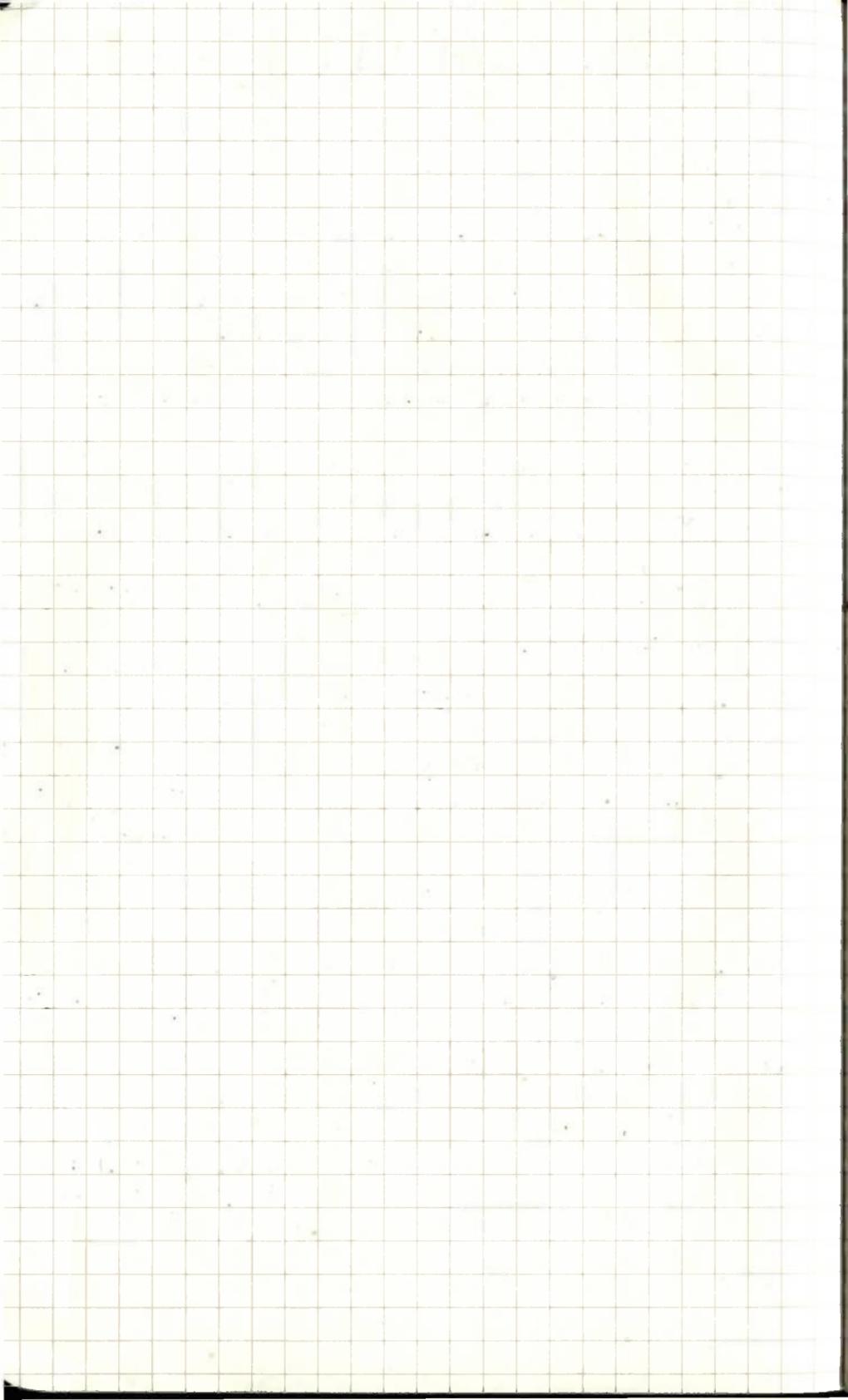


DEEP LEARNING 6



[Binary Classification]

As we know it's an algorithm for binary classification, is or is not ($y=1$ or $y=0$).

We might have an input of an image and determine if something is there; an image is stored by 3 separate matrices one for each color channel, i.e. $R \in \mathbb{R}^{n \times n}$, $G \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times n}$. Then we extend the matrices into a long vector (unroll), putting all values in a column. For example if the image is 64×64 pixels, its real dimension is $64 \times 64 \times 3 = 12288$ values, unrolled make a vector of 12288 pixels.

Notation to be used during the course:

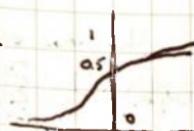
- Single training example: represented by a pair (x, y) where $x \in \mathbb{R}^{n_x}$, $y \in \{0, 1\}$
- Training sets: m training examples $(x^{(1)}, y^{(1)})$, $(x^{(2)}, y^{(2)})$, \dots , $(x^{(m)}, y^{(m)})$
↳ m_{train} m_{test}
- The training set inputs $X = \begin{bmatrix} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & x^{(3)} \\ 1 & 1 & 1 \end{bmatrix}^T \in \mathbb{R}^{n_x \times m}$ ↳ easier than transposing
- For output labels Y : stack them in columns $Y = [y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(m)}]^T \in \mathbb{R}^{1 \times m}$

[Logistic Regression]

Learning algorithm for classification, given ' x ', output the prediction ' \hat{y} ' what is the probability that ' y ' equals one with input x , $\hat{y} = P(y=1|x)$; $x \in \mathbb{R}^{n_x}$

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

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



$$\sigma(z) = \frac{1}{1+e^{-z}}$$

If z is large, $\sigma(z) \approx 1$
If z is small, $\sigma(z) \approx 0$

When we program neural networks, we usually keep the parameter w and parameter b separate, where here b correspond to a intercept.

In some conventions an extra feature is defined as $x_0 = 1$, so that now $x \in \mathbb{R}^{n+1}$ and then define $\hat{y} = g(\theta^T x)$, in this alternate convention we have pvector parameters θ from $\theta^{(0)}$ to $\theta^{(n)}$, and so $\theta^{(0)}$ plays the role of b and the rest as w :

$$x_0 = 1, x \in \mathbb{R}^{n+1}$$

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

$$\theta = \begin{bmatrix} \theta^{(0)} \\ \theta^{(1)} \\ \vdots \\ \theta^{(n)} \end{bmatrix} \quad b$$

however we want use this notation as it's more convenient to keep b and W separated.

[Logistic Regression Cost Function]

To recall: $\hat{y} = \sigma(w^T x + b)$ where $\sigma(z) = \frac{1}{1+e^{-z}}$, what $\hat{y} \approx y$

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, what $\hat{y}^{(i)} \approx y^{(i)}$

$$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b), \text{ where } \sigma(z) = \frac{1}{1+e^{-z}}, \text{ what } \hat{y}^{(i)} \approx y^{(i)}$$

$$x^{(i)} / y^{(i)} / z^{(i)}$$

Loss error function: $J(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$ (but not suitable for logistic regression)

so:

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

600¢ for logistic regression for a single example

This equation is balanced to determine a value for $y=1$ or $y=0$; by replacing y with one $\{0, 1\}$, one of the terms is eliminated and it gives a big or small value. If $y=1$ we want the value as big as possible, but if $y=0$, we want it as small as possible.

The cost function applies the loss error function for all the samples

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

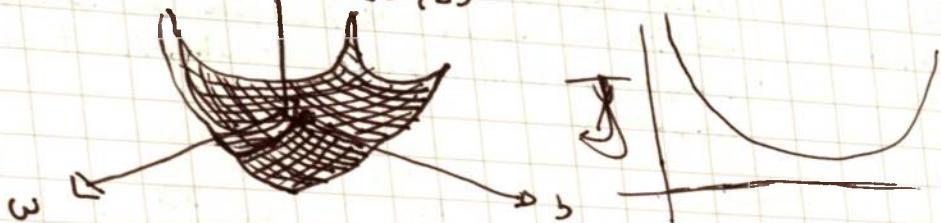
applied to parameters

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

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

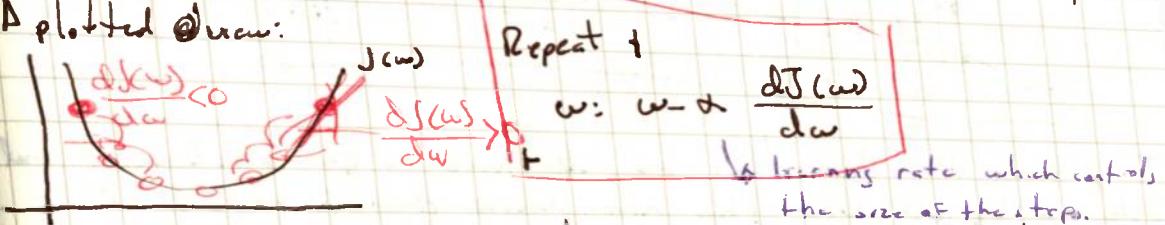
gradient]

Used to learn the parameters (w and b) of the cost function. So we want to find the combination of w and b that minimizes the value of $J(w, b)$.



After iteration, we tend to converge at the ~~local~~ global optima.

Plotted @ $J(w)$:



From calculus, we know the derivative tells us the slope of a point in a graph, thus by dividing the height and width of the point we get a triangle and it's possible to compute it's current slope; this applies on both ends, if the derivative is lower than zero, then it will move to the right \rightarrow bottom on each iteration; whereas, if the derivative is higher than zero, it will move to the bottom left after each iteration.

On the previous equation we defined for J of w , but our cost function is a function of both $J(w$ and b), so the inner loop of gradient descent ends up as

$J(w, b)$

$$\boxed{\begin{aligned} w &:= w - \alpha \frac{dJ(w, b)}{dw} \\ b &:= b - \alpha \frac{dJ(w, b)}{db} \end{aligned}}$$

how much the slope function slopes into w direction.

(Computation Graph)

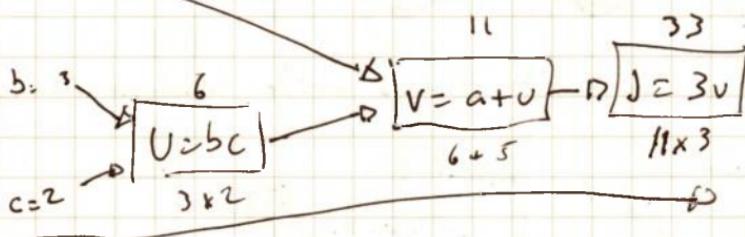
$$J(a, b, c) = 3(a + bc)$$

$$U = bc$$

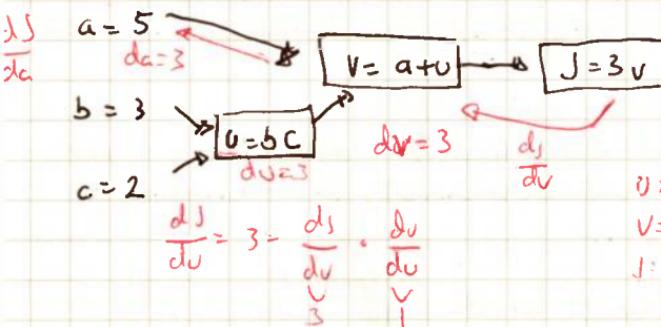
$$V = a + U$$

$$J = 3V$$

$$a=5$$



To perform the derivatives, we have to go backward and determine by how much the value changes with respect to a modification in the previous step. As values change one after another, the results are carried out by that fact as well the derivatives. This enables us to use the chain rule of derivatives. $a \rightarrow v \rightarrow J \quad \frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da}$



$$\begin{aligned} U &= 6 \rightarrow 6.001 \\ V &= 11 \rightarrow 11.001 \\ J &= 33 \rightarrow 33.003 \end{aligned}$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial b}$$

$$\begin{aligned} b &= 3 \rightarrow 3.001 \\ U &= bc \rightarrow 6 - 0.6.002 = c = 2 \\ J &= 33.001 \rightarrow V = 11.002 \\ &\quad \vdots \\ J &= 3V \dots \end{aligned}$$

[Logistic Regression Gradient Descent] - derivatives computation for G.D.

Logistic Regression recap:

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

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

$$\bullet L(a, y) = -(y \log(a) + (1-y) \log(1-a))$$

Supposing we have two features x_1 and their weights

$$\begin{aligned} & w_1 \\ & x_1 \\ & x_2 \\ & w_2 \\ & b \end{aligned} \quad z = w_1 x_1 + w_2 x_2 + b$$
$$\hat{y} = a = \sigma(z) = L(a, y)$$
$$\frac{\partial L}{\partial w_i} = \frac{dL}{dz} = \frac{\partial L(a, y)}{\partial z}$$
$$\frac{\partial a}{\partial z} = \frac{d\sigma}{da} = \frac{-y}{a} + \frac{1-y}{1-a}$$

$$\frac{\partial L}{\partial w_i}, \quad \frac{\partial L}{\partial w_i} = x_i \frac{\partial L}{\partial z},$$

$$\frac{\partial L}{\partial z} = a - y$$

$$\begin{cases} \frac{\partial L}{\partial w_i} = x_i \frac{\partial L}{\partial z} \\ \frac{\partial L}{\partial b} = \frac{\partial L}{\partial z} \end{cases}$$

[Logistic Regression on "m" examples]

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

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

And the derivative of the weight w_j is also equal to the sum of the derivatives

~~$$\frac{\partial J}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_j} (L(a_i^{(i)}, y^{(i)}))$$~~

Algorithm: ~~compute over 2~~ assuming $n=2$

$J=0$; $d\omega_1=0$, $d\omega_2=0$; $db=0 \rightarrow$ commutator

for $i=1$ to m

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

$$\alpha^{(i)} = g(z^{(i)})$$

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

$$d\omega^{(i)} = \alpha^{(i)} - y^{(i)}$$

$$d\omega_1 += x_1^{(i)} d\omega^{(i)} \quad \left. \begin{array}{l} n=2 \\ \text{for loop over the} \end{array} \right.$$

$$d\omega_2 += x_2^{(i)} d\omega^{(i)} \quad \left. \begin{array}{l} \text{features} \end{array} \right.$$

$$db += d\omega^{(i)}$$

$$J / m$$

$$d\omega_1 / m; d\omega_2 / m; db / m$$

At the end of the computation:

$$d\omega_1 = \frac{\partial J}{\partial \omega_1} \quad d\omega_2 = \frac{\partial J}{\partial \omega_2} \quad \text{and} \quad db =$$

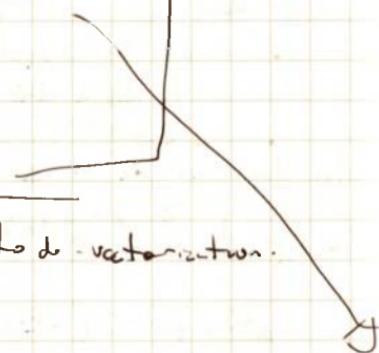
$$\omega_1 := \omega_1 - \alpha d\omega_1$$

$$\omega_2 := \omega_2 - \alpha d\omega_2$$

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

The for loops aren't optimal, so best is to do vectorization.

one step
of
gradient
descent



[Vectorization]

columns, F | Cs

As we try to optimize the calculation, taken in advantage the modern CPU architectures that allow the use of SIMD (Single Instruction Multiple Data), it's possible to take advantage of matrix operations, several of the steps can be greatly optimized.

It becomes possible to eliminate the use of for loops and increase the speed by magnitudes . Order of magnitudes.

To ~~first~~ compute $dZ^{(i)}$, we need to do $C^{(i)} - Y^{(i)}$, and as we have vectors representing the whole array of $dZ = [dZ^{(1)} \ dZ^{(2)} \ dZ^{(m)}]$ we can just do ~~with~~ ~~numpy~~:

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

$$dZ = A - Y = [a^{(1)} - y^{(1)}, a^{(2)} - y^{(2)}, \dots, a^{(m)} - y^{(m)}]$$

So with one line of code $\{dZ = A - Y\}$ we can compute it all at the same time.

Now in the case of W , we had to multiply $(X^{(i)})$ with $dZ^{(i)}$, this for all the features " n ". We can as well get rid of the for loop to vectorize it, applying matrix multiplication, which multiply and do the same sum.

Similar (the same) happens with b ! $db = \frac{1}{m} \sum_{i=1}^m dZ^{(i)}$, and in python we can do: ~~\star np.dot~~ $\boxed{1/m \cdot \text{np.sum}(dZ)}$,

$$\text{for } dw = \boxed{\frac{1}{m} X dZ^T} = \frac{1}{m} \begin{bmatrix} x^{(1)} & \dots & x^{(n)} \end{bmatrix} \begin{bmatrix} dZ^{(1)} \\ dZ^{(2)} \\ \vdots \\ dZ^{(m)} \end{bmatrix} = \frac{1}{m} [x^{(1)} dZ^{(1)} + x^{(2)} dZ^{(2)} + \dots + x^{(n)} dZ^{(n)}]$$

Vectorized Version:

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

$$= \text{np.dot}(W^T, X) + b$$

$$A = \boxed{Z}$$

$$dZ = A - Y$$

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

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

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

* Single iteration of gradient descent

Forward Prop and backward Prop.

Broadcasting in Python

When we do operations between matrices and real numbers, Python uses a feature called broadcasting, and what it does is to replicate the value (real, horizontal or vertical vector) to match the shape of the (m,n) matrix.

If we have for example:

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

or $(1,n)$

$$\begin{bmatrix} 123 \\ 456 \end{bmatrix} + [100 \ 200 \ 300] \text{ it becomes } \begin{bmatrix} 123 \\ 456 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

$$\begin{bmatrix} 123 \\ 456 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} \text{ becomes } \begin{bmatrix} 123 \\ 456 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

$(m,1)$

General principle

$(m,1)$ $\xrightarrow{\text{+}}$ $(1,n)$ $\Rightarrow (m,n)$ it will copy n times
 $(m,1)$ $\Rightarrow (m,n)$ it will copy m times.

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

$$\underline{[1 \ 2 \ 3] + 100 = [101 \ 102 \ 103]}$$

Notes on Python and vectors

Do not use one rank vectors, always define the dimensions, i.e.

$a = np.random.rand(5,1)$ = Column vector $\left| \text{assert(a.shape} == (5,1)) \right.$
 $a = np.random.rand(1,5)$ = Row Vector $\left| \text{to verify the shape is correct} \right.$

Logistic Regression test in Jupyter

train-set-x has shape (209, 64, 64, 3)

train-set-y has shape (209,) - float64
to 12cat
or non-cat

A = rows
B = height
C = width
D = R6B → 3 color representation

test-set-x
test-set-y

(50, 64, 64, 3)
(50,)

confirmed by
`str(... .shape)`

From (209, 64, 64, 3), `shape[n]` gives
`shape[0] = 209`

[1] = 64 ∴ on using shape reduces (0, 1, 2, 3)
[2] = 64
[3] = 3

I might be wrong but seems numpy handles the columns
and rows in an inverted way (n, m) instead
of mind as in Matlab

From (209, 64, 64, 3), we can flatten the dataset to get a
normal array of shape (64 × 64 × 3, 2), where each column will
represent a flattened image

Flattened = ~~normal.shape[0]~~

So $\text{normal.reshape}(\text{normal.shape[0]}, -1).T$
 $\text{normal.reshape}(209, 64, 64, 3).reshape(209, -1).T - 1$ = lots more to figure out
what to do with the rest.
in this case, one dimension with everything concatenated

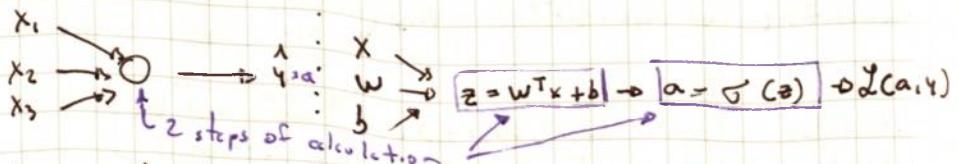
After the images in the arrays have been flattened, then we need to
standardize the dataset, meaning that we subtract the mean of the
whole numpy for each example and then divide each example by the
standard deviation of the whole numpy array. But for picture datasets,
it's simpler and more convenient to just divide every row of the dataset by
255, which is the maximum value for a pixel channel (0:255; 65535; 1655)

train-set-x = train-set-x.flatten / 255

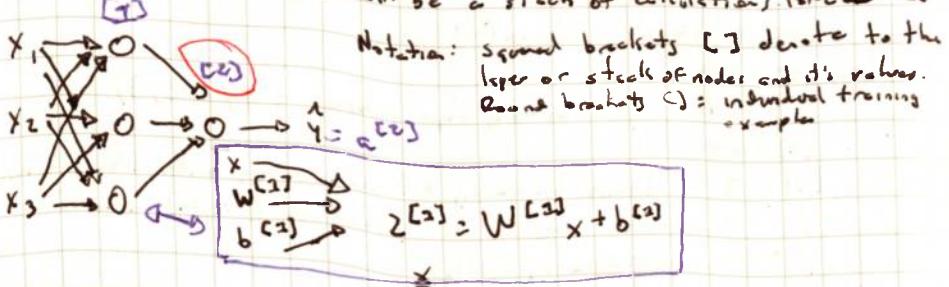
test-set-x = test-set-x.flatten / 255

[Shallow Neural Networks]

Last module we took logistic regression and we've seen how the following model worked:



In a neural networks, z will be a stack of calculations for all examples.



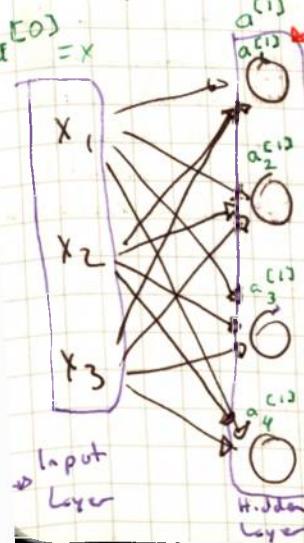
$$z^{[2]} = W^{[2]}x + b^{[2]} = a^{[2]} = \sigma(z^{[2]}) \rightarrow z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} = a^{[2]} = \sigma(z^{[2]})$$

$$\frac{\partial L}{\partial W^{[2]}} = \frac{\partial L}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial W^{[2]}}$$

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

The calculation w forward propagated across each layer $^{[l]}$, and also the back propagation will perform the derivatives across.

NN Representation



$$a^{[1]} = \begin{bmatrix} a^{[1]}_1 \\ a^{[1]}_2 \\ a^{[1]}_3 \\ a^{[1]}_4 \end{bmatrix}$$

Later on, we'll see that $a^{[1]}$ is a 4 by 1 vector, with four nodes

$$W^{[1]} = (4, 3) \quad b^{[1]} = (4, 1)$$

$$a^{[2]} = \begin{bmatrix} a^{[2]}_1 \\ a^{[2]}_2 \\ a^{[2]}_3 \\ a^{[2]}_4 \end{bmatrix}$$

Later on, we'll see that $a^{[2]}$ is a 4 by 1 vector, with four nodes

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

This is a 2 Layer NN, as the input layer isn't created.

As we know, each node performs two operations, it computes the weight vs feature multiplication + bias, and then the activation sigmoid function. That said in vectorized approach would be:

$$z = w^T x + b \rightarrow \sigma(z).$$

As this repeats on each node on the NN, taking in consideration as an example, the previous NN drawn, we get:

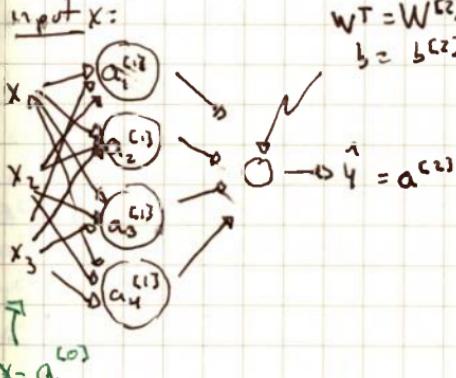
$$\begin{aligned} z_1^{(1)} &= w_1^{(1)T} x + b_1^{(1)}, & a_1^{(1)} &= \sigma(z_1^{(1)}) \\ z_2^{(1)} &= w_2^{(1)T} x + b_2^{(1)}, & a_2^{(1)} &= \sigma(z_2^{(1)}) \\ z_3^{(1)} &= w_3^{(1)T} x + b_3^{(1)}, & a_3^{(1)} &= \sigma(z_3^{(1)}) \\ z_4^{(1)} &= w_4^{(1)T} x + b_4^{(1)}, & a_4^{(1)} &= \sigma(z_4^{(1)}) \end{aligned}$$

Computing the above in for loops would be efficient, so we can vectorize them.

$$W^{(1)}$$

$$\begin{aligned} z^{(1)} &= \begin{bmatrix} w_1^{(1)T} \\ w_2^{(1)T} \\ w_3^{(1)T} \\ w_4^{(1)T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix} \rightarrow \begin{bmatrix} w_1^{(1)T} x + b_1^{(1)} \\ w_2^{(1)T} x + b_2^{(1)} \\ w_3^{(1)T} x + b_3^{(1)} \\ w_4^{(1)T} x + b_4^{(1)} \end{bmatrix} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix} \\ a^{(1)} &= \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ a_3^{(1)} \\ a_4^{(1)} \end{bmatrix} = \sigma(z^{(1)}) \end{aligned}$$

Given that, what we see for the first layer of the neural network, given the input x :



$$W^T = W^{(2)}$$

$$b = b^{(2)}$$

$$z^{(2)} = W^{(2)} x + b^{(2)}$$

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

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

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

[Vectorizing across multiple examples]

For $i = 1$ to m :

$$z^{(2)}(i) = W^{(2)}x^{(i)} + b^{(2)}$$

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

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

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

or $A^{(2)}$

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

$$A^{(2)} = G(Z^{(1)})$$

$$Z^{(2)} = W^{(2)}A^{(1)} + b^{(2)}$$

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

Capital X contains the matrix X with our training examples stacked in columns.

$$X = \begin{bmatrix} & 1 & \\ x^{(1)} & x^{(2)} & x^{(m)} \\ & 1 & \end{bmatrix} : \star$$

We went to lowercase to uppercase by stacking in columns:

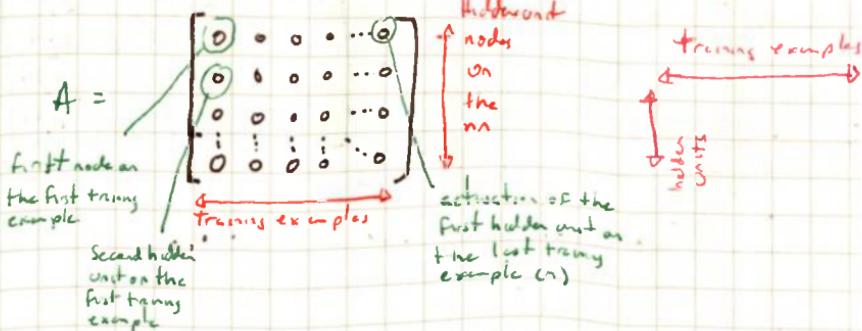
$$Z^{(1)} = \begin{bmatrix} z^{(1)(1)} & z^{(1)(2)} & z^{(1)(3)} & \dots & z^{(1)(m)} \end{bmatrix}$$

One of the properties of this notation

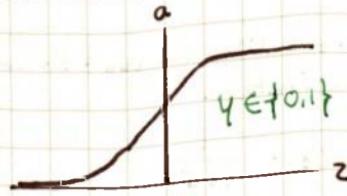
is that these matrices, horizontally contains

$$A^{(2)} = \begin{bmatrix} a^{(2)(1)} & a^{(2)(2)} & a^{(2)(3)} & \dots & a^{(2)(m)} \end{bmatrix}$$

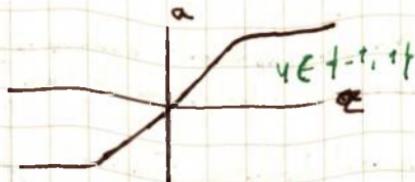
training examples and vertically the nodes in the NN.



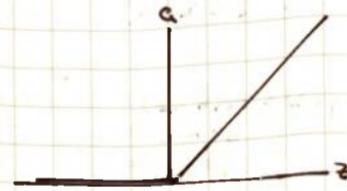
[Activation Functions]



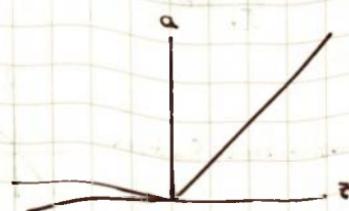
$$\text{Sigmoid: } a = \frac{1}{1 + e^{-z}}$$



$$\tanh = a_z \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



$$\text{ReLU } a = \max(0, z)$$



$$\text{Leaky ReLU } a = \max(0.01z, z)$$

Recommendations:

- Sigmoid: Don't use it, except at the first layer if doing binary classification.
- Tanh: much better, slightly superior
- ReLU: the default, the most common
- Leaky ReLU: Values greater or equal to zero (no negatives)

[Why do we need an activation function?] (non-linear)

To compute interesting functions with a neural network, we do need to pick a non-linear activation function.

If we decide to avoid the activation function and pass the values straight into the next layer, we are just composing another simple variation of the equation (linear), which:

$$\begin{aligned}
 z^{(l+1)} &= W^{(l+1)}x + b^{(l+1)} \\
 a^{(l+1)} &= g(W^{(l+1)}z^{(l+1)})z^{(l+1)} \Rightarrow a^{(l+1)} = z^{(l+1)} = W^{(l+1)}x + b^{(l+1)} \\
 z^{(l+2)} &= W^{(l+2)}a^{(l+1)} + b^{(l+2)} \\
 c^{(l+2)} &= g(W^{(l+2)}a^{(l+1)} + b^{(l+2)}) = g(W^{(l+2)}(W^{(l+1)}x + b^{(l+1)}) + b^{(l+2)}) = \\
 &= (W^{(l+2)}W^{(l+1)})x + (W^{(l+2)}b^{(l+1)} + b^{(l+2)})
 \end{aligned}$$

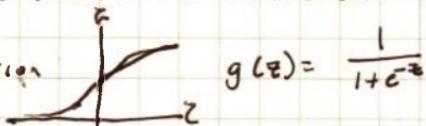
A linear hidden layer is more or less useless because the composition is itself a linear function ($y = mx + b$)

-Exception: If doing for a regression problem -

[Derivatives of activation functions]

When we implement back propagation for the NN, we need to compute the slope or derivative of the activation functions.

- Sigmoid Activation Function



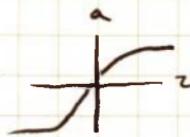
$$g(z) = \frac{1}{1+e^{-z}}$$

Alternative Notation:

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

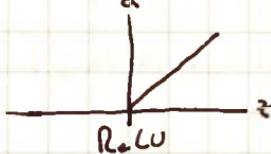
$$g'(z) = \frac{d}{dz} g(z)$$

• Tanh activation function



$$g(z) = \tanh(z) \\ = \frac{e^z - e^{-z}}{e^z + e^{-z}} \\ = 1 - (\tanh(z))^2 \\ g'(z) = 1 - a^2$$

• ~~ReLU~~ ReLU and Leaky ReLU



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

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \\ \text{undef if } z=0 \end{cases}$$



$$g(z) = \max(0.001z, z)$$

$$g(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

[Gradient descent for Neural Networks]

Forward Propagation:

$$\begin{aligned} z^{(1)} &= W^{(1)}X + b^{(1)} \\ A^{(1)} &= g^{(1)}(z^{(1)}) \\ z^{(2)} &= W^{(2)}A^{(1)} + b^{(2)} \\ A^{(2)} &= g(z^{(2)}) \end{aligned}$$

Back propagation: $y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$

$$\begin{aligned} dz^{(2)} &= A^{(2)} - Y \\ dW^{(2)} &= \frac{1}{m} dz^{(2)} A^{(1)T} \\ db^{(2)} &= \frac{1}{m} \text{np.sum}(dz^{(2)}, \text{axis}=1, \text{keepdims=True}) \\ dz^{(1)} &= W^{(2)T} dz^{(2)} * g'(z^{(1)}) \quad \begin{matrix} \text{sum across } \\ \text{the vertical axis} \\ (\text{axis}=1) \end{matrix} \\ dW^{(1)} &= \frac{1}{m} dz^{(1)} X^T \quad \begin{matrix} \text{element wise prod} \\ (\text{axis}=1) \end{matrix} \\ db^{(1)} &= \frac{1}{m} \text{np.sum}(dz^{(1)}, \text{axis}=1, \text{keepdims=True}) \quad \begin{matrix} \text{axis=1} \\ (\text{axis}=1) \end{matrix} \end{aligned}$$

[Random Initialization]

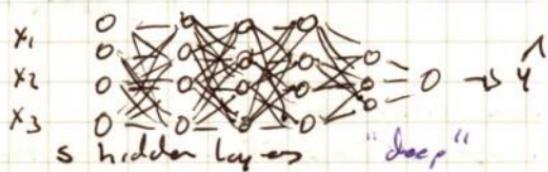
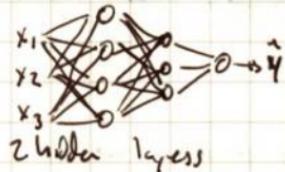
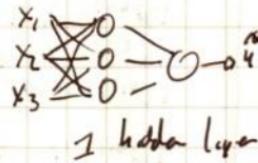
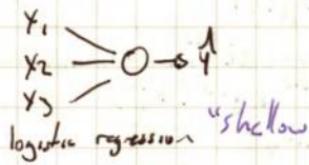
It's important for a NN to initialize the weights to some random values instead of zeros. Zeros works fine for linear regression but for NN's is a problem. The main issue is that the hidden units will compute the same values, which by the way, by multiplying by zero, they get annihilated.

To do it, we can use numpy \hookrightarrow

$$\begin{aligned} W^{(1)} &= \text{np.random.randn}(2, 2) * 0.01 \quad \begin{matrix} \leftarrow \text{very small random value} \\ \rightarrow \text{not go to the } \text{intelligence} \\ \text{of the model} \end{matrix} \\ b^{(1)} &= \text{np.zeros}(2, 1) \quad \text{as long as } W \text{ is random, this won't break strictly} \end{aligned}$$

[Week 4]

What is a deep neural network



Is good to try one by one, perhaps adding layers

Notation

L = Number of layers

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

In the sample: $a^{(25)}$ -units in layer I

L=4

$$n^{(0)} = n^{(1)} = 5, \quad n^{(2)} = 5, \quad n^{(3)} = 3, \quad n^{(4)} = n^{(5)} = 3$$

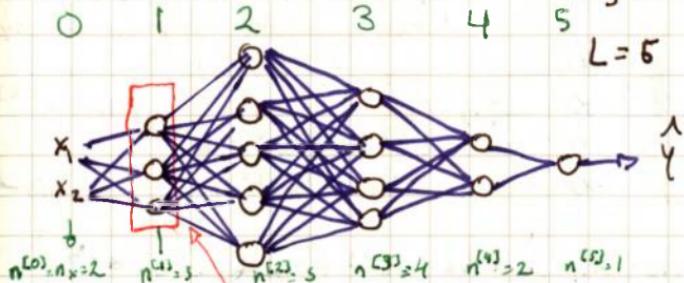
[Forward propagation in a Deep Network]

$$Z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)} \quad | \text{ } l=1 \text{ layer}$$

$$\alpha^{(k)} = g^{(k)}(\hat{z}^{(k)}) \quad \rightarrow \text{ok to use for loops on forward prop to iterate between layers}$$

[Getting the Matrix dimensions right]

When implementing a DNN, one of the tools to use is to check correctness of the code is to pull a piece of paper and just work thru the dimensions of matrix we are working with.



$$\text{Taking in consideration } z^{(l)} = W^{(l)} \cdot x + b^{(l)}$$

$$\begin{matrix} \frac{\partial}{\partial} \\ (3, 2) \\ (n^{(1)}, 2) \end{matrix} \quad \begin{matrix} \frac{\partial}{\partial} \\ (2, 1) \\ (n^{(2)}, 1) \end{matrix}$$

$$\begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \rightarrow \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \quad + \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

$$\begin{matrix} (3, 2) \\ (n^{(1)}, n^{(2)}) \end{matrix} \quad \begin{matrix} \swarrow \\ \uparrow \end{matrix}$$

$$W^{(1)}: (n^{(1)}, n^{(0)}) \quad \frac{1}{2} \quad \frac{1}{2}$$

$$W^{(2)}: (5, 3) = (n^{(2)}, n^{(1)})$$

$$z^{(2)}: W^{(2)} \cdot a^{(1)} + b^{(2)}$$

$$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \text{and so on...}$$

$$(2, 1) \rightarrow (5, 3) \leftarrow (3, 1)$$

$$w^{(2)} = (4, 5)$$

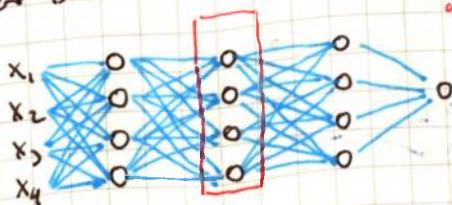
$$w^{(4)} = (2, 4), w^{(5)} = (1, 2)$$

$$W^{(l)} = (n^{(l)}, n^{(l-1)})$$

$$\text{and } b^{(l)} = (n^{(l)}, 1)$$

[Building Blocks of Deep Neural Networks]

We've already seen the basic building blocks of forward propagation and back propagation, the key components to build a deep net. Here is a network with a few layers. Let's pick a layer and look at the computations only for that layer.



Layer l : $W^{(l)}, b^{(l)}$

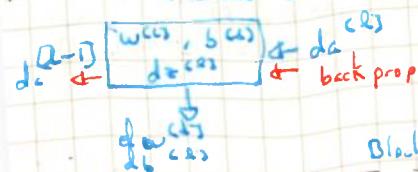
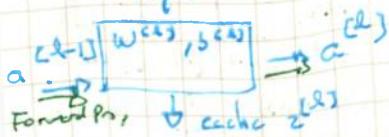
Forward: Input $a^{(l-1)}$, output $a^{(l)}$

$$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}$$

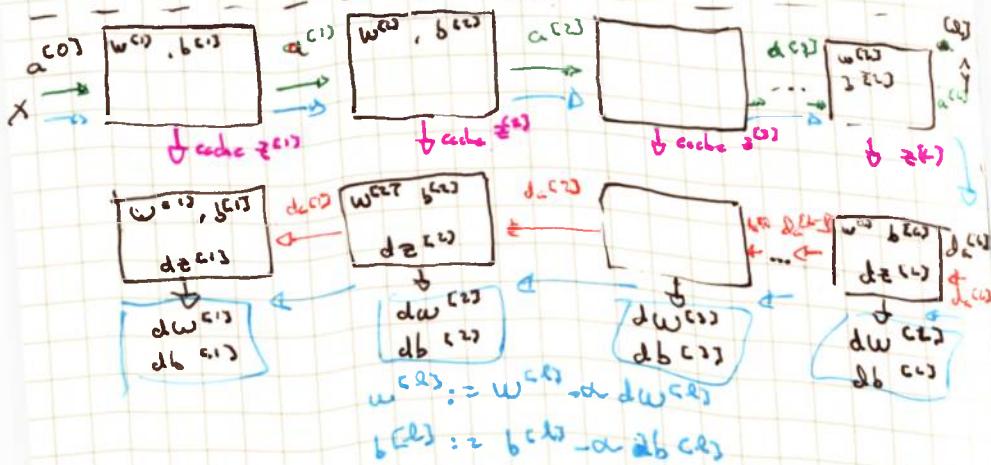
$$a^{(l)} = g^{(l)}(z^{(l)})$$

Turns out that for later use, would be useful to cache $z^{(l)}$, because would be useful for the backward propagation step. For the backward prop step, we will implement a function that inputs $d a^{(l)}$ and outputs $\frac{d a^{(l)}}{d a^{(l-1)}}$
 $cache(z^{(l)})$

Layer l



Block of the DNN



[Forward and Backward Propagation] - Implementation

→ Input $a^{[l-1]}$

→ Output $a^{[l]}$, cache ($z^{[l]}$)

$$\begin{aligned} z^{[l]} &= W^{[l]} \cdot a^{[l-1]} + b^{[l]} \\ a^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned}$$

Forward Pass

Vectorized:

$$\begin{aligned} z^{[l]} &= W^{[l]} \cdot A^{[l-1]} + b^{[l]} \\ A^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned}$$

Backward propagation for layer l

→ Input: $d_a^{[l]}$

→ Output: $da^{[l-1]}, dW^{[l]}, db^{[l]}$

$$da^{[l-1]} = da^{[l]} * g'(z^{[l]}) \quad \text{element wise prod}$$

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

$$db^{[l]} = dA^{[l-1]} \cdot a^{[l-1]T}$$

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

Vectorized:

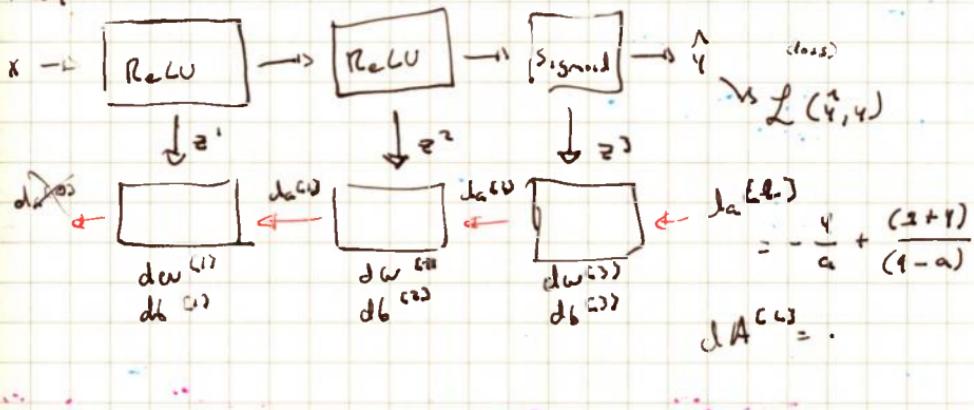
$$dA^{[l]} = dA^{[l-1]} * g'(z^{[l]})$$

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

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

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

Summary:



[Parameters vs hyperparameters]

Parameters are those given from the data: $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, \dots$

Hyperparameters are those that control the ultimate parameters W and b because they are controlling the result, they determine the final values of W and b :

Parameters: $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, \dots$

Hyperparameters:

- Learning rate α
- # iterations
- # hidden layers L
- # hidden units $n^{(1)}, n^{(2)}, \dots$
- choose activation function

There are other parameters in deep learning such as:

- Momentum
- Mini batch size
- Regularization parameters.

[Course 2: Improving Deep Neural Networks]

- Hyperparameter tuning,
- Regularization and
- Optimization

[Setting up the Machine Learning Application]

• Train / Dev. / Test sets

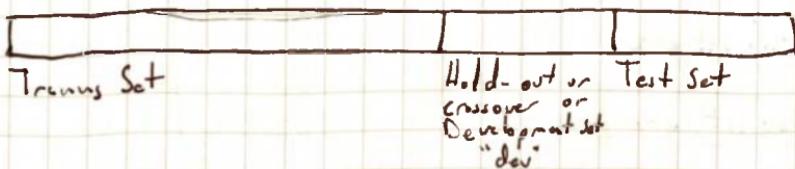
When training a NN, we have to make a lot of decisions such as, the amount of layers, nodes, activations for hidden layers, so in practice applying ML is a highly iterative process: code -> code -> experiments.

As there are more different fields of applications, is not just possible to bring the experience or intuitions from one field over to the next, for example, bringing knowledge from NLP into Computer Vision or from speech recognition into Ads, or from security into logistics.

One of the things that determine how quickly we can make progress is, how efficiently we can iterate the cycle code -> code -> test, and as the processing power is somehow bounded into the CPU's or GPU's capabilities, we want to make it fast.

Setting up our datasets well, in terms of train, development and tests, can help out on that step. At least to execute it faster. Dividing the data into 3 sets:

Dsets:



In the previous era of ML was common to split in 70/30% or 60/20/20, this was considered reasonable when the amount of data was about 1000 or ~~10K~~ 1K or 10K samples. In the current era of Big Data, where the amount of samples is extremely big, (ie: 1000000), the dev and test set has become a small percentage.

For example if we have 1M of samples, perhaps 1K (1%) is enough to have as dev and test purposes, i.e., 98/1/1%, there are cases where 99.5/0.25/0.25% or 99.5/0.4/0.1%.

(Mismatched train/test distribution)

It might happen that we train our NN with very high quality samples, for example, if the ML algorithm is supposed to recognize cats, we might get samples not only of professional pictures with perfectly framed cats, while the users of our app might take pics of their cats with low res phones, or the cats being blurry, out of frame, etc. So it's good idea to find realistic samples.

It's good idea to find realistic samples. Make sure dev and test sets come from the same distribution.

(Bias / Variance)

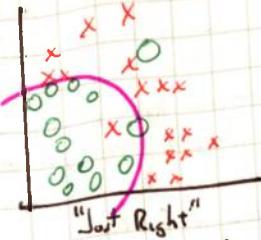
Good practitioners of ML have tended to be very sophisticated in understanding of bias and variance. It sounds simple but gets complex to master.



High Bias

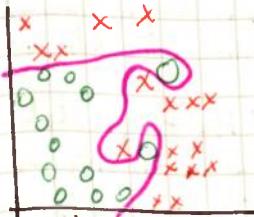
underfit

Doesn't capture grad anch



Captures good enough

"Just Right"



"high Variance"

Over Fit

too forced to capture all of it

Continuing with the cat classification, if we have a set of train/dev/exemplars. We would look at the loss numbers:

Train set error

Case 1
7%

Case 2
15%

Case 3
15%

Case 4
0.5%.

Dev set error

11%

16%

30%

1%

Has high
Variance

High
Bias

High Bias
and
High Variance

low bias
and
low Variance

This analysis is predicated on the assumption that human level performance gets near 0% error, or more generally, that the optimal error is nearly 0%. If the optimal error is 0.5%, then the case 2 would have $\frac{15}{0.5} = 30$ times worse error than the best case.

To take away, by looking at the training set error, we can get a sense of how well we are fitting at least the training data, and then we can look at how much higher is the error in dev set, that gives us sense how bad the variance is going.

[Basic Recipe for ML]

Ask

- Does the algorithm has high bias? \rightarrow the bigger net from longer

$\downarrow N$

• (the another NN architecture)

- High Variance? \rightarrow (overfit performance)

\downarrow

- More data \times

• Regularization

• (another NN architecture)

\downarrow
Done

the Bias-Variance tradeoff, there were not tools to reduce one without increasing the other.

- Using a bigger network, almost always reduces the Bias without hurting the Variance
- Getting more data, always reduces the variance and doesn't hurt the bias very much

[Regularization]

If we suspect the NN is overfitting the data (high Variance), we can use regularization, as training data might be too noisy.

Taking from the cost function

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i) + \frac{\lambda}{2m} \|w\|_2^2$$

L2 regularization $\rightarrow \|w\|_2^2 = \sum_{j=1}^{n_2} (w_j)^2 = w^T w$

omitted b, b^T
more subtle impact every
all the other parameters

Lambda λ is the regularization parameter, and in python we'll refer to it as lambd because the word lambda is a reserved word.

How we apply regularization in a neural network
In the NN, what has the cost function which is the function of all parameters:

$$J(w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i^{(2)}, y_i^{(2)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|^2$$

$\|w^{(l)}\|^2 = \sum_{i=1}^n \sum_{j=1}^{n_{l+1}} (w_{ij}^{(l)})^2$

From backprop:

$$dw^{(l)} := (\text{from backprop}) + \frac{\lambda}{m} w^{(l)}$$

$$w^{(l)} := w^{(l)} - \alpha dw^{(l)}$$

"Weight decay"

Because slowly we decay by reducing the dw by a little less than 1

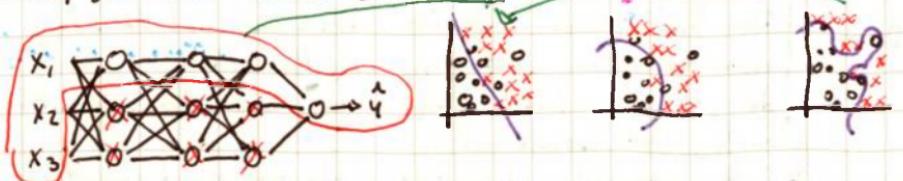
$$w^{(l)} = w^{(l)} - \alpha [(\text{from backprop}) + \frac{\lambda}{m} w^{(l)}]$$

$$= [w^{(l)} - \frac{\alpha \lambda}{m} w^{(l)}] - \alpha (\text{from backprop})$$

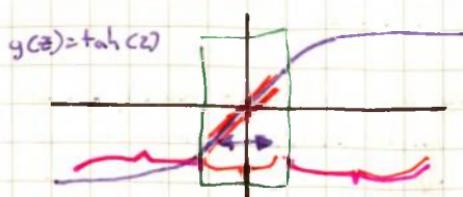
$$= 1 - \frac{\alpha \lambda}{m} \alpha$$

[How does regularization prevent overfitting?]

If we set the λ to be a high value, then $\frac{\lambda}{2} \sum_{l=1}^L \|W^{(l)}\|_F^2$ it will force the weight matrices W to be close to zero! $W^{(l)} \approx 0$ and that will cancel out several neurons, thus simplifying the NN, taking us from a complex function with high variance (i.e., complex n-degree function) into a simpler and lower-grade function, reducing the impact of a lot of hidden units, as if we were using logistic regression. They, what happens is that they just have smaller effect.



Another visualization with the tanh function:

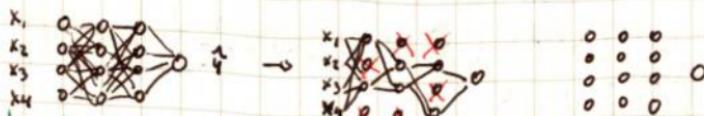


As long as z is quite small, so if z takes a small set of values, then it's using only the linear regime of the tanh function. If z uses larger values it starts to become less linear.

So, if t takes a small set of feature values, $z^{(l)} = W^{(l)}a^{(l-1)} + b^{(l)}$ or $g(z)$ will be roughly linear. And if the error is linear, even a DNN, it will only calculate a linear function.

[Dropout Regularization]

Dropout consists into randomly drop out (inactivated) hidden nodes, so it simplifies the network. Then apply back prop on them and repeat.



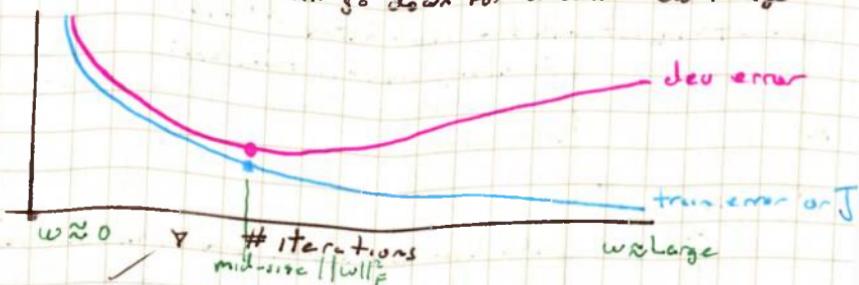
In code: For example, with layer $l=3$, and "keepprob" = 0.8 (80% to keep 20% drop).
 $d3 = np.random.rand(a3.shape[0], a3.shape[1]) < keepprob$
 $a3 = np.multiply(a3, d3)$ # $a3 \neq d3$
 $a3 /= keepprob$

↓ bajamos los de los valores al verificar cuáles son mejores máquinas a keep-prob, esto da una matriz booleana. Para conservar los valores en $a^{(l-1)}$, subimos los valores de las radios "keep" en 80%.

At test time we should not use dropout, otherwise we will get just noise on the predictions.

[Early Stopping]

As we run gradient descent, plot the cost function J , with early stop, we also plot our dev set error, and what we find is that the dev set error will go down for a while and then go up.



We stop halfway when we see it was performing better, by stopping halfway, we will have a mid-size value of w and similar to L2 regularization, it's likely at that point it will overfit less.

It does have a downside (early stopping), because we should consider that lowering the cost of function J is one thing and regularizing the dev set is completely different thing. This is called **orthogonalization** and its point is about thinking or doing one task at a time.

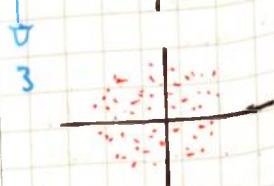
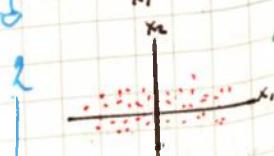
Doing early stopping, breaks whatever we are doing to reduce the cost function J and also simultaneously try to not overfit.. instead of using two tools to solve different problems, we are mixing them into one.

[Normalizing test sets]

One of the techniques in NN, that will speed up the trains, is if we normalize the inputs. Let's see what that means:

Imagine a training set with two input features $x > [x_1 \ x_2]$

Normalizing consists in two steps:
 1: Subtract mean: $\bar{x} = \frac{1}{N} \sum_{i=1}^N x^{(i)}$, then $x' = x - \bar{x}$ for every example x
 2: Normalize



We have now x' , every example has zero mean,
 then the second step is:
 2: Normalize the variance: $\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x'^{(i)})^2$ element wise square
 now σ^2 is a vector with the variances for each of the features.

Now we take every example and divide it by sigma squared vector σ^2 : $x'' = x'/\sigma^2$ and we end up with

IMPORTANT: If we use this to scale the training data, then use it as well on the test set, the same values of N and σ^2 . They have to be normalized/scaled in the same way because we want them defined by the same transformation.

[Vanishing / exploding gradients]

When we have a very deep NN, we have checked several multiplications which are by one will grow exponentially.

If the values of the weights are greater than 1 (3.0, 1.5, etc) then after several layers, the value will grow greatly. Similar happens when the weight is less than 1, 0.9 or less, exponentially, 1 will decrease to a value so small. Therefore the gradients might explode or vanish.

$$\begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ x_1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ x_2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix}$$

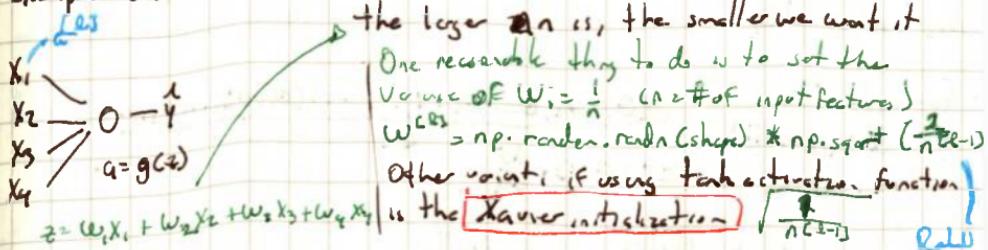
$$d \cdot n \cdot w$$

This makes training slow and difficult, but it can be partially solved correctly picking the correct choice of the random initialization.

[Weight Initialization for Deep Networks]

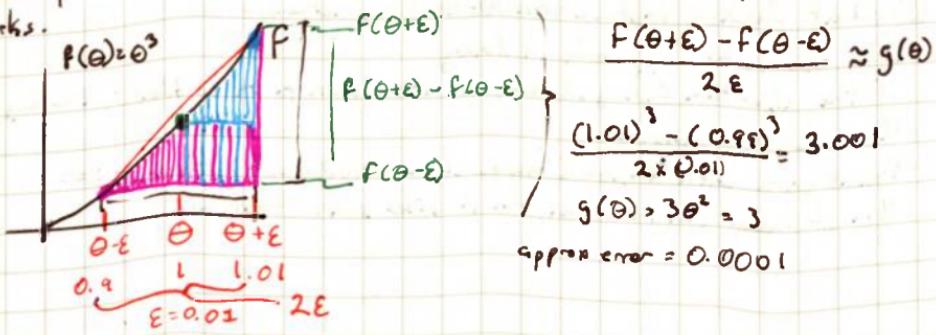
To minimize the effects of the vanishing/exploding gradients, it's possible to adjust the weights in each layer of the DNN by taking in consideration the values in the previous layer.

Example with a single neuron.



[Numerical Approximation of Gradients]

When we implement back propagation, we will find that there is a test called gradient checking that can help us check that the back propagation is correct. This is because sometimes we write all the equations and we might not be 100% sure if we got all the details right. So in order to build up the gradient checking, let's see how numerical approximation works.



The value is very close to the real derivative of $F(\theta)$, the real value is 3.0301, thus we end with a real error of 0.03, which just for checking is good enough but not for real calculations. Also an important note is that running it together with the back prop, will duplicate the time required to compute and converge, so the test should be disabled after the calculations are verified to be implemented properly.

[Gradient check for a neural network]

The NN will have some set of parameters $W^{(1)}, b^{(1)}, \dots, W^{(L)}, b^{(L)}$. To take them and reshape them into a big vector θ^* , so we need $W^{(1)}$ matrix and make it vector, same for $b^{(1)}$, same for $W^{(L)}$ and $b^{(L)}$ by taking concatenates them. so instead of $J(w_1, b_1, \dots, w_L, b_L)$, we get $J(\theta^*)$ and then Then do the same for $dW^1, db^1, \dots, dW^L, db^L$ and reshape them into a big vector $d\theta^*$.

$$\text{for each } i = \\ d\theta_{\text{approx}}[i] = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i, \dots)}{2\epsilon} \\ \approx d\theta^{(i)} \cdot \frac{dJ}{d\theta^i}$$

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

in practice $\epsilon = 10^{-7}$, so if the formula gives 10^{-7} or less then $\approx 0\%$
 $10^{-5} \dots 10^{-4}$
 10^{-3} - warning!

[Gradient Checking Implementation Notes] - practical.

- Don't use in training, only in debug
- If the algorithm fails gradcheck, look at components to try to identify the bug
- Remember the regularization
- Doesn't work with dropout (because elements random sets), turn it off
- Run at random initialization, run it after some iterations.

Week 6

[Optimization Algorithms] [Mini-Batch Gradient Descent]

As trains on large datasets is slow, and considering that we might need to run the model several times to find the right set of parameters and optimization algorithms. With the big data in our time, deep learning and its deep layers, even by using vectorization, still takes a long while to complete the training, and as it's empirical, we need to optimize it.

As we know, vectorization allows us to efficiently compute on an example, so regarding the training set we would have:

$$\begin{aligned} X &= \left[\underbrace{X^{(1)}, X^{(2)}, X^{(3)}, \dots, X^{(1000)}}_{X^{1000}} \right] \left[\underbrace{X^{(1001)}, \dots, X^{(2000)}}_{X^{1000}} \right] \dots \\ &\quad \left[\dots \underbrace{X^{(m)}}_{X^{1000}} \right] \\ Y &= \left[\underbrace{y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(1000)}}_{Y^{1000}} \right] \left[\underbrace{y^{(1001)}, \dots, y^{(2000)}}_{Y^{1000}} \right] \dots \\ &\quad \left[\dots \underbrace{y^{(m)}}_{Y^{1000}} \right] \end{aligned}$$

If we have a dataset of $m=5000,000$, we could split it in groups for example of 1000 samples, thus we end up with 5000 smaller groups or samples; for this we will introduce a new notation, "curly braces" and each group of 1000 samples are referring to; batches are separated as well as X only. Getting for example a minibatch $t = \{X^{t+1}, Y^{t+1}\}$

The process to follow is as usual with the difference if it were all the training set but in this case it's limited to one pass on the batch only, one

for $t=1, \dots, 5000$

$$\begin{aligned} \text{Forward prop on } X^{t+1} & \sim \\ Z^{(t+1)} &= W^{(t+1)} X^{t+1} + b^{(t+1)} \\ A^{(t+1)} &= g(Z^{(t+1)}) \end{aligned}$$

Vectorized implementation
for 1000 examples

$$A^{(t+1)} = g^{(t+1)}(Z^{(t+1)})$$

for X^{t+1}, Y^{t+1}

$$\text{Compute cost } J^{t+1} = \frac{1}{1000} \sum_{i=1}^{1000} L(Y_i, A^{(t+1)}) + \frac{\lambda}{2 \cdot 1000} \sum_{i=1}^{1000} \|W^{(t+1)}\|^2$$

$$\text{Backprop to compute gradients with respect to } J^{t+1} \text{ using } (X^{t+1}, Y^{t+1})$$

$$W^{(t+1)} = W^{(t+1)} - \alpha dW^{(t+1)}, \quad b^{(t+1)} = b^{(t+1)} - \alpha db^{(t+1)}$$

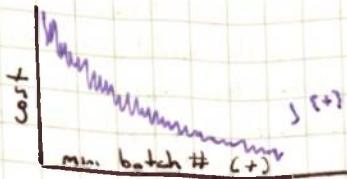
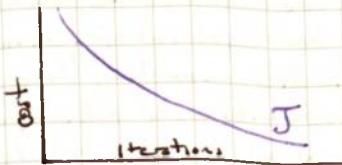
The code above is 1 pass to the training set, this is called epoch and it means a single pass thru the training set.

With mini-batch gradient descent a single pass thru the training set (1 epoch) allows us to take "1" gradient descent steps (1000 i.e.). Where as normal training, a gradient descent step would be done after every 1000 steps.

[Preliminary content]

[Understanding batch gradient descent]

Training a mini-batch gradient descent, would look like:



The reason of the bumpy costs on minibatches, is that each minibatch has different samples and some of them might be more expensive than others.

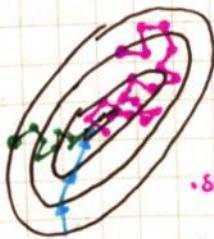
[Choosing the minibatch size]

There are two extreme cases:

If the minibatch size is m (all samples): We end up with batch gradient descent

If the minibatch size is equal to 1: Stochastic gradient descent (using sample i is own)

Let's look what these cases would do while trying to optimize the cost function, on the contour functions.



- Will take large steps towards the global minima.
- Will start randomly and will jump wildly, it will never converge but will end close.

Stochastic Gradient Descent
↳
Looses speed from vectorization

- In between (batch size depends)
 - fastest learning
 - Vectorization ✓
 - Progress without waiting for long

Batch Gradient Descent
↳
too long iteration

• In practice: Something in between • and ♡

[Guidelines:]

- if training set is small: batch gradient descent
($m \leq 2000$)

- big training set: $64, 128, 256, 512$ (powers of $2, 2^6, 2^7, 2^8, 2^9$)

- Make sure the minibatch fits in CPU/GPU

Thus, however, are not the best approaches, not bad, but there are better than gradient descent.

[Exponentially Weighted Averages]

This is a key component of several optimization algorithms.

Let's take a sample of the average temperature during a year

$$\theta_1 = 4.4^\circ C$$

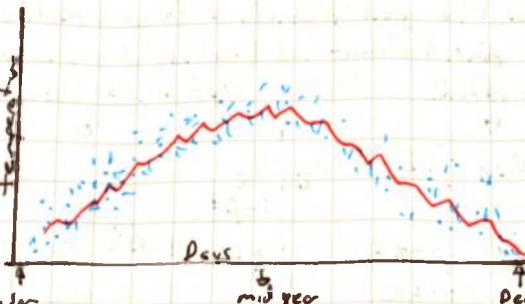
$$\theta_2 = 9.4^\circ C$$

$$\theta_3 = 7.2^\circ C$$

\vdots

$$\theta_{100} = 15.6^\circ C$$

$$\theta_{101} = 13.33^\circ C$$



$$V_0 = 0$$

$$V_1 = 0.9V_0 + 0.1\theta_1$$

$$V_2 = 0.9V_1 + 0.1\theta_2$$

$$V_3 = 0.9V_2 + 0.1\theta_3$$

$$\vdots$$

$$V_t = 0.9V_{t-1} + 0.1\theta_t$$

$$+ \dots$$

The data looks a bit noisy and if we want to compute the trends of the local average or moving average we initialize $V_0 = 0$ and then average it with a weight of 0.9 times whatever was the previous value, plus 0.1 times that day temperature.

The formula then becomes $V_t = 0.9V_{t-1} + 0.1\theta_t$, which generalizing we can call the 0.9 weight as beta β , (being previously equal to 0.9), so we have:

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

Averaging, we can think as $V_t \approx \frac{1}{1-\beta}$ days (samples), so for example

$$\beta = 0.9 : \approx 10 \text{ days temp}$$

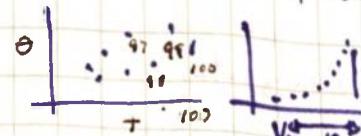
$$\beta = 0.98 : \approx 50 \text{ days temp} - \text{giving priority to the past values and reducing the current active's weight}$$

$$\beta = 0.5 : \approx 2 \text{ days temp} - \text{less influence from the past and more to come}$$

$$V_{100} = 0.1\theta_{100} + 0.9(V_{99}) \cdot (0.1\theta_{99} + 0.9V_{98}) \dots 0.1\theta_{98} + 0.9V_{97} \dots \\ = 0.1\theta_{100} + 0.1 \times 0.9 \cdot \theta_{99} + 0.1(0.9)^2 \theta_{98} + 0.1(0.9)^3 \theta_{97} + 0.1(0.9)^4 \theta_{96} + \dots \text{and all of this coefficients add up to one up to a detail called bias correction.}$$

$$0.9^{10} = 0.35 \approx \frac{1}{e}$$

$$\frac{(1-\varepsilon)}{0.9} = \frac{1}{e}$$



Implementation:

$$\theta_0 = 0$$

$$= \beta V_0 + (1-\beta)\theta_1$$

$$= \beta V_1 + (1-\beta)\theta_2$$

$$= \beta V_2 + (1-\beta)\theta_3$$

\dots

In programming

$$V = 0$$

for (each) :

get next θ_t

$$V = \beta V + (1-\beta)\theta_t$$

}

[Bias Correction]

If we use exponentially weighted averages, our start will be impacted greatly on the values, as the $V_0 = 0$, ~~over~~ over the first iteration and subsequent iterations will be greatly penalized. Because multipliers by β will cancel out the first element and the second one is very small (the one that is added).

Optimized (Bias correction)

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

$$V_0 = 0$$

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

exponent

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

$$\frac{V_2}{0.0396} = \frac{0.0196(\theta_1) + 0.02\theta_2}{0.0396} = 0.0396$$

[Gradient Descent with Momentum]

The basic idea is to use the Exponential Weighted Averages from the gradients and then use that gradient to update the ~~weight~~ ^{weight} itself.



If we start gradient descent here and after some iteration of gradient descent either batch or minibatch, it will take large steps and slowly oscillate toward the minima, and if we end up using much larger learning rate, we may end up overshooting and diverging.

Gradient descent with momentum:

On iteration t:

Compute dW, db on the current batch

$VdW = \beta VdW + (1-\beta)dW$ - Moving Avg

$Vdb = \beta Vdb + (1-\beta)db$

$W = W - \alpha VdW, b = b - \alpha Vdb$

$\beta = 0.9$ (usually) in practice, Bias Correction is not used while momentum can

Another goal to look at it, is that we want in the vertical axis to have a slower learning, whereas on the horizontal axis, we want faster learning. By using exponential weighted averages, we will smooth down and average the steps of gradient descent.

A way to rationalize the effect of the exponential weighted averages's impact on gradient descent, is to see the derivatives dW, db as providing ~~gradient~~ ^{acceleration} to the descent, while β gives friction and VdW, Vdb provide speed:

$$VdW = \beta VdW + (1-\beta)dW$$

[RMS Prop]

Root-mean-square Prop has very similar behavior to gradient descent with momentum, however it makes use of squares and square roots to penalize the weights and aid or reward the biweights, with the same purpose, to optimize the descent and approximation to the global minima, while avoiding to bounce laterally.

On iteration t :

Compute dW , db on current mini-batch element

$$SdW = \beta_2 SdW + (1 - \beta_2) dW^2$$

$$Sdb = \beta_2 Sdb + (1 - \beta_2) db^2$$

$$w = w - \alpha \frac{dW}{\sqrt{SdW} + \epsilon}, b = b - \alpha \frac{db}{\sqrt{Sdb} + \epsilon}$$

$\beta_2 = \text{to not mix it with } \rho \text{ from EWA}$

$\epsilon = \text{Very small value to avoid numerical error}$
 $\text{if } SdW \text{ or } Sdb \text{ become too small } \approx \epsilon = 10^{-8}$

[Adam Optimizer Algorithm] = ADapted MOMent ESTIMATION

The ADAM optimization algorithm comes from the fusion of Gradient Descent with Momentum and RMS Prop. This algorithm has proven to be very efficient in several applications of ML. The implementation is as follows:

$$VdW=0, SdW=0, Vdb=0, Sdb=0$$

On iteration $t=1$

Compute dW , db using current mini-batch

$$VdW = \beta_1 VdW + (1 - \beta_1) dW, SdW = \beta_2 SdW + (1 - \beta_2) dW^2 \quad \leftarrow \text{"momentum"} \quad \beta_1$$

$$Sdb = \beta_2 Sdb + (1 - \beta_2) db^2, Vdb = \beta_1 Vdb + (1 - \beta_1) db \quad \leftarrow \text{"RMSProp"} \quad \beta_2$$

$$VdW^{corrected} = VdW / (1 - \beta_1^t), Vdb^{corrected} = Vdb / (1 - \beta_1^t)$$

$$SdW^{corrected} = SdW / (1 - \beta_2^t), Sdb^{corrected} = Sdb / (1 - \beta_2^t)$$

$$w = w - \alpha \frac{VdW^{corrected}}{\sqrt{SdW^{corrected}} + \epsilon}, b = b - \alpha \frac{Vdb^{corrected}}{\sqrt{Sdb^{corrected}} + \epsilon}$$

Hyperparameters choice

α : needs to be tuned

$\beta_1 : 0.9$

$\beta_2 : 0.999$

$\epsilon : 10^{-8}$

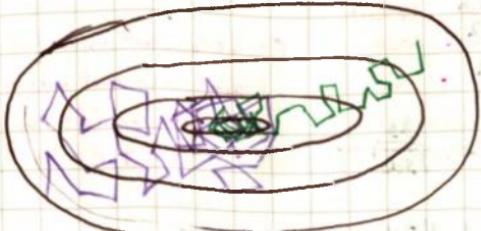
Adam Paper
Author's recommended
Values

(dW)

$(dW^2) \quad db$

[Learning Rate Decay]

One of the things that may help to speed up our learning algorithm is to slowly reduce the learning rate over time, this is called learning rate decay. If we reduce slowly the learning rate α , during the initial phases, the α will be yet big thus allowing ~~longer~~ faster learning, while decreasing, the GD will approach to the minima, and with smaller oscillations, or not, that happens with a constant α . This will allow to enable the algorithm to approach very confidently to the minima, with micro oscillations that can be negligible.



- With constant learning rate α
- With decaying learning rate α

Remembering: 1 epoch = 1 pass on all the training set data

$$\alpha = \frac{1}{1 + \text{decay rate} \cdot \text{epoch}} \text{ with } \alpha_0 = 0.2 \text{ and decay rate } = 1$$

epoch #	α
1	0.7
2	0.67
3	0.5
4	0.4



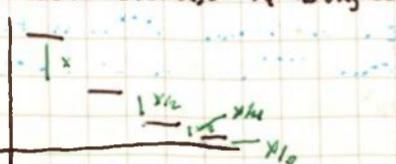
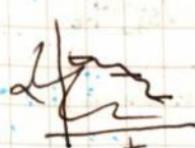
Other learning rate decay methods

$$\text{Exponential decay } (\alpha \propto e^{-t}) : \alpha^{epoch_num} \cdot \alpha_0$$

$$\text{or } \alpha = \frac{k}{\sqrt{\text{epoch_num}}} \cdot \alpha_0$$

$$\text{or } \alpha = \frac{k}{\sqrt{t}} \cdot \alpha_0$$

or discrete staircase by halving each iteration



[Hyperparameter Tuning]

On DNN, as we know, we have to tune up ~~some~~ of the ~~most~~ hyperparameters, which are:

Learning rate α

Momentum	$\beta_1 \approx 0.9$
Adam's	β_1, β_2 and ϵ ≈ 0.999 $\times 10^{-8}$
# layers	
# hidden units	
Learning rate decay	
mini-batch size	

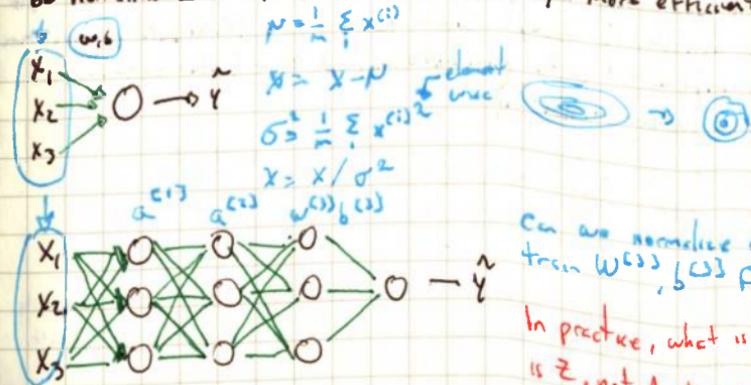
In order of recommendation

- 1 Learning rate α
- 2 Momentum β
- 2 Minibatch size
- 2 # hidden units
- 3 Number of layers
- 3 Learning rate decay
- 4 Adam parameters (defaults are ok)

- try random parameters

[Batch Normalization] - batch norm

- Makes the hyperparameter search much easier, much more robust
- In the view of deep learning, one of the most important ideas has been an algorithm called batch normalization, created by two researchers, Sergey Ioffe and Christian Szegedy. Batch normalization makes our hyperparameter search problem much ~~harder~~ easier, makes our NN much more robust.
- The choice of hyperparameters is a much bigger range of hyperparameters that work well and will also enable us to much easily train even very deep neural networks.
- Similar as it was done on input parameters, when we normalized them in order to turn the elongated contours of our learning problem into more 'round', which also let the gradient descent algorithm to be optimized
- The same principle applies on each hidden layer, where the activations can be normalized as well, to make the next layer more efficiently.



Can we normalize $a^{(2)}$ so as to train $W^{(2)}, b^{(2)}$ faster ($L+1$)

In practice, what is commonly normalized
is Z , not A . Are on the other hand Z is

[Implementing Batch Norm]

Given some intermediate values in the NN $\tilde{z}^{(2)}$... $\tilde{z}^{(n)}$ \downarrow ^{lower sample}

$$\mu = \frac{1}{n} \sum z^{(i)} \text{ - mean } (\mu)$$

$$\sigma^2 = \frac{1}{n} \sum (z^{(i)} - \mu)^2 \text{ - variance}$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad \text{to avoid NaN if } \sigma^2 = 0$$

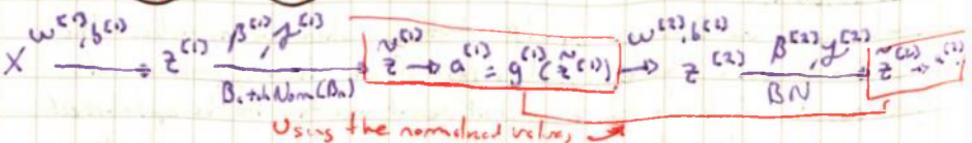
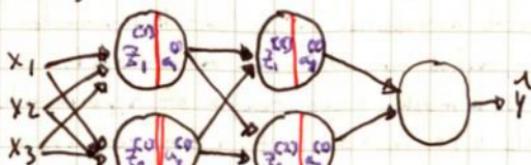
So now we have every component of \tilde{z} with mean μ and variance σ^2 , but we don't want the hidden units to always have mean μ and variance σ^2 . Maybe it makes sense for hidden units to have a different distribution, so for this we'll compute \tilde{z} like:

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta \quad \begin{matrix} \text{learnable parameters} \\ \uparrow \quad \uparrow \end{matrix}$$

the effect of γ and β is that allows us to set the mean of \tilde{z} to be whatever we want it to be.

- So, if $\gamma = \sqrt{\sigma^2 + \epsilon}$ and $\beta = \mu$, they would equal the denominator and numerator, then the effect of $\gamma z_{\text{norm}}^{(i)} + \beta$ will be inverting the equation, then $\tilde{z}^{(i)} = z^{(i)}$, it's computing essentially the identity function.
- To fit it into the neural network, instead of $\tilde{z}^{(i)}$, we use $\tilde{z}^{(i)}$ for the later computations of the neural network.

[Adding Batch Norm to a network]



- In practice we don't need usually to implement Batch norm, because the ML packages already have this available, for example with TensorFlow: tf.nn.batch_normalization
- In practice it is implemented with mini batches, and the batch normalization looks only on the current minibatch.

$$X \xrightarrow{\substack{w^{(1)}, b^{(1)}}} \tilde{z}^{(1)} \xrightarrow{\substack{\gamma^{(1)}, \beta^{(1)}}} z^{(1)} \xrightarrow{\substack{w^{(2)}, b^{(2)}}} \tilde{z}^{(2)} \xrightarrow{\substack{\gamma^{(2)}, \beta^{(2)}}} z^{(2)} \xrightarrow{\substack{w^{(3)}, b^{(3)}}} \tilde{z}^{(3)} \xrightarrow{\substack{\gamma^{(3)}, \beta^{(3)}}} z^{(3)} \dots$$

The term b gets cancelled as we apply normalization and all z 's have similar distributed values, so the bias element becomes a constant and we know it won't change the outcome.

[Implement Gradient Descent with Batch Norm]

For $t \geq 1$... num minibatches

compute forward prop on $X^{(t)}$

in each hidden layer use batch norm to replace $\tilde{z}^{(l)}$ with $\hat{z}^{(l)}$

use back prop to compute $dW^{(l)}$, ~~$d\beta$~~ , $d\beta^{(l)}$, $d\gamma^{(l)}$

update parameters:

$$\left. \begin{array}{l} W^{(l)} := W^{(l)} - \alpha dW^{(l)} \\ \beta^{(l)} := \beta^{(l)} - \alpha d\beta^{(l)} \\ \gamma^{(l)} := \gamma^{(l)} - \alpha d\gamma^{(l)} \end{array} \right\} \text{G.D.}$$

Works with momentum, RMSProp, Adam

[Batch Norm as regularization]

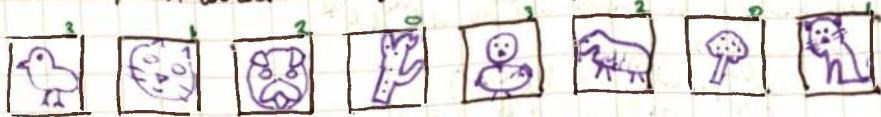
- Each minibatch is scaled by the mean / variance computed on just that mini-batch.
- Adds some noise to the values of $\tilde{z}^{(l)}$ within that minibatch. So similar to dropout, it adds some noise to ~~that~~ each hidden layer's activations.
- This has a slight regularization effect.

[Batch Norm at Test Time]

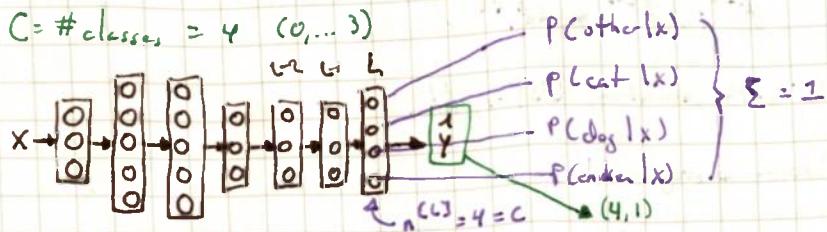
On the test samples we won't have ~~the~~ access to the training minibatches to compute the original mean and variance on each minibatch. To circumvent it, we can compute the exponential weighted average, to keep track of the latest average value, which in turn will be used as μ and σ^2 , but with running average at test.

[Multiclass Classification]

In difference of logistic regression, that was used to do binary classification, there is a generalization of logistic regression called **Softmax** regression, that lets make predictions when we want to recognize 1 to C classes, for example if we want to recognize on a set cats¹, dogs², chickens³ and others⁴.



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



The standard way model for getting your network to do this, uses what's called a Softmax layer in order to generate these outputs.

- In the final layer of the neural networks, we are going to compute as usual the linear part of the layer: $z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]}$. Now having computed z , we need to apply what's called the Softmax activation function.
- The activation function is a bit unusual for the softmax layers, but this is what it does:

Activation function: compute temp "t" s

$$t = e^{z^{[L]}} \quad (z = (4,1)) \therefore t = (4,1) \text{ [element wise]}$$

$$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_i e^{z^{[L]}}} \quad a_i^{[L]} = \frac{t_i}{\sum_i t_i}$$

Applied for example:

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ 4 \\ 3 \end{bmatrix}$$

$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^4 \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix} \quad \sum_j t_j = 176.3$$

$$a^{[L]} = \frac{t}{176.3}$$

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

0	$\frac{e^5}{176.3} = 0.942 \rightarrow 80\%$
0	$\frac{e^2}{176.3} = 0.042 \rightarrow 4\%$
0	$\frac{e^4}{176.3} = 0.002 \rightarrow 0.1\%$
0	$\frac{e^3}{176.3} = 0.114 \rightarrow 6.4\%$

[Deep Learning Frameworks]

At this moment we've learnt how to implement almost from scratch deep learning algorithms. Now, it turns out that is not practical to implement all of them by scratch ourselves but it's a better approach to use frameworks that are out there, already optimised and growing month by month. A list of the leading frameworks available

- Caffel Caffel2
 - CNTK
 - DL4J
 - Keras
 - Lasagne
 - mxnet
 - Paddle Paddle
 - Tensor Flow
 - Theano
 - Torch
 - Accord (.Net)

Criteria to choose:

- Ease of programming: Deck, testben, deployment
 - Running speed :
 - Truly Open : Not only open source, but also good governance

[Code example for Tensor Flow]

```
import numpy as np  
import tensorflow as TF
```

```
coefficients = np.array([[1], [-20], [25]])
```

w = tf.Variable((0), dtype=tf.float32)

`x = tf.placeholder(tf.float32, [3, 1])`

$\text{cost} = x[0] \cdot 0 \cdot w * * 2 + x[1,0] \cdot w + x[2] \cdot 0 \# (w-5) * * 2$

```
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
```

```
init = tf.global_variables_initializer()
```

Section = Section.

```
session.run(init)  
init(session, graph)
```

} with `t.F.Session()` as `session`:
`session.ran (cont)`
`print (session.ran[0])`

for i in range(1000):

```
session.run(train, feed_dict={x:coefficients})
```

point (session . run (w))

$$X[0](6) \rightarrow [X[0](6) + w^2] \rightarrow \dots \rightarrow \text{cost} \rightarrow$$

$w = w^2$

$X[2](6)$

Course 3 : ML Strategy

[Orthogonalization]

Orthogonal means 90° , like on a hyperplane \Leftrightarrow

Change of assumptions on ML

- Fit training set well on cost function (\propto human-level performance) - Adam
- Fit dev set well on cost function - Regularization
- Fit test set well on cost function - Bigger train set
- Performs well in real world - changes the dev set or cost function

Orthogonalization in ML: it's about tweaking the right hyperparameters to obtain the right results. Mixing or altering one HP to reduce or maximize a particular effect should be avoided. Each HP is designated to both fit/fail a specific situation and problem, therefore is not good to mix them to get a better result. (i.e., tuning the regularization parameter and increasing the regularization optimization algorithm depends on a single "knob".)

[Single Numeric Evaluation Metrics]

Having a single measure unit to evaluate the performance of the model \Rightarrow to have of course some iterations and evaluate them, for example using Precision and Recall:

Precision: Of the correctly identified samples, what percentage of them are From the whole set, (i.e., if a set has 100 pictures, 95 are cats and 5 are something else). If the 95 cats are properly recognized, then we have a precision of 98% because the 95% of the set was properly recognized.

Recall: From the sample above, if only 90% of the cats were recognized, (that's 90% of the 95 cats in the set), then the percentage of actual cats recognized would be 94.7%.

A way to combine them into one single number to measure the performance is to use the F1 score, or "harmonic mean".

F1 Score will allow us to choose the best tradeoff between Prec. and Rec.

Classifier	Precision	Recall	F1 Score	
A	95%	90%	92.4%	$\frac{2}{\frac{1}{P} + \frac{1}{R}}$
B	98%	85%	91.0%	

Course 4: Convolutional Neural Networks

Computer Vision Problems

Image classification: Cat or not?

Object detection: Figure out the position

Neural Style transfer: With content image and style image, merge them.

Challenges: Inputs can get really big! in the previous problems we have worked on scales of $64 \times 64 \times 3$ images (12288 pixels), but that's very small. Working with larger images (i.e., 1000×1000), we end up with 3 million pixels in to process ($1000 \times 1000 \times 3$)

Passing the 3M parameters to a 1000 units hidden layer, would imply to have a $(1000 \times 3M)$ matrix which means the matrix will have 3 billion params.

This causes overfitting and computational requirements become very high.

$$\begin{matrix} 3M & 1000 & [1000, 3M] \\ x_1 & 0 & \xrightarrow{\text{...}} 3B, 1, 1000 \\ x_2 & 0 & 0 \\ x_3 & 0 & 0 \\ x_n & 0 & 0 \end{matrix}$$

EDGE Detection

Vertical edge detection

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

convolution

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

3x3

filter or Kernel

Python: conv FORWARD

TF : `tf.nn.conv2d`

Keras : `Conv2D`

-5	-4	0	0
-1	-2	2	3
0	-2	-4	-7
-3	-7	-3	-16

4x4

Vertical
edge
detection

6x6

$$3x1 + 1x1 + 2x1 + \\ 0x0 + 5x0 + 7x0 + \\ 1x-1 + 8x1 + 2x-1 = -5$$

$$0x1 + 3x1 + 7x1 + \\ 0x0 + 8x0 + 2x0 + \\ 2x-1 + 9x-1 + 3x-2 = -4$$

$$1x1 + 8x1 + 2x1 + \\ 2x0 + 9x0 + 5x0 + \\ 7x1 + 3x-1 + 1x-2 = 0$$

Simplified example:

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

1 0 -1

=

1 0 +1

0 +1 0

0 0 +1

0 0 0 +1

Dark to light

transition

0 30 30 0

0 30 30 0

0 30 30 0

0 30 30 0



bright
edge
dark

Vertical and horizontal edge detection

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

Vertical

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

Horizontal

$$\begin{array}{ccccccc} 10 & 10 & 10 & 0 & 0 & 0 & \\ 10 & 10 & 10 & 0 & 0 & 0 & \\ 10 & 10 & 10 & 0 & 0 & 0 & \\ 0 & 0 & 6 & 10 & 10 & 10 & \\ 0 & 0 & 0 & 10 & 10 & 10 & \\ 0 & 0 & 0 & 10 & 10 & 10 & \\ \end{array}$$

0 0 0 0

30 20 -10 -30

30 20 -10 -30

0 0 0 0

Learning to detect edges

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

Sobel filter

$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

$$\begin{bmatrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{bmatrix}$$

Sobel filter

There has been large debate among Computer Vision researchers about the values on the filters to use. But now, instead of

had picking the values, with the rise of Deep Learning, one of the things we've learned is that when we really want to detect edges on some complicated edge, we might not need a researcher to hand pick the nine values of the filter. Maybe we just can learn them and treat the nine numbers of this matrix as parameters which we can then learn using backpropagation.

$$\begin{bmatrix} 3 & 0 & 1 & 2 & 7 & 4 \\ 1 & 5 & 8 & 9 & 3 & 1 \\ 2 & 7 & 2 & \boxed{5} & 1 & 3 \\ 0 & 1 & 3 & 1 & 7 & 9 \\ 4 & 2 & 1 & 6 & 2 & 8 \\ 2 & 4 & 5 & 2 & 3 & 9 \end{bmatrix}$$

[Padding]

As we've seen, doing convolution does shrink the data, after each hidden layer so it happens ~~more~~ on a ratio $n-f+1$; thus can become problematic if because after several layers the data would be shrunk vastly. Also there is the problem that taking convolutions on the edge pixels, there are not many pixels around, so it can't have greater influence on the convolution, for example the pixel in red has only pixels down and to the right, whereas the green pixel has neighbors pixels. So to not throw away info from the image's edges.

To fix the problems, we can add padding to the image, for example adding one pixel around. Like in blue color. So after $\hat{x} = (n+p-f+1, n+p-f+1)$.

convolution

$$\begin{array}{c} w_1 \ w_2 \ w_3 \\ * \quad w_4 \ w_5 \ w_6 \\ w_7 \ w_8 \ w_9 \end{array}$$

$(F \times F)$

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

$(n-F+1 \times n-F+1)$

Same size as original

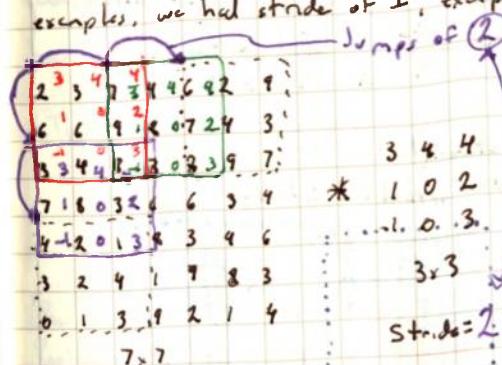
Valid and Same convolutions:

"Valid": no-padding, $n \times n * f \times f \rightarrow (n-f+1) \times (n-f+1)$
 $6 \times 6 * 3 \times 3 \rightarrow (4 \times 4)$

"Same": Pad, so that the output size is the same as the input size
 $(n+2p-f+1) \times (n+2p-f+1) * (f \times f) = n \times n$
 $p = \frac{f-1}{2}$ by convention n. (V, p should be odd)

[Strided Convolutions]

It's another piece of the basic building blocks of convolutions as used in Convolutional Neural Networks. What basically means is that we define as stride, the number of steps to "jump"; on the previous convolution examples, we had stride of 1, except



$$\frac{n+2p-f+1}{2} \times \frac{n+2p-f+1}{2}$$
$$= \begin{matrix} 91 & 100 & 83 \\ 69 & 91 & 127 \\ 44 & 72 & 74 \end{matrix}$$
$$\frac{7+0-3}{2} \approx \frac{4}{2} + 1 = 3$$

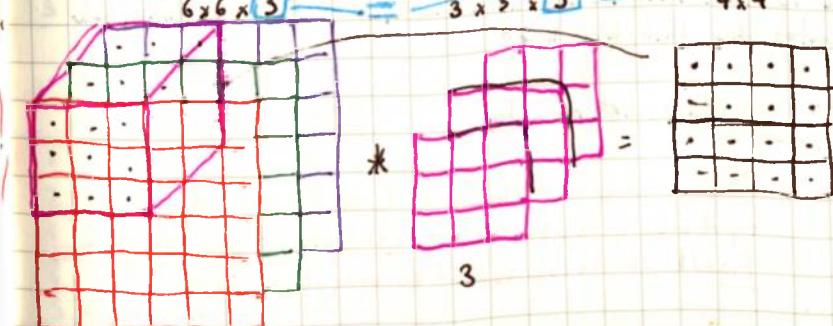
There is a convention to [floor] the dimensions, so they are not decimal.

[Summary]

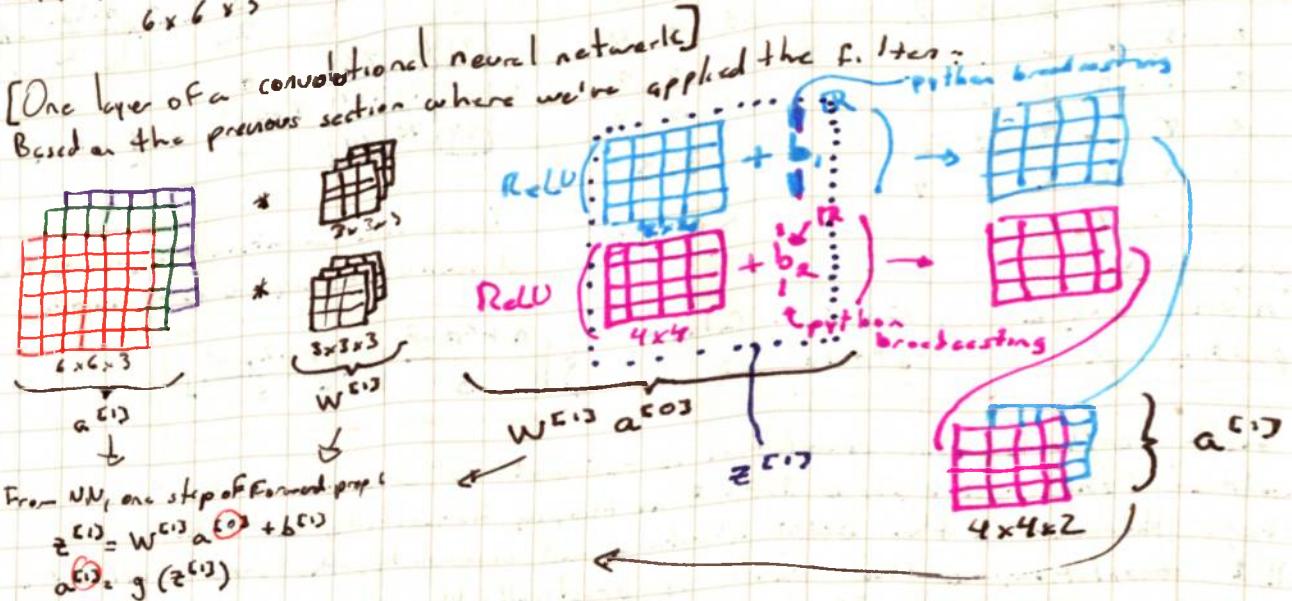
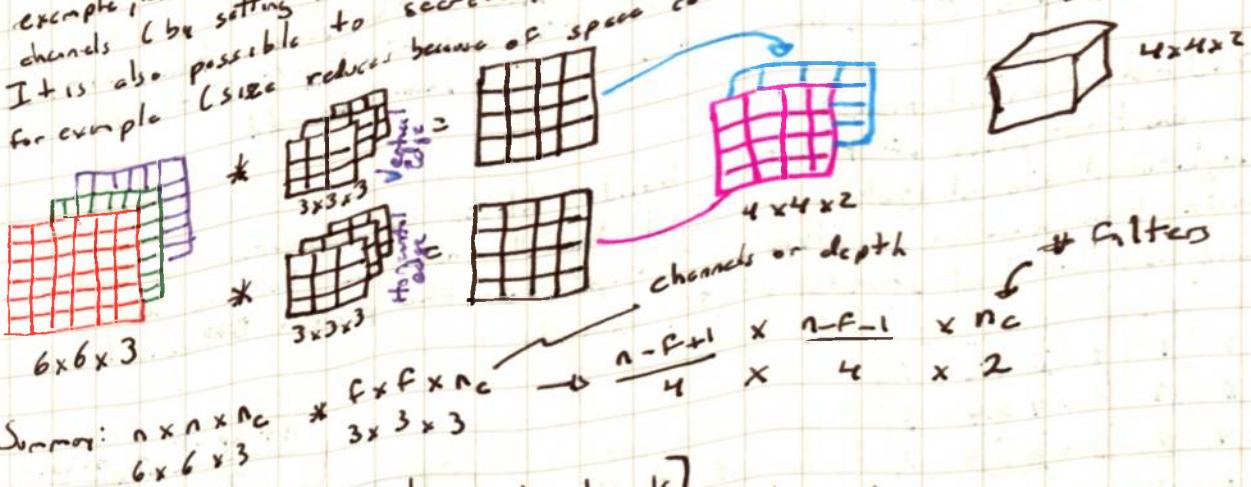
$n \times n$: image	$f \times f$: filter	$\left\lfloor \frac{n+2p-f}{s} \right\rfloor \times \left\lfloor \frac{n+2p-f}{s} \right\rfloor$
padding: p	stride: s	

In most textbooks, this convolution is referred as "cross-correlation", and convolution (the real one), is defined as flipping horizontally and vertically the filter. However in the DNN literature and community, the "cross-correlation" is called convolution, as the result is practically the same and it saves coding time.

(Convolutions over volume): [Convolutions on RGB Image] some number of channels



[Multiple Filters]
 As we learnt before, it's possible to use the filters horizontally, diagonal or in a specific "degree". The filters can apply to one color channel (on a 3-channel image), or to all channels (by setting the appropriate channel's values on the filter). It is also possible to search for various edges (if we stack the filters for example (size reduces because of spec constraints)).



- Numbers of parameters in one layer
- If we have 10 Filters that are $3 \times 3 \times 3$ in one layer of a neural network, how many parameters does that layer have?

$3 \times 3 \times 3 + 1 \times b^{(1)}$

$28 \times 10 = 280 \text{ parameters.}$

One nice thing about this, is that doesn't matter how big the input image is, $5K \times 5K, 1K \times 1K$; the number of parameters remain fixed! It can be applied to very large images.

[Summary of notation]

$f^{[L]}$ = filter size
 $p^{[L]}$ = padding
 $s^{[L]}$ = stride
 n_C = number of filters

Input : $n_h^{[L-1]} \times n_w^{[L-1]} \times n_C^{[L-1]}$
 Output : $n_h^{[L]} \times n_w^{[L]} \times n_C^{[L]}$
 Volume : $n_h^{[L]} = \left\lfloor \frac{n_h^{[L-1]} + 2p^{[L-1]} - f^{[L-1]}}{s^{[L-1]}} + 1 \right\rfloor$

Each filter is: $f^{[L]} \times f^{[L]} \times n_C$

$$n_w^{[L]} = \left\lfloor \frac{n_w^{[L-1]} + 2p^{[L-1]} + f^{[L-1]}}{s^{[L-1]}} + 1 \right\rfloor$$

Activations: $a^{[L]} \rightarrow n_h^{[L]} \times n_w^{[L]} \times n_C^{[L]}$

→ vectorized

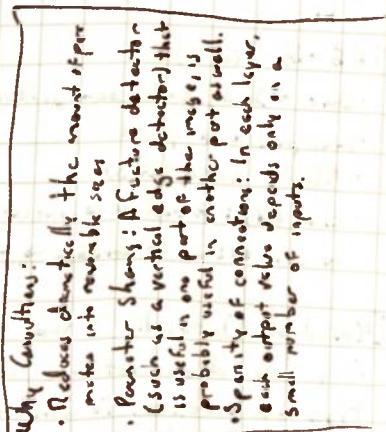
$$A^{[L]} \rightarrow m \times n_h^{[L]} \times n_w^{[L]} \times n_C^{[L]}$$

Weights: $f^{[L]} \times f^{[L]} \times n_C^{[L]} \rightarrow$ Filters in layer L

$$\text{bias} : n_C = (1, 1, 1, n_C^{[L]})$$

[Types of layers in a convolutional network:]

- Convolution [Conv]
- Pooling [pool]
- Fully Connected [Fc]



[Pooling Layers]

Pooling is used to shrink the output from very large channels; the principle is very similar to Conv, where a filter is passed over the pixels, however instead of multiplying and adding, only the "max" is taken. As Conv, it also has a filter size, and a stride; the padding is practically never used (but is possible to add).

An example

Max Pooling

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

4x4



9	2
6	3

2x2

$F=2$

$S=2$

There is another version called Average pooling, which does what we can deduct, instead of taking the max.

A common setting is:

$$F=2, S=2 \text{ or } F=3, S=2$$

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

→

9	9	5
9	4	5
8	6	9

$F=3$

$S=1$

$2 \times 3 \times n_C$

[Classic networks]

- LeNet-5 : trained on grayscale images ($32 \times 32 \times 1$)
 - 60k parameters.
 - backpropagation
 - conv -> pool -> conv -> pool FC, FC, output

- sigmoid + tanh
- LeCun et al., 1998
- read chapter II and III

AlexNet: Trained in color, RGB ($227 \times 227 \times 3$)

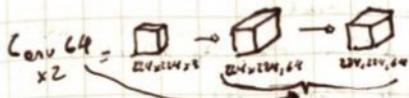
• ReLU - layout

- Much bigger than LeNet, b.s. more, 60 million params
- conv - pool - conv - pool - conv - conv - conv - pool - FC - FC - FC - softmax
- Multiple GPUs
- Local response normalization (LRN)

• Krizhevsky et al., 2012

VGG-16 (or 19)

- Conv = 3×3 , s=1, same padding
- max pool = 2×2 , s=1



- RGB ($224 \times 224 \times 3$)

• number of channels doubles in each stack of layers

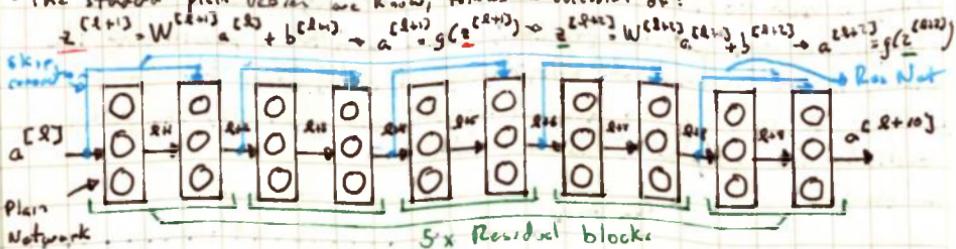
- Conv $64 \times 2 \rightarrow$ pool $64 \rightarrow$ conv $128 \times 2 \rightarrow$ pool $128 \rightarrow$ conv $256 \times 2 \rightarrow$ pool $256 \rightarrow$ conv $512 \times 3 \rightarrow$ pool 512
- \rightarrow conv $512 \times 3 \rightarrow$ pool $512 \rightarrow$ FC \rightarrow FC \rightarrow Softmax

• Simonyan & Zisserman 2015, Very Deep Conv.

[Residual Networks - Res Nots]

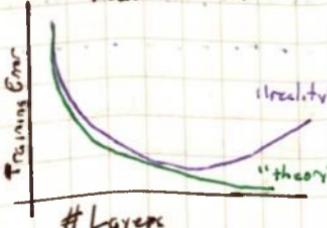
Very Very deep networks are difficult to train because of vanishing and exploding types of problems. To avoid that, it's possible to take the activation from one layer and suddenly fed it into another layer even much deeper in the NN.

- The standard "plain" neurons are known, follows a succession of:

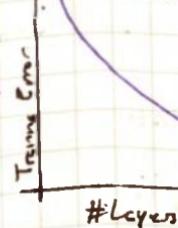


Behaviour in training:

Plain



ResNet



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

$$= g(w^{[L+1]} a^{[L+1]} + b^{[L+1]} + \dots + z^{[L+1]})$$

the dimensions should match

- with "zero" padding or
- reducing the dimension in the activation by adding other matrix weights ("multiplication") W_s (experiment)

$W_s = R$

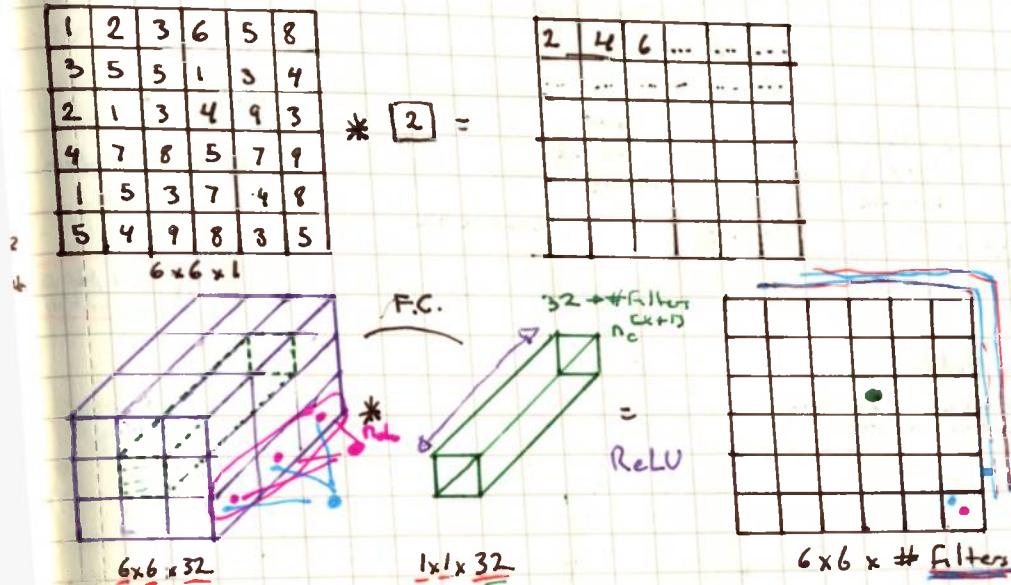
↳ learned parameters or

↳ zero initial

He et al., 2015. Deep Residual networks for image recognition

- Lin et al., 2013. Networks in Network

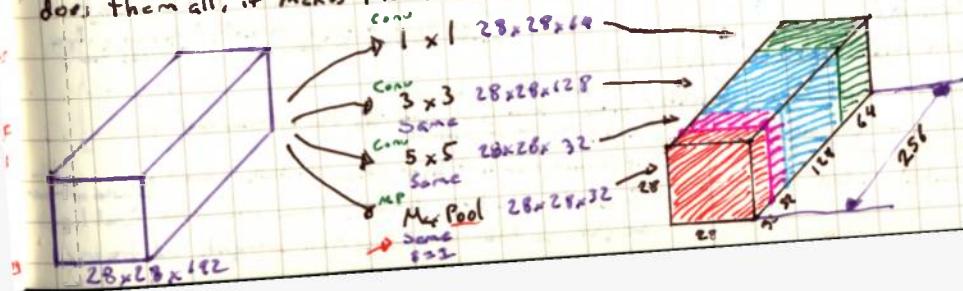
Filters as we have seen are larger, i.e., 2×2 , 3×3 , and thus one is just a real number $(1, 1)$; on first glance it doesn't seem useful as what it does is just multiply the original value times the filter, however it makes sense while working with several channels; what it will actually do, is to look at each of the different positions and take an element wise product with the number of channels of the original volume and the volume of the filter.



This idea has been very influential on several NN architectures, including the Inception Network.

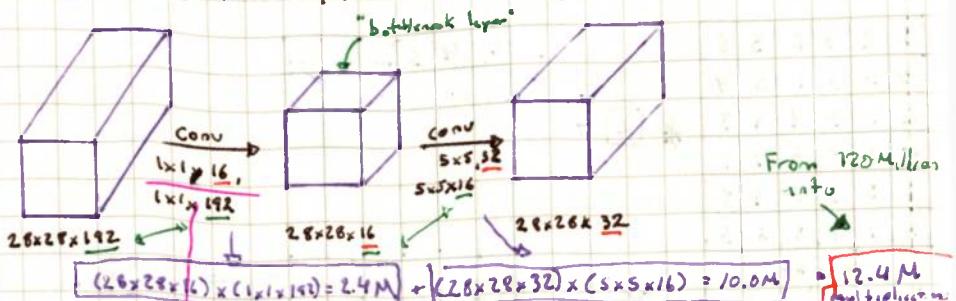
If we want to reduce the height and width, we can use a pooling layer, but if the number of channels get big, then we can use a lot of 1x1 filters to shrink the number of channels, if needed, otherwise can be the n_c kept equal.

- Inception Network Motivation: Szegedy et al., 2014: Going Deeper with convolution
- When designing a layer for a ConvNet, we might have to pick if we want a 1×1 or 3×3 or 5×5 filter, or if we want a pooling layer... ; What the inception network does, is that does them all; it makes the network architecture more complicated but works well!



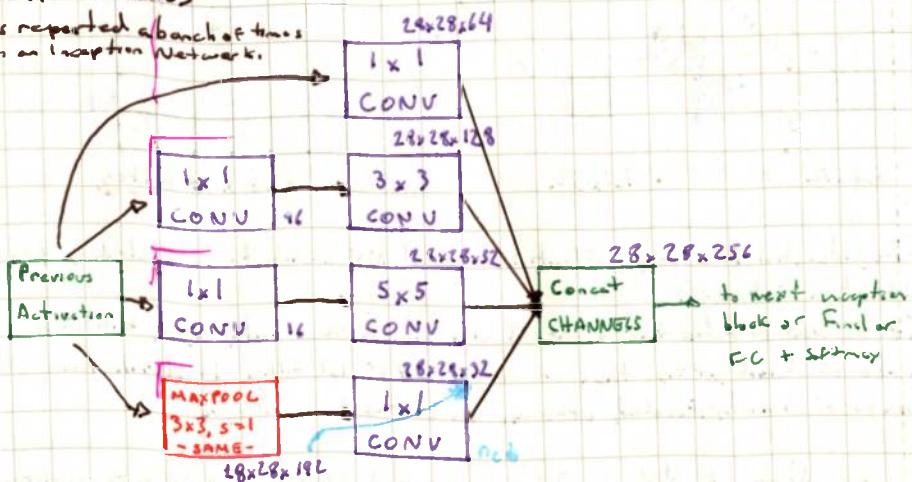
- The problem of the computational cost.
- Taking from the previous example the convolution $5 \times 5, 32$, where the original parameter space of $28 \times 28 \times 192$ we have
 $\text{Input} = 28 \times 28 \times 192 \rightarrow \text{Conv } 5 \times 5, 32 \text{ (same)} \rightarrow 28 \times 28 \times 32$; this mean we have:
 32 filters $5 \times 5 \times 192$; which gives up the multiplication of $(28 \times 28 \times 32) \times (5 \times 5 \times 192)$, which in turn equals to $120,422,400 = \underline{\text{120 Million}}$ of output nodes multiplications.

We can use a 1×1 convolution to reduce the count of computation needed; We input the volume, use a $1 \times 1 \times 2$ convolution, to reduce the volume, and then run the 5×5 convolution.



[Inception Module]

is reported a bunch of times on an Inception Network.

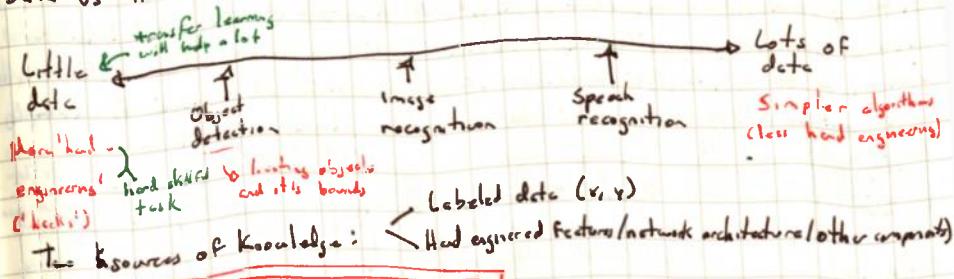


On the final layers, we added softmax layers that help to regularize the outputs, helps to ensure the features computed in hidden/intermediate layers are not too bad to predict the class on the output layer. An inception network is called GoogLeNet.

[Practical Advice]

- Use open source implementations; contribute!
- Transfer learning; use pre-trained weights
 - freeze pretrained layers
 - continue training by removing the softmax layer
 - set trainable parameters
 - freeze
 - prune the last layer adde to avoid processing the whole network
- to seriously encouraged unless we have a huge dataset
- PCA color augmentation (RGB)
- if we have a huge dataset, then use the first layer weight as the initialization weights and train.

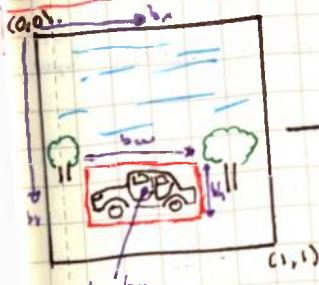
Data vs Hard Engineering



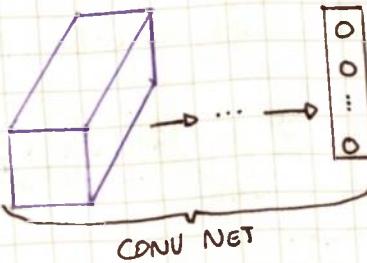
-Week 3 - Object Detection -

It's one of the areas of computer vision that is exploding and working much better than just a couple years ago. We already know and do object classification (telling if a picture is a cat, dog, etc), but now in the detection, it's needed to locate one or many objects and their location in bounding boxes, where is the picture are.

-Classification with localization



(x_i, y_i)



CONV NET

$| \text{bbox} | \text{or} | \text{car} | \geq 1$

Softmax (4)
 -pedestrian
 -car
 -motorcycle
 -background

b_x } midpoint
 b_y } Y-position
 b_h } height
 b_w } width

The idea is to output from the convnet a softmax predicting the object, but also additional parameters identifying the object location (b_x, b_y = midpoint of the object, from origin) and its height and width (b_h, b_w), so we can calculate its bounding box.

Need to output: $b_x, b_y, b_h, b_w, \text{class label}$ (e.g.)

-Defining the target label. (y_i)

Class : 1 = pedestrian 2 = car 3 = motorcycle 4 = other

P_c	Is there an object?		x_i		y_i
b_x					
b_y					
b_h					
b_w					
C_1	pedestrian				
C_2	car				
C_3	motorcycle				

Loss Function:

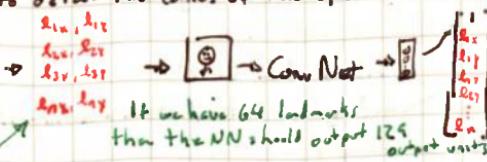
$$L(\hat{y}_i, y_i) = (y_{i1} - \hat{y}_{i1})^2 + (y_{i2} - \hat{y}_{i2})^2 + \dots + (y_{iC} - \hat{y}_{iC})^2$$

, if $y_i = 1$

$$(y_{i1} - \hat{y}_{i1})^2 \quad (\text{if } y_i = 0)$$

[Landmark Detection]

On the previous section we saw how we can get a NN to output 4 numbers ((x_1, y_1, x_2, y_2)) to specify the bounding box of an object we want the NN to localize. For landmarks, we want a NN to output (x_i, y_i) + numbers of important points of an image. For example, if we want to detect the corners of the eyes on a face:



{ 1
1
1
1
1
1
1
1
1
1
1
1 } = Face?
 128 } Used for APR
CGI, detection.

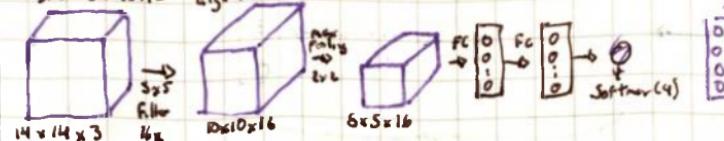
- Landmarks should most be consistent across different images -

[Sliding Window Object Detection]

Although it's computationally inefficient, it is possible to detect objects on a picture by sliding small to big "squares" which are one by one fed into a ConvNet to identify if there is an object or not. For this the ConvNet should be trained with close-ups of the object we want to detect . This sliding windows are not recommended, too slow.

[Convolutional Implementation of Sliding Windows]: Turning FC layers into convolutional layers

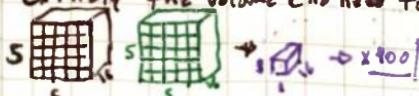
Original detection algorithm:



Turned into convolutional layers



Each 5×5 filter has 16 channels, and applying it into the 5×5 volume, it will cover entirely the volume (no need to stride), so the output is going to be $1 \times 1 \times n_c$

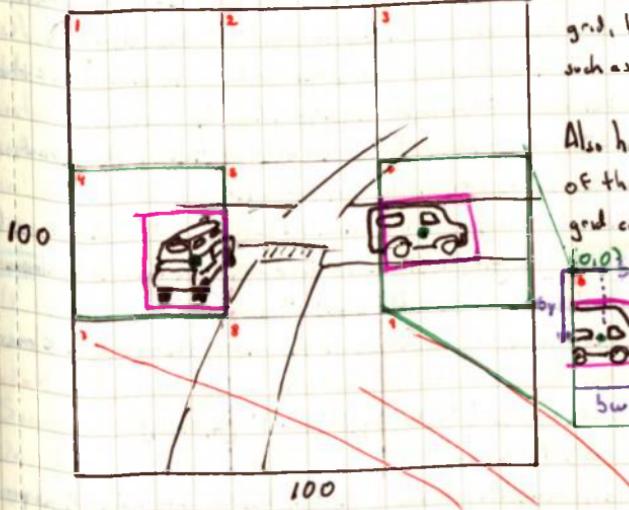


If we input larger images, the conv net will ~~take~~ scan it with stride, the stride defined or ruled by the Network size. So with a MP of 2, the conv net will slide over the image with strides of 2, and at the last filter, the convnet will mark on the output volume, in which "window" the object was found. (See stepled paper).

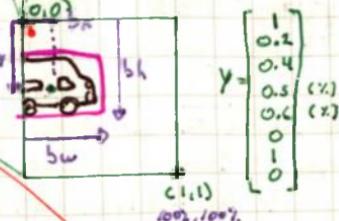
• Sermanot et al., 2014, OverFeat: Integrated recognition, localization using ConvNets

[Bounding Box Predictions] : YOLO algorithm: You Only Look Once, Redmon et al 2015
 with the convolutional implementation of sliding windows, we've found a computationally efficient way to output bounding boxes, however it could not be always accurate; because it might happen that the sliding window's stride skips a crucial element of the object, or it was too small to fit it, or the object was wider.
 - It works by placing a grid over an image, and the basic idea is to take the image classification and localization algorithm, and apply it on each of the grids of the image. For labeling the image, we specify a label Y which is a vector containing as before: $[1_{\text{is present}}, b_x, b_y, b_h, b_w, c_1, c_2, c_3]^T$, this for each grid cell.

For this sample we assume a (9×9) grid, but finer grids are used in reality such as (19×19) .



Also here we look at the midpoint of the object, but Relative to the grid cell where it is located!

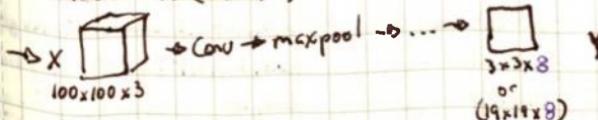


y_1	y_2	y_3	y_4	y_5	y_6	y_7	y_8	y_9	y_{10}	y_{11}
0	0	0	1	0	1	0	0	0	0	0
?	?	?	bx	?	bx	?	?	?	?	?
?	?	?	by	?	by	?	?	?	?	?
?	?	?	bh	?	bh	?	?	?	?	?
?	?	?	bw	?	bw	?	?	?	?	?
?	?	?	0	?	0	?	?	?	?	?
?	?	?	1	?	1	?	?	?	?	?
?	?	?	0	?	0	?	?	?	?	?

Labels for training for each grid cell.

? = numbers that will be ignored (noise)

The target output (based on this sample) will be a $(3 \times 3 \times 8)$ volume. With grids of 19×19 , then it would be $(19 \times 19 \times 8)$



This is one of the harder paper to understand, even for senior researchers

[Intersection over Union: IoU]

How to tell if the object detection algorithm is working well?

Having a match with the ground truth bounding box and if our algorithm outputs this bounding box ... is this a good or bad outcome?



$$\text{IoU} = \frac{\text{Size of } \cap}{\text{Size of } \cup}$$

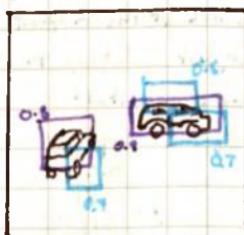
"Correct" if $\text{IoU} \geq 0.5 \dots$ | or 0.6
"perfect" if $\text{IoU} = 1$ | 0.7

↓ is a measure of the overlap between

↓ IoU: two bounding boxes.

[Non-max Suppression]

With finer grids, it's almost always true that the algorithm may find multiple detections of the same objects. Non-max Suppression is a way to make sure the algorithm detects each object only once. This happens while running the identifiers, each cell of the grid might 'think' that it has the center of the object and output a high P_c .



What Non-max suppression does, is that it cleans up those up the detections so that we end up with one detection per object.

It looks at the probabilities associated with each of the detections P_c . It takes the largest one (the most confident one), then looks at the remaining rectangles with a high overlap (IoU) and suppress them.

Non-max means, it's going to suppress the rest of the rectangles & max

Multiple detections of each object Note: P_c is evaluated per each class, in our example, P_c has 3 object steps:

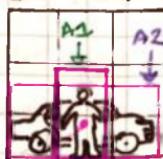
1: Each output prediction is $[P_c, b_x, b_y, b_h, b_w]^T$ $[C_1, C_2, C_3]$

2: Discard all boxes with $P_c < 0.6$

3: While there are any remaining boxes: a) Pick the box with largest P_c as prediction
b) Discard any remaining box with $\text{IoU} \geq 0.5$ with the box output in the previous step

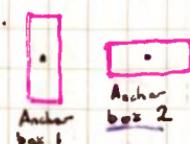
[Anchor Boxes] (Note: choose between 5-10 classes for known to choose representative shapes)

What if a grid cell wants to detect multiple objects? We can use the idea of anchor boxes. Below a simplified 3×3 grid for illustration purposes. We can see that



both objects share the same midpoint on the same cell.

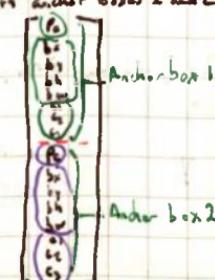
If we want to detect multiple objects, we can use anchor boxes and define the class label containing as well with anchor boxes:



Normal: with anchor boxes 1 and 2 Anchor Box algorithm:

P_c
b _{x1}
b _{y1}
b _{h1}
b _{w1}
c ₁
c ₂
c ₃

Cont defect
more than
one object
at the same
time.



Previously

Each object in training image is assigned to grid cell that contains that object's midpoint.

Output y :

$8 \times 3 \times 3 \times 8$

With two anchor boxes

Each object in training image is assigned to grid cell that contains the object's midpoint and anchor box for the grid cell with highest IoU. Grid cell, anchor box, and anchor

Output

$3 \times 3 \times (8 \times 2)$

[Course 4 - Week 4 : Face Recognition and Neural Style Transfer]

[Face Recognition]

- 1) Face Verification: - Input image θ , name / ID 1:1
 Building Block For • Output: whether the input image is that of the claimed person
- 2) Recognition :- Has a database of K persons 1:K
 - Get an input image
 - Output ID if the image is any of the K persons (or not recognized)

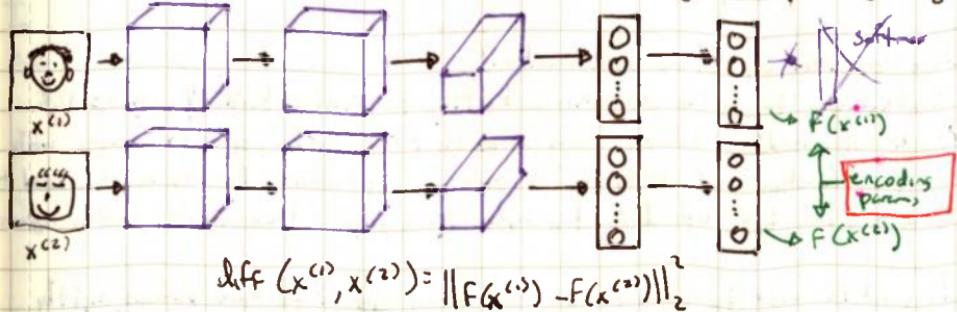
[One Shot Learning]

What it means, is that for most Face recognition applications, we need to be able to recognize a person given just one example of that person's face. And historically, Deep learning algorithms didn't work well if we have only one training example. To address this, we can't rely on a DNN with softmax output layer, as this would require to retrain and add the additional unit. The approach in this case is to have a function to compute the similarity between two images, that means, compare the image in the database vs the image taken in the moment. If the difference is small, then the person is the same, whereas a large number would indicate the person is not the same.

[Siamese Network]: the similarity function.

We are used to seeing pictures of ConvNets where we input an image and through a sequence of convolutional, pooling and Fully Connected layers, we end up with a feature vector, and sometimes is fed up to a softmax unit to make classification. On the case of Siamese Networks, we focus on the last vector of let's say 128 numbers computed by some FC layer that is deeper in the network.

We can refer to this vector as $F(x^{(c1)})$, or as "encoding of $x^{(c1)}$ ", encoding of image ' $c1$ '



The goal is to find a set of parameters that we can then use to find the difference. If the difference is small between images of the same subject, we have it!

As we vary the parameters, we end up with different encodings, and what we can do is use back propagation, to adjust all the parameters to ensure the the conditions are satisfied.

Taigman et al. 2014. Deep Face closing the gap to human level performance

[Triplet Loss]

One way to learn the parameters of the neural network so that it gives a good encoding for pictures or faces, is to define and apply gradient descent on the triplet loss function.

To apply the triplet loss, is needed to compare pairs of images. For example given a picture **A** as Anchor, compare it with another image of the same object which will be **P** positive, and another from someone else, which will then be **N** negative sample. The goal is to have pairs **(A, P)** and **(A, N)**, we want the encodings to be very similar for AP, whereas for AN, to be larger.

Loss function: Given 3 images



α - margin

$$\cdot L(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2, \phi)$$

Cost function:

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

Of course it's needed to have multiple pictures of the same person, i.e., 10 persons so we can have for example:
Training set of 10k pictures of 1k persons

After being trained with enough data, then it can be possible to use just one shot learning

Choosing the triplets **A, P, N**

- Avoid randomization: During training, if APN are chosen randomly, then is very likely that $d(A, P) + \epsilon \leq d(A, N)$, so the condition will be easily satisfied
- Choose triplets that are 'hard' to train on, choose those which similarity is very close so that $d(A, P) + \epsilon \leq d(A, N)$ is forced to have better competition and push the quantities further away.

[Face Verification and Binary Classification]

Another way to do face verification is to take the final embeddings, maybe 128 dimensional embeddings and use them as an input for logistic regression, where with a sigmoid function we'll verify if the person is valid or not. For this is not needed to have triplets, but just different samples of a person; then apply logistic regression.

To avoid the need to compute on real time, ~~this~~ is to have perhaps nightly precomputed embeddings and run the network with the precomputed values.

$$\hat{Y} = \sigma \left(\sum_{i=1}^n w_i [f(x^{(i)})_K - f(x^{(j)})_K] + b \right)$$

Schroff et al., 2015, FaceNet: A unified embedding for face recognition and clustering

[Neural Style Transfer]

The goal is to generate an image based on the style of one of them and when it's being applied to the content (source) image.

- [Content Cost Function]: -

$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

- Use a hidden layer l to compute content cost
- 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.

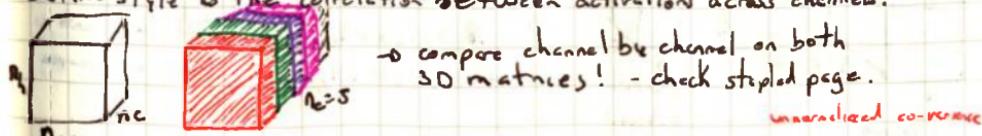
$$\hookrightarrow J_{\text{content}}(C, G) = \frac{1}{2} \| (a^{[L](C)} - a^{[L](G)})^2 \|_F^2 \rightarrow b \quad \begin{cases} \text{elements are} \\ \text{unrolled into vectors} \end{cases}$$

} sum of square differences between content in layer l for C and G

- [Style Cost Function]: -

~~SIMILAR~~ we are using layer l 's activations to measure "style"

Define style as the correlation between activations across channels.



unnormalized covariance

Style Matrix: Let $a^{[L]}_{ijk, k'} = \text{activation at } (i, j, k, k')$. $G^{[L]}$ is $n_w \times n_c$ $G^{[L]}_{kk'}$,
for later: $G^{[L]}_{kk'} = \sum_{i=1}^{n_w} \sum_{j=1}^{n_h} a^{[L]}_{ijk} a^{[L]}_{ijk'}$

$$G^{[L]CS}_{kk'} = \sum_{i=1}^{n_w} \sum_{j=1}^{n_h} a^{[L]}_{ijk} a^{[L]}_{ijk'}$$

In linear algebra Gram matrix

$$J_{\text{style}} = (S, G) = \frac{\| G^{[L]CS} - G^{[L]CS] \|_F^2}{\sqrt{n_w \cdot n_h \cdot n_c}} \cdot \frac{1}{2 \cdot n_w \cdot n_h \cdot n_c}$$

$$+ \beta \cdot \frac{\sum_{K, K'} (G^{[L]CS}_{KK'} - G^{[L]CS}_{KK'})^2}{\sum_{K, K'} G^{[L]CS}_{KK'}} \cdot \frac{1}{2 \cdot n_w \cdot n_h \cdot n_c}$$

To get more visually attractive results, works better to use the style cost func. from multiple different layers.

$$J_{\text{style}}(S, G) = \sum_L J_{\text{style}}^{[L]}(S, G)$$

- Overall cost function:

$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

1D and 3D Generalizations

There are sources of 1 dimensional data, such as electrocardiograms, which also can be seen as a time-series problem. However it is possible to compute convolutions as well - For example:

$$\begin{array}{c} \text{Diagram of a 2D input image} \\ \text{2D } 14 \times 14 \end{array} * \begin{array}{c} \text{Diagram of a 2D filter} \\ 2D \\ 5 \times 5 \\ \text{Filter} \end{array} \rightarrow \begin{array}{l} 14 \times 14 \times 3c * 5 \times 5 \times 3c \\ \rightarrow 10 \times 10 \times 16 \end{array} \left. \right\} \text{Convolution step.}$$

$$\begin{array}{c} \text{Diagram of a 1D signal} \\ \text{1D } 14 \end{array} * \begin{array}{c} \text{Diagram of a 1D filter} \\ 1D \\ 5 \end{array} \rightarrow \begin{array}{l} 14 * 5 \\ \rightarrow 10 \\ 5 \times 16 \end{array} \left. \right\} \text{1D data convolved with 1D filter and also with 16 channels.}$$

$$\begin{array}{c} \text{Diagram of a 3D volume} \\ 3D \\ 14 \times 14 \times 14 \end{array} * \begin{array}{c} \text{Diagram of a 3D filter} \\ 3D \\ 5 \times 5 \times 5 \\ \text{Filter} \\ 3 \text{ channel} \\ 16 \text{ channel} \end{array} \rightarrow \begin{array}{l} 14 \times 14 \times 14 \times 16 \\ * 5 \times 5 \times 5 \times 16 \\ \rightarrow 10 \times 10 \times 10 \times 16 \end{array} \left. \right\} \text{3D data such as frames or a movie}$$

- COURSE 5 - Sequence Models -

Sequence Models are one of the most exciting areas in Deep Learning; Models like recurrent Neural Networks or RNNs have transformed speech recognition, Natural Language Processing and other areas. Sequence Models can be useful in:

	input	output	
Speech Recognition	WAV file	Text	- sequence
Music Generation	# or an integer	Music	- sequence
Sentiment classification	"this is not that like"	Star rating	
DNA seq. analysis	A G C C C T G T G A G G A C T A G	A G C C C T G T G A G G A C T A G	
Machine Translation	"hi how are you"	"Hello, como estas?	
Video activity recognition	Image	action on the video	
Name Entity Recognition	Yesterday, Harry Potter met Hermione Granger	Yesterday, Harry Potter met Hermione Granger	
Just to mention some...			

All of these problems can be addressed as supervised learning with labeled data x, y ; In some, both the input x and the output y can have different lengths, and in some of these examples only either x or y is a sequence.

Notation: As a motivating example, let's say we want to build a sequence model to input a sentence like this:

$x = \underline{\text{Harry}} \ \underline{\text{Potter}} \ \underline{\text{and}} \ \underline{\text{Hermione}} \ \underline{\text{Granger}}$ invented a new spell.

- And we want a model to tell us where are the **people's names** in this sentence. So this is a problem called Name-Entity recognition, and this is used for example by search engines to index, let's say, the last 24 hours of news, and looks for people names, so they can be indexed appropriately.
- Name-Entity recognition systems can be used to find people's names, company names, times, locations, country names, currency names and so on, on diff. types of text.

- Now, given the input x , let's say we want a model to output y that has 1 (zeros) per input words, and the target output, the desired y tells for each of the input words if that's part of a person's name! Technically, this is not the best output representation, there are more sophisticated output representations that tells not only if it's a name, but tells you where are the start and ends of people's names in the sentence!

$$^1 Y = [1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0]$$

$$^2 Y = [(1 \ 1) \ 0 \ (1 \ 1) \ 0 \ 0 \ 0 \ 0]$$

For this motivating example, we will stick with this simpler output representation. The input is a sequence of nine words, and we will use superscript $\langle t \rangle$ to denote or index positions. t implies are those are temporal sequences (whether the sequence is temporal or not). On this example we have $x^{\langle 1 \rangle} \dots x^{\langle 9 \rangle}$, and we have also will have for the output the same notation $y^{\langle 1 \rangle} \dots y^{\langle 9 \rangle}$.

$X = \text{Harry Potter and Hermione Granger invented a new spell.}$

$x^{\langle 1 \rangle} x^{\langle 2 \rangle} x^{\langle 3 \rangle} x^{\langle 4 \rangle} x^{\langle 5 \rangle} x^{\langle 6 \rangle} x^{\langle 7 \rangle} x^{\langle 8 \rangle} x^{\langle 9 \rangle}$ $Tx=9$

$Y = \begin{matrix} 1 & \emptyset & 0 & 1 & 1 & 0 & 0 & 0 & 0 \end{matrix}$ $Ty=9$

As we know, the length of X and Y , don't necessarily have to match! But on this example, both $Tx=9$ and $Ty=9$. Finally, to denote the current sample, we continue using $\langle i \rangle$, to refer to the i th training example. So:

$x^{\langle i \rangle \langle t \rangle} = t^{\text{th element of the training example } i}$; $Tx^{\langle i \rangle} = \text{sequence length}$
 $y^{\langle i \rangle \langle t \rangle} = t^{\text{th element of the output sequence of train example } i}$; $Ty^{\langle i \rangle} = \text{sequence length of } i$

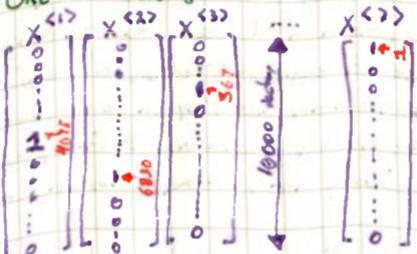
This is our first serious foray into NLP or Natural Language Processing, and one of the things we need to decide is, how to represent individual words in the sequence.

To represent a word in a sentence, the first thing we do, is to come up with a vocabulary, sometimes also called a dictionary, and that means, making a list of the words that we will use in our representations. Starting with A... till the end...

Vocabulary

aaron	1
and	2
...	367
...	...
harry	4075
potter	6830
zulu	10,000

One hot encoding:



Real world dictionaries contain from 30 to 50k or 100k words.
Big internet companies have dictionaries even of 1 million words

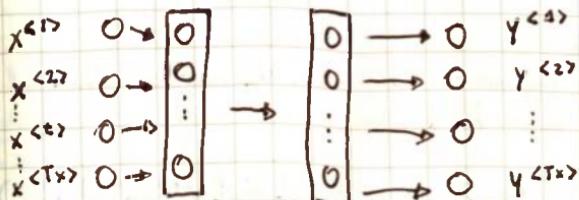
A way to come up with our own dictionary is to take the 'top n ' most common words in our sentences and use only those.

If we encounter a word that is not present in our vocabulary, then we can create a nowtaker or fake word, such as unknowm, more later or

[Recurrent Neural Networks]

Why not a standard network?

At first glance we could think on using a standard/naive NN to do the work, to learn the mappings between input words and entities on the output; even we could pad the input and output to match a consistent expected length of input data, however, this doesn't work well



Problems:

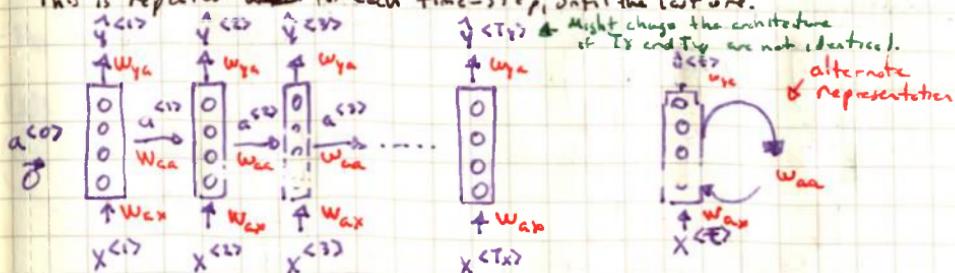
- Inputs and outputs can be different in different examples.
- Doesn't share features learned across different positions of text.

And the recurrent NN doesn't have these disadvantages.

[Recurrent Neural Networks]

In simple words, in a RNN, each $X^{<t>}$ is being fed into a layer, together with the activation of the layer before. On the first layer, a vector of zeros \emptyset or a randomly initialised vector w is used.

This is repeated until for each time-step, until the last one.

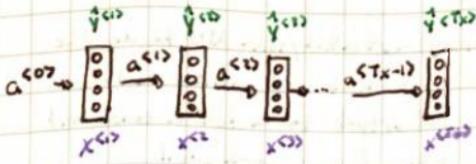


W_{ax} = It's the sum from across each time step (shared), governs the move from $X^{<t>}$ to the $O^{<t>}$
 W_{aa} = Activation are governed as well by a shared parameter.

One weakness of the RNN is that uses only information that is earlier in the sequence to make a prediction, it uses past information, but doesn't use information later in the sequence

- He said, "Teddy Roosevelt was a great President" & Should be a name
- He said, "Teddy bears are on sale!" & Should not be a Name
— would solve it later with Bidirectional RNNs

Forward Propagation



$$\begin{aligned} a^{t-1} &= \vec{0} \\ a^{t-1} &= g(W_{aa} a^{t-1} + W_{ax} x^{t-1} + b_a) \\ y^{t-1} &= g(W_{ya} a^{t-1} + b_y) \end{aligned}$$

\vdots

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

Simplified RNN Notation:

A way to deduce or simplify the notation is to use a W_a matrix which will contain both the a , or the x values.

$$a^{t-1} = g(W_{aa} a^{t-1} + W_{ax} x^{t-1} + b_a) \rightarrow a^t = g(W_a [a^{t-1}, x^{t-1}] + b_a)$$

$$\hat{y}^t = g(W_{ya} a^t + b_y) \rightarrow \hat{y} = g(W_y a^t + b_y)$$

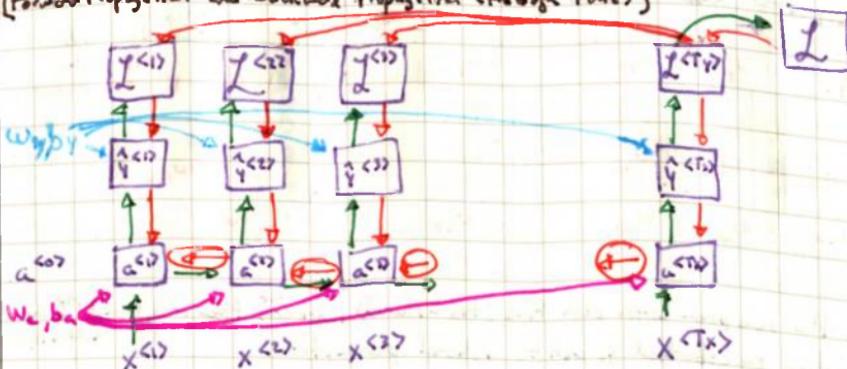
The simplification comes from stacking W_{aa} and W_{ax} ($100,100$) and x^{t-1} ($100,10000$) resp., so

$$\begin{matrix} 100 \\ \downarrow \\ \xrightarrow[100]{\left[\begin{matrix} W_{aa} : W_{ax} \end{matrix} \right]} \end{matrix} \xrightarrow[10000]{\left[\begin{matrix} a^{t-1}, x^{t-1} \end{matrix} \right]} = W_a$$

$$\left[\begin{matrix} a^{t-1}, x^{t-1} \end{matrix} \right] = \left[\begin{matrix} a^{t-1} \\ \dots \\ x^{t-1} \end{matrix} \right] \quad \begin{matrix} 100 \\ \downarrow \\ 10K \\ \downarrow \\ 100 \end{matrix}$$

$$[W_{aa} : W_{ay}] = \begin{bmatrix} a^{t-1} \\ x^{t-1} \end{bmatrix} = [W_{aa} a^{t-1} + W_{ax} x^{t-1}]$$

[Forward Propagation and Backward Propagation (Through time)]

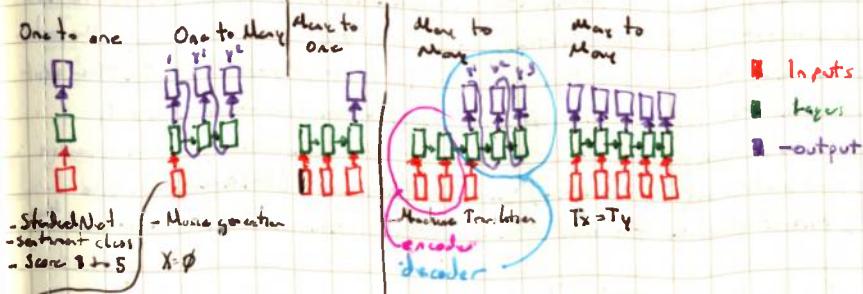


- Forward Propagation
- Backward Propagation
- Backward Propagation Through Time
- Parameters shared for a
- Parameters shared for y

$$\begin{aligned} \mathcal{L}^{t-1} (\hat{y}^{t-1}, y^{t-1}) &= -y^{t-1} \log \hat{y}^{t-1} - (1-y^{t-1}) \log (1-\hat{y}^{t-1}) \\ \mathcal{L}(y_i, \hat{y}_i) &= \sum_{t=1}^T \mathcal{L}^{t-1} (\hat{y}^{t-1}, y^{t-1}) \end{aligned}$$

[Different Types of RNNs]

We've seen so far RNNs where the number of input parameters match the number of output parameters, but that's not the case for different applications.



[Language Model and Sequence Generation]

What is language modelling? **Def:** What a language model does, is that given a sentence its job is to tell you what's the probability of a sentence, for example, given sentence:

- The apple and pear saled [pear and pear sound similar but a good language model would say "pear" is the correct word.]
- The apple and pear saled $P(\text{the apple and pear saled}) = 3.2 \times 10^{-13}$

$$P(\text{the apple and pear saled}) = 5.7 \times 10^{-10} \quad P(y^{(1)}, y^{(2)}, \dots, y^{(T)})$$

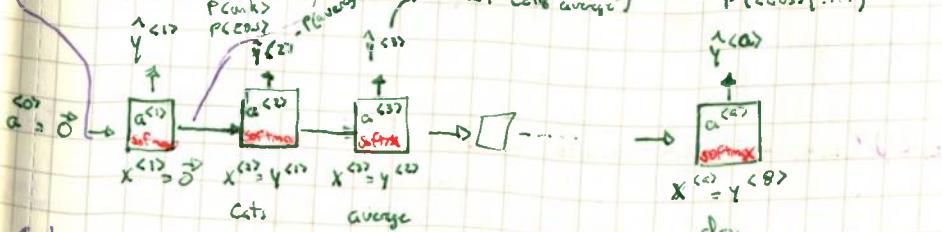
fundamental for speech recognition and machine translation

- Training set: Large corpus of X language text. \rightarrow corpus = NLP terminology for a very large body of text sentences.
→ Tokenize \rightarrow Gets Average 15 hours of sleep a day. $\langle \text{EOS} \rangle$ = end of sentence
 $y^{(1)} \ y^{(2)} \ y^{(3)} \ y^{(4)} \ y^{(5)} \ y^{(6)} \ y^{(7)} \ y^{(8)} \ y^{(9)} \ y^{(10)}$ Tokenizer

- Flag unknown words in the dictionary. The Elephant \rightarrow CUNKY \rightarrow $\langle \text{EOS} \rangle$ is a breed of cats, $\langle \text{EOS} \rangle$

RNN Model: Start with zeros vector, predict, and next layer has as input the first token.

$$P(a^{(1)}, P(a^{(2)}), \dots, P(a^{(t)}), \dots, P(a^{(n)})) = P(a^{(1)}) P(a^{(2)} | a^{(1)}) \dots P(a^{(t)} | a^{(1)}, \dots, a^{(t-1)}) \dots P(a^{(n)} | a^{(1)}, \dots, a^{(n-1)})$$



- Gets average 15 hours of sleep a day. $\langle \text{EOS} \rangle$

- each word is fed after the first layer, then probabilities are calculated in sequence softmax of the first word, then softmax of the second word (including the previous), and soon, cumulating each word all together to output the probability of a sentence.

$$\text{Cost function: } L = \sum_i y^{(i)} \log p_i \quad \leftarrow \text{softmax loss function}$$

Sample for simple sentence $y^{(1)}, y^{(2)}, y^{(3)}$:

$$\text{Loss: } L = \sum_i y^{(i)} (p_i - p_{y^{(i)}})$$

$p_{y^{(1)}} = \frac{p(y^{(1)}, y^{(2)}, y^{(3)})}{p(y^{(1)}, y^{(2)}) + p(y^{(1)}, y^{(3)}) + p(y^{(2)}, y^{(3)})}$

[Sample Novel Sequences]

After we train a model, one of the ways we can informally get a sense of what's learned is to have sample novel sequences.

To sample, we just want the NN to generate sentences based on the training, to do so, we do a normal sampling scan on the softmax the first word we will use. Means, to take the first word that comes out (based on probabilities), then we for example np.random.choice, on the vector, and sample according to the distribution defined by the vector of probabilities. On the next time step, pass the first prediction, and continue in the same fashion until the end.

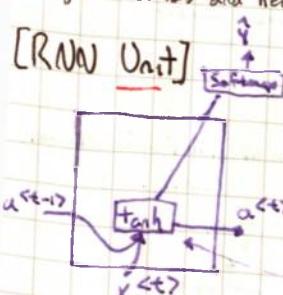
- A character level language model can also be used, however is computationally more expensive vs the word model, and the vocabulary would be the alphabet {0-9A-Z}. One benefit is that it eliminates the need of using `<UNK>` tokens, as all letters would be in the list.

[Vanishing Gradients with RNNs]

Recurrent Neural Networks, as we've seen so far, are not very good at capturing very long-term dependencies. For example, in the sentence, where there should be coherence, i.e. ^{long-term} ~~local~~ ^{global} parts,
The cat which already ate was full. The cats wh were full.
To address it, we'll address some options, while as for exploding gradients, the procedure is much simpler/easier, by using "Gradient Clipping".

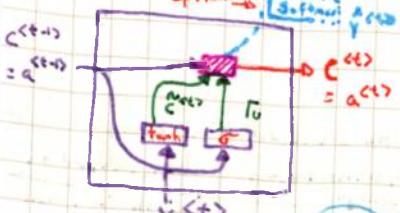
Gated Recurrent Unit (GRU)

Is a modification to the RNN hidden layer that makes it much better capturing long range connections and helps a lot with the vanishing gradient problems.



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

GRU (Simplified) - New variable $\rightarrow c = \text{cell}$ (memory cell)
optimization



$$c = \text{memory cell}$$

$$c^{(t+1)} = a^{(t+1)}$$

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

Candidate to
replace c

$$\tilde{c}^{(t+1)} = \sigma(W_u [c^{(t+1)}, x^{(t+1)}] + b_u) \quad 0 < \sigma < 1$$

"update"

$C^{(t)}$ can be a vector, then also
 $\hat{C}^{(t)}$, and F_0 as well. In
fact one the multiplication
are element wise!

[Full GRU]:

Modifying a bit the previous equations to account for relevance r_t :

$$\begin{aligned} \tilde{c}^{t+1} &= \tanh(W_c[r_t * c^{t+1}, x^{t+1}] + b_c) \\ r_t &= \sigma(W_r[c^{t+1}, x^{t+1}] + b_r) \\ f_t &= \sigma(W_f[c^{t+1}, x^{t+1}] + b_f) \\ c^{t+1} &= r_t * \tilde{c}^{t+1} + (1 - r_t) * c^{t+1} \end{aligned}$$

r_t = relevance of word

in academic literature
 h
 u
 r
 h

[LSTM: Long Short Term Memory]

This is even more powerful than the GRU; This was a really seminal paper that had a huge impact on sequence modeling. This paper is one of the most difficult to read.

GRU and LSTM

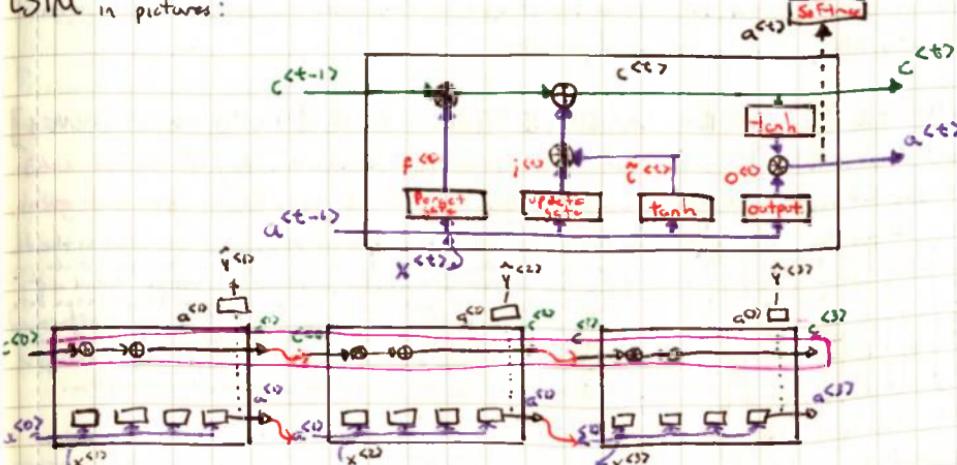
GRU

$$\begin{aligned} \tilde{c}^{t+1} &= \tanh(W_c[r_t * c^{t+1}, x^{t+1}] + b_c) \\ r_t &= \sigma(W_r[c^{t+1}, x^{t+1}] + b_r) \\ f_t &= \sigma(W_f[c^{t+1}, x^{t+1}] + b_f) \\ c^{t+1} &= r_t * \tilde{c}^{t+1} + (1 - r_t) * c^{t+1} \\ a^{t+1} &= c^{t+1} \end{aligned}$$

LSTM

$$\begin{aligned} \tilde{c}^{t+1} &= \tanh(W_c[a^{t+1}, x^{t+1}] + b_c) \\ r_t &= \sigma(W_r[a^{t+1}, x^{t+1}] + b_r) \quad \text{update} \\ f_t &= \sigma(W_f[a^{t+1}, x^{t+1}] + b_f) \quad \text{forget} \\ o_t &= \sigma(W_o[a^{t+1}, x^{t+1}] + b_o) \quad \text{output} \\ c^{t+1} &= r_t * \tilde{c}^{t+1} + f_t * c^{t+1} \\ a^{t+1} &= r_t * c^{t+1} \end{aligned}$$

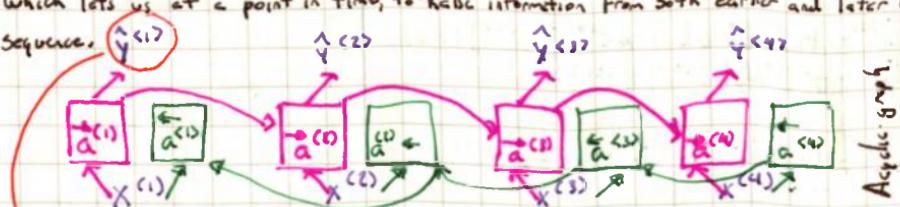
LSTM in pictures:



A common variation is the "peephole connection", meaning that the gate values may depend not just on a^{t+1} and x^{t+1} , but also on the previous memory cell value c^{t+1} .

Bidirectional RNN (BRNN)

By now we've seen most of the key building blocks of RNNs, one is bidirectional RNNs which lets us, at a point in time, to track information from both earlier and later in the sequence.



$$Y^{(t+1)} = g(w_t[\vec{a}, \vec{x}] + b_t)$$

- Forward in time
 - backward in time - is forward prop! not back prop!
 - computations are done with network activation =

For example, at time 3:

- Production flows from $X^{(t)} \rightarrow a^{(t)} \rightarrow c^{(t)} \rightarrow a^{(t+1)} \rightarrow y^{(t+1)}$ so information from $X^{(t)}, X^{(t+1)}$ and $a^{(t)}$ counts.
 - While info. from $X^{(t+1)} \rightarrow a^{(t+1)} \rightarrow c^{(t+1)} \rightarrow y^{(t+1)}$
 - So this allows at T3 to take account its input to take info from the past, present and future
 - / the blocks can be, not only the standard RNN blocks, but also GRU or LSTM blocks; for NLP or text problems, a BRNN with BiLSTM is commonly used.

The disadvantage of the Bidirectional Recurrent Neural Networks is that we do need the entire sequence of data before we can make predictions anywhere. For example if we are building a speech recognition system, then the BRNN will let us take into account the entire speech utterance, but if we use the straightforward implementation, then we need the person to stop talking to make the actual prediction.

WEEK 2: NATURAL LANGUAGE PROCESSING & WORD EMBEDDINGS

[Word Representation]: So far we've been representing words in a vocabulary $M = \{a, \text{aaaa}, \dots, d\}$ as one-hot vectors. It trivially encodes nothing into itself and doesn't allow the algo. to generalize.

The relationship between apple and orange is not any closer to the other words like star, women, king, queen, or "orange juice" is also common context as "apple juice", and the NN has reached us at a popular phase, it won't realize the similarity. ~~With~~ ^{With} you see both fruits are similar things, that's because the inner product between their two hot vectors is zero, and there's no distance, so it doesn't know that apple and orange are more similar than apple and green orange and king, so finally we can use neither a Factorized representation,

~~08391 09653 04914 0247 0456 06257~~ tuning features on the front end side

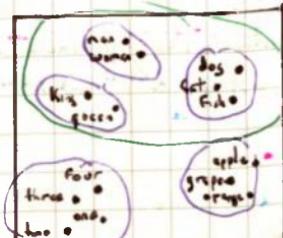
[Factored Representation]

Gender	Men	Women	King	Queen	Apple	Orange
total	5371	9853	4914	7187	456	6257
Race	0.1	1	-0.95	0.97	0.00	0.01
age	-0.01	0.02	0.93	0.95	-0.01	0.00
Food	0.03	0.02	0.7	0.69	0.02	-0.02
size	0.04	0.01	0.18	0.01	0.65	0.97
cost	1	*	*	*	*	*
altru	*	*	*	*	*	*

• Sandarity

i.e. 300 days \downarrow \rightarrow

Visualizing word Embeddings ($300D \rightarrow 2D$)



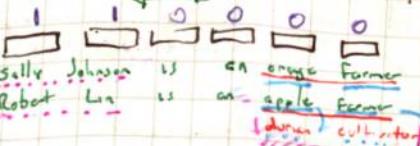
t-SNE

6m

the strong
similarity of
many Features
and rare
Features even
in common

[Word Embeddings]

Once we have feature representations of words, the embedding vectors, we can 'easily' design an algorithm that can generalize for example, similarities like "farmer" and "cultivator" or names, persons. But for that, we should have seen these words on the training set, so it can infer that "cultivator" is also a farmer or "dunin" is a fruit, but also even that a "dunin cultivator" is a person.



For this, the algorithm should examine 3 Billion words \rightarrow 100 Billion \rightarrow huge but reasonable task to learn embeddings that group them together. Take the word embeddings and apply it to some other recognition task. Transfer learning to train a smaller "100K" related training set.

(Transfer Learning and Word Embeddings)

1: Learn Word Embedding from large text corpus (2-100B words) (or download pre-trained embeddings online)

2: Transfer embedding to new task with smaller training set (Sep 100K words)

3: Optional: Continue fine-tuning the word embeddings with new data.

Men:Women : Boy:Girl
Others:Gods : Vampire:Krone
Year:Japan : Table:Russia
Big:Russia : Tall:Taller

[Properties of Word Embeddings]: Help with analogy reasoning and convey a sense of what these words embeddings are doing. Using the dot product, and assuming a D vector only:

	Men	Woman	King	Queen	Apple	Orange	6257
Men	0.851	0.853	0.814	0.817	0.781	0.777	0.821
Woman	0.01	0.02	0.03	0.02	0.01	0.01	0.01
King	0.03	0.02	0.10	0.09	0.03	0.02	0.03
Queen	0.01	0.01	0.02	0.01	0.01	0.01	0.01
Apple	-0.95	0.97	0.00	0.00	-0.95	0.97	-0.95
Orange	0.97	-0.95	-0.01	0.01	0.97	-0.95	0.97
6257	0.00	0.00	0.00	0.00	0.00	0.00	0.00

$$\text{Man} - \text{Woman} + \text{King} = ?$$

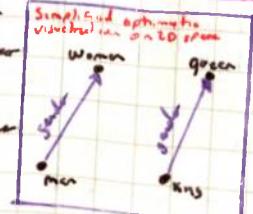
$$\text{Queen} - \text{Woman} + \text{King} = ?$$

$$\text{Orange} - \text{Apple} + \text{Orange} = ?$$

so to find the word P , represented as e_P :

$$e_P = \text{similarity}(e_W, e_B - e_A + e_C)$$

Previously we talked about using t-SNE to visualize words, which t-SNE does, but it takes 300D data and maps it in a very non-linear way to 2D space and so the mapping with t-SNE learns many, we can't expect these types of relationships to hold true, but in higher dimensional space, analogies hold.



The most common similarity function is the cosine similarity:

$$\text{Sim}(e_U, e_V) = \text{Sim}(U, V) = \frac{U \cdot V}{\|U\|_2 \|V\|_2}$$

excluding $\|U\|_2 \|V\|_2$ if similar the lengths, product will be large

Another option, but far "dissimilarity" is to use square distances, $\|U - V\|^2$

Dist. is how it anomalies between vectors U and V

[Embedding Matrix] 10000×100 (100D)

0.0000 | 0 = Features

$E \cdot O_j = e_j$ = embedding forward for the vocabulary

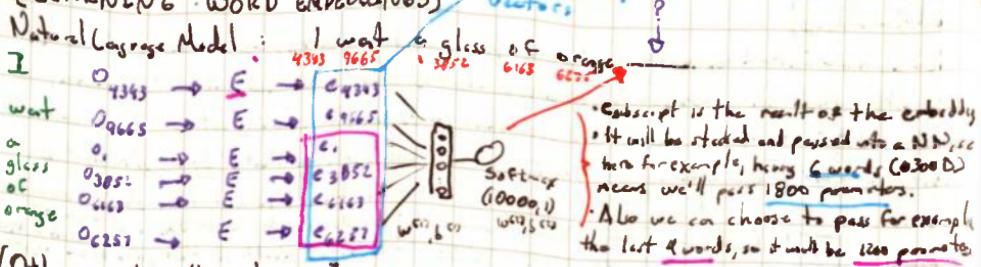
Note: On the relatively high dimensional matrix doing the vector multiplication is not so efficient because the count of zeros, so we will better use a specialized and more efficient way like with Keras function "Embedding"

$$E \cdot O_{6257} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = e_{6257}$$

Screening

Embedding matrix E , initialized randomly and use gradient descent to learn all the parameters of this 10000×10000 parameter matrix.

[LEARNING WORD EMBEDDINGS]



[Other context/Target pairs]

$$X = I \text{ want a } \xrightarrow{\text{glass of orange}} \text{ juice} \xrightarrow{\text{target}} \text{to go along with my cereal.}$$

Content: last 4 words

- 4 words on left and right: $\text{a glass of orange} \xrightarrow{\text{?}} \text{to go along with}$
- Last 2 words: $\text{orange} \xrightarrow{\text{?}}$
- Nearby 1 word: $\text{gloss} \xrightarrow{\text{?}}$ } Skip Gram Model

[Word2vec] model: Simpler and computationally efficient to learn ~~this~~ embeddings (Skip-grams). Mikolov et al., 2013. Efficient estimation of word representations in vector space.

In the Skip-gram model, what are we going to do is to come up with a few context to target pairs to create our supervised learning problem. So rather than having the context be always the last 4 words or last 2 words immediately after or before the target word, randomly pick a word to be the context word; let's say, the word orange and randomly select another word within some window. Say ± 5 or ± 10 words of the context words. So it's happening that we're choosing the immediate neighbors of the target word.
Obviously this is not an easy learning problem, because between $\pm n$ words, there could be there are many words on both sides

$$X = I \text{ want a } \xrightarrow{\text{glass of orange}} \text{ juice} \xrightarrow{\text{target}} \text{to go along with my cereal.}$$

Content	Target
Orange	juice
Orange	glass
Orange	my

Model: Vocabulary 10000 words
Context C("orange") \rightarrow Target t ("juice")
 $6257 \xrightarrow{4834}$

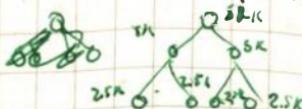
$$O_c \rightarrow E \rightarrow c_t \xrightarrow{\text{softmax}} Y_t$$

$$\text{Softmax} \Rightarrow p(c|t) = \frac{e^{E_t^T c_t}}{\sum_{i=1}^{10000} e^{E_i^T c_t}} \Rightarrow Q_t = \text{counter associated with output } t$$

$$Y_t = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \text{ weight } 4834$$

Paper:
Mikolov et al., 2013: Efficient estimation of word representation in vector space.

This 10K is already too slow, even worse with vectors of 1M, or more! An alternative is the hierarchical softmax, which uses a balanced tree with binary classification. In theory, the thresholds are balanced, however in this application, the thresholds are set by the most common at the top, to the less common at the bottom.

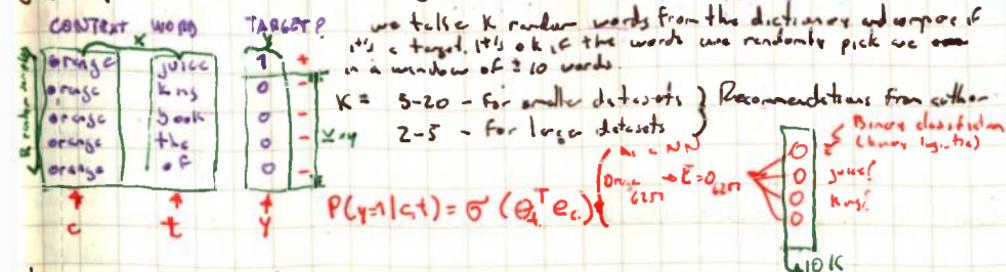


To sample "G": If unmodified, it will sample very common words like "the", "is", "and", so in practice the distribution is skewed. P(c) isn't uniform, therefore diff. frequencies to below it with less

Demand diff. frequencies to below it with less

[Negative Sampling]: Mikolov et al., 2013: Distributed representation of words and phrases and their compositionality
 We already saw how the Skip-Gram model allows us to construct a supervised learning task. So we map from context to target and how that allows us to learn useful word embeddings. But the downside of that, was that the softmax objective was too slow to compute. Now we will address Negative sampling that will allow us to do something similar to the Skip-Gram Model, but with a much more efficient algorithm.

What are we going to do to create a new supervised learning problem: given a pair of words, we are going to predict "is this a context-target pair?". $x = \text{I want a glass of orange juice to go along with my cereal}$



In comparison to the straightforward softmax, we end up with a 10K logistic regression classification problem, which is going to classify if the word is juice, book, etc. It's having 10K but on each iteration is going to train only $k+1$ logistic units.

(Selecting Negative Samples)

Sample proportionally by the frequency of the word: $P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{1000} f(w_j)^{3/4}}$) Author recommended and validated

[GloVe word vectors]: Pennington et al., 2014. 6BBo: Global Vectors for word representations

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

co-context c, t
 target $x_{ij} = \# \text{ times } i \text{ appears in the context of } j$
 co-target x_{ji}
 word

Might happen that depending on the choice of symmetry, x_{ij} would be before and after over a window of K so $x_{ij} \neq x_{ji}$; if restricted to be before only then it won't happen.

Model

$$\text{minimize}_{\theta} \sum_{i=1}^{1000} \sum_{j=1}^{1000} f(x_{ij}) \cdot (\theta_i^T e_j + b_i + b_j - \log x_{ij})^2$$

$\phi \log \phi = 0$ $f(x_{ij}) = 0 \text{ if } x_{ij} = 0$ $\text{if } x_{ij} \neq 0 \text{ then add } \beta \text{ to sum over the } x_{ij} \text{ per row}$

$\theta_i \text{ and } e_j$ are symmetric
 $e_w^{(final)} = \frac{e_w + \theta_w}{2}$

$f(x_{ij})$ also helps to give proper weights to words that appear less frequently, but also gives (but not unduly) weights to more common words. (weighting factor)
 - frequent words: this, a, or a
 - less frequent: dinner ...
 check the paper for implementation of $f(x_{ij})$

The conclusion

[Sentiment Classification]

The idea is to identify the meaning whether positive or negative, for example with rating where we have examples like "good at taste" $\rightarrow \star\star\star\star\star$. "The dessert is excellent" $\rightarrow 5$ stars. "The dessert is excellent" $\rightarrow 5$ stars.

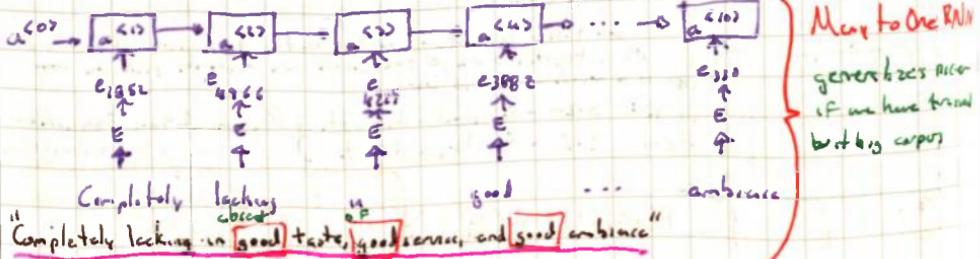
the $\rightarrow 08928 \rightarrow E \rightarrow e2228$
 dessert $\rightarrow 08928 \rightarrow E \rightarrow e2468$
 excellent $\rightarrow 04644 \rightarrow E \rightarrow e4644$
 $\rightarrow 03180 \rightarrow E \rightarrow e3180$

$\star\star\star\star\star$ \rightarrow [AUG] \rightarrow 0 softmru \rightarrow 5 stars

\rightarrow 20000 \rightarrow 4

Simple model, works reasonably well
 it averages the meaning of words, +
 ignores word order, so might give
 false positives

[RNN For Sentiment Classification]



[Debiasing Word Embeddings]: Bolukbasi et al., 2016. This is to computer programmer as housewife is to homemaker?
 The algorithms interpret the data as it is given them, so they can make stereotypical or biased assumptions based on the word embeddings; Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of the text used to train the model:
 - Men : Women, King : Queen
 - Men : Computer programmer, Women : Homemaker
 - Father : Doctor, Woman : Nurse : Nurse

[Addressing bias in word embeddings]

1: Identify bias direction

$$E_{he} - E_{she}$$

$$E_{male} - E_{female}$$

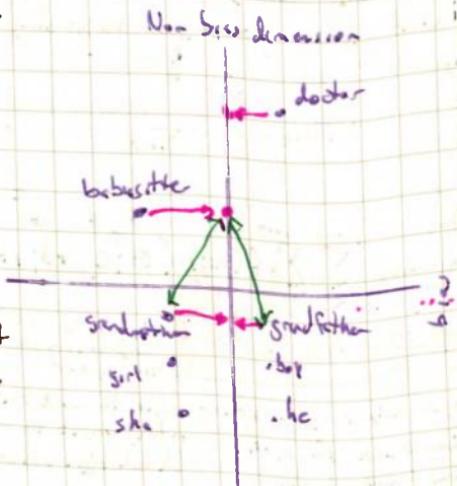
E...:

Average

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

3: Equate pairs

$$\begin{aligned} \text{grandmother} &= \text{grandfather} \\ \text{girl} &= \text{boy} \end{aligned}$$



- - WORK 3 : Sequence Models and Attention Mechanism]]
for Machine translation to speech recognition.

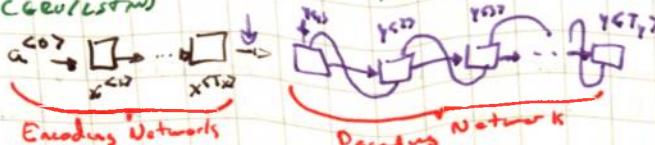
Model 3

versus to sequence model

The sentence is French

+ there's ~~CARROT~~ (GRILLED)

at a vector that
sits the is part
not

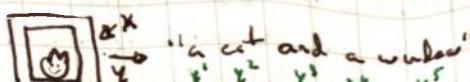


Sutcliffe et al.,
2014. Seq2Seq
learning with NN

Chap. 1, 2014

Learning phase representation using LNU encoder-decoder for statistical machine translation.

se Captain on



Like a ConvNet, here it learns and encodes a lot of features of the picture

at end of the softmax layer and via the
process vector that represents the image

Find the vector $\vec{v_{AB}}$ to a RW

and et al., 2014. Deep Captioning with multimodal RNNs

Engels et al., 2020. Show and tell: Neural image caption generation

Korbari and Li, 2015. Deep visual-semantic alignments for generating image descriptions

Picking the Most Likely Sentence

There is some similarity with the sequence-to-sequence machine translation model and the language models that we have worked on the fifth course, but there are some significant differences as well.

Language Model:  Allows us to estimate the probability of a sentence or generate novel sentences.

- Machine Translation:
 - Encoding Network
 - Decoding Network

Thus becomes into what Andrew calls "Conditional language model", as the parameters already encoded and representing the input sequence, are passed to the decoding network. This then becomes a calculation between the probability of an english translation, conditioned on a French sentence

[Finding the most likely translates]

Given the following sentence, the model will give us the probability of different corresponding english translations.

Jane visite l'Afrique en Septembre

$$P(Y^{(1)}, \dots, Y^{(T)} | X)$$

What we don't want is to sample translation at random! For example:

- 1: Jane is visiting Africa in September ✓
- 2: Jane is going to be visiting Africa in September ✗✓
- 3: In September, Jane will visit Africa ✗
- 4: Her African friend welcomed Jane in September ✗!

$$\arg \max_{y^{(1)}, \dots, y^{(T)}} P(y^{(1)}, \dots, y^{(T)} | x)$$

By chance we might sample a bad translation, so rather than using this model we want to find the sentence that maximizes the conditional probability.

An algorithm to do find the best sentence that maximizes the conditional probability, is called "Beam Search". But still important to mention other approaches that

• Greedy search: Search the best possible word in our language model, and iteratively, do the - **NO!** - same over and over. This picks only the most likely words. It will end up choosing by "common" rather than doing it by semantic way:

- Jane is visiting Africa in September } "going" might be more common
- Jane is going to be visiting Africa in September } than "visiting, thus."
 $P(\text{Jane is going} | x) > P(\text{Jane is visiting})$

[BEAM SEARCH]

(For drawing, check the stepped process)

Beam search differentiates itself from greedy search because this algorithm takes "B" samples at a time. If we set B=2, then it will become essentially a greedy algorithm.

The first thing it does, is to try to pull the first words of the translation; with B=3, then it will pick the "B" or 3 ~~most likely~~ maximum possibilities at a time (beam width).

So from the vocabulary, of let's say (in the example) of 10000 words, the 3 (or 2) most likely words will be chosen.

2 Pass each of the first $\frac{B}{2}$ words into an encoder-decoder network to get the first "D" possible words, (from the vocabulary) $P(y^{(1)} | x)$

3 Having picked the 3 top words, is for each of these choices, what should be the next word, for that, also using the same vocabulary, so for each of them, they will be passed to the network to predict the probability of the first word with a new "top word" from the dictionary, then, by passing each combination to 3 different instances of the network (3 in this case):

$$P(y^{(2)} | x, "a") , P(y^{(2)} | x, "b") \text{ or } P(y^{(2)} | x, "Jane") / P(y^{(2)} | x, "September")$$

On this example, that would mean that there will be 3×10000 comparisons, and the search would need to pull the top "B" again

4 Repeat again the process, now pass the top "3" two-word sentences to the 3 different instances of the model, and once each pass them "handwritten" and then the probability (conditional) should be now: $P(y^{(3)} | x, "a", "in September")$ and so on.

5 Repeat again with the next word, until $\langle \text{EOS} \rangle$ - end of sentence -

[Refinements on Beam Search]

[Beam search normalization]: Is a small change to the basic beam search algorithm that can help us to get much better results. Beam search is about maximizing the probabilities and with the current algorithm, the probabilities multiplied on each step, tends to reduce the values to a numerical underflow.

$$P(y^{(1)} \dots y^{(T_f)} | x) = P(y^{(1)} | x) \cdot P(y^{(2)} | x, y^{(1)}) \dots P(y^{(T_f)} | x, y^{(1)}, \dots, y^{(T_f-1)})$$

$$\arg \max \prod_{t=1}^{T_f} P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

Probabilities are from 0 to 1 thus multipliers
several values < 1 will gradually reduce the number
leading to underflow.

In practice, instead of maximizing this product we will take logs, as the log of a product becomes the sum of a log, and maximizing this sum of log probabilities, should give the same result in terms of selecting the most likely sentence y . We end with this, with a more numerically stable algorithm, that is less prone to numerical rounding errors or underflows.

$$\arg \max \sum_{t=1}^{T_f} \log \{P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})\}$$

However this still has a problem, and it's that it will prefer or "reward" shorter sentences as the value will be higher with less sentences than with more, because of the nature of the multiplication or the sum of logs, for that this needs to normalize it, by doing it's will reduce the penalty:

$$\frac{1}{T_f} \sum_{t=1}^{T_f} \log \{P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})\}$$

~~penalty hyperparameter~~
~~would be not precisely~~

$$\alpha = 0.7$$

$\alpha = 1$ - complete normalization
 $\alpha = 0$ - no normalization

To wrap up: As we run beam search, we will see a lot of sentences with length ℓ_1 , a lot with $\text{len}(\ell)=2$, a lot with $\text{len}(\ell)=3$, and so on. Maybe we run the search for 30 steps and we consider output sentences up to 30 (let's say). And so with $B=3$, we will be keeping track of the top 3 possibilities for each of those possible sentence lengths, $\ell=2, 3, 4$ and so on up to 30. Then we will look at all the output sentences and score them against the score. And so we can take our top sentences and just compute the objective function onto sentences that we have seen through the beam search process. And finally, of all of these sentences that we validate this way, we pick the one that achieves the highest value on the normalized log probability objective. Sometimes it's called a Normalized log likelihood objective. And then that would be the final translation your outputs.

(Beam Search discussion)

- Beam width B_f : The larger it is, the more possibilities we consider, but also computing requirements increase rapidly producing heavy Research papers.
- Large B : better results, slower
- Small B : worse results, faster

1 → 3 → 10 → 100 → 1000 → 3000

[Error analysis or Beam Search]

It's given for example a fragment translation by a human, but the model outputs a wrong one. It's important to determine if the problem is on the RNN or the BeamSearch algorithm. We can compute $P(\hat{Y}|x)$ and $P(\hat{Y}|X)$ and determine which is worse.

6. f. the training set and look at the training example and make the error reduction process

Human	Algorithm	$P(C z)$	$P(\bar{C} z)$	At Fault?
—	—	2×10^{-10}	1×10^{-10}	B
1	1	:	:	R
1	1	:	:	R

Case 1: $P(\hat{Y}|x) > P(Y|x)$

Better choose \hat{Y} . But F prefers "higher" $P(Y|x)$.

Case 2: $P(\hat{Y}|x) \leq P(Y|x)$

\hat{Y} is better translation than Y . But RNN pred.

$P(\hat{Y}|x) < P(Y|x)$: RNN is at fault

[Blew Score]: Bilingual Education Understudy

[Papineni et. al., 2012. Bleu: A method for automatic evaluation of machine translation
- Evaluating Machine translation]

To match equally good translations (here), that might be written similarly but never or express exactly the same. For example:

French: Le chat est sur le tapis

Reference 3. The cat is on the mat

Reference 2: There is a ~~act~~ on the mat.

MT : the the the the the the the the } Texts

We will give a word credit
up to the maximum times
it appears in the reference
sentences

Precision: $\frac{1}{7}$ Modified Precision: $\frac{2}{7}$

Blow Score on 5, from

	Score on signs	Count	Count clip
1 the cat	2x	1	:
2 cat the	1x	0	
3 cat on	1x	1	$\rightarrow \frac{4}{6} = \frac{2}{3}$
4 on the	1x	1	
5 the mat	1x	1	

- Ref 1: the cat is on the mat
- Ref 2: There is a cat on the mat
- MT: the cat is on the mat

MT : The cat the cat on the mat

- Bleu Scene on Vayours

$$P_{\frac{1}{k}} = \frac{\sum_{\text{Unigrams}} \hat{y} \text{ Count}_{\text{clip}}(\text{Unigram})}{\sum_{\text{Unigrams}} \hat{y} \text{ Count}_{\text{unigram}}}$$

\downarrow
Single words in isolation

$$P_n = \frac{\sum_{\text{n-gram} \in T} \text{Count}_{clip}(n\text{-gram})}{\sum_{\text{n-gram} \in T} \text{Count}(n\text{-gram})}$$

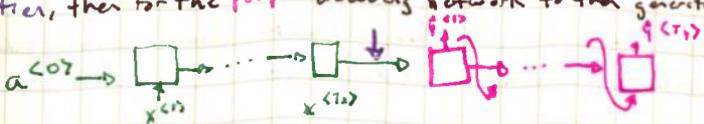
Details: Short translations give high scores, so we add BP (Brevity Penalty), which is a factor for adjustment, that penalizes translation systems that output translations that are too short.

$$B_p = \begin{cases} 1 & \text{if MT output length} > \text{reference output length} \\ \exp(1 - \text{MT output length} / \text{reference output length}) & \text{else} \end{cases}$$

[Attention Model Intuition]

-The problem of long sentences:

When we input a very long sentence, what we are asking to the green encoder to read in the whole sentence and then memorize it, the whole is stored in the activations, then for the purple decoding network to then generate the translation.



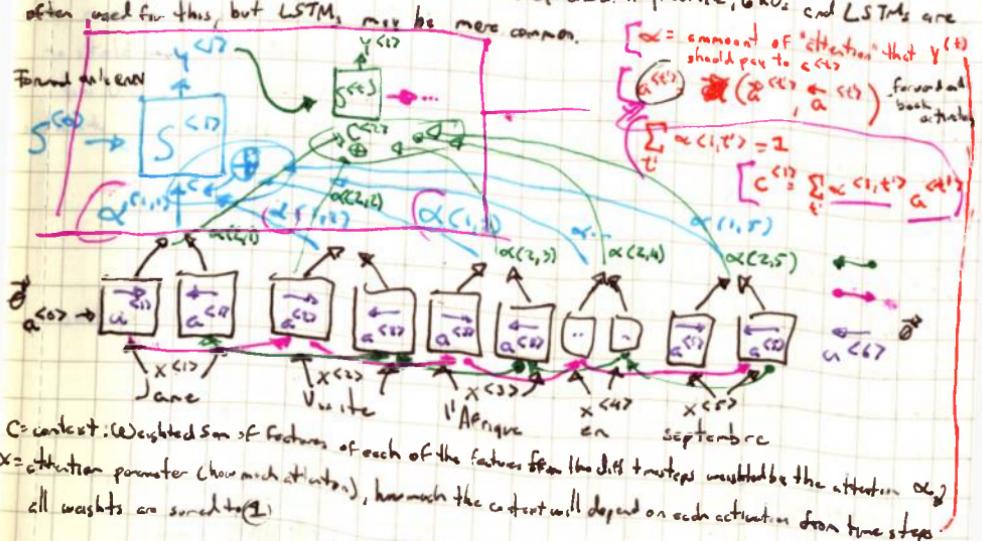
Now, the way a human translator would translate a long sentence is not by reading the whole thing, memorize it, and then regenerate the whole translated version from scratch. Instead the human translator would read the first part of the sentence and generate some words, read more and generate more words, and so on.



What we see for the encoder-decoder architecture shown above, is that it works quite well for short sentences, so we might achieve relatively high Bleu score, but for very long sentences, maybe larger than 30-40 words, the performance comes down. Later on we will see the improvements, but the Attention Model solves this, working as follows.

[Attention Model]: Pay attention to parts of the sentence.

Let's assume we have an input sentence, we use a RNN, or Bidirectional GRU or bidirectional LSTM to compute features of every word. In practice, GRUs and LSTMs are often used for this, but LSTMs may be more common.



(Bahdanau et al., 2014. Neural Machine Translation by learning to align and translate)

(Xu et al., Show, attend and tell: Neural image caption generation with visual attention)

[Speech Recognition - Audio Data]: Sequence to Sequence models applied to audio data
What is the speech recognition problem? Gives an audio clip $W \times H$ matrix a text transcript.

Audio is taken over three dimensions: time, frequency and features or spectrograms, thus we carry pre-processors.

Academic has research with 300 to 3000 hours of transcripts, whereas commercial grade systems have 10K or 100K hours of audio data and its transcript.

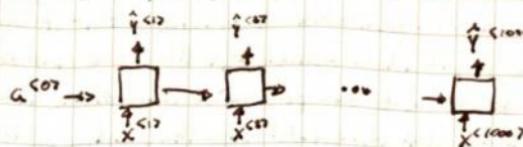
Having similar inputs as the attention model requires, it can be used as well.

CTC Cost for speech recognition [Graves et al., 2006. Connectionist Temporal Classification: segmenting audio labeling assignment (Connectionist Temporal Classification)]

"The quick brown fox"

Simplified for demonstrating purposes, as in reality we would need to use: Bidirectional LSTM +

segmented data with RNNs.



Bidirectional GRU, and deeper model

++ $_$ h $_$ eee $_$ $_$ l $_$ $_$ qqq $_$ $_$ the quick

Basic rule: collapse repeated characters not separated by "blank".

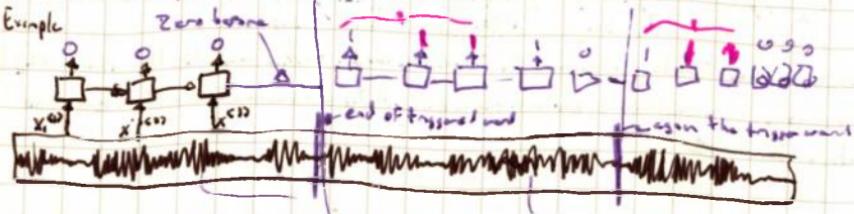
Notice the number of steps is very large (1000), it's fixed, and it's because as shown above, the model can output repeated characters, it will be filled with spaces (blanks) or the correct characters. So, the input sentence, even if it has 19 characters, the NN will output 1000, but once the repeated characters are collapsed and blanks removed, it will have succeeded. THIS REQUIRES LOADS OF DATA

[Trigger Word Detection]: Still evolving technology

The tech behind smart assistants like Alexa, Google Home, Siri, etc.

It requires less data to train and it can be used to trigger specific activities, like turning on a lamp, denoting a command, etc.

Example



Having to output so many zeros ~~or other values~~ might create imbalance, so what we can do is to make the model to output a bit more "one", so the ratio 1:0 gets more even.

