

BATCH NORM

- avoids the problem of internal co-variate shifts
- when performing Gr update, for a given weight matrix, our Gr update assume that for a given W , the other ones are held constant. This is not true in general & therefore gradient step in earlier layers may affect the statistics of the neurons for the next layer. BN helps normalize these statistics so the Gr's are scaled and are drastically different after each Gr step.
- BN prevents outputs of each layer being too large or too small

BN

- per neuron per batch that is mean & variance calculated for all instances in a batch per channel
- per feature across batch & spatial dims.
- Large batches for stable statistics
- CNNs
- uses batch stats. for training & running averages during inference
- Image classification tasks

LN

- mean & variance calculated for a single instance across a layer.
- per sample across all features
- independent works with a single sample
- RNNS, Transformers
- Sample stats. for both training & testing
- NLP tasks, variable length sequences
- For deeper networks & CNNs
- A lot less variance batch size is reduced, for sequential data

BN → linear transformation

- accelerates training by requiring fewer iterations
- introduces noise to h_n activation bc mean & SD are estimated with a mini-batch of data
- uses running mean & variance from train. stats
- Reduces I.C.S
- Two learnable params (γ, β) - to scale, shift, & normalize values
- Regularization, bc adds noise to each mini batch
- During inference - it uses population stats computed during training on the entire data rather than the mini-batch
- For deeper networks & CNNs
- can introduce dependency between examples in a mini-batch

This shift in input distribution to hidden layers is called Internal Cov. Shift.
→ This is where batch normalization helps

layer norm → Statistics of hidden layers
batch norm → Statistics for mini-batch

Batch normalization → address the problem of ICS by normalizing the input of each layer

Let $x \in B$ (batch) denote an input to a BN layer
Then it is transformed as follows:

training sample \rightarrow minibatch

$$BN(x) = \gamma \frac{x - \mu_B}{\sigma_B} + \beta$$

Statistics of mini batch

$$\mu_B = \frac{1}{|B|} \sum_{x \in B} x$$

Minibatch sample mean
in mini-batch

\rightarrow W/out the scaling & shifting parameters, BN would always follow a unit variance & 0 mean distribution.

$$\sigma_B^2 = \frac{1}{|B|} \sum_{x \in B} (x - \mu_B)^2 + \epsilon$$

Minibatch sample variance
in mini-batch

γ → Scaling parameter
 β → Shifting parameter

\rightarrow $\beta \neq 0$.
 \rightarrow Learnable parameters

\rightarrow $\beta \neq 0$ & γ helps prevent constraint allowing the network to flexibly adjust feature representations.

how we apply it

$$h_i = f(BN(W_i h_{i-1} + b_i))$$

Matrix \times vector \times vector

$$\arg \min_{\theta} \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \theta^T \hat{x}^{(i)})^2 + \frac{\lambda}{2} \|\theta\|_2^2$$

$$L = \frac{1}{2} \sum_{i=1}^N (y^{(i)} - \theta^T \hat{x}^{(i)})^T (y^{(i)} - \theta^T \hat{x}^{(i)}) + \frac{\lambda}{2} \|\theta\|_2^2$$

$$= \frac{1}{2} \sum_{i=1}^N (y^{(i)\top} y^{(i)} - y^{(i)\top} \theta^T \hat{x}^{(i)} - \theta \hat{x}^{(i)\top} y^{(i)} + \theta \hat{x}^{(i)\top} \theta^T \hat{x}^{(i)}) + \frac{\lambda}{2} \|\theta\|_2^2$$

$$= \frac{1}{2} \sum_{i=1}^N (y^{(i)\top} y^{(i)} - 2 y^{(i)\top} \theta^T \hat{x}^{(i)} + \theta^T \hat{x}^{(i)\top} \hat{x} \theta) + \frac{\lambda}{2} \|\theta\|_2^2$$

Let, $X = [x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(n)}] \rightarrow$ data

$y = [y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(n)}] \rightarrow$ output, be matrix

$$= \frac{1}{2} (y^T y + 2 y^T X \theta + \theta^T X^T \theta) + \frac{\lambda}{2} \|\theta\|_2^2$$

$$\nabla_{\theta} L = \frac{1}{2} (-Z X^T y + Z X^T X \theta) + \lambda \theta$$

$$= -X^T y + X^T X \theta + \lambda \theta$$

$$X^T y = \theta (X^T X + \lambda I)$$

$$\theta = (X^T X + \lambda I)^{-1} X^T y$$

a) Regularization

- Helpful for better generalization
- Improves test set prediction

b) Mini-batch SGD

- It's a noisy approximation because only uses a few samples
- Only a part of the training data is involved

c) Adagrad Optimizer

- has equal or smaller learning LRs.
- The gradient history stays the same or increases in each dimension for every iteration of Adagrad

d) SGID $E = 5 \times 10^{-6}$

- Smaller LR's usually lead to getting stuck at a local minima
- Slower convergence
- Poor generalization

$$L(w, b) = \frac{1}{K} \sum_{i=1}^K \text{hinge}_y(x^{(i)}) + \lambda \|w\|_2^2$$

$$\nabla_w L(w, b) = \frac{1}{K} \sum_{i=1}^K \nabla_w \text{hinge}_y(x^{(i)}) + \lambda \nabla_w \|w\|_2^2$$

2-tensor to compute the gradients for

$$\text{hinge}_y(x^{(i)}) = \max(0, 1 - y^{(i)}(w^T x^{(i)} + b))$$

Piece-wise linear function

$$= \begin{cases} 0 & , y^{(i)}(w^T x^{(i)} + b) > 1 \\ 1 - y^{(i)}(w^T x^{(i)} + b), y^{(i)}(w^T x^{(i)} + b) \leq 1 \end{cases}$$

$$\text{If } y^{(i)}(w^T x^{(i)} + b) > 1, \text{ then } \nabla_w \text{hinge}_y(x^{(i)}) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^d$$

Yellow tensor gradients

$$\text{If } y^{(i)}(w^T x^{(i)} + b) \leq 1, \text{ then } \nabla_w \text{hinge}_y(x^{(i)}) = -y^{(i)} x^{(i)}$$

$$\text{Putting all together, } \nabla_w \text{hinge}_y(x^{(i)}) = \prod_{j=1}^d \begin{cases} y^{(i)}(w^T x^{(i)} + b) < 1 & 0 \\ y^{(i)}(w^T x^{(i)} + b) \geq 1 & -y^{(i)} x^{(i)} \end{cases}$$

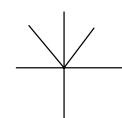
$$\|w\|_1 = \sum_{i=1}^d |w_i|$$

$$\begin{bmatrix} y^{(i)}(w^T x^{(i)} + b) > 1 \\ \vdots \\ y^{(i)}(w^T x^{(i)} + b) > 1 \end{bmatrix}$$

vector of all 1's

$$\nabla_w \|w\|_1 = \begin{bmatrix} \nabla_w |w_1| \\ \nabla_w |w_2| \\ \vdots \\ \nabla_w |w_d| \end{bmatrix}$$

$$|w_i| = \begin{cases} w_i, & w_i > 0 \\ -w_i, & w_i \leq 0 \end{cases}$$



$$\nabla_{w_i} |w_i| = \begin{cases} 1, & w_i > 0 \\ -1, & w_i \leq 0 \end{cases}$$

Sub-gradient
- anything between
1 & -1



$$\text{If } w_i = 0, \nabla_{w_i} (w_i) = 0$$

$$\nabla_{w_i} |w_i| = \text{Sgn}(w_i)$$

$$\nabla_w \|w\|_1 = \text{Sgn}(w) \in \mathbb{R}^d$$

$$\nabla_w L(w, b) = \frac{1}{K} \sum_{i=1}^K \prod_{j=1}^d \begin{cases} y^{(i)}(w^T x^{(i)} + b) < 1 & 0 \\ y^{(i)}(w^T x^{(i)} + b) \geq 1 & \text{Sgn}(w) \end{cases}$$

$$\nabla_w L(w, b) = \frac{1}{K} \sum_{i=1}^K \text{hinge}_y(x^{(i)})$$

$$= \max(0, 1 - y^{(i)}(w^T x^{(i)} + b))$$

We can treat it as a piece-wise function

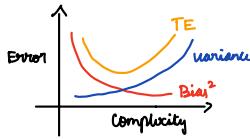
$$\text{hinge}_y(x^{(i)}) = \begin{cases} 0 & , \text{If } y^{(i)}(w^T x^{(i)} + b) \geq 1 \\ 1 - y^{(i)}(w^T x^{(i)} + b), \text{ If } y^{(i)}(w^T x^{(i)} + b) < 1 \end{cases}$$

$$\text{When } y^{(i)}(w^T x^{(i)} + b) < 1,$$

$$\nabla_w \text{hinge}_y(x^{(i)}) = -y^{(i)} x^{(i)}$$

$$\therefore \nabla_w L = \frac{1}{K} \sum_{i=1}^K \prod_{j=1}^d \begin{cases} y^{(i)}(w^T x^{(i)} + b) < 1 & 0 \\ y^{(i)}(w^T x^{(i)} + b) \geq 1 & -y^{(i)} x^{(i)} \end{cases}$$

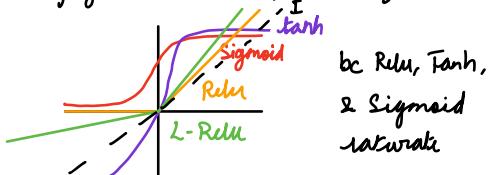
$$\nabla_b L = \frac{1}{K} \sum_{i=1}^K \prod_{j=1}^d \begin{cases} y^{(i)}(w^T x^{(i)} + b) < 1 & 0 \\ y^{(i)}(w^T x^{(i)} + b) \geq 1 & -y^{(i)} \end{cases}$$



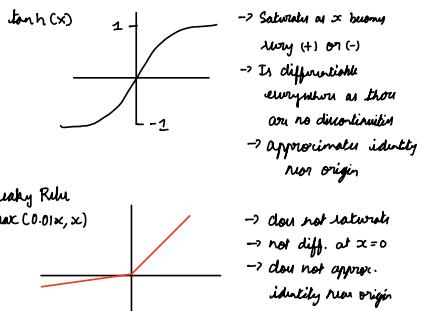
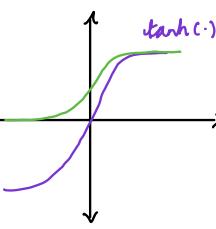
- 2↑ underfit training data → ↑ bias - ↓ variance (less complex)
 2↓ overfit training data → ↓ bias - ↑ variance (more complex)

Dropout (an approx. to bagging - where all models share weights)
 → During each forward pass, some activations may be zeroed out
 → But units somehow get included in the training process
 → ∴ they have the same no. of parameters as the vanilla network
 → The weights dropped during the FP are not zeroed out but they are simply not updated in subsequent BP steps. ∴ Scaling by p does not reduce the size of the model

Vanishing gradients: ReLU, Tanh, Sigmoid



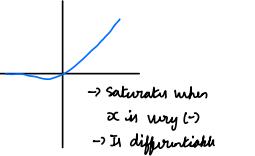
bc ReLU, Tanh,
 & Sigmoid
 saturate



- Saturates as $x \rightarrow \infty$ (very (+) or (-))
- Is differentiable everywhere as there are no discontinuities
- Approximates identity near origin

Gradient Error Linear Unit (GELU)

$$GELU(x) = 0.5x(1 + \tanh[\sqrt{2/\pi}(x + 0.045x^3)])$$



- Does not saturate
- not diff. at $x=0$
- does not approximate identity near origin

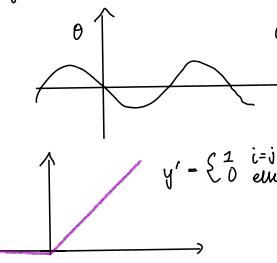
- Saturates when x is very (-)
- Is differentiable everywhere
- Does not approximate identity near origin

Leaky ReLU

$$\max(0.01x, x)$$

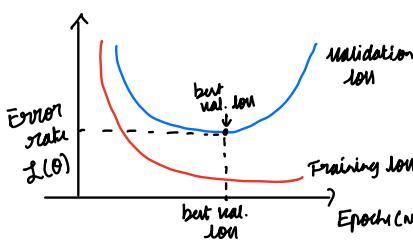
$$\max(\alpha \cdot g, g), 0 < \alpha < 2$$

$$y = \sin(\omega x)$$

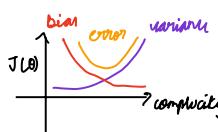


Pro: ① non-linear / differentiable everywhere

Con: Periodic $\sin(\omega x)$ may equal to non-unique outputs for different inputs. The oscillatory gradient can result in vanishing gradients → leading to unstable training



↑ Train → Overfitting → Do not add layers!
 accuracy ↓ bias → Regularize
 bad testing error ↑ variance → augmentation
 testing error ↑ variance → dropout
 very complex



Initializations

- d. Constant values may not be a good idea for a fully connected NN. This can lead to symmetry breaking issues where neurons in the layers start to behave identically limiting the network capacity to learn diverse features. Each unit gets signal equal to the sum of inputs. The gradients of all the neurons will be the same in the layer.
- Weights initialized at 0 will result in 0 gradients, ∴ no gradients
- Large initializations can lead to vanishing gradients (saturation range for tanh & sigmoid).
- “0 initialization kills learning”
- Init. with the same weight will not lead to the weights being equal
- Symmetry can cause σ to be equal for weights sharing the same input feature, but weights connected to different input features can break symmetry
- Random init. → exploding / vanishing gradients & worse I.C.S.
- For the right batch size, random init. does not affect BN strength bc the $\mu \approx 0$ of the output of BN is constrained & learned

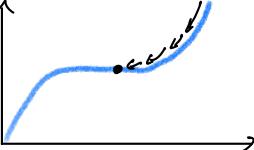
Name	Plot	Equation	Derivative (with respect to x)
Identity		$f(x) = x$	$f'(x) = 1$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a. Sigmoid or Soft step)		$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
TanH		$f(x) = \tanh(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})}$	$f'(x) = 1 - f(x)^2$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$
ElliotSig ^{[9][10][11]} Softsign ^{[12][13]}		$f(x) = \frac{x}{1 + x }$	$f'(x) = \frac{1}{(1 + x)^2}$
Inverse square root unit (ISRU) ^[14]		$f(x) = \frac{x}{\sqrt{1+ax^2}}$	$f'(x) = \left(\frac{1}{\sqrt{1+ax^2}}\right)^3$
Inverse square root linear unit (ISRLU) ^[14]		$f(x) = \begin{cases} \frac{x}{\sqrt{1+ax^2}} & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \left(\frac{1}{\sqrt{1+ax^2}}\right)^3 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Square Nonlinearity (SQNL) ^[11]		$f(x) = \begin{cases} 1 & : x > 2.0 \\ x - \frac{x^2}{4} & : 0 \leq x \leq 2.0 \\ x + \frac{x^2}{4} & : -2.0 \leq x < 0 \\ -1 & : x < -2.0 \end{cases}$	$f'(x) = 1 \mp \frac{x}{2}$
Rectified linear unit (ReLU) ^[15]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Leaky rectified linear unit (Leaky ReLU) ^[16]		$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0.01 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Parametric rectified linear unit (PReLU) ^[17]		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Randomized leaky rectified linear unit (RReLU) ^[18]		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
SoftPlus ^[23]		$f(x) = \ln(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$
Bent identity		$f(x) = \sqrt{x^2 + 1} - 1 + x$	$f'(x) = \frac{x}{2\sqrt{x^2 + 1}} + 1$
SoftExponential ^[26]		$f(\alpha, x) = \begin{cases} -\frac{\ln(1 - \alpha(x + \alpha))}{\alpha} & \text{for } \alpha < 0 \\ x & \text{for } \alpha = 0 \\ \frac{e^{\alpha x} - 1}{\alpha} & \text{for } \alpha > 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \frac{1}{1 - \alpha(x + \alpha)} & \text{for } \alpha < 0 \\ e^{\alpha x} & \text{for } \alpha \geq 0 \end{cases}$
Soft Clipping ^[27]		$f(\alpha, x) = \frac{1}{\alpha} \log \frac{1 + e^{\alpha x}}{1 + e^{\alpha(x-1)}}$	$f'(\alpha, x) = \frac{1}{2} \sinh\left(\frac{p}{2}\right) \operatorname{sech}\left(\frac{px}{2}\right) \operatorname{sech}\left(\frac{p}{2}(1-x)\right)$
Sinusoid ^[28]		$f(x) = \sin(x)$	$f'(x) = \cos(x)$
Sinc		$f(x) = \begin{cases} 1 & \text{for } x = 0 \\ \frac{\sin(x)}{x} & \text{for } x \neq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x = 0 \\ \frac{\cos(x)}{x} - \frac{\sin(x)}{x^2} & \text{for } x \neq 0 \end{cases}$
Gaussian		$f(x) = e^{-x^2}$	$f'(x) = -2xe^{-x^2}$

- GD $\theta_{t+1} \leftarrow \theta_t + \epsilon \nabla_{\theta} J(\theta)$
 SGD $\theta_{t+1} \leftarrow \theta_t + \epsilon \nabla_{\theta} J_i(\theta)$
 SGD \rightarrow faster than GD when dataset is large
 \rightarrow sensitive to noise in the data
 \rightarrow more efficient in terms of memory usage
 GD + momentum update: $v_{t+1} = \epsilon v_t - \epsilon \nabla_{\theta} L(\theta_t)$
 \rightarrow do not adjust step size per parameter
 \rightarrow overshoot in \uparrow curvature regions
 GD + Nesterov update: $v_{t+1} = \epsilon v_t - \epsilon \nabla_{\theta} L(\theta_t + \epsilon v_t)$
 Weight update: $\theta_{t+1} = \theta_t + v_{t+1}$

Momentum - avg. the zigzagging components
 $v \leftarrow \epsilon v - \epsilon g$ v is the running avg. of historical gradients
 $\theta \leftarrow \theta + v$
 - helpful in larger networks
 - SGD + momentum \rightarrow helps find shallow local optimum
 \rightarrow sensitive to noise

Nesterov momentum (anticipates upcoming gradients)
 \rightarrow gradient is calculated at parameter setting after taking a step direction along the direction of momentum history
 $v \leftarrow \epsilon v - \epsilon \nabla_{\theta} J(\theta + \epsilon v)$
 $\theta \leftarrow \theta + v$
 $v_{\text{new}} = \epsilon v_{\text{old}} - \epsilon \nabla_{\theta} J(\theta_{\text{old}})$
 $v = \epsilon v - \epsilon \nabla_{\theta} L(\theta + \epsilon v)$
 $\theta_{\text{new}} = \theta_{\text{old}} + \epsilon v_{\text{new}}$
 $= \theta_{\text{old}} + v_{\text{new}} + \epsilon v_{\text{new}}$
 $= \theta_{\text{old}} - \epsilon v_{\text{old}} + v_{\text{new}} + \epsilon v_{\text{new}}$
 $= \theta_{\text{old}} + v_{\text{new}} + \epsilon(v_{\text{new}} - v_{\text{old}})$
 $\theta_{\text{new}} = \theta_{\text{old}} + v_{\text{new}} + \epsilon(v_{\text{new}} - v_{\text{old}})$
 Don't have to evaluate or at $\theta + \epsilon v$

Nesterov momentum will help with faster convergence after a steep slope bc the accumulated momentum will be large. The weight update is always opposite to the direction of the gradient, \therefore in our case it will be a large (\downarrow) update on weights overshooting in the case of GD + momentum



Adaptation LR

Adagrad (accumulates squared gradients)
 $a \leftarrow a + g \otimes g$ \rightarrow 2nd moment takes variance
 $\theta \leftarrow \theta - \frac{\epsilon}{\sqrt{a} + \epsilon} \otimes g$ into account \rightarrow reduces effective step size to approx. annealing LR where G's are large
 $\rightarrow a \leftarrow a + g \otimes g$ \rightarrow accum. sum of past G's squared
 \rightarrow global minimum reached!
 \rightarrow division by a larger number \rightarrow small step
 \rightarrow No momentum calculated, only update θ on past G.
 Adagrad / RMS-prop + momentum = Adam

1st moment $v \leftarrow \beta_1 v + (1-\beta_1) g$ (momentum-like)
 $a \leftarrow \beta_2 a + (1-\beta_2) g \otimes g$ (G normalization)
 $\theta \leftarrow \theta - \frac{\epsilon}{\sqrt{a} + \epsilon} \otimes v$

Adam with bias correction

$$v \leftarrow \beta_1 v + (1-\beta_1) g$$

$$a \leftarrow \beta_2 a + (1-\beta_2) g \otimes g$$

$$\tilde{v} = \frac{1}{1-\beta_1^t} v; \quad \tilde{a} = \frac{1}{1-\beta_2^t} a;$$

$$\theta \leftarrow \theta - \frac{\epsilon}{\sqrt{\tilde{a}} + \epsilon} \otimes \tilde{v}$$

at t gets large, $\beta_1^t \rightarrow 0$, $\beta_2^t \rightarrow 0$,
 convolution disappears & we have
 $\tilde{v} = v, \quad \tilde{a} = a$

RMS Prop $a \leftarrow \beta a + (1-\beta) g \otimes g$
 $\theta \leftarrow \theta - \frac{\epsilon}{\sqrt{a} + \epsilon} \otimes g$
 \rightarrow makes a an exponentially weighted moving avg.
 \rightarrow does not get slower towards the end like Adagrad bc of the β parameter
 \rightarrow can take bigger steps
 \rightarrow leads to global minimum.

RMS prop + momentum

\rightarrow Although it makes a large excursion, gets to optima quickly

$$a \leftarrow \beta a (1-\beta) g \otimes g$$

$$v \leftarrow \epsilon v - \frac{\epsilon}{\sqrt{a} + \epsilon} \otimes g$$

$$\theta \leftarrow \theta + v$$

$$g_{t+1} = \beta g_t + (1-\beta) (\nabla_{\theta} L(\theta_t))^2$$

$$v_{t+1} = \epsilon v_t - \frac{\epsilon}{\sqrt{g_{t+1}} + \delta} \nabla_{\theta} L(\theta_t)$$

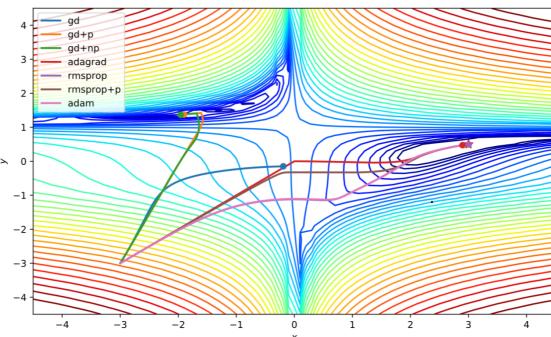
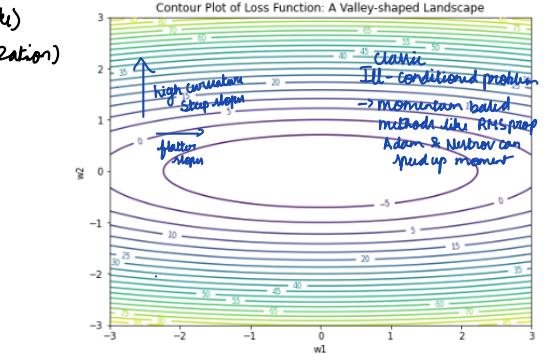
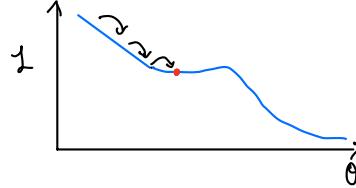
$$\theta_{t+1} = \theta_t + v_{t+1}$$

$$g_t = \nabla_{\theta} L(\theta_{t-1})$$

i) Small step size

GD + momentum / Adam

will be good option to get out of the local minimum bc they would be able to gather momentum starting up-hill on the left which will eventually drive the optimizer out of the plateau where the G's can vanish.



Second Order methods

- \rightarrow use second derivative (or Hessian) to avoid curvatures
 \rightarrow can capture curvature of the loss surface
 \rightarrow But computationally expensive

Newton's method

$$\theta \leftarrow \theta - H^{-1} g$$

$$\theta \leftarrow \theta - (H + \delta I)^{-1} \nabla_{\theta} J(\theta_0)$$

$$\rightarrow$$
 memory storage is expensive

$$\rightarrow$$
 Inverting $H \rightarrow O(n^3)$

$$\rightarrow H$$
 requires very large batch sizes.

Update magnitude order

GD + momentum > Adam > Adagrad
 momentum will help increase the update magnitude at the plateau while Adagrad will shrink it. Combining both, Adam will give an update magnitude in the middle of the 2 extremes

Vanishing Gradients

update magnitude order

Newton's method > Adam + bias correction > Adam
 > Nesterov Momentum > GD + momentum > RMSprop + momentum > SGD > GD > RMSprop > Adagrad

Layer Dimension:

$$\frac{W - W_f + 2 \cdot \text{pad}}{\text{Stride}} + 1 = \frac{h - h_f + 2 \cdot \text{pad}}{\text{Stride}} + 1$$

Total params: $(3 \times 3 \times 3 + 1) \times 5 = 140$

\uparrow n_f
 \uparrow bias

MAXPOOL

→ max pool layer does not have any params.

→ to reduce invariance due to shifts

→ concept of saliency maps

→ Filter has its own bias term (+1 per filter)

of neurons → Output height × Output width × n_f

of params → (filter size + 1) × n_f

Newton Momentum + GD

$$V_t = \alpha V_{t-1} - \epsilon \nabla f(\theta_{t-1} + \alpha V_{t-1})$$

$$\theta_t = \theta_{t-1} + v_t$$

$$v_0 = 0$$

$$v_1 = -\epsilon \nabla f(\theta_0)$$

$$v_2 = \alpha v_1 - \epsilon \nabla g_1$$

$$v_3 = \alpha v_2 - \epsilon \nabla f g_2$$

$$v_t = -\epsilon \sum_{i=1}^t \alpha^{t-i}$$

$$\theta_t = \theta_{t-1} + v_t ; \quad \theta_t = \theta_0 + \sum_{j=1}^t v_j$$

$$\theta_t = \theta_0 - \epsilon \sum_{j=1}^t \sum_{i=1}^j \alpha^{j-i} \nabla f(\theta_{i-1} + \alpha v_{i-1})$$

$$\text{OR } \theta_t = \theta_0 - \epsilon \alpha \sum_{i=1}^t \alpha^{t-i} g_i - \epsilon \sum_{j=1}^t \left(\sum_{i=1}^j \alpha^{j-i} g_i \right)$$

RMSprop + Momentum

The update equations are:

$$g_t = \beta g_{t-1} + (1 - \beta)(\nabla_\theta L(\theta_{t-1}))^2$$

$$v_t = \alpha v_{t-1} - \frac{\epsilon}{\sqrt{g_t} + \delta} \nabla_\theta L(\theta_t)$$

$$\theta_t = \theta_{t-1} + v_t$$

Recursion Expansion

Expanding v_t recursively, assuming $v_0 = 0$:

$$v_t = -\sum_{i=1}^t \alpha^{t-i} \frac{\epsilon}{\sqrt{g_i} + \delta} g_i$$

Now, substituting into the weight update equation:

$$\theta_t = \theta_0 + \sum_{j=1}^t v_j$$

$$\theta_t = \theta_0 - \sum_{j=1}^t \sum_{i=1}^j \alpha^{j-i} \frac{\epsilon}{\sqrt{g_i} + \delta} g_i$$

Adagrad + Momentum

The update equations are:

$$g_t = g_{t-1} + (\nabla_\theta L(\theta_{t-1}))^2$$

$$v_t = \alpha v_{t-1} - \frac{\epsilon}{\sqrt{g_t} + \delta} \nabla_\theta L(\theta_t)$$

$$\theta_t = \theta_{t-1} + v_t$$

Recursion Expansion

Expanding v_t recursively:

$$v_t = -\sum_{i=1}^t \alpha^{t-i} \frac{\epsilon}{\sqrt{g_i} + \delta} g_i$$

Now, substituting into the weight update equation:

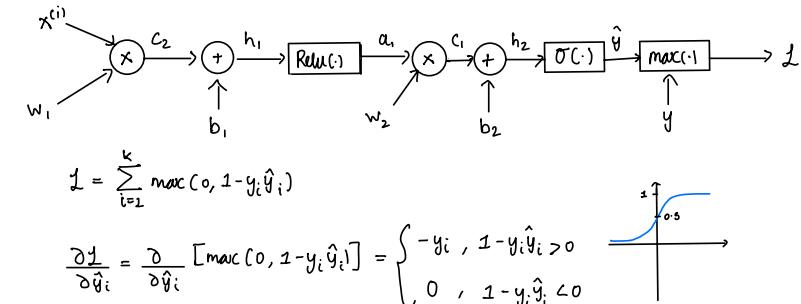
$$\theta_t = \theta_0 + \sum_{j=1}^t v_j$$

$$\theta_t = \theta_0 - \sum_{j=1}^t \sum_{i=1}^j \alpha^{j-i} \frac{\epsilon}{\sqrt{g_i} + \delta} g_i$$

Consider the neural network with the following architecture

$$\begin{aligned} h_1 &= W_1 x^{(i)} + b_1 \\ a_1 &= \text{ReLU}(h_1) \\ h_2 &= W_2 a_1 + b_2 \\ \hat{y} &= \sigma(h_2) \\ L &= \sum_{i=1}^K \max(0, 1 - y_i \hat{y}_i) \end{aligned}$$

where $x^{(i)} \in \mathbb{R}^{D_x \times 1}$, $W_1 \in \mathbb{R}^{D_{a_1} \times D_x}$, $W_2 \in \mathbb{R}^{K \times D_{a_1}}$. σ is the sigmoid activation function and y is a K -dimensional vector of 1's and 0's. y_i represents the i^{th} element of y and \hat{y}_i represents the i^{th} element of \hat{y} .



\hat{y} is the output of a sigmoid operator, so $\hat{y} < 1$.

from problem statement $y_i \leq 1$. Hence $1 - y_i \hat{y}_i > 0$

$$\frac{\partial L}{\partial \hat{y}_i} = -y_i \quad , \quad \frac{\partial L}{\partial y_i} = -y$$

$$\frac{\partial L}{\partial h_2} = \frac{\partial \hat{y}_i}{\partial h_2} \frac{\partial L}{\partial \hat{y}_i} \quad \text{softmax is defined elementwise}$$

$$\frac{\partial \hat{y}_i}{\partial h_2} = \sigma(h_{2,i}) [1 - \sigma(h_{2,i})] = \text{diag}[\sigma(h_{2,1}), \sigma(h_{2,2}), \dots, \sigma(h_{2,K})]$$

$$\frac{\partial L}{\partial h_2} = -\hat{y}_i \sigma(1 - \hat{y}_i) \sigma(y_i) = \delta_{h_2}$$

$$\oplus \quad \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial h_2} = \delta_{h_2} \quad \frac{\partial L}{\partial c_1} = \frac{\partial L}{\partial h_2} = \delta_{h_2}$$

$$\otimes \quad \frac{\partial L}{\partial w_2} = \frac{\partial c_1}{\partial w_2} \frac{\partial L}{\partial c_1} = \delta_{h_2} a_1^T$$

$$\frac{\partial L}{\partial a_1} = \frac{\partial c_1}{\partial a_1} \frac{\partial L}{\partial c_1} = w_2^T \delta_{h_2}$$

$$a_1 = \text{ReLU}(h_1)$$

$$\frac{\partial L}{\partial h_1} = \frac{\partial a_1}{\partial h_1} \frac{\partial L}{\partial a_1} = \prod_{i=1}^K (h_i > 0) \odot \delta_{a_1}$$

$$\oplus \quad \frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial h_1} = \delta_{h_1}$$

$$\frac{\partial L}{\partial c_2} = \frac{\partial L}{\partial h_1} = \delta_{h_1}$$

$$\otimes \quad C_2 = w_1 \times^{(i)}$$

Tensor derivative trick

$$\frac{\partial L}{\partial w_1} = \frac{\partial c_2}{\partial w_1} \frac{\partial L}{\partial c_2} = \delta_{h_1} x^{(i)T}$$

$$\mathbb{R}^{D_w \times D_x \times D_{a_1}}$$

$$2. a) \nabla_{\hat{y}^{(i)}} \mathcal{L}^{(i)} \rightarrow \delta_{\hat{y}^{(i)}}$$

$$\mathcal{L}^{(i)} = y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})$$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}^{(i)}} = \frac{y^{(i)}}{\hat{y}^{(i)}} - \left[\frac{1-y^{(i)}}{1-\hat{y}^{(i)}} \right]$$

$$b) \nabla_{z_3} \mathcal{L}^{(i)} \rightarrow \delta_{z_3}$$

$$\frac{\partial \mathcal{L}}{\partial z_3} = \frac{\partial \hat{y}}{\partial z_3} \frac{\partial \mathcal{L}}{\partial \hat{y}} = \sigma(z_3)[1-\sigma(z_3)] \cdot \delta_{\hat{y}^{(i)}} = \delta_{z_3}$$

$$c) \nabla_{w_2} \mathcal{L}^{(i)} \rightarrow \delta_{w_2}$$

$$m = w_2 z \quad \frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial z}{\partial w_2} \frac{\partial \mathcal{L}}{\partial z} = z^T \cdot \delta_{z_3}$$

$$d) \nabla_z \mathcal{L}^{(i)} \rightarrow \delta_z$$

$$m = w_2 z \quad \frac{\partial \mathcal{L}}{\partial z} = \delta_{z_3} \cdot w_2^T$$

$$e) \nabla_{h_q} \mathcal{L}^{(i)} \rightarrow \delta_{h_q}$$

$$\frac{\partial \mathcal{L}}{\partial h_q} = \frac{\partial z_2}{\partial h_q} \frac{\partial \mathcal{L}}{\partial z_2} = \left[\frac{e^{\beta h_q}}{1+e^{\beta h_q}} \right] \cdot \delta_z$$

$$f) \nabla_{h_p} \mathcal{L}^{(i)} \rightarrow \delta_{h_p}$$

$$\frac{\partial \mathcal{L}}{\partial h_p} = \frac{\partial z_1}{\partial h_p} \frac{\partial \mathcal{L}}{\partial z_1} = \left[\frac{e^{\beta h_p}}{1+e^{\beta h_p}} \right] \cdot \delta_z$$

$$g) \nabla_b \mathcal{L}^{(i)}$$

$$\nabla_b \mathcal{L}^{(i)} = \delta_{h_p} + \delta_{h_q}$$

$$h) \nabla_{w_i} \mathcal{L}^{(i)}$$

$$\nabla_{w_i} \mathcal{L}^{(i)} = \delta_{h_p} x_p^{(i)T} + \delta_{h_q} x_q^{(i)T}$$

$$i) \tilde{\mathcal{L}} = \left(-\frac{1}{m} \sum_{i=1}^m \mathcal{L}^{(i)} \right) + \ln \left(\sum_{i=1}^m e^{x_i} \right) + \|w_i\|_\infty$$

Since
gradients
in a linear
operator

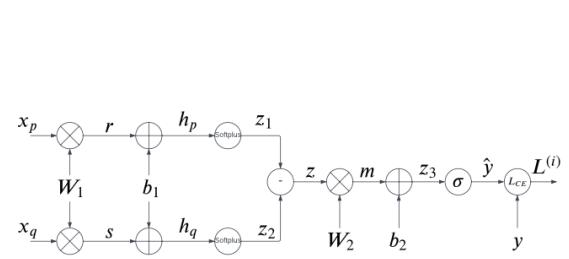
$$\nabla_{w_i} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial w_i} = -\frac{1}{m} \sum_{i=1}^m \nabla_{w_i} \mathcal{L}^{(i)} + \nabla_{w_i} \|w_i\|_\infty$$

$$\nabla_{w_2} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial w_2} = -\frac{1}{m} \sum_{i=1}^m \nabla_{w_2} \mathcal{L}^{(i)} + \nabla_{w_2} \ln \left(\sum_{i=1}^m e^{x_i} \right)$$

$$\nabla_{w_i} \|w_i\|_\infty = \begin{bmatrix} 0 \\ 0 \\ \text{Sign}(w_j) \\ 0 \end{bmatrix} \quad \nabla_{w_2} \text{LSE} = \begin{bmatrix} \frac{e^{x_1}}{\sum_{i=1}^m e^{x_i}} \\ \frac{e^{x_2}}{\sum_{i=1}^m e^{x_i}} \\ \vdots \end{bmatrix}^T = \text{softmax}$$

$$\therefore \nabla_{w_i} \mathcal{L} = -\frac{1}{m} \sum_{i=1}^m \nabla_{w_i} \mathcal{L}^{(i)} + \begin{bmatrix} 0 \\ 0 \\ \text{Sign}(w_j) \\ 0 \end{bmatrix}$$

$$\nabla_{w_2} \mathcal{L} = -\frac{1}{m} \sum_{i=1}^m \nabla_{w_2} \mathcal{L}^{(i)} + \begin{bmatrix} \frac{e^{x_1}}{\sum_{i=1}^m e^{x_i}} \\ \frac{e^{x_2}}{\sum_{i=1}^m e^{x_i}} \\ \vdots \end{bmatrix}^T$$



$$h_p = W_1 x_p^{(i)} + b_1$$

$$z_1 = \text{softplus}(h_p)$$

$$h_q = W_2 x_q^{(i)} + b_2$$

$$z_2 = \text{softplus}(h_q)$$

$$z = z_1 - z_2$$

$$z_3 = W_2 z + b_3$$

$$\hat{y}^{(i)} = \sigma(z_3)$$

$$L^{(i)} = L_{CE}(y^{(i)}, \hat{y}^{(i)})$$

$$L = -\frac{1}{m} \sum_{i=1}^m L^{(i)}$$