$$\frac{\partial \mathcal{L}}{\partial W_{ij}^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial W_{ij}^{[2]}}$$

$$= \frac{\partial \mathcal{L}}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial W_{ij}^{[2]}}$$

$$= \left( a^{[3]} - y \right) \cdot W^{[3]} \; \mathrm{diag}\left[ g'\left(z^{[2]}\right) \right]$$

$$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ a_j \\ \vdots \end{bmatrix} \leftarrow \; i^{th} \; \text{position}$$

element
wise product

$$= \left[ \left( a^{[3]} - y \right) \; W^{[3]} \underset{\underbrace{\qquad X \qquad}}{\odot} g'\left[ z^{[2]} \right] \right] \begin{bmatrix} 0 \\ 0 \\ a_j^{[1]} \\ 0 \end{bmatrix} \_\; i^{th}$$

$$= \quad X_i \, a_j^{[1]}$$

$[\ i^{th}$ element of vector $X$

$$\frac{\partial \mathcal{L}}{\partial W^{[2]}} = X \, a^{[1]T}$$

# Network View

$$f : \mathbb{R}^n \rightarrow \mathbb{R}^m$$
$$\frac{\partial f}{\partial x} \in \mathbb{R}^{m \times n}$$

$\longleftarrow$ Row vector $\delta^{[2]}$ $\longrightarrow$

$\overbrace{\qquad}^{\text{diag } g' \, w^{[l-1]} \dots}$

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^{[2]}} = \frac{\partial \mathcal{L}}{\partial a^{[L]}} \underbrace{\frac{\partial a^{[L]}}{\partial z^{[L]}}}_{1 \times 1} \underbrace{\frac{\partial z^{[L]}}{\partial a^{[L-1]}}}_{1 \times m} \dots \underbrace{\frac{\partial a^{[2]}}{\partial z^{[2]}}}_{m \times m} \underbrace{\frac{\partial z^{[2]}}{\partial W_{ij}^{[2]}}}_{m \times 1}$$

$\underbrace{\phantom{xxxx}}_{[\times]}$  $\underbrace{\phantom{xx}}_{1 \times 1}$

Logistic Regression
$$[y^{(i)} - h_\theta(x^{(i)})] x^{(i)} \leftarrow \frac{\partial h}{\partial X}$$
$$(y^{(i)} - a^{[L]}) w^{[L]} \leftarrow \frac{\partial \mathcal{L}}{\partial a}$$

$\underset{\text{column vector}}{\longleftrightarrow}$

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^{[1]}} = \frac{\partial \mathcal{L}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial w_{ij}^{[1]}}$$

$$1 \times 1 \qquad 1 \times m \qquad m \times m \qquad m \times m \qquad m \times 1$$

$$\frac{\partial a^{[L-1]}}{\partial z^{[L-1]}} = \begin{bmatrix} g' & & \\ & g' & O \\ & & g' \\ O & & g' \\ & & g' \end{bmatrix} \Bigg\updownarrow a$$

$$\longleftarrow z \longrightarrow$$

Diagonal

$$a_i = g(z_i)$$

$$\frac{\partial z^{[2]}}{\partial w_{ij}^{[2]}} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ a_j \end{bmatrix} \leftarrow i^{th} \text{ position}$$

$$z_i = \sum_j w_{ij} a_j + b_i$$

$$\frac{\partial \lambda}{\partial w_{ij}^{[2]}} = \delta^{[2]} \cdot \begin{bmatrix} 0 \\ \vdots \\ a_j \\ \vdots \end{bmatrix}$$

$$= \frac{\partial \lambda}{\partial z^{[2]}} \begin{bmatrix} 0 \\ \vdots \\ a_j \\ \vdots \end{bmatrix} = \left[\frac{\partial \lambda}{\partial z^{[2]}}\right]_i a_j$$

$$\frac{\partial L}{\partial w^{[2]}} = \left(\frac{\partial \lambda^T}{\partial z^{[2]}}\right)\left(a^{[1]T}\right) \longrightarrow \text{think as} \atop \text{outer product}$$

$$w^{[l]} := w^{[l]} - \alpha \frac{\partial L}{\partial w^{[l]}}$$

$$J(w, b) = \sum_{i=1}^{B} \mathcal{L}\left(y^{(i)}, \hat{y}^{(i)}\right)$$

$B$ is size of batch

$$\underset{m}{\underline{Z^{[1]}}} = \underset{m \times d}{\underline{w^{[1]}}} \underset{d}{\underline{x^{(i)}}} + \underset{m}{\underline{b^{[1]}}}$$

$$Z^{[1]} = \underset{m \times d}{\underline{w^{[1]}}} \underbrace{\left[ x^{(1)} \; x^{(2)} \cdots x^{(B)} \right]}_{d \times B} + \underset{m \times 1}{\underline{b^{[1]}}}$$

$\xleftarrow{\hspace{4cm}}$ $m \times B$ $\xrightarrow{\hspace{1cm}}$ $\xleftarrow{} m \times 1 \xrightarrow{}$

Broadcasting

# Stochastic Gradient Descent

$^{(i)}$

$$\alpha$$

$$\sum_t \alpha^{(t)} = \infty$$

$$\sum_t \alpha^{(t)^2} < \infty$$

$\left.\begin{array}{c}\end{array}\right\}$ if cond$^n$ satisfied local optima SGD

NN. : learnable feature map
$\underbrace{\text{linear model}}_{\text{GLM}}$

You learning representation [like $x \to \phi(x)$] & then you apply GLM

Associated with each n.n. is a kernel,

$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j), \quad \phi \text{ takes you have}$$

input to represent learnt in $(l-1)^{th}$ layer

So you can learn kernel & combine it with G.P.

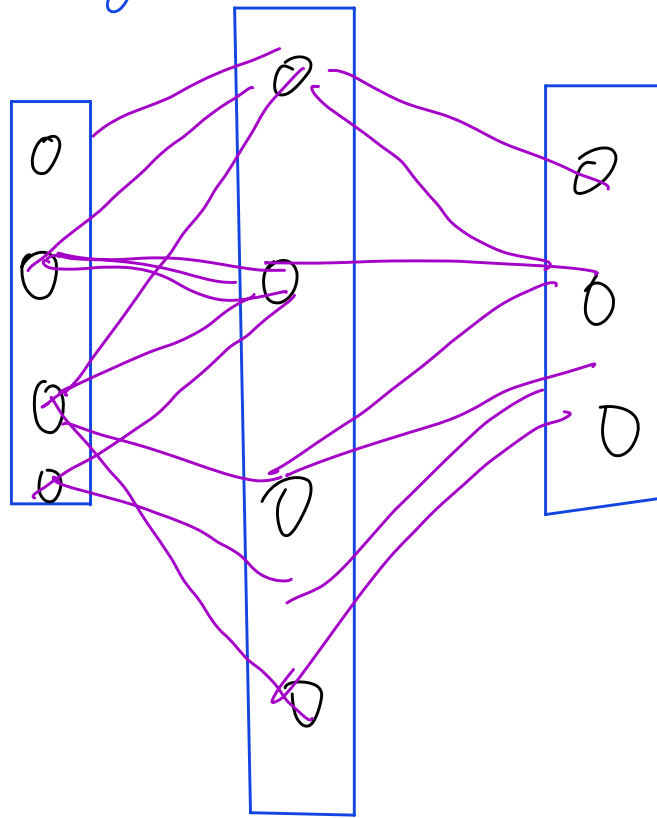N.N. becomes G.P. $\Leftarrow$

Also, 2 layer NN as no. of neuron in hidden layer $\to \infty$

# Universal Approximation Theorem

$y = f(x)$

$x \in \mathbb{R}^d$

$y = \in \mathbb{R}^k$

$\varepsilon = 10^{-6}$

bounded region of $x$

Exists a NN of 1 hidden layer.



that mimics function to an arbitrary precision assuming function is continuous