

Logistic Regression

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

Notation:

m : sample size

n_x : # of feature.

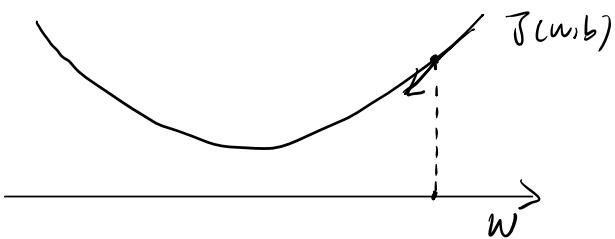
$w: n_x, 1$

$x: n_x, m$ (column/sample)

$$\text{Loss fx: } l(y, \hat{y}) = -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

$$\text{Cost fx: } J(w, b) = \frac{1}{m} \sum_{i=1}^m [-y \log \hat{y} - (1-y) \log(1-\hat{y})] \quad - \text{ convex fx}$$

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



Chain rule & Computation graph:

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

$$a = 5 \quad v = a + u \quad J = 3v$$

$$b = 3 \quad u = bc \quad du = dv \cdot 1 = 3$$

$$c = 2 \quad v = a + u \quad dv = 3$$

$$da = dv \cdot 1 = 3$$

$$db = b = 3$$

$$dc = 9$$

$$du = dv \cdot 1 = 3$$

$$u = bc ; v = a + u ; J = 3v$$

For logistic regression (1 sample)

$$x_1 \quad d\hat{z} = \frac{d\ell(a, y)}{da} \cdot \boxed{\frac{da}{d\hat{z}} = a - y} \quad \frac{d\ell(a, y)}{da} = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$w_1 \quad \boxed{Z = w_1 x_1 + w_2 x_2 + b} \rightarrow \boxed{\hat{y} = a = \sigma(Z)} \rightarrow \boxed{\ell(a, y)}$$

$$x_2$$

$$w_2$$

$$b$$

$$d\hat{z} = x_1 \cdot d\hat{z} ; d\hat{z} = x_2 \cdot d\hat{z} ; db = d\hat{z}$$

For logistic regression (m sample) - One iteration of batch gradient descent

$$\hat{z} \leftarrow \frac{1}{m} \sum_{i=1}^m \hat{z}_{i, \text{old}}$$

$$\overbrace{\sum_{i=1}^m \mathcal{L}(w_i; y_i)}^J(w, b) = \overbrace{\sum_{i=1}^m \mathcal{L}(w_i; y_i)}^{dw_i^{(i)} = (x^{(i)}, y^{(i)})}$$

$$J=0; \quad \boxed{dw_1=0; dw_2=0; db=0} \quad dw = np.zeros((n_x, 1))$$

For $i=1$ to m : For loop I.

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

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

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

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

$$dw_1 += X_1^{(i)} dz^{(i)}$$

$$dw_2 += X_2^{(i)} dz^{(i)}$$

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

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

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

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

For loop II.

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

assuming $n \geq 2$

$$J /= m; \quad dw_1 /= m; \quad dw_2 /= m; \quad db /= m$$

Vectorization: To get rid of explicit for loops.

Vectorizing Logistic Regression

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

$$x.shape = (n_x, m)$$

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

$$Z = np.dot(w.T, x) + b \quad \text{broadcasting}$$

$\begin{matrix} \text{Broadcasting Rule:} \\ (m, n) \xrightarrow{\text{+}} (1, n) \rightsquigarrow (m, n) \\ \text{matrix} \xrightarrow{\times} (m, 1) \rightsquigarrow (m, n) \end{matrix}$	
--	--

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

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

$$A = \sigma(Z)$$

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

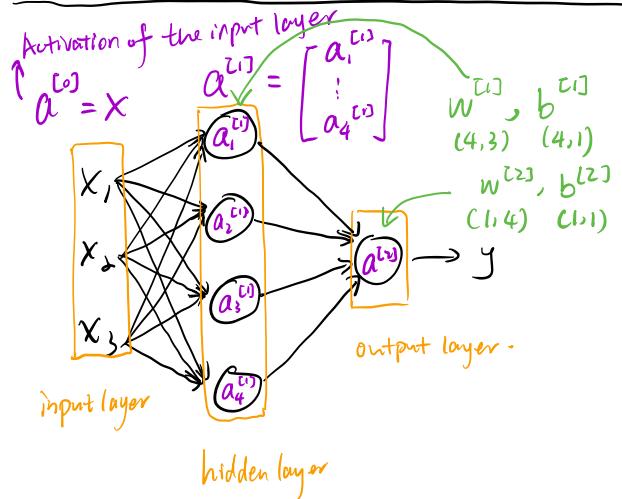
$$dZ = A - Y$$

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

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

$$\text{Ldb} = \frac{1}{m} \text{np.sum(Ldz)}$$

Neural Network Representation



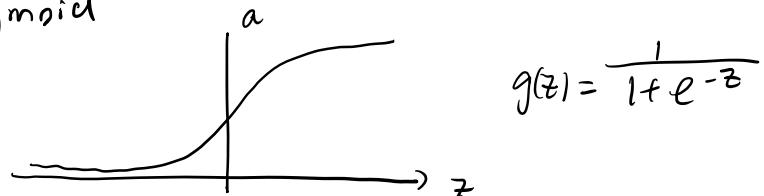
Vectorizing across multiple examples

$$X = \begin{bmatrix} | & | \\ X^{(1)} & \dots & X^{(m)} \\ | & | \end{bmatrix}_{n_x \times m}$$

$$\begin{aligned} z^{(i)} &= W^{(i)} X + b^{(i)} \\ A^{(i)} &= \sigma(z^{(i)}) \\ z^{(T)} &= W^{(T)} A^{(i)} + b^{(T)} \\ A^{(T)} &= \sigma(z^{(T)}) \end{aligned}$$

Derivatives of Activation functions

Sigmoid



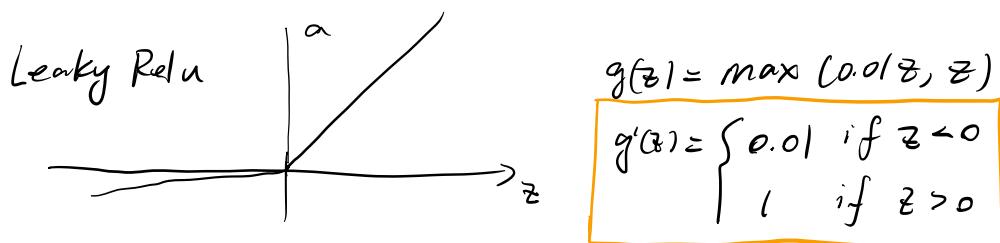
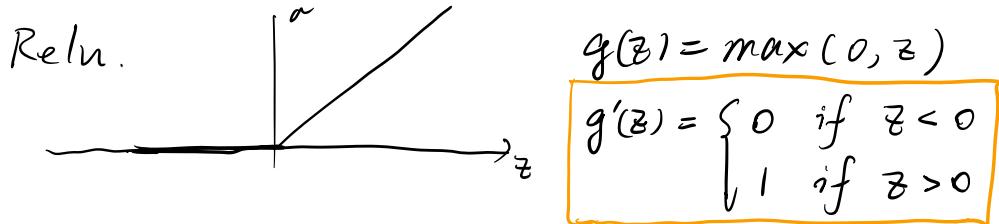
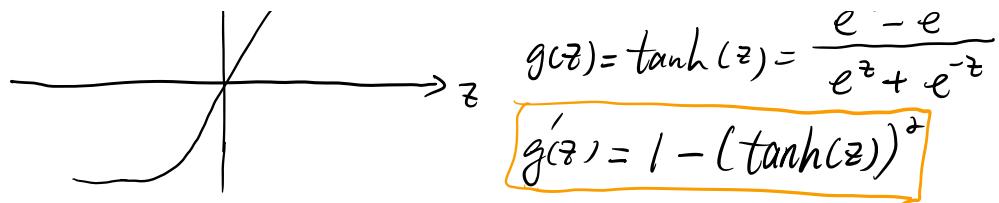
$$\frac{d}{dz} d(z) = \frac{d}{dz} (1 + e^{-z})^{-1} = -(1 + e^{-z})^{-2} e^{-z} \cdot (-1)$$

$$= \frac{e^{-z}}{(1 + e^{-z})^2} = \frac{e^{-z}}{(1 + e^{-z})(1 + e^{-z})} = \boxed{g(z) \cdot (1 - g(z))}$$

tanh

$$1 - \frac{e^{-z}}{e^z}$$

$$-z \quad -z$$



Gradient descent for neural Networks

Params: $w^{[l]}, b^{[l]}$, $n_x = n^{[0]}$, $n^{[l]}$, $n^{[l]} = 1$
 $(n^{[0]}, n^{[0]})$ $(n^{[0]}, 1)$ $(n^{[l]}, n^{[l]})$ $(n^{[l]}, 1)$ $n^{[l]} = 1$
 input feature hidden unit output

Cost function: $J(w^{[l]}, b^{[l]}, w^{[l]}, b^{[l]}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i)$

Gradient descent:

Repeat $\left\{ \begin{array}{l} \text{Compute predict } (\hat{y}^{(i)}, i=1 \dots m) \\ \frac{d w^{[l]}}{d w^{[l]}} = \frac{\partial J}{\partial w^{[l]}} ; \quad \frac{d b^{[l]}}{d b^{[l]}} = \frac{\partial J}{\partial b^{[l]}} \\ \text{update } w^{[l]}, b^{[l]}, w^{[l]}, b^{[l]} \end{array} \right\}$

Forward propagation:

$$z^{[l]} = w^{[l]} X + b^{[l]}$$

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

$$\bar{z}^{[l]} = w^{[l]} A^{[l]} + b^{[l]}$$

$$A^{[l]} = \underbrace{\sigma}_{\text{sigmoid}} \bar{z}^{[l]}$$

Backward propagation:

$$\frac{(n^{[l]}, m)}{d A^{[l]}} = -\frac{Y}{A^{[l]}} + \frac{1-Y}{1-A^{[l]}}$$

$$\frac{(n^{[l]}, m)}{d z^{[l]}} = d A^{[l]} \sigma'(z^{[l]}) = A^{[l]} - Y$$

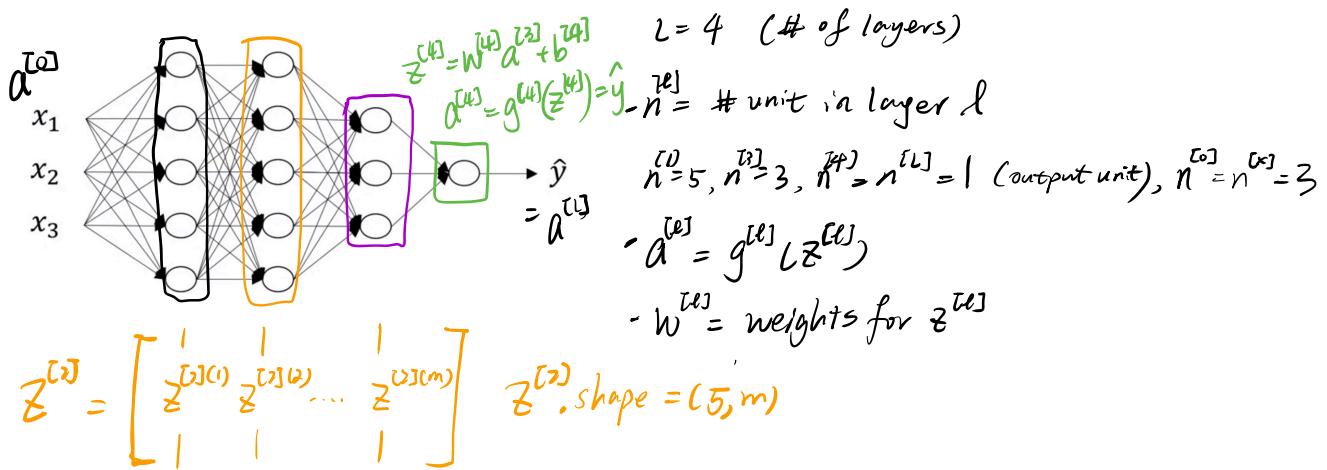
$$\frac{(n^{[l]}, n^{[l]})}{d w^{[l]}} = \frac{1}{m} \frac{(n^{[l]}, m)}{d z^{[l]}} A^{[l] T}$$

$$\frac{(n^{[l]}, m)}{d b^{[l]}} = \frac{1}{m} \text{np.sum}(d z^{[l]}, \text{axis}=1, \text{keepdims=True})$$

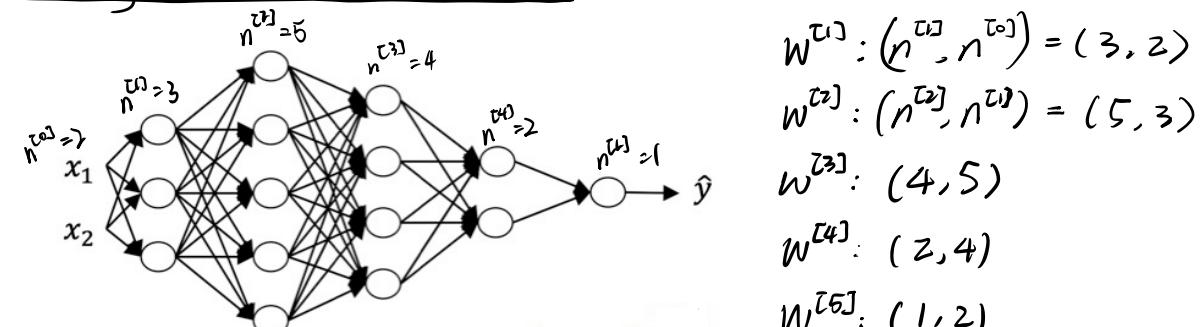
$$\frac{(n^{[l]}, m)}{d A^{[l]}} = W^{[l] T} \frac{(n^{[l]}, m)}{d z^{[l]}}$$

$$\begin{cases}
 \frac{\partial Z^{[l]}}{\partial Z^{[l]}} = (n^{[l]}, m^{[l]}) \\
 \frac{\partial A^{[l]}}{\partial Z^{[l]}} = \frac{1}{m} (n^{[l]}, m^{[l]}) X^T \\
 \frac{\partial b^{[l]}}{\partial Z^{[l]}} = \frac{1}{m} \text{np.sum}(dZ^{[l]}, axis=1, keepdims=True)
 \end{cases}$$

Deep neural network notation



Getting matrix dimension right



$$\begin{aligned}
 (3,1) &\leftarrow (3,2) \quad (2,1) \\
 z^{[l]} &= W^{[l]} \cdot x + b^{[l]} \\
 (n^{[0]}, m) &(n^{[0]}, n^{[1]}) (n^{[1]}, m)
 \end{aligned}$$

$$\begin{aligned}
 W^{[1]} : (n^{[0]}, n^{[1]}) &= (3, 2) \\
 W^{[2]} : (n^{[1]}, n^{[2]}) &= (5, 3) \\
 W^{[3]} : (4, 5) & \\
 W^{[4]} : (2, 4) & \\
 W^{[5]} : (1, 2) &
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial w}{\partial w} / \frac{\partial w}{\partial w} &: n^{[l]}, n^{[l-1]} \\
 \frac{\partial b}{\partial b} / \frac{\partial b}{\partial b} &: n^{[l]}, 1 \\
 \frac{\partial z}{\partial z} / \frac{\partial A}{\partial A} &: n^{[l]}, m
 \end{aligned}$$

Forward and Backward functions

For layer l : $W^{[l]}, b^{[l]}$

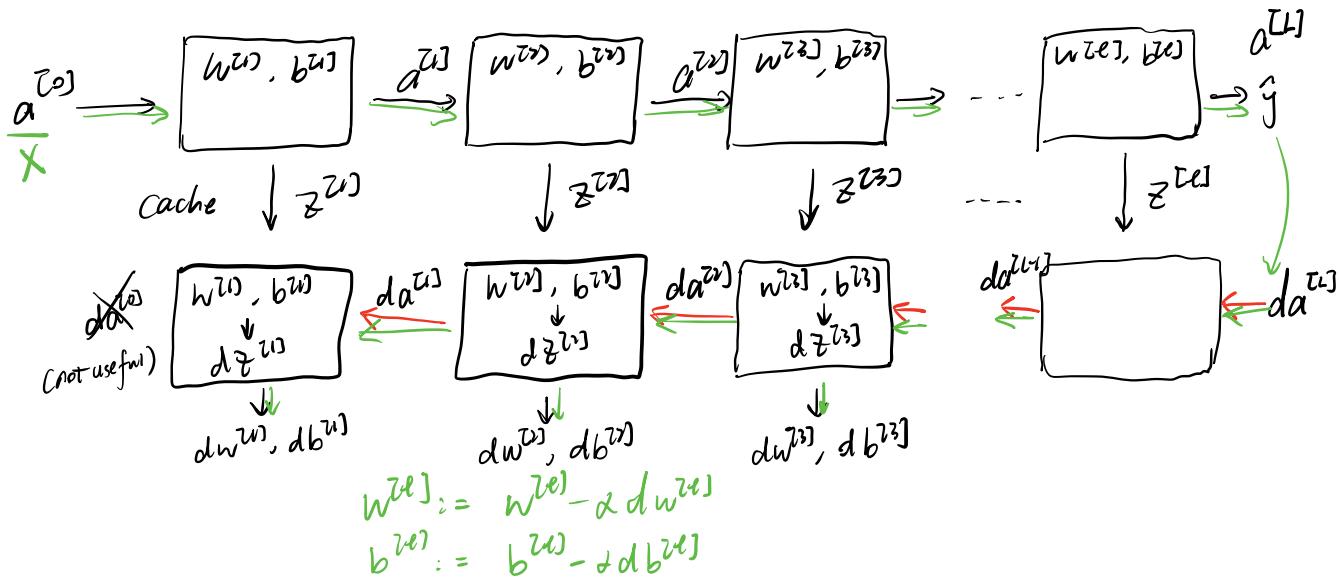
Forward: Input $a^{[l-1]}$, Output $a^{[l]}$

$$z^{[l]} = W^{[l]} \cdot a^{[l-1]} + b^{[l]} \rightarrow \text{cache } Z^{[l]}$$

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

Backward: Input $da^{[L]}$, Output $da^{[L-1]}$

$$z^{[L]} \quad dW^{[L]}, db^{[L]}$$



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

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

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

$$A^{[L-1]} = W^{[L]T} dZ^{[L]}$$

$$dZ^{[L-1]} = A^{[L-1]} * g'([L-1])(Z^{[L-1]})$$

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

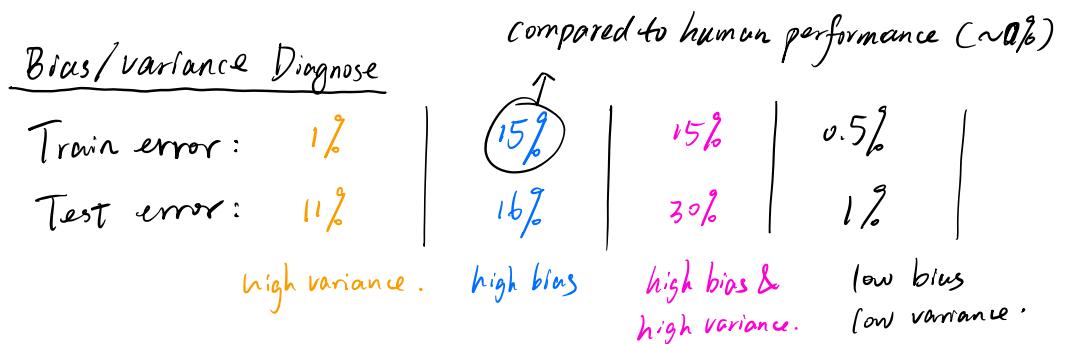
$$dW^{[1]} = \frac{1}{m} dZ^{[1]} A^{[0]T}$$

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

Course 2: Improving DNN - Hyperparameter tuning, Regularization and Optimization

Train / Dev / Test sets

- If sample size is small, could do 70/30 or 60/20/20 split
 - If sample size is very large (e.g., $n > 1M$), could do 99/0.6/0.4 etc.
 - Make sure dev/test sets come from the same distribution.
 - Not having a test set might be okay



* Human error $\approx 0\%$ \Rightarrow compared to train error to decide bias.

Optimal (Bayes) error : 0%

⇒ To diagonalize whether the algorithm's problem comes from bias/variance or both

Basic "Recipe" for machine Learning: (Bias-Variance tradeoff) (not necessarily)

① High bias or not? ————— try: bigger network, train longer.
(training data performance) CNN architecture)

② High variance or not? → try: More data, Regularization,
(CNN architecture)
Dev set performance
↓
Done!

When high variance : Add regularization to the network.

Logistic regression:

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2 + \underbrace{\frac{\lambda}{2m} b^2}_{\text{omit}}$$

(Euclidean norm square)

$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2$$

$$\rightarrow L_2 \text{ regularization } \frac{\lambda}{2m} \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 \quad (\text{Euclidean norm square}) \quad \| \cdot \|_2^2$$

$$L_1 \text{ regularization } \frac{\lambda}{2m} \|w\|_1 = \sum_{j=1}^{n_x} |w_j| \Rightarrow w \text{ will be sparse (most 0s)}$$

λ = regularization parameter, decided by dev set.

Neural Network:

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

$$\text{where } \|w^{[l]}\|^2 = \sum_{i=1}^{n^{[l+1]}} \sum_{j=1}^{n^{[l]}} (w_{ij}^{[l]})^2 \quad \text{given } w: (n^{[0]}, n^{[L+1]})$$

↳ "Frobenius norm" $\| \cdot \|_F$

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

$$w^{[l]} = w^{[l]} - \alpha dJ/dw^{[l]} = \underbrace{w^{[l]} - \frac{\alpha \cdot \lambda}{m} w^{[l]}}_{\text{L2: weight decay.}} - \alpha (\text{from backprop})$$

L2: weight decay.

Dropout regularization (Inverted dropout)

Illustrate with layer $l=3$ ↑ chance of keeping the unit.

$d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \underline{\text{keep_prob}}$ $a_3 = \text{np.multiply}(a_3, d_3)$ $a_3 = a_3 / \text{Keep_prob}$ (scale it up, keep $[a^{(i)}]$ the same.)	$\left\{ \begin{array}{l} d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \underline{\text{keep_prob}} \\ a_3 = \text{np.multiply}(a_3, d_3) \\ a_3 = a_3 / \text{Keep_prob} \quad (\text{scale it up, keep } [a^{(i)}] \text{ the same.}) \end{array} \right.$
--	---

Test time: Not use dropout

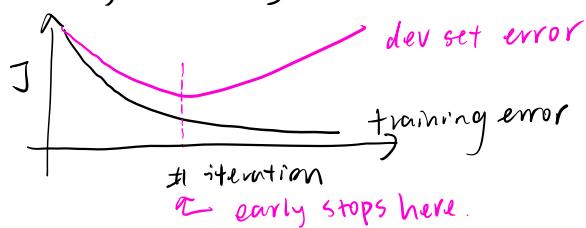
Intuition: Cannot rely on any one feature, so have to spread out weights.

⇒ Shrink weights just like L2

Choosing keep-prob : smaller keep-prob when for w that has many parameters.
 No drop-out (or very close to 1) for input layer.

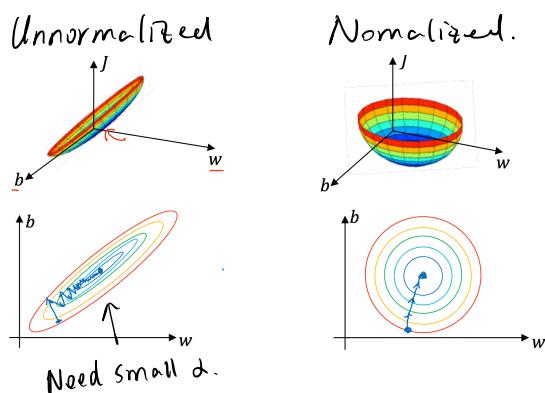
Other regularization methods

- Data augmentation. — creating more training data by modifying available training samples. (Flip horizontally, rotation)
- Early stopping.

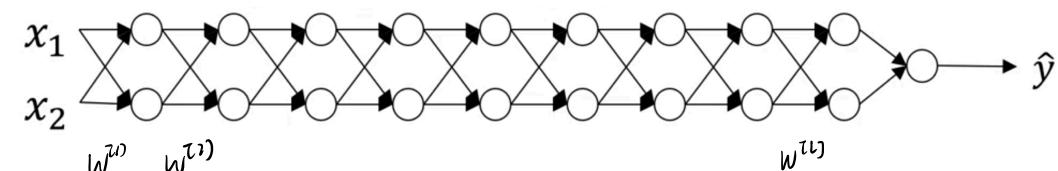


Normalizing Inputs

- Subtract mean with μ
 - Normalize variance with σ^2
- } use the same set of μ & σ^2 for train & test sets.



Vanishing / Exploding gradient



$$\text{Let } g(z) = z, \quad b^{[k]} = 0, \quad n^{[k]}$$

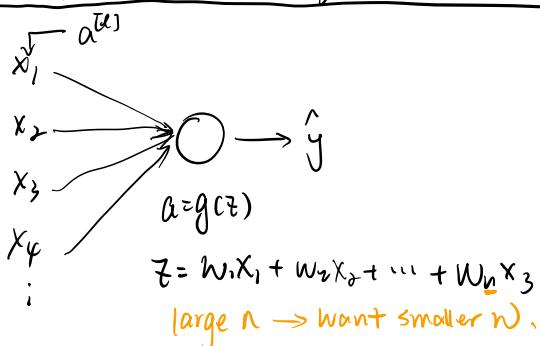
$$L-1 \quad 1.5^{L-1} \rightarrow \infty$$

$$\hat{y} = w^{[L]} w^{[L-1]} \dots \underbrace{w^{[2]} w^{[1]} x}_{z^{[1]} = a^{[1]}}^{[1]} \Rightarrow \hat{y} = w^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} x \quad 0.5^{L-1} \rightarrow 0$$

if $w^{[1]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$

Problem: If the weights are bigger/smaller than identity matrix, the activation could vanish/explode in a deep neural network.

Partial Solution: Weight Initialization



$$\text{Var}(w_i) = \frac{1}{n} \quad (\text{if } g(\cdot) \text{ is relu})$$

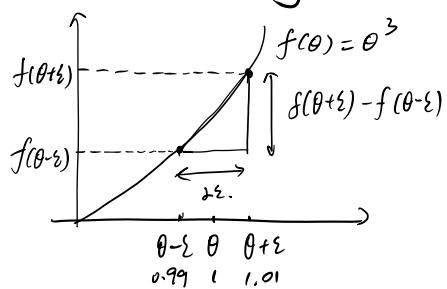
$$w^{[l]} = \text{np.random.randn(shape)} * \boxed{\text{np.sqrt}(\frac{1}{n^{(l-1)}})}$$

equivalent to setting
var of the gaussian dist = 2

Other variants:
tanh $\sqrt{\frac{1}{n^{(l-1)}}}$ Xavier .

$$\sqrt{\frac{2}{n^{(l-1)} + n^{(l)}}}$$

Gradient checking



$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx g(\theta) = 3\theta^2.$$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001 \approx 3. \quad \text{approx error: } 0.0001$$

Take $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ and reshape into a big vector θ .
Concatenate $\Rightarrow J(\theta)$

Take $dW^{[1]}, dB^{[1]}, \dots, dW^{[L]}, dB^{[L]}$ and reshape into a big vector $d\theta$
concatenate \Rightarrow Is $d\theta$ the gradient of $J(\theta)$?

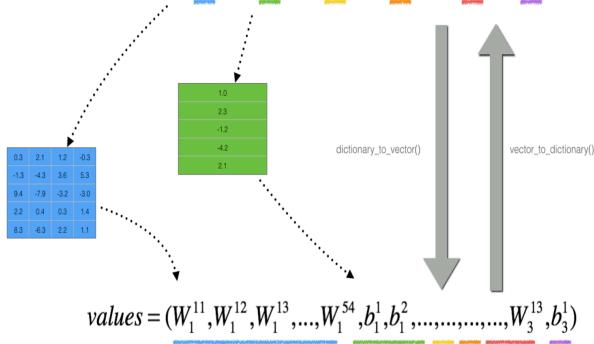
For each i :

$$J_{\theta^{[i]}} = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon} \approx J_{\theta^{[i]}}$$

$$\text{uv approx} = \frac{2\sum_{i=1}^{10^7} \text{uv}}{\text{uv}} = \overline{\delta\theta_i} \quad | \quad \text{uv approx} \approx \text{uv}$$

check $\frac{\|\delta\theta_{\text{approx}} - \delta\theta\|_2}{\|\delta\theta_{\text{approx}}\|_2 + \|\delta\theta\|_2} \approx \begin{cases} 10^{-7} & - \text{great} \\ 10^{-5} & \\ 10^{-3} & - \text{worry} \end{cases}$

`parameters = {"W1": ..., "b1": ..., "W2": ..., "b2": ..., "W3": ..., "b3": ...}`



Batch vs. Mini-batch gradient descent

$$X = \underbrace{[x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(1000)}]}_{(N_x, m)} \mid \underbrace{x^{(1001)} \ \dots \ x^{(2000)}}_{(N_x, 1000)} \mid \dots \mid \underbrace{x^{(m)}}_{(N_x, 1000)}$$

$$Y = \underbrace{[y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(1000)}]}_{(1, m)} \mid \underbrace{y^{(1001)} \ \dots \ y^{(2000)}}_{(1, 1000)} \mid \dots \mid \underbrace{y^{(m)}}_{(1, 1000)}$$

if $m = 5,000,000$, then 5000 minibatches \rightarrow minibatch t : $X^{(t)}, Y^{(t)}$

for epoch = ... {

for $t = 1, \dots, 5000$ {

forwardProp on $X^{(t)}$

$$\begin{aligned} Z^{(l)} &= W^{(l)} X^{(l)} + b^{(l)} \\ A^{(l)} &= g^{(l)}(Z^{(l)}) \end{aligned} \quad \left. \begin{array}{l} \text{Vectorized implementation} \\ (\text{1000 example}) \end{array} \right.$$

$$A^{(l)} = g^{(l)}(Z^{(l)}) \quad \text{Indicate # of sample in minibatch}$$

$$\text{Compute cost } J = \frac{1}{1000} \sum_{i=1}^{1000} \ell(g^{(l)}, y^{(i)}) + \frac{\lambda}{2000} \sum_l \|W^{(l)}\|_F^2$$

cost of this minibatch

from minibatch $X^{(t)}, Y^{(t)}$ # of layers.

backprop to compute gradients wrt. $J^{(t)}$ (using $X^{(t)}, Y^{(t)}$)

One epoch.

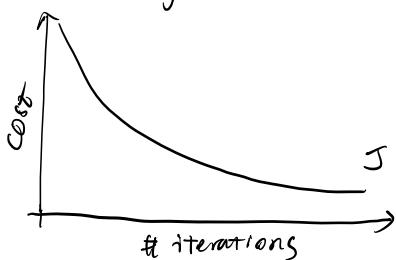
$$w^{[t+1]} := w^{[t]} - \alpha \nabla w^{[t]} ; \quad b^{[t+1]} := b^{[t]} - \alpha \nabla b^{[t]}$$

}

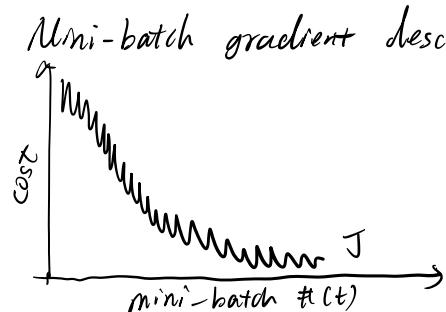


}

Batch gradient descent



Mini-batch gradient descent



Choosing mini-batch size:

- If mini-batch size = m : batch gradient descent $(X^{[t]}, Y^{[t]}) = (X, Y)$ Takes too long per iteration
- If mini-batch size = 1: Stochastic gradient descent $(X^{[t]}, Y^{[t]}) = (x^{(1)}, x^{(2)}) \dots$ lose speedup from vectorization

\Rightarrow In between: Fastest learning \rightarrow Vectorization speed up.

\Rightarrow Make progress without processing entire set

\Rightarrow Won't always converge, would bounce around a small region, could be solved by lr decay.

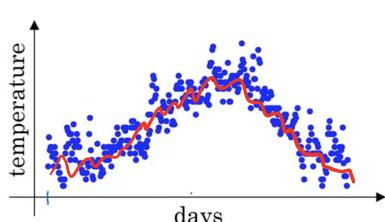
If $(m \leq 2000)$ small training set: batch gradient descent.

If big training set: minibatch sizes: 64-512 (better be 2^n)

Make sure $X^{[t]}, Y^{[t]}$ fits CPU/GPU memory.

★ High β : More weights on previous data points (Averaged across more previous data points). Line will become smoother.

Exponentially weighted averages



$$V_0 = 0$$

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

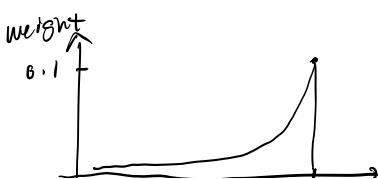
$$V_2 = 0.9V_1 + 0.1\theta_2$$

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

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

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

V_t is approximately averaging over $\approx \frac{1}{1-\beta}$ days' temperature.



$$V_\theta = 0$$

for $t \in 0 \dots 100 \{$

$$V_\theta := \beta V_\theta + (1-\beta)\theta_t$$

bias
correction \Rightarrow

$$V_\theta = 0$$

for $t \in 0 \dots 100 \{$

$$V_\theta := \beta V_\theta + (1-\beta)\theta_t$$

$$V_\theta := \frac{V_\theta}{n}$$

$t=100$

$3 \quad 1-\beta^t$

Gradient descent with momentum

Momentum:

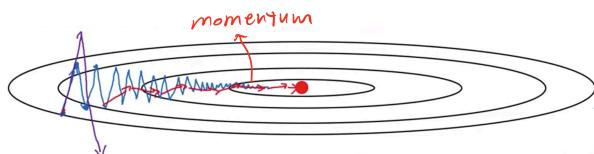
On iteration t :

Compute dw, db on current mini-batch

$$V_{dw} = \beta V_{dw} + (1-\beta) dw \quad \Rightarrow \text{moving average of } dw \text{ for each mini-batch}$$

$$V_{db} = \beta V_{db} + (1-\beta) db \quad \text{acceleration} \quad \beta = 0.9$$

$$w := w - \alpha V_{dw}; \quad b := b - \alpha V_{db}$$



RMSprop: (root-mean-square prop)

On iteration t :

Compute dw, db on the current mini-batch

$$S_{dw} = \beta S_{dw} + (1-\beta) dw^2 \quad \text{element-wise} \quad * \text{Decay learning rate faster for steeper dimensions than flatter ones.}$$

$$S_{db} = \beta S_{db} + (1-\beta) db^2$$

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw} + \epsilon}}; \quad b := b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}}$$

* Minimize oscillation
* Faster learning with diverging.

Adam: Adaptive momentum

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

On iteration t :

Compute dw, db using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dw; \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) db \quad \leftarrow \text{momentum}$$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2; \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \quad \leftarrow \text{RMSprop}$$

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

$$S_{dw}^{\text{corrected}} = \frac{\dot{S}_{dw}}{1-\beta_2^t}; S_{db}^{\text{corrected}} = \frac{\dot{S}_{db}}{1-\beta_2^t}$$

$$w := w - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon} ; b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

Hyperparameter choices:

α : Needs to be tuned

β_1 : Moving weighted average of dw/db 0.9

β_2 : Moving weighted average of d^2w/d^2b 0.999

$\epsilon: 10^{-8}$

Learning Rate Decay (Separate from Adam)

$$\alpha = \frac{1}{1 + \text{decayRate} * \text{epoch_Num}} \cdot \alpha_0$$

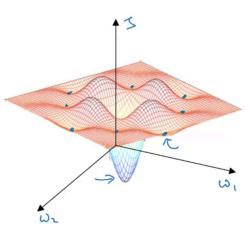
Epoch	α
1	0.1
2	0.06
3	0.05
4	0.04

if $\alpha_0 = 0.2$

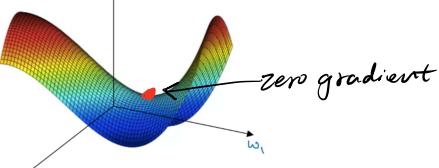
$$\alpha = 0.95^{\text{epoch_num}} \cdot \alpha_0 \quad \text{or} \quad \alpha = \frac{1}{\text{epoch_Num}} \cdot \alpha_0$$

Local Optima in neural networks

used to think:



Local optima in high dimension NN
are usually saddle points



⇒ Unlikely to get stuck in a bad local optima.

⇒ Plateaus can make learning slow (Where Adam can help)

Hyperparameter tuning process

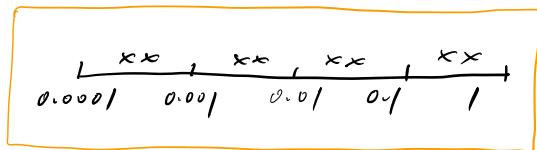
$\alpha, \beta_1, \beta_2, \gamma, \# \text{ layers}, \# \text{ hidden units}, \text{learning rate decay}, \text{mini-batch size} \dots$
 ~ 0.9

→ Try random values: Don't use a grid.

→ Coarse to fine sampling scheme.

Appropriate scale for hyperparameters: log scale

e.g. $\alpha = 0.0001 \dots 1$



$$r = -4 * np.random.rand() \leftarrow r \in [-4, 0] \leftarrow \text{sample } r \text{ uniformly}$$

$$\alpha = 10^r \leftarrow 10^{-4} \dots 10^0$$

Hyperparameters for exponentially weighted averages:

$$\beta = 0.9 \dots 0.999 \Rightarrow r \in [-3, -1] \quad 1-\beta = 10^r \quad \beta = 1-10^r$$

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

Batch Normalization

Recall: Normalizing inputs to speed up learning.

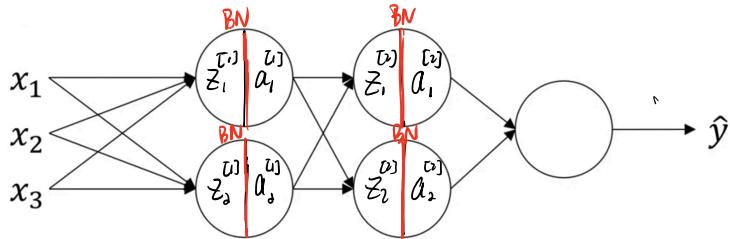
→ Can we normalize $a^{[l]}$ or $z^{[l+1]}$ so as to train $w^{[l+1]}, b^{[l+1]}$ faster?

Given some intermediate values in NN, $z^{[l+1]}, \dots, z^{[L+1]}$

$$\mu = \frac{1}{m} \sum_i z^{[l+1]} ; \sigma^2 = \frac{1}{m} \sum_i (z^{[l+1]} - \mu)^2$$

$$z_{\text{norm}}^{[l+1]} = \frac{z^{[l+1]} - \mu}{\sqrt{\sigma^2 + \epsilon}} ; \hat{z}^{[l+1]} = \gamma z_{\text{norm}}^{[l+1]} + \beta$$

learnable parameters of model.



$$x \xrightarrow{W^{[l]}, b^{[l]}} z^{[l]} \xrightarrow{\gamma^{[l]}, \beta^{[l]}} \hat{z}^{[l]} \rightarrow a^{[l]} = g^{[l]}(\hat{z}^{[l]}) \xrightarrow{W^{[l+1]}, b^{[l+1]}} z^{[l+1]} \xrightarrow{\gamma^{[l+1]}, \beta^{[l+1]}} \hat{z}^{[l+1]} \rightarrow a^{[l+1]}$$

will center at 0 before BN either way, not useful if using BN.

parameters: $W^{[l]}, b^{[l]}, W^{[l+1]}, b^{[l+1]}, \dots, W^{[L]}, b^{[L]} \}$ $\beta^{[l]} = \beta^{[l]} - \alpha d\beta^{[l]}$
 $\gamma^{[l]}, \beta^{[l]}, \gamma^{[l+1]}, \beta^{[l+1]}, \dots, \gamma^{[L]}, \beta^{[L]} \}$
 $(n^{[l]}, 1) \quad (n^{[l]}, 1)$

Implementing gradient Descent for Batch Normalization:

for $t=1 \dots \text{num-mini-batches}$:

Compute forward prop on x^{it} have normalized mean/variance

In each hidden layer, use BN to replace $\hat{z}^{[l]}$ with $\tilde{z}^{[l]}$

Use backprop to compute $dW^{[l]}$, $d\beta^{[l]}$, $d\gamma^{[l]}$

update parameters $W^{[l]} = W^{[l]} - \alpha dW^{[l]} \}$ Can use other optimizers
 $\beta^{[l]} = \beta^{[l]} - \alpha d\beta^{[l]} \}$ as well.
 $\gamma^{[l]} = \gamma^{[l]} - \alpha d\gamma^{[l]} \}$

Why does batch normalization work?

→ Learning on shifting input distribution

→ Reduces the amount that the distribution of the hidden units' activation changes.

→ Mean/Variance within each mini-batch add some noise to the value $\hat{z}^{[l]}$ and $\tilde{z}^{[l]}$ → regularization effect.

Batch Norm at test time

Estimate μ , σ^2 using exponentially weighted average (across mini-batches)

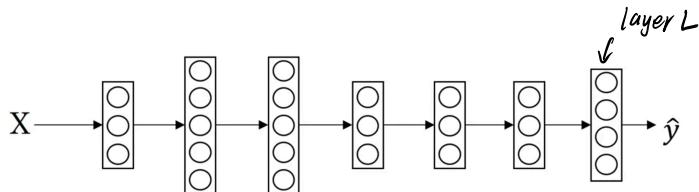
$x^{[1]}, x^{[2]}, x^{[3]}, \dots$ | At test (with each testing sample).

$$\begin{array}{c} \downarrow \quad \downarrow \quad \downarrow \\ \mu^{[1][t]} \quad \mu^{[2][t]} \quad \mu^{[3][t]} \\ \sigma^{[1][t]} \quad \sigma^{[2][t]} \quad \sigma^{[3][t]} \end{array} \rightarrow \begin{array}{c} \mu \\ \sigma^2 \end{array}$$

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

$$\hat{z} = f_{\text{norm}} + \beta$$

Multi-class classification (Softmax layer)



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

$$\text{Softmax activation fx: } t = e^{(z^{[L]})} ; a^{[L]} = \frac{e^{(z^{[L]})}}{\sum_{i=1}^c t_i} ; a_i^{[L]} = \frac{t_i}{\sum_{j=1}^c t_j}$$

$$\text{e.g. } z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix} \quad \sum_{i=1}^c t_i = 176.3 \quad a^{[L]} = \frac{t}{176.3} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} = \hat{y}$$

Loss function:

$$f(\hat{y}, y) = - \sum_{j=1}^m y_j \log \hat{y}_j \Rightarrow \text{Make } \hat{y}_j \text{ as big as possible.}$$

$$J = \frac{1}{m} \sum_{j=1}^m f(\hat{y}, y)$$

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

$$= \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad \text{One hot}$$

(4, m) if 4 classes in total

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

$$= \begin{bmatrix} 0.3 \\ 0.12 \\ 0.1 \\ 0.4 \end{bmatrix} \quad \dots$$

(4, m)

$$\text{Backprop: } d\hat{z}^{[L]} = \hat{y} - y$$

Class 3: Structuring Machine Learning project

Chain of assumptions in ML:

- ① Fit training set well on cost function (\approx human level performance)
- ↓
② Fit Dev set well on cost function \Rightarrow regularization
bigger training set
- ↓
③ Fit test set well on cost function \Rightarrow bigger dev set
- ↓
④ Performs well in real world. \Rightarrow change dev set
change cost function

bigger network
optimization (e.g. Adam)

→ Using a single number evaluation metric

- Instead of using precision & recall \rightarrow use F1 score = $\frac{2}{P+R}$ "harmonic mean"
- Choosing the "best" model from a bunch. \rightarrow Speed up iterating

→ Satisficing and Optimizing metric

e.g. classifier	Optimizing	Satisficing.	metric:
	↑	↑	
A	Accuracy	running time	Maximize Accuracy subject to
B	90%	80ms	running time \leq 100ms
C	92%	95ms	
	95%	1500ms	

- If N metrics : 1 optimizing ; N-1 satisfying.

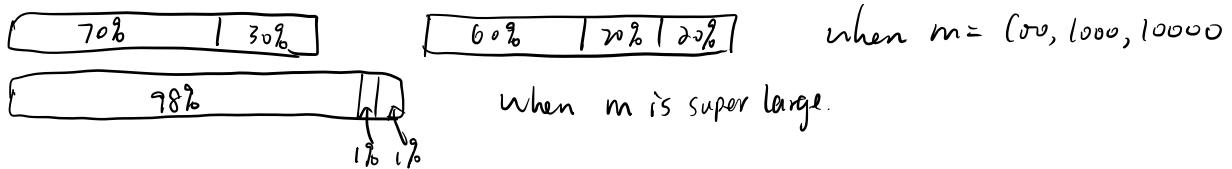
→ Train/dev/test distributions

- Make dev & test set from the same distribution

e.g., Regions: US ; UK ; Other Europe ; South America ; India ; China ; Other Asia
randomly shuffle into dev/test

→ Size of dev and test sets

- Old way of splitting data



→ When to change Dev/Test sets and metrics?

- If doing well on your metric + dev/test set does not correspond to doing well on your application, change your metric and/or dev/test sets.

Summary of bias/variance with human-level performance.

- O%
- Human level error (proxy for bayes error)
- Training error
- Dev error

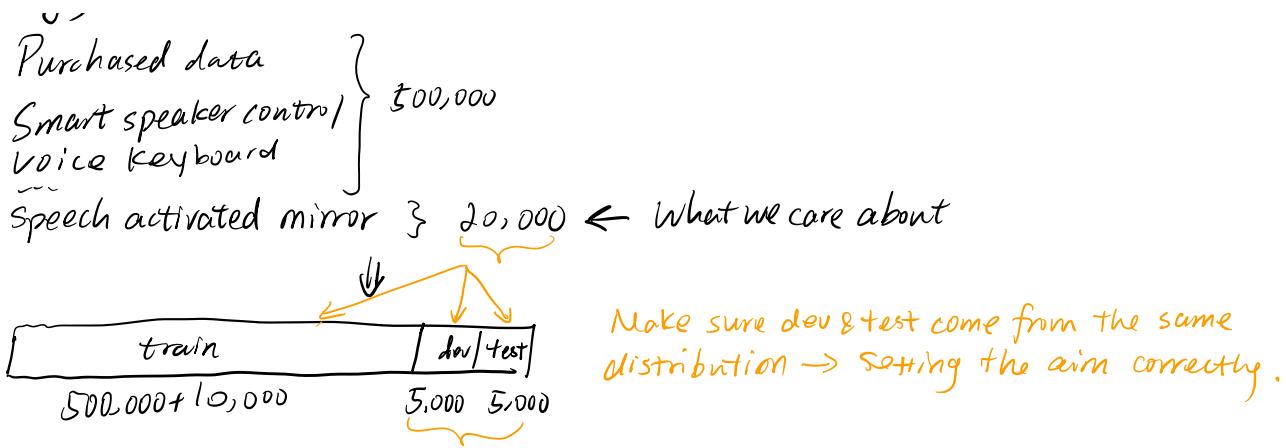


The 2 fundamental assumptions of supervised learning.

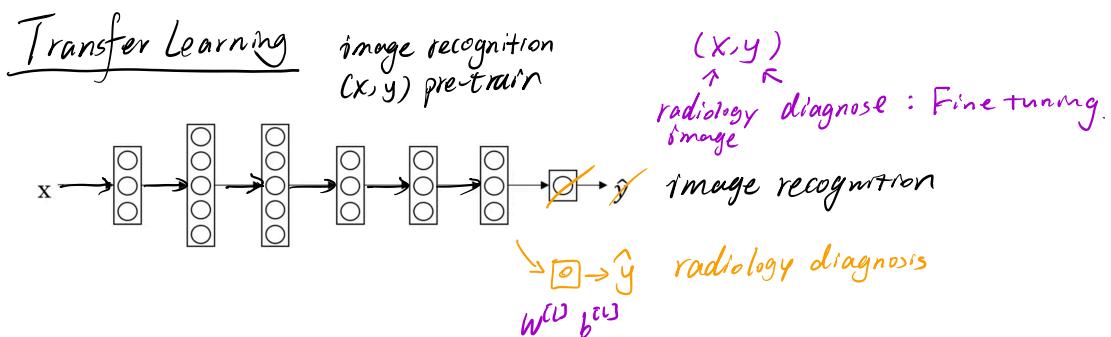
1. You can fit the training set pretty well. ~ low avoidable bias
 - Train bigger model
 - Train bigger/better optimization algorithms (momentum; RMSprop ...)
 - NN architecture/hyperparameters search.
2. The training set performance generalizes pretty well to the dev/test set ⇒ More training data / regularization (Dropout; L2)
 - More data
 - Regularization (L2; dropout; data augmentation)
 - NN architecture/hyperparameters search

When training & testing data do not come from the same distribution:

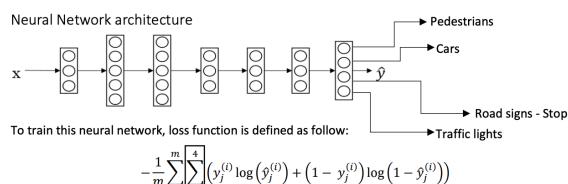
e.g., Speech recognition problem:



Bias/Variance on mismatched training and dev/test sets



Multiple task learning

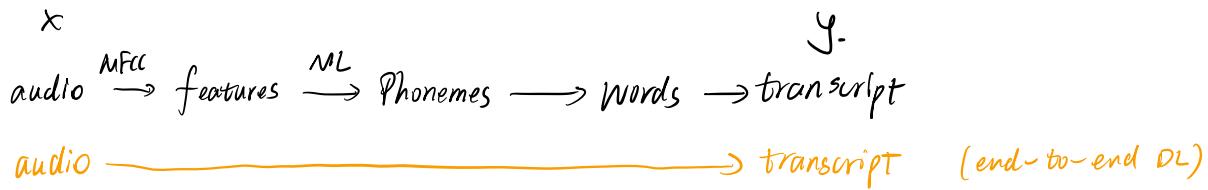


$$-\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^4 \left[y_j^{(i)} \log(\hat{y}_j^{(i)}) + (1 - y_j^{(i)}) \log(1 - \hat{y}_j^{(i)}) \right]$$

Also, the cost can be compute such as it is not influenced by the fact that some entries are not labeled.
Example:

$$Y = \begin{bmatrix} 1 & 0 & ? & ? \\ 0 & 1 & ? & 0 \\ 0 & 1 & ? & 1 \\ ? & 0 & 1 & 0 \end{bmatrix}$$

End-to-end learning



- Pros: Let data speak
Less hand-designing of components needed

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

Class 4: Convolution neural Network

Conv
 Pool
 FC

Computer vision problems: Image classification, Object detection, Neural style transfer

Convolution: Vertical edge detection

$\rightarrow 3 \times 1 + 1 \times 1 - 2 \times 1 + 0 \times 0 + 5 \times 0 + 7 \times 0 + 1 \times 1 + 8 \times -1 + 2 \times -1 = -5$

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

n $\xrightarrow{5 \times 6}$
 convolution
 $\xrightarrow{*}$
 $f \xrightarrow{3 \times 3}$
 → filter kernel
 $n-f+1=4$
 $\xrightarrow{4 \times 4}$
 Sobel filter (vertical)

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
10	10	10	0	0	0

\times
 $\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$
 =
 $\begin{bmatrix} 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \end{bmatrix}$

→ back prop to learn filters / kernels.

Padding:

To not ① Shrink output

② Throw away info from edges

"Valid" convolution: No padding. $n \times n \times f \times f = n-f+1 \times n-f+1$

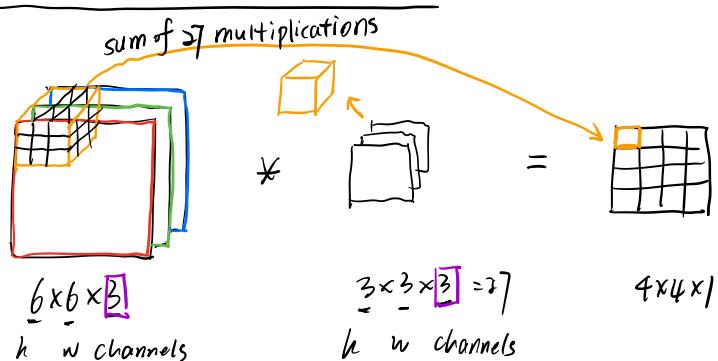
"Same" convolution : Pad so that output size is the same as the input size.
 $P = \frac{f-1}{2}$ (f is usually odd)

Strided Convolution:

$$n \times n * f \times f \xrightarrow{\text{w/ Padding } P \text{ & Stride } S} \left\lfloor \frac{n+2P-f}{S} + 1 \right\rfloor \times \left\lfloor \frac{n+2P-f}{S} + 1 \right\rfloor$$

$\lfloor z \rfloor = \text{floor}(z)$

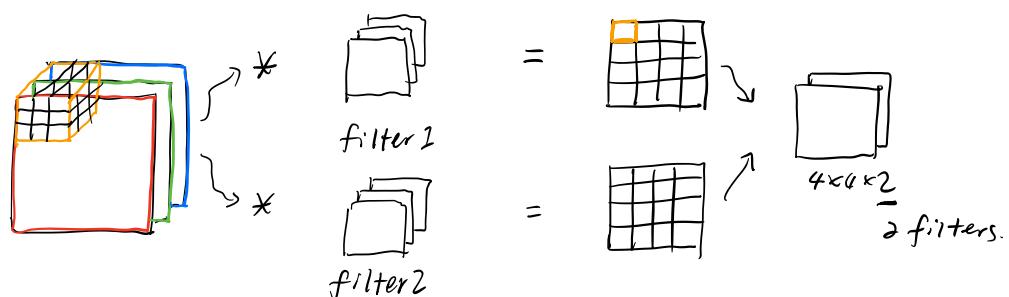
Convolutions over volumes:



filter examples :

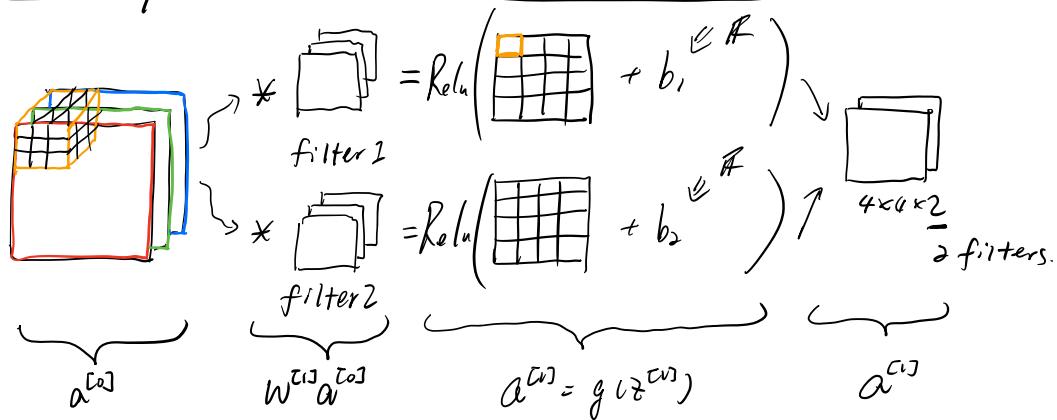
$$R \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} G \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} B \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \rightarrow 3 \times 3 \times 3 : \text{A filter that detects edges in the Red channel of the image.}$$

$$R \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} G \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} B \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} : \text{A filter that detects edges in any color}$$



$$\text{Summary: } n \times n \times \underline{n_c} * f \times f \times \underline{n_c} \xrightarrow{\text{channel/Depth}} n_f \times n_f \times \underline{n_c'} \xrightarrow{\text{channel/Depth}} \# \text{ of filters.}$$

One layer of Convolutional Neural Network.

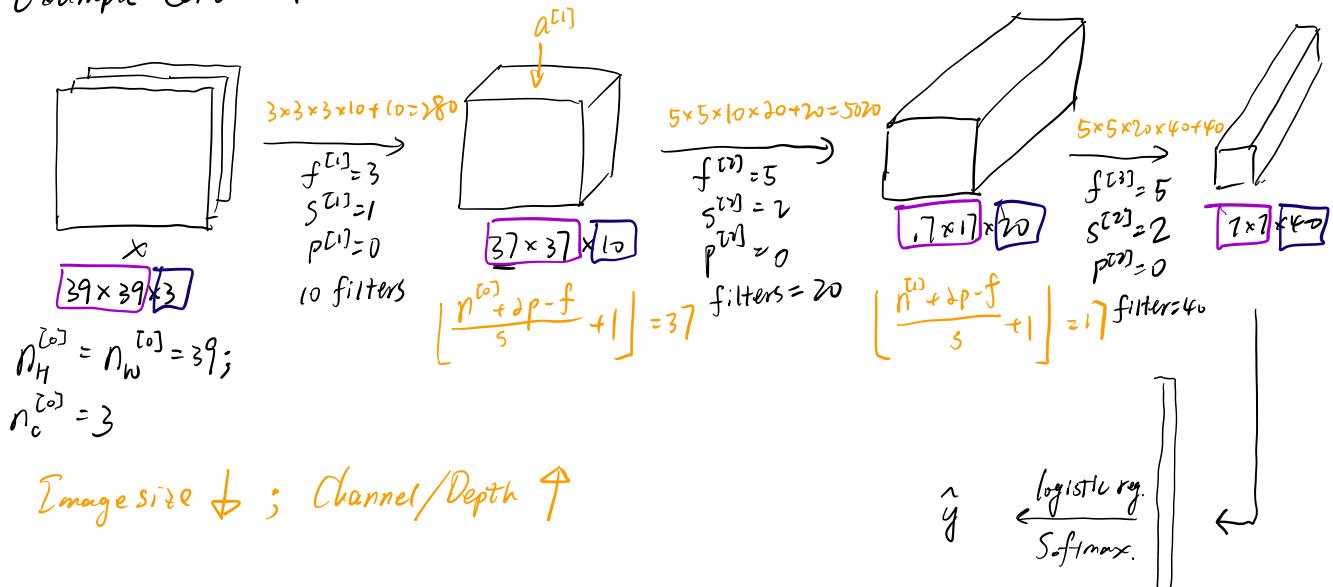


- If you have 10 filters that are $3 \times 3 \times 3$ in one layer of a CNN, # of parameters? $(3 \times 3 \times 3 + 1) \times 10 = 280$.

Notations: If layer l is a convolution layer:

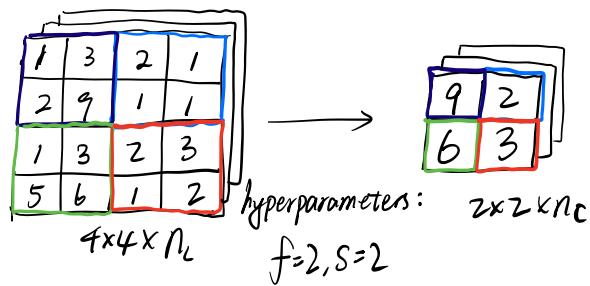
$f^{[l]}$ = filter size	Each filter is: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$	Input = $n_H^{[l-1]} \times n_W^{[l-1]} \times n_c^{[l-1]}$
$p^{[l]}$ = padding	Activations: $a^{[l-1]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$	Output = $n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$
$s^{[l]}$ = stride	$A^{[l-1]} \rightarrow M \times n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$	$n_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$
$n_c^{[l]}$ = number of filters	Weights: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$	# of filters in layer l
	bias: $1 \times 1 \times 1 \times n_c^{[l]}$	

Example Conv Net:

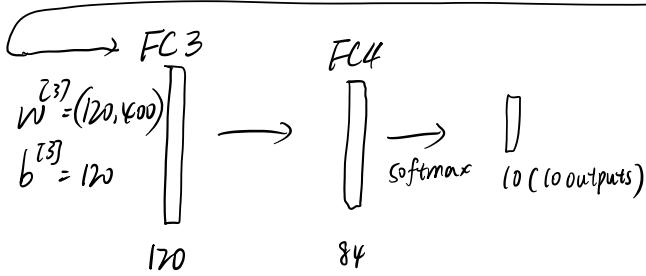
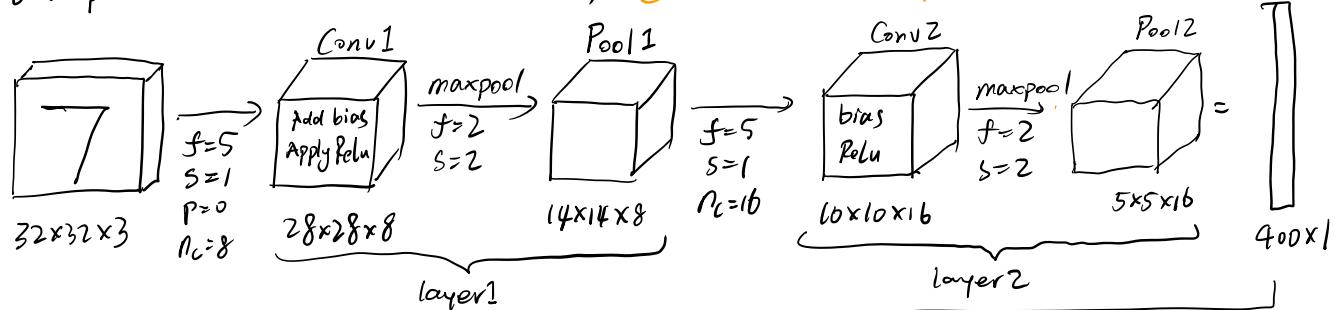


Pooling layer: Max pooling.

vectorize 1960



Example Conv. Net: (LeNet-5-ish) Conv - pool - Conv - pool - Fc - Fc - Softmax.



Params:

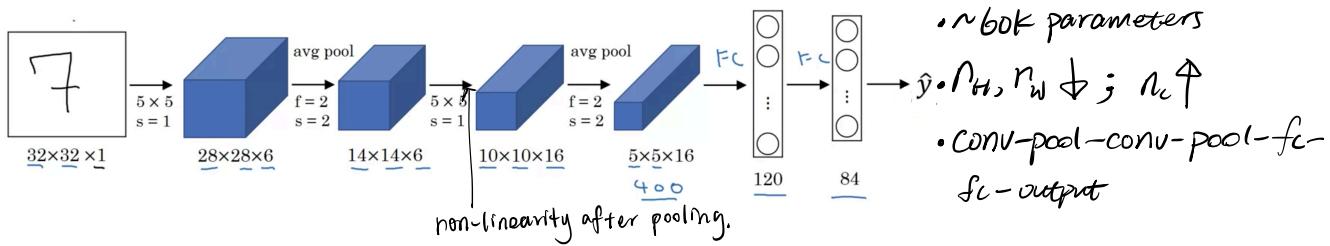
layer 1:	softmax: $10 \times 84 + 10 = 850$
$5 \times 5 \times 3 \times 8 + 8 = 608$	
layer 2:	$5 \times 5 \times 8 \times 16 + 16 = 3216$
layer 3:	$120 \times 400 + 120 = 48120$
layer k:	$84 \times 120 + 84 = 10160$

Why Convolutions:

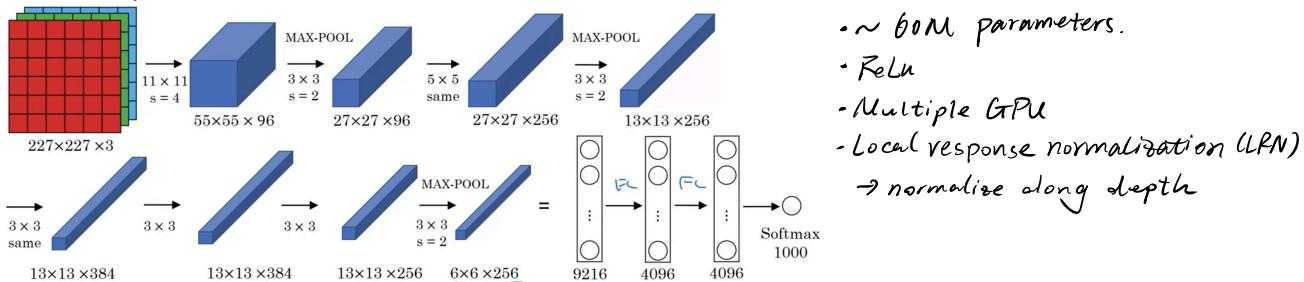
1. **Parameter Sharing**: A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.
2. **Sparsity of Connections**: In each layer, each output value depends only on a small number of inputs (prevent overfitting) \rightarrow translation invariance

CNN architectures:

LeNet-5 (Document recognition on grayscale image)

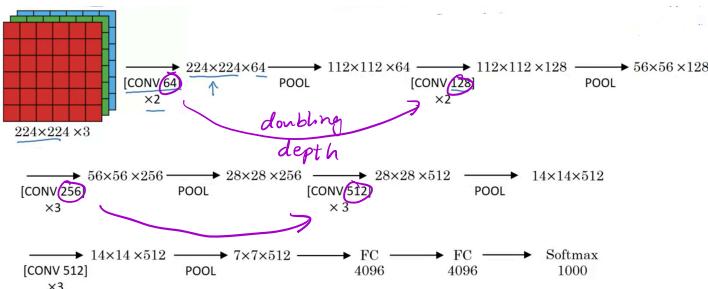


AlexNet

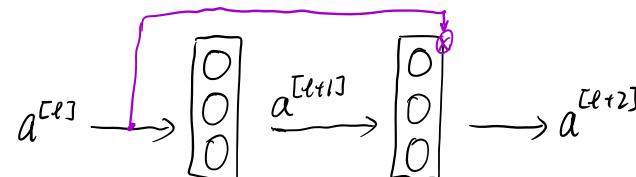


VGG-16 : Conv = 3×3 filter, $s=1$, padding = "same"

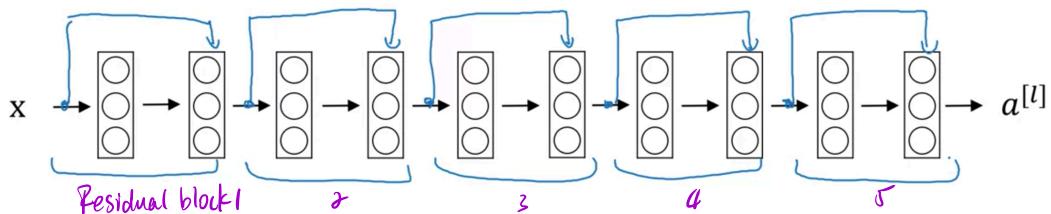
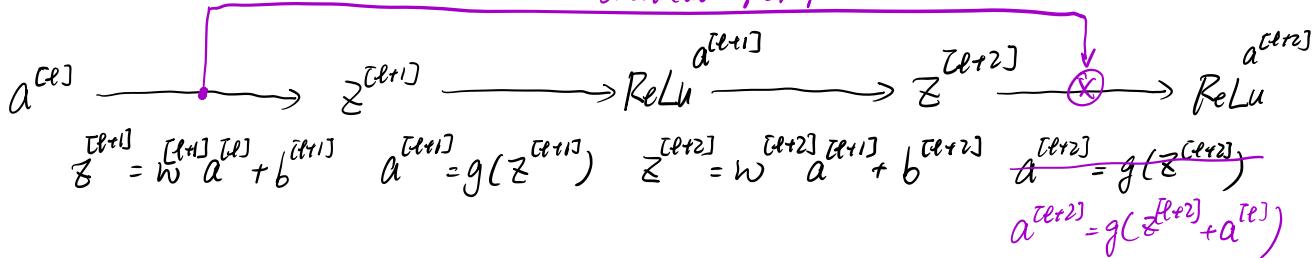
Maxpool = 2×2 , $s=2$

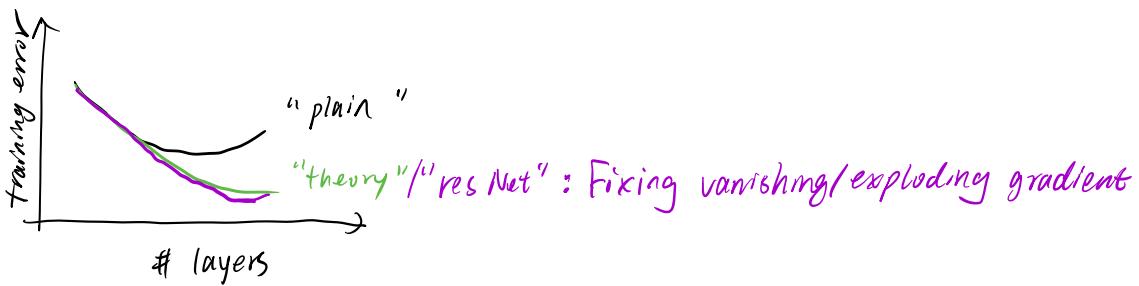


Residual Network (skip connections)



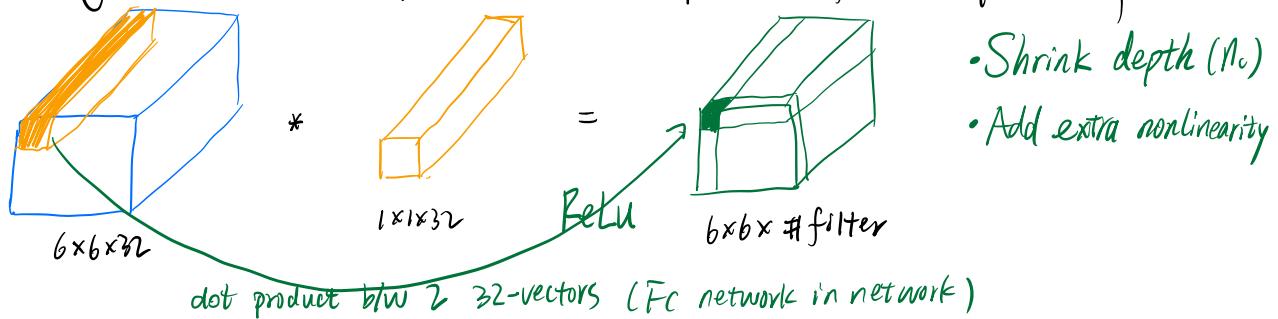
"short cut" / skip connection





Networks in networks and 1×1 convolutions

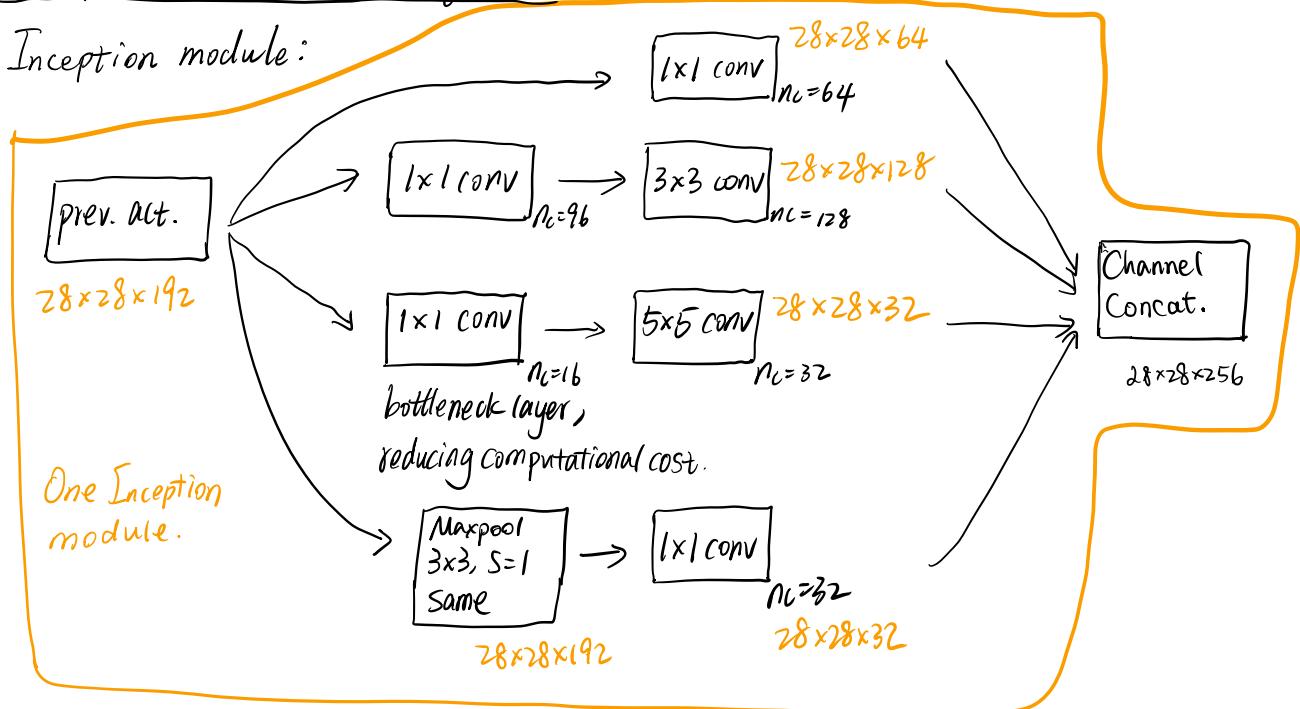
- Make sense when depth > 1
- Having FC networks applied to each position of the input image



- Shrink depth (n_c)
- Add extra nonlinearity

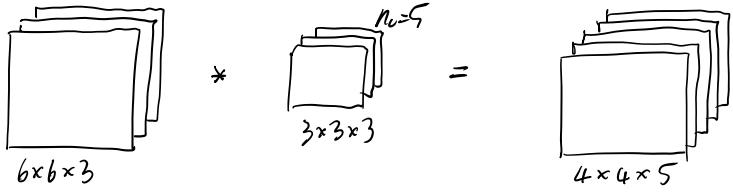
Inception Conv. Net: GoogleNet

Inception module:



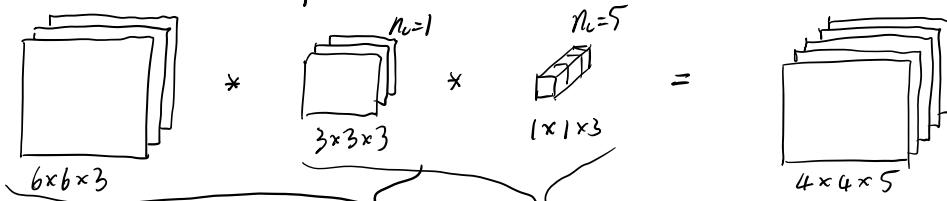
MobileNet

Key idea: Normal vs. depthwise-separable convolutions.



Normal convolution = # of filter parameters \times # of filter positions $\times n_c$
computation cost

$$= 27 \times 16 \times 5 = 2160$$



depthwise-separable conv = $\underline{27 \times 16} + \underline{3 \times 16 \times 5} = 672$
computation cost

Object classification with localization

Define the target label y : Need to output b_x, b_y, b_h, b_w , class label (1-4)

1. pedestrian
2. car
3. motorcycle
4. background

$$y = \begin{bmatrix} P_e \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \rightarrow \begin{array}{l} \text{Is there an object?} \\ 1: \text{If there is } 0: \text{background.} \end{array}$$

} which object

$$\ell(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + (\hat{y}_3 - y_3)^2 & \text{if } y_1 = 1 \\ \dots + (\hat{y}_8 - y_8)^2 & \\ (\hat{y}_1 - y_1)^2 & \text{if } y_1 = 0 \end{cases}$$

Object detection : Could have multiple objects in one image

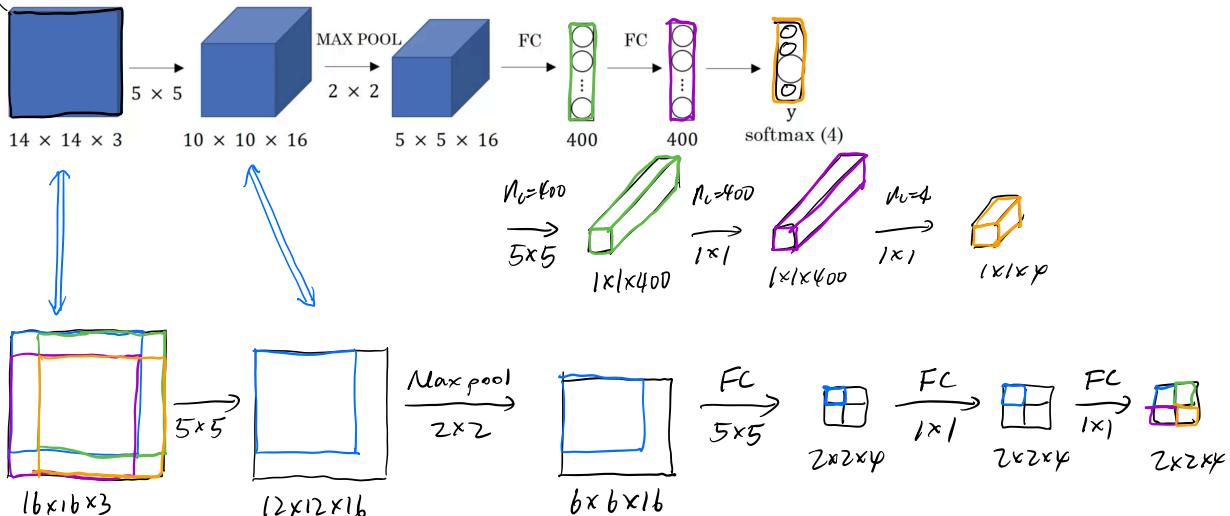
\Rightarrow Sliding window method

↙ closely cropped image



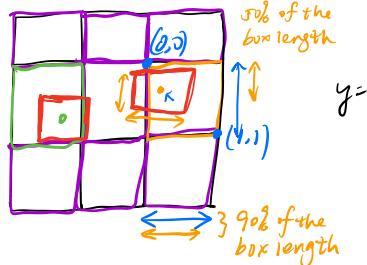
Convolutional Implementation of sliding window

Turning FC layer into convolutional layers



Bounding box predictions

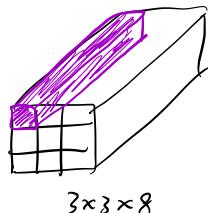
YOLO algorithm



labels for training
for each grid cell

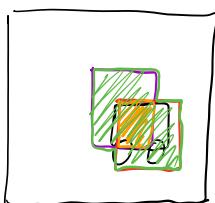
$$y = \begin{bmatrix} y \\ bx \\ by \\ bh \\ bw \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

target output volume



object assigned based on
mid-point location

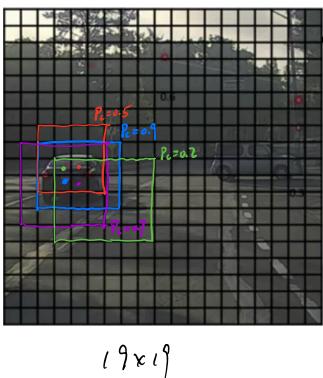
Intersection over union : Evaluating object localization



$$\text{Intersection over union} = \frac{\text{Size of } \text{[orange box]}}{\text{Size of } \text{[green box]}}$$

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

Non-max suppression: Make the algorithm detect an object only once.



Each output prediction is:

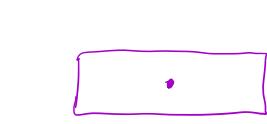
$$\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$$

- Discard all boxes with $P_c \leq 0.6$
- While there are any remaining boxes:
 - Pick the box with the largest P_c , output that as pred.
 - Discard any remaining box with $\text{IoU} \geq 0.5$ with the box output in the previous step.

Anchor boxes: Letting 1 grid cell to detect multiple objects.



Anchor box 1: Anchor box 2:



$$y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ P_c \\ b_x \\ \vdots \end{bmatrix} \left\{ \begin{array}{l} \text{anchor box 1 (pedestrian)} \\ \text{anchor box 2 (car)} \end{array} \right.$$

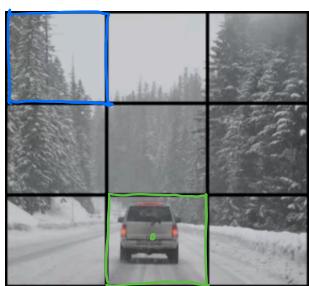
With two anchor boxes: Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU

YOLO Algorithm

Training: (3 class labels)

if using 2 anchor boxes: $y.\text{shape} = (3, 3, 2, 8)$ or $(3, 3, 16)$

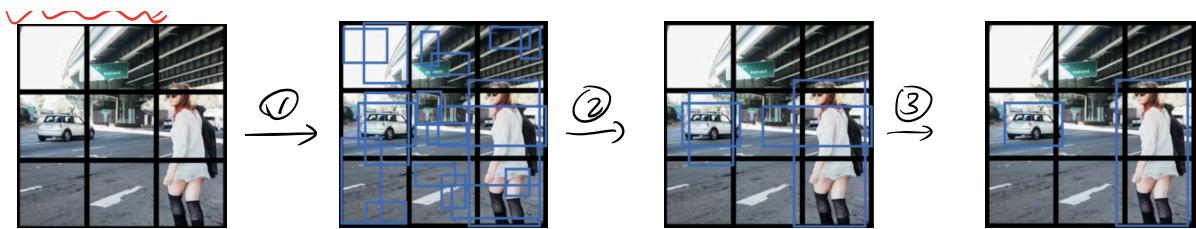
\downarrow
3x3 grid cell PC, box, 3 classes



$$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ 0 \\ ? \\ ? \\ ? \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ 1 \\ b_x \\ b_y \\ b_w \\ b_h \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Predicting



- ① For each grid cell, get 2 predicted bounding boxes, 1 for each anchor.
- ② Get rid of low P_c boxes
- ③ For each class (Pedestrian, car, motorcycle) use non-max suppression to generate final pred.

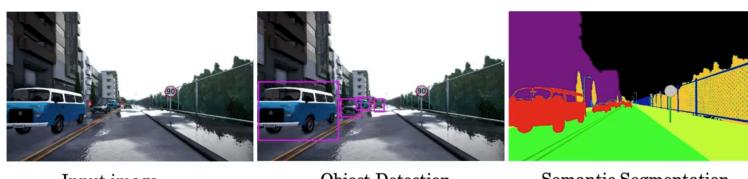
R-CNN: Region proposal CNN.

R-CNN: Propose regions. Classify proposed regions one at a time.
Output label + bounding box

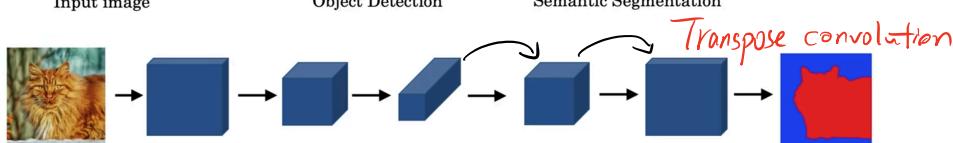
Fast R-CNN: Propose regions. Use convolution implementation of sliding windows to classify all the proposed regions

Faster R-CNN: Use convolutional network to propose regions.

Semantic segmentation: label every pixel

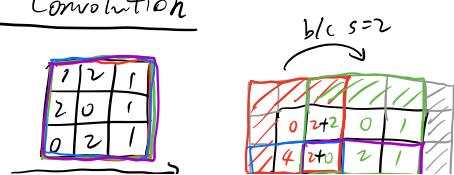


Motivation: locating tumor in a MR scan



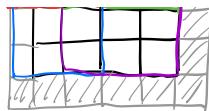
→ Instead of getting smaller and smaller, width and height go back to normal and P_c decreases to # class

Transpose Convolution

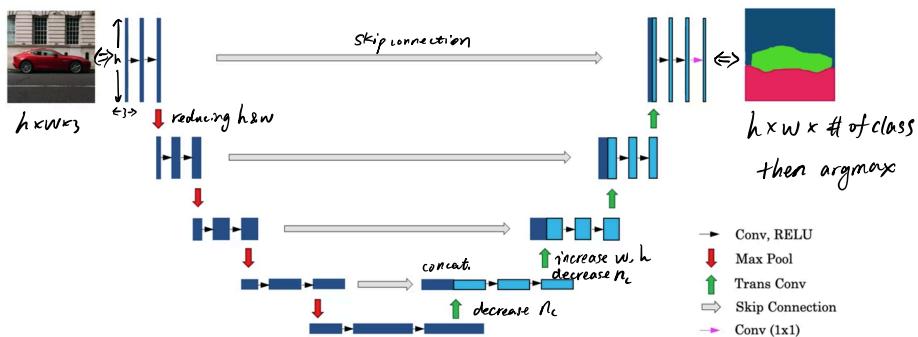


Filter $f \times f = 3 \times 3$
padding $P = 1$
 $s = 2$

32



U-Net Architecture



Face Verification vs. Face recognition

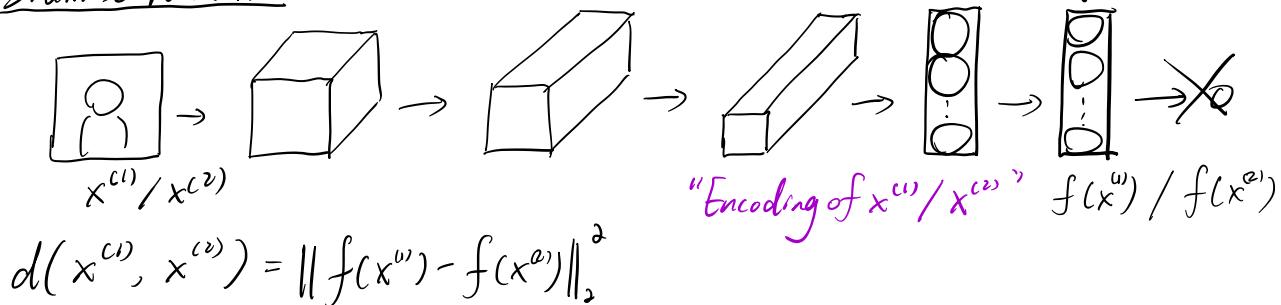
verification: Whether the ID matches the face.

recognition: Output ID if the face is any of the K person in the database.

One-Shot learning: Learning from 1 example to recognize the person again.

→ Learning a "similarity" function: $d(\text{img1}, \text{img2}) = \text{degree of difference}$ b/w images
 if $d(\text{img1}, \text{img2}) \leq T$ same person
 if $d(\text{img1}, \text{img2}) > T$ different person } verification

Siamese Network:



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

• Learn parameters so that: if $x^{(i)}, x^{(j)}$ are the same person, $\|f(x^{(i)}) - f(x^{(j)})\|^2$ is small

Objective function: Triplet loss (anchor, positive, negative images)

$$\text{Goal: } \underbrace{\|f(A) - f(P)\|^2}_{d(A,P)} + \underbrace{\|f(A) - f(N)\|^2}_{d(A,N)} \leq 0$$

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha \leq 0 \quad (\text{preventing the NN from encoding everything the same})$$

Given 3 images: A, P, N

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

$$J = \sum_{i=1}^m L(A^{(i)}, P^{(i)}, N^{(i)}) \leq 0$$

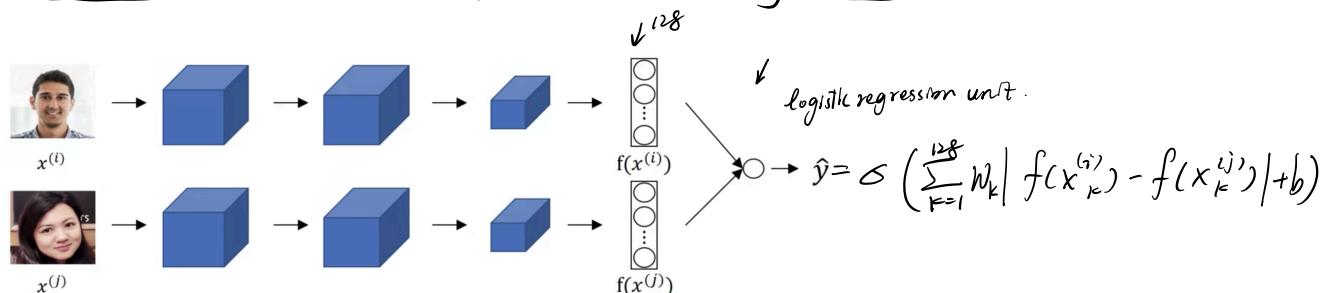
E.g., Training set: $10k$ pictures of $1k$ persons (i.e., need multiple images of a person)

\Rightarrow The trained parameters might be able to solve the one-shot learning problem.

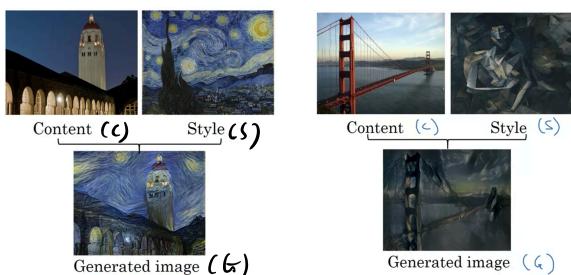
How to choose triplets: If A, P, N are randomly chosen, then

$$d(A, P) + \alpha \leq d(A, N) \text{ is easily satisfied.}$$

Another Approach: Face Verification and Binary classification



Neural Style Transfer



Visualizing what a deep network is learning:

- Pick a unit in layer 1. Find the nine image patches that maximize the unit's activation ; Repeat, for other units



Layer 1



Layer 2



Layer 3



Layer 4



Layer 5

Grams et al., 2015. A neural algorithm of artistic style.

Cost function: Content component + Style component.

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

similarity b/w contents similarity b/w styles

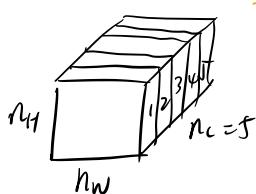
1. Initiate G randomly
2. Use gradient descent to minimize $J(G)$

Content cost function :

- Say you use hidden layer l to compute content cost
 - Use pretrained ConvNet (e.g., VGG network)
 - Let $a^{[l]^{(c)}}$ and $a^{[l]^{(G)}}$ be the activation of layer l on the images
- $$J_{\text{content}}(C, G) = \|a^{[l]^{(c)}} - a^{[l]^{(G)}}\|^2$$
- (element-wise sum of square differences)

Style cost function :

- Say you are using layer l 's activation to measure "style"
- Define style as correlation b/w activations across channels at layer l
 \rightarrow high correlation indicates the high probability of co-occurrence of certain features.



• Style matrix:

Let $a_{i,j,k}^{[l]}$ = activation at (i, j, k) . $G^{[l]}$ is $n_c^{[l]} \times n_c^{[l]}$ correlation matrix b/w each pair of channel.

$$G_{kk'}^{[l](S)} = \sum_i^H \sum_j^W a_{ijk}^{[l](S)} a_{ijk'}^{[l](S)}$$

Some type of co-variation

$$G_{kk'}^{[l](G)} = \sum_i^H \sum_j^W a_{ijk}^{[l](G)} a_{ijk'}^{[l](G)}$$

(Gram matrix)

$$J_{\text{style}}^{[l]}(S, G) = \|G^{[l](S)} - G^{[l](G)}\|_F^2 = \frac{1}{2(n_h^{[l]}n_w^{[l]}n_d^{[l]})^2} \sum_k \sum_{k'} (G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)})^2$$

↓ sum over all layers,

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

Class 5: Recurrent Neural Network

Use cases: Speech recognition; Music generation; Sentiment classification;
DNA sequence analysis; Machine translation; Video activity recognition
Name entity recognition.

Notations (using Name entity recognition as an example)

x : Harry Potter and Hermione Granger invented a new spell

$x^{<1>}$	$x^{<2>}$	$x^{<3>}$	---	$x^{<t>}$	---	$x^{<9>}$	$T_x = 9$
-----------	-----------	-----------	-----	-----------	-----	-----------	-----------

y : 1 1 0 1 1 0 0 0 0

$y^{<1>}$	$y^{<2>}$	$y^{<3>}$	---	$y^{<t>}$	---	$y^{<9>}$	$T_y = 9$
-----------	-----------	-----------	-----	-----------	-----	-----------	-----------

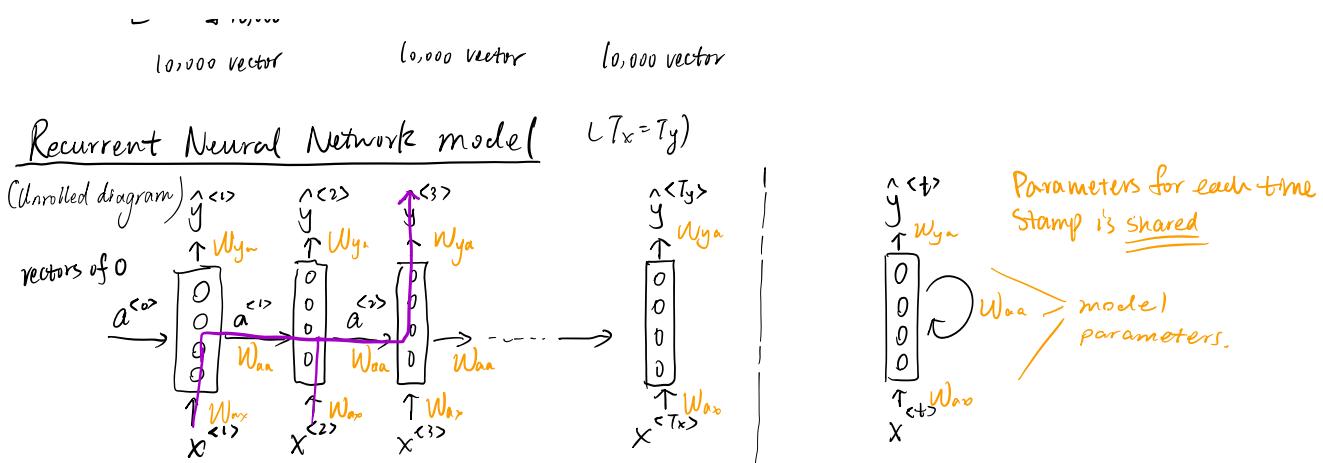
$x^{(i)<t>}$: t^{th} element in the sequence of training example i

$T_x^{(i)}$ (examples in the training set can have different length):

Input sequence length for training example i

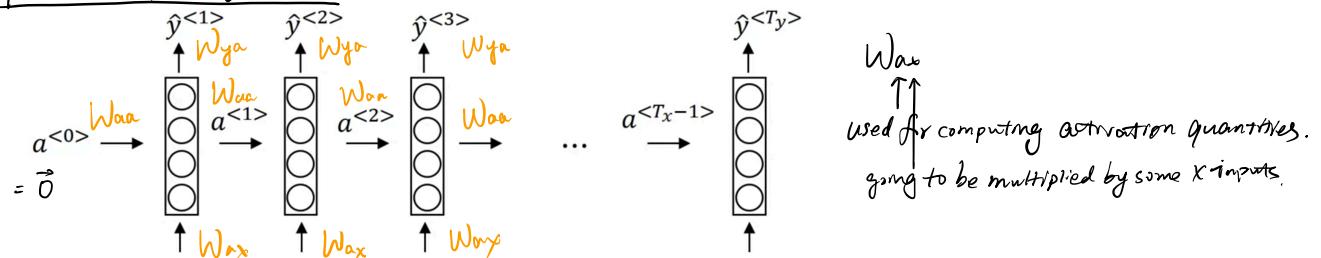
Representing words:

Vocab:	$x^{<1>} (\text{char})$	$x^{<2>} (\text{potter})$	one-hot representation.
$\begin{bmatrix} \text{a} \\ \text{and} \\ \text{harry} \\ \text{potter} \\ \vdots \\ \text{zulu} \end{bmatrix}^T \begin{bmatrix} 1 \\ 367 \\ 4075 \\ 6830 \\ \vdots \\ 10000 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}^T \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}^T \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}$	



Limitation: Prediction for a certain timepoint relies on the information earlier in the sequence only
 → Bidirectional RNN (BRNN)

Forward Propagation:



$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a) \quad g = \tanh / \text{ReLU}$$

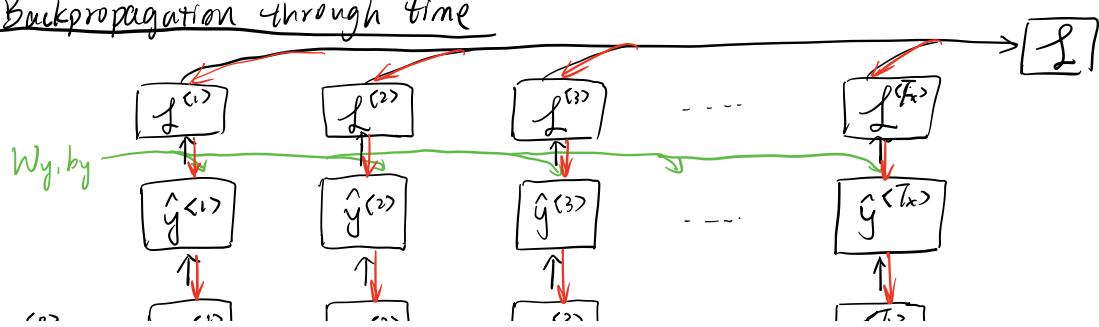
$$\hat{y}^{<t>} = g(W_{ya}a^{<t>} + b_y) \quad g \text{ depends on the type of output.}$$

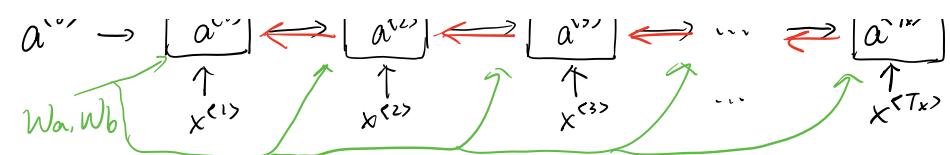
$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a) \rightarrow a^{<t>} = g(W_a[a^{<t-1>}, x^{<t>}] + b_a)$$

$$W_a = \begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix}; [a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}$$

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

Backpropagation through time



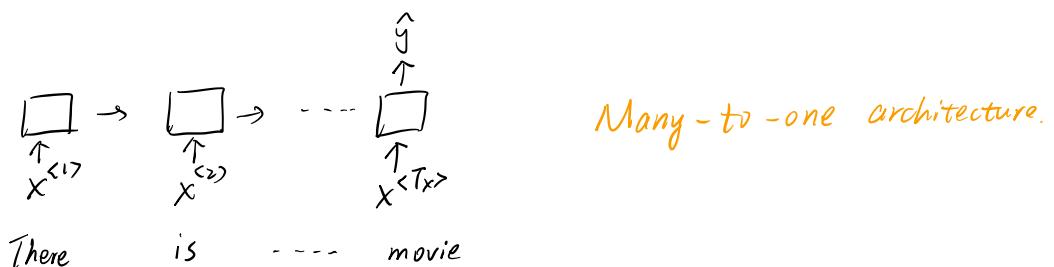


$$L^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1-y^{(t)}) \log (1-\hat{y}^{(t)})$$

$$L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

Different types of RNNs (T_x may not equal to T_y)

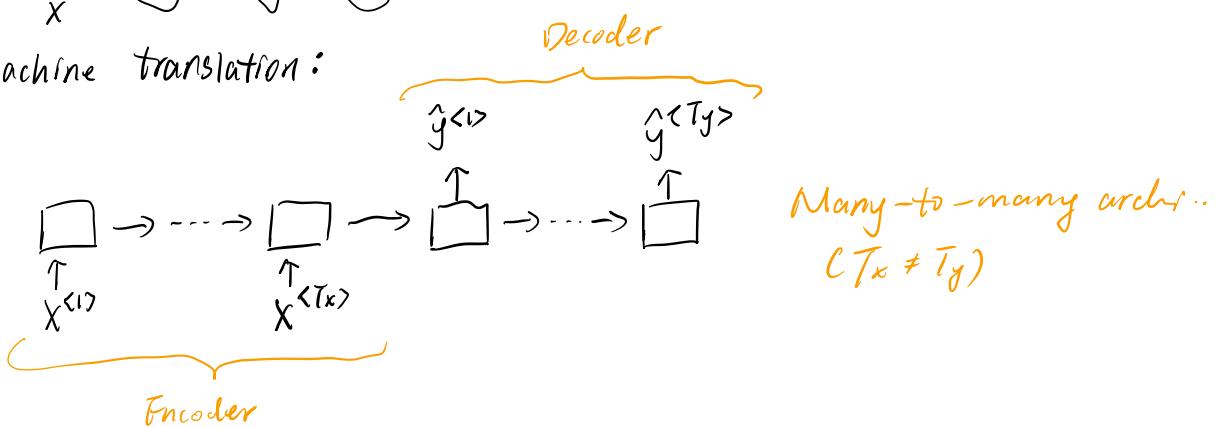
- Sentiment classification : $X = \text{text}$; $Y = 0/1$ or $1 \dots 5$



- Music generation : $X = \text{genre/first note}$ $Y = \text{sequence of music notes}$



- Machine translation :



Language Model and sequence generation

- Speech recognition : Model $P(\text{sentence}) = ? P(y^{(1)}, y^{(2)}, \dots, y^{(T_y)})$

Training set : Large corpus of english text.

Tokenize :

"Cats average 15 hours of sleep a day." <EOS> (end of sentence)

$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad \dots \quad y^{<8>} \quad y^{<9>}$

"The Egyptian Mau is a breed of cat." <EOS>

<UNK> (unknown token if this word is not in the dictionary)

RNN model (For language model)

"Cats average 15 hours of sleep a day."

- outputting the probability
of the first word being $y^{<1>}$
each of the words in the dictionary

$\alpha^{<0>} = \vec{0}$ $\xrightarrow{\text{softmax}}$ $a^{<1>}$
 $x^{<1>} = \vec{0}$ $\xrightarrow{\text{softmax}}$ $y^{<1>} = \text{"Cat"}$

- Model $P(y_i | \text{cat})$

for i in dict

$y^{<2>} = \text{"Average"}$

$y^{<3>} = \text{"Sleep"}$

$y^{<4>} = \text{"Hours"}$

$y^{<5>} = \text{"Ave"}$

$y^{<6>} = \text{"Rea"}$

$y^{<7>} = \text{"Bri"}$

$y^{<8>} = \text{"Dre"}$

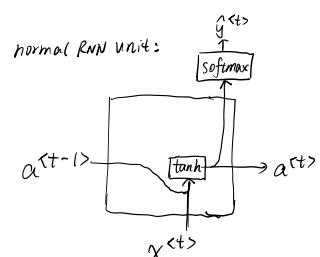
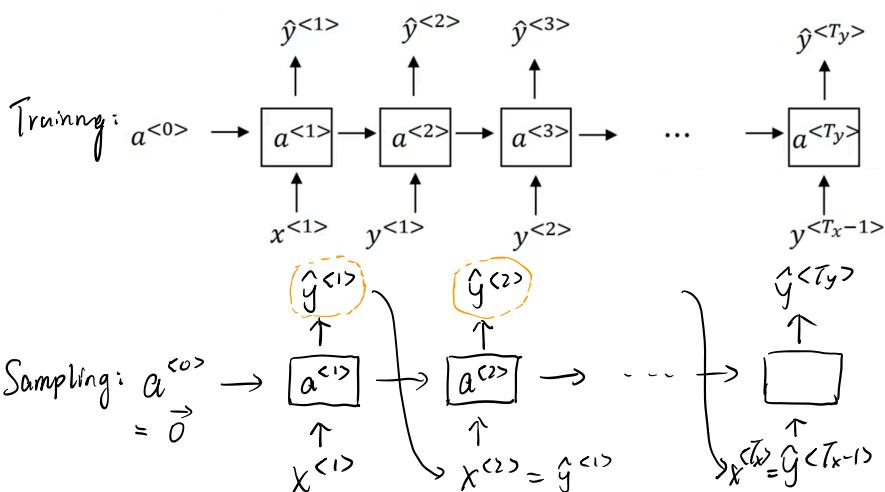
$y^{<9>} = \text{"Ave"}$

- If dict size is 1000, $\hat{y}^{<t>}$ is a 1000-dimensional vector. (<EOS>; <UNK>)

$$\mathcal{L}(\hat{y}^{<t>}, y^{<t>}) = -\sum_i y_i^{<t>} \log \hat{y}_i^{<t>} ; \quad \mathcal{L} = \frac{1}{T} \sum_t \mathcal{L}(\hat{y}^{<t>}, y^{<t>})$$

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

Sampling a Novel sequence from a trained RNN (Word level)

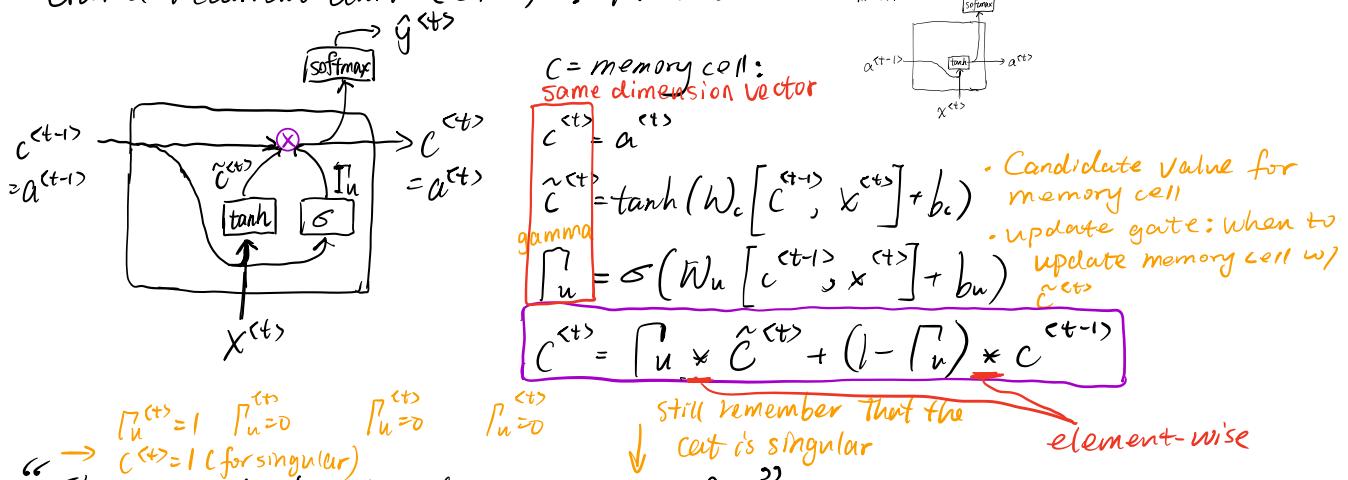


$\rightarrow p(a) \cdot p(a|aaron) \cdots p(zulu) \sim \text{np.random.choice}(5, 3, p=[0.1, 0, 0.3, 0.6, 0])$
 $\rightarrow \text{array([3, 3, 0])}$

Vanishing gradient with RNN

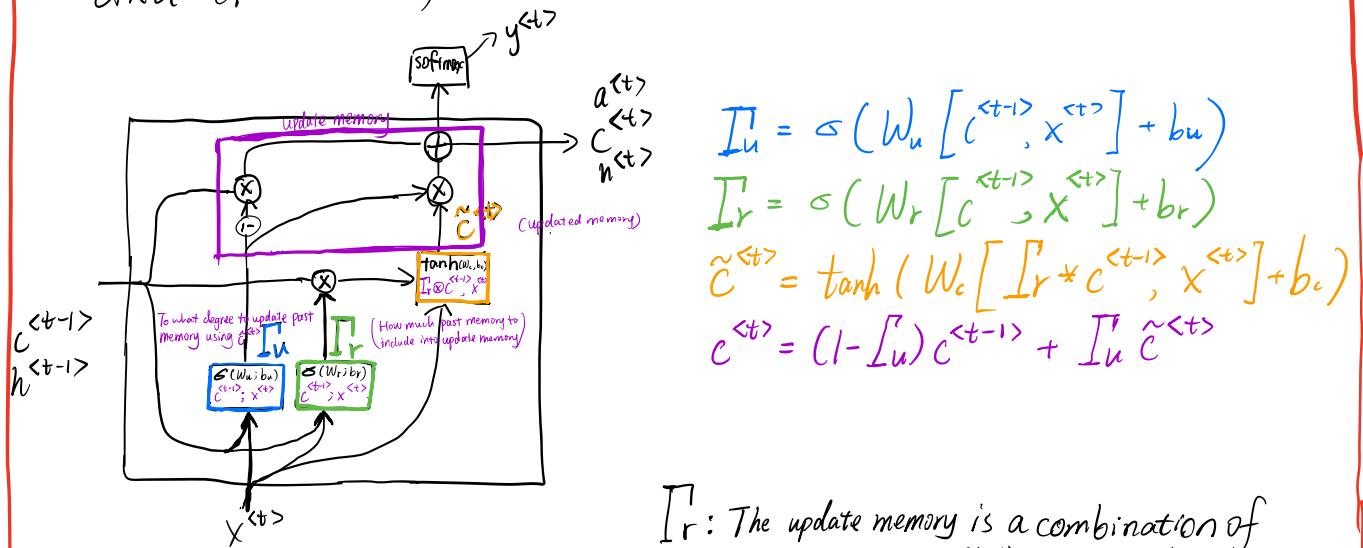
$\hat{y}^{(t)}$ is mainly influenced by values close to $\hat{y}^{(t)}$ (Very local influence)

- Gated Recurrent Unit (GRU) : RNN unit



"The cat, which already ate...., was full"

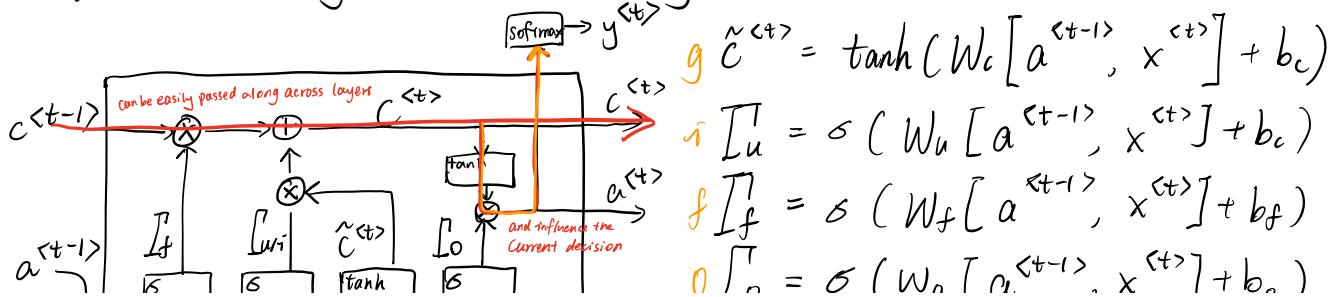
- GRU (Full version)

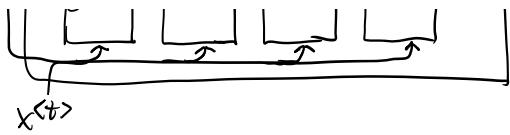


I_u : giving memory ($C^{(t-1)}$) and current input ($x^{(t)}$); How much should we update memory, resulting $C^{(t-1)}$

I_r : The update memory is a combination of past memory ($C^{(t-1)}$) and $x^{(t)}$. How much past memory should be included in the update memory?
 How relevant is $C^{(t-1)}$ to compute $C^{(t)}$

- LSTM : Long Short term memory Unit.





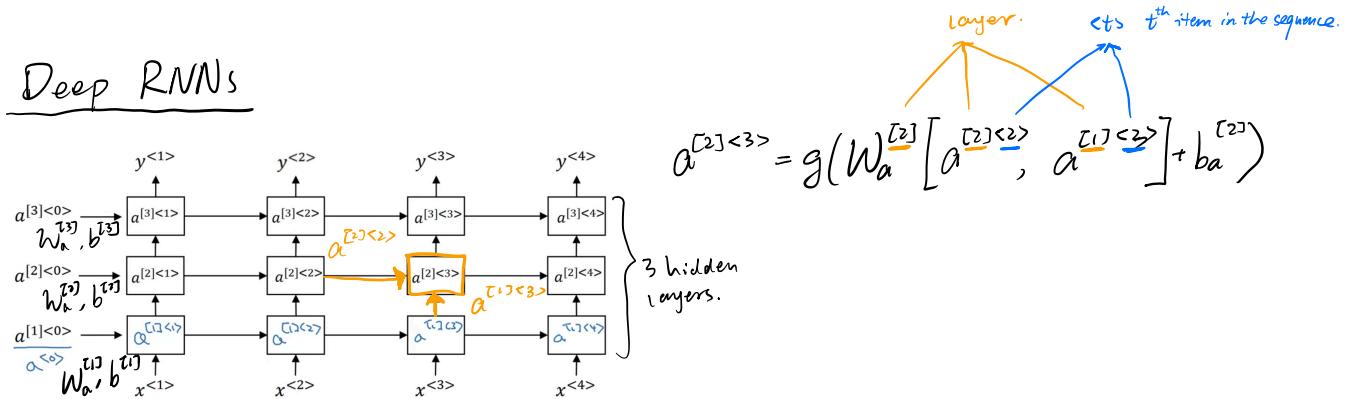
$$C^{(t)} = I_u * \hat{C}^{(t)} + I_f * \hat{C}^{(t-1)}$$

$$a^{(t)} = O_o * \tanh(C^{(t)})$$

Bidirectional RNN: Getting information from the future.

- He said, “Teddy bears are on sale” whether this is a human name depends on future information
- He said, “Teddy Roosevelt was a great president”

Deep RNNs



$$a^{[2]<3>} = g(W_a^{[2]} [a^{[2]<2>}, a^{[1]<2>}] + b^{[2]})$$

Word Representation

- 1-hot representation : Inner product b/w any 2 1-hot representations = 0
thus correlation matrix doesn't convey any semantic related info.

- Featurized representation: Word embedding embedding matrix

word feature	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
food	0.09	0.01	0.02	0.01	0.95	0.97
size	;	;	;	;		
cost	;	;	;	;		
	e_{5391}	e_{9853} (300-vector)				

300D → 2D
sklearn.manifold.tSNE

Using Word Embeddings (Like the encoding of faces in the Siamese network)

- ① Learn/Download word embedding from a large sample of unlabeled data (~1B)

Transfer learning

② Train, say named entity recognition model, using word embedding to represent each word ($\sim \text{unk}$)

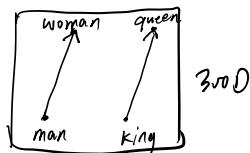
Note:

Cosine Similarity & Pearson Correlation look very similar. But:

- ① Cosine similarity computes the similarity b/w 2 samples from the same dist.
- ② Pearson correlation computes the relationship b/w 2 jointly distributed random variables

Properties of word embeddings

$$e_{\text{man}} - e_{\text{woman}} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$



$$e_{\text{king}} - e_{\text{queen}} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

Can learn Man: Woman \equiv Boy: Girl
Big: Bigger \equiv Tall: Taller

$$\underset{w}{\operatorname{argmax}} \operatorname{Sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$$

$$\operatorname{Sim}(u, v) = \frac{\langle u, v \rangle}{\|u\|_2 \|v\|_2}$$

Learning word embeddings (For language model)

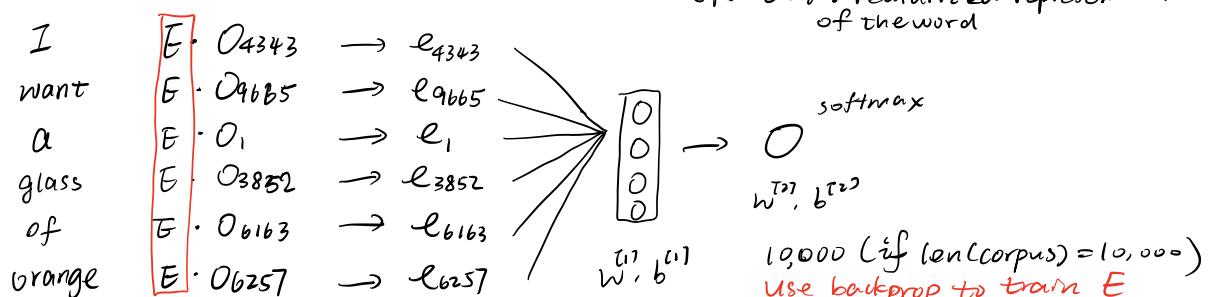
I want a glass of orange ?
4343 9665 1 3852 6163 6257

i^{th} word in the corpus.

O_i : One-hot representation of the word.

E : embedding matrix

$e_i = E O_i$: Featureized representation of the word



10,000 (if len(corpus)=10,000)
use backprop to train E

For constructing Embedding matrix

Skip Gram model: Learn context \rightarrow target mapping, targets are ± 10 from context

$$O_{\text{context}} E \rightarrow e_{\text{context}} \rightarrow O \xrightarrow{\text{softmax}} \hat{y} \quad P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}} \quad L(\hat{y}, y) = -\log \hat{y}_t$$

Negative Sampling To solve this:

choose context choose target.

"I want a glass of orange juice to go along with my cereal" Corpus=10000

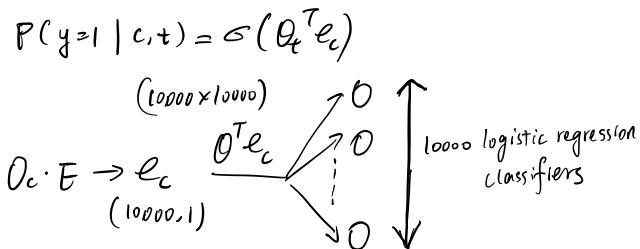


Context	word	target
orange	Juice	1
orange	King	0
orange	book	0
K orange	the	0
orange of	of	0

negative examples

$c \uparrow$ $t \uparrow$ $y \uparrow$

$K=5-20$ for small dataset
 $= 2-5$ for large dataset



Turning 10000 way softmax to 10000 logistic regressions

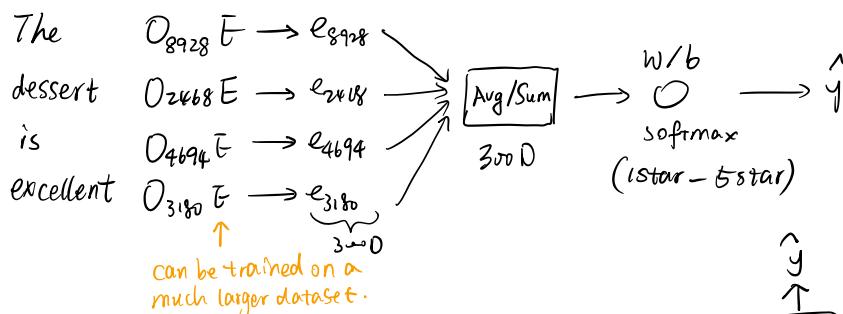
GloVe (global vector for word representation) model

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

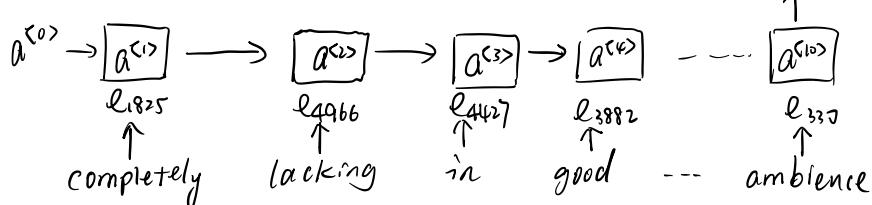
minimize $\sum_{i=1}^{10000} \sum_{j=1}^{10000} f(X_{ij})(\underbrace{\theta_i^T e_j + b_i + b_j}_{\text{weighting term measures how related is } i \& j} - \underbrace{\log X_{ij}}_c)^2$

Sentiment classification problem

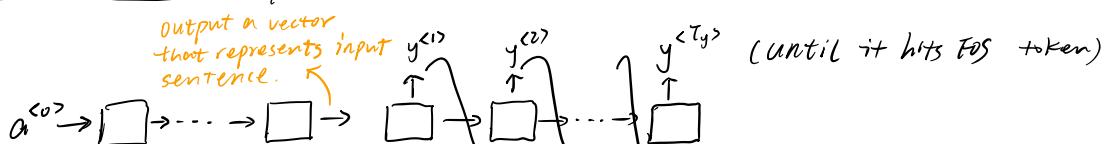
"The dessert is excellent" → ⭐⭐⭐⭐

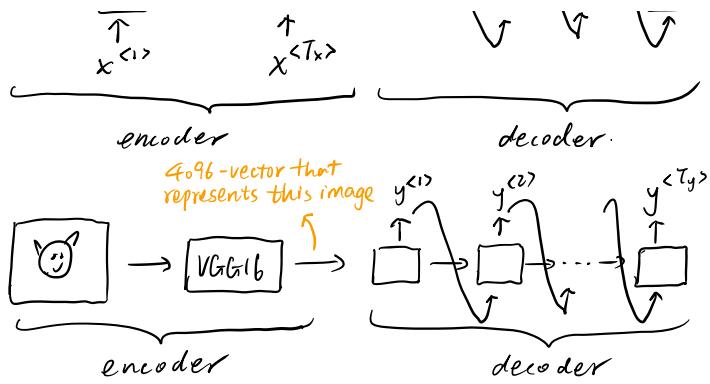


Many-to-One architecture:



Sequence-to-sequence model (e.g. Machine translation / Image captioning)



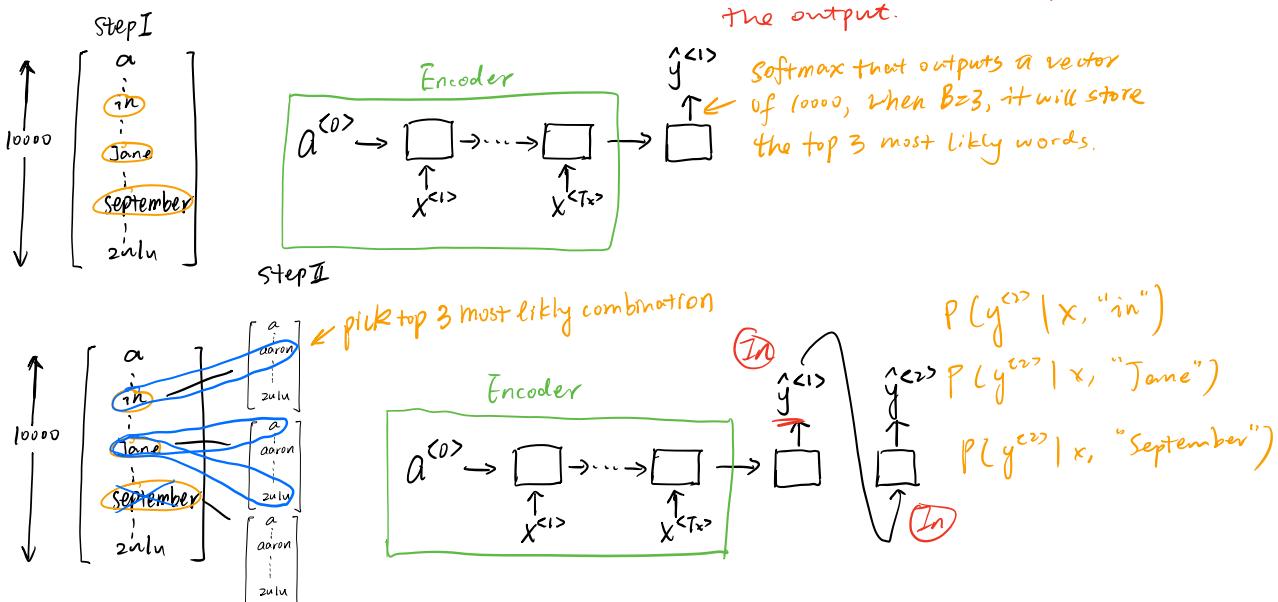


Picking the most likely sentence

- In a regular language model, we train a RNN to model $P(y^{<1>} | y^{<2>} \dots y^{<T_y>})$, and then randomly sample y from this distribution.
- In the case of machine translation, the encoder-decoder model $P(y^{<1>} \dots y^{<T_y>} | x)$, but we do not want to randomly sample our output from this distribution.

Instead $\underset{y^{<1>} \dots y^{<T_y>}}{\operatorname{argmax}} P(y^{<1>} \dots y^{<T_y>} | x) \leftarrow$ Find y that maximize this conditional joint prob.

Beam Search algorithm ($B=3$; beam width): At each step, instantiate 3 copies of the network to evaluate partial sentence of the output.



- Length Normalization:

$$\underset{y}{\operatorname{argmax}} \frac{1}{T_y} \sum_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>} \dots y^{<t-1>}) \rightarrow \text{Problems: } \begin{array}{l} \text{Can be too small (} \textcircled{1} \text{)} \\ \text{long sentence would have } \textcircled{2} \end{array}$$

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

small probability.

- Large B : better result, slower.

Small B : worse result, faster.

encoder + decoder.

Error analysis in Beam search: Problem due to beam search or RNN model?

e.g., Jane visits Africa in September (\hat{Y}^*)

Jane visited Africa last September (\hat{Y}) → Decide whether the problem is the beam width (B)?

→ Test $P(Y^*|x)$ vs. $P(\hat{Y}|x)$

↳ Meaning first pass x through encoder then use decoder softmax outputs to compute $P(Y^*|x)$ & $P(\hat{Y}|x)$

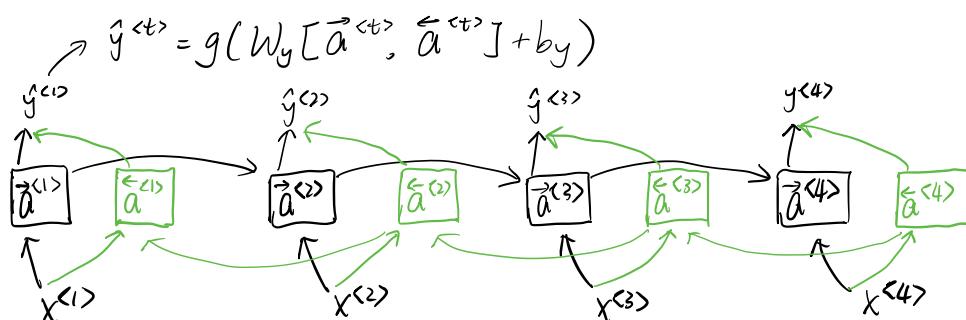
- If $P(Y^*|x) > P(\hat{Y}|x)$: problem with beam width.

- If $P(Y^*|x) \leq P(\hat{Y}|x)$: RNN model is at fault. (After length normalization)

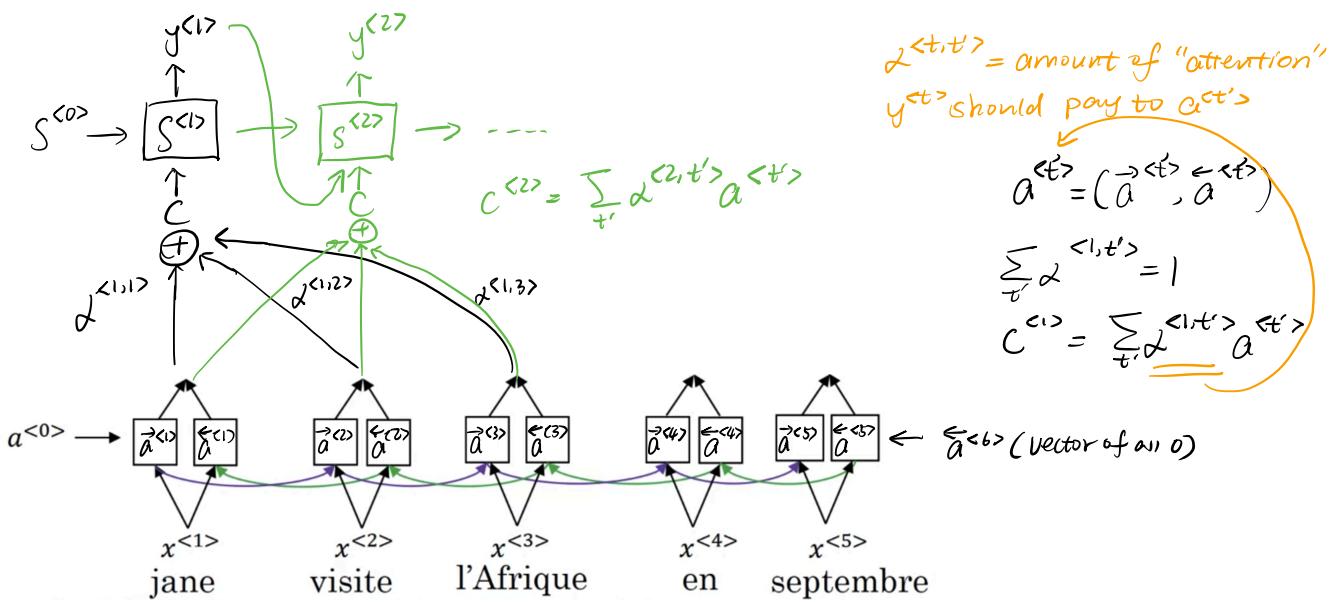
Attention: Improving sequence to sequence model

- In the case of machine translation, the encoder needs to remember the whole sentence before passing it to the decoder, which is not the case for humans.
- With the Attention model, the RNN could look at part of the sentence at each time.

Recall BRNN: First compute all hidden layer activations, then make predictions



Attention model:



- Computing attention $\alpha^{t,t'}$

$\alpha^{t,t'}$ = amount of attention $y^{t'}$ should pay to $a^{t'}$

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

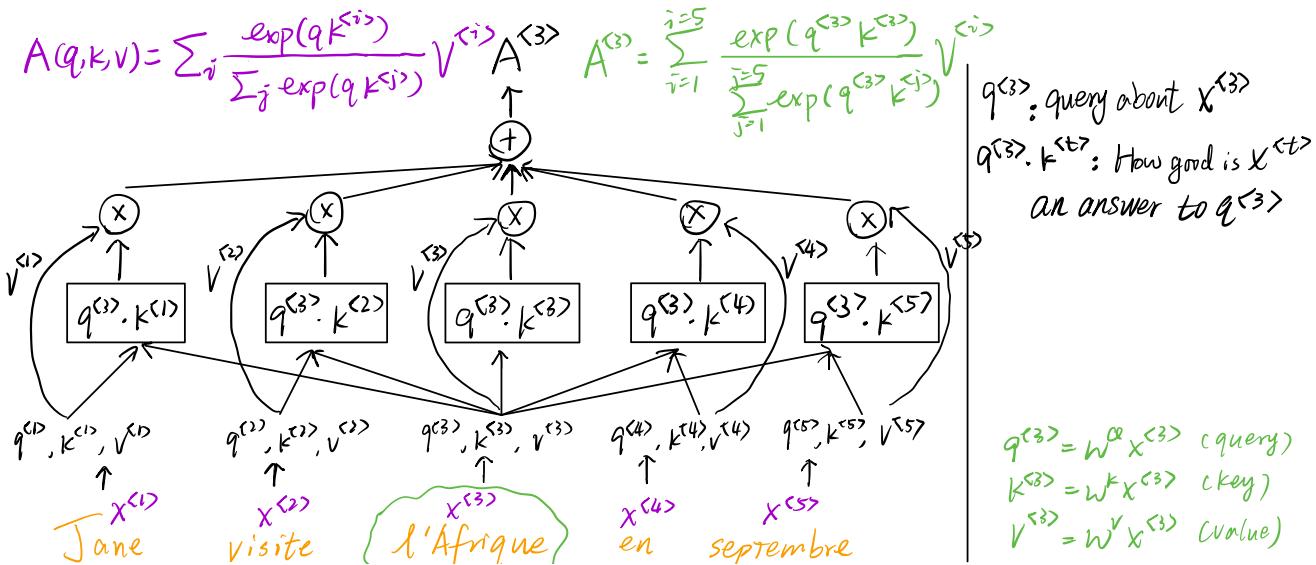
$$S^{t-1} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow e^{t,t'}$$

Transformer Network Intuition: Attention + CNN.

Vaswani et al. 2017 Attention Is All You Need.

→ Self Attention (Machine Translation as an example)

- Calculate for each input word attention-based vector representation



$$q^{(3)} = W^Q x^{(3)} \quad (\text{query})$$

$$k^{(3)} = W^K x^{(3)} \quad (\text{key})$$

$$v^{(3)} = W^V x^{(3)} \quad (\text{value})$$

$q^{(3)}$: query about $x^{(3)}$
 $q^{(3)}.k^{(t)}$: How good is $x^{(t)}$
 an answer to $q^{(3)}$

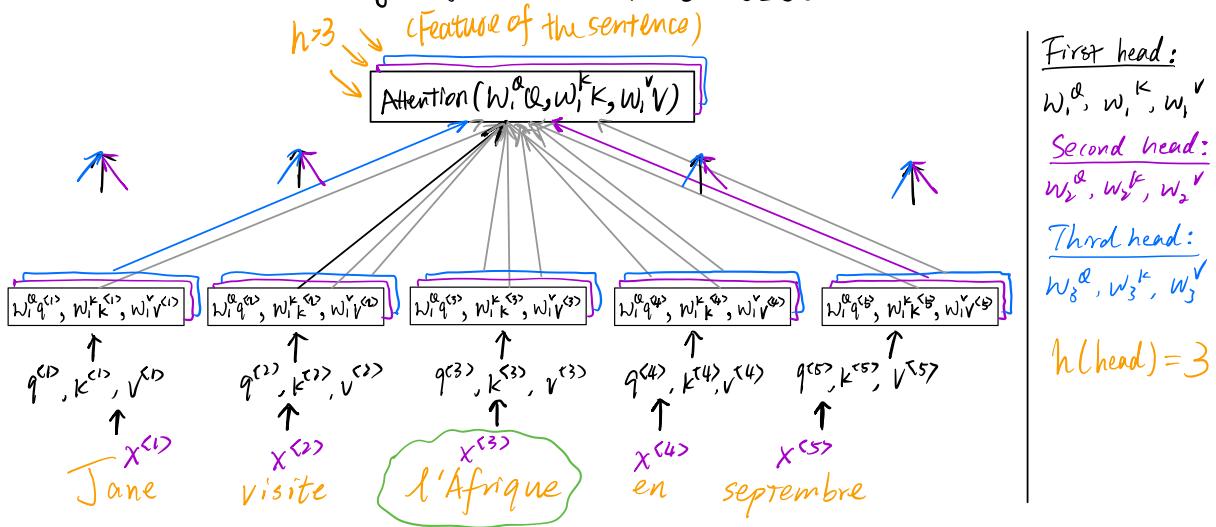
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{dk}}\right)V$$

scale parameter

→ The word would not be a fixed word embedding, but defined by the sequence.

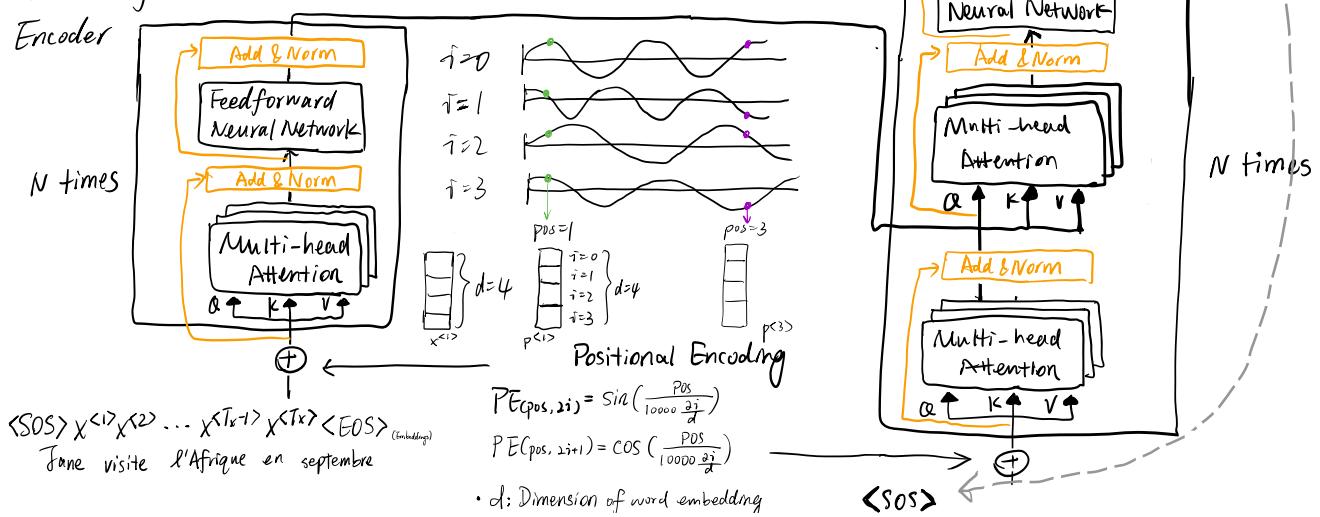
→ Multi-head Attention :

- Each time calculating self-Attention for a sequence is called a head



- Attention (Q, K, V) = $\text{softmax}\left(\frac{\alpha K^T}{\sqrt{d_k}}\right)V$
 - Multihead (Q, K, V) = $\text{concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h)W_0$
 - $\text{head}_i = \text{Attention}(w_i^Q Q, w_i^K K, w_i^V V)$

→ Transformer Network.



Intuition of Decoder:

- $p^{(i)}$: positional encoding of the i^{th} word.
- i : different dimension
 - Thus same word's embedding would be different given different pos.

- Ask questions (query) from the current translation and get context from encoder
- Residue connection: Pass along positional encoding to the whole network.