbf{X})}{\partial \mathbf{X}} && \text{by (22)} \\ &= -2\mathbf{A} + 2\mathbf{X} + \mathbf{U} - \mathbf{U}^T && \text{by (13)}\end{aligned}$$

CALCULATION

Third -- set the derivative equal to 0:

$$\frac{\partial g(\mathbf{X})}{\partial \mathbf{X}} = -2\mathbf{A} + 2\mathbf{X} + \mathbf{U} - \mathbf{U}^T = 0$$

$$2\mathbf{X} = 2\mathbf{A} + \mathbf{U}^T - \mathbf{U}$$

$$\mathbf{X} = \mathbf{A} + \frac{\mathbf{U}^T - \mathbf{U}}{2}$$

However, now note that: $\mathbf{X} = \mathbf{X}^T$

$$\mathbf{X}^T = \left(\mathbf{A} + \frac{\mathbf{U}^T - \mathbf{U}}{2} \right)^T$$

$$\mathbf{X}^T = \mathbf{A}^T + \frac{\mathbf{U} - \mathbf{U}^T}{2}$$

RESULT

Fourth -- add \mathbf{X}^T to both sides and solve for \mathbf{X} :

$$\mathbf{X} + \mathbf{X}^T = \mathbf{A} + \frac{\mathbf{U}^T - \mathbf{U}}{2} + \mathbf{A}^T + \frac{\mathbf{U} - \mathbf{U}^T}{2}$$

$$\mathbf{X} + \mathbf{X} = \mathbf{A} + \mathbf{A}^T + \frac{\mathbf{U}^T - \mathbf{U}}{2} - \frac{\mathbf{U}^T - \mathbf{U}}{2}$$

$$2\mathbf{X} = \mathbf{A} + \mathbf{A}^T$$

$$\hat{\mathbf{A}} = \mathbf{X} = \frac{\mathbf{A} + \mathbf{A}^T}{2}$$

Therefore, to approximate \mathbf{A} with a Symmetric Matrix, the **best** matrix (according to the Least Squares Criterion) is the **average** of the elements of \mathbf{A} and the elements of \mathbf{A}^T .

DEFINITION

The function we want to maximize:

$$u = |\mathbf{U}^{-1}(\mathbf{R} - \mathbf{F}\mathbf{F}^T)\mathbf{U}^{-1}|$$

- ① \mathbf{R} is a correlation matrix.
 - $\text{diag}(\mathbf{R}) = \mathbf{I}$
- ② \mathbf{F} is a factor pattern matrix of **uncorrelated** common factors.
- ③ \mathbf{U}^2 is a covariance matrix of **uncorrelated** unique factors.
 - $\mathbf{U}^2 = \text{diag}(\mathbf{U}^2) = \mathbf{I} - \text{diag}(\mathbf{F}\mathbf{F}^T)$

DEFINITION

The function we want to maximize:

$$u = |\mathbf{U}^{-1}(\mathbf{R} - \mathbf{F}\mathbf{F}^T)\mathbf{U}^{-1}|$$

The function u is a likelihood ratio criterion for a test of independence after the common factors have been partialled out of the covariance matrix.

$\mathbf{U}^{-1}(\mathbf{R} - \mathbf{F}\mathbf{F}^T)\mathbf{U}^{-1}$ should be close to \mathbf{I} , so u should be close to 1.

We want to find the \mathbf{F} (and consequently the \mathbf{U}^2) that results in a determinant as close to 1 as possible.

DEFINITION

To make the derivatives simpler, let:

$$u_1 = |\mathbf{U}^{-2}| \quad \text{and} \quad u_2 = |\mathbf{R} - \mathbf{F}\mathbf{F}^T|$$

Note that:

$$u_1 u_2 = |\mathbf{U}^{-2}| |\mathbf{R} - \mathbf{F}\mathbf{F}^T| = |\mathbf{U}^{-1}| |\mathbf{R} - \mathbf{F}\mathbf{F}^T| |\mathbf{U}^{-1}| = |\mathbf{U}^{-1}(\mathbf{R} - \mathbf{F}\mathbf{F}^T)\mathbf{U}^{-1}|$$

Thus, we can use the product rule to find the derivative:

$$\begin{aligned} \frac{\partial u}{\partial \mathbf{F}} &= \frac{\partial(u_1 u_2)}{\partial \mathbf{F}} \\ &= \frac{\partial u_1}{\partial \mathbf{F}} u_2 + u_1 \frac{\partial u_2}{\partial \mathbf{F}} \end{aligned} \quad \text{by (7)}$$

CALCULATION: $\frac{\partial u_1}{\partial \mathbf{F}}$

Let's find the derivative of the first part.

$$\begin{aligned}
 \frac{\partial u_1}{\partial \mathbf{F}} &= \frac{\partial |\mathbf{U}^{-2}|}{\partial \mathbf{F}} = \frac{\partial |\mathbf{U}^2|^{-1}}{\partial \mathbf{F}} && \text{(by determinant rules)} \\
 &= \frac{\partial \text{tr}(|\mathbf{U}^2|^{-1})}{\partial \mathbf{F}} && \text{(since } \text{tr}(|\mathbf{X}|) = |\mathbf{X}|) \\
 &= -\frac{\text{tr}(|\mathbf{U}^2|^{-2} |\mathbf{U}^2|)}{\partial \mathbf{F}} && \text{by (16)} \\
 &= -|\mathbf{U}^2|^{-2} \frac{\partial |\mathbf{U}^2|}{\partial \mathbf{F}} && (24)
 \end{aligned}$$

Our next objective is to find the derivative of the highlighted part.

CALCULATION: $\frac{\partial u_1}{\partial \mathbf{F}}$

Continuing:

$$\begin{aligned} \frac{\partial |\mathbf{U}|}{\partial \mathbf{F}} &= \frac{\partial |\mathbf{I} - \text{diag}(\mathbf{F}\mathbf{F}^T)|}{\partial \mathbf{F}} && \text{(by definition)} \\ &= \frac{\partial \text{tr} \left(\mathbf{Q}_c [\mathbf{I} - \text{diag}(\mathbf{F}\mathbf{F}^T)] \right)}{\partial \mathbf{F}} && \text{by (21)} \\ &= \frac{\partial \text{tr}(\mathbf{Q}_c)}{\partial \mathbf{F}} - \frac{\partial \text{tr}[\mathbf{Q}_c \text{diag}(\mathbf{F}\mathbf{F}^T)]}{\partial \mathbf{F}} && \text{by (1) \& (2)} \\ &= 0 - \frac{\partial \text{tr}[\mathbf{Q}_c \text{diag}(\mathbf{F}\mathbf{F}^T)]}{\partial \mathbf{F}} \\ &= - \frac{\partial \text{tr}[\mathbf{Q}_c \text{diag}(\mathbf{F}\mathbf{F}^T)]}{\partial \mathbf{F}} \end{aligned}$$

CALCULATION: $\frac{\partial u_1}{\partial \mathbf{F}}$

Based on (21), \mathbf{Q}_c is the Adjoint of $[\mathbf{I} - \text{diag}(\mathbf{F}\mathbf{F}^T)] = \mathbf{U}^2$.

Therefore:

$$\begin{aligned} |\mathbf{U}^2|^{-1} \mathbf{Q}_c &= (\mathbf{U}^2)^{-1} && \text{by (17)} \\ \mathbf{Q}_c &= |\mathbf{U}^2|(\mathbf{U}^{-2}) \end{aligned}$$

Because \mathbf{U}^2 is a diagonal matrix, $\mathbf{Q}_c = |\mathbf{U}^2|(\mathbf{U}^{-2})$ is a diagonal matrix.

And:

$$\frac{\partial |\mathbf{U}|}{\partial \mathbf{F}} = -\frac{\partial \text{tr}[\mathbf{Q}_c \text{diag}(\mathbf{F}\mathbf{F}^T)]}{\partial \mathbf{F}} = -\frac{\partial \text{tr}[\text{diag}(\mathbf{Q}_c \mathbf{F}\mathbf{F}^T)]}{\partial \mathbf{F}} = -\frac{\partial \text{tr}(\mathbf{Q}_c \mathbf{F}\mathbf{F}^T)}{\partial \mathbf{F}}$$

Because the trace only **operates on the diagonal**, the trace of the diagonal of a matrix is the same as the trace of the original matrix.

CALCULATION: $\frac{\partial u_1}{\partial \mathbf{F}}$

Continuing:

$$\begin{aligned}
-\frac{\partial \operatorname{tr}(\mathbf{Q}_c \mathbf{F} \mathbf{F}^T)}{\partial \mathbf{F}} &= -\frac{\partial \operatorname{tr}(\mathbf{F}^T \mathbf{Q}_c \mathbf{F} \mathbf{I})}{\partial \mathbf{F}} && \text{by (4)} \\
&= -(\mathbf{Q}^T \mathbf{F} \mathbf{I}^T + \mathbf{Q} \mathbf{F} \mathbf{I}) && \text{by (14)} \\
&= -(\mathbf{Q}^T + \mathbf{Q}) \mathbf{F} \\
&= -2\mathbf{Q} \mathbf{F} && (\mathbf{Q} \text{ is symmetric}) \\
&= -2|\mathbf{U}^2| \mathbf{U}^{-2} \mathbf{F} && \text{by (17)}
\end{aligned}$$

And, thus:

$$\begin{aligned}
\frac{\partial u_1}{\partial \mathbf{F}} &= -|\mathbf{U}^2| \frac{\partial |\mathbf{U}|^2}{\partial \mathbf{F}} && \text{by (24)} \\
&= -|\mathbf{U}^2|^{-2} (-2|\mathbf{U}^2| \mathbf{U}^{-2} \mathbf{F}) \\
&= 2|\mathbf{U}^2|^{-1} \mathbf{U}^{-2} \mathbf{F} \\
&= 2|\mathbf{U}^{-2}| \mathbf{U}^{-2} \mathbf{F} && (25)
\end{aligned}$$

CALCULATION: $\frac{\partial u_2}{\partial \mathbf{F}}$

Now, let's find the derivative of the second part.

$$\begin{aligned}
 \frac{\partial u_2}{\partial \mathbf{F}} &= \frac{\partial |\mathbf{R} - \mathbf{F}\mathbf{F}^T|}{\partial \mathbf{F}} \\
 &= \frac{\partial \text{tr}[\mathbf{Q}_c(\mathbf{R} - \mathbf{F}\mathbf{F}^T)]}{\partial \mathbf{F}} && \text{by (21)} \\
 &= \frac{\partial \text{tr}(\mathbf{Q}_c \mathbf{R})}{\partial \mathbf{F}} - \frac{\partial \text{tr}(\mathbf{Q}_c \mathbf{F}\mathbf{F}^T)}{\partial \mathbf{F}} && \text{by (1) \& (2)} \\
 &= \mathbf{0} - \frac{\partial \text{tr}(\mathbf{F}^T \mathbf{Q}_c \mathbf{F} \mathbf{I})}{\partial \mathbf{F}} && \text{by (4)} \\
 &= -(\mathbf{Q}^T + \mathbf{Q})\mathbf{F} && \text{by (14)} \\
 &= -2\mathbf{Q}\mathbf{F} && (\mathbf{Q} \text{ is symmetric}) \\
 &= -2|\mathbf{R} - \mathbf{F}\mathbf{F}^T|(\mathbf{R} - \mathbf{F}\mathbf{F}^T)^{-1}\mathbf{F} && (26)
 \end{aligned}$$

CALCULATION: ENTIRE THING

Putting the pieces together:

$$\frac{\partial u}{\partial \mathbf{F}} = \frac{\partial u_1}{\partial \mathbf{F}} u_2 + u_1 \frac{\partial u_2}{\partial \mathbf{F}}$$

Which implies that

$$\begin{aligned} \frac{\partial u}{\partial \mathbf{F}} = & (2|\mathbf{U}^{-2}| \mathbf{U}^{-2} \mathbf{F}) |\mathbf{R} - \mathbf{F}\mathbf{F}^T| \\ & + |\mathbf{U}^{-2}| (-2|\mathbf{R} - \mathbf{F}\mathbf{F}^T| (\mathbf{R} - \mathbf{F}\mathbf{F}^T)^{-1} \mathbf{F}) \end{aligned}$$

CALCULATION: FINDING MAXIMUM

To find the maximum of this function, we must set it equal to 0 and solve.

$$0 = 2|\mathbf{U}^{-2}|(\mathbf{U}^{-2}\mathbf{F})|\mathbf{R} - \mathbf{F}\mathbf{F}^T| - 2|\mathbf{U}^{-2}||\mathbf{R} - \mathbf{F}\mathbf{F}^T|(\mathbf{R} - \mathbf{F}\mathbf{F}^T)^{-1}\mathbf{F}$$

$$0 = 2|\mathbf{U}^{-2}||\mathbf{R} - \mathbf{F}\mathbf{F}^T|(\mathbf{U}^{-2}\mathbf{F} - (\mathbf{R} - \mathbf{F}\mathbf{F}^T)^{-1}\mathbf{F})$$

$$0 = \mathbf{U}^{-2}\mathbf{F} - (\mathbf{R} - \mathbf{F}\mathbf{F}^T)^{-1}\mathbf{F}$$

CALCULATION: FINDING MAXIMUM

Finishing the calculation:

$$\begin{aligned}0 &= \mathbf{U}^{-2}\mathbf{F} - (\mathbf{R} - \mathbf{F}\mathbf{F}^T)^{-1}\mathbf{F} \\(\mathbf{R} - \mathbf{F}\mathbf{F}^T)^{-1}\mathbf{F} &= \mathbf{U}^{-2}\mathbf{F} \\ \mathbf{F} &= (\mathbf{R} - \mathbf{F}\mathbf{F}^T)\mathbf{U}^{-2}\mathbf{F} \\ \mathbf{F} &= \mathbf{R}\mathbf{U}^{-2}\mathbf{F} - \mathbf{F}\mathbf{F}^T\mathbf{U}^{-2}\mathbf{F} \\ \mathbf{R}\mathbf{U}^{-2}\mathbf{F} - \mathbf{F} &= \mathbf{F}\mathbf{F}^T\mathbf{U}^{-2}\mathbf{F} \\ (\mathbf{R}\mathbf{U}^{-2} - \mathbf{I})\mathbf{F} &= \mathbf{F}(\mathbf{F}^T\mathbf{U}^{-2}\mathbf{F})\end{aligned}\tag{27}$$

RESULT

Based on the previous slide, we have

$$(\mathbf{R}\mathbf{U}^{-2} - \mathbf{I})\mathbf{F} = \mathbf{F}(\mathbf{F}^T\mathbf{U}^{-2}\mathbf{F})$$

Let $\mathbf{\Lambda} = (\mathbf{F}^T\mathbf{U}^{-2}\mathbf{F})$ be diagonal. Then

$$(\mathbf{R}\mathbf{U}^{-2} - \mathbf{I})(\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_n) = (\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_n) \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \vdots & \vdots \\ \vdots & \cdots & \ddots & 0 \\ 0 & \cdots & 0 & \lambda_n \end{pmatrix}$$

$$(\mathbf{R}\mathbf{U}^{-2} - \mathbf{I})(\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_n) = (\mathbf{f}_1\lambda_1, \mathbf{f}_2\lambda_2, \dots, \mathbf{f}_n\lambda_n)$$

$$(\mathbf{R}\mathbf{U}^{-2} - \mathbf{I})(\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_n) = (\lambda_1\mathbf{f}_1, \lambda_2\mathbf{f}_2, \dots, \lambda_n\mathbf{f}_n)$$

is an implicit eigenproblem.

REFERENCES

- ▶ Schönemann, P. H. (1985). On the formal differentiation of traces and determinants. *Multivariate Behavioral Research*, *20*, 113–139.
- ▶ Schönemann, P. H. (2001). Better never than late: Peer review and the preservation of prejudice. *Ethical Human Sciences and Services*, *3*, 7–21.