48.) **Deep mult layer Perceptrons : 1980's to 2018's**

1980's, 1990's, 2006 ⎤ → 2-3 layerd n/w
2010 ⎟→ Vanishing gradient → ①
⎟→ too little data → overfit → ②
⎦→ too little compute → too much time → ③

Early 2010 ÷ → we had lots of data
↳ we had labelled data (because of internet
Companies)
↳ we have a new type of Compute infrastructure
↳ Gpu
↳ v.v. good & super suitable for deep learning

↳ new ideas & new algorithms


classical ML → SVM (90's)        { ① ②
→ RF (2000's)         Theory, Experiments.


**modern ml**

Experiment & Theory
①            ②

## 48.2 Dropout layers & Regularization

Deep nn → overfitting ← $L_1, L_2$
↳ many layers ⇒ many weights

RF: we are doing randomization of features to create regularization.

RFR (Randomization for Regularization) → nn (MLP's)
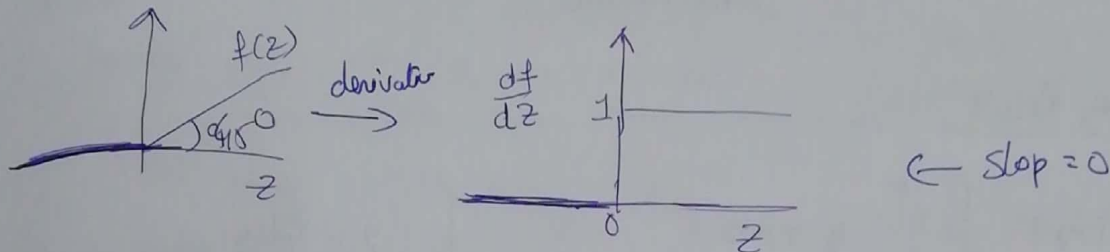
dropout rate : $0 \leq p \leq 1$

→ If I say P=0.2, in every layers 0.2 neurons will be dropped out

-) dropout ≈ random subset of features

48.3 Rectified Linear Units (ReLU) → Converge faster

As of 2018 → Best activation
↳ default activation in MLP

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



ReLU : not differentiable at zero

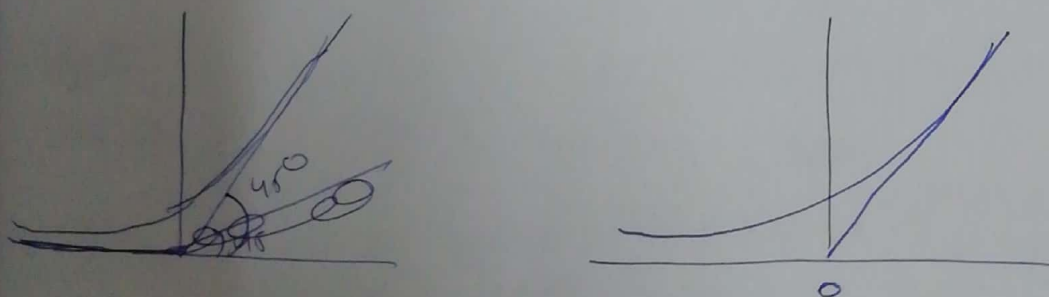$$\frac{d f_{ReLU}}{dz} \in \{0, 1\} \rightarrow$$ we don't have any value called Exploding gradient or vanishing gradient

→ As we have zero we will have dead activations

A smooth approximation of ReLU is the **Softplus** function

$$f(x) = \log(1 + \exp x).$$

The derivative of softplus is $f'(x) = \dfrac{\exp(x)}{1 + \exp x} = \dfrac{1}{1 + \exp(-x)}$ ⇒

which is a logistic function



$$\frac{d f_{ReLU}}{dz} = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \end{cases}$$
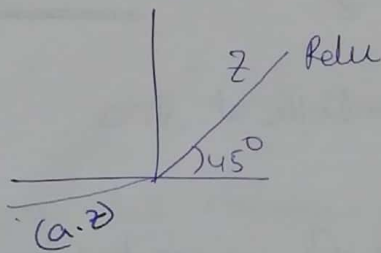
soft plus

# Noisy ReLU's

~~let good~~ $f(x) = max(0, x + y)$ with $y \sim N(0, \sigma(x))$

↑

adding a random value from a gaussian distribution with some variance

# Leaky ReLu's

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

$a$ = hyper parameter

↳ This can lead to vanishing gradient

48.4 Weight Initialization : Deep MLP

→ Logistic Regression → inialize weight randomly
$\quad\quad\quad\quad\quad\quad\quad\hookrightarrow W_{ij}^k \sim N(0, \sigma)$

uniform    Gaussian random
random.

Idea 1

initialize $w_{ij}^k = 0 \; \forall \; i,j,k \rightarrow$ V. V. bad.
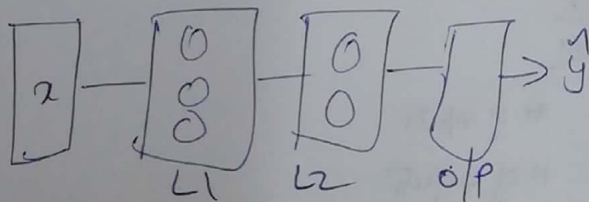$\quad\quad\quad\quad\quad\quad$ for all
→ $f_{ij}$ : Relu
→ all neurons compute the same thing
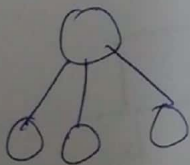→ Same gradient updates happen to all neurons  } Symmetry

If we initialize either 0 & 1 & 2 same problem as below (Symmetry)

→ we want models to be asymetry



$\quad\quad\quad$ L1 $\quad\quad$ L2 $\quad\quad$ o/p

Ensemble : more different the base models are, the better will be
the output of Ensembling

Stacking



Idea 2

Initialize $w_{ij}^k$ = large -ve numbers $\quad\quad$ relu

Relu $\quad w^T x = z \equiv$ large -ve value $\approx f(z) = 0$
$\quad\quad\quad\quad\quad\uparrow$
$\quad\quad\quad\quad$ normalized $\xleftarrow{\quad}$ mean centry
$\quad\quad\quad\quad\quad\quad\quad\quad$ variance scaling
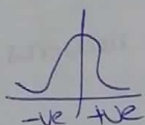
* data normalization is mandatory

# Solutions

## Idea 1

→ weights should be small (not too small)

→ not all zero

→ good-variance ← $Var(w^k_{ij})$

⟶ All weights come from a ~~star~~ Gaussian distribut

$$w^k_{ij} \sim N(0, \sigma) \quad (\sigma: \text{is reasonably small})$$
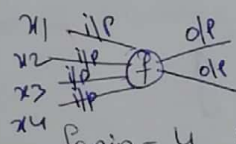
① $\sigma$ ✓

② +ve, -ve

Centered at 0

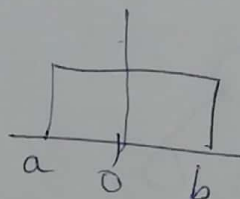## better init strategy

→ fan - in

→ fan - out

fanin = 4   # of inputs

fanout = 2   # of outputs

## uniform initialization (work well for Sigmoid)

$$w^k_{ij} \sim unif.\left[\underbrace{\frac{-1}{\sqrt{fanin}}}_{a}, \underbrace{\frac{1}{\sqrt{fanout}}}_{b}\right]$$

min value   max value

→ no concrete agreement amongst all researchers about best initialize

## Xavier/Glorot init (works well for sigmoid)

<u>normal</u>

a) $w^k_{ij} \sim N(0, \sigma)$      $\sigma_j = \sqrt{\dfrac{2}{fanin + fanout}}$

<u>unif</u>

b) $w^k_{ij} \sim u\left[\dfrac{-\sqrt{6}}{\sqrt{fanin + fanout}}, \dfrac{+\sqrt{6}}{\sqrt{fanin + fanout}}\right]$

for Each neuron fanin and fanout will change

He -Iniatilizatia (2015) (Good for ReLu)

a) Normal

$$w_{ij}^k = \sqrt{\frac{2}{fanin}}$$

b) uniform

$$w_{ij}^k \sim u\left[\frac{-\sqrt{6}}{\sqrt{fanin}}, \frac{+\sqrt{6}}{\sqrt{fanx}}\right]$$
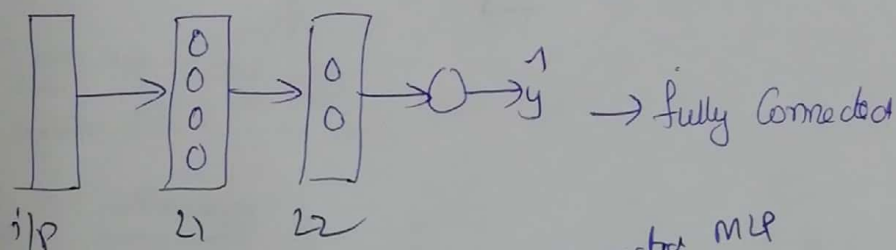
## 48.5 Batch normalization (2015)

$D = \{x_i, y_i\}$

Preprocessing : Data normalization $(x_i)$
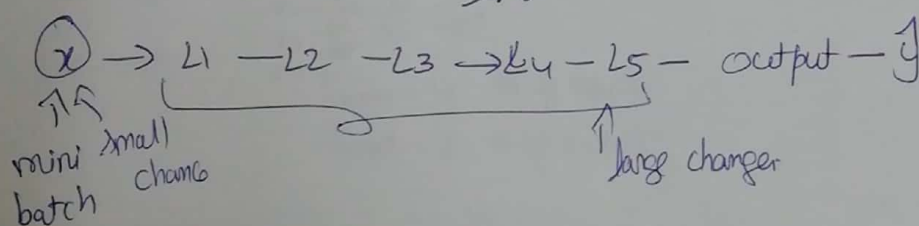   - mean centering
   - Var. scaling.

$$\tilde{x}_i = \frac{x_i - \mu}{\sigma}$$

$\mu = $ mean $(x_i)$
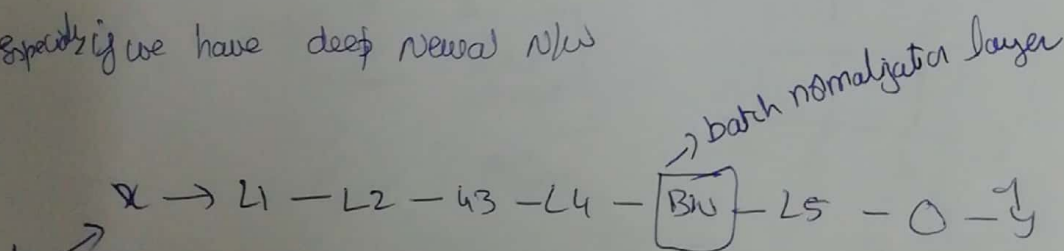
$\sigma = $ stddev $(x_i)$



i/p      $z_1$      $z_2$      → fully connected

→ fully connected MLP

$(x) \to z_1 - z_2 - z_3 \to z_4 - z_5 -$ output $- \hat{y}$

mini batch      small change      large changes

A small change in i/p will lead to a large change in the output especially if we have deep neural n/w

$x \to z_1 - z_2 - z_3 - z_4 - \boxed{BN} - z_5 - O - \hat{y}$

→ batch normalization layer

$b_1 \to$
$b_2 \to$
$b_3 \to$

$$b_1 = \frac{\{x_1, x_2 \cdots x_k\}}{\Downarrow}$$

normalization



$b_1 \quad b_2 \to b_3$

→ Batch normalization will work better when it is placed deep in the n/w.

$$B_i : \boxed{\begin{matrix} BN \\ x_i \longrightarrow \gamma,\beta \end{matrix}} \longrightarrow \tilde{x}_i \; (\beta)$$

$\mu_\beta \leftarrow \frac{1}{m} \sum\limits_{i=1}^{m} x_i$  // mini batch mean.

$\tilde{\sigma}_\beta^2 \leftarrow \frac{1}{m} \sum\limits_{i=1}^{m} (x_i - \mu_\beta)^2$  // mini batch variance.

$\hat{x}_i \leftarrow \dfrac{x_i - \mu_\beta}{\sqrt{\sigma^2 + \epsilon}}$ → small value to avoid division by zero Errors.

$\tilde{x}_i \rightarrow \gamma \hat{x}_i + \beta = BN_{\gamma,\beta}(x_i)$  // scale & shift

scale ↑   shift ↑

## Advantage

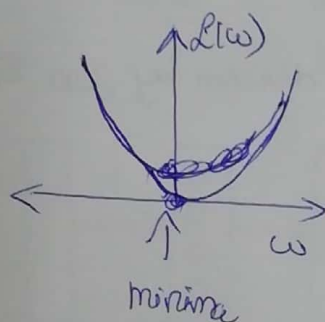→ Helps to have faster convergence

→ weak regularizer

$\boxed{BN + dropout}$

→ Internal covariate shift → we can train deep NN.

# 48.6 Optimizers: Hill-descent analogy in 2D

Ls. Reg & optimisation → GD, SGD, mini batch-SGD

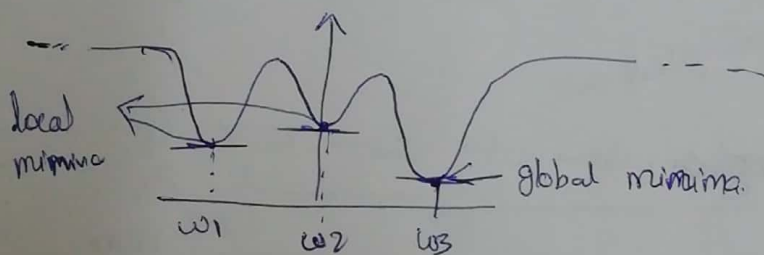$$\min_{w} \mathcal{L}(w)$$

① If $w$ is scalar



— There is no maxima

→ There is only one minima

minima



local minima

global minima.

$w_1$   $w_2$   $w_3$

At both minima/maxima $\exists\; df/dw = 0$
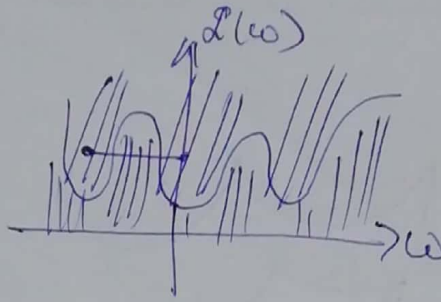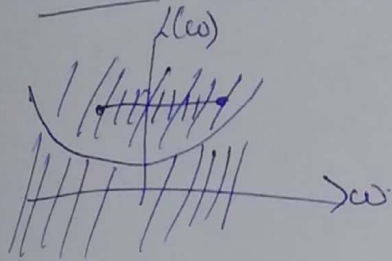
If $\dfrac{d\mathcal{L}}{dw} = 0$  we can be either at minima/maxima/Saddle point

$$w_{new} = w_{old} - \eta \frac{d\mathcal{L}}{dw}$$

We will keep updating $w$ until $\dfrac{d\mathcal{L}}{dw}$ become zero, and this can be zero at minima/maxima/Saddle point

→ SGD/mini batch SGD can get stuck at a saddle point

convex functions & non convex functions



Convex functions has only one minima & one maxima.
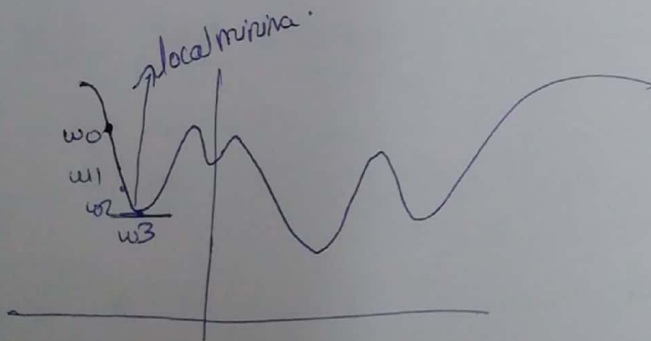
↳ local minima = Global minima

**√√Imp**

Lr. reg, lr. reg, SVM all of them the loss function can be shown as convex function. → local minima = global minima.

DL & MLP
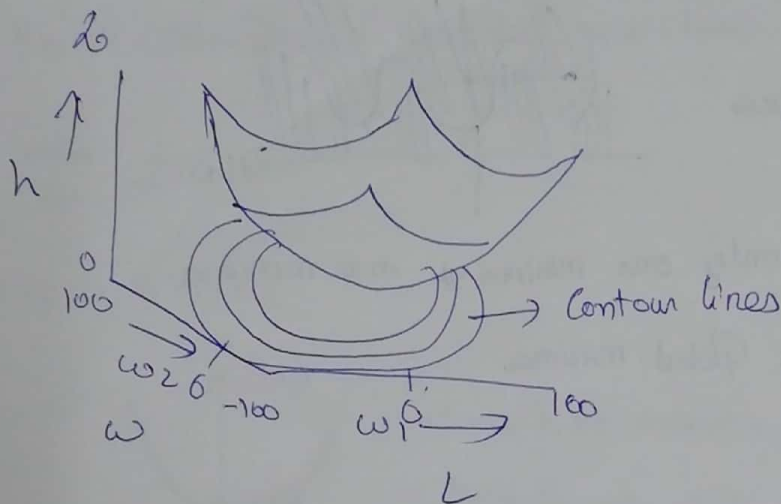
↳ non convex function → multiple local minima & can stuck at
                                                                            saddle point



→ Base on initial weights we can land up at different minima.

$\rightarrow$ Contour lines

L

$$(\omega^k_{ij})_{new} = (\omega^k_{ij})_{old} - \eta \left[\dfrac{\partial L}{\partial \omega^k_{ij}}\right]$$



update function

$$\omega \longrightarrow \omega^k_{ij}$$

$$\omega_t = \omega_{t-1} - \eta \left[\dfrac{\partial L}{\partial \omega}\right]_{\omega_{t-1}}$$

with $\uparrow^{0.01}$ above $\eta$

$$\mathcal{D} = \{x_i, y_i\}_{i=1}^{n}$$

$\dfrac{\partial L}{\partial \omega}$ 

$\rightarrow$ using all the n pts in $\mathcal{D}$ $\rightarrow$ GD

$\hookrightarrow$ using only one pt $x_i$ @ random $\rightarrow$ SGD

$\hookrightarrow$ using a random subset of $k$ pts in $\mathcal{D}$

$\hookrightarrow$ mini-batch SGD

mini batch SGD
_____

(noisy) $\rightarrow \left(\dfrac{\partial L}{\partial \omega}\right)_{min-SGD} \overset{approx}{\approx} \left(\dfrac{\partial L}{\partial \omega}\right)_{GD}$

local minima.

mini batch SGD



local minima

→ Each of the update are more noisy in SGD & MBSGD than GD

→ denoising gradient from SGD to Converge faster

## 48.9  Batch SGD with Momentum

$t=1 \quad t=2 \quad t=3 \quad t=4$ { value at specific time.

$a_1 \quad a_2 \quad a_3 \quad a_4 \quad a_5 \quad a_6 \cdots$ { number

(t=1)  $V_1 = a_1$

$t=2 \quad V_2 = \gamma V_1 + a_2 \quad (\gamma=1)$

$\qquad V_2 = V_1 + a_2 \qquad (0 \le \gamma \le 1)$

$\qquad = a_1 + a_2$

if $\gamma = 0$

$V_2 = a_2$

$\gamma = 0.5$

① $V_2 = 0.5 V_1 + a_2$

$\qquad = 0.5 a_1 + a_2$

$\qquad \uparrow$ ↑ giving more weight to the new points

low weight to old points

$V_1 = a_1$

$V_2 = \gamma V_1 + a_2$

$V_3 = \gamma V_2 + a_3$

$\qquad = \gamma(\gamma V_1 + a_2) + a_3$

$\qquad = \gamma^2 V_1 + \gamma a_2 + a_3$

if $\gamma=0.5$ ⟹ $0.25 a_1 + 0.5 a_2 + a_3$

$\qquad\qquad\qquad \uparrow \qquad \uparrow \qquad \uparrow$

$\qquad\qquad\qquad t=1 \qquad t=2 \qquad t=3$

$V_4 = 0.125 a_1 + 0.25 a_2 + 0.5 a_3 + a_4$

$$v_1 = a_1$$
$$v_t = \gamma v_{t-1} + a_t$$

⟶ Exponentially weighted averages.

⟶ $v_t \approx$ denoised Estimate

$\uparrow$ change over time.

$v_{t-1}\alpha_t + \gamma a_{t-1} + \gamma^2 a_{t-2}$
- - -

$1 \geq \gamma > \gamma^2 \geq \gamma^3 \geq \gamma^4 \gamma^5 \dots$

MB-SGD $\longrightarrow$ $w_t = w_{t-1} - \eta \left( \dfrac{\partial L}{\partial w} \right)_{w_{t-1}}$

$g_t = \left( \dfrac{\partial L}{\partial w} \right)_{w_{t-1}}$  ⟶ gradient at time $t$

$$\boxed{w_t = w_{t-1} - \eta (g_t)}$$ ⟶ mini batch SGD

$$v_t = \gamma v_{t-1} + \eta g_t$$  ⟶ using exponential weight avg.

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

Case1 $\gamma = 0$ ; $v_t = \eta g_t \Rightarrow w_t = w_{t-1} - \eta g_t$

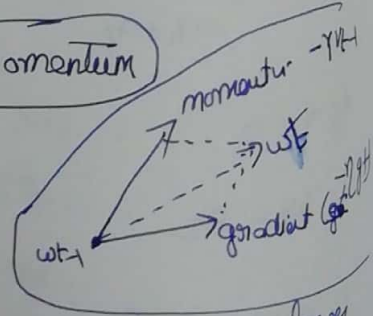$$\boxed{\text{People typically use } \boxed{\gamma = 0.9}}$$

Case 2 $\div \gamma = 0.9 \Rightarrow v_t = 0.9 v_{t-1} + \eta g_t$ ; $w_t = w_{t-1} - \left( 0.9(v_{t-1}) + \eta g_t \right)$

when we use exponential weight to denoise your SGD gradient

$\downarrow$

SGD + momentum

$$w_t = w_{t-1} - \left[ \underbrace{\gamma v_{t-1}}_{\text{Momentum}} + \underbrace{\eta g_t}_{\text{gradient}} \right]$$
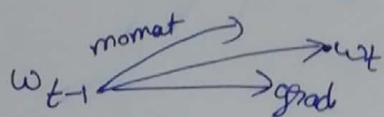
⟶ momentum will make longer steps to reach minima

$$v_t = 1(\eta g_t) + \gamma(\eta g_{t-1}) + \gamma^2 (\eta g_{t-2}) + \gamma^3 (\eta g_{t-3}) + \dots$$

SGD + momentum ⟶ Speep up Convergence

48.10    Nesterov Accelarate Gradient (NAG)



$\eta g_t$ -- gradient term

$\gamma V_{t-1}$ = momentum

## N AG

1) First Compute momentum and move in that direction, & Compute $\omega'$
and Compute gradient at $\omega'$



$$\omega_t = \omega_{t-1} - (\gamma v_t + \eta g_\phi')$$

$$\uparrow$$

$$g' = \left(\frac{\partial \mathcal{L}}{\partial \omega}\right)_{\omega'}$$

$$\omega' = \omega_{t-1} - \gamma V_{t-1}$$

$$\boxed{V_t = \gamma V_{t-1} + \eta g'}$$

$$\boxed{\omega_t = \omega_{t-1} - V_t}$$

48.11    Optimizers : Ada Grad

In SGD & SGD + Momentum → learning rate $(\eta = 0.01)$ → same for each weight

Key idea of Ada Grad

Each (weight) has a different $\eta$
↳ parameter
⤷ why?

feature → Sparse → (BOW) → few times we can see non zero
↳ dense → (W2V)

Simple SGD
→ Same for all weights

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

Adagrad

$$w_t = w_{t-1} - \left(\eta_t'\right) g_t$$

↳ different $\eta$ for each weight @ each iteration $t$
it changes

$$\eta_t' = \frac{\eta^{\nearrow 0.01}}{\sqrt{\alpha_{t-1} + \epsilon}}$$

↳ small +ve number to avoid division by '0'

$$\alpha_{t-1} = \sum_{i=1}^{t=\text{all}} g_i^{\nu} \longrightarrow \left(\frac{\partial L}{\partial w}\right)_{(w_{t-1})}$$

↳ always +ve as it is a square

Computing $g_1, g_2, g_3 \cdots g_{t-1}$

$$\alpha_t \geq \alpha_{t-1}$$

@ as $t \uparrow \rightarrow \alpha_{t-1} \uparrow \rightarrow \eta_t' \downarrow$

As iteration increases, learning rate for that weight is decreasing

It uses all the information in the previous gradients

⊕ no need to tune $\eta$ for each iterate

⊕ Sparse & dense feature

⊖ $\alpha_{t-1}$ can become very large as $t \uparrow$

Optimizers : Ada delta & RMS Prop

Adagrad : $\alpha_{t-1}$ V. large $\rightarrow$ slow convergence

$$\eta'_t = \frac{\eta \, (=0.01)}{\sqrt{\alpha_{t-1}) + \epsilon}} \quad ; \quad \alpha_{t-1} = \sum_{i=1}^{t} g_i^{\sim} \rightarrow \left(\frac{\partial L}{\partial w}\right)_{(\omega \, t-1)}$$

$\hookrightarrow$ gradient squares

## Ada delta

$$w_t = w_{t-1} - \eta'_t \, g_t$$

$$\eta'_t = \frac{\eta}{\sqrt{Eda_{t-1} + \epsilon}}$$

$\uparrow$ $\longrightarrow$ Controll the growth (It shoud grow but it should not become large)

Exponential decaying average.

$$edat_1 = \gamma \, edat_{-2} + (1-\gamma) g_{t-1}^{\sim}$$

$$\boxed{\text{Typically } \gamma = 0.95}$$

$$edat_1 = \underbrace{0.95 \, edat_{-2} + 0.05 g_{t-1}^{\sim}}_{1}$$

$$= 0.95 * \left[0.05 \, g_{t-2}^{\sim} + 0.95 \, Eda \, t_{-3}\right] + 0.05 \, g_{t-1}^{\sim}$$

$$= 0.05 \, g_{t-1}^{\sim} + (0.95 * 0.05) \, g_{t-2}^{\sim} + (0.95)^2 \, edat_{-2}$$

## Ada delta : Exp. weighted avg of $g_i^{\sim}$ instead of sum. of $g_i^{\sim}$ (Adagrad) so as to avoid large denominator in $\eta'_t$

$\downarrow$

Some avoid slow convergence.

48.13 : __Adam__ (Adaptive moment Estimation) (Best algorithm)

idea ÷ Adadelta $\rightarrow$ storing Exp. weight. avg of $g_t^\gamma$ $\rightarrow$ learning rate $(n't)$

└ what if we store Eda of $g_t$

In statistics

mean $\rightarrow$ 1st order momentum

Var $\rightarrow$ 2nd order momentum

$$\underset{\substack{? \\ \text{mean at,} \\ \text{time } t}}{m_t} = \overset{\rightarrow 0.9}{\beta_1} m_{t-1} + (1-\beta_1)g_t \qquad 0 \le \beta_1 \le 1$$
$$\rightarrow \text{①}$$

$$V_t = \underset{\substack{\nearrow \\ 0.99}}{\beta_2} m_{t-2} + (1-\beta_2) g_t^\vee \qquad 0 \le \beta_2 \le 1$$
$$\rightarrow \text{②}$$

$$m'_t = \frac{m_t^\bullet}{1-(\beta_1)^t} \quad ; \quad v'_t = \frac{v_t^\bullet}{1-(\beta_2)^t} \qquad \Big\} \text{ Bias correction.}$$

$$\boxed{ w_t = w_{t-1} - \frac{\alpha \, m'_t}{\sqrt{v'_t + \epsilon}} }$$
$$\nearrow \text{①}$$

If $\beta_1 = 0 \Rightarrow$ Adadelta

If $\beta_1 = \beta_2 = 0 \Rightarrow g_t^\vee$ will remain.

## 48.14 Which algorithm to choose when ?

$$\underset{\underset{\text{Small/shallow}}{\underset{\text{NN}}{\underbrace{\text{MB-SGD}}}} \to \underset{\underset{\text{works well}}{\underset{\text{in most}}{\underset{\text{cases}}{\underset{\text{(Slower)}}{\downarrow}}}}}{\overset{\text{Momentum}}{\text{NAG}}} \to \underset{}{\text{Adagrad}} \to \underset{\underset{(\text{Sparse})}{\downarrow}}{\text{Adadelta}} \to \underset{\underset{\underset{\text{fastest}}{\text{most.favorite}}}{\downarrow}}{\text{Adam}}}$$

## 48.15 Gradient monitoring & clipping
↳ update weights

$x_i \to$  $\to \hat{y_i}$

→ monitor gradients & updates
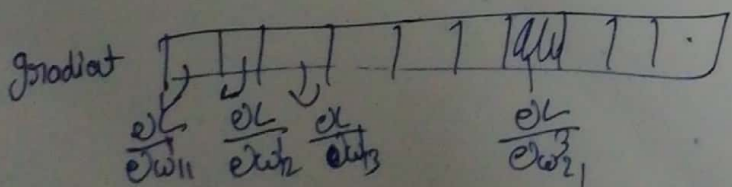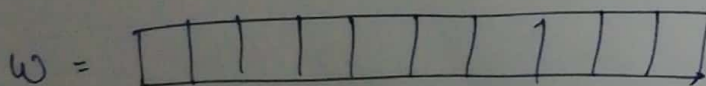
→ Each epoch

→ layer

→ If we found vanishing Gradient problem, immediately we should change the activation function to ReLu

→ **Exploding gradient**

**soln** clipping

$w = $ [ grid of cells ]

gradient [ grid of cells ]

$\dfrac{\partial L}{\partial w_{11}} \quad \dfrac{\partial L}{\partial w_{12}} \quad \dfrac{\partial L}{\partial w_{13}} \qquad \dfrac{\partial L}{\partial w_{21}^{3}}$

L2 n8m clipping

$$G_{new} = \frac{G}{||G||_2} * \tau$$

$\uparrow$ threshold



$$G = \begin{array}{|c|c|c|c|c|} \hline G_1 & G_2 & G_3 & G_4 & G_5 \\ \hline 2 & 5 & 10 & 100 & 05 \\ \hline \end{array}$$

$$||G||_2 = \sqrt{G_1^v + G_2^v + G_3^v + \cdots}$$

$$\frac{G}{||G||} = \begin{array}{|c|c|c|c|c|} \hline <1 & <1 & <1 & <1 & <1 \\ \hline \end{array} * \tau \qquad (\tau = 2)$$

$\underbrace{\qquad\qquad\qquad}$

$< 1$

---

48.16 : Softmax and Cross Entropy for multi class classification

Softmax - classifier

$\underset{\downarrow}{\text{logistic reg}} \rightarrow$ binary classification

$\underset{\downarrow}{\text{multiclass}} \rightarrow$ one vs Rest

$\{$ Log. reg + multiclass $\} =$ Softmax.



$$z = w^T x_i$$

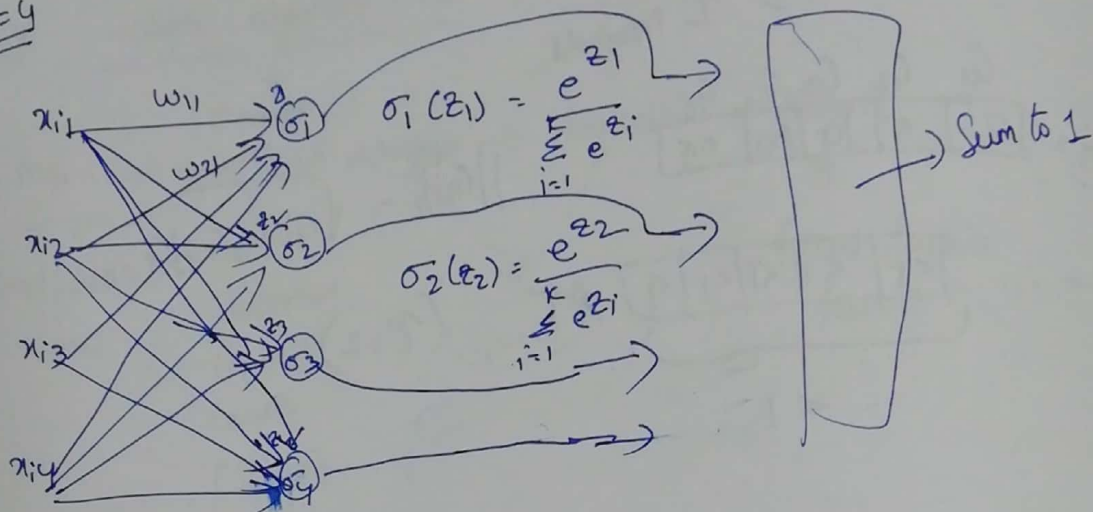$$P(y = 1|x_i) = \hat{y_i} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

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

Soft max

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

$$y_i \in \{1, 2, 3, 4, 5, 6, 7, 8, 9 \cdots K\}$$



$x_i \rightarrow \boxed{m} \rightarrow P(y_i = 1|x_i)$

$P(y_i = 2|x_i)$

$\vdots$

$P(y_i = K|x_i)$

$\}$ Sum to 1

$K = y$



$$\sigma_1(z_1) = \frac{e^{z_1}}{\sum\limits_{i=1} e^{z_i}}$$

$$\sigma_2(z_2) = \frac{e^{z_2}}{\sum\limits_{i=1}^{K} e^{z_i}}$$

$\rightarrow$ Sum to 1

$$z_1 = \sum_{j=1}^{d} x_{ij} w_{j1} + b, \quad z_2 = \sum_{j=1}^{d} x_{ij} w_{j2} + b$$

## Softmax

$\hookrightarrow$ generalizatio to LR to "multi class" setting.

$\rightarrow$ minimize binary class log loss.

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

$$\text{Softmax}: \sigma_1(z_1) = \frac{e^{z_1}}{\sum\limits_{i=1}^{K} e^{z_i}} \quad ; \quad (\sigma_1, \sigma_2 \dots \sigma_K)$$

minimize multiclass log loss

$\hookrightarrow$ cross entropy

## Ex

no pts

K class

$$-\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{K} y_{ij} \log p_{ij}$$

$$\begin{cases} 1 & \text{if } y_i \in j \\ 0 & \text{otherwise} \end{cases}$$

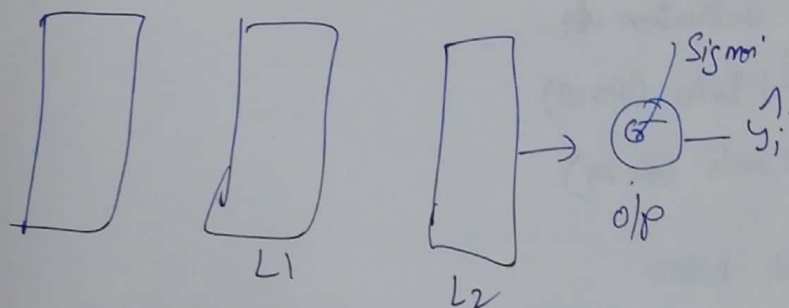$\hookrightarrow P(y_i \in j | x_i)$

### 2 class

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

$x_i, y_i = 3$



|   | 1 | 2 | 3 | 4 | ... | k |
|---|---|---|---|---|-----|---|
|   | 0 | 0 | 1 | 0 | ... | 0 |

$y_{i1}$ $y_{i2}$ $y_{i3}$ $y_{i4}$ ... $y_{ik}$

$x_i \rightarrow \boxed{m}$ ⟶ $P(y_i=1|x_i) = 0.2$
⟶ $P(y_i=2|x_i) = 0.1$
⟶ $P(y_i=3|x_i) = 0.7$

2-class



$L_1$   $L_2$   Sigmoid
o/p   $\hat{y}_i$

k-class



$x_i \rightarrow L_1 \rightarrow L_2 \rightarrow$ soft max $\rightarrow$
$\begin{cases} P(y_i=1|x_i) \\ \vdots \\ P(y_i=k|x_i) \end{cases}$

Vector of size k

Regression



$L_1$   $L_2$   Nr.num
i/p   o/p
$w^Tx$   $w^Tx$   $\hat{y}_i$

How to train a deep MLP?

1) Preprocess data
 ↳ Data normalization.

2) weight Init
 ↳ Xavier/Glorot → Sigmoid/tanh
 ↳ He → ReLu
 ↳ Gaussian (small $\sigma$)

3) Choose activation fn
 ↳ ReLu (2018)
 ↳ Selu (2017)

4) Batch Norm
 ↳ Esp for deep MLP (later layers)
 dropout → Ⓟ

5) Optimizer
 ↳ Adam (2018) (Adaptive fast)

6) hyper parameters
 ↳ Architecture ⇒ # layers
        # neurons
 ↳ drop out
 ↳ Adam: $\beta_1$, $\beta_2$, $\alpha$.

7) Loss function : 2-class classificat → log loss
      k-class classifical → multiclass LL
      regression → square loss

8) Monitor your Gradient
 ↳ Gradient clipping (if needed)

9) plots

Train
loss



epoch

10) avoid overfitting
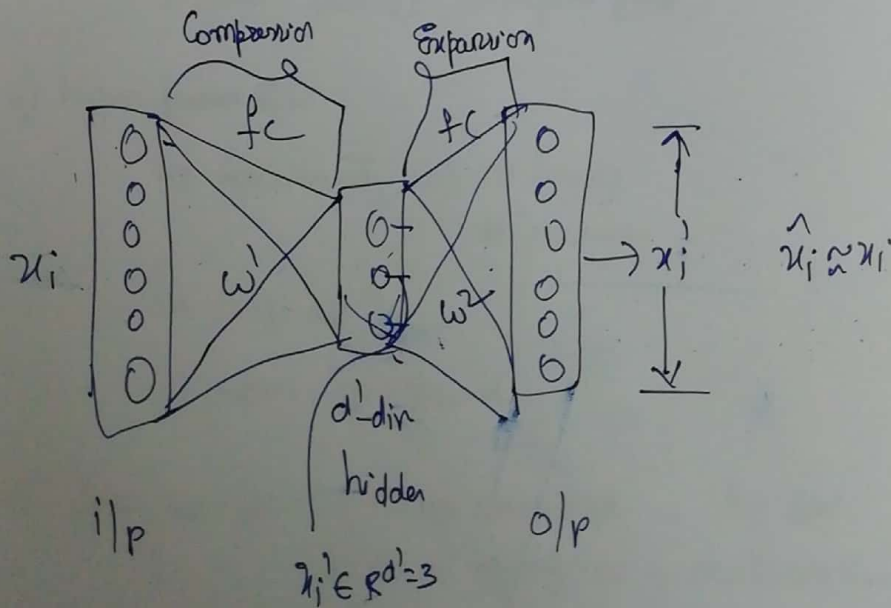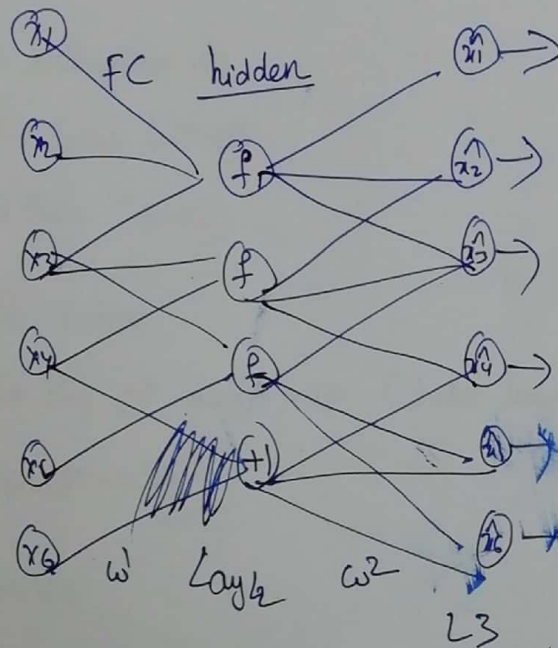{train & test losses}

48.18  Auto Encoders

→ It is a NN, which performs dimensionality reduction [PCA, TSNE]
    → weighting
    → variance

→ Some times better than PCA, TSNE

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

$D' = \{x_i'\}_{i=1}^n$    $x_i' \in \mathbb{R}^{d'}$    $\boxed{d' < d}$

Ex



Compression    Expansion



$x_i$    $\omega^1$    $\omega^2$    $\to x_i'$    $\hat{x}_i \simeq x_i$

i/p    d'-dim hidden    o/p

$x_i' \in \mathbb{R}^{d'=3}$

$$\mathcal{L}(x_i, \hat{x}_i) = \|x_i - \hat{x}_i\|$$

$\boxed{\mathcal{L} = 0}$  → distance loss

# Denoising Auto Encoder :

$$\mathcal{D} = \{x_1, x_2, \dots x_n\} \leftarrow \text{actual data}$$

target $\rightarrow$ $\tilde{\mathcal{D}} = \{\tilde{x}_1, \tilde{x}_2, \tilde{x}_3 \dots \tilde{x}_n\} \rightarrow$ noisy & corrupted
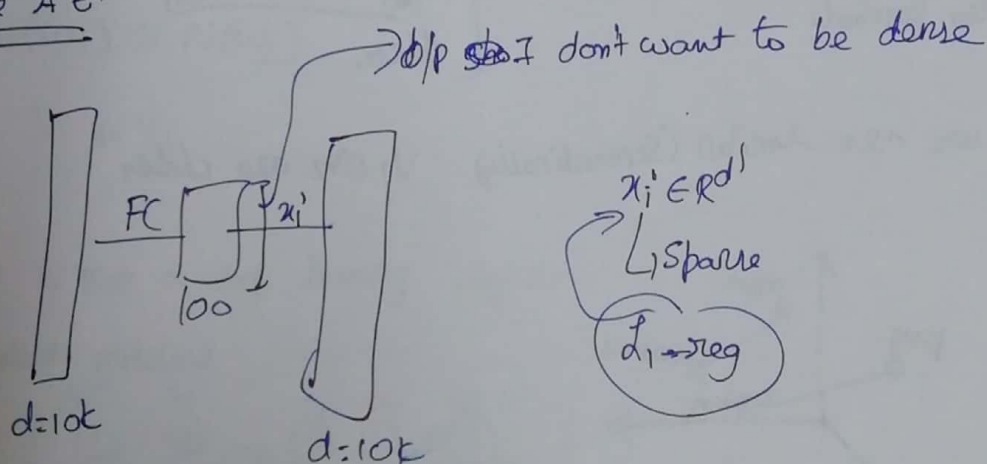
$\rightarrow$ AE

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

$$\boxed{x_i \rightarrow \tilde{x}_i = x_i + N(0, \sigma)^2}$$

## Sparse A E

$\rightarrow$ o/p set I don't want to be dense



FC
100

d=10k

d=10k

$x_i \in \mathbb{R}^{d'}$
$\rightarrow$ Sparse
$l_1 \rightarrow$ reg

$\rightarrow$ If all linear activation are used & only a single sigmoid hidden layer, then the optimal solution to an autoEncoder is strongly related to principal component analysis (PCA)

48.19   W8d2vec: CBOW

→ give it a w8d it return a vect8r

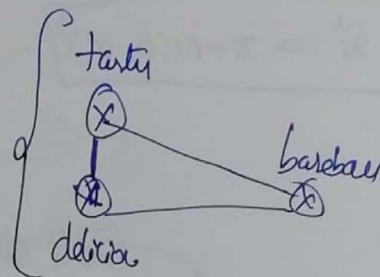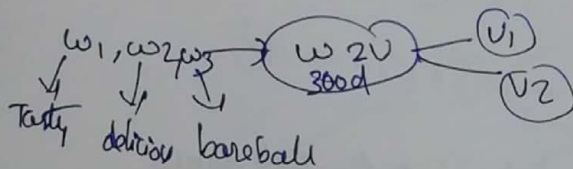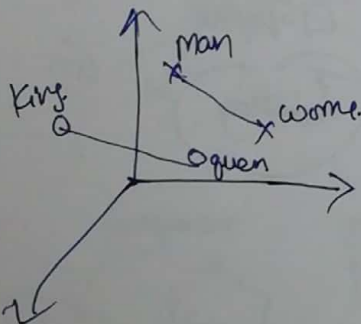→ Semantic meaning of w8d's into consideration.

✆ w8d → vect8r

w8d → ▢ → d-dime vect
           → not a spance veet

d typically

$(50, 100, 200, 300)$
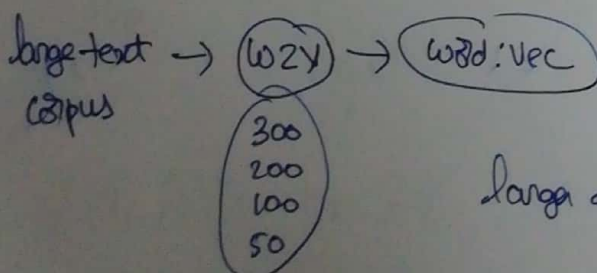
$w_1, w_2, w_3$ → $\boxed{W\ 2V}$ $\underset{300d}{}$ → $\begin{cases} (v_1) \\ (v_2) \end{cases}$

↓    ↓    ↓
Tasty delicion bareball



w1 & w2 are similar (Semantically)    $v_1$ & $v_2$ are closer



$(V_{man} - V_{woman}) \parallel (V_{king} - V_{queen})$

(W2V) → learning relationships automatically from raw-text

large-text → (W2V) → w8d:vec
corpus

$\boxed{\begin{matrix} 300 \\ 200 \\ 100 \\ 50 \end{matrix}}$      larger dimension → more info such the vect8r b
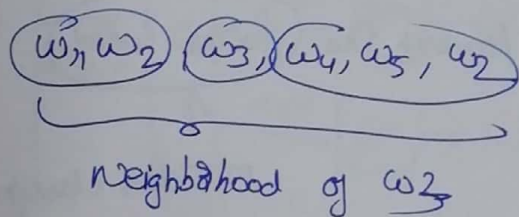
data corpus ↑ → d↑

Google -news $\rightarrow$ W2V $\rightarrow$ (300-dim)

review
text $\rightarrow$ W2V $\rightarrow$ word:vec
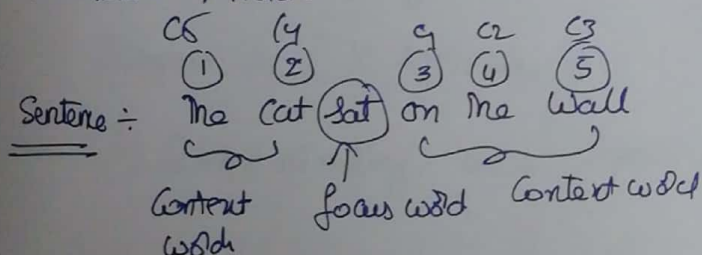500k

**core ÷ W2V**

will look at the sequence of words,

$(\omega_1 \omega_2) (\omega_3 )(\omega_4, \omega_5, \omega_2)$

neighbourhood of $\omega_3$

$W2V (\omega_3) =$

$N(\omega_i) \approx N(\omega_j)$

$V_i \approx V_j$

→ W2V is not a deep learning algorithm.

→ mixtav mikolav

Sentence ÷

| c5 | c4 | | c1 | c2 | c3 |
|----|----|----|----|----|----|
| ① | ② | | ③ | ④ | ⑤ |
| The | cat | (sat) | on | the | wall |

Content     focus word     Context word
word

→ Context words are very useful in understanding the focus word. and
vice versa

2-algo $<$ CBOW (Continuous Bag of words)
          skipgram