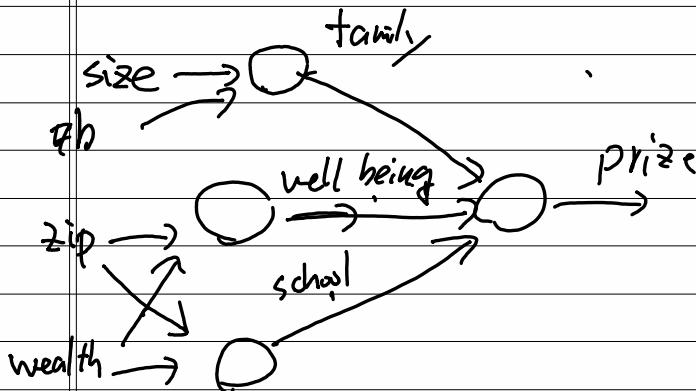


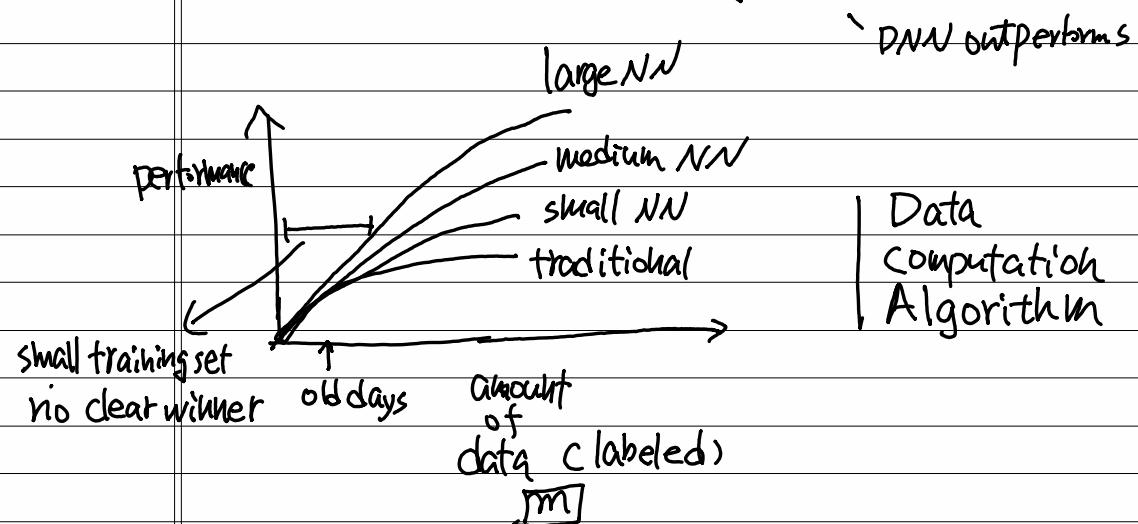
Coursera.



good intuition

but usually full-connected

structured data / unstructured data
table image, video, audio...



W2

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

$$X = \begin{bmatrix} | & | \\ x^{(1)} & \dots & x^{(m)} \\ | & | \end{bmatrix}_{n \times m} \text{ features}$$

training sample

$$X \in \mathbb{R}^{n \times m}$$

$$X.\text{shape} = (n, m)$$

$$Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}_{1 \times m}$$

$$Y \in \mathbb{R}^m$$

stack in column
for easier implementation

logistic regression

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

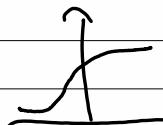
$$X \in \mathbb{R}^{n \times m}$$

parameter: $w \in \mathbb{R}^n$, $b \in \mathbb{R}$

sigmoid function

0

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



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

$$\Rightarrow \sigma(z) = \frac{1}{1+0} = 1$$

$$\begin{aligned} z \downarrow \sigma(z) &= 1 \\ \text{negative} &= \frac{1}{1+e^{\infty}} \\ &= 0 \end{aligned}$$

COST FUNC

cost
func

$$\hat{y} = \sigma(w^T x + b) \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

want $\hat{y} \approx y^{(i)}$

loss:

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

loss
func

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

$$y=1 \rightarrow L(\hat{y}, y) = -\log(\hat{y}) \quad \text{want } -\log(\hat{y}) \downarrow$$

\hat{y} large

$$y=0 \rightarrow L(\hat{y}, y) = -\log(1-\hat{y}) \quad \text{want } \log(1-\hat{y}) \uparrow$$

\hat{y} small

$$\hat{y} = (0, 1)$$

cost func

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

GRADIENT DESCENT

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

convex non-convex

$J(w)$

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

learning rate



$J(w, b)$

$$w := w - \alpha \frac{d J(w, b)}{dw}$$

$$\frac{\partial J(w, b)}{\partial w}$$

$$b := b - \alpha \frac{d J(w, b)}{db}$$

(partial derivative symbol)

∂ = partial (for multiple var)

d = derivative

DERIVATIVE

$$\text{slope} = \frac{h}{w}$$

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

$$f(a) = 3a$$

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

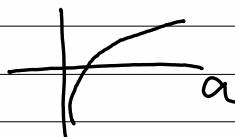
increase a by 0.1

$$y \uparrow 2 \times 0.1 = 0.2$$

$$\frac{d f(a)}{da} = 3$$

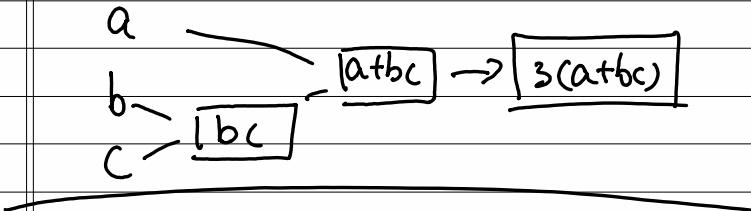
nudge $a \rightarrow$, y increase $3a$

$$f(a) = \log(a) \quad \frac{d}{da} f(a) = \frac{1}{a}$$

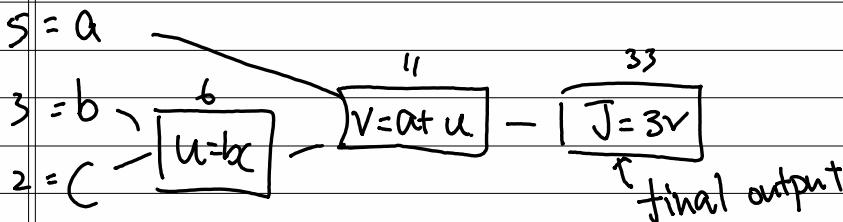


Computation Graph

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



Derivatives with computation graph



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

$$\frac{dJ}{da} = ?$$

$$v = 11 \rightarrow 11.001$$

$$a = 5 \rightarrow 5.001 \rightarrow \frac{dv}{da} = 1$$

$$J = 33 \rightarrow 33.003$$

$$v = 11 \rightarrow 11.001$$

$$dt/dt = 0.001 \times 3 = 0.003$$

$$J = 33 \rightarrow 33.003$$

$$\frac{dJ}{da} = \frac{dJ}{dv} \frac{dv}{da} \quad \text{chain rule}$$

Coding
convention

$$\frac{d \text{finaloutputvar}}{d \text{var}} \rightarrow \underline{d \text{var}}$$

- derivative of final output var.

backpropagation

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

$$S = a$$

$$3 = b$$

$$2 = c$$

$$u = bc$$

$$\frac{du}{dc} = 3$$

$$11$$

$$33$$

$$V = a + u$$

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

$$J = 3V$$

$$u = 6 \rightarrow 6.001$$

$$v = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$\frac{dJ}{du} = \frac{dJ}{dv} \cdot \frac{dv}{du} = 3.$$

$$3 \quad 1$$

$$b = 3 \rightarrow 3.001$$

$$\frac{dJ}{db} = \frac{dJ}{dv} \cdot \frac{dv}{du} \cdot \frac{du}{db} = 6$$

$$(\begin{array}{ccc} 3 & 1 & 1 \\ & 1 & 2 \end{array})$$

$$u = 6 \rightarrow 6.002$$

$$v = 11 \rightarrow 11.002$$

$$J = 33 \rightarrow 33.006$$

$$\frac{dJ}{dc} = \frac{dJ}{du} \cdot \frac{du}{dc} = 9$$

$$(\begin{array}{c} 1 \\ 3 \end{array})$$

Logistic regression gradient descent

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

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

single sample

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

$$\begin{aligned} Z &= w_1 x_1 + w_2 x_2 + b \rightarrow \alpha = \sigma(Z) \rightarrow L(a, y) \\ \frac{dZ}{dz} = \frac{dL}{dz} &= \frac{dL(a,y)}{da} \\ da &= \frac{dL(a,y)}{da} \\ &= -\frac{y}{a} + \frac{1-y}{1-a} \\ &= \underbrace{\alpha(1-\alpha)}_{\text{sigmoid derivative}} \end{aligned}$$

$$\begin{aligned} \frac{dL}{dw_i} &= \frac{dL}{dz} \cdot \frac{dz}{dw_i} \\ &= x_i \cdot \frac{dL}{da} \end{aligned}$$

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

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

To calculate $\rightarrow d_{w_1}^{(i)}, d_{w_2}^{(i)}, db^{(i)}$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} L(a^{(i)}, y^{(i)})$$

$$\frac{\partial}{\partial w_2} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_2} L(a^{(i)}, y^{(i)})$$

▷ Vectorization

CODE

$$\frac{dL(a,y)}{dz} = a - y$$

$$J(w,b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)}) \quad \frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \frac{\partial J}{\partial b}$$

for each sample for $i = 1 \text{ to } m$

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

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

$$dw_1 += x_1 \cdot (a - y)$$

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

$$dw_2 += x_2 \cdot (a - y)$$

$$J += \{y^{(i)} / \log(o^{(i)}) + (1-y^{(i)}) / \log(1-o^{(i)})\}$$

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

$$dw_1 += x_1 dz^{(i)}$$

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

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

$$J /= m$$

$$dw_1 /= m$$

$$dw_2 /= m$$

$$db /= m$$

Vectorization / vectorized logistic regression

$$\underline{z}^{(1)} = \underline{w}^T \underline{x}^{(1)} + b \quad \underline{z}^{(2)} = \underline{w}^T \underline{x}^{(2)} + b$$

$$\alpha^{(1)} = \sigma(\underline{z}^{(1)}) \quad \alpha^{(2)} = \sigma(\underline{z}^{(2)}) \quad \dots$$

$$\begin{matrix} X & X \in \mathbb{R}^{n \times m} \\ \begin{bmatrix} \underline{z}^{(1)} & \dots & \underline{z}^{(m)} \end{bmatrix} = \underline{w}^T \underline{X} + [b \dots b] \\ 1 \times n & n \times n & 1 \times m \end{matrix}$$

$$= \begin{bmatrix} \underline{w}^T \underline{x}^{(1)} + b & \underline{w}^T \underline{x}^{(2)} + b & \underline{w}^T \underline{x}^{(3)} + b & \dots \end{bmatrix}$$

$$\underline{z} = \text{np.dot}(w^T, x) + b$$

$\underbrace{\phantom{\text{np.dot}(w^T, x) + b}}$

$\alpha = \sigma(\underline{z})$

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

$\underbrace{\phantom{\alpha^{(i)} - y^{(i)}}}$

CODE

one iter of gradient descent

$$\underline{z} = \underline{w}^T \underline{x} + b$$

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

$$\alpha = \sigma(\underline{z})$$

$$d\underline{z} = A - Y$$

$$dw = \frac{1}{m} \underline{x} d\underline{z}^T$$

$$db = \frac{1}{m} \text{np.sum}(d\underline{z})$$

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

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

Broadcasting

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 \rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \dots$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} \xrightarrow{(1, n) \sim \cancel{(m, n)}} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \dots$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \\ 300 \end{bmatrix} \xrightarrow{(m, n) \quad (m, 2)} \begin{bmatrix} (m, n) - m \cdot n \end{bmatrix} = \dots$$

$$\begin{array}{rcl} (m, n) & + & (l, n) \sim (m, l) \\ \times & & \text{broadcast} \\ \hline \div & & (m, 1) \sim (m, n) \\ & & \text{broadcasting} \end{array}$$

DEBUG

Avoid "rank-1" array.

randn(5) - DON'T USE

randn(5, 1) - column vector

randn(1, 5) - row vector

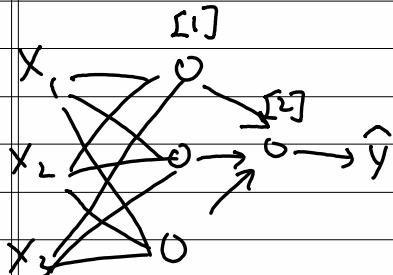
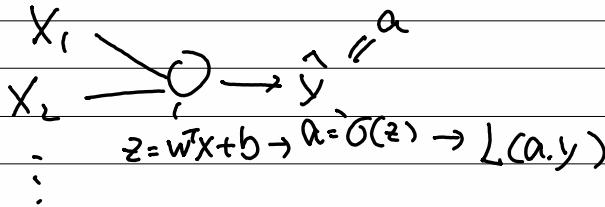
* → assert(a.shape == (5, 1))

↳ (5, 1).reshape((5, 1))

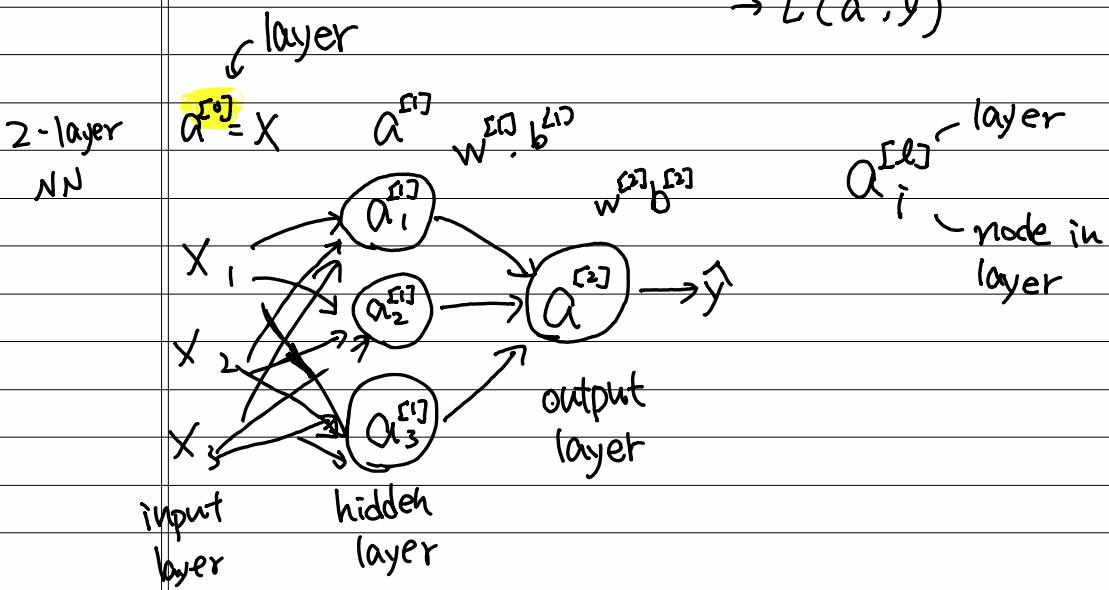
NN

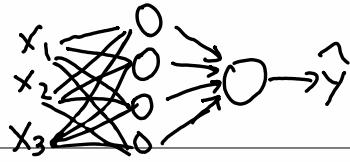
WEEK 3

Logistic regressionish



$$\begin{aligned} z^{[1]} &= w^{[1]} x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \\ &\rightarrow a^{[2]} = \sigma(z^{[2]}) \\ &\rightarrow L(a^{[2]}, y) \end{aligned}$$





Computing
NN output

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

hen roh

vectorize

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

$$\begin{bmatrix} -w_1^{(1)} \\ -w_2^{(1)} \\ \vdots \\ -w_4^{(1)} \end{bmatrix}_{4 \times 3} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} b_1^{(1)} \\ \vdots \\ b_4^{(1)} \end{bmatrix} = \begin{bmatrix} w_1^{(1)T} \cdot x + b_1^{(1)} \\ \vdots \\ w_4^{(1)T} \cdot x + b_4^{(1)} \end{bmatrix}$$

4×3

$$W^{(1)}$$

$$b^{(1)}$$

$$a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ \vdots \\ a_4^{(1)} \end{bmatrix} = \sigma(z^{(1)})$$

CODE

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

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

m samples

$$x \rightarrow a^{(1)} = \hat{y}$$

$$x^{(1)} \rightarrow a^{(1)(1)} = \hat{y}^{(1)}$$

$$\vdots$$

$$x^{(m)} \rightarrow a^{(1)(m)} = \hat{y}^{(m)}$$

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

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

for $i = 1$ to m

$$z^{(1)(i)} = W^{(1)T} x^{(i)} + b^{(1)}$$

\vdots

vectorize
in training

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

$n \times m$.

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

$$A^{[1]} = G\{Z^{[1]}\}$$

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

$$A^{[2]} = G\{Z^{[2]}\}$$

$$Z^{[1]} = \begin{bmatrix} | & | & | \\ Z^{1} & Z^{[1](2)} & \dots & Z^{[1](m)} \\ | & | & | \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} | & | & | \\ a^{1} & \dots & a^{[1](m)} \\ | & | & | \end{bmatrix}$$

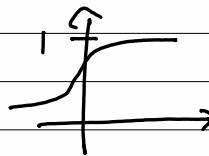
\nwarrow different
Node
in
layers

\nearrow different
sample

$$Z^{1} = W^{[1]} X^{(1)}$$

ACTIVATION

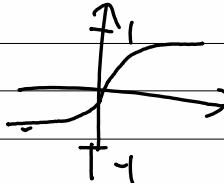
$$G = \frac{1}{1+e^{-z}}$$



prob only for output
 $[0, 1]$

$$g = \tanh(z)$$

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



almost always better
 for hidden layers

on:

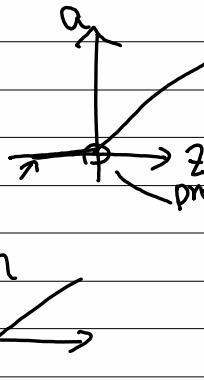
Δ large \Rightarrow small gradient \rightarrow slow

DEFAULT

ReLU

Rectified Linear Unit

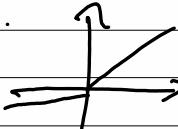
$$\alpha = \max(0, z)$$



pro: learn faster!
 - bcⁿ not affected by
 small slope
 - usually much faster
 w/ $z > 0$

Leaky ReLU

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



Derivative of Activation Function

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

prime

$$\begin{aligned} g'(z) &= \frac{1}{1+e^{-z}} \left(-\frac{1}{1+e^{-z}} \right) \\ g \text{ prime of } z &= g(z) \cdot (1-g(z)) \end{aligned}$$

$$\begin{aligned} \frac{d}{dz} e^{-z} \\ = -z \cdot e^{-z} \end{aligned}$$

* Power rule applies when numerator is
the power of derivative

$$\frac{d}{dx} x^n = n \cdot x^{n-1}$$

$$\begin{aligned} \frac{d}{dx} f(x)^n &= \frac{d}{df(x)} f(x)^n \cdot \frac{df(x)}{dx} \\ &= n \cdot f(x)^{n-1} \cdot f'(x) \end{aligned}$$

$$\begin{aligned} &= -\frac{d}{dz} (1+e^{-z}) \times (1+e^{-z})^2 \\ &= -e^{-z} \cdot (1+e^{-z})^{-2} \\ &= \frac{-e^{-z}}{(1+e^{-z})^2} \end{aligned}$$

$$\begin{aligned} \frac{d}{dx} f(x)^{-1} &= -f(x)^{-2} \cdot f'(x) \\ &= \frac{-f'(x)}{f(x)^2} \end{aligned}$$

$$\begin{aligned} \frac{d}{dx} \frac{\theta(x)}{f(x)} \\ = \frac{d}{dx} \theta(x) \cdot f(x)^{-1} \end{aligned}$$

$$\tanh z = \tanh(z)$$

$$\frac{d}{dz} \tanh(z) = 1 - (\tanh(z))^2$$

ReLU

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

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

implement NN

Forward Prop

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

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

$$z^{[2]} = w^{[1]} A^{[1]} + b^{[2]}$$

$$\begin{aligned} A^{[2]} &= g^{[2]}(z^{[2]}) \\ &= \sigma(z^{[2]}) \end{aligned}$$

Back Prop

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

$$dw^{[2]} = \frac{1}{m} dz^{[2]} X^T$$

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

$$dz^{[1]} = w^{[1]T} dz^{[2]} * g'(z^{[1]})$$

$$dw^{[1]} = \frac{1}{m} dz^{[1]} X^T$$

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

Initialization

\times all $w = 0$

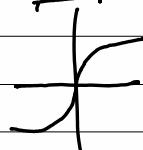
- all hidden unit will be the same

$$w^{[1]} = \text{np.random.randn}(C_{2,2}) \xrightarrow{\text{under}} 0, 0.1 \xrightarrow{\text{in this range}}$$

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

:

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



W3 Assignment

planar

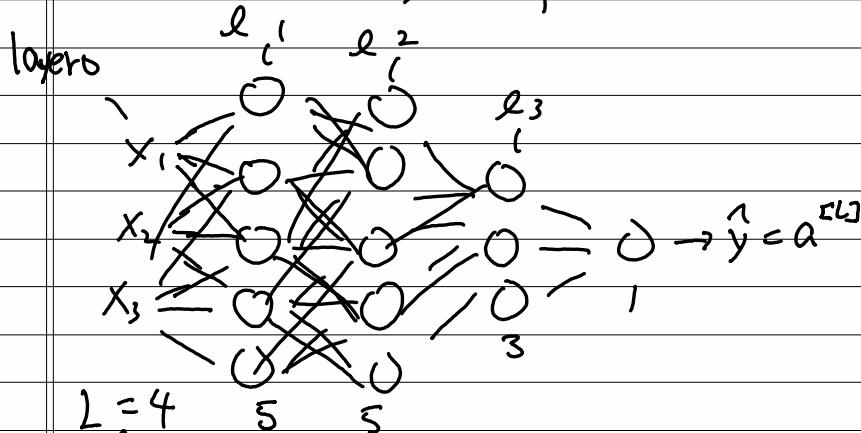
$$X \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad 2 \times m$$

$$Y \in \mathbb{R}^m \quad 1 \times m \text{ (label)}$$

w4

DNN - Deep NN

shallow vs deep



$$L = 4 \quad 5 \quad 5$$

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

$$h^{[1]}=5, h^{[2]}=5, h^{[3]}=3, n^{[4]}=1$$

$$h^{[1]} = n_x = 3$$

$$a^{[l]} = \text{activation in layer } l$$

$$a^{[l]} = g^{[l]}(z^{[l]}) \quad w^{[l]} = \text{weights for } z^{[l]}$$

z

Forward Prop

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

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

↓

vectorize

$$z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

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

$$z^{[l]}$$

$$\{ z^{[0]}, \dots, z^{[L-1]} \}$$

$$z^{[l]}$$

$$\hat{y} = g(z^{[L]}) = A^{[L]}$$

DEBUG DNN

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

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

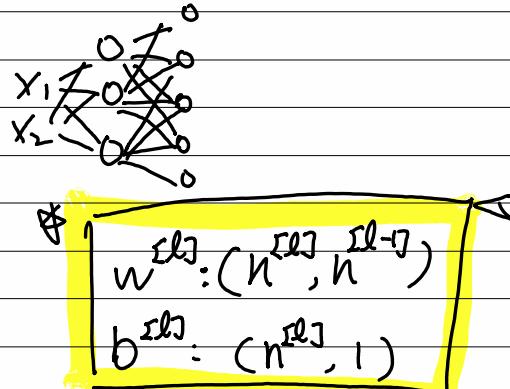
so, w should be $(3, 2)$

$$\Rightarrow (h^{[1]}, 1) = (h^{[0]}, h^{[1]})(h^{[0]}, 1)$$

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

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

$$\therefore (5, 3)$$



$$w^{[l]} : (n^{[l]}, n^{[l-1]})$$

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

$$dw^{[l]} : (n^{[l]}, n^{[l-1]})$$

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

⋮
⋮
⋮

VECTORIZED

$$z^{[l]} = W^{[l]} \cdot X + b^{[l]}$$

($n^{[l-1]} \times 1$) ($n^{[l]} \times n^{[l-1]}$) ($n^{[l]} \times 1$) ($n^{[l]} \times 1$)

$$\sum^{[l]} = W^{[l]} \cdot X + b^{[l]}$$

($n^{[l]} \times m$) ($n^{[l]} \times n^{[l-1]}$) ($n^{[l]} \times m$) ($n^{[l]} \times 1$)

Python broadcasting

\downarrow
($n^{[l]} \times m$)

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

$$l=0, A^{[0]} = X = (n^{[0]}, m)$$

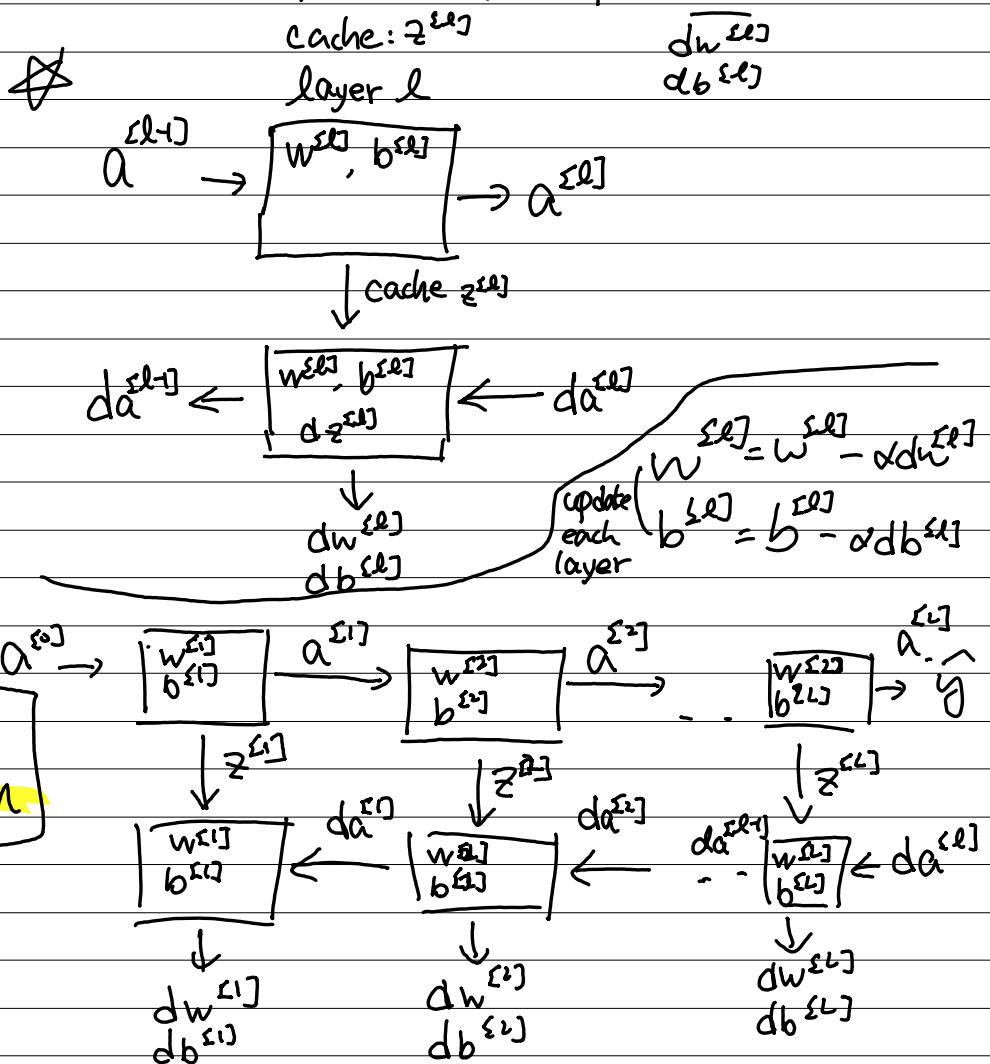
$$dz^{[l]} - dA^{[l]} : (n^{[l]}, m)$$

Forward: input $a^{[l-1]}$, output $a^{[l]}$

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

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

Backward: input $da^{[l]}$, output $da^{[l-1]}$



IMPLEMENTATION

BACKPROP

input: $d\alpha^{[L]}$

output: $d\alpha^{[L-1]}, dW^{[L]}, db^{[L]}$

$$d\alpha^{[L]} = d\alpha^{[L]} * g^{[L]}(z^{[L]})$$

$$dW^{[L]} = d\alpha^{[L]} \cdot a^{[L-1]T}$$

$$db = d\alpha^{[L]}$$

$$d\alpha^{[L-1]} = W^{[L]T} \cdot d\alpha^{[L]}$$

$$d\alpha^{[L-1]} = W^{[L]T} d\alpha^{[L]} * g^{[L-1]}(z^{[L]})$$

$$\Delta d\alpha^{[L]} = -\frac{y}{a} + \frac{(1-y)}{1-a}$$

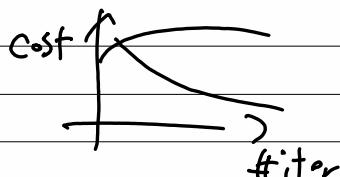
$$dA^{[L]} = \left[-\frac{y}{a} + \frac{(1-y)}{1-a}, \dots, -\frac{y^{(m)}}{a^{(m)}} + \frac{(1-y^{(m)})}{1-a^{(m)}} \right]$$

HYPERPARAMETER

- learning rate
- iteration
- hidden layer L
- hidden unit, $n^{[1]}, n^{[2]}$
- activation func

parameters
 $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots$

look at cost function graph to decide



hyperparameter

trial-and-error.

Course 2

Hyperparameter

Tuning

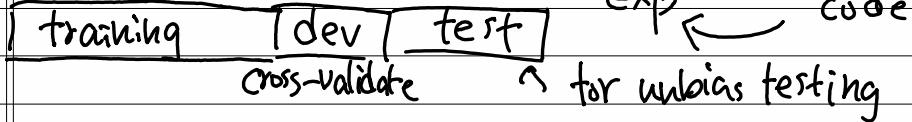
COURSE 2

- Iterative process

for hyperparameter tuning

idea

DATA



Previous

70 / 30 or 60 / 20 / 20.

100 ~ 10000 samples

big data

1,000,000 data, 10,000 test might be enough
10,000

→ 98% / 1% / 1%

Mismatch train/test distribution

training

e.g. cat on web

dev / test

cat on app

- try making

dev / test

with same distribut

Not having test set might be OK

BIAS & VARIANCE

- under fitting — high bias (not fitting well)
- just right
- over fitting — high variance (too flexible)

training set error: 1% 15% 15%

dev set error : 11% 16% 30%

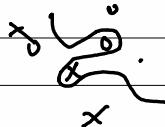
→ overfitting -underfit high bias
high variance high bias var

0.5%

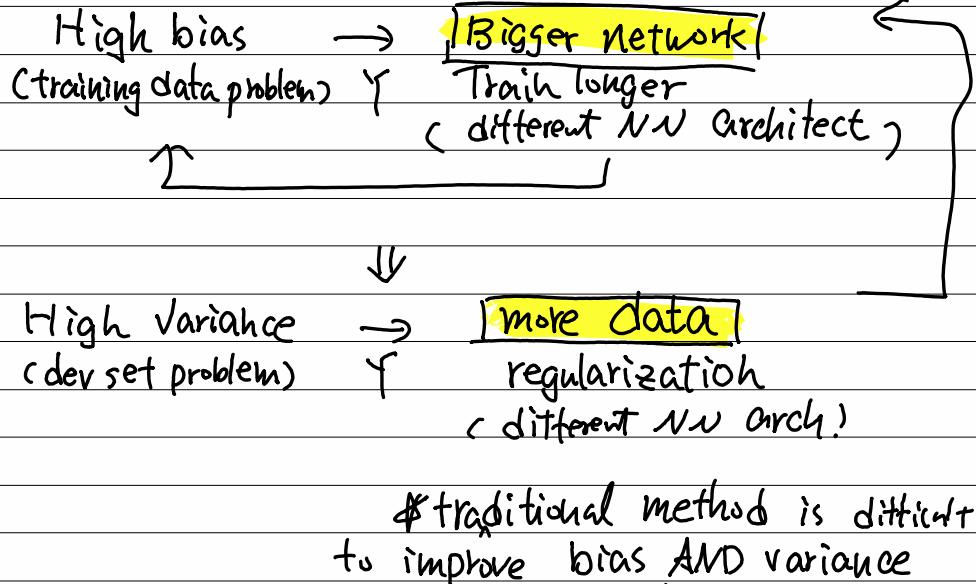
↙ high bias + high variance

1%

↙ low bias var



RECIPE for ML



high bias → more training data is not useful

high variance → regularization
more training data.

Regularization

Logistic Regre..

$$\min_{w, b} J(w, b)$$

$$w \in \mathbb{R}^{k_x}, b \in \mathbb{R}$$

λ : regularization para

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

$$\Delta L_2 \text{ regu } \|w\|_2 = \sqrt{\sum_{j=1}^{k_x} w_j^2} = \sqrt{w^T w}$$

omit

$$L_1 \text{ regularization } \|w\|_1, \frac{\lambda}{2m} \sum_{j=1}^{k_x} |w_j| = \frac{1}{2m} \|w\|_1$$

Neural Network

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

$$\|w^{[l]}\|^2_F = \sum_{i=1}^{n^{[l+1]}} \sum_{j=1}^{n^{[l]}} (w_{ij}^{[l]})^2 \quad w = (w^{[0]}, \dots, w^{[L]}) \quad \text{square norm}$$

Frobenius form sum of square norm of matrix

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

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

$$= w^{[l]} - \underbrace{\frac{\alpha \lambda}{m} w^{[l]}}_{\text{weight decay}} - \alpha (\text{from back prop})$$

weight decay

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

$\lambda \uparrow \Rightarrow w \downarrow$ (bcz we penalize high cost)
 \Rightarrow simpler NN

Dropout regularization

- dropout random nodes in NN
(hidden unit)

INVERTED DROPOUT

e.g. at layer 3,

$$\text{keep-prop} = 0.8$$

$$d_3 = \text{np.random.rand}(a_3.\text{shape}[1], a_3.\text{shape}[0])$$

$$a_3 = \text{np.multiply}(a_3, d_3) \quad < \text{keep-prop}$$

* $\Rightarrow a_3 \neq \text{keep-prop}$

e.g. total 50 units, start off 10 units (20%)

$$\overset{(20)}{\tilde{z}} = \overset{(20)}{W} \cdot \overset{(20)}{a} + \overset{(20)}{b}$$



bump up by 20%

* no drop-out at test time

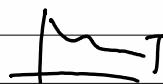
why drop-out works?

* intuition \rightarrow can't rely on one feature, so we spread out weights.

* keep-drop can vary in layers \rightarrow shrink weights similar to L2

- Computer Vision uses drop out b/c high dimension

- Downside. J might looks weird in the plot

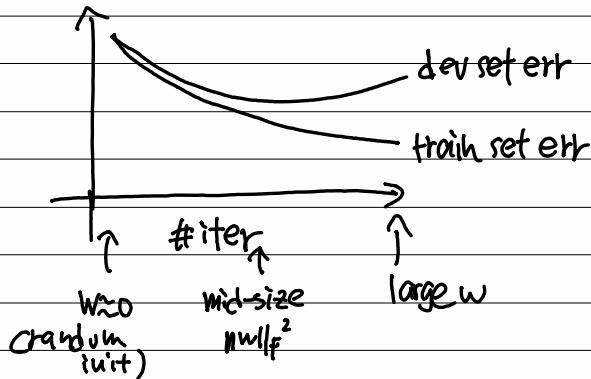


Reguralization methods

Data augmentation

- flipping, random crop images. (flipping cats)
- Adding distortion to characters etc

Early stopping



↳ breaks orthogonalization

Normalizing Input

subtract mean

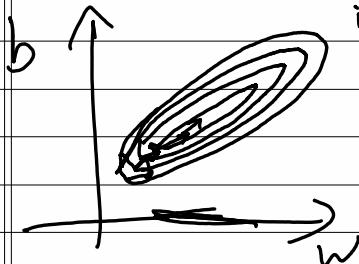
normalize variance

$$X = X - \mu$$

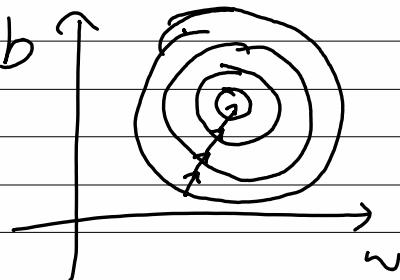
$$X \sim \sigma^2$$

* use same μ & σ^2 for
testing / eval / train sets

unnormalized b might be in very diff range normalized



require small learning rate
for gradient descent



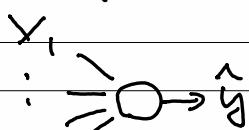
better!
sater to take larger step!

vanishing / exploding gradient

(assuming linear activation) $\hat{y} = w \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}^{L-1} x$ - exploding

back-prop $\hat{y} = w \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^{L-1} x$ - vanishing

solution (initialization)



$$z = w_1 x_1 + \dots + w_n x_n + b$$

large $n \rightarrow$ small w_i

$$x_n \propto g(z)$$

$$\text{Var}(w_i) = \frac{1}{n} \quad \text{if ReLU}, \quad \text{Var}(w_i) = \frac{2}{n}$$

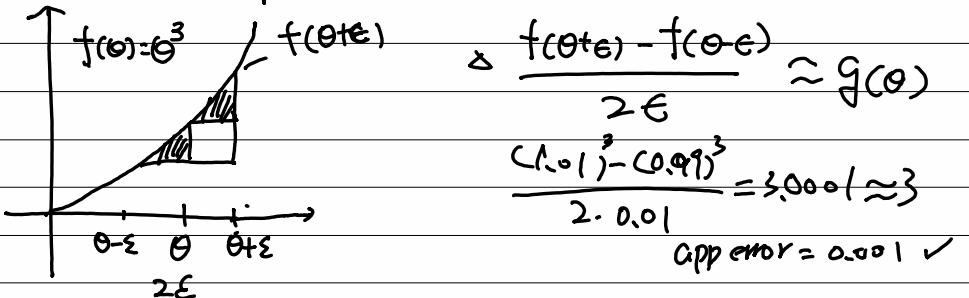
$$w^{[L]} = \text{np.random.rand(slope)} \times \text{np.sqrt}\left(\frac{2}{n^{[L]}}\right)$$

ReLU

$$\tanh = \sqrt{\frac{1}{1+e^{-x}}}$$

Xavier initialization $\sqrt{\frac{2}{n^{[L-1]} + n^{[L]}}}$

Numeric Approx of Gradients



$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \quad O(\epsilon^2)$$

$$f'(\theta) \approx \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \quad \text{err: } O(\epsilon)$$

Δ we use two-sided for gradient checking (more accurate) instead of one-sided checking

DEBUG TOOL

Gradient Checking (Grad Check)

$$J(w^{(1)}, b^{(1)}, \dots, w^{(m)}, b^{(m)}) = J(\theta)$$

$$J(dw^{(1)}, db^{(1)}, \dots, dw^{(m)}, db^{(m)}) = J(d\theta)$$

for each i:

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

$$\approx d\theta^{(i)} = \frac{\partial J}{\partial \theta_i} \Rightarrow \text{check } d\theta_{\text{app}} \approx d\theta?$$

$$\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2} \approx 10^{-7} - \text{great}$$

$$\approx 10^{-5}$$

$$10^{-3} - \text{worry! bug!}$$

Grad Check Implementation

- Don't use in training
- look at component to try to identify the bug
- remember regularization
- doesn't work with dropout
 - keep-prop=1 for debug
- run at random initialization

mini-batch gradient descent

split [training set] into batches

$$X^{t \in S}, Y^{t \in S}$$

{

△ notation

$$X^{(i)}$$

batch v.s. mini-batch

$$X^{[l]}$$

(all) (small batches)

$$X^{[t \in S]}$$

curly bracket.

training set batch

EXAMPLE

5,000,000 data \rightarrow 5000 mini-batch, each has 1000 data.

for $t = 1$ to 5000

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

$$\text{cost } J^{(t)} = \frac{1}{1000} L(y^{(t)}, y_{(t)}) + \frac{\lambda}{2 \cdot 1000} \|w^{(t)}\|_F^2$$

backprop using $J^{(t)}$

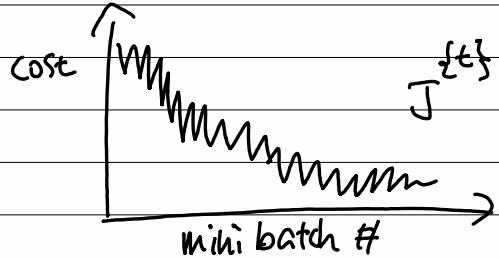
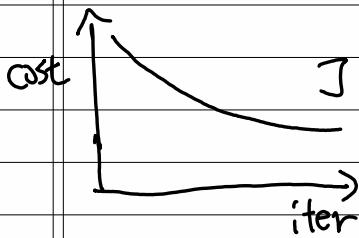
$$w^{(t+1)} = w^{(t)} - \alpha d w^{(t)} \quad b^{(t+1)} = b^{(t)} - \alpha d b^{(t)}$$

1 epoch - a single pass of 'All' data

BATCH - 1 step descent

MINIBATCH - 5000 step descent.

WHY MINI-BATCH WORK?



noisy, but trend down
bcz batch might be difficult

mini-batch size

size = $m \rightarrow$ batch gradient descent

size = 1 \rightarrow stochastic gradient descent.

BATCH - too long per iteration

STOCHASTIC - lose the benefit of vectorization

GUIDE:

Small training set:

(~2000) use batch gradient

typical size

64 ~ 512

make sure mini-batch fits CPU/GPU memory.

Exponentially Weighted Average

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

moving

△

V_t is approx average
over $\frac{1}{1-\beta}$ samples

$$\beta = 0.9$$

$$V_{100} = 0.1 \theta_{100} + 0.9 V_{99}$$

e.g. $\beta = 0.9 \rightarrow 10$ samples

$$V_{99} = 0.1 \theta_{99} + 0.9 V_{98}$$

$\beta = 0.98 \rightarrow 50$ samples

$$V_{98} = 0.1 \theta_{98} + 0.9 V_{97}$$

:

$$V_{100} = 0.1 \theta_{100} + 0.9 V_{99}$$

$$= 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + 0.9 V_{98})$$

$$= 0.1 \theta_{100} + 0.1 \cdot 0.9 \theta_{99} + 0.9^2 (0.1 \theta_{98} + 0.9 V_{97})$$

$$= 0.1 \theta_{100} + 0.1 \cdot 0.9 \theta_{99} + 0.1 \cdot (0.9)^2 \theta_{98} + (0.9)^3 V_{97} + \dots$$

- - -

$(0.9)^{10} \approx 0.35 \approx \frac{1}{e}$ after $\frac{1}{e}$, it drops faster.

$$(0.98)^{50} \approx \frac{1}{e}$$

Bias Correction in EWMA

bias because $V_0 = 0$

$$U_t = \beta V_{t-1} + (1-\beta) \Theta_t$$

incorrect V in the beginning of series

$$V_0 = 0$$

$$U_1 = 0.9V_0 + 0.1\Theta_1 = 0.1\Theta_1$$

$$V_2 = 0.9U_1 + 0.1\Theta_2$$

$$= 0.09\Theta_1 + 0.1\Theta_2$$

⋮

Correction

$$\frac{U_t}{1-\beta^t}, \text{ correct in the beginning}$$

$t \uparrow; \beta^t \sim 0$, correction does not affect U_t .

Gradient Descent with Momentum

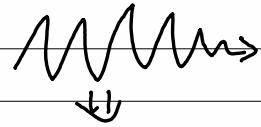
Compute EWA of gradient, then update

on iter t

compute dw db on current mini-batch

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

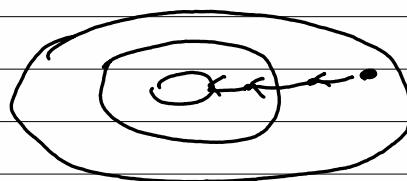
$$V_{db} = \beta V_{db} + (1-\beta)db$$



$$w = w - \alpha V_{dw}$$

$$B = B - \alpha V_{db}$$

dampen oscillation
overshoot.



ball rolling.

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

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

Acceleration

velocity

friction

hyper parameters

α - learning rate

β

0.9 usually works

RMSprop (root mean square prop)

On iter t.

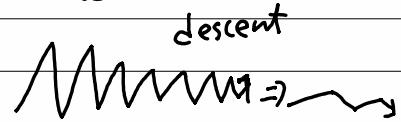
compute dw, db on mini-batch

$$Sdw = \beta Sdw + (1-\beta)dw^2 \quad \text{- element-wise square}$$

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

$$w = w - \alpha \frac{dw}{\sqrt{Sdw} + \epsilon} \quad \text{P.G.}$$

$$b = b - \alpha \frac{db}{\sqrt{Sdb} + \epsilon} \quad \text{avd divided by 0}$$



$$Sdw \downarrow \quad w \uparrow$$

$$Sdb \uparrow \quad b \downarrow$$

Adam Optimization Algorithm

Adaptive
moment
estimation

Momentum + RMSdrop

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

On iter t:

compute dw, db of mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dw \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) db$$

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

bias correction

$$V_{dw}^{\text{corrected}} = \frac{V_{dw}}{1-\beta_1^t} \quad V_{db}^{\text{corr}} = \frac{V_{db}}{1-\beta_1^t}$$
$$S_{dw}^{\text{corrected}} = \frac{S_{dw}}{(1-\beta_2^t)} \quad S_{db}^{\text{corr}} = \frac{S_{db}}{(1-\beta_2^t)}$$

RMSprop

RMSprop

$$w = w - \alpha \frac{V_{dw}^{\text{corr}}}{\sqrt{S_{dw}^{\text{corr}}} + \epsilon} \quad b = b - \alpha \frac{V_{db}^{\text{corr}}}{\sqrt{S_{db}^{\text{corr}}} + \epsilon}$$

hyperparameter α - tube

$$\beta_1: 0.9 \rightarrow dw$$

first momentum

$$\beta_2: 0.999 \rightarrow dw^2$$

second momentum

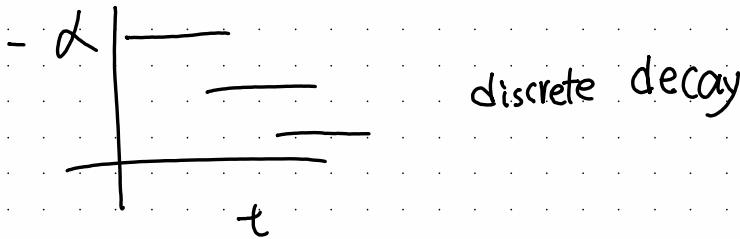
$$\epsilon: 10^{-8}$$

Learning Rate Decay

slowly reduce α (avoid overshoot when close to solution)

$$\alpha = \frac{1}{\text{lr decay rate} \times \text{epoch num}} \alpha_0 \Rightarrow \begin{matrix} \text{para.} \\ \alpha_0 \\ \text{decay rate} \end{matrix}$$

- $\alpha = 0.95^{\frac{1}{\text{epoch num}}} \cdot \alpha_0$ - exponential decay
 - $\alpha = \frac{k}{\sqrt{\text{epoch num}}} \cdot \alpha_0$ or $\frac{k}{\sqrt{t}} \cdot \alpha_0$



- manual decay (manually update ...)

Local optima

most local optima happens at saddle point.

- unlikely to stuck in bad local



- plateaus can make learning slow.



Hyperparameter tuning

α . β . $\beta_1, \beta_2, \epsilon$ * α

layers

1 / # hidden units second important

2nd - learning rate decay

mini-batch

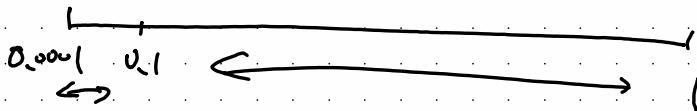
- use "random" search, not grid

- coarse to fine

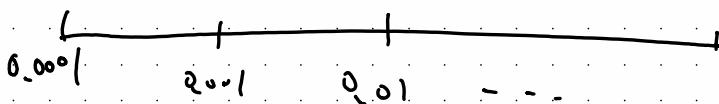
appropriate scale to pick hyper-parameter.

$n^{[l]}$, L uniform sample works

α .



in log scale



$$\gamma = -4 * \text{np.random.rand} \leftarrow \gamma \in [-4, 0]$$

$$\alpha = 10^\gamma \leftarrow \gamma \in [-4, 0]$$

Hyperparameter for EWA (exponentially weighted averages)

$$\beta = 0.9 \dots 0.999$$

$$r \in [-1, 0] \quad \begin{matrix} 0 \\ \nearrow \\ 0.1 \quad 0.01 \quad 0.001 \end{matrix}$$

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

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

$$\frac{1}{1-\beta} \leftarrow \text{when } \beta \rightarrow 1.$$

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

Small changes have big

impact.

$$\beta = 1 - 10^{(r-1)}$$

Hyperparameter tuning

Babysitting one model = fewer resource
Training many models in parallel

large model

Panda
Cavier.

Batch Normalization (single layer)

can we normalize $a^{(l)}$ so that training $w^{(l+1)}, b^{(l+1)}$ faster

$$\mu = \frac{1}{m} \sum z^{(l)}$$

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

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

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

$$\text{if } \gamma = \sqrt{\sigma^2 + \epsilon}$$

$$\beta = \mu$$

$$\text{then } \tilde{z}^{(l)} = z^{(l)}$$

$$\underbrace{z^{(l)} \dots z^{(n)}}_{z^{(l), c_i}} \quad \star$$

learnable hyperparameter

⇒ use $\tilde{z}^{(l)}$ instead of $z^{(l)}$
activation

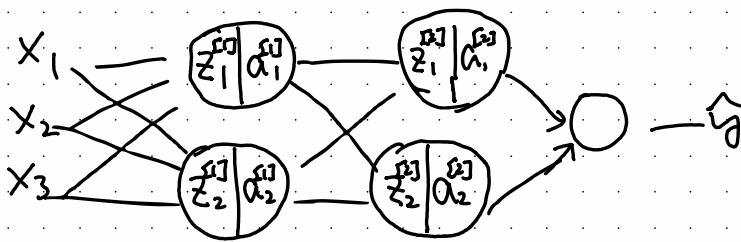


normalization

avoids value squeeze in small area

(by spreading the value)

Batch norm in a NN



$$X \xrightarrow{w^{[l]}, b^{[l]}} z^{[l]} \xrightarrow{\text{BatchNorm}, \beta^{[l]}, \gamma^{[l]}} \tilde{z}^{[l]} \xrightarrow{\alpha^{[l]}} \hat{y}$$

$$w^{[l+1]}, b^{[l+1]} \xrightarrow{z^{[l+1]}} \tilde{z}^{[l+1]} \xrightarrow{\text{BatchNorm}, \beta^{[l+1]}, \gamma^{[l+1]}} \hat{y}$$

Parameters:

$$(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}, \beta^{[1]}, \gamma^{[1]}, \dots, \beta^{[L]}, \gamma^{[L]})$$

this β is different
from β in Momentum

ADAM
etc

(tf.nn.batch-normalization)

WORKING with MINI-BATCH

$$X^{[1]} \\ X^{[2]} \\ \vdots \\ X^{[n]}$$

parameter: $w^{[l]}, b^{[l]}, \beta^{[l]}, \gamma^{[l]}$ → $z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$

Batch-Normalization

$(n, 1) \quad (h, w)$ → $\tilde{z}^{[l]} = w^{[l]} \tilde{a}^{[l-1]}$ eliminate

$\tilde{a}^{[l]} = \gamma^{[l]} \tilde{z}^{[l]} + \beta^{[l]}$

IMPLEMENTATION.

for $t = 1 \dots \# \text{MINI BATCH}$

compute forward prop on $X^{[t]}$

in each layer, use BN to replace $Z^{[e]}$ with $\tilde{Z}^{[e]}$

use backprop to compute $dW^{[e]}$, $b^{[e]}$, $ds^{[e]}$, $dr^{[e]}$

update para $W^{[e]} = W^{[e]} - \alpha dW^{[e]}$

$B^{[e]} = B^{[e]} - \alpha ds^{[e]}$

work with momentum, RMSprop, ADAM etc

Batch norm work bcz.

- learning on shifting of distribution

'covariance shift' e.g. training on black cat

✗ Batch norm reduce "cor shift" hard to generalize to yellow cat.

△ Batch norm as regularization

- mean/var of each mini-batch differ

- Add noise to value $Z^{[e]}$

- similar to dropout.

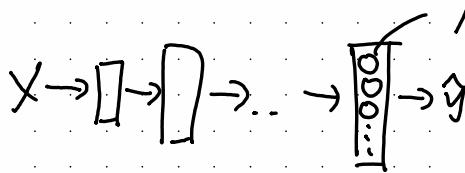
if add noise to each hidden layer's activation

→ slight regulation

Batch norm at test time

M, S^2 for testing: exponentially weighted average across all mini-batch

softmax regression



$$p(\text{class}; | x)$$

$$N^{\text{class}} = \# \text{ class}$$

e.g. 4 class

$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)} \quad (4, 1)$$

Activation Function

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

$$a^{(l)} = \frac{e^{z^{(l)}}}{\sum_{j=1}^4 t_j}, \quad a_i^{(l)} = \frac{t_i}{\sum_{j=1}^4 t_j}$$

hard max

$$\begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

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

$$(4, 1) \quad (4, 1)$$

softmax activation.

vector \rightarrow vector

others output real number.

$C=2 \rightarrow$ softmax reduced to logistic regression

Loss function

$$C=4, \quad y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

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

$$\text{e.g. } y_2 = -y_2 \log \hat{y}_2$$

reduce loss $\rightarrow \hat{y}_2 \uparrow$

in sample

maximum likelihood estimation

$$\mathcal{J}(w^{(i)}, b^{(i)}, \dots) = \frac{1}{m} \sum L(y^{(i)}, \hat{y}^{(i)})$$

Deep Learning Frameworks

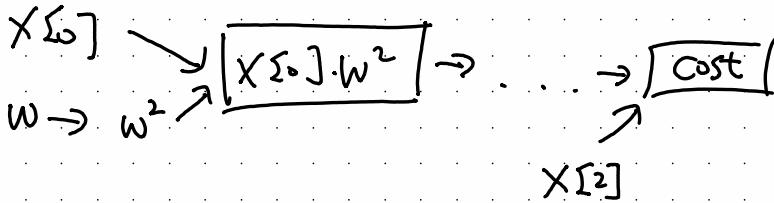
TensorFlow

$$\leftarrow \{w=s\}^2, w-s=0$$

$$J(w) = w^2 - 10w + 25$$

↖ only need to calculate forward prop.

Gradient Tape for back prop



Course 3

STRUCTURING ML Project

Orthogonalization

- many potential tools

Fit training set well on cost function

- bigger network

- Adam

- :

↗
X early stopping
↖ X.

Fit dev set well on cost function

- Regularization

- bigger training set

Fit test set well on cost function

- Bigger dev set

perform well in real-world

- change dev set or cost function

Single Number Value Evaluation

precision → recognized cat, % of them as cat

recall → % of actual cat got recognized

Competing... hard to evaluate

⇒ F1 score ~ average of P & R

$$\frac{2}{P+R} \text{ "Harmonic Mean"}$$

Satisficing & Optimizing Metric

e.g.

maximizing accuracy

subject to running time < 100ms

Optimizing

Satisficing

e.g. wake word / trigger word

'hi siri'

maximize accuracy

s.t. ≤ 1 false positive every 24 hrs.

Train / Dev / Test set

Dev / Test should have same distribution

choose Dev/Test that reflect your future target

Size of dev/test sets

Old

170%	30%
Train	Test
60%	20%
Train	D

100
10,000

1,000,000 data	98%	1.1%
Train	Dev	test

high

size of test set ~ big enough to give confidence on performance

- [train | dev] ← sometimes not recommended
- [train | test] X not good naming...

When to change Dev/Test and Metrics

real data
different from expectation

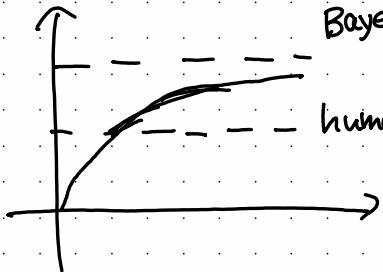
- [metric + dev prefer A → time to change
your user prefer B e.g. bigger penalty on
undesired image]

orthogonality
- place target (define metrics)

- shoot at it (change cost func... etc)

human-level performance

Bayes optimal error \leftarrow best possible



→ progress slows when
surpass human

if ML is worse

- get more labels
- gain insight from manual analysis
- better analysis of bias / variance

avoidable bias

▷ assume human-level
error \approx bayes error
(esp. CV tasks)

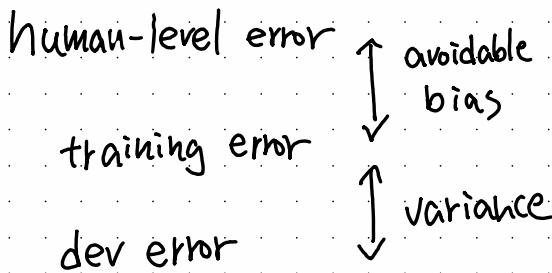
✗ human	1%	7.5%
training error	8%	8%
dev error	10%	10%
	↑	↑
reduce bias		reduce variance

Avoidable
bias

define human-level error

human-level error as proxy to bayes error.

e.g. human doctor
theoretical upper bound
group of doctors ← closer to bayes error



surpassing human-level error

- online ad
- product recommend
- logistic

speech recog
medical data
:

learned from structure data
not involve perception
lots of data

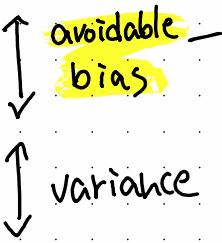
fundamental assumption of supervised learning

- one can fit training set well
- training set generalizes to dev/test.

human-level error

training error

dev error



bigger model
better optimize
NN Arch / hyper-para change

More data
regularization
NN Arch / Hyper-para search

Error Analysis

cat vs dog. if find mislabelled ones are dog.

Analysis

1. get ~100 mislabeled dev set example
2. count how many are dogs
 - if low, probably not worth it.

↳ build table for incorrectly labelled dev set.

Image Dog big cat blurry .. ↪ ih

1	✓
:	
n	✓

%



based on percentage, decide what's next step.

Incorrect Labelled Data

- robust at "random" error in training data
- for dev/test, label accuracy might be less important.
(based on the err analysis table)

Correcting dev/test labels

- apply same process for both dev/test set
(ensure same distri)
- consider examples your algorithm got "right" as well as "wrong"
- note training and dev/test might come from different distri

Build system quickly then iterate

- set dev/test set early
- build
- bias/var/error analysis
 - esp for new applications

Training / Dev from different distri

e.g.

Cat from webpage web from mobile (blur etc)

~200,000

~10,000

what we care

option:

1. $200k + 10k = 210k$ (X)

train dev/test

$$2,500 \rightarrow 2500 \times \frac{10}{210} = 119 \text{ mobile image.}$$

2.

train dev/test (v)

200k + 5k
web mobile

2500 2500
mobile photos

Bias & variance with mismatched data dist.

training

train dev/test

train on this

[training-dev set]

human err 0%

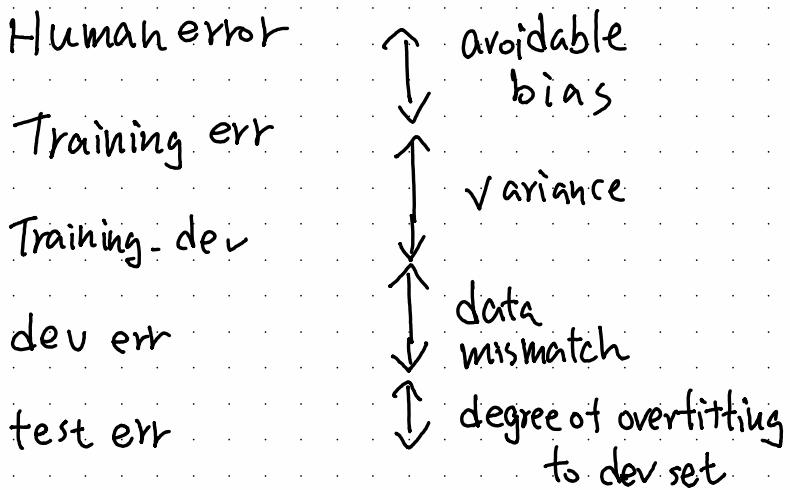
Train ↑ 1% 1% 10% 10%

same dist as training set

train-dev ↓ 9% 1.5% ↑ 11% 11%

dev 6% 10% ↓ 12% 20%

Variance data bias bias
problem mismatch problem + data mismatch



Sometimes dev performs better than training

- e.g. dev/test are of different distribution
- additional label on the dev/training might better understand human performance

Addressing Data Mismatching

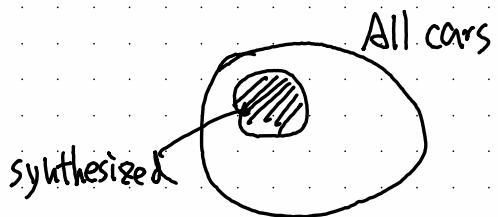
- manual error analysis to understand the diff between training, dev/test.
 - making training data more similar

Artificial Data Synthesis

e.g. speech + car noise = synthesised in-car speech

note ↑ ↑ ↑
10,000 hrs 1 hr might be overfitting
the 1 hr car noise

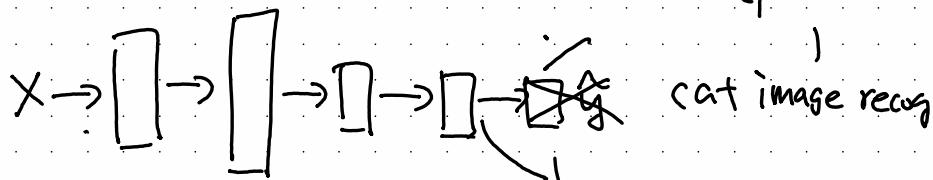
e.g. car photo



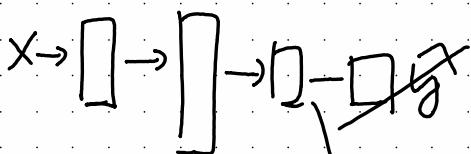
8 if using synth data, make sure it capture enough variety.

Transfer Learning

(pre-training)



or

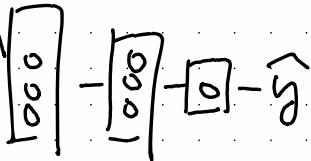


\hat{y} radiology

$w^{(4)}, b^{(4)}$

(fine-tuning)

Working bcz low-level
feature might be
similar.



Transfer from A \rightarrow B, useful when.

- task A & B have the same input
- much more data in A
- lower level feature of A is helpful for B

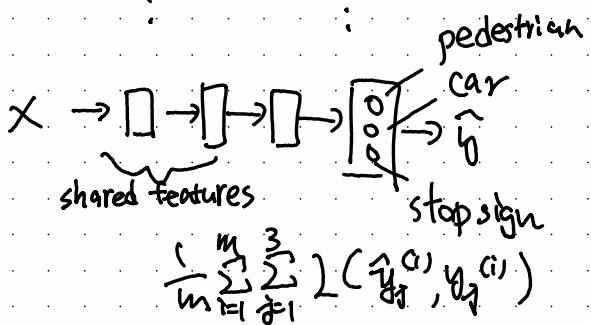
Multi-task Learning

(less often used)
except CV. compared to TL

e.g. autonomous driving

pedestrian	1
Cars	0
stop sign	0
:	:
:	:

▷ one image with
multiple labels



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

▷ might still work when some images only have
partial labels

MAKE SENSE WHEN

- tasks share similar low-level features
- amount of data for each feature are similar
- can train a big enough NN to do well on all tasks

End-to-End DL

X Audio $\xrightarrow{\text{MFCC}}$ features $\xrightarrow{\text{ML}}$ phonemes \rightarrow words \rightarrow transcript
audio \longrightarrow transcript

need a lot of data!

e.g.

Face recognition

image $\xrightarrow{\text{tracking}}$ face $\xrightarrow{\text{ML}}$ recognition
face detection

& have lot data for each of 2 sub-tasks.
even more data for end-to-end data.

Machine Translation

X Eng \rightarrow text analysis $\rightarrow \dots \rightarrow$ French

✓ Eng \longrightarrow French

Estimate child age

✓ image $\xrightarrow{\text{①}}$ bone $\xrightarrow{\text{②}}$ age

X image \longrightarrow age

wheather to use End-to-End DL

Pro :

- let the data speak
- less hand-design features

Con :

- data hunger
- exclude hand-designed features

↳ hand-design feature helps when training set is small

KEY Q:

sufficient data to learn a function maps X to Y ?

Course 4

Convolutional Neural Networks

CoVNet

Computer Vision

Edge Detection

larger number
= edge.

Convolution

$\begin{matrix} 3 & 0 & 1 & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \end{matrix}$

\times

filter (vertical edge)

$=$

$\begin{matrix} 6 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{matrix}$

6×6 img.

3×3

4×4

Python: conv-forward

tensor : tf.nn.conv2d

$n-j+1$

Sobel filter

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

Scharr filter.

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

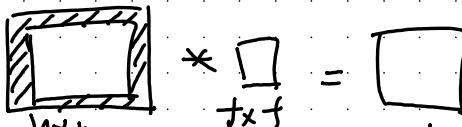
\times

$\begin{matrix} w_1 & w_2 & \cdot \\ w_3 & w_4 & \cdot \\ w_5 & w_6 & \cdot \\ w_7 & w_8 & w_9 \end{matrix}$

\uparrow

use NN to learn filter

Padding



valid -
(no padding)

$$(n+p) \times (n+p) \times P_{fxf} = n-f+2p+1 \times n-f+2p+1$$

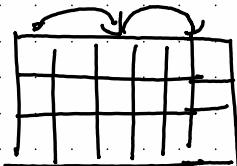
same
(padding)

$$P = \frac{f-1}{2}$$

$$\text{e.g. } f=3 \rightarrow p=1, \quad f=5 \rightarrow p=2$$

f is almost always odd $\begin{matrix} 3 \times 3 \\ 5 \times 5 \\ \vdots \end{matrix}$

Strided convolution



$$n \times n \times f \times f$$

$$\text{pad: } p, \quad \text{stride: } s \Rightarrow L \left\lceil \frac{n+2p-f}{s} + 1 \right\rceil \times L \left\lceil \frac{n+2p+f}{s} + 1 \right\rceil$$

Convolution in math

$$\begin{array}{|c|c|c|} \hline 3 & 4 & 5 \\ \hline 1 & 0 & 1 \\ \hline 9 & 7 & 8 \\ \hline \end{array}
 \xrightarrow{\text{flip before conv.}}
 \begin{array}{|c|c|c|} \hline 8 & 7 & 9 \\ \hline 1 & 0 & 1 \\ \hline 5 & 4 & 3 \\ \hline \end{array}$$

cross-corr : no flip filter
 conv : flip filter

ML community just call cross-corr as convolution

Conv in volume

$$\begin{array}{c} \text{input} \\ \text{6x6x3} \\ \text{---} \end{array} * \begin{array}{c} \text{filter} \\ \text{3x3x3} \\ \text{---} \end{array} = \begin{array}{c} \text{output} \\ \text{4x4} \\ \text{---} \end{array}$$

↑ same channels

multiple filters

$$\begin{array}{c} \text{input} \\ \text{6x6x1} \\ \text{---} \end{array} * \begin{array}{c} \text{filter} \\ \text{3x3x3} \\ \text{---} \end{array} = \begin{array}{c} \text{intermediate output} \\ \text{4x4x1} \\ \text{---} \end{array}$$

horizontal

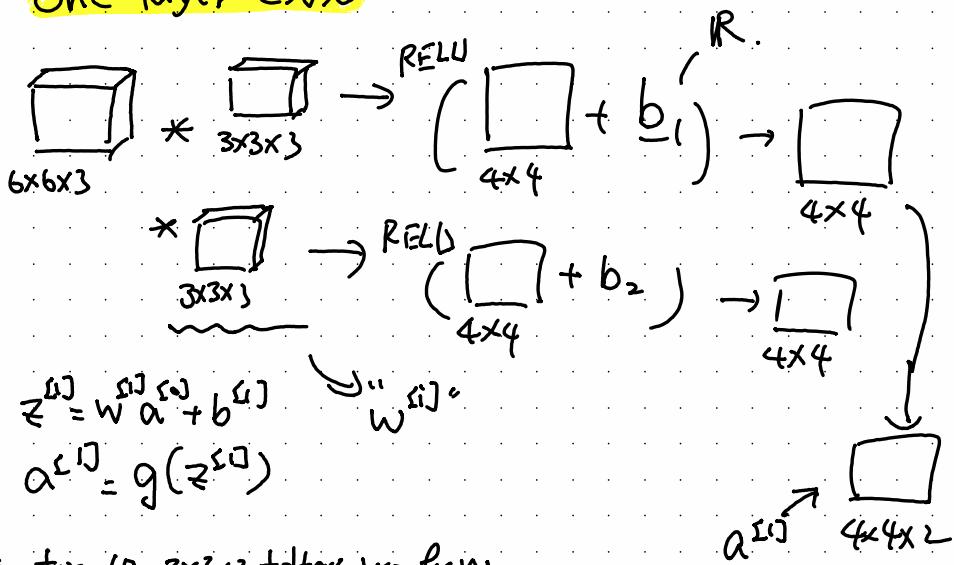
$$\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \end{array} * \begin{array}{c} \text{filter} \\ \text{3x3x3} \\ \text{---} \end{array} = \begin{array}{c} \text{intermediate output} \\ \text{4x4x1} \\ \text{---} \end{array}$$

vertical

$$\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \end{array} = \begin{array}{c} \text{final output} \\ \text{4x4x2} \\ \text{---} \end{array}$$

$$n \times n \times n_c \times f \times f \times n_c := n-f+1 \times n-f+1 \times n_c$$

One-layer CNN



e.g. for 10 $3 \times 3 \times 3$ filters, we have

$$(3 \cdot 3 \cdot 3 + 1) \times 10 = 280 \text{ parameters}$$

* parameters are independent from image resolution

Notation

$f^{[l]}$ - filter

$p^{[l]}$ - padding

$s^{[l]}$ - stride

$n_c^{[l]}$ - number of filters
each filter.

$$f^{[l]} \times f^{[l]} \times n_c^{[l]}$$

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

output: $n_H^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

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

Activation

$$\alpha^{[l]} \rightarrow n_H^{[l]} \times n_w^{[l]} \times n_c^{[l]}$$

$$A^{[l]} \rightarrow m \times n_H^{[l]} \times n_w^{[l]} \times n_c^{[l]}$$

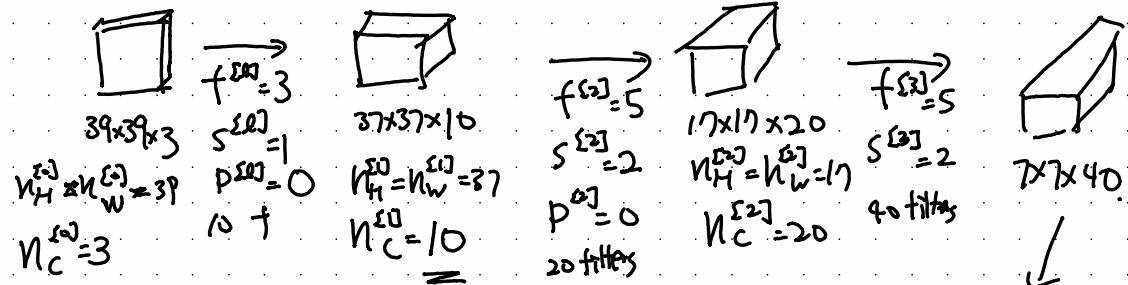
weights

$$f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$$

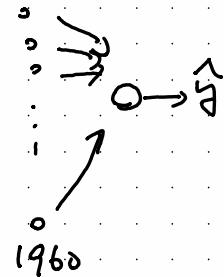
bias:

$$n_c^{[l]} = (1, 1, 1, n_c^{[l]})$$

CNN



$$\frac{W+2p-f}{S} + 1$$

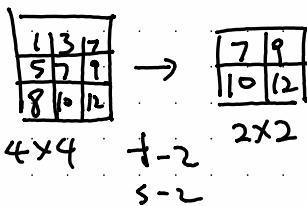


TYPE of LAYERS in CNN

- convolution (Conv)
- pooling (Pool)
- fully connected (FC)

POOLING LAYER

Max pooling - preserve feature?



△ no parameters to learn △

Average pooling - not often used

hyperparameters

f - filter size $f=2, s=2$

s - stride $f=3, s=2$

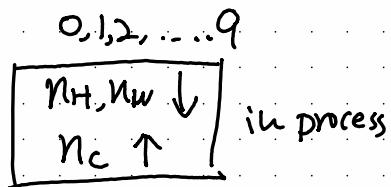
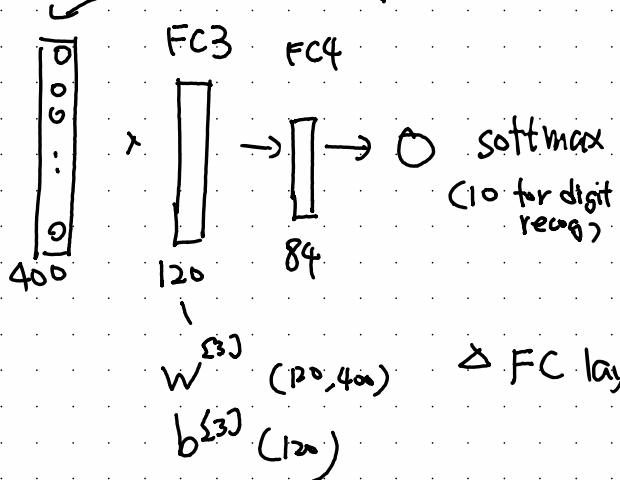
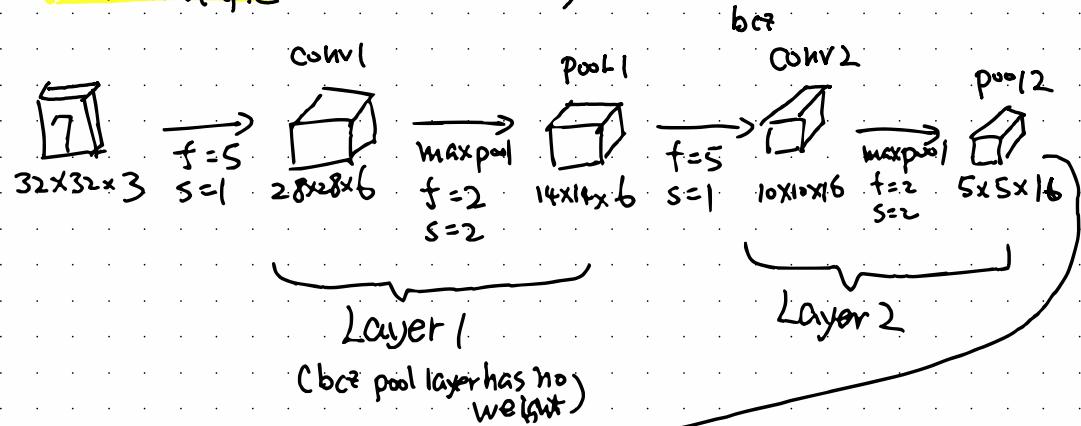
type of pooling layer

$$\left\lfloor \frac{n-f}{s} + 1 \right\rfloor$$

* No parameter to learn

NN example

(LeNet - 5)



△ FC layers have lots parameters!

* Why Conv? *

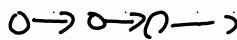
Reduce #parameters compared

- parameter sharing - feature detector useful for all parts of image
- sparsity of connections - each output value depends on small part of image

△ translation invariant.
(pooling layer)

Keras

- sequential API

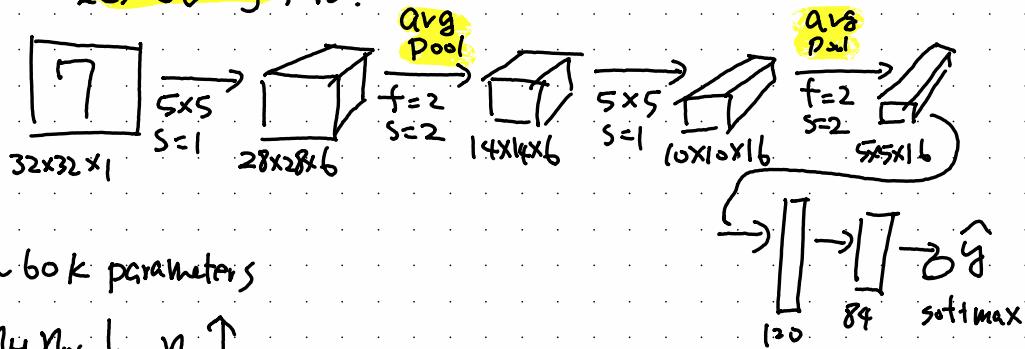


- functional API



CASE STUDY

LeNet-5, 98.



~60k parameters

$n_H n_W \downarrow n_c \uparrow$

conv pool / conv pool + fc pool

original

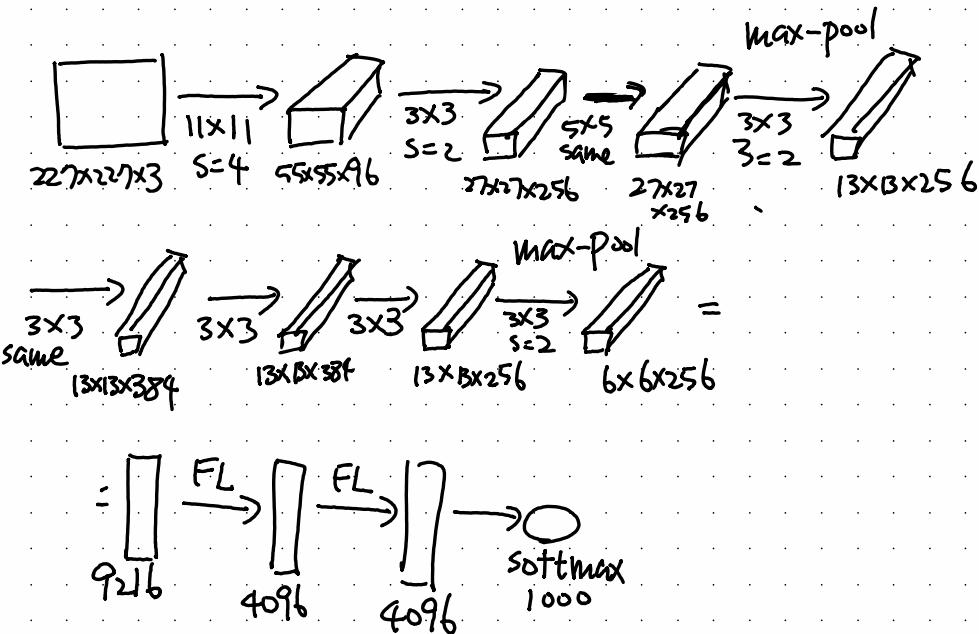
Sigmoid/tanh

different filter

nonlinearity after pooling

section II, III

AlexNet. 2012



$\sim 60M$ parameters

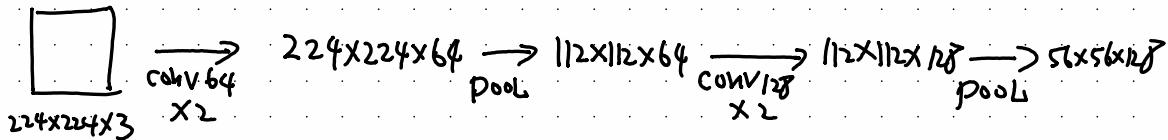
- ReLU

Orig

- multiple GPU
- local response normalization (LRN)
 - normalize along axis
 - less often used.

VGG-16, 2015

CONV = 3x3, S=1, same, MAX-Pool = 2x2, S=2



CONV 256 → POOL → CONV 512 - POOL?

... → FC → FC → softmax
4096 4096 1000

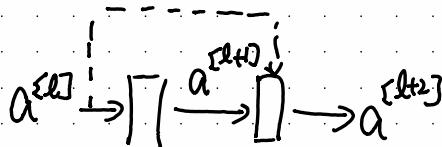
~138M parameters LARGE!

~Uniform

Δ w h n w ↓ n c ↑

Residual Network (ResNets) 2015

Residual Block.

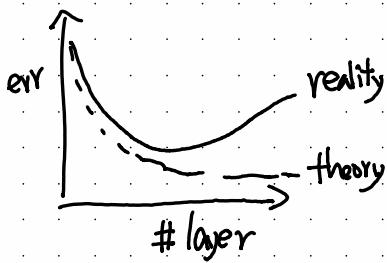


SHORT CUT / SKIP CONNECTION

* $a^{[l]}$ $\xrightarrow{\text{linear}} \text{ReLU} \xrightarrow{\text{linear}} \text{ReLU} \xrightarrow{\oplus} \text{ReLU} \rightarrow a^{[l+2]}$

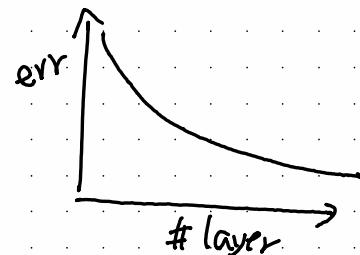
$$z^{[l+1]} = w^{[l+1]} a^{[l]} + b^{[l+1]}, a^{[l+1]} = g(z^{[l+1]}), z^{[l+2]} = w^{[l+2]} a^{[l+1]} + b^{[l+2]}, a^{[l+2]} = g(z^{[l+2]})$$

plain



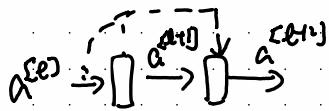
$$\text{ResNet} \Rightarrow a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

Res



△ help training deeper network.

Why ResNet Work?



$$a^{l+2} = g(z^{l+2} + a^{l+1})$$

$$= g(w^{l+2} a^{l+1} + b^{l+2} + a^{l+1})$$

IF $\begin{matrix} \downarrow \\ 0 \end{matrix}$ $\begin{matrix} \downarrow \\ 0 \end{matrix} \Rightarrow a^{l+2} = g(a^{l+1})$

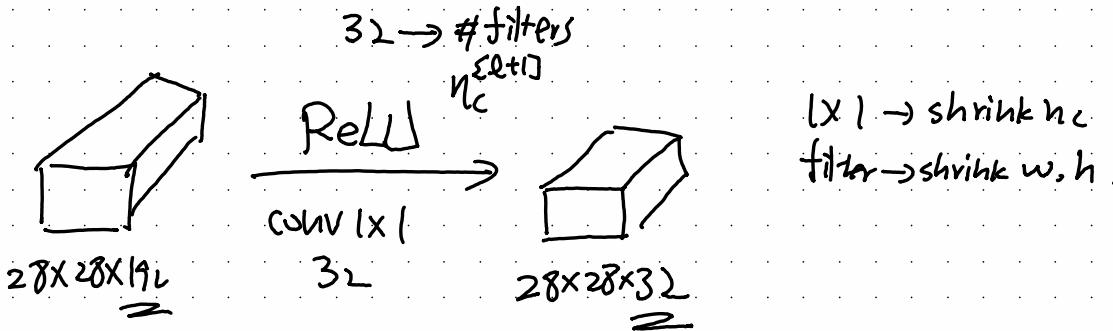
→ identity function is easy to learn

→ ① res block doesn't hurt.

→ ② prevent deeper layers making result worse.

1×1 convolution / network in network

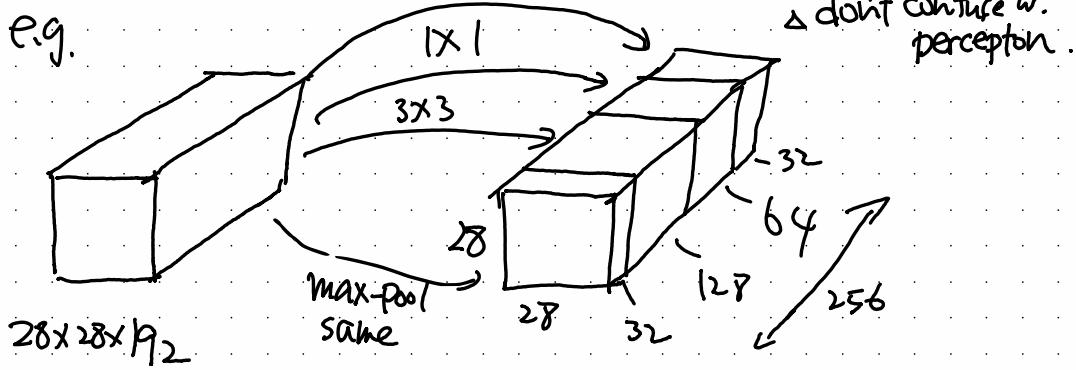
$$\begin{array}{c} \text{3D tensor} \\ 6 \times 6 \times 32 \end{array} * \begin{array}{c} \text{filter} \\ 1 \times 1 \times 32 \end{array} \xrightarrow{\text{ReLU}} = \begin{array}{c} \text{3D tensor} \\ 6 \times 6 \times \# \text{filters} \end{array}$$



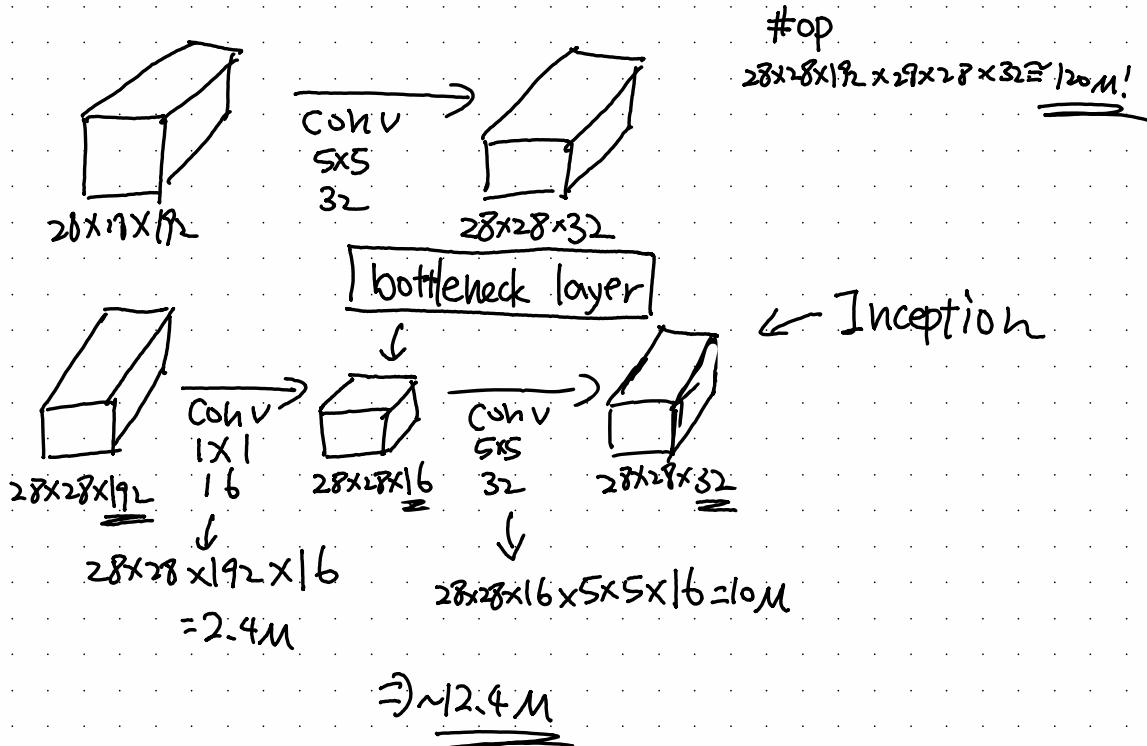
shrink n_c , if we like.

Inception Network, 2014

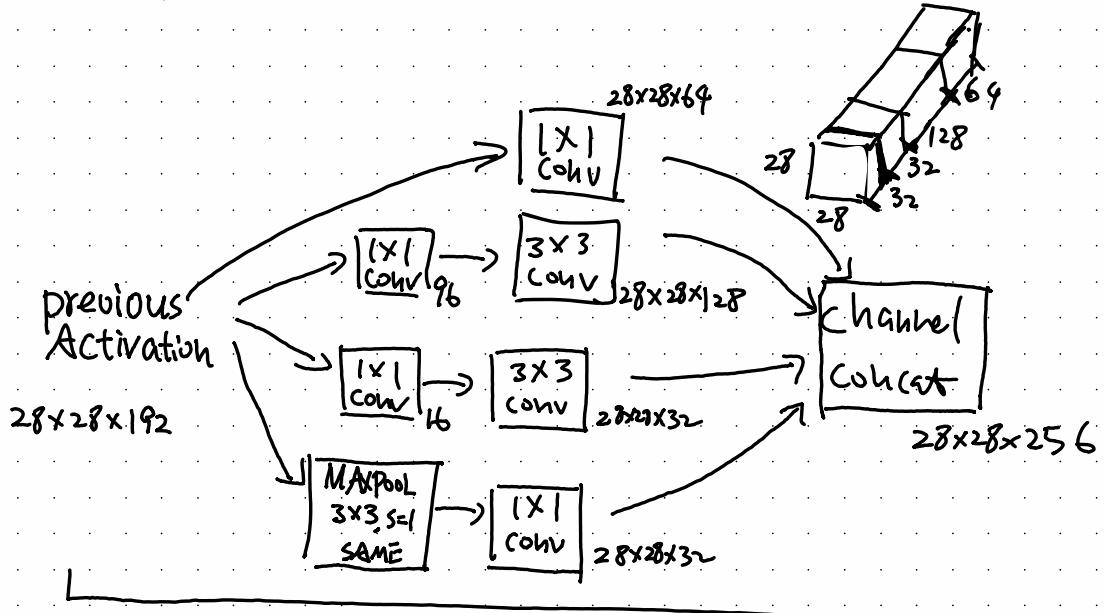
e.g.



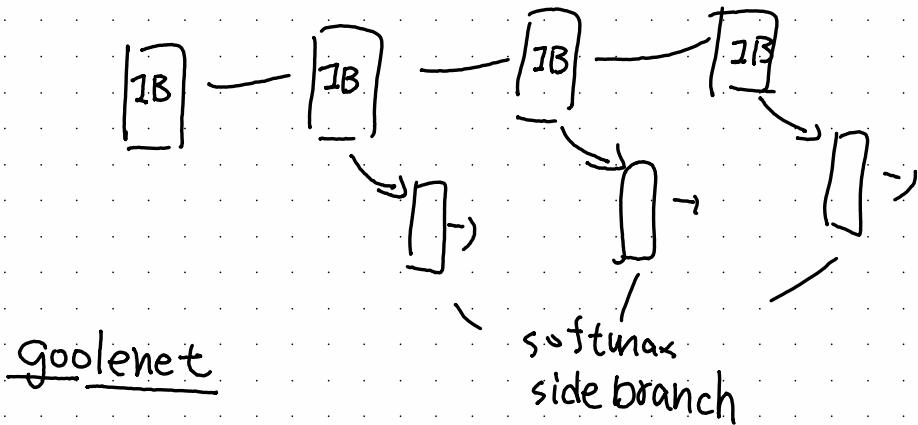
△ try all different size at once!



Inception Module



Inception Block



goolenet

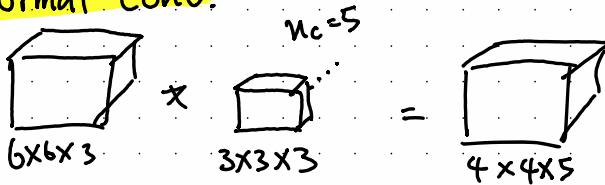
"we need to go deeper"

Inception

MobileNet, 2017

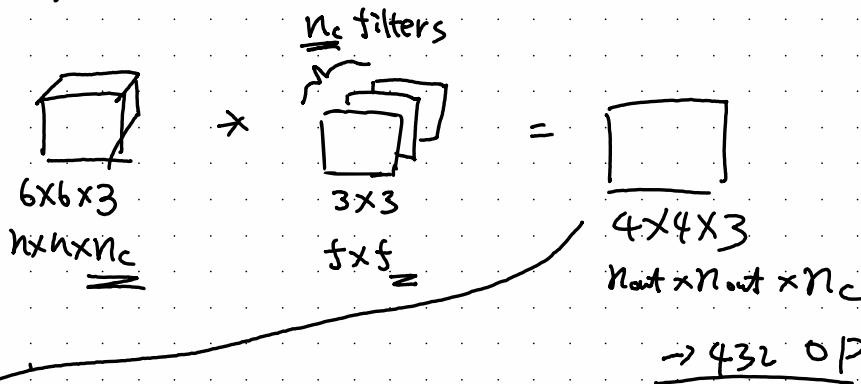
- low computational cost.
- Normal v.s depthwise - separable convolution
 \Rightarrow depth-wise + point-wise.

Normal Conv.

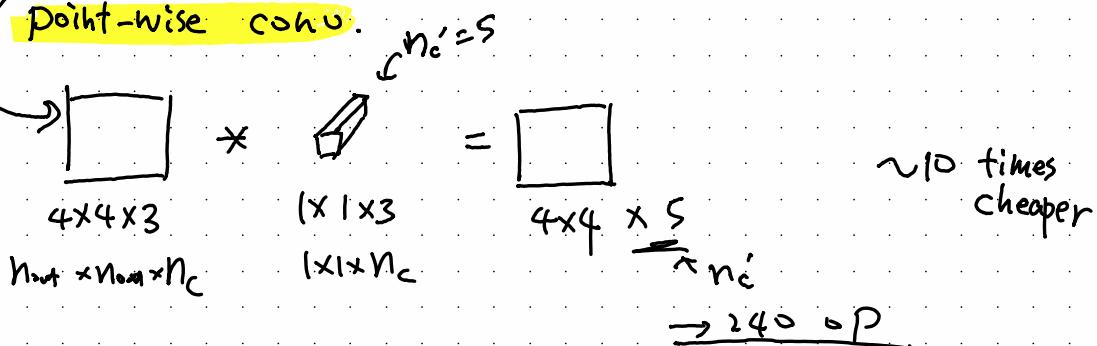


$$\text{cost} = \# \text{filter para} \times \# \text{filter pos} \times \# \text{filter} \\ = 2160 \text{ op.}$$

Depth-wise Conv.

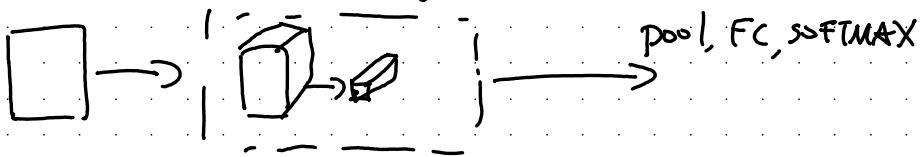


Point-wise conv.



MobileNet V1, 2017

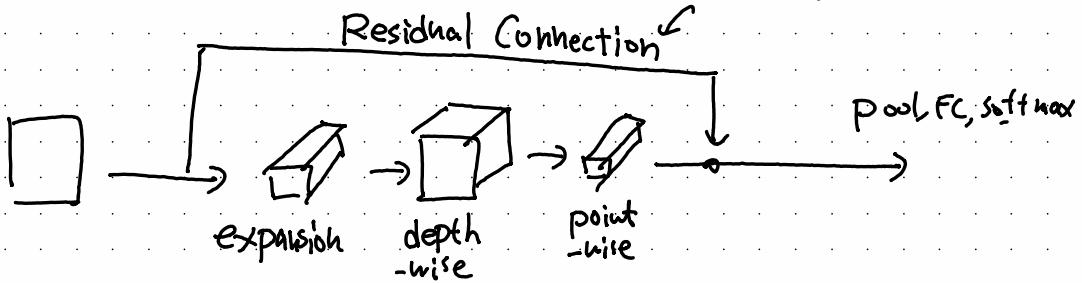
13 times



MobileNet V2, 2019

bottleneck blocks

17 times



eg. $n \times h \times 3 \rightarrow n \times h \times 18 \rightarrow n \times h \times 18 \rightarrow n \times h \times 3$

expansion

(1×3)
18 filters

depth
-wise



point-wise

$1 \times 1 \times 18$
3 filters

(projection)

- expansion → increase representation in bottleneck block

- projection → reduce bandwidth between layers

EfficientNet. 2019

- automatic scale up/down model

scale r resolution

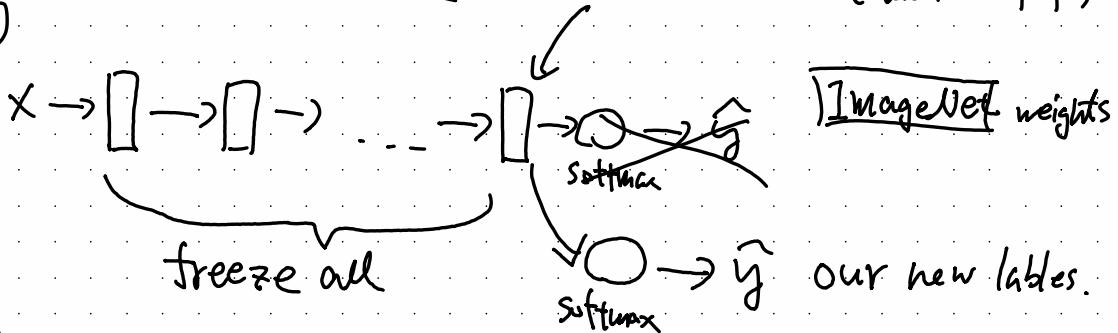
d depth

w width the model selects r,d,w that suits device.

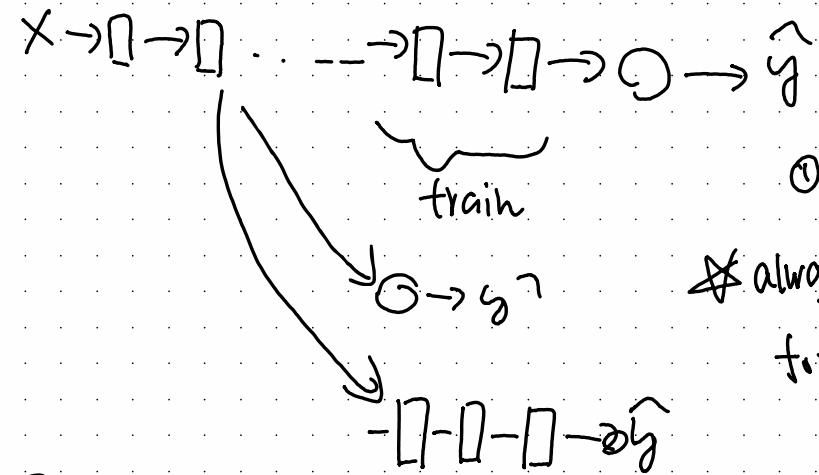
Transfer Learning

precompute this, to accelerate
(Avoid forward prop.)

①



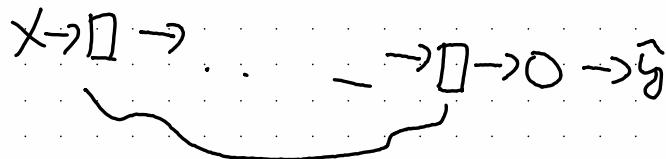
②



①, ② if small dataset.

* always start from existing weights.
for CV problems

③

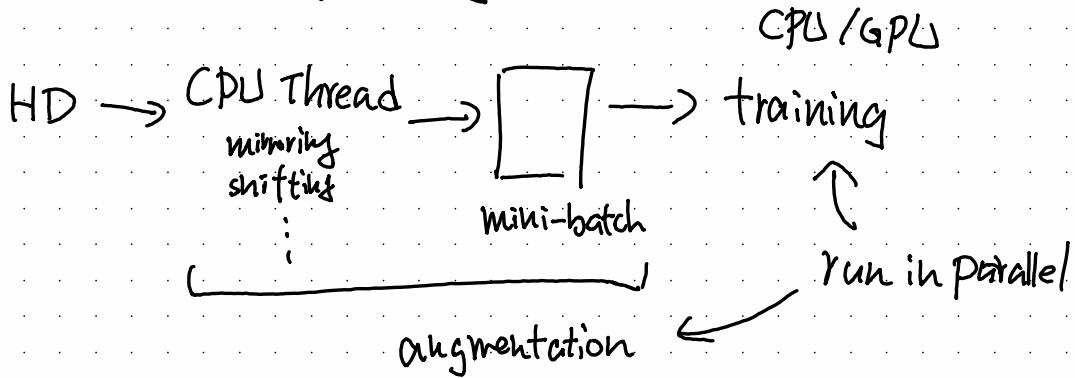


use previous weights
as initialization
(if we have a large data set)

Data Augmentation

- mirroring
- random cropping
 - might crop some less informative area
- color shifting
 - PCA color augmentation
 - simple RGB shifting, e.g. different sunlight.

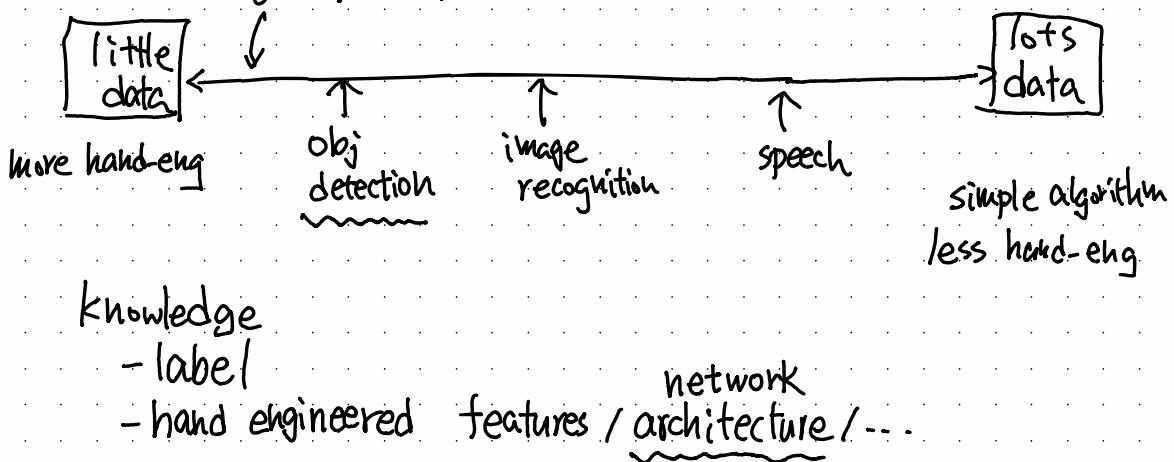
Distortion during training



State for CV.

data vs hand-engineering
(hack)

your specific problem (TL)



Tips for doing well on benchmarks

- Ensembling (3~15 networks)
 - train several networks and average outputs (\hat{y})
Class useful in practical)
- Multi-crop at test time
 - run classifier on multiple vers of test images and avg results
e.g. (10-crops)

Use Open source code

- use network arch.
- use implementation
- use pre-trained model

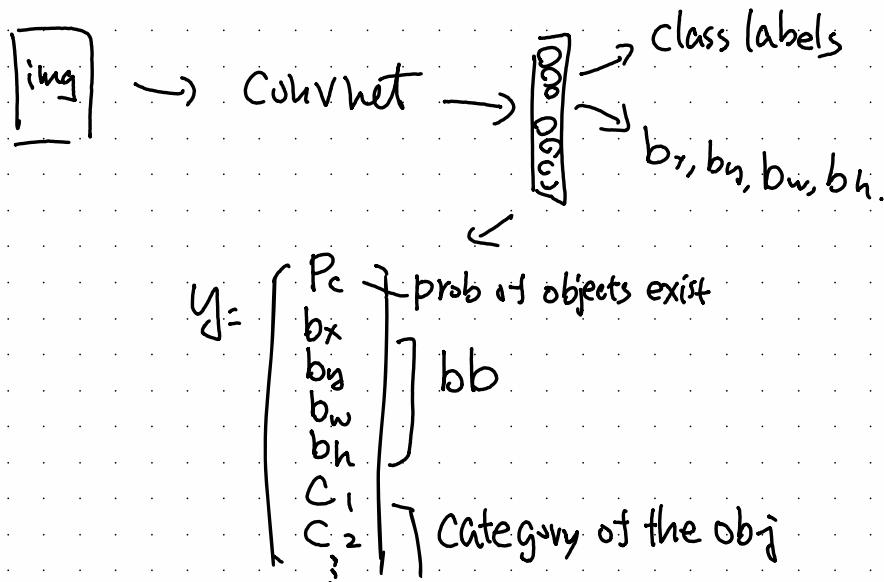
Object detection

classification - car, cat etc

localization - bounding box

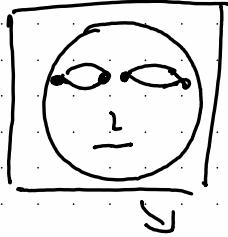
detection - multiple objs of different categories

classification with localization



$$L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + \dots + (\hat{y}_j - y_j)^2 & \text{if } y_i = 1 \\ (\hat{y}_1 - y_1)^2 & \text{if } y_i = 0 \end{cases}$$

Landmark detection



eyes: l_{1x}, l_{1y}
 l_{2x}, l_{2y}
 l_{4x}, l_{4y}

additional landmarks e.g. mouse.
nose
:

Cuhvel



each out a landmark on face

- same for pose detection

skeleton landmarks in a photo.

Obj detection

sliding window detection

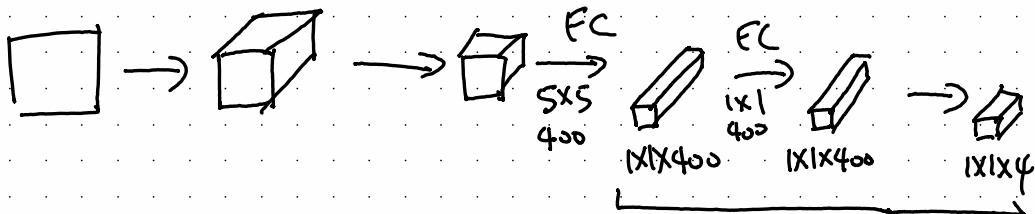
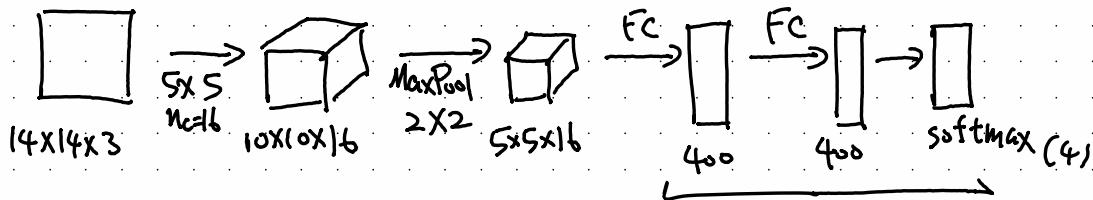


send each window through ConvNet.
from small to large window size

- high computation cost.

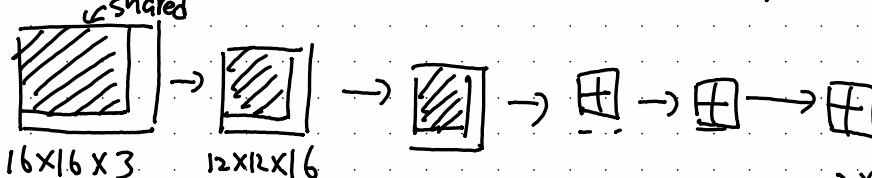
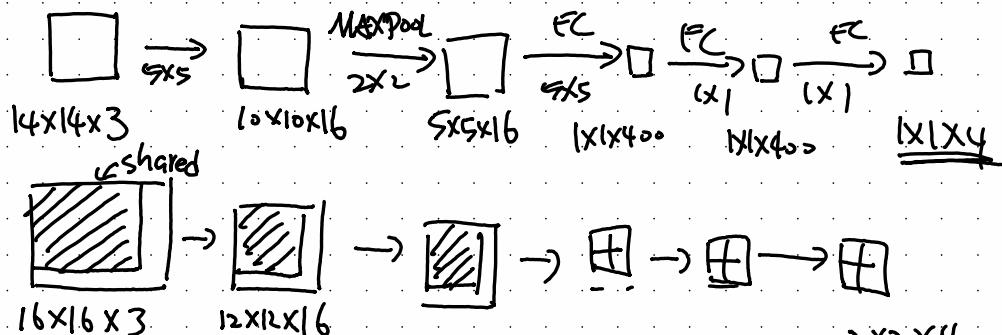
Convolutional Implementation

turn FC to convolutional layer



OverFeat. 2014.

Conv implementation

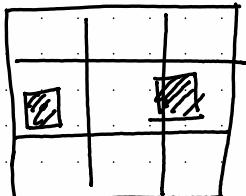


- reuse overlapped computation result sequentially
- instead of send each sliding window through ConvNet, simultaneously do prediction for all window position
- the window position is not precise though.

Output accurate bounding boxes

YOLO - You only look once, 2015 (a difficult paper to read)

(100)



100

target output

grid 3x3x8
y P, b, c

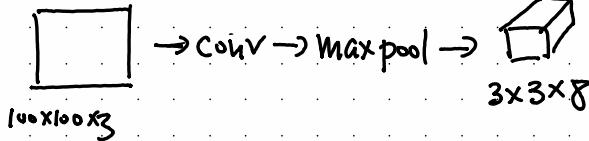
for each grid cell

$$y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_w \\ b_h \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$$

- relative coord

to the grid cell

- b_x, b_y - between 0 and 1
 b_w, b_h - could be > 1



⇒ increase grid resolution if need higher precision.

~ still one ConvNet implementation

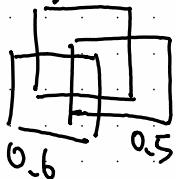
Intersection over Union (IoU)



correct if $\text{IoU} > 0.5$

Non-max Suppression

0.7



find max \rightarrow suppress others with high IoU

each grid cell $\begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$

discard all boxes with $p_c \leq 0.6$

for remaining boxes

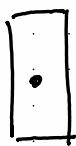
- pick largest p_c

discard remaining with $\text{IoU} \geq 0.5$ with other boxes.

Anchor boxes



predefined shapes



box 1



box 2

Anchor
box 1

Anchor
box 2

$$\Rightarrow y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

Anchor box 1

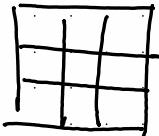
Anchor box 2

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

Anchor box 1

each obj assigned to cell
contain mid-point.
- Anchor box for the cell with high IoU

Yolo algorithm

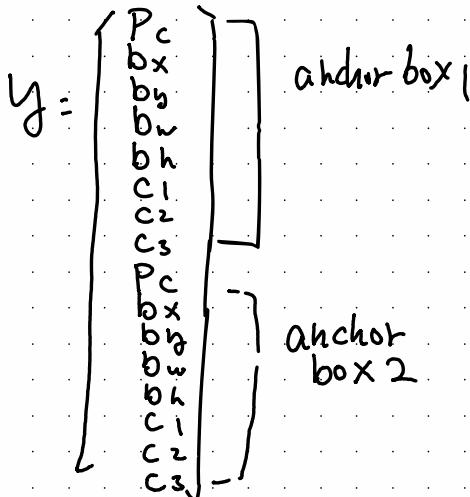


C

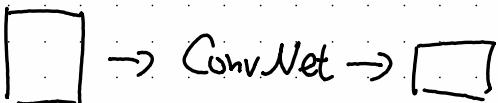
- 1 - pedestrian
- 2 - car
- 3. - motorcycle

$$y \quad 3 \times 3 \times 2 \times 8$$

anchor boxes 5 + #classes

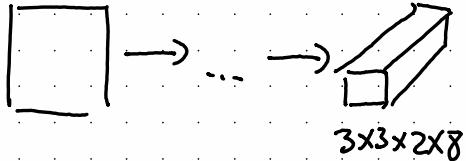


Training



$$3 \times 3 \times 16$$

Prediction



$$3 \times 3 \times 2 \times 8$$

$$y = \begin{bmatrix} 1 \\ 0.2 \\ 0.3 \\ 0.2 \\ 0.1 \\ 0 \\ 0 \\ \vdots \end{bmatrix}$$

non-max suppression

for each cell,

- get 2 predicted bb (for each anchor boxes)
- get rid of low p_c prediction
- for each class, run non-max suppression

Region Proposals. R-CNN. 2013

- pick interesting region to run CNN

segmentation → run CNN on interesting regions

R-CNN: propose region, CNN, output

Fast R-CNN : propose region, Convolution implementation, output
C20/5)

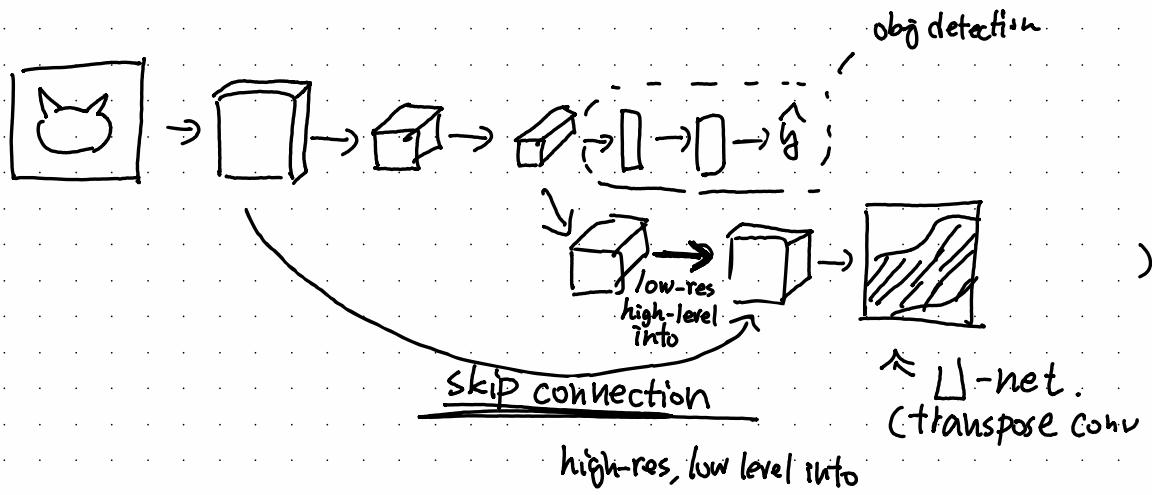
Faster R-CNN : conv net to propose regions (2016)

still slower than YOLO...

Semantic Segmentation U-net, 2017

obj detection - bb

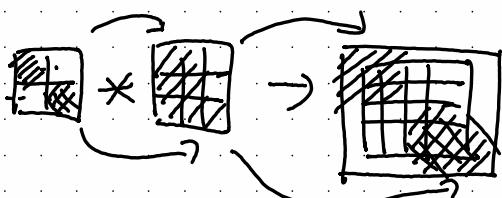
semantic seg - per-pixel labeling



Transpose Conv, deconvolution, upscale conv

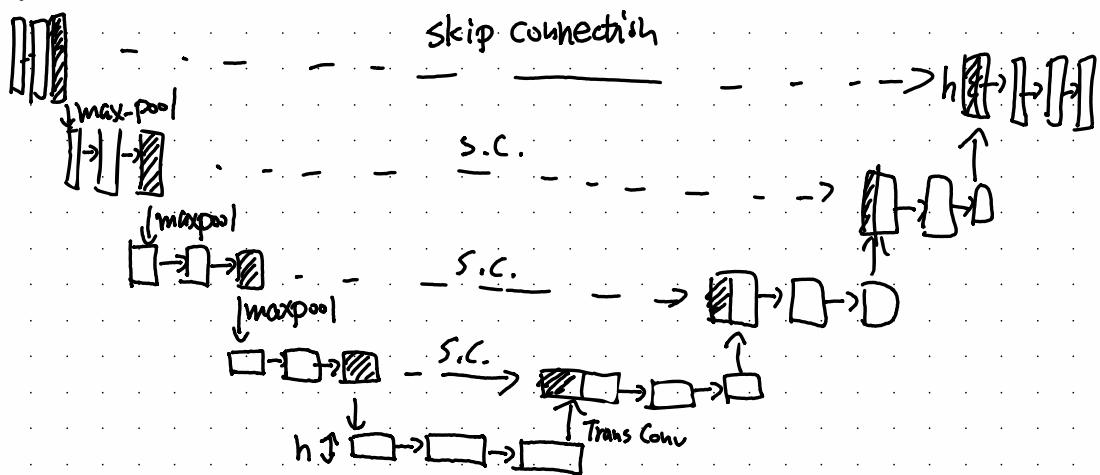
$$\begin{matrix} \text{2x2} & * & \begin{matrix} \text{3x3} \end{matrix} & \rightarrow & \begin{matrix} \text{4x4} \end{matrix} \end{matrix}$$

$$f=3, p=1, s=2$$



if overlap, add values together.

U-Net, 2015 (on biomedical image seg.)



Face Recognition

- verification

- input image, name
- output whether image/name match
(1 to 1)

- Recognition

- database of k people
- input image
- output ID of the image among k.
(1 to k)

one-shot learning

learn from one example to recognise the person again

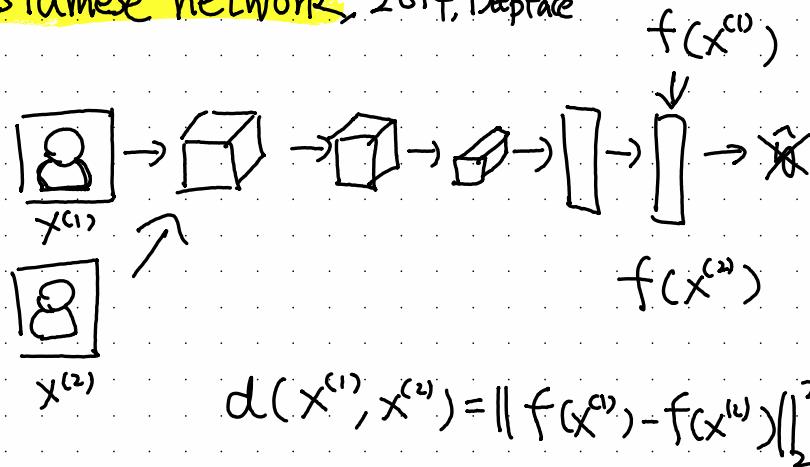
learn a "similarity" function

$d(\text{img1}, \text{img2})$ = degree of difference between imgs.

if $d(\text{img1}, \text{img2}) > \tau$, same person

$< \tau$, no match

Siamese network, 2014, DeepFace



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

parameters of NN define $f(x^{(i)})$

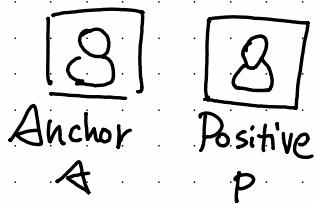
learn parameters that

if $x^{(i)}, x^{(j)}$ are same person $d(x^{(i)}, x^{(j)})$ small

different

large

Triplet Loss FaceNet, 2015



$$\text{WANT: } \|f(A) - f(P)\|^2 + \alpha \leq \|f(A) - f(N)\|^2$$

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

$\underbrace{}$
margin

to prevent trivial solution

$$\text{e.g. } f(\text{img}) = 0$$

$$\text{or } f(\text{img}^1) = f(\text{img}^2),$$

Loss Function

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

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

training set: 10K pictures of 1K persons

Choosing A, P, N

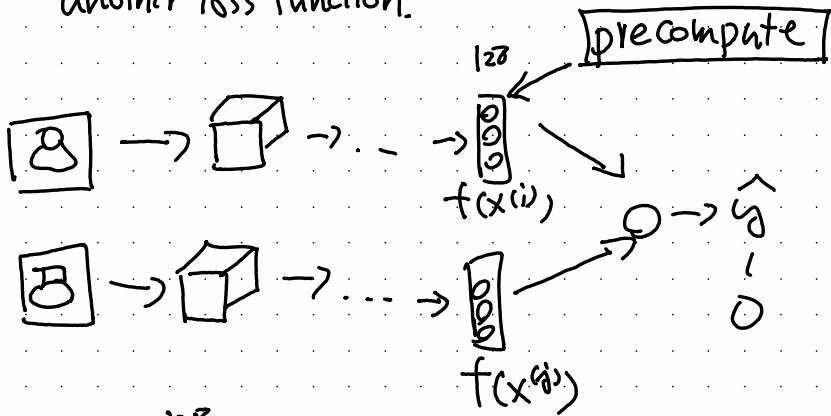
if chosen randomly, $d(A, P) + \alpha \leq d(A, N)$ is too easy.

choose A, P, N that are "hard" (NN learns little)
(check paper)

Current K = 1M images
commercial

DeepFace, 2014

another loss function.



$$\hat{y} = \sigma \left(\sum_{k=1}^{128} w_k |f(x^{(i)})_k - f(x^{(g)})_k| + b \right)$$

- kth feature or x^2

Neural Style Transfer, 2015

Content + style \rightarrow generated
(C) (S) (G)

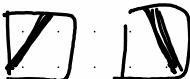
Visualize ConvNet 2013

hidden unit covers larger patch
 \rightarrow

[8] \rightarrow conv \rightarrow conv $\rightarrow \dots \rightarrow$ FC \rightarrow FC $\rightarrow \hat{y}$

pick a unit in layer 1, find image patches that maximize
unit activation
repeat for all units
patch size \uparrow , to deeper layer.

lower level -



higher level



Neural Style Transfer. 2015

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

1. init G randomly (like noise)
2. gradient descent to minimize $J(G)$

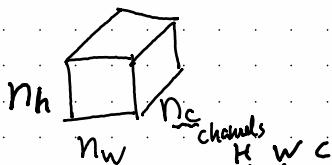
$J_{\text{content}}(C, G)$

- use hidden layer l to compute content cost (usually middle of network)
- use pre-trained ConvNet.
- Let a^{enc} and $a^{\text{enc}(G)}$ be activation of layer l

$$J_{\text{cont}}(C, G) = \frac{1}{2} \| a^{\text{enc}} - a^{\text{enc}(G)} \|^2$$

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

at layer l , define "style" as correlation between activation across channels



intuition: each channel corresponds

Q_{ijk} = activation at (i, j, k) . $G^{\{l\}} : n_c^{\text{enc}} \times n_c^{\text{enc}}$ to hidden units that capture different levels of features

$$G_{kk'}^{\text{enc}(S)} = \sum_i \sum_j Q_{ijk}^{\text{enc}} Q_{ijk'}^{\text{enc}} \leftarrow \text{style} \quad \text{"Gram Matrix"}$$

$$G_{kk'}^{\text{enc}(G)} = \dots \leftarrow \text{generated}$$

$$J_{\text{style}}(S, G) = \| G^{\text{enc}(S)} - G^{\text{enc}(G)} \|_F^2 = \frac{1}{2n_H n_w n_c} \sum_k \sum_{k'} (G_{kk'}^{\text{enc}(S)} - G_{kk'}^{\text{enc}(G)})^2$$

Course 5

Sequence Model

Sequence Model

Sequence Model.

X: "Harry Potter" and "Hermione Granger" invented a spell

$x^{<1>} \quad x^{<2>}$

$x^{<9>}$

$T_x = 9$

y: 1 1 0 1 1 0 0 0

$y^{<1>} \quad y^{<2>}$

$y^{<9>}$

i-th training
example

$x^{(i)<t>} \quad T_x^{(i)}$

$T_y = 9$

$y^{(i)<t>} \quad T_y^{(i)}$

Representation

a		$x^{<1>}$	$x^{<2>}$
Aaron	1	0	0
:	2	:	:
and	:	0	0
harry	4075	1	0
:		0	1
potter	6830	1	0
:		0	0
Zulm	20000	0	0

10,000

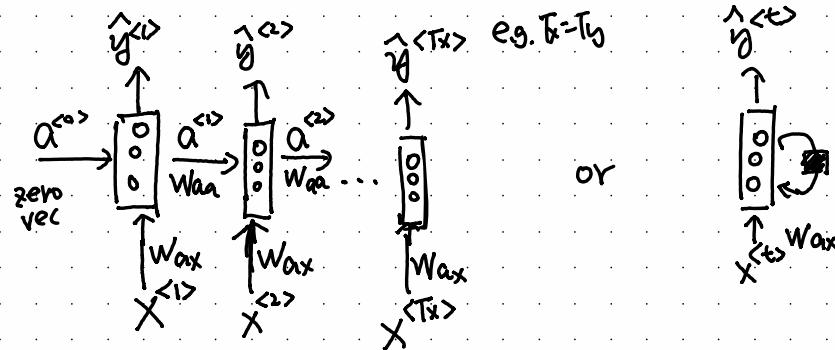
one-hot.

~30,000 for
commercial software

Recurrent NN

why not standard network

- input / output could be of different length
- doesn't share features across different pos.



/imitation: only use info "before" t

Δ B-RNN

$$a^{(0)} = 0 \quad a^{(1)} = g(W_{aa}a^{(0)} + W_{ax}x^{(1)} + b_a) \leftarrow \text{tanh, ReLU}$$
$$\hat{y}^{(1)} = g(W_{ya}a^{(1)} + b_y) \leftarrow \text{sigmoid}$$

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

$$\hat{x}^{(t)} = g(W_{aa}\hat{a}^{(t-1)} + W_{ax}\hat{x}^{(t)} + b_a)$$

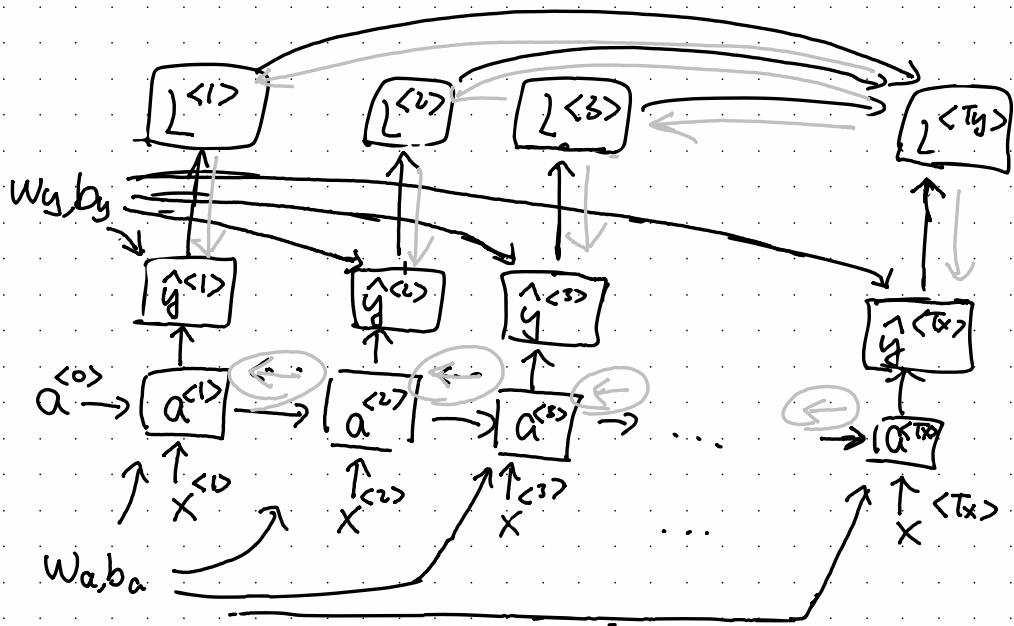
$$\hat{y}^{(t)} = g(W_{ya}\hat{a}^{(t)} + b_y)$$

\Downarrow simplify

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

$$\underbrace{\begin{bmatrix} 100 \\ W_{aa}; W_{ax} \\ 100 \\ 10000 \end{bmatrix}}_{\underline{W_a}} = \underline{W_a} \quad \begin{bmatrix} \hat{a}^{(t-1)} \\ \hat{x}^{(t)} \end{bmatrix} \begin{matrix} \uparrow 100 \\ \downarrow 10,000 \end{matrix} \quad \begin{bmatrix} 10,100 \\ 10,100 \end{bmatrix}$$

Backprop. through time

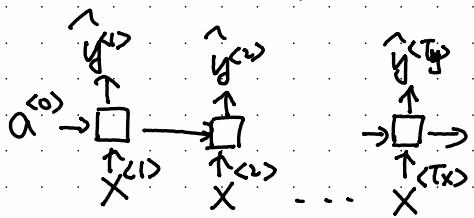


$$L^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -\hat{y}^{<t>} \log \hat{y}^{<t>} - (1 - \hat{y}^{<t>}) \log (1 - \hat{y}^{<t>})$$

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

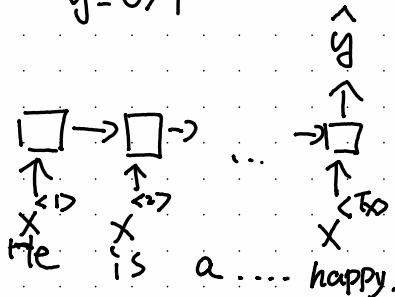
Different RNN types

$$T_x = T_y$$



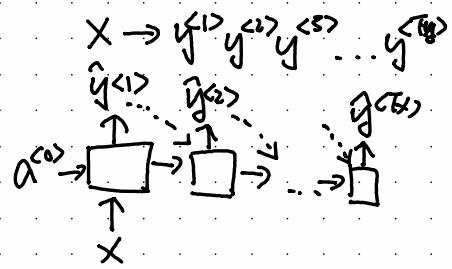
many-to-many

e.g. sentiment detection
 $X = \text{text}$
 $Y = 0/1$



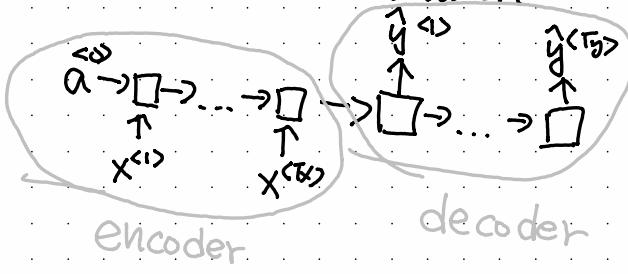
many-to-one

music generation



One-to-many

machine translation



many-to-many

Language Model

Speech Recognition

the apple and pair salad

the apple and pear salad \leftarrow higher p, although both sound the same

language model: $P(y^{<1>} \dots y^{<n)})$?

Training Set: large corpus of text.

- An apple a day ... doctor away. <EOS>

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

- The Egyptian Mau is a type of cat

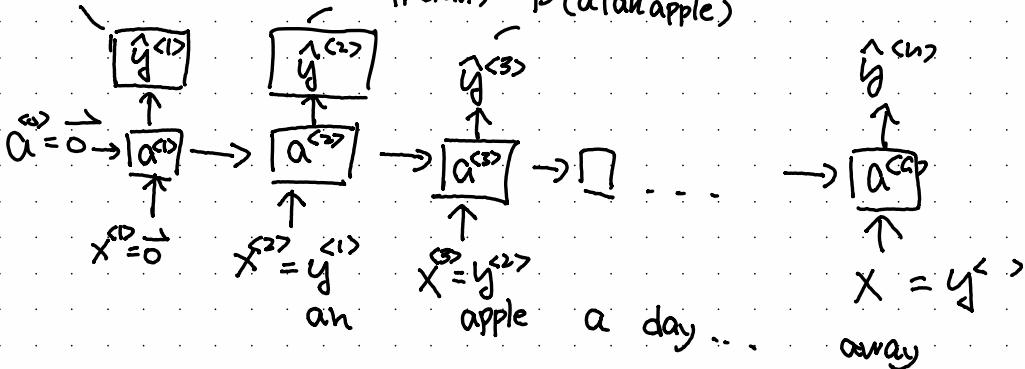
corpus

0	a
:	ab
0	acat
:	z
0	zebra

\downarrow replace
<UNK>

unique token for unknown word.

$$P(\text{apple}|\text{an}) \quad P(\text{an}|\text{apple})$$



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

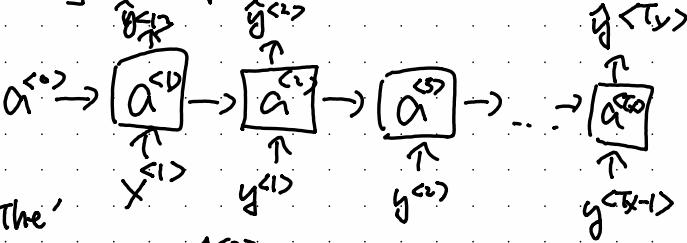
$$\mathcal{L}(y^{<1>} | y^{<2>} | y^{<3>}) = -\sum_t y_i^{<t>} \log p_i^{(t)}$$

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

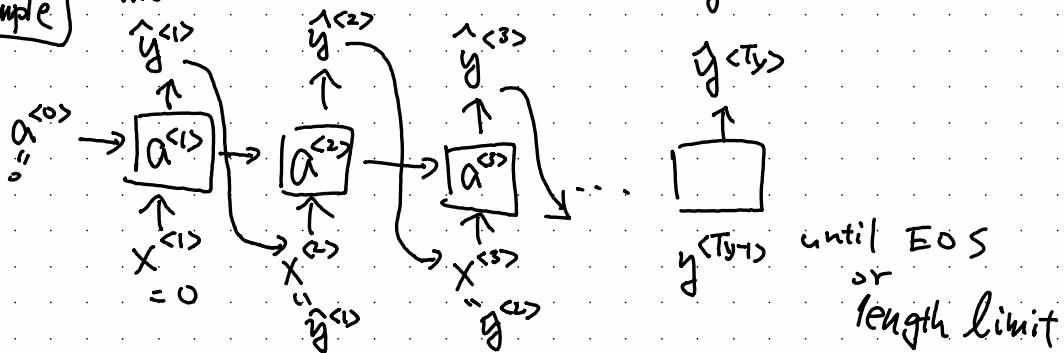
- generate random sentence

Sampling a seq from trained RNN

Trained



sample



< np.random.choice >

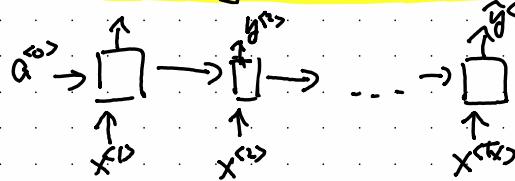
until EOS
or
length limit

- can also do "character-level" language model

- very long
- slower to train

- word level model is more popular.

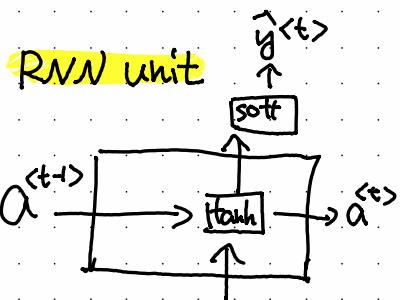
Vanishing gradients with RNN



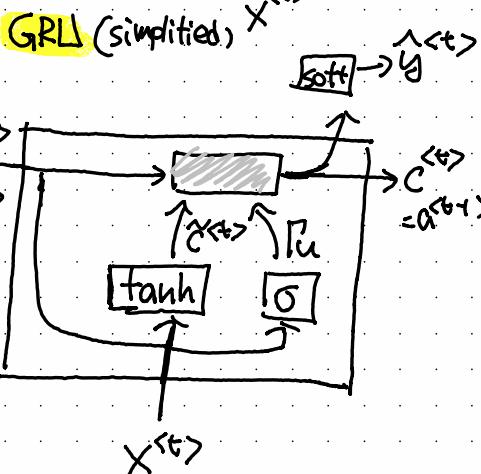
e.g. cats, which ~~are~~ are cute
is cute
- RNN cannot memorize long-term influence

exploding gradient - solve with clipping

GRU - Gated Recurrent Unit, 2014



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



C = memory cell

$$\underline{C}^{(t)} = a^{(t)}$$

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

$$\Gamma_u^{(t)} = \sigma(W_u[C^{(t-1)}, x^{(t)}] + b_u)$$

update

$$\underline{C}^{(t)} = \Gamma_u \times \tilde{C}^{(t)} + (1 - \Gamma_u) \times \underline{C}^{(t-1)}$$

(gamma)

the cat, which ate .. . is full

$$\begin{aligned} C^{(t)} &= 1 \dots = \\ \Gamma_u &= 0, 0 \dots , 1 \end{aligned}$$

GRL (Full)

simpler, and can build bigger model

$$\tilde{h}^{(t)} = \tanh (W_c[\tilde{r}_r \times C^{(t-1)}, x^{(t)}] + b_c)$$

$$l_u = \sigma (W_u[C^{(t-1)}, x^{(t)}] + b_u)$$

$$l_r = \sigma (W_r[C^{(t-1)}, x^{(t)}] + b_r) \quad - \text{relevance}$$

$$h^{(t)} = l_u \times \tilde{h}^{(t)} + (1 - l_u) \times C^{(t-1)}$$

$$C^{(t)} = a^{(t)}$$

LSTM, 1997 more complex, but more powerful.

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

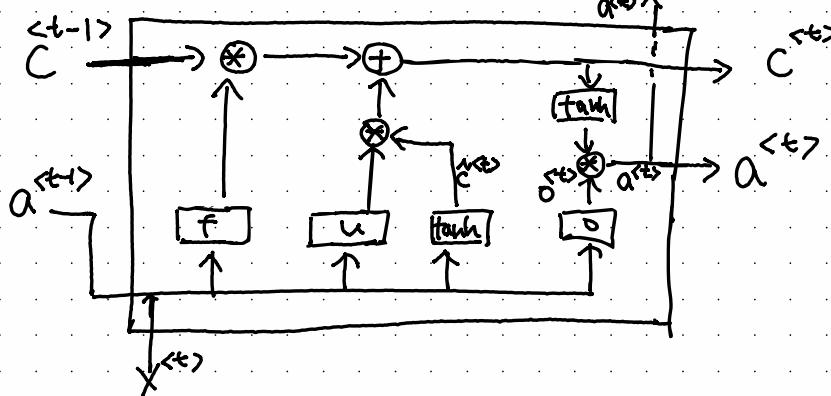
$$l_u = \sigma (W_u[a^{(t-1)}, x^{(t)}] + b_u) \quad - \text{update} \quad \begin{matrix} \text{peephole connection} \\ C^{(t-1)} \text{ affects} \end{matrix}$$

$$l_f = \sigma (W_f[a^{(t-1)}, x^{(t)}] + b_f) \quad - \text{forget}$$

$$l_o = \sigma (W_o[a^{(t-1)}, x^{(t)}] + b_o) \quad - \text{output}$$

$$C^{(t)} = l_u \times \tilde{C}^{(t)} + l_f \times C^{(t-1)}$$

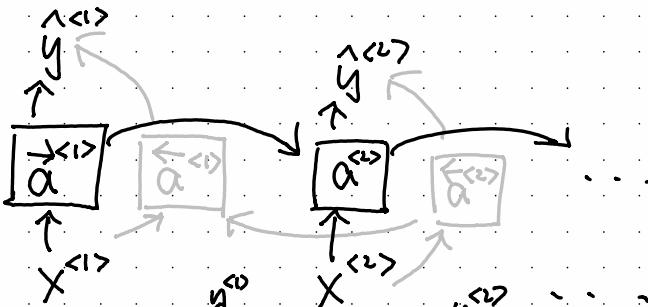
$$a^{(t)} = l_o \times \tanh (C^{(t)})$$



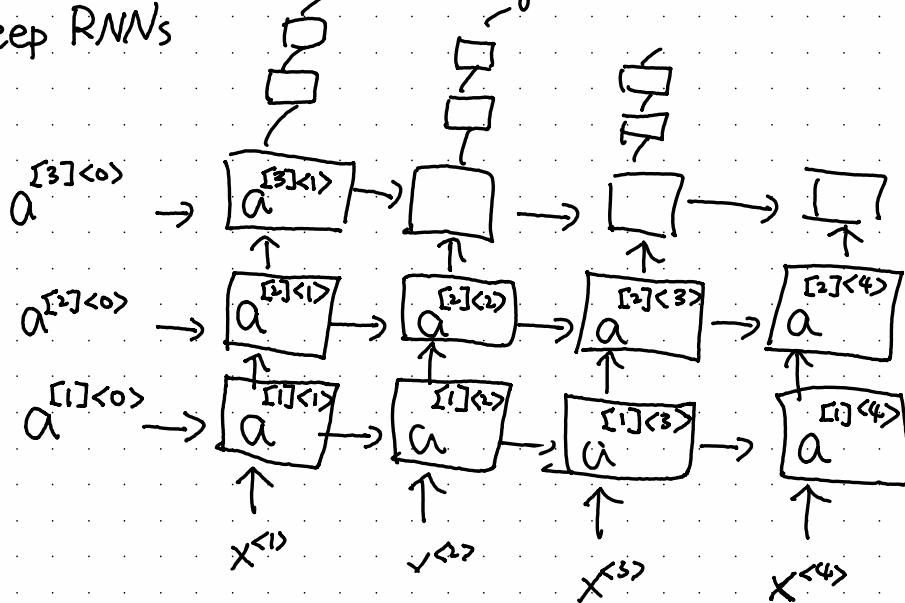
$\tilde{C}^{(t-1)}$ affects l_u, l_f, l_o

Bidirectional RNN (BRNN)

He said Teddy bears are cute
getting into from future
He said Teddy Roosevelt was a good man



Deep RNNs



$$a^{[i,j]<3>} = g(W_a(a^{[i,j]<2>}, \underline{\underline{\underline{a}}}) + b_{a^{[i,j]}})$$

Word Embeddings

one-hot representation

$$32 \rightarrow \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 0 \\ ? \\ ? \end{bmatrix} \leftarrow \text{dot prodct is } 0, \text{ doesn't reflect semantic distance}$$

apple orange

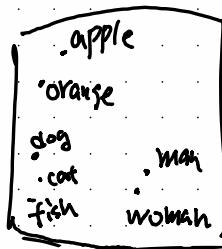
featurized representation: word embedding

	man	woman	apple	orange	...
gender	-1	1	0	0	
age	0.02	0.03	0.02	0.03	
size	0.4	0.4	0.01	0.02	
food	0.01	0.02	0.95	0.94	
e.g.					
word					

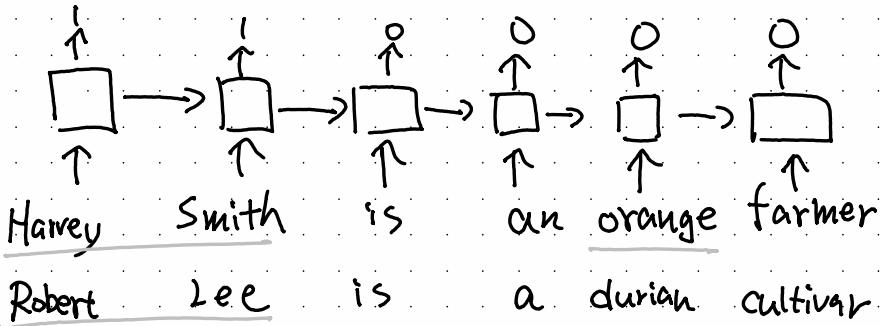
Vishalize word embedding
embed 3D to 2D for viz

H-SNE

"embed" words into feature space



ex: Name entity recognition



↑
name not corporate name

↑
what if we have fruit
not trained before?

Transfer Learning

1. learn from large text corpus
(or download pre-trained embedding online)
2. transfer embedding to new tasks w. smaller set
3. (optional) continuously finetune word embedding

properties of word embedding, 2013 Mikolov

e.g. man → woman

king → ?

$$E_{\text{man}} - E_{\text{woman}} \approx E_{\text{king}} - E_{?}$$

find word w

$$\arg \max_w \underline{\text{sim}}(E_w, E_{\text{king}} - E_{\text{man}} + E_{\text{woman}})$$

30~75% accuracy in papers.

	man	woman	king	queen
gender	:	:	:	:
age	:	:	:	:
royal	:	:	:	:
:	:	:	:	:

Cosine Similarity

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

$$\text{man:woman} = \text{boy:girl}$$

$$\text{RMB:China} = \text{Rupee:India}$$

Embedding Matrix

300
(features)

E

1,000 x
(words)
corpus

$E \cdot O = []_{(300, 1)}$

The diagram illustrates the matrix multiplication $E \cdot O$. On the left, there is a large matrix labeled E with dimensions 300 (features) by 1,000 (words in corpus). A bracket indicates that the columns represent words. An arrow points from this bracket to the product $E \cdot O$, which is shown as a 300x1 vector with dimensions (300, 1).

$E \cdot O_j = e_j$ - embedding for word j
(in practise, AP2 extract the column)

Learning Word Embeddings

history: difficult \rightarrow simple

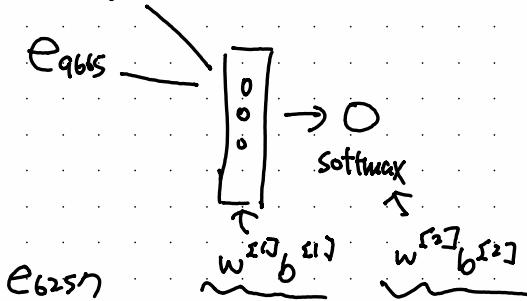
Neural Language Model, 2003, Bengio

$$I \quad O_{4343} \rightarrow E \rightarrow e_{4343}$$

$$\text{want } O_{9665} \rightarrow E \rightarrow e_{9665}$$

want
a
glass
of
Orange

Juice



if use previous 4 words

$$\text{learn } E, W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]} \quad 4 \times 300 = 1200 \text{ dim}$$

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

Context:

last 4 words

last / next 4 words

last 1 word

"nearby" 1 word (skip gram)

Word2Vec, 2d3, Mikolov

skip-gram

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

<u>context</u>	<u>target</u>	<u>randomly chose</u>
orange	juice	
Orange	glass	
Orange	my	e.g. vocab size = 10000.

Model

$$x \rightarrow y$$

Context c → Target t
 "orange" "juice"

$$o_c \rightarrow E \rightarrow e_c \rightarrow o \rightarrow \hat{y}$$

$$\text{softmax} \quad e^{o_t^T e_c}$$

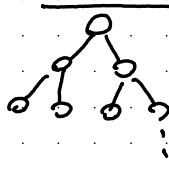
$$\text{Softmax: } P(t|c) = \frac{e^{o_t^T e_c}}{\sum_{j=1}^{10000} e^{o_j^T e_c}}$$

o_t : parameter associates with output t
 i.e. target t being the label

slow, so

$$L(\hat{y}, y) = - \sum_i y_i \log \hat{y}_i \quad y = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix}$$

hierarchical softmax



how to sample context c ?

too ~~common~~, the, a, of, ...

→ orange, apple...

$P(c)$ based on frequency

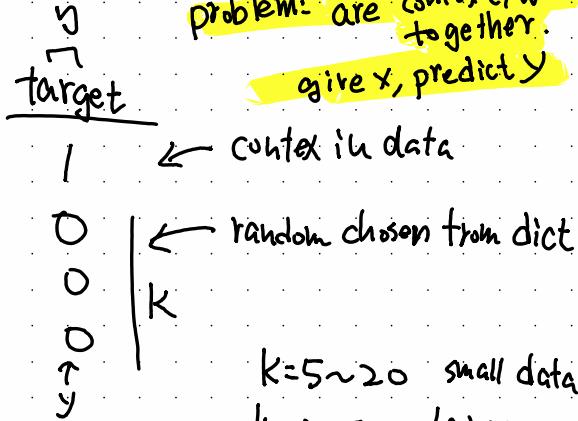
Negative Sampling

<u>Context</u>	<u>word</u>
Orange	juice
Orange	king
Orange	book
Orange	of
c ↑ t	t

2013, Mikolov.

problem: are context/word together.

give x, predict y



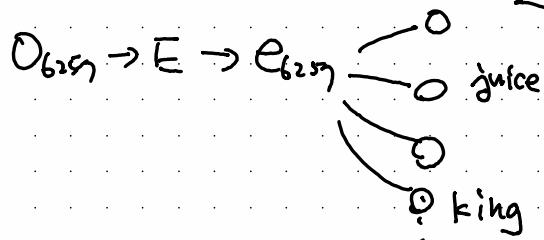
Model

softmax:

$$P(c|t) = \frac{e^{\theta_c^T e_t}}{\sum_j e^{\theta_c^T e_j}} \quad \leftarrow 10,000 \text{ softmax, slow!}$$

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

orange: O_{6257}



10,000 binary classification

each iteration, train K+1 classifier

Sampling:

freq. $P(w_i)$

$$\Rightarrow P(w_i) = \frac{f(w_i)^{\frac{3}{4}}}{\sum f(w_i)^{\frac{3}{4}}} \quad \text{heuristic}$$

uniform $\frac{1}{|V|}$

Glove Word, 2014. Pennington.
(global vector)

C, t

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

$$x_{ij}^* = X_{ij}$$

$$\min \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) (\Theta_i^T e_j + b_i + b_j' - \log x_{ij})^2$$

t c
↓ ↓
weighting

$$f(x_{ij}) = 0 \text{ if } X_{ij}=0 \quad \text{if } \log 0 = 0$$

Θ_i, e_j are symmetric

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

this is...
during
zulu

note:

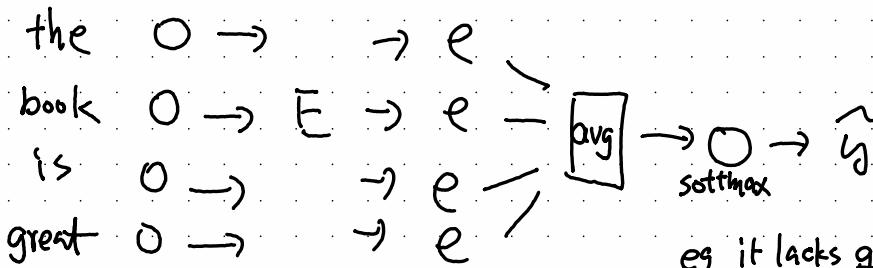
~~✓~~ e_i might not align with defined interpretable feature

axis of the embedding matrix

Sentiment classification

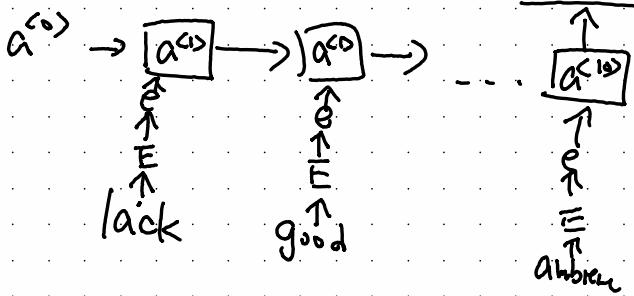
X		Y	
the book is great		****	data set is smaller.
the book lacks depth		*	
reads Ok		**	
:			

Model 1.



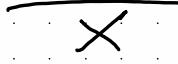
eg. it lacks good insight,
good depth, good taste
 problem!

Model 2 - RNN



many-to-one

Debiasing Word Embeddings, 2016.

Man = Computer Programmer as Woman = Homemaker 

Word embedding reflect gender, ethnicity used to train the model.

1. identify

identify bias axis

2. neutralize

project onto axis except definitive terms

3. equalize

Move definitive pairs to same distance around axis

Seq to Seq Model 2014

e.g. translation

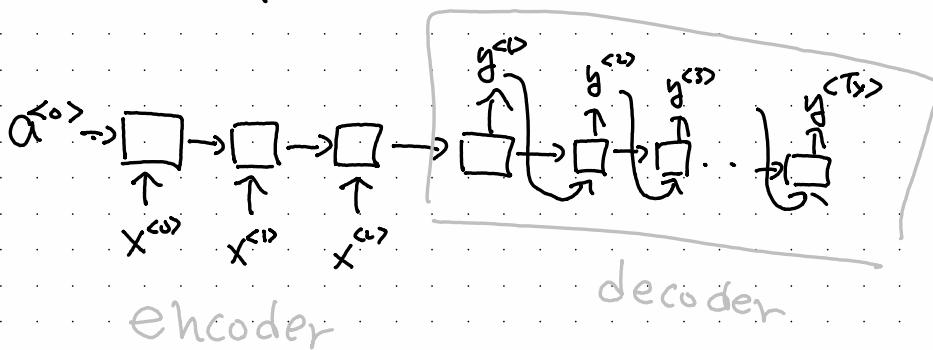
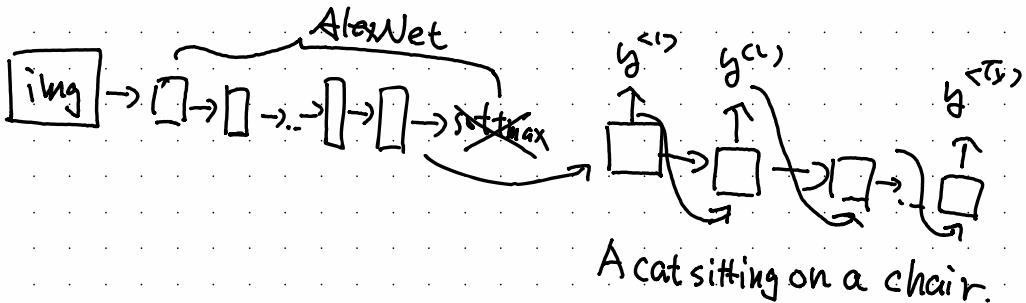


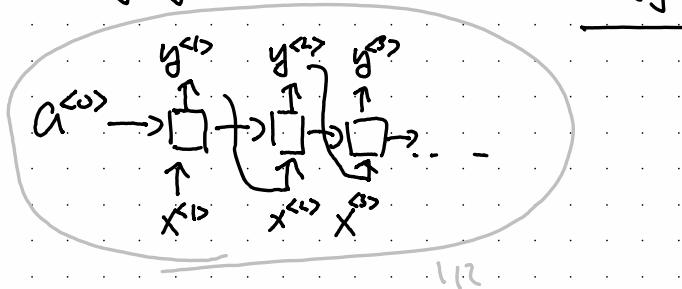
Image Captioning, 2014, 2015



seq2seq vs language model

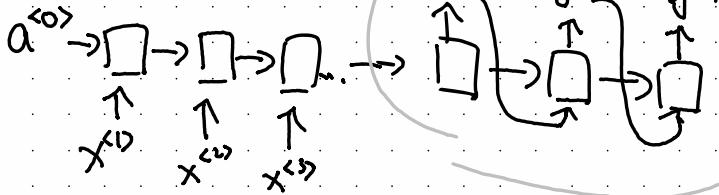
machine trans as building conditional language model.

Language Model:



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

Machine Trans.



Conditional
language model

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

e.g. $P(\text{English} | \text{French})$

$$\underset{y^{<1>} \dots y^{<n>}}{\operatorname{argmax}} P(y^{<1>} \dots y^{<n>} | x)$$

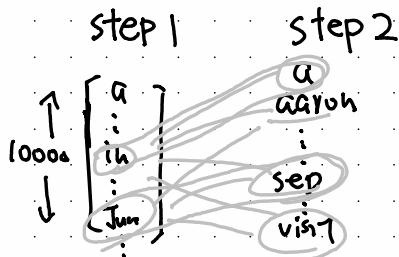
French sentence

(
beam search

△ greedy search usually result in worse sentence

e.g. picking common word like "going", "the", ... etc

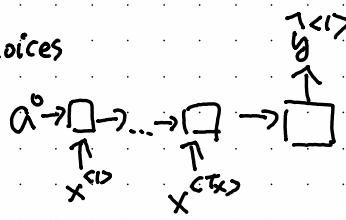
Beam Search



$B=3$ beam width \curvearrowright keep track of top B choices

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

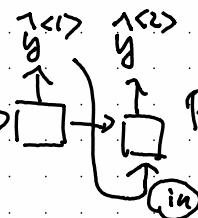
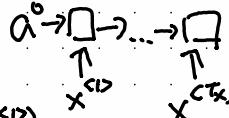
Memorize 3 good choices



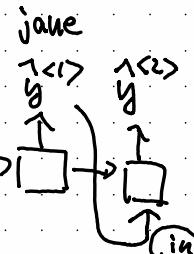
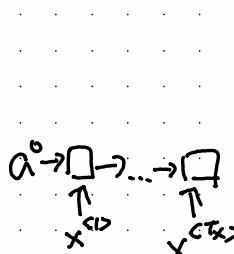
Step 2

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

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



$$P(y^{<2} | x, 'in')$$



$$P(y^{<2} | x, 'jane')$$

beam width

eventually pick $B=3$ most likely output
then moves on to next step.

step 3

in september



keep B NNs.

jane is

jane visit



$B=1$ ~ beam search
"greedy" search.

Refinement: Length normalization

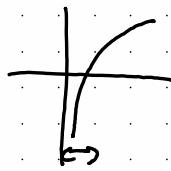
$$\underset{y}{\operatorname{argmax}} \prod_{t=1}^{T_y} P(y^{ct} | x, y^{<1>} \dots y^{<t-1>})$$

- multiplication of p mitigate value
- tend to pick short sentence

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

$$\Rightarrow \frac{1}{T_y} \sum_{t=1}^{T_y} \log P(y^{ct} | x, y^{<1>} \dots y^{<t-1>})$$

$$\alpha = 0.7$$



/ large $B \rightarrow$ slow, better
small $B \rightarrow$ fast, faster.

Production $B \sim 10$

research $= B \sim 100 \sim 1000$

to squeeze performance

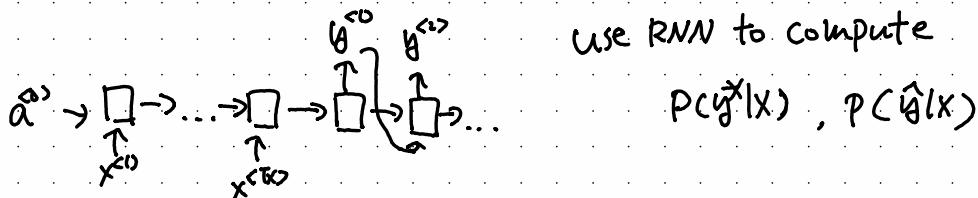
△ Compared to BFS, DFS, beam search doesn't guarantee best solution
- heuristic search

Error analysis in Beam Search.

Jane visite l'Afrique en septembre

Human: Jane visit Africa in September (y^*)

Machine: Jane visit Africa last September (\hat{y})



use RNN to compute

$$P(y^*|x), P(\hat{y}|x)$$

Case 1.

$$P(y^*|x) > P(\hat{y}|x)$$

beam search can improve

Case 2.

$$P(y^*|x) \leq P(\hat{y}|x)$$

RNN model is at fault.

→ y^* is better, but RNN says \hat{y} is better, thus RNN is at fault.



do the above analysis for all dev set,

then decide what to do next. - increase B

or

update RNN

Bleu score: bilingual evaluation underway, 2002

- deal with equally good translations

Bleu: Unigram

Sentence: Le chat est sur le tapis

ref 1 : the cat is on the mat

ref 2 : There is a cat on the mat.

:

clipping # occurrence

MT output: the the the ...

Count_{clip}('the')

precision: $\frac{7}{7}$

modified: $\frac{2}{7}$

Count('the')

Bleu score on bigrams

MT: The cat the cat on the mat.

	Count	Count _{clip}
the cat	2	
cat the	1	
the cat	1	
cat on	1	
on the	1	
the mat	1	

$$P_i = \frac{\sum_{\text{eg}} \text{Count}_{\text{clip}}(\text{n-gram})}{\sum_{\text{eg}} \text{Count}(\text{n-gram})}$$

n-gram

$$P_n = \frac{\sum_{\substack{\text{n-gram} \\ \text{eg}}} \text{Count}_{\text{clip}}(\text{n-gram})}{\sum_{\text{n-grams}} \text{Count}(\text{n-gram})}$$

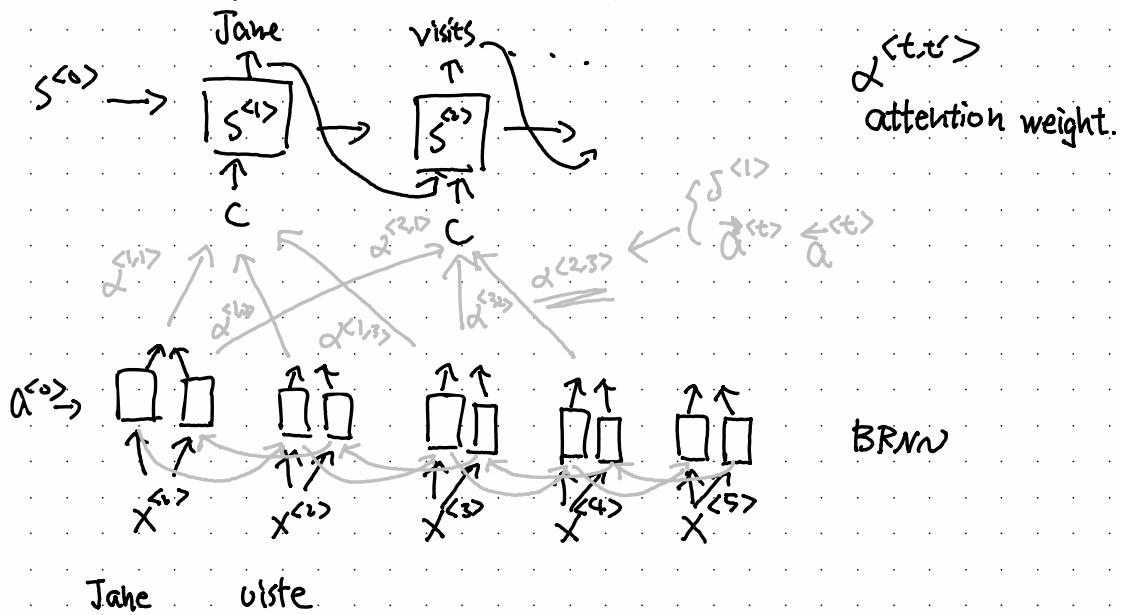
$$\text{Combined: } P = \exp\left(\frac{1}{n} \sum_{i=1}^n P_i\right)$$

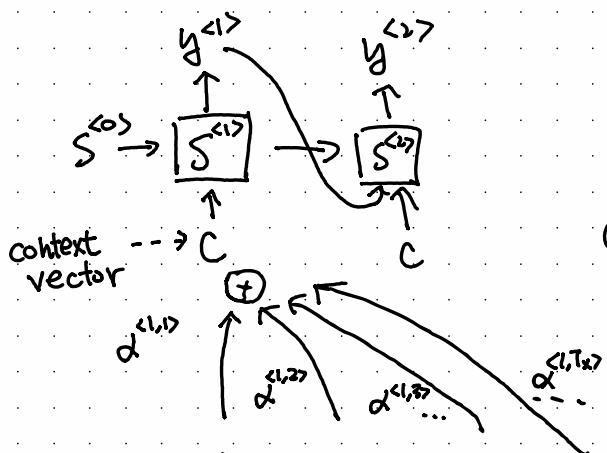
BP, brief penalty \rightarrow penalize short trans

$$B_P = \begin{cases} 1 & \text{if MT output longer than human} \\ e^{(1 - \frac{\text{human_len}}{\text{MT_len}})} & \end{cases}$$

Attention Model, 2014

- cannot handle long sequences



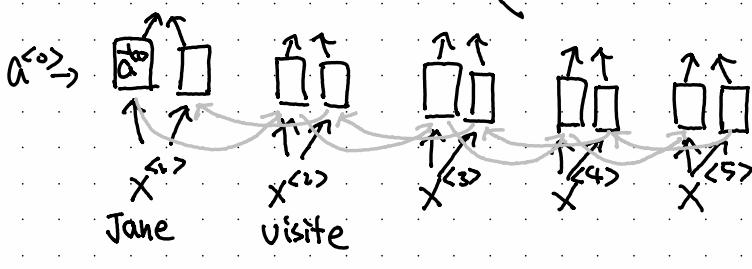


$$\alpha^{(t)} = (\alpha^{(t)}, \alpha^{(t)})$$

$$\sum_t \alpha^{(t)} = 1$$

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

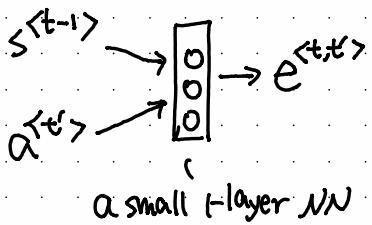
$\alpha^{(t,t')}$ is amount of attention
 $y^{(t)}$ pay to $a^{(t')}$



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

similar to softmax

runtime



$T_x \cdot T_y$ parameters
 $\alpha^{(t,t')}$
 usually T_x, T_y are not huge.

2015, image captioning, \rightarrow pay attention to a part of the picture

speech recognition

$X \rightarrow Y$
audio transcript

data: 300h ~ 3000h

Commercial: 100,000 h

- Attention model works well

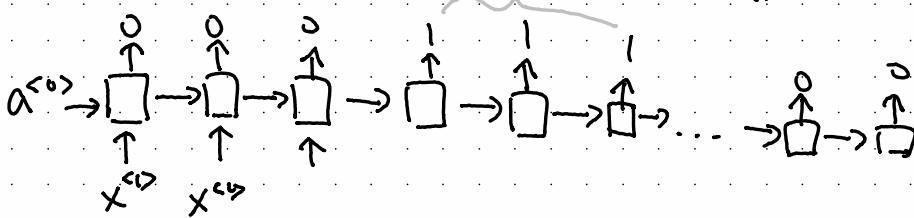
- CTC cost [2006]

Connectionist temporal classification tt...h_eee...w...ggg...
 $\Delta T_x \gg T_y$, e.g. 1000 Hz in audio \rightarrow the quick brown fox

\rightarrow collapse repeated chars not separated by "blank"

trigger word detection

multiple / for balanced
data set.



Transformer. 2017

RNN → GRU → LSTM

coh:

increase complexity
sequential

transformer:

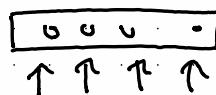
- attention + CNN

- self-attention. $A^{(1)} A^{(2)} \dots$

- multi-head attention. multiple copies

words

↑



parallel.

self-attention

transformer attention

$A(q, k, v)$ - attention-based vector representation

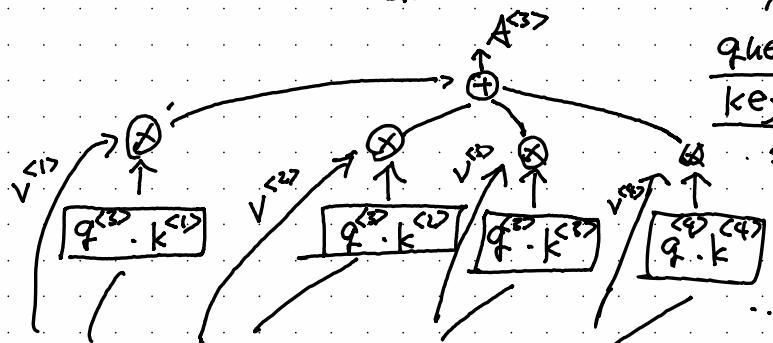
$$= \sum_i \frac{e^{q \cdot k^{(i)}}}{\sum_j e^{q \cdot k^{(j)}}} v^{(i)}$$

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

$$q^{(3)} = W^Q x^{(3)}$$

$$k^{(3)} = W^K x^{(3)}$$

$$v^{(3)} = W^V x^{(3)}$$



query - ask question

key - used to calculate

similarity between word and query

value - for plugging into the representation.

$q^{(1)}, q^{(2)}, q^{(3)}, q^{(4)}, v^{(1)}, v^{(2)}, v^{(3)}, v^{(4)}$
 $q \cdot k^{(1)}, q \cdot k^{(2)}, q \cdot k^{(3)}, q \cdot k^{(4)}$
 $q \cdot k^{(1)}, q \cdot k^{(2)}, q \cdot k^{(3)}, q \cdot k^{(4)}$
 $q \cdot k^{(1)}, q \cdot k^{(2)}, q \cdot k^{(3)}, q \cdot k^{(4)}$
 $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$
 Jane visite l'Afrique en Septembre

more flexible than fixed word embedding

Multi-head Attention

for example

#heads = 8 and? head 1 w_i^Q, w_i^K, w_i^V - what's happening

head 2 w_2^Q, w_2^K, w_2^V - when?

; ; - how?

MultiHead(Q, K, V)

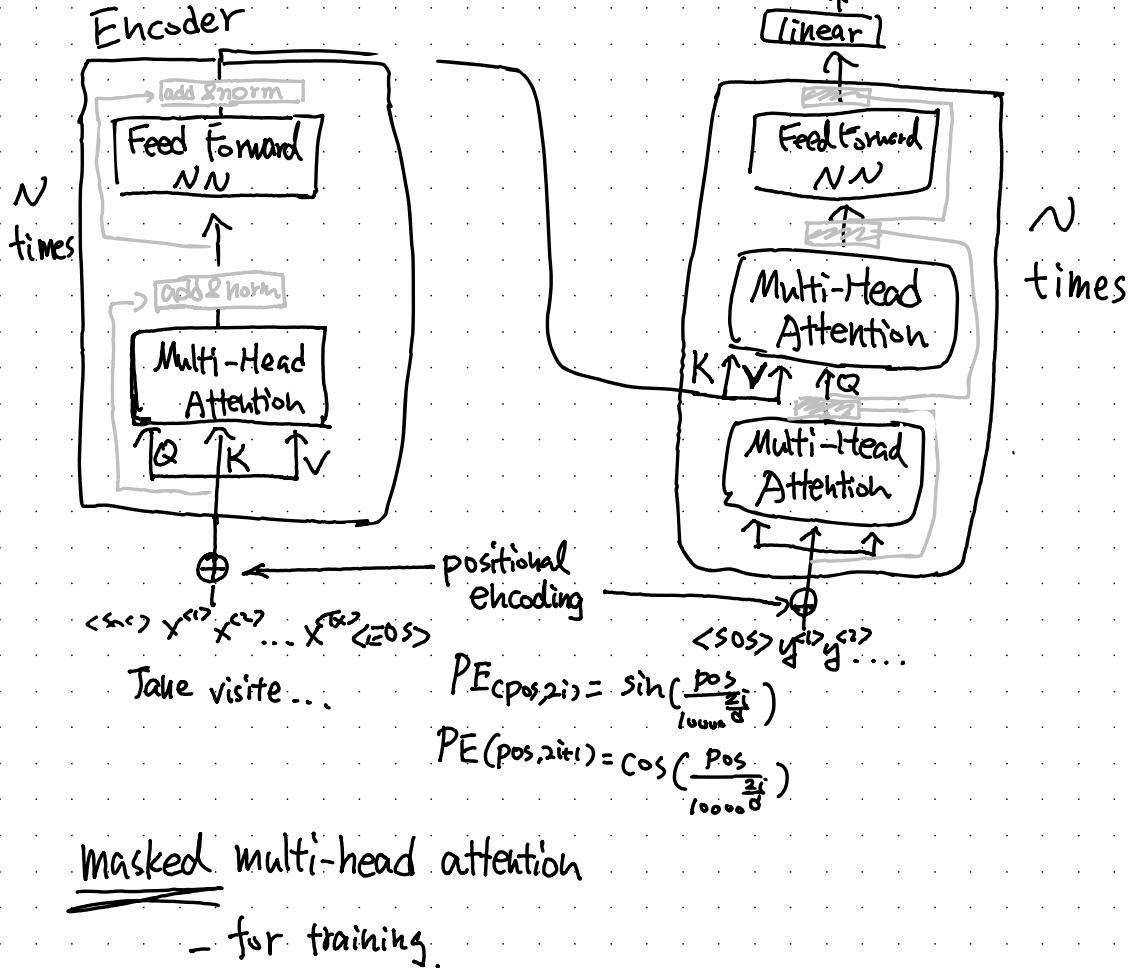


Attention($w_i^Q Q, w_i^K K, w_i^V V$)

$w_i^Q \leftarrow, w_i^K \leftarrow, w_i^V \leftarrow$

$Q \leftarrow, K \leftarrow, V \leftarrow$ $Q, K, V \rightarrow - - q, k, v \rightarrow$
 $X^{(1)} \rightarrow$ $X^{(2)} \rightarrow$ $X^{(3)} \rightarrow$
Jane visite

Transformer Network



masked multi-head attention

- for training.