



Lecture 7

Math Foundations Team



BITS Pilani

Pilani | Dubai | Goa | Hyderabad



- ▶ In last lecture, we discussed about differentiation of univariate functions, partial differentiation, gradients and gradients of vector valued functions.
- ▶ Now we will look into gradients of matrices and some useful identities for computing gradients.
- ▶ Finally, we will discuss back propagation and automatic differentiation.



The gradient of an $m \times n$ matrix A with respect to a $p \times q$ matrix B , the resulting Jacobian would be an $(m \times n) \times (p \times q)$, i.e., a four-dimensional tensor J , whose entries are given as

$$J_{ijkl} = \frac{\partial A_{ij}}{\partial B_{kl}}$$

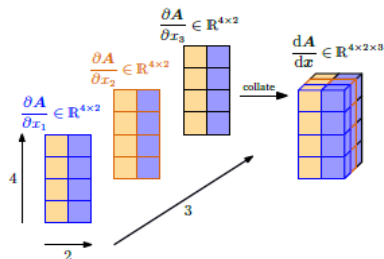
Since, we can consider $\mathbb{R}^{m \times n}$ as \mathbb{R}^{mn} , we can shape our matrix into vectors of length mn and pq respectively. The gradient using mn vectors results in a Jacobian of size $mn \times pq$

Gradients of Matrices

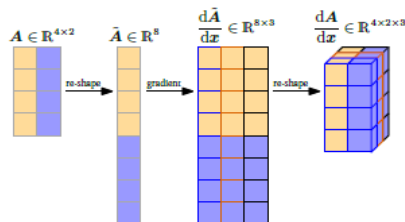


$$A \in \mathbb{R}^{4 \times 2} \quad x \in \mathbb{R}^3$$

Partial derivatives:



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



Let $f = Ax$ where $A \in \mathbb{R}^{m \times n}$, and $x \in \mathbb{R}^n$, then

$$\frac{\partial f}{\partial A} \in \mathbb{R}^{m \times (m \times n)}$$

By definition

$$\frac{\partial f}{\partial A} = \begin{bmatrix} \frac{\partial f_1}{\partial A} \\ \vdots \\ \frac{\partial f_m}{\partial A} \end{bmatrix}, \frac{\partial f_i}{\partial A} \in \mathbb{R}^{1 \times (m \times n)}$$



Now, we have

$$f_i = \sum_{j=1}^n A_{ij}x_j, i = 1, \dots, m.$$

Therefore, by taking partial derivatives with respect to A_{iq}

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

Hence, i^{th} row becomes

$$\frac{\partial f_i}{\partial A_{i,:}} = x^T \in \mathbb{R}^{1 \times 1 \times n}$$

$$\frac{\partial f_i}{\partial A_{k,:}} = 0^T \in \mathbb{R}^{1 \times 1 \times n}, \text{ for } k \neq i$$

Hence, by stacking the partial derivatives, we get

$$\frac{\partial f_i}{\partial A_{k,:}} = \begin{bmatrix} 0^T \\ \vdots \\ x^T \\ \vdots \\ 0^T \end{bmatrix} \in \mathbb{R}^{1 \times m \times n}$$



Let $B \in \mathbb{R}^{m \times n}$ and $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{n \times n}$ with

$$f(B) = B^T B =: K \in \mathbb{R}^{n \times n}$$

Then, we have

$$\frac{\partial K}{\partial B} \in \mathbb{R}^{(n \times n) \times (m \times n)}.$$

Moreover

$$\frac{\partial K_{pq}}{\partial B} \in \mathbb{R}^{1 \times (m \times n)}, \text{ for } p, q = 1, \dots, n$$

where K_{pq} is the $(p, q)^{th}$ entry of $K = f(B)$



Let i^{th} column of B be b_i , then

$$K_{pq} = r_p^T r_q = \sum_{l=1}^m B_{lp} B_{lq}$$

Computing the partial derivative, we get

$$\frac{\partial K_{pq}}{\partial B_{ij}} = \sum_{l=1}^m \frac{\partial}{\partial B_{ij}} B_{lp} B_{lq} = \partial_{pqij}$$



Clearly, we have

$$\partial_{pqij} = B_{iq} \quad \text{if } j = p, p \neq q$$

$$\partial_{pqij} = B_{ip} \quad \text{if } j = q, p \neq q$$

$$\partial_{pqij} = 2B_{iq} \quad \text{if } j = p, p = q$$

$$\partial_{pqij} = 0 \quad \text{otherwise}$$

where $p, q, j = 1, \dots, n$ $i = 1, \dots, m$

Useful Identities for Computing Gradients



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

Useful Identities for Computing Gradients



- ▶ $\frac{\partial a^T X^{-1} b}{\partial X} = -(X^{-1})^T a b^T (X^{-1})^T$
- ▶ $\frac{\partial x^T a}{\partial x} = a^T$
- ▶ $\frac{\partial a^T x}{\partial x} = a^T$
- ▶ $\frac{\partial a^T X b}{\partial X} = a b^T$
- ▶ $\frac{\partial x^T B}{\partial x} = x^T (B + B^T)$
- ▶ $\frac{\partial}{\partial s} (x - As)^T W (x - As) = -2(x - As)^T W A$

for symmetric W .

Consider the function

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2))$$

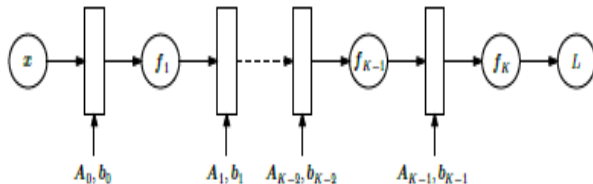
Taking derivatives

$$\begin{aligned}\frac{df}{dx} &= \frac{2x + 2x\exp(x^2)}{2\sqrt{x^2 + \exp(x^2)}} - \sin(x^2 + \exp(x^2))(2x + 2x\exp(x^2)) \\ &= 2x\left(\frac{1}{2\sqrt{x^2 + \exp(x^2)}} - \sin(x^2 + \exp(x^2))\right)(1 + \exp(x^2))\end{aligned}$$



- ▶ The implementation of the gradient could be significantly more expensive than computing the function, which imposes unnecessary overhead where we get such lengthy expressions.
- ▶ We need an efficient way to compute the gradient of an error function with respect to the parameters of the model.
- ▶ For training deep neural network models, the backpropagation algorithm is one such method.

In neural networks with multiple layers



$$f_i(x_{i-1}) = \sigma(A_{i-1}x_{i-1} + b_{i-1})$$

where x_{i-1} is the output of layer $i - 1$ and σ is an activation function.



To train these model, the gradient of the loss function L with respect to all model parameters $\theta_j = \{A_j, b_j\}, j = 1, \dots, K$ and inputs of each layer needs to be computed. Consider,

$$f_0 := x$$

$$f_i := \sigma_i(A_{i-1}f_{i-1} + b_{i-1}), i = 1, \dots, K.$$

We have to find $\theta_j = \{A_j, b_j\}, j = 1, \dots, K - 1$ such that

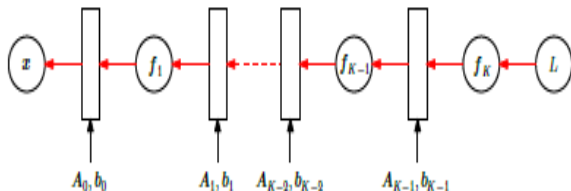
$$L(\theta) = ||y - f_K(\theta, x)||^2$$

is minimum where $\theta = \{A_0, b_0, \dots, A_{K-1}, b_{K-1}\}$

Using the chain rule, we get

$$\begin{aligned}\frac{\partial L}{\partial \theta_{K-1}} &= \frac{\partial L}{\partial f_K} \frac{\partial f_K}{\partial \theta_{K-1}} \\ \frac{\partial L}{\partial \theta_{K-2}} &= \frac{\partial L}{\partial f_K} \frac{f_K}{f_{K-1}} \frac{\partial f_{K-1}}{\partial \theta_{K-2}} \\ \frac{\partial L}{\partial \theta_{K-3}} &= \frac{\partial L}{\partial f_K} \frac{f_K}{f_{K-1}} \frac{\partial f_{K-1}}{\partial f_{K-2}} \frac{\partial f_{K-2}}{\partial \theta_{K-3}} \\ \frac{\partial L}{\partial \theta_i} &= \frac{\partial L}{\partial f_K} \frac{f_K}{f_{K-1}} \dots \frac{\partial f_{i+2}}{\partial f_{i+1}} \frac{\partial f_{i+1}}{\partial \theta_i}\end{aligned}$$

Backpropagation



If the partial derivatives $\frac{\partial L}{\partial \theta_{i+1}}$ are computed, then the computation can be reused to compute $\frac{\partial L}{\partial \theta_i}$.

Example

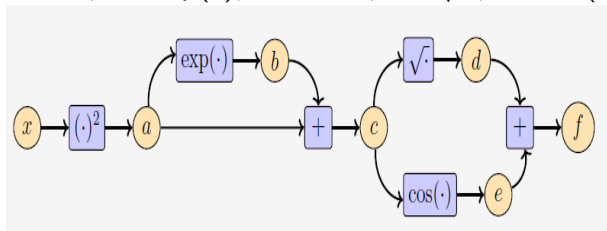


Consider the function

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2))$$

Let

$$a = x^2, b = \exp(a), c = a + b, d = \sqrt{c}, e = \cos(c) \Rightarrow f = d + e$$



Example



$$\begin{aligned}\Rightarrow \frac{\partial a}{\partial x} &= 2x \\ \frac{\partial b}{\partial a} &= \exp(a) \\ \frac{\partial c}{\partial a} &= 1 = \frac{\partial c}{\partial b} \\ \frac{\partial d}{\partial c} &= \frac{1}{2\sqrt{c}} \\ \frac{\partial e}{\partial c} &= -\sin(c) \\ \frac{\partial f}{\partial d} &= 1 = \frac{\partial f}{\partial e}\end{aligned}$$

Example



Thus, we have

$$\begin{aligned}\frac{\partial f}{\partial c} &= \frac{\partial f}{\partial d} \frac{\partial d}{\partial c} + \frac{\partial f}{\partial e} \frac{\partial e}{\partial c} \\ \frac{\partial f}{\partial b} &= \frac{\partial f}{\partial c} \frac{\partial c}{\partial b} \\ \frac{\partial f}{\partial a} &= \frac{\partial f}{\partial b} \frac{\partial b}{\partial a} + \frac{\partial f}{\partial c} \frac{\partial c}{\partial a} \\ \frac{\partial f}{\partial x} &= \frac{\partial f}{\partial a} \frac{\partial a}{\partial x}\end{aligned}$$

Example



Substituting the results, we get

$$\frac{\partial f}{\partial c} = 1.(\frac{1}{2\sqrt{c}} + 1).(-\sin(c))$$

$$\frac{\partial f}{\partial b} = \frac{\partial f}{\partial c}.1$$

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial b}\exp(a) + \frac{\partial f}{\partial c}.1$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial a}2x$$

Thus, the computation for calculating the derivative is of similar complexity as the computation of the function itself.



Let x_1, \dots, x_d : input variables.

x_{d+1}, \dots, x_{D-1} : intermediate variables.

x_D : output variable, then we have,

$$x_i = g_i(x_{Pa(x_i)})$$

Note that g_i s are elementary functions and are also called as forward propagation function and $x_{Pa(x_i)}$ is the set of parent nodes of variable x_i in the graph.



Now,

$$f = x_D \Rightarrow \frac{\partial f}{\partial_D} = 1$$

For other variables, using chain rule, we get

$$\frac{\partial f}{\partial x_i} = \sum_{x_j: x_i \in Pa(x_j)} \frac{\partial f}{\partial x_j} \frac{\partial x_j}{\partial x_i} = \sum_{x_j: x_i \in Pa(x_j)} \frac{\partial f}{\partial x_j} \frac{\partial g_i}{\partial x_i}$$

The last equation is the back propagation of the gradient through the computation graph. For neural network training, we back propagate the error of the prediction with respect to the label.