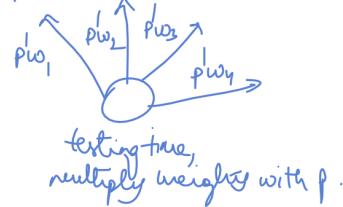
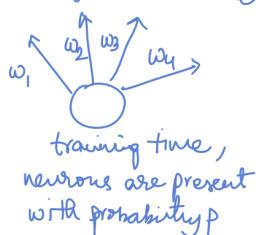


Deep Multi Layered Perceptrons : 1980's to 2010s (2.1) :-

- Upto 2010 people were using NN with 2-3 layers. Because of Vanishing Gradient & lack of data. Also no computing power.
- From 2010, we started having lots of labelled data. We also started having GPUs. New ideas & algs as well.
- Approach was also changed. Instead of going theory first, experiments were performed more which took lot less time than coming up with new theory.

Dropout Layers & Regularization (2.2) :-

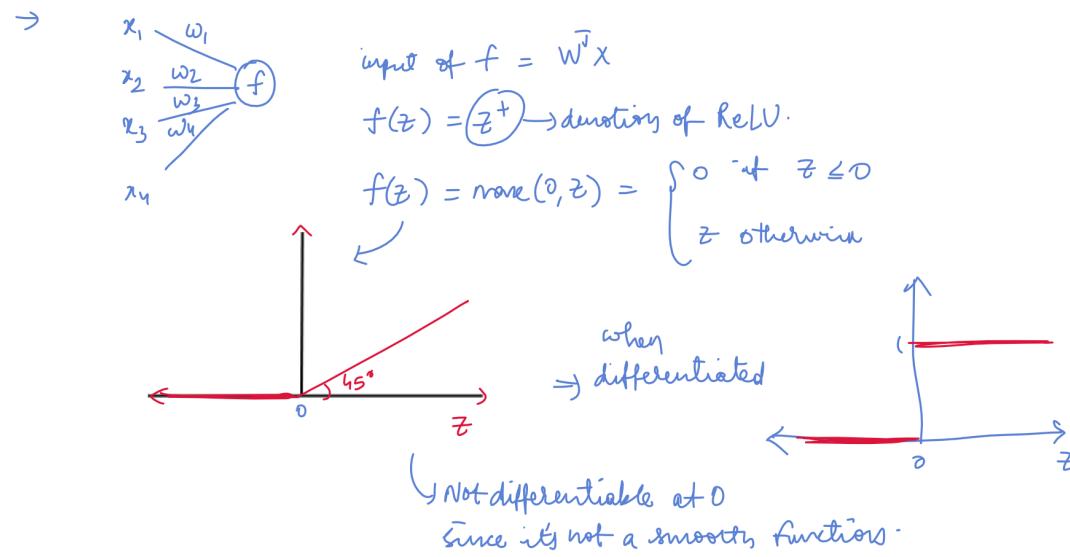
- one of the major issues of Deep NN is overfitting. We use regularization to overcome this.
- In Random Forest, we get many trees. But we only select subset of trees & subset of features. We use randomization of features to create regularization.
- Core idea of dropout layer is to use this philosophy for MLP's
- In each iterations randomly remove some neurons. As a result, some incoming & outgoing connections will be lost. dropout rate $0 \leq p \leq 1$. At any iteration, in each layer $p \neq 1$. & neurons shall be inactive.
- Dropout is very similar to having random subset of features in Random Forest -
- During testing we multiply w with p



- When we are overfitting with many weights & few data points, dropout rate should be large.
- Dropouts can also be thought of as a layer that takes inputs from one layer but does not send all of them to the next layer.

Rectified Linear Unit (ReLU) (2.3) :-

- One of the main problems with classical NN is vanishing gradient problem.
- ReLU is the most common activation functions for MLP these days.



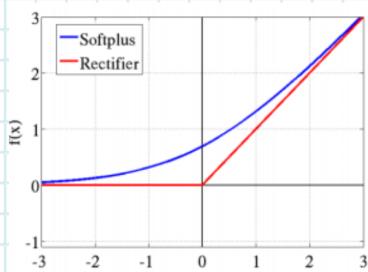
→ $\frac{\partial f_{\text{ReLU}}}{\partial z} \neq 0$ \Rightarrow As a result no vanishing gradient, no exploding gradient -

But this 0 could cause dead activations.

→ ReLU converges faster than tanh/sigmoid because there's no vanishing gradient/exploding gradient.

→ ReLU is sometimes also called simply as rectifiers.

→ softplus function is like an approximation of ReLU



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

$$\frac{\partial f}{\partial x} = \frac{\exp(x)}{(1 + \exp(x))} \rightarrow \text{logistic function.}$$

→ Variation of ReLU:-

Noisy ReLUs :- $f(x) = \max(0, x + \gamma)$, $\gamma \sim N(0, \sigma(x))$

Leaky ReLUs :- $z = w^T x$, if weights are -ve,
then $\gamma = -ve$.

if $\gamma = -ve$, $\frac{\partial f}{\partial z} = 0$.

→ the product of this derivative with others in chain rule = 0

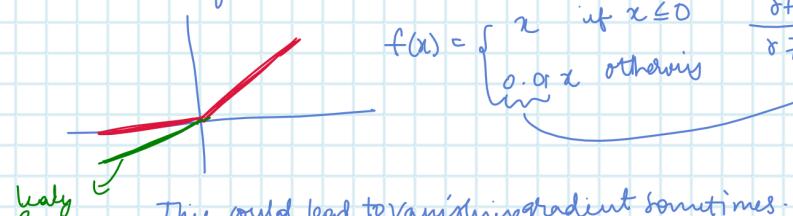
$$\rightarrow w_{new} = w_{old} - \eta \underbrace{\left(\frac{\partial f}{\partial w} \right)}_0$$

→ $w_{new} = w_{old}$. - weights are not updated.

This state is called dead activation state.

We must always keep track of this when using ReLU

Leaky ReLU fixes this



This could lead to vanishing gradient sometimes.

$$f(x) = \begin{cases} x & \text{if } x \leq 0 \\ \alpha x & \text{otherwise} \end{cases}$$
$$\frac{\partial f}{\partial z} = \begin{cases} 1 & \text{if } z \text{ is +ve} \\ \alpha & \text{if } z \text{ is -ve} \end{cases}$$

its the hyperparameter typical $0 \leq \alpha \leq 1$

Weight Initialization (2/4):-

→ In SGD, we initialize randomly (uniform or Gaussian) i.e. $w_{ij} \sim N(0, \sigma)$

→ Problem of Symmetry :- when all weights are initialized to 0/1/- this is bad.

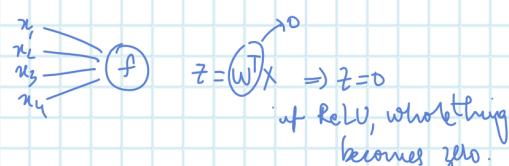
if f is also same, then

(i) All neurons compute the same thing.

(ii) Same gradient update

This is called problem of symmetry.

We want stuff to be asymmetric because
we want each neuron of layer to learn something different.



→ if we initialize to large -ve values, then $f(z) = 0$ if ReLU is used. if Sigmoid or tanh is used, leads to vanishing gradient problem.

Idea ① :- → weight should be small

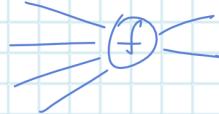
→ not all zero

→ good variance

We can do $w_{ij} \sim N(0, \sigma)$ σ is small.

This is called Gaussian/Normal initialization

As a result we get good σ and there will be both +ve and -ve weights.



for f, fan-in = 4
fan-out = 2

→ After lots of experiments, people concluded, weights should be based on fan-in, fan-out.

Idea 2 :- Uniform Initialization :-

$$w_{ij}^k \sim \text{Uniform dist} \left[\frac{1}{\sqrt{\text{fan-in}}}, \frac{1}{\sqrt{\text{fan-out}}} \right] \text{ for any neuron}$$

This works well for sigmoid functions.

Idea 3 :- Xavier / Glorot :-

two variations.

(a) Xavier / Glorot Normal :-

$$w_{ij}^k \sim N(0, \sigma_{ij}) \quad \sigma_{ij} = \sqrt{\frac{2}{\text{fan-in} + \text{fan-out}}} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{ we use fanin + fanout in this. In previous one we only used fan-in.}$$

(b) Xavier / Glorot Uniform :-

$$w_{ij}^k \sim U\left[\frac{-\sqrt{6}}{\sqrt{\text{fan-in} + \text{fan-out}}}, \frac{+\sqrt{6}}{\sqrt{\text{fan-in} + \text{fan-out}}}\right]$$

Works well for sigmoid too.

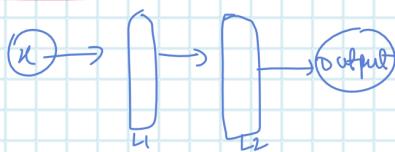
Idea 4 :- He-init :-

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

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

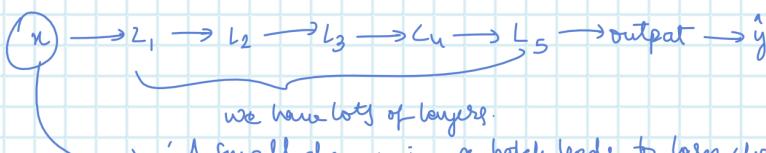
Work well for ReLU.

Batch Normalization (2-5) :



$$\hat{x}_i = \frac{x_i - \mu}{\sigma} \quad \begin{array}{l} \sigma = \text{std-dev} \\ \mu = \text{mean} \end{array}$$

→ In minibatch SGD, we only send batches of points



∴ A small change in x batch leads to big change in w.

batch 1 } → L1 won't see much difference.
batch 2 } But L5 will because of chain rule before it.
⋮ } This is called internal covariance shift.

→ In Batch Normalization, we add a new layer that normalizes the data of that batch only.

→ γ & β are to be learned in backpropagation

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad \boxed{\hat{x}_i = \gamma \hat{x}_i + \beta} \quad \text{small value.}$$

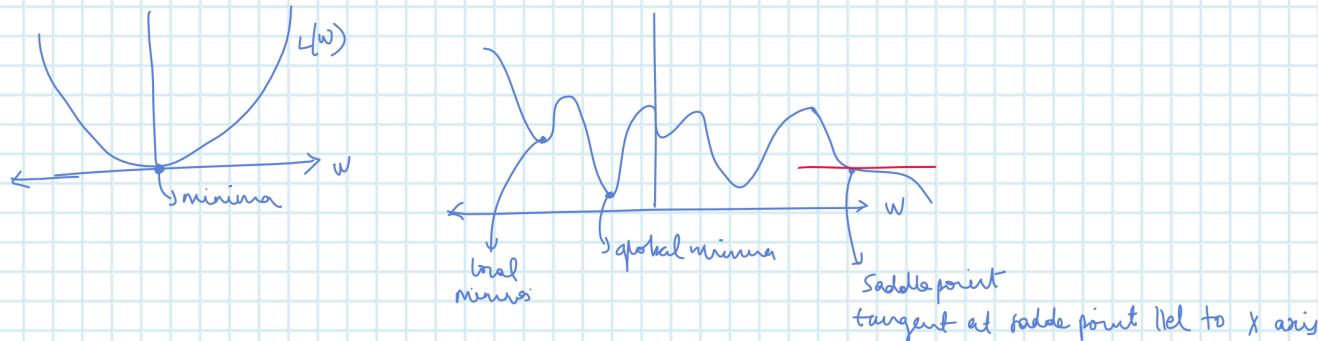
→ batch Normalization is differentiable.

→ During testing, we take the mean & variance of the whole dataset. □

→ Generally, if we use higher learning rate with no normalization, it tends to overshoot minima & convergence is much slower.

Optimizers: Hill Descent analogy in 2D :- (3.6)

→ Simple SGD in deep learning cannot solve the optimization functions properly.



Simple SGD would get stuck at minima/maxima or saddle point.

→ Convex functions :-

shortest dist b/w any two points of one regions i.e. the same regions its a convex functions

don't lie in the same regions - !. Not convex.

→ Convex functions only have one minimum & 1 maxima

local minima = global minima -

non convex have multiple local minima & 1 global minima.

→ In log reg, lin reg & SVM: all are convex funcns.

→ In deep learning: non convex functions i.e., we have local minima, saddle points etc;
Depending on w initializations we could end up with local minima -

rajab_chaudhary

47 Votes

Hi, the way I see this, it is upto us to choose which ever loss function we can choose for our algorithm. So that one of the reason we use squared loss in linear regression, as the squared loss is a convex function, hence no issue of getting stuck at local minima or local maxima or a saddle point. just wanted to confirm my understanding here.

Reply



Mar 03, 2019 08:20 AM

Applied_AI

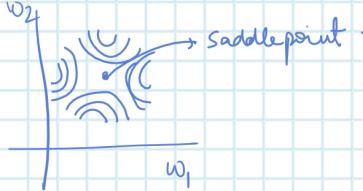
Good question especially for the interviews. You are asked why to use the squared loss and power 4 or modulus. One of the major reason is squared loss upon derivative produces a single value of parameter set and hence gives one unique solution while others don't. However it is not true for one value of the parameter, the loss would be minimal. Loss can be minimal for multiple values also. Hence you can't tell that squared loss would only give you the least solution. Just because you have multiple minima doesn't mean that they need to be different in values. Honestly in Square and power 4 both generates minimal solution but the squared loss gives one solution while power 4 gives 3 solutions.

Reply



Mar 03, 2019 10:36 AM

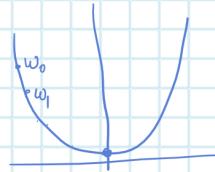
contour plot of a saddle point :-



Simple mini batch SGD will get stuck at saddle point.

SGD Recap (2.8) :-

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



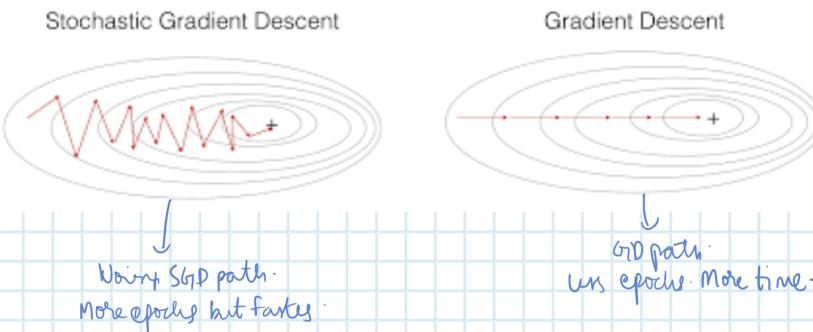
$$w_t = w_{t-1} - \eta \left[\frac{\partial L}{\partial w} \right]_{w_{t-1}}$$

$\rightarrow \frac{\partial L}{\partial w}$ when calculated using all points \rightarrow gradient descent (computationally expensive)
using one random point in each iteration \rightarrow SGD
using a random subset of points each time \rightarrow MiniBatch SGD (most common in DL)

Minibatch SGD is erroneous, noting it an approximation of GD

i.e.) SGD with lots of iterations \approx GD

We are estimating the gradient in SGD. So the updates in each iterations are more noisy in SGD than in GD.



Batch SGD with Momentum (2.9) :-

\rightarrow Assume we are calculating values

at $t=1 \quad 2 \quad 3$

$$\begin{matrix} a_1 & a_2 & a_3 \\ v_1 & v_2 & v_3 \end{matrix}$$

Let γ be a value b/w 0 & 1

$$v_1 = a_1$$

$$v_2 = \gamma v_1 + a_2$$

$$v_3 = \gamma v_2 + a_3 = \gamma(\gamma v_1 + a_2) + a_3$$

if $\gamma = 1$

$$v_1 = a_1$$

$$v_2 = a_1 + a_2$$

$$v_3 = a_1 + a_2 + a_3$$

$\gamma = 0$

$$v_1 = a_1$$

$$v_2 = a_2$$

$$v_3 = a_3$$

$\gamma = 0.5$

$$v_1 = a_1$$

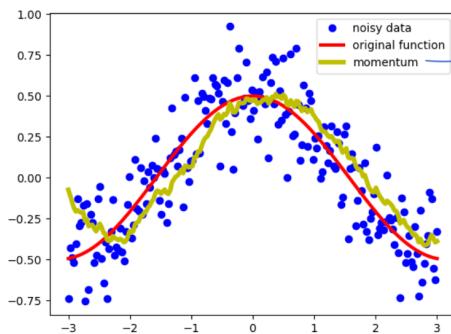
$$v_2 = 0.5a_1 + a_2$$

$$v_3 = 0.25a_1 + 0.5a_2 + a_3$$

most weightage to a_3
less to a_2
lower to $a_1 \dots$

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

v_t at any time is called a derived value. These are called exponentially weighted averages.



→ obtained after exponentially weighted average. Use approximation to the original function.

→ taking this concept & applying it to mini batch SGD :-

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

$$\text{Let } g_t = \left(\frac{\partial L}{\partial w} \right)_{w_{t-1}} \Rightarrow w_t = w_{t-1} - \eta g_t$$

Applying exponential weighting.

$$v_t = (\gamma v_{t-1}) + \eta g_t \quad \text{momentum}$$

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

In practice typical value of $\gamma = 0.9$

dealing with exponential weighting = SGD + momentum

(With momentum, SGD starts taking longer steps towards minima.)

∴ SGD + momentum speeds up convergence.

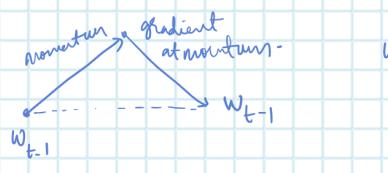
→ As we are moving closer to minima SGD point in the best direction. Therefore we use exponential average & not simple average.

Nesterov Accelerated Gradient (2.10) :-

$$\text{in SGD + momentum } \hat{g}_t = \left[\frac{\partial L}{\partial w} \right]_{w_{t-1}}$$



But in NAG:-



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

$$g' = \left(\frac{\partial L}{\partial w} \right)_{w_{t-1} + \gamma v_{t-1}} \quad \text{= gradient calculated at momentum.}$$

→ if we can skip minia in an iterations, NAG will take fewer steps till convergence than SGD with momentum.

Adagrad (2.11) :-

→ In SGD & SGD+Momentum, learning rate (η) is constant.

feature → sparse } same learning rate for both is bad.
→ dense }

→ Adagrad = Adaptive Learning rate.

$$\text{SGD} : w_t = w_{t-1} - \eta g_t \quad \text{same for all weights.}$$

Adagrad :- $w_t = w_{t-1} - \eta'_t g_t$ η'_t is different for each weight at each iteration

$$\eta'_t = \frac{\eta}{\sqrt{\alpha_t + \epsilon}} \quad \epsilon \text{ is a small value to prevent division error.}$$

η is a normal learning rate value.

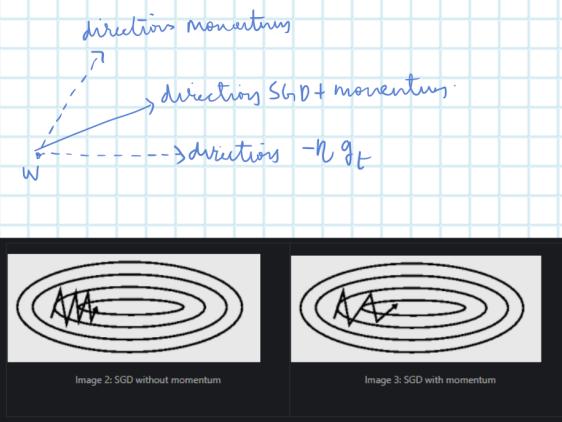


Image 2: SGD without momentum

Image 3: SGD with momentum

$$\alpha_{t+1} = \sum_{j=1}^{t+1} \left(\frac{g_j^2}{\gamma_j} \right) \rightarrow \text{always true} \Rightarrow \alpha \text{ is always +ve}$$

$$\alpha_t > \alpha_{t-1} > \alpha_{t-2} \dots$$

$$\alpha_t \uparrow, \alpha_t \uparrow, \eta_t \downarrow$$

as iterations increase, learning rate decreases adaptively -

Advantages :-

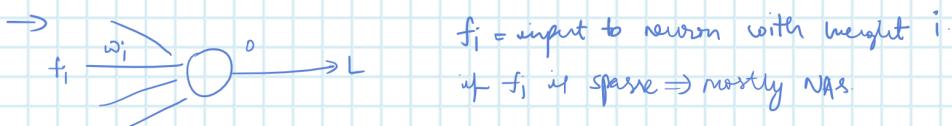
→ No need of manually tuning η .

→ if there are some sparse & some dense features, adagrad automatically chooses the best learning rate.

Disadvantage :-

→ As $t \uparrow$, α_t would become very large & η_t becomes very small $\Rightarrow \eta_t \approx \alpha_{t-1}$ leads to slower convergence.

→ to avoid getting stuck at saddle points, add small value ϵ to optimization problem.



in minibatch SGD, lot of batches contains mostly NAs & only few valid values.

∴ let loss be squared loss

$$\Rightarrow \left(\frac{\partial L}{\partial w} \right)_{w_{t-1}} = \sum_{i=1}^n (-2x_i)(y_i - w_{t-1}x_i)$$

↳ will be 0 for sparse features.

$$\Rightarrow \left(\frac{\partial L}{\partial w} \right)_{w_{t-1}} \text{ will be much smaller for sparse than dense features.}$$

∴ Adagrad will help with this.

Optimizers : Adadelta and RMSprop (2.12) :-

→ In adagrad :- if α_{t+1} becomes very large, η_t becomes low, convergence becomes slower.

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

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

→ In Adadelta, we replace α with exponentially decaying average so as to avoid small η_t values & avoid slow convergence.

$$\Rightarrow \eta_t = \frac{1}{\sqrt{cda_t + \epsilon}}$$

$$cda_t = \sqrt{cda_{t-1} + (1-\gamma)g_t^2}$$

$$\begin{aligned} \text{if } \gamma = 0.95, \text{ then } cda_t &= 0.95 cda_{t-1} + 0.05 g_t^2 = 0.95 g_{t-1}^2 + [0.95 \{0.05 g_{t-1}^2 + 0.95 cda_{t-2}\}] \\ &= 0.95 g_{t-1}^2 + (0.95 \cdot 0.05) g_{t-1}^2 + (0.95)^2 cda_{t-2}^2 \end{aligned}$$

We are controlling the growth of α

Adam (2.13) :-

→ In Adadelta, we are taking the cda of g_t^2 .

In Adam, we take cda of g_t .

Adam = Adaptive Moment Estimation.

In Statistics : mean is aka 1st order moment (Adam more like 1st order)

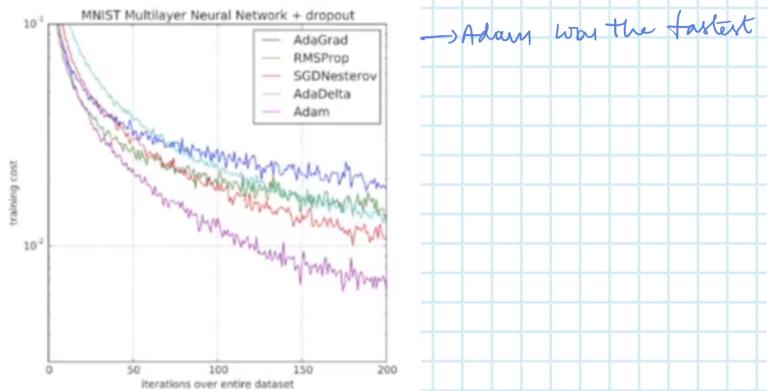
Variance is aka 2nd order moment (Adadelta more like 2nd order)

$$\begin{aligned}
 \text{mean at time } t, \quad \hat{m}_t &= \beta_1 m_{t-1} + (1-\beta_1) q_t \quad 0 \leq \beta_1 \leq 1 \quad \rightarrow \text{eda } q_t \\
 v_t &= \beta_2 v_{t-1} + (1-\beta_2) q_t^2 \quad 0 \leq \beta_2 \leq 1 \quad \rightarrow \text{eda } q_t^2 \\
 \hat{m}_t &= \frac{m_t}{1-(\beta_1)_t} \quad \hat{v}_t = \frac{v_t}{1-(\beta_2)_t} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{ very similar to} \\
 w_t &= w_{t-1} - \alpha \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}
 \end{aligned}$$

- β_1 is typically 0.9
- β_2 is typically 0.99
- α is typically 0.001

if $\beta_1 = \beta_2 = 0 \Rightarrow$ Adoletta
 $\beta_1 = \beta_2 \neq 0 \Rightarrow$ not exactly Adoletta but similar

\rightarrow Adam has all the advantages as that of AdaDelta

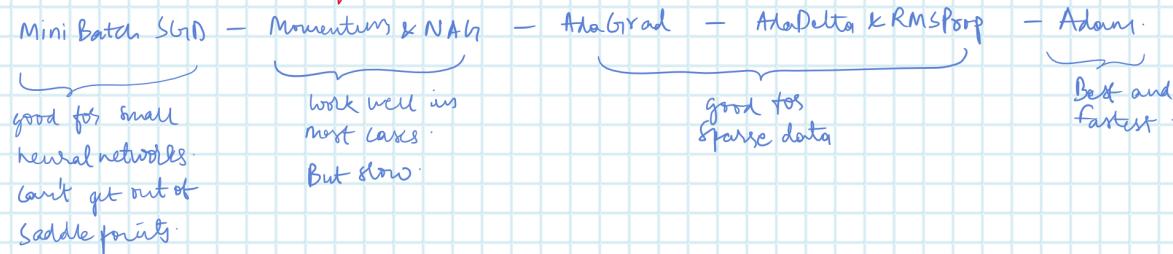


→ Adam has the problem of overshooting the minima

$$\text{NADAM} = \text{NAG} + \text{Adam} + \text{RMSProp}$$

→ As of 2021 AdaBound is Adam Variant that employs dynamic bounds on learning rates to achieve gradual & smooth transitions to SGD, is faster.

When to choose which algorithm (21h)!



Gradient checking & Clipping (2.15):-

→ good habit to monitor weight & gradients for each epoch, weight & losses.

→ helps detect vanishing gradient/exploding gradient

→ Solution to exploding gradient = clipping:

$$w = \boxed{ \dots - } \quad w_{in}$$

$$G_1 = \frac{\partial L}{\partial \dot{L}} \quad \frac{\partial L}{\partial \dot{L}_1}$$

$$L_2 \text{ Normal Clipping: } g_{\text{new}} = \frac{g}{\|g\|_2} * Z \quad Z = \text{threshold}$$

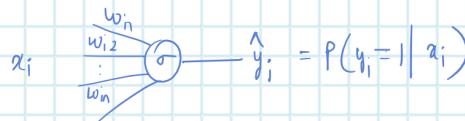
All values of G_{new} will be < 2

Softmax and Cross Entropy for Multi-Class Classification (2.16) :-

→ Logistic Regression on its own works only for binary classification
for multi-class → one v/s rest.

Instead of one v/s rest : log reg for multi-class = softmax.

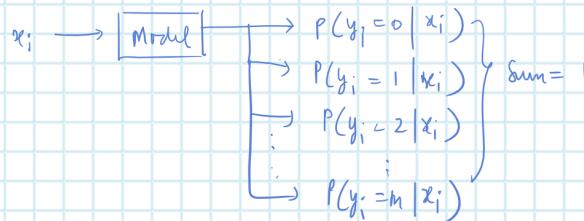
→ Logistic Regression from neural net perspective



$$z = \mathbf{x}_i^T \mathbf{w}$$

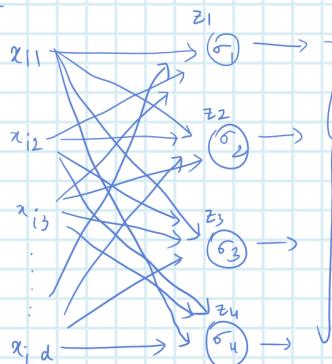
$$P(y_i=1|x_i) = \hat{y}_i = \sigma(z) = \frac{1}{1+e^{-z}} = \frac{e^z}{e^z + 1}$$

→ $D = \{(x_i, y_i)\} \quad y_i \in \{0, 1, 2, \dots, m\}$



soft-mode :-

let $m=4$



$\underbrace{z_1, z_2, \dots, z_m}_{\text{inputs to neurons}}$

$$\sigma_1(z_1) = \frac{e^{z_1}}{\sum_{i=1}^m e^{z_i}}$$

$$\sigma_2(z_2) = \frac{e^{z_2}}{\sum_{i=1}^m e^{z_i}}$$

$$\sigma_1(z_1) + \sigma_2(z_2) + \dots = \frac{z_1}{\sum_{i=1}^m e^{z_i}} + \frac{z_2}{\sum_{i=1}^m e^{z_i}} + \dots + \frac{z_m}{\sum_{i=1}^m e^{z_i}}$$

$$= \frac{\sum_{i=1}^m z_i}{\sum_{i=1}^m e^{z_i}} = 1 \quad \text{sum} = 1$$

→ soft mode minimizes multi-class log loss aka cross-entropy.

$$-\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K y_{ij} \log p_{ij} \quad y_{ij} = 1 \text{ if } x_{ij} \in y_{ij} \quad p_{ij} = \text{prob of belonging to } y_{ij} \text{ otherwise}$$

How to train MLP (2.17) :-

① Preprocess : Data Normalization.

② weight initialization : (i) Xavier / Glorot → Sigmoid / Tanh
(ii) He → ReLU
(iii) Gaussian → Small σ

③ Choose activation functions :- ReLU as of now

④ Perform batch Normalization (especially for deep MLP (later layers))

Dropout is compulsory.

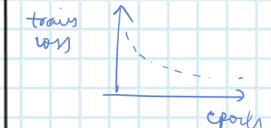
⑤ Optimizer : Adam (best one rn)

⑥ hyperparameters : lots of hyperparameters in DL (hyperas if you have GPU)

⑦ loss_fn → 2 classes → log loss
K classes → multi-class log loss
regression → square loss

⑧ Monitor Gradients - Clipping

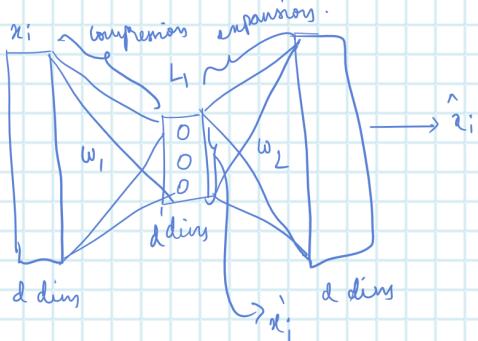
⑨ Decay -



⑩ Avoid overfitting - Easy to overfit in deep learning because lots of weights.

Autoencoder :- (2.18)

→ Neural network which can perform dim-red.

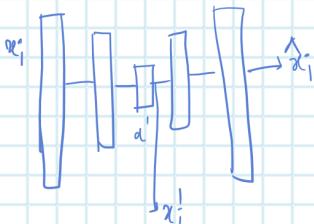


if the \hat{x}_i outputs from L_2 e.g. \hat{x}_i accurately represent x_i
then x_i encodes all of the info of x_i

$$L(x_i, \hat{x}_i) = \|x_i - \hat{x}_i\|^2$$

$$\text{if } x_i = \hat{x}_i, L=0$$

→ just like MLP, we can have deep autoencoders i.e., multiple hidden layers



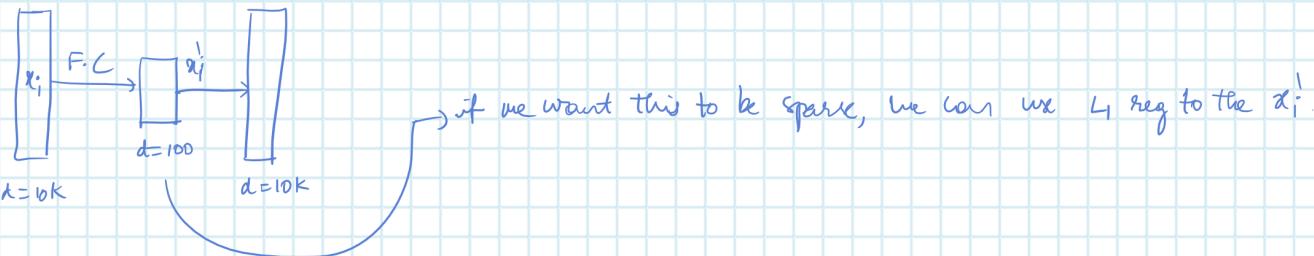
Denoising Autoencoder :-

$$x_i = \{\hat{x}_1, \dots, \hat{x}_n\} \rightarrow \text{noisy dataset}$$

$$\text{AutoEncoder}(x_i) = \hat{x}_i \quad \text{noise free, robust. Because it reduces the number of dimensions}$$

→ if original dataset is already noise free, we add noise to it ourselves and denoise it. Because test data points may not be noise free. Model can learn this.

Sparse Auto Encoder :-



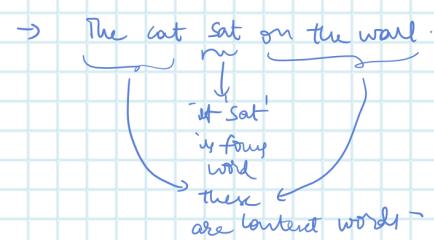
→ If there's only one hidden layer that only uses sigmoid act. function, the AE is strongly related to PCA.

→ AE can also be seen as unsupervised feature extraction process.

e.g. - MNIST 784 → 50 dimensions features.

Word2Vec : CBOW (2.19) :-

→ Not a neural network.



→ Word2Vec uses 2 algos

(i) CBOW :- Continuous Bag of Words.

(ii) Skipgram.

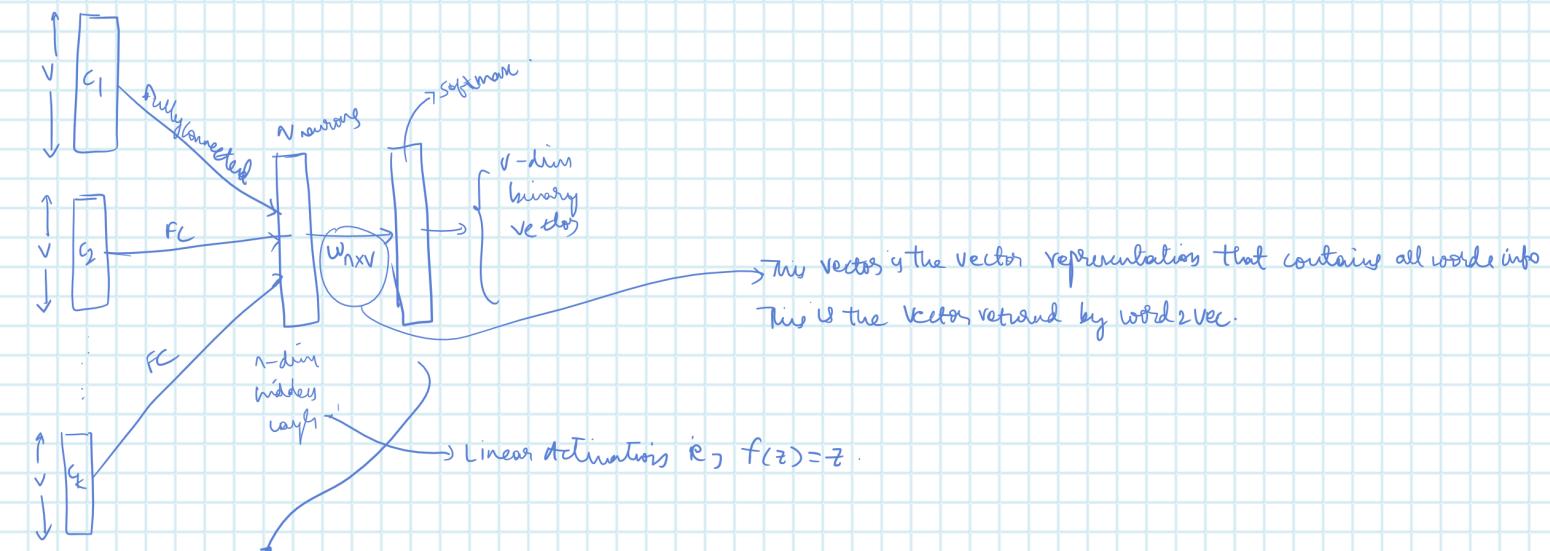
CBOW :-

① dictionary / vocab of words = V

$v = \text{length / size of vocab}$

② one hot encoding each word

$$w_i \in \mathbb{R}^V$$



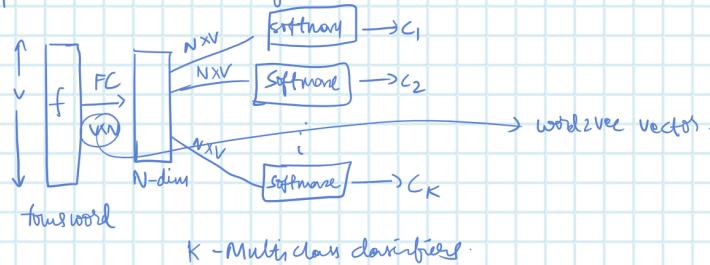
We index : predict focus word using content words -

\Rightarrow Multi class classification problem.

Training is done by taking all possible combinations of four words & content words -

word2Vec (skipgram) (2.20) :-

\rightarrow predict content words given four words.



K - Multi class classifier.

\rightarrow CBOW : 1 softmax

skipgram : K softmax

\therefore skipgram is computationally more expensive.

\rightarrow CBOW \rightarrow faster than skipgram

\rightarrow better for frequently occurring words

skipgram \rightarrow can work with smaller amounts of data.

\rightarrow works well for infrequent words.

$K = \# \text{ of content words}$
 $K \uparrow \Rightarrow$ more content \Rightarrow N -dim vector for each point is better.
 N is typically 200/300.

\rightarrow # of weights to be trained is large in both $\approx (K+1)(N \times V)$

So we need to find algorithmic optimizations.

word2vec (Algorithmic Optimizations) (1.21) :-

\rightarrow lots of weights to be trained.

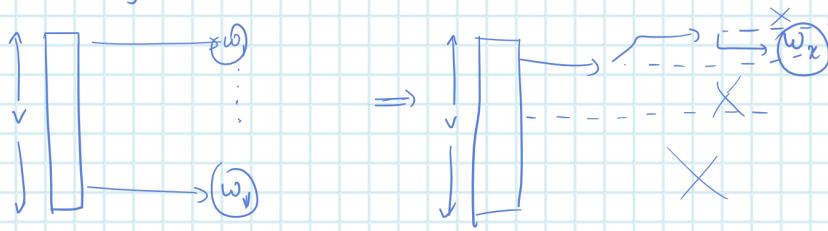
\rightarrow 2 optimization

(i) hierarchical softmax (algs)

(ii) Negative Sampling (statistical)

Hierarchical :-

→ idea is to modify v-softmaxes to make it optimal (skipgram)



binary tree like structure. We find correct word by backpropagation as correct half route in less loss.

= $\log(v)$ softmax classifiers required.

Negative Sampling :-

→ update only a sample of words per iterations

↳ always keep the target

↳ sample non target words with $P(w_i) = \frac{1}{\text{freq}(w_i)} \cdot 10^{-5}$

