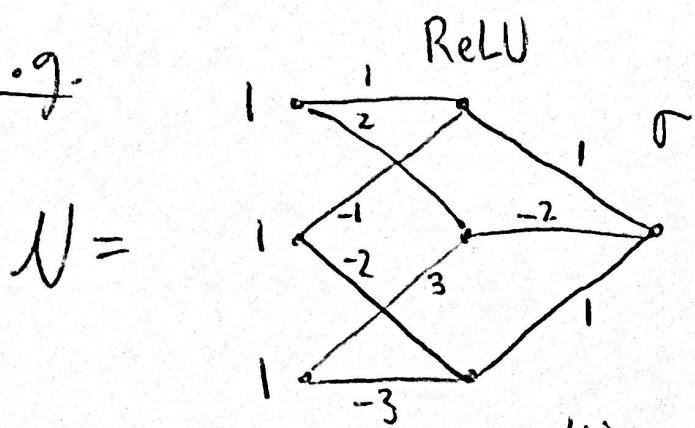


E.g.



$$b^{(0)} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad b^{(1)} = 0$$

$$N\begin{pmatrix} 1 \\ 1 \end{pmatrix} = \sigma F^{(1)} \text{ReLU} F^{(0)}\begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$F^{(0)}\vec{x} = \begin{pmatrix} 1 & -1 & 0 \\ 2 & 0 & 3 \\ 0 & -2 & -3 \end{pmatrix} \vec{x} + \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$$= \sigma F^{(1)} \text{ReLU} \begin{pmatrix} 1 \\ 6 \\ -4 \end{pmatrix}$$

$$F^{(1)}\vec{x} = (1 \ -2 \ 1) \vec{x}$$

$$= \sigma F^{(1)} \begin{pmatrix} 1 \\ 6 \\ 0 \end{pmatrix}$$

$$= \sigma(-1)$$

$$= \frac{1}{1 + e^{-1}}$$

$$\approx .000051 \quad (\text{small})$$

E.g. Suppose N has depth D and $G_1^{(1)}, G_2^{(1)}, \dots, G_B^{(1)}$ are all linear activations ($G_i^{(1)}(\vec{x}) = \vec{x} \ \forall i$)

$$\text{Then } N\vec{x} = f_{D,1}^{(D-1)} f_{D,2}^{(D-2)} \dots f_1^{(1)} f_0^{(0)} \vec{x}$$

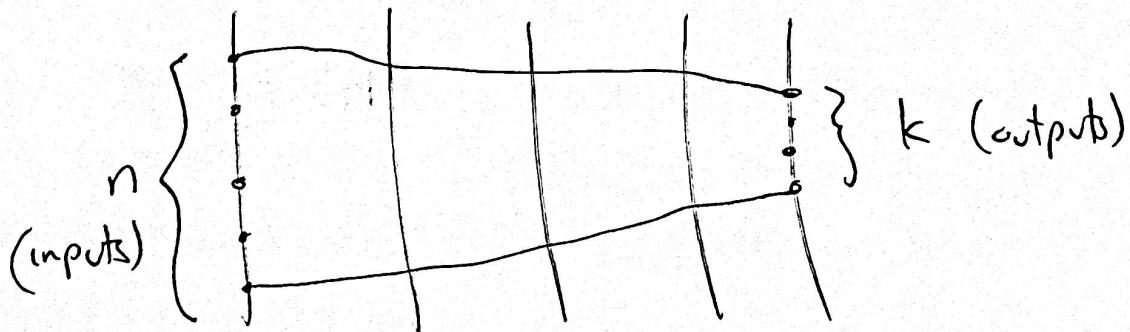
Remark The composition of affine maps is an affine map, so $\exists W, \vec{b}$ s.t.

$$N\vec{x} = W\vec{x} + \vec{b}$$

(so this NN can only fit affine functions)
 i.e. this is essentially linear regression

Neural Nets as Learning Algorithms

Suppose there is a function $f^*: \mathbb{R}^n \rightarrow \mathbb{R}^k$ (often $k=1$) and we wish to "approximate" or "model" f^* . Let N be a neural network $N = (H, F, G)$ with depth D & $H = (P, \subseteq, \varphi, \leq')$
 s.t. $\#P^{(0)} = n$ and $\#P^{(D)} = k$



$$F^{(i)}(\vec{x}) = W^{(i)}\vec{x} + b^{(i)}$$

where nonzero entries
 of $W^{(i)}, b^{(i)}$ are
 treated as parameters

Let $\mathbf{W} = (W^{(0)}, \dots, W^{(D-1)})$

$\mathbf{B} = (\vec{b}^{(0)}, \dots, \vec{b}^{(D-1)})$

& define $f_{\text{model}}(\vec{x}; \cancel{\mathbf{W}, \mathbf{B}}) = N\vec{x}$
(feedforward
f \vec{x})

How to find appropriate values of \mathbf{W}, \mathbf{B} ?

Given a dataset $X \in \mathbb{R}^{m \times n}, Y \in \mathbb{R}^{m \times k}$ we define a cost function

~~MB~~

$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^m L(f_{\text{model}}(\vec{x}^{(i)}; \mathbf{W}, \mathbf{B}), \cancel{y^{(i)}})$$

$y^{(i)}$

per-example error
function or
"loss function"

& try to find values for \mathbf{W}, \mathbf{B} which minimize this using some optimization algorithm

Note: We will usually need to add a regularization term (soon)

Note: A common choice for L is

$$L(f_{\text{model}}(\vec{x}; W, B), \vec{y}) = \|f_{\text{model}}(\vec{x}; W, B) - \vec{y}\|_2^2$$

As a maximum likelihood estimator:

Let's assume X, Y are sampled from a distribution $p_{\text{data}}(\vec{x}, \vec{y})$ under I.I.D. assumptions.

Let $p_{\text{model}}(\vec{y} | \vec{x}; \theta)$ be a model for $p_{\text{data}}(\vec{y} | \vec{x})$, parameterized by $\theta \in \mathbb{R}^N$ some N .

$$\text{Then } \theta_{\text{ML}} = \underset{\theta}{\operatorname{argmin}} \left(-\sum_{i=1}^m \log(p_{\text{model}}(\vec{y}^{(i)} | \vec{x}^{(i)}; \theta)) \right)$$

E.g. Suppose we choose p_{model} to be Gaussian:

$$p_{\text{model}}(\vec{y} | \vec{x}; W, B) = N(\vec{y}; N\vec{x}, \Sigma)$$

Normal distribution Neural network

$$= \frac{1}{\sqrt{2\pi^n \det \Sigma}} \exp \left(-\frac{1}{2} (\vec{y} - N\vec{x}) \Sigma^{-1} (\vec{y} - N\vec{x}) \right)$$

$$\text{where } \mathcal{N} = (H, F, G) \quad F_{\vec{x}}^{(i)} = W^{(i)} \vec{x} + \vec{b}^{(i)}$$

$$W = (W^{(0)}, \dots, W^{(D-1)})$$

$$B = (\vec{b}^{(0)}, \dots, \vec{b}^{(D-1)})$$

Then finding $\theta_{ML} = (W, B)_{ML}$ is equivalent to minimizing

$$J(W, B) = - \sum_{i=1}^m \log \left(\frac{1}{\sqrt{2\pi^n \det \Sigma}} \exp \left(-\frac{1}{2} (\vec{y}^{(i)} - \vec{N}_{\vec{x}^{(i)}}) \Sigma^{-1} (\vec{y}^{(i)} - \vec{N}_{\vec{x}^{(i)}})^T \right) \right)$$

$$= \frac{1}{m} \sum - \log \left(\frac{1}{\text{const}} \right) + \frac{1}{2} (\vec{y}^{(i)} - \vec{N}_{\vec{x}^{(i)}}) \Sigma^{-1} (\vec{y}^{(i)} - \vec{N}_{\vec{x}^{(i)}})^T$$

↑ Does not affect argmin
(throw it away)

$$J(W, B) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (\vec{y}^{(i)} - \vec{N}_{\vec{x}^{(i)}}) \Sigma^{-1} (\vec{y}^{(i)} - \vec{N}_{\vec{x}^{(i)}})^T$$

(Note: can usually take ~~Σ~~ $\Sigma = \text{id}$ by
first normalizing the dataset if needed)

in which case

$$J(W, B) = \frac{1}{m} \sum_{i=1}^m \left\| \vec{y}^{(i)} - W \vec{x}^{(i)} \right\|_2^2$$

(Naive)

so we recover our cost function from before.

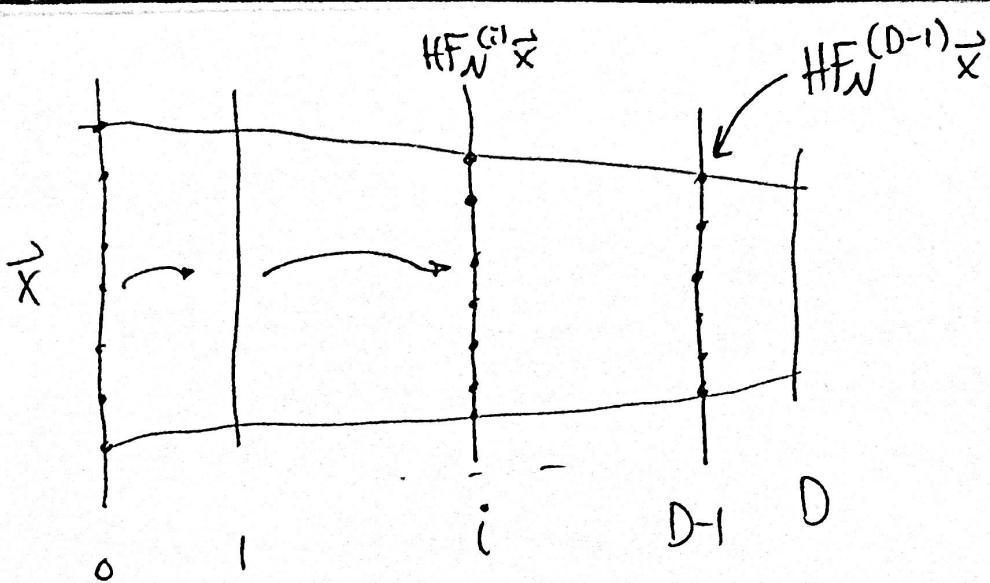
(Note: We see we can either specify a cost or specify a model for pdata and can usually go back & forth between those)

Output Neurons

Def: Given a neural network $N = (H, F, G)$ with depth D , for $\vec{x} \in \mathbb{R}^{\# P_0} \cong V_P^{(0)}$ (an input vector), the hidden features associated to \vec{x} at layer i are

$$HF_N^{(i)} \vec{x} = G_i F_{i-1} \dots F_1 G_1 F_0 \vec{x}$$

$$\left(\text{also set } HF_N^{(0)} \vec{x} = \vec{x} \right)$$



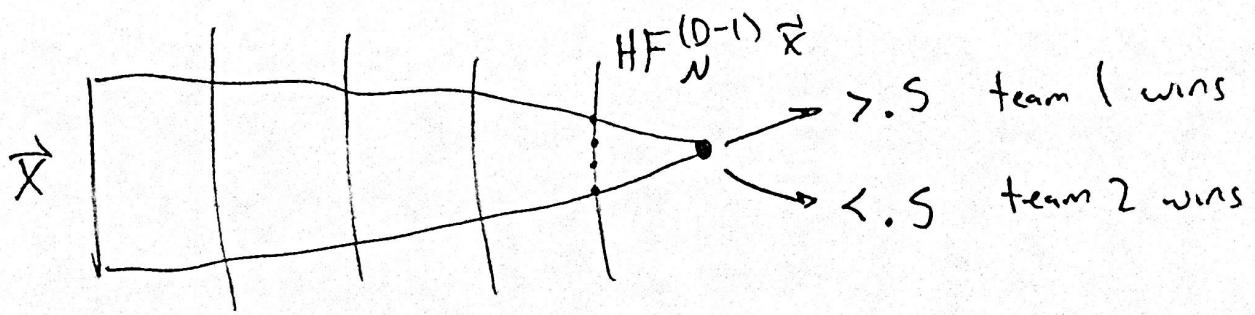
So the output neurons are responsible for predicting the target from $HF_N^{(0-1)} \vec{x}$. The choice of activation at the output layer depends on the type of targets possible
 (i.e. on $p_{\text{data}}(\vec{y})$)

E.g. Input \vec{x} = (some statistics on two sports teams)
 $y = \begin{cases} 0 & \text{if team 1 wins} \\ 1 & \text{if team 2 wins} \end{cases}$

Dataset X, \vec{y} of previous results. Goal:
 predict winners of upcoming games.

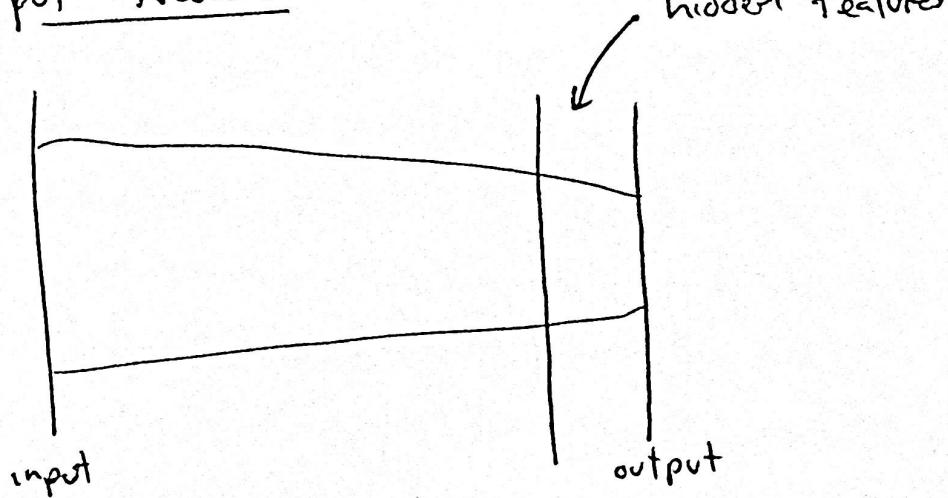
Idea: Apply a neural network N so that

$$N\vec{x} = p(y=1 | \vec{x})$$



Sigmoid activation makes sense here.

Output Neurons

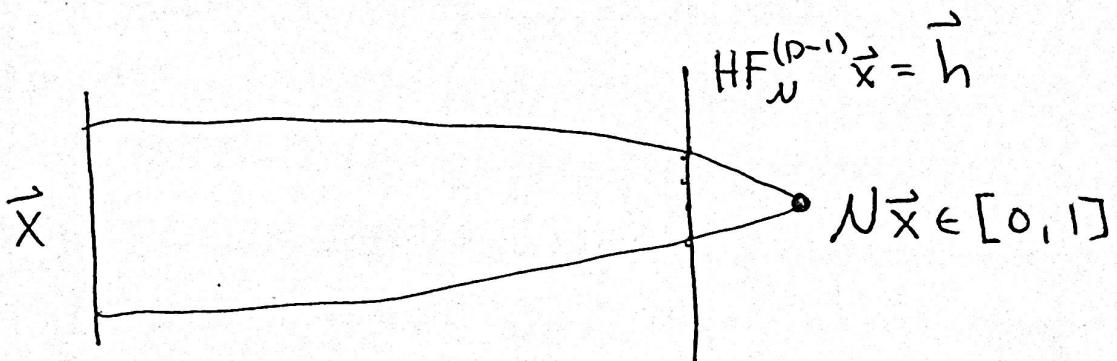


$$\cancel{\vec{x}^{(i)} = HF_N^i(\vec{x}) = G^{(i)} F^{(i-1)} \dots G^{(1)} f_{\alpha}^{(0)} \vec{x}}$$

E.g. Input \vec{x} statistics on sports teams 1 & 2

$$y = \begin{cases} 1 & \text{team 1 wins} \\ 0 & \text{team 2 wins} \end{cases}$$

Idea: want $N\vec{x} = p(y=1 | \vec{x})$

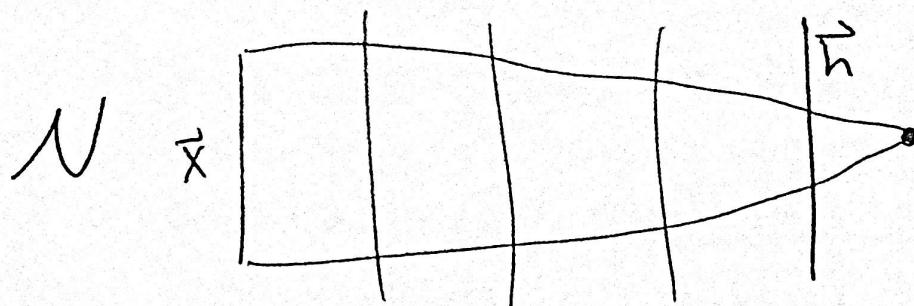


Makes sense to use sigmoid activation for $G^{(D)}$

Note: This is equivalent to transforming \vec{x} to new features \vec{h} and performing logistic regression on \vec{h} .

E.g. \vec{x} statistics about a company
 y stock price in 1 week

(y has
 continuous distribn.
 e.g. Gaussian)



linear activation $G^{(D)}(y) = y$ makes sense here.

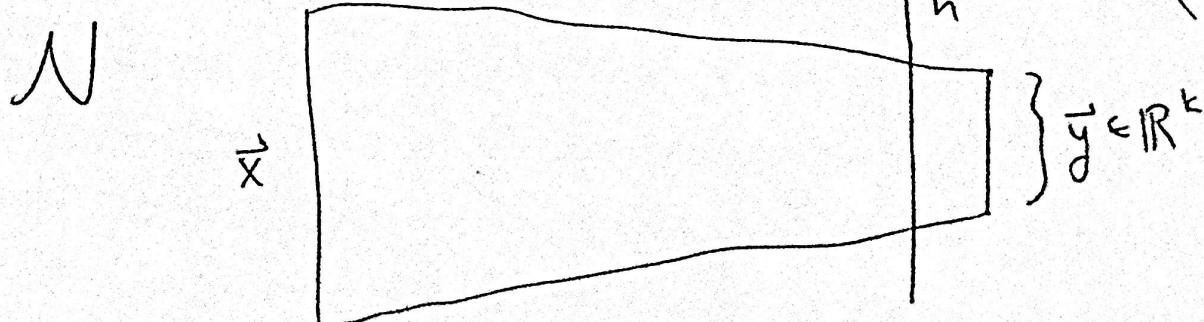
Note: this is equivalent to performing linear regression on $H\mathcal{F}_N^{(D-1)}(\vec{x})$

E.g. \vec{x} = vectorized image of an animal's types

$$y = \begin{cases} e_1 & \text{if cat} \\ e_2 & \text{if dog} \\ \vdots \\ e_k & \text{if llama} \end{cases} \in \mathbb{R}^k = \langle e_1, \dots, e_k \rangle$$

(called a one-hot encoding)

want $N\vec{x} = \hat{y}$ where $\hat{y}_i = p(y=e_i | \vec{x})$ $\left(\sum_{i=1}^k \hat{y}_i = 1 \right)$



Define $G^{(D)} = \text{softmax}$ where $\text{softmax}(h_i) = \frac{\exp(h_i)}{\sum_{j=1}^k \exp(h_j)}$

Softmax makes sense here.

Remarks In each case we are choosing an activation function which will give the gradient nice properties. We want the gradient to have the property that when W, B are far from a local min, the gradient stays sufficiently large, (so GD does not "get stuck").

Optimization for Neural Networks

Let N be a neural network and $J(W, B)$ a cost function to measure "proximity" between $N\vec{x}$ and target \vec{y} on a dataset X, Y .

To run gradient descent to find good values for W, B we need to compute

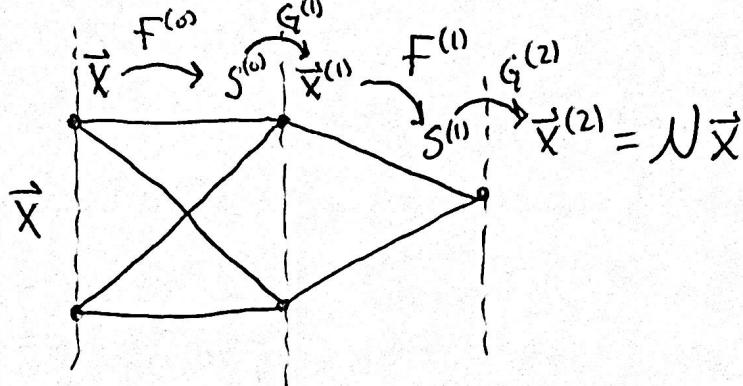
$$\frac{\partial J}{\partial W^{(i)}}, \frac{\partial J}{\partial b^{(i)}} \quad 0 \leq i \leq D-1$$

Definition The signal out of layer i for input \vec{x} is

$$S^{(i)} = F^{(i)} G^{(i)} \dots G^{(1)} F^{(0)} \vec{x} \quad S^{(0)} = F^{(0)} \vec{x}$$

Note: $G^{(i+1)} S^{(i)} = H F_N^{(i+1)}(\vec{x}) = \vec{x}^{(i+1)}$

E.g.



$$W^{(0)} = \begin{pmatrix} w_{11}^0 & w_{12}^0 \\ w_{21}^0 & w_{22}^0 \end{pmatrix} \quad \vec{b}^{(0)} = \begin{pmatrix} b_1^0 \\ b_2^0 \end{pmatrix}$$

$$W^{(1)} = (w'_{11} \quad w'_{12}) \quad \vec{b}^{(1)} = (b'_1)$$

$$S^{(0)} = W^{(0)} \vec{x} + \vec{b}^{(0)}$$

$$\vec{x}^{(1)} = G^{(1)}(S^{(0)}) = G^{(1)}(W^{(0)} \vec{x} + \vec{b}^{(0)})$$

$$S^{(1)} = W^{(1)} G^{(1)}(S^{(0)}) + b^{(1)} = W^{(1)} \vec{x}^{(1)} + b^{(1)}$$

$$\vec{x}^{(2)} = G^{(2)}(S^{(1)}) = N \vec{x}$$

Cost function $J(W, B) = \sum_{i=1}^m L(W\vec{x}_i, \vec{y}_i)$

Need to find derivatives of $L(W\vec{x}, \vec{y})$

Note: $\frac{\partial}{\partial v} (f \circ g)(v) = \frac{\partial f(w)}{\partial w} \Big|_{w=g(v)} \frac{\partial}{\partial v} g(v)$

write $\frac{\partial f(w)}{\partial w} \Big|_{w=g(v)} = \frac{\partial f}{\partial g(v)}$

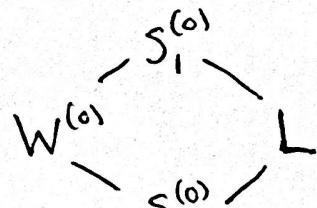
$$\frac{\partial L}{\partial b^{(1)}} = \frac{\partial L}{\partial S^{(1)}} \frac{\partial S^{(1)}}{\partial b^{(1)}} = \frac{\partial L}{\partial S^{(1)}} := \delta^{(1)}$$

$$\frac{\partial L}{\partial W^{(1)}} = \frac{\partial L}{\partial S^{(1)}} \frac{\partial S^{(1)}}{\partial W^{(1)}} = \frac{\partial L}{\partial S^{(1)}} (\vec{X}^{(1)})^T = \delta^{(1)} (\vec{X}^{(1)})^T$$

$$\frac{\partial L}{\partial b^{(0)}} = \frac{\partial L}{\partial S^{(0)}} \frac{\partial S^{(0)}}{\partial b^{(0)}} = \frac{\partial L}{\partial S^{(0)}} := \delta^{(0)}$$

$$\frac{\partial L}{\partial W^{(0)}} = \frac{\partial L}{\partial S^{(0)}} \frac{\partial S^{(0)}}{\partial W^{(0)}} = \frac{\partial L}{\partial S_1^{(0)}} \frac{\partial S_1^{(0)}}{\partial W^{(0)}} + \frac{\partial L}{\partial S_2^{(0)}} \frac{\partial S_2^{(0)}}{\partial W^{(0)}}$$

problem: this
is a 3-tensor



$$= \frac{\partial L}{\partial S^{(0)}} \left[\frac{\partial W^{(0)}}{\partial w_1} (w_1^0 x_1 + w_2^0 x_2) \right] + \frac{\partial L}{\partial S^{(0)}} \left[\frac{\partial}{\partial w_1} (w_1^0 x_1 + w_2^0 x_2) \right]$$

$$= \frac{\partial L}{\partial S^{(0)}} \begin{bmatrix} x_1 & x_2 \\ 0 & 0 \end{bmatrix} + \frac{\partial L}{\partial S^{(0)}} \begin{bmatrix} 0 & 0 \\ x_1 & x_2 \end{bmatrix}$$

$$= \left(\frac{\partial L}{\partial S^{(0)}} \begin{pmatrix} x_1 & x_2 \\ 0 & 0 \end{pmatrix} \right) \begin{pmatrix} x_1 & x_2 \\ 0 & 0 \end{pmatrix} = \cancel{\left(\frac{\partial L}{\partial S^{(0)}} \begin{pmatrix} x_1 & x_2 \\ 0 & 0 \end{pmatrix} \right)} = \cancel{\left(\frac{\partial L}{\partial S^{(0)}} \begin{pmatrix} x_1 & x_2 \\ 0 & 0 \end{pmatrix} \right)^T} = \cancel{\text{outer}}(\cancel{S^{(0)}}, \cancel{x}) = (S^{(0)})^T x = S^{(0)} x$$

$$\text{Still need } S^{(0)} \text{ and } S^{(1)}$$

$$S^{(0)} = \frac{\partial L}{\partial S^{(0)}} = \frac{\partial L}{\partial S^{(0)}} \frac{\partial S^{(1)}}{\partial S^{(0)}} = S^{(1)} W^{(1)} G^{(0)} = S^{(1)} \frac{\partial G^{(0)}}{\partial S^{(0)}}$$

$$S^{(1)} = \frac{\partial L}{\partial S^{(1)}} = \frac{\partial L}{\partial x} \frac{\partial N \bar{x}}{\partial S^{(1)}} = \frac{\partial L}{\partial x} \frac{\partial G^{(1)}}{\partial S^{(1)}} =$$

Backpropagation for $P(n_0, \dots, n_D)$

Cost function $J(W, B) = \frac{1}{m} \sum_{i=1}^m L(\vec{U}\vec{x}_i, \vec{y}_i)$

(each layer
fully connected to
adjacent layers)

$$\begin{aligned}\delta^{(D-1)} &= \frac{\partial L}{\partial s^{(D-1)}} = \frac{\partial L}{\partial U\vec{x}} \frac{\partial U\vec{x}}{\partial s^{(D-1)}} \\ &= \frac{\partial L}{\partial U\vec{x}} \frac{\partial G^{(D)}}{\partial s^{(D-1)}}\end{aligned}$$

$$U\vec{x} = G^{(D)}(s^{(D-1)})$$

and for $0 \leq l < D-1$

$$\begin{aligned}\delta^{(l)} &= \frac{\partial L}{\partial s^{(l)}} = \frac{\partial L}{\partial s^{(l+1)}} \frac{\partial s^{(l+1)}}{\partial s^{(l)}} \\ &= \delta^{(l+1)} W^{(l+1)} \frac{\partial G^{(l+1)}}{\partial s^{(l)}}\end{aligned}$$

$$s^{(l+1)} = W^{(l+1)} G^{(l+1)}(s^{(l)}) + b^{(l+1)}$$

$$\begin{aligned}\frac{\partial L}{\partial W^{(l)}} &= \frac{\partial L}{\partial s^{(l)}} \frac{\partial s^{(l)}}{\partial W^{(l)}} \\ &= \sum_{k=1}^{n_{l+1}} \frac{\partial L}{\partial s_k^{(l)}} \frac{\partial s_k^{(l)}}{\partial W^{(l)}} \\ &= \sum_{k=1}^{n_{l+1}} \delta_k^{(l)} \sum_{j=1}^{n_l} E_{kj} x_j^{(l)} \\ &= \sum_{k,j} \delta_k^{(l)} x_j^{(l)} E_{kj} \\ &= \text{outer}(\delta^{(l)}, x^{(l)})\end{aligned}$$

$$\begin{aligned}s^{(l)} &= W^{(l)} x^{(l)} + b^{(l)} \\ s_k^{(l)} &= W_{k+}^{(l)} x^{(l)} + b_k^{(l)} \\ &= \sum_{j=1}^{n_l} W_{kj}^{(l)} x_j^{(l)} + b_k^{(l)}\end{aligned}$$

Note $\frac{\partial}{\partial w_{ui}} W_{kj}^{(l)} = E_{kj}$

$$\frac{\partial L}{\partial b^{(l)}} = \frac{\partial L}{\partial s^{(l)}} \frac{\partial s^{(l)}}{\partial b^{(l)}} = \frac{\partial L}{\partial s^{(l)}} = \delta^{(l)}$$

matrix with
1 in (k, j)
0 elsewhere

Gradient for N with architecture $P(n_0, \dots, n_D)$

Dataset X, Y with $X \in \mathbb{R}^{m \times n}$ $Y \in \mathbb{R}^{m \times k}$

$$J(W, B) = \frac{1}{m} \sum_{i=1}^m L(N\vec{x}_i, \vec{y}_i)$$

$$\frac{\partial J}{\partial W^{(l)}} = \frac{1}{m} \sum_{i=1}^m \text{outer}(\delta_i^{(l)}, \vec{x}_i^{(l)})$$

$$\frac{\partial J}{\partial b^{(l)}} = \frac{1}{m} \sum_{i=1}^m \delta_i^{(l)}$$

where $\delta_i^{(D-1)} = \frac{\partial L}{\partial N\vec{x}_i} \frac{\partial G^{(D)}}{\partial S_i^{(D-1)}}$

$$\delta_i^{(l)} = \delta_i^{(l+1)} W^{(l+1)} \frac{\partial G^{(l+1)}}{\partial S_i^{(l)}} \quad 0 \leq l < D-1$$

Regularization (See HW 4)

A Stochastic Gradient Descent Step Pick $\alpha > 0$

Pick a batch $B_a \subseteq \{1, \dots, m\}$ and approximate

$$J(W, B) \approx \frac{1}{|B_a|} \sum_{i \in B_a} L(N\vec{x}_i, \vec{y}_i) \quad (\text{add to list of costs})$$

for $i \in B_a$:

$$\vec{x}_i$$

$$N\vec{x}_i$$

compute $[(x_i^{(0)}, s_i^{(0)}), (x_i^{(1)}, s_i^{(1)}), \dots, (x_i^{(D-1)}, s_i^{(D-1)}), (x_i^{(0)})]$

compute $\delta_i^{(D-1)}, \delta_i^{(D-2)}, \dots, \delta_i^{(0)}$

$$\text{compute } \frac{\partial J}{\partial W^{(l)}} \approx \frac{1}{|\mathcal{B}_l|} \sum_{i \in \mathcal{B}_l} \text{outer}(\delta_i^{(l)}, x_i^{(l)})$$

$$\frac{\partial J}{\partial b^{(l)}} \approx \frac{1}{|\mathcal{B}_l|} \sum_{i \in \mathcal{B}_l} \delta_i^{(l)}$$

for l in range(D):

$$W^{(l)} = W^{(l)} - \cancel{\alpha} \alpha \frac{\partial J}{\partial W^{(l)}}$$

$$b^{(l)} = b^{(l)} - \alpha \frac{\partial J}{\partial b^{(l)}}$$

Common Cost (loss) functions : $J(W, B) = \frac{1}{m} \sum_{i=1}^m L(\vec{N}\vec{x}_i, \vec{y}_i)$

Square Error Cost $L_{SE}(\vec{N}\vec{x}, \vec{y}) = \|\vec{N}\vec{x} - \vec{y}\|_2^2$ (MLE for gaussian Pmodel)
 Used for regression
 (mostly)

Cross Entropy Cost $L_{CE}(\vec{N}\vec{x}, \vec{y}) = -\vec{y}^T \text{Log}(\vec{N}\vec{x})$
 Used for multiclass
 classification with
 softmax in last layer

$$= - \sum_{i=1}^{n_D} y_i \log((\vec{N}\vec{x})_i)$$

Binary Cross Entropy Cost $L_{BCE}(\vec{N}\vec{x}, \vec{y}) = -y \log(\vec{N}\vec{x}) - (1-y) \log(1-\vec{N}\vec{x})$
 Used for binary classification

with sigmoid output neurons

(one output neuron $n_D=1$)

(MLE \Leftrightarrow for
 ? Pmodel
 Bernoulli)

How do these affect backprop?

They only affect $\delta^{(D-1)}$

$$\delta^{(D-1)} = \frac{\partial L}{\partial N\vec{x}} \frac{\partial N\vec{x}}{\partial S^{(D-1)}}$$

for SE: ~~$\frac{\partial L_{SE}}{\partial N\vec{x}}$~~ $\frac{\partial L_{SE}}{\partial N\vec{x}} = 2(N\vec{x} - \vec{y})^T$

for CE: $\frac{\partial L_{CE}}{\partial N\vec{x}} = \frac{\partial}{\partial N\vec{x}} (-\vec{y}^T \text{Log}(N\vec{x}))$

$$= -\vec{y}^T \begin{pmatrix} \frac{1}{(N\vec{x})_1}, & 0 \\ 0, & \ddots \\ & \frac{1}{(N\vec{x})_{n_D}} \end{pmatrix}$$

for BCE: $\frac{\partial L_{BCE}}{\partial N\vec{x}} = \frac{-y}{N\vec{x}} + \frac{1-y}{1-N\vec{x}}$

These lead to nice expressions for $\delta^{(D-1)}$ when you add specific output layer neurons (activations)

(see HW 4)