

Deep learning

— Neural networks and deep learning

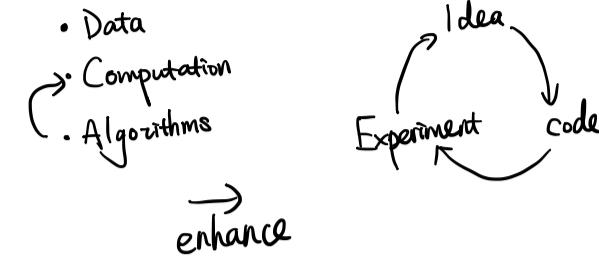
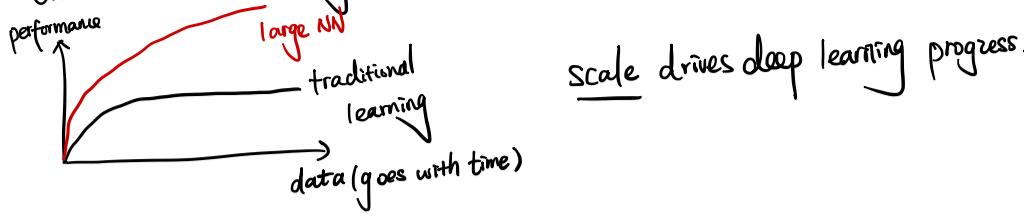
— Coursera
Andrew Ng



Neural networks and deep learning.

Structured data: like an excel list, very well defined meaning of each input

Unstructured data: Images, Audio.



Binary classification:

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

$$X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(n)} \\ | & | & | & | \end{bmatrix} \quad X.\text{shape} = (n, m)$$

$$Y = \begin{bmatrix} y^{(1)}, \dots, y^{(n)} \end{bmatrix} \quad Y.\text{shape} = (1, m)$$

Logistic regression

Given x , want $\hat{y} = P(y=1 | x; \theta)$

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

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

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

Computation graph:

$$\begin{array}{l} a=5 \\ b=3 \\ c=2 \end{array} \rightarrow u=bc \rightarrow v=a+u \rightarrow J=3v$$

In code $\frac{d \text{Find output}}{d \text{var}}$ is represented by 'dvar'

e.g.

$$\begin{array}{l} x_1 \\ w_1 \\ x_2 \\ w_2 \\ b \end{array} \rightarrow z = w_1 x_1 + w_2 x_2 + b \rightarrow a = \sigma(z) \rightarrow J(a, y) \\ dz = da \cdot (1-a)a = a-y \quad da = -\frac{y}{a} + \frac{1-y}{1-a}$$

For m examples, implement average of gradients and update.

SIMD - single instruction multiple data.

In python, Vectorized implementation:

```
for i in range(1000):  
    z = w^T x + b = np.dot(w.T, x) + b  
    A = sigma(z)  
    dZ = A - Y  
    dw = 1/m * x * dZ^T  
    db = 1/m * np.sum(dZ)  
    w := w - alpha * dw  
    b := b - alpha * db
```

Broadcasting in python:

$$cal = A.sum(axis=0) \downarrow^0$$

$$\text{percentage} = 100 * A / cal$$

$$(m, n) \xrightarrow{*} (1, n) \sim (m, n) \\ (m, 1) \sim (m, n) \\ \mathbb{R}$$

use $\text{np.random.rand}(5, 1)$ this will make your code easy to understand

assert(a.shape == (5, 1))

a = a.reshape((5, 1))

Note:

$$P(y|x) = \hat{y}^y (1-\hat{y})^{1-y} \quad \text{max } p \leftarrow \text{maximize likelihood estimation}$$

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

min loss

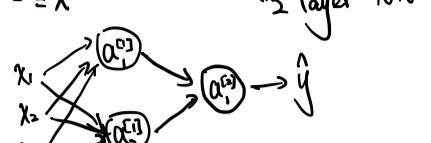
Neural networks and deep learning

Week 3 shallow neural network

$w^{(l)}$ the l layer

$x^{(i)}$ the i the x example

$$a^{[0]} = x$$

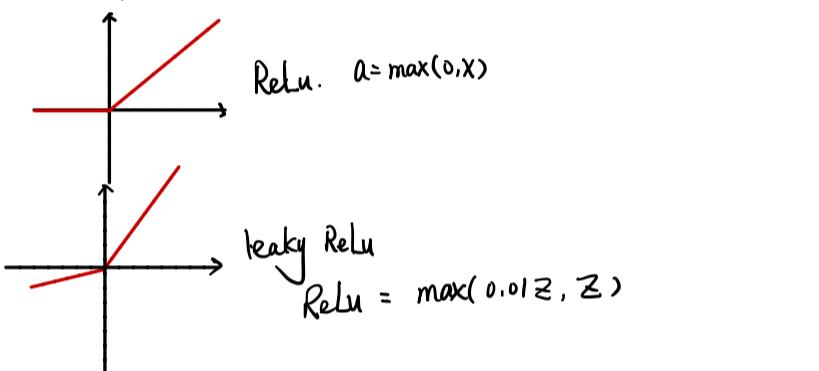
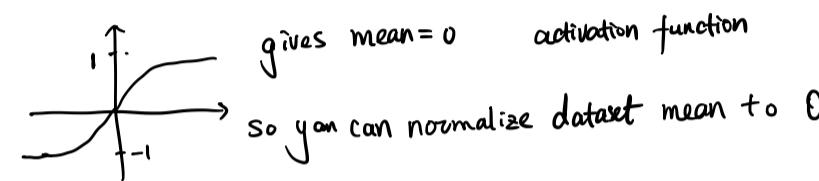


$$\begin{bmatrix} -w_1- \\ -w_2- \\ -w_n- \end{bmatrix} \begin{bmatrix} x_1^{(n)} & \dots & x_m^{(n)} \\ x_2^{(n)} & \dots & x_m^{(n)} \\ x_3^{(n)} & \dots & x_m^{(n)} \end{bmatrix} = \begin{bmatrix} z_1^{(n)} & \dots & z_m^{(n)} \\ z_2^{(n)} & \dots & z_m^{(n)} \\ z_n^{(n)} & \dots & z_m^{(n)} \end{bmatrix}$$

horizontal stands for different training example.
vertically stands for different features/nodes

Activation function:

$$\tanh = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad \text{tanh is nearly always better than sigmoid.}$$



rule: if want to output a probability only. use sigmoid.

normally use ReLU because it learns a lot faster than sigmoid/tanh
because the derivative does not go to zero if x very large.

sigmoid < tanh < ReLU ≈ leaky ReLU. way: try them all.
default.

derivative of activation functions:

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

$$\tanh'(x) = 1 - (\tanh(x))^2 \quad \text{tanh}$$

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

gradient descent algorithm in python:

$$dZ^{[2]} = A^{[2]} - Y \rightarrow \text{loss}' \cdot g'(z)$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T} \quad W^{[2]} \Rightarrow A^{[1]} \rightarrow Z^{[2]}$$

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

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g'(Z^{[1]}) \quad x \xrightarrow{W^{[1]} g'(Z^{[1]})} a_1 \xrightarrow{W^{[2]} g'(Z^{[1]})} a_2$$

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

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

Random initialization:

to zero will not work, break symmetry.

$$W^{[i]} = \text{np.random.randn}(z, z) * 0.01$$

$$b^{[i]} = \text{np.zeros}(z, 1)$$

Week 4 deep neural network.

Input is not a layer.

$$\begin{aligned} l &= 4 \text{ (layers)} \\ n^{[l]} &= \# \text{ units in layer } l \\ a^{[l]} &= \text{activations in layer } l \\ a^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned}$$

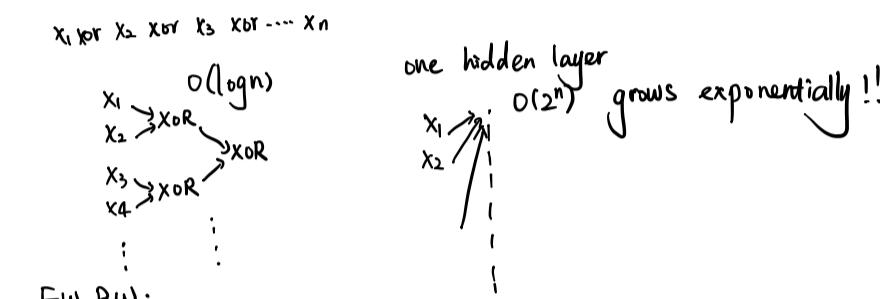
$$\begin{aligned} w^{[l]} &= \text{weight for } z^{[l]} \\ n^{[0]} &= n_x \quad x = a^{[0]} \\ a^{[l]} &= y \end{aligned}$$

$$\begin{aligned} z^{[l]} &= w^{[l]} a^{[l-1]} + b^{[l]} \quad w^{[l]} : (n^{[l]}, n^{[l-1]}) \\ a^{[l]} &= g^{[l]}(z^{[l]}) \quad b^{[l]} : (n^{[l]}, 1) = d b^{[l]} \\ z^{[l]}, A^{[l]} &: (n^{[l]}, m) \end{aligned}$$

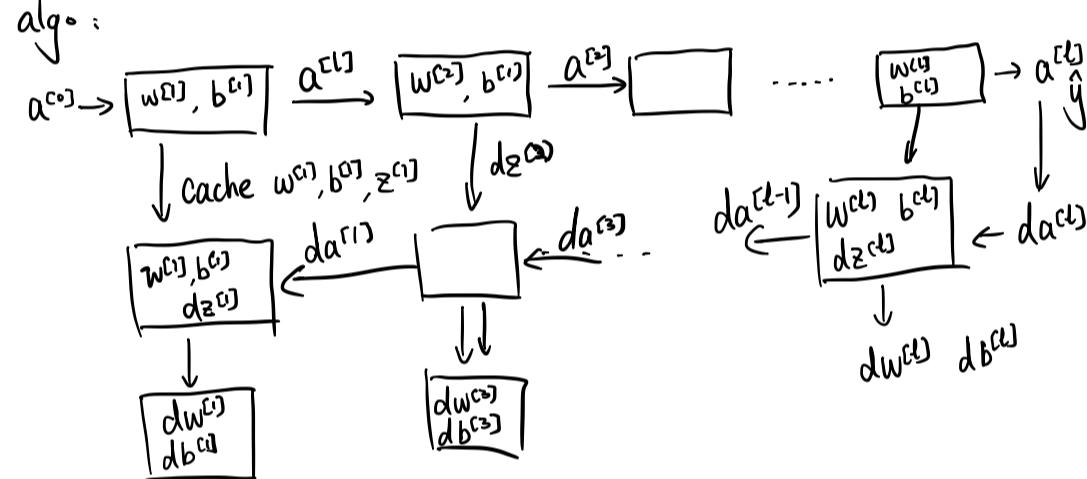
use a for loop to loop through all code.

getting matrix dimension right.

Circuit Theory and deep learning



FW BW.



Implementation:

Forward prop

Input $a^{[l-1]}$
Output $a^{[l]}$, cache $[z^{[l]}]$

$$z^{[l]} = w^{[l]} \cdot A^{[l-1]} + b^{[l]}$$

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

Backward prop for layer l

Input $da^{[l]}$
Output $da^{[l-1]}, d w^{[l]}, d b^{[l]}$

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

$$d w^{[l]} = \frac{1}{m} d z^{[l]} \cdot A^{[l-1] T}$$

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

$$d A^{[l-1]} = W^{[l] T} \cdot d z^{[l]}$$

if sigmoid function at last layer.

$$d a^{[l]} = \left[\frac{-y^{[l]}}{a} + \frac{(1-y)^{[l]}}{1-a} \right]$$

Applied deep learning is a very empirical process.

Improving Deep Neural Networks:

Hyperparameter tuning.

Regularization and

Optimization

— Andrew Ng

Setting up Machine learning application.

60% training set 20% valid set 20% test set small dataset

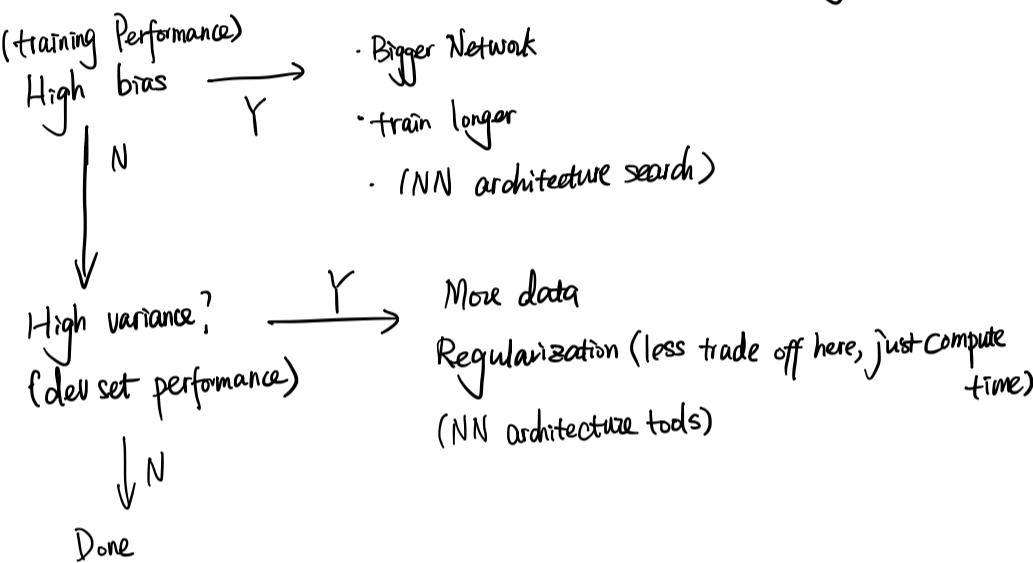
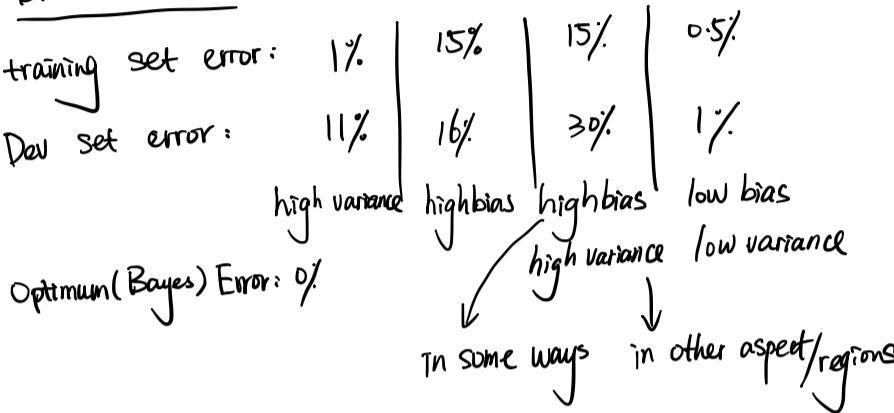
large dataset

1,000,000 train 10,000 valid 10,000 test

make sure dev/test set from same distribution.

training data can be download from internet.

Bias variance:



Regularization:

$$L_2 \text{ regularization: } \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

$$L_1 \text{ regularization: } \frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1, \quad w \text{ will be sparse.}$$

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

$$\|w^{(l)}\|_F^2 = \sum_{i=1}^{n^{(l-1)}} \sum_{j=1}^{n^{(l)}} (w_{ij}^{(l)})^2 \quad \text{"Frobenius norm"}$$

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

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

regularization: ① use all the weights, each weights
② activation will be more linear,

more linear assumption.



Dropout regularization:

Implement dropout:

```
d3 = np.random.rand(a3.shape[0], a3.shape[1]) < keep_prob
```

```
a3 = np.multiply(a3, d3)
```

```
a3 /= keep_prob.
```

also zero out at back prob

No drop out at test time!

Intuition: Can't rely on any feature. so have to spread out weights.
→ computer vision, apply a lot. because have a lot of inputs and each works.

When plotting $J(\theta)$ first set $\text{keep_dim} = 1$ to find J is decreasing because dropout may affect J .

other techniques:

1. Data augmentation

2. early stopping



bad side of early stopping:

Orthogonalization first:

- optimize cost function - Gradient descent / momentum
- overfitting solving - R_2 regularization

but early stopping make us to do these two steps together.

Or we can set a good R_2 regularization parameter λ .

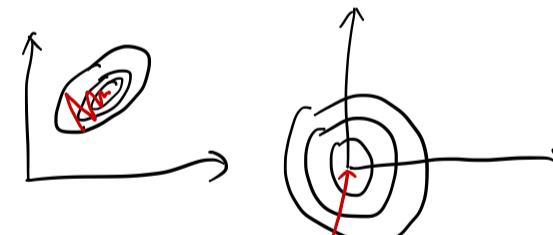
So we can train as much as we want.

but need to try several times.

faster your training:

Normalizing training set. store μ_{train} , σ_{train}
and apply it on test set.

Why normalize:



Vanishing/exploding gradient.

when initialize W , random initialize weights.

$[W]^L = X$ grows exponentially.

Random initialization:

$$W^{(l)} = \text{np.random.randn(shape)} + \text{np.sqrt}\left(\frac{2}{n^{(l-1)}}\right) \text{ReLU (Kaiming he)}$$
$$\text{np.sqrt}\left(\frac{1}{n^{(l-1)}}\right) \tanh (\text{Xavier initialization})$$
$$\text{np.sqrt}\left(\frac{2}{n^{(l-1)} + n^{(l-1)}}\right) \text{ (also can)}$$

checking derivative: numerical checking.

$$\frac{f(\theta + \varepsilon) - f(\theta - \varepsilon)}{2\varepsilon} = g(\theta)$$

turn $[w_1, b_1, \dots, w_n, b_n] \rightarrow [\theta_1, \theta_2, \dots, \theta_n]$

for each i :

$$\frac{d\theta_i^{approx}}{d\theta_i} = \frac{J(\theta_1, \dots, \theta_{i+\varepsilon}) - J(\theta_1, \dots, \theta_{i-\varepsilon}, \dots)}{2\varepsilon}$$

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

$$\begin{aligned} & \varepsilon = 10^{-7} \approx 10^{-7} - \text{correct!} \\ & 10^{-5} \\ & 10^{-3} - \text{worry.} \end{aligned}$$

notes:

- Don't use in training - only to debug
- If algorithm fails, debug.
- Remember regularization
- Doesn't work with dropout keep dims = 1.
- Run at random initialization, check again after some training (less common).

week2: optimization algorithm:

Batch vs mini batch gradient descent.

$x^{fit}, y^{fit} \Rightarrow$ first + mini batch.

minibatch - run faster

minibatch \rightarrow stochastic gradient descent: longer training time.
update parameter each time. make progress everytime.

choose minibatch size:

If small training set ≤ 2000 : use batch gradient descent.

Typical minibatch sizes: 64, 128, 256, 512,

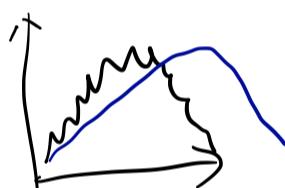
depend on GPU.

Exponentially weighted averages

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

$$\text{average over steps: } \frac{1}{1-\beta}$$

β bigger, the less adaptive. got latent changes



$$V_{100} = (1-\beta)\beta \theta_{99} + (1-\beta)\beta^2 \theta_{98} + \dots + (1-\beta)\theta_{100}$$

$$(1-\beta) + (1-\beta)\beta + (1-\beta)\beta^2 + \dots \approx 1.$$

$$0.9^{10} \approx \frac{1}{e} \quad 0.98^{50} \approx \frac{1}{e}$$

memory saving.

$$V_\theta := 0$$

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

:

:

bias correction in initial phase:

$$V_\theta = \frac{(1-\beta)V_{\theta-1} + \theta_t}{1-\beta^t} \quad (\text{usually do not bother do it})$$

Gradient descent with momentum.

can get off local min.

On iteration t:

Compute dW, db , on current mini batch

$$V_{dw} = \beta V_{dw} + (1-\beta) dW \quad \rightarrow \text{have this term}$$

all need to rescale
& everytime you
change β

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

$$W = W - \alpha V_{dw}, b = b - \alpha V_{db}$$

RMS Prop:

On iteration t:

Compute dW, db on current minibatch.

$$S_{dw} = \beta_1 S_{dw} + (1-\beta_1) dW^2 \leftarrow \text{element wise}$$

$$S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \leftarrow \text{element wise}$$

$$W := W - \alpha \frac{dW}{\sqrt{S_{dw} + \epsilon}}, b := b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}} \rightarrow 10^{-8}$$

slow learning rate
when gradient very

large

for numeric stability.

Adam optimization: Adaptive moment estimation

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

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

On iteration t,

Compute dW, db using correct minibatch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dW, V_{db} = \beta_2 V_{db} + (1-\beta_2) db$$

$$S_{dw} = \beta_1 S_{dw} + (1-\beta_1) dW^2, S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2$$

$$V_{dw}^{\text{corrected}} = V_{dw} / (1-\beta_1^t), V_{db}^{\text{corrected}} = V_{db} / (1-\beta_2^t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1-\beta_1^t), S_{db}^{\text{corrected}} = S_{db} / (1-\beta_2^t)$$

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

α : need to tune

$$\beta_1 = 0.9 \quad (dW)$$

$$\beta_2 = 0.999 \quad (dW^2)$$

$$\epsilon = 10^{-8}$$

learning rate decay:

$$1. \quad \alpha = \frac{1}{1 + \text{decay rate} * \text{epoch num}} \alpha_0$$

$$2. \quad \alpha = 0.95^{\text{epoch num}} \alpha_0$$

$$3. \quad \alpha = \frac{k}{\sqrt{\text{ep-num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \alpha_0$$

$$4. \quad \alpha = \frac{1}{t}$$

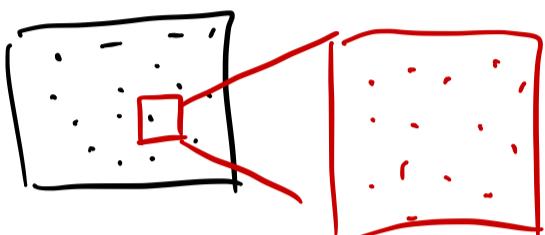
5. Manual decay.

Local optima in neural network:

most saddle point.

hyperparameters tuning:

- α
- $\beta \sim 0.9$
- # hidden units
- # minibatch size
- layers
- learning rate decay.
- random parameter, do not use grid..



Appropriate scale for hyperparameters

$$r = -4 * \text{np.random.rand}()$$

$$\alpha = 10^r \quad \leftarrow \log \text{ scale.}$$

$$\beta = 0.9 \dots 0.9999$$

$$(-\beta = 0.1 \dots 0.001)$$

focus on.

One model



Batch normalization:

usually normalize $z^{(i)}$

Implement:

Given some intermediate values in NN.

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

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

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

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

learn γ & β .
use nonlinearity better

$$\text{Parameters} = [w^{(0)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}]$$

$$[\beta^{(0)}, \gamma^{(0)}, \dots, \beta^{(L)}, \gamma^{(L)}]$$

tf.nn.batch_normalization.
if use batch norm, we can ignore $b^{(l)}$.

Implementing gradient descent

for $t=1 \dots \text{num Mini batches.}$

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

In each hidden layer, use BN replace $z^{(l)}$ with $\tilde{z}^{(l)}$

Use back prop to find $dW^{(l)}, dB^{(l)}, d\gamma^{(l)}$

update parameters

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

$$\beta^{(l)} := \beta^{(l)} - \alpha dB^{(l)}$$

$$\gamma^{(l)} := \gamma^{(l)} - \alpha d\gamma^{(l)}$$

This works with adam/momentum as well.

Why batch norm work?

(1) combat covariance shift of each layer, each neural network can still perform task independently even though the input from previous layer change.

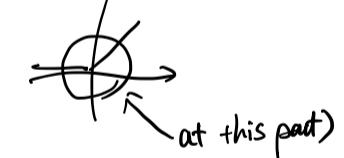
(2) regularization effect:
each minibatch is scaled by mean/variance computed on just that minibatch.

This add noise to value $z^{(l)}$ within that minibatch.
So similar to dropout. it slightly force nn not to rely on individual models.

(bigger minibatch \rightarrow smaller regularization)

(3) make nn to make better use of activation.

morvan taught instead of Ng.



Batch Norm at test time:

keep track of moving average μ and σ^2 in training time and use it at test time.

multi class classification:

softmax layer:

activation function.

$$t = e^{(z^{[l]})}$$

$$a^{[l]} = \frac{e^{z^{[l]}}}{\sum_i t_i} \rightarrow \text{if only one layer, softmax decision boundary will be linear.}$$

softmax is soft hard max.

logistic is one example of softmax.

$$\mathcal{L}(\hat{y}, y) = - \sum_{j=1}^{\# \text{category}} y_j \log \hat{y}_j$$

$$J() = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}, y) \quad \# \text{category} = 4.$$

Back prop : $d\mathcal{L}^{[l]} = \hat{y} - y \in 4 \text{ by } 1 \text{ vector.}$

Deep learning frame works:

- Caffe / caffe2 choosing criteria:
 - 1. easy to code (deploy)
 - 2. run time.
 - 3. open source (in the future)
maintenance.
- CNTK
- DL4J
- Keras
- Lasagne
- mxnet
- PaddlePaddle
- Tensorflow
- Theano
- Torch.

Tensorflow:

see code

$+$ $*$ $/$ is overloaded in tensorflow so they are = `tf.add...`
makes your code cleaner.

Structuring machine learning

Project

Andrew Ng.



ML strategy

Orthogonalization:

separately each factor that only will affect one factor/function of algo.
 speed → angle
 driving a car → performance

Chain of ML assumptions

Fit training set well on cost function → bigger network
 ↓ (≈ human level performance) Adam
 Fit dev set well on cost function → regularization
 ↓ bigger training set
 Fit test set well on cost function → bigger dev set
 ↓
 real world performance well. → change dev set or cost function.

Single number evaluation metric. (at start project)

so we can pick out a better solution.

e.g. recall, precision, F1 score.
 speed up iterating.

How to set evaluation metric.

1 Optimizing metric: accuracy

N-1 satisfying metric: running time.

set dev/test sets:

Come from same source. can reflect real world problem. test. you care about.

Size of test set:

set your test set to be big enough to give high confidence in overall performance of your system.

• When to change dev/test sets and metrics.

e.g. pornographic images weight more in loss function.

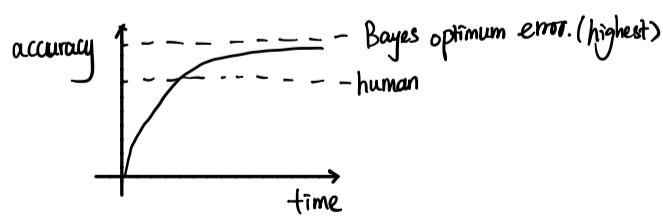
define new metric perform your favourite.

1. set metric to evaluate classifiers (do not worry it is not very accurate, set early)

2. Worry separately about how to do well on this metric.

$$J = \frac{1}{\sum w^{(i)}} \sum_{i=1}^m w^{(i)} L(\hat{y}^{(i)}, y^{(i)}) \quad w^{(i)} = \begin{cases} 1 & \text{norm image} \\ 0 & \text{image you do not want.} \end{cases}$$

Human-level performance.



Why compare to human level performance:

so long as ML is worse than humans:

- Get labelled data from humans
- Gain insight from manual error analysis
- Better analysis on bias/variance.

bias or not?

humans (≈ Bayes)	1%	:	75%	↑ avoidable bias ↓ variance
Training error	8%	:	8%	
Dev error	10%	:	10%	

focus on bias, focus on variance.

Human level as a proxy for Bayes error.

after Surpass human level performance. (team of experts)

- Online advertising
- Product recommendation
- Logistic (predict transit time)
- Loan Approvals.

Structured data,
 Not natural perception
 lots of data

Bayes error (human)

↓ avoidable bias → { bigger model
 train longer optimization algo
 NN architecture/hyperparameter search.

Training error

↓ variance → { more data
 regularization - L2, dropout, data augmentation, etc
 NN architecture, hyperparameter search.

Dev error

machine learning idea:

Flight simulator: train with scenarios that will take time to encounter in real experience

error analysis:

Look at dev examples to evaluate ideas:

Evaluate multiple ideas in parallel.

Ideas for cat detection:

- Fix pictures dogs being recognized as cats
- Fix great cats (lions, panther) misclassification
- blurry images affect performance.

draw a sheet of error!

Image	Dog	Great cats	Blurry	Comments
1	✓			pitbull.
2		✓	✓	
of total	8%	43%	61% (circled)	12%

incorrect labels :

DL algorithms are quite robust to random errors in training set.

it is not robust to systematic errors in training set.

If In test/test/valid set: add another column in sheet of errors
to see how much.

Is it worth while to go and fix incorrect labels:

(whether incorrect labels makes up a large percent of errors)

Correct incorrect dev/test sets:

- apply same process to your dev and test sets to make sure they continue to come from same distribution.
- consider what you got right as well as you got wrong.

(sometimes only look at what you got wrong make you over feature engineering and turn some correct cases to wrong)

- train & dev/test sets now may come from different distribution.

Build your first system quickly, and iterate.

speech recognition example:

ways to improve:

- Noisy background
different noises (car)
- accented
- Far from microphone.

work flow:

① Set up dev/test set and metric

② build initial system quickly

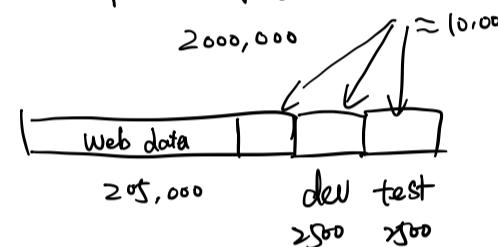
③ Use bias/variance analysis & Error analysis to prioritize next steps.

Mismatched training and dev/test set:

training and testing on different distribution:

Data from other source (big) Data matters (small)

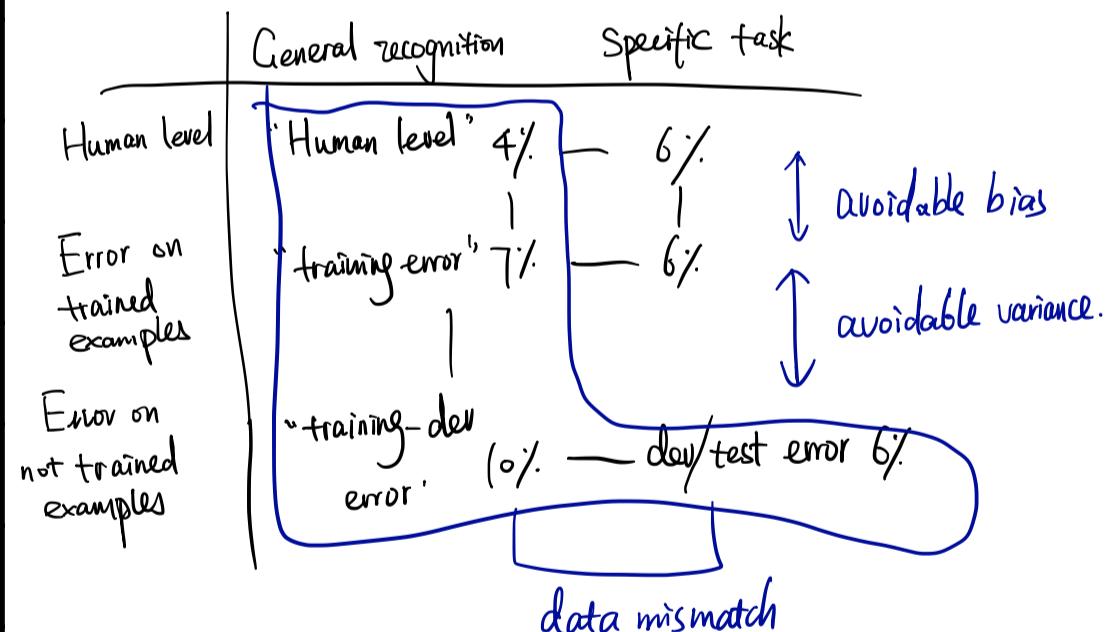
Data from webpage: Data matters:



if they come from different distribution:

training-dev set: A part in training set your training data is test on.
to see whether it is high bias or data mismatch.

Sheet of analysis



Addressing data mismatch:

- carry out manual error analysis to find difference between training & dev/test sets.

- Make training data more similar or recollect data.

Artificial data analysis

(make sure the way you synthesize data, do not overfit)
It is hard to know sometimes.

Transfer training:

When transfer learning make sense:

- Task A and B have same input type.
- A have more data than B.
- A and B shares common small features.

Multiclass learning:

assign multiple labels to one example.

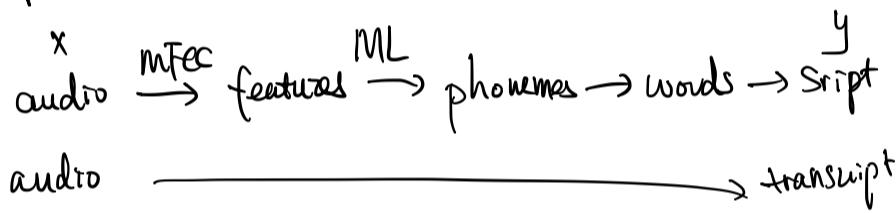
sometimes this performs better than single class.

When multiclass make sense:

- training on a sets of task that could benefit from shared - low level features.
- Usually amount of data you have for each task is similar.
- can train big enough network do well on all tasks.
(usually used in object detection)

End-to-End deep learning:

speech recognition example:



training size matters a lot.

Pros and cons of end-to-end deep learning

Pros:

- let data speak.
- Less hand-designing of component needed.

Cons:

- need a large data. (have enough data for subtask)
- cannot add hand designed features.

Should I try end-to-end:

key question: do you have sufficient data to learn
a function of complexity needed to map x to y?

Convolutional NNs

— Andrew Ng



CNN:

tf.nn.conv2d

keras : Conv2D

Convolutional operation: \star (cross-correlation)

Edge detector:

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

sobel filter for vertical line:
if there is a line: extreme larger or extrem negative

Padding:

- shrink output
- throw away info from edge (there is no overlapping at edge)
- padding zeros.

valid convolution (no padding):

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

same: pad to make sure output same

$$\frac{p-1}{2}, p \text{ is usually odd.}$$

Stride convolution:

$$n \times n \star f \times f \text{ strides } \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor$$

Convolution over volume:

e.g.

$$n \times n \times 3 \star f \times f \times 3 \rightarrow n-f+1 \star n-f+1 \times n_c' \quad \text{\# filters.}$$

CNN layer:

$$\text{ReLU}(\boxed{} + b)$$

notation:

$f^{[l]}$ = filter size,

$p^{[l]}$ = padding size

$s^{[l]}$ stride

$$\text{Input: } n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$$

$$\text{Output: } n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$$

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

Each filter: $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$

order is not
consistent in different
networks

Activations: $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$$A = m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$$

$$\text{weights: } f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]}$$

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

Pooling layer:

Max pooling

$$\left\lfloor \frac{n_H-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W-f}{s} + 1 \right\rfloor$$

$$5 \times 5 \times n_C \rightarrow 3 \times 3 \times \frac{n_C}{2}$$

↓ does not change
No parameter to learn.

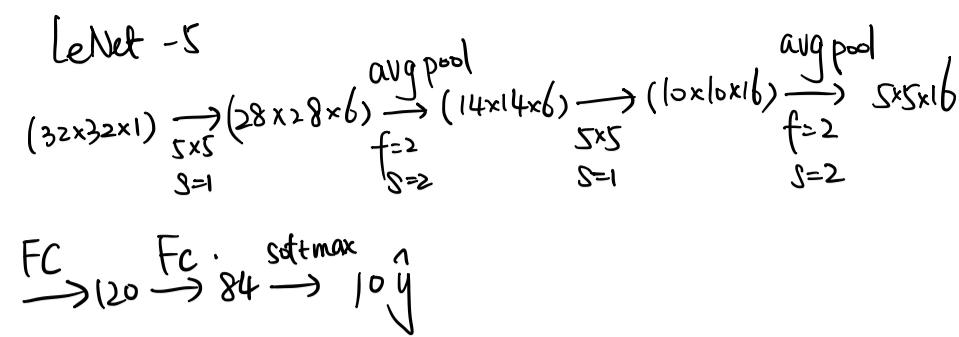
see hyperparameters works for others will also work for you.

why CNN work:

Parameter sharing

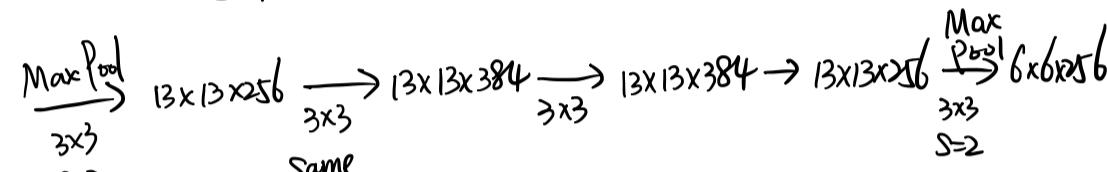
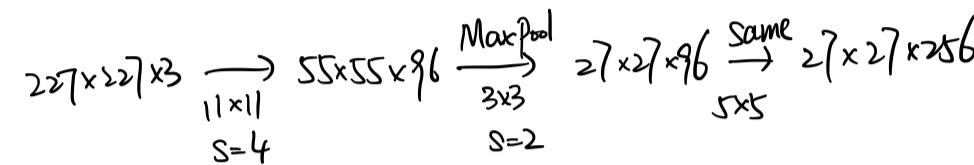
Sparsity of connections

Classic Networks:



$n_H, n_W \downarrow n_C \uparrow$

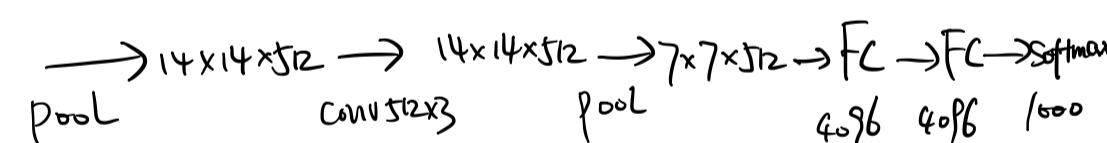
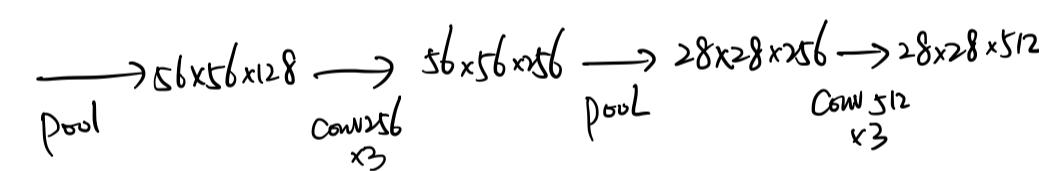
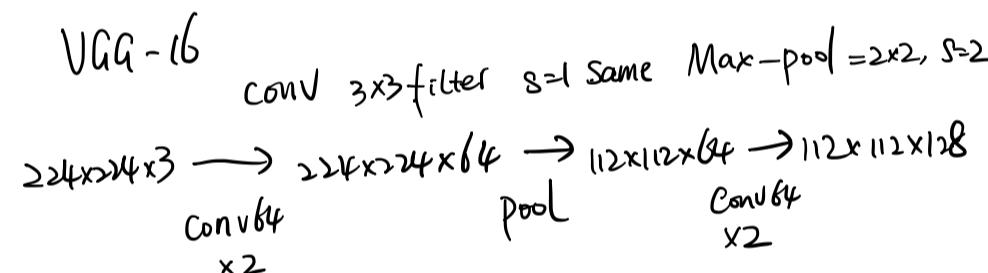
conv pool conv pool fc fc output
activation (sigmoid) at that time.



$$= 9216 \xrightarrow{\text{FC}} 4096 \xrightarrow{\text{FC}} 4096 \xrightarrow{\text{softmax}} 1000$$

— Similar to LeNet, but much bigger $\sim 60M$ params

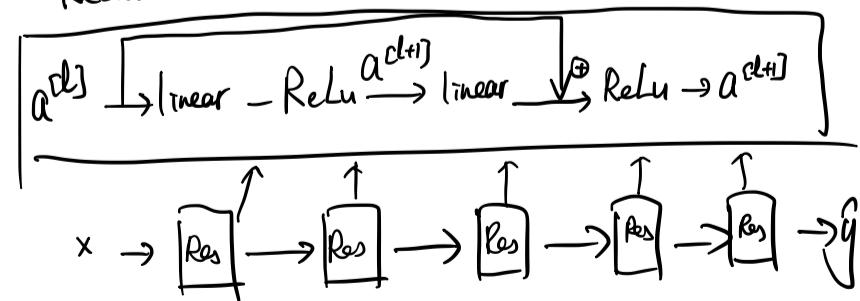
— ReLU



$\sim 138M$.

ResNets:

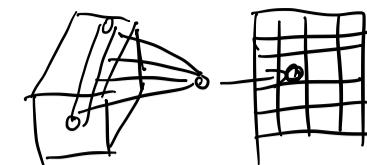
Residual blocks:



Identity function is easier for Residual block to learn.

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

1×1 Conv do?

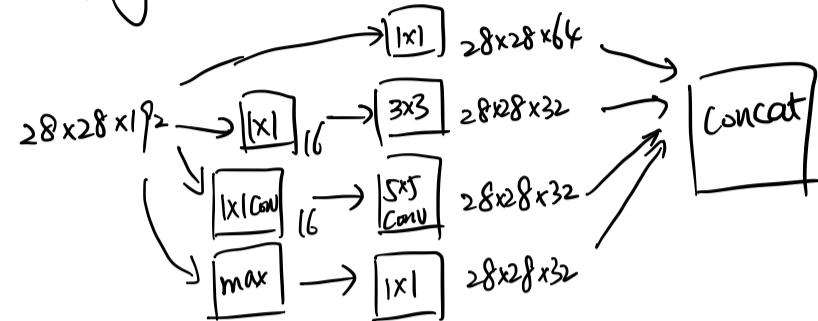


network in networks.

add linearity

Inception network:

using 1×1 convolution to reduce parameters to $\frac{1}{10} M$.



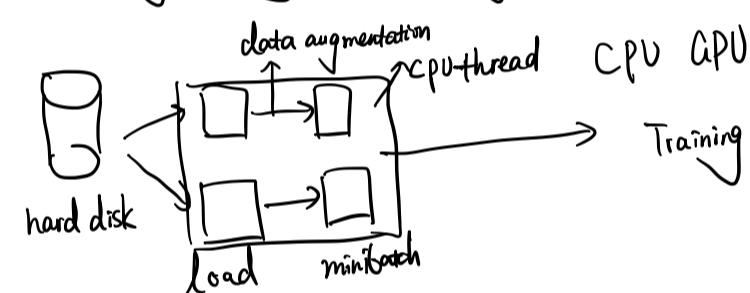
(also side output, to regularize and make sure intermediate layers do not go wrong too much)

Transfer learning:

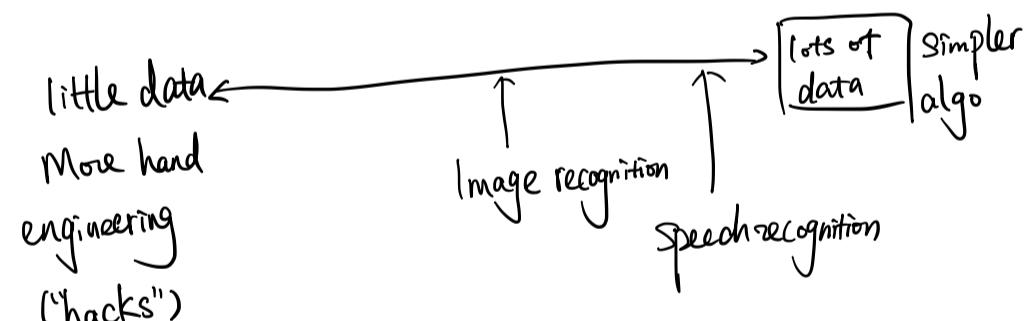
if data really large: can freeze more modules.

Data augmentation:

color shifting, mirroring. "PCA color augmentation".



Data vs hand-engineering



Two source of knowledge:

① labelling

② complex nn and hand engineering feature

Tips on breaking benchmarks:

• Ensembling: train several network independently and average their output.

• Multi crop at test time.

test data augmentation.

object localization:

output

$$y = \begin{cases} p_x \\ p_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{cases}$$

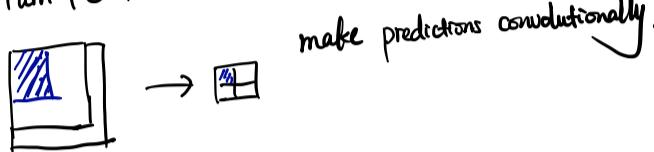
Same nn.

landmark detection.

object detection:

sliding window detection.

Turn FC into convolutional layers:



YOLO algorithm:

assign a classification label [7] to one cell.
(b_h, b_w can be bigger than 1)

Intersection over Union (IoU) : "correct if $\frac{\text{detected area has object}}{\text{cell area}} > 0.5$.

Non max suppression:

① discard < 0.5 p.

while there is cell:

② highlight the biggest probability

③ discard ones have high overlap with the first one.

anchor box: have two part.



choose window box
by kmeans to cover all objects?

Region proposal:

segmentation algorithm

R-CNN: Propose region, Classify proposed region one at a time.

Output label + Bounding box

Fast-Rcnn: Propose region, Use convolution implementation of sliding windows to classify all proposed regions.

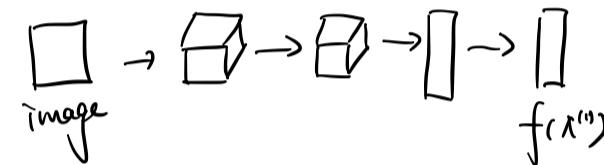
Faster-Rcnn: Use convolutional neural network to propose regions.

Face verification:

one shot learning: have only 1 training image.

learn $d(\text{image}_1, \text{image}_2) \leq T$. is same body.

Siamese network



If $x^{(i)}, x^{(j)}$ is same person, $\|f(x^{(i)}) - f(x^{(j)})\|$ is small.

Triplet loss function:

Given 3 images Anchor, positive, negative.

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

margin.

training set. 10k pictures of 1k person.

choose triplets that are "hard" to train on. (paper: FaceNet).

$$J = \sigma \left(\sum_{k=1}^{128} w_k (||f(x^{(i)})_k - f(x^{(j)})_k|| + b) \right) \quad \begin{array}{l} \circ \text{ True} \\ \downarrow \text{ or } \\ \circ \text{ False} \end{array}$$

$$\frac{(f(x^{(i)})_k - f(x^{(j)})_k)^2}{(x^{(i)})_k + f(x^{(j)})_k}$$

Neural transfer:

[Zeiler and Fergus 2013. Visualizing and understanding CNNs]

ways to visualize each layers of CNN.

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

[Gatys et al. 2015. A neural algorithm of artistic style]

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

Content C Style S
 ↓
 Generate image G

$$J_{\text{content}}(G, C) = \frac{1}{2} \|a^{[L]}(C) - a^{[L]}(G)\|^2 \text{ some middle layer.}$$

Style matrix:

$$a_{ijk}^{[L]} = \text{activation at } (ijk) \quad A^{[L]} = n_c^{[L]} \times n_c^{[L]} \quad J_{\text{style}}(S, G) = \frac{1}{2n_H^{[L]} n_W^{[L]} n_c^{[L]}} \sum_{k,k'} \sum_{i,i'} (a_{ik}^{[L]} - a_{ik'}^{[L]})$$

$$G^{[L]} = \sum_{i=1}^{n_H^{[L]}} \sum_{j=1}^{n_W^{[L]}} a_{ijk}^{[L]} a_{ijk'}^{[L]} \text{ covariance.}$$

convolution in 1D.

ECG signal using 1D filter.

3D conv. on CT scan.

movie data using 3D filter.

Sequence model

- Andrew Ng coursera



RNNs

Notation

$x^{(i),t}$ The t -th word in training sample i

$y^{(i),t}$

$T_x^{(i)}$ The length of i -th training example.

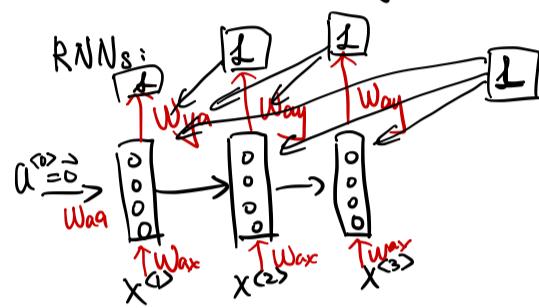
$T_y^{(i)}$

<UNK> for word not in the model.

Why not standard network:

Problems:

Inputs, outputs can input use various length.
don't share weights across different position of text.



$$a^{<t>} = \vec{0} \quad | \quad a^{<t>} = g_1(W_a a^{<t-1>} + W_{ax} x^{<t>} + b_a) \leftarrow \text{tanh . ReLu}$$

$$g^{<t>} = g_2(W_y a^{<t>} + b_y) \leftarrow \text{can be sigmoid depend on output.}$$

$$a^{<t>} = g_1(W_a [a^{<t-1>}, x^{<t>}] + b_a)$$

$$g^{<t>} = g_2(W_y a^{<t>} + b_y) \leftarrow \text{sigmoid if output 0/1}$$

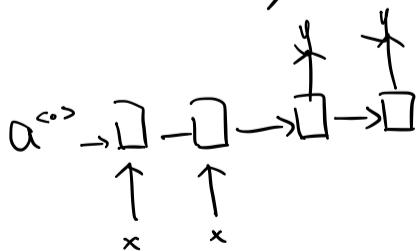
Backprop through time (BPTT)

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

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

Different RNNs

Many-to-Many (NER machine translation)	Many-to-one (movie review)	one-to one (standard neural network)	one-to Many (music generation)
--	-------------------------------	--	-----------------------------------



language model and sequence generation.

$$P(\text{sentence}) = P(y^{<1>} \cdot y^{<2>} \cdot \dots \cdot y^{<T>} \mid \text{EOS})$$

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

$$y^{<n-1>} = x^{<n>}$$

Sampling novel sequences. feed any random word to next cell to create a novel sentence.

Vanishing gradient /

Exploding gradient - gradient clipping, clip over a threshold.

GRU (Gated Recurrent unit)

$$C^{<t>} = \tilde{C}^{<t>}$$

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

$$\tilde{I}_u = \sigma(W_u [C^{<t-1>}, x^{<t>}] + b_u) \text{ gate}$$

$$C^{<t>} = \tilde{I}_u * \tilde{C}^{<t>} + (1 - \tilde{I}_u) * C^{<t-1>}$$

$\tilde{I}_u \approx 0$, does not contribute very much to

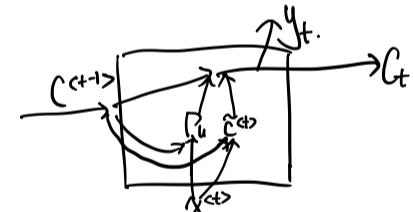
or the full version:

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

$$\tilde{I}_u = \sigma(W_u [C^{<t-1>}, x^{<t>}] + b_u)$$

$$\tilde{R}_r = \sigma(W_r [C^{<t-1>}, x^{<t>}] + b_r)$$

$$C^{<t>} = \tilde{I}_u * \tilde{C}^{<t>} + (1 - \tilde{I}_u) * \tilde{C}^{<t-1>}$$



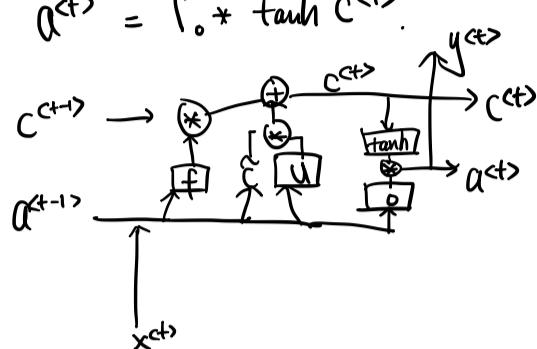
Papers:

Cho et al. 2014 On the properties of neural machine translation encoder-decoder approach

Chung et al. 2014 Empirical evaluation of Gated RNN on sequential Modeling.

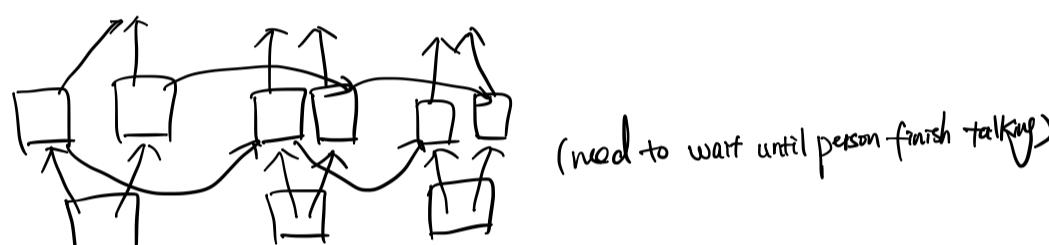
LSTM:

$$\begin{aligned}\tilde{c}^{(t)} &= \tanh(W_c[a^{(t-1)}, x^{(t)}] + b_c) \\ r_u &= \sigma(W_u[a^{(t-1)}, x^{(t)}] + b_u) \quad \text{update gate} \\ r_f &= \sigma(W_f[a^{(t-1)}, x^{(t)}] + b_f) \quad \text{forget gate} \\ r_o &= \sigma(W_o[a^{(t-1)}, x^{(t)}] + b_o) \quad \text{output gate} \\ c^{(t)} &= r_u * \tilde{c}^{(t)} + r_f * c^{(t-1)} \quad (c^{(t-1)} \text{ sometimes}) \\ a^{(t)} &= r_o * \tanh(c^{(t)})\end{aligned}$$

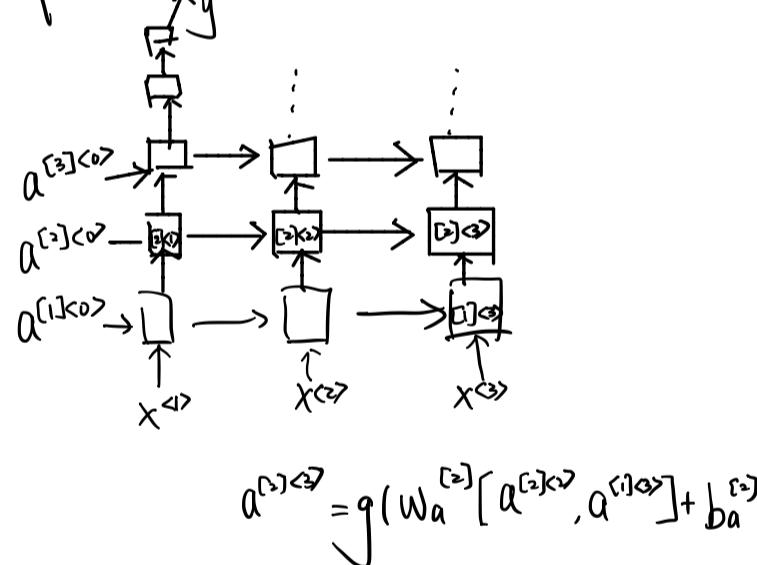


[Hochreiter & Schmidhuber 1997. Long short term memory]

b) - RNN:



Deep RNN:



Word representation.

[van der Maaten and Hinton, 2008. visualizing data t-SNE]

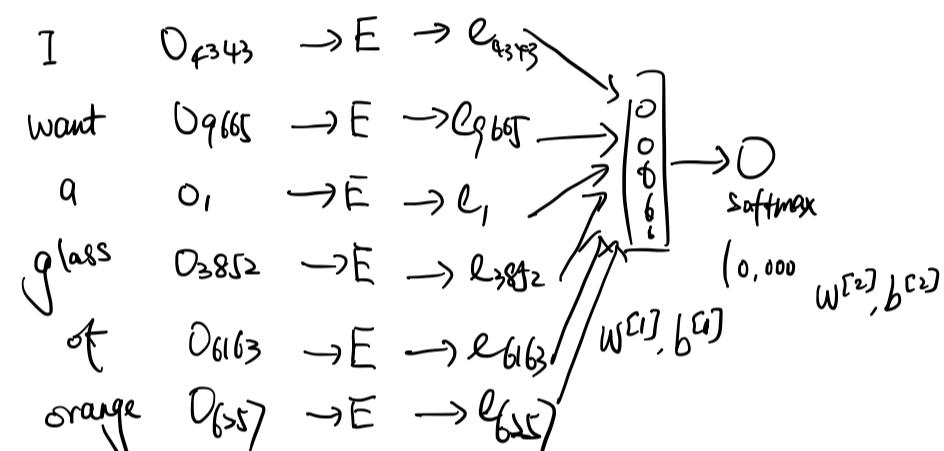
Transfer learning and word embedding.

(face encoding have unlimited faces while word embedding is limited vocab)

[Mikolov et al 2013 Linguistic regularities in continuous space word representations)

$$\begin{aligned}E \cdot o_j &= e_j \\ &\quad \downarrow \\ \vdots &= \vdots \quad \text{embedding for word } j\end{aligned}$$

I want a glass of orange —



4 words on left & right

Context: last 4 words:

embedding is nearby 4 words:
previous 1 words!
nearby words.

Word vector model [Mikolov et al 2013 Efficient estimation of word embeddings in vector space. skip grams:

Context Target.
I want a orange juice
+10 window.

$$\text{softmax: } p(t | c) = \frac{e^{\theta_t^T c}}{\sum_{j=1}^{10,000} e^{\theta_j^T c}} \quad \theta_t = \text{parameter associated with output } t.$$

$$J(\hat{y}, y) = - \sum_{i=1}^{100,000} y_i \log \hat{y}_i$$

Hierarchical softmax can reduce computational cost to:
 $O(\text{vocab}) \nrightarrow O(\log |\text{vocab}|)$

Negative Sampling:

[Mikolov et al 2013 Distributed representation of words and phrases and their compositionality]

orange juice 1 choose nearby k words and replace with random words. to predict True or False.
 orange king 0 common word is placed with false.
 orange book 0 randomly chosen (of and)
 $K = 5-20$ small dataset
 $K = 2-5$ large dataset

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$$

each iteration.

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

$K+1$. binary classification problem

Selecting negative examples:

$$P(w_i) = \frac{f(w_i)^{\frac{3}{4}}}{\sum_{j=1}^{10000} f(w_j)^{\frac{3}{4}}} \quad \begin{array}{l} \text{in middle of from} \\ \text{empirical frequency} \\ \text{to standard distribution} \end{array}$$

Glove algo:

[Pennington et al. 2014 GloVe Global vector for word representation]

x_{ij} = # times i appears in context of j.

$$\min \left(\sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) (\theta_i^T e_j + b_i + b_j - \log x_{ij})^2 \right)$$

$f(x_{ij}) = 0$ if $x_{ij} = 0$

stop words low weights

common words big weights

θ_i, e_j are symmetric

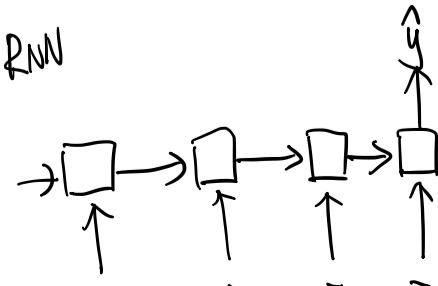
$$e_w^{(final)} = \frac{\theta_w + \theta_w}{2}$$

Sentiment analysis algorithm.

The $\rightarrow e$
 desert $\rightarrow e$
 is $\rightarrow e$
 excellent $\rightarrow e$

Avg $\rightarrow O$ softmax $\rightarrow y$

RNN



The problem of bias in word embeddings

Word embedding can reflect gender, age.

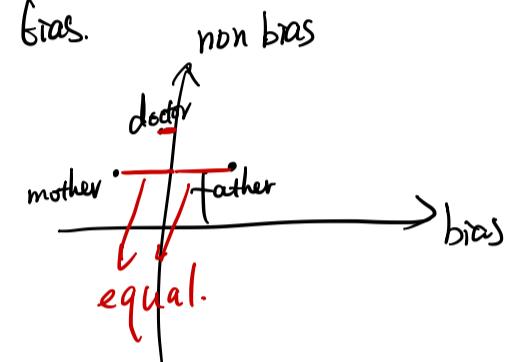
[Bolukbasi et al 2016. Man is programmer as women is homemaker Debiasing word embedding]

1. identify bias direction \rightarrow (gender)

Singular value decomposition to find gender
 $\left\{ \begin{array}{l} \text{female} - \text{female} \\ \text{male} - \text{female} \end{array} \right.$

2. for word that are not definitional.

project to get rid of bias.



3. equalize pairs.

sequence to sequence model.

[Sutskever et al. 2014 Sequence to Sequence learning with neural networks]

[cho et al 2014 Learning phrase representations using RNN encoder-decoder for statistical machine translation

machine translation

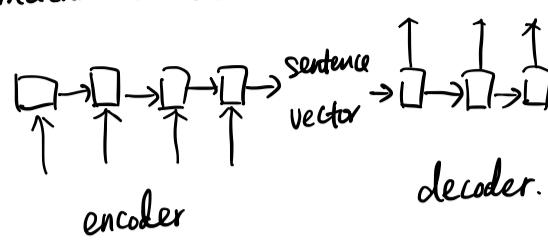
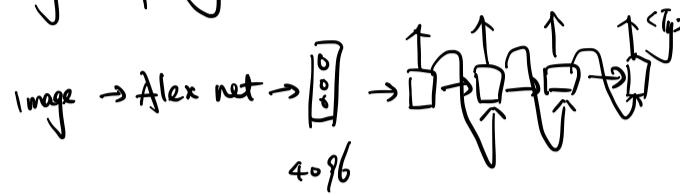


Image captioning:

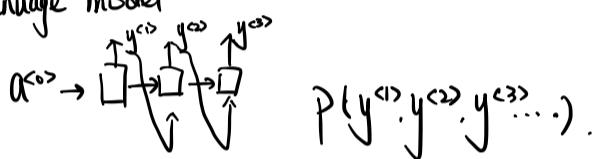


[Mao et al 2014 Deep captioning with multimodal recurrent neural networks]

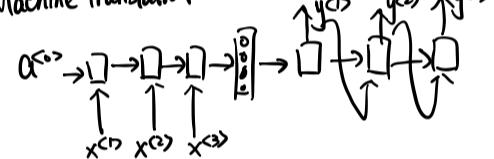
[Vinyals et al 2014 Show and tell : Neural image caption generator]

[Karpathy and Li 2015 Deep visual semantic alignments for generating image descriptions]

Language model



Machine translation



"conditional model" $P(y^{(0)}, y^{(1)}, \dots | x^{(0)}, x^{(1)}, \dots)$

$\text{argmax}_{y^{(0)}, \dots, y^{(T)}} P(y^{(0)}, \dots | x^{(0)}, \dots, x^{(T)})$

Beam search (not greedy search)

B = 3 (beam width)

try first B most possible candidate

$P(y^{(0)}, y^{(1)} | x) = P(y^{(0)} | x) \cdot P(y^{(1)} | x, y^{(0)})$

every step evaluate 3 candidate. in total.

B=1, beam search = greedy search.

$\text{argmax}_y \sum_{t=1}^T \log P(y^{(t)}, | x, y^{(0)}, \dots, y^{(t-1)})$

This makes do not prefer very short length.

Unlike exact search algo BFS (Breadth first Search)
or DFS (Depth first search), Beam search runs faster but not
guaranteed to find exact maximum for $\text{argmax}_y P(y^{(0)}, \dots, y^{(T)} | x)$

Error analysis in Beam search.

RNN compute $P(y^* | x) \cdot P(y^* | x)$. if human's RNN.
beam is bad increase B

y^* human annotate see which one taking up more error
 y RNN output.

Bleu Score (Bilingual evaluation under study)

Ref 1 . A go to B to do C

Ref 2 A visit B to do C,
each word appear in both sentences (bigram)

$$P_1 = \sum_{\text{Unigram } i \in y} \frac{\text{Count}_{\text{clip}}(\text{unigram})}{\text{Count}(\text{unigram})}$$

P_n compute on only gram.

$$\text{combined Bleu score} = BP \exp\left(-\frac{1}{4} \sum_{n=1}^4 P_n\right)$$

BP = brevity penalty

$$BP = \begin{cases} 1 & \text{if MT-length} > \text{human_len} \\ \exp\left(1 - \frac{\text{MT_len}}{\text{reference_len}}\right) & \text{otherwise} \end{cases}$$

[Papineni et al. 2002. Bleu : A method for automatic evaluation of machine translation]

Attention Model intuition.

[Bahdanau et al 2014 Neural machine translation by jointly learning to align and translate]

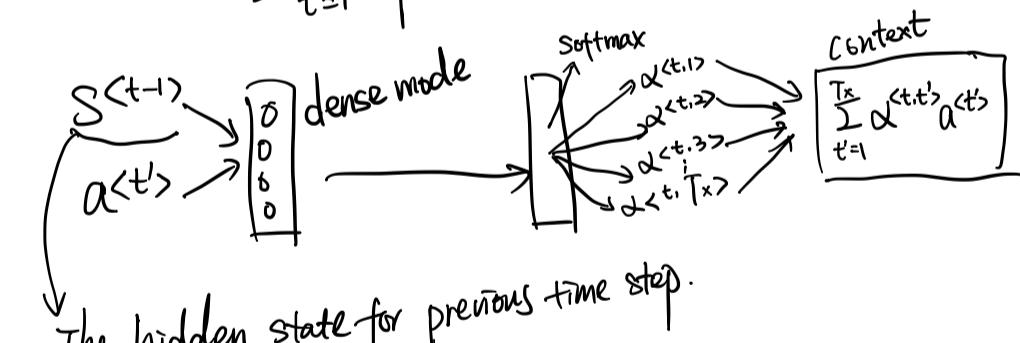
Attention weights:

$$\sum_{t'} \alpha^{(1, t')} = 1$$

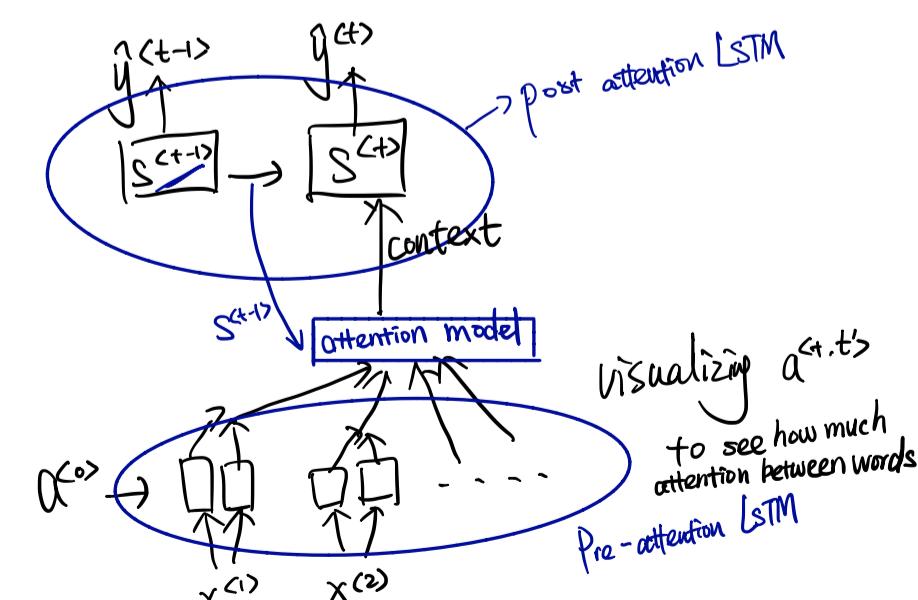
$$c^{(1)} = \sum_{t'} \alpha^{(1, t')} a^{(t')} \quad (a^{(t')} = (\vec{a}^{(t')}, \hat{a}^{(t')}))$$

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

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



The hidden state for previous time step.



The cost is quadratic cost, relatively high.

image captioning use attention model:

[Xu et al 2015 show, attend and tell : Neural image caption generation with visual attention]

speech recognition:

audio clip → transcript

previously: phonemes:

now: end-to-end system with 1000,000 h. audio

CTC loss [Graves et al 2006 Connectionist Temporal Classification: Labelling unsegmented sequence data with Recurrent neural networks.]

Trigger word detection System.

-programming exercise.