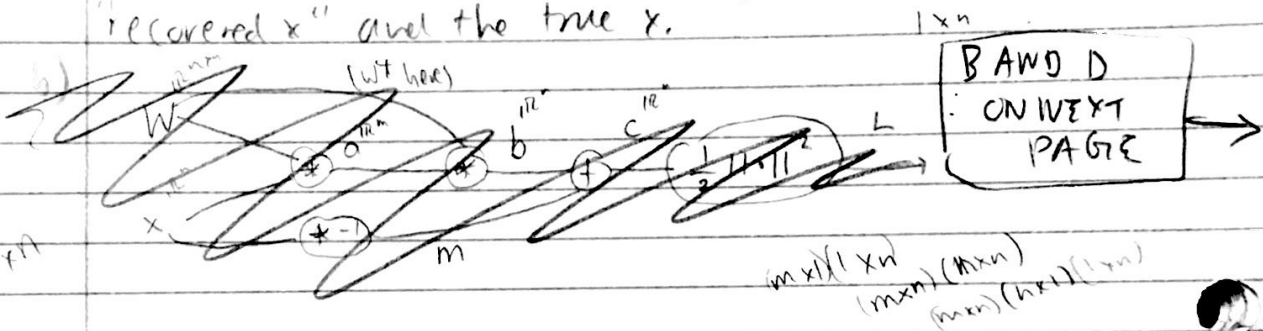


Scalar chain: $\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$ $(W^T W x - x)(W x + W^T x)$

$\frac{dz}{dx} = \frac{dy}{dx} \frac{dz}{dy}$ $a = m \times 1$ $c = n \times 1$ $W = m \times n$

- a) We know that the transformation Wx would preserve information perfectly if we "reverse" this transform via $W^T(Wx)$ and get an approximation a that is "close" to x , for some closeness metric. Essentially, if we approximate x by $W^T W x$ where Wx reduces the dimensionality of x , then W would have to be selected such that it learns important features about x , so that it can be recovered. In order to do this, we'd want to minimize a cost such as the L2 squared distance between our "recovered x " and the true x .



- c) The two paths to b should be summed. This is because we can think of the 2 paths out of a as separate functions that w affects. For example, we could assign each f and g as functions coming out of the w 's node. Then, f and g are both functions of w , and contribute some incoming gradient $\nabla_w f$ and $\nabla_w g$ to the node w , so these gradients should be summed: $\nabla_w f + \nabla_w g$

- d) Let $\frac{dL}{dc} = 1$.

$\frac{dL}{dc}$ (scalar w.r.t to vector) $= \frac{1}{2} \|c\|_2^2 = c$ $c \in \mathbb{R}^n$

$\frac{dL}{db} = \frac{dc}{db} \frac{dL}{dc} = \frac{dc}{db} c = \frac{d(b+m)(c)}{db} = c$

$\frac{dL}{dw^T}$ (from b) $= \left(\frac{db}{dw^T} \frac{dL}{db} \right) = \frac{dW^T a(c)}{dw^T} = ac^T$

$\frac{dL}{da} = \frac{db}{da} \frac{dL}{db} = \frac{dW^T a(c)}{da} = Wc$

$\frac{dL}{dw}$ (from a) $= \frac{da}{dw} \frac{dL}{da} = \frac{d(Wx)(Wc)}{dw} = Wcx^T$

$\frac{dL}{dw} = ac^T + Wcx^T$
 $= Wx(W^T W x - x)^T + W(W^T W x - x)x^T$

h y m

$$w_{n \times m} \frac{dx}{(n \times 1)(1 \times n)}$$

$$\frac{h \times l}{12^m}$$



$m \times 1$

$$m \times n \times m$$

4612

$$A_{sm}^+ \quad A_{cm}^+ \quad I_m^T$$

$$m \times l$$

we can infer that $dL/dw^T = \boxed{Ca^T}$

$$\rightarrow Wx(W^TWx - x)^T + W(W^TWx - x)x^T$$