

Deep MLP (1980's to 2018)

①

1980's, 1990's — 2006 - 2012

→ 2,3 layered network

①

②

③

→ vanishing gradient

→ too little data (easy to overfit)

→ too little computation power. (hard to run for too many epoch because require lot of time).

problems with growth of Deep MLP.

2010 → lots of data → Internet

↳ labelled data (Image net)

→ Computers (GPU, CPU)

↳ NVidia (super useful for solving deep learning algorithm)
→ v.v. good for deep learning.

→ new Idea's & Algorithm

Modern DL

classical ML

SVM (90's) old new.
RF (2000's) { Theory, Experimental }

Random forest)

* Building new theory is very hard.

So they do lot of experiments and if it is successful then build the theory.

great result ↙

lot of new idea's in DL → University of toronto theory
↳ Stanford
↳ Companies. ↳ experiment.

② Dropout & regularization

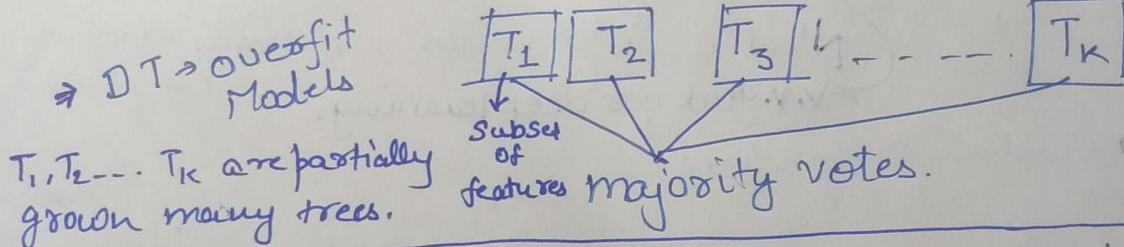
Deep NN \rightarrow overfit easily
many layers \rightarrow many weights

* So we must avoid overfitting we use regularizer like L1, L2 or DROPOUT.

\hookrightarrow extremely simple & elegant

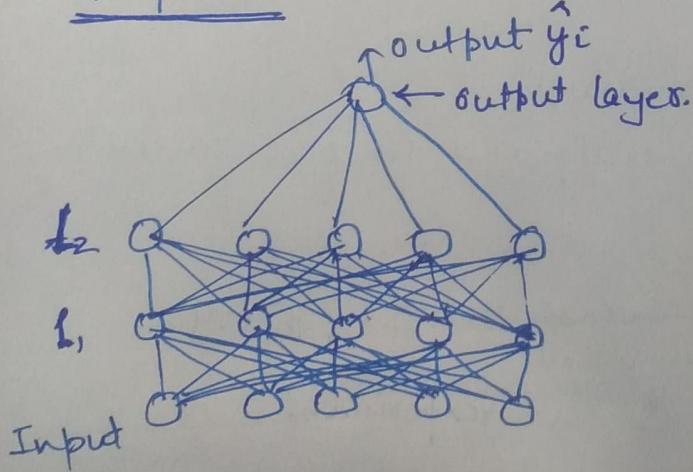
2012 \vdash Nitesh S. Hinton
 \hookrightarrow Master's thesis

Random forest \rightarrow sample a subset of columns/features
 \hookrightarrow many trees (DT) \rightarrow fully grown

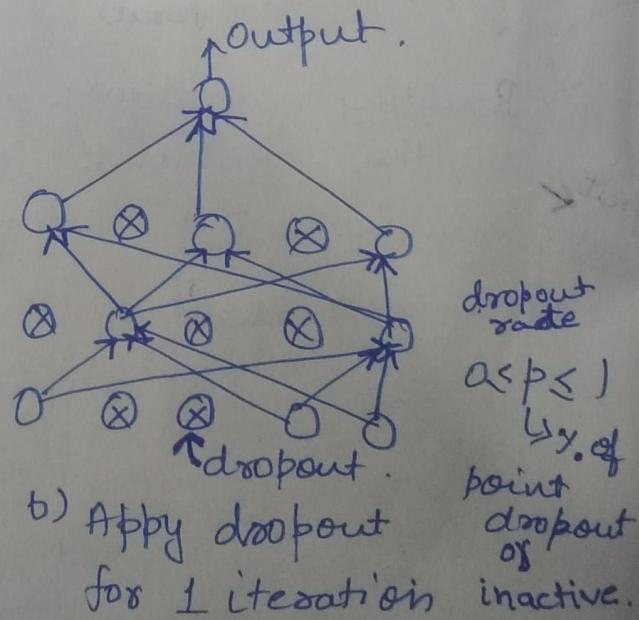


① Random subset of feature and randomization of regularization through randomization.

Dropouts \rightarrow



a) standard NN

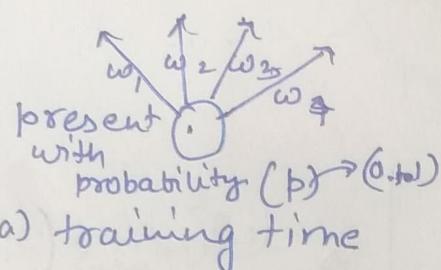


b) Apply dropout for 1 iteration

⇒ dropout is random subset of features in random forest but in NN we have dropout {random subset of neurons active}.

Q training is done but now test will done.

probability of dropout.



» p is multiplied because this neuron was active only for p numbers of times.

→ p is multiplied when completing the model for final testing.

b) test time

p is hyperparameter generally keep low if overfitting.

③ Rectified Linear Unit (ReLU) ~~Activation function.~~^{new}

→ Classical NN → vanishing gradient (Sigmoid or tanh)

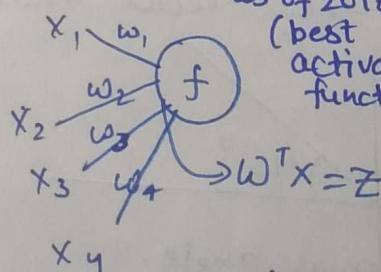
→ convergence slow

$$(\omega_{ij}^k)_{\text{new}} = (\omega_{ij}^k)_{\text{old}} - \eta \left[\frac{\partial L}{\partial \omega_{ij}^k} \right] \rightarrow \text{very small.}$$

Very slow changes. When training large network cause lot of time.

To solve this comes ReLU

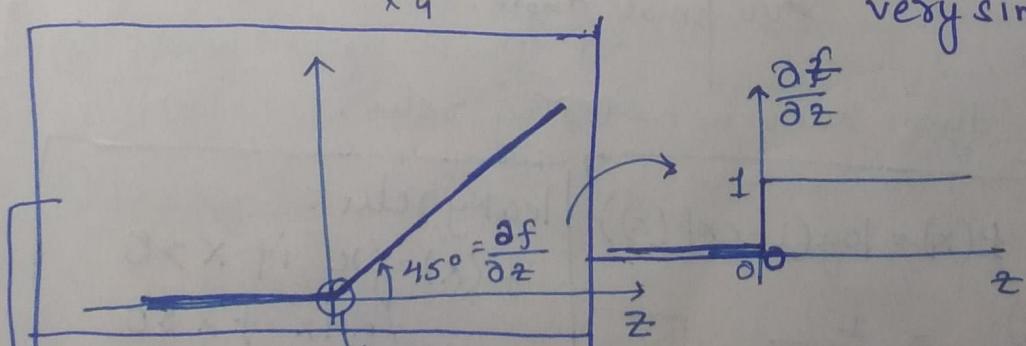
as of 2018
(best activation function)



$$f(z) \text{ is said } z^+ \text{ which is } \max(0, z)$$

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

very simple. still powerful.



this is graph of
ReLU activation

function i.e. $f(z) = \max(0, z)$

discontinuous function.
i.e. not differentiable at 0. means $\frac{\partial f}{\partial z}$ at 0 is not defined.
which will be optimized in later chapter.

Computing derivative relu.

easy $\left\{ \begin{array}{l} \frac{df}{dz} = 0 \text{ if } z < 0 \\ \frac{df}{dz} = 1 \text{ if } z > 0 \end{array} \right.$

- easy to diff.
- NO problem of vanishing gradient because now derivative will be 0 or 1.
- NO problem of exploding gradient

hard
$$\frac{\partial f_{\text{sigma}}}{\partial z} = \frac{1}{1 + \exp(-z)}$$
, $\frac{\partial f_{\text{tanh}}}{\partial z} = 1 - \tanh^2(z)$

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

$$(w_{ij}^k)_{\text{new}} = (w_{ij}^k)_{\text{old}}$$

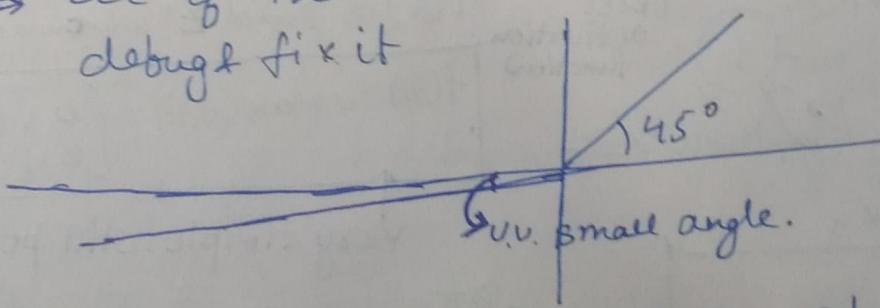
If z is negative every $\frac{\partial L}{\partial w_{ij}^k}$ this will 0. \uparrow dead activation.

$$0 * 1 + 1 * 1 * 1 \Rightarrow 0.$$

* ReLU also converge faster for same number of epochs than tanh function. (because it does not have vanishing gradient problem).

This is because if z is negative means weights are highly negative and derivation is always 0. the we can use variation called leaky ReLU.

→ See if there are more dead than active.
debug & fix it



This is known as leaky ReLU.

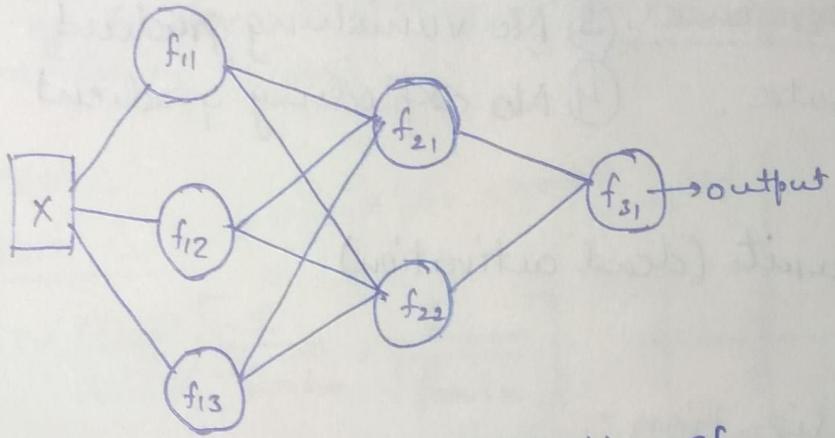
softplus function: $f(x) = \log(1 + \exp(x))$

derivative of softplus function $f'(x) = \frac{\exp(x)}{1 + \exp(x)} = \frac{1}{1 + \exp(-x)}$

leaky ReLU.

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0.01x & \text{if } x \leq 0 \end{cases}$$

$$f(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0.01 & \text{if } x \leq 0 \end{cases}$$

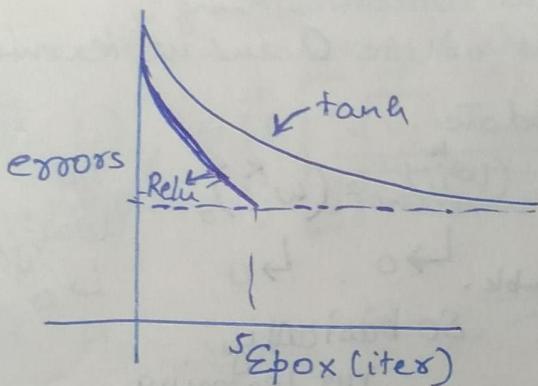


No problem of.

$$\frac{\partial f}{\partial z} \in [0, 1] \rightsquigarrow \begin{array}{l} \times \text{ exploding} \\ \times \text{ vanishing} \end{array}$$

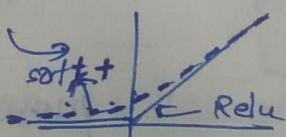
but can cause dead-activation

Experimentally how fast is ReLU.



so, for equal number of epochs ReLU converge faster than tanh.

To solve diff at 0 problem solution was softplus + function. $f(x) = \log(1 + \exp(x))$



This can also solve dead activation to certain level.

Variants of ReLU

① Noisy ReLU's this is like adding some bias
 $f(x) = \max(0, x + Y)$, with $Y \sim N(0, \sigma(x))$

② Leaky ReLU (avoid Dead activation) like mixture of Vanishing & ReLU.

$$\frac{\partial f_{LR}}{\partial z} = \begin{cases} 1 & \text{if } z > 0 \\ a & \text{if } z < 0 \end{cases}$$

$\hookrightarrow (0, 0.01)$

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

Advantage of ReLU.

- ① Speed up convergence.
- ② No vanishing gradient
- ③ easy to compute.
- ④ No exploding gradient

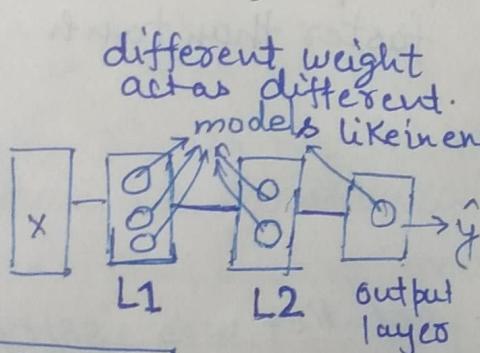
⑤ disadvantage

- ① dead ReLU units (dead activation).

⑥ Weight initialization:

Logistic reg → ① initialize the weight randomly
 or ② $w_{ij} \sim N(0, \sigma)$
 (gaussian)

- ① init $w_{ij}^k = 0 \quad \forall i, j, k \rightarrow$ v.v. bad idea.
 ↑ all
 → all Neuron will start with same thing.
 and always input will be 0 and no learning.
 → same gradient update.



$$(w_{ij}^k)_n = (w_{ij}^k)_0 - \eta \left[\frac{\partial L}{\partial w_{ij}^k} \right]$$

↳ 0 ↳ 0 ↳ 0

So basically,
 no learning.

RF, GBDT

ensemble → more different models are better
 the output, the better is output

- ② if $w_{ij}^k = \text{large -ve number}$ of ensembling.

ReLU $w^T x = z = -\text{ve large}$

normalized.
data ↳ Vanishing gradient
 or
 dead activation.

Solution

properties of weights.

- ideal ①-
 - ① weight should be small (not too small)
 - ② not all zero

③ good variance $\leftarrow \text{Var}(w_{ij}^k)$

$w_{ij}^k \sim N(0, \sigma)$ this σ is also very small but not too small.
 Gaussian/Normal initialization.

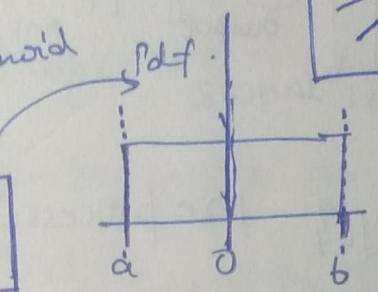
better init strategies.

- ↳ lots of experiments (strong) ✓ on how to initialise weight.
↳ some theory (weak)

Technique -

idea ② Uniform init \rightarrow v good for sigmoid

$w_{ij}^k \sim \text{Unif distib} \left[\frac{-1}{\sqrt{\text{fanin}}}, \frac{1}{\sqrt{\text{fanin}}} \right]$



fanin = 4
tout = 2
No. of output

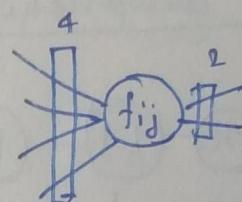
[no concrete agreement amongst all researchers.] \rightarrow that which is best initialized.

idea ③ Xavier/Glorot init (2010) \rightarrow v good for sigmoid.

a) $w_{ij}^k \sim N(0, \sigma_i)$ $\sigma_i = \sqrt{\frac{2}{\text{fanin} + \text{fanout}}}$

(normal) \downarrow
normal distribution

uniform distribution



b) $w_{ij}^k \sim U \left[-\frac{\sqrt{6}}{\sqrt{\text{fanin} + \text{fanout}}}, +\frac{\sqrt{6}}{\sqrt{\text{fanin} + \text{fanout}}} \right]$

(uniform)

idea ④ He-init. (2015) \rightarrow v. good for ReLU.

a) $w_{ij}^k \sim N(0, \sigma)$ $\sigma = \sqrt{\frac{2}{\text{fanin}}}$

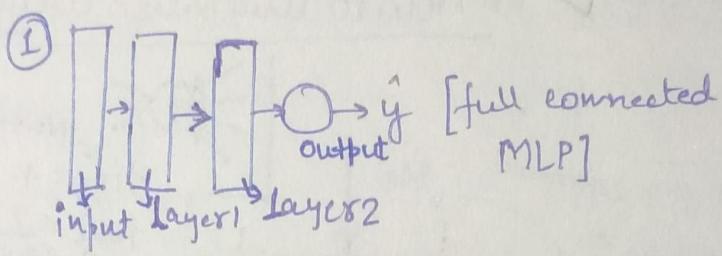
(normal)

b) $w_{ij}^k \sim U \left[-\sqrt{\frac{6}{\text{fanin}}}, +\sqrt{\frac{6}{\text{fanin}}} \right]$

(uniform)

He init is good for relu.

5) Batch Normalization :- (2015)



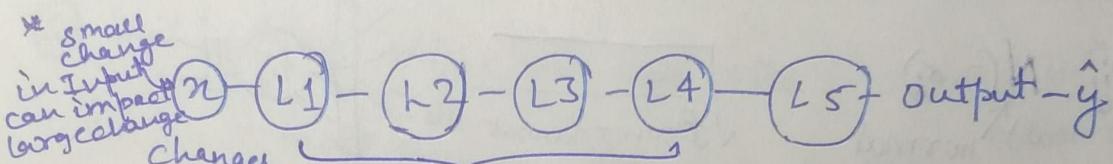
$D = \{x_i, y_i\}$ pre processing step \rightarrow data normalization on x_i i.e. mean centring and var-scaling

Dataset

$$x_i = \frac{x_i - \mu}{\sigma} \quad \begin{matrix} \text{mean}(x_i) \\ \sigma \end{matrix}$$

$$\text{std dev}(x_i)$$

ex of deep network



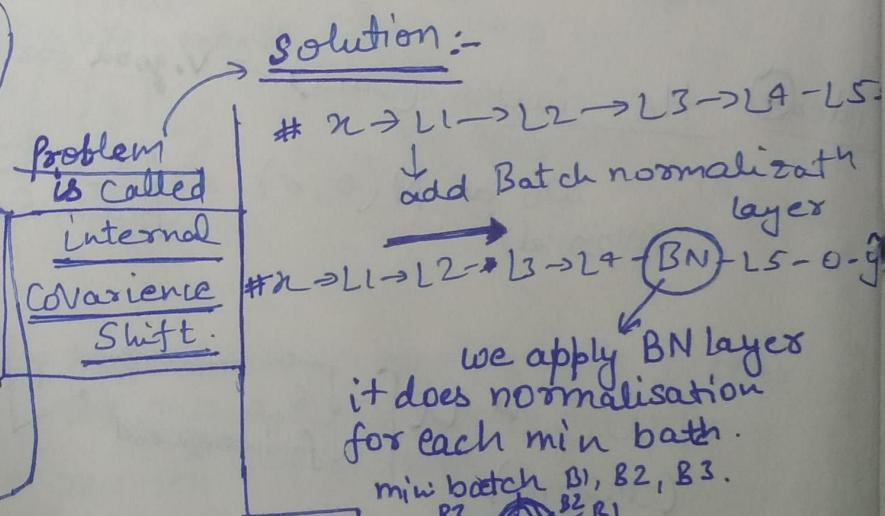
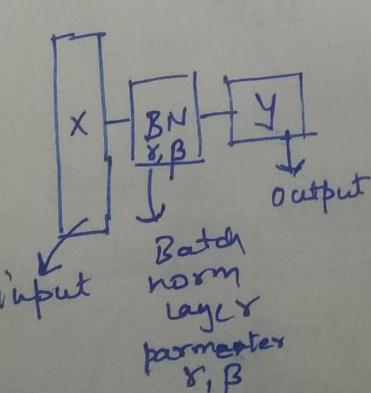
SGD

mini-batch SGD

Example

So many things are multiplying so small change in x and become large change in the Network at Layer 5.

1 \rightarrow 1.01 \leftarrow slight change
 1.01 \rightarrow 1.02 \rightarrow 1.04 after 1 layer
 1.04 \rightarrow 1.08 \rightarrow 1.16 \rightarrow 1.24 \rightarrow 1.32 \rightarrow 1.37 after 5 layers
 (1) will come to 1.375
 37% with slight change on sq function neuron at 5 layers



Input in BN x_i

Output: $\tilde{x}_i = BN_{\gamma, \beta}(x_i) \rightarrow$ It just mean input $z = w\tilde{x} + b$ not x .

$$\mu_{BN} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma_{BN}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{BN})^2$$

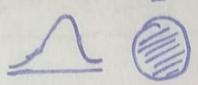
$$\hat{x}_i \leftarrow \frac{x_i - \mu_{BN}}{\sqrt{\sigma_{BN}^2 + \epsilon}}$$

$$\tilde{x}_i \leftarrow \gamma \hat{x}_i + \beta = BN_{\gamma, \beta}(x_i)$$

small value so that dividing by 0 error does not come.

Learning of γ, β is done just like weights learning it can be done using backpropagation.

- + Advantages of Batch Normalisation.
- + faster convergence \rightarrow LS



- larger learning rate.
 η

Can be used as $\frac{\partial L}{\partial w_k}$ will be small

- + act as weak regularizer
use BN + Dropout

and $\frac{\partial L}{\partial w_{ij}}$ will be small because all the values come from same distribution.

- + Avoid internal Covariate Shift \rightarrow train deep NN easy

⑥ Optimizers :- [Hill descend]

L.R. Reg & Optimization \rightarrow GD, SGD, mini batch-SGD

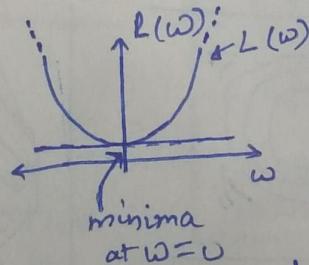
$$w^* = \min_w L(w)$$

W Loss function

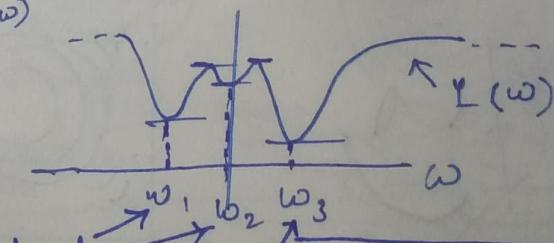
Deep learning \rightarrow simple SGD, GD \rightarrow advances??

$$\min_w L(w)$$

a) w is scalar
↑ 1 dimension



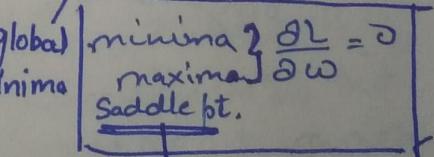
2 D visualisation



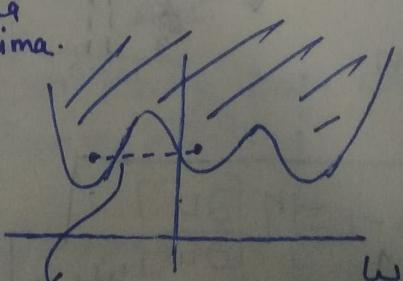
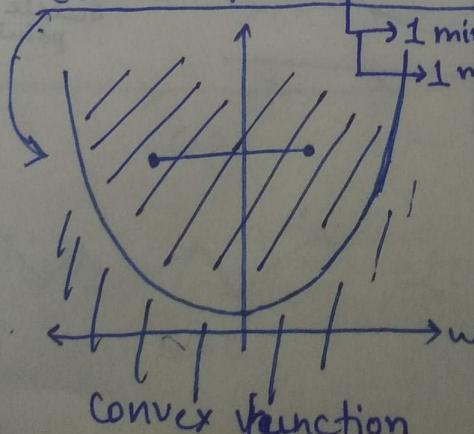
$$w_{new} = w_{old} - \eta \left(\frac{\partial L}{\partial w} \right)$$

Stop when $\frac{\partial L}{\partial w} = 0$

So it could stuck at saddle pt.



Convex function and Non-Convex function

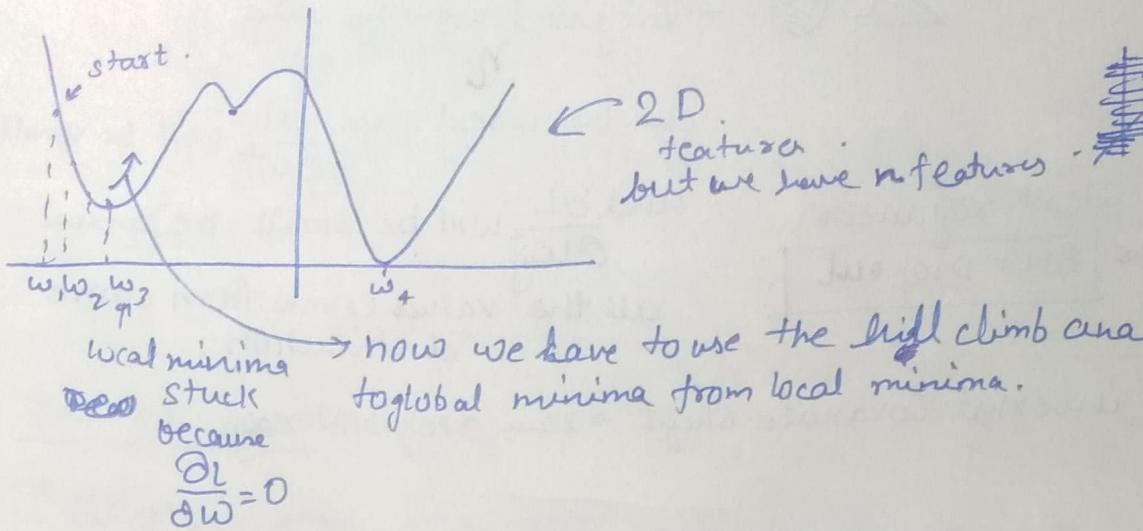


In Saddle pt
also $\frac{\partial L}{\partial w} = 0$

Some pts like in other region ie. it is not convex function.

property \rightarrow Convex function have only one minima or maxima.
 Local minima = global minima

In real life we have non convex functions.



Now we have to use the hill climb analogies to go to global minima from local minima.

So based on init value of w we can land of diff local minima but we have seen what is global minima and all depends on what initial weights we have initialised.

⑦ 3D and Contours (Hill descent)

$$w = [w_1, w_2]$$

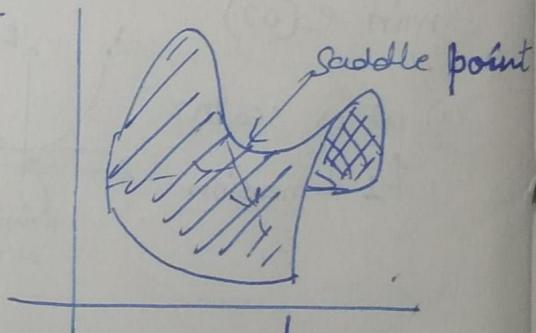


Contours ptot.



Contour pt. is 2 D

projection of 3D plot.



Saddle point

⑧ SGD & Momentum

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

Update function

$$w \leftarrow w_{ij}^k$$

notation for
simlicity

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

t^{th} iteration $t-1^{th}$ iteration



Saddle point.

$\eta \leftarrow 0.01$ generally.

but $\frac{\partial L}{\partial w}$?

$\rightarrow D\{x_i, y_i\}_{i=1}^n$ size of data

\rightarrow using all the n points in $D \Rightarrow$ G.D. \rightarrow consume so much time & memory.

\rightarrow using only one point $x_i @$ random? this is ~~batch~~

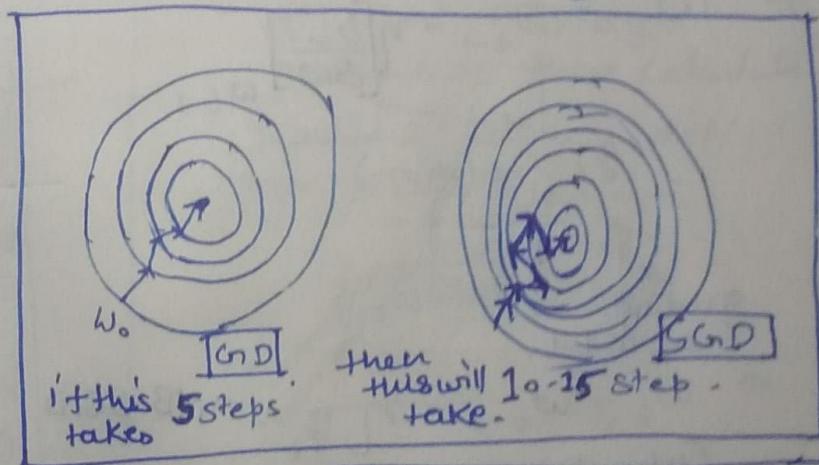
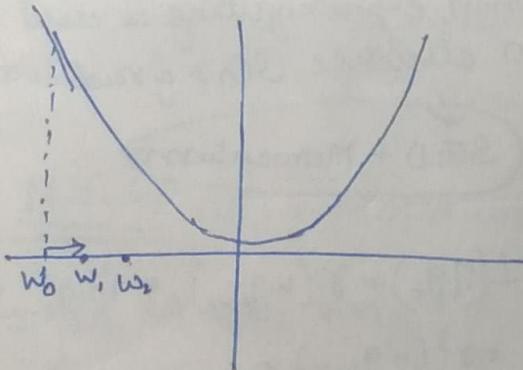
\rightarrow using a random subset of K points in D . SGD. What we have discussed earlier for backpropagation.

Minibatch SGD.

While applying minibatch SGD

$$\left[\frac{\partial L}{\partial w} \right]_{\text{minibatch SGD}} \approx \left[\frac{\partial L}{\partial w} \right]_{\text{G.D.}}$$

These are noisy but not equal to $\left[\frac{\partial L}{\partial w} \right]_{\text{G.D.}}$



Basically what we are doing is estimating gradient/des using SGD.

G.D.

w^* SGD result will be same as w^* G.D. but take more iterations.

Each of updates are more insipid than G.D.

\rightarrow G.D.

Eventually both will reach to the same point.

\rightarrow SGD

(Noisy) \rightarrow but we can remove this noise by using the concept of momentum.

How to denoising gradient in SGD?

⑨ SGD with momentum

change in gradient term
No change in learning rate term
just add a term.

denoise
iter $\rightarrow t_1, t_2, t_3, t_4$
val $\rightarrow a_1, a_2, a_3, a_4, a_5, a_6, \dots$

take averages.

PTCEZCUC

$$\left\{ \begin{array}{l} \text{at } t_1, v_1 = a_1 \\ \text{at } t_2, v_2 = \gamma v_1 + a_2 \end{array} \right. \quad \begin{array}{l} \text{recent point.} \\ \gamma = 1 \\ 0 \leq \gamma < 1 \end{array}$$

recursive equation!

recent point has more weight.

$$t_3 = \gamma^2 a_1 + \gamma a_2 + a_3$$

oldest point middle point newest point

also called exponential moving average.

Back to SGD: Mini-Batch

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

g_t (notation)

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

gradient at w_{t-1} .

$$w_t = w_{t-1} - \eta g_t \rightarrow \text{MB-SGD}$$

(update)

No. what is this.

\Rightarrow (exp-weighting)

$$\left\{ \begin{array}{l} v_t = \gamma v_{t-1} + \eta g_t \\ w_t = w_{t-1} - v_t \end{array} \right.$$

new momentum in SGD which is added.

v_t is $(0, 1)$ generally 0.9

gradient as old

gradient + momentum

Case 1: $\gamma = 0$

$$v_t = \eta g_t$$

$$w_t = w_0 - \eta g_t$$

* When exp-weighting is used to denoise SGD gradient

SGD + Momentum

$$v_t = \underbrace{1(\eta g_t)}_{\text{SGD}} + \underbrace{\gamma(v_{t-1})}_{\text{momentum}} + \underbrace{\gamma^2(g_{t-2})}_{\text{old}} + \underbrace{\gamma^3(g_{t-3})}_{\text{old}} + \dots$$

Core-concept is \rightarrow removing noise using exp-weight.

Case 2: $\gamma = 0.9 v_{t-1}$

$$w_t = w_{t-1} - (0.9 v_{t-1} + \eta g_t)$$

v_{t-1} old

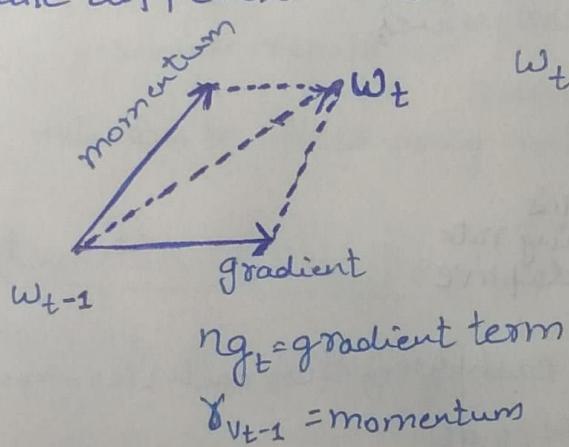
momentum (γv_{t-1}) \leftarrow it has all the information of past movements of $w_{t-2}, w_{t-3} \dots w_1$.
 $w_t = w_{t-1} + (-\gamma v_{t-1} - \eta g_t)$
 w_t \leftarrow this is the main movement all this point but it has only current information about w_{t-1} .

* SGD + momentum \rightarrow Speedup convergence.

and overcome the drawback of slow convergence in

Nesterov Accelerated Gradient (NAG) \leftarrow No change in learning rate, no change in gradient just add another term

(10) It is the similar technique as SGD + momentum but with subtle difference and better than SGD + momentum.



NAG. (Who people write)

$$\begin{cases} v_t = \gamma v_{t-1} + \eta g'_t \\ w_t = w_{t-1} - v_t \end{cases}$$

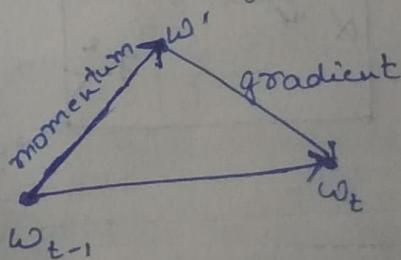
g'_t is at w_{t-1}

In NAG g'_t is at $w_{t-1} + \gamma v_{t-1}$.

$$g_t = \frac{\partial L}{\partial w}|_{w_{t-1}}$$

$$g'_t = \frac{\partial L}{\partial w}|_{w_{t-1} + \gamma v_{t-1}}$$

Nesterov AG
 If γ says don't make resultant of momentum + gradient move in momentum then calculate gradient' which is diff. from earlier gradient.

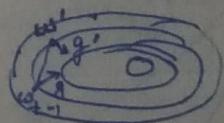


$$w_t = w_{t-1} - (\gamma v_{t-1} + \eta g'_t)$$

$$g'_t \neq g_t$$

$$g' = \frac{\partial L}{\partial w}$$

$$(w_{t-1} - \gamma v_{t-1})$$



Because we have move in that direction already.

⑪ Adagrad (Adaptive gradient) No change in gradient change in learning rate

→ In SGD & SGD with momentum.

: learning rate ($\eta = 0.01$) → same for each weight parameter.

In Adagrad → Adaptive gradient.

idea → each weight/parameter has different η .
why?

features → sparse? (BOW) most of them are zero
dense → (less zero values)

$$SGD: \quad w_t = w_{t-1} - \eta g_t$$

Adagrad: $w_t = w_{t-1} - \eta' t g_t$ Not this learning rate is adaptive.

diff η' for each weight @ each iteration.

$\eta' t = \eta$ ← constant like (0.01)

$\alpha_{t-1} = \sum_{i=1}^{t-1} g_i^2$ small +ve number (to avoid divide by zero error)

α_{t-1} is +ve $\frac{\partial L}{\partial w}$ is +ve that's why squaring

$$g_1 \ g_2 \ g_3 \ g_4 \dots \ g_{t-1} \ g_t$$

as $t \uparrow \rightarrow \alpha_t \uparrow \rightarrow (\eta'_t) \downarrow$ as t increases α increases because we are summing the square

Hence learning rate will reduce constantly with iterations.

So, with the increasing iteration, learning rate is decreasing. Convergence in start is fast and slow to avoid overshooting problem.

Advantage & Disadvantage

(+) no need for manually tuning η

$\eta' t \rightarrow$ weight, time iter

(+) Sparse & dense feature \rightarrow Adagrad

disadvantage

(-) α_{t-1} can become very large as t increase.

$$\text{And } w_t = w_{t-1} - \eta g_t \quad \text{↑ release}$$

very small because α_{t-1} very large.

which result in $w_t \approx w_{t-1}$.

But this can be fixed using next algo. amount of movement is so small.

(12) Adadelta & RMS prop \rightarrow only change in learning rate
 \rightarrow NO change in gradient

Adagrad α_{t-1} v. large \rightarrow slow convergence

$$\eta' t = \frac{\eta (=0.01)}{\sqrt{\alpha_{t-1} + \epsilon}} ; \quad \alpha_{t-1} = \sum_{i=1}^{t-1} g_i^2$$

$\eta' t$ is becoming small because α_{t-1} is becoming large

α_{t-1} is becoming large because it is summation of g_i^2 .

Idea: here is take exponentially decaying rate. (eda)

$$\text{Adadelta: } w_t = w_{t-1} - \eta' t g_t$$

$$\eta' t = \frac{\eta}{\sqrt{\text{eda}_{t-1} + \epsilon}}$$

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

exponential decay average

$\Rightarrow \gamma$ is hyperparameter typically 0.95.

$$\text{so } \text{eda}_{t-1} = 0.95 \text{eda}_{t-2} + 0.05 g_{t-1}^2$$

~~Adam: Adaptive Moment Estimation~~

$$\text{eda}_{t-1} = (0.05) g_{t-1}^2 + 0.95 \text{eda}_{t-2} \quad (1)$$

$$\text{eda}_{t-2} = (0.05) g_{t-1}^2 + 0.95 \text{eda}_{t-3} \quad (2)$$

$$\text{eda}_{t-1} = (0.05) g_{t-1}^2 + [0.95 * \{0.05 g_{t-2}^2 + 0.95 \text{eda}_{t-3}\}]$$

~~Now eda_{t-1}~~
 do not grow
 v.v. large and
 v.v. slow η'_t .
 $\text{eda}_{t-1} = 0.05 g_{t-1}^2 + (0.95 * 0.05) g_{t-2}^2 + (0.95)^2 \text{eda}_{t-3}$
 5% of present weight
 5% times to prev. top rev.
 Weight
 95% times of present weight

But doing this we can control the growth of

g_{t-1} or named here eda_{t-1} , (denominator)
 (like in adagrad) $(\text{ada}\delta\text{elta})$

RMS prop is same as Adadelta but some thing are diff.
 but idea is same.

(13) Adam: Adaptive Moment Estimation. (2018) best
 ↳ Both learning rate & gradient changed & made something else.
 Idea: Adadelta → storing $\frac{\text{exp weight avg of } g_t^2}{\text{eda}}$ → learning rate (hit)

In statistic mean, → 1st order moment

Var., → 2nd order moment

$$\begin{cases} m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t & (1) \quad 0 \leq \beta_1 \leq 1 \quad \text{typically } \beta_1 = 0.9 \\ V_t = \beta_2 V_{t-1} + (1-\beta_2) g_t^2 & (2) \quad 0 \leq \beta_2 \leq 1 \quad \text{typically } \beta_2 = 0.99 \\ \text{eda}_t \end{cases}$$

$$\text{eda}_t = \sqrt{m_t / V_t}$$

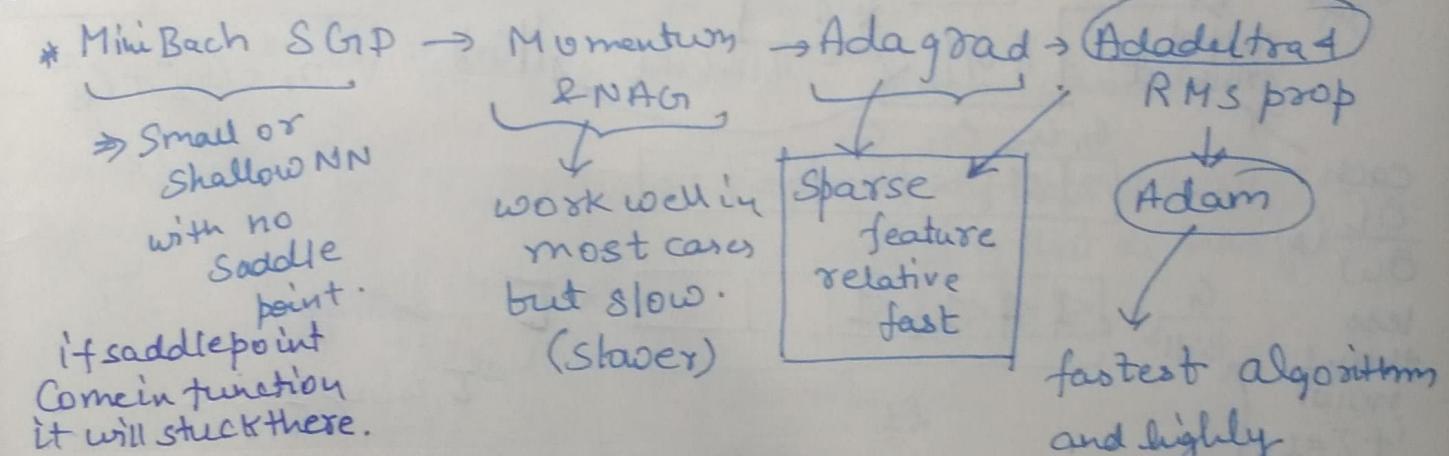
$$\left\{ \begin{array}{l} \hat{m}_t = \frac{m_t}{1 - (\beta_1)^t} ; \quad \hat{v}_t = \frac{v_t}{1 - (\beta_2)^t} \\ w_t = w_{t-1} - \frac{\alpha \hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}} \end{array} \right.$$

} Adam is best that comes by lot of experiment.
} NO need to remember formula.

If $\beta_1 = 0$
 $m_t = g_t$ } Adadelta

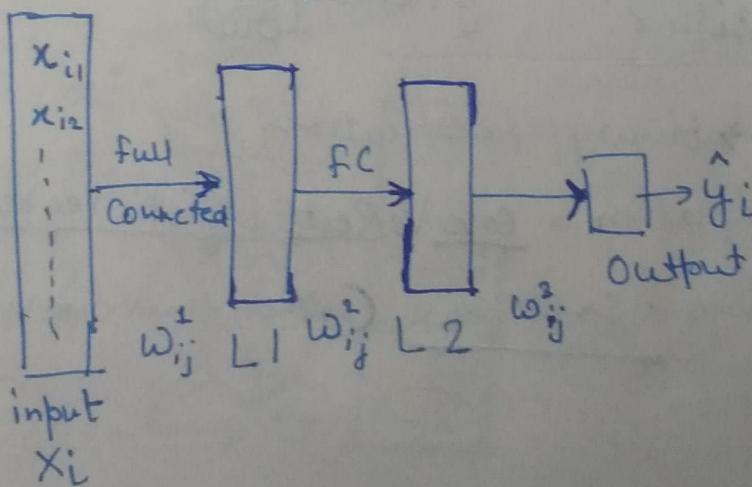
If $\beta_1 = \beta_2 = 0$

14) Which optimizer to use.



15) Gradient Monitoring & clipping.

used to update weights.



* no need to specify η .

→ good habit to monitor gradient and update

↓
to detect prob of vanishing gradient

If gradient is large then prob. come will exploding gradient. Then we can visualize it can be used as change the optimization by Relu.

Another solution is Clipping, w_{ij}^k

$$W = \begin{bmatrix} w_{11}^1 & w_{12}^1 \\ \vdots & \vdots \\ w_{11}^3 & \vdots \end{bmatrix}$$

G_i = Gradient: $\frac{\partial L}{\partial w_{11}^1}, \frac{\partial L}{\partial w_{12}^1}, \dots, \frac{\partial L}{\partial w_{11}^3}$ some are very large then we can do clipping otherwise it will create exploding gradient.

L2 norm clipping

$$G_{i, \text{new}} = \frac{G_i}{\|G_i\|_2} * T$$

clipped gradient

$$G_i = \begin{bmatrix} G_{i1} & G_{i2} & G_{i3} & G_{i4} \\ 2 & 5 & 10 & 10000 & 0.5 & 50 \end{bmatrix}$$

$$\|G_i\|_2 = \sqrt{G_{i1}^2 + G_{i2}^2 + G_{i3}^2 + G_{i4}^2 + \dots}$$

$$\frac{G_i}{\|G_i\|_2} = \begin{bmatrix} 1 & 1 & 1 & 1 & \dots \\ \underbrace{\quad}_{<1} \end{bmatrix} * T$$

if $T=1$ or less than T if T is anything else.

this is called
L2 norm clipping of Gradient

and it is done to avoid V.V. Large gradient ($\frac{\partial L}{\partial w}$).

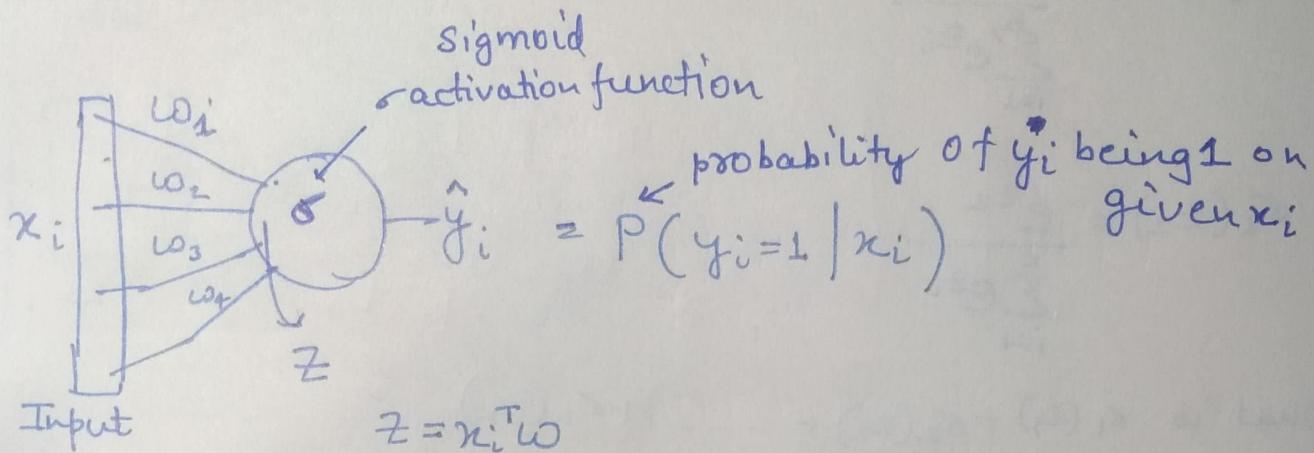
16) Softmax - classifier:

logistic regression \rightarrow binary classification

for multi-class classification \rightarrow One Vs Rest Logistic regression

log. reg + Multi-class = Softmax (rather than One Vs Rest)

~~Logistic~~ regression recap.

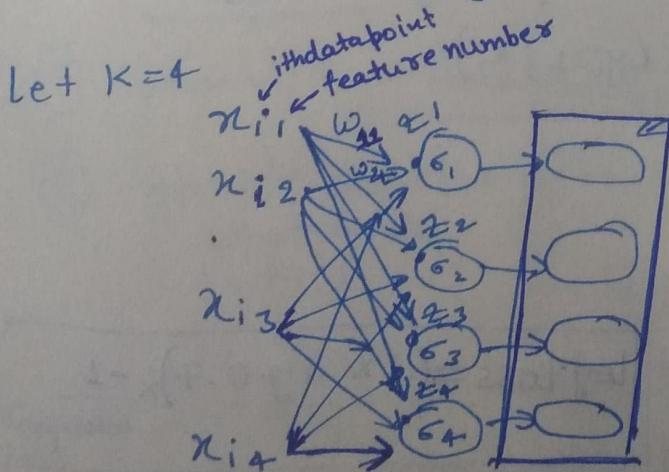
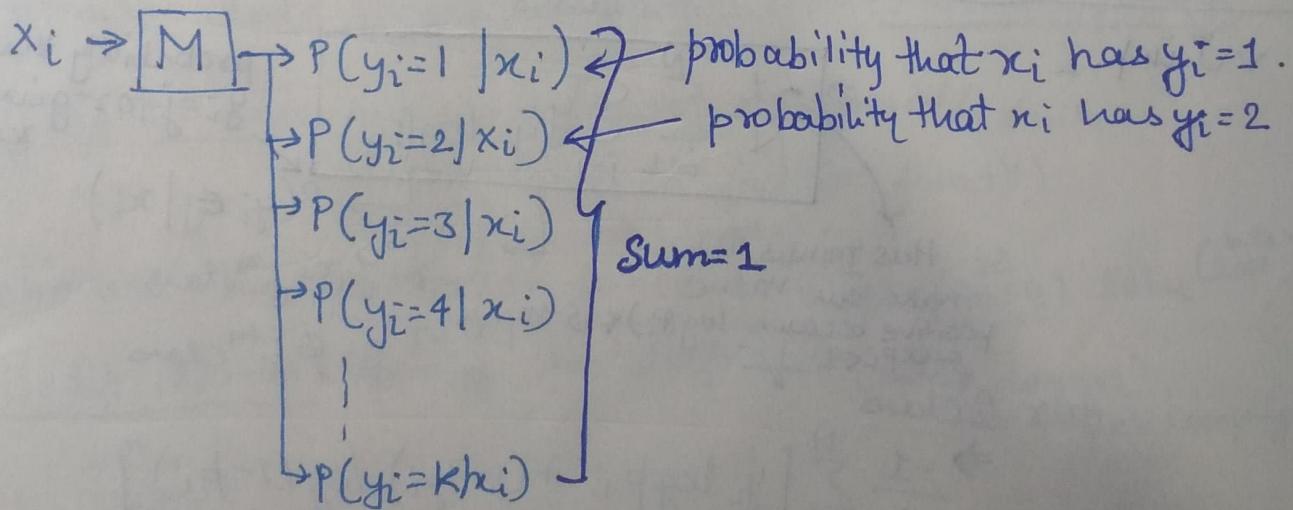


$$P(y_i=1 | x_i) = \hat{y}_i = \sigma(z) = \frac{1}{1-e^{-z}} = \frac{e^z}{e^z + 1}$$

Softmax classifier

$\mathcal{D} = \{x_i, y_i\}$; $y_i \in \{0, 1\}$ in Log. reg.

$y_i \in \{1, 2, 3, \dots, K\}$ in softmax classifier.



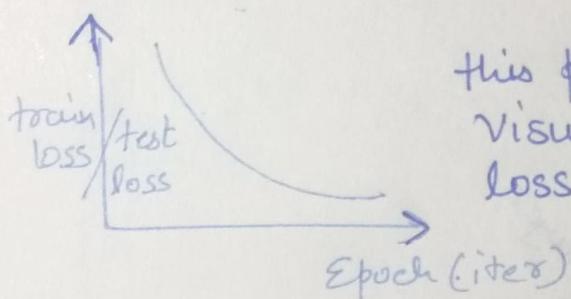
Sum to 1.

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

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

⑧ monitor gradients \rightarrow apply clipping

⑨ plots



this plot is very imp to visualise the train/test loss with number of epochs.

⑩ Avoid overfitting \rightarrow Deep MLP are easy to overfit in multilayer perceptron.

⑪ Auto-encoder (AE) \rightarrow revisit dimension-reduction.

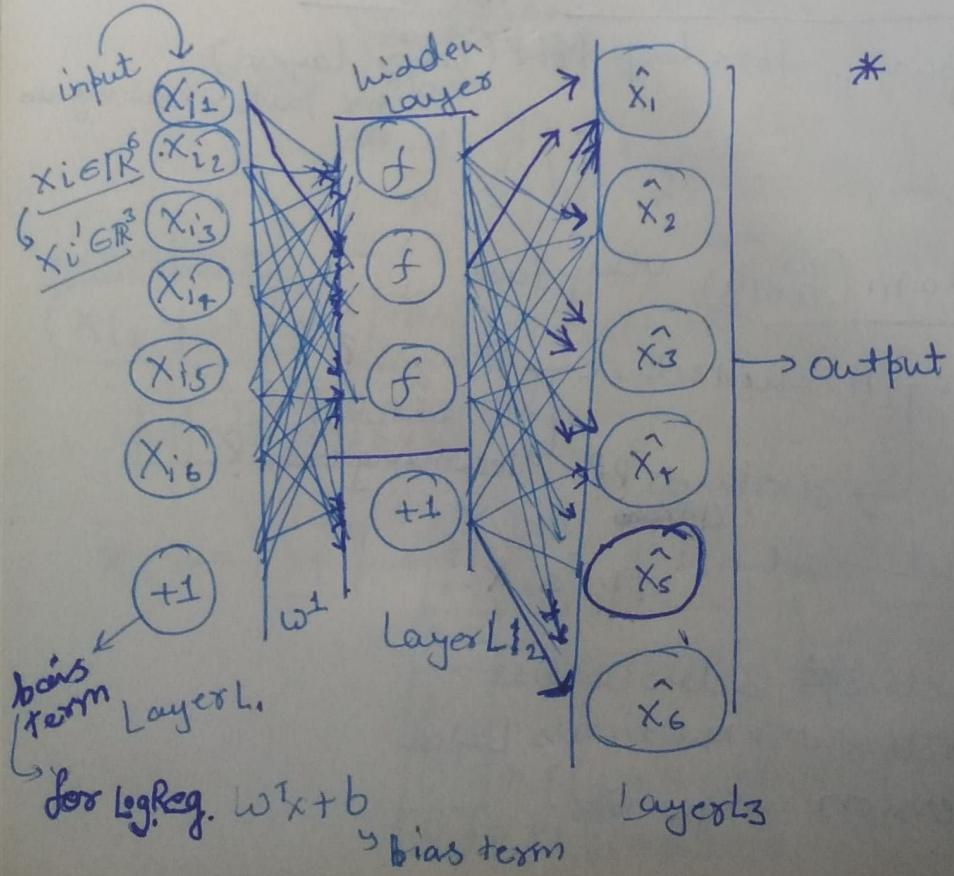
* Neural network which perform dim-red. like [PCA, tSNE.]

* Simple Auto Encoder

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

$$\mathcal{D}' = \{x'_i\}_{i=1}^n \quad x'_i \in \mathbb{R}^{d'} \quad d' < d$$

Variance
Neighbourhood.



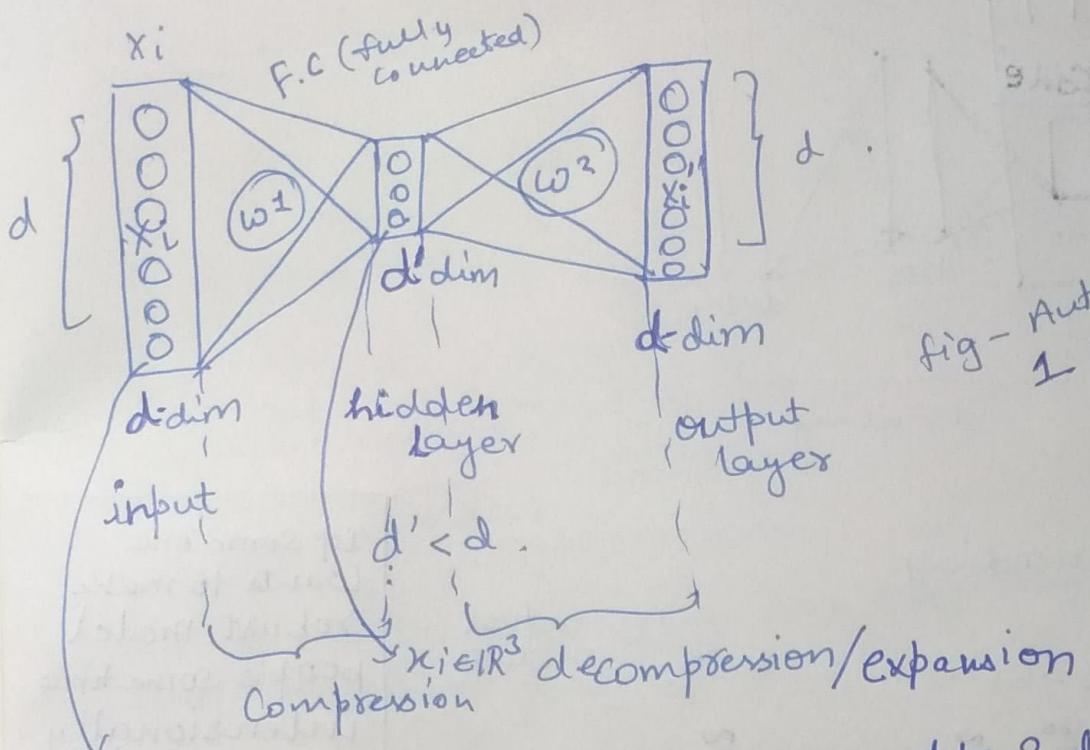
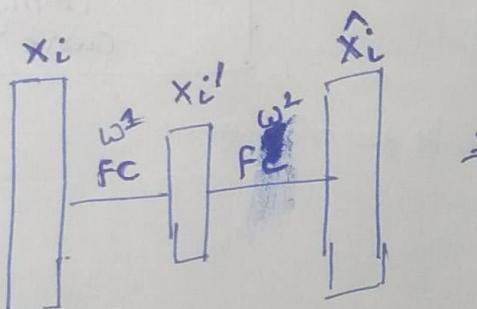


fig - Auto Encoder with 1 hidden layer

$x_i \in \mathbb{R}^d$ * Here d dim are compressed to 3 dim & then all these 3 dim can be obtained from 3 dims.

Actual NN of Auto Encoder



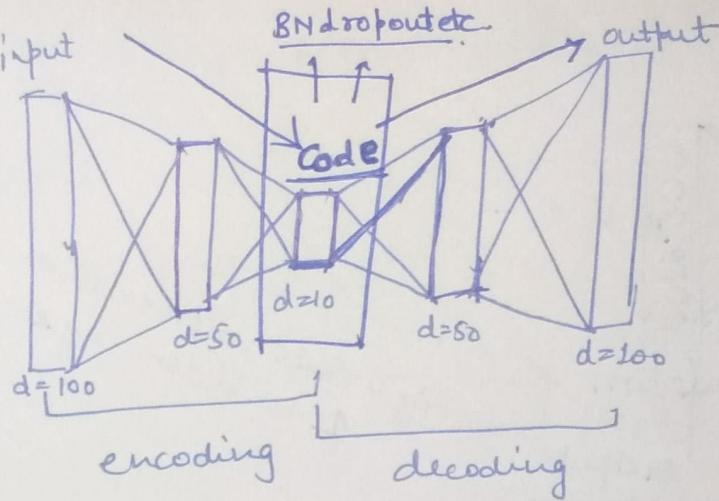
Q. So what we want from autoencoder??
→ want $\hat{x}_i \approx x_i$ → decomparsision after Comprasion

$$L(x_i, \hat{x}_i) = \|x_i - \hat{x}_i\|^2 \quad \text{lets say squared loss / distance loss.}$$

① differential, so we can minimise it or do optimisation.

fun part is $x_i' \ll x_i$ so we need not create v.v. deep model.

Layer wise pretraining very popular but now Adam is used.



Denoising Auto Encoding

$D = \{x_1, x_2, x_3, \dots, x_n\}$ actual data
 we get $\tilde{D} = \{\tilde{x}_1, \tilde{x}_2, \tilde{x}_3, \dots, \tilde{x}_m\}$ noisy & corrupt

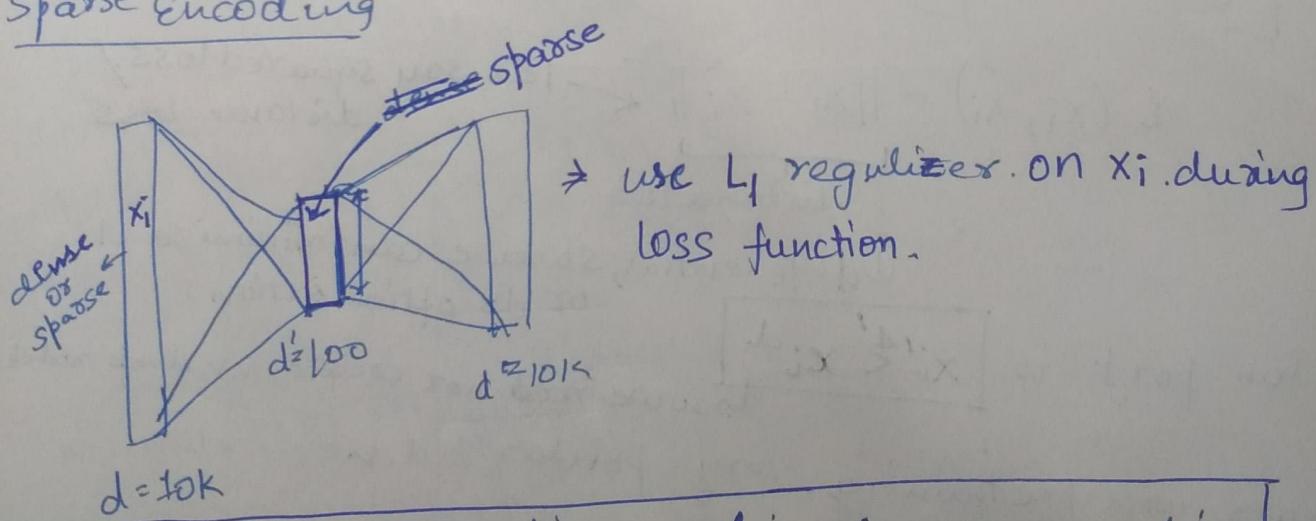
$$\tilde{x}_i \rightarrow \boxed{AE} \rightarrow x_i \in \mathbb{R}^{d'}$$

$$d' < d$$

automatically noise will be removed.

If some one wants to make robust model people sometimes intentionally add noise
 $\tilde{x}_i = x_i + \tilde{N}(0, \sigma)$
random value from gaussian curve,

Sparse Encoding



* if Activation function used is sigmoid output is much similar to PCA with 1 hidden layer.

19 Word 2. Vec → Amazon fine food reviews.

$\{w2v(w) \rightarrow \text{vectors}\}$

not a deep learning algo. but see the NN explanation.

def $\begin{array}{ccccccc} C_5 & C_4 & f & C_1 & C_2 & C_3 \\ \text{The} & \text{cat} & \underline{\text{sat}} & \text{on} & \text{the} & \text{wall.} \\ \text{Context} & \text{words} & \text{focus word} & & \text{Context words} & \end{array}$

Core Id is context words are very useful in understanding focus word.

2-algorithm $\begin{cases} \text{C-Bow (continuous bag of words)} \\ \text{(popular)} \quad \text{skip gram} \end{cases}$

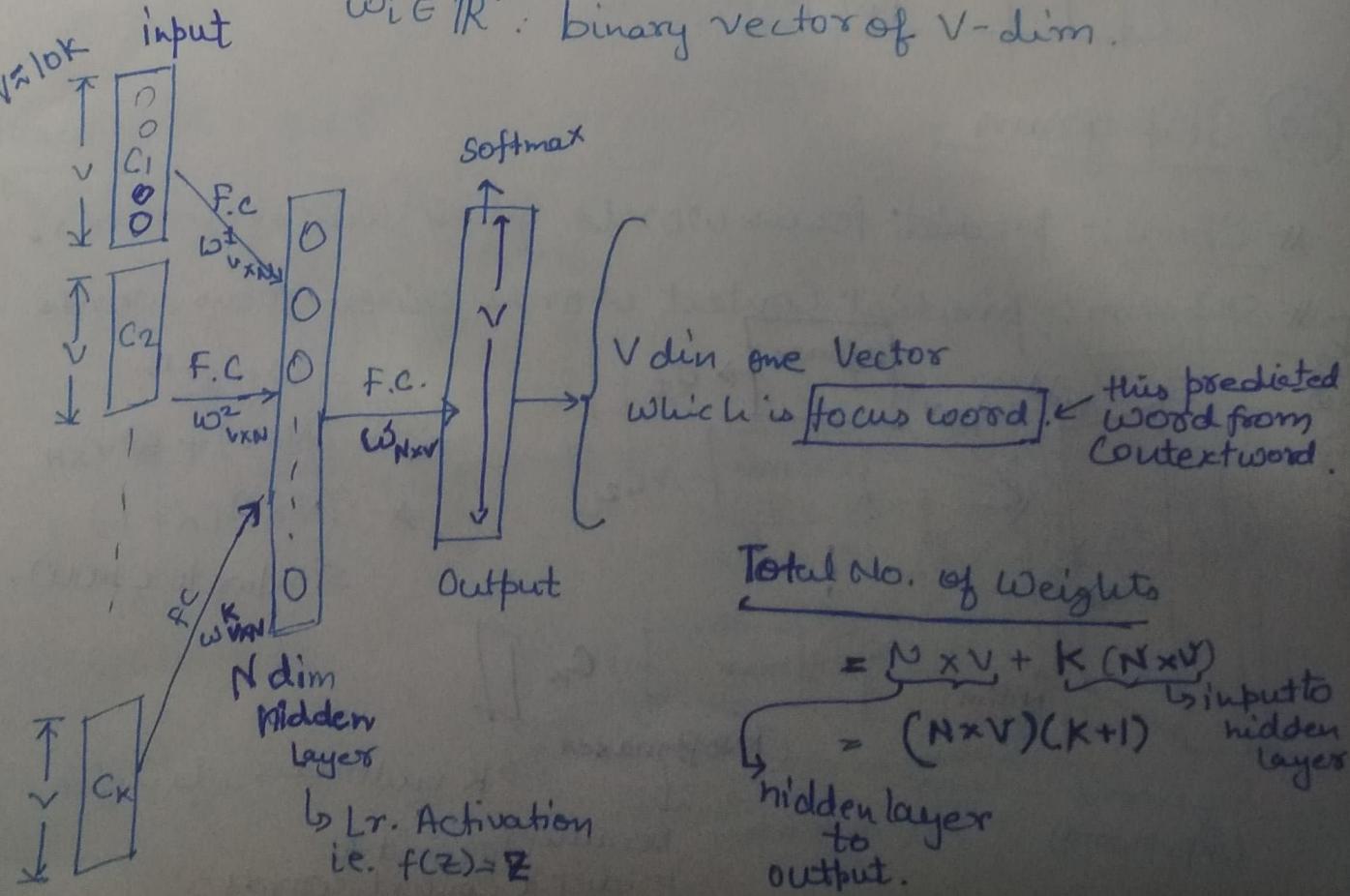
C-Bow → given context word predict focus word.

① dictionary/Vocabulary of words = V

$V = \text{length}/\text{size of words}$

② One-hot encoding

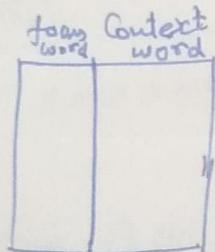
input $w_i \in \mathbb{R}^V$: binary vector of V -dim.



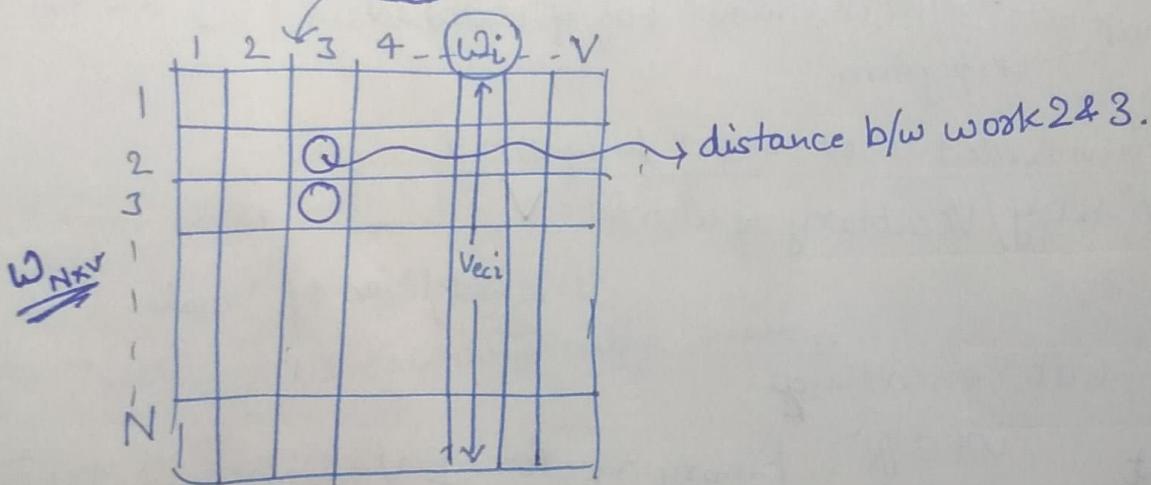
Train? CBOW

text

- * Make combination of focus words & context words. that becomes training data.



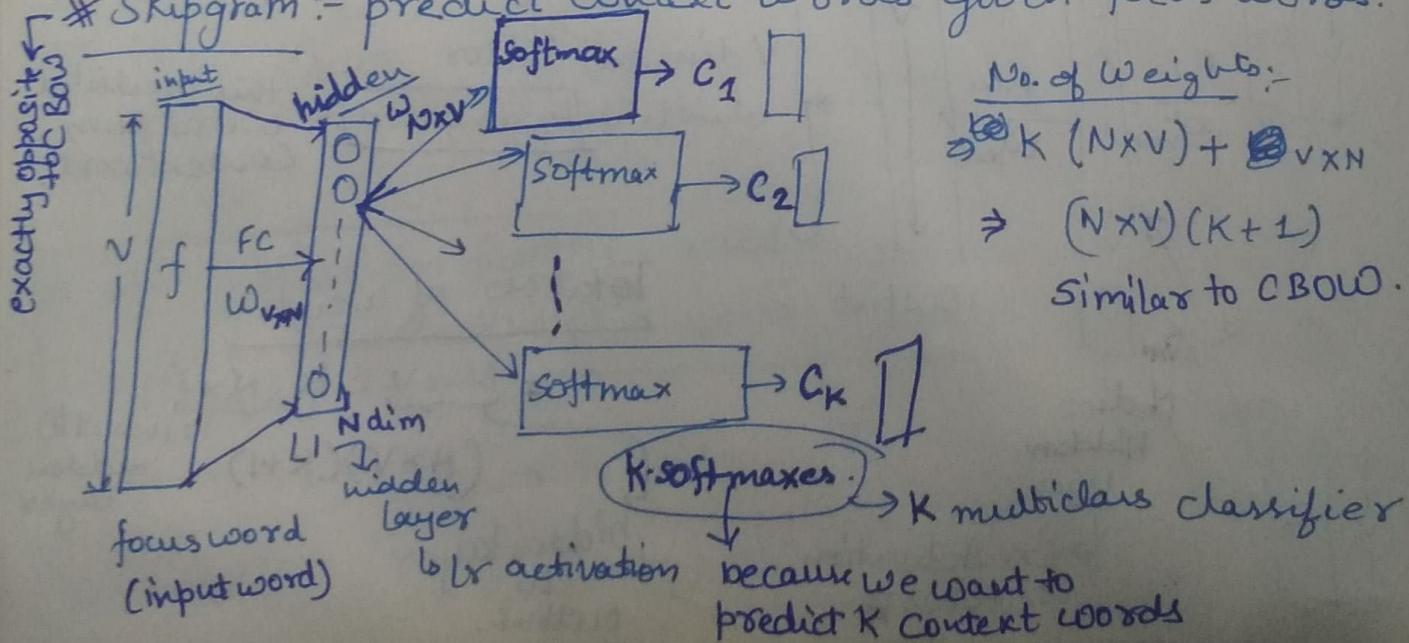
The vector $w_{N \times V}$ has summary of $w_{V \times N}^1 + w_{V \times N}^2 + \dots + w_{V \times N}^K$



20) Skip gram

* CBOW:- predict focus words given Context words.

* Skipgram:- predict Context words given focus words.



$$\begin{aligned} \text{No. of Weights: } & \\ \Rightarrow & K(N \times V) + V \times N \\ \Rightarrow & (N \times V)(K+1) \end{aligned}$$

Similar to CBOW.

because we want to predict K context words

In both w. weight are same but

In CBOW \rightarrow 1-softmax ?

In Skipgram \rightarrow K-softmax } Skipgram is computationally more expensive (takes more time)

CBOW \oplus faster to train than skipgram

\oplus better for frequently occurring word (more variety of data for same focus word so better training)

Skipgram \rightarrow \oplus can work with less amount of data

\oplus better with less freq words also.

For both skipgram & CBOW \rightarrow

lets say 'k' \rightarrow # context words

$K \uparrow \rightarrow$ more context $\rightarrow N \cdot \text{dim}$ for each pt is better.

How hard is this?

$$\# \text{weights} = (K+1)(N \times V)$$

\downarrow \downarrow \downarrow
 ≈ 6 ≈ 200 $\approx 10k$

$$= 1.2k \times 10k$$

= 12 Million (impossible) for normal computers.

So we need to optimise this??
we see how they use algorithms to optimise this?

21 Algorithm Optimization

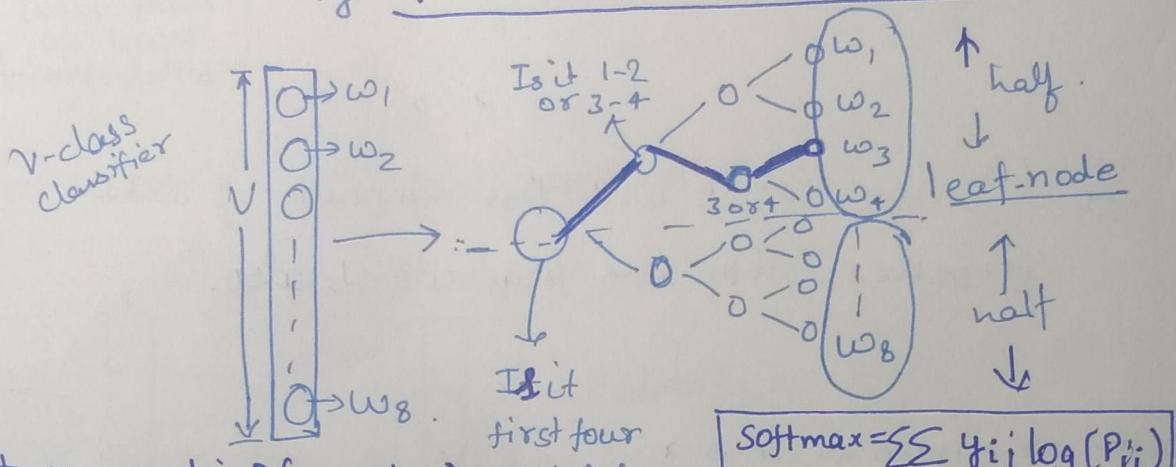
Algorithms
CBow & Skipgram \rightarrow Million of weights.

① Hierarchical softmax (algo)

② Negative Sampling (stat)

① Hierarchical softmax.

Some how modify V-softmax to make it optimal

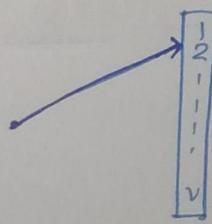


What we want is $P(y_i=3|x_i)$
rather than checking for all
the 8 units we use some
tree kind of structure which
goes to path 1 \rightarrow 3 \rightarrow 4 \rightarrow 4 i.e. ③ computations only.

So instead of 8 activation we need.

3 activation i.e. $\lceil \log(V) \rceil$ complexity now.

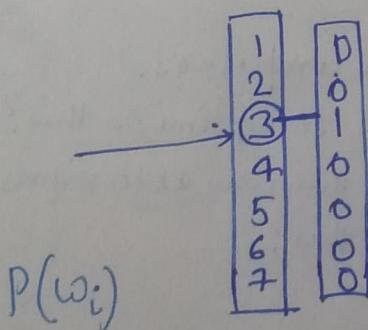
② Negative Sampling (stats)



{ update only a sample of words per iteration }

① always keep the target word

② Amongst non-target word update sample of them



$$P(w_i) = 1 - \sqrt{\frac{1}{\text{freq}(w_i)}}^{10^{-5}}$$

for more study \rightarrow

<https://blog.acolyer.org/2016/04/21/the-amazing-power-of-word-vectors/>

if word is very freq. $\sqrt{\frac{1}{\text{freq}}}$ is small $P(w_i)$ is large