

Matrix Calculus

Jack David Carson - January 31, 2025

In machine learning, the fundamental operation of backpropagation is to take the derivative of a matrix with respect to a vector or another matrix. We may seek to calculate the gradient of a matrix $A_{m \times n}$ with respect to $B_{p \times q}$ and accumulate the partial derivatives in a Jacobian **tensor** J whose enteries are given as $J_{ijkl} = \partial A_{ij} / \partial B_{kl}$.

0.1 Gradient of a Matrix with Respect to a Matrix

For the matrix expression $A_{m \times p} = B_{m \times n} C_{n \times p}$, we can coordinate expand into $A_{ij} = \sum_{r=1}^n B_{ir} C_{rj}$ by definition of matrix multiplication such that the derivative

$$\begin{aligned} \frac{\partial A_{ij}}{\partial B_{kl}} &= \frac{\partial}{\partial B_{kl}} \left(\sum_{r=1}^n B_{ir} C_{rj} \right) \\ &= \sum_{r=1}^n \frac{\partial B_{ir}}{\partial B_{kl}} C_{rj} \end{aligned}$$

and we can define the **Kronecker delta** δ_{ij} as

$$\delta_{i,j} := \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

which effectively “selects” the term of $i = j$ from the sum. Since for an isolated matrix, a value B_{ij} does not depend on any other term but itself, its gradient with respect to other terms will always be 0 or 1. Hence

$$\begin{aligned} \frac{\partial A_{ij}}{\partial B_{kl}} &= \sum_{r=1}^n \frac{\partial B_{ir}}{\partial B_{kl}} C_{rj} \\ &= \sum_{r=1}^n \delta_{ik} \delta_{rl} C_{rj} \\ &= \delta_{ik} C_{lj}. \end{aligned}$$

Here, the first Kronecker delta δ_{ik} is not affected by the summation and is pulled out, and the second δ_{rl} sets every term except for C_{lj} to 0 where $r = l$.

0.2 Gradient of a Vector with Respect to a Matrix

In the context of an expression $\mathbf{y} = \mathbf{A}\mathbf{x}$, we can take the derivative starting again in coordinate form $y_i = \sum_{j=1}^n A_{ij} x_j$ to find

$$\frac{\partial y_i}{\partial A_{kl}} = \frac{\partial A_{ij}}{\partial A_{kl}} x_j = \delta_{ik} x_l$$

0.2.1 Simple Loss Function

If we have simple loss \mathcal{L} for a one-layer neural network then we can say

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial A_{kl}} &= \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial y_i} \frac{\partial y_i}{\partial A_{kl}} \\ &= \sum_{i=1}^m \sum_{j=1}^n \frac{\partial \mathcal{L}}{\partial y_i} \frac{\partial A_{ij}}{\partial A_{kl}} x_j \\ &= \sum_i \frac{\partial \mathcal{L}}{\partial y_i} \delta_{ik} x_l \\ &= \frac{\partial \mathcal{L}}{\partial y_k} x_l \end{aligned}$$

which is a *very* useful result in machine learning.

Single Layer Network

EXAMPLE 0.2.1

For a simple network with weights W , target vector \mathbf{t} , output vector $\mathbf{y} = W\mathbf{x}$ and loss function $\mathcal{L}(W) = \frac{1}{2} \|W\mathbf{x} - \mathbf{t}\|^2$, we seek to find the gradient $\nabla_W \mathcal{L}$.

In coordinate form

$$(W\mathbf{x})_i - t_i = \sum_{j=1}^n W_{ij} x_j - t_i$$

and using the chain rule to compute

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W_{ij}} &= \frac{\partial}{\partial W_{ij}} \left[\frac{1}{2} \left(\sum_{k=1}^n W_{ik} x_k - t_i \right)^2 \right] \\ &= \underbrace{\left(\sum_{k=1}^n W_{ik} x_k - t_i \right)}_{\text{define as } r_i} \underbrace{x_j}_{\frac{\partial}{\partial W_{ij}(\cdot)}} \\ &= ((W\mathbf{x})_i - t_i) x_j \end{aligned}$$

Now we connect this the vector form to say $\mathbf{r} = W\mathbf{x} - \mathbf{t}$ such that

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial W_{ij}} &= r_i x_j \\ (\nabla_W \mathcal{L})_{ij} &= (W\mathbf{x} - \mathbf{t})_i x_j \end{aligned}$$

Then by the **outer product**

$$\mathbf{r}\mathbf{x}^T = (W\mathbf{x} - \mathbf{t})\mathbf{x}^T$$

$$\nabla_W \mathcal{L} = (W\mathbf{x} - \mathbf{t})\mathbf{x}^T$$

□

Identities for Computing Gradients

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{f}(\mathbf{X})^\top = \left(\frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \right)^\top$$

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

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

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

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{x}^\top \mathbf{a} = \mathbf{a}^\top$$

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{a}^\top \mathbf{x} = \mathbf{a}^\top$$

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{a}^\top \mathbf{X} \mathbf{b} = \mathbf{a} \mathbf{b}^\top$$

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{x}^\top \mathbf{A} \mathbf{x} = \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top)$$