

Your Paper

You

November 18, 2023

Abstract

Your abstract.

1 Introduction

Remembering that

$$\frac{\partial g(M)}{\partial t} = \sum_{i,j} \frac{\partial g(M)}{\partial M_{i,j}} \frac{\partial M_{i,j}}{\partial t},$$

chain rule allows us to say the following:

$$\left[\frac{\partial L}{\partial W} \right]_{m,n} = \frac{\partial g(Y)}{\partial W_{m,n}}$$

$$\left[\frac{\partial L}{\partial W} \right]_{m,n} = \frac{\partial g(Y)}{\partial W_{m,n}} \quad \text{because of } L = g(Y) \quad (1)$$

$$= \sum_{i,j} \frac{\partial g(Y)}{\partial Y_{i,j}} \frac{\partial Y_{i,j}}{\partial W_{m,n}} \quad \text{because of property above + chain rule} \quad (2)$$

$$= \sum_{i,j} \frac{\partial g(Y)}{\partial Y_{i,j}} \frac{\partial}{\partial W_{m,n}} \left([XW^T]_{i,j} + \cancel{B_{i,j}} \right) \quad \text{Since } B_{i,j} \text{ is not a function of } W \quad (3)$$

$$(4)$$

Expanding $[XW^T]_{i,j}$ we get:

$$[XW^T]_{i,j} = \sum_k X_{i,k} W_{k,j}^T \quad (5)$$

$$= \sum_k X_{i,k} W_{j,k} \quad \text{flipping rows and columns for } W \text{ transpose} \quad (6)$$

Which now allows us to state:

$$\frac{\partial [XW^T]_{i,j}}{\partial W_{m,n}} = \sum_k X_{i,k} \frac{\partial W_{j,k}}{\partial W_{m,n}} \quad (7)$$

$$= \sum_{\cancel{k}} X_{i,k} \delta_{j,m} \delta_{\cancel{k},n} \quad \delta_{k,n} = 1 \text{ iff } k = n \quad (8)$$

$$= X_{i,n} \delta_{j,m} \quad (9)$$

Plugging in this result in the main equation above yields:

$$= \sum_{i,j} \frac{\partial g(Y)}{\partial Y_{i,j}} \frac{\partial}{\partial W_{m,n}} \left([XW^T]_{i,j} \right) \quad (10)$$

$$= \sum_{i,j} \frac{\partial g(Y)}{\partial Y_{i,j}} X_{i,n} \delta_{j,m} \quad (11)$$

$$= \sum_i \sum_{\cancel{j}} \frac{\partial g(Y)}{\partial Y_{i,j}} X_{i,n} \delta_{\cancel{j},m} \quad \delta_{k,n} = 1 \text{ iff } j = m \quad (12)$$

$$= \underbrace{\sum_i X_{i,n}}_{\text{loop over the rows}} \left[\frac{\partial L}{\partial Y} \right]_{i,m} \quad (13)$$

$$= X^T \frac{\partial L}{\partial Y} \quad (14)$$

again, since $L = g(Y)$:

$$\frac{\partial L}{\partial b_k} = \sum_{i,j} \frac{\partial g(Y)}{\partial Y_{i,j}} \frac{\partial Y_{i,j}}{\partial b_k} \quad (1)$$

We can unpack $Y_{i,j}$ just like we did above:

$$\frac{\partial Y_{i,j}}{\partial b_k} = \frac{\partial}{\partial b_k} \left(\cancel{[XW^T]_{i,j}} + B_{i,j} \right) \quad XW^T \text{ independent of } b_k \quad (2)$$

$$= \frac{\partial b_j}{\partial b_k} \quad \text{due to } B_{i,j} = b_j \quad (3)$$

$$= \delta_{j,k} \quad (4)$$

Plugging the above result into the main equation yields:

$$\frac{\partial L}{\partial b_k} = \sum_i \sum_j \cancel{\frac{\partial g(Y)}{\partial Y_{i,j}} \delta_{j,k}} \quad \delta_{j,k} = 1 \text{ iff } j = k \quad (5)$$

$$= \sum_i \frac{\partial L}{\partial Y_{i,k}} \quad (6)$$

$$\frac{\partial L}{\partial \mathbf{b}} = \sum_i \frac{\partial L}{\partial \mathbf{y}_i} \quad (7)$$

Once again, since $L = g(Y)$:

$$\left[\frac{\partial L}{\partial X} \right]_{m,n} = \frac{\partial g(Y)}{\partial X_{m,n}} \quad (1)$$

$$= \sum_{i,j} \frac{\partial g(Y)}{\partial Y_{i,j}} \frac{\partial Y_{i,j}}{\partial X_{m,n}} \quad \text{because of property above + chain rule} \quad (2)$$

$$= \sum_{i,j} \frac{\partial g(Y)}{\partial Y_{i,j}} \frac{\partial}{\partial X_{m,n}} \left([XW^T]_{i,j} + \cancel{B_{i,j}} \right) \quad \text{Since } B_{i,j} \text{ is not a function of } X \quad (3)$$

Expanding $[XW^T]_{i,j}$ we get:

$$[XW^T]_{i,j} = \sum_k X_{i,k} W_{k,j}^T \quad (4)$$

$$= \sum_k X_{i,k} W_{j,k} \quad \text{flipping rows and columns for } W \text{ transpose} \quad (5)$$

Which now allows us to state:

$$\frac{\partial [XW^T]_{i,j}}{\partial X_{m,n}} = \sum_k W_{j,k} \frac{\partial X_{i,k}}{\partial X_{m,n}} \quad (6)$$

$$= \sum_{/k} W_{j,k} \delta_{i,m} \delta_{k,n} \quad \delta_{k,n} = 1 \text{ iff } k = n \quad (7)$$

$$= W_{j,k} \delta_{i,m} \quad (8)$$

Plugging in this result in the main equation above yields:

$$= \sum_{i,j} \frac{\partial g(Y)}{\partial Y_{i,j}} \frac{\partial}{\partial X_{m,n}} \left([XW^T]_{i,j} \right) \quad (9)$$

$$= \sum_{i,j} \frac{\partial g(Y)}{\partial Y_{i,j}} W_{j,k} \delta_{i,m} \quad (10)$$

$$= \sum_{/i} \sum_j \frac{\partial g(Y)}{\partial Y_{i,j}} W_{j,k} \delta_{i,m} \quad \delta_{i,m} = 1 \text{ iff } i = m \quad (11)$$

$$= \underbrace{\sum_j W_{j,k}}_{\text{loop over the columns}} \left[\frac{\partial L}{\partial Y} \right]_{m,j} \quad (12)$$

$$= \frac{\partial L}{\partial Y} W \quad (13)$$

And the dimensions match due to $Y \in \mathbb{R}^{S \times N}$, resulting in $\frac{\partial L}{\partial Y} \in \mathbb{R}^{S \times N}$. Since $W^T \in \mathbb{R}^{M \times N}$, $W \in \mathbb{R}^{N \times M}$, thus the matrix product results in a $S \times M$ matrix, which is the original size of X .

Once again, since $L = g(Y)$ and $\mathbf{Y} = h(\mathbf{X}) \rightarrow Y_{i,j} = h(X_{i,j})$:

$$\left[\frac{\partial L}{\partial X} \right]_{m,n} = \frac{\partial g(Y)}{\partial X_{m,n}} = \frac{\partial g(h(X))}{\partial X_{m,n}} \quad (1)$$

$$= \sum_{i,j} \frac{\partial g(h(X))_{i,j}}{\partial h(X)_{i,j}} \frac{\partial h(X)_{i,j}}{\partial X_{i,j}} \frac{\partial X_{i,j}}{\partial X_{m,n}} \quad \text{chain rule} \quad (2)$$

$$= \sum_{i,j} \frac{\partial g(h(X))_{i,j}}{\partial h(X)_{i,j}} \frac{\partial h(X)_{i,j}}{\partial X_{i,j}} \delta_{i,m} \delta_{j,n} \quad (3)$$

$$= \sum_i \sum_j \frac{\partial g(h(X))_{i,j}}{\partial h(X)_{i,j}} \frac{\partial h(X)_{i,j}}{\partial X_{i,j}} \delta_{i,m} \delta_{j,n} \quad \delta_{i,m}, \delta_{j,n} = 1 \text{ iff } i = m, j = n \quad (4)$$

$$= \frac{\partial L}{\partial Y_{m,n}} \frac{\partial h(X)_{m,n}}{\partial X_{m,n}} \quad (5)$$

Thus, rewriting the result in matrix notation provides us with the result

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial Y} \circ \frac{\partial h(X)}{\partial X},$$

where \circ is the Hadamard product.

Since in a local minimum the Hessian is positive definite,

$$\mathbf{x}^\top \mathbf{H} \mathbf{x} > 0 \quad \forall \mathbf{x} \neq \mathbf{0}$$

By definition, an eigenvalue of \mathbf{H} is the one scalar λ for which

$$\mathbf{H} \mathbf{x} = \lambda \mathbf{x}$$

Left multiplying both sides by \mathbf{x}^\top we get

$$\mathbf{x}^\top \mathbf{H} \mathbf{x} = \mathbf{x}^\top \lambda \mathbf{x}$$

Since $\mathbf{x}^\top \mathbf{H} \mathbf{x}$ is greater than 0, focusing on the RHS we can see

$$\lambda \mathbf{x}^\top \mathbf{x}$$

and given that λ is a scalar, commutative, it follows that

$$= \lambda \|\mathbf{x}\|^2$$

The norm of a non-zero vector is greater than zero by definition, which leads us to the desired result of λ being greater than zero for the equality to hold.