

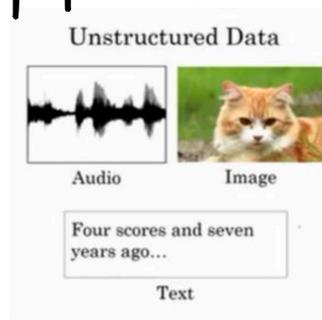
Deep Learning 学习笔记

Lecture One: Neural Network and Deep Learning:

Introduction to Deep Learning:

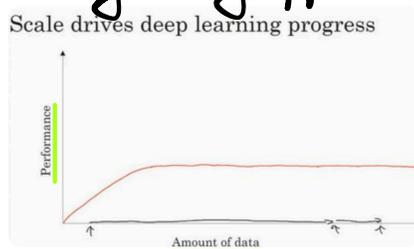
① Neural Network has high performance on:

Unstructured Data:



② Why is deep learning taking off?

1' Big Data



2' From Sigmoid to ReLU: FASTER! and avoid Vanishing Gradient.

Basic of Neural Network programming:

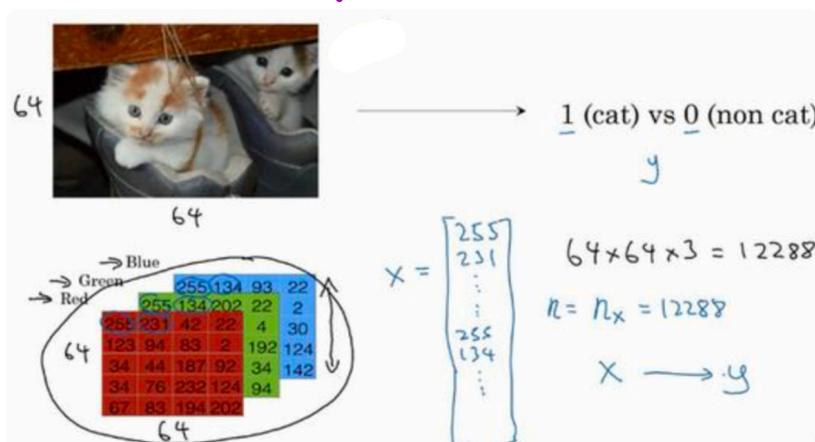
① Binary Classification:

1' Logistic Regression:

A image of cat is defined as: $(64 * 64 * 3)^{RGB}$

The feature vector will be: $n \times 1$.

The dimension of feature vector is: 12,288



Input: X.shape

(n_x, m)

Dimension quantity

Label: Y.shape

$(1, m)$

0 or 1

Via sigmoid function ($\delta(z) = 1 / (1 + e^{-z})$)

$$\hat{y} = \delta(b + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n), \text{ let } x_0 = 1.$$

$$\hat{y} = \delta(\theta^T X), \text{ where } \theta = [b, \theta_1, \dots, \theta_n]^T$$

\downarrow $\underbrace{\hspace{2cm}}$
 b w

② Logistic Regression Cost Function:

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

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

③ Gradient Descent in Logistic Regression:

Given two features: x_1, x_2 , we have $z = w_1 x_1 + w_2 x_2 + b$, $\hat{y} = a = \delta(z)$

Cost of EACH sample is: $L(a, y) = -y \log(a) - (1-y) \log(1-a)$

$$\Rightarrow \frac{\partial L}{\partial a} = -\frac{y}{a} + (1-y)/(1-a), \quad \frac{\partial L}{\partial z} = a - y$$

$$\frac{\partial L}{\partial w_1} = x_1 \cdot \frac{\partial L}{\partial z} = x_1 (a - y), \quad \frac{\partial L}{\partial w_2} = x_2 \cdot \frac{\partial L}{\partial z} = x_2 (a - y), \quad \frac{\partial L}{\partial b} = \frac{\partial L}{\partial z} = (a - y)$$

$$\Rightarrow w_1' = w_1 - \alpha \cdot \frac{\partial L}{\partial w_1}, \quad w_2' = w_2 - \alpha \cdot \frac{\partial L}{\partial w_2}, \quad b' = b - \alpha \cdot \frac{\partial L}{\partial b}$$

The disaster of "For loop": to traverse m samples is only one step, you also need to traverse each feature!

④: Vectorization: for accelerating computation.

Without vectorization: $z = w^T X + b$

With vectorization:

$$z = 0;$$

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

for i in range($n-x$):

$$z += W[i] * X[i]$$

$$z += b$$

Conclusion: Try not to use "FOR loop"

Homework:

① In numpy: $\begin{cases} 1 \text{ means ROW} \\ 0 \text{ means COLUMN} \end{cases}$ and $\begin{cases} a = \text{np.random.randn}(5) \\ a = \text{np.random.randn}(5, 1) \end{cases}$

Potential Risk \rightarrow

② ASSERT is always useful: `assert(a.shape == (5, 1))`
for reducing risk.

③ L2 normalization: $\text{np.sum}(\hat{y} - y)$ VS $\text{np.dot}((y - \hat{y}), (y - \hat{y}).T)$

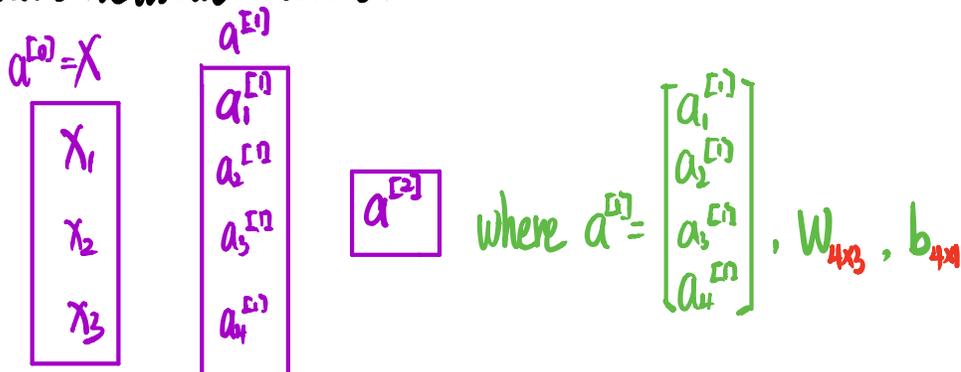
④ `np.squeeze()`: delete shape where dimension is one. `[[[1]]] \Rightarrow 1`

⑤ Remember Reshape your training set!

★⑥ Trick: `X.flatten(X.shape[0], -1).T \Rightarrow (a, b, c, d) \Rightarrow (b*c*d, a)`

Shallow neural networks:

Neural Network Overview:



Activation Functions:

sigmoid: Only used in output layer of two-classification

· tanh: widely used.

Relu, Leak Relu: widely used.

Random Initialization:

Without Random Initialization, Gradient Descent will not work!

Deep Neural Network:

Forward and backward propagation:

① Forward

$$1' \ z^{[L]} = W^{[L]} \cdot \underline{a^{[L-1]}} + b^{[L]}$$

save as cache! input

$$2' \ \underline{a^{[L]}} = g^{[L]}(z^{[L]})$$

output

② Backward:

$$1' \ dz^{[L]} = \underline{da^{[L]}} \cdot g^{[L]'}(z^{[L]})$$

input

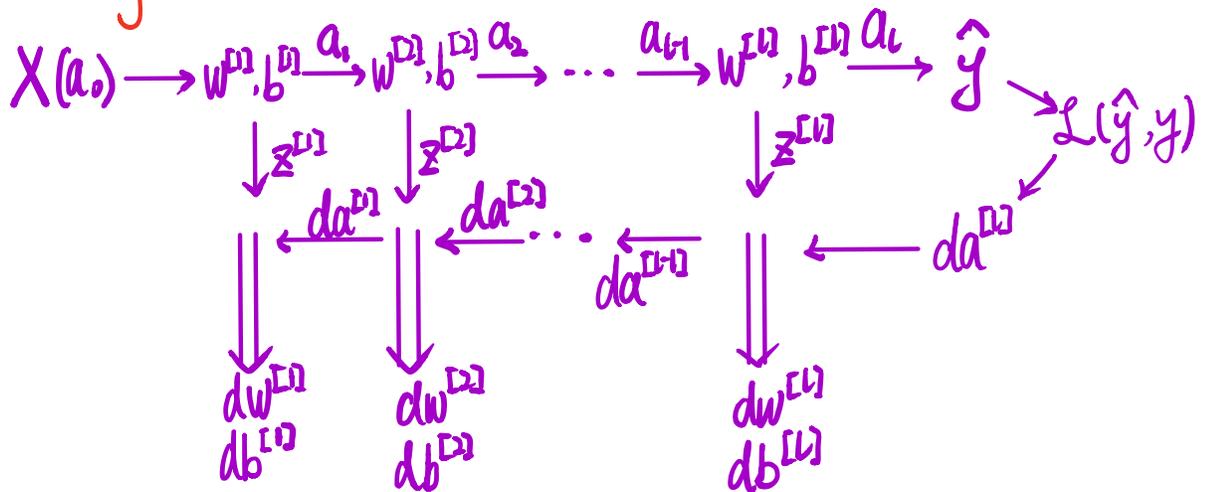
$$2' \ \underline{dw^{[L]}} = dz^{[L]} \cdot a^{[L-1]}$$

$$3' \ \underline{db^{[L]}} = dz^{[L]}$$

$$4' \ \underline{da^{[L-1]}} = W^{[L]T} \cdot dz^{[L]}$$

output

③ Summary



Lecture Two: Improving Deep Neural Networks: Hyperparameter tuning, Regularization, and Optimization.

Practical aspects of Deep Learning:

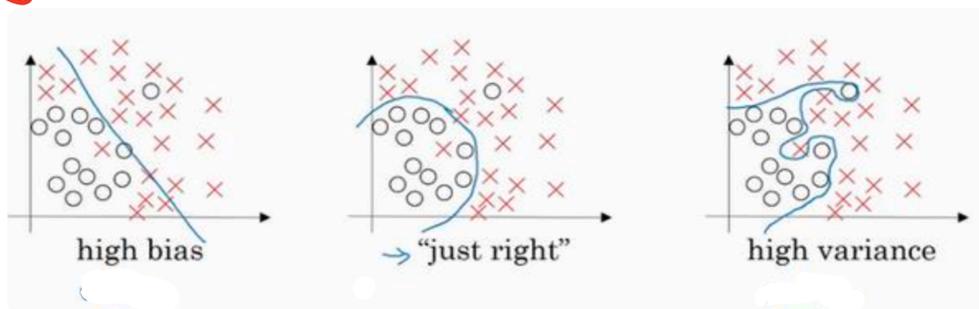
① Divide data set into: Train / Dev / Test

small data: 0.6 0.2 0.2

1m big data: 0.98 0.01 0.01

over 1m big data: 0.995 0.0025 0.0025

② Bias / Variance:



high bias: underfitting \Rightarrow change parameter or model

high variance: overfitting \Rightarrow increase data size or Regularization

③ Regularization:

\star L2 Regularization: $J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|^2$ sum the square of each element

L1 Regularization: $J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_1$

where $\|w\|^2 = \sum_i \sum_j (w_{ij})^2$

④ Dropout Regularization:

⑤ Normalizing inputs:

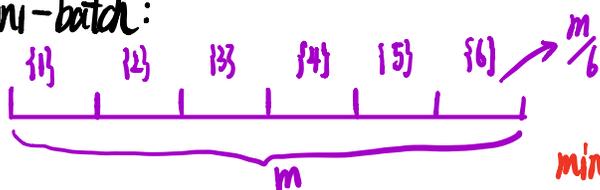
⑥ Vanishing / Exploding gradients:

Weight Initialization: For certain layer:

$np.random.randn(\text{shape}) * np.sqrt(\frac{1}{n^{[l-1]}})$ 2 is better for ReLU
number of l-1 layer's neuron.

Optimization algorithms:

① mini-batch:



one epoch only means one traverse of data
 mini-batch makes you do more gradient-descent

② Gradient descent with Momentum: \rightarrow usually set 0.9

$$\text{let } dw = V_{dw} = \beta V_{dw} + (1-\beta) dW \quad db = V_{db} = \beta V_{db} + (1-\beta) db$$

we obtain: $W = W - \alpha \cdot V_{dw}$, $b = b - \alpha V_{db}$ learning rate

③ RMS prop:

$$W := W - \alpha \cdot \frac{dw}{\sqrt{S_{dw}}} \quad b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$

* Adam: Widely Used !!!

⑤ Learning rate decay:

$$\text{set } \alpha = 1 / (1 + \text{decayrate} \times \text{epoch-num}) \times \alpha_0$$

⑥ The problem of local optimal): This may not that important!

Hyperparameter tuning:

① Batch normalization: \Rightarrow normalization " $z^{[L]}$ " * Before it came to " $a^{[L]}$ " *

$$\text{For certain layer: } \mu = \frac{1}{m} \sum z^{(i)}, \quad \sigma^2 = \frac{1}{m} \sum (z_i - \mu)^2$$

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

② Softmax regression: (multi-classify)

$$1' \text{ activation: } a_i^{[L]} = \frac{e^{z_i^{[L]}}}{\sum_j e^{z_j^{[L]}}}$$

$$2' \text{ loss: } \mathcal{L}(\hat{y}, y) = - \sum_i y_i \log \hat{y}_i$$

Lecture 3: Structuring Machine Learning Projects.

Orthogonalization: resolve complicated thing into single function

Error analysis:

Transfer Learning: (Serial) VS Multi-Task Learning (parallel)

End-to-end deep learning: DIRECT! better with lots of data.

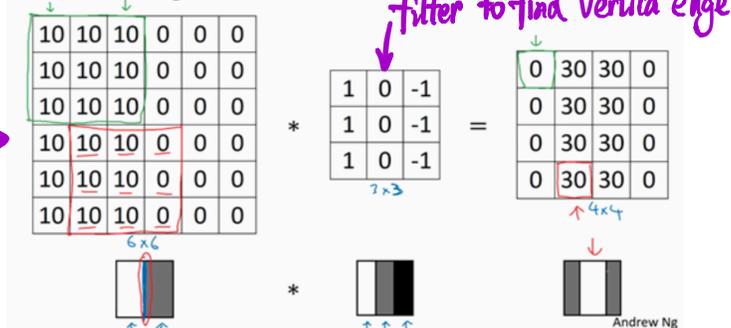
Lecture 4: Convolutional Neural Network

Foundations of Convolutional Neural Networks:

Big Picture \Rightarrow High Dimension $\not\rightarrow$ Low memory

Edge detection

Vertical edge detection



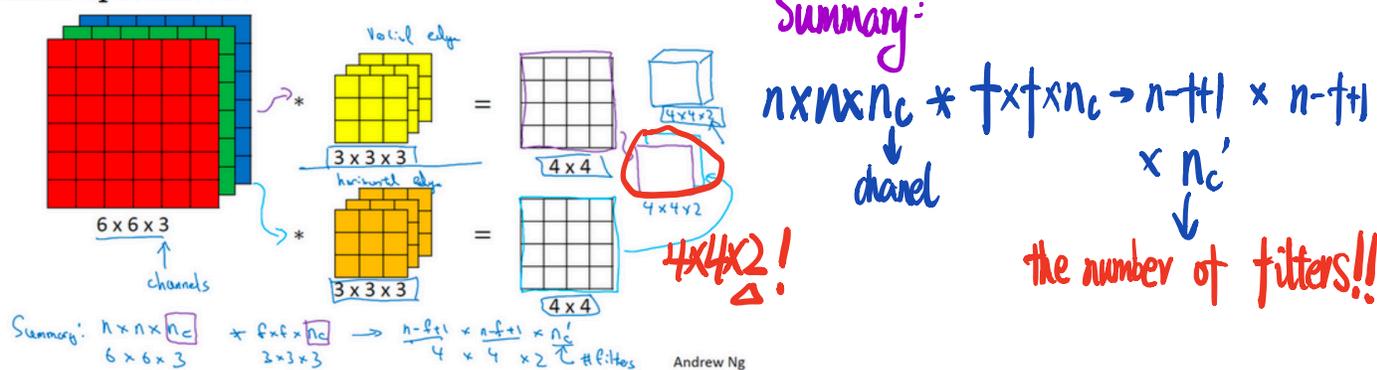
Padding:

Prevent image from becoming "1x1" size, padding image with "0".

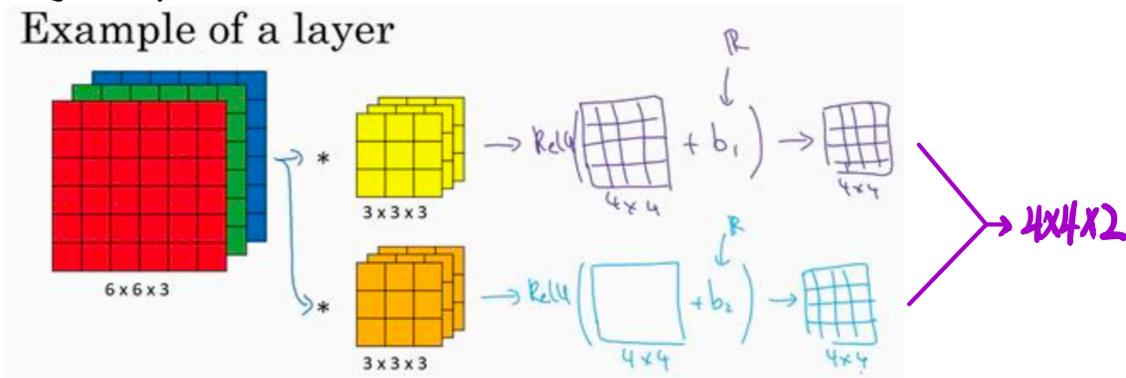
Stride convolution:

Convolutions over volumes:

Multiple filters



One layer of a convolutional network:



\therefore the number of parameter: $(3 \times 3 \times 3 + 1) \times 2 = 56$

each filter \downarrow bias \downarrow filter_num

Parameter description: if layer L is a convolution layer

$f^{[L]}$ = filter size

$p^{[L]}$ = padding

$s^{[L]}$ = stride

$n_c^{[L]}$ = number of filters

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

Activations: $a^{[L]} \rightarrow n_H^{[L]} \times n_W^{[L]} \times n_c^{[L]}$

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

Bias: $n_c^{[L]} \rightarrow (1, 1, 1, n_c^{[L]})$

Input: $n_H^{[L-1]} \times n_W^{[L-1]} \times n_c^{[L-1]}$

Output: $n_H^{[L]} \times n_W^{[L]} \times n_c^{[L]}$

Pooling layers:

Further reduce the size of the model.

① max pooling:

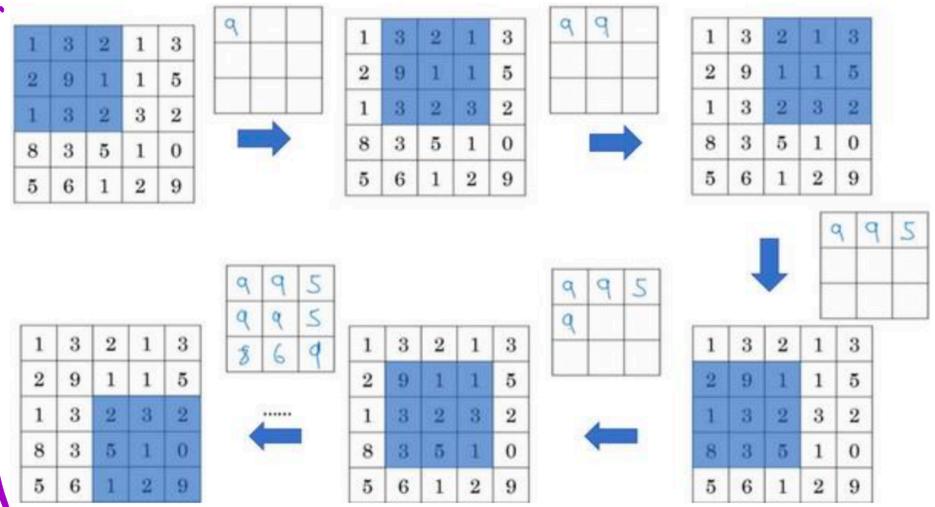
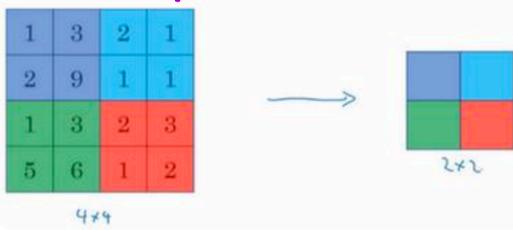
One example:

$f=2$
 $s=2$

widely used

another example

$f=3$
 $s=1$



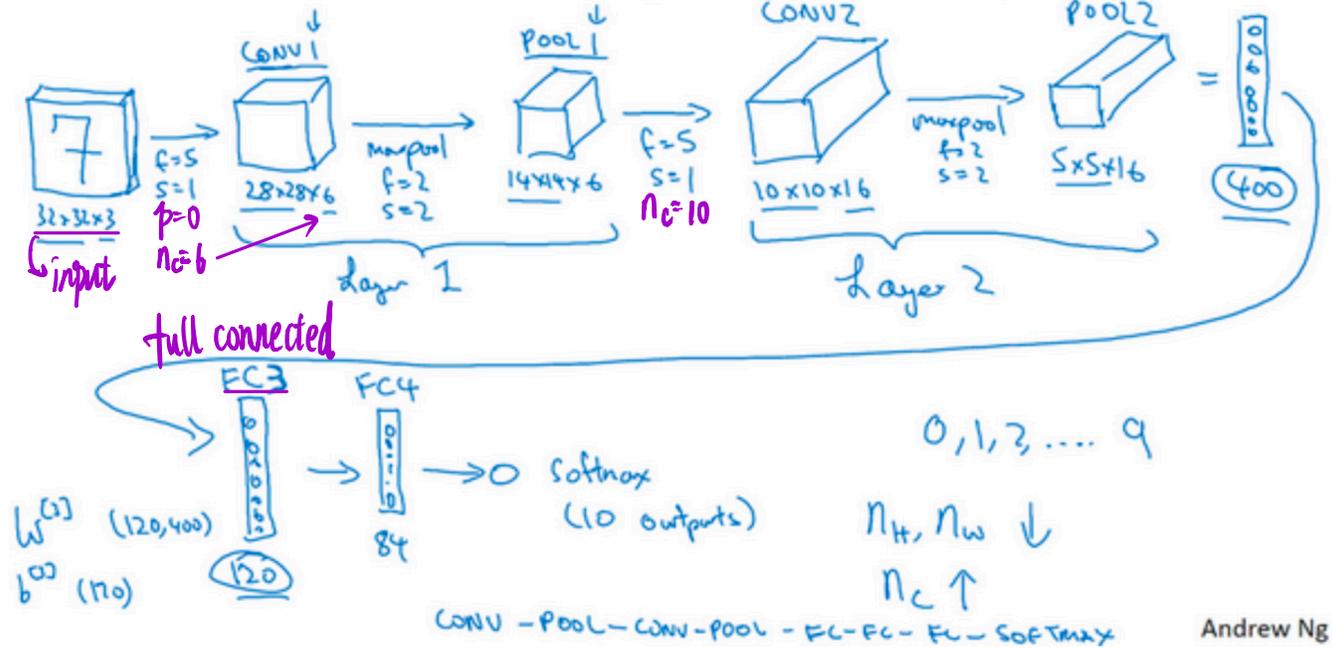
② formula:

max pooling output size: $\lfloor \frac{n_H - f}{s} + 1 \rfloor \times \lfloor \frac{n_W - f}{s} + 1 \rfloor \times n_c$

convolution output size: $\frac{n + 2p - f}{s} + 1$

Convolutional neural network example

Neural network example (LeNet-5)



Lower height and weight & higher channels.

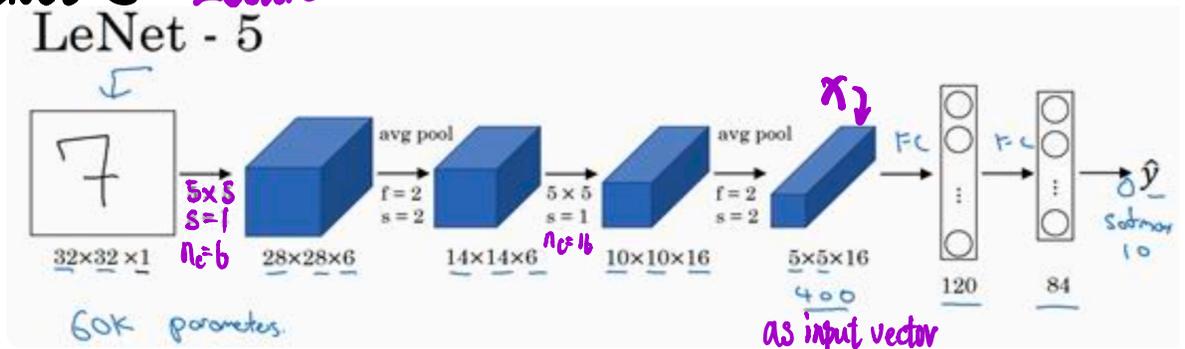
Why convolutions?

- ① Shared Parameters.
- ② Sparse Connection.

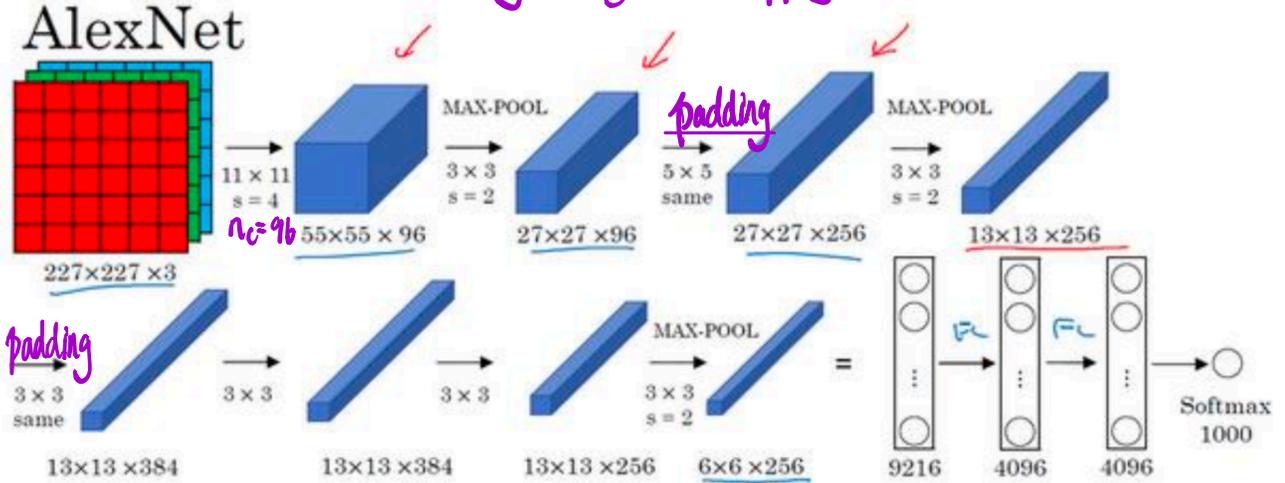
Deep Convolutional models: case studies.

Classic networks:

① LeNet-5: Lecun



② AlexNet: Alex Krizhevsky, Ilya, Geoffrey Hinton.

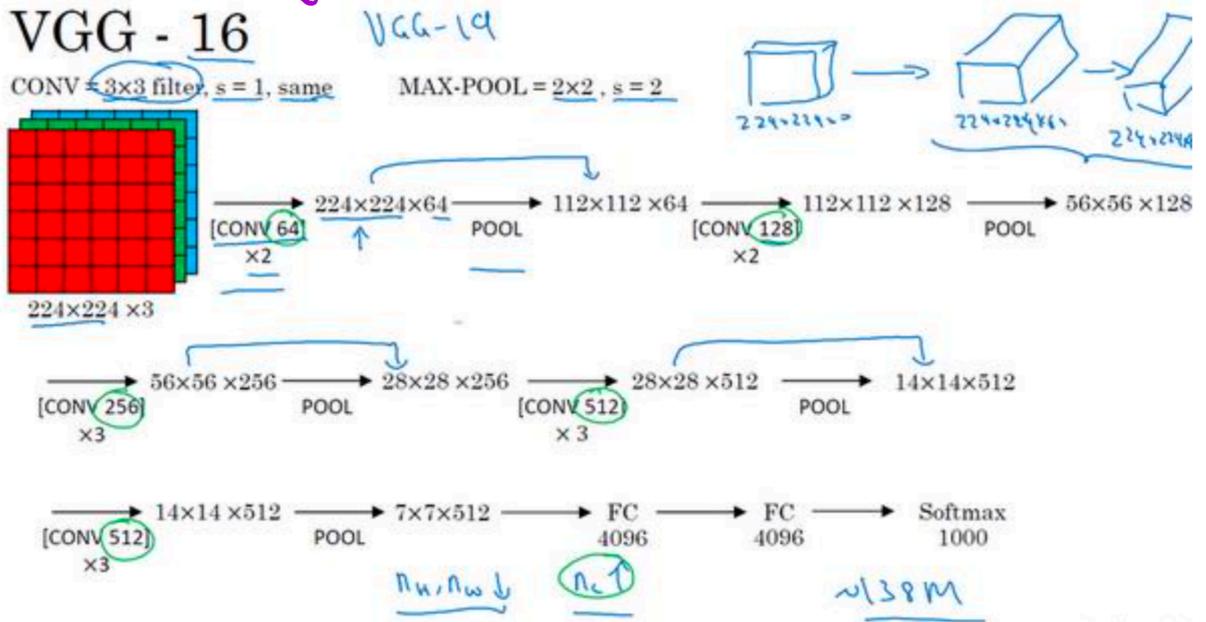


Activation is: ReLU compare to Sigmoid in LeNet-5

③ VGGNet:

filter_num: 64 \rightarrow 128 \rightarrow 256 \rightarrow 512 more and more

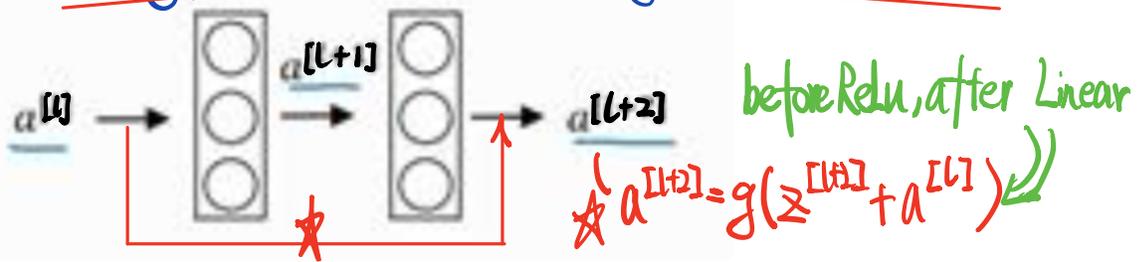
Convolutional_layer_num: 16



Residual Network [ResNets]

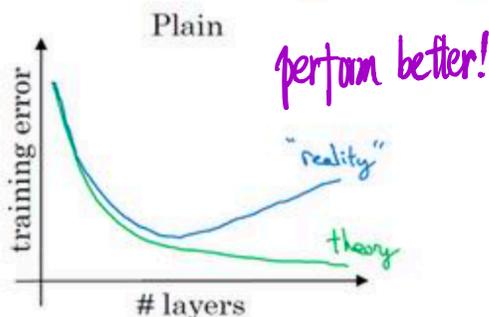
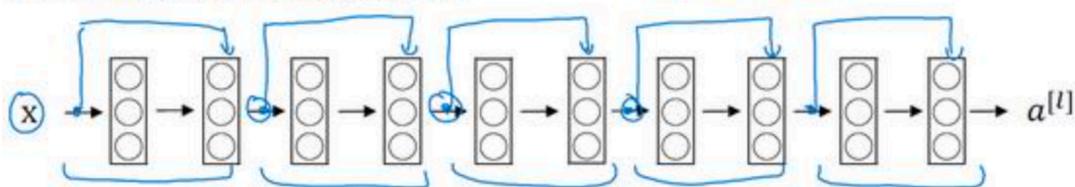
① Residual Block: message send to deeper layer directly!

~~$a^{[L+1]} = g(w^{[L+1]} \cdot a^{[L]} + b^{[L+1]})$ $a^{[L+2]} = g(w^{[L+2]} \cdot a^{[L+1]} + b^{[L+2]})$~~



② Residual Network with 5 blocks:

Residual Network



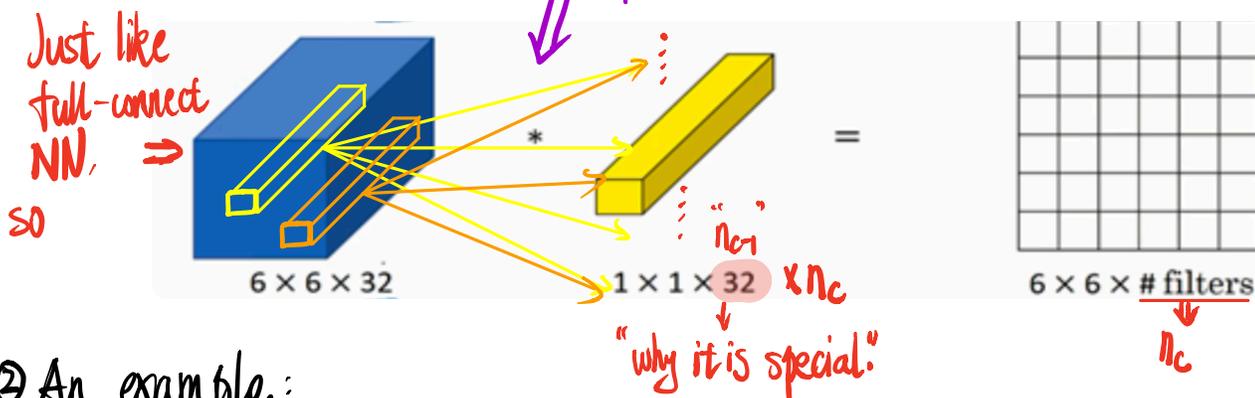
③ Why ResNets work?

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

ReLU, L2 regularization / Weight Decay make them to 0
 $\therefore a^{[L+2]} = a^{[L]}$

1x1 convolutions / (Network in network) ←

① Real size: $1 \times 1 \times n_{c1}$, this part like FNN, so it also called

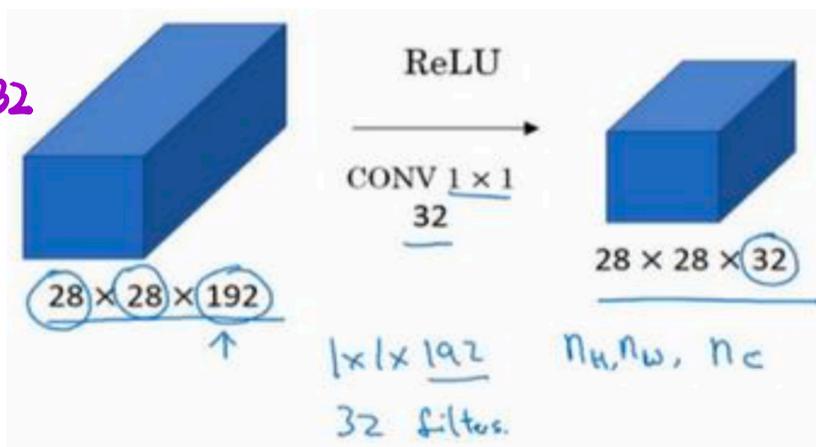


② An example:

How could we convert $28 \times 28 \times 192$ to $28 \times 28 \times 32$

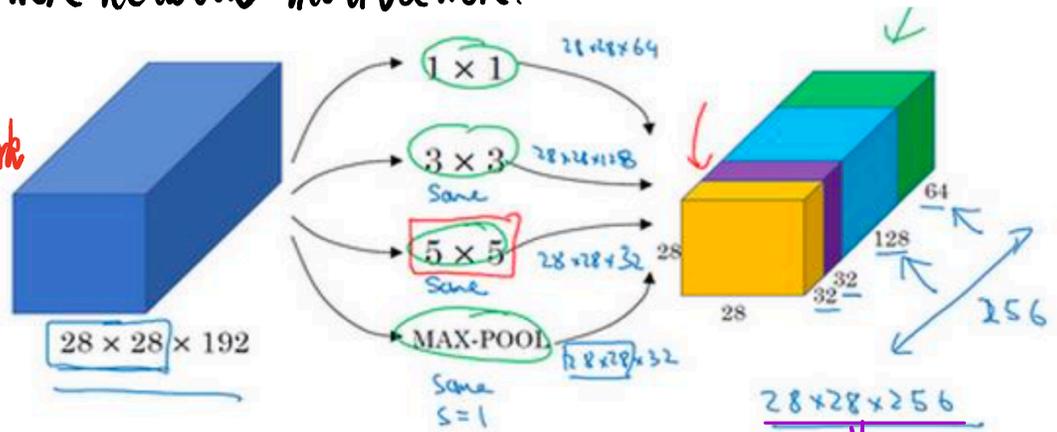
with 1x1 convolutions

$$28 \times 28 \times 192 \xrightarrow{n_{c1}} 1 \times 1 \times 192 \times 32 \xrightarrow{n_{c1}, n_c} 28 \times 28 \times 32$$



Google Inception network motivation.

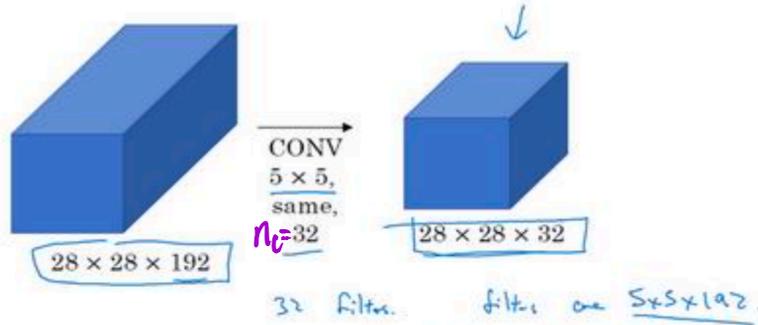
Let neural network learn by itself



Problem: COST!

The problem of computational cost

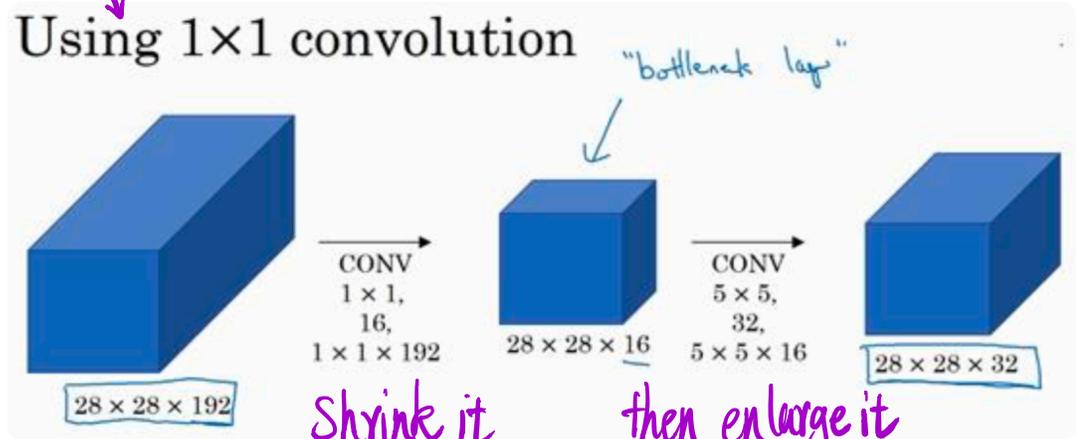
it is called Inception Module



Calculation: $(28 \times 28 \times 192 \times 5 \times 5 \times 32) = 120M$

Solution:

Using 1x1 convolution

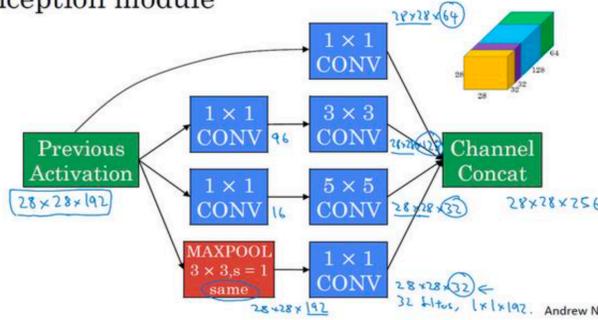


Shrink it then enlarge it

Calculation: $28 \times 28 \times 192 \times 16 + 28 \times 28 \times 16 \times 32 \times 5 \times 5 = 12.4M$

Inception network:

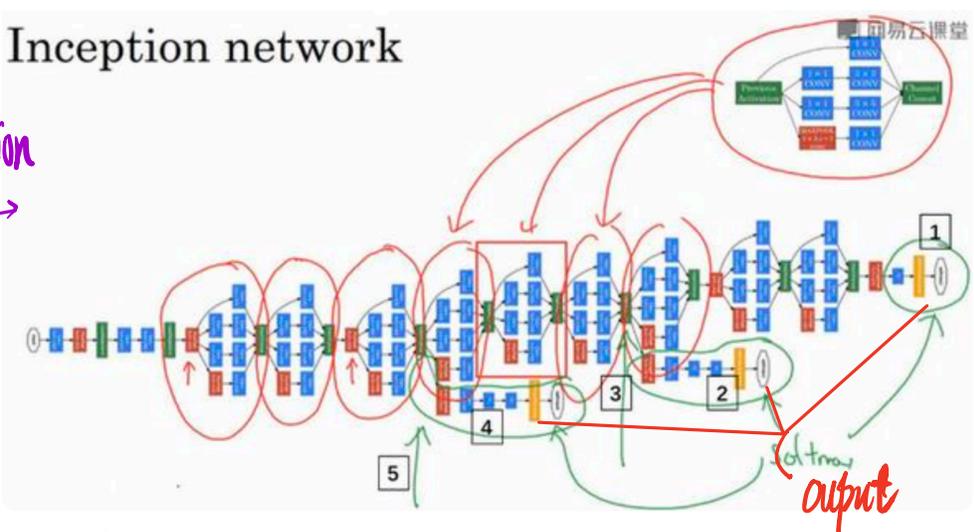
Inception module



with 1x1 conv to reduce calculation.

combine inception module into inception network

Inception network



If you want to go deeper: you can jump (ResNet) or you can quit everywhere (Inception)

Transfer Learning: Use other weights as Prior Knowledge.

Tricks for doing well on benchmarks / winning competitions.

- ① integration: mean the output of LOTS OF neural networks.
- ② enlarge dataset: Multi-crop

Object detection:

Object localization:

① Symbolic representation:

1' Top Left corner: $(0, 0)$; Bottom Right corner $(1, 1)$

2' $b_x, b_y, b_h, b_w \rightarrow$ central point (b_x, b_y) ; $b_h, b_w \in [0, 1]$

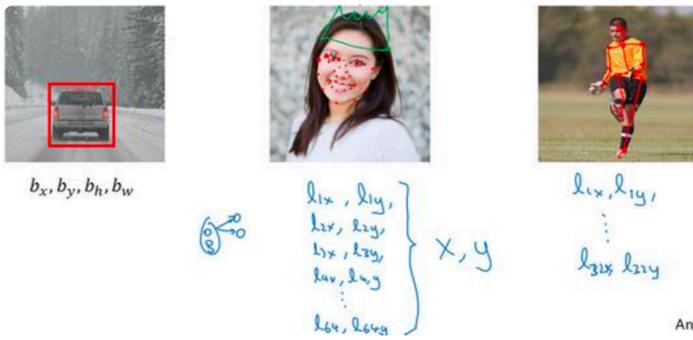
② How to define label "Y":

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

\rightarrow have object or not $\{0, 1\}$
 \rightarrow coordinate $\in [0, 1]$
 \rightarrow what object is $\{0, 1\}$

③ Loss function: $L(\hat{y}, y) = \sum_i (\hat{y}_i - y_i)^2$

Landmark detection:

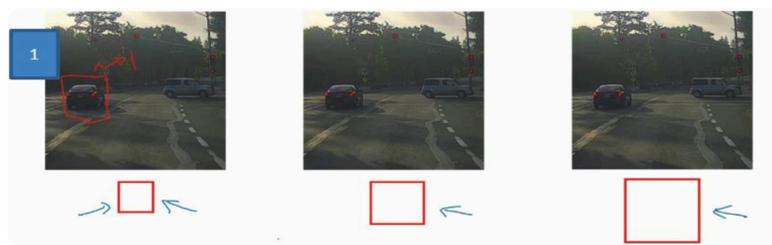


Andrew Ng

Object detection:

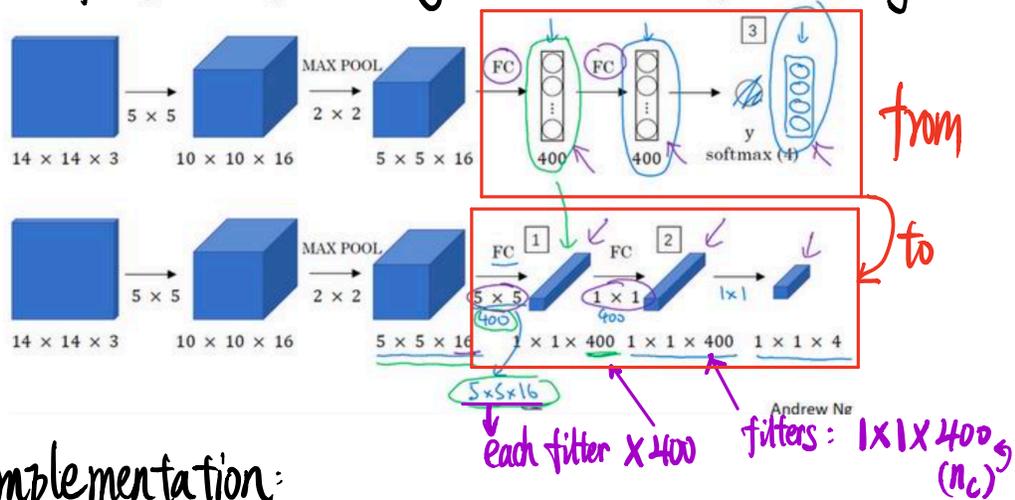
Sliding Window: \rightarrow low efficiency.

- 1' Train a CNN to detect car
- 2' Sliding windows in different size

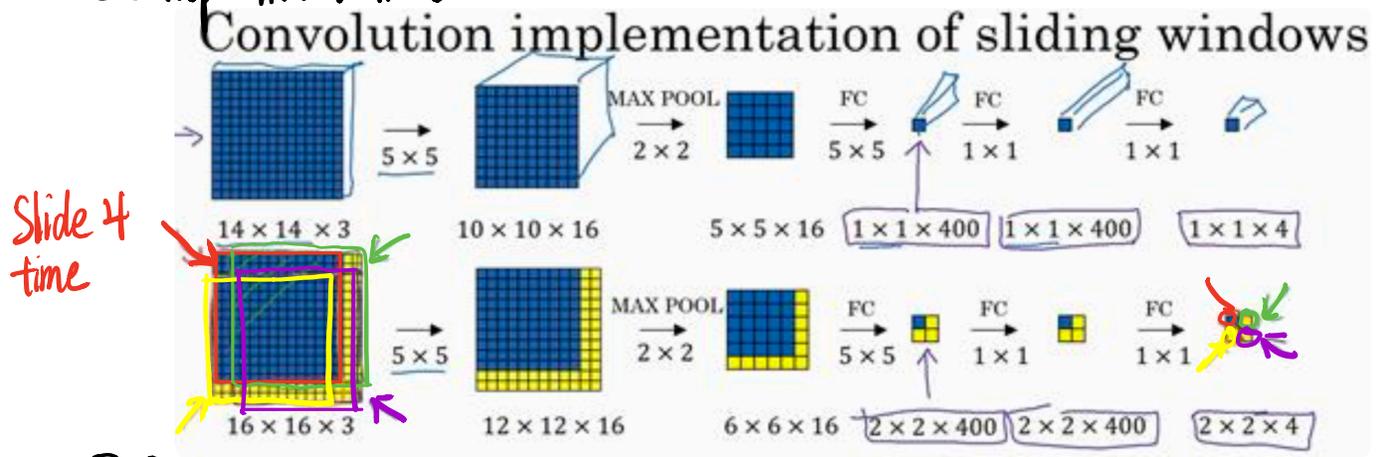


Convolutional implement of Sliding windows:

① Turn full-connected layer into Convolutional layer:



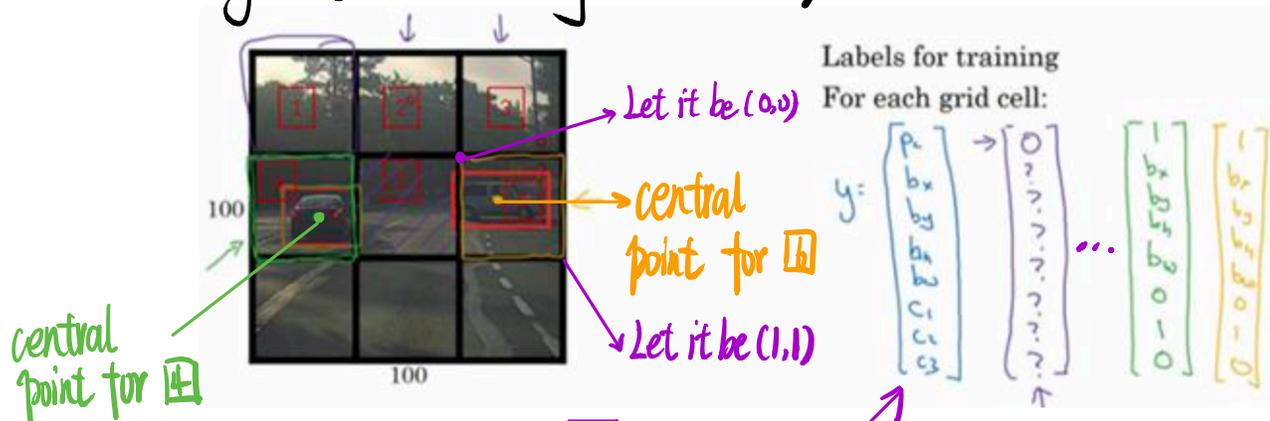
② Implementation:



③ Summary: sub_image \rightarrow CNN \rightarrow Output

Bounding Box predictions: (How to be more accurate)

① Yolo Algorithm: (You Only Look Once)



output is: $3 \times 3 \times 8$  8 vectors.

⇒ Let \square central point be $(0.4, 0.3)$

then determine b_h & b_w

Intersection over union: evaluate detection algorithm.

$$IoU = \frac{\text{Size of intersection}}{\text{Size of Union}} \quad \left\{ \begin{array}{l} \geq \text{threshold value} \quad \checkmark \\ \text{else} \quad \times \end{array} \right.$$

Non-max suppression:

make sure your algorithm only detect object once a time.

→ Only the classification results with the highest probability are output.

Anchor Boxes:

Let one box to detect multi-objects.

→ pre-define different kind of anchor boxes

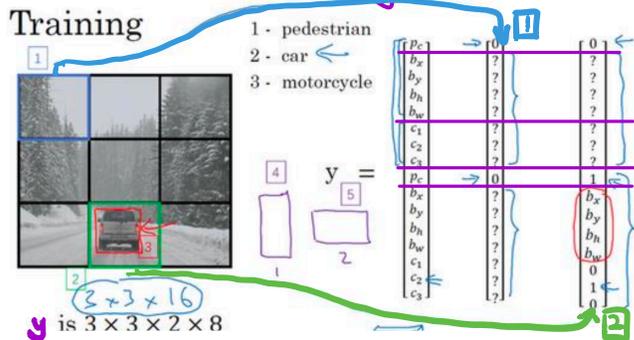
$$\therefore y_{\text{old}} = [p_c \ b_x \ b_y \ b_h \ b_w \ c_1 \ c_2 \ c_3]^T$$

$$y_{\text{new}} = [p_c \ b_x \ b_y \ b_h \ b_w \ c_1 \ c_2 \ c_3 \ p_c \ b_x \ b_y \ b_h \ b_w \ c_1 \ c_2 \ c_3]^T$$

Yolo Algorithm: Put all component above into algorithm.

① Build Training Set:

Let anchor box_num = 2, $y = 3 \times 3 \times (2 \times 8)$



② Predicting

③ Outputting the non-max suppressed outputs

Special applications: Face recognition & Neural style transfer

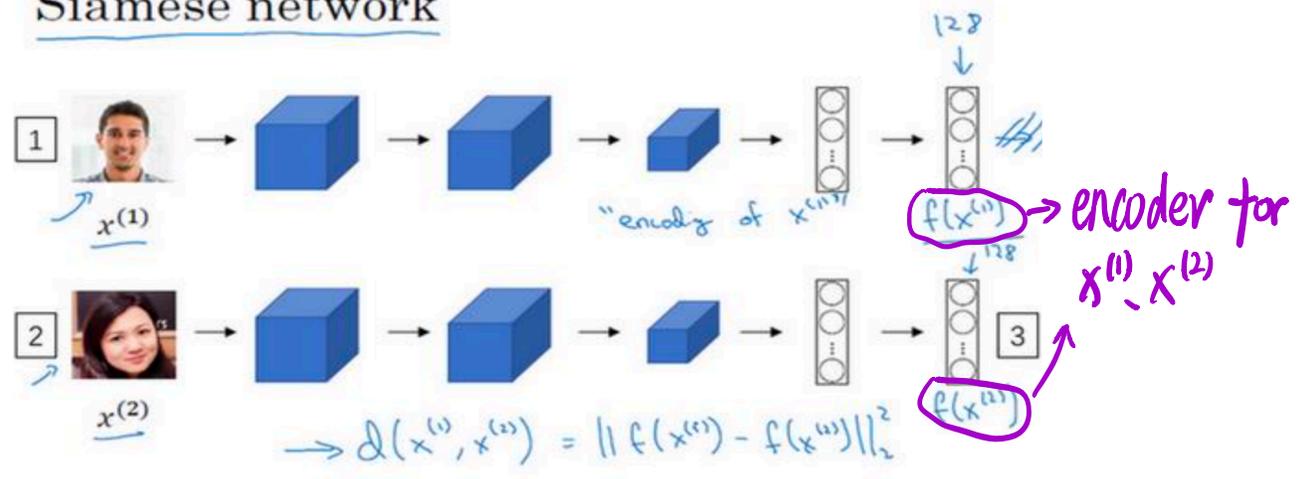
One-shot learning:

One sample to recognize a person.

→ to learn similarity function is a good idea.

Siamese network:

Siamese network



same person: low $\|f(x^{(i)}) - f(x^{(j)})\|^2$

different person: large $\|f(x^{(i)}) - f(x^{(j)})\|^2$

Triplet loss: $(A, P, N) \equiv (\text{Anchor, Positive, Negative})$

① Loss function

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

margin to prevent "0 ≤ 0"

$$\therefore \mathcal{L}(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0)$$

② Choosing the triplets A, P, N

1' A, P must be same person.

2' A, N must be different person.

③ Choosing triplets that're "hard" to train on.

$$d(A, P) \approx d(A, N)$$

Face verification and binary classification:

① other parameter learning method:

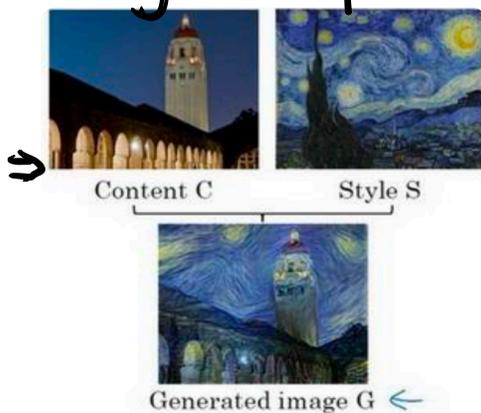
1' logistic regression

② face verification supervised learning:



data & label for face verification.

Neural Style Transfer:



We define Loss function as:

$$\textcircled{1} J_{\text{content}}(C, G)$$

$$\textcircled{2} J_{\text{style}}(S, G)$$

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

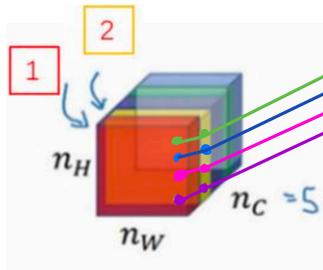
Content cost function:

If you use hidden layer l to compute content cost:

$$J_{\text{content}}(C, G) = \frac{1}{2} \| a^{[l][C]} - a^{[l][G]} \|^2 \rightarrow \text{are they similar enough?}$$

Style cost function:

The style of a picture is correlation index across different channels.



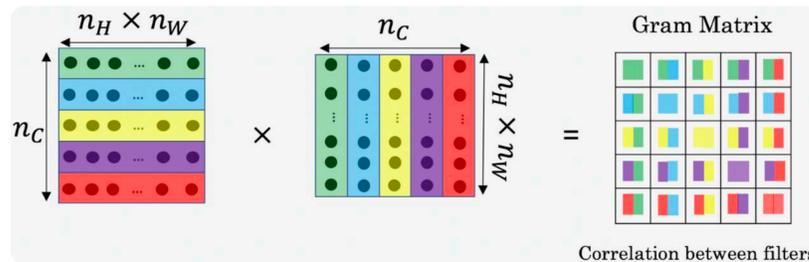
correlation means the uniformity probability of their own content is very high.

\Rightarrow if you measure style cost on layer l :

Let $a_{i,j,k}^{[l]}$ = activation at (i,j,k) . $G^{[l][S]}$ is $n_C \times n_C$ $\xrightarrow{\text{layer } l}$ style image

$$\text{we obtain: } G_{kk'}^{[l](S)} = \sum_{i=1}^{n_H} \sum_{j=1}^{n_W} a_{i,j,k}^{[l](S)} \cdot a_{i,j,k'}^{[l](S)}$$

$$\text{And } G_{kk'}^{[l](G)} = \sum_{i=1}^{n_H} \sum_{j=1}^{n_W} a_{i,j,k}^{[l](G)} \cdot a_{i,j,k'}^{[l](G)}$$



$$\Rightarrow \text{We finally get: } J_{\text{style}}(S, G) = \frac{1}{\beta} \| G^{[l](S)} - G^{[l](G)} \|^2$$

Lecture 5: Sequence Model

Recurrent Neural Network:

Notations:

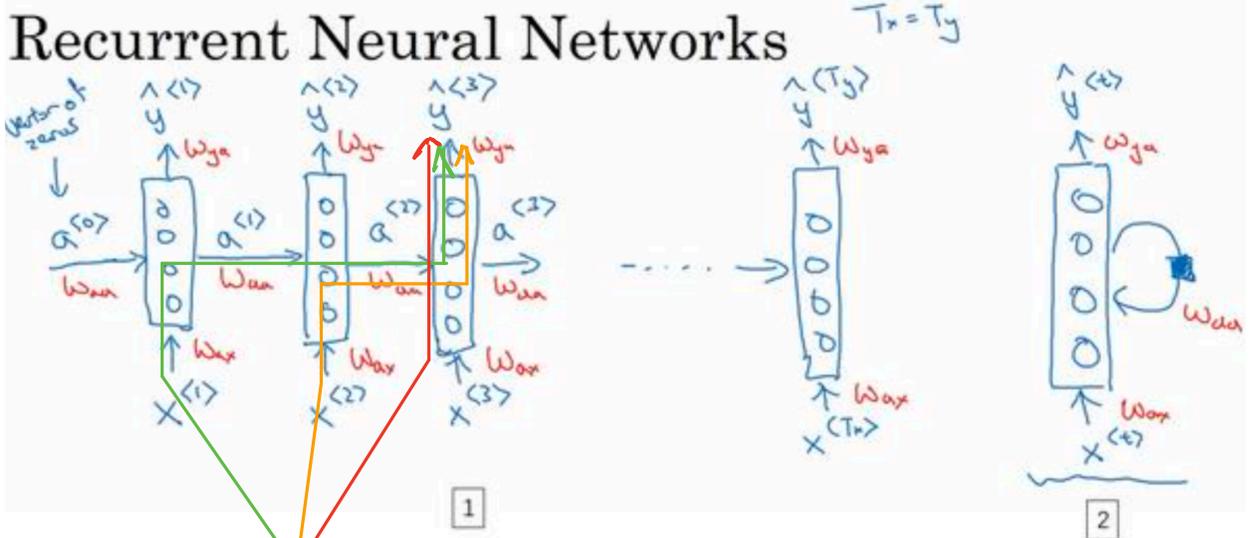
Input: $x^{(1)}, x^{(2)}, \dots, x^{(i)}$

Output: $y^{(1)}, y^{(2)}, \dots, y^{(j)}$

Vocabulary: [a, Aaron, ..., and, ..., hary, ..., potter, ..., Zulu]
 Index: 1 2 ... 367 ... 4075 ... 6830 ... 10,000

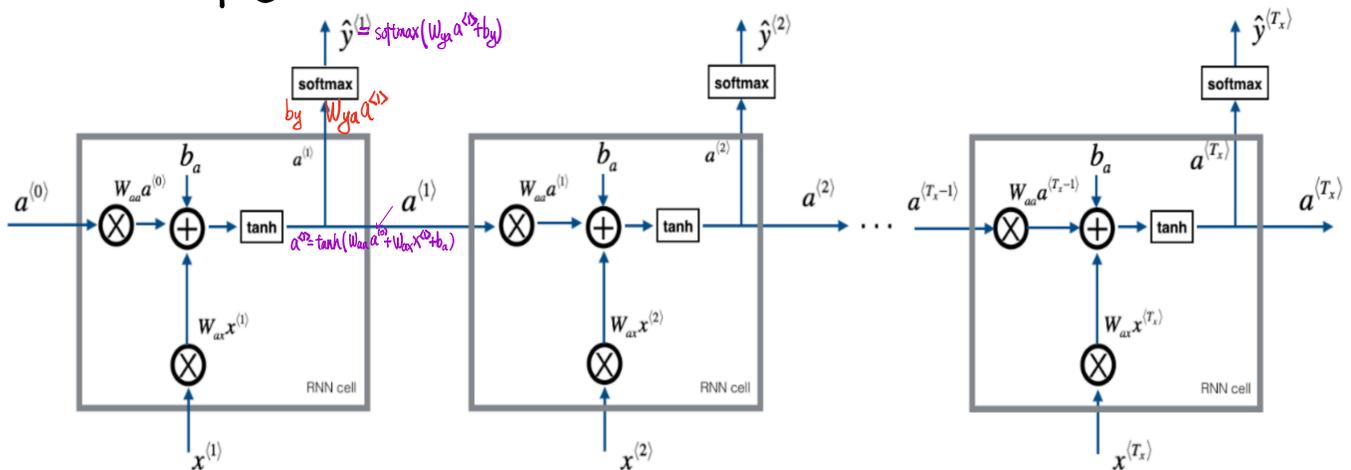
One-hot vector to represent words: lots of zero and an one.

Recurrent Neural Network Model:



$x^{(1)}, x^{(2)}, x^{(3)}$ both effect on $y^{(3)}$, BUT without other nodes X
 \implies Solution: BRNN (Bidirection)

Forward Propagation:

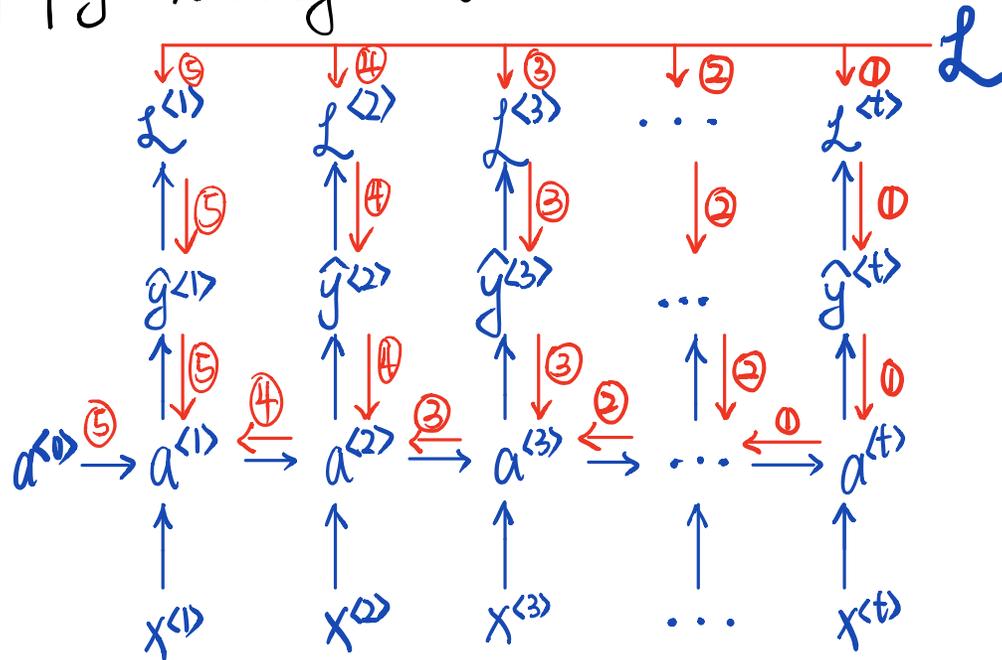


$$a^{(i)} = g_1(W_{aa} a^{(i-1)} + W_{ax} x^{(i)} + b_a)$$

$$y^{(i)} = g_2(W_{ya} a^{(i)} + b_y)$$

→ only one!

Backpropagation through time:



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

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

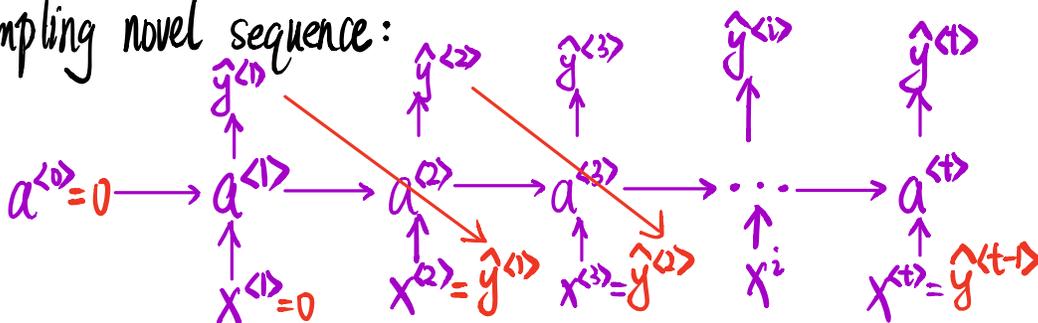
Language model

Language model: decide better sentence among context, according to probability.

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

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

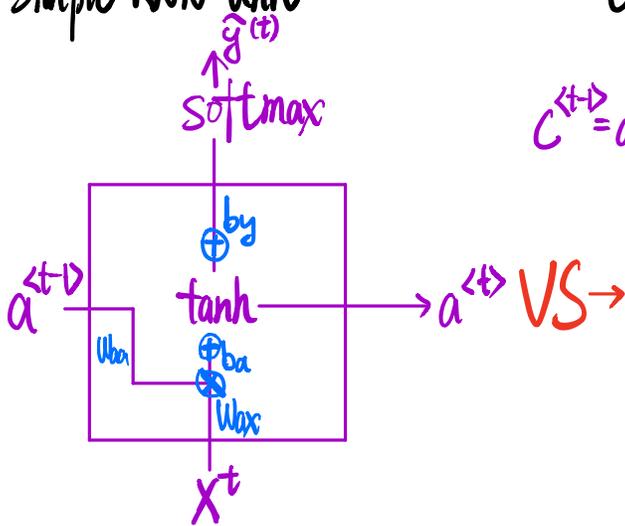
Sampling novel sequence:



Vanishing gradients with RNNs:

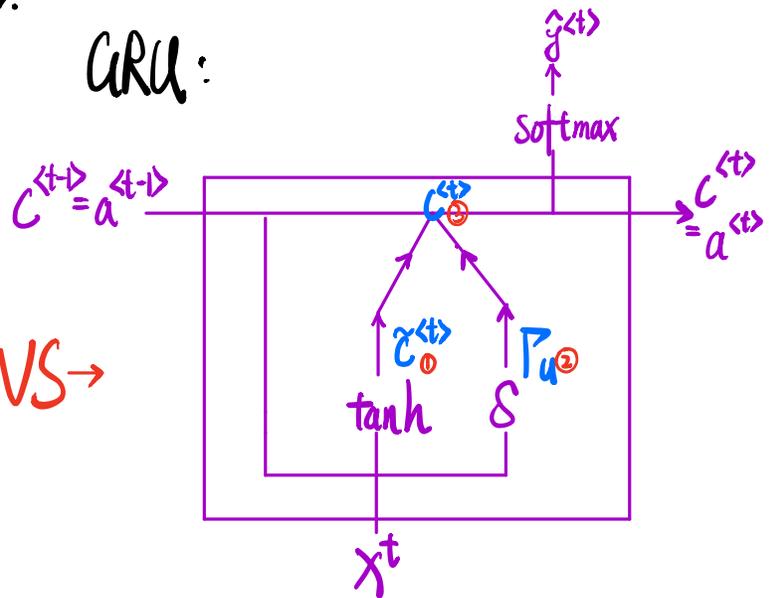
GRU: Gated Recurrent Unit.

Simple RNN Unit



$$a^{(t)} = \tanh(W_a [a^{(t-1)}, x^{(t)}] + b_a)$$

GRU:



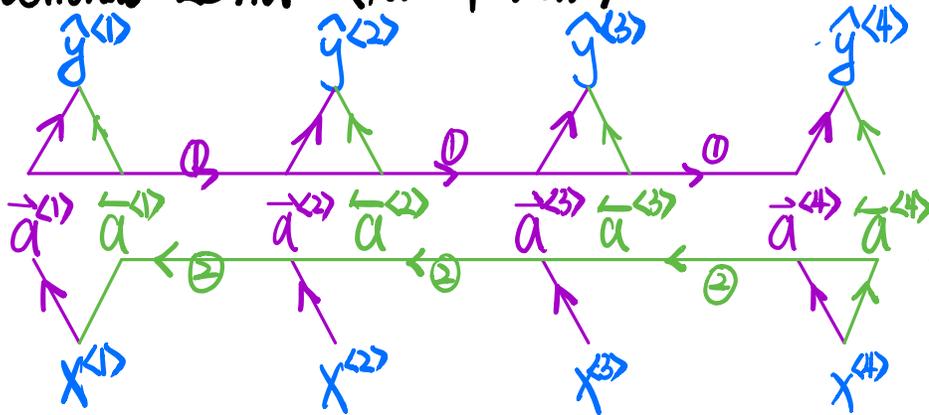
$$\text{Let } \textcircled{1} \tilde{c}^{(t)} = \tanh(W_c [c^{(t-1)}, x^{(t)}] + b_c)$$

$$\textcircled{2} r_u = \sigma(W_u [c^{(t-1)}, x^{(t)}] + b_u)$$

$$\Rightarrow \textcircled{3} c^{(t)} = r_u * \tilde{c}^{(t)} + (1 - r_u) * c^{(t-1)}$$

LSTM: Previous blog contains this part.

Bidirectional LSTM: (For 4 word)



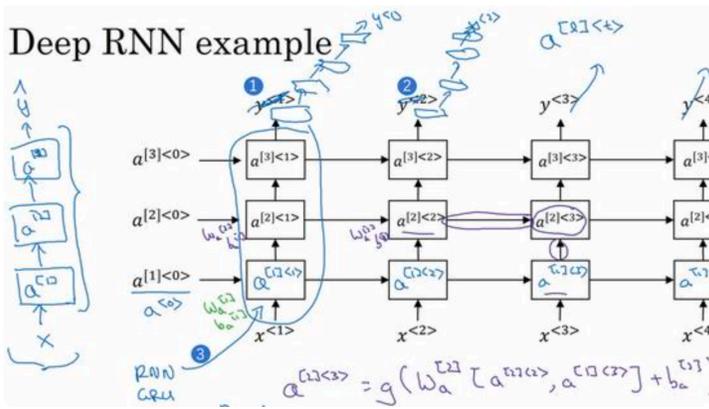
start from a^{(4)}

$$* \hat{y}^{(t)} = g(W_y [a^{(t)}, \overleftarrow{a}^{(t)}] + b_y)$$

Deep RNNs:

something like this:

Deep RNN example



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

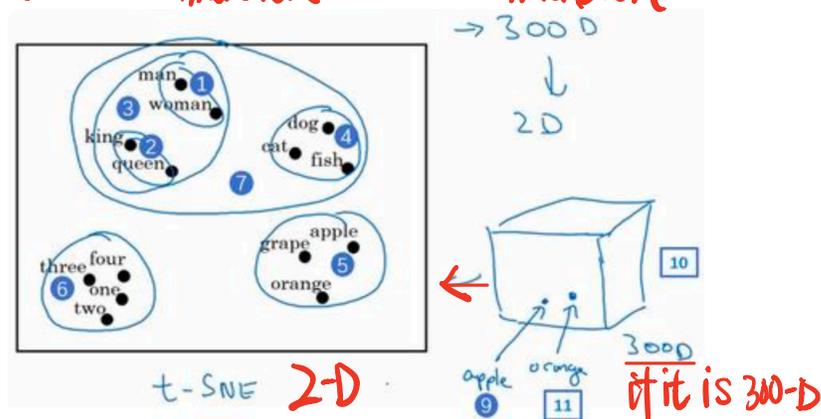
Natural Language Processing and Word Embedding:

Word Representation:

① Featurized representation: word embedding (to replace one-hot)

	Man	Women	King	Queen	Apple	Orange
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
⋮	⋮	⋮	⋮	⋮	⋮	⋮

② Visualizing word embedding: t-SNE Algorithm.
From 300-Dimension to 2-Dimension



Properties of Word Embeddings:

① Let $e_{man} = [-1, 0.01, 0.03, 0.09]^T$, $e_{woman} = [1, 0.02, 0.02, 0.09]^T$

$e_{king} = [-0.95, 0.93, 0.70, 0.02]^T$, $e_{queen} = [0.97, 0.95, 0.69, 0.01]^T$

then $e_{man} - e_{woman} \approx [-2, 0, 0, 0]^T$

& $e_{king} - e_{queen} \approx [-2, 0, 0, 0]^T$

Analogical reasoning: $e_{man} - e_{woman} \approx e_{king} - e_{?}$

So let us define cosine similarity:

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

Embedding Matrix:

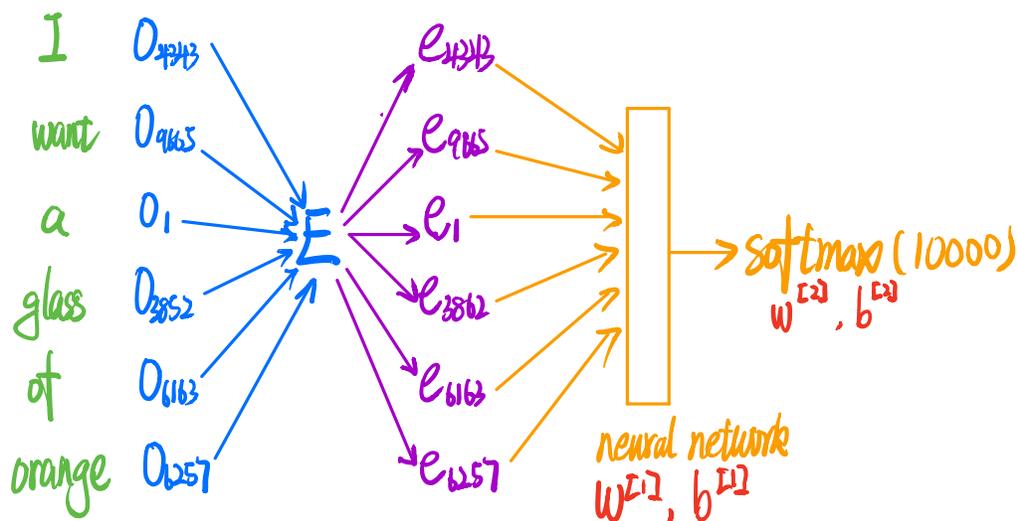
$E(\text{Dimension}, \text{Word_num})$ like $E(300, 10000)$

Learning Word Embedding:

Neural Language Model:

I want a glass of orange _____?

index: 4343 9665 1 3852 6163 6257



Other context / target pairs:

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

Context: $\left\{ \begin{array}{l} \text{Last 4 words.} \\ \text{Last 1 word.} \\ \text{Nearby 1 word.} \end{array} \right.$

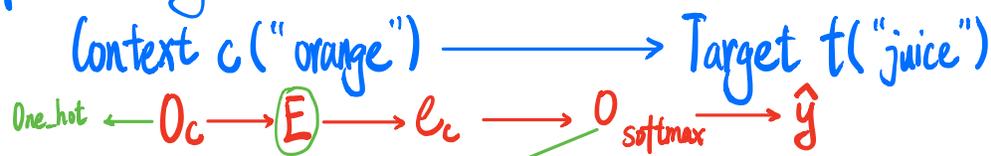
Word2Vec:

① Skip-Grams Model

abstract context: randomly around target word.

match target word.

if Vocabulary Size = 10,000 k



We need to minimize the softmax loss function:

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

$$L(\hat{y}, y) = - \sum_{i=1}^{10000} y_i \log \hat{y}_i$$

(one-hot)

② Summary:
 { CBOW: origin sentence \rightarrow target word
 | Skip-Gram: target word \rightarrow origin sentence

Negative Sampling:

New Problem: Give a pair of words (like orange and juice), let us to predict whether they are context-target or not.

{ orange - juice - 1 \rightarrow positive sample
 | orange - king - 0 \rightarrow negative sample.

	Input: x	Label: y	
	context word	target	
K	orange - juice	1	$\Rightarrow P(y=1 t,c) = \sigma(\theta_t^T e_c)$
	orange - king	0	
	orange - book	0	
	orange - the	0	
	orange - of	0	

K=5-20 if small datasets

K=2-5 if large datasets

GloVe Word Vectors:

Given X_{ij} is the number of i is in context of j .

After traverse whole corpus. you will found $X_{ij} = X_{ji}$

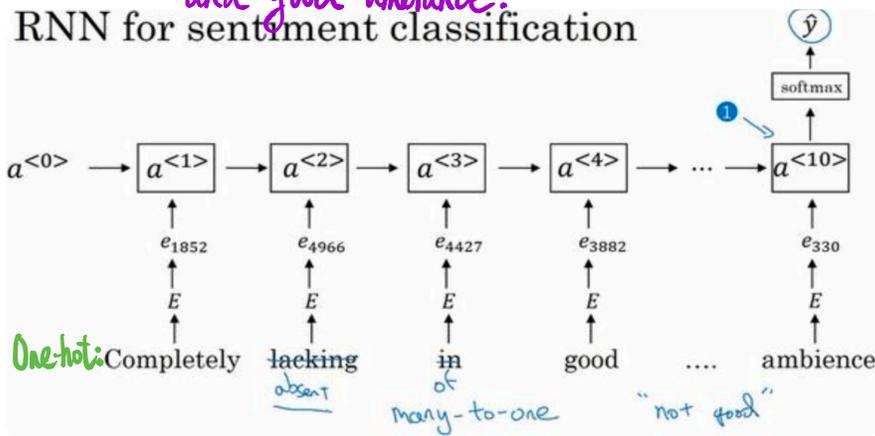
$$\therefore \text{minimize } \sum_{i=1}^{(n)} \sum_{j=1}^{(n)} t(X_{ij}) (\theta_i^T e_j + b_i + b_j' - \log X_{ij})^2$$

The weight function to prevent huge weight of "a, the, ..."

just like t&C, a measurement of how many relationship between them.

Sentiment Classification:

Given a sentence: *completely lacking in good taste, good service and good ambiance.*



Sequence model & Attention mechanism.

Various sequence to sequence architectures:

① Sequence to sequence.

Sequence to sequence model

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$
 Jane visite l'Afrique en septembre

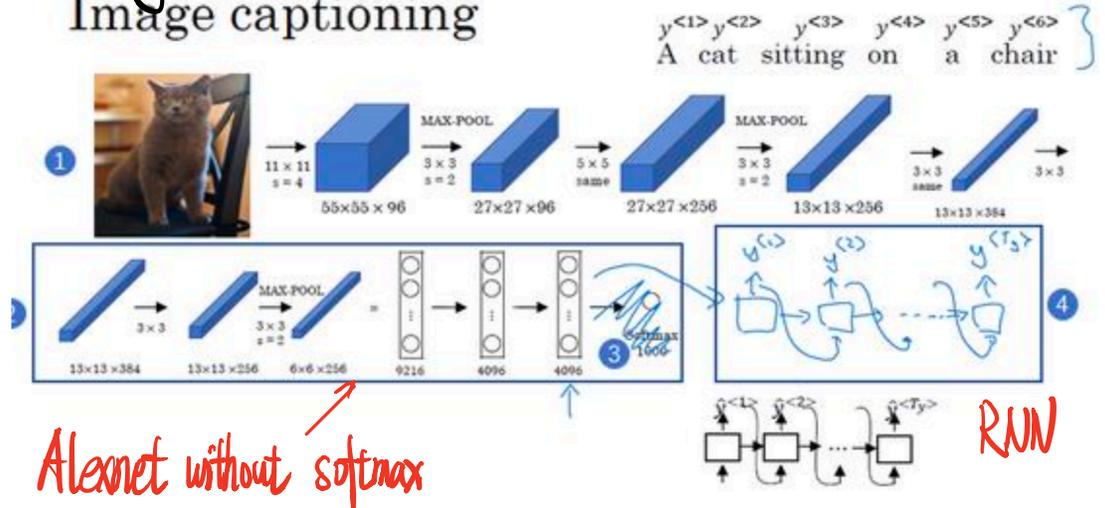
→ Jane is visiting Africa in September.

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

RNN (encoder network) → RNN (decoder network)

② Image captioning

Image captioning



Picking the most likely sentences:

Jane visite l'Afrique en septembre.
(french)

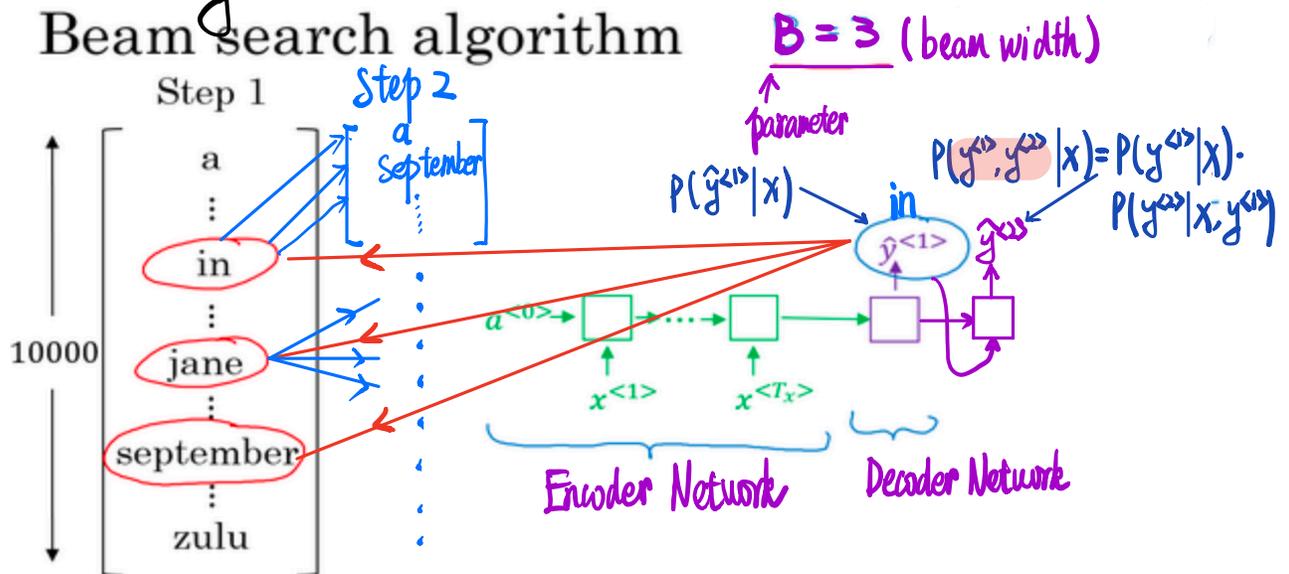
- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.
- In September, Jane will visit Africa.
- Her African friend welcomed Jane in September.

$$\text{argmax}_{y_1, \dots, y_T} P(y_1, \dots, y_T | x)$$

English French

Beam searching:

Beam search algorithm



Output after two step: "in september", "jane is", "jane visits"

Refinements to Beam Search:

① Length normalization:

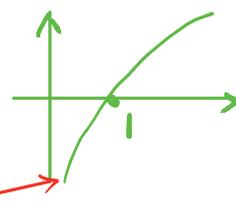
target function can write as:

$$\operatorname{argmax}_y \prod_{t=1}^T P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

log formation:

$$\operatorname{argmax}_y \prod_{t=1}^T \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

to prevent numerical underflow:

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


Error analysis in beam search:

① beam search can seem as an approximate search algorithm or a heuristic search algorithm.

② Process:

- 1' traverse dataset to find errors
- 2' find why: algorithm or model?

Bleu score: (bilingual evaluation under study).

Given a translation from machine and automatically compute a score to measure this translation.

① Bleu score on word:

Evaluating machine translation

French: Le chat est sur le tapis.

→ Reference 1: The cat is on the mat. ← (2 overlaps)

→ Reference 2: There is a cat on the mat. ←

→ MT output: the the the the the the.

✗ Precision: $\frac{7}{7}$ ✓ Modified precision: $\frac{2}{7}$ ← Count("the")

② Bleu score on bigrams:

like a sliding window

$$P = \frac{\sum_{\text{unigram}} \text{Count}_{\text{clip}}(\text{Unigram})}{\sum_{\text{Unigram}} \text{Count}(\text{Unigram})}$$

the cat	2 ←
cat the	1 ←
cat on	1 ←
on the	1 ←
the mat	1 ←

Example: Reference 1: The cat is on the mat. ←

Reference 2: There is a cat on the mat. ←

MT output: The cat the cat on the mat. ←

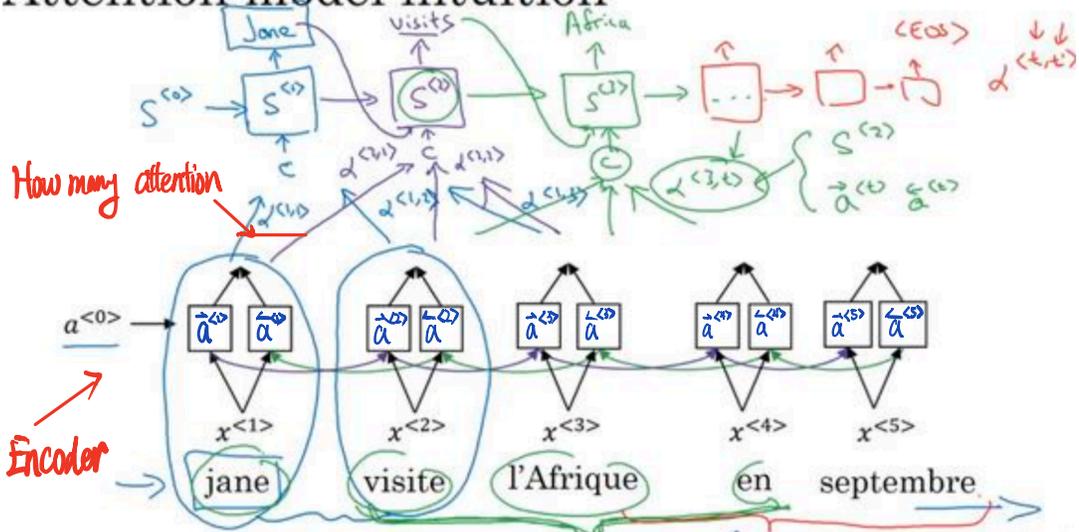
Count	Count _{clip}	
2 ←	1 ←	= $\frac{4}{6}$
1 ←	0	
1 ←	1 ←	
1 ←	1 ←	
1 ←	1 ←	

Attention Model Intuition:

① How human-being work?

Only read a **part** of sentence once a time.

② Attention model intuition



Key point: a set of attention weights: $\alpha^{<i,j>}$
 output | input

③ Compute attention $a^{<t,t'>}$

$a^{<t,t'>}$ = amount of attention $y^{<t>}$ should pay to $x^{<t'>}$

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