

Backdrop 4

Proof by my own

Vector = x, y, b, a, z

Matrix = W

Scalar = $x_m, y_n, w_{mn}, L, M, N$

Function = $f_{NN}, g, \sigma, C, \text{len}$

Def: Neural Network is a composed function with minimizes a loss function C

$$f_{NN}(x) = f \left(f^{L-1} \left(\dots f^1 \left(f^0(x) \right) \right) \right) = \hat{x}$$

LAYERS L

Length of 0

$$\text{Dataset} = D = \{(x^i, y^i)\}$$

$$\min_C C(f_{NN}(x^i), y^i)$$

WHAT WE WANT

$$\hat{a}^L = x \quad \hat{a}^{L-1} = \hat{x} \quad C = (\hat{x} - x)^2$$

$$g^L(z) = z$$

$$g = g$$

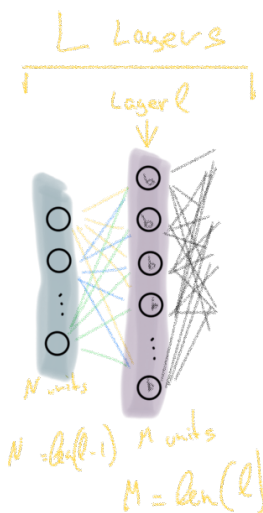
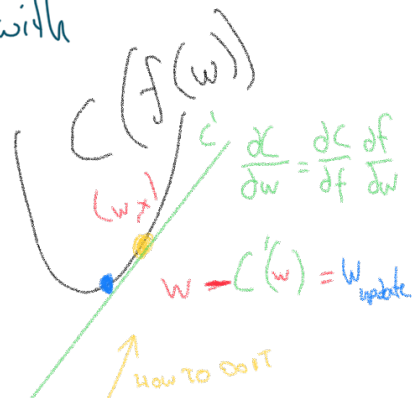
$$f^L(a^{L-1}) = a^L = g^L \left(W^L \cdot a^{L-1} + b^L \right)$$

activation

$$\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_M \end{bmatrix} = z^L = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1N} \\ w_{21} & w_{22} & \dots & w_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ w_{M1} & w_{M2} & \dots & w_{MN} \end{bmatrix} \cdot \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_N \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_M \end{bmatrix}$$

any

$$z_m^L = (w_{m1}^L \cdot a_1^{L-1} + w_{m2}^L \cdot a_2^{L-1} + \dots + w_{mN}^L \cdot a_N^{L-1} + b_m^L)$$



$$\frac{\partial C}{\partial a_n^{l-1}} = \begin{cases} \text{if } l = L-1 \\ \frac{\partial C}{\partial \hat{x}} = 2(\hat{x}_n - y_n)^2 \\ \text{else} \\ \sum_{m=0}^{M-1} \frac{\partial C}{\partial a_m^l} \frac{\partial a_m^l}{\partial z_m^l} \frac{\partial z_m^l}{\partial a_n^{l-1}} = \sum_{m=0}^{M-1} \frac{\partial C}{\partial a_m^l} g'^l(z_m^l) \cdot W_{mn}^l \end{cases}$$

HEIGHT OF LAYER l

$$\frac{\partial C}{\partial W_{mn}^l} = \frac{\partial C}{\partial a_m^l} \frac{\partial a_m^l}{\partial z_m^l} \frac{\partial z_m^l}{\partial W_{mn}^l} = \frac{\partial C}{\partial a_m^l} g'^l(z_m^l) \cdot a_n^{l-1}$$

$$\frac{\partial C}{\partial b_m^l} = \frac{\partial C}{\partial a_m^l} \frac{\partial a_m^l}{\partial z_m^l} \frac{\partial z_m^l}{\partial b_m^l} = \frac{\partial C}{\partial a_m^l} g'^l(z_m^l)$$

THESE EQUATIONS DESCRIBE GRADIENT DESCENT

We need to update the weights and biases once we have $\frac{\partial C}{\partial a_n^l}$ and $\frac{\partial C}{\partial b_m^l}$ but this correction is computed by looking at a single datapoint (x_i, y_i) . To get the true correction we must average it over all datapoints.

$$W_{mn}^l \leftarrow W_{mn}^l - \rho \cdot \frac{1}{D} \sum_{d=0}^{D-1} \frac{\partial C_d}{\partial W_{mn}^l}$$

Updated Learning rate

$$b_m^l \leftarrow b_m^l - \rho \cdot \frac{1}{D} \sum_{d=0}^{D-1} \frac{\partial C_d}{\partial b_m^l}$$

Error for datapoint d

