

Deep Learning Basics

①

Basic - Neuron

$$\text{Output} = f \left(\sum_{i=1}^n w_i x_i \right)$$

f - activation function

w_i - weight

x_i - input

perception

a Neuron where,

$$f(x) = \begin{cases} 1, & \text{if } w^T x + b = 0 \\ 0, & \text{else} \end{cases}$$

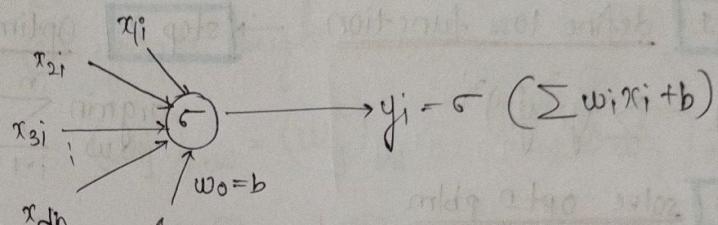
perception = LR (diff: no squashing $f^2(\cdot)$)

How Log-Reg = Neuron

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

In neuron,
let activation $f_n = \sigma$

$$\hat{y}_j = \sigma \left(\sum_{i=1}^n w_i x_{ij} + b \right)$$

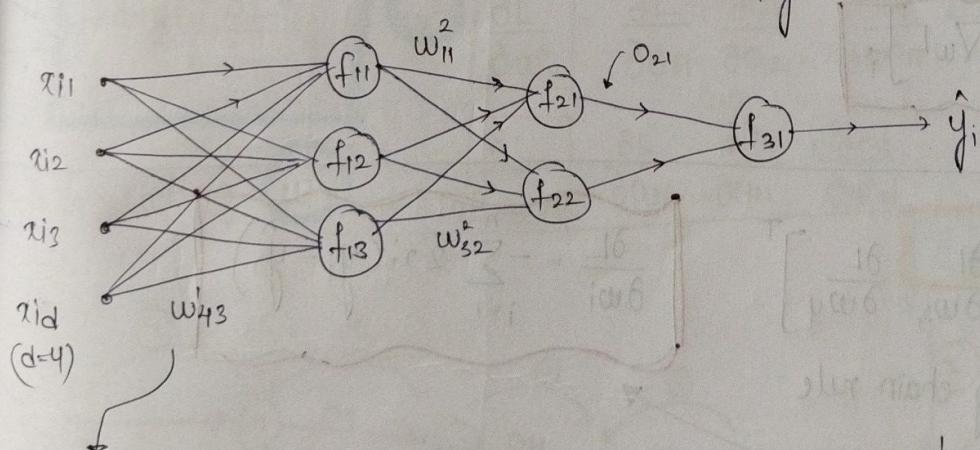


d = dimensions of data

last x_d

Notation

$$D = \{x_i, y_i\}_{i=1}^n, x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$$



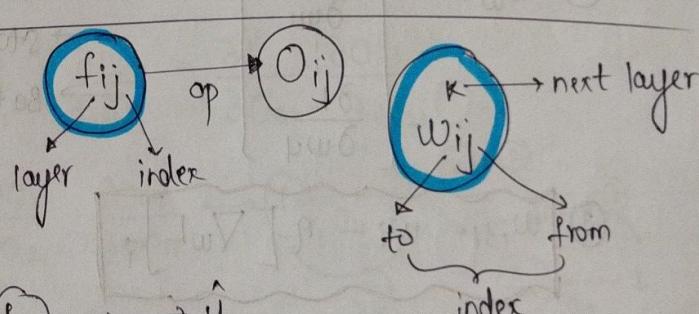
$$w^1 \rightarrow \begin{bmatrix} w_{11}^1 & w_{12}^1 & w_{13}^1 \\ w_{21}^1 & | & | \\ w_{31}^1 & | & | \\ w_{41}^1 & \dots & w_{43}^1 \end{bmatrix}_{4 \times 3}$$

similarly

$$w^1 = \frac{\text{dim}}{4 \times 3}$$

$$w^2 = 3 \times 2$$

$$w^3 = 2 \times 1$$



② Train Single Neuron model

Find best weights.

→ Linear regression as neural model

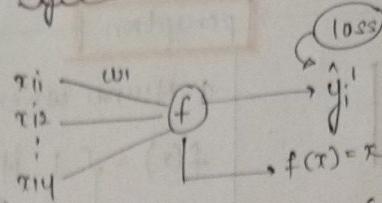
$$\hat{y} = w^T x_i (w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b)$$

loss fn

$$\min_{w_i} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{reg}$$

Classification
Regression

- Log-reg / perceptron
- linear regression



for Linear regression
 $f(x) = \text{identity fn}$

$$\therefore \hat{y}_i = w^T x_i = (w_1 x_1 + w_2 x_2 + \dots + w_n x_n)$$

we ignore reg for simplicity and $y_i \neq \hat{y}_i$

→ step 1 define loss function

$$L = \sum (y_i - \hat{y}_i)^2$$

→ step 2 Optimization

$$\underset{w}{\operatorname{argmin}} \sum_{i=1}^n (y_i - f(w^T x_i))^2$$

Linear-R : identity
Log-R : sigmoid

→ step 3 solve optn pbmlm

① Initialize w_i 's, n

$$\nabla_w L = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_4} \end{bmatrix}$$

Calculating $\nabla_w L$

→ Gradient descent — all pts (x_i, y_i)

→ Stochastic GD — one pt

→ Batch-SGD — some pts $< n$

$$\text{③ } w_{i+1} = w_i - \eta [\nabla_w L]_i$$

Calculate $\nabla_w L$

$$\nabla_w L = \left[\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial w_3}, \frac{\partial L}{\partial w_4} \right]^T$$

$$\frac{\partial L}{\partial w_i} = - \sum_{j=1}^n 2x_j(y_j - \hat{y}_j)$$

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial w_i}$$

using chain rule

$$\frac{\partial L}{\partial f} = \frac{\partial \sum (y_i - f)^2}{\partial f} = -2 \sum_{i=1}^n (y_i - f(w^T x_i))$$

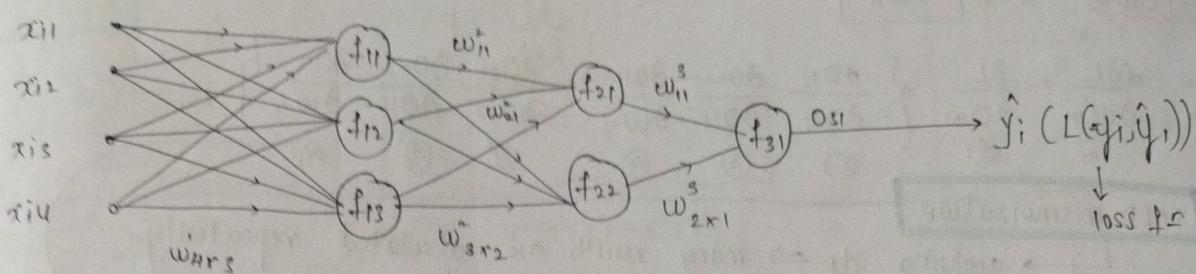
$$\frac{\partial f}{\partial w_i} = \frac{\partial w^T x_i}{\partial w_i} = x_i$$

$$\therefore \left[\frac{\partial L}{\partial f} = -2 \sum (y_i - f(w^T x_i)) \quad \frac{\partial f}{\partial w_i} = x_i \right]$$

Train multi-layer perceptron

→ task? find $w^1, w^2, w^3 = 12 + 6 + 2 = 20$ weights

(3)



→ step 1 Loss fn → step 2

$$\text{Loss fn} = \sum (y_i - \hat{y}_i)^2 + \text{reg}$$

$$\text{let } L_i = (y_i - \hat{y}_i)^2$$

$$L = \sum_{i=1}^n L_i + \text{reg}$$

optimization

① Initialize weights, n

② update weights,

$$(w_{ij}^k)_{\text{new}} = (w_{ij}^k)_{\text{old}} - \eta \left(\frac{\partial L}{\partial w_{ij}^k} \right)_{\text{old}}$$

③ repeat ① and ② until convergence

→ Δw is small

$$\min_{w_{ij}^k} L$$

① find w^3

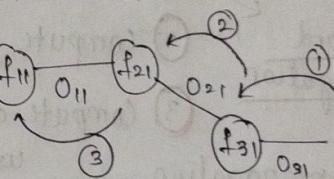
$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{11}^3}; \quad \frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{21}^3}$$

$$\text{finding } \frac{\partial L}{\partial w_{ij}^k}$$

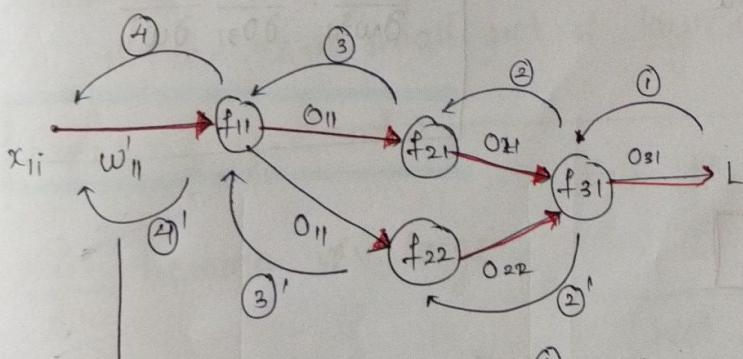
$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial w_{11}^2}$$

$$\frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial w_{21}^2}$$

→ calculating w^1



Special Case of chain rule



$$\frac{\partial h}{\partial x} = \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}$$

$$\frac{\partial L}{\partial x} = \underbrace{\frac{\partial L}{\partial h}}_{1} * \underbrace{\frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}}_{2}$$

$$\frac{\partial h}{\partial x} = \frac{\partial L}{\partial h} \cdot \frac{\partial h}{\partial x}$$

$$\frac{\partial L}{\partial w_{11}^i} = \frac{\partial L}{\partial o_{31}} \left(\frac{\partial o_{31}}{\partial w_{11}^i} \right) \rightarrow \frac{\partial o_{31}}{\partial w_{11}^i} = \left(\frac{\partial o_{31}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}^i} \right) + \left(\frac{\partial o_{31}}{\partial o_{22}} \cdot \frac{\partial o_{22}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}^i} \right)$$

$$\therefore \frac{\partial L}{\partial w_{11}^i} = \frac{\partial L}{\partial o_{31}} \left\{ \begin{array}{l} \text{(1)} \quad \frac{\partial o_{31}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}^i} \\ \text{(2)} \quad \frac{\partial o_{31}}{\partial o_{22}} \cdot \frac{\partial o_{22}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w_{11}^i} \end{array} \right. \quad \begin{array}{l} \text{(3)} \\ \text{(4)} \\ \text{(5)} \\ \text{(6)} \end{array}$$

MLP memoization

→ applying dp → many results are calculated repeatedly

$$\text{Ex: } \frac{\partial o_{11}}{\partial w_{11}^i}$$

You can use chain rule

Back Propagation

Back-Prop = chain rule + memoization

$$D = \{x_i, y_i\}$$

Inp

* Activation f is differentiable

Algorithm

Step 1 : Initialize w_{ij}^k 's

Step 2 : for each (x_i) in D :

 ① pass x_i to forward n/w

 ② Compute $L(y_i, \hat{y}_i)$

 ③ Compute all derivatives

 using chain-rule and
 memoization

Back-propagation

 ④ update weights from end of n/w to start

Step 3 : repeat ② until convergence

i.e. $(w)_{\text{old}} \approx (w)_{\text{new}}$

If !convergence:

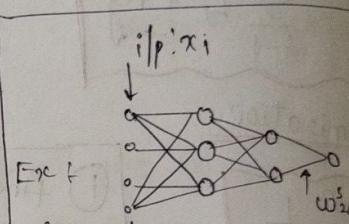
repeat algo for multiple epochs

mini-batch Back-prop

instead of sending (x_i) 1-by-1,

send them in batches of K.

forward-prop



- ① send x_{ii}
- ② calculate $L(y_i, \hat{y}_i)$

- ③ Back-prop
update weight

$$(w_{11}^3)_{\text{new}} = (w_{11}^3)_{\text{old}} - \eta \left(\frac{\partial L}{\partial w_{11}^3} \right)$$

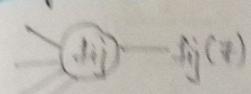
start

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial o_{31}} \frac{\partial o_{31}}{\partial w_{11}^3}$$

Activation Functions

(5)

Sigmoid

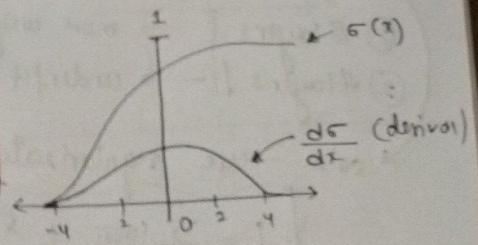


$$\text{ip} : w_i x_i + \sum w_i x_i = z$$

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

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

$$\begin{aligned}\frac{\partial \sigma}{\partial z} &= \sigma(z)(1 - \sigma(z)) \\ \sigma'(z) &= \sigma(z)(1 - \sigma(z))\end{aligned}$$



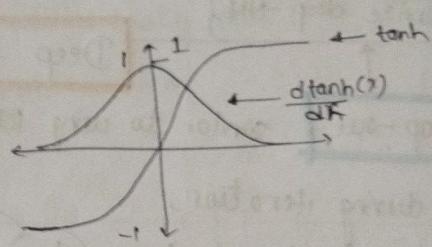
tanh

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

$$\frac{da}{dz} = 1 - a^2$$

sigmoid → old
tanh

ReLU
Leaky ReLU → new



Vanishing Gradients

→ slows down convergence

$$\rightarrow w_{\text{new}} = w_{\text{old}} - n \left(\frac{\partial L}{\partial w} \right)$$

$$\frac{\partial o_{31}}{\partial o_{21}} = \frac{\partial \sigma(o_{21})}{\partial o_{21}}$$

$$\text{w.r.t. } \frac{\partial \sigma}{\partial x} = 0 < x < 1$$

as value of $\left(\frac{\partial L}{\partial w} \right)$ is very small

each values are 0 ↔ 1

→ when u multiply,

the answer is very small ≪ 1

$w_{\text{old}} - n \cdot \frac{\partial L}{\partial w}$ → negligible

$\therefore w_{\text{new}} \approx w_{\text{old}}$ → very small value → so this cause

$w_{\text{new}} = w_{\text{old}}$ and

it is assumed to converge

\therefore due to vanishing grad,

u can train small amt of layers $\leq 2, 3$, so u can't build deep NN.

Exploding gradients

$\frac{\partial L}{\partial w}$ become very large

$\therefore \Delta w$ varies massively at end of each iterations

NEVER CONVERGES

as $\Delta w \not\rightarrow 0$

$$\text{Ex: } \frac{\partial L}{\partial w} = \left(\frac{\partial o_{31}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial o_{11}} \dots \right)$$

$> 1 \dots = \text{large number}$

(Not with sigmoid activatⁿ fⁿ)

⑥ Bias-variance tradeoff - (Multi-layer perceptron)

- ① # layers ↑ → more weights → overfit / (high variance) $L_1 \text{ reg}$
- ② # layers ↓ → underfit (high-bias)

soln - use regularization

$$L = \sum_{i=1}^n \text{loss}_i + \lambda \cdot L_2/L_1 \text{ regularization on weight}$$

$L_2 \text{ reg}$

$L_1 \rightarrow$ sparsity,
some $w_i = 0$

regularize deep-NN

Deep Learning

Drop-out

similar to using RF to do our regularization

* during iteration,

at each layer drop some neurons

remove all connections

(drop-out rate = p) ($0 \leq p \leq 1$)

$p = 0.2$, 20% of neurons
in layer are dropped

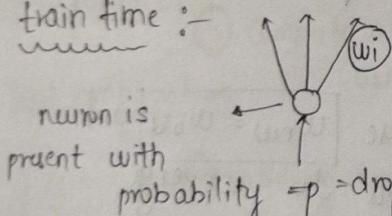
→ in each iteration,

drop different neurons

* similar to selecting some features in random-trees

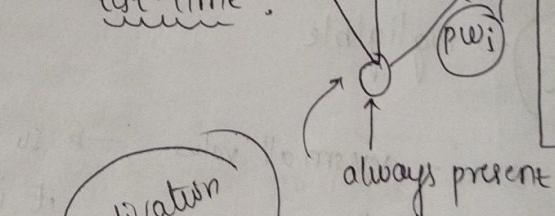
∴ dropout ≈ random subset of features / inputs

train time :-



neuron is present with probability p = dropout rate

test-time :-



why pw_i ?

because during train
0 was present
only $p(p)$ times

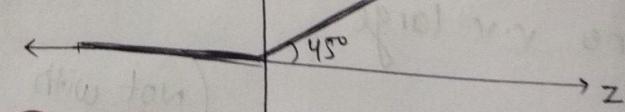
Rectified Linear Unit

$$f(z) = z^+ = \max(0, z)$$

$$f(z) = \begin{cases} z, & \text{if } z > 0 \\ 0, & \text{otherwise} \end{cases}$$

ReLU

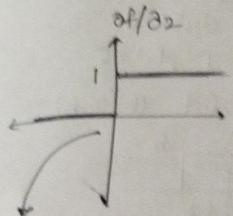
$$f(z)$$



$$\frac{df(z)}{dz} = \begin{cases} 0, & \text{if } z \leq 0 \\ 1, & \text{otherwise} \end{cases}$$

$$\tan(45^\circ) = 1$$

$$\frac{\partial f(z)}{\partial z} = \begin{cases} 0, & z \leq 0 \\ 1, & \text{otherwise} \end{cases}$$



discontinuous at $z=0$
∴ not differentiable at 0

$$\therefore \frac{\partial f(z)}{\partial z} \in \{0, 1\}$$

i) never > 1 ,
∴ !exploding gradient

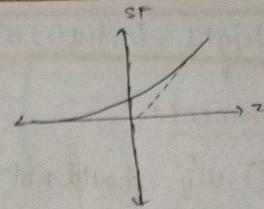
ii) can't be decimal like 0.001
∴ !vanishing gradient

X iii) if $\frac{\partial f(z)}{\partial z} = 0$, then all other terms don't matter

softplus fn

$$f(z) = \log(1 + e^z)$$

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



⑦

ReLU advantages

- 1) speeds up convergence
- 2) easy to compute derivatives - just if

disadvantages

- 1) dead activations $\xrightarrow{\text{soln}} \text{Leaky relu (init strat)}$
↳ when weight $\rightarrow -\infty$

$$w_{\text{new}} = w_{\text{old}} - n(0 \dots)$$

∴ weights won't change

(Dead activation)

Variants of ReLU

1) Noisy ReLU

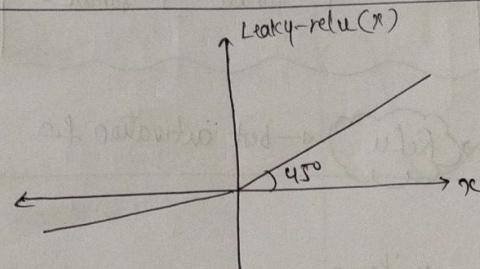
→ add gaussian noise

$$f(z) = \max(0, z + \gamma)$$

where,

$$\gamma \sim N(0, \sigma^2)$$

→ γ is picked from a gaussian distribution with small noise \equiv variance



$$\frac{d f(z)}{d z} = \begin{cases} 1, & z > 0 \\ a, & z < 0 \end{cases}$$

this may lead to vanishing gradient

because $0 < a < 1$ and if a is multiple times it may lead to vanishing

2) Leaky ReLU

→ solve dead activation problem

$$f(x) = \begin{cases} x, & x > 0 \\ a \cdot x, & x \leq 0 \end{cases}$$

$$[a = 0.01] \rightarrow \text{hyperparam / fixed}$$

⑧ Weight Initialization

Requirements

- ① w_{ij}^k should not be zero
- ② w_{ij} should be small (\neq too small)
- ③ good variance $\text{var}(w_{ij}^k)$

more different
better learning

① Gaussian/Normal Init

$$W_{ij}^k \sim N(0, \sigma^2)$$

↑ pick ↓ small var

Xavier/Glorot → sigmoid outⁿ

→ gaussian

$$W_{ij}^k \sim N(0, \sigma_{ij}^2)$$

$$\sigma_{ij}^2 = \frac{2}{\text{fanin} + \text{fanout}}$$

He-init?

→ Normal

$$W_{ij}^k \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{\text{fanin}}}$$

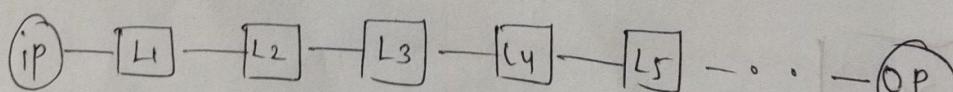
→ uniform

$$W_{ij}^k \sim U\left(-\sqrt{\frac{6}{\text{fanin}}}, +\sqrt{\frac{6}{\text{fanin}}}\right)$$

Batch-Normalization

→ consider inputs are normalized ie $\text{mean}=0, \text{sd}=1$

→ we send x_i as mini-batchy



b_1
 b_2
 b_3
 b_4

① new E_{b1}
② new E_{b2}

both are same as normalized

①
②

① and ② vary due to
OP performed by
layers L_2, L_3

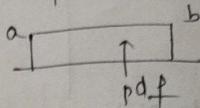
- picking weights for step ① of optimization

* $N(0, \sigma^2)$

normal distⁿ

mean = 0, std = σ

* uniform dist $U(a, b)$



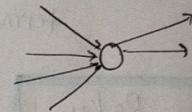
② Uniform init

$$W_{ij}^k \sim U\left(\frac{-1}{\sqrt{\text{fan-in}}}, \frac{1}{\sqrt{\text{fan-in}}}\right)$$

good for σ activation fn

fan-in : no of inputs

fan-out : no of outputs



fan-in = 4
fan-out = 2

ReLU ← best activation fn

He-init?

→ Normal

$$W_{ij}^k \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{\text{fanin}}}$$

→ uniform

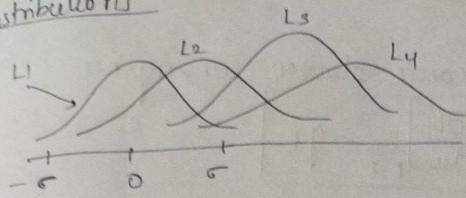
$$W_{ij}^k \sim U\left(-\sqrt{\frac{6}{\text{fanin}}}, +\sqrt{\frac{6}{\text{fanin}}}\right)$$

fan-in

fan-out

Prob

distributions



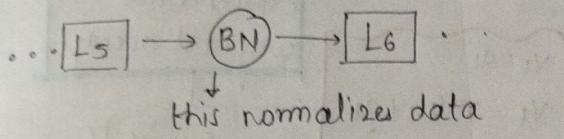
Internal Covariance
shifting

(9)

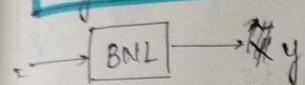
Can be avoided using
batch normalization

∴ Deeper nets can be
trained

[soln] → add a batch layer in b/w deep layer



Algorithm for Batch-Norm



advantages

(1) faster convergence

(2) works as weak regularization

use $(BN + drop)$, always

input : values of x over minibatch

$$B = \{x_1, \dots, x_m\}$$

output :- params to learn, γ, β → hyper params

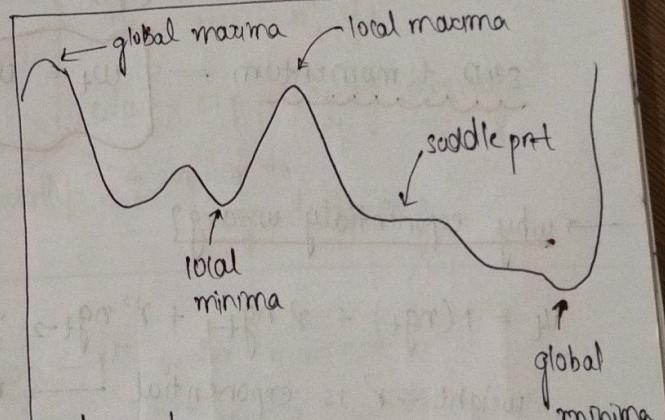
$$\{y_i = BN_{\gamma, \beta}(x_i)\}$$

$$\text{# } \mu_B \rightarrow \frac{1}{m} \sum_{i=1}^m x_i \quad (\text{Batch-mean})$$

$$\text{# } \sigma_B^2 \rightarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \quad (\text{min-batch variance})$$

$\hat{x}_i \rightarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$ normalize
in case $\sigma_B^2 = 0$,
 $\epsilon \approx 0$ very small value

o/p $\rightarrow y_i = \gamma \hat{x}_i + \beta$ scale and shift



slope at $= 0$

minima

slope (saddle) $= 0$, $\leftrightarrow 0 = 0$
 $\tan 0 = 0$

problems

stuck at saddle point

⑩ SGD with momentum

$\text{SGD} \approx \text{GD}$, \therefore the gradients are noisy
 update : $w_t = w_{t-1} - n \left[\frac{\partial L}{\partial w} \right]_{t-1}$

which cause slower convergence

$$\text{let } \left[\frac{\partial L}{\partial w} \right]_{t-1} = g_t$$

* Denoise data

$$\begin{array}{cccc} t_1 & t_2 & t_3 & t_4 \\ \rightarrow a_1 & a_2 & a_3 & a_4 \end{array} \quad \begin{array}{c} \gamma = 0 \leftrightarrow 1 \\ \therefore t_1 = a_1 \\ t_2 = \gamma t_1 + a_2 \end{array}$$

$$(t_1) \rightarrow v_1 = a_1$$

$$(t_2) \rightarrow v_2 = a_2 + \gamma v_1$$

$$(t_3) \rightarrow v_3 = a_3 + \gamma v_2$$

$$= a_3 + \gamma (a_2 + \gamma v_1) = a_3 + \gamma a_2 + a_1 \gamma^2$$

↓ less weight to older points
more weight to current point

denoise SGD

momentum

$$\text{initially, } v_1 = (g_1 \cdot n) \quad \text{learning rate}$$

$$v_t = \gamma v_{t-1} + (n g_t)$$

$$\Rightarrow w_t = w_{t-1} - v_t$$



$$0 \leq \gamma \leq 1$$



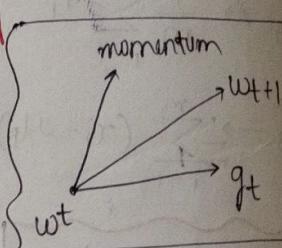
$$\text{if } \gamma = 0,$$

$$\text{SGD + momentum} = \text{SGD}$$

SGD + momentum =

$$w_t = w_{t-1} - \left[\gamma v_t + n g_t \right]$$

↓ exp-weight ↓ gradient
 org (momentum)



→ why exponentially w-avg?

$$v_t = 1(n g_t) + \gamma n g_{t-1} + \gamma^2 n g_{t-2} \dots$$

weight $\equiv \gamma$ is exponential → more weight to currently discovered points and less weight to past points

change LR(n) dynamically

ADAGRAD

→ In SGD the learning rate (n) is constant for all weights

[idea] have different (n) for different weight

why?

when data has sparse features it's helpful

Adagrad = adaptive gradient

$$w_t = w_{t-1} - n_t^{-1} \cdot g_t$$

(n_t) is different for each weight and each iteration

$$\eta_t = \frac{n}{\sqrt{\alpha_{t-1} + \epsilon}}$$

$\alpha_t \geq 0$ &
 $\alpha_t > \alpha_{t-1}$

$$\alpha_{t-1} = \sum_{i=1}^{t-1} g_i^2$$

↓ always true

sum of all gradients² upto $t-1$

$$\left[\frac{\partial L}{\partial w} \right]_{w_{t-1}}$$

(12)

[Note] t - iteration number
 $\alpha_t \uparrow, \alpha_{t-1} \uparrow \rightarrow \eta_t \downarrow$ ∵ as iterations increase, learning rate decreases adaptively

[uses] η_t tunes automatically

works if D has sparse and dense

[problem] as $t \uparrow, \alpha_{t-1}$ may become v.v.large
if $\alpha_{t-1} \rightarrow \infty, \eta_t \rightarrow 0$ which causes slower convergence (very less Δ)

Adadelta & RMS prop → soln to α_t becoming v.v.large

Adadelta

$$\eta_t = \frac{n}{\sqrt{eda_{t-1} + \epsilon}}$$

$$eda_{t-1} = \gamma \cdot eda_{t-2} + (1-\gamma) g_{t-1}^2$$

$\gamma = 0.95$ usually

recursive

$$w_t = w_t - \eta_t g_t$$

eda → exp decaying average

similar to exp-decay-avg

use:

→ control growth of eda_t

In stats

mean - 1st order moment

variance - 2nd order moment

Adam (adaptive moment estimate)

idea: along with using (g_t^2) for eda_t also use (g_t) → mean

$$eda_{gt} \quad m_t = \beta_1 \cdot m_{t-1} + (1-\beta_1) g_t$$

$$V_t = \beta_2 \cdot V_{t-1} + (1-\beta_2) g_t^2$$

usually,
 $\beta_1 = 0.95 \quad 0 \leq \beta_1, \beta_2 \leq 1$
 $\beta_2 = 0.99$

$$\hat{m}_t = \frac{m_t}{1-(\beta_1)^t} \quad \hat{V}_t = \frac{V_t}{1-(\beta_2)^t} \rightarrow \beta_2 \approx t$$

$$w_t = w_{t-1} - \alpha \cdot \frac{\hat{m}_t}{\sqrt{\hat{V}_t + \epsilon}}$$

if $\beta_1 = 0$,
adadelta (similar)

(12) Which optimizer to use?

mini-Batch SGD
gets stuck in saddle point
& the good for smaller NN

momentum & NAG
works well in most cases
but slower

Adagrad
good for sparse data

adadelta
rmsprop

Adam
better overall

gradient monitor and clipping

monitor gradients and update for each epoch, layer and each weight

why? vanishing gradient, exploding gradient

monitoring

clipping

bad clipping

$$W = \begin{bmatrix} | & | & | & | & | \end{bmatrix}$$

$\downarrow w_{11} \quad \downarrow w_{12}$

$$g = \begin{bmatrix} | & | & | & | & | \end{bmatrix}$$

$\frac{\partial L}{\partial w_{11}}$

$\rightarrow W$ & g are vector of all gradients and weights.

\rightarrow if weights don't change \rightarrow vanishing gradient

clipping

$$\vec{g}_{\text{new}} = \frac{\vec{g}_{\text{old}}}{\|\vec{g}\|_2}$$

$$\|\vec{g}\|_2 = \sqrt{g_1^2 + g_2^2 + \dots + g_n^2}$$

\therefore in g_{new} each term ≤ 1

\therefore it won't cause Δ exploding

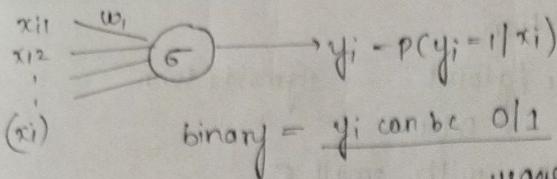
Softmax classifier

\rightarrow Logistic regn. + multi class = Soft-max

expanded upon on Logistic regression

* By default LR is for Binary class,
use 1 vs rest for multi class

Logistic regression or MN



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

$$z = w^T x$$

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

softmax

$x_i \rightarrow \text{model} \rightarrow p(y=1|x_i)$

$p(y=2|x_i)$

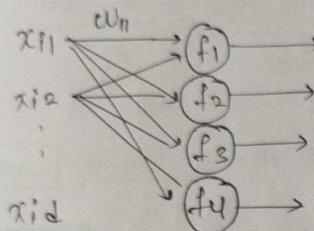
$p(y=n|x_i)$

sum = 1

K-class

$x_i \rightarrow [] \rightarrow \dots + [] \rightarrow \sum_j p_j$

softmax layer



→ let f_i be (z_i) , $z_i = w^T x$

$$z_1 = \sum_{j=1}^d w_{j1} x_{ij}$$

$$z_2 = \sum_{j=1}^d w_{j2} x_{ij}$$

sum of $f(z_i)$

$$\frac{e^{z_1} + e^{z_2} + e^{z_3} + \dots}{\sum e^{z_i}} = \frac{\sum e^{z_i}}{\sum e^{z_i}} = 1$$

$f_i(x)$

$$f_i(z_i) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

$$f_i(z_i) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

where, $K = \text{no of } f_i \text{ s / classes}$

$$\therefore (\text{LR}) \rightarrow f(z) = \frac{e^z}{e^z + 1}; \quad \text{softmax} \rightarrow f(z_i) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

minimizing 2-class

minimize

multi-class log-loss

Crossentropy

multi-class log loss

N points, K classes

$$\text{mc-log loss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K y_{ij} \log(p_{ij})$$

2-class log loss

$$\frac{1}{N} \sum_{i=1}^N [y_{ij} \log p_i + (1-y_{ij}) \log(1-p_i)]$$

→ (y_{ij}) : binary indicator → $\begin{cases} 1, & \text{if } y_i \in \text{class } j \\ 0, & \text{otherwise} \end{cases}$

$$\rightarrow (P_{ij}) = p(y_i = j | x_i)$$

simple Interpret

$$-\frac{1}{N} \sum_{i=1}^N \left(\underset{\in \text{class } j}{\underset{\text{does } x_i}{\left(\right)}} \times \text{prob}(x_i \in j) \right)$$

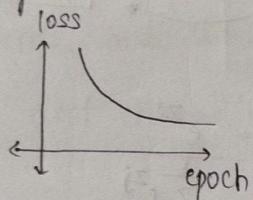
(14)

Steps to train MLP

- ① Pre-process : data normalization
- ② weight initialization : Xavier / glorot \rightarrow sigmoid / tanh
He \rightarrow ReLU
gaussian with small σ
- ③ activation f_D : ReLU
- ④ Batch Normalization \rightarrow deep layer (other layers)
Dropout \rightarrow dropout rate (p)
- ⑤ optimizers : adam (2018) \rightarrow fast to converge
- ⑥ hyper-params
 - \rightarrow architecture \rightarrow # layers
neurons
 - \rightarrow dropout rate (p)
 - \rightarrow Adam (B_1, B_2, α)
- ⑦ Loss function :
 - 2-class \rightarrow log-loss
 - K-class \rightarrow multiclass log-loss
 - regression \rightarrow squared error

- ⑧ monitor gradients : clipping

- ⑨ plots :



- ⑩ ! avoid overfitting

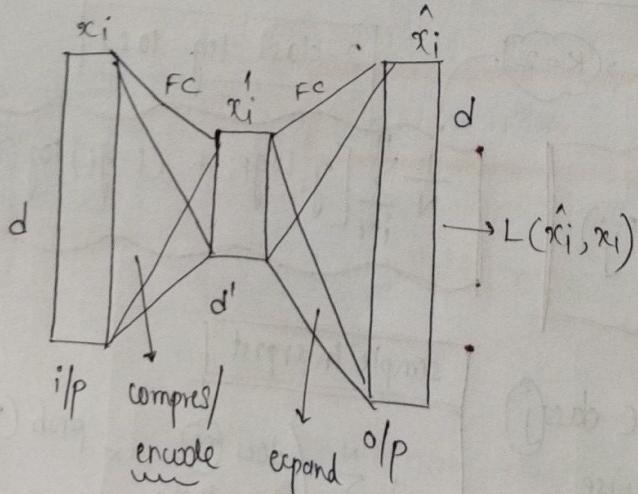
\hookrightarrow diff b/w train / test loss

$\left. \begin{array}{l} AE \approx PCA \\ AE \text{ with 1 layer} \end{array} \right\}$ AE with 1 layer + fc

dimensionality reduction

$$x_i \rightarrow R^d \rightarrow x'_i \rightarrow R^{d'}$$

Auto Encoders



lets consider, $d=6$

\rightarrow hidden layer $d'=3$

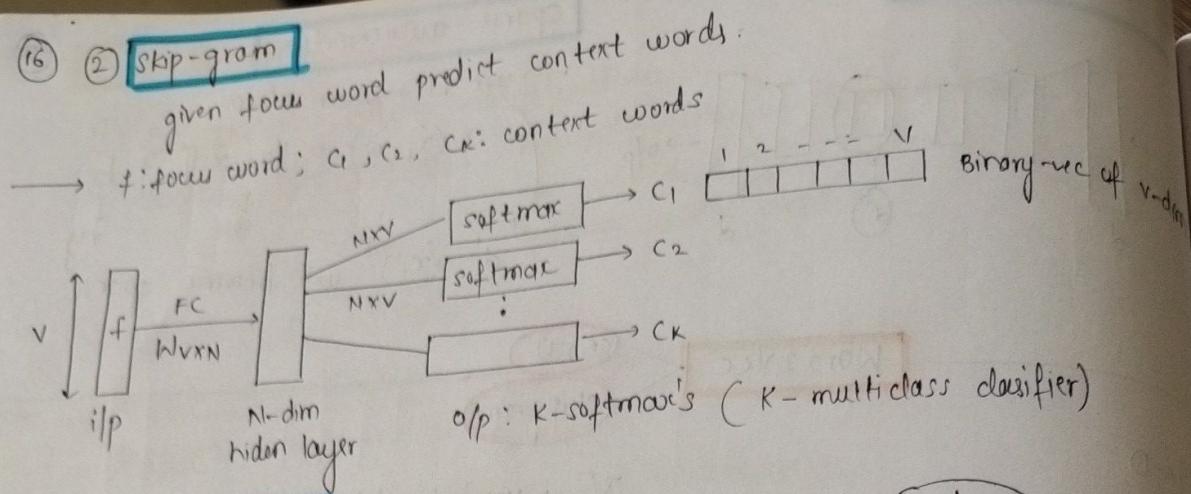
\therefore 3 o/p's from hidden layer = x'_i

$$x_i \rightarrow R^6 \quad \hat{x}_i = R^6$$

$$\boxed{x'_i \rightarrow R^3}$$

dimensionality reduced

⑯ ② **skip-gram**
given focus word predict context words.



skipgram vs CBOW

$$\# \text{weights in CBOW \& skipgram} = (K+1)(V \times N) \quad \text{vv large}$$

CBOW → 1 softmax to train \therefore skipgram is computationally expensive
skipgram → K softmax

CBOW

- faster
- better for frequently occurring words

SKIP

- ⊕ smaller amt of data
- ⊕ infrequently occurring words, works well
- ⊖ slower

optimizations for w2v

due to large weights

Hierarchical softmax (algo-based)

Negative sampling (stats-based)

Neg-sampling

idea: update only sample of words (weights)

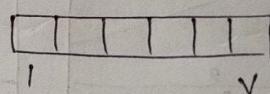
- ① always keep target word (which u need to predict)
- ② sample from non-target words

$$P(w_i) = 1 - \sqrt{\frac{c}{\text{freq}(w_i)}} \rightarrow f^{-r}$$

$\rightarrow P(\text{you pick word } w_i)$

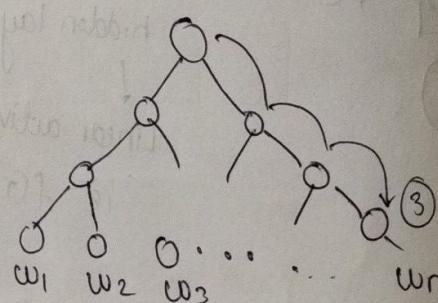
Hierarchical softmax

→ optimize softmax o/p.



→ you need to check $i=1 \rightarrow y$

↓ we binary tree to search



incase (wn) is answer (most probable)

u only need ③ iterations

Convolutional Neural Networks

(17)

Sobel Kernel → Edge detection

$$G_x = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 0 & 0 \\ -1 & -2 & -3 \end{bmatrix}$$

(detects edges)

$$G_{xy} = \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 3 & 0 & -3 \end{bmatrix}$$

$$\text{Input } 6 \times 6 \xrightarrow{\text{convolution}} \text{Output } 3 \times 3$$

$$\text{Image } 4 \times 4$$

with horizontal
edges detected

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

output

Padding

$$6 \times 6 \xrightarrow{n} 4 \times 4 \xrightarrow{k=3}$$

we want out shape = input shape

$$n=6, k=3$$

$$\rightarrow o/p = n-k+1 = 6-3+1 = 4$$

we can get o/p = 6 if $n=8$,

$$8-3+1 = 6$$

$$\begin{matrix} x & x \\ x & x \end{matrix}_{2 \times 2} \xrightarrow{p=1} \begin{matrix} & & & \\ & & & \\ & & & \\ & & & \end{matrix}_{4 \times 4}$$

$$\rightarrow \text{add 2 col (R L)} + 2 \text{ rows (T B)}$$

$$\begin{matrix} & & & \\ & & & \end{matrix}_{(n \times n)} \xrightarrow{\text{padding with some value}} \begin{matrix} & & & \\ & & & \end{matrix}_{(n+2) \times (n+2)}$$

$$\begin{matrix} & & & \\ & & & \end{matrix}_{n \times n} \xrightarrow[\text{Pad}=p]{(K \times K)} \begin{matrix} & & & \\ & & & \end{matrix}_{(n+2p-k+1) \times (n+2p-k+1)}$$

padding + convolution

stride

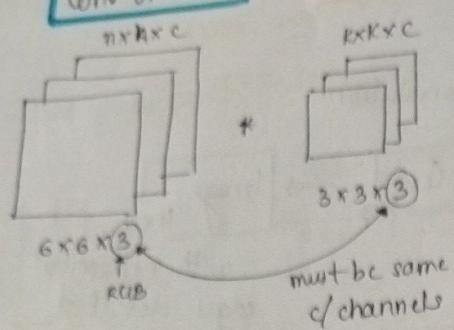
→ By how many values u shifted
the conv-kernel (by default = 1)

$$n \times n \xrightarrow[\text{stride}=s]{K \times K} \left(\frac{n-K+1}{s} \right) \times \left(\frac{n-K}{s} \right) + 1$$

Padding + Stride + Conv ?

$$\begin{matrix} & & & \\ & & & \end{matrix}_{(n \times n)} \xrightarrow[\text{input}]{\text{padding } p} \left(\frac{n-k+2p}{s} + 1 \right) \times \left(\frac{n-k+2p}{s} + 1 \right)$$

(18) Conv on RGB

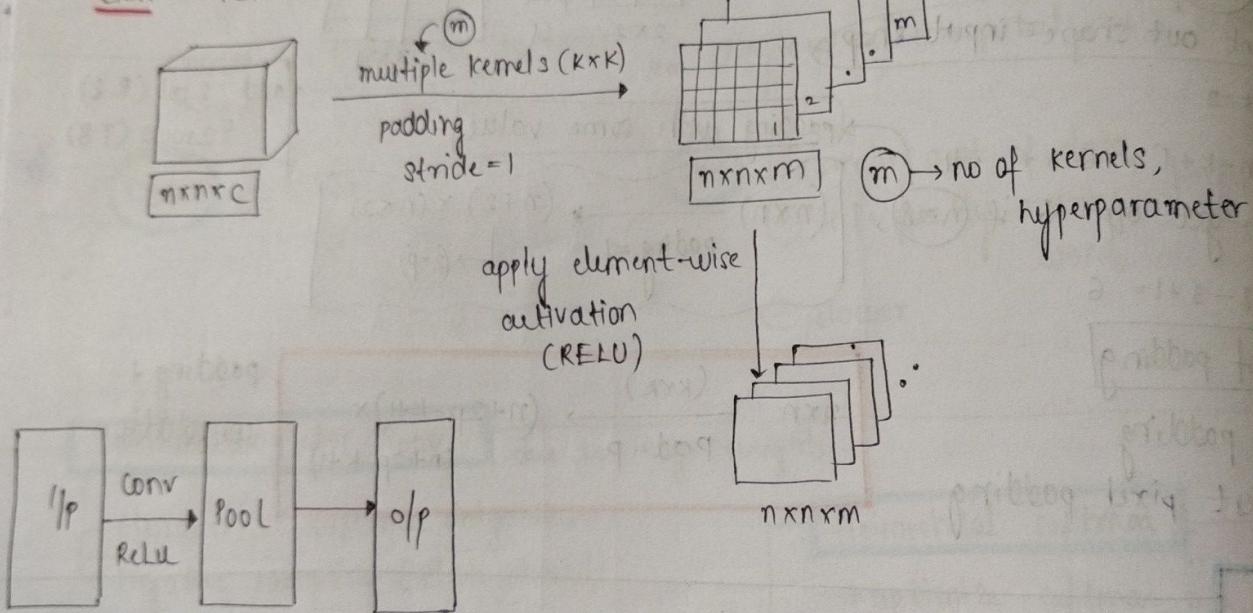


$(n-k+1) \times (n-k+1) \times 1 \rightarrow o/p$ is always a matrix/2D o/p

Convolution layer in NN

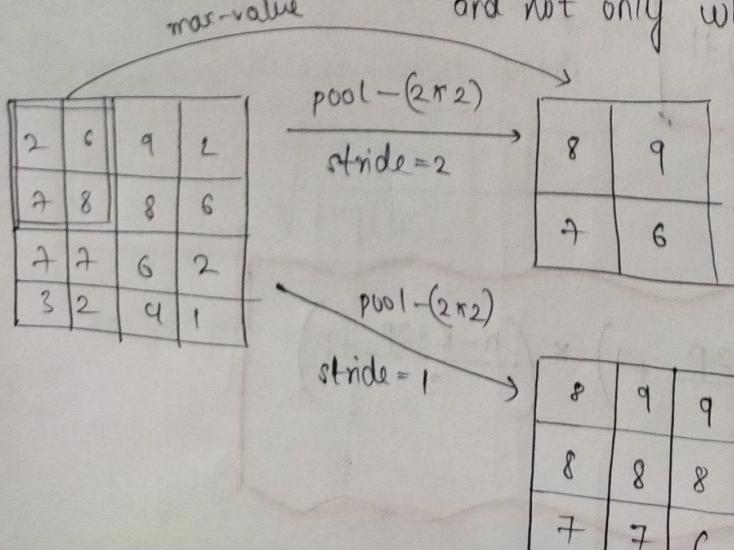
MLP \rightarrow learn weights

CNN \rightarrow learn Kernel matrix



max-pooling \rightarrow same as conv \rightarrow but pick max value

aim: location invariance: - the feature could be present anywhere in image and not only where it was seen



we can say if val \uparrow , there is atleast one occurrence of feature in image

CNN Train

+ in backprop

condition

- { 1) conv-
2) more-

must be

derivative

pool \rightarrow

$$f'(x) = \frac{(df(x))}{(dx)}$$

LeNet

CNN Training

* in backprop, w.r.t loss fn, must be differentiable and all operations

- condition
 - { 1) conv layer ✓
 - { 2) max-pool layer
 - must be differentiable
- conv + [ReLU] → differentiable
- element wise mult + addition
→ similar to dot product
 $w^T x \rightarrow$ is differentiable

derivative of max-pooling

$$\text{pool} \rightarrow f(x) = \text{max}(x)$$

elements which occur in max pool kernel]

0	1	2
6	2	9
1	3	4

mp = 2x2

6	9
6	7

6 depends only on (2,1)

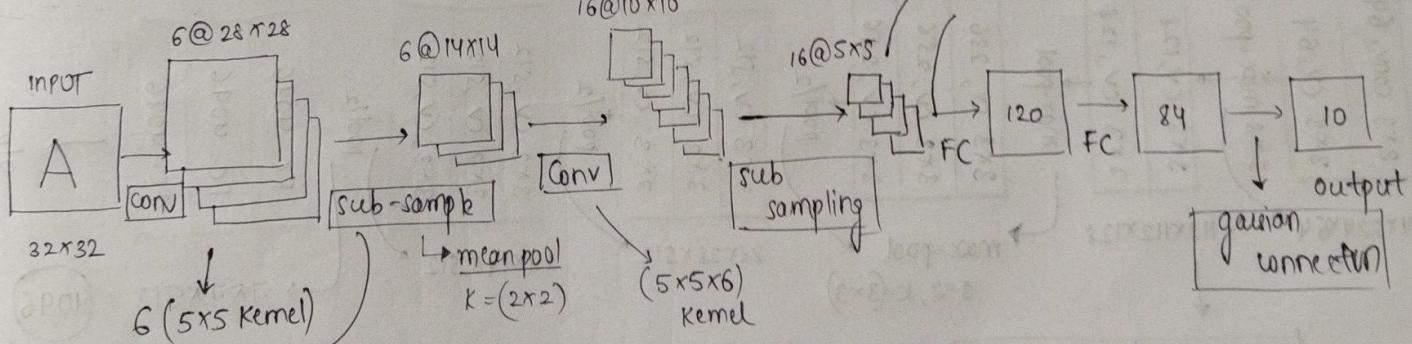
$$\therefore \frac{\partial \text{mp}}{\partial x_i} = 1 \text{ for } (x_{21})$$

for others

$$\frac{\partial \text{mp}}{\partial x_i} = 0 \text{ b/c mp doesn't depend on } x_i$$

LeNet

→ sigmoid activation



Sub-Sampling = Mean Pooling

$$\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array} \rightarrow 10/4$$

Data Augmentation

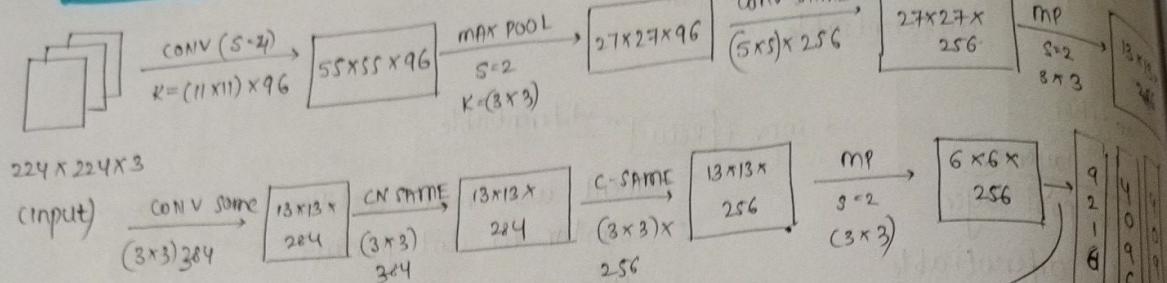
→ modify some data to make different variations
→ rotate, flip, skew, noise, reshape ..etc, zoom, shape

[use]

→ Reduces various invariances.

→ convert small dataset to large dataset

(20)

AlexNet (2012)

new ideas:
ReLU
Dropout
GPUs

→ Local Response Normalization (LRN)

VGG Net

→ Simplified improved version of alexnet

convolutions

Conv → (3x3)

S=1

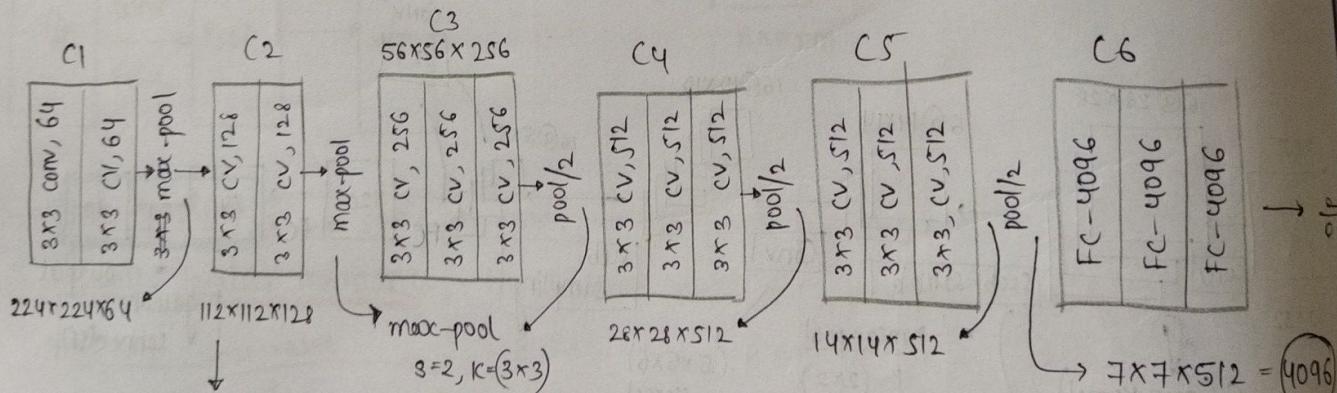
p=same

and max-pool → (2x2)
(S=2)

VGG-16

VGG-19

Convnet architecture of **VGG-16**



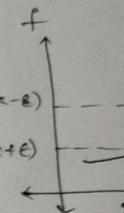
output dim
after (c2) completes

i/p (224x224x3)

→ o/p (4096) → 1000-class

ResNets**Gradient**

* Numerical
of de



grad-ch

Ex: W =
dp/dW

(du)

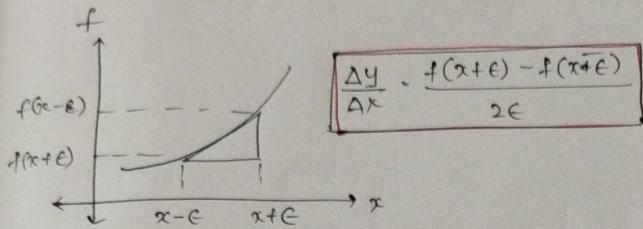
Grad

ResNets | Residual Networks

Gradient Checking

* Numerical approximation of derivative of f : $\lim_{\epsilon \rightarrow 0} \frac{f(x+\epsilon) - f(x-\epsilon)}{2\epsilon}$

* Actual der \rightarrow calculated by formula or backprop.



because ϵ is small,
difference also is small

$$\text{grad-check} = \frac{\|d\mathbf{w} - d\mathbf{w}_{\text{approx}}\|_2}{\|d\mathbf{w}\|_2 + \|d\mathbf{w}_{\text{approx}}\|_2}$$

Euclidean distance

$\begin{cases} \text{very small } < 10^{-7} \rightarrow \text{correct} \\ > 10^{-3} \rightarrow \text{false/not correct} \end{cases}$

Ex: $\mathbf{w} = w_1, w_2, \dots, x = x_1, x_2 \quad f(w_1, w_2, x_1, x_2) = w_1^2 x_1 + w_2^2 x_2$

$$\frac{df}{dw_1} = 2 \cdot w_1 x_1 \quad \frac{df}{dw_2} = 2 \cdot w_2 x_2; \quad \left(\frac{df}{dw_1} \right)_{\text{approx}} = \frac{f(w_1 + \epsilon_1, w_2, \dots) - f(w_1 - \epsilon_1, \dots)}{2\epsilon}$$

$$(d\mathbf{w}) = 2w_1x_1 = 6 \rightarrow \textcircled{1}$$

$$(d\mathbf{w}_1)_{\text{approx}} = 5.9999999 \rightarrow \textcircled{2}$$

$$\text{grad-check} = \left(\frac{6 - 5.9999999}{6 + 5.9999999} \right) = 4.251268e^{-13}$$

\therefore as $< e^{-13} \rightarrow$ backprop
is correct