

MLP Backpropagation math

```
def forward(self, x):  
    x1 = x @ w1.T  
    x2 = sigmoid(x1)  
    x3 = x2 @ w2.T  
    x4 = softmax(x3)  
    return x4
```

Basic matrix function derivatives

1. Notice that, for scalar function derivatives, we can directly use chain rules to get the target variable derivative. However, for scalar-matrix function derivative, if we use chain rules, we may face to very complicated matrix-matrix function derivative. To avoid such dilemma, we will use the total derivative to get the final results.

2. Some basic rules:

$$\begin{aligned}d(X \pm Y) &= dX \pm dY, d(XY) = dXY + XdY \\dX \odot Y &= dX \odot Y + X \odot dY \\d(\sigma(X)) &= \sigma'(X) \odot dX \\tr(AB) &= tr(BA) \\tr(A^T(B \odot C)) &= tr((A \odot B)^T C), A, B, C \in \mathbb{R}^{n \times n} \\df &= tr\left(\frac{\partial f}{\partial X}^T dX\right)\end{aligned}$$

3. We will use formula (6) as the very basic beginning for getting our derivatives:

Backpropagation starts at loss function

$$\begin{aligned}l &= -y \log \text{softmax}(\sigma(XW_1^T)W_2^T)^T \\X &\in \mathbb{R}^{1 \times \dim 1}, W_1 \in \mathbb{R}^{\dim 2 \times \dim 1}, W_2 \in \mathbb{R}^{ncl \times \dim 2} \\a_1 &= XW_1^T, h_1 = \sigma(a_1), a_2 = h_1W_2^T \\&\text{Find: } \frac{\partial l}{\partial W_1} \text{ and } \frac{\partial l}{\partial W_2}\end{aligned}$$

In first way, we get to $\frac{\partial l}{\partial a_2}$, in a very traditional way:

$$\begin{aligned}
l &= -y \log \text{softmax}(a_2)^T \\
l &= -y(\log(\exp(a_2)) - \mathbf{1}_{ncl \times ncl} \log(\exp(a_2) \mathbf{1}_{ncl \times 1}))^T \\
l &= -y a_2^T + \log(\mathbf{1}_{ncl \times 1}^T \exp(a_2^T)) \\
dl &= -y dW_2 h_1^T + \frac{\mathbf{1}_{ncl \times 1}^T (\exp(W_2 h_1^T) \odot (dW_2 h_1^T))}{\mathbf{1}_{ncl \times 1}^T \exp(W_2 h_1^T)} \\
dl &= -y dW_2 h_1^T + \frac{\exp(h_1 W_2^T) dW_2 h_1^T}{\mathbf{1}_{ncl \times 1}^T \exp(W_2 h_1^T)} \\
dl &= \text{tr}(h_1^T (-y) + h_1^T \frac{\exp(h_1 W_2^T)}{\mathbf{1}_{ncl \times 1}^T \exp(W_2 h_1^T)} dW_2) \\
dl &= \text{tr}(-h_1^T y + h_1^T \text{softmax}(h_1 W_2^T) dW_2) \\
\frac{\partial l}{\partial W_2} &= (-y^T + \text{softmax}(W_2 h_1^T)) h_1
\end{aligned}$$

Or, we can have:

$$\begin{aligned}
\frac{\partial l}{\partial a_2} &= \text{softmax}(a_2) - y \quad (*1) \\
dl &= \text{tr}(\frac{\partial l}{\partial a_2}^T da_2) = \text{tr}(\frac{\partial l}{\partial a_2}^T dh_1 W_2^T + \frac{\partial l}{\partial a_2}^T h_1 dW_2^T) \\
dl &= \text{tr}((\frac{\partial l}{\partial a_2} W_2)^T dh_1) + \text{tr}((\frac{\partial l}{\partial a_2}^T h_1)^T dW_2) \\
\frac{\partial l}{\partial W_2} &= \frac{\partial l}{\partial a_2}^T h_1 = (-y^T + \text{softmax}(W_2 h_1^T)) h_1 \quad (*2)
\end{aligned}$$

The results are the same. Obviously, the later is much simpler.

On top of that, we need to get further results:

$$\begin{aligned}
dl_2 &= \text{tr}((\frac{\partial l}{\partial a_2} W_2)^T dh_1) \text{ ignore } W_2 \text{ part for further formula} \\
\frac{\partial l}{\partial h_1} &= \frac{\partial l}{\partial a_2} W_2 \quad (*3) \\
dl_2 &= \text{tr}((\frac{\partial l}{\partial h_1})^T d\sigma(a_1)) \\
dl_2 &= \text{tr}((\frac{\partial l}{\partial h_1})^T \sigma'(a_1) \odot da_1) \\
dl_2 &= \text{tr}((\frac{\partial l}{\partial h_1} \odot \sigma'(a_1))^T da_1) \\
\frac{\partial l}{\partial a_1} &= \frac{\partial l}{\partial h_1} \odot \sigma'(a_1) = \frac{\partial l}{\partial a_2} W_2 \odot \sigma'(a_1) \quad (*4) \\
dl_2 &= \text{tr}((\frac{\partial l}{\partial a_1})^T da_1) = \text{tr}((\frac{\partial l}{\partial a_1})^T X dW_1^T) \\
dl_2 &= \text{tr}(((\frac{\partial l}{\partial a_1})^T X)^T dW_1) \\
\frac{\partial l}{\partial W_1} &= \frac{\partial l}{\partial a_1}^T X \quad (*5)
\end{aligned}$$

Matrix Function Derivatives Methods:

1. We need a universal formula for matrix-matrix function derivatives:

$$\text{vec}(dF) = \frac{\partial F}{\partial X}^T \text{vec}(dX)$$

2. Composition of matrix – matrix function derivatives :

$$\text{vec}(dF) = \frac{\partial F}{\partial Y}^T \text{vec}(dY) = \frac{\partial F}{\partial Y}^T \frac{\partial Y}{\partial X}^T \text{vec}(dX)$$

And we could get use the chain rules for the complicated gradient computations

Properties:

$$\text{vec}(A + B) = \text{vec}(A) + \text{vec}(B)$$

$$\text{vec}(AXB) = (B^T \otimes A) \text{vec}(X)$$

$$\text{vec}(A^T) = K_{mn} \text{vec}(A), A \in \mathbb{R}^{m \times n}, K_{mn} \in \mathbb{R}^{mn \times mn}, K_{mn} \text{ is commutation matrix}$$

$$\text{vec}(A \odot X) = \text{diag}(\text{vec}(A)) \text{vec}(X), \text{diag}(A) \in \mathbb{R}^{mn \times mn} \text{ is a diagonal use elements in } A, \text{ ordered by columns}$$

$$(A \otimes B)^T = A^T \otimes B^T$$

$$\text{vec}(ab^T) = b \otimes a$$

$$(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$$

$$K_{mn} = K_{nm}^T, K_{mn} K_{nm} = I$$

$$K_{pm}(A \otimes B)K_{nq} = B \otimes A, A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{p \times q}$$