

## Operators

Define the operator  $\mathbf{vec}_{l,m,n} : \mathbb{R}^{l \times m \times n} \rightarrow \mathbb{R}^{lmn}$ , which flattens a  $l \times m \times n$  tensor into a  $lmn$ -dimensional vector. Clearly,  $\mathbf{vec}_{l,m,n}$  is linear and invertible. We will refer to the inverse of  $\mathbf{vec}_{l,m,n}$  as  $\mathbf{ten}_{l,m,n}$ , which reshapes a  $lmn$ -dimensional vector into a  $l \times m \times n$  tensor.

## Jacobian

Let  $f : \mathbb{R}^{l \times m \times n} \rightarrow \mathbb{R}^{p \times q \times r}$  be a function of the form

$$f(\mathbf{x}) = \begin{array}{|c|} \hline f_{p,1,1}(\mathbf{x}) \cdots f_{p,1,r}(\mathbf{x}) \\ \hline f_{1,1,1}(\mathbf{x}) \cdots f_{1,1,r}(\mathbf{x}) \\ \vdots \quad \ddots \quad \vdots \\ f_{1,q,1}(\mathbf{x}) \cdots f_{1,q,r}(\mathbf{x}) \\ \hline \end{array}(\mathbf{x})$$

where the component functions  $f_{i,j,k} : \mathbb{R}^{l \times m \times n} \rightarrow \mathbb{R}$  are scalar-valued. It follows that the difference  $f(\mathbf{x} + \epsilon) - f(\mathbf{x})$  is approximated by

$$\begin{array}{|c|} \hline \left( \frac{\partial}{\partial \mathbf{x}} f_{p,1,1}(\mathbf{x}) \right)^T \mathbf{vec}_{l,m,n}(\epsilon) \cdots \left( \frac{\partial}{\partial \mathbf{x}} f_{p,1,r}(\mathbf{x}) \right)^T \mathbf{vec}_{l,m,n}(\epsilon) \\ \hline \left( \frac{\partial}{\partial \mathbf{x}} f_{1,1,1}(\mathbf{x}) \right)^T \mathbf{vec}_{l,m,n}(\epsilon) \cdots \left( \frac{\partial}{\partial \mathbf{x}} f_{1,1,r}(\mathbf{x}) \right)^T \mathbf{vec}_{l,m,n}(\epsilon) \\ \vdots \quad \ddots \quad \vdots \\ \left( \frac{\partial}{\partial \mathbf{x}} f_{1,q,1}(\mathbf{x}) \right)^T \mathbf{vec}_{l,m,n}(\epsilon) \cdots \left( \frac{\partial}{\partial \mathbf{x}} f_{1,q,r}(\mathbf{x}) \right)^T \mathbf{vec}_{l,m,n}(\epsilon) \\ \hline \end{array}(\epsilon) = \mathbf{ten}_{p,q,r} \left( \left( \frac{\partial}{\partial \mathbf{x}} f(\mathbf{x}) \right)^T \mathbf{vec}_{l,m,n}(\epsilon) \right)$$

where  $\frac{\partial}{\partial \mathbf{x}} f(\mathbf{x})$  is the Jacobian of  $f$  with respect to  $\mathbf{x}$  — that is, the  $lmn \times pqr$  matrix with columns given by the gradients  $\frac{\partial}{\partial \mathbf{x}} f_{i,j,k}(\mathbf{x})$ .

## Gradient descent

If  $f$  is scalar-valued, then the above approximation becomes

$$\left( \frac{\partial}{\partial \mathbf{x}} f(\mathbf{x}) \right)^T \mathbf{vec}_{l,m,n}(\epsilon) = \left\| \frac{\partial}{\partial \mathbf{x}} f(\mathbf{x}) \right\|_2 \|\mathbf{vec}_{l,m,n}(\epsilon)\|_2 \cos(\theta)$$

where  $\theta$  is the angle between the gradient and  $\mathbf{vec}_{l,m,n}(\epsilon)$ . Hence, the approximation is most negative when  $\mathbf{vec}_{l,m,n}(\epsilon) = -\gamma \frac{\partial}{\partial \mathbf{x}} f(\mathbf{x})$  for some positive constant  $\gamma$  (i.e., when  $\cos(\theta) = -1$ ). It follows that the direction from  $\mathbf{x}$  in which  $f$  decreases most rapidly is the direction of  $\epsilon = -\gamma \mathbf{ten}_{l,m,n} \left( \frac{\partial}{\partial \mathbf{x}} f(\mathbf{x}) \right)$ . Thus, the gradient descent update rule is given by

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \mathbf{ten}_{l,m,n} \left( \frac{\partial}{\partial \mathbf{x}} f(\mathbf{x}_t) \right)$$

## Chain rule

Consider the composition  $h = f \circ g$  for functions  $g : \mathbb{R}^{l \times m \times n} \rightarrow \mathbb{R}^{i \times j \times k}$  and  $f : \mathbb{R}^{i \times j \times k} \rightarrow \mathbb{R}^{p \times q \times r}$ . From the results above,

$$\begin{aligned}
h(\mathbf{x} + \boldsymbol{\epsilon}) - h(\mathbf{x}) &= f(g(\mathbf{x} + \boldsymbol{\epsilon})) - f(g(\mathbf{x})) \\
&= f(g(\mathbf{x}) + g(\mathbf{x} + \boldsymbol{\epsilon}) - g(\mathbf{x})) - f(g(\mathbf{x})) \\
&\approx f\left(g(\mathbf{x}) + \mathbf{ten}_{i,j,k}\left(\left(\frac{\partial}{\partial \mathbf{x}}g(\mathbf{x})\right)^T \mathbf{vec}_{l,m,n}(\boldsymbol{\epsilon})\right)\right) - f(g(\mathbf{x})) \\
&\approx \mathbf{ten}_{p,q,r}\left(\left(\frac{\partial}{\partial g(\mathbf{x})}f(g(\mathbf{x}))\right)^T \mathbf{vec}_{i,j,k}\left(\mathbf{ten}_{i,j,k}\left(\left(\frac{\partial}{\partial \mathbf{x}}g(\mathbf{x})\right)^T \mathbf{vec}_{l,m,n}(\boldsymbol{\epsilon})\right)\right)\right) \\
&= \mathbf{ten}_{p,q,r}\left(\left(\frac{\partial}{\partial g(\mathbf{x})}f(g(\mathbf{x}))\right)^T \left(\frac{\partial}{\partial \mathbf{x}}g(\mathbf{x})\right)^T \mathbf{vec}_{l,m,n}(\boldsymbol{\epsilon})\right) \\
&= \mathbf{ten}_{p,q,r}\left(\left(\frac{\partial}{\partial \mathbf{x}}g(\mathbf{x})\frac{\partial}{\partial g(\mathbf{x})}f(g(\mathbf{x}))\right)^T \mathbf{vec}_{l,m,n}(\boldsymbol{\epsilon})\right)
\end{aligned}$$

Hence, the Jacobian of  $h$  is given by

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

### Example

Consider the operation of cross-correlating a  $C \times H \times W$  image  $\mathbf{x}$  with a  $C \times H' \times W'$  kernel  $\mathbf{k}$ . The dimension of the cross-correlation is given by

$$\left\lfloor \frac{H - H'}{s_H} + 1 \right\rfloor \times \left\lfloor \frac{W - W'}{s_W} + 1 \right\rfloor$$

where  $s_H$  and  $s_W$  denote the strides in the vertical and horizontal spatial dimensions, respectively. The  $(i, j)^{\text{th}}$  element of the cross-correlation is given by

$$(\mathbf{x} * \mathbf{k})[i, j] = \sum_{c=0}^{C-1} \sum_{h=0}^{H'-1} \sum_{w=0}^{W'-1} \mathbf{x}[c, s_H i + h, s_W j + w] \mathbf{k}[c, h, w]$$

Defining the operation  $\mathbf{flat}_{d_1, \dots, d_n} : \mathbb{N}^n \rightarrow \mathbb{N}$  by

$$\mathbf{flat}_{d_1, \dots, d_n}(i_1, \dots, i_n) = \sum_{j=1}^n i_j \prod_{k=1}^{j-1} d_k$$

which takes an  $n$ -dimensional index of a  $d_1 \times \dots \times d_n$  tensor and returns the corresponding 1-dimensional index, it follows that rows  $\mathbf{flat}_{C,H,W}(c, s_H i + h, s_W j + w)$  of column  $\mathbf{flat}_{\left\lfloor \frac{H-H'}{s_H} + 1 \right\rfloor, \left\lfloor \frac{W-W'}{s_W} + 1 \right\rfloor}(i, j)$  in  $\frac{\partial}{\partial \mathbf{x}}(\mathbf{x} * \mathbf{k})$  are given

by  $\mathbf{k}[c, h, w]$  and that all other rows are zero. Hence, we have the following algorithm for computing  $\frac{\partial}{\partial \mathbf{x}}(\mathbf{x} * \mathbf{k})$ .

```

Initialize  $\frac{\partial}{\partial \mathbf{x}}(\mathbf{x} * \mathbf{k})$  with all zeros
for  $i = 0, \dots, \left\lfloor \frac{H - H'}{s_H} + 1 \right\rfloor - 1$  :
    for  $j = 0, \dots, \left\lfloor \frac{W - W'}{s_W} + 1 \right\rfloor - 1$  :
         $\text{col} \leftarrow \text{flat}_{\left\lfloor \frac{H-H'}{s_H} + 1 \right\rfloor, \left\lfloor \frac{W-W'}{s_W} + 1 \right\rfloor}(i, j)$ 
        for  $c = 0, \dots, C - 1$  :
            for  $h = 0, \dots, H' - 1$  :
                for  $w = 0, \dots, W' - 1$  :
                     $\text{row} \leftarrow \text{flat}_{C, H, W}(c, s_H i + h, s_W j + w)$ 
                     $\frac{\partial}{\partial \mathbf{x}}(\mathbf{x} * \mathbf{k})[\text{row}, \text{col}] \leftarrow \mathbf{k}[c, h, w]$ 

```

The above algorithm can be used to compute cross-correlation itself. Since cross-correlation is linear,

$$\mathbf{x} * \mathbf{k} = ((\mathbf{x} + \mathbf{x}) * \mathbf{k}) - (\mathbf{x} * \mathbf{k}) = \text{ten}_{\left\lfloor \frac{H-H'}{s_H} + 1 \right\rfloor, \left\lfloor \frac{W-W'}{s_W} + 1 \right\rfloor} \left( \left( \frac{\partial}{\partial \mathbf{x}}(\mathbf{x} * \mathbf{k}) \right)^T \text{vec}_{C, H, W}(\mathbf{x}) \right)$$

Note that in the above computations, we consider elements of a tensor in channel, row, column order. For example, the operator **flat** returns a 1-dimensional index by counting elements of a tensor front to back (channel), top to bottom (row), left to right (column). The order doesn't matter as long as one is consistent.