

Natural Language Processing

CSE 325/425



Sihong Xie

Lecture 12:

- Neural networks revisit

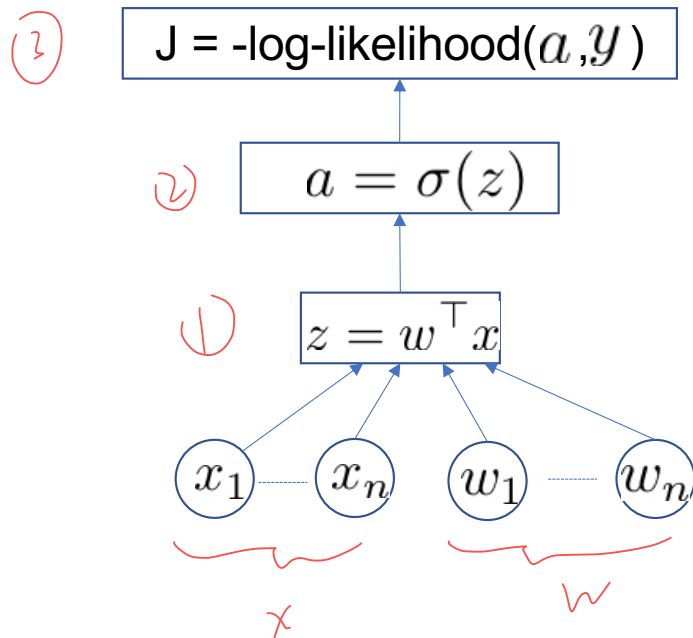
Revisiting neural networks

POS tagging goes neural!

- use neural networks to predict sequences:
 - word sequences: language models
 - POS tag sequences: POS-taggers
- upgrade HMM/MEMM to RNN
 - local normalization to distributions over all words/tags.
- upgrade CRF to CRF-LSTM
 - normalize to distribution over all sequences.

Logistic regression is a neural network

A computation graph is a differentiable system for **evaluation** and **differentiation**.



Forward pass:

- Compute the value of the hidden, output units, and the loss.

(1) $z = w^T x = \sum_{i=1}^n w_i x_i \in \mathbb{R}$

(2) $a = \sigma(z) = \frac{1}{1 + \exp(-z)} \in \mathbb{R}$

(3) $J = -y \log a - (1-y) \log (1-a) \in \mathbb{R}$

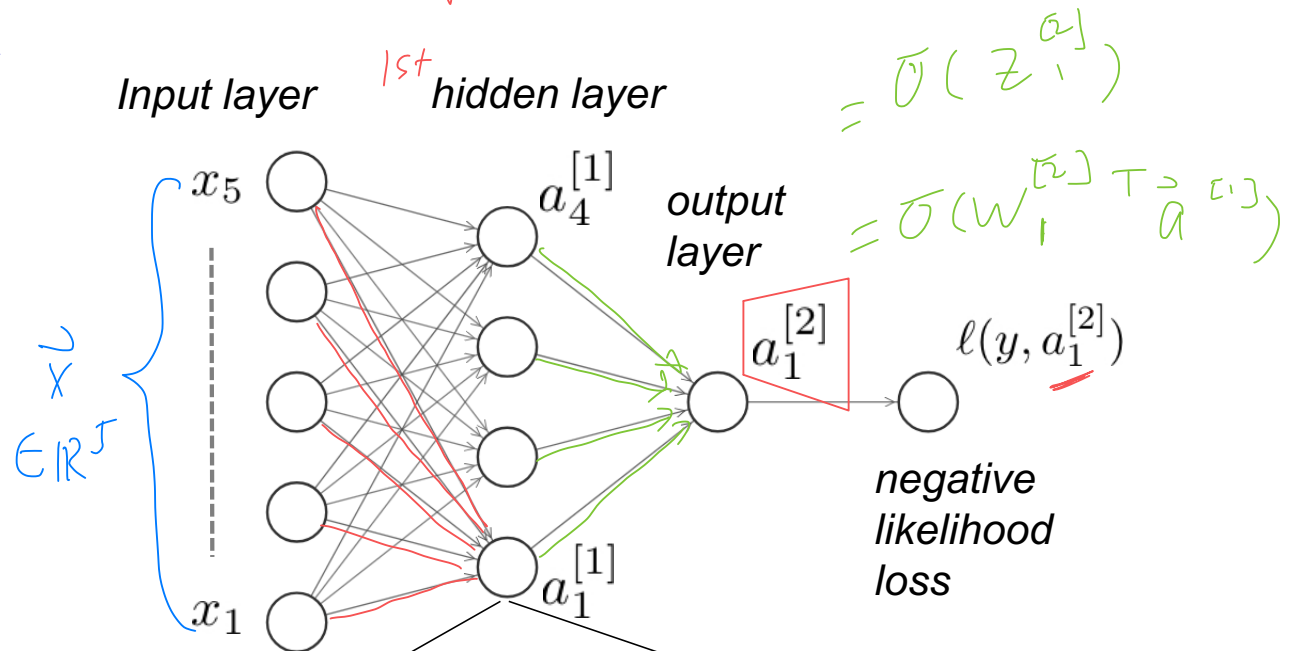
Back-propagation:

- compute the gradients using the chain rules.
- more later after we review matrix calculus.

Neural Network

Multiple Layer Perceptron (MLP)

A neural network is a computation graph that stacks many Logistic regression models.



$$W_1^{[1]} \in \mathbb{R}^5$$

One Logistic regression model

$$z_1^{[1]} = W_1^{[1]T} x \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

the vector of weights on edges connecting x to $z_1^{[1]}$

activation (σ is an activation function)

$$z_2^{[1]} = W_2^{[1]T} x$$

$$a_2^{[1]} = \sigma(z_2^{[1]})$$

tanh, ReLU

$$z^{[1]} = W^{[1]} x \in \mathbb{R}^{4 \times 1}$$

$$a^{[1]} = \sigma(z^{[1]}) = \begin{bmatrix} \sigma(z_1^{[1]}) \\ \vdots \\ \sigma(z_4^{[1]}) \end{bmatrix} \in \mathbb{R}^4$$

$$W^{[1]} = \begin{bmatrix} \text{---} W_1^{[1]T} \text{---} \\ \text{---} W_2^{[1]T} \text{---} \\ \text{---} W_3^{[1]T} \text{---} \\ \text{---} W_4^{[1]T} \text{---} \end{bmatrix} \in \mathbb{R}^{4 \times 5}$$

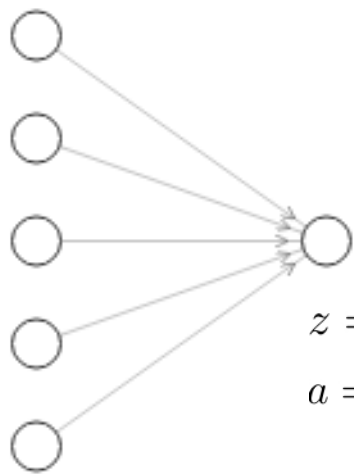
{ Vectorization
of W_N }

Forward propagation

In general, for the j -th neural on the first layer:

$$z_j^{[1]} = W_j^{[1]T} x + b_j^{[1]}$$

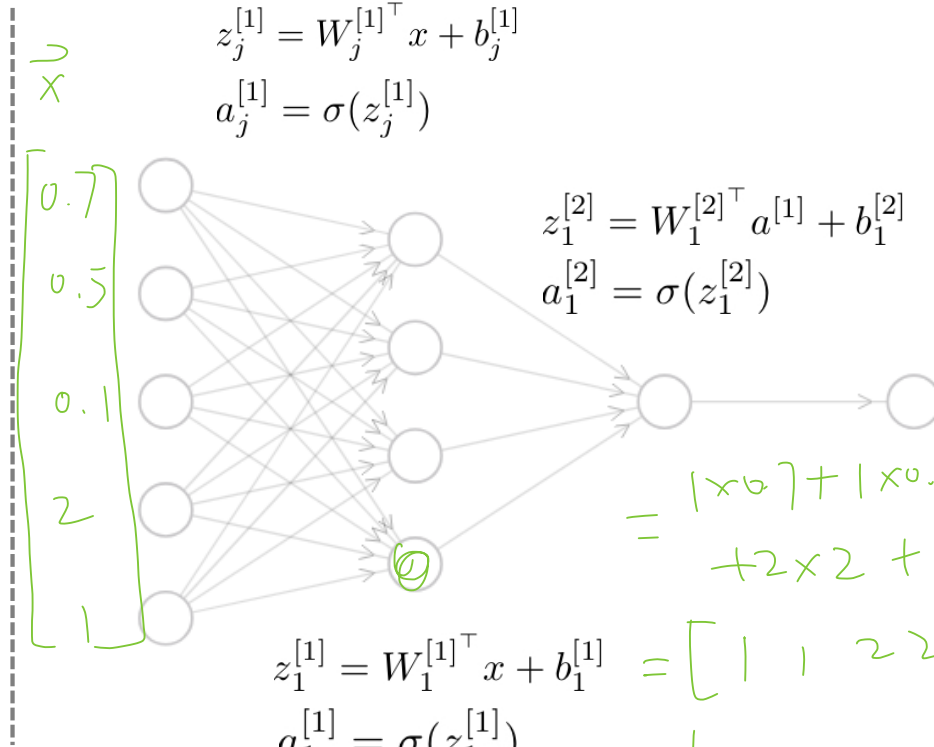
$$a_j^{[1]} = \sigma(z_j^{[1]})$$



$$z = w^T x + b$$

$$a = \sigma(z)$$

$\in \mathbb{R}$



$$z_1^{[2]} = W_1^{[2]T} a^{[1]} + b_1^{[2]}$$

$$a_1^{[2]} = \sigma(z_1^{[2]})$$

$$= 1 \times 0.7 + 1 \times 0.5 + 2 \times 0.1 + 2 \times 2 + 4 + 1 + 0.4 = 9.8$$

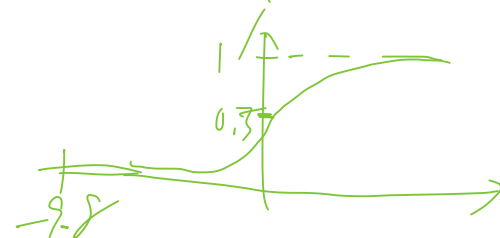
$$z_1^{[1]} = W_1^{[1]T} x + b_1^{[1]} = [1 \ 1 \ 2 \ 2 \ 4] \begin{bmatrix} 0.7 \\ 0.5 \\ 0.1 \\ 2 \\ 1 \end{bmatrix} + 0.4$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$= \frac{1}{1 + \exp(-9.8)} \approx 1$$

$$W_1^{[1]} = \begin{bmatrix} 1 \\ 1 \\ 2 \\ 2 \\ 4 \end{bmatrix}$$

$$b_1^{[1]} = 0.4$$



Matrix Calculus

- To train a neural network with parameters w (may include the biases), usually we minimize a scalar loss with respect to the parameters

$$\min_{w \in \mathbb{R}^d} L(w)$$

- Most optimization algorithms need the gradient of the loss with respect to the parameters w :

$$\frac{\partial L}{\partial w}$$

- Logistic regression optimization is a simple case:
 - can you recall what loss function we used?
 - and what's the gradient?

Matrix Calculus

- Basic definitions. Given a differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$

↗ set of all parameters
(Not the data space)

- differential $f(x+h) = f(x) + d_x f(h) + o(h)$ $d_x f : \mathbb{R}^n \rightarrow \mathbb{R}$

- gradient: explicit form of $d_x f$

$$\lim_{h \rightarrow 0} \frac{o(h)}{h} = 0$$

$$d_x f(h) = (\nabla_x f)^T h$$

- partial derivative

$$\nabla_x f : \mathbb{R}^n \rightarrow \mathbb{R}^n \quad \nabla_x f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

- Examples:

$$f([x_1, x_2]) = 3x_1 + x_2^2$$

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$\nabla_x f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 3 \\ 2x_2 \end{bmatrix}$$

$$\nabla_x f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

$$\nabla_x f|_{x=\begin{bmatrix} 0 \\ 1 \end{bmatrix}} = \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

Matrix Calculus

- Generalization of gradient to higher dimensional output space $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$

- Jacobian

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} [\partial f_1 / \partial \vec{x}]^T \\ \vdots \\ [\partial f_m / \partial \vec{x}]^T \end{bmatrix} : \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$$

- Examples:

$$g([y_1, \dots, y_{100}]) \quad n=3$$

$$g([y_1, y_2, y_3]) = \begin{bmatrix} y_1 + 2y_2 + 3y_3 \\ y_1 y_2 y_3 \end{bmatrix} \quad m=2$$

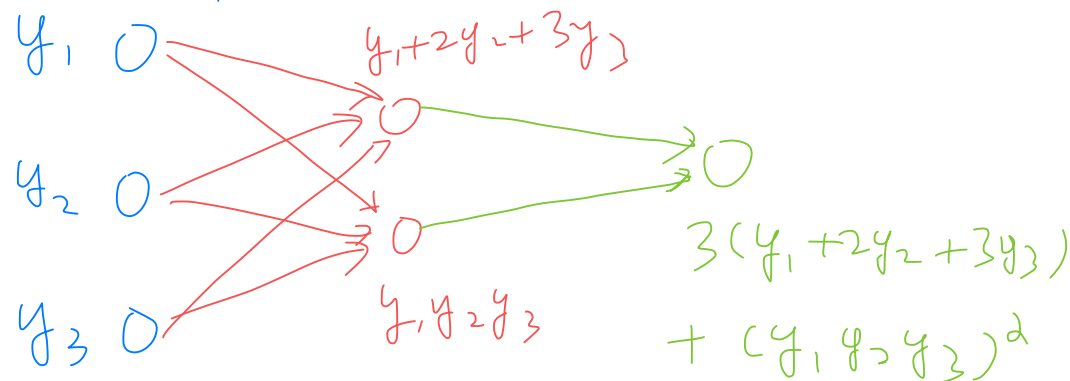
$$\frac{\partial g}{\partial y} = \begin{bmatrix} g_1(y_1, \dots, y_{100}) \\ g_2(y_1, \dots, y_{100}) \\ g_3(y_1, \dots, y_{100}) \end{bmatrix}$$

$$\frac{\partial g}{\partial y} \in \mathbb{R}^{3 \times 100}$$

$$\frac{\partial g}{\partial y} = \begin{bmatrix} 1 & 2 & 3 \\ y_2 y_3 & y_1 y_3 & y_1 y_2 \end{bmatrix} \in \mathbb{R}^{2 \times 3}$$

Matrix Calculus

Input Layer



$$= f \circ g(y_1, y_2, y_3)$$

- Chain rule for derivatives:

- Given two differentiable functions $g : \mathbb{R}^p \rightarrow \mathbb{R}^n$ $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
- Composition of two functions $f \circ g : \mathbb{R}^p \rightarrow \mathbb{R}^m$
- Jacobian of $f \circ g$ is the product of the Jacobians

$$d_x(f \circ g) = d_{g(x)}(f) \circ d_x(g)$$

$$J_{f \circ g} = J_f \circ J_g$$

- Examples:

output from g at x

$$g([y_1, y_2, y_3]) = \begin{bmatrix} y_1 + 2y_2 + 3y_3 \\ y_1y_2y_3 \end{bmatrix}$$

$$f([x_1, x_2]) = 3x_1 + x_2^2$$

$$\frac{\partial f}{\partial y} = \begin{bmatrix} 1 & 2 & 3 \\ y_2y_3 & y_1y_3 & y_1y_2 \end{bmatrix}$$

$$\frac{\partial f}{\partial \vec{x}} = \begin{bmatrix} 3 \\ 2x_2 \end{bmatrix}$$

$$d_y(f \circ g)|_{y=[0,1,2]}$$

$$= d_{g(y)}(f)|_{g(y)=g(0,1,2)} \circ d_y(g)|_{y=[0,1,2]} = [3, 0] \times \begin{bmatrix} 1 & 2 & 3 \\ 2 & 0 & 0 \end{bmatrix}$$

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

Matrix Calculus

- Common examples

- Vector inner product

$$\frac{\partial}{\partial w} (w^\top x) = x$$

- Matrix-vector (w.r.t. vector)

$$\frac{\partial}{\partial x} Wx = W \quad \text{(use def of Jacobian)}$$

- Vector-matrix (w.r.t. vector)

$$\frac{\partial}{\partial x} x^\top W = W^\top$$

- Matrix-vector (w.r.t. matrix)

$$\frac{\partial}{\partial W} Wx \quad \leftarrow \begin{array}{l} \text{of the shape of} \\ (\text{\#rows of } W \times \text{shape of } x) \end{array}$$

$J = f(z) \quad z = Wx$
 UV Params

- Usually don't directly find this Jacobian: why?
- Rather, embed this in chain-rule for some scalar function
- Example in the next two slides.

$$\frac{\partial J}{\partial W} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial W}$$

of the same shape of W

$$\frac{\partial}{\partial w} (w_1 x_1 + w_2 x_2 + \dots + w_n x_n) = \begin{bmatrix} \frac{\partial}{\partial w_1} (w_1 x_1 + \dots + w_n x_n) \\ \vdots \\ \frac{\partial}{\partial w_n} (w_1 x_1 + \dots + w_n x_n) \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$$\textcircled{1} \quad z = w^T x + b$$

Prob. Interpretation:

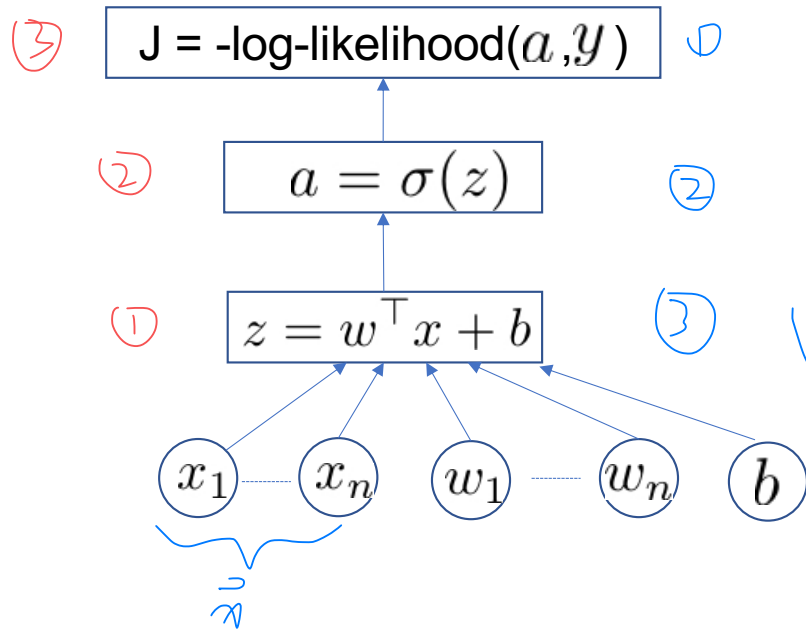
$$\textcircled{2} \quad a = \sigma(z) = \frac{1}{1 + \exp(-z)} = \Pr(y=1|x)$$

forward
prop.

$$\textcircled{3} \quad J = -y \log a - (1-y) \log (1-a)$$

Backpropagation for binary Logistic Regression

going Backwards



$$\textcircled{1} \quad \frac{\partial J}{\partial a} = -y \frac{1}{a} + (1-y) \frac{1}{1-a}$$

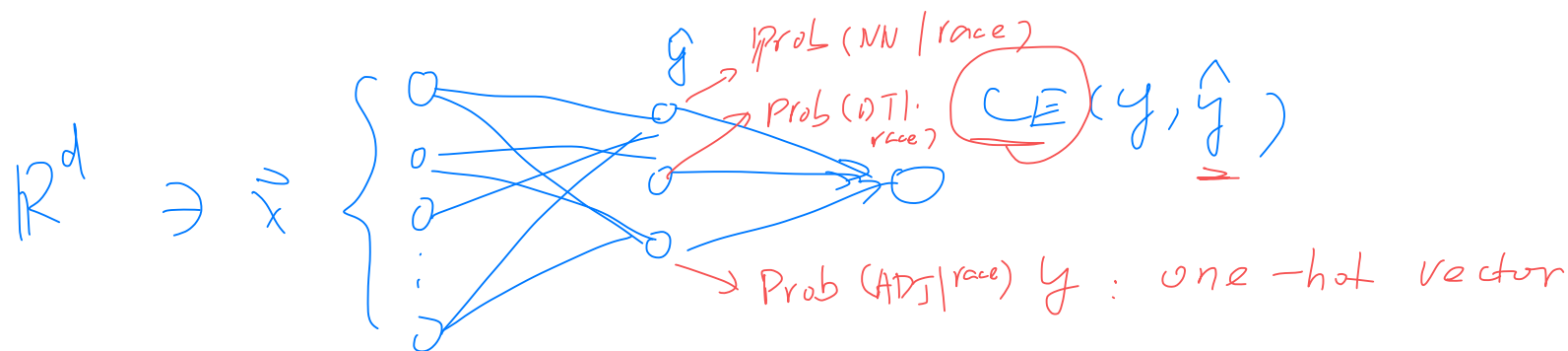
$$\begin{aligned} \textcircled{2} \quad \frac{\partial J}{\partial z} &= \frac{\partial J}{\partial a} \frac{\partial a}{\partial z} \quad (\text{chain rule}) \\ &= \left(-\frac{y}{a} + \frac{1-y}{1-a} \right) \times (a) \times (1-a) \\ &= -y(1-a) + a(1-y) \\ &= a - y \end{aligned}$$

gradient descent:

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \times \frac{\partial J}{\partial w} = w^{(t)} + \eta (y \vec{x} - a \vec{x})$$

$$\begin{aligned} \textcircled{3} \quad \frac{\partial J}{\partial w} &= \frac{\partial J}{\partial z} \times \frac{\partial z}{\partial w} \\ &= (a - y) \times \vec{x} \end{aligned}$$

MEMM



Backpropagation for multi-class Logistic Regression

- Example: predicting more than two classes.

$$y \in \{\text{NN}, \text{DT}, \text{ADJ}\}$$

- Cross-entropy loss (negative log-likelihood loss)
- Used in logistic regression and neural networks for classification.

$$z = Wx = \begin{bmatrix} -w_{\text{NN}} \\ -w_{\text{DT}} \\ -w_{\text{ADJ}} \end{bmatrix} \quad \text{--- "race" ---}$$

$$\hat{y} = \text{softmax}(z)$$

$$\hat{y}_i = \frac{\exp(z_i)}{\sum_{j=1}^3 \exp(z_j)}$$

$$\text{CE}(y, \hat{y}) = - \sum_{i=1}^3 y_i \log \hat{y}_i$$

$\hat{y}_i \triangleq$ Prob of class i predicted by the model

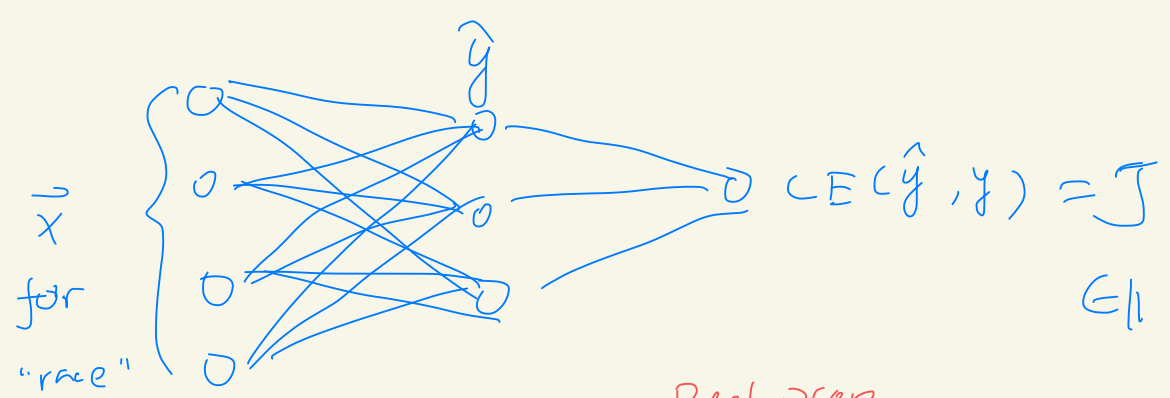
$$y = \begin{bmatrix} 1, 0, 0 \\ 0, 1, 0 \\ 0, 0, 1 \end{bmatrix}$$

$$= -1 \times \log \hat{y}_1$$

$$- 0 \times \log \hat{y}_2$$

$$- 0 \times \log \hat{y}_3$$

$$= -\log \hat{y}_1 \triangleq \text{negative log-likelihood of NN give "race"}$$



Back prop.

$$\textcircled{1} z = Wx$$

$$\textcircled{2} \hat{y} = \text{softmax}(z)$$

$$\textcircled{3} J = CE(\hat{y}, y)$$

$$\textcircled{1} \frac{\partial J}{\partial \hat{y}}$$

$$\textcircled{2} \frac{\partial J}{\partial z} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z}$$

$$\textcircled{3} \frac{\partial J}{\partial W} = \left(\frac{\partial J}{\partial z} \right) \left(\frac{\partial z}{\partial W} \right)^T = \hat{y} - y$$

(Swapped due to matrix shape compatibility.)

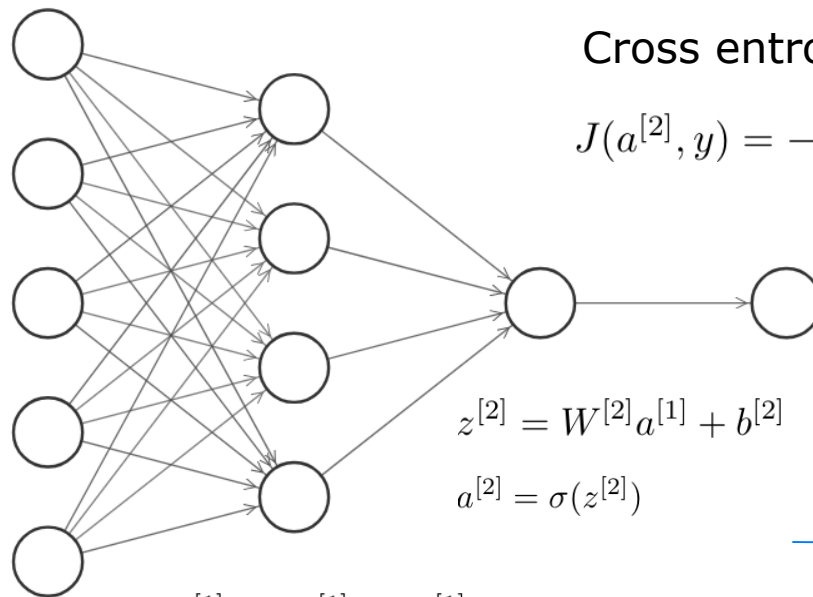
predicted prob vector

one hot vector

$$= (\hat{y} - y) x^T$$

$$\begin{bmatrix} -w_1 \\ -w_2 \\ -w_3 \end{bmatrix} \quad \text{(outer product)}$$

Backpropagation for Neural Network



Cross entropy loss for binary classification:

$$J(a^{[2]}, y) = -y \log a^{[2]} - (1 - y) \log(1 - a^{[2]})$$

Based on the gradient of Logistic

Regression, $\frac{\partial J}{\partial z^{[2]}} = a^{[2]} - y$ (error term)

To train $w^{[2]}$ & $b^{[2]}$, we need

$$\frac{\partial J}{\partial w^{[2]}} = \frac{\partial J}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial w^{[2]}} \leftarrow \text{too large to evaluate}$$

$$= (a^{[2]} - y) (a^{[1]})^T \leftarrow \text{this is correct if verified each element of the outer product.}$$

$$\frac{\partial J}{\partial b^{[2]}} = \frac{\partial J}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial b^{[2]}}$$

$$= (a^{[2]} - y) \leftarrow z^{[2]} = w^{[2]} a^{[1]} + (I) b^{[2]} \text{ (Identity matrix)}$$

Input vector x
 $z^{[1]} = W^{[1]}x + b^{[1]}$
 $a^{[1]} = \sigma(z^{[1]})$

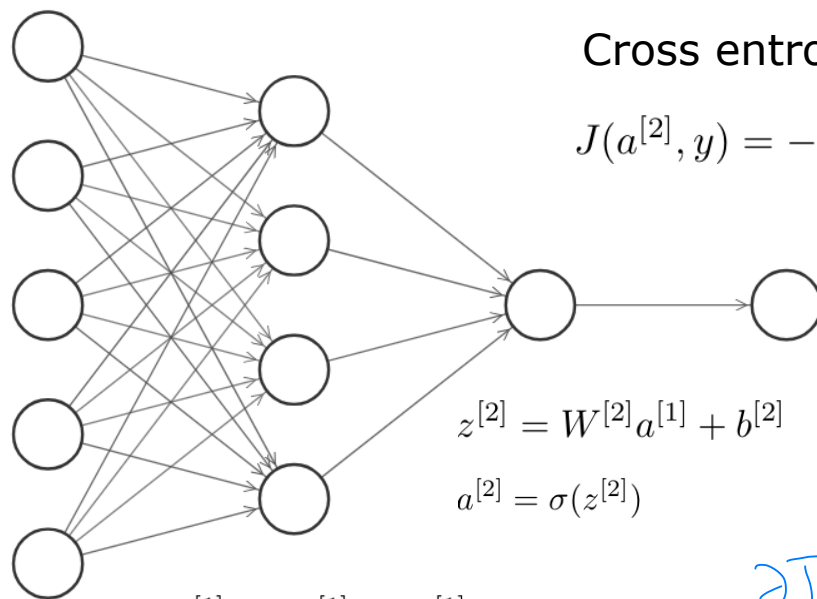
Get the dimensions right!

$$\frac{\partial J}{\partial w^{[1]}} = \frac{\partial J}{\partial z^{[1]}} x^T$$

($z^{[1]}$ has the same number of rows as $w^{[1]}$, & x^T has the same number of columns of $w^{[1]}$)

$$\frac{\partial J}{\partial b^{[1]}} = \frac{\partial J}{\partial z^{[1]}}$$

Backpropagation for Neural Network



Cross entropy loss for binary classification:

$$J(a^{[2]}, y) = -y \log a^{[2]} - (1 - y) \log(1 - a^{[2]})$$

Need $\frac{\partial J}{\partial z^{[1]}}$ since that will help finding $\frac{\partial J}{\partial w^{[1]}}$ & $\frac{\partial J}{\partial b^{[1]}}$ for training

$w^{[1]}$ & $b^{[1]}$

$$\frac{\partial J}{\partial a^{[1]}} = \left(\frac{\partial z^{[2]}}{\partial a^{[1]}} \right)^T \frac{\partial J}{\partial z^{[2]}} = (w^{[2]})^T \frac{\partial J}{\partial z^{[2]}} = (w^{[2]})^T (a^{[2]} - y)$$

Input vector x
 $z^{[1]} = W^{[1]}x + b^{[1]}$
 $a^{[1]} = \sigma(z^{[1]})$

Get the dimensions right!

$$a^{[1]} = \sigma(z^{[1]}) = [\sigma(z_1^{[1]}), \dots, \sigma(z_4^{[1]})]$$

$$\frac{\partial a^{[1]}}{\partial z^{[1]}} = \text{diag} \left(\left[\frac{\partial a_1^{[1]}}{\partial z_1^{[1]}}, \dots, \frac{\partial a_4^{[1]}}{\partial z_4^{[1]}} \right] \right)$$

$$\frac{\partial J}{\partial z^{[1]}} = \frac{\partial J}{\partial a^{[1]}} \times \frac{\partial a^{[1]}}{\partial z^{[1]}}$$

Continue at the top

Jacobian of the vector function $a^{[1]}$ with input $z^{[1]}$

$$W^{[2]} = \begin{bmatrix} w_{11}^{[2]} & w_{12}^{[2]} & \boxed{w_{13}^{[2]}} & w_{14}^{[2]} \\ w_{21}^{[2]} & \dots & \boxed{w_{23}^{[2]}} & \vdots \\ w_{31}^{[2]} & \dots & \boxed{w_{33}^{[2]}} & 1 \end{bmatrix}$$

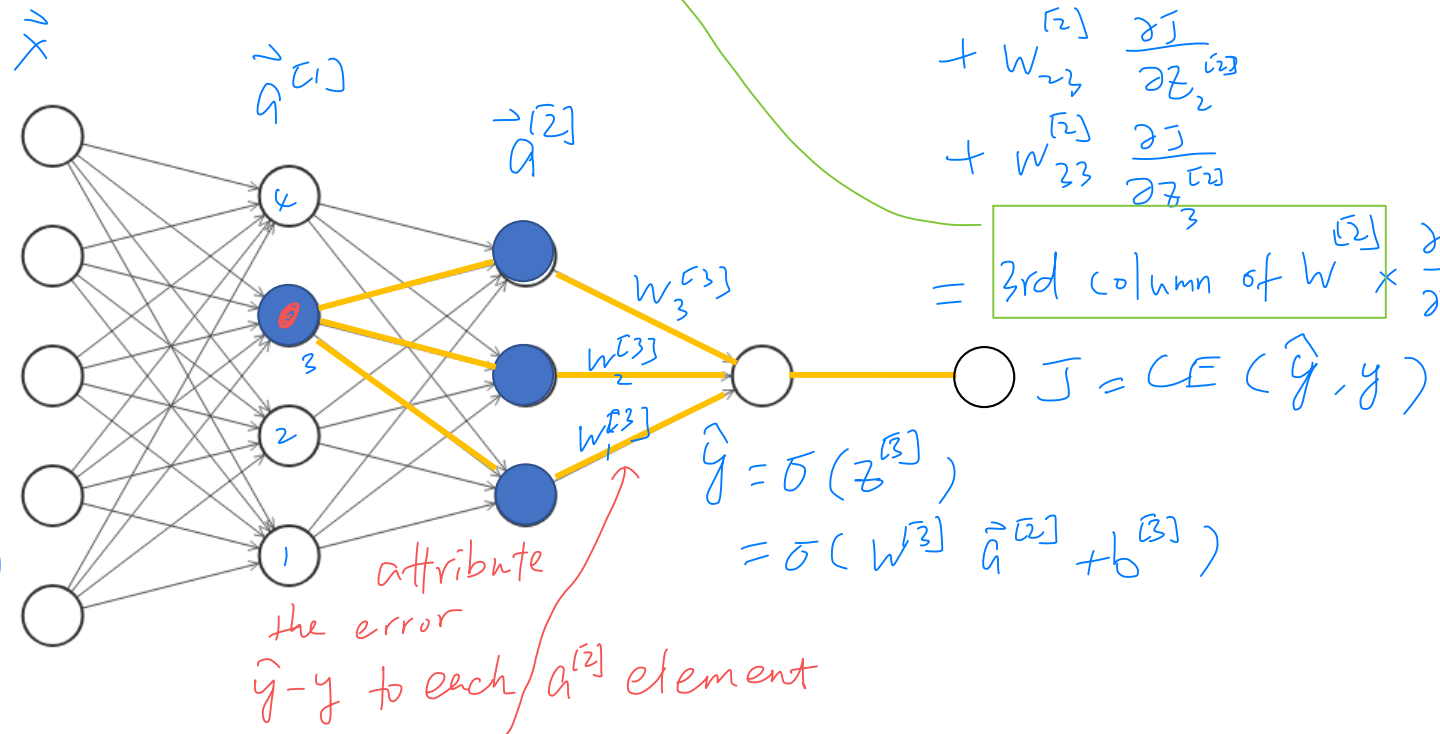
Since $a_3^{[1]} \in \mathbb{R}$ contributes to J through $z_1^{[2]}, z_2^{[2]},$ & $z_3^{[2]}$, with weights $w_{13}^{[2]}, w_{23}^{[2]}$ & $w_{33}^{[2]}$, so

Backpropagation intuition

How much J will change if $a_3^{[1]}$ is slightly perturbed to $a_3^{[1]} + h$, with $h \rightarrow 0$

$$\frac{\partial J}{\partial a_3^{[1]}} = \lim_{h \rightarrow 0} \frac{J - J(a_3^{[1]} + h)}{h}$$

$$\begin{aligned} \frac{\partial J}{\partial a_3^{[1]}} &= w_{13}^{[2]} \frac{\partial J}{\partial z_1^{[2]}} + w_{23}^{[2]} \frac{\partial J}{\partial z_2^{[2]}} + w_{33}^{[2]} \frac{\partial J}{\partial z_3^{[2]}} \\ &= \boxed{\text{3rd column of } W^{[2]}} \times \frac{\partial J}{\partial z^{[2]}} \end{aligned}$$



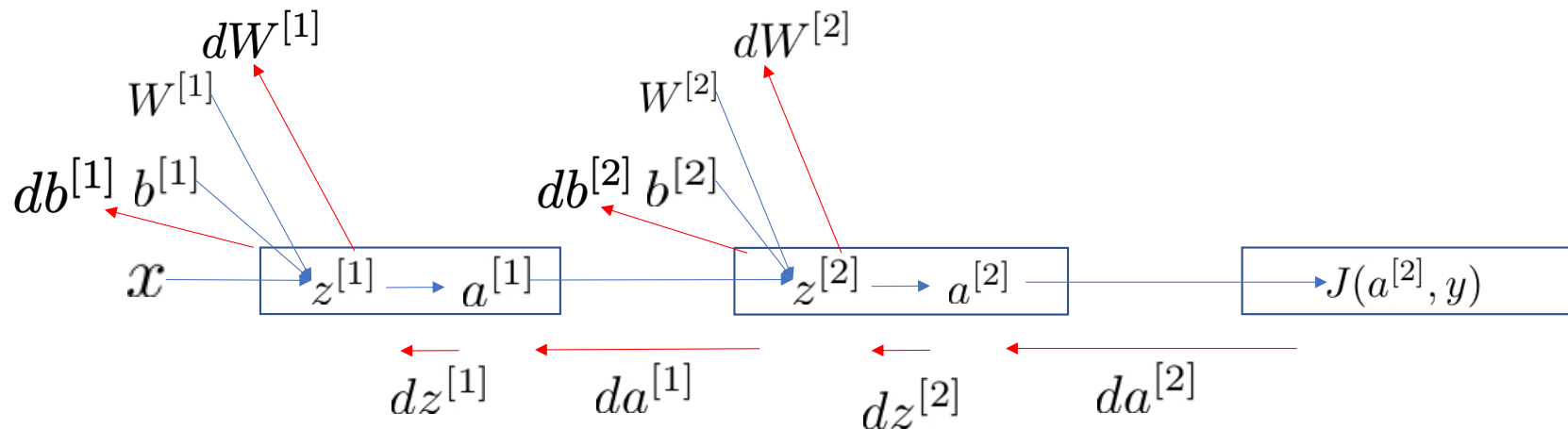
Backprop: $\frac{\partial J}{\partial z^{[3]}} = \hat{y} - y$, $\frac{\partial J}{\partial a^{[2]}} = (W^{[3]})^T \frac{\partial J}{\partial z^{[3]}}$

$$\begin{aligned} &= (W^{[3]})^T (\hat{y} - y) \\ &= \begin{bmatrix} w_1^{[3]} \times (\hat{y} - y) \\ \vdots \\ w_3^{[3]} \times (\hat{y} - y) \end{bmatrix} \end{aligned}$$

$$\frac{\partial J}{\partial z^{[2]}} = (W^{[3]})^T (\hat{y} - y) \times \text{diag}\left(\frac{\partial a^{[2]}}{\partial z^{[2]}}\right)$$

Backpropagation for Neural Network

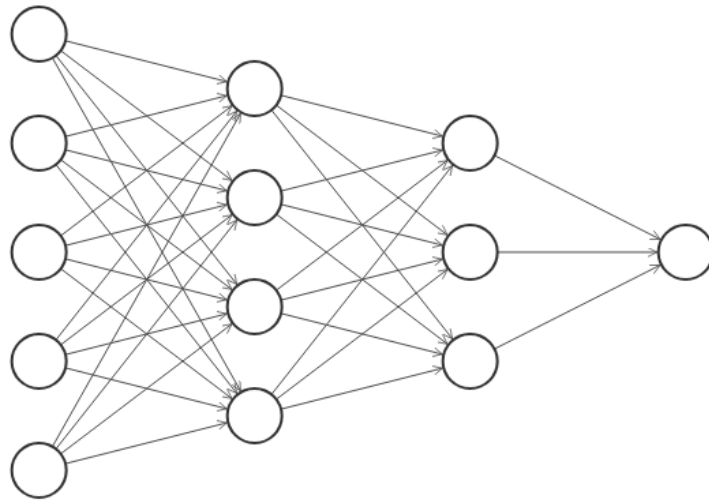
The big picture



Get the dimensions right!

Gradient descent for network training

1. Initialize network weights randomly.
2. Compute gradients through backpropagation.
3. Update weights (learning rate?).
4. Repeat 2 – 3 until convergence.



$$w^{(t+1)} \leftarrow w^{(t)} - \eta \left. \frac{\partial J}{\partial w} \right|_{w = w^{(t)}}$$
$$w^{(t+2)} \leftarrow w^{(t+1)} - \eta \left. \frac{\partial J}{\partial w} \right|_{w = w^{(t+1)}}$$

$$\vdots$$
$$w^{(*)} \leftarrow w^{(*)} - \eta \left. \frac{\partial J}{\partial w} \right|_{w = w^{*}}$$
$$= w^{(*)}$$

(Note: The boxed equation above is circled in red in the original image, with a red '0' written below it.)