



Coursera Deep Learning Specialization

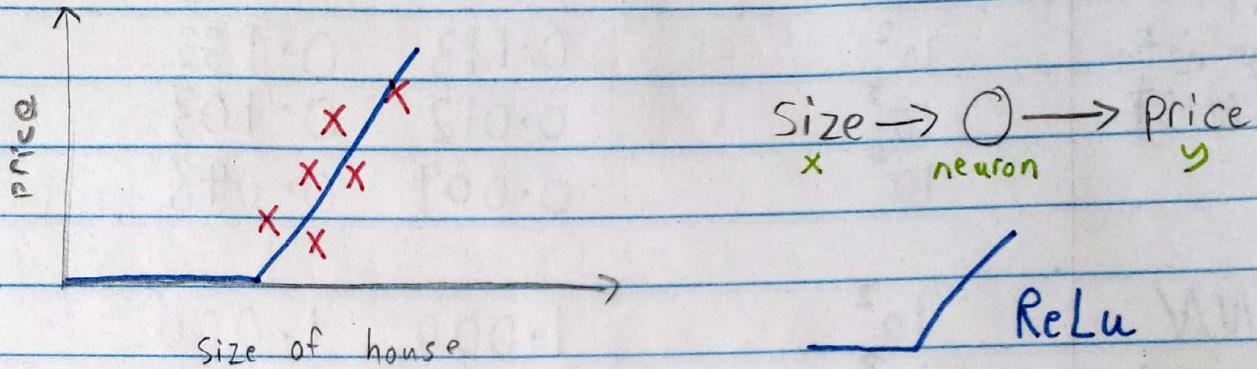
deeplearning.ai
by Andrew Ng

Mar 2020

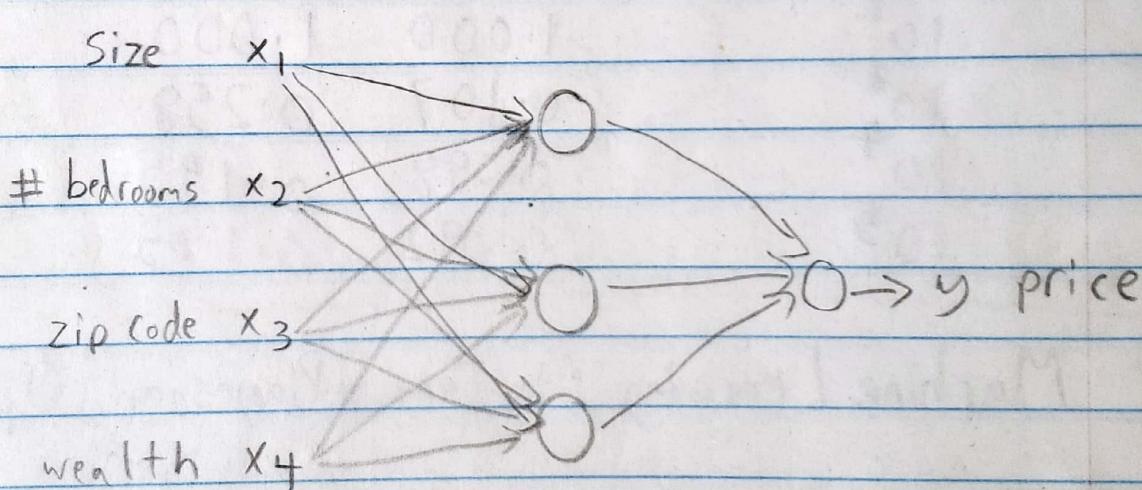
Neural Networks & Deep Learning

Intro

- Example: Housing Price Prediction



→ With more features



Layers are **densely** connected as input features connected to every circle in the middle

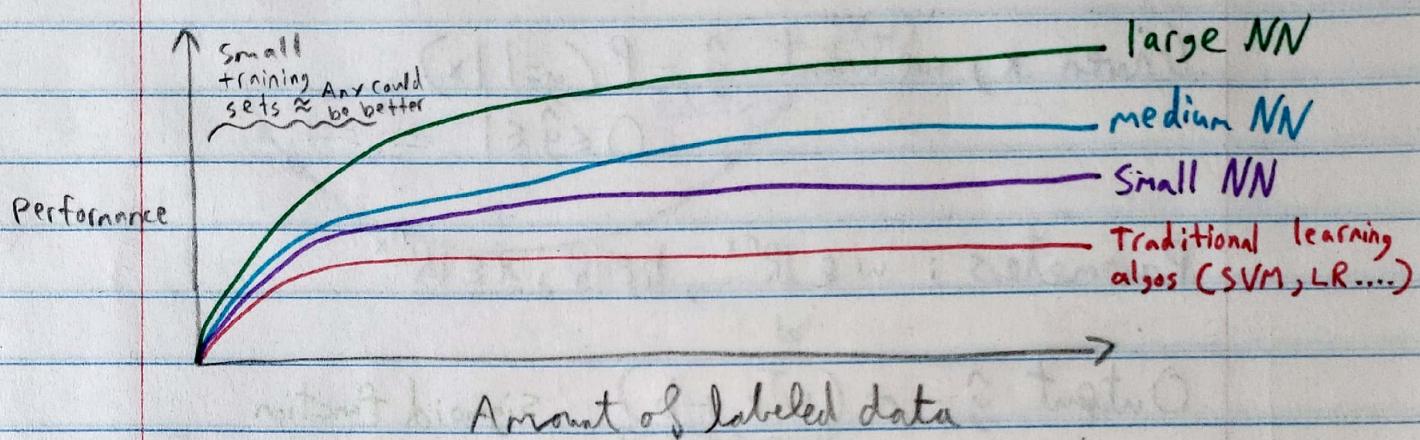
- Supervised Learning
Structured Data

Size	#bedrooms	Price
2104	3	400
1600	3	330
2400	3	369
:	:	:
3000	4	540

Unstructured Data

- Audio
- Image
- Text

- Scale drives deep learning progress



- Data
- Computation
- Algorithms

~~Week 2~~ Binary Classification

Week 2 - Logistic Regression as a Neural Network

Binary Classification

Image $\rightarrow 1$ (cat) vs 0 (non cat)

64x64

y

$$64 \times 64 \times 3 \quad x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad n = n_x = 12,288$$

$x \rightarrow y$

m training examples: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$M = M_{train}$ $M_{test} = \# \text{ test examples}$

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}_{n \times m}^T \quad X \in \mathbb{R}^{n \times m}$$

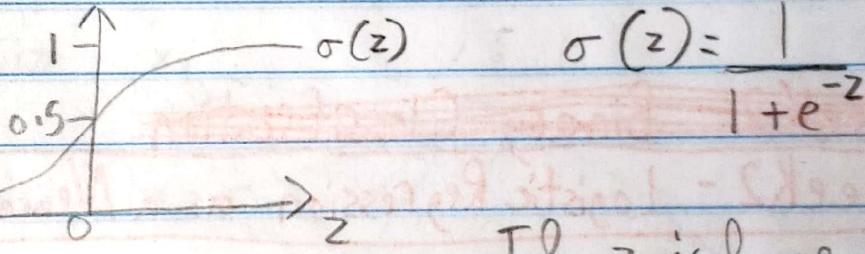
$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]_{1 \times m} \quad Y \in \mathbb{R}^{1 \times m}$$

Logistic Regression

Given x , we want $\hat{y} = P(y=1|x)$
 $0 \leq \hat{y} \leq 1$

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

Output $\hat{y} = \sigma(\underbrace{w^T x + b}_z)$ sigmoid function



If z is large $\sigma(z) \approx \frac{1}{1+0} = 1$

If z is small $\sigma(z) \approx \frac{1}{1+big\ number} \approx 0$

Loss (error) function:

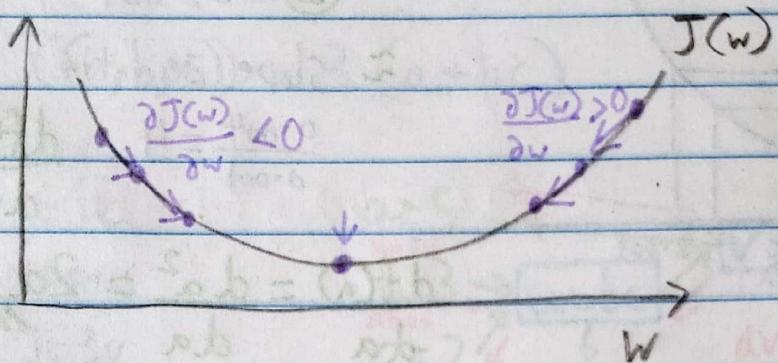
$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$$

Case 1 If $y=1$: $L(\hat{y}, y) = -\log \hat{y}$ ← make \hat{y} as large as possible

Case 2 If $y=0$: $L(\hat{y}, y) = -\log(1-\hat{y})$ ← make \hat{y} as small as possible

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

Gradient Descent



Repeat {

$$w := w - \alpha \frac{\partial J(w)}{\partial w}$$

}

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

[could be written as]

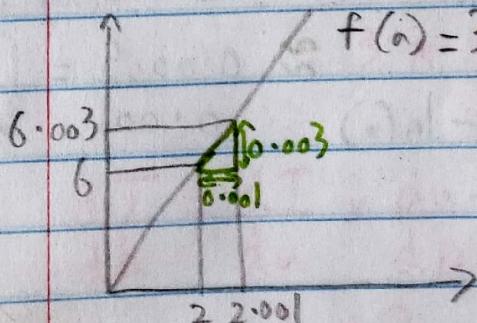
just dw

$$b := b - \alpha \frac{\partial J(w)}{\partial b}$$

$$\frac{\partial b}{\partial b}$$

[could be written as in our code]
just db

Derivatives



$$\textcircled{1} \quad a=2 \quad f(a)=6$$

$$\textcircled{2} \quad a=2.001 \quad f(a)=6.003$$

\approx slope (derivative) of $f(a)$

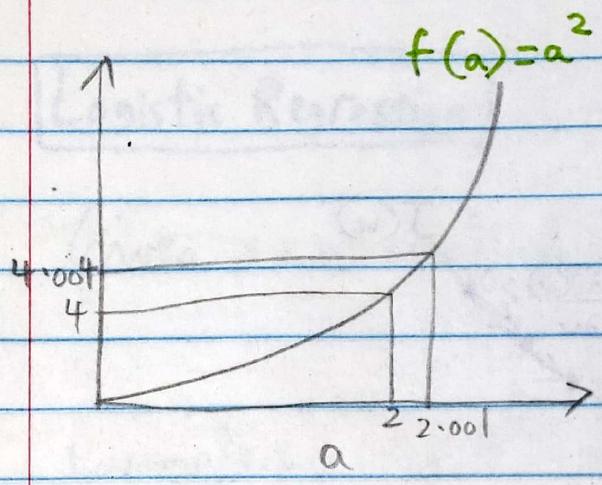
$$\frac{h}{w} = \frac{0.003}{0.001} \quad \text{at } a=2 \text{ is 3}$$

$$f(a) = 3a$$

$$\frac{df(a)}{da} = 3$$

Anywhere on the
graph the slope is 3

* Not the case for every $f(a)$!



$$f(a) = a^2$$

- ① $a=2$ when $f(a)=4$
 ② $a=2.001$ $f(a) \approx 4.004$

\approx Slope (derivative)

$$\frac{0.004}{0.001} = 4 \quad \frac{df(a)}{da} = 4 \text{ when } a=2$$

$$\frac{df(a)}{da} = \frac{d a^2}{da} = 2a$$

Anywhere on the graph, the derivative varies on the curve at diff points

- Formal definition of a derivative: What happens when we move a to the right by an infinitesimal amount $0.000\dots$

$$f(a) = a^3$$

$$\frac{df(a)}{da} = 3a^2$$

$$\textcircled{1} \quad a=2 \quad f(a)=8$$

$$\textcircled{2} \quad a=2.001 \quad f(a) \approx 8.012$$

$$\approx \frac{0.012}{0.001} = 12 \text{ // slope when } a=2$$

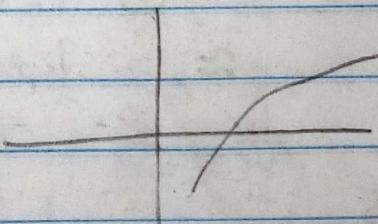
$$f(a) = \log_e(a) \quad \text{or } \ln(a)$$

$$\frac{df(a)}{da} = \frac{1}{a}$$

$$\textcircled{1} \quad a=2 \quad f(a) \approx 0.69315$$

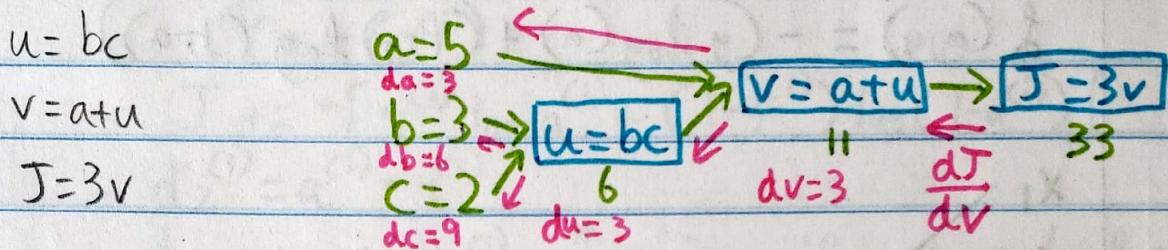
$$\textcircled{2} \quad a=2.001 \quad f(a) \approx 0.69365$$

$$\approx \frac{0.0005}{0.001} = \frac{1}{2} \text{ // slope when } a=2$$



Computational Graph

$$J(a, b, c) = 3(a + bc)$$



$$\frac{dJ}{dv} = ? = 3$$

$$J = 3V$$

$$V = 11 \rightarrow 11.\underline{001} \quad \frac{0.003}{0.001} = 3$$

$$J = 33 \rightarrow 33.\underline{003}$$

$$\frac{dJ}{da} = 3 = \frac{dJ}{dv} \times \frac{dv}{da}$$

(Chain rule)

$$= 3 \times 1$$

$a = 5 \rightarrow 5.\underline{001}$
 $v = 11 \rightarrow 11.\underline{001}$
 $J = 33 \rightarrow 33.\underline{003}$

$$\frac{dJ}{du} = 3 = \frac{dJ}{dv} \times \frac{dv}{du}$$

$$= 3 \times 1$$

$u = 6 \rightarrow 6.\underline{001}$
 $v = 11 \rightarrow 11.\underline{001}$
 $J = 33 \rightarrow 33.\underline{003}$

$$\frac{dJ}{db} = \frac{dJ}{du} \times \frac{du}{db}$$

$$= 3 \times 2$$

$$= 6$$

$b = 3 \rightarrow 3.\underline{001}$
 $u = 6 \rightarrow 6.\underline{002}$
 $J = 33 \rightarrow 33.\underline{006}$

$$\left[\frac{dJ}{dv} \times \frac{dv}{du} \times \frac{du}{db} \right] = [3 \times 1 \times 2] = [6]$$

$$\frac{dJ}{dc} = 3 \times 3$$

$$= 9$$

Logistic Regression Gradient Descent

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

$$\hat{y} = a = \sigma(z)$$

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

$$x_1, w_1, x_2, w_2, b \rightarrow z = w_1 x_1 + w_2 x_2 + b \rightarrow a = \sigma(z) \rightarrow \mathcal{L}(a, y)$$
$$\frac{dz}{d\mathcal{L}} = \frac{d\mathcal{L}}{dz} = \frac{d\mathcal{L}(a, y)}{da}$$
$$\frac{da}{d\mathcal{L}} = \frac{d\mathcal{L}(a, y)}{d\mathcal{L}}$$
$$= a - y$$
$$= \frac{d\mathcal{L}}{da} \cdot \frac{da}{dz}$$
$$= \left[\frac{-y}{a} + \frac{1-y}{1-a} \right] \times [a(1-a)] = -y(1-a) + a(1-y)$$
$$= -y + ay + a - ay$$
$$= a - y //$$

$$\frac{d\mathcal{L}}{dw_1} = "d_{w_1}" = x_1 \cdot dz \quad dw_2 = x_2 \cdot dz \quad db = dz$$

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

$$\frac{dz}{d\mathcal{L}} = \frac{d\mathcal{L}}{dz}$$

The simplification. Could be written as just that in our code.

$$J=0 ; dw_1=0 ; dw_2=0 ; db=0$$

For $i=1$ to m

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

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

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

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

$$dw_1 += x_1^{(i)} dz^{(i)} \quad \uparrow n=2 \quad \text{if more, need a for loop}$$

$$dw_2 += x_2^{(i)} dz^{(i)} \quad \downarrow$$

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

$\therefore 2 \text{ 'For' Loops}$

$$J / m$$

$$dw_1 / m ; dw_2 / m ; db / m$$

$$\rightarrow w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

To make this all more efficient Vectorize //

Vectorization

"Art of getting rid of 'for' loops"

Numpy
Python

$$z = \underbrace{\text{np.dot}(w, x)}_{w^T x} + b \rightarrow \text{way faster than non-vectorized version}$$

Can be done on both CPU & GPU

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

import numpy as np
u = np.exp(v)

Vectorized Logistic Regression (Incomplete)

$$J=0, dw=np.zeros((n_x, 1)), db=0$$

of features

for i=1 to m:

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

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

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

$$dz^{(i)} = a^{(i)}(1-a^{(i)})$$

$$dw += x^{(i)} dz^{(i)}$$

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

$$J /= m, dw /= m, db /= m$$

$$z^{(i)} = w^T x^{(i)} + b \rightarrow Z = np.dot(w^T, X) + b$$

"Broadcasting"

$$\begin{array}{c} \downarrow \\ [z^1, z^2, \dots, z^m] \end{array} \quad \begin{array}{c} \downarrow \\ [w^1, \dots, w^m] \end{array} \quad \begin{array}{c} \downarrow \\ [x^1] \end{array} \quad \begin{array}{c} \downarrow \\ [b^1, \dots, b^m] \end{array}$$

size (n^x, m)

The same thing for every other line

Complete Vectorized Logistic Regression

$$J=0, dw=np.zeros((n-x, 1)), db=0$$

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

$$A = \sigma(Z)$$

$$dZ = A - Y$$

(I believe same thing for J)

$$dw = \frac{1}{m} X dZ^T$$

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

$$w := w - \alpha dw$$

$$b := b - \alpha db$$

Broadcasting in Python

expands to match array sizes

ex1

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

ex2

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

$1 \times 3 \rightarrow 2 \times 3$

ex3

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

$2 \times 1 \rightarrow 2 \times 3$

A.sum(axis=0) sums vertically, $\downarrow^0 \rightarrow = 1$ horizontally
cal.reshape(1, 4) reshape to 1×4

- $a = np.random.rand(5)$
- $\text{print}(a, \text{shape})$ $\xrightarrow{\text{result}}$ $(5,)$ Called Rank 1 array [avoid]
- $\text{print}(a) \rightarrow [\dots \dots \dots \dots \dots] \leftarrow \text{Only 1}$
- $a = np.random.rand(5, 1)$
- $\text{print}(a, \text{shape}) \rightarrow (5, 1)$ 5×1 matrix
- $\text{print}(a) \rightarrow [[\dots \dots \dots \dots \dots]] \leftarrow 2 \text{ square brackets}$
- $\text{assert}(a, \text{shape} == (5, 1))$ To double check dimensions

Further explaining Logistic Regression Cost function

$$\text{If } y=1 \quad p(y|x) = \hat{y}$$

$$\text{If } y=0 \quad p(y|x) = 1 - \hat{y}$$

$$\rightarrow p(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

$$\text{If } y=1 \quad p(y|x) = \hat{y}$$

$$\text{If } y=0 \quad p(y|x) = 1 - \hat{y}$$

$$\log p(y|x) = \log \hat{y}^y (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log (1-\hat{y})$$

$$= -f(\hat{y}, y) \downarrow \text{minimize loss}$$

maximize log of probability
 \approx maximum likelihood estimation

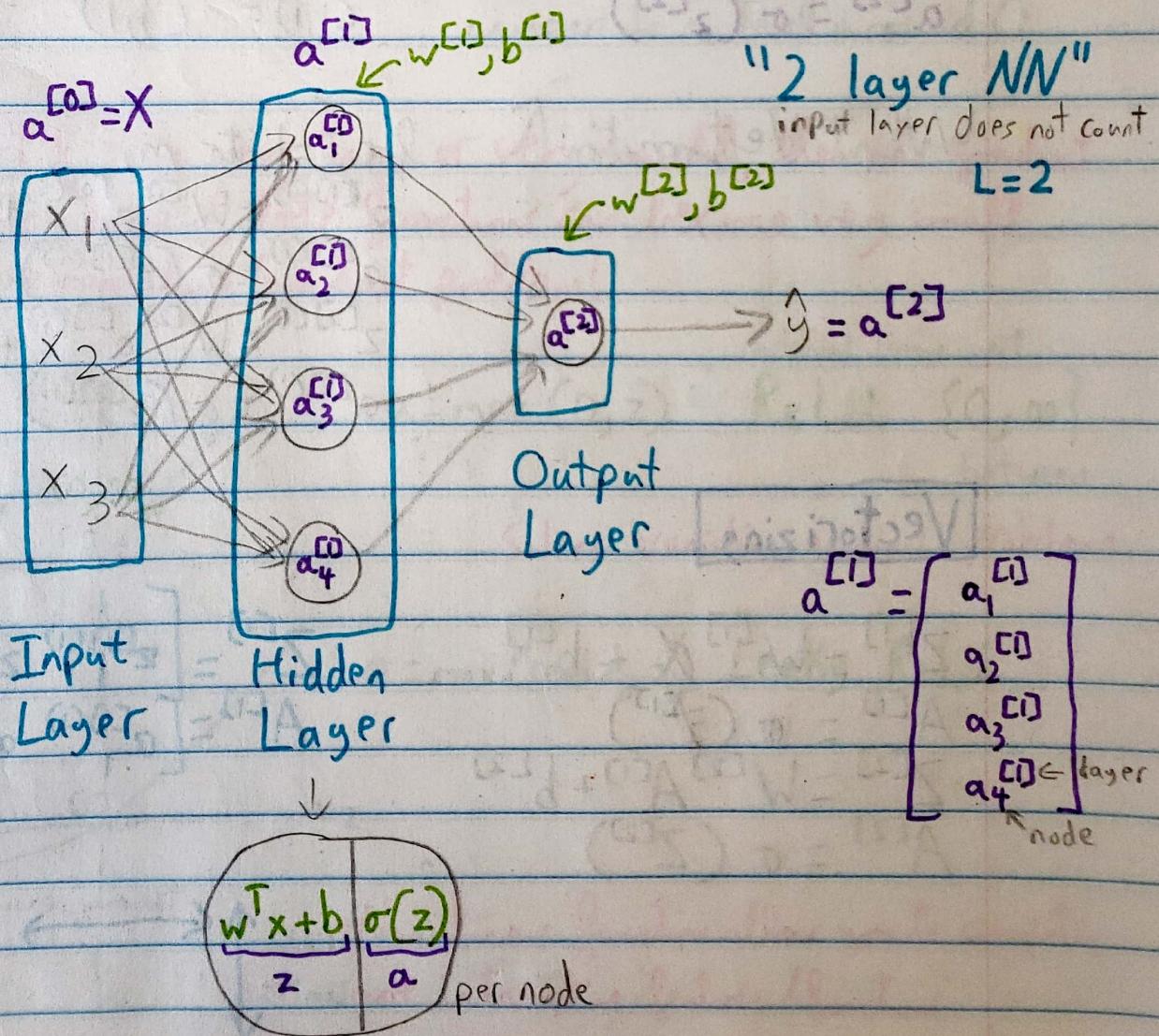
$$\text{Cost: } J(w, b) = \frac{1}{m} \sum_{i=1}^m f(\hat{y}^{(i)}, y^{(i)})$$

(minimize)

Useful tips

- A neuron computes a linear function ($z = w^T x + b$) followed by an activation function
- Column vector ($x, 1$) Row vector ($1, x$)
- $\text{np.dot}(a, b) \rightarrow$ matrix multiplication
- $a^* b \rightarrow$ element-wise multiplication. Sizes must match!
(can broadcast too!)

Week 3 - Shallow Neural Networks



$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]} \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]} \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

(4,1) (4,3) (3,1) (4,1) (4,1) (4,1)

$$\begin{bmatrix} -w_1^{[1]T} - \\ -w_2^{[1]T} - \\ \vdots \\ -w_4^{[1]T} - \end{bmatrix}$$

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

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

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \quad (1,1) = (1,4)(4,1) + (1,1)$$

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

$$(1,1) = \sigma(1,1)$$

Without Vectorization

for i=1 to m

$$z^{[1](i)} = W^{[0]} a^{[0](i)} + b^{[1]}$$

$$a^{[0](i)} = \sigma(z^{[1](i)})$$

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

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

Vectorizing

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

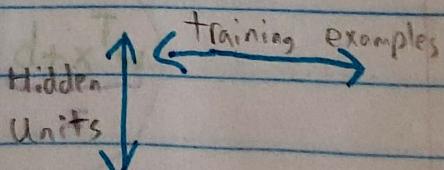
$$A^{[1]} = \sigma(Z^{[1]})$$

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

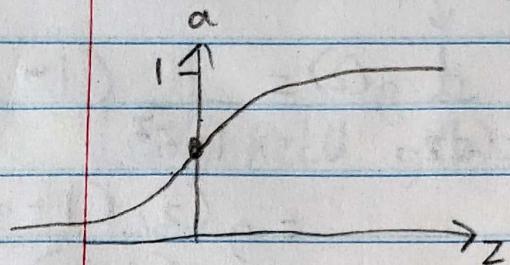
$$A^{[2]} = \sigma(Z^{[2]})$$

$$Z^{[1]} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ z_1^{1} & z_1^{[1](2)} & \dots & z_1^{[1](n)} \\ \vdots & \vdots & \ddots & \vdots \\ a_1^{[1]} & a_2^{[1]} & \dots & a_m^{[1]} \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ a_1^{[0](1)} & a_1^{[0](2)} & \dots & a_1^{[0](n)} \\ \vdots & \vdots & \ddots & \vdots \\ a_1^{[1]} & a_2^{[1]} & \dots & a_m^{[1]} \end{bmatrix}$$



Activations

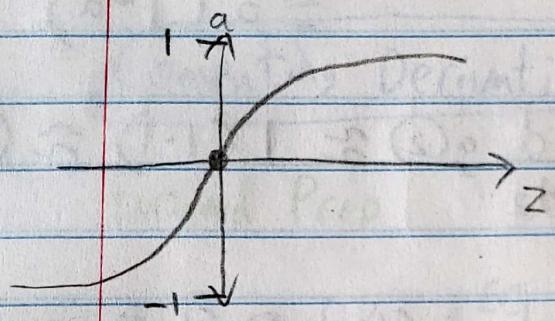


$$a = \frac{1}{1 + e^{-z}}$$

Sigmoid $\{0, 1\}$

between

- never use as tanh superior
- maybe only for binary classification



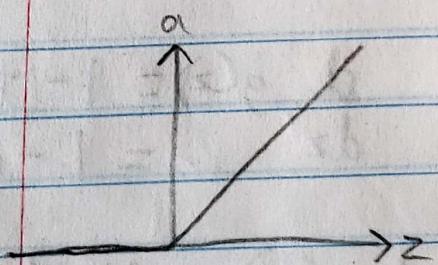
$$a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Tanh $\{-1, 1\}$

between

- Almost always work better than sigmoid as mean is centered around 0

- If z is very large or very small, then the gradient or slope of both functions can become very small
 \rightarrow Vanishing Gradient problem!

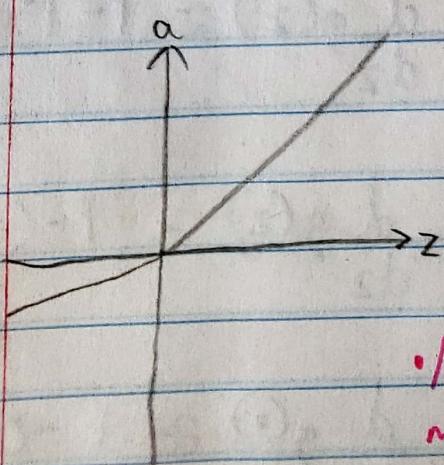


$$a = \max(0, z)$$

rectified linear unit
ReLU $\{0, \infty\}$

between

Solves vanishing gradient problem



$$a = \max(0.01z, z)$$

Leaky ReLU

- Non linear functions allow you to compute more interesting features / functions

Activation Derivatives

$$\text{Sigmoid } g(z) = \frac{1}{1+e^{-z}}$$

(slope)



$$\begin{aligned} \frac{d}{dz} g(z) &= \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}} \right) \\ &= g(z)(1-g(z)) \\ &= a(1-a) \end{aligned}$$

$$\rightarrow z=10 \quad g(z) \approx 1$$

$$\frac{d}{dz} g(z) \approx 1(1-1) \approx 0$$

$$\rightarrow z=-10 \quad g(z) \approx 0$$

$$\frac{d}{dz} g(z) \approx 0(1-0) \approx 0$$

$$\rightarrow z=0 \quad g(z) \approx \frac{1}{2}$$

$$\frac{d}{dz} g(z) \approx \frac{1}{2}(1-\frac{1}{2}) \approx \frac{1}{4}$$

$$\text{Tanh } g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\begin{aligned} \frac{d}{dz} g(z) &= 1 - (\tanh h(z))^2 \\ &= 1 - a^2 \end{aligned}$$

$$\rightarrow z=10 \quad g(z) \approx 1$$

$$\frac{d}{dz} g(z) \approx 1 - 1^2 \approx 0$$

$$\rightarrow z=-10 \quad g(z) \approx -1$$

$$\frac{d}{dz} g(z) \approx 1 - (-1)^2 \approx 0$$

$$\rightarrow z=0 \quad g(z) \approx 0$$

$$\frac{d}{dz} g(z) \approx 1 - 0^2 = 1$$

• ReLU $g(z) = \max(0, z)$ $g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$

• Leaky ReLU $g(z) = \max(0.01z, z)$ $g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$

Computing Derivatives

Forward Prop

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

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

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(Z^{[2]}) = \sigma(Z^{[2]})$$

Back prop

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dZ^{[2]}, \text{axis}=1)$$

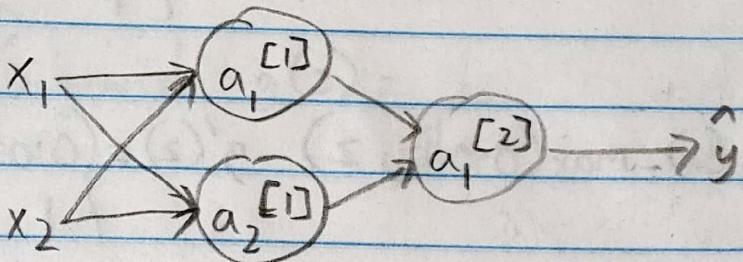
Keraspkins = True

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[2]'}(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \dots)$$

Random Initialization



$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$W^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]}$$

$$d_{z1}^{[1]} = d_{z2}^{[1]}$$

$dW = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$ ↪ Symmetric ∵ Both hidden units computing same thing. Not useful! Cost does not really decrease

Solution:

$$W^{[1]} = np.random.rand(2, 2) * 0.01$$

$$b^{[1]} = np.zeros(2, 1)$$

← does not have the same symmetry breaking problem

It's okay to have them as 0's

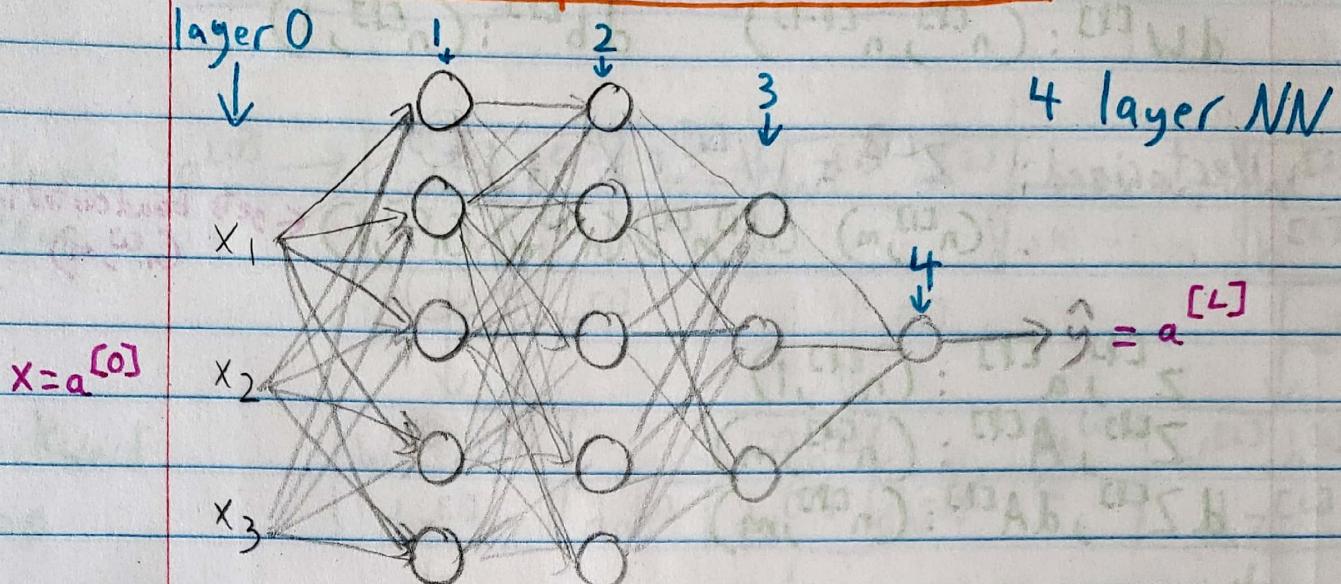
$$W^{[2]} = \dots$$

$$b^{[2]} = 0$$

W too large: Cost starts very high, e.g. Sigmoid activation outputs results

that are very close to 0 or 1 for some examples, and when it gets it wrong it incurs a very high loss for that example. $\log(a^{[2]}) = \log(0)$ Loss goes to ∞

Week 4 - Deep Neural Networks



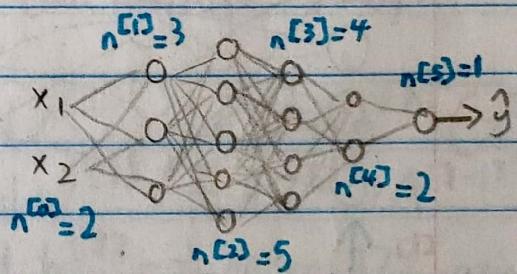
$$L=4$$

$n^{[l]}$ = # units in layer l

$$\begin{aligned} n^{[1]} &= 5, n^{[2]} = 5, n^{[3]} = 3. \\ n^{[4]} &= n^{[L]} = 1, n^{[0]} = n_x = 3 \end{aligned}$$

- Forward Propagation same way as we did before but now over 4 layers

Matrix Dimensions



Note:
z & x not
vectorized here.
Check next pg

$$z^{[1]} = W^{[1]} \cdot x + b^{[1]}$$

$$(3,1) \quad (3,2) \quad (2,1) \quad (3,1)$$

$$(n^{[0]},1) \quad (n^{[1]}, n^{[0]}) \quad (n^{[0]}, 1)$$

$$z^{[2]} = W^{[2]} \cdot a^{[1]} + b^{[2]}$$

$$(5,1) \quad (5,3) \quad (3,1) \quad (5,1)$$

$$(n^{[1]},1) \quad (n^{[2]}, n^{[0]}) \quad (n^{[1]}, 1)$$

$$\begin{aligned} W^{[1]} &: (n^{[1]}, n^{[0]}) \\ W^{[2]} &: (n^{[2]}, n^{[1]}) \\ \dots & \\ W^{[L]} &: (n^{[L]}, n^{[L-1]}) \\ b^{[1]} &: (n^{[1]}, 1) \end{aligned}$$

dW & db should be the same dimensions as W and b

$$dW^{[l]} : (n^{[l]}, n^{[l-1]}) \quad db^{[l]} : (n^{[l]}, 1)$$

Vectorized: $Z^{[l]} = W^{[l]} \cdot X + b^{[l]}$

$(n^{[l]}, m) \quad (n^{[l]}, n^{[0]}) \quad (n^{[0]}, m) \quad (n^{[l]}, 1) \leftarrow \text{gets broadcasted in python}$
 $(n^{[0]}, m)$

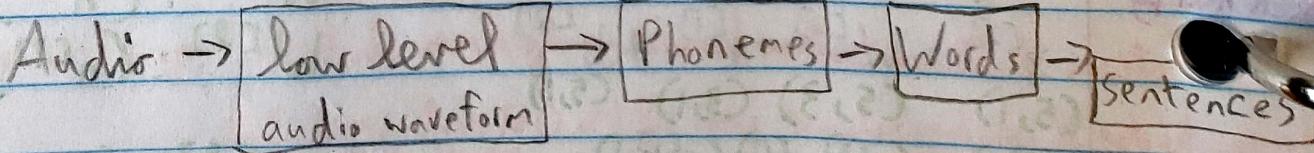
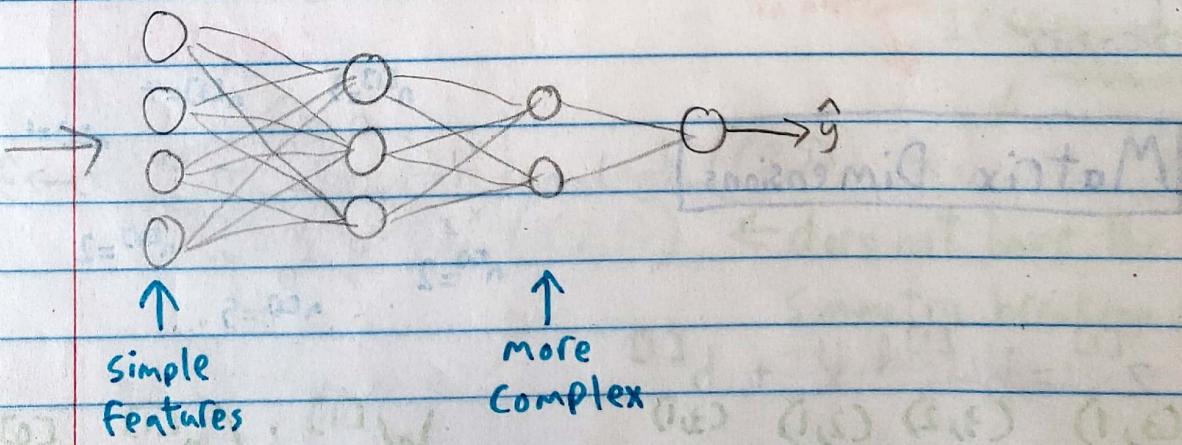
$$Z^{[l]}, a^{[l]} : (n^{[l]}, 1)$$

$$Z^{[l]}, A^{[l]} : (n^{[l]}, m)$$

$$dZ^{[l]}, dA^{[l]} : (n^{[l]}, m)$$

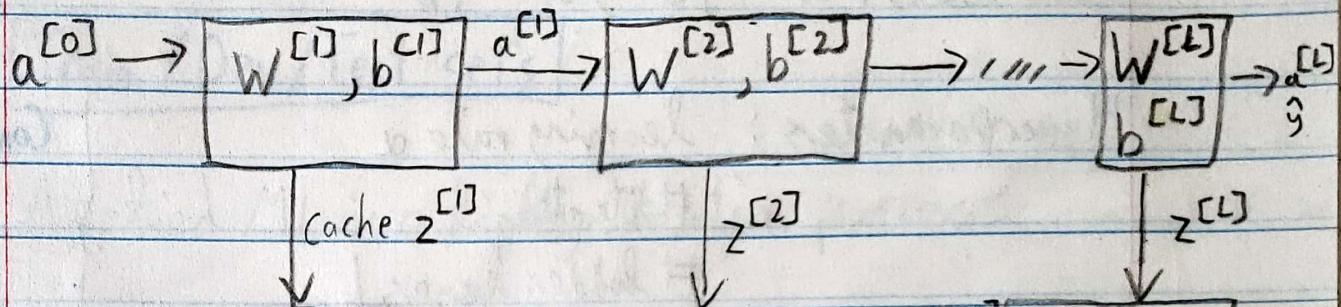
Why deep representation?

There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden units to compute

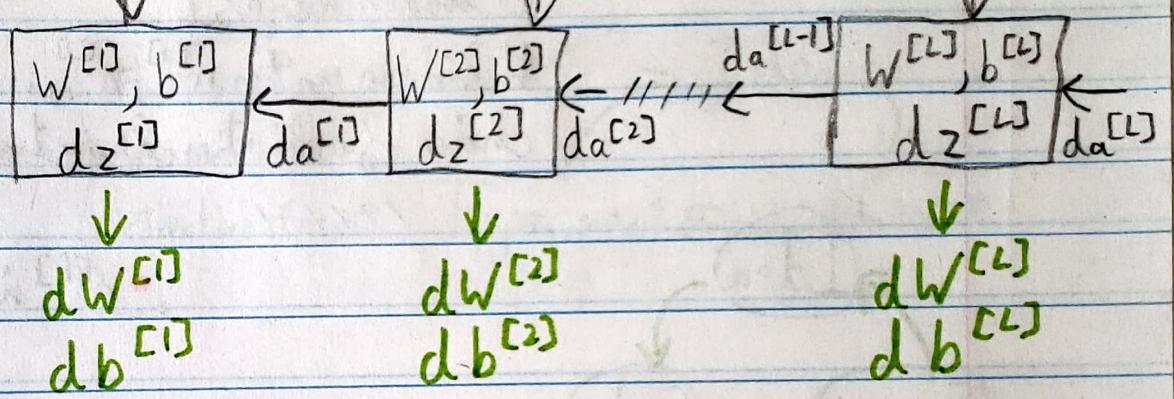


Deep Learning Building Blocks

Forward Prop



Backward Prop



$$W^{[l]} := W^{[l]} - \alpha dW^{[l]}$$

$$b^{[l]} := b^{[l]} - \alpha db^{[l]}$$

Forward Prop

$$Z^{[l]} = W^{[l]} \cdot A^{[l-1]} + b^{[l]}$$

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

Backward Prop

$$dZ^{[l]} = dA^{[l]} \cdot g^{[l]}'(Z^{[l]})$$

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

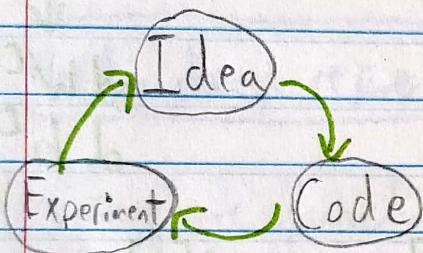
$$db^{[l]} = \frac{1}{m} \text{np.sum}(dZ^{[l]}), \text{axis}=1$$

Keepdims=True

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

Parameters vs Hyperparameters

- Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots$
 - Hyperparameters:
 - learning rate α
 - # iterations
 - # hidden layers L
 - # hidden units $n^{[1]}, n^{[2]}, \dots$
 - choice of activation function
- $a^{[L]}$ is not a hyperparameter



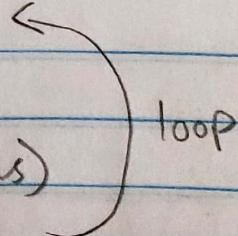
Vectorization Note:

→ Vectorization does not allow you to compute forward propagation in an L -layer neural network without an explicit for-loop over layers $l=1, 2, \dots, L$

NN Steps

1. Initialize parameters
2. Compute forward prop
3. Compute cost
4. Compute backward prop (get grads)
5. Update parameters

Optimize cost



Coursera Deep Learning Specialization: Course #2

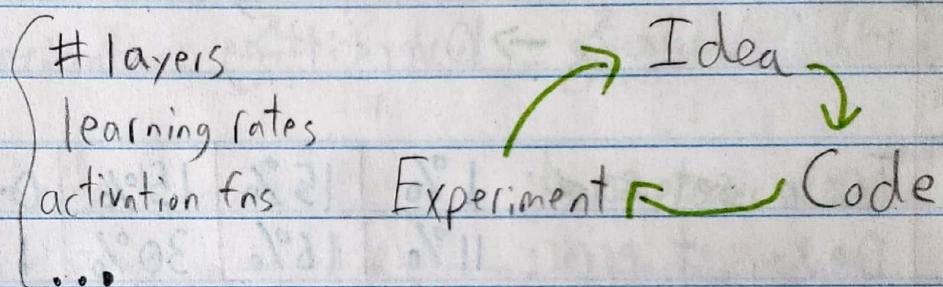
Improving Deep NN: Hyperparameter tuning, Regularization, Optimization

Week 1 - Practical Aspects of Deep Learning

Setting up
ML

Train / Dev / Test sets

- Applied ML is a highly iterative process



Data

training set	- Hold out CV	test set
--------------	---------------	----------

- Previous era: 70/30 or 60/20/20

when data > 1,000,000

- Big data: 1,000,000

98/1/1

- Development set "dev"

To evaluate algorithms

• We don't need that much for dev & test

Given your final classifier, give a confident estimate of how well it is doing [unbiased]

Mismatched train / test distribution

e.g. Training set: cat pics

from web

Dev / Test set: cat pics

from user apps

→ Make sure dev & test come from same distribution

- Not having a test set might be okay (only dev set)
 - If you only want to iterate & evaluate to build the best model

Bias / Variance

- High Bias → Underfitting
- High Variance → Overfitting

Train set error:	1%	15%	15%	0.5%
Dev set error:	11%	16%	30%	1%
	high variance	high bias	high bias & high variance	low bias & low variance

→ Assumption Human \approx 0% error
 Optimal (Bayes) error \approx 0%

Basic recipe for ML

- High bias? some fixes
 - Get additional features
 - Add polynomial features
 - Decrease λ
 - NN architecture search
- High variance? some fixes
 - More training examples
 - Smaller set of features
 - Increase λ
 - NN architecture search

- Modern era: No need for Bias / Variance trade off due to the amount of tools we have

Regularization

Regularization

λ = regularization parameter

w_1, w_2, w_3

(np.sum(np.square(w)))

+ ... + w_2
+ ... + w_2
 $\times [2/m]$

→ Weights end
in smaller
(Weight Decay)

$$L_2 \text{ regularization: } \|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w \quad \begin{array}{l} \rightarrow \text{drive } w \text{ to smaller values} \\ \rightarrow \text{used much more often} \end{array}$$

$$L_1 \text{ regularization: } \frac{\lambda}{2m} \sum_{i=1}^n |w_i| = \frac{\lambda}{2m} \|w\|_1, \quad \begin{array}{l} w \text{ will be sparse} \\ \rightarrow \text{it will have a lot of 0's} \end{array}$$

• we don't regularize b as it's only a number, whereas w is a high dimensional parameter vector

→ compresses model ∴ uses less memory

- Neural networks

$$J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \ell(g(i), y(i)) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{[l]})^2 \quad w: (n^{[0]} \ n^{[1]} \ \dots \ n^{[L-1]})$$

"Frobenius norm" $\|\cdot\|_2^2 \quad \|\cdot\|_F^2$

$$\rightarrow dw^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[l]} \quad \begin{array}{l} \frac{d}{dw} \left(\frac{1}{2} \frac{\lambda}{m} w^2 \right) = \frac{\lambda}{m} w \\ \text{if } \lambda < 0 \\ (1 - \alpha \frac{\lambda}{m}) \end{array}$$

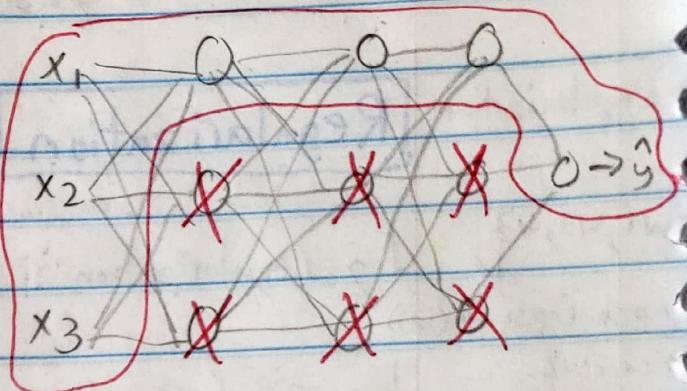
$$w^{[l]} = w^{[l]} - \alpha dw^{[l]}$$

→ λ : penalizes weight matrices from being too large

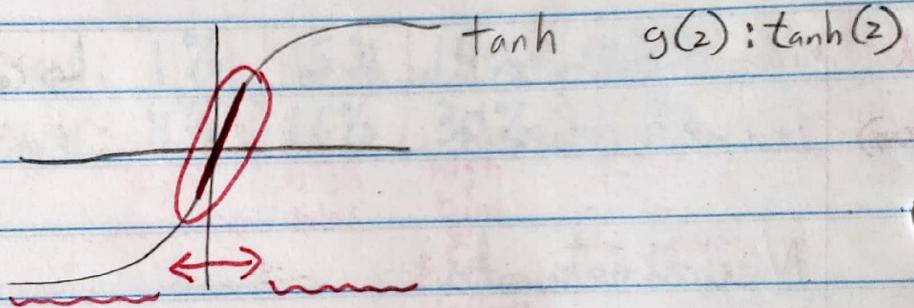
regularization

Why regularization reduces overfitting

- Intuition #1 $w^{[l]} \approx 0$



- Intuition #2



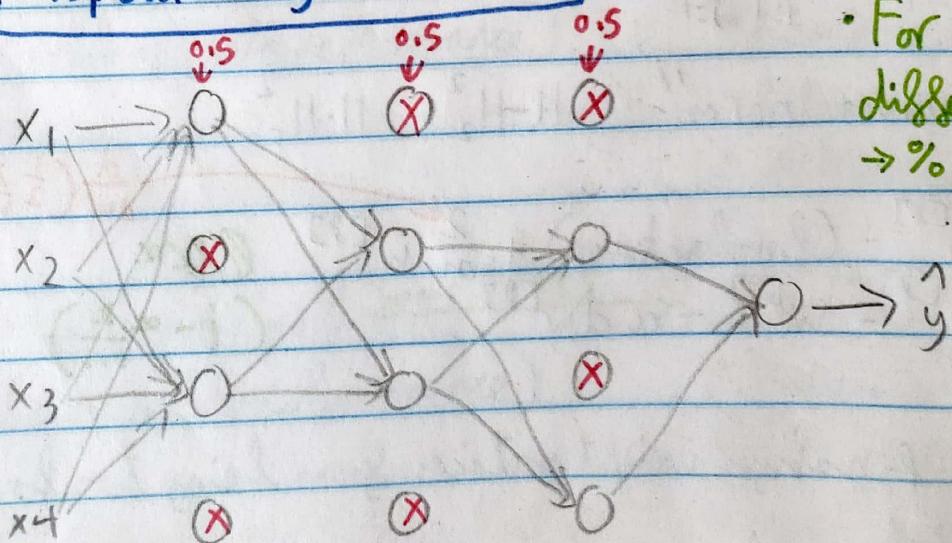
$$z \uparrow \quad w^{[l]} \downarrow \quad z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

Every layer \approx linear

regularization

Dropout Regularization

Reg technique to avoid overfitting



- For every example
different neurons activated
 \rightarrow % depends on dropout factor

- Implementing "Inverted Dropout" Most common way

Illustrate with layer $l=3$ Keep-prob = 0.8

$$\rightarrow d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \text{Keep-prob}$$
$$a_3 = \text{np.multiply}(a_3, d_3)$$

$a_3 /= \text{Keep-prob}$ 50 units \leadsto 10 units shut off
 \downarrow Scale value of neurons $\rightarrow a_3$ reduced by 20%
that haven't been shut down + assure cost result = same expected val as without drop out

\rightarrow Makes test time easier as there is no scaling problem

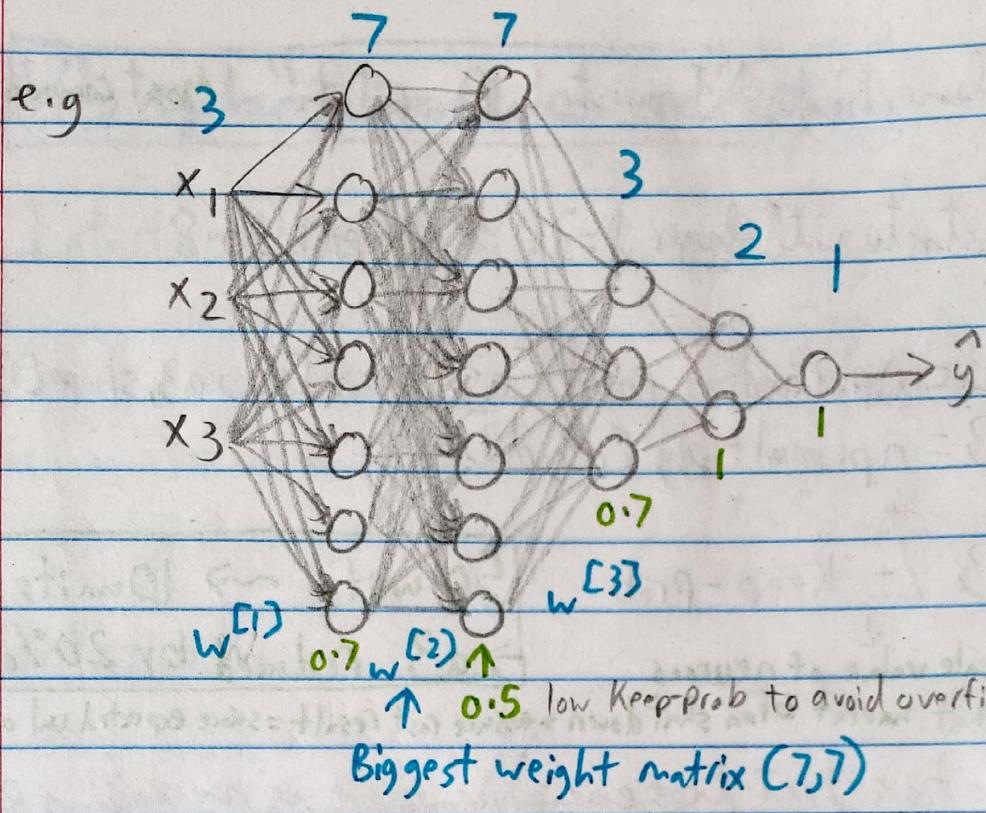
- At test time \rightarrow No dropout as it will just add noise to predictions (we don't want output to be random)
 \rightarrow Computationally inefficient to average result over many iterations if dropout was used

regularization

Understanding Dropout

- Why does dropout work?

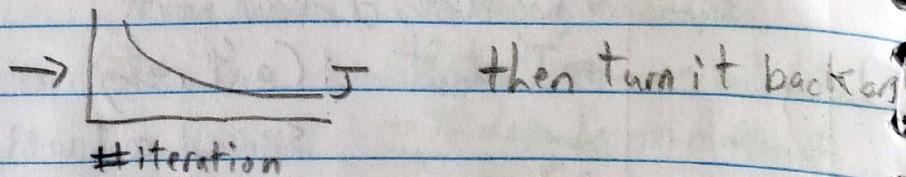
\rightarrow Intuition: Can't rely on any one feature, so have to spread out weights across all features.
 \therefore Shrinks weights!



→ e.g. In computer vision due to low amount of data dropout almost always used by default to avoid overfitting

Downside: J(C) not well defined due to the random activated neurons

→ Solution: Turn off dropout, make sure J(C) keeps decreasing

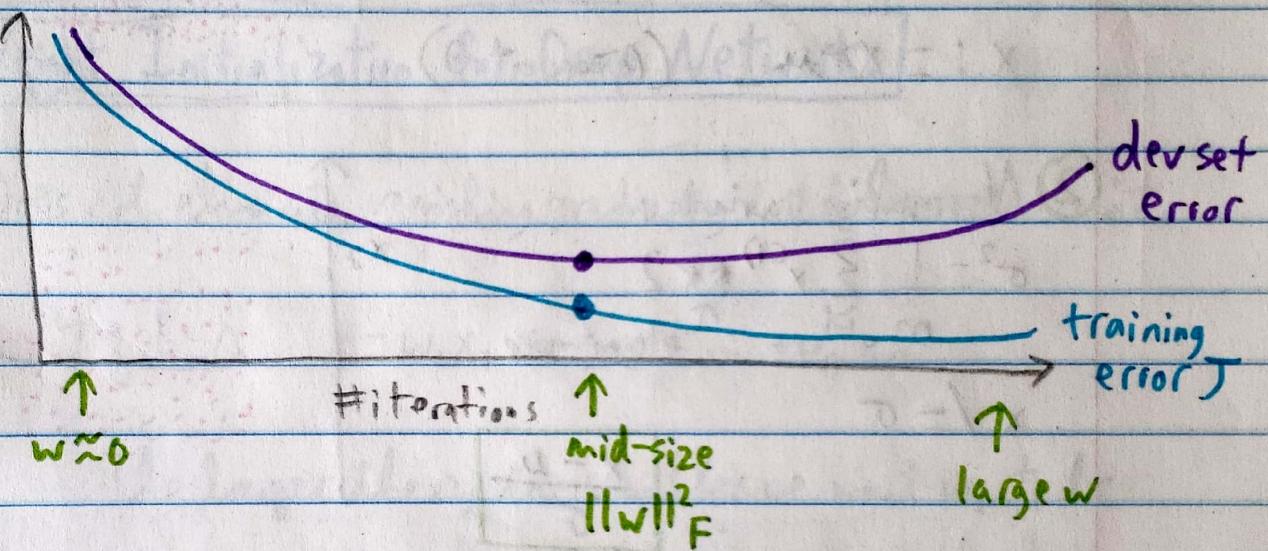


• Apply dropout both during forward & backward propagation

Regularization

Other Regularization Methods

- Data Augmentation e.g. flipping, zooming
- Early Stopping



→ Neural network was doing best during mid-size so we can stop then

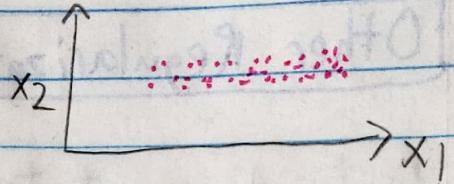
Downside those 2 are separate tasks

- Optimize cost function J
- GD, Adam, ...
- Not overfit
- Regularization, ...

Early stopping couples those 2 tasks. Does not work optimally for both. Need to use different tools for each task.

Setting up
your optimization
problem

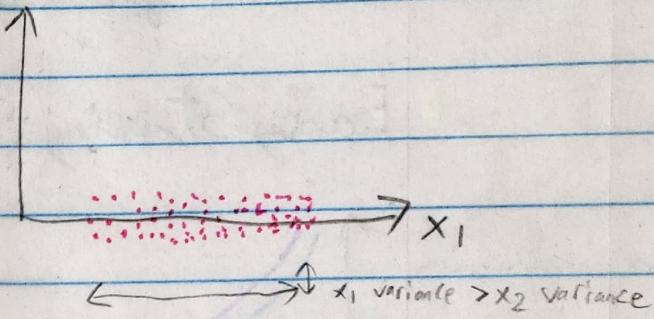
Normalizing Inputs



- ① Subtract mean:

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x := x - \mu \quad (\text{centers at } 0)$$

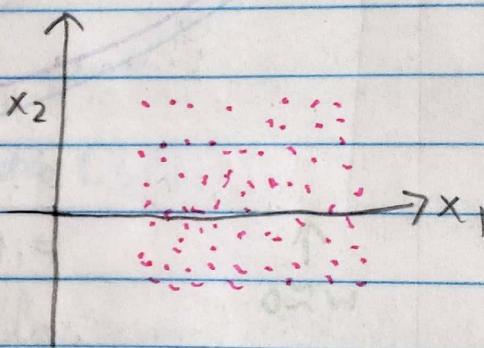


- ② Normalize variance:

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} \star \star 2 \quad \text{element-wise}$$

$$x / \sigma$$

$$\boxed{\frac{x - \mu}{\sigma}}$$



Important: We use same μ and σ to normalize test set

- Why normalize inputs?

- To make features on same scale, making cost function J easier and faster to optimize

Setting up your
optimization
problem

Vanishing/Exploding Gradients

- When derivatives/slopes get really big or really small, making training difficult.

If linear activation

If $W^{[L]} > 1$ Activations will explode exponentially
e.g 1.5^L where L is the number of layers

If $W^{[L]} < 1$ Activations will decrease exponentially
e.g 0.5^L

Setting up your optimization problem

Weight Initialization for Deep Networks

- Does not eliminate vanishing gradient, but helps reduce it

e.g If $b=0$ $z = w_1x_1 + w_2x_2 + \dots + w_nx_n$

The larger then n , the smaller we want w_i to be

$\rightarrow \text{Var}(w_i) = \frac{1}{n}$ if ReLU $\frac{2}{n}$ works better

$$W^{[l]} = \text{np.random.randn(Shape)} * \text{np.sqrt}\left(\frac{1}{n^{[l-1]}}\right)$$

Xavier Initialization $\frac{1}{\sqrt{n^{[l-1]}}}$ if tanh

Other variants $\frac{2}{n^{[l-1]} + n^l}$

He $\frac{2}{n^{[l-1]}}$

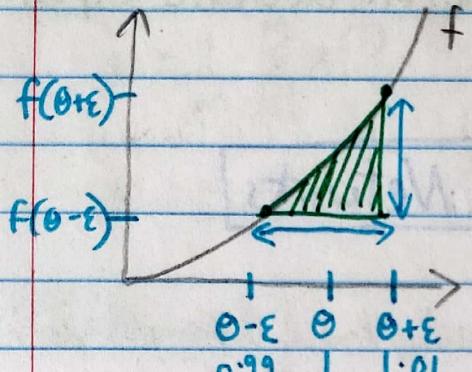
- These could be hyperparameters that can be tuned if we would like
 \rightarrow Low priority compared to other hyperparameters

Setting up
your optimization
problem

Numerical Approximations of Gradients

& Checking your
derivative computation

$$\text{e.g. } f(\theta) = \theta^3$$



$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx g(\theta)$$

$$\frac{1.01^3 - 0.99^3}{2(0.01)} = \frac{3.0001}{0.02} \approx 3$$

$$\epsilon = 0.01$$

$$g(\theta) = 3\theta^2 = 3$$

$$\text{approx error: } 0.0001$$

- Get a much better gradient estimate when taking $(\theta - \epsilon)$ into account

Gradient Checking

- To verify that backpropagation is correct
- Take $W^{[0]}, b^{[0]}, \dots, W^{[L]}, b^{[L]}$ and reshape into a big vector θ
 $J(W^{[0]}, b^{[0]}, \dots, W^{[L]}, b^{[L]}) = J(\theta)$

Take $dW^{[0]}, db^{[0]}, \dots, dW^{[L]}, db^{[L]}$ and reshape into a big vector $d\theta$
Is $d\theta$ the gradient of $J(\theta)$?

for each i :

$$\text{grad approx} = d\theta_{\text{approx}}[i] = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\tilde{d}\theta[i] = \frac{\partial J}{\partial \theta_i} \quad | \quad d\theta_{\text{approx}} \tilde{d}\theta$$

Check $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}$

$\approx 10^{-3}$ worry
 $\approx 10^{-2}$ great!
if $\epsilon = 10^{-7}$

Gradient Checking Implementation Notes

- Don't use in training - only use to debug because it's slow.
- If algorithm fails grad check, look at components (e.g. $db^{(l)}$ or $dW^{(l)}$) to try to identify bug.
- Remember regularization.
- Doesn't work with dropout. (Keep-prob=1)
- Run at random initialization; perhaps again after some training

$$\frac{\partial J}{\partial \theta} = \lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

Assignments

Initialization Assignment Conclusions:

- Different initializations lead to different results
- Random initialization is used to break symmetry & for units to learn different things
- Don't initialize to values that are too large
- He initialization works well for networks with ReLU activations

	Train acc	
O's initialization	50%	Fails to break symmetry
Random	83%	Too large weights
He	99%	Recommended method

Week 2 - Optimization Algorithms

Mini-Batch Gradient Descent

$$X = \left[\begin{array}{cccc|ccccc} (1) & (2) & (3) & \dots & (1000) & | & (1001) & \dots & (2000) \\ X & X & X & \dots & X & | & \dots & \dots & \dots \\ (n_x, m) & & & & & | & & & & (n_x, m) \end{array} \right] \quad X^{\{1\}}_{(n_x, 1000)} \quad X^{\{2\}}_{(n_x, 1000)} \quad \dots \quad X^{\{5000\}}_{(n_x, 1000)}$$

$$Y = \left[\begin{array}{cccc|ccccc} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} & | & y^{(1001)} & \dots & y^{(2000)} \\ \dots & \dots & \dots & \dots & \dots & | & \dots & \dots & \dots \\ (1, m) & & & & & | & & & & (1, m) \end{array} \right] \quad Y^{\{1\}}_{(1, 1000)} \quad Y^{\{2\}}_{(1, 1000)} \quad \dots \quad Y^{\{5000\}}_{(1, 1000)}$$

- If $m=5,000,000$
- 5000 mini-batches of 1,000 each

for $t=1, \dots, 5000$

→ Forward prop on $X^{[t]}$

$$Z^{[0]} = W^{[1]} X^{[t]} + b^{[0]}$$

$$A^{[0]} = g^{[0]}(Z^{[0]})$$

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

Vectorized

→ Compute Cost $J^{[t]} = \frac{1}{1000} \sum_{i=1}^l f(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \|W\|^2$

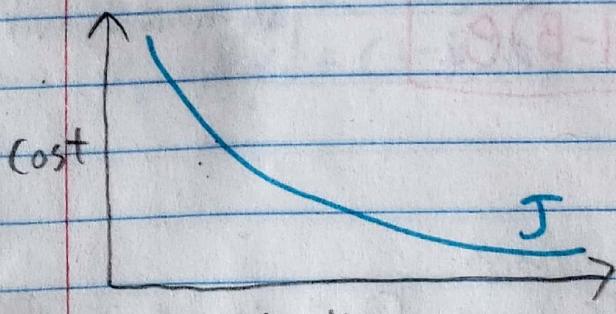
→ Backward prop to compute gradients w.r.t $J^{[t]}$

$$W^{[l]} := W^{[l]} - \alpha dW^{[l]}, \quad b^{[l]} := b^{[l]} - \alpha db^{[l]}$$

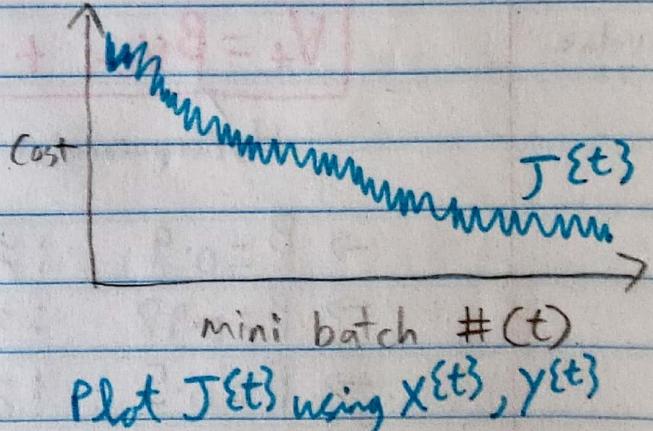
}" "1 epoch" \approx 1 pass through training set

→ Much faster than Batch Gradient Descent

Batch Gradient Descent



Mini-Batch Gradient Descent



Plot $J^{[t]}$ using $X^{[t]}, y^{[t]}$

Choosing Mini-Batch Size

- If mini-batch size = m : Batch Gradient Descent
- If mini-batch size = 1 : Stochastic Gradient Descent

Stochastic GD

→ Lose speedup
from vectorization

Mini-Batch GD

→ Value between 1 & m
→ Fastest learning
• Vectorization
• Doesn't process entire training set

Batch GD

→ Too long per iteration

- It's another hyperparameter that can be tuned
- If small training set ($m < 2000$), use batch GD
- Typical mini-batch sizes: 64, 128, 256, 512
- Make sure mini-batch fit in CPU/GPU memory

Exponentially Weighted Averages

(moving average)

- Reduces noise and locates trends via local weighted average

V = value

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

Hyperparameter

$$\rightarrow \beta = 0.9 \quad ; \approx 10 \text{ days}$$

$$\rightarrow \beta = 0.98 \quad ; \approx 50 \text{ days} \quad [\text{least noisy but graph shifts}]$$

$$\rightarrow \beta = 0.5 \quad ; \approx 2 \text{ days}$$

- Averaging over short window

- More noise

- More susceptible to outliers, interacts quickly

$$V_t \text{ averaging over } \approx \frac{1}{1-\beta}$$

e.g. $\beta = 0.9$

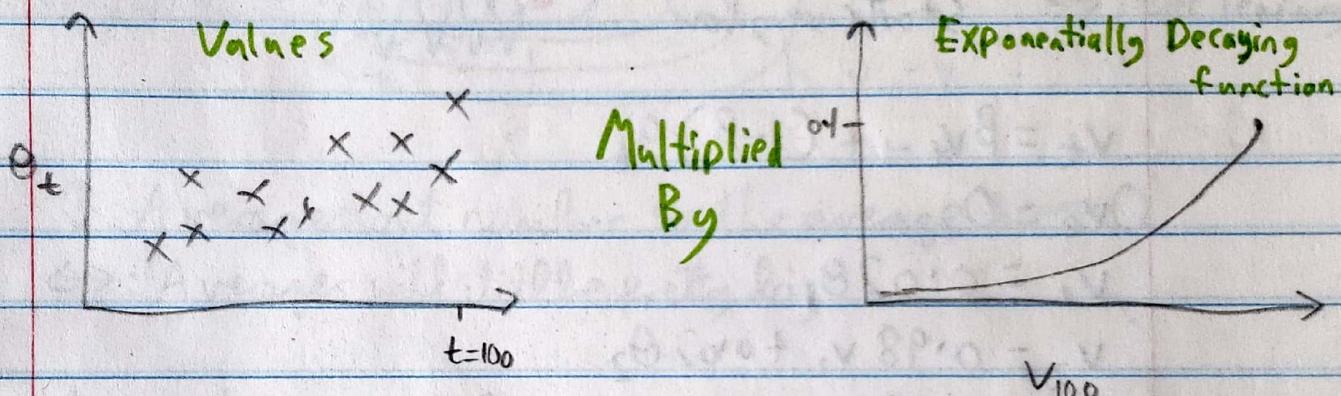
$$v_{100} = 0.9 v_{99} + 0.1 v_{100}$$

$$v_{99} = 0.9 v_{98} + 0.1 v_{99}$$

$$v_{98} = 0.9 v_{97} + 0.1 v_{98}$$

$$v_{100} = 0.1 v_{100} + 0.9 \cancel{v_{99}} \xrightarrow{(0.1 v_{99} + 0.9 v_{98})} 0.1 v_{99} + 0.9 v_{97}$$

$$= 0.1 v_{100} + 0.1 \times 0.9 v_{99} + 0.1 \times (0.9)^2 v_{98} + 0.1 (0.9)^3 v_{97} + \dots$$



Why 10 days? $0.9^{10} \approx 0.35 \approx \frac{1}{e} = \underbrace{(1-\epsilon)^{\frac{1}{\epsilon}}}_{0.9} = \frac{1}{e}$

$$\text{e.g. } 0.98 \quad \underbrace{(1-\epsilon)^{\frac{1}{\epsilon}}}_{0.98} = 0.98^{\frac{1}{0.02}} = 0.98^{50} = \frac{1}{e} \therefore 50 \text{ days}$$

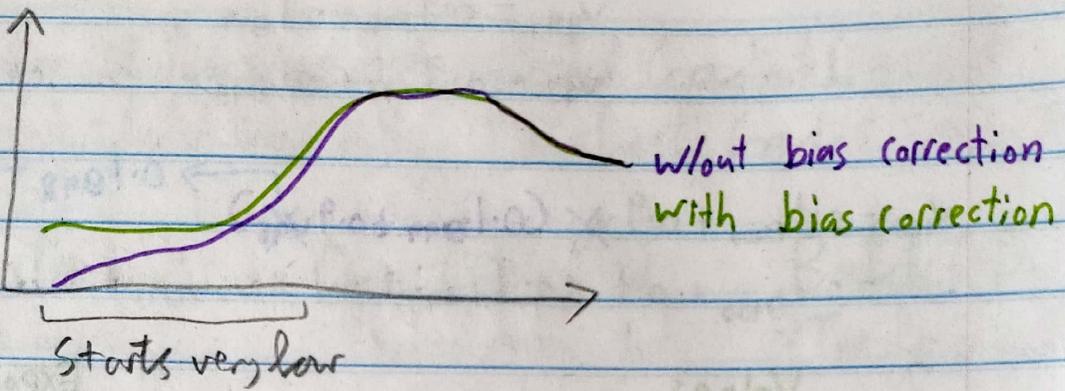
$$v_0 = 0$$

$$v_1 = \beta v_0 + (1-\beta) \theta_1$$

$$v_2 = \beta v_1 + (1-\beta) \theta_2$$

⋮

Bias Correction with Exponentially Weighted Averages



$$v_t = \beta v_{t-1} + (1-\beta) \theta_t$$

$$v_0 = 0$$

$$v_1 = 0.02 \theta_1$$

$$v_2 = 0.98 v_1 + 0.02 \theta_2$$

$$= 0.98 \times 0.02 \times \theta_1 + 0.02 \theta_2$$

$$= 0.0196 \theta_1 + 0.02 \theta_2$$

Not quite accurate

Introduce bias correction $\frac{v_t}{1-\beta^t}$

$$t=2: \quad 1-\beta^t = 1-(0.98)^2 = 0.0396$$

$$\frac{v_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$

numerator = denominator

When t is large $\rightarrow \beta^t = 0$

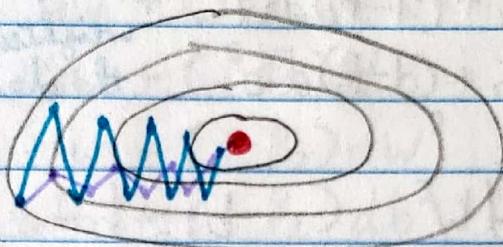
\rightarrow Help obtain better estimate! (during initial phase of learning)

Gradient Descent with Momentum

- Smoothes out the steps of gradient descent

GD

(GD w/ Momentum)



We want ↓ slower learning
↔ faster learning

- ↓: Averages out number so the average is ≈ 0
- ↔: Averages will still be pretty big

Momentum

On iteration t :

Compute dW, db on current mini-batch

$$V_{dw} = \beta V_{dw} + (1-\beta) dW$$

initialized to 0, same size as W

$$V_{db} = \beta V_{db} + (1-\beta) db$$

initialized to 0, same size as b

Friction Velocity Acceleration

$$W := W - \alpha V_{dw}$$

$$b := b - \alpha V_{db}$$

Hyperparameters: α, β β usually = 0.9

Average \uparrow over last ≈ 10 gradients

- In practice, bias correction is not really implemented. If $\beta = 0.9$, algorithm would have warmed up after 10 iterations & there isn't a need for bias estimate.

RMSprop

Root-Mean-Squared Prop

- On iteration t :

Compute dW, db on current mini-batch

$$S_{dW} = \beta_2 S_{dW} + (1 - \beta_2) dW^2 \quad \text{element-wise. Keeps exponentially weighted average of the squares of the derivatives}$$

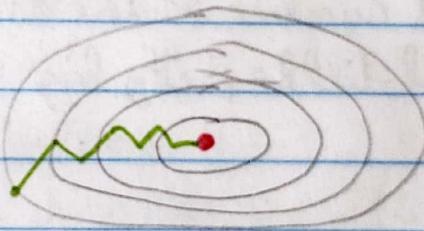
$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) d b^2$$

$$W := W - \alpha \frac{dW}{\sqrt{S_{dW} + \epsilon}}$$

$$b := b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}}$$

Ensures not dividing by zero

Just for illustration
in real life,
very high dimensions



What we are hoping for

$S_{dW} \rightarrow$ small, in order to speedup $\frac{dW}{\sqrt{S_{dW}}}$
 $S_{db} \rightarrow$ large, in order to slow down $\frac{db}{\sqrt{S_{db}}}$

Adam

Adaptive moment estimation

Adam = RMSprop + Gradient Descent with momentum

Algorithm →

dW 1st moment
 dW^2 2nd moment

- $V_{dw} = 0, S_{dw} = 0$ $V_{db} = 0, S_{db} = 0$

On iteration t :

Compute dW, db using current mini-batch

$$\begin{aligned} V_{dw} &= \beta_1 V_{dw} + (1-\beta_1) dW \\ V_{db} &= \beta_1 V_{db} + (1-\beta_1) db \end{aligned} \quad] \text{ "Momentum"}$$

$$\begin{aligned} S_{dw} &= \beta_2 S_{dw} + (1-\beta_2) dW^2 \\ S_{db} &= \beta_2 S_{db} + (1-\beta_2) db^2 \end{aligned} \quad] \text{ "RMSprop"}$$

$$V_{dw}^{\text{corrected}} = \frac{V_{dw}}{1-\beta_1^t} \quad V_{db}^{\text{corrected}} = \frac{V_{db}}{1-\beta_1^t}$$

Bias
Correction

$$S_{dw}^{\text{corrected}} = \frac{S_{dw}}{1-\beta_2^t} \quad S_{db}^{\text{corrected}} = \frac{S_{db}}{1-\beta_2^t}$$

$$W := W - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon} \quad b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon} \quad] \text{ Updates}$$

- Hyperparameters choice:

α : needs to be tuned

β_1 : 0.9

(dw) Commonly Used

β_2 : 0.999

(dw²) Adam's authors recommendation

ϵ : 10⁻⁸

Adam's authors recommendation

Learning Rate Decay

To learn fast in the beginning then converge to local minimum

$$\alpha = \frac{1}{1 + \text{decay_rate} * \text{epoch_num}} \cdot \alpha_0$$

e.g. $\alpha_0 = 0.2$

decay-rate = 1

Epoch	α
1	0.1
2	0.067
3	0.05
4	0.04
:	

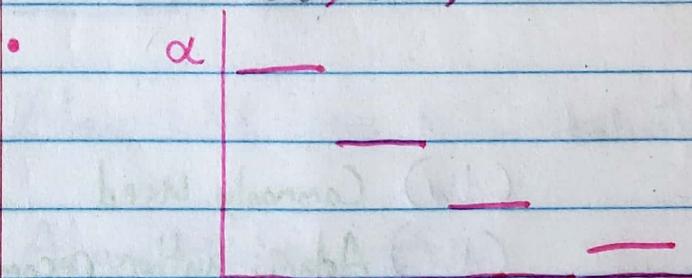
Other formulas

• $\alpha = 0.95^{\text{epoch_num}} \cdot \alpha_0$

exponential decay

• $\alpha = \frac{K}{\sqrt{\text{epoch_num}}} \cdot \alpha_0$

~~manually decay~~



discrete staircase

↑
mini batch
size

• Manually decay

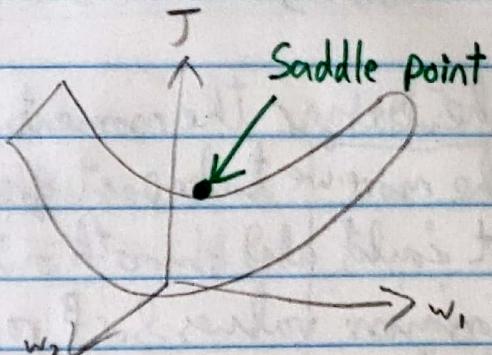
Local Optima Problem

In very high dimensions

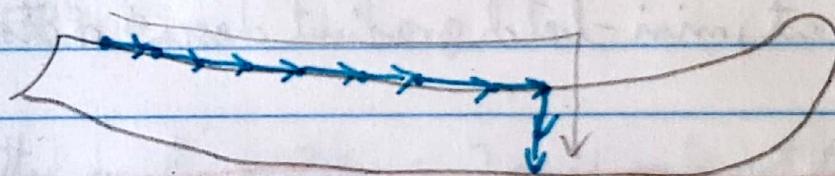
→ Could get stuck at
a saddle point

→ 0 gradients

(In low dimensions it gets stuck in local optimas)



→ Plateaus can make learning slow



→ This is where momentum, RMS prop, Adam etc help with speeding up the rate at which you could get down the plateau and move out of the plateau

Assignment

- When training set is large, SGD can be faster than BGD, but parameters will oscillate towards minimum rather than converge smoothly.

SGD

→ 3 for loops total

- Over the number of iterations

- Over the m training examples

- Over the layers to update params from $(W^{[l]}, b^{[l]})$ to $(W^{[l]}, b^{[l]})$

Mini-batch Shuffling & Partitioning are 2 required steps for mini-batch G-D

Momentum

- The larger the momentum β is, the smoother the update because the more we take past gradients into account. But if β is too big, it could also smooth out the updates too much.
- Common values for β range from 0.8 to 0.999. If you don't feel inclined to tune this, $\beta=0.9$ is reasonable by default.
- Momentum takes past gradients into account to smooth out the steps of gradient descent. It can be applied to batch gradient descent, mini-batch gradient descent or stochastic gradient descent.

Week 3 - Hyperparameter Tuning, Regularization and Optimization

Hyperparameter
Tuning

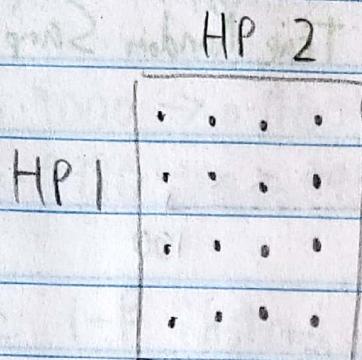
Tuning Process

- So far, here are our hyperparameters:

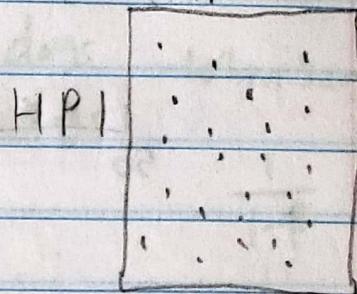
- α ← highest importance
- β
- $\beta_1, \beta_2, \epsilon$ $\sim 0.9, 0.999, 10^{-8}$ (Mr Andrew doesn't tune this)
- # layers
- # hidden units
- learning rate decay
- mini-batch size

- Try random values: Don't use a grid

e.g.



HP2

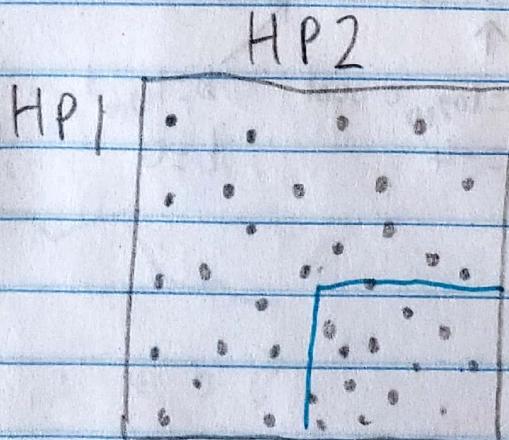


e.g. $HP1 \propto \epsilon$
 $HP2 \propto \epsilon$

α matters much more than ϵ . In grid, only got to try 4 values of α , whereas in random, we got to try 16 different values. In practice, there are much more hyperparameters.

- It's usually hard to know in advance which hyperparameters matter more than others for the application in hand

- Another common practice: Coarse to fine

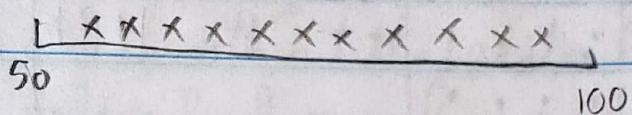


- that lie
- Found best samples in region
- Sample further values in that region (smaller square)
- Pick value that works best on train/dev set

Hyperparameter
Tuning

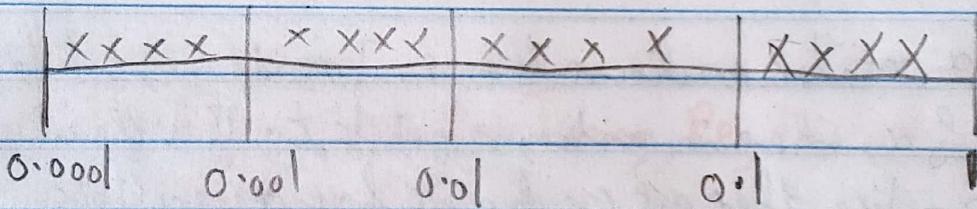
Using an appropriate scale to pick hyperparameters

- $n^{[x]} = 50, \dots, 100$ e.g. True Random Sampling works



#layers L: 2-4 2,3,4 True Random Sampling Works

- $\alpha = 0.001, \dots, 1$ Use Logarithmic Scale



$$r = -4 * \text{np.random.rand}() \quad \leftarrow r \in [-4, 0]$$
$$\alpha = 10^r \quad \leftarrow 10^{-4}, \dots, 10^0$$

Generally if $\alpha = 10^a \dots 10^b$

$$a = \log_{10} 0.0001 \quad b = \log_{10} 1$$
$$a = -4 \quad b = 0$$
$$r \in [a, b] \quad \leftarrow [-4, 0] \quad \alpha = 10^r$$

- $\beta = 0.9 \dots 0.999$ \leftarrow Bad idea to use for sampling
 \downarrow \downarrow
 10 days 1000 days

$0.9000 \rightarrow 0.9005$] 10 extra days

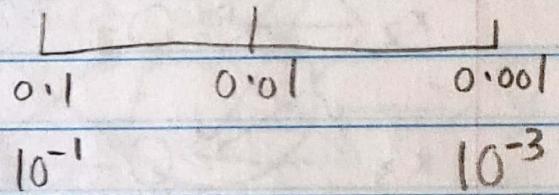
$0.9990 \rightarrow 0.9995$] 1000 extra days

when same increase

$$\frac{1}{1-\beta}$$

\therefore use $1-\beta$ instead

$$1-\beta = 0.1 \dots 0.001$$



$$r \in [-3, -1]$$

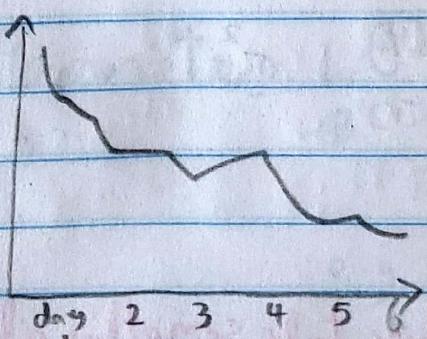
$$1-\beta = 10^r$$

$$\beta = 1 - 10^r$$

Hyperparameter
Tuning

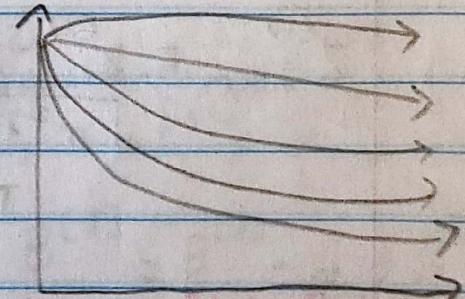
Hyperparameter tuning in practice

Babysitting one mode
(low on resources)



"Panda" approach

Training many models in
parallel



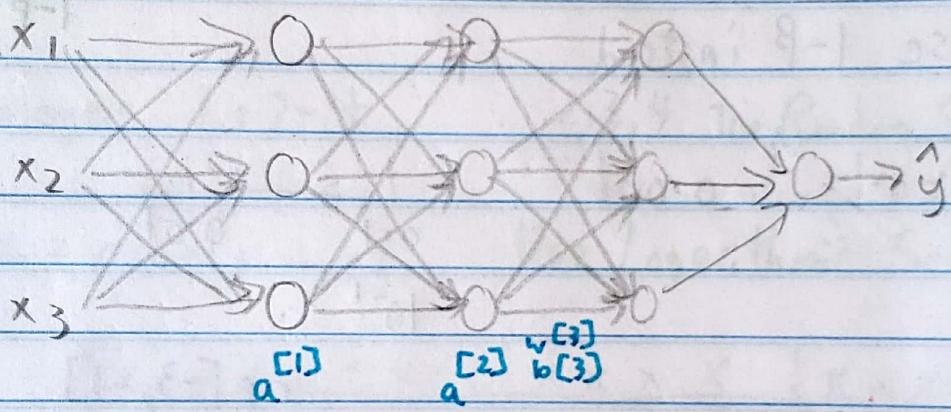
"Caviar" approach

Batch
Normalization

Normalizing activations in a network

- Normalizing inputs speeds up learning.

What about in deeper network?



Can we normalize $a^{[2]}$ so we can train $w^{[3]}, b^{[3]}$ faster?
→ Normalize $z^{[2]}$

- Implementing Batch Norm

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (z_i - \mu)^2$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\mu = 0 \quad \sigma^2 = 1$$

To avoid 0's

What if we want μ and σ^2 values to be something else?

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

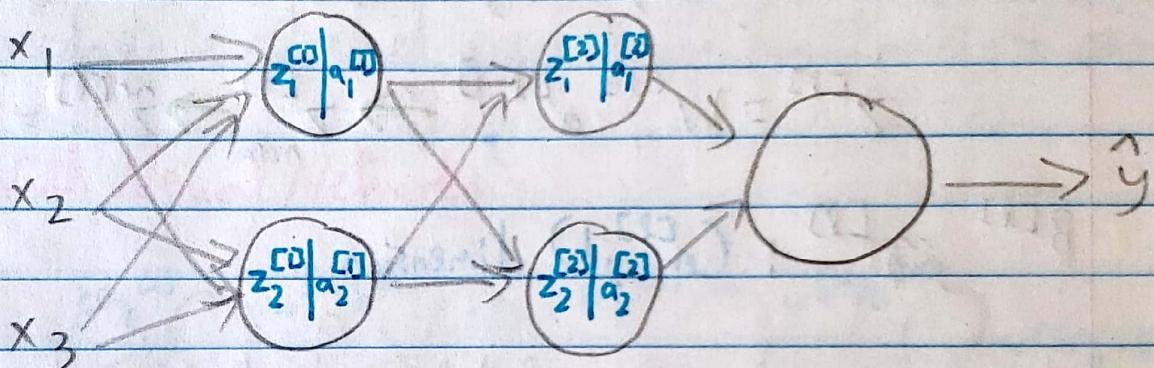
Learnable parameters

Use $\tilde{z}^{[l]}(i)$ instead of $z^{[l]}(i)$

→ To better take advantage of the non-linearity of the activation function

Batch
Normalization

Fitting Batch Norm into a neural network



$$\begin{aligned}
 & \text{If mini-batches } \leftarrow \\
 & x_{\Sigma 1}, x_{\Sigma 2}, \text{ all} \\
 & x_{\Sigma 3} \text{ separately} \\
 & x_{\Sigma 3} \\
 & \times \xrightarrow{w^{[0]}, b^{[0]}} z^{[0]} \xrightarrow[\text{Batch Norm (BN)}]{\beta^{[0]}, \gamma^{[0]}} \tilde{z}^{[0]} \rightarrow a^{[0]} = g^{[0]}(\tilde{z}^{[0]}) \\
 & a^{[0]} \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\text{BN}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \rightarrow a^{[1]} \rightarrow \dots
 \end{aligned}$$

Parameters: $W^{[0]}, b^{[0]}, W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}$
 $\beta^{[0]}, \gamma^{[0]}, \beta^{[1]}, \gamma^{[1]}, \dots, \beta^{[L]}, \gamma^{[L]}$

$$\text{e.g. } d\beta^{[l]} \quad \beta^{[l]} = \beta^{[l]} - \alpha d\beta^{[l]}$$

- If using batch-norm: $b^{[l]}$ can be eliminated

$$z^{[l]} = w^{[l]} a^{[l-1]} \left[+ b^{[l]} \right]$$

\uparrow

subtracting and adding constants to find mean ends up equaling 0 for batch-norm

Parameters: $w^{[l]}, \cancel{b^{[l]}}, \beta^{[l]}, \gamma^{[l]}$

$$z^{[l]} = w^{[l]} a^{[l-1]} \rightarrow z_{\text{norm}}^{[l]} \rightarrow \tilde{z}^{[l]} = \gamma^{[l]} z_{\text{norm}}^{[l]} + \beta^{[l]}$$

This now controls bias terms

- $\beta^{[l]}, \gamma^{[l]} : (n^{[l]}, 1)$ dimension

- Implementing mini-batch gradient descent

for $t = 1, \dots, \text{num Mini-Batches}$

Compute forward pass on $X^{\{t\}}$

In each hidden layer, use BN to replace $z^{[l]}$ with $\tilde{z}^{[l]}$
 Use backprop to compute $dW^{[l]}, d\beta^{[l]}, d\gamma^{[l]}$

Update params

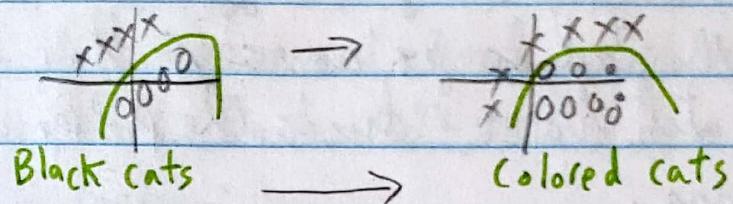
$$\begin{aligned} W^{[l]} &:= W^{[l]} - \alpha dW^{[l]} \\ \beta^{[l]} &:= \beta^{[l]} - \alpha d\beta^{[l]} \\ \gamma^{[l]} &:= \gamma^{[l]} - \alpha d\gamma^{[l]} \end{aligned}$$

Works with momentum, RMSprop, Adam....

Batch Normalization

Why does Batch Norm work?

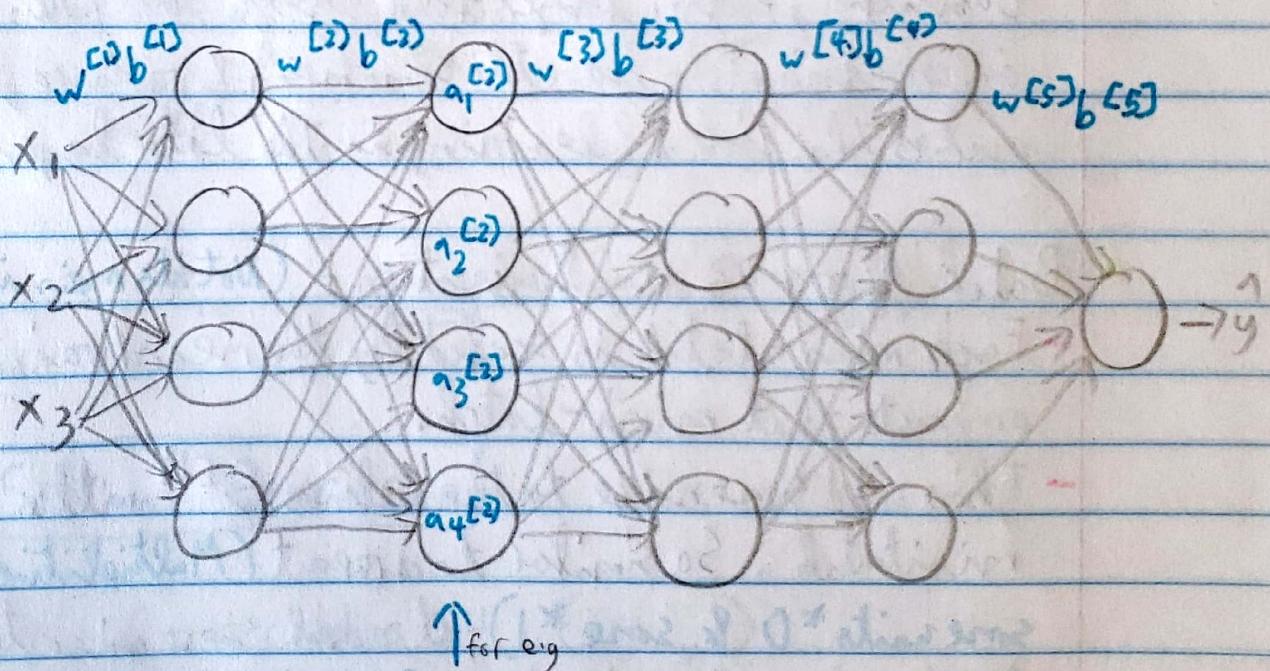
e.g.



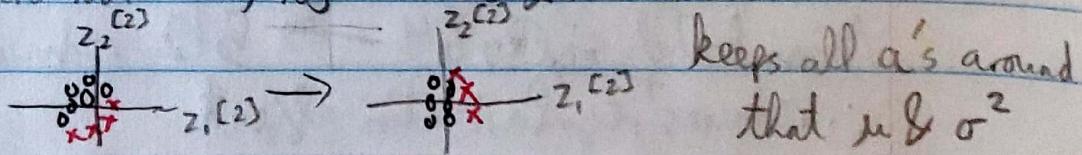
"Covariate shift" $x \rightarrow y$

If the distribution of x changes, then the algorithm needs to be retrained, even if the function remains unchanged

With Neural Networks



- $w^{[0]}, b^{[0]}$ & $w^{[1]}, b^{[1]}$ keep on changing, so the $a^{[1]}$ values keep on changing \therefore suffers from covariate shift problem
- Batch norm reduces the distribution that the $a^{[1]}$ values shift around. e.g. if mean 0 & variance 1



- Limits the amount to which updating the parameters in the earlier layers affect the distribution of values that the 3rd layer now sees & \therefore has to learn on. So Batch Norm reduces the problem of the input values changing.
- Causes these values to become more stable so the later layers in the network have a firm ground to stand on.
- Weakens coupling between earlier layers & later layers allows each layer to learn more independently, speeds up learning of the whole network
- Basically, earlier layers don't shift that much as they are constrained to have the same mean & variance, which makes the job of learning in the later layers easier

• Batch Norm as regularization : (Not its main intent though)

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.
- This adds some noise to the values $z^{[l]}$ within that minibatch. So similar to dropout (**Multiplicative noise as some units * 0 & some * 1**), it adds some noise to each hidden layer's activations (**Additive noise due to mean & multiplicative due to variance**).
- This has a slight regularization effect. (Doesn't rely on one hidden unit by downstream hidden units)
- Mini batch size $\approx 64 \xrightarrow{\text{increase}} 512$ reduces Batch Norm noise \therefore reduces regularization effect

Batch
Normalisation

Batch Norm at test time

- During training time, μ and σ^2 computed over all examples in 1 mini-batch. Each mini-batch has a different μ and σ^2 .
 - But during testing:
- Estimate μ, σ^2 using exponentially weighted average across all mini-batches

$$X^{[1]}, X^{[2]}, X^{[3]}, \dots$$

$$\begin{matrix} \downarrow & \downarrow & \downarrow \\ \mu^{[1][l]} & \mu^{[2][l]} & \mu^{[3][l]} \\ \theta_1 & \theta_2 & \theta_3 \end{matrix} \rightarrow \underline{\mu}$$

$$\sigma^{2[1][l]} \rightarrow \sigma^2$$

$$\therefore z_{\text{norm}} = \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad \hat{z} = \gamma z_{\text{norm}} + \beta$$

Multi-class
Classification

Softmax Regression

- e.g. $C = \# \text{ classes} = 4 \quad (0, \dots, 3)$

$$n^{[L]} = 4 = C$$

$$x \rightarrow 1/11 \rightarrow \begin{array}{|c|c|c|c|} \hline & 0 & 0 & 0 & 0 \\ \hline \end{array} \rightarrow \hat{y} (4, 1)$$

$$\uparrow n^{[L]} = 4$$

$$\hat{y} \text{ sum} = 1$$

- Activation function: $t = e^{z^{[L]}}$ $(4,1)$ dim

$$(g) \quad a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{i=1}^4 t_i}, \quad a_i^{[L]} = \frac{t_i}{\sum_{i=1}^4 t_i} \quad (4,1) \text{ dim}$$

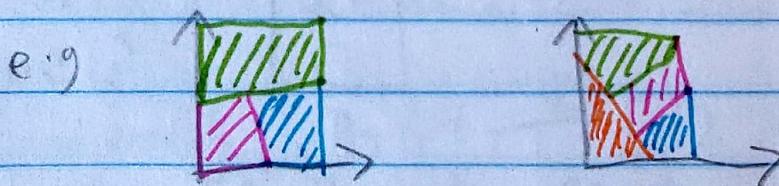
e.g. $z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$ $t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$

$$\sum_{i=1}^4 t_i = 176.3 \quad a^{[L]} = \frac{t}{176.3}$$

\circ	$\rightarrow e^5/176.3 = 0.842$
\circ	$\rightarrow e^2/176.3 = 0.042$
\circ	$\rightarrow e^{-1}/176.3 = 0.002$
\circ	$\rightarrow e^3/176.3 = 0.114$

$a^{[L]} = g^{[L]}(z^{[L]})$
 $(4,1)$

- The decision boundary between classes is linear



Multi-class
Classification

Training a softmax classifier

- Softmax $\begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$ Hard max $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

- Softmax regression generalizes logistic regression to C classes
 - If $C=2$, softmax reduces to logistic regression

- Loss function $f(\hat{y}, y) = -\sum_{j=1}^4 y_j \log \hat{y}_j$

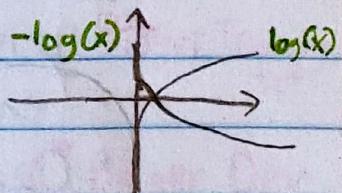
e.g. $y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ $\hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$

$$f(\hat{y}, y) = -y_2 \log \hat{y}_2 = -\log \hat{y}_2$$

↑
To minimize this ↑
we need to maximize this

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}] \quad (4, m)$$

$$\hat{Y} = [\hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(m)}] \quad (4, m)$$



- Backprop/Gradient Descent

$$X \rightarrow \dots \rightarrow \boxed{\quad} \rightarrow \hat{y}$$

$$L_2^{[L]} \rightarrow a^{[L]}$$

$$d_z^{[L]} = \hat{y} - y \quad (4, 1)$$

$$\left(\frac{\partial J}{\partial z^{[L]}} \right)$$

Programming Frameworks

- Caffe / Caffe2
- NTK
- DL4J
- Keras
- Lasagne
- mxnet
- PaddlePaddle
- TensorFlow
- Theano
- Torch

Choosing deep learning frameworks:

- Ease of programming (development and deployment)
- Running speed
- Truly open (open source with good governance)

Tensorflow

- Coding session

- Two main object classes:

- Tensors (Variables, Placeholders, ...)
- Operations (tf.matmul, tf.add,)

Structuring Machine Learning Projects

Week 1 - Why ML Strategy

Intro

Orthogonalization

- Eg. radio tuning: if the radio sound was too noisy, we want 1 knob to tune in order to adjust the frequency of the radio channel. We don't want to have to carefully adjust five different knobs that all affect different things.

- Chain of assumptions in ML

If it doesn't

Fit training set well on cost function ↗ bigger network

Adam



Fit dev set well on cost function ↗ Regularization

Bigger Training set



Fit test set well on cost function ↗ Bigger Dev set



Performs well in real world ↗ Change Dev set or cost function or metric

Setting up your goal

Single Number Evaluation Metric

- Idea
Experiment
Code

Classifier	Precision	Recall	F1 Score
A	95%	90%	92.4%
B	98%	85%	91.0%

(well defined)

→ Derive + Single real number evaluation metric
Speeds up iterating

e.g.	Algorithm	US	China	India	Other	Average
	A	3%	7%	5%	9%	6%
	B	5%	6%	5%	10%	6.5%
	C	2%	3%	4%	5%	3.5%
	D	5%	8%	7%	2%	5.25%
	E	4%	5%	2%	4%	3.75%
	F	7%	11%	8%	12%	9.5%

Satisficing and Optimizing Metrics

- e.g. Classifier

Optimizing Accuracy Satisficing Running Time

A	90%	80ms
B	92%	95ms
C	95%	150ms

← Best Classification

→ Maximize accuracy subject to runningTime $\leq 100\text{ms}$

- N metrics :
 - 1 optimizing
 - N-1 satisfying
- e.g. voice controlled devices
 - Maximize accuracy subject to ≤ 1 false pos every 24 hours
 - Optimizing: Accuracy
 - Satisficing: False Positives

setting up
your goal

Train / dev / test distributions

- Make sure dev and test set data come from the same distribution
 - Distributions e.g. Regions, high/low income individuals
- Choose a dev set and test set to reflect data you expect to get in the future and consider important to do well on.

setting up
your goal

Size of the dev and test sets

- Earlier era 70/30 or 60/20/20 for $< 10K$ data
Now 98/1/1 for $> 1000K$ data
- Set your test set to be big enough to give high confidence in the overall performance of your system.

setting up
your goal

When to change dev/test sets and metrics

- Metric: classification error

Algorithm A: 3% error → Has "unacceptable" images

Algorithm B: 5% error → Doesn't

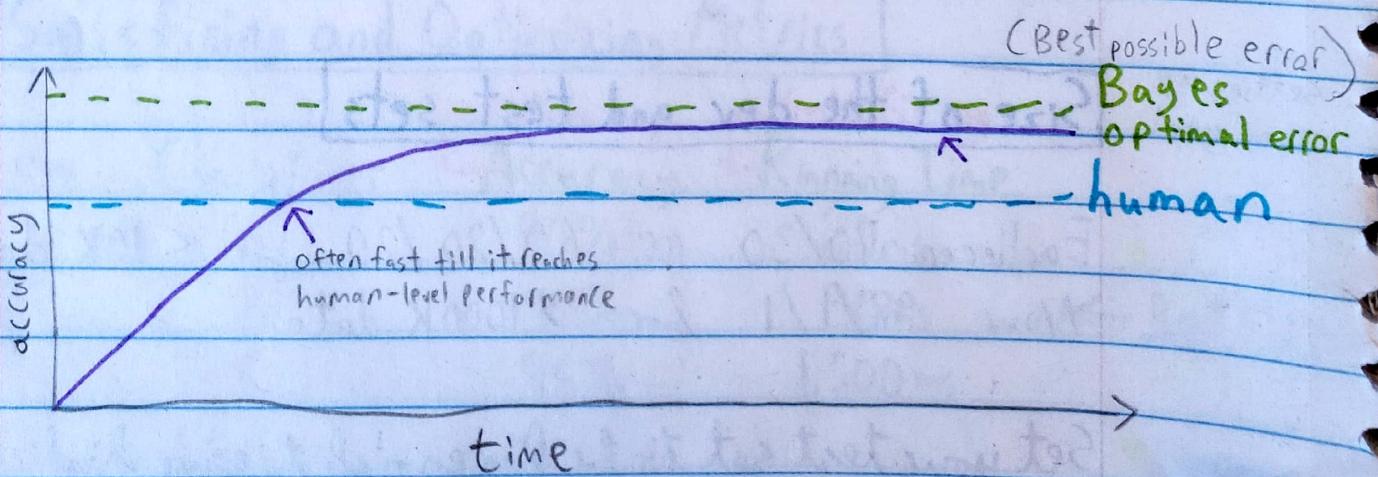
- Metric + Dev: Prefer A
- You/users: Prefer B

$$\text{Error} = \frac{1}{\sum w^{(i)}} \cancel{\frac{1}{m_{dev}}} \sum_{i=1}^{m_{dev}} w^{(i)} \delta \{y_{pred}^{(i)} \neq y^{(i)}\}$$

$$\rightarrow w^{(i)} = \begin{cases} e.g. 1 & \text{if acceptable image} \\ e.g. 10 & \text{if unacceptable image} \end{cases}$$

Comparing
to human
level performance

Why human-level performance?

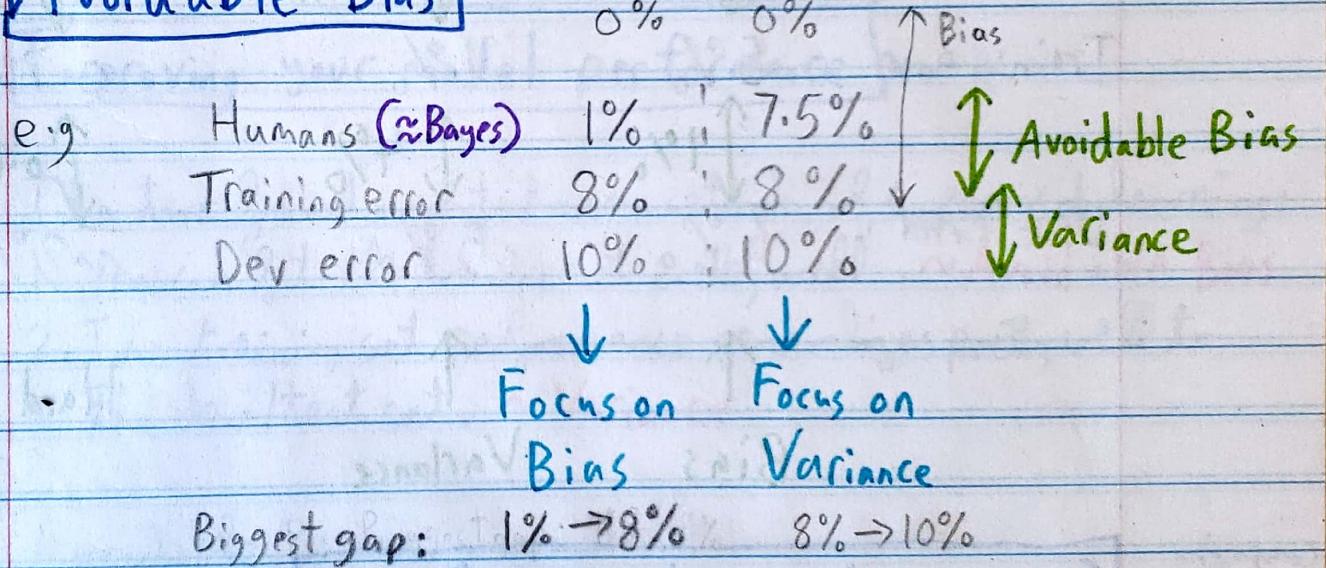


Why progress often slows down:

- Human-level performance for many tasks not far off from Bayes optimal error
- If ML is worse than humans, we can:
 - Get labeled data from humans
 - Gain insight from manual error analysis
 - Better analysis of bias/variance

Comparing to
human-level
performance

Avoidable Bias



Human-level error acts as a proxy/estimate for Bayes error

Comparing to
human-level
performance

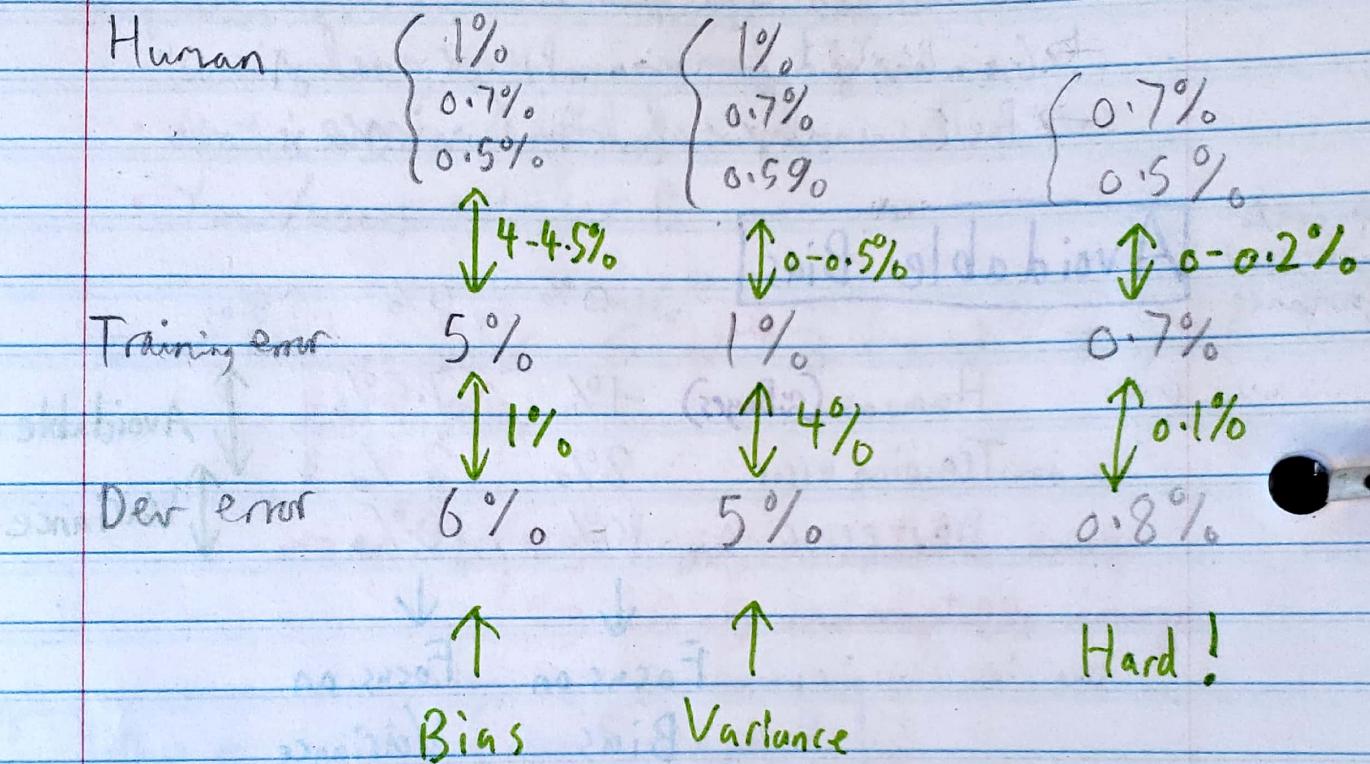
Understanding human-level performance

- Human-level error as a proxy for Bayes error

e.g. Medical Image Classification example:

- (a) Typical human 3% error
 - (c) Experienced doctor 0.7% error
 - (b) Typical doctor 1% error
 - (d) Team of experienced doctors 0.5% error
- What is "human-level" error? Bayes error $\leq 0.5\%$

- As machine learning algorithm is doing really well, ML project progress gets harder. Pro human-level performance becomes harder to tease out the bias and variance effects



Comparing
to human
level performance

Surpassing human-level performance

Team of humans	0.5%
Training error	0.3%
Dev error	0.4%

What is avoidable bias? Very unclear. It's hard for us to now know if we should focus on teasing out bias or variance effects.

- Online advertising
 - Product recommendations
 - Logistics (predicting transit time)
 - Loan approvals
- These all learn from structural data. Nothing to do with natural perception, which we humans are really good at.
- Lots of data

Improving your model performance

- Comparing to
human-level
performance
- The two fundamental assumptions of supervised learning:
 1. You can fit the training set pretty well ~ **Avoidable Bias**
 2. The training set performance generalizes pretty well to the dev/test set. ~ **Variance**
 - Avoidable Bias:
 - Train bigger model
 - Train longer / better optimization algorithms
 - NN architecture / hyperparameters search (RNN, CNN etc)
 - Variance:
 - More data
 - Regularization (L2, dropout, data augmentation)
 - NN architecture / hyperparameters search

Week 2 - ML Strategy (2)

Error Analysis

Error Analysis

- Evaluate multiple ideas in parallel. Ideas:
 - Fix pictures of dogs being recognized as cats.
 - Fix great cats (lions, panthers, etc...) being misrecognized
 - Improve performance on blurry images

Image	Dog	Great Cats	Blurry	Comments
1	✓			Pitbull
2			✓	
3		✓	✓	Rainy day at 300
:				
% of total	8%	43%	61%	<p>It's okay to start adding new filters or categories</p> <p>Can have separate teams working on them</p>

Error Analysis

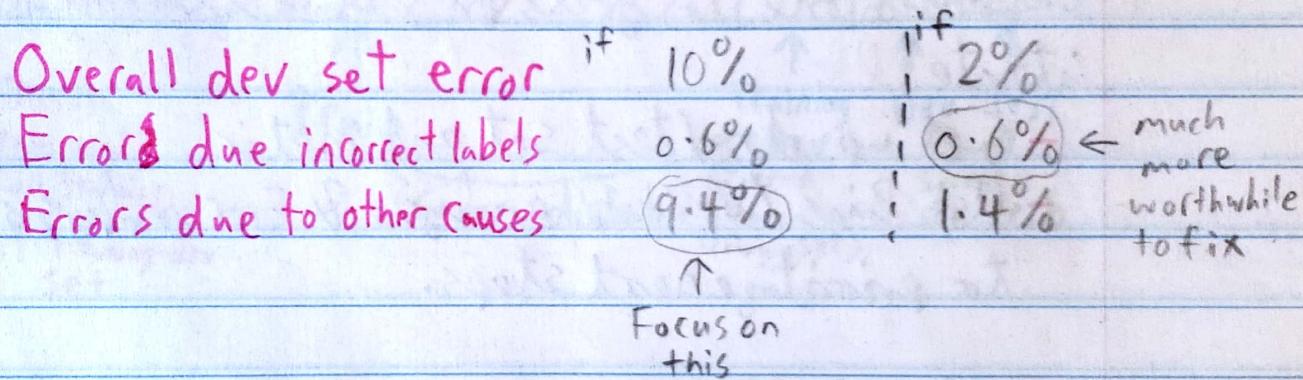
Cleaning up incorrectly labeled data

- DL algorithms are quite robust to random errors in the training set.
 - Less robust to systematic errors.

- I'll Devset

5

Image	Dog	Cat	Blurry	Incorrectly Labeled	Comments
98	X	-		✓	Labels missed cat in background
99	-	X			
100		X		✓	Drawing of cat, not real
% of total	8%	43%	61%	6%	



- Goal of devset is to help you select between two classifiers A & B. To know which one to trust, it's good to first eliminate incorrect labels.

- Correcting incorrect dev/test set examples:

- Apply same process to your dev and test sets to make sure they continue to come from the same distribution
- Consider examining examples your algorithm got right as well as ones it got wrong to avoid biased advantages.
- Train and dev/test data may now come from slightly different distributions.

Error Analysis

Build your first system quickly, then iterate

- Speech recognition example. Algorithm could be built to be more robust to:

- Noisy background
- Accented speech
- Far from microphone
- Young children's speech
- Stuttering

But which one to focus on? Build system quickly, then iterate!

- Set up dev/test set and metric
- Use Bias/Variance analysis & Error analysis to prioritize next steps.

- Less strongly applicable to domains where experts are or in domains with huge literature available.

Mismatched training and dev/test set

Training and Testing on different distributions

Cat app example

- Data from webpages (web crawl) $\approx 200,000$
- Data from mobile app (users what we care about) $\approx 10,000$

Good option

Train	Dev	Test
200,000 web + 5,000 app	2,500 app	2,500 app

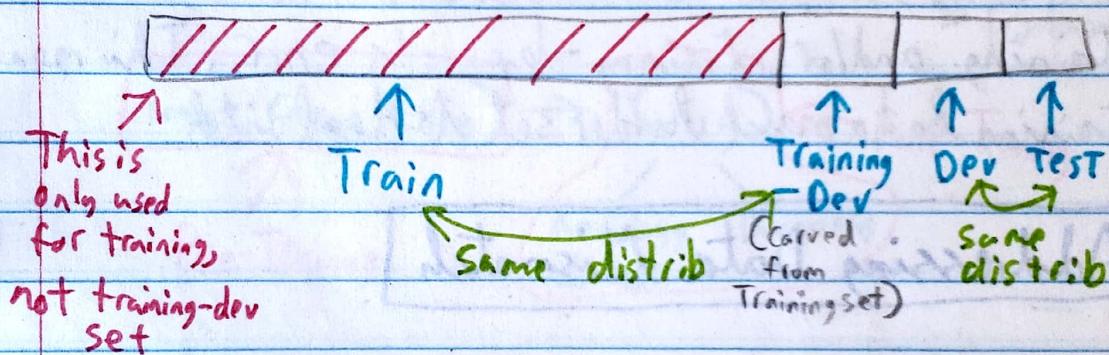
Optimizes on the mobile app data as it's the whole dev set

Mismatched
Training and
dev/test set

Bias and Variance with mismatched data distributions

- Example training error 1% ↓ 9% Is it a variance problem dev error 10% or mismatched distributions

→ Create Training-Dev set: Same distribution as training set but not used for training



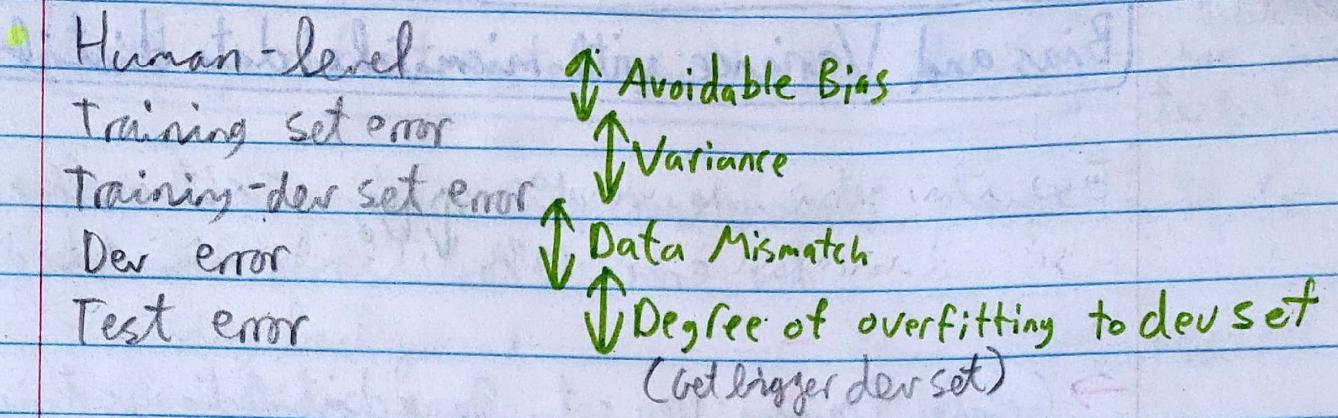
For example Training error 1% ↑ Variance error!
Training-dev err 9%
dev error 10%

Training error 1%
Training-dev err 1.5% ↓ Data mismatch problem!
dev error 10%

Key
Quantities
to look at

Human level	0%	↓ Avoidable bias
Training error	10%	
Training-dev error	11%	
dev error	12%	

Could be a combination of both



- It could happen that dev/test error is lower than training and/or training-dev sets error. This means training had a much harder set to deal with

Mismatched
train &
dev/test
set

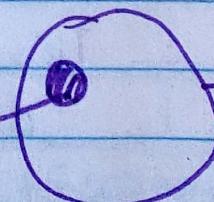
Addressing Data Mismatch

- Carry out manual error analysis to try to understand differences between training and dev/test sets. e.g. noisy data
 - Make training data more similar; or collect more data similar to dev/test sets e.g. simulate noisy in-car data
- Artificial Data Synthesis (When you need more specific data)

e.g. Man talking + Car Noise = Synthesized in-car audio

Need to be careful

Synthesized



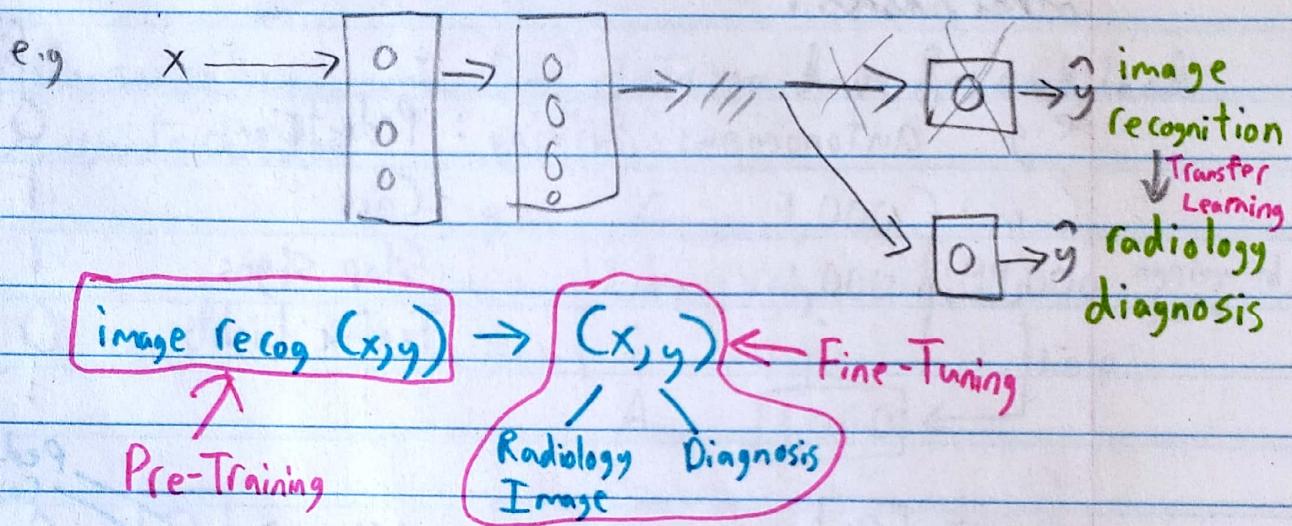
All cars or
all audio in car
(or all possibilities)

Learning from
multiple tasks

Transfer Learning

(Done Sequentially)

- Take knowledge the neural network has learned from one task and apply that knowledge to a separate task.



- Transfer learning makes sense when you have a lot of data for the problem you're transferring from and usually relatively less data for the problem you're transferring to.

e.g. image recognition 1,000,000 data
radiology diagnosis 100 data

- When task A and B have the same input x . (both are images)
- You have a lot more data for Task A than Task B.
- Low level features from A could be helpful for learning B.

- Used very often

Learning from
multiple tasks

Multi-task learning

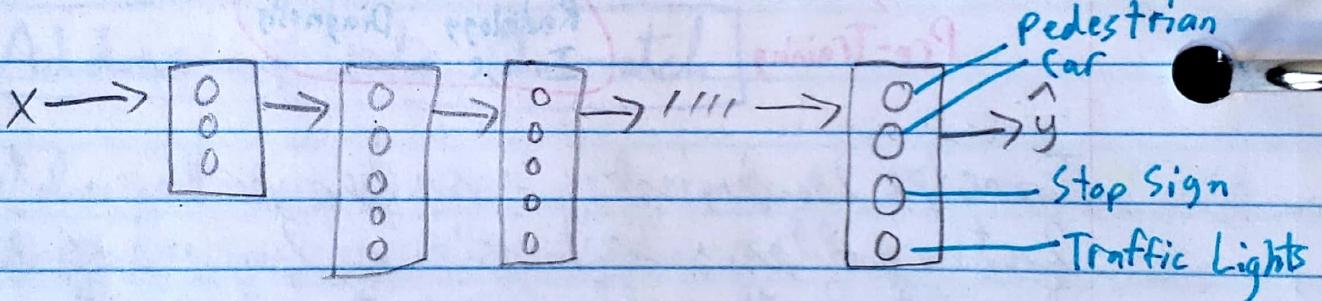
- Simultaneously have one neural network do several things at the same time, and then each of these tasks helps all the other tasks.

e.g. autonomous driving: Pedestrians,

$y^{(i)}$	(4, 1)
0)
1)
1)
0)

Stop Signs

Traffic Lights



Loss: $\hat{y}^{(i)}$
 $(4, 1)$

$$\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^4 l(\hat{y}_j^{(i)}, y_j^{(i)})$$

main difference $-y_j^{(i)} \log \hat{y}_j^{(i)} - (1-y_j^{(i)}) \log (1-\hat{y}_j^{(i)})$

~Image can have multiple labels

[Unlike softmax regression where one image can have only 1 label]

- If some of the earlier features in neural network can be shared between these different types of objects, then you find that training one neural network to do few things result in better

performance than training four completely separate neural networks to do the four tasks separately.

- Works even when some $y_i^{(1)}$ values are unlabeled -
- Makes sense when amount of data you have for each task is quite similar

e.g. $\begin{array}{l} \left\{ \begin{array}{ll} A_1 & 1,000 \\ A_2 & 1,000 \\ \vdots & \vdots \\ A_{100} & 1,000 \end{array} \right\} \\ \text{99,000 aggregated data} \end{array}$ help

- Makes sense when you can train a big enough neural network to do well on all the tasks.
→ It actually hurts performance when network not big enough.

End-to-end
deep learning

What is end-to-end deep learning?

- Traditional Pipeline: audio $\xrightarrow{\text{MFCC}}$ features $\xrightarrow{\text{ML}}$ Phonemes \rightarrow Words \rightarrow Transcript
- End-to-end learning: audio $\xrightarrow{\text{works really well when you have lots of data}}$ transcript
→ Sometimes breaking things into 2 sub-problems works better than pure end-to-end DL approach
- Nowadays works well with machine translation e.g. English \rightarrow French, not so well with estimating child's age from hand X-ray so
(lack of data) Image \rightarrow Bones \rightarrow Age instead of Image \rightarrow Age

Whether to use end-to-end deep learning

Pros:

- Let the data speak (when you have enough data)... it does its own $x \rightarrow y$ mapping rather than being forced to human misconceptions
- Less hand-designing of components needed

Cons:

- May need large amount of data
- Excludes potentially useful hand-designed components

Key question: Do you have sufficient data to learn a function of the complexity needed to map x to y ?

$$\begin{matrix} \text{input} & \text{output} \\ \text{end} & \text{end} \\ x & \downarrow & y \end{matrix}$$

Coursera
Deep
Learning
Specialization
Course #4

Coursera: Convolutional Neural Networks

Week 1 : Foundations of CNNs

★ Vertical edge detection

$$\begin{array}{|c|c|c|c|c|c|} \hline
 3 & 0 & 1 & 2 & 7 & 4 \\ \hline
 1 & 5 & 8 & 9 & 3 & 1 \\ \hline
 2 & 7 & 2 & 5 & 1 & 3 \\ \hline
 0 & 1 & 3 & 1 & 7 & 8 \\ \hline
 4 & 2 & 1 & 6 & 2 & 8 \\ \hline
 2 & 4 & 5 & 2 & 3 & 9 \\ \hline
 \end{array}$$

\star

"convolution"

6×6
 $n \times n$

$*$

"convolution"

$$\begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array}$$

3×3
Filter/Kernel
 $f \times f$

$$\begin{array}{|c|c|c|c|} \hline
 -5 & -4 & 0 & 8 \\ \hline
 -10 & -2 & 2 & 3 \\ \hline
 0 & -2 & 4 & -7 \\ \hline
 -3 & -2 & 3 & -16 \\ \hline
 \end{array}$$

4×4

$n-f+1 \times n-f+1$

e.g.

$$\begin{array}{|c|c|c|c|c|c|} \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 10 & 10 & 10 & 0 & 0 & 0 \\ \hline
 \end{array}$$

$*$

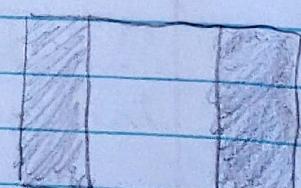
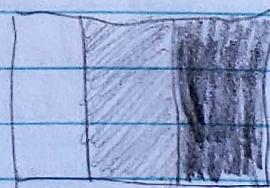
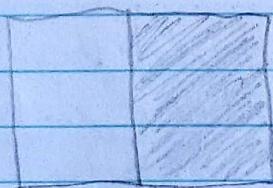
$$\begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array}$$

Filter/Kernel

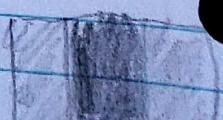
$$\begin{array}{|c|c|c|c|} \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 0 & 30 & 30 & 0 \\ \hline
 \end{array}$$

4×4

6×6



if



↑
Indicates high likelihood of v edge in that area

• Horizontal Edge Filter/Kernel

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$

• Different weights & filters could be used for vertical & horizontal edge detections

→ Sometimes those numbers can be learned

→ Can be also used to detect edges at angles e.g. $45^\circ, 70^\circ$

- ★ • To ensure numbers on the edge of an image gets factored in as much as numbers in other places of the image, we introduce **padding**
- Solves shrinking output
 - Solves throwing away lots of info from edges

e.g. padding = 1

$$\begin{array}{ccccccc} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & & & & & & \\ 0 & & & & & & \\ 0 & & & & & & \\ 0 & & & & & & \\ 0 & & & & & & \\ 0 & & & & & & \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{array}$$

$$* \quad 3 \times 3 = 4 \times 4 \rightarrow 6 \times 6$$

$$\begin{aligned} n+2p-f+1 \\ 6+2-3+1 \\ = 6 \end{aligned}$$

$$6 \times 6 \rightarrow 8 \times 8$$

• "Valid Convolution": $n \times n * f \times f \rightarrow n-f+1 \times n-f+1$
 $6 \times 6 \quad 3 \times 3 \quad 4 \times 4$

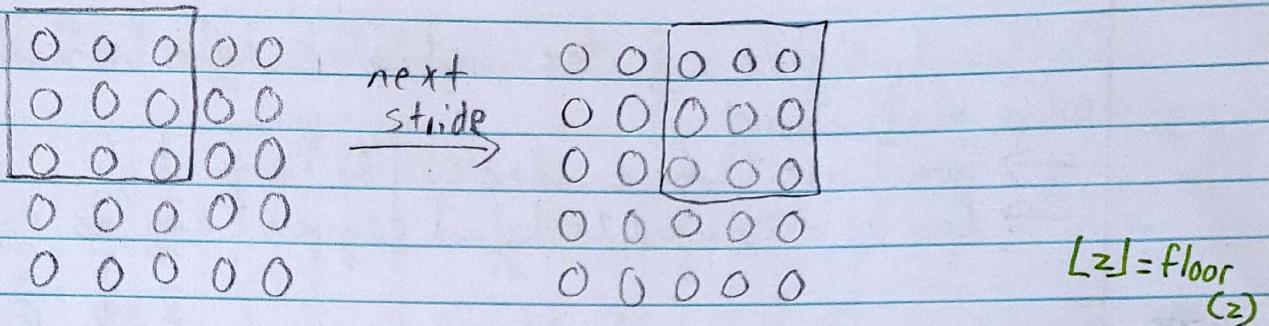
• "Same Convolution": Pad so that output size is the same as the input size

$$p = \frac{f-1}{2} \quad f \text{ is usually odd}$$

$$\begin{aligned} f = 3 \times 3 & \quad p = 1 \\ f = 5 \times 5 & \quad p = 2 \end{aligned}$$

★ Stided Convolution

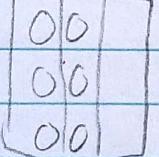
e.g. stride = 2



$n \times n$ * $f \times f$
padding p strides

$$\left\lceil \frac{n+2p-f+1}{s} \right\rceil \times \left\lceil \frac{n+2p-f+1}{s} \right\rceil$$

e.g. 7×7 3×3
 $p=0$ $s=2$ $\frac{7+0-3+1-3}{2} = 3$ $\therefore 3 \times 3$

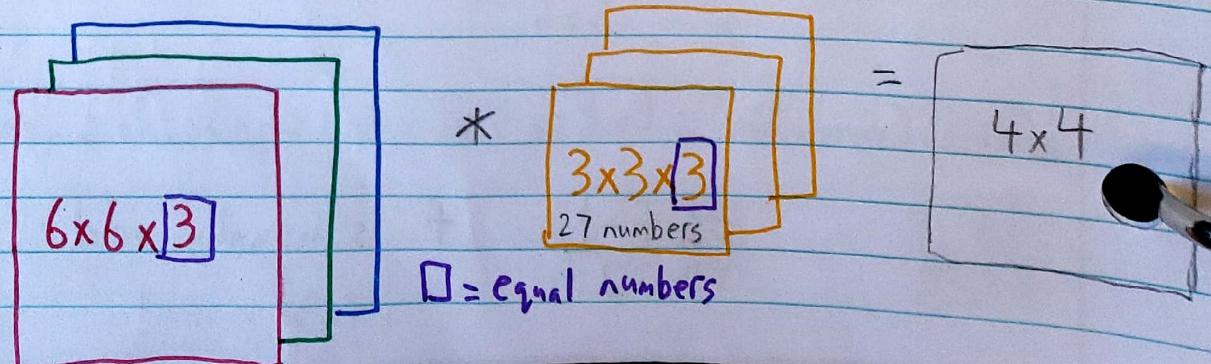
if  x filter must lie completely in image

Convolution in maths is to

$$\begin{bmatrix} 3 & 4 & 5 \\ 1 & 0 & 2 \\ -1 & 9 & 7 \end{bmatrix} \rightarrow \begin{bmatrix} 7 & 9 & -1 \\ 2 & 0 & 1 \\ 5 & 4 & 3 \end{bmatrix}$$

What we do in ML is called cross-correlation in maths, but ML results are the same so flipping isn't usually done

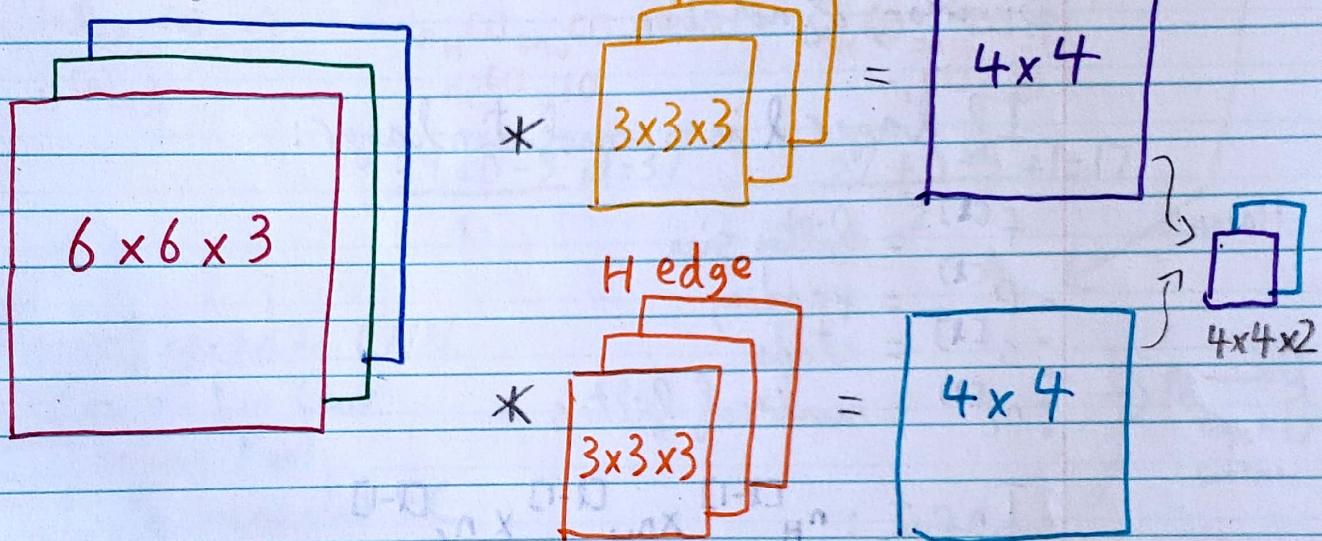
★ Convolutions on RGB images (volumes)



e.g. if we want to detect only Red v edges

$$R \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} \quad G \& B \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

- Multiple filters



$$n \times n \times n_c \times f \times f \times n_c \rightarrow n-f+1 \times n-f+1 \times n_c' \quad \# \text{filters}$$

$$6 \times 6 \times 3 \quad 3 \times 3 \times 3 \quad 4 \times 4 \times 2$$

- Example of a layer

(Using above figures)

$$a^{[0]}$$

$$w^{[1]}$$

$$z^{[1]} = w^{[1]} a^{[0]} + b^{[1]}$$

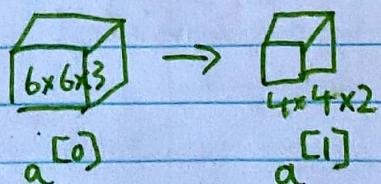
$$a^{[1]} = g(z^{[1]})$$

in $^{[]}$ indicates layer

$$\underbrace{g}_{= z^{[1]}} \underbrace{\underbrace{w^{[1]} a^{[0]}}_{= z^{[1]}}} = \text{ReLU} \left(\underbrace{[4 \times 4]}_{\text{ReLU}} + b_1 \right) \rightarrow [4 \times 4]$$

$$\text{ReLU} \left(\dots + b_2 \right) \rightarrow [4 \times 4]$$

$$= \underbrace{\begin{bmatrix} & \\ & \end{bmatrix}}_{4 \times 4 \times 2} \underbrace{a^{[1]}}_{a^{[1]}}$$



• Number of parameters in one layer

e.g. 10 filters $3 \times 3 \times 3$ in one NN layer params?

$$3 \times 3 \times 3 + 1 = 28 \quad \begin{matrix} \uparrow \\ \text{Bias} \end{matrix} \quad 28 \times 10 = 280 \text{ /}$$

• Summary of notation

If layer ℓ is a convolution layer:

$f^{[\ell]}$ = filter size

$p^{[\ell]}$ = padding

$s^{[\ell]}$ = stride

$n_c^{[\ell]}$ = number of filters

Input: $n_H^{[\ell-1]} \times n_w^{[\ell-1]} \times n_c^{[\ell-1]}$

Output: $n_H^{[\ell]} \times n_w^{[\ell]} \times n_c^{[\ell]}$

$$n_{\text{HorW}}^{[\ell]} = \left\lfloor \frac{n^{[\ell-1]} + 2p^{[\ell]} - f^{[\ell]}}{s^{[\ell]}} + 1 \right\rfloor$$

Each filter is: $f^{[\ell]} \times f^{[\ell]} \times n_c^{[\ell-1]}$

Activations: $a^{[\ell]} \rightarrow n_H^{[\ell]} \times n_w^{[\ell]} \times n_c^{[\ell]}$ $A \rightarrow m \times a^{[\ell]}$

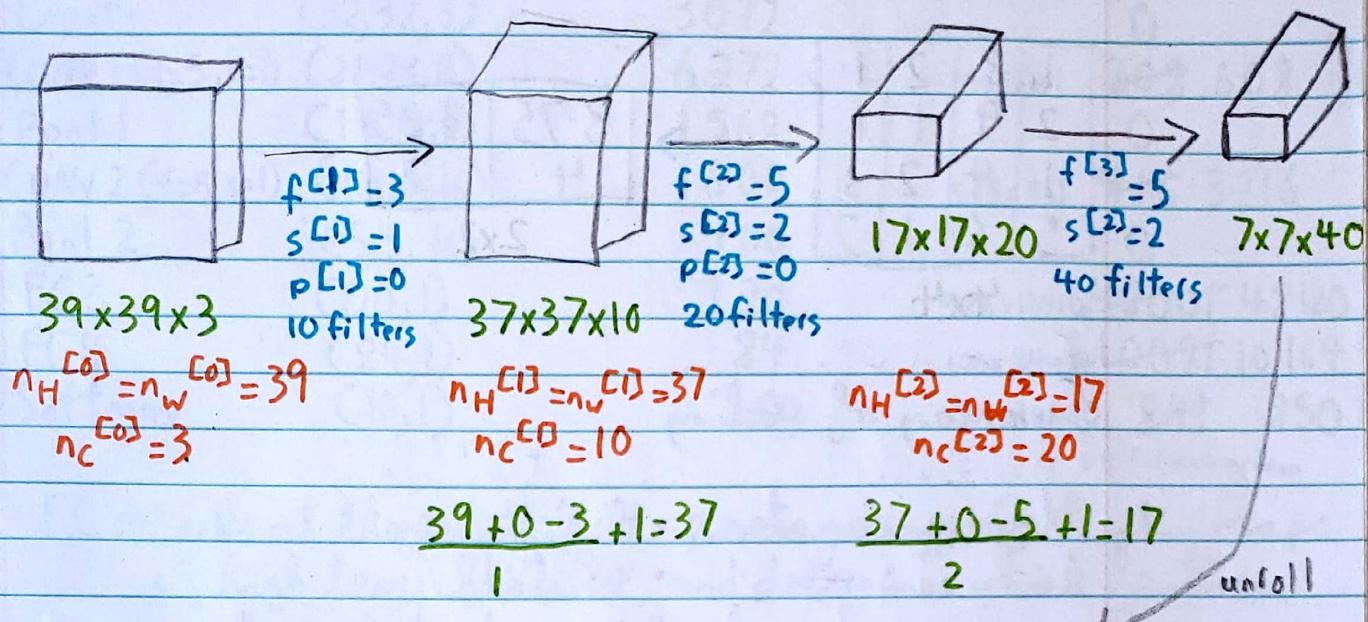
Weights: $f^{[\ell]} \times f^{[\ell]} \times n_c^{[\ell-1]} \times n_c^{[\ell]}$

Bias: $n_c^{[\ell]} - (1, 1, 1, n_c^{[\ell]})$

→ The n_c can sometimes be $n_c \times n_H \times n_W$

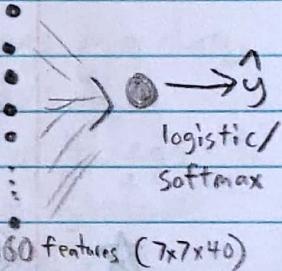
As long as we are consistent

Example ConvNet



Types of layers in NN

- Convolution Conv
- Pooling Pool
- Fully connected FC



Pooling layers

- Max pooling (works well... proven in many experiments)

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2

$f=2$

$s=2$

$4 \times 4 \rightarrow 2 \times 2$

Basically the max from each region

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

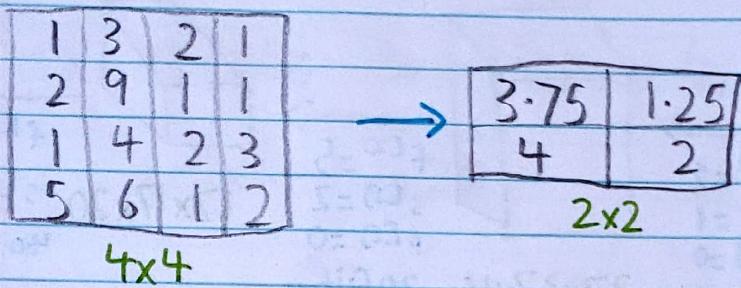
$f=3$

$s=1$

$5 \times 5 \rightarrow 3 \times 3$

$\frac{n+2p-f}{s} + 1 = 3 \times 3$

Average Pooling



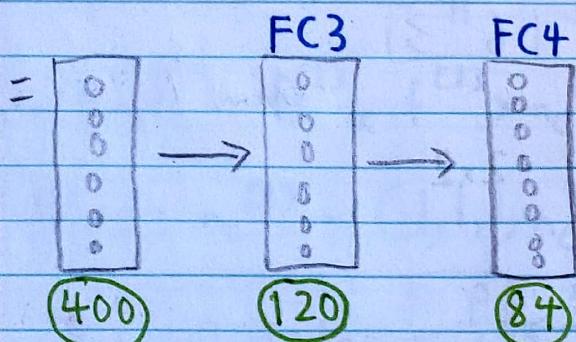
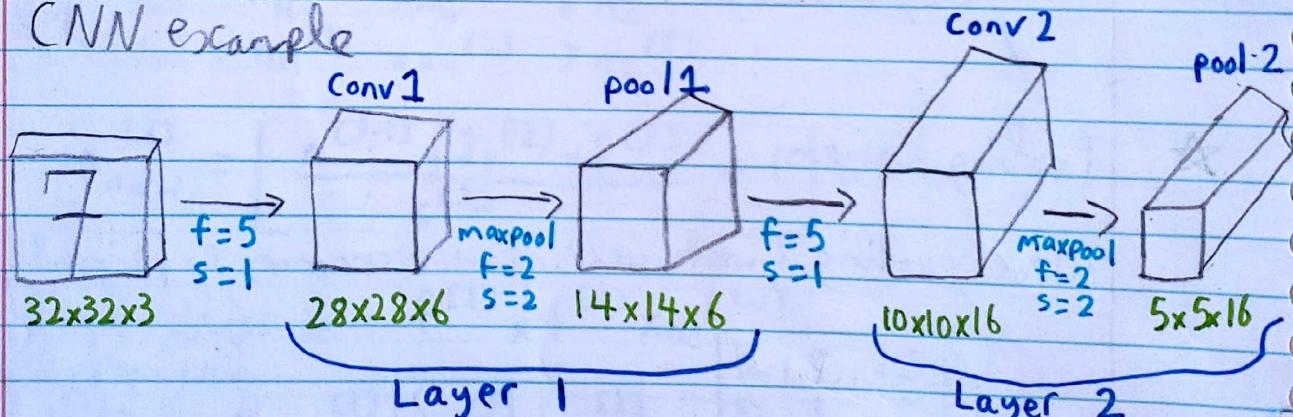
- Summary of pooling

- Hyperparameters :
- f : filter size
 - s : stride
 - Max or average pooling
 - p : padding (not usually with pooling)

No parameters to learn!

- ★ CNN example

Digit
Recognition



• softmax
if output $0 \rightarrow 9$
10 outputs

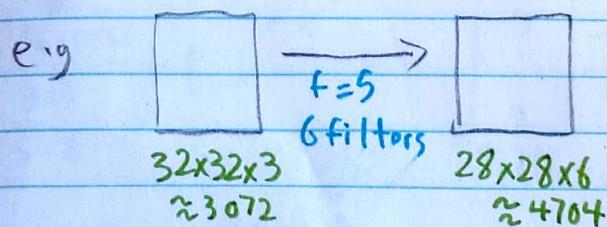
$n_H, n_W \downarrow$
 $n_C \uparrow$

	Activation Shape	Activation Size	# parameters
Input:	(32, 32, 3)	3072	0
CONV1 ($f=5, s=1$)	(28, 28, 8)	6272	$5 \times 5 \times 3 + 1 = 256$ 608
Pool 1	(14, 14, 8)	1568	0
CONV2 ($f=5, s=1$)	(10, 10, 16)	1600	$5 \times 5 \times 8 + 1 = 400$ 3216
Pool 2	(5, 5, 16)	400	0
FC 3	(120, 1)	120	$400 \times 120 + 120 = 48000$ 48120
FC 4	(84, 1)	84	$120 \times 84 + 84 = 10080$ 10164
Softmax	(10, 1)	10	$84 \times 10 + 10 = 850$ 850

Bias for each neuron

- FC \Rightarrow Looks at the output of the previous layer (activation maps of high level features) and determines which features most correlate to a particular class.
 - \rightarrow e.g. if the program is predicting that some image is a dog, it will have high values in the activation maps that represent high level features like a paw or 4 legs etc.
 - \rightarrow Basically a FC layer looks at what high level features most strongly correlate to a particular class and has particular weights so that when you compute the products between the weights and previous layer, you get the correct probabilities for the different classes.
- (Backpropagation makes all this work)

★ Why convolutions?



Instead of $3072 \times 4704 \approx 14M$
only $(5 \times 5 \times 3 + 1) \times 6 = 456$
 $f^{(l)} f^{(l)} n^{(l)} b^{(l)} n^{(l)}$

- \rightarrow Parameter Sharing: A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

e.g. 3x3 filter / parameter

1	0	-1
1	0	-1
1	0	-1

across each part of the image

→ Sparsity of connections : In each layer, each output value depends only on a small number of inputs

→ Very good at capturing translation invariance (a shifted image)
≈ still will be able to detect features

- Less prone to overfitting
- Can work with smaller training sets

• Training set $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$

$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m L(g^{(i)}, y^{(i)})$$

Use SGD or RMSprop or ADAM or something to optimize params
& reduce cost J

Week 2 - Deep CNN: Case Studies

case studies

Why look at case studies?

- The same way we learn how to code is by looking at other people's code, we can learn from or use other built CNN models for our task in hand

• Classic Networks:

- LeNet-5
- AlexNet
- VGG

• ResNet

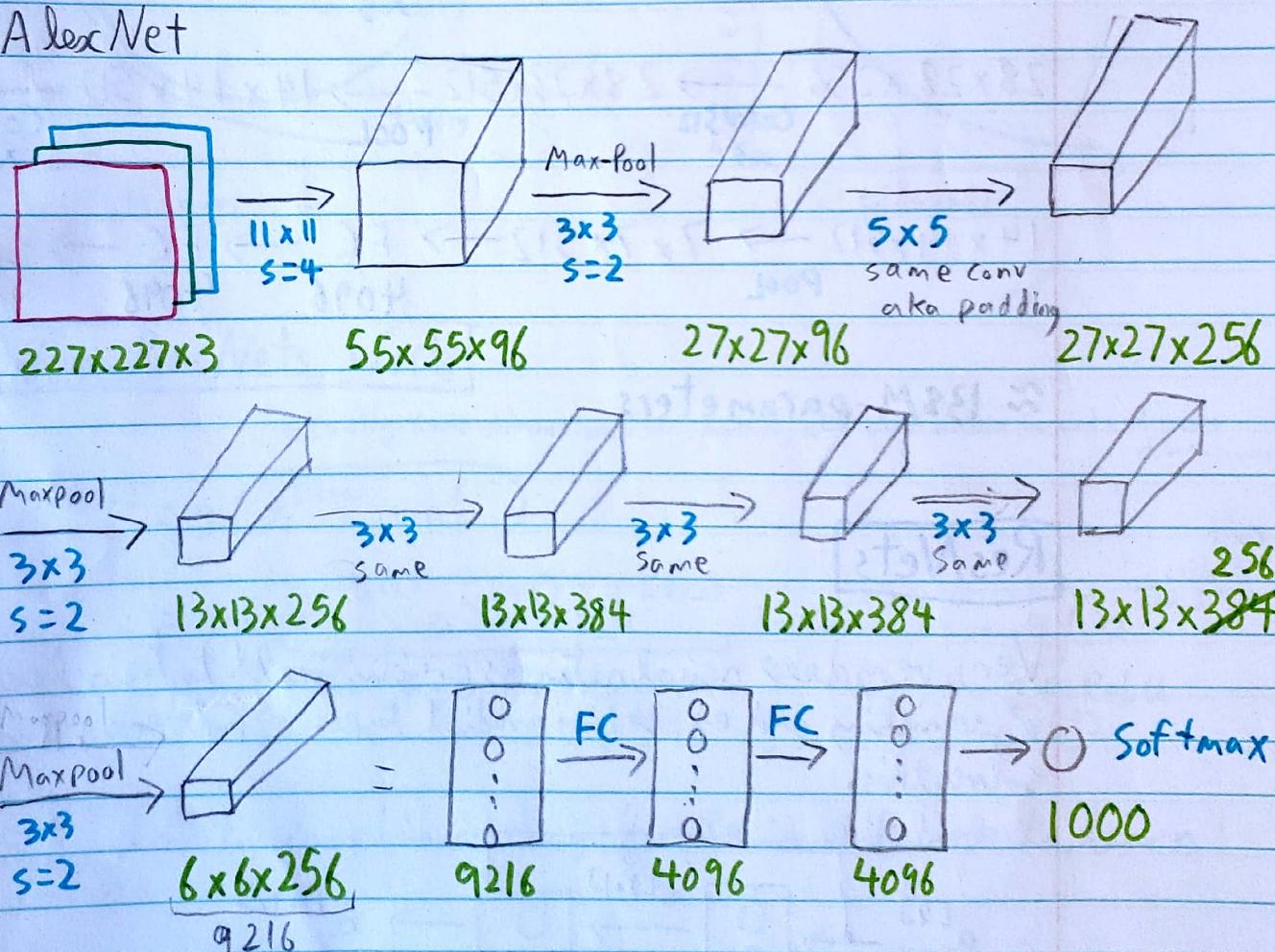
• Inception

Classic Networks

- LeNet-5 : Used originally to detect handwritten digits (1998)

- Same exact model as the CNN model 2 pages prior, but instead of max pool avg pool was used as it was more common back then and input was $32 \times 32 \times 1$ as opposed to $32 \times 32 \times 3$, as image was grayscale.
- $\approx 60K$ parameters
- Sigmoid/tanh instead of ReLU

- AlexNet

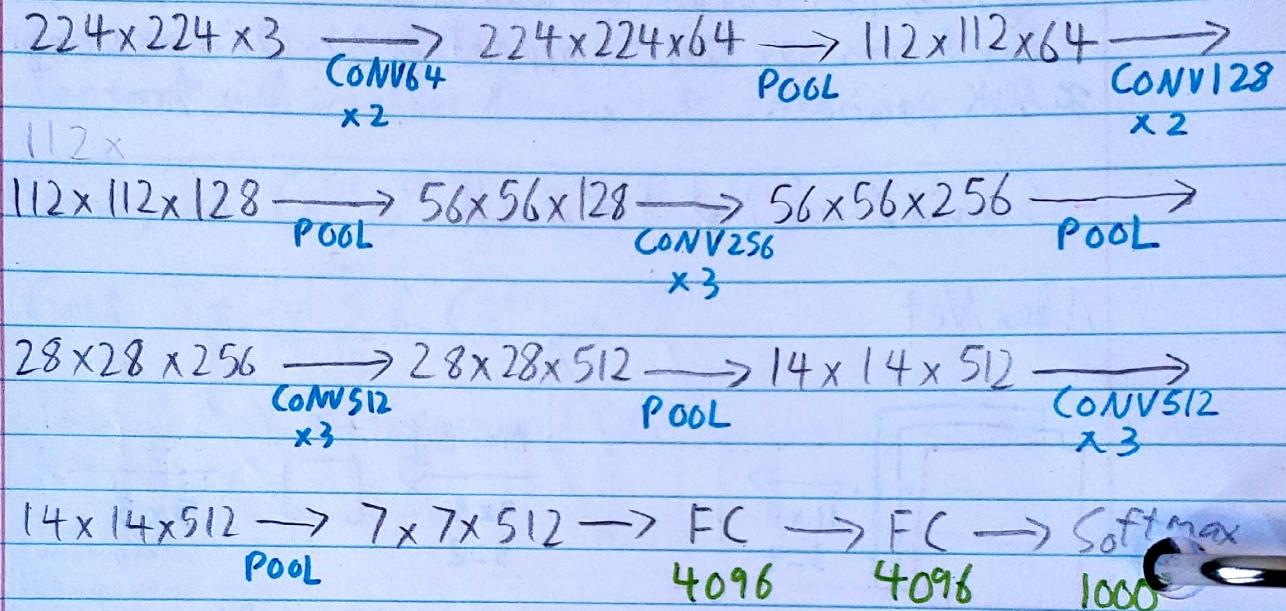


≈ 60 million parameters

- ReLU
- Local Response Normalization (not used often nowadays)

VGG-16

- $\text{CONV} = 3 \times 3 \text{ filter}, s=1, \text{ same}$
 $\text{MAX-POOL} = 2 \times 2, s=2$

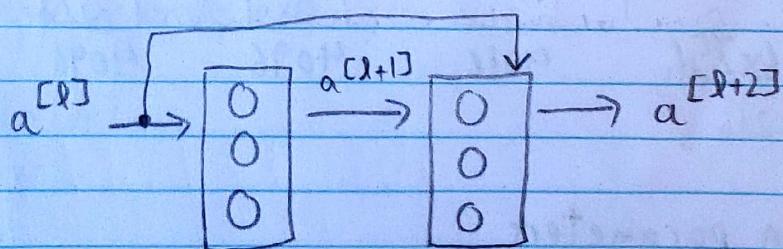


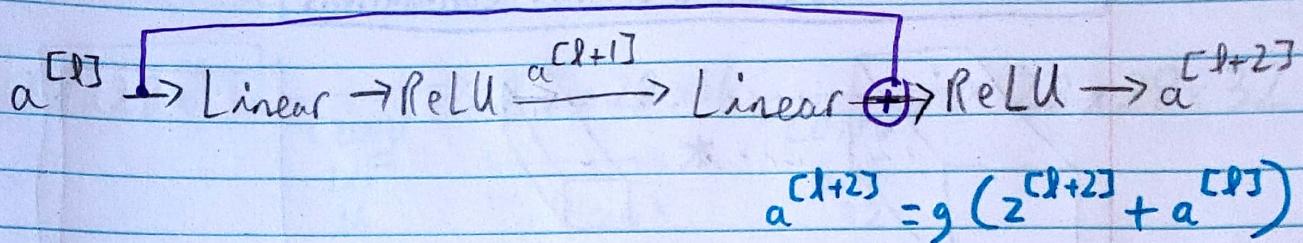
$\approx 138 \text{M parameters}$

Case Studies

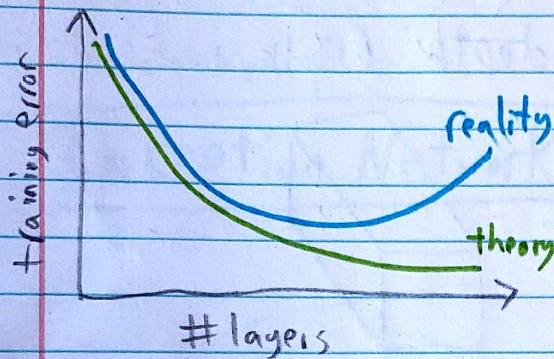
ResNets

- Very, very deep neural networks are difficult to train because of vanishing and exploding gradient types of problems. ResNets solves this.

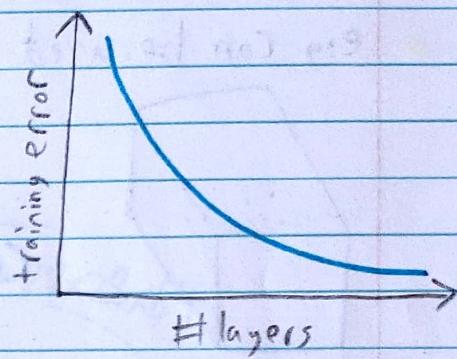




Plain Network



ResNet



Why ResNets Work

usually same dimensions If not, multiply $a^{[l]}$ by a matrix to make it same dimensions

- $a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$

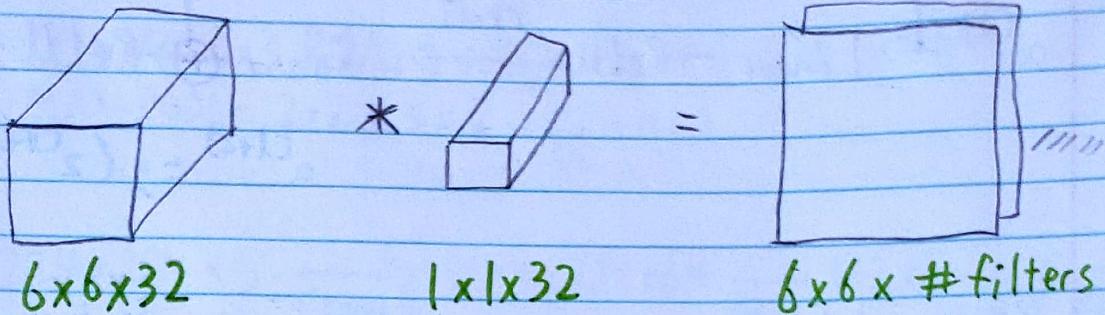
$$= g(W^{[l+2]} a^{[l+1]} + b^{[l+2]} + a^{[l]}) = g(a^{[l]})$$

If $= 0$ e.g. from weight decay

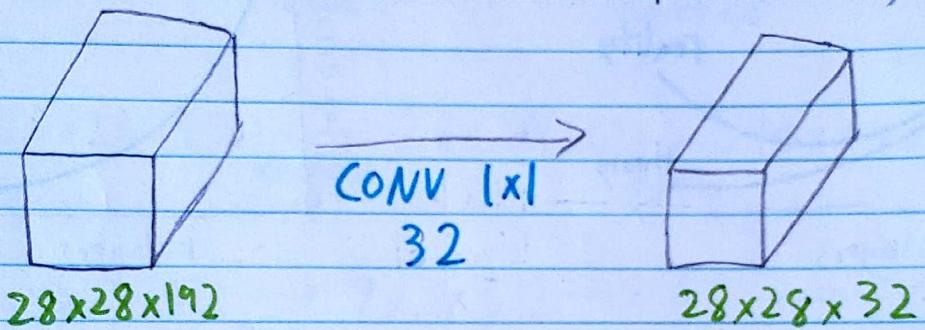
If ReLU
 $= a^{[l]}$

→ Identity function is easy for Residual Block to learn!

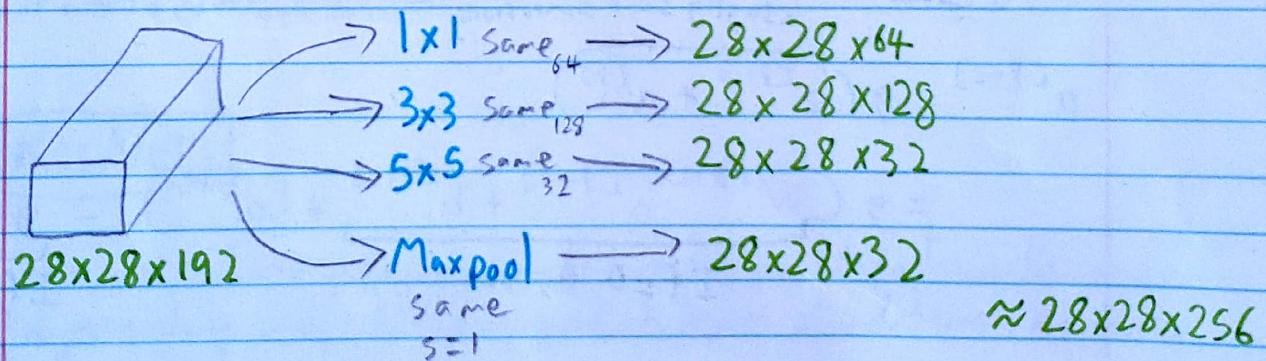
Networks in Networks & 1×1 Convolutions



- e.g. can be used to shrink depth (Pooling only shrinks H & W)



Inception Network Motivation



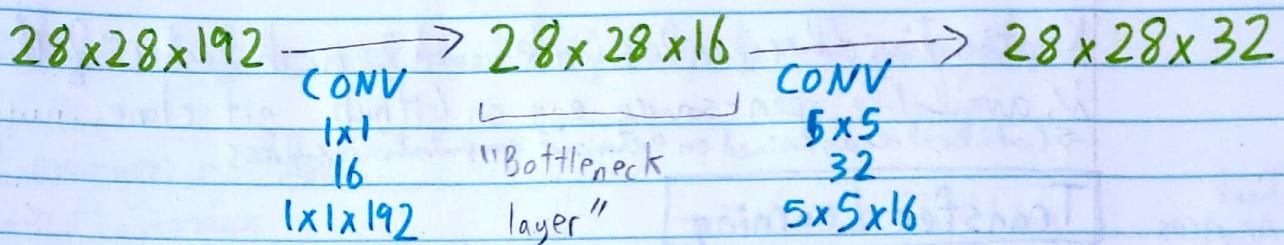
- Instead of deciding what to do, perform different computations and concatenate outputs. Let model decide then what to do later.

- Computational cost problem

$$\text{e.g. } 28 \times 28 \times 192 \xrightarrow{\begin{array}{l} \text{Conv } 5 \times 5 \\ \text{Same } 32 \end{array}} 28 \times 28 \times 32$$

$28 \times 28 \times 32 \times 5 \times 5 \times 192$
 $= 120 \text{M parameters}$

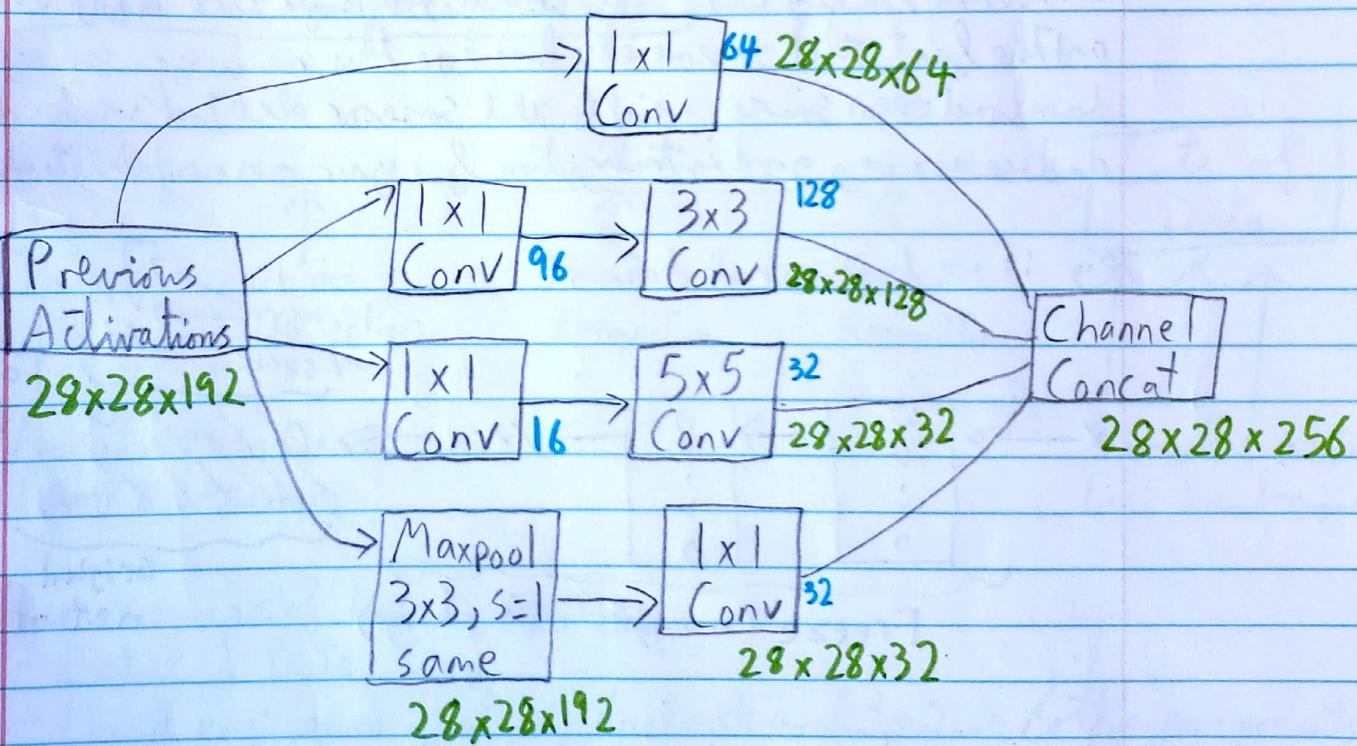
To reduce it : Use 1×1 Convolution



$$28 \times 28 \times 16 \times 1 \times 1 \times 192 = 2.4M \quad] \rightarrow 12.4M \text{ parameters}$$
$$28 \times 28 \times 32 \times 5 \times 5 \times 16 = 10M \quad]$$

Case studies

Inception Network



↑ 1 inception module. The network puts these modules together.

Practical
Advices

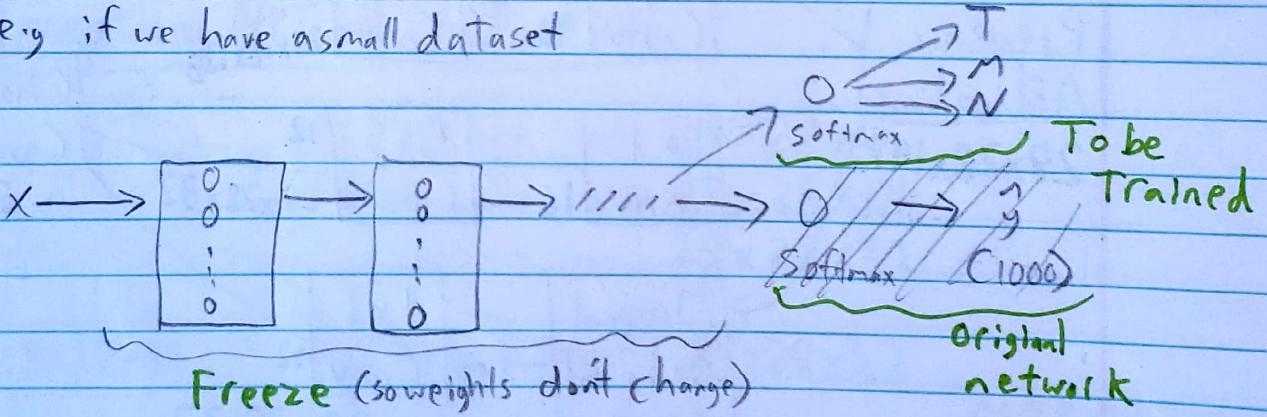
Using open-source implementations

- Replicating other people's work would be much faster if their code is available open-source e.g. on Github git clone...
→ Could be already trained on extensive computation resources

Practical
Advices

Transfer Learning

- Often make much faster progress if you download weights that someone else has already trained on the network architecture and use that as pre-training and transfer that to a new task that you might be interested in.
- Sometimes training could take several and might take many GPUs and the fact that someone else has done this means you can often download open source weights that someone else took weeks or months and use as a very good initialization for your own neural network.
- e.g. if we have a small dataset



- e.g. if we have a larger dataset

- Freeze fewer layers
- Number of layers to be trained on top increases

- e.g. if we have a HUGE dataset

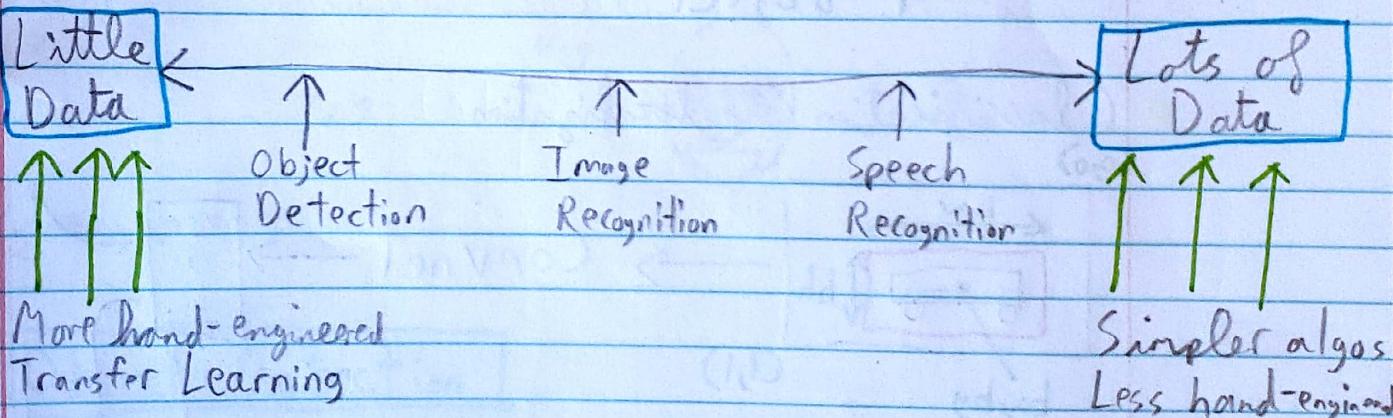
- We don't have to freeze any layers and use weights as initialization

Data Augmentation

- Used to improve performance of computer vision systems
- Common methods:
 - Mirroring
 - Random Cropping
 - Rotation
 - Shearing
 - Local Warping
 - ...

→ Color Shifting

State of Computer Vision

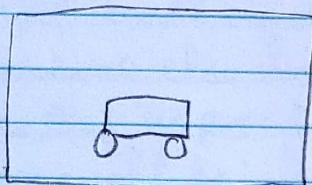


- Two sources of knowledge:
 - Labeled Data (x, y)
 - Hand engineered features / network architecture / other components
- Tips for doing well on benchmarks / winning competitions
 - Ensembling
 - Train several networks [3-15 networks] independently and average their outputs
 - Multi-crop at test time
 - Run classifiers on multiple versions

Week 3 - Detection Algorithms

Object Localization

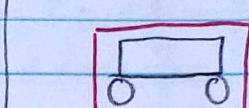
- Image Classification



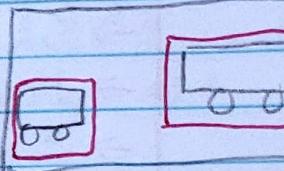
"Car"

1 object

- Classification with localization



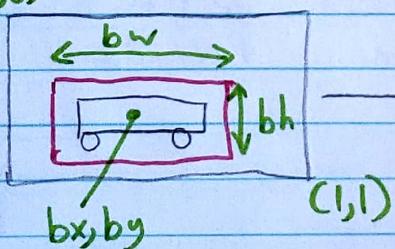
"Car"



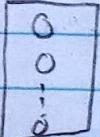
multiple objects

- Classification with localization

(c_0, c_1)



Convnet



Softmax
(4)

bounding box
 bx, by, bh, bw

1-Pedestrian

2-Car

3-Motorcycle

4-Background

$$y = \begin{bmatrix} P_c \\ bx \\ by \\ bh \\ bw \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

is there any obj?

e.g. car	1
bx	0
by	0
bh	0
bw	0
c ₁	1
c ₂	0
c ₃	0

background	0
?	?
?	?
?	?
?	?
?	?
?	?
?	?

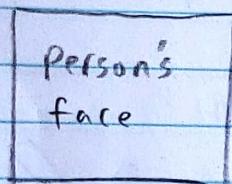
"Don't care"

if squared loss

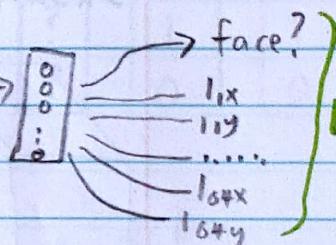
$$L(\hat{y}, y) =$$

$$\begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2 & \text{if } y_1 = 1 \\ (\hat{y}_1 - y_1)^2 & \text{if } y_1 = 0 \end{cases}$$

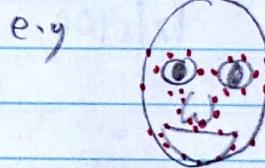
Landmark Detection



So people's \rightarrow Convnet \rightarrow faces



$$\left. \begin{matrix} l_{1x}, l_{1y} \\ l_{2x}, l_{2y} \\ \vdots \\ l_{64x}, l_{64y} \end{matrix} \right\} x, y$$



Object Detection

- e.g. Car Detection example

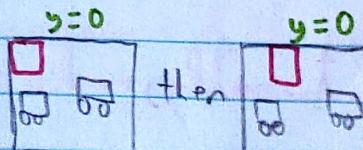
Training set:

x	y
closely	1-car
cropt pics	0-no
of cars	car
or background	

image \rightarrow Convnet \rightarrow y

Sliding Windows Detection

Then bigger windows

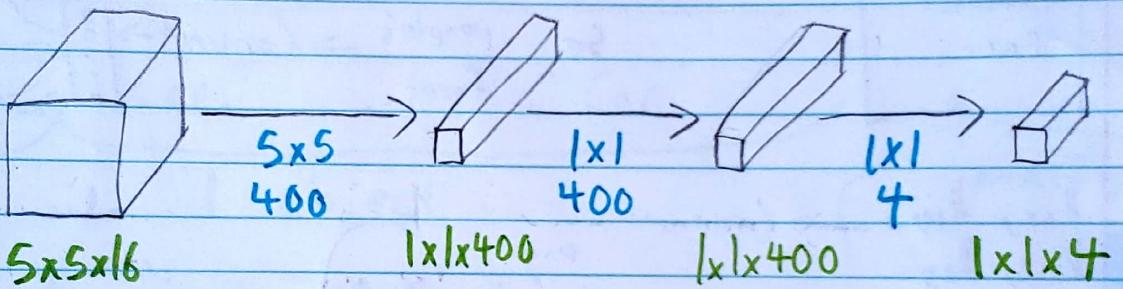
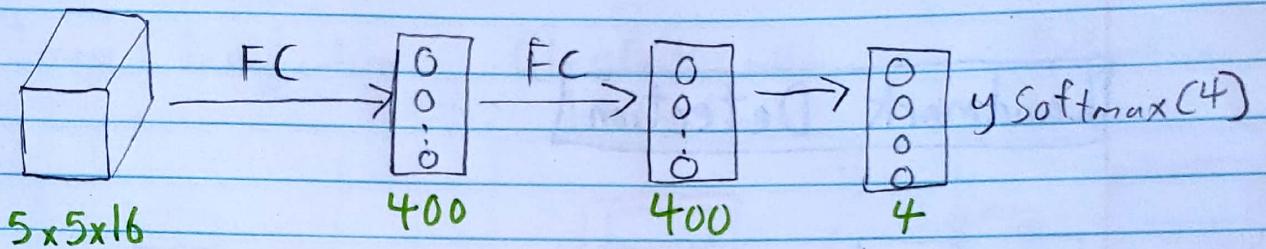


& so on till window
slid across every
position in the image

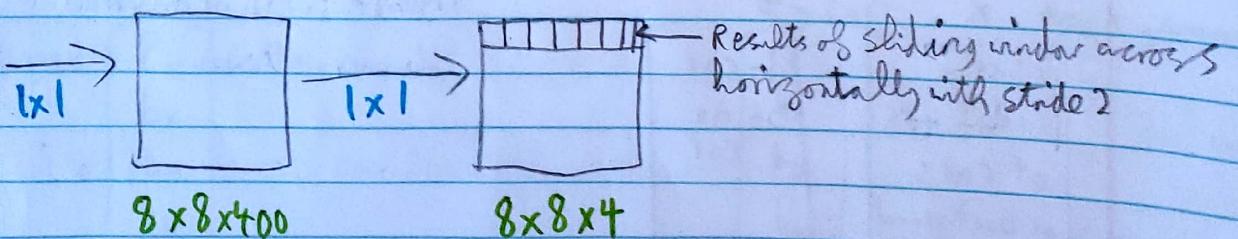
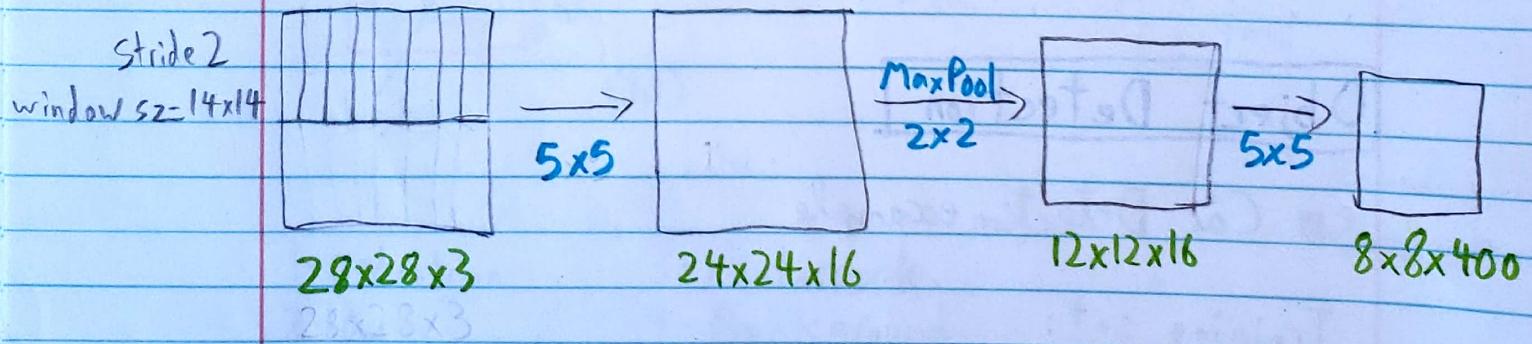
Computation cost :: HIGH!

Convolutional Implementation of Sliding Windows

- Turning FC layer into convolutional layers



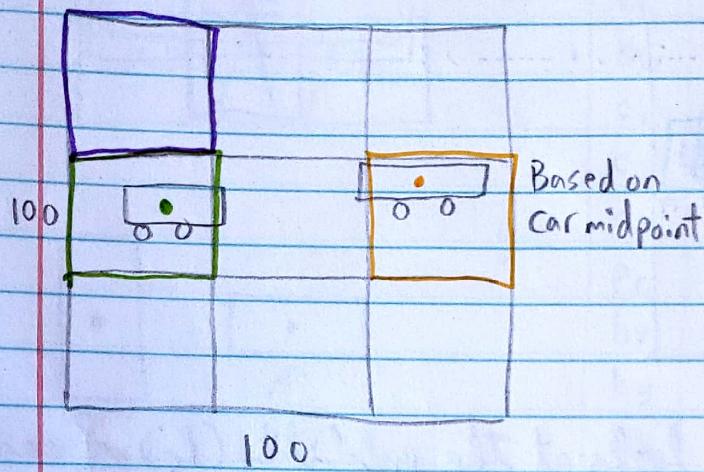
- Convolutional implementation of sliding windows



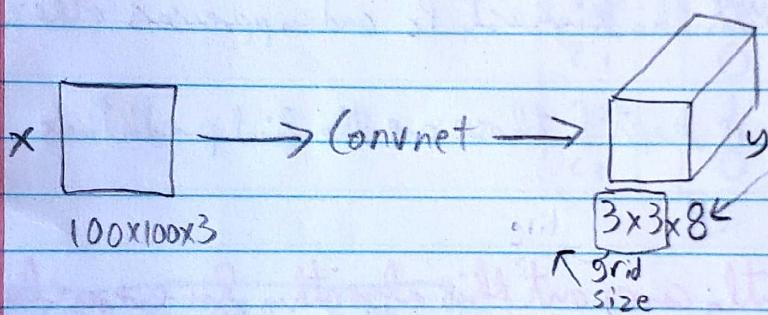
→ Convolutionally, make all predictions at the same time

Bounding Box Predictions

- YOLO Algorithm [You Only Look Once]



$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \quad \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

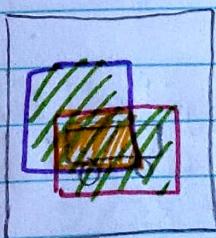


Usually finer grid to avoid having multiple objects in one box

b_x, b_y, b_h, b_w relative to grid size

Intersection Over Union [IoU]

- Evaluating object localization



- Predicted
- Union
- Intersection

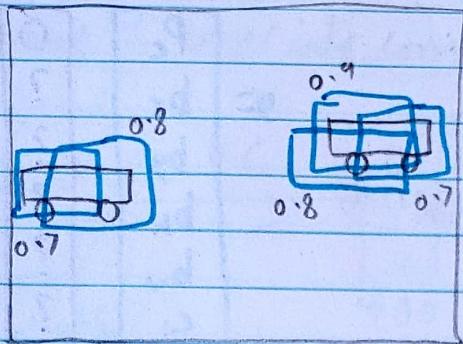
$$\text{Intersection over union} = \frac{\text{Size of intersection}}{\text{Size of union}}$$

Correct if $\text{IoU} \geq 0.5$

More generally, IoU is a measure of the overlap between 2 bounding boxes

Non-max Suppression

- e.g. 19x19 grid, can have multiple detections per car



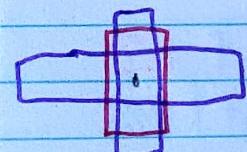
- Non-max suppression looks at the probability, (P_c) of each box
- Highlights the one with the highest P_c and suppresses others with high IoU
- Highlighted ones are kept & those are the final predictions
∴ 1 box per car

Note: Independently, carry out this algorithm for every class
e.g. cars, pedestrians, busses

Anchor Boxes

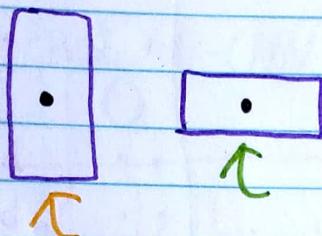
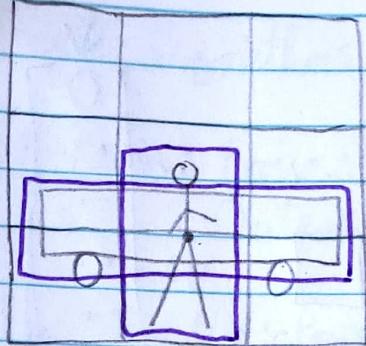
- Previously:

Each object in training image is assigned to grid cell that contains that object's midpoint. $3 \times 3 \times 8$



With two anchor boxes:

Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU. $3 \times 3 \times 16$
 $(3 \times 3 \times 2 \times 8)$

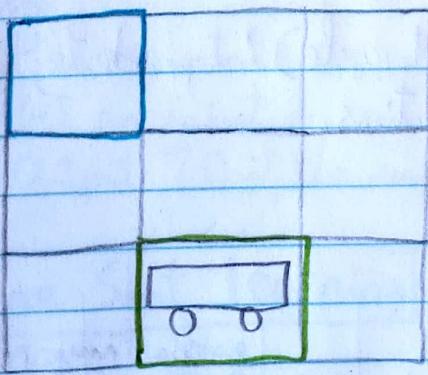


	Car only?	
P _c	1	0
b _x	b _x	?
b _y	b _y	?
b _h	b _h	?
b _w	b _w	?
c ₁	1	?
c ₂	0	?
c ₃	0	?
P _c	1	1
b _x	b _x	b _x
b _y	b _y	b _y
b _h	b _h	b _h
b _w	b _w	b _w
c ₁	0	0
c ₂	1	1
c ₃	0	0

YOLO Algorithm

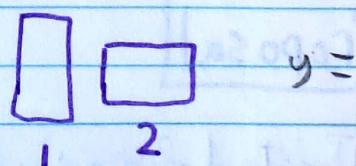
Putting all the components learnt together

Training

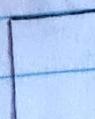


y $3 \times 3 \times 16$

(Bounding box needs to be provided)



$3 \times 3 \times 16$

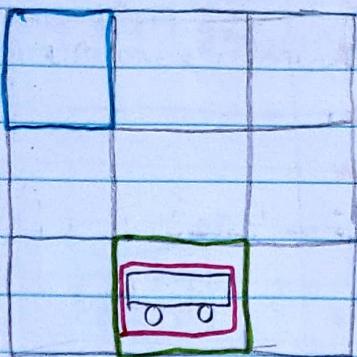


\rightarrow ConvNet \rightarrow

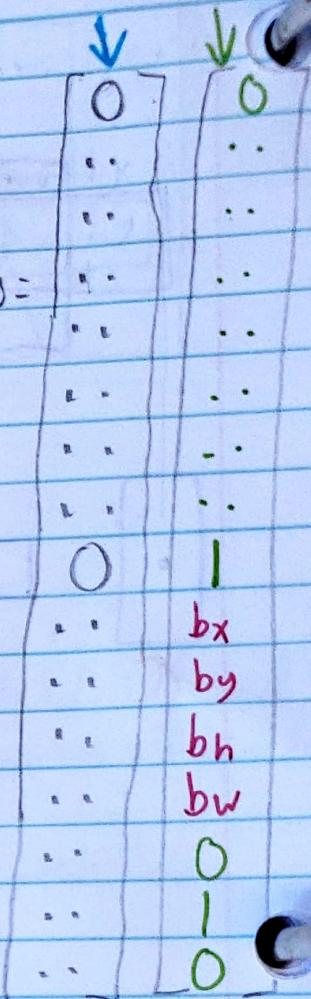
$100 \times 100 \times 3$

1- pedestrian	P _c	0	0
2- car	b _x	?	?
3- motorcycle	b _y	?	?
	b _h	?	?
	b _w	?	?
	c ₁	?	?
	c ₂	?	?
	c ₃	?	?
	P _c	0	1
	b _x	?	b _x
	b _y	?	b _y
	b _h	?	b _h
	b _w	?	b _w
	c ₁	?	0
	c ₂	?	1
	c ₃	?	0

- Making predictions



$y = \dots$
 $3 \times 3 \times 2 \times 8$



- Outputting the non-max suppressed outputs

- For each grid cell, get 2 predicted bounding boxes
- Get rid of low probability predictions
- For each class (pedestrian, car, motorcycle) independently, use non-max suppression to generate final predictions.

Region Proposal

- R-CNN : Instead of running sliding window, only run detection algo on classified regions across all image parts

→ Regions selected based on segmentation algo (similar to an image heatmap)

Faster algorithms:

- R-CNN: Propose regions. Classify proposed regions one at a time.
Output label + bounding box.
- Fast R-CNN: Propose regions. Use convolutional implementation of sliding windows to classify all the proposed regions.
- Faster R-CNN: Use convolutional network to propose regions.

Week 4 - Face Recog & Neural style transfer

Face
Recog

What is Face Recog?

- Face verification:
 - Input image, name / ID
 - Output whether the input image is that of the claimed person
- Face recognition:
 - Has a database of K persons
 - Get an input image
 - Output ID if the image is any of the K persons (or "not recognized")

Face
Recog

One Shot learning

- Learning from one example to recognize the person again
- Learn a "similarity" function

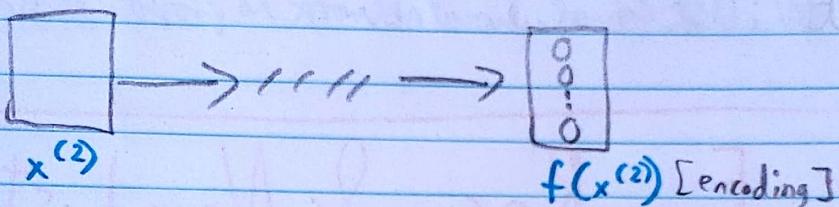
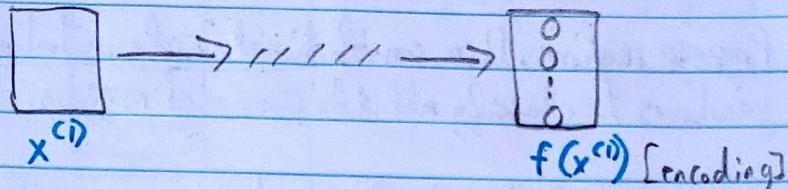
$$\rightarrow d(\text{img 1}, \text{img 2}) = \text{degree of difference between images}$$

$$\begin{cases} \text{If } d(\text{img 1}, \text{img 2}) \leq T \\ \text{If } d(\text{img 1}, \text{img 2}) > T \end{cases}$$

"same"
"different" } Verification

Siamese Network

- The idea of running two identical CNN on two different inputs and then comparing them.



$$d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|_2^2$$

- Goal of learning

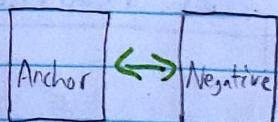
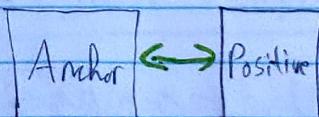
→ Parameters of NN define an encoding $f(x^{(i)})$

→ Learn parameters so that:

- If $x^{(i)}, x^{(j)}$ are the same person, $\|f(x^{(i)}) - f(x^{(j)})\|^2$ is small
- If $x^{(i)}, x^{(j)}$ are different persons, $\|f(x^{(i)}) - f(x^{(j)})\|^2$ is large

Triplet Loss

- Learning Objective



$$\text{Want: } \frac{\|f(A) - f(P)\|^2}{d(A, P)} + \frac{\alpha}{\text{margin}} \leq \frac{\|f(A) - f(N)\|^2}{d(A, N)}$$

Pushes the Pairs
Further away from each other

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha \leq 0$$

Loss function:

Given 3 images A, P, N

$$l(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0)$$

$$J = \sum_{i=1}^m l(A^{(i)}, P^{(i)}, N^{(i)})$$

Choosing the triplets A, P, N

During training if A, P, N are chosen randomly, $d(A, P) + \alpha \leq d(A, N)$ is easily satisfied.

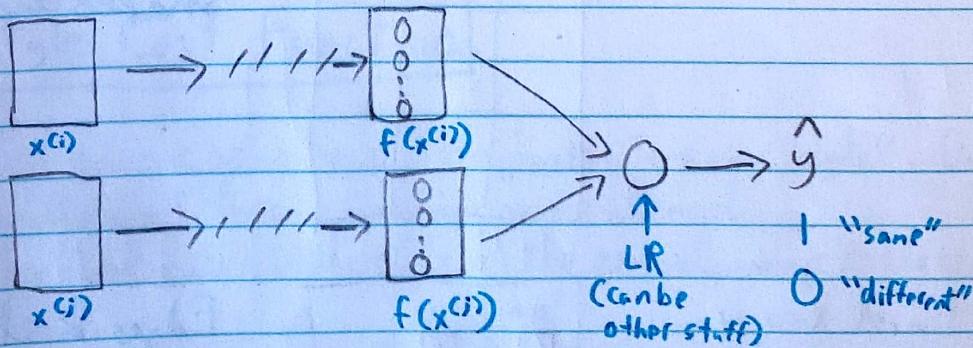
$\rightarrow \therefore$ Choose triplets that are hard to train on

2015 FaceNet paper Schroff et al

Face
Recog

Face Verification and Binary Classification

Learning the similarity function



$$\hat{y} = o\left(\sum_{K=1}^{128} w_K |f(x^{(i)})_K - f(x^{(j)})_K| + b\right)$$

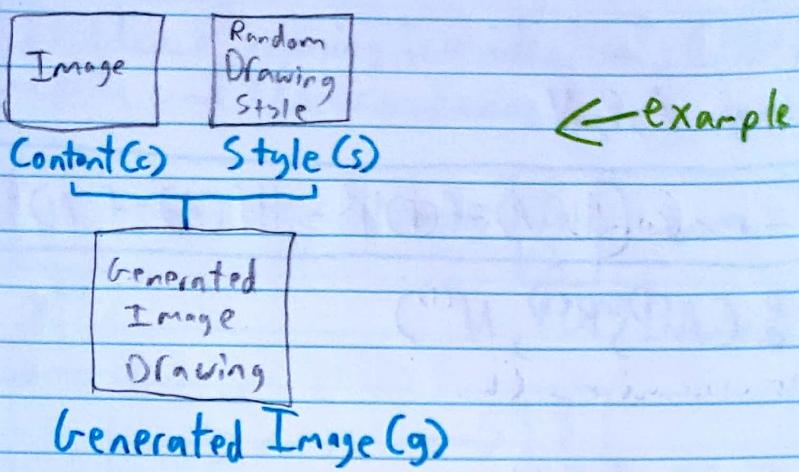
(could be other formulas)

$$\text{e.g. } \frac{f(x^{(i)})_K - f(x^{(j)})_K}{f(x^{(i)})_K + f(x^{(j)})_K}$$

- If image is from database, encoding could be precomputed to avoid running the convNets everytime and avoid storing the raw images, saving significant computation

Neural Style Transfer

What is Neural Style Transfer? (Optimization Technique)



Neural Style Transfer

What are deep ConvNets learning?

- Pick a unit in layer 1. e.g Find the nine image patches that maximize the unit's activation. Repeat for other units. Then move to other
- Usually in earlier layers, small image patches are examined and network learns about edges or shadows of colors
- In deeper layers, larger image patches are examined \therefore more complex shapes and patterns

e.g In last layer

9 Keyboard Images		
	words	
	flowers	dogs

More complex shapes and patterns

First layer

<input type="checkbox"/>	<input checked="" type="checkbox"/>
	learns green
	learns orange

Edges & shadows of colors

Cost Function

$$J(G) = \alpha J_{\text{Content}}(C_G) + \beta J_{\text{Style}}(S_G)$$

- 1) Initiate G randomly

$G: 100 \times 100 \times 3$ (Basically random noise)

- 2) Use GD to minimize $J(G)$

$$G := G - \frac{\partial}{\partial G} J(G)$$

Content Cost Function

- Say you use hidden layer l to compute content cost. l should not be too deep or too shallow. e.g. $\square \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \hat{g}$ \uparrow good choice
- Use pre-trained ConvNet. (e.g. VGG network)
- Let $a^{[l]}_{[C]}$ and $a^{[l]}_{[G]}$ be the activation of layer l on the images.
- If $a^{[l]}_{[C]}$ and $a^{[l]}_{[G]}$ are similar, both images have similar content

$$J_{\text{Content}}(C_G) = \frac{1}{2} \|a^{[l]}_{[C]} - a^{[l]}_{[G]}\|^2$$

Style Cost Function

- Say you are using layer l 's activation to measure "style". Define style as correlation between activations across channels.
- See how similar is the style of the generated image to the style of the input image

Style matrix (gram matrix)

Let $a^{[l]}_{i,j,k} = \text{activation at } (i,j,k)$. $G^{[l]}$ is $n_c^{[l]} \times n_c^{[l]}$

$$G_{KK'}^{[l]} = \sum_{i=1}^{n_h^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l]} a_{ijk'}^{[l]}$$

$$G_{KK'}^{[l]} = \sum_{i=1}^{n_h^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l]} a_{ijk'}^{[l]}$$

$K = 1, \dots, n_c^{[l]}$

$$J_{\text{style}}^{[1]}(S, G) = \frac{1}{(2n_H^{[1]} n_W^{[1]} n_C^{[1]})^2} \sum_K \sum_{K'} \underbrace{(G_{KK'}^{[1](S)} - G_{KK'}^{[1](G)})^2}_{\|G^{[1](S)} - G^{[1](G)}\|_F^2}$$

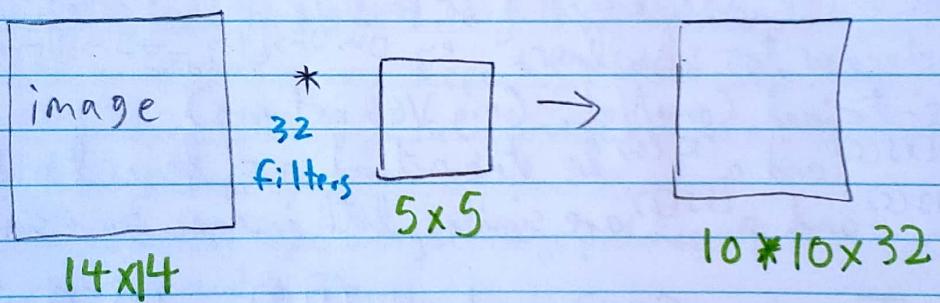
$$J_{\text{style}}(S, G) = \sum \lambda^{[l]} J_{\text{style}}^{[l]}(S, G)$$

↑↑

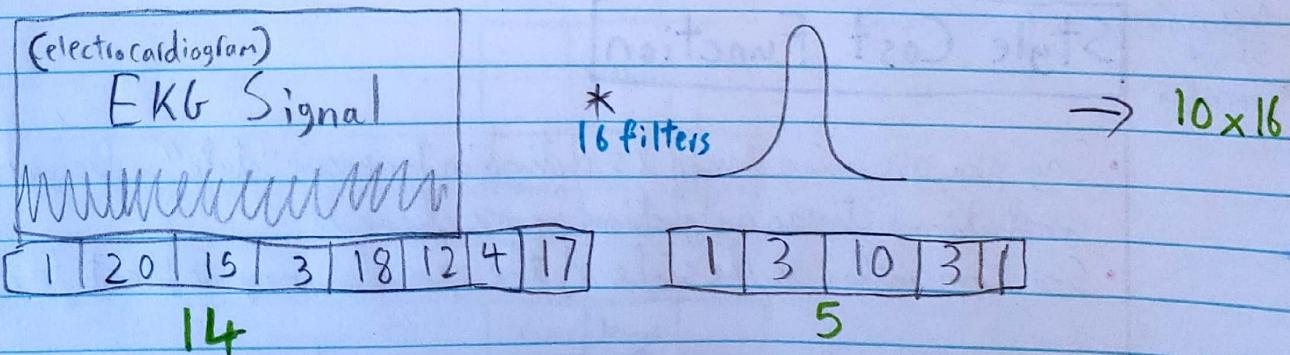
Take both low level and high level correlation into account
when computing style

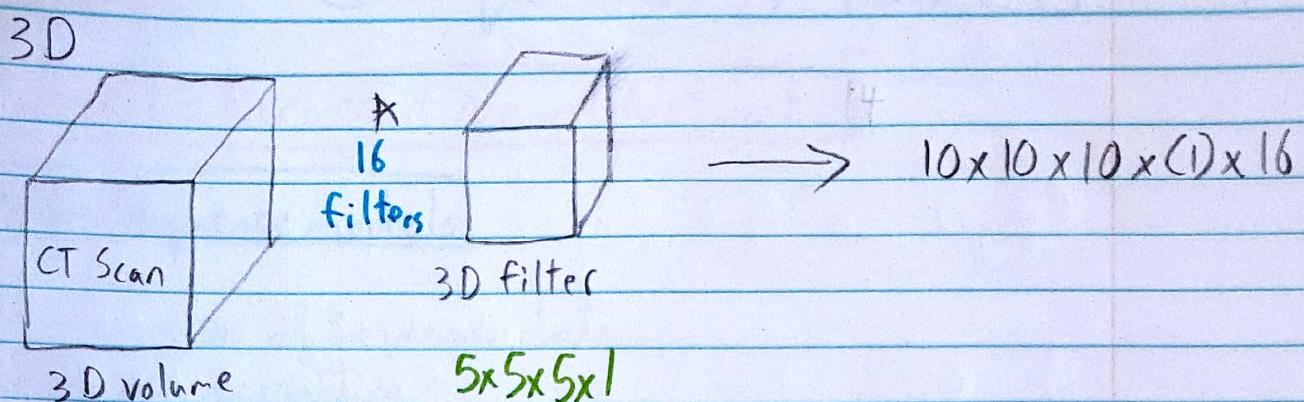
1D and 3D Generalizations

2D



1D





$14 \times 14 \times 14 \times 1$

$\uparrow \uparrow \uparrow \uparrow$

H W D channels

Sequence Models

Week 1 - Recurrent Neural Networks

RNNs

Why sequence models

- Examples of sequence data:

- Speech recognition
- Music generation
- Sentiment classification (e.g. "Movie sucks" → ★★☆☆☆)
- DNA sequence analysis
- Machine translation (e.g. French → English)
- Video activity recognition
- Name entity recognition

RNNs

Notation

- Motivating example (name entity recognition)

x : Harry Potter and Hermione Granger invented a new spell

$$\begin{array}{ccccccccc} x^{(1)} & x^{(2)} & x^{(3)} & \dots & \dots & x^{(t)} & \dots & x^{(9)} & T_x = 9 \\ y^{(1)} & y^{(2)} & y^{(3)} & \dots & \dots & \dots & \dots & y^{(9)} & T_y = 9 \end{array}$$

$$\begin{array}{l} x^{(i)} \langle t \rangle \\ y^{(i)} \langle t \rangle \end{array} \quad \begin{array}{l} T_x^{(i)} \\ T_y^{(i)} \end{array}$$

- Representing words (One-hot)

Not in vocab? <UNK>

Vocabulary

a	1
aaron	2
and	367
harry	4075
potter	6830
zula	10000

$$x^{(1)} \quad \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \leftarrow 4075$$

$$x^{(2)} \quad \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

← 6830

$$x^{(3)} \quad \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

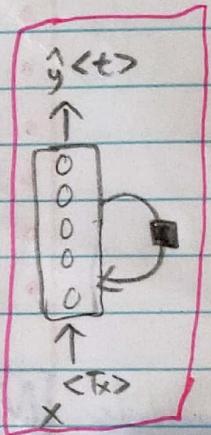
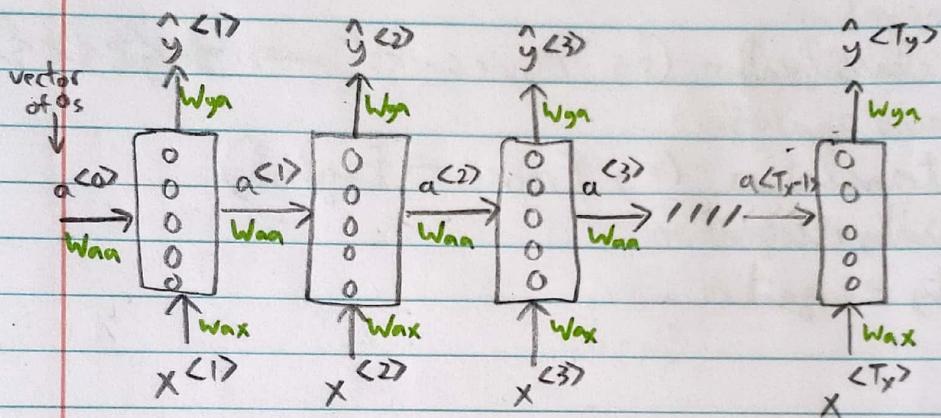
← 367

10,000

RNNs

Recurrent Neural Network Model

- Problems of a standard network:
 - Inputs, outputs can be different lengths in different examples
 - Doesn't share features learned across different positions of text
- Simple unidirectional RNN



Limitation:

- He said "Teddy Roosevelt was a great President"
- He said "Teddy bears are on sale!"

Need bidirectional RNN //

- Forward Propagation

$$a^{<0>} = \vec{0}$$

more generally

$$\begin{cases} a^{<1>} = g_1(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a) \\ \hat{y}^{<1>} = g_2(W_{ya} a^{<1>} + b_y) \end{cases} \quad \begin{matrix} \leftarrow \text{tanh/ReLU} \\ \leftarrow \text{Sigmoid} \end{matrix}$$

$$\begin{cases} a^{<t>} = g(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a) \\ \hat{y}^{<t>} = g(W_{ya} a^{<t>} + b_y) \end{cases}$$

$$W_a = [W_{aa} : W_{ax}]$$

$$a^{<t>} = g(W_a [a^{<t-1>}], x^{<t>}) + b_a$$

Back propagation through time

- Arrows go in opposite direction of forward propagation: Scan from right to left

$$\delta^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{(t)} \log \hat{y}^{<t>} - (1-y^{<t>}) \log (1-\hat{y}^{<t>})$$

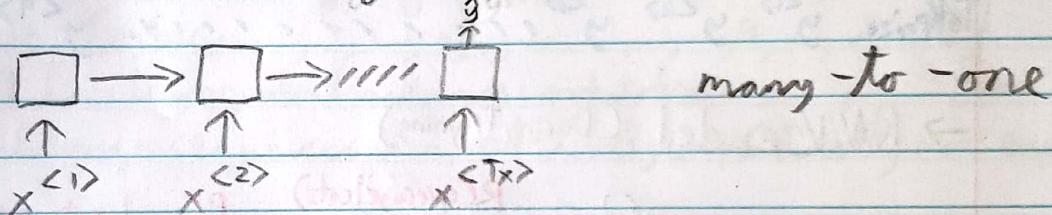
$$\delta(\hat{y}, y) = \sum_{t=1}^{T_y} \delta^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

RNNs

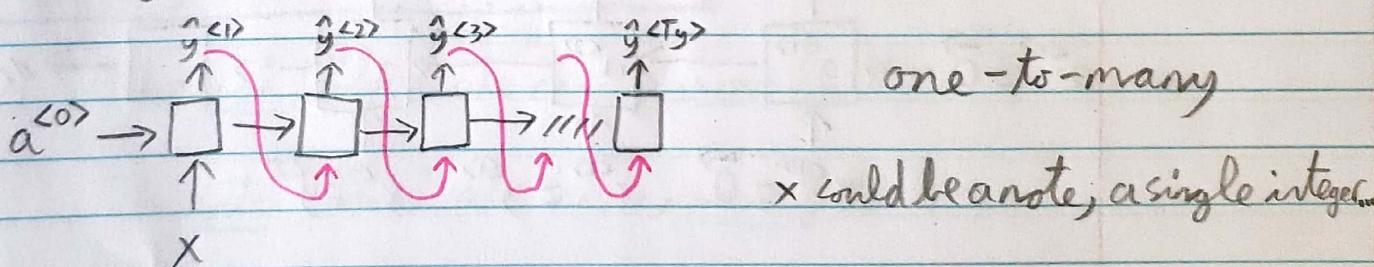
Different types of RNNs

- The RNN we have seen so far is many-to-many

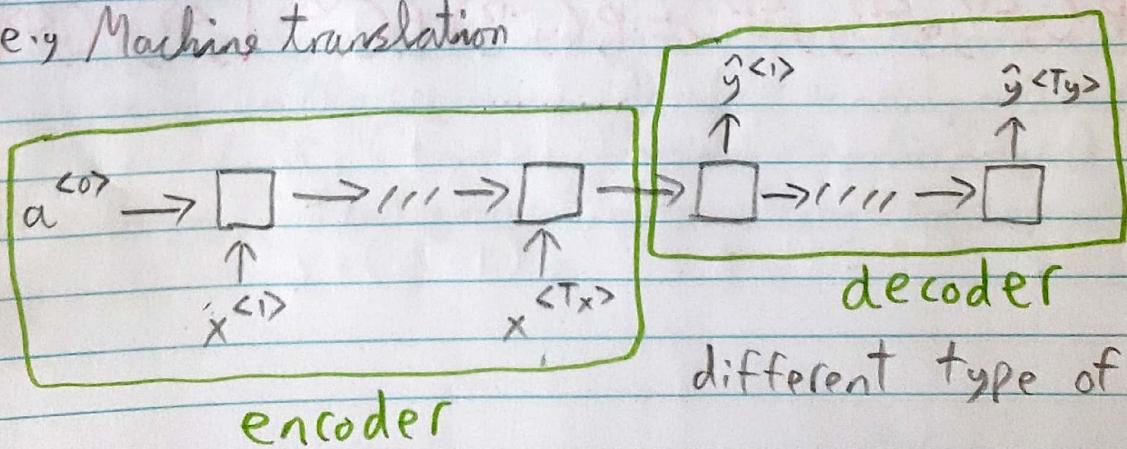
- e.g. in Sentiment Classification $x = \text{text}$ $y = 0/1 \text{ or } 1 \dots 5$



- e.g. Music generation $x \rightarrow y^{<1>} y^{<2>} \dots y^{<Ty>}$



- e.g. Machine translation



RNNs

Language Model and Sequence Generation

- What is language modelling? e.g. in Speech Recognition

$$P(\text{The apple and } \underline{\text{pair}} \text{ salad}) = 3.2 \times 10^{-13}$$

$$P(\text{The apple and } \underline{\text{pear}} \text{ salad}) = 5.7 \times 10^{-10}$$

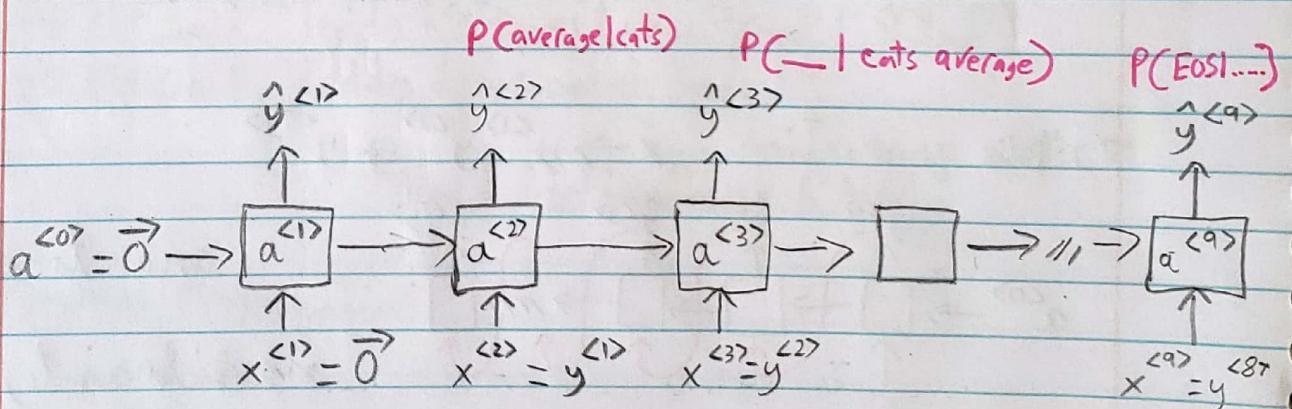
$$P(y^{<1>} y^{<2>} \dots y^{<T_y>})$$

- Training set: large corpus of English text

e.g. Cats average 15 hours of sleep a day. $\langle \text{EOS} \rangle$
end of sentence

Tokenize $y^{<1>} y^{<2>} y^{<3>} \dots \dots \dots y^{<9>}$

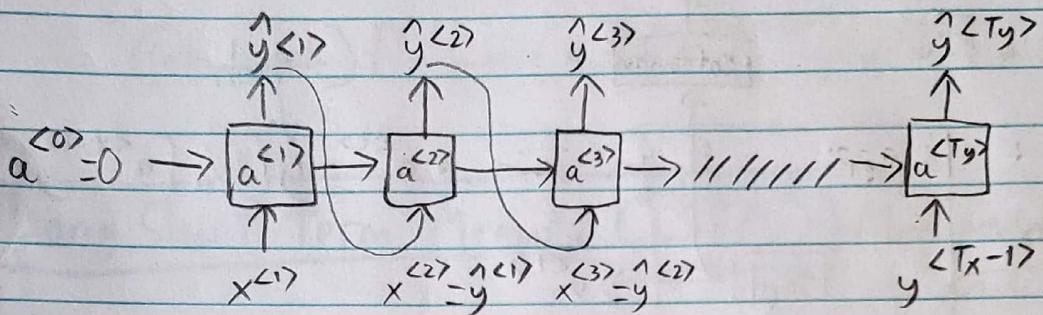
→ RNN model (During Training)



$$P(y^{<1>} y^{<2>} y^{<3>}) = P(y^{<1>}) P(y^{<2>} | y^{<1>}) P(y^{<3>} | y^{<1>} y^{<2>}) \dots \dots$$

Sampling novel sequences

- ## Sampling a sequence from a trained RNN



- Word-level language model
Vocab = [a, aaron, ..., zulu, <UNK>]

Character-level language model

Vocab = [a, b, c, ..., z, ʌ, ɒ, ɒ, ..., ə, ..., ɪ, ..., A, ..., Z]

- Don't have to worry about unknown word token
 - Much longer sequences than word-level, computationally more expensive to train

RNNs

Vanishing Gradients with RNNs

- e.g. The cat which ate was full

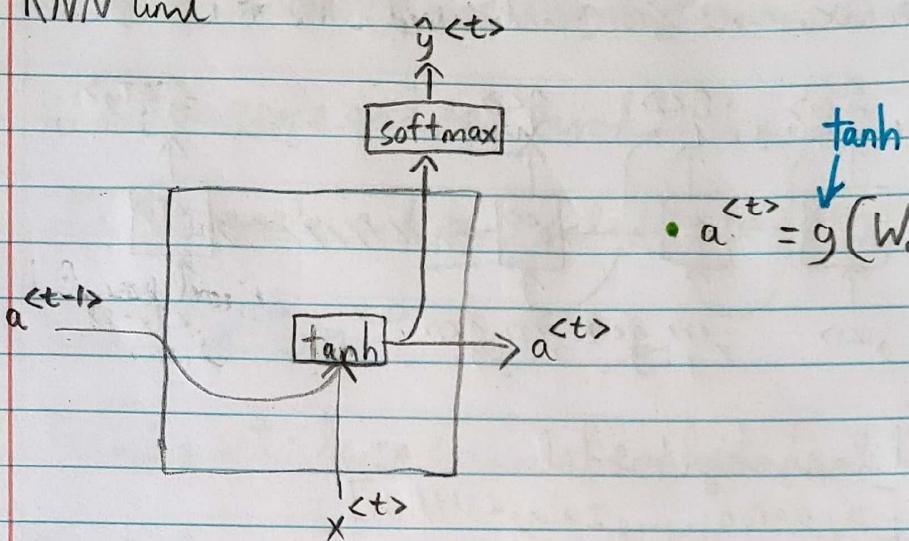
The cats, which ate raw, wet fish

- Hard for $\hat{y}^{(t_y)}$ to backpropagate all the way to the beginning of the sequence and modify how the RNN is doing the computation earlier in the sequence. An output is mainly influenced by values close to it. (ie modify the weights of earlier layers)
- Basic RNN not good at capturing long-term dependencies

RNNs

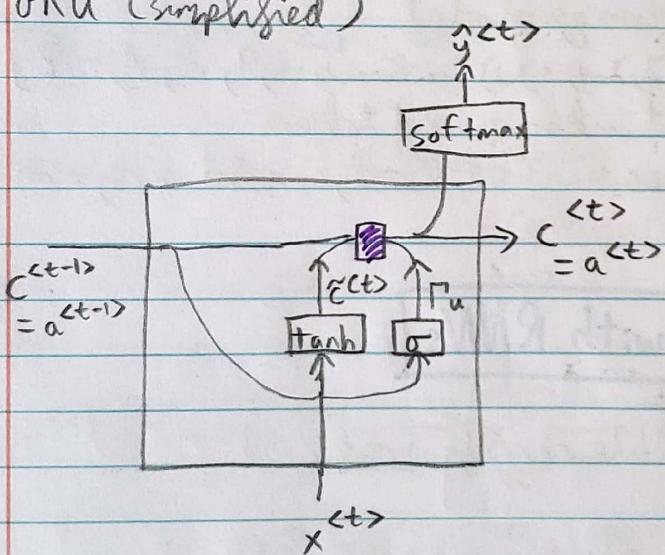
Gated Recurrent Unit (GRU)

- RNN unit



$$\bullet a^{<t>} = g(W_a [a^{<t-1>}, x^{<t>}] + b_a)$$

- GRU (simplified)



• $c = \text{memory cell}$

$$c^{<t>} = a^{<t>}$$

$$\tilde{c}^{<t>} = \tanh (W_c [c^{<t-1>}, x^{<t>}] + b_c)$$

candidate to
replace c

Helps with
Vanishing gradient



update gate $\bullet \Gamma_u = \sigma (W_u [c^{<t-1>}, x^{<t>}] + b_u)$

$$\boxed{\therefore c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}}$$

element-wise

Full GRU

The difference:

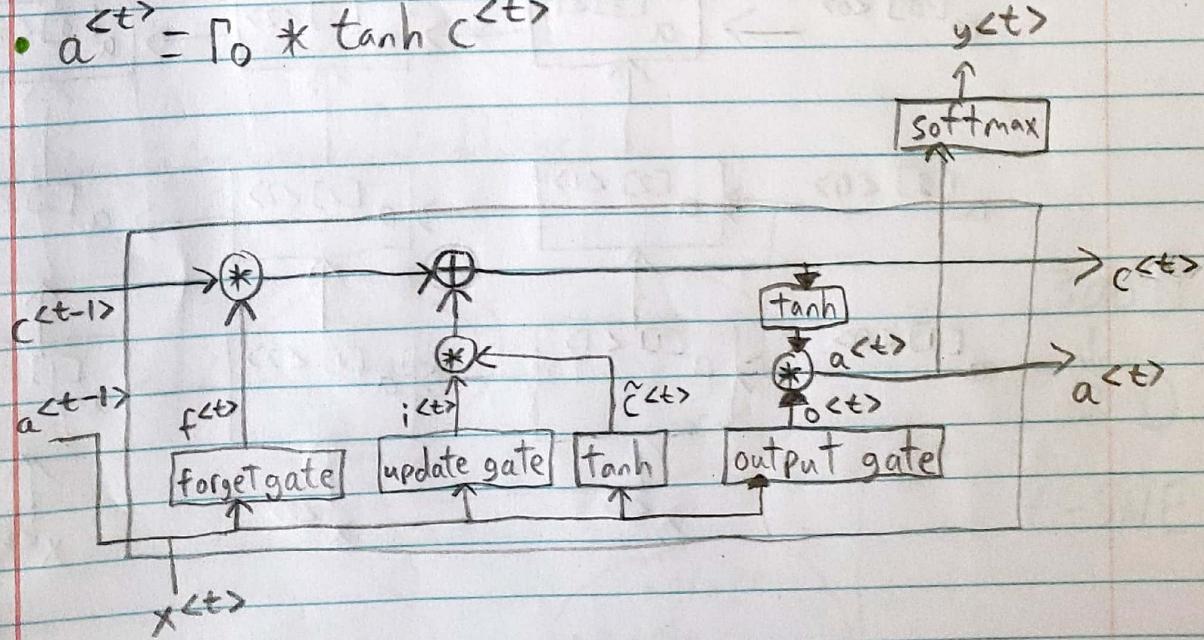
- $\tilde{c}^{<t>} = \tanh(W_c[c^{<t-1>}, x^{<t>}] + b_c)$
- (Relevance gate) $\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$

RNNs

Long Short Term Memory (LSTM)

More powerful & general than GRU

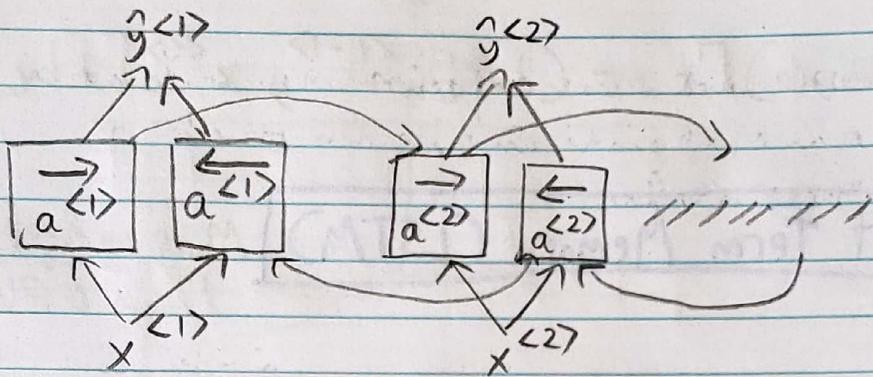
- $\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$
- $\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$ update gate
- $\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$ forget gate
- $\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$ output gate
- $c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$
- $a^{<t>} = \Gamma_o * \tanh(c^{<t>})$



RNNs

Bidirectional RNN (BRNN)

- It's an Acyclic graph



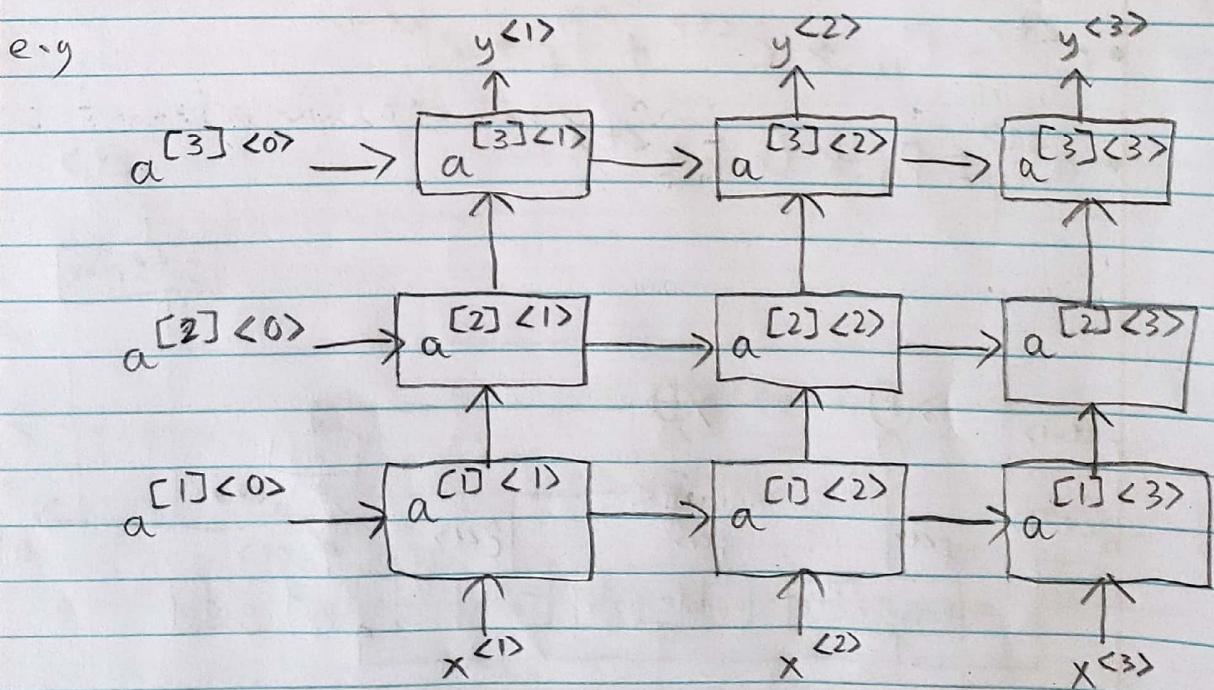
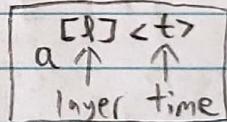
These blocks
can be GRU
LSTM

$$\hat{y}^{<t>} = g(W_y [\vec{a}^{<t>} \wedge \vec{a}^{<t>}] + b_y)$$

Disadvantage: Need entire sequence before being able to do a prediction

RNNs

Deep RNNs



Week 2 - Natural language processing & word embeddings

Intro
to
word
embeddings

Word Representation

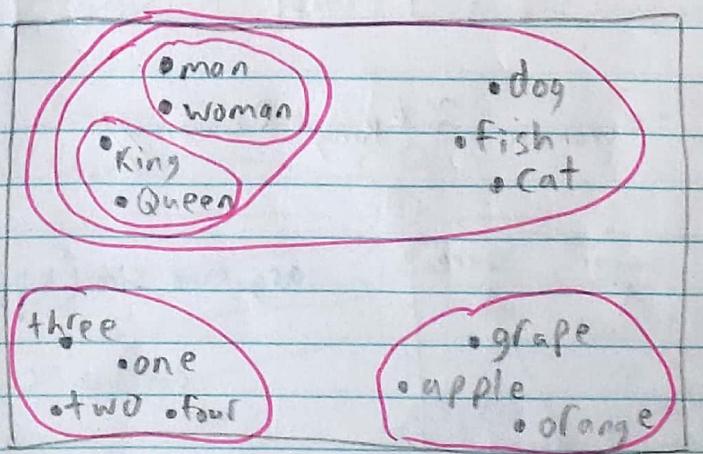
Instead of one-hot representation where words are not related to each other and are either 0/1 in the vocab, we can learn featureized representation

Vocab position	5391	9853	4914	7157	456	6257
	Man	Woman	King	Queen	Apple	Oranges
Gender	-1	1	-0.95	0.97	0.02	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
i	⋮					

300 dimensional vector to represent man
 e_{5391} (embedding vector)

→ Allows better generalizations across different words and better relationships

- Visualizing word embeddings



Learn similar features for concepts that feel like they should be related

Using word embeddings

- Transfer learning and word embeddings

- Learn word embeddings from large text corpus (1-100B words)
(or download pre-trained embedding online)
- Transfer embedding to new task with smaller set
e.g. Name entity recognition task, 100K words
- Optional: Continue to finetune the word embeddings with new data

Related to face encoding/embedding of siamese network of CNN
difference: lets say we have a fixed vocab of 10,000 words,
we'll learn fixed encodings for each of 10K words in
the vocab. In CNN, can learn encoding for any
input image

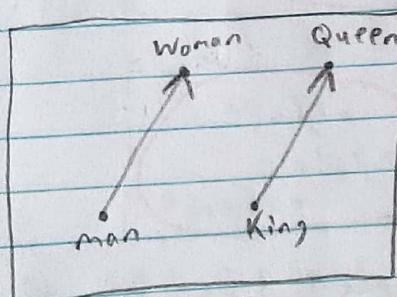
Properties of word embeddings

- Man → Woman as King → ?

$$e_{\text{man}} - e_{\text{woman}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix}$$

$$e_{\text{King}} - e_{\text{Queen}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix}$$

$$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{King}} - e_{\text{Queen}} //$$



300D

$$\arg \max \text{Sim}(e_?) e_{\text{King}} - e_{\text{man}} + e_{\text{woman}}$$

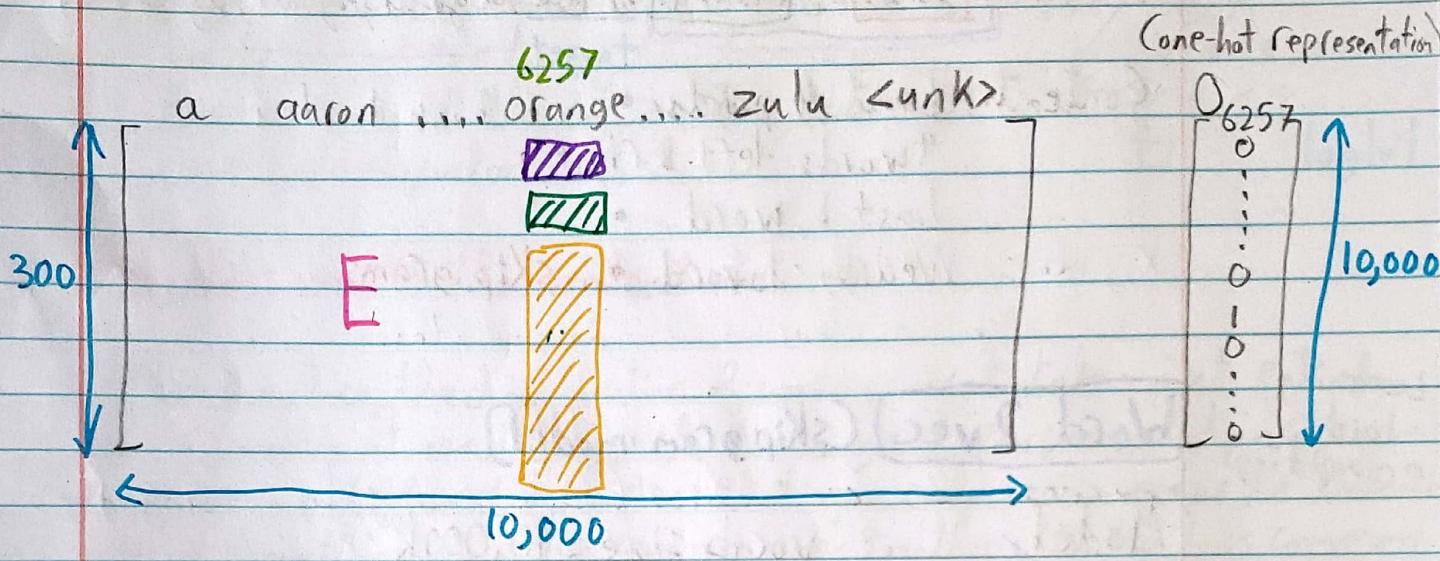
↓
could be cosine similarity

$$\text{Sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

cosθ

Intro to
word
embeddings

Embedding Matrix



$$E \cdot O_{6257} = \begin{bmatrix} \text{purple} \\ \text{green} \\ \text{orange} \end{bmatrix} = e_{6257} \quad (\text{embedding for orange})$$

(300,
10K)
1)

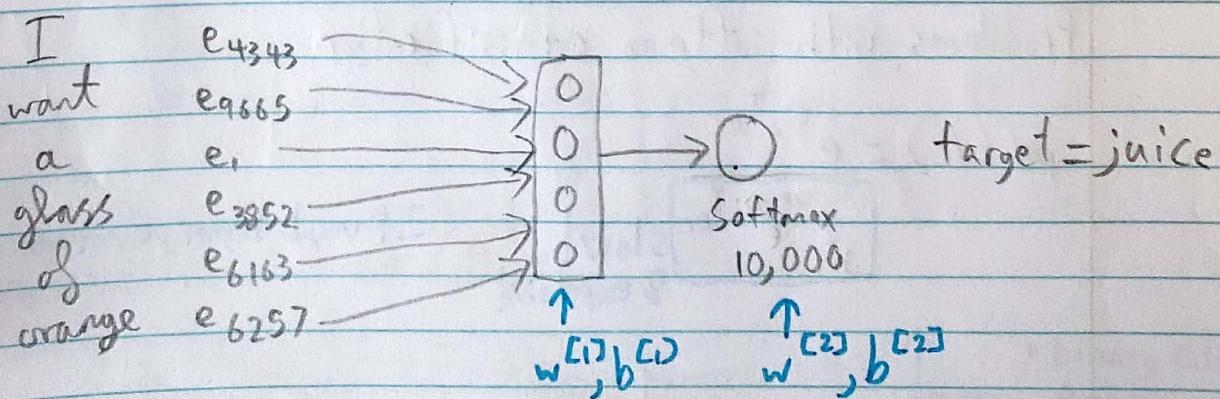
(300,
1)

→ In practice, use specialized function to look up an embedding as it's way more efficient

Learning
word
embeddings

Learning Word Embeddings

- Yoshua Bengio : A neural probabilistic language model 2003



- Other context/target pairs

e.g.
I want a glass of orange juice to go along with my cereal
target

Context: Last 4 words •

4 words left & right •

Last 1 word •

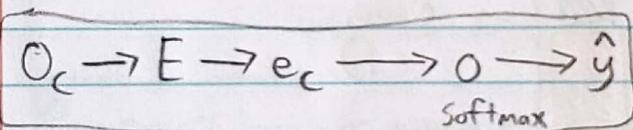
Nearby 1 word • skip grams

Learning
word
embeddings

Word 2 vec (Skipgram model)

- Model vocab size = 10,000K

We would like to learn mapping context $c \rightarrow$ target t



$$\text{Softmax: } p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

θ_t = parameter associated
with output t
(given by softmax)

$$\mathcal{L}(\hat{y}, y) = - \sum_{i=1}^{10,000} y_i \log \hat{y}_i$$

- Problems with softmax classification

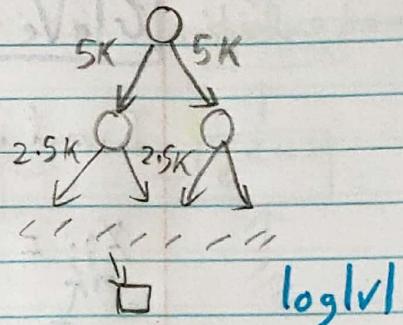
$$p(t|c) = e^{\theta_t^T e_c}$$

$\sum_{j=1}^{10,000} e^{\theta_j^T e_c}$ slow!
& expensive

If vocab bigger, even worse!

Potential solution: Hierarchical Softmax

→ Can be developed such that common words are at the top of the tree, and less common at the bottom. They don't have to be equal on both sides.



How to sample the context c?

→ In practice, the distribution of words pc isn't taken just entirely uniformly at random for the training set purpose, but instead there are different heuristics that you could use in order to balance out something from the common words together with the less common words.

common → to, the, a, from
uncommon → orange, durian, apple

Learning word embeddings

Negative Sampling

I want a glass of orange juice to go along with my cereal

<u>context</u>	<u>word</u>	<u>target</u>
orange	juice	1
orange	King	0
orange	book	0
orange	the	0
orange	of	0

x → context
y → target
K → number of negative samples

$$P(y=1 | c, t) = \sigma(\theta_t^T e_c)$$

$$\text{orange} \rightarrow \theta_{6257} \rightarrow E \rightarrow e_{6257} \rightarrow \begin{matrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{matrix}$$

• binary classification
• Trained on K+1
• 10,000 Problem

Learning word embeddings

GloVe word vectors

Global Vectors for word Representations

I want a glass of orange juice to go along with my cereal

$$X_{ij} = \# \text{ of times } j \text{ appears in context of } i$$

X_{ij} → count that captures how often do words i and j appear with each other or close to each other

$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\theta_i^T e_j + b_i + b_j - \log X_{ij})^2$$

symmetric

$\swarrow \downarrow \searrow$

weighting term

makes sure it doesn't give the frequent words too much weight & the less frequent too little weight

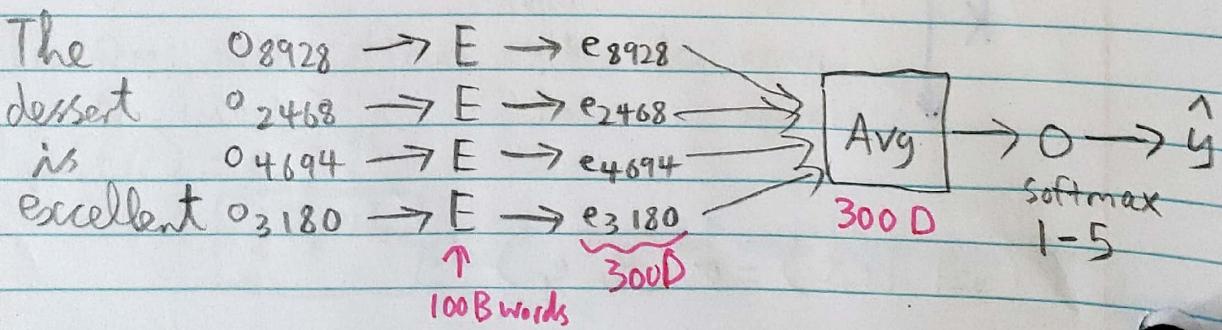
Applications using word embeddings

Sentiment Classification

Simple sentiment classification model

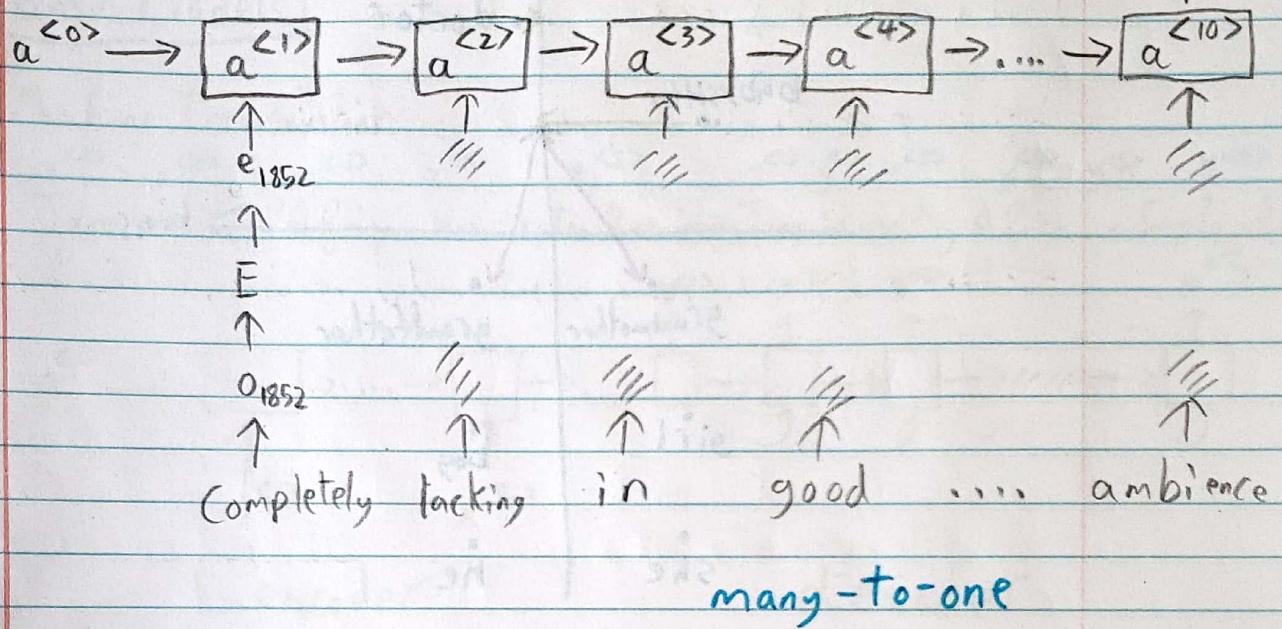
The dessert is excellent

8928 2468 4694 3180



Does not take word order into consideration

RNN for sentiment classification



→ Generalize much better towards words that were not in the training set used to train the word embeddings, if word was in the 100B word corpus used to train word embeddings.

Applications using word embeddings

Debiasing word embeddings

The problem of bias in word embeddings

Man:Woman as King:Queen

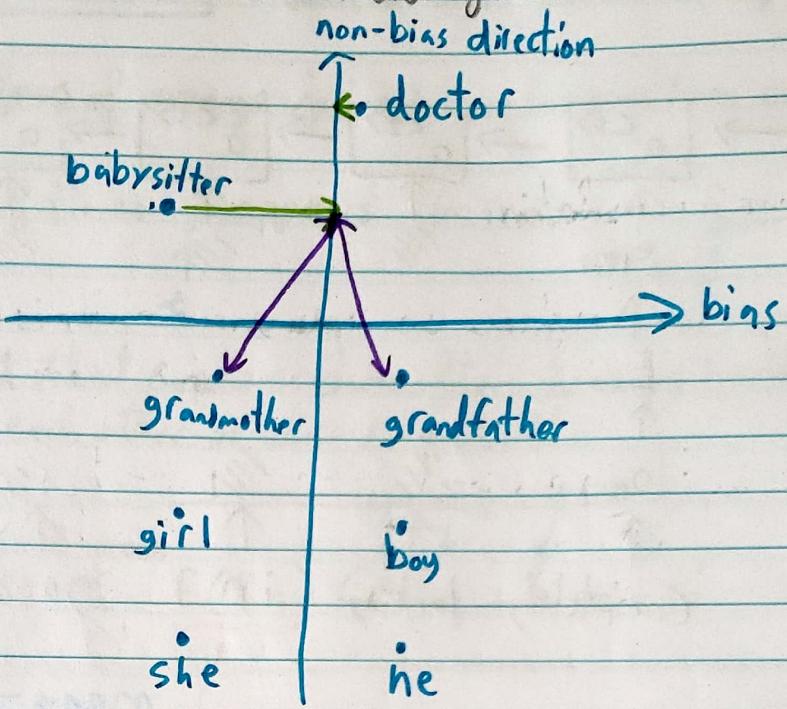
Man:Computer Programmer as Woman:Homemaker

Father:Doctor as Mother:Nurse

X

Word embeddings can reflect gender, ethnicity, age, sexual orientation and other biases of the text used to train the model.

Addressing bias in word embeddings



1. Identify bias direction

$e_{\text{he}} - e_{\text{she}}$
 $e_{\text{male}} - e_{\text{female}}$
 ...
 average

2. Neutralize : For every word that is not definitional, project to get rid of bias

3. Equalize pairs : grandmother - grandfather
 girl - boy
 Make sure
 Same distance

Week 3 - Sequence models & Attention mechanism

Basic Models

Machine Translation

$x^{<1>} x^{<2>} x^{<3>} x^{<4>} x^{<5>} y^{<1>} y^{<2>} y^{<3>} y^{<4>} y^{<5>} y^{<6>}$

Jane visite l'Afrique en septembre → Jane is visiting Africa in September

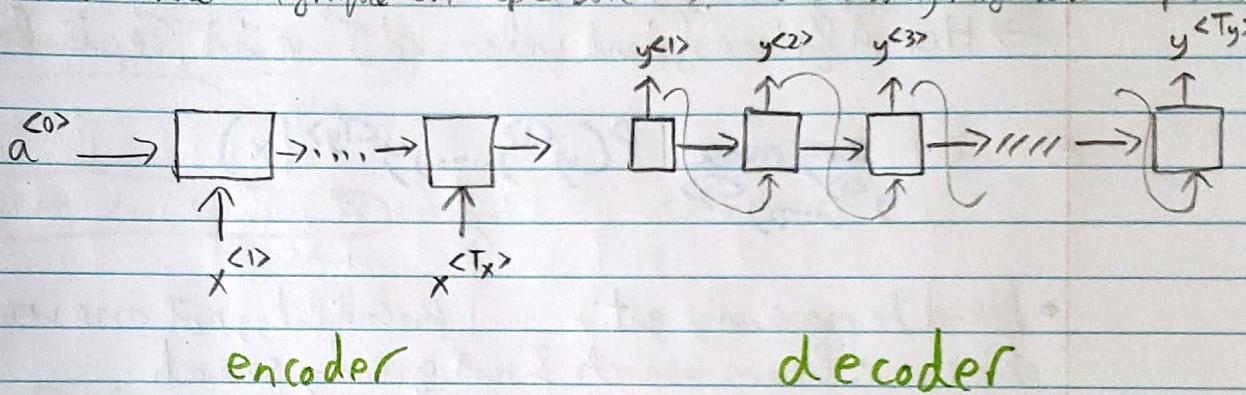
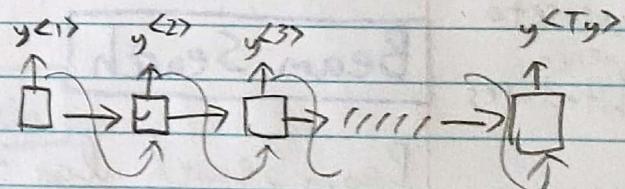
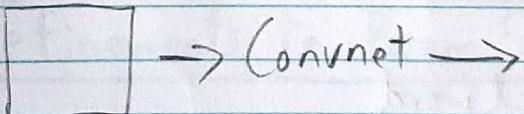


Image Captioning



Cat image

X

A cat sitting on a chair

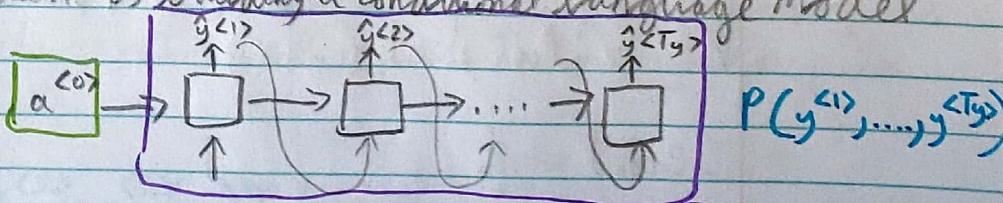
$y^{<1>} y^{<2>} \dots \dots \dots y^{<6>}$

Sequence
to
sequence
architectures

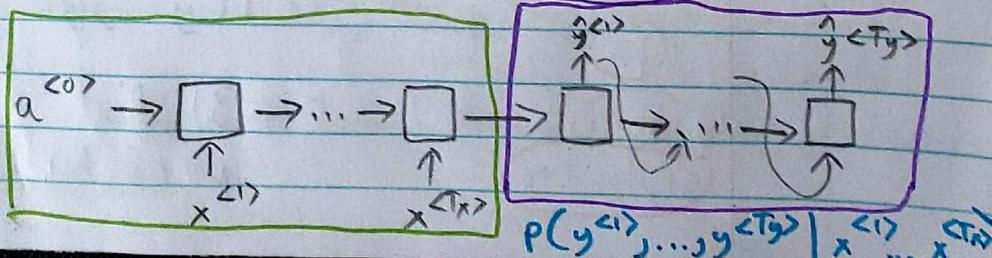
Picking the most likely sentence

Machine translation as building a conditional language model

Language model:



Machine translation:



- Finding the most likely translation

Jane visite l'Afrique en septembre

English French

$$P(y^{<1>} \dots y^{<T_y>} | x)$$

- Jane is visiting Africa in September
- Jane is going to be visiting Africa in September
- In September, Jane will visit Africa
- Her African friend welcomed Jane in September

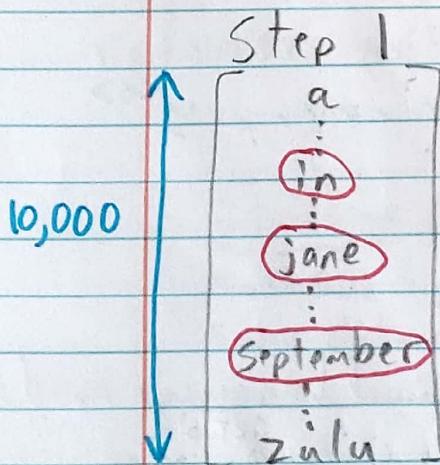
$$\arg \max_{y^{<1>} \dots y^{<T_y>}} P(y^{<1>} \dots y^{<T_y>} | x)$$

- Need to maximize entire probability, not one word at a time. ∴ Beam Search & not greedy search

Sequence-to-sequence architectures

Beam Search

- Beam Search algo - First word

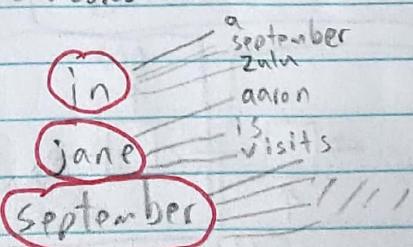


$B=3$ (BeamWidth)
Keeps the highest 3 possibilities in memory//

$$P(y^{<1>} | x)$$

(if $B=1$, this is equivalent to greedy search)

- Second word



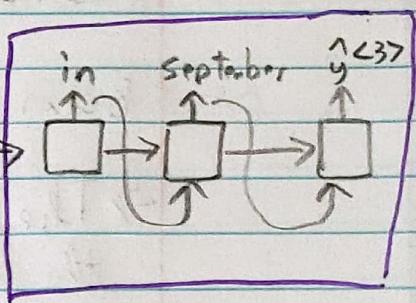
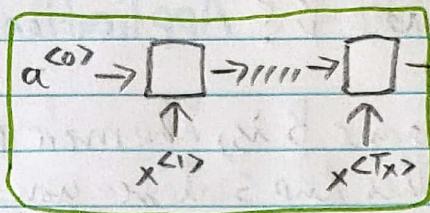
$$P(y^{<1>} y^{<2>} | x) = P(y^{<1>} | x) P(y^{<2>} | y^{<1>})$$

3 Highest possibilities: in september
jane is
jane visits

Third word

in september

a qalon
zulu



jane is

a visiting
zulu

|||||

jane visits

a africa
zulu

|||||

Outcome: jane visits africa in september. <EOS>

Refinements to Beam Search

Length normalization

$$\begin{aligned} \text{Right now } \arg \max P(y^{<1>} | \dots | x) \\ = P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>}) \dots P(y^{<Ty>} | x, y^{<1>} \dots y^{<Ty-1>}) \end{aligned}$$

→ Probabilities are all numbers less than 1. Often much less than 1. Multiplying them all together will result in a tiny, tiny number, which can result in numerical underflow. Meaning it's too small for the floating point representation in your computer to store accurately.

objective function

$$\therefore \arg \max \sum_{t=1}^{Ty} \log P(y^{<t>} | x, y^{<1>} \dots y^{<t-1>}) \text{ more numerically stable}$$

Normalization

Avg of the log of the prob of each word

The more terms, the more negative it becomes

- Beam search discussion

Beam width B ? Application & Domain dependent. 10, 100, 1000, 3000

- The larger B is, the more possibilities you're considering and the better the sentence you find. But the larger B is, the more computationally expensive the algorithm is, because you're also keeping a lot more possibilities around.
- Unlike exact search algorithms like BFS (Breadth First Search) or DFS (Depth First Search), Beam Search runs faster but it is not guaranteed to find exact maximum for $\arg \max P(y|x)$

Sequence to
sequence
architectures

Error Analysis in Beam Search

- Example

Jane visits l'Afrique en septembre

- Human: Jane visits Africa in September (y^*)
- Algorithm: Jane visited Africa last September (\hat{y})

Which one is at fault? → RNN?
→ Beam search?

$$\textcircled{1} \quad P(y^*|x) > P(\hat{y}|x) \quad \text{Beam search at fault}$$

$$\textcircled{2} \quad P(y^*|x) \leq P(\hat{y}|x) \quad \text{RNN model is at fault}$$

Then figure out what fraction of errors are "due to" beam search vs. RNN model

Bleu Score

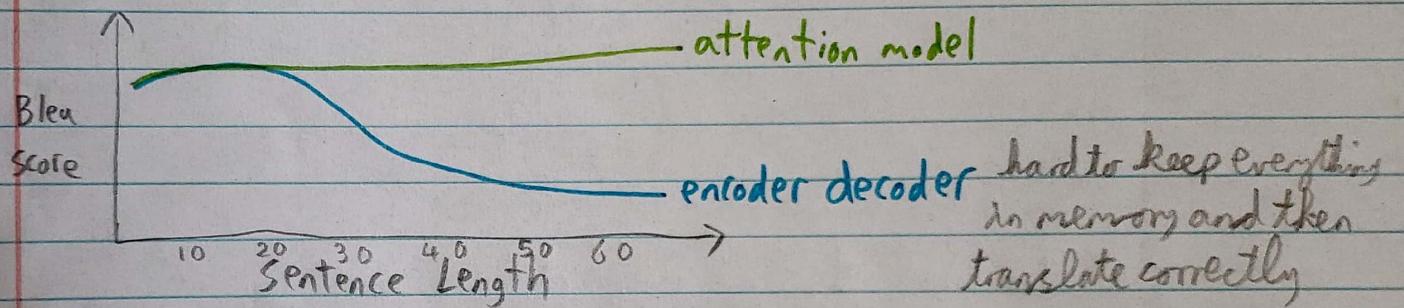
Bilingual Evaluation Understudy

- Measuring accuracy when there are a couple of great answers

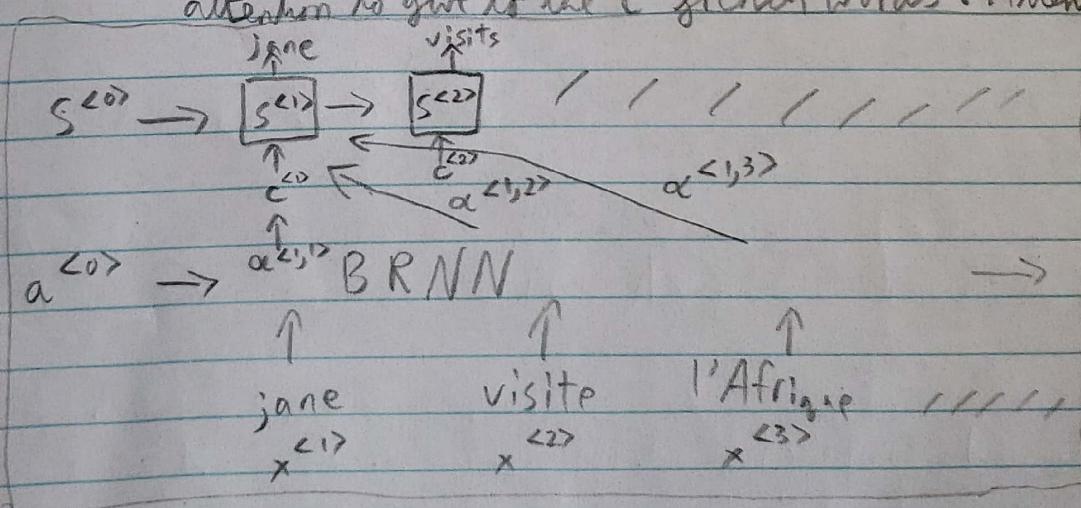
Videoooooo noisy to hear!

Attention Model Intuition

- The problem of long sequences



- $\alpha^{<t,t'>}$ when you're trying to generate t English words, how much attention to give to the t' French words. Attention weights!



$$c^{<1>} = \sum_{t'} \alpha^{<1,t'>} a^{<t'>}$$

$$c^{<2>} = \sum_{t'} \alpha^{<2,t'>} a^{<t'>}$$

$$\alpha^{<t,t'>} = \frac{\exp(e^{<t,t'>})}{\sum_{t'=1}^T \exp(e^{<t,t'>})}$$

Quadratic cost time!

