

Week 1

Friday, August 25, 2017 8:47 AM

Specialization Outline

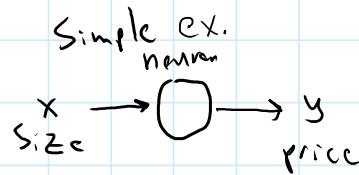
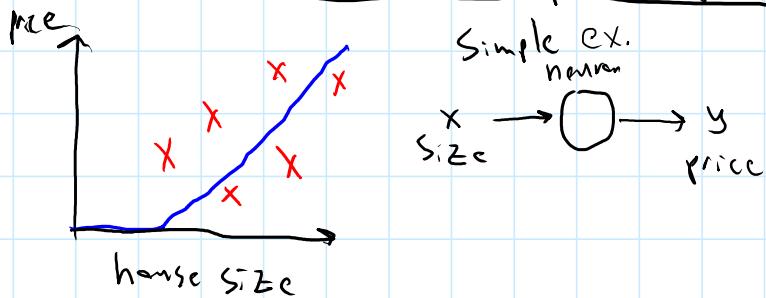
Courses

- 1) Neural Networks & Deep Learning (This one)
- 2) Improving Deep NNs: Hyperparameter Tuning, Regularization, Optimization
- 3) Structuring an ML project
- 4) Convolutional NNs
- 5) Natural Language Processing: Building Sequence Models

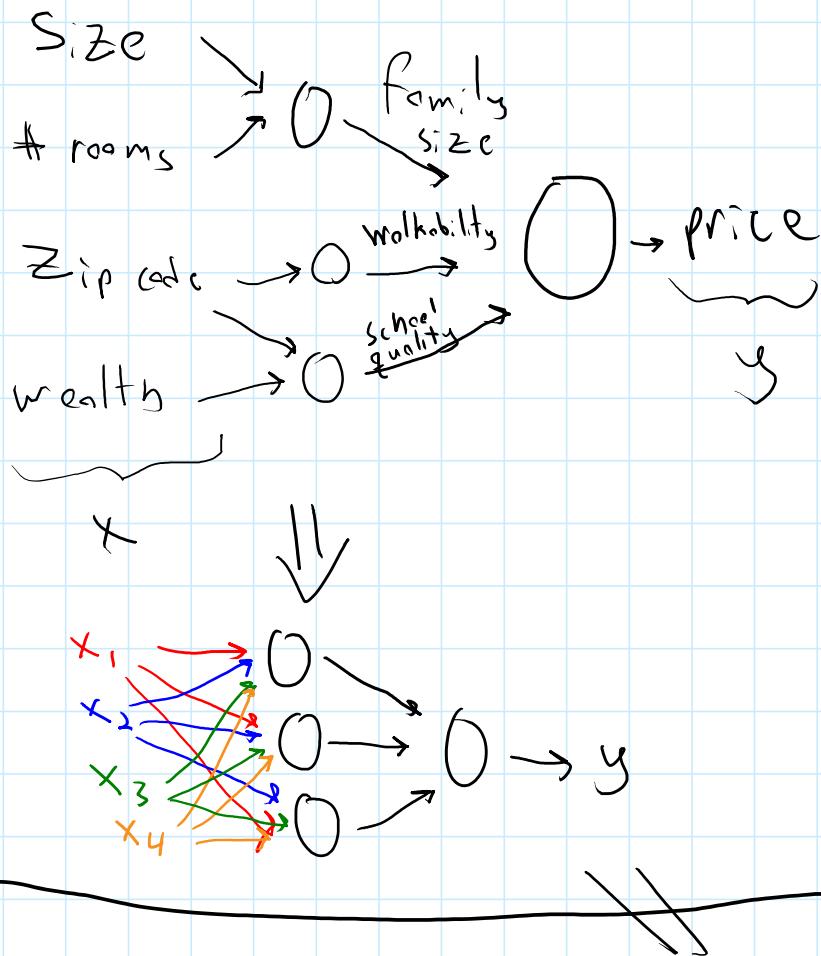
Intro to Deep Learning

Friday, August 25, 2017 9:02 AM

What is a Neural Network?



Rectified
Linear
Unit



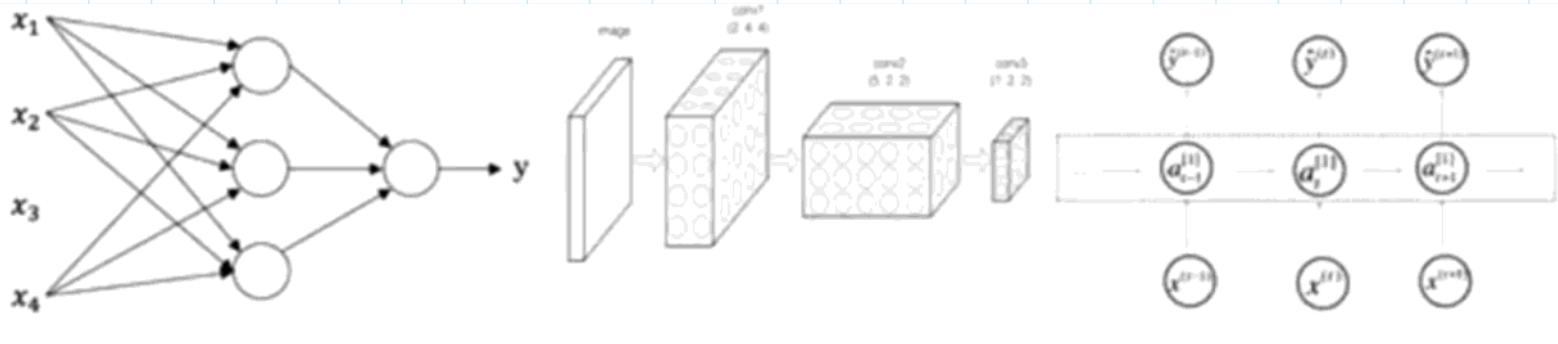
Supervised Learning w/ Neural Networks

Examples

| Input (x) | Output (y) | Application |
|--|----------------------------------|-------------------------|
| Home features $x = [x_1, \dots, x_n]$ | Price $y = [y_1, \dots, y_m]$ | Real estate {Sd. NN} |

| | | |
|---------------|-----------------------|---------------------|
| Home features | Price | Real estate |
| Ad, user info | Click? (0,1) | Ads |
| Image | Object (1, -1000) | Photo tagging |
| Audio | Text | Speech recognition |
| English | Chinese | Machine Translation |
| Image, Radar | Position & other cars | Autonomous Driving |

Std. NN CNN RNN Custom hybrid



Standard NN

Convolutional NN

Recurrent NN

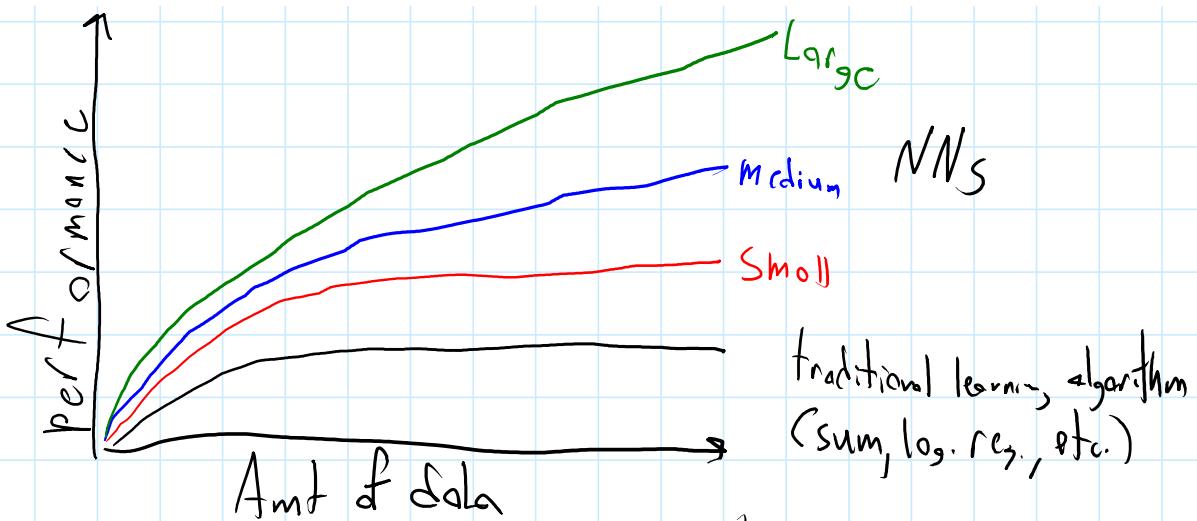
Structured Data

Tables

Unstructured Data

Audio Images
Text

Why have NNs become so popular



Employing RCLU instead of Sigmoids in processing,
has sped up training

Week 2

Monday, August 28, 2017 9:41 AM

Basics of NN programming

Vectorization
Forward & back propagation

Logistic Regression as a Neural Network

Monday, August 28, 2017 9:42 AM

Binary Classification

Image \rightarrow (0 or 1)

An $n \times m$ picture has 3 (R , G , B) $n \times m$ matrices
unroll into 1 feature vector x

$$\begin{bmatrix} x_R(:) \\ x_G(:) \\ x_B(:) \end{bmatrix}$$

$$x \in \mathbb{R}^{1 \times (n \cdot m \cdot 3)}$$

Notation

$$(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$$

m training examples : $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ python notation

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \quad X \in \mathbb{R}^{n_x \times m} \quad X.shape = (n_x, m)$$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \quad Y \in \mathbb{R}^{1 \times m} \quad Y.shape = (1, m)$$



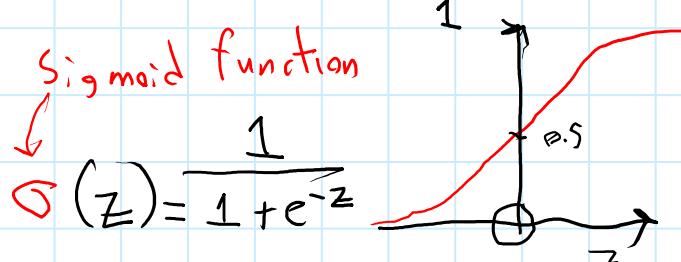
Logistic Regression

$$\text{Given } x \quad \text{want } \hat{y} \stackrel{\text{prediction}}{=} P(y=1|x) \quad x \in \mathbb{R}^{n_x}$$

Given X , want $\hat{y} \xleftarrow{\text{prediction}} P(y=1 | X) \quad x \in \mathbb{R}^{n_x}$

Parameters: $w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$

Output: $\hat{y} = \sigma(w^T x + b)$ where $\sigma(z) = \frac{1}{1+e^{-z}}$



Alt Notation: (used in Andrew Ng's ML course, but NOT here)

$$x_0 = 1, \quad x \in \mathbb{R}^{n_x+1}, \quad \hat{y} = \sigma(\Theta^T x), \quad \Theta = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \vdots \\ \Theta_{n_x} \end{bmatrix} \quad \left\{ \begin{array}{l} \Theta_0 \\ \Theta_1 \\ \vdots \\ \Theta_{n_x} \end{array} \right\} w$$



Logistic Regression Cost Function

Want $\hat{y}^{(i)} \approx y^{(i)}$ ($\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$)

Loss (error function): $L(\hat{y}, y)$

Sq. error $[\frac{1}{2}(y - \hat{y})^2]$ struggles to find global minima w/ gradient descent
when used w/ logistic regression

$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

Intuition: $y=1: L(\hat{y}, y) = -\log(\hat{y}) \leftarrow \text{want } \log(\hat{y}) \text{ large, want } \hat{y} \text{ large}$

$y=0: L(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow \text{want } \log(1-\hat{y}) \text{ large, want } \hat{y} \text{ small}$

Cost function (loss for whole training set):

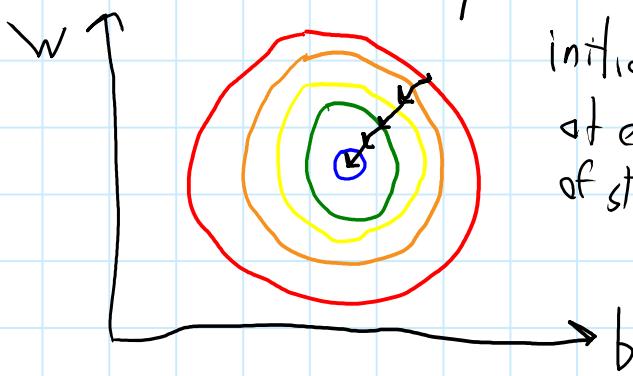
Cost function (loss for whole training set):

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(y^{(i)}, \hat{y}^{(i)})$$

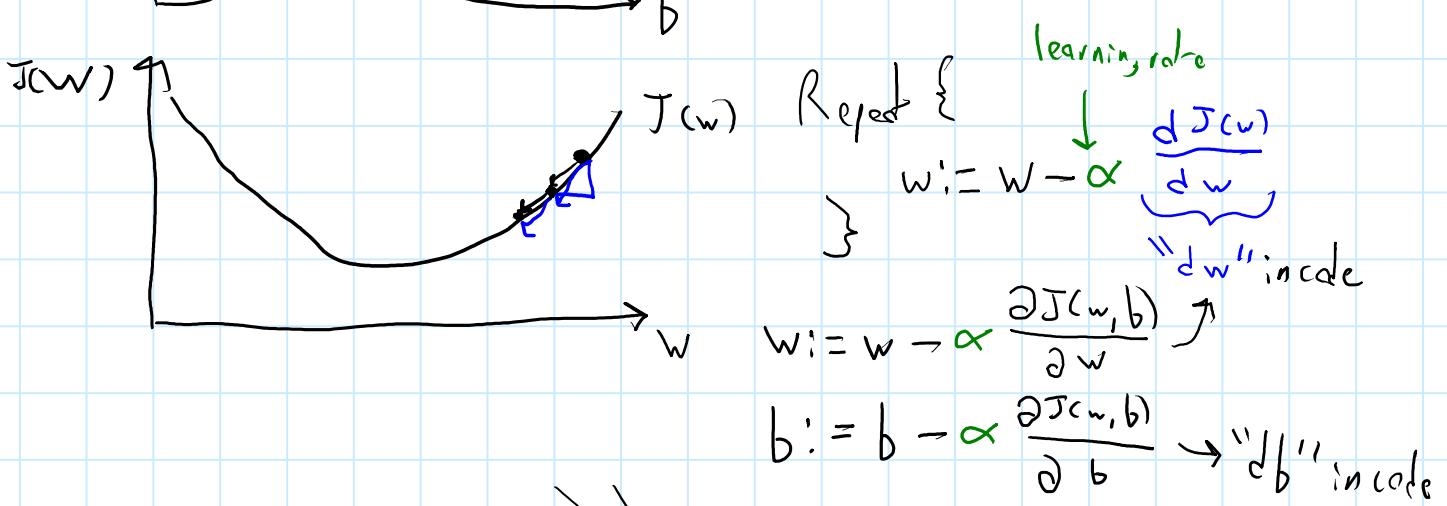
$$= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$$

Gradient Descent (learning w+b parameters)

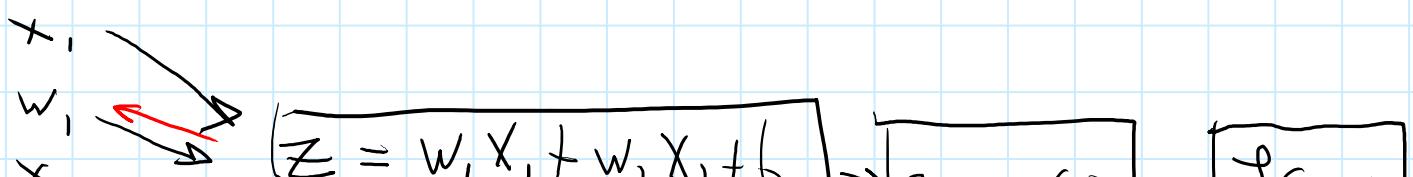
Want to find w, b that minimize $J(w, b)$

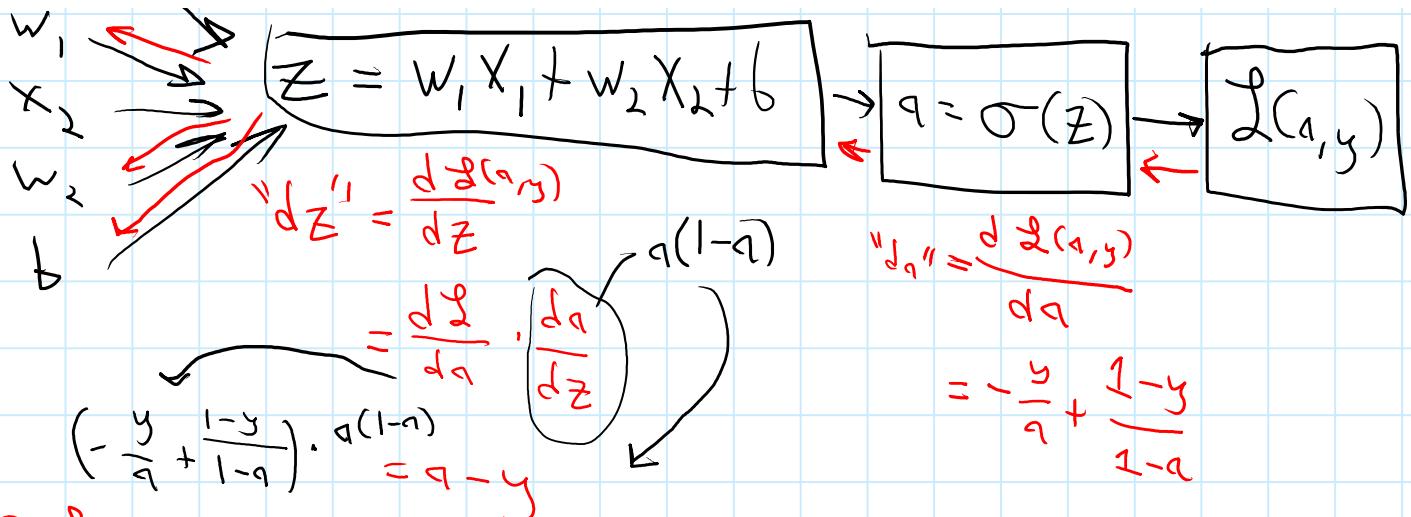


initialize @ 0 or random value,
at each step, move in direction
of steepest descent



Logistic Regression & Gradient Descent





$$\frac{\partial \mathcal{L}}{\partial w_n} = "d_{w_n}" = x_n d_Z ; \quad d_b = d_Z$$

Gradient Descent over a Training Set

$$\frac{\partial J}{\partial w_i} = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_i} \underbrace{\mathcal{L}(a^{(i)}, y^{(i)})}_{d w_i^{(i)} - (x^{(i)}, y^{(i)})}$$

Initialize $J, d_{w_1}, d_{w_2}, db = 0$

for i in range($1, m+1$):

$$Z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(Z^{(i)})$$

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

$$d_Z^{(i)} = a^{(i)} - y^{(i)}$$

$$d_{w_1} += x_1^{(i)} d_Z^{(i)}$$

$$d_{w_2} += x_2^{(i)} d_Z^{(i)}$$

$$db += d_Z^{(i)}$$

\uparrow can do as for loop,
 \downarrow but inefficient, use vectorization

$$J /= m$$

$$\delta w_1 / = m$$

$$\delta w_2 / = m$$

$$\delta b / = m$$

Python and Vectorization

Tuesday, August 29, 2017 9:02 AM

Vectorization

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

$$w = [\cdot] \quad x = [\cdot] \quad w \in \mathbb{R}^{n_x} \quad x \in \mathbb{R}^{n_x}$$

Non-vectorized

$$Z = 0$$

for i in range(n_x):

$$Z += w[i] * x[i]$$

$$Z += b$$

~ 450 ms

vectorization $\sim 300 \times$ faster !! Use over for loops whenever poss.

Exponentiation

$$v = [v_1, \dots, v_n] \rightarrow u = [e^{v_1}, \dots, e^{v_n}]$$

non-vectorized

$$u = np.zeros((n, 1))$$

for i in range(n):

$$u[i] = \text{math.exp}(v[i])$$

Vectorized

$$z = \underbrace{\text{np.dot}(w, x)}_{w^T x} + b$$

iterations over 10^7 values

~ 1.5 ms

vectorized

$$u = \text{np.exp}(v)$$

Revisiting the Logistic Regression for loop

$$d_w = \text{np.zeros}((n_x, 1))$$

~~$$J = 0, d_w[1] = 0, d_w[2] = 0, d_b = 0$$~~

For i in range(1, m+1):

$$Z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(Z^{(i)})$$

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

$$dZ^{(i)} \approx a^{(i)} - y^{(i)}$$

$$\frac{dw_1}{dw} += x_1^{(i)} dZ^{(i)}$$

$$\frac{dw_2}{dw} += x_2^{(i)} dZ^{(i)}$$

$$db += dZ^{(i)}$$

$$J /= m, \quad dw1 /= m, \quad dw2 /= m, \quad db /= m$$

$dw /= m$

Vectorizing Logistic Regression

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} \quad X.\text{shape} = (n_x, m)$$

$$w^T X + [b \dots b] = [z^{(1)} z^{(2)} \dots z^{(m)}]$$

$$[w^T x^{(1)} + b \quad w^T x^{(2)} + b \quad \dots \quad w^T x^{(m)} + b] \Rightarrow \boxed{Z = np.\text{dot}(w.T, X) + b}$$

real number
"broadcast" into $1 \times n$ matrix

$$A = \sigma(Z)$$

Vectorizing Logistic Regression, Gradient Descent

$$Z = [dZ^{(1)} \dots dZ^{(m)}], \quad A = [a^{(1)} \dots a^{(m)}], \quad Y = [y^{(1)} \dots y^{(m)}]$$

$$dZ = A - Y$$

$$db = \frac{1}{m} \sum_{i=1}^m dZ^{(i)} \rightarrow \left(\frac{1}{m}\right) \text{np.sum}(dZ)$$

$$dw = \frac{1}{m} \sum_{i=1}^m X^{(i)} dZ^{(i)} \rightarrow \left(\frac{1}{m}\right) X dZ^T \rightarrow \left(\frac{1}{m}\right) \text{np.dot}(X, dZ.T)$$

Implementation of code w/o for loops in process.

for iter in range(iterations):

$$\bar{Z} = \text{np.dot}(w.T, X) + b$$

$$\bar{A} = \sigma(\bar{Z})$$

$$dZ = \bar{A} - Y$$

$$dw = \left(\frac{1}{m}\right) \text{np.dot}(X, dZ.T)$$

$$db = \left(\frac{1}{m}\right) \text{np.sum}(dZ)$$

$$w -= \alpha dw$$

$$b -= \alpha db$$

Broadcasting Example

calories dist in 100g of diff foods

$$\begin{matrix} & \text{Apples} & \text{Beef} & \text{Eggs} & \text{Potatoes} \\ \text{Carbs} & 56.0 & 0.0 & 4.4 & 68.0 \\ \text{Protein} & 1.2 & 104.0 & 52.0 & 8.0 \\ \text{Fat} & 1.8 & 135.0 & 99.0 & 0.9 \end{matrix} = A$$

(Set % of calories of each nutrient class for each food w)

\tilde{n}_9 for loops

$$csl = A \cdot \text{sum}(x_1, 3=0) \quad \# \text{ sum of colones per food}$$

printL(col) [59. 231. 155. 4 76. 9]

$$\text{percenta,c} = 100 * A / \text{col.reshape}(1, 4)$$

printL(percenta,c)

| | | | |
|----|--------|--------|----|
| | ✓(3,4) | ✓(1,4) | |
| 94 | 2 | 2 | 88 |
| 2 | 43 | 33 | 10 |
| 3 | 56 | 63 | 1 |

redundant here,
but not costly

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

broadcasting,

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} \stackrel{(1,n) \rightsquigarrow (m,n)}{\sim} \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

General Principle

$$(m, n) \stackrel{+}{=} (1, n) \rightsquigarrow (m, n)$$

$$(m, n) \stackrel{\times}{=} (m, 1) \rightsquigarrow (m, n)$$

$$(m, n) \stackrel{t}{\times} \mathbb{R} \rightsquigarrow \mathbb{R} \Rightarrow \text{np.ones}(m, n)$$

$$(m, n) \xrightarrow{f} \mathbb{R} \rightsquigarrow \mathbb{R} \otimes \text{np.ones}(m, n)$$

Note: when implementing arrays, rank 1 ($a.shape = (n,)$) arrays behave in strange ways, use the reshape command to fit calculations e.g. $a = \text{random.randn}(5) \rightarrow \text{array}([a_1, a_2, a_3, a_4, a_5])$

$$a.reshape(5, 1) = \begin{bmatrix} a_1 \\ \vdots \\ a_5 \end{bmatrix} \quad a.reshape(1, 5) = [a_1, \dots, a_5]$$



Homework Notes

Thursday, August 31, 2017 9:59 AM

4.3 - Forward + Back prop

$$A = \sigma(w^T X + b) \rightarrow \text{Sigmoid(np.dot}(w, T, X) + b)$$

$$J = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(a^{(i)}) + (1 - y^{(i)}) \log(1 - a^{(i)})$$

$$\left(-(\text{np.sum}(Y * \text{np.log}(A) + (1 - Y) * \text{np.log}(1 - A)) / m \right)$$

$$\frac{\partial J}{\partial w} = \frac{1}{m} X (A - Y)^T \rightarrow \text{np.dot}(X, (A - Y), T) / m$$

$$\frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)}) \rightarrow \text{np.sum}(A - Y) / m$$

4.3-2 - Optimize

$$\theta := \theta - \alpha d\theta$$

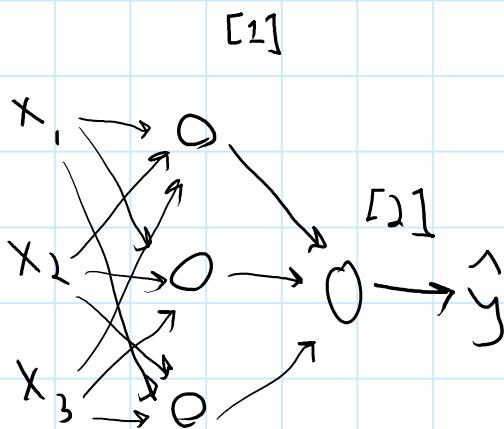
4.3-3 - Predict

Week 3

Wednesday, September 13, 2017 2:53 PM

Shallow NNs

What NNs are



$$w^{[1]} x + b^{[1]} = z^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

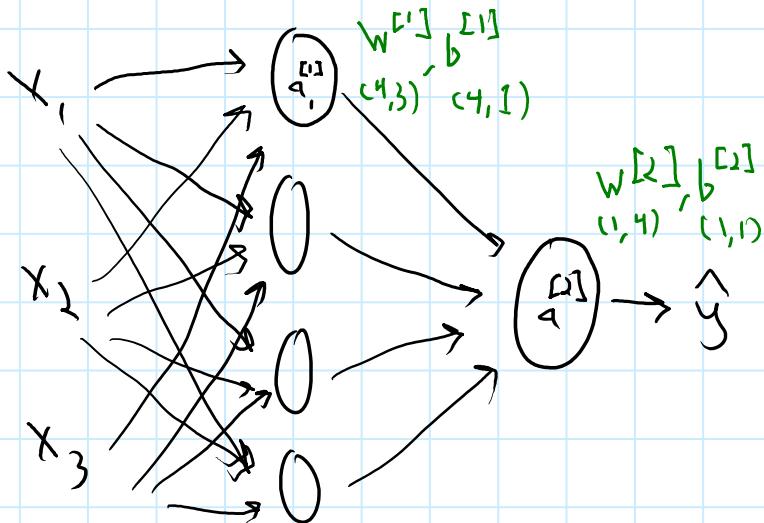
$$w^{[2]} a^{[1]} + b^{[2]} = z^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

Shallow NNs

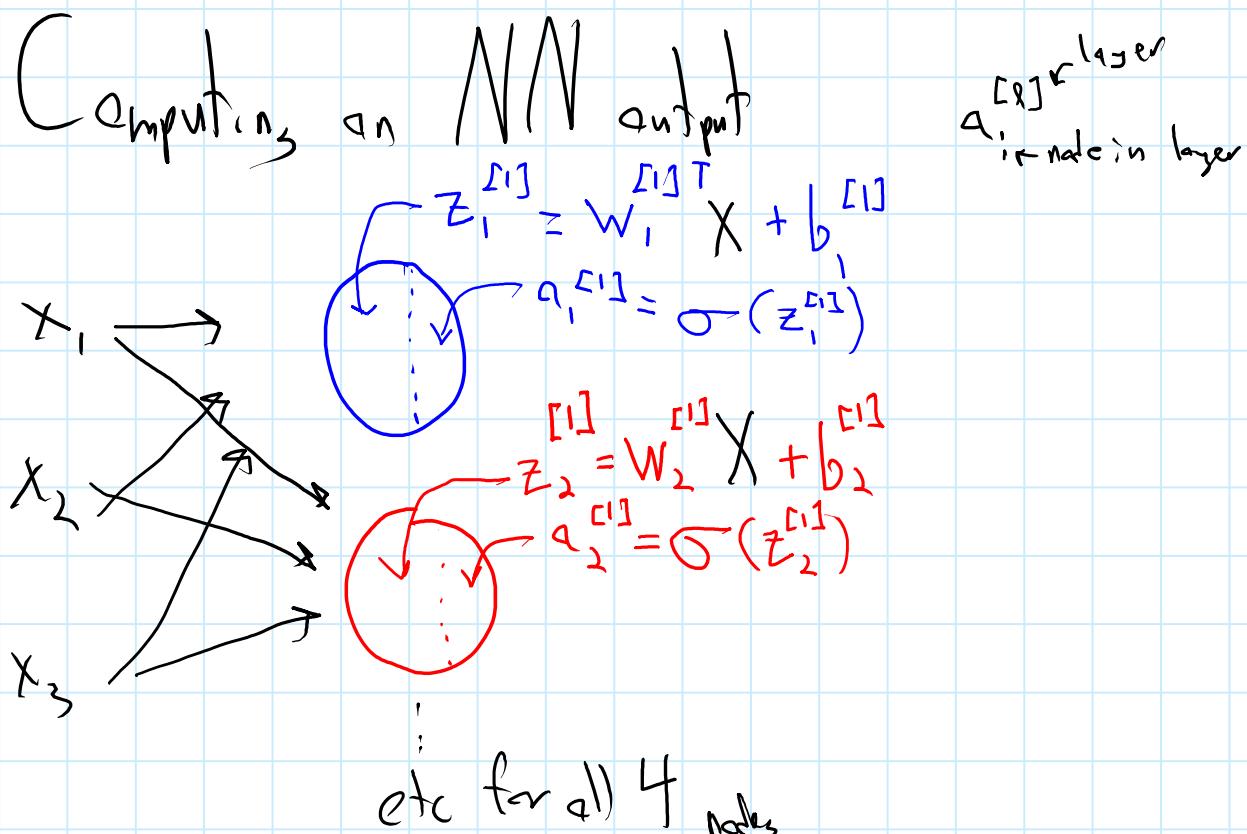
Wednesday, September 13, 2017 2:54 PM

Neural Network Representation



$$q = \begin{bmatrix} q_1 \\ q_2 \\ q_3 \\ q_4 \end{bmatrix}$$

input layer Hidden layer output layer



etc for all 4 nodes

Vectorizing

Inefficient! for $i \in \text{range}(0, \text{nodes})$:

$$\begin{aligned} z_i^{[1]} &= w_i^{[1]T} X + b_i^{[1]} \\ a_i^{[1]} &= \sigma(z_i^{[1]}) \end{aligned}$$

Vectorize:

$$\left[\begin{array}{c} -w_1^{[1]T} \\ -w_2^{[1]T} \\ \vdots \\ -w_4^{[1]T} \end{array} \right] \left[\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right] + \left[\begin{array}{c} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{array} \right] = \left[\begin{array}{c} w_1^{[1]T} X + b_1^{[1]} \\ \vdots \\ w_4^{[1]T} X + b_4^{[1]} \end{array} \right] = \left[\begin{array}{c} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{array} \right] = Z^{[1]}$$

$W^{[1]}$

$b^{[1]}$

$$\left[\begin{array}{c} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{array} \right]$$

$$w^{[2]}, b^{[2]} \\ (4, 4) \quad (1, 1)$$

$$\begin{aligned} z^{[2]} &= W^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} &= y = \sigma(z^{[2]}) \end{aligned}$$

Vectorizing Over Multiple Examples

Given m training examples

$$\begin{aligned} X^{(1)} &\rightarrow q^{[2](1)} = \hat{y}^{(1)} \\ X^{(2)} &\rightarrow q^{2} = \hat{y}^{(2)} \end{aligned}$$

$$X^{(m)} \rightarrow q^{[2](m)} = \hat{y}^{(m)}$$

$$\begin{bmatrix} X^{(1)} & \dots & X^{(m)} \end{bmatrix}$$

inefficient:

for $i = 1 \rightarrow m$,

$$Z^{[1](i)} = W^{[1]} X^{(i)} + b^{[1]}$$

$$q^{[1](i)} = \sigma(Z^{[1](i)})$$

$$Z^{[2](i)} = W^{[2]} X^{(i)} + b^{[2]}$$

$$q^{[2](i)} = \sigma(Z^{[2](i)})$$

Better:

$$Z^{[1]} = W^{[1]} X + b^{[1]} = \begin{bmatrix} Z^{1} & \dots & Z^{[1](m)} \end{bmatrix}$$

$$A^{[1]} = \sigma(Z^{[1]})$$

$$Z^{[2]} = W^{[2]} X + b^{[2]}$$

$$A^{[2]} = \sigma(Z^{[2]}) = \begin{bmatrix} q^{[2](1)} & \dots & q^{[2](m)} \end{bmatrix}$$



Justification for Vectorization

$$\begin{aligned} Z^{1} &= W^{[1]} X^{(1)} + b^{[1]} \\ Z^{[1](2)} &= W^{[1]} X^{(2)} + b^{[1]} \\ Z^{[1](3)} &= W^{[1]} X^{(3)} + b^{[1]} \end{aligned}$$

$$X = \begin{bmatrix} | & | & | \\ X^{(1)} & X^{(2)} & X^{(3)} \\ | & | & | \end{bmatrix}$$

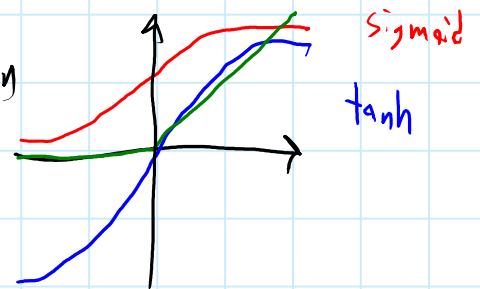
$$W^{[1]} X = Z^{[1]} = \begin{bmatrix} | & | & | \\ Z^{1} & Z^{[1](2)} & Z^{[1](3)} \\ | & | & | \end{bmatrix}$$

Activation Functions

Friday, September 15, 2017 10:05 AM

Sigmoid function not always best option

$$q^{[1]} = \sigma(z^{[1]}) = g(z^{[1]})$$



$\tanh\left(\frac{e^z - e^{-z}}{e^z + e^{-z}}\right)$ "centers" mean closer to 0

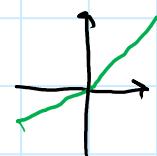
ReLU ($a = \max(0, z)$)

If using diff activation functions for diff layers, label w/ layer
e.g., $a^{[1]} = \sigma(z^{[1]})$, $a^{[2]} = g(z^{[2]})$

(General) guide to use:

Sigmoid: output is binary (output layer)

ReLU: pretty much everything else, or Leaky ReLU



Why Non-linear Activation functions?

$$a^{[2]} = W^{[2]}(W^{[1]}x + b^{[1]}) + b^{[2]}$$

$$= (\underbrace{W^{[2]}W^{[1]}}_{W'x} x + \underbrace{(W^{[2]}b^{[1]} + b^{[2]})}_{b'})$$

This is no more expressive than logistic regression

Note: if $y \in \mathbb{R}$ (like housing prices) output layer may be ok to be linear

Derivatives of Activation Functions

$$\text{Sigmoid}: g'(z) = a(1-a)$$

$$\tanh: g'(z) = 1 - a^2$$

$$\text{ReLU}: g'(z) \begin{cases} a & z < 0 \\ 1 & z > 0 \\ \text{undefined} & z = 0 \end{cases}$$

$$\text{Leaky ReLU}: g'(z) \begin{cases} 0.01 & z < 0 \\ 1 & z > 0 \\ \text{undefined} & z = 0 \end{cases}$$

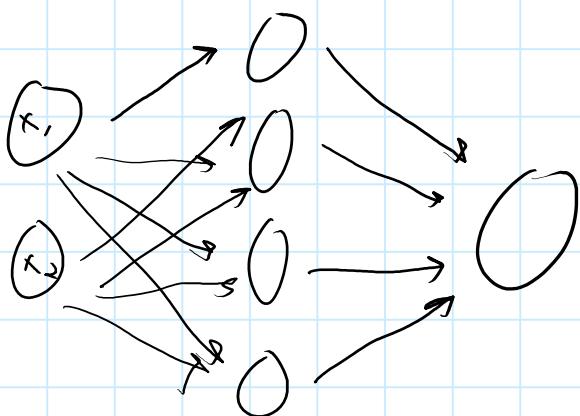
Homework Notes

Thursday, September 14, 2017 10:19 AM

Goals: Develop intuition of back-prop

Recognize more hidden layers = more complex IDim,
Build all helper functions

"Flower dataset"



$$y^{(i)} \begin{cases} 1 & a^{[2](i)} > 0.5 \\ 0 & \text{else} \end{cases}$$

$$J = -\frac{1}{m} \sum_{i=0}^m \left[y^{(i)} \log(a^{[2](i)}) + (1-y^{(i)}) \log(1-a^{[2](i)}) \right]$$

layer_sizes(X, Y):

n_x = size of input layer (parameters)

n_h = # of nodes in hidden layer

n_y = size of output layer

initialize_parameters(n_x, n_h, n_y)

$W_1 = \text{randn}(n_h, n_x)$

$b_1 = \text{zeros}(n_h, 1)$

$W_2 = \text{randn}(n_y, n_h)$

$$W_2 = \text{random}(n_y, n_h)$$
$$b_2 = \text{zeros}(n_y, 1)$$

forward_prop(X, params):
params = W_1, b_1, W_2, b_2 as dict entries
unpack params
param = params["param"]

+ forward

$$Z_1 = np.dot(W_1, X) + b_1$$

$$A_1 = \text{sigmoid}(Z_1)$$

$$Z_2 = np.dot(W_2, A_1) + b_2$$

$$A_2 = \text{sigmoid}(Z_2)$$

Week 4

Thursday, September 21, 2017

9:19 AM

Deep Neural Networks (L -layers, $L \gg 2$)

Notation

$$\begin{array}{ccccc}
 x_1 & 0 & 0 & & L=4 \\
 & 0 & 0 & 0 & \\
 x_2 & 0 & 0 & 0 & \xrightarrow{\quad} y \\
 & 0 & 0 & 0 & \\
 x_3 & 0 & 0 & 0 & \\
 & 0 & 0 & 0 & \\
 n^{[0]} = n_x & n^{[1]} = 5 = n^{[2]} & n^{[3]} = 3 & n^{[L]} = 1 &
 \end{array}$$

$$a^{[l]} = g^{[l]}(z^{[l]}) \quad W^{[l]}, b^{[l]} = \text{weights for } z^{[l]}$$

Deep NNs

Thursday, September 21, 2017 9:26 AM

Forward Prop in a deep NN

$$\begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} \quad \begin{matrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix} \quad \begin{matrix} 9 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix} \quad \begin{matrix} 0 \\ 0 \\ 0 \end{matrix}$$

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

$$z^{[1]} = W^{[1]} x + b^{[1]}, \quad a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}, \quad a^{[2]} = g^{[2]}(z^{[2]})$$

$$z^{[3]} = W^{[3]} a^{[2]} + b^{[3]}, \quad a^{[3]} = g^{[3]}(z^{[3]})$$

Getting Matrix Dimensions Right

$$\begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} \quad \begin{matrix} 0 & 9 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \quad \rightarrow \quad \hat{y}$$

$$\text{General: } W^{[l]} = (n^{[l]}, n^{[l-1]}) = \boxed{W^{[l]}}$$

$$b^{[l]} = (n^{[l]}, 1) = \boxed{b^{[l]}}$$

$$b^{Lx-1} = \begin{pmatrix} n^{Lx-1}, 1 \end{pmatrix} = db^{-1}$$
$$Z^{[x]}, A^{[x]} = \begin{pmatrix} n^{[x]}, M \end{pmatrix} = dZ^{[x]}, dA^{[x]}$$

Why Deep Networks?

Layers split up tasks among neurons
early layers may do best low level things (image audio edges, waveform)

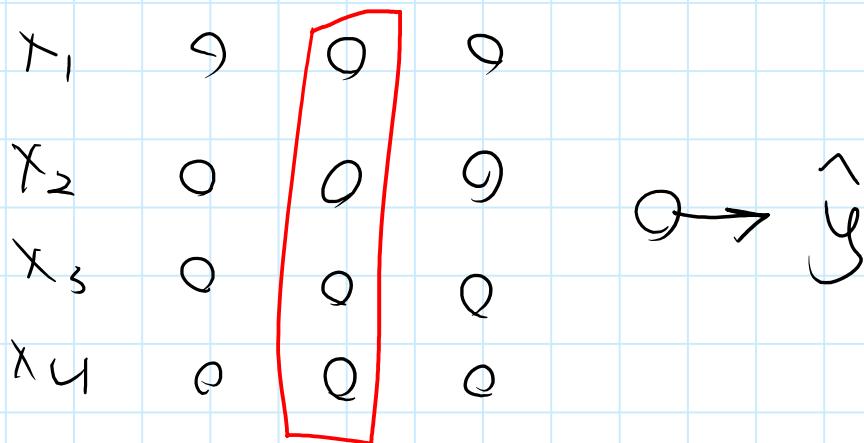
More layers can negate need for more hidden units in a given layer

Building Blocks

Thursday, September 21, 2017

10:14 AM

Forward & Backward Functions

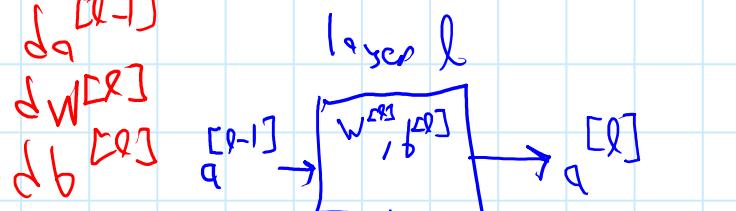


$$\text{Forward: } Z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]} \quad \text{"cache" } Z^{[l]} \text{ for later use}$$

$$a^{[l]} = g^{[l]}(Z^{[l]})$$

$$\text{Backward: Input: } d_a^{[l]}, \text{ output: } d_a^{[l-1]}$$

cache $Z^{[l]}$



$$d_a^{[l-1]} \leftarrow d_a^{[l]}$$

$$d_w^{[l]} \leftarrow d_w^{[l]} - \alpha d_w^{[l]}$$

$$d_b^{[l]} \leftarrow d_b^{[l]} - \alpha d_b^{[l]}$$

Forward & Backward Prop

Forward: input $a^{[l-1]}$: output $a^{[l]}$, cache($Z^{[l]}$)

Forward : input $a^{[l-1]}$, output a^l , cache($z^{[l]}$)

backward : input $d_a^{[l]}$, output $d_a^{[l-1]}$, $dW^{[l]}$, $db^{[l]}$

$$d_z^{[l]} = d_a^{[l]} \star g^{[l]}'(z^{[l]}) = W^{[l+1]T} d_z^{[l+1]} \star g^{[l]}'(z^{[l]})$$

$$dW^{[l]} = d_z^{[l]} \cdot a^{[l-1]}$$

$$db^{[l]} = d_z^{[l]}$$

$$d_a^{[l-1]} = W^{[l]T} \cdot d_z^{[l]}$$

Vectorized:

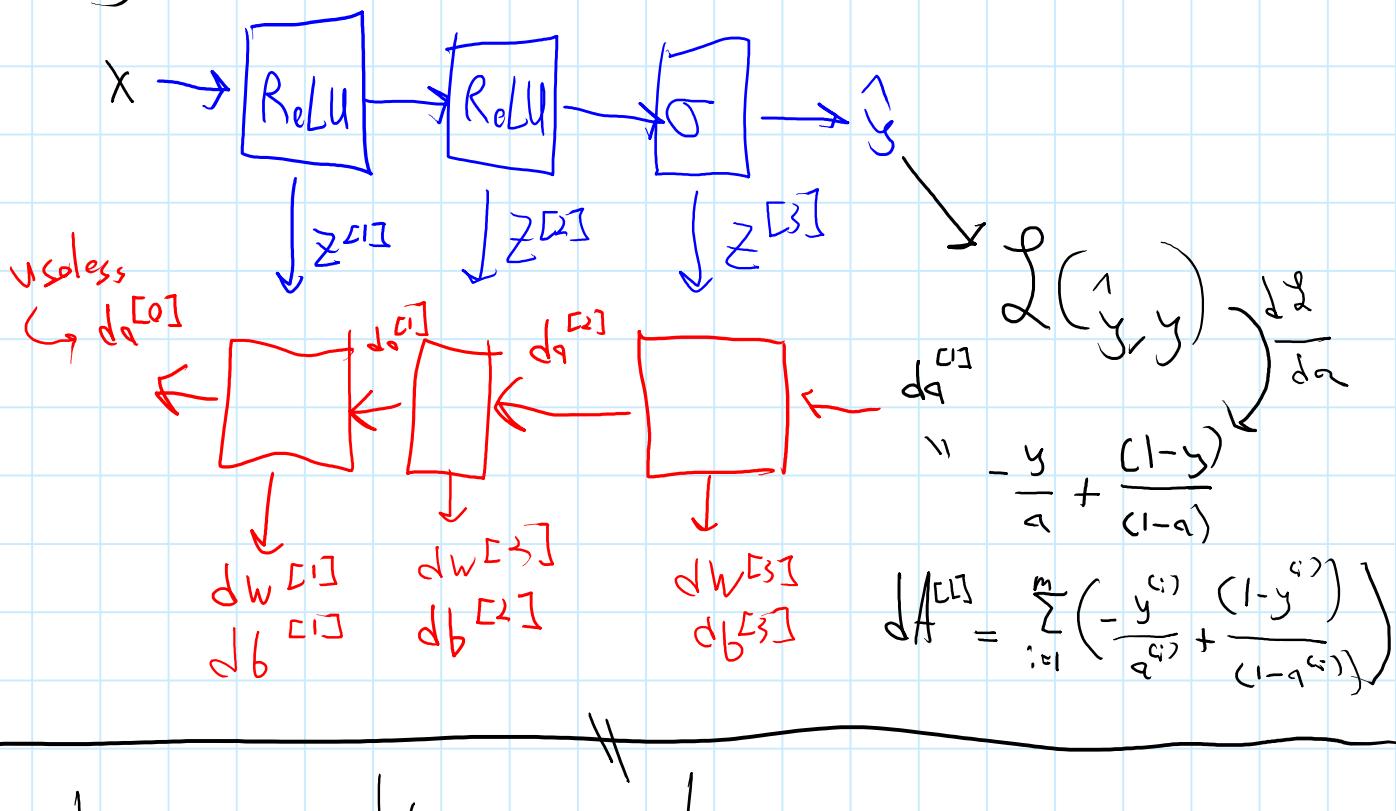
$$d_z^{[l]} = d_f^{[l]} \star g^{[l]}'(z^{[l]})$$

$$dW^{[l]} = \frac{1}{m} d_z^{[l]} \cdot A^{[l-1]T}$$

$$db^{[l]} = \frac{1}{m} \text{np.sum}(d_z^{[l]}, axis=1, keepdims=True)$$

$$dA^{[l-1]} = W^{[l]T} \cdot d_z^{[l]}$$

Summary



Parameters vs Hyperparameters

Parameters: $W^{[l]}, b^{[l]}$

Hyperparameters: learning rate (α)

of iterations

of hidden layers (L)

hidden units ($n^{[1]} n^{[2]} \dots$)

Choice of activation functions

Layer: Momentum, minibatch size, regularizations, ...