


Matrix & Array Derivatives

Mark van der Wilk

Department of Computing
Imperial College London

@markvanderwilk
m.vdwilk@imperial.ac.uk

October 12, 2021

Motivation: Next level curve fitting

You have now solved Linear Regression. A key design choice was which **basis functions** to use, e.g.:

$$\phi(x)^T = [x^3 \ x^2 \ x \ 1] \quad (1)$$

Motivation: Next level curve fitting

You have now solved Linear Regression. A key design choice was which **basis functions** to use, e.g.:

$$\boldsymbol{\phi}(x)^{\top} = [x^3 \ x^2 \ x \ 1] \quad (1)$$

Instead, can we learn the basis functions?

Motivation: Next level curve fitting

You have now solved Linear Regression. A key design choice was which **basis functions** to use, e.g.:

$$\boldsymbol{\phi}(x)^\top = [x^3 \ x^2 \ x \ 1] \quad (1)$$

Instead, can we learn the basis functions?

A **neural network** parameterises functions:

$$f_\ell(\mathbf{x}) = \sigma(\mathbf{A}_\ell \mathbf{x} + \mathbf{b}_\ell) \quad (2)$$

$$f_{NN}(\mathbf{x}) = f_L(f_{L-1}(\dots f_1(\mathbf{x}) \dots)) \quad (3)$$

Motivation: Next level curve fitting

You have now solved Linear Regression. A key design choice was which **basis functions** to use, e.g.:

$$\boldsymbol{\phi}(x)^\top = [x^3 \ x^2 \ x \ 1] \quad (1)$$

Instead, can we learn the basis functions?

A **neural network** parameterises functions:

$$f_\ell(\mathbf{x}) = \sigma(\mathbf{A}_\ell \mathbf{x} + \mathbf{b}_\ell) \quad (2)$$

$$f_{NN}(\mathbf{x}) = f_L(f_{L-1}(\dots f_1(\mathbf{x}) \dots)) \quad (3)$$

- Parameters are $\boldsymbol{\theta} = \{\mathbf{A}_\ell, \mathbf{b}_\ell\}_{\ell=1}^L$.

Motivation: Next level curve fitting

You have now solved Linear Regression. A key design choice was which **basis functions** to use, e.g.:

$$\boldsymbol{\phi}(x)^\top = [x^3 \ x^2 \ x \ 1] \quad (1)$$

Instead, can we learn the basis functions?

A **neural network** parameterises functions:

$$f_\ell(\mathbf{x}) = \sigma(\mathbf{A}_\ell \mathbf{x} + \mathbf{b}_\ell) \quad (2)$$

$$f_{NN}(\mathbf{x}) = f_L(f_{L-1}(\dots f_1(\mathbf{x}) \dots)) \quad (3)$$

- ▶ Parameters are $\boldsymbol{\theta} = \{\mathbf{A}_\ell, \mathbf{b}_\ell\}_{\ell=1}^L$.
- ▶ How do we differentiate w.r.t. matrices?

Derivatives of matrices/arrays

How should we find derivatives like $\frac{d}{d\theta} \mathbf{x}^\top \mathbf{A}(\theta) \mathbf{x}$ or $\frac{d}{d\mathbf{A}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2$?

Derivatives of matrices/arrays

How should we find derivatives like $\frac{d}{d\theta} \mathbf{x}^\top \mathbf{A}(\theta) \mathbf{x}$ or $\frac{d}{d\mathbf{A}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2$?

- By the same arguments as before (can look at differences of any array, and directional derivative arguments), we can just collect the gradients of all outputs w.r.t. all inputs:

$$\frac{\partial f_{ijkl}}{\partial x_{abc}} \quad (4)$$

Derivatives of matrices/arrays

How should we find derivatives like $\frac{d}{d\theta} \mathbf{x}^\top \mathbf{A}(\theta) \mathbf{x}$ or $\frac{d}{d\mathbf{A}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2$?

- By the same arguments as before (can look at differences of any array, and directional derivative arguments), we can just collect the gradients of all outputs w.r.t. all inputs:

$$\frac{\partial f_{ijkl}}{\partial x_{abc}} \quad (4)$$

Wouldn't it be nice if there was a chain rule?

$$\frac{df}{d\theta} = \frac{df}{d\mathbf{A}} \frac{d\mathbf{A}}{d\theta} \quad \text{or} \quad \frac{df}{d\mathbf{A}} = \frac{df}{d\mathbf{g}} \frac{d\mathbf{g}}{d\mathbf{A}}?$$

Chain rule

A function of a matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ is **just a multivariate function**:

$$f(\mathbf{A}) = \|\mathbf{Ax} - \mathbf{y}\|^2$$
$$f(A_{11}, A_{21}, \dots, A_{M1} \dots A_{MN}) = \sum_i \left(\sum_{ij} A_{ij} x_j - y_i \right)^2$$

Chain rule

A function of a matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ is **just a multivariate function**:

$$f(\mathbf{A}) = \|\mathbf{Ax} - \mathbf{y}\|^2$$
$$f(A_{11}, A_{21}, \dots, A_{M1} \dots A_{MN}) = \sum_i \left(\sum_{ij} A_{ij} x_j - y_i \right)^2$$

We just **arrange** the numbers in a different way.

Chain rule

A function of a matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ is **just a multivariate function**:

$$f(\mathbf{A}) = \|\mathbf{Ax} - \mathbf{y}\|^2$$
$$f(A_{11}, A_{21}, \dots, A_{M1} \dots A_{MN}) = \sum_i \left(\sum_{ij} A_{ij} x_j - y_i \right)^2$$

We just **arrange** the numbers in a different way. So the chain rule is the same!

$$f(\mathbf{g}) = \|\mathbf{g}\|^2, \quad \mathbf{g}(\mathbf{A}) = \mathbf{Ax} - \mathbf{y} \quad (5)$$

$$\frac{\partial f}{\partial A_{ij}} = \sum_k \frac{\partial f}{\partial g_k} \frac{\partial g_k}{\partial A_{ij}} \quad (6)$$

Chain rule

A function of a matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ is **just a multivariate function**:

$$f(\mathbf{A}) = \|\mathbf{Ax} - \mathbf{y}\|^2$$
$$f(A_{11}, A_{21}, \dots, A_{M1} \dots A_{MN}) = \sum_i \left(\sum_{ij} A_{ij} x_j - y_i \right)^2$$

We just **arrange** the numbers in a different way. So the chain rule is the same!

$$f(\mathbf{g}) = \|\mathbf{g}\|^2, \quad \mathbf{g}(\mathbf{A}) = \mathbf{Ax} - \mathbf{y} \quad (5)$$

$$\frac{\partial f}{\partial A_{ij}} = \sum_k \frac{\partial f}{\partial g_k} \frac{\partial g_k}{\partial A_{ij}} \quad (6)$$

$$f(\mathbf{A}) = \mathbf{x}^\top \mathbf{Ax} \quad (7)$$

$$\frac{\partial f}{\partial \theta} = \sum_{jk} \frac{\partial f}{\partial A_{jk}} \frac{\partial A_{jk}}{\partial \theta} \quad (8)$$

Chain rule is not straightforward

$$f(\mathbf{A}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$$

$$\frac{\partial f}{\partial \theta} = \sum_{jk} \frac{\partial f}{\partial A_{jk}} \frac{\partial A_{jk}}{\partial \theta}$$

$$f(\mathbf{g}) = \|\mathbf{g}\|^2, \mathbf{g}(\mathbf{A}) = \mathbf{A} \mathbf{x} - \mathbf{y}$$

$$\frac{\partial f}{\partial A_{ij}} = \sum_k \frac{\partial f}{\partial g_k} \frac{\partial g_k}{\partial A_{ij}}$$

Can we find a convenient notation like earlier?

$$\frac{df}{d\theta} = \frac{df}{d\mathbf{A}} \frac{d\mathbf{A}}{d\theta}?$$

$$\frac{df}{d\mathbf{A}} = \frac{df}{d\mathbf{g}} \frac{d\mathbf{g}}{d\mathbf{A}}? \quad (9)$$

Chain rule is not straightforward

$$f(\mathbf{A}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$$

$$\frac{\partial f}{\partial \theta} = \sum_{jk} \frac{\partial f}{\partial A_{jk}} \frac{\partial A_{jk}}{\partial \theta}$$

$$f(\mathbf{g}) = \|\mathbf{g}\|^2, \mathbf{g}(\mathbf{A}) = \mathbf{A} \mathbf{x} - \mathbf{y}$$

$$\frac{\partial f}{\partial A_{ij}} = \sum_k \frac{\partial f}{\partial g_k} \frac{\partial g_k}{\partial A_{ij}}$$

Can we find a convenient notation like earlier?

$$\frac{df}{d\theta} = \frac{df}{d\mathbf{A}} \frac{d\mathbf{A}}{d\theta}?$$

$$\frac{df}{d\mathbf{A}} = \frac{df}{d\mathbf{g}} \frac{d\mathbf{g}}{d\mathbf{A}}? \quad (9)$$

- ▶ NOT matrix multiplication, even though both $\frac{df}{d\mathbf{A}}$ and $\frac{d\mathbf{A}}{d\theta}$ look like matrices.

Chain rule is not straightforward

$$f(\mathbf{A}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$$

$$\frac{\partial f}{\partial \theta} = \sum_{jk} \frac{\partial f}{\partial A_{jk}} \frac{\partial A_{jk}}{\partial \theta}$$

$$f(\mathbf{g}) = \|\mathbf{g}\|^2, \mathbf{g}(\mathbf{A}) = \mathbf{A} \mathbf{x} - \mathbf{y}$$

$$\frac{\partial f}{\partial A_{ij}} = \sum_k \frac{\partial f}{\partial g_k} \frac{\partial g_k}{\partial A_{ij}}$$

Can we find a convenient notation like earlier?

$$\frac{df}{d\theta} = \frac{df}{d\mathbf{A}} \frac{d\mathbf{A}}{d\theta}?$$

$$\frac{df}{d\mathbf{A}} = \frac{df}{d\mathbf{g}} \frac{d\mathbf{g}}{d\mathbf{A}}? \quad (9)$$

- ▶ NOT matrix multiplication, even though both $\frac{df}{d\mathbf{A}}$ and $\frac{d\mathbf{A}}{d\theta}$ look like matrices. Check the shapes!
- ▶ Shape of $\frac{d\mathbf{g}}{d\mathbf{A}}$ isn't even a matrix!

Derivatives with Respect to Matrices

- ▶ Recall: A function $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ has a gradient that is an $M \times N$ -matrix with

$$\frac{df}{dx} \in \mathbb{R}^{M \times N}, \quad df[m, n] = \frac{\partial f_m}{\partial x_n}$$

Gradient dimension: # target dimensions \times # input dimensions

Derivatives with Respect to Matrices

- ▶ Recall: A function $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ has a gradient that is an $M \times N$ -matrix with

$$\frac{df}{dx} \in \mathbb{R}^{M \times N}, \quad df[m, n] = \frac{\partial f_m}{\partial x_n}$$

Gradient dimension: # target dimensions \times # input dimensions

- ▶ This generalizes to when the inputs (N) or targets (M) are **matrices**

Derivatives with Respect to Matrices

- Recall: A function $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ has a gradient that is an $M \times N$ -matrix with

$$\frac{df}{dx} \in \mathbb{R}^{M \times N}, \quad df[m, n] = \frac{\partial f_m}{\partial x_n}$$

Gradient dimension: # target dimensions \times # input dimensions

- This generalizes to when the inputs (N) or targets (M) are **matrices**
- Function $f : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{P \times Q}$, has a gradient that is a $(P \times Q) \times (M \times N)$ object (array)

$$\frac{df}{dX} \in \mathbb{R}^{(P \times Q) \times (M \times N)}, \quad df[p, q, m, n] = \frac{\partial f_{pq}}{\partial X_{mn}}$$

Derivatives with Respect to Matrices

- Recall: A function $f : \mathbb{R}^N \rightarrow \mathbb{R}^M$ has a gradient that is an $M \times N$ -matrix with

$$\frac{df}{dx} \in \mathbb{R}^{M \times N}, \quad df[m, n] = \frac{\partial f_m}{\partial x_n}$$

Gradient dimension: # target dimensions \times # input dimensions

- This generalizes to when the inputs (N) or targets (M) are **matrices**
- Function $f : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{P \times Q}$, has a gradient that is a $(P \times Q) \times (M \times N)$ object (array)

$$\frac{df}{dX} \in \mathbb{R}^{(P \times Q) \times (M \times N)}, \quad df[p, q, m, n] = \frac{\partial f_{pq}}{\partial X_{mn}}$$

Autodiff packages have similar consistency of shapes.

Example 1: Derivatives with Respect to Matrices

$$\mathbf{f} = \mathbf{A}\mathbf{x}, \quad \mathbf{f} \in \mathbb{R}^M, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{x} \in \mathbb{R}^N$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

$$\frac{d\mathbf{f}}{d\mathbf{A}} \in \mathbb{R}^?$$

Example 1: Derivatives with Respect to Matrices

$$\mathbf{f} = \mathbf{A}\mathbf{x}, \quad \mathbf{f} \in \mathbb{R}^M, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{x} \in \mathbb{R}^N$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

$$\frac{d\mathbf{f}}{d\mathbf{A}} \in \mathbb{R}^{\# \text{ target dim} \times \# \text{ input dim}} = M \times (M \times N)$$

$$\frac{d\mathbf{f}}{d\mathbf{A}} = \begin{bmatrix} \frac{\partial f_1}{\partial \mathbf{A}} \\ \vdots \\ \frac{\partial f_M}{\partial \mathbf{A}} \end{bmatrix}, \quad \frac{\partial f_i}{\partial \mathbf{A}} \in \mathbb{R}^{1 \times (M \times N)}$$

Example 2: Derivatives with Respect to Matrices

$$f_i = \sum_{j=1}^N A_{ij}x_j, \quad i = 1, \dots, M$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

$$\frac{\partial f_i}{\partial A_{iq}} = ?$$

$$\frac{\partial f_i}{\partial A_{i,:}} = ?$$

$$\frac{\partial f_i}{\partial A_{k \neq i,:}} = ?$$

$$\frac{\partial f_i}{\partial \mathbf{A}} = ?$$

Example 2: Derivatives with Respect to Matrices

$$f_i = \sum_{j=1}^N A_{ij}x_j, \quad i = 1, \dots, M$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

$$\frac{\partial f_i}{\partial A_{iq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_i}{\partial A_{i,:}} = ? \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = ? \quad \frac{\partial f_i}{\partial \mathbf{A}} = ?$$

Example 2: Derivatives with Respect to Matrices

$$f_i = \sum_{j=1}^N A_{ij}x_j, \quad i = 1, \dots, M$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

$$\frac{\partial f_i}{\partial A_{iq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_i}{\partial A_{i,:}} = \underbrace{\mathbf{x}^\top}_{\in \mathbb{R}^{1 \times 1 \times N}} \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = ? \quad \frac{\partial f_i}{\partial \mathbf{A}} = ?$$

Example 2: Derivatives with Respect to Matrices

$$f_i = \sum_{j=1}^N A_{ij}x_j, \quad i = 1, \dots, M$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

$$\frac{\partial f_i}{\partial A_{iq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_i}{\partial A_{i,:}} = \underbrace{\mathbf{x}^\top}_{\in \mathbb{R}^{1 \times N}} \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = \underbrace{\mathbf{0}^\top}_{\in \mathbb{R}^{1 \times N}} \quad \frac{\partial f_i}{\partial \mathbf{A}} = ?$$

Example 2: Derivatives with Respect to Matrices

$$f_i = \sum_{j=1}^N A_{ij}x_j, \quad i = 1, \dots, M$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_M \end{bmatrix} = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_i(\mathbf{x}) \\ \vdots \\ f_M(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots + A_{1N}x_N \\ \vdots \\ A_{i1}x_1 + A_{i2}x_2 + \cdots + A_{iN}x_N \\ \vdots \\ A_{M1}x_1 + A_{M2}x_2 + \cdots + A_{MN}x_N \end{bmatrix}$$

$$\frac{\partial f_i}{\partial A_{iq}} = \underbrace{x_q}_{\in \mathbb{R}} \quad \frac{\partial f_i}{\partial A_{i,:}} = \underbrace{\mathbf{x}^\top}_{\in \mathbb{R}^{1 \times N}} \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = \underbrace{\mathbf{0}^\top}_{\in \mathbb{R}^{1 \times N}} \quad \frac{\partial f_i}{\partial \mathbf{A}} = \underbrace{\begin{bmatrix} \mathbf{0}^\top \\ \vdots \\ \mathbf{x}^\top \\ \vdots \\ \mathbf{0}^\top \end{bmatrix}}_{\in \mathbb{R}^{1 \times (M \times N)}}$$

Chain rule

- ▶ We now understand how gradients involving matrices are arranged in “multidimensional arrays” or “tensors”.

Chain rule

- ▶ We now understand how gradients involving matrices are arranged in “multidimensional arrays” or “tensors”.
- ▶ How do we perform the chain rule? Can we find a meaning for convenient notation like:

$$\frac{df}{d\theta} = \frac{df}{d\mathbf{A}} \frac{d\mathbf{A}}{d\theta} ? \qquad \frac{df}{d\theta} = \frac{df}{d\mathbf{g}} \frac{d\mathbf{g}}{d\mathbf{A}} ? \qquad (10)$$

Chain rule

- ▶ We now understand how gradients involving matrices are arranged in “multidimensional arrays” or “tensors”.
- ▶ How do we perform the chain rule? Can we find a meaning for convenient notation like:

$$\frac{df}{d\theta} = \frac{df}{d\mathbf{A}} \frac{d\mathbf{A}}{d\theta} ? \qquad \frac{df}{d\theta} = \frac{df}{d\mathbf{g}} \frac{d\mathbf{g}}{d\mathbf{A}} ? \qquad (10)$$

- ▶ Recall: Function $f : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{P \times Q}$, has a gradient that is a $(P \times Q) \times (M \times N)$ object (tensor)

$$\frac{df}{d\mathbf{X}} \in \mathbb{R}^{(P \times Q) \times (M \times N)}, \qquad df[p, q, m, n] = \frac{\partial f_{pq}}{\partial X_{mn}}$$

Chain rule (2)

- ▶ Start from index notation chain rule (**always correct!**):

$$\frac{\partial f_{pq}}{\partial X_{mn}} = \sum_{rs} \frac{\partial f_{pq}}{\partial A_{rs}} \frac{\partial A_{rs}}{\partial X_{mn}}$$

Chain rule (2)

- ▶ Start from index notation chain rule (**always correct!**):

$$\frac{\partial f_{pq}}{\partial X_{mn}} = \sum_{rs} \frac{\partial f_{pq}}{\partial A_{rs}} \frac{\partial A_{rs}}{\partial X_{mn}}$$

- ▶ Like matrix multiplication, but with **vectorised** (vectors stacked column-by-column) matrices!

$$\frac{\text{dvec}(f)}{\text{dvec}(X)} = \frac{\text{dvec}(f)}{\text{dvec}(A)} \frac{\text{dvec}(A)}{\text{dvec}(X)}$$

Chain rule (2)

- ▶ Start from index notation chain rule (**always correct!**):

$$\frac{\partial f_{pq}}{\partial X_{mn}} = \sum_{rs} \frac{\partial f_{pq}}{\partial A_{rs}} \frac{\partial A_{rs}}{\partial X_{mn}}$$

- ▶ Like matrix multiplication, but with **vectorised** (vectors stacked column-by-column) matrices!

$$\frac{\text{dvec}(f)}{\text{dvec}(X)} = \frac{\text{dvec}(f)}{\text{dvec}(A)} \frac{\text{dvec}(A)}{\text{dvec}(X)}$$

- ▶ Keep track of grouping, sum over grouped indices:

$$\underbrace{\frac{df}{dX}}_{(P \times Q) \times (M \times N)} = \underbrace{\frac{df}{dA}}_{(P \times Q) \times (R \times S)} \cdot \underbrace{\frac{dA}{dX}}_{(R \times S) \times (M \times N)}$$

Summary: Matrix differentiation

We saw:

- ▶ Principle is the same for matrix and vector differentiation

Summary: Matrix differentiation

We saw:

- ▶ Principle is the same for matrix and vector differentiation
- ▶ Difference: Management of the numbers. It's about **convention**

Summary: Matrix differentiation

We saw:

- ▶ Principle is the same for matrix and vector differentiation
- ▶ Difference: Management of the numbers. It's about **convention**
- ▶ Mathematical principle is index notation, convention is defined

Summary: Matrix differentiation

We saw:

- ▶ Principle is the same for matrix and vector differentiation
- ▶ Difference: Management of the numbers. It's about **convention**
- ▶ Mathematical principle is index notation, convention is defined

You should be able to:

- ▶ Do the bookkeeping of matrix derivative shapes
- ▶ Compute derivatives of matrices
- ▶ Abstract complex derivatives into the well-defined chain rule.
- ▶ Describe the detailed index-wise summation for the chain rule.