# CSE463: Neural Networks

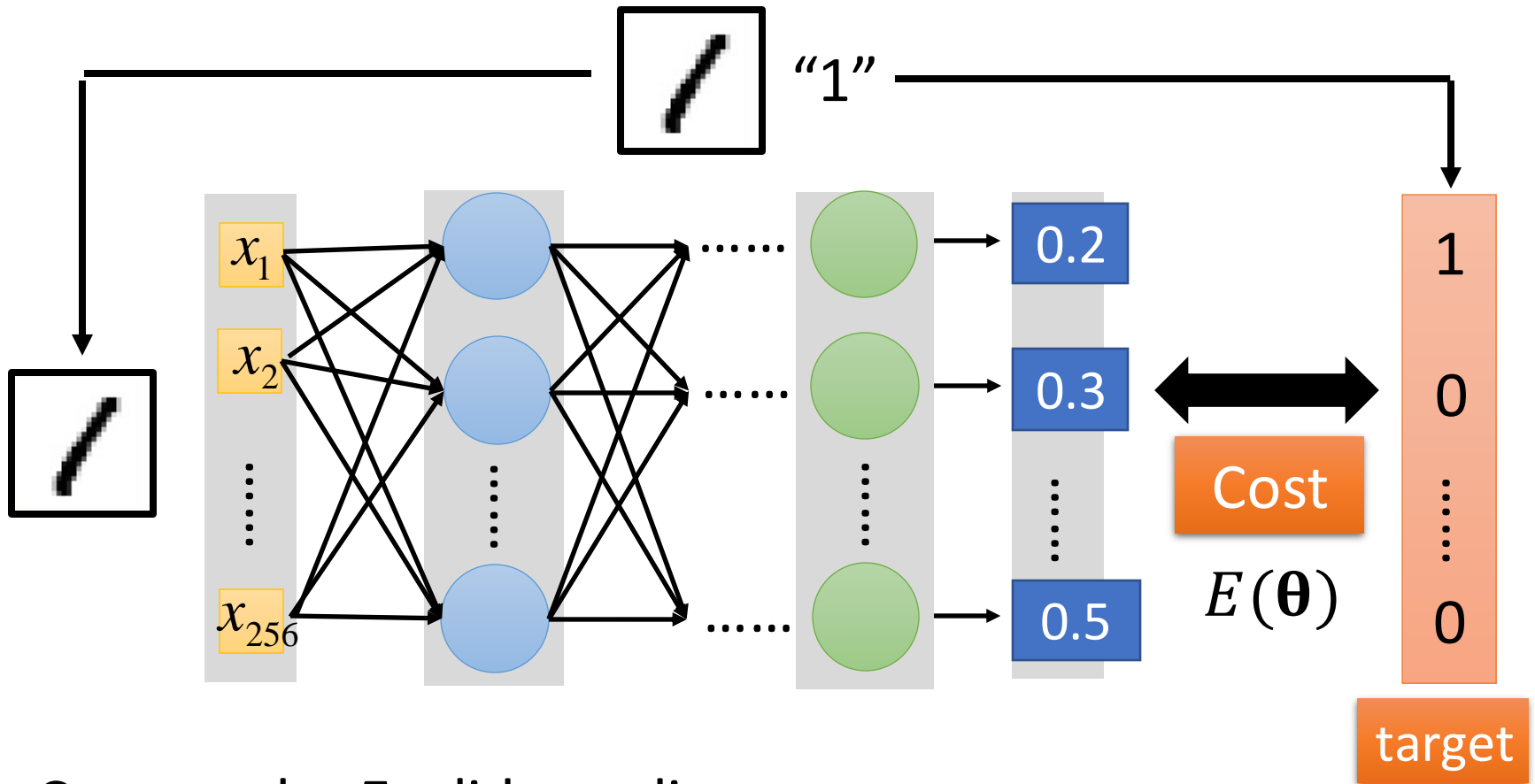# Backpropagation

by:

**Hossam Abd El Munim**

**Computer & Systems Engineering Dept.,**

**Ain Shams University,**

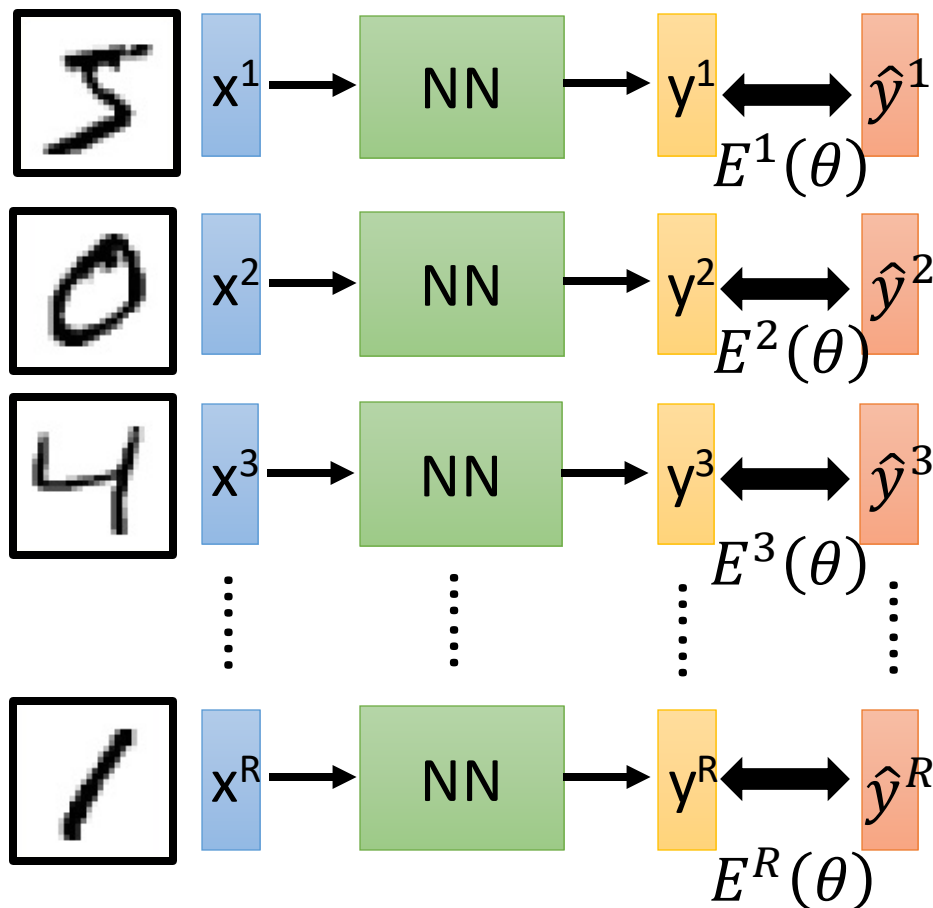**1 El-Sarayat Street, Abbassia, Cairo 11517**

# Cost

Given a set of network parameters $\theta$, each example has a cost value.



"1"

$$x_1$$
$$x_2$$
$$x_{256}$$

0.2
0.3
0.5

$E(\mathbf{\theta})$

Cost

1
0
0

target

Cost can be Euclidean distance or cross entropy of the network output and target

# Total Cost

For all training data …



Total Cost:

$$E(\boldsymbol{\theta}) = \sum_{r=1}^{R} E^r(\boldsymbol{\theta})$$

How bad the network parameters $\boldsymbol{\theta}$ is on this task

Find the network parameters $\boldsymbol{\theta}^*$ that minimize this value

# Problem Formulation

Given a data set, D = {($\mathbf{X}_1$, $\mathbf{Y}_1$), ($\mathbf{X}_2$, $\mathbf{Y}_2$)... ($\mathbf{X}_N$, $\mathbf{Y}_N$)} of labelled feature vectors where $\mathbf{Y}$ represents the target vectors designed as shown before. We to estimate the vector $\boldsymbol{\theta}$ through the conventional gradient descent:-

$$\boldsymbol{\theta}(t + 1) = \boldsymbol{\theta}(t) - \eta \frac{\partial E(\boldsymbol{\theta}(t))}{\partial \boldsymbol{\theta}}$$

# Let us simplify the problem:-

$$E(\boldsymbol{\theta}) = \frac{1}{2} ||\widehat{\mathbf{Y}}(\boldsymbol{\theta}) - \mathbf{Y}||^2$$

$$E(\boldsymbol{\theta}) = \frac{1}{2}(\widehat{\mathbf{Y}}(\boldsymbol{\theta}) - \mathbf{Y})^{\mathrm{T}}(\widehat{\mathbf{Y}}(\boldsymbol{\theta}) - \mathbf{Y})$$

$$\frac{\partial E}{\partial \boldsymbol{\theta}} = (\widehat{\mathbf{Y}}(\boldsymbol{\theta}) - \mathbf{Y})^{\mathrm{T}} \frac{\partial \widehat{\mathbf{Y}}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

# Vector Calculus

## 5.2 Partial Differentiation and Gradients

Differentiation as discussed in Section 5.1 applies to functions $f$ of a scalar variable $x \in \mathbb{R}$. In the following, we consider the general case where the function $f$ depends on one or more variables $x \in \mathbb{R}^n$, e.g., $f(x) = f(x_1, x_2)$. The generalization of the derivative to functions of several variables is the *gradient*.

We find the gradient of the function $f$ with respect to $x$ by *varying one variable at a time* and keeping the others constant. The gradient is then the collection of these *partial derivatives*.

**Definition 5.5** (Partial Derivative). For a function $f : \mathbb{R}^n \rightarrow \mathbb{R}, x \mapsto f(x), x \in \mathbb{R}^n$ of $n$ variables $x_1, \ldots, x_n$ we define the *partial derivatives* as

$$\frac{\partial f}{\partial x_1} = \lim_{h \to 0} \frac{f(x_1 + h, x_2, \ldots, x_n) - f(x)}{h}$$

$$\vdots \tag{5.39}$$

$$\frac{\partial f}{\partial x_n} = \lim_{h \to 0} \frac{f(x_1, \ldots, x_{n-1}, x_n + h) - f(x)}{h}$$

and collect them in the row vector

$$\nabla_x f = \operatorname{grad} f = \frac{\mathrm{d}f}{\mathrm{d}x} = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} & \frac{\partial f(x)}{\partial x_2} & \cdots & \frac{\partial f(x)}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{1 \times n}, \tag{5.40}$$

### 5.2.1 Basic Rules of Partial Differentiation

In the multivariate case, where $x \in \mathbb{R}^n$, the basic differentiation rules that we know from school (e.g., sum rule, product rule, chain rule; see also Section 5.1.2) still apply. However, when we compute derivatives with respect to vectors $x \in \mathbb{R}^n$ we need to pay attention: Our gradients now involve vectors and matrices, and matrix multiplication is not commutative (Section 2.2.1), i.e., the order matters.

Here are the general product rule, sum rule, and chain rule:

$$\text{Product rule:} \quad \frac{\partial}{\partial x}\left(f(x)g(x)\right) = \frac{\partial f}{\partial x}g(x) + f(x)\frac{\partial g}{\partial x} \qquad (5.46)$$

$$\text{Sum rule:} \quad \frac{\partial}{\partial x}\left(f(x) + g(x)\right) = \frac{\partial f}{\partial x} + \frac{\partial g}{\partial x} \qquad (5.47)$$

$$\text{Chain rule:} \quad \frac{\partial}{\partial x}(g \circ f)(x) = \frac{\partial}{\partial x}\left(g(f(x))\right) = \frac{\partial g}{\partial f}\frac{\partial f}{\partial x} \qquad (5.48)$$

## 5.2.2 Chain Rule

Consider a function $f : \mathbb{R}^2 \to \mathbb{R}$ of two variables $x_1, x_2$. Furthermore, $x_1(t)$ and $x_2(t)$ are themselves functions of $t$. To compute the gradient of $f$ with respect to $t$, we need to apply the chain rule (5.48) for multivariate functions as

$$\frac{\mathrm{d}f}{\mathrm{d}t} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial x_1(t)}{\partial t} \\ \frac{\partial x_2(t)}{\partial t} \end{bmatrix} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}, \qquad (5.49)$$

where d denotes the gradient and $\partial$ partial derivatives.

If $f(x_1, x_2)$ is a function of $x_1$ and $x_2$, where $x_1(s, t)$ and $x_2(s, t)$ are themselves functions of two variables $s$ and $t$, the chain rule yields the partial derivatives

$$\frac{\partial f}{\partial s} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial s} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial s}, \qquad (5.51)$$

$$\frac{\partial f}{\partial t} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}, \qquad (5.52)$$

If $f(x_1, x_2)$ is a function of $x_1$ and $x_2$, where $x_1(s,t)$ and $x_2(s,t)$ are themselves functions of two variables $s$ and $t$, the chain rule yields the partial derivatives

$$\frac{\partial f}{\partial s} = \frac{\partial f}{\partial x_1}\frac{\partial x_1}{\partial s} + \frac{\partial f}{\partial x_2}\frac{\partial x_2}{\partial s} \,, \tag{5.51}$$

$$\frac{\partial f}{\partial t} = \frac{\partial f}{\partial x_1}\frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2}\frac{\partial x_2}{\partial t} \,, \tag{5.52}$$

and the gradient is obtained by the matrix multiplication

$$\frac{\mathrm{d}f}{\mathrm{d}(s,t)} = \frac{\partial f}{\partial \boldsymbol{x}}\frac{\partial \boldsymbol{x}}{\partial(s,t)} = \underbrace{\begin{bmatrix} \dfrac{\partial f}{\partial x_1} & \dfrac{\partial f}{\partial x_2} \end{bmatrix}}_{=\frac{\partial f}{\partial \boldsymbol{x}}} \underbrace{\begin{bmatrix} \dfrac{\partial x_1}{\partial s} & \dfrac{\partial x_1}{\partial t} \\[2mm] \dfrac{\partial x_2}{\partial s} & \dfrac{\partial x_2}{\partial t} \end{bmatrix}}_{=\frac{\partial \boldsymbol{x}}{\partial(s,t)}} . \tag{5.53}$$

## 5.3 Gradients of Vector-Valued Functions

Thus far, we discussed partial derivatives and gradients of functions $f : \mathbb{R}^n \to \mathbb{R}$ mapping to the real numbers. In the following, we will generalize the concept of the gradient to vector-valued functions (vector fields) $f : \mathbb{R}^n \to \mathbb{R}^m$, where $n \geqslant 1$ and $m > 1$.

For a function $f : \mathbb{R}^n \to \mathbb{R}^m$ and a vector $x = [x_1, \dots, x_n]^\top \in \mathbb{R}^n$, the corresponding vector of function values is given as

$$f(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_m(x) \end{bmatrix} \in \mathbb{R}^m . \tag{5.54}$$

Writing the vector-valued function in this way allows us to view a vector-valued function $f : \mathbb{R}^n \to \mathbb{R}^m$ as a vector of functions $[f_1, \dots, f_m]^\top$, $f_i : \mathbb{R}^n \to \mathbb{R}$ that map onto $\mathbb{R}$. The differentiation rules for every $f_i$ are exactly the ones we discussed in Section 5.2.

Therefore, the partial derivative of a vector-valued function $f : \mathbb{R}^n \to \mathbb{R}^m$ with respect to $x_i \in \mathbb{R}$, $i = 1, \ldots n$, is given as the vector

$$\frac{\partial f}{\partial x_i} = \begin{bmatrix} \frac{\partial f_1}{\partial x_i} \\ \vdots \\ \frac{\partial f_m}{\partial x_i} \end{bmatrix} = \begin{bmatrix} \lim_{h \to 0} \frac{f_1(x_1, \ldots, x_{i-1}, x_i + h, x_{i+1}, \ldots x_n) - f_1(x)}{h} \\ \vdots \\ \lim_{h \to 0} \frac{f_m(x_1, \ldots, x_{i-1}, x_i + h, x_{i+1}, \ldots x_n) - f_m(x)}{h} \end{bmatrix} \in \mathbb{R}^m .$$

(5.55)

From (5.40), we know that the gradient of $f$ with respect to a vector is the row vector of the partial derivatives. In (5.55), every partial derivative $\partial f / \partial x_i$ is a column vector. Therefore, we obtain the gradient of $f : \mathbb{R}^n \to \mathbb{R}^m$ with respect to $x \in \mathbb{R}^n$ by collecting these partial derivatives:

$$\frac{\mathrm{d} f(x)}{\mathrm{d} x} = \begin{bmatrix} \boxed{\dfrac{\partial f(x)}{\partial x_1}} & \cdots & \boxed{\dfrac{\partial f(x)}{\partial x_n}} \end{bmatrix}$$

(5.56a)

$$= \begin{bmatrix} \boxed{\dfrac{\partial f_1(x)}{\partial x_1}} & \cdots & \boxed{\dfrac{\partial f_1(x)}{\partial x_n}} \\ \vdots & & \vdots \\ \boxed{\dfrac{\partial f_m(x)}{\partial x_1}} & \cdots & \boxed{\dfrac{\partial f_m(x)}{\partial x_n}} \end{bmatrix} \in \mathbb{R}^{m \times n} .$$

(5.56b)

**Definition 5.6** (Jacobian). The collection of all first-order partial derivatives of a vector-valued function $f : \mathbb{R}^n \to \mathbb{R}^m$ is called the *Jacobian*. The Jacobian $J$ is an $m \times n$ matrix, which we define and arrange as follows:

$$J = \nabla_x f = \frac{\mathrm{d}f(x)}{\mathrm{d}x} = \begin{bmatrix} \dfrac{\partial f(x)}{\partial x_1} & \cdots & \dfrac{\partial f(x)}{\partial x_n} \end{bmatrix} \tag{5.57}$$

$$= \begin{bmatrix} \dfrac{\partial f_1(x)}{\partial x_1} & \cdots & \dfrac{\partial f_1(x)}{\partial x_n} \\ \vdots & & \vdots \\ \dfrac{\partial f_m(x)}{\partial x_1} & \cdots & \dfrac{\partial f_m(x)}{\partial x_n} \end{bmatrix}, \tag{5.58}$$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad J(i,j) = \frac{\partial f_i}{\partial x_j}. \tag{5.59}$$

As a special case of (5.58), a function $f : \mathbb{R}^n \to \mathbb{R}^1$, which maps a vector $x \in \mathbb{R}^n$ onto a scalar (e.g., $f(x) = \sum_{i=1}^{n} x_i$), possesses a Jacobian that is a row vector (matrix of dimension $1 \times n$); see (5.40).

**Example 5.9 (Gradient of a Vector-Valued Function)**
We are given

$$f(x) = Ax, \qquad f(x) \in \mathbb{R}^M, \quad A \in \mathbb{R}^{M \times N}, \quad x \in \mathbb{R}^N.$$

To compute the gradient $df/dx$ we first determine the dimension of $df/dx$: Since $f : \mathbb{R}^N \to \mathbb{R}^M$, it follows that $df/dx \in \mathbb{R}^{M \times N}$. Second, to compute the gradient we determine the partial derivatives of $f$ with respect to every $x_j$:

$$f_i(x) = \sum_{j=1}^{N} A_{ij} x_j \implies \frac{\partial f_i}{\partial x_j} = A_{ij} \qquad (5.67)$$

We collect the partial derivatives in the Jacobian and obtain the gradient

$$\frac{df}{dx} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \cdots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} = \begin{bmatrix} A_{11} & \cdots & A_{1N} \\ \vdots & & \vdots \\ A_{M1} & \cdots & A_{MN} \end{bmatrix} = A \in \mathbb{R}^{M \times N}. \quad (5.68)$$

**Example 5.11 (Gradient of a Least-Squares Loss in a Linear Model)**
Let us consider the linear model

$$y = \Phi\theta \,, \tag{5.75}$$

where $\theta \in \mathbb{R}^D$ is a parameter vector, $\Phi \in \mathbb{R}^{N \times D}$ are input features and $y \in \mathbb{R}^N$ are the corresponding observations. We define the functions

$$L(e) := \|e\|^2 \,, \tag{5.76}$$
$$e(\theta) := y - \Phi\theta \,. \tag{5.77}$$

We seek $\frac{\partial L}{\partial \theta}$, and we will use the chain rule for this purpose. $L$ is called a *least-squares loss* function.

Before we start our calculation, we determine the dimensionality of the gradient as

$$\frac{\partial L}{\partial \theta} \in \mathbb{R}^{1 \times D} \,. \tag{5.78}$$

The chain rule allows us to compute the gradient as

$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial e}\frac{\partial e}{\partial \theta}, \tag{5.79}$$

where the $d$th element is given by

$$\frac{\partial L}{\partial \theta}[1, d] = \sum_{n=1}^{N} \frac{\partial L}{\partial e}[n]\frac{\partial e}{\partial \theta}[n, d]. \tag{5.80}$$

We know that $\|e\|^2 = e^\top e$ (see Section 3.2) and determine

$$\frac{\partial L}{\partial e} = 2e^\top \in \mathbb{R}^{1 \times N}. \tag{5.81}$$

Furthermore, we obtain

$$\frac{\partial e}{\partial \theta} = -\mathbf{\Phi} \in \mathbb{R}^{N \times D}, \tag{5.82}$$

such that our desired derivative is

$$\frac{\partial L}{\partial \theta} = -2e^\top \mathbf{\Phi} \overset{(5.77)}{=} -\underbrace{2(y^\top - \theta^\top \mathbf{\Phi}^\top)}_{1 \times N}\underbrace{\mathbf{\Phi}}_{N \times D} \in \mathbb{R}^{1 \times D}. \tag{5.83}$$

*Remark.* We would have obtained the same result without using the chain rule by immediately looking at the function

$$L_2(\boldsymbol{\theta}) := \|\boldsymbol{y} - \boldsymbol{\Phi}\boldsymbol{\theta}\|^2 = (\boldsymbol{y} - \boldsymbol{\Phi}\boldsymbol{\theta})^\top (\boldsymbol{y} - \boldsymbol{\Phi}\boldsymbol{\theta}). \qquad (5.84)$$

This approach is still practical for simple functions like $L_2$ but becomes impractical for deep function compositions. $\diamondsuit$
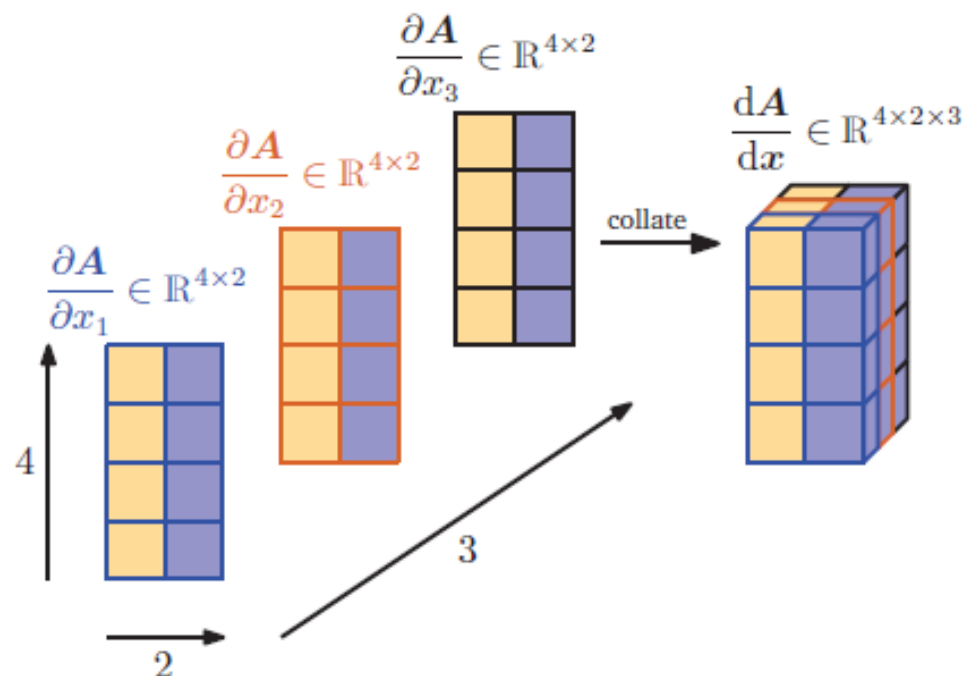
## 5.4 Gradients of Matrices

We will encounter situations where we need to take gradients of matrices with respect to vectors (or other matrices), which results in a multidimensional tensor. We can think of this tensor as a multidimensional array that
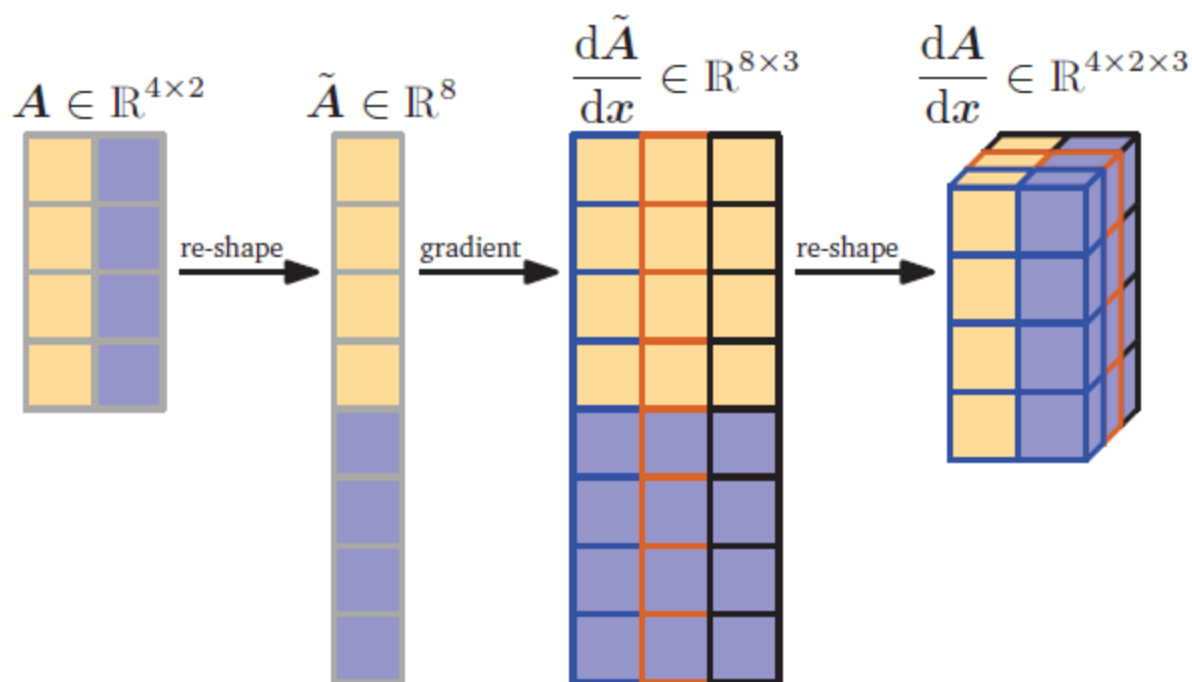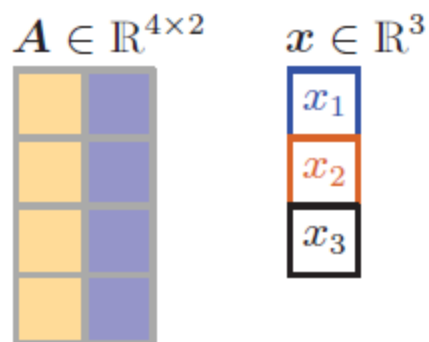
collects partial derivatives. For example, if we compute the gradient of an $m \times n$ matrix $\boldsymbol{A}$ with respect to a $p \times q$ matrix $\boldsymbol{B}$, the resulting Jacobian would be $(m \times n) \times (p \times q)$, i.e., a four-dimensional tensor $\boldsymbol{J}$, whose entries are given as $J_{ijkl} = \partial A_{ij} / \partial B_{kl}$.

$A \in \mathbb{R}^{4 \times 2}$  $x \in \mathbb{R}^3$

$x_1$
$x_2$
$x_3$

Partial derivatives:

$\dfrac{\partial A}{\partial x_3} \in \mathbb{R}^{4 \times 2}$

$\dfrac{\partial A}{\partial x_2} \in \mathbb{R}^{4 \times 2}$

$\dfrac{\mathrm{d} A}{\mathrm{d} x} \in \mathbb{R}^{4 \times 2 \times 3}$

$\dfrac{\partial A}{\partial x_1} \in \mathbb{R}^{4 \times 2}$

collate

4

2

3

(a) Approach 1: We compute the partial derivative $\frac{\partial A}{\partial x_1}, \frac{\partial A}{\partial x_2}, \frac{\partial A}{\partial x_3}$, each of which is a $4 \times 2$ matrix, and collate them in a $4 \times 2 \times 3$ tensor.

(b) Approach 2: We re-shape (flatten) $A \in \mathbb{R}^{4 \times 2}$ into a vector $\tilde{A} \in \mathbb{R}^8$. Then, we compute the gradient $\frac{\mathrm{d}\tilde{A}}{\mathrm{d}x} \in \mathbb{R}^{8 \times 3}$. We obtain the gradient tensor by re-shaping this gradient as illustrated above.

**Example 5.12 (Gradient of Vectors with Respect to Matrices)**
Let us consider the following example, where

$$f = Ax, \quad f \in \mathbb{R}^M, \quad A \in \mathbb{R}^{M \times N}, \quad x \in \mathbb{R}^N \qquad (5.85)$$

and where we seek the gradient $df/dA$. Let us start again by determining the dimension of the gradient as

$$\frac{df}{dA} \in \mathbb{R}^{M \times (M \times N)} . \qquad (5.86)$$

By definition, the gradient is the collection of the partial derivatives:

$$\frac{df}{dA} = \begin{bmatrix} \frac{\partial f_1}{\partial A} \\ \vdots \\ \frac{\partial f_M}{\partial A} \end{bmatrix}, \quad \frac{\partial f_i}{\partial A} \in \mathbb{R}^{1 \times (M \times N)} . \qquad (5.87)$$

To compute the partial derivatives, it will be helpful to explicitly write out the matrix vector multiplication:

$$f_i = \sum_{j=1}^{N} A_{ij} x_j, \quad i = 1, \ldots, M , \qquad (5.88)$$

and the partial derivatives are then given as

$$\frac{\partial f_i}{\partial A_{iq}} = x_q.$$ (5.89)

This allows us to compute the partial derivatives of $f_i$ with respect to a row of $A$, which is given as

$$\frac{\partial f_i}{\partial A_{i,:}} = x^\top \in \mathbb{R}^{1 \times 1 \times N},$$ (5.90)

$$\frac{\partial f_i}{\partial A_{k \neq i,:}} = 0^\top \in \mathbb{R}^{1 \times 1 \times N}$$ (5.91)

where we have to pay attention to the correct dimensionality. Since $f_i$ maps onto $\mathbb{R}$ and each row of $A$ is of size $1 \times N$, we obtain a $1 \times 1 \times N$-sized tensor as the partial derivative of $f_i$ with respect to a row of $A$.

We stack the partial derivatives (5.91) and get the desired gradient in (5.87) via

$$\frac{\partial f_i}{\partial A} = \begin{bmatrix} 0^\top \\ \vdots \\ 0^\top \\ x^\top \\ 0^\top \\ \vdots \\ 0^\top \end{bmatrix} \in \mathbb{R}^{1 \times (M \times N)}.$$ (5.92)

## 5.6 Backpropagation and Automatic Differentiation

In many machine learning applications, we find good model parameters by performing gradient descent (Section 7.1), which relies on the fact that we can compute the gradient of a learning objective with respect to the parameters of the model. For a given objective function, we can obtain the gradient with respect to the model parameters using calculus and applying the chain rule; see Section 5.2.2. We already had a taste in Section 5.3 when we looked at the gradient of a squared loss with respect to the parameters of a linear regression model.

Consider the function

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos\left(x^2 + \exp(x^2)\right). \tag{5.109}$$

By application of the chain rule, and noting that differentiation is linear, we compute the gradient

$$\frac{df}{dx} = \frac{2x + 2x \exp(x^2)}{2\sqrt{x^2 + \exp(x^2)}} - \sin\left(x^2 + \exp(x^2)\right)\left(2x + 2x \exp(x^2)\right)$$

$$= 2x \left(\frac{1}{2\sqrt{x^2 + \exp(x^2)}} - \sin\left(x^2 + \exp(x^2)\right)\right)\left(1 + \exp(x^2)\right). \tag{5.110}$$

Writing out the gradient in this explicit way is often impractical since it often results in a very lengthy expression for a derivative. In practice, it means that, if we are not careful, the implementation of the gradient could be significantly more expensive than computing the function, which imposes unnecessary overhead. For training deep neural network models, the *backpropagation* algorithm (Kelley, 1960; Bryson, 1961; Dreyfus, 1962; Rumelhart et al., 1986) is an efficient way to compute the gradient of an error function with respect to the parameters of the model.
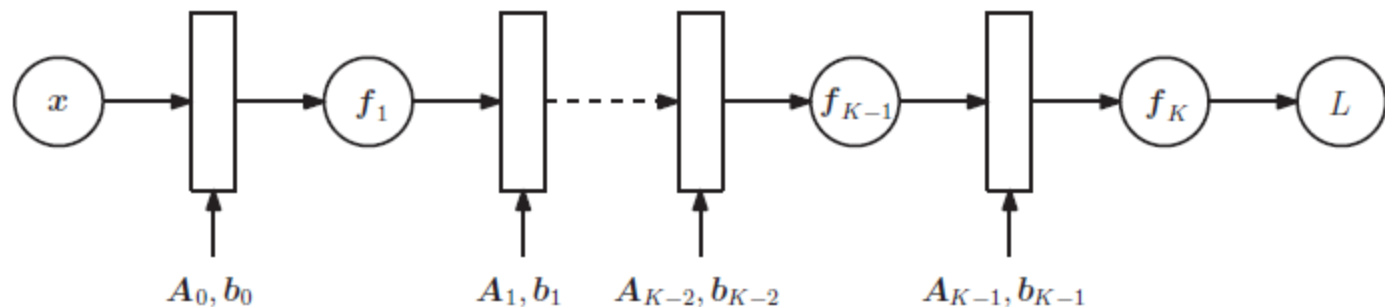
### 5.6.1 Gradients in a Deep Network

An area where the chain rule is used to an extreme is deep learning, where the function value $y$ is computed as a many-level function composition

$$y = (f_K \circ f_{K-1} \circ \cdots \circ f_1)(x) = f_K(f_{K-1}(\cdots(f_1(x))\cdots)), \quad (5.111)$$

where $x$ are the inputs (e.g., images), $y$ are the observations (e.g., class labels), and every function $f_i, i = 1, \ldots, K$, possesses its own parameters.

**Figure 5.8** Forward pass in a multi-layer neural network to compute the loss $L$ as a function of the inputs $x$ and the parameters $A_i, b_i$.

In neural networks with multiple layers, we have functions $f_i(x_{i-1}) = \sigma(A_{i-1}x_{i-1} + b_{i-1})$ in the $i$th layer. Here $x_{i-1}$ is the output of layer $i - 1$ and $\sigma$ an activation function, such as the logistic sigmoid $\frac{1}{1+e^{-x}}$, tanh or a rectified linear unit (ReLU). In order to train these models, we require the gradient of a loss function $L$ with respect to all model parameters $A_j, b_j$ for $j = 1, \ldots, K$. This also requires us to compute the gradient of $L$ with respect to the inputs of each layer. For example, if we have inputs $x$ and observations $y$ and a network structure defined by

$$f_0 := x \tag{5.112}$$

$$f_i := \sigma_i(A_{i-1}f_{i-1} + b_{i-1}), \quad i = 1, \ldots, K, \tag{5.113}$$

see also Figure 5.8 for a visualization, we may be interested in finding $A_j, b_j$ for $j = 0, \ldots, K - 1$, such that the squared loss

$$L(\theta) = \|y - f_K(\theta, x)\|^2 \tag{5.114}$$

is minimized, where $\theta = \{A_0, b_0, \ldots, A_{K-1}, b_{K-1}\}$.

To obtain the gradients with respect to the parameter set $\boldsymbol{\theta}$, we require the partial derivatives of $L$ with respect to the parameters $\boldsymbol{\theta}_j = \{\boldsymbol{A}_j, \boldsymbol{b}_j\}$ of each layer $j = 0, \ldots, K - 1$. The chain rule allows us to determine the partial derivatives as

$$\frac{\partial L}{\partial \boldsymbol{\theta}_{K-1}} = \frac{\partial L}{\partial \boldsymbol{f}_K} \frac{\partial \boldsymbol{f}_K}{\partial \boldsymbol{\theta}_{K-1}} \tag{5.115}$$

$$\frac{\partial L}{\partial \boldsymbol{\theta}_{K-2}} = \frac{\partial L}{\partial \boldsymbol{f}_K} \boxed{\frac{\partial \boldsymbol{f}_K}{\partial \boldsymbol{f}_{K-1}} \frac{\partial \boldsymbol{f}_{K-1}}{\partial \boldsymbol{\theta}_{K-2}}} \tag{5.116}$$

$$\frac{\partial L}{\partial \boldsymbol{\theta}_{K-3}} = \frac{\partial L}{\partial \boldsymbol{f}_K} \frac{\partial \boldsymbol{f}_K}{\partial \boldsymbol{f}_{K-1}} \boxed{\frac{\partial \boldsymbol{f}_{K-1}}{\partial \boldsymbol{f}_{K-2}} \frac{\partial \boldsymbol{f}_{K-2}}{\partial \boldsymbol{\theta}_{K-3}}} \tag{5.117}$$

$$\frac{\partial L}{\partial \boldsymbol{\theta}_i} = \frac{\partial L}{\partial \boldsymbol{f}_K} \frac{\partial \boldsymbol{f}_K}{\partial \boldsymbol{f}_{K-1}} \cdots \boxed{\frac{\partial \boldsymbol{f}_{i+2}}{\partial \boldsymbol{f}_{i+1}} \frac{\partial \boldsymbol{f}_{i+1}}{\partial \boldsymbol{\theta}_i}} \tag{5.118}$$
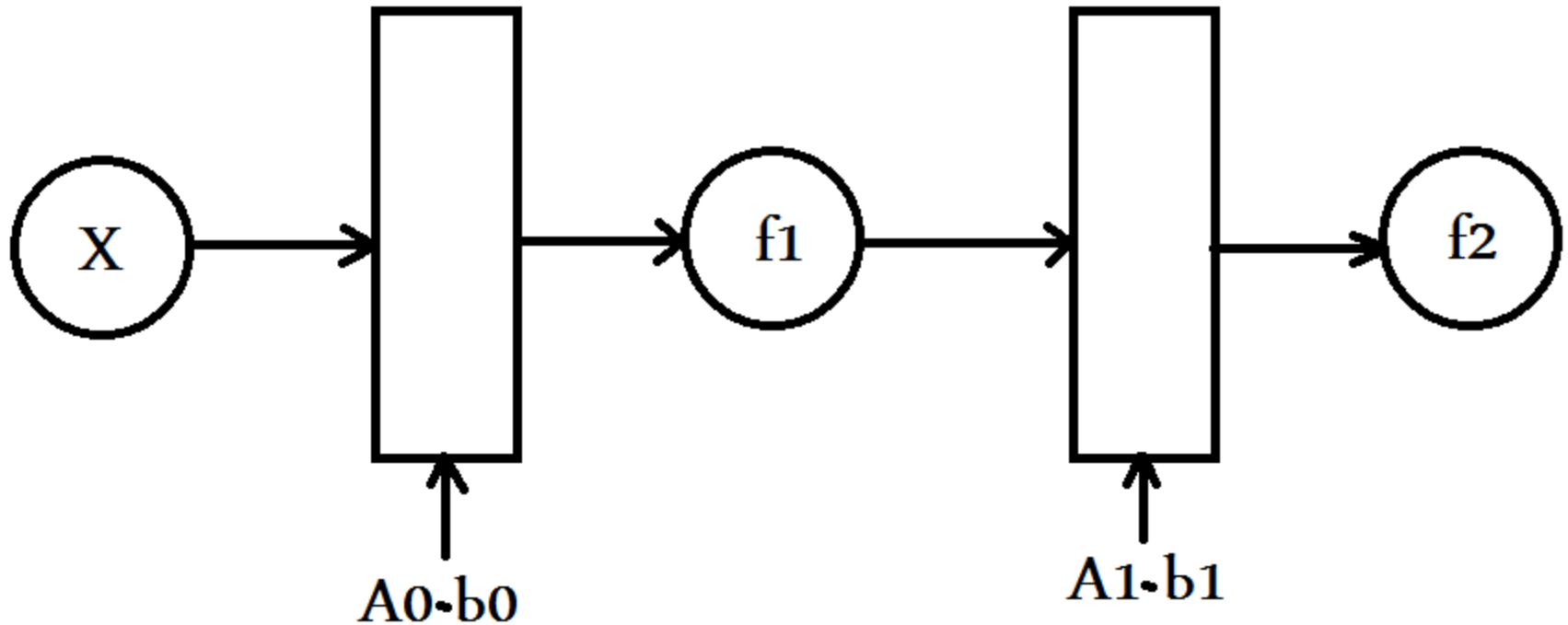
The **orange** terms are partial derivatives of the output of a layer with respect to its inputs, whereas the **blue** terms are partial derivatives of the output of a layer with respect to its parameters. Assuming, we have already computed the partial derivatives $\partial L / \partial \boldsymbol{\theta}_{i+1}$, then most of the computation can be reused to compute $\partial L / \partial \boldsymbol{\theta}_i$. The additional terms that we

**Figure 5.9**
Backward pass in a multi-layer neural network to compute the gradients of the loss function.

need to compute are indicated by the boxes. Figure 5.9 visualizes that the gradients are passed backward through the network.

# Specific Example

# Forward Path

$$\mathbf{f}_0 = \mathbf{X}$$

$$\mathbf{Z}_1 = \mathbf{A}_0\mathbf{f}_0 + \mathbf{b}_0$$

$$\mathbf{f}_1 = \sigma(\mathbf{Z1})$$

$$\mathbf{Z}_2 = \mathbf{A}_1\mathbf{f}_1 + \mathbf{b}_1$$

$$\mathbf{f}_2 = \sigma(\mathbf{Z}_2)$$

$$E = (1/2)\,(\mathbf{f}_2 - \mathbf{Y})^\top(\mathbf{f}_2 - \mathbf{Y})$$
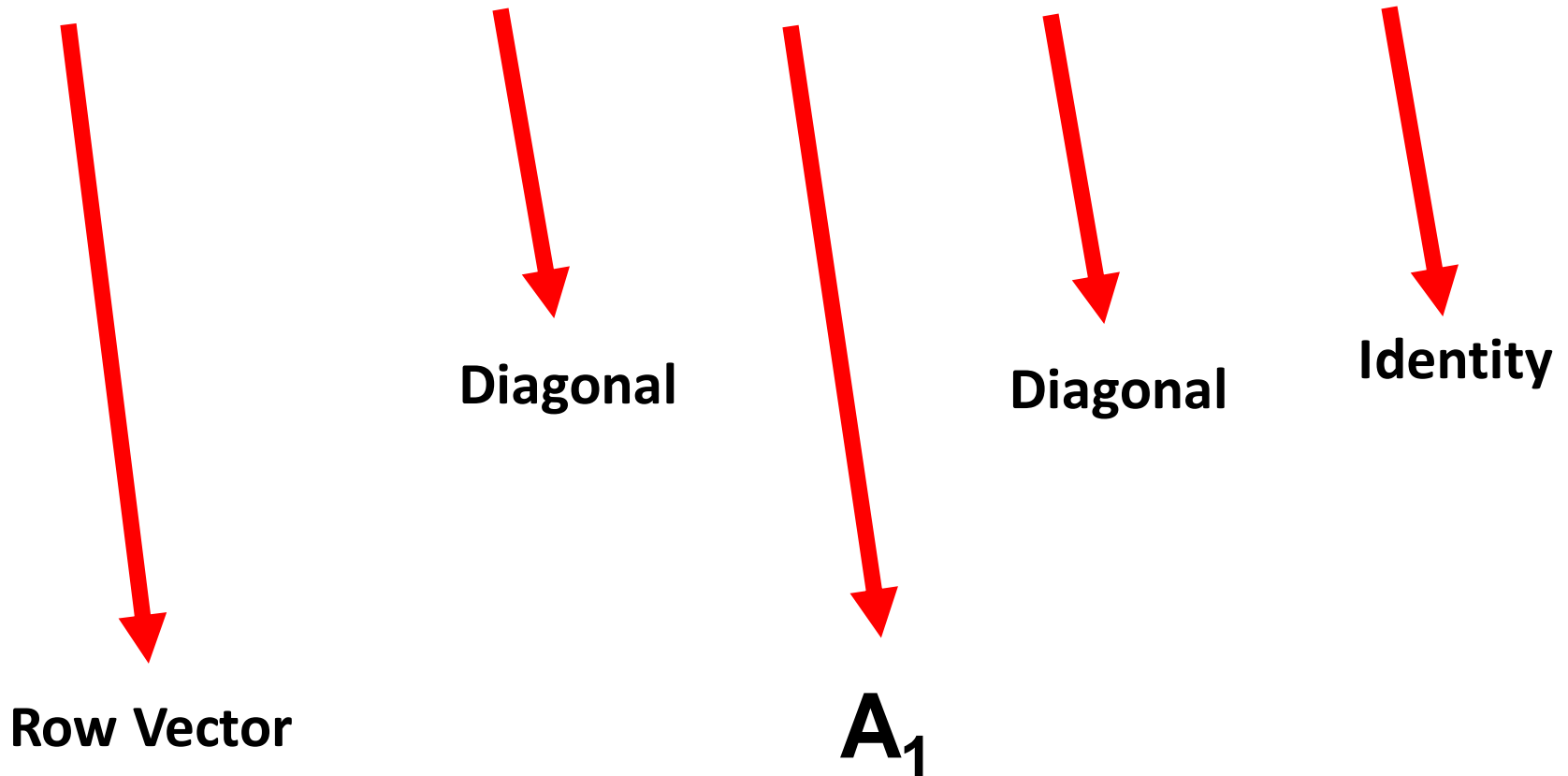
# Backward Path (1)

$$\partial E/\partial \mathbf{b_1} = (\mathbf{f_2} - \mathbf{Y})^\top (\partial \mathbf{f_2}/\partial \mathbf{Z_2})(\partial \mathbf{Z_2}/\partial \mathbf{b_1})$$

$(\partial \mathbf{Z2}/\partial \mathbf{b_1})$ = Identity Matrix

$(\partial \mathbf{f_2}/\partial \mathbf{Z2})$ = Diagonal matrix of σ'

# Backward Path (2)

$$\partial E / \partial \mathbf{b_0} = (\mathbf{f_2} - \mathbf{Y})^{\top} (\partial \mathbf{f_2} / \partial \mathbf{Z_2})(\partial \mathbf{Z_2} / \partial \mathbf{f_1})\ (\partial \mathbf{f_1} / \partial \mathbf{Z_1})(\partial \mathbf{Z_1} / \partial \mathbf{b_0})$$

**Diagonal**

**Diagonal**

**Identity**

**Row Vector**

$$\mathbf{A_1}$$

# Backward Path (3): Simplifying the derivative w.r.t a matrix

$$\partial E/\partial \alpha = (\mathbf{f}_2 - \mathbf{Y})^\top (\partial \mathbf{f}_2/\partial \mathbf{Z}_2)(\partial \mathbf{Z}_2/\partial \alpha)$$

$$\partial \mathbf{Z}_2/\partial \alpha = [\partial(\mathbf{A}_1)/\partial \alpha]\mathbf{f}_1$$

$$\alpha \in \mathbf{A}_1$$

# Backward Path (4): Simplifying the derivative w.r.t a matrix

$$\partial E / \partial \beta = (\mathbf{f}_2 - \mathbf{Y})^\top (\partial \mathbf{f}_2 / \partial \mathbf{Z}_2)(\partial \mathbf{Z}_2 / \partial \mathbf{f}_1)\ (\partial \mathbf{f}_1 / \partial \mathbf{Z}_1)(\partial \mathbf{Z}_1 / \partial \beta)$$

$$\partial \mathbf{Z}_1 / \partial \beta = [\partial (\mathbf{A}_0) / \partial \beta]\mathbf{f}_0$$

$$\beta \in \mathbf{A}_0$$