

## 1 Matrix Calculus

### 1.1 Derivative of Vector wrt Matrix

(1.1) We have vector  $\mathbf{r} = \begin{bmatrix} r_1 \\ r_2 \end{bmatrix}$ , matrix  $\mathbf{A} = \begin{bmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \end{bmatrix}$

(1.2) We have that  $\frac{\partial \mathbf{r}}{\partial \mathbf{A}} = \begin{bmatrix} \frac{\partial r_1}{\partial \mathbf{A}} \\ \frac{\partial r_2}{\partial \mathbf{A}} \end{bmatrix}$ , where each  $\frac{\partial r_i}{\partial \mathbf{A}}$  is a Scalar-Matrix derivative

(2.1) If we view  $\frac{\partial \mathbf{r}}{\partial \mathbf{A}} = \begin{bmatrix} \frac{\partial r_1}{\partial \mathbf{A}} \\ \frac{\partial r_2}{\partial \mathbf{A}} \end{bmatrix} = \begin{bmatrix} \frac{\partial r_1}{\partial w_1} & \frac{\partial r_1}{\partial w_2} & \frac{\partial r_1}{\partial w_3} & \frac{\partial r_1}{\partial w_4} & \frac{\partial r_1}{\partial w_5} & \frac{\partial r_1}{\partial w_6} \\ \frac{\partial r_2}{\partial w_1} & \frac{\partial r_2}{\partial w_2} & \frac{\partial r_2}{\partial w_3} & \frac{\partial r_2}{\partial w_4} & \frac{\partial r_2}{\partial w_5} & \frac{\partial r_2}{\partial w_6} \end{bmatrix}$ , this is termed as 'flattening' the matrix

(2.2) Notice that if  $\mathbf{r} \in \mathbb{R}^{m \times 1}$  and  $\mathbf{A} \in \mathbb{R}^{n \times k}$ , we have that  $\frac{\partial \mathbf{r}}{\partial \mathbf{A}} \in \mathbb{R}^{m \times (n \times k)}$

(3a) Alternatively, we could view  $\frac{\partial \mathbf{r}}{\partial \mathbf{A}} = \begin{bmatrix} \frac{\partial r_1}{\partial \mathbf{A}} \\ \frac{\partial r_2}{\partial \mathbf{A}} \end{bmatrix}$  as a rank-3 tensor

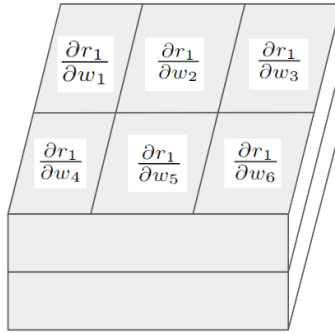


Figure 1: Rank-3 tensor

### 1.2 Derivative of Matrix wrt Matrix

(1a) We have matrix  $\mathbf{A} = \begin{bmatrix} a_1 & a_2 \\ a_3 & a_4 \end{bmatrix}$ , matrix  $\mathbf{W} = \begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \end{bmatrix}$

(1b) We have that  $\frac{\partial \mathbf{A}}{\partial \mathbf{W}} = \begin{bmatrix} \frac{\partial a_1}{\partial \mathbf{W}} & \frac{\partial a_2}{\partial \mathbf{W}} \\ \frac{\partial a_3}{\partial \mathbf{W}} & \frac{\partial a_4}{\partial \mathbf{W}} \end{bmatrix}$ , where each  $\frac{\partial a_i}{\partial \mathbf{A}}$  is a Scalar-Matrix derivative

(1c) Notice that if  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{W} \in \mathbb{R}^{k \times l}$ ; then  $\frac{\partial a_i}{\partial \mathbf{W}} \in \mathbb{R}^{k \times l}$ ,  $\frac{\partial \mathbf{A}}{\partial \mathbf{W}} \in \mathbb{R}^{(k \times l \times m) \times (n)}$

### 1.3 Derivative of Vector wrt Matrix - Alternative methods

Suppose that  $\mathbf{W} \in \mathbb{R}^{n \times m}$  and  $\mathbf{x} \in \mathbb{R}^{m \times 1}$ . How do we calculate  $\frac{d\mathbf{W}\mathbf{x}}{d\mathbf{W}}$ ?

We know that the quantity in question is a  $3^{rd}$  order tensor.

#### 1.3.1 Index Notation

(1a) We know that  $\mathbf{f} = \mathbf{W}\mathbf{x}$

(1b) Define  $f_i = W_{ij}x_j$ ; Note that we are using Einstein summation notation

$$(1c) \quad \frac{\partial f_i}{\partial W_{mn}} = \frac{\partial f_i}{\partial W_{ij}} \frac{\partial W_{ij}}{\partial W_{mn}} = \frac{\partial f_i}{\partial W_{ij}} x_j = \delta_{im} \delta_{jn} x_j = \delta_{im} x_n$$

#### 1.3.2 Vectorization

$$\begin{aligned} (1) \quad \mathbf{f} &= \mathbf{W}\mathbf{x} \\ &= \mathbf{I}\mathbf{W}\mathbf{x} \\ &= (\mathbf{x}^T \otimes \mathbf{I}) \text{vec}(\mathbf{W}) \\ &= (\mathbf{x}^T \otimes \mathbf{I}) \mathbf{w} \end{aligned}$$

$$(2) \quad \text{Thus, } \frac{\partial \mathbf{f}}{\partial \mathbf{w}} = (\mathbf{x}^T \otimes \mathbf{I})$$

#### 1.3.3 Special Case

See matrixDifferentiation.pdf for more

## 2 2-layer NN

### 2.1 Backpropagation - Bias terms

Note that we are considering backpropagation w.r.t a single **layer**, where a single layer may encapsulate **more than one neuron**

(1a) Input data  $\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \end{bmatrix} \in \mathbb{R}^{N \times D}$  ;  $N = \text{numSamples}$ ,  $D = \text{dataDimension}$

(1b) Weights  $\mathbf{W} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix} \in \mathbb{R}^{D \times C}$  ;  $C = \text{number of classes}$

(1c) Bias vector  $\mathbf{b} = [b_1 \quad b_2 \quad b_3] \in \mathbb{R}^C$

(1d) Bias matrix  $\mathbf{B} = \begin{bmatrix} b_1 & b_2 & b_3 \\ b_1 & b_2 & b_3 \end{bmatrix} \in \mathbb{R}^{N \times C}$ ; Bias vector is 'broadcast' once for each data sample

(1e) Upstream term  $\mathbf{Z} = \mathbf{XW} + \mathbf{B}$

$$= \begin{bmatrix} xw_{11} & xw_{12} & xw_{13} \\ xw_{21} & xw_{22} & xw_{23} \end{bmatrix} + \begin{bmatrix} b_1 & b_2 & b_3 \\ b_1 & b_2 & b_3 \end{bmatrix}$$

$$= \begin{bmatrix} z_{11} & z_{12} & z_{13} \\ z_{21} & z_{22} & z_{23} \end{bmatrix} \in \mathbb{R}^{N \times C}$$

(2a)  $\frac{\partial L_i}{\partial \mathbf{b}} = \frac{\partial L_i}{\partial \mathbf{Z}_i} \cdot \frac{\partial \mathbf{Z}_i}{\partial \mathbf{b}}$  ;  $\mathbf{Z}_i$  stands for  $i^{th}$  row of  $\mathbf{Z}$  (corresponding to  $i^{th}$  sample),  $L_i$  stands for loss of  $i^{th}$  sample

$$= \frac{\partial L_i}{\partial \mathbf{Z}_i} \cdot \frac{\partial}{\partial \mathbf{b}} \begin{bmatrix} z_{i1} & z_{i2} & z_{i3} \end{bmatrix}$$

$$= \frac{\partial L_i}{\partial \mathbf{Z}_i} \cdot \begin{bmatrix} \frac{\partial z_{i1}}{\partial b_1} & \frac{\partial z_{i1}}{\partial b_2} & \frac{\partial z_{i1}}{\partial b_3} \\ \frac{\partial z_{i2}}{\partial b_1} & \frac{\partial z_{i2}}{\partial b_2} & \frac{\partial z_{i2}}{\partial b_3} \\ \frac{\partial z_{i3}}{\partial b_1} & \frac{\partial z_{i3}}{\partial b_2} & \frac{\partial z_{i3}}{\partial b_3} \end{bmatrix} = \frac{\partial L_i}{\partial \mathbf{Z}_i} \cdot \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$= \frac{\partial L_i}{\partial \mathbf{Z}_i} \in \mathbb{R}^{1 \times C}$$

(2b) Thus,  $\frac{\partial L}{\partial \mathbf{b}} = \sum_i \frac{\partial L_i}{\partial \mathbf{b}} = \sum_i \frac{\partial L_i}{\partial z_i}$