

RNN Systems POV: FIR Filters

- Convolution-like form; some maximum receptive field (accepts finite inputs)
- Constant z-transform
- Response characteristics filter
- All CNNs/GNNs

2) IIR Filters operate on potentially infinite length signals

- State (evolution over time) encodes dynamics; complex z-transform
- Multiple poles exist due to recurrent dynamics
- View as streaming/online algo

Ex Kalman Filtering. Given known dynamic of linear sys. driven by Gauss. noise (known covariance)

$$\begin{aligned} x_{t+1} &= Ax_t + B_u u_t + w_t \\ y_{t+1} &= Cx_{t+1} + v_t \end{aligned}$$

- Assm. w, v iid. jointly Gauss. & x, y, w, v all jointly
- Following these being linear map $y \rightarrow (w, v)$

- Gauss. min. $\mathbb{E}[\|x_t - \hat{x}_t\|^2]$

$$\hat{x}_{t+1} = \hat{A}\hat{x}_t + \hat{B}_u u_t + \hat{B}_y y_{t+1}$$

- Int. \hat{x}_t should capture stable dyn. could be learnable, or O/random hoping requirements will get us along

RNN seq2seq challenges:

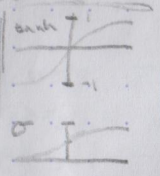
- Requires sequential processing (slow)
- Fixed length context (doesn't work well for long-range dependencies)
- Target/source word ordering may differ but can't "align" to fix
- Motivates encoder-decoder seq2seq

Non-linearity: Can output into, motivated skip

- Non-saturating: If $\lim_{x \rightarrow \pm\infty} f(x) = \pm\infty$
- Ex. ReLU (no)

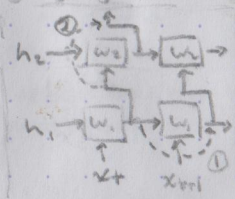
- Saturating: Stabilizes dynamics by preventing explosions

- Ex. tanh & sigmoid
- Apply to recurrent state
- here $z = \sigma(Ah + b)$
- could also layer norm h_t or x_t
- Not batch, as this wouldn't affect recurrent aspect
- Leads to dying grads for extreme values



SKIPS: Tackle dying grads

- Leads to too much info/long-range dependence
- LSTM solved by gating
- Attention solves by using soft hashmap



- Capture hierarchical relations based on similarity measure between states
- Weighted avg. so grads aren't harsh
- Alleviates bottleneck from encoder; enhance feat. expressiveness/reduce
- Cross-attn: (k, v) from enc. regularization q from decoder
- Self-attn: (k, q, v) from same source
- Decoder often made causal by masking
- Multi-head: lookup l items over l tables; concatenation
- Multi-query: Lookup l items over 1 table; (k, v) shared across queries
- Less params; faster @ eval
- Softmax induces 1 primary feature per head
- Reception field whole seq.

Auto Encoders / Linear: $\mathcal{L} = \|X - W_e W_d X\|_F^2$

- Projects $W_e W_d$ onto k largest S.V. of X
- $\rightarrow W_e W_d$ should look like I for $k \ll n$
- i.e. $W_e \approx W_d^T \equiv \Sigma_k \approx \Sigma_k^T$
- Regularizer: $\lambda(\|W_e\|_F^2 + \|W_d\|_F^2)$
- Induces orthonorm $\lambda(\Sigma_k^T \Sigma_k + \Sigma_k \Sigma_k^T) = \lambda(\Sigma_k^T + \Sigma_k)$
- cells W_e $0 = \nabla[\dots] \equiv 2\sigma_i = 2/\sigma_i$
- rows W_d $\Rightarrow \sigma_i \equiv 1$
- asym Unique $\rightarrow [n] \rightarrow \infty$
- λ small enough to decouple terms

$$\begin{aligned} \text{New loss } \mathcal{L} &= \|X - WX\|_F^2 + \lambda\|W\|_F^2 \\ &= \sum_{i=1}^n \|x_i - Wx_i\|_2^2 + \lambda\|W\|_F^2 \end{aligned}$$

$$H_g = ABH + A\|b\|_2^2$$

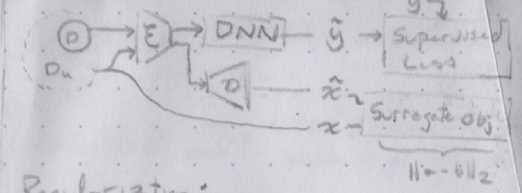
$$\begin{aligned} \text{For } y &= XT \\ b &= WT \\ A &= XT \\ \Rightarrow \|x - x^T W\|_2^2 + \lambda\|W\|_F^2 \\ &= (ATA + \lambda I)^{-1} ATy \\ W^T &= (KT^T \lambda I)^{-1} XT^T \\ W &= XX^T (XX^T + \lambda I)^{-1} \\ &= (U \Sigma U^T) (U \Sigma^T U^T + \lambda I)^{-1} \\ &= U \Sigma \Sigma^T U^T (\Sigma^T + \lambda I)^{-1} \\ &= U (\Sigma \Sigma^T / (\Sigma \Sigma^T + \lambda I)) U^T \end{aligned}$$

$$\begin{aligned} \text{Attrn}(k, q, v): \quad e_i &= q^T v_i / \sqrt{d} \sim N(\mu, \sigma) \\ d_i &= \exp(e_i) / \sum_j \exp(e_j) \\ r_i &= \sum_j d_i v_j \end{aligned}$$

- Scaling: $q^T k_i \approx \frac{1}{\sqrt{d}} \sum_j q_j k_{ij} \rightarrow N(\mu, \sqrt{d})$
- Ex $X \sim N(\mu, \sqrt{d})$ and sample x_1, \dots, x_n
- $\mathbb{E}[x_i - x_j] = \sqrt{d}/n$ so dist. between pts $\propto \sqrt{d}$, enlarged by softmax (acts like hardmax at init b.c. winner wins by scale/dumb luck)

- Had we normalized by d , we'd constrain $q, k \in [0, 1]$ so sim. hardly differ (all small)
- \rightarrow Prior, aligned case magnitude $\frac{d}{\sqrt{d}}$, here it's $\frac{d}{d} = 1$
- Improper scaling \Rightarrow large loss causing softmax to give small grads (slow learning)

$$\text{Aside: } T = 0 \equiv \text{pick best}$$



Regularization:

- Adding noise, $x_i + \eta$ for $\eta \sim N(0, \sigma^2)$ makes network de-noise, which requires understanding of implicit structure of x .
- Making \equiv subtractive noise

$$\begin{aligned} \text{VAE } x &= \mathbb{E}_{z \sim N(\mu, \Sigma)} [x] \\ &= \mathbb{E}_{z \sim N(\mu, \Sigma)} [x] - \mathbb{E}_{z \sim N(0, I)} [z] \end{aligned}$$

$$VT: O((n/p \times w/p)^2 D)$$

MCA only computes $XW_k = XWV$ over saving $2Md^2(h-1)$ operations.

For input seq = M , output = N , h heads $W_k, W_q, W_v \in \mathbb{R}^{d \times d}$

Encoder SA takes $3Md^2$ For $\{XW_k\}_{k=1}^h$ compute attn mat. $M^2 d$ and $O(M^2)$ mul. w/ value $M^2 d$ over h heads

- Enc. SA/MCA: $O((Md + M^2)dh)$
- Cross Attn: $O((Md + Md + MN)dh)$
- Decoder SA: $O((Nd + N^2)dh)$

$$\text{PE} \text{ have } k, q \sim N(\mu, \sigma^2 I); \mu = 0, \sigma = 1$$

$$\begin{aligned} \text{Var}(q^T k) &= \text{Var}(\sum_i q_i k_i) = \sum_i \text{Var}(q_i k_i) \\ &= d \cdot (\mathbb{E}[q_i^2] \text{Var}(k_i) + \mathbb{E}[k_i^2] \text{Var}(q_i)) \\ [\text{Attn}] &= \mathbb{E}[(q^T k)^2] - \mathbb{E}[q^T k]^2 = \text{Var}(q_i) \text{Var}(k_i) \\ &= M^T \Sigma_q M + M \Sigma_k M + \text{tr}(\Sigma_k \Sigma_q) \\ &= 2\sigma^2 \|MM^T\|_F + d\sigma^2 \end{aligned}$$

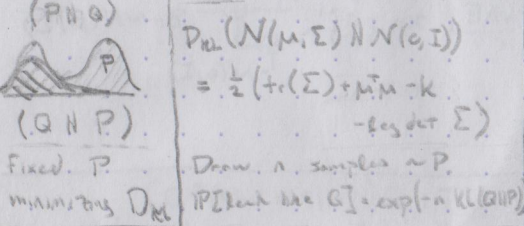
$$\text{So } \text{Var}(q^T k / d) = d/d^2 = 1/d \Rightarrow S = \sqrt{d}$$

$$\begin{aligned} &\text{sample } T_i \sim p(T) \\ &\text{For all } T_i: \text{Calc grad over } h \text{ samples} \\ &\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{T_i}(F_{\theta}) \\ &\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_i \mathcal{L}_{T_i}(F_{\theta_i}) \end{aligned}$$

$$\lim_{g \rightarrow 0} \frac{f(x)}{g(x)}; h(x) = f'(x)g(x) - f(x)g'(x)$$

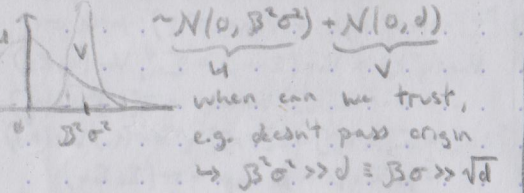
D. Fission $H(Y) = \mathbb{E}_Y[-\log(p(Y=K))]$
 $H(P, Q) = -\sum_K P(Y=K) \cdot \log(Q(K))$
 $J = \mathbb{E}_{P(x)}[\log(1/q(x))]$
 Negative log likelihood $= -\sum_x p(x) \log(q(x))$

$D_{KL}(P \parallel Q) = \sum_x p(x) \log\left(\frac{p(x)}{q(x)}\right)$
 $= \mathbb{E}_{P(x)}[\log p(x)] - \mathbb{E}_{P(x)}[\log q(x)]$
 $= -H(P) + H(P, Q)$



If data is "clean" (ie high SNR)
 (categories fairly distinct) then
 classifier has lots of margin; not need
 to pick out true characteristics over noise

Adversarial Ex $x \in \mathbb{R}^{d+1}$ st. $x_i \sim N(0, \sigma^2)$
 $w / \text{class}(x) = \text{sgn}(x_1)$ $x_1 \sim N(0, 1)$
 learned cls = $\text{sgn}(g^T x)$ $V_i = 2 \cdot d \cdot 1$
 $x \sim \beta, d: 1 \Rightarrow \beta x_1 + \sum x_i$



Consider $x' = x + Sx$ (S is adversary noise)
 st. $\text{score}(x) = x^T x + S x^T x$
 $= U + V + S(d^2 + d)$
 optm. $\sigma \gg \beta$ so nature of x not
 changed to x'
 Can flip IF $d \gg \beta \sigma = \text{Var}(U)$
 Sng $\sigma = d^a, \beta = d^b$ sng
 $a+b > 1/2$ so $a > 1/2$
 $a > b$ so $b > 1/4$
 $1 > a+b$

PE Applied to queries + heads
Relative: Each head has learnable
 bias term $\langle q_i, K_i \rangle + b_{j,i}$
 • Makes nearby cls bigger
 • Doesn't increase dimensionality
 like RoPE; no shw down

RoPE: $q = W_q x \rightarrow q = \Theta_e W_q x$
 where $\Theta_e = \text{diag}(\exp(i t \theta_1), \dots)$
 $\theta_1 = 10, \dots, 2\pi/d$ $\dots \exp(i t \theta_d/2)$
 • No learnable params
 • Depend just on relative locations
 Enl. Given $x \in \mathbb{R}^{2n}$
 Represent complex: $z = [x_1 + i x_{1+n}, \dots, x_n + i x_{n+n}]$
 $\text{RoPE}(z; t) = [z, \exp(i t \theta_1), \dots, \exp(i t \theta_n)]$
 $\rightarrow n$ learnable params

$x_1 \sim N(0, 1-\beta)$
 $x_2 \sim N(0, \beta)$
 βx_1 returns signal
 $x_t = \sqrt{\beta} x_{t-1} + \sqrt{1-\beta} \epsilon_{t-1}$
 $= \sqrt{\beta} (\sqrt{\beta} x_{t-2} + \sqrt{1-\beta} \epsilon_{t-2}) + \sqrt{1-\beta} \epsilon_{t-1}$
 $= \sqrt{\beta^2} x_{t-2} + \sqrt{\beta(1-\beta)} \epsilon_{t-2} + \sqrt{1-\beta} \epsilon_{t-1}$
 $= (\dots) + \sqrt{1-\beta} x_{t-1}$
 $= \sqrt{\prod_{s=1}^{t-1} \beta} x_0 + \sqrt{1-\prod_{s=1}^{t-1} \beta} \epsilon_0$
 $x_t \sim N(0, \beta^{t-1} (1-\beta) + (1-\beta^t))$
 $\therefore q(x_1, x_2) = N(x_1; \sqrt{\beta} x_2, (1-\beta)I)$

Adv $q^T K_i = (x + p)^T (x_2 + p)$
 $= x^T x_2 + p^T p_2 + x^T p_2 + x_2^T p$
 • only compare
 words/words pos/pos
 • Position term not
 dim or constructed st
 nearby positions not

Form Vectors: Let $z = f(x) = x - \mu \mathbb{1}$
 For $M = \frac{1}{2} \mathbb{1} \mathbb{1}^T$
 $g(x) = \frac{1}{2} x^T x = \frac{\partial x}{\partial x} \cdot \frac{\partial z}{\partial x} = (I - \frac{1}{2} \mathbb{1} \mathbb{1}^T) g(x)$
 Let $y = h(z) = \frac{z}{\sigma} = \sqrt{d} z / \|z\|_2$
 for $\sigma = \sqrt{d \|z\|_2^2 / d}$
 so $g(z) = \frac{\partial x}{\partial y} \cdot \frac{\partial y}{\partial z} = \sqrt{d} \left(\frac{\mathbb{1} \mathbb{1}^T z}{\|z\|_2} - \frac{z z^T}{\|z\|_2^3} \right) g(x)$
 For $z = \alpha \cdot z$ $= \sqrt{d} \left(I - \frac{z z^T}{\|z\|_2^2} \right) g(x)$
 $\|f(z)\|_2 = \alpha \|f(z)\|_2$
 $\|g(z)\|_2 = \|g(z)\|_2$
 $\|h(z)\|_2 = \|h(z)\|_2$
 $\|g(z)\|_2 = \frac{1}{2} \|g(z)\|_2$
 { Speed bump preventing
 exploding grads

Forgetting mitigation: Boils down to
 (1/step or better here) regularization
 1] Freezing. Early layers for if
 we're in new data domain
 Later layers if new task
 \rightarrow Could also lower LR on
 per-basis level

2] Rerun
 Combine CG data with new data
 If new task, graft head and
 keep OG for supervision

Linear Probing. Off-diagonal
 slope in CG from origin

RNN backprop: $h_t = W(u_t + h_{t-1})$

 $\frac{\partial y}{\partial W} = \sum_{t=1}^T \frac{\partial y}{\partial h_t} \frac{\partial h_t}{\partial W}$
 $y = W \cdot t = W u_2 + W^2 u_2 + W^3 u_1$
 $\frac{\partial y}{\partial W} = u_2 + 2 W u_2 + 3 W^2 u_1$
 $= t + r \frac{\partial y}{\partial S} + P \frac{\partial y}{\partial q}$
 $= u_2 + W u_2 + W^2 u_1 + (u_2 + q)(W)$
 $+ (u_1)(W^2)$
 $\frac{\partial y}{\partial t} = W = \frac{\partial y}{\partial S}$
 $\frac{\partial y}{\partial r} = \frac{\partial y}{\partial S} \cdot W; \frac{\partial y}{\partial q} = \frac{\partial y}{\partial t} \cdot 1$

Finetuning
 1] Prompt Tuning. Essentially make
 first layers keys learnable
 • Vectors more free than tokens (english)
 • Works best on large models

$p(x) = \int p(x, z) dz$ intractable
 $= p(x, z) / p(z|x)$ unknown,
 surrogate model: what we want
 to learn!
 $p(x|z) \approx q(x|z)$

ELBO maximize
 $\log p(x) = \log p(x) \int q(z|x)$
 $= \mathbb{E}_{q(z|x)} [\log p(x)]$
 $= \mathbb{E} \left[\log \frac{p(x, z)}{p(z|x)} \cdot \frac{q(z|x)}{q(z|x)} \right]$
 $= \mathbb{E}_q \left[\log \frac{p(x, z)}{q(z|x)} \right] + D_{KL}(q(z|x) \parallel p(z|x))$
 $\geq \mathbb{E}_q \left[\log \frac{p(x, z)}{q(z|x)} \right] \geq 0$
 $= \mathbb{E} \left[\log p(x|z) \right] - D_{KL}(q(z|x) \parallel p(z|x))$
 reconstruction quality learned dist. diff



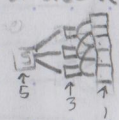
If blurry Q and $q(z|x)$ is far from $p(z)$
 learned map eval sampler

Conv vs Input: $C \times H \times W$ ($C \times K \times K$)
 Filter size (K), padding (P)
 Stride (S), number Filters (F)
 Weights: $K^2 \times F \times C$ Learnable
 # biases: F parameters

$W' = \lfloor (W - K + 2P) / S \rfloor + 1$, $H' = \dots$
 $C' = F$

Receptive Field: $L(K-1)+1 = O(LN)$

→ L successive conv layers of size K
 → Formula for stride=1
 → Stride makes increase faster than linear



→ Then 3x3, same receptive F. → 7x7
 • BUT former has non-lin; more
 • $3 \times (C \times (3 \times 3 \times C)) = 27C^2$ expensive
 v.s. $C \times (7 \times 7 \times C) = 49C^2$ parameters

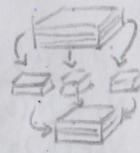
• General: $r_{i+1} = S_i \times r_i + k_{i+1} - S_i$

Weight Sharing: Rather than per-pixel weights, we have per-patch weights that we apply over whole img.

- Less parameters / slow blowup
- Introduces Translation Invariance

Data Augmentation: Introduce insignificant (wrt classification) perturbations to data s.t. model learns these invariances
 • SGD moves along weighed avg. of img patches, so augmentation can be thought of as regularizing noise

Depthwise: Filters unique to each channel of input
 • Less expensive, now quite common



Peeling! $W' = (W - K) / S + 1$
 Receptive Field $\sim K$ (asm $S=1$)
 Spatial resolution $\sim 1/K^2$ (in 2D)
 → Why we double output channels

Global Avg: Over all channels

- Done prior to SoftMax, rather than Flatten: reduce from $H \times W \times C$ to C
- Loses spatial info; okay for classification (present or not, binary) but not segmentation
- Enables us to use sized input: output solely depends on C .

Max: Take highest value of $K \times K$ pool

- Induces learning "invariant" features
- High-pass: retains "strongest" signal

$\frac{\partial x}{\partial x_{ij}} = \sum_{yab} dy_{ab} \cdot \frac{\partial y_{ab}}{\partial x_{ij}}$
 $\frac{\partial y_{ab}}{\partial x_{ij}} = dy_{ab} \times \mathbb{1}\{x_{ij} = \max\{x_{ij}, \dots\}\}$

Avg: Reduces spatial domain by downsampling; retains some info about distribution of patches values

- Better at capturing smooth transitions / subtle changes than Max pool

$\frac{\partial x}{\partial x_{ij}} = \sum_{yab} dy_{ab} \cdot \frac{\partial y_{ab}}{\partial x_{ij}} = dy_{ab} \cdot \frac{1}{4}$

UNET For pixel-level classification

- Double channels when downsampling s.t. "total info" is const.
- Low-res, high-d (middle) integrates contextual info from whole img
- Up-sampling turns low-res feature vectors into higher res per-pixel preds
- Combine with details (fine info) from earlier

Normalizing! IF we move somewhere, we want a reason - not b.c. it's big.
 • Note: Adam already normalizes grad. over every param. by the avg. across many batches

- Go to $N(0,1)$: $(x - \mathbb{E}[x]) / \sigma[x]$
 → operations can be in-place by weights + biases; no reduction in expressive power



- 1) Batch: Norm. each channel in a batch by μ, σ of each feature (all channel pixels)



- Treats channels as unrelated
- Backprop grad over whole batch; non-trivial implement.
- Test time: Use stored μ, σ
- Suitable for convnets, not FC layers (no symm/ spatial)
- Tackles internal covariate shift



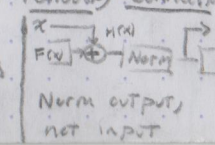
- 2) Layer: Norm. over all pixels in each channel
- E.g. later layers in network when features should be "learned" and thus sizes should be alike

- Can group-norm (only some of channels) too
- Faster than Batch

- 3) Instance: Norm across each channel in each img/P
- Agnostic to contrast of input

- Practically: Inits + Norm + Conv. isn't enough for training
- Grads become too small
- Backprop. signal noised
- Solved w/ skip (or residual) connections

- Decouple behavior: enables fixing paths
- Tackles internal covariate shift



Inits! Metric: Don't want dead ReLU
 → Wasted computation
 • Usually unable to recover weights:

- 0: Can't backprop more than 1 layer; always zero gradient
- Ideally, inputs follow $N(0,1)$
- 1) Xavier: $\sigma^2 = 1/d$, d = Fan-in
- Appropriate for sigmoid/tanh; not ReLU
- Saturating fctns; go to some $\neq 0$ const
- 2) He: $\sigma^2 = 2/d$

- Before, half of ReLUs begin off so Fan-in is halved
- Solves Xavier's vanishing grad problem

Biases:

- Xavier with $b=1$: treat as weight
- 0 (optimizer to set within first cycles)
- C.o. Empirically better than 0
- Any small random #. (doesn't matter as we expect optimizer to set)

convnet Expand to higher dim, then compress/shrink down

- Each layer modifies input, rather than completely transforming it.
- Classification obj. Features inherently low-dim. (e.g. $[class, x, y]$)
- So it can "fit" in few # channels
- Don't ReLU output: destroys too much info over low-dim spaces

- Weight Decay: Account for batch selection variance from denoising
- Effectively smoothing over batches
- Const. is good, cosine better
- Label smooth: Cross-Entropy of SoftMax will rarely have '1' for label

- LS relaxes constraint; ok w/ 0.9
- s.t. fail cases are focuses

$\frac{\partial H}{\partial x} = \frac{\partial F}{\partial x} + I$

Ensures never too small

Optimization Algos | GD

$$x_{t+1} = x_t - \eta \nabla_x \frac{1}{2} \|Ax - y\|_2^2$$

$$= (I - \eta A^T A) x_t + \eta A^T y$$

$$x_{t+1} - x^* = (I - \eta A^T A)^t (x_0 - x^*)$$

$$\|x_{t+1} - x^*\| \leq \|I - \eta A^T A\|_2^t \|x_0 - x^*\|$$

$$\sigma_{\min}\{I - \eta A^T A\}^t$$

$$w_t = (I - \eta A^T A)^t w_0 + \sum_{i=0}^{t-1} (I - \eta A^T A)^i x^T y$$

- Learning rate:
 - σ_{\min} limits how fast we can converge without oscillations/instability
 - σ_{\min} requires large learning rate
 - No osc. if $(1 - 2\eta \sigma_{\min}^2) > -1$
 - $\Rightarrow \eta < 1/\sigma_{\min}^2$
- Fastest rate: $(1 - 2\eta \sigma_{\min}^2) = (-1 - 2\eta \sigma_{\max}^2)$
- $\sigma = \sigma_{\max}/\sigma_{\min} \Rightarrow \eta^* = 1/\sigma_{\min}^2 \sigma_{\max}^2$
- Rate = $(\sigma^2 - 1)/(\sigma^2 + 1)$

Momentum: We want our direction to not solely depend on singular values, but also on past steps.

- By dampening oscillations, we effectively downscale LR for largest sing. values, enabling higher η

$$w_{t+1} = w_t - \eta \alpha_{t+1} \nabla \mathcal{L}(w_t)$$

$$\alpha_{t+1} = (1 - \beta) \alpha_t + \beta (\nabla \mathcal{L}(w_t))$$

\Rightarrow Exponential moving avg. of grad

\Rightarrow Still represents just $\sigma_{\min}/\sigma_{\max}$

ADAM: Pretend coords are indep directions

eg. using different η each direction

$$\alpha_{t+1} = (1 - \beta) \alpha_t + \beta (\nabla \mathcal{L}(w_t))$$

$$v_{t+1} = (1 - \beta') v_t + \beta' \times \left\{ \frac{1}{\sqrt{2}} \frac{\partial}{\partial w_i} \mathcal{L}(w_t) \right\}$$

$$w_{t+1} = w_t - \eta \alpha_{t+1} / \sqrt{v_{t+1}} + \epsilon$$

- $\beta' \ll \beta$ s.t. variation in normalization doesn't drive dynamics; v_t decays slowly
- Size of grads = V ; ensured by small steps
- Ideally wld find 2nd order deriv; practically infeasible

Training/Debug

- Train err. of bigger network is worse than smaller / slower
- \hookrightarrow Underfitting: grads too small
 - \square Inits to He/Xavier; non-zero
 - \square Skip connections
 - \square Normalizing layers

For \dots in range (epochs):

batch-loss = []

For batch in dataset:

[]

optimizer.zero_grad()

loss.backward()

optimizer.step()

batch-loss.append(loss.item())

training-loss.append(np.mean(batch-loss))

Adam introduces bias - not good w/ many possible solns

AdamW - Normalize st. sharp

$$\hat{m}_t = m_t / (1 - \beta^t)$$

\hookrightarrow smoother

$$\hat{v}_t = v_t / (1 - \beta'^t)$$

\hookrightarrow goes away after many iters

GNNs | Generalization of CNNs

- Graph-level: Look at all graph vts and make statement/pred. about whole
- Node-level: "statement about nodes"
 - \hookrightarrow Duality: Can be defined over edges

Weight Sharing: Must not have dependency on neighbor cnt or order

Aggregate operations: tend to blur, rather than sharpen

$\hookrightarrow \Sigma / \max / \min / \text{TT} / \text{LS} / \text{Variance}$

$f_{w_1}(m_e, \sum_{w_2} S_{w_2}(m_e, \text{them}) g_{w_3}(\text{them}))$

Attention Importance

Report | Input standardization is to

- batch-layer norm as data avg. is to dropout
- Fully-connected layers have no spatial awareness (unlike conv. nets) so normalization less effective
- Implementation: Randomly "kill" a unit (set to 0, don't upd weights)
 - \hookrightarrow Eval: scale weights by $(1 - p)$ [if not done when training]
- Ensembles: wholeistic network: no one unit should be solely responsible
 - \hookrightarrow Encourage redundant activations, i.e. lower-dimension; discourage sparsity (e.g. one large singular value)
- Stochastic Depth Regularization:
 - Drop residual/skip blocks
 - \hookrightarrow Used in CNNs
- Drop channels: Image already captures some redundancy of pixels
 - \hookrightarrow Causes jagged grids
- Dropout params: don't depend on fine-grain details of pixels
 - \hookrightarrow Only for Fully-connected layers

- Global: Inter-layer communication
 - \hookrightarrow Notion of time/cycle invariance
 - Messages either go from node to neighbors, or to all [Larger receptive field]
 - \hookrightarrow Distinguish w/ global weights
- Many global nodes \equiv single node, multiple channels

Peeling: Task specific

- Content-specific: Alter topology, i.e. add edges till FC
- Cluster: Eg random walk, finds "central" nodes

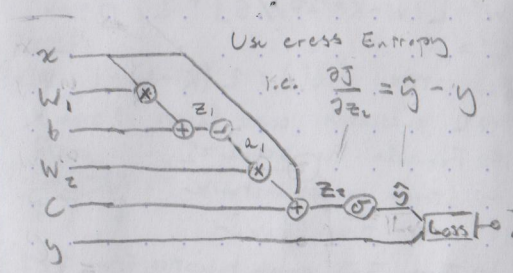
Vector Calc

$$y = W^{(l)} \text{ReLU}(W^{(l-1)} x + b) + c$$

$$\nabla_{W^{(l)}} y = \text{ReLU}(0, W^{(l-1)} x + b)$$

$$\nabla_{b^{(l)}} y = \mathbb{1}\{W^{(l-1)} x + b > 0\} \times W^{(l)}$$

$$\nabla_{W^{(l-1)}} y = \mathbb{1}\{W^{(l-1)} x + b > 0\} \times W^{(l)} x$$



$$\tilde{y} = \sigma(W_2 \text{ReLU}(W_1 x + b) + b_2 + x)$$

$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial z_2} \frac{\partial z_2}{\partial w_2} = (\tilde{y} - y) \cdot a_1$$

$$\frac{\partial J}{\partial c} = \frac{\partial J}{\partial z_2} \frac{\partial z_2}{\partial c} = \tilde{y} - y$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial z_1} \frac{\partial z_1}{\partial w_1} = \tilde{y} - y$$

$$= [(\tilde{y} - y) W_2 \frac{\partial z_1}{\partial w_1}] x$$

$$\frac{\partial J}{\partial x} = [N] W_1 + (\tilde{y} - y)$$

$$\sigma(x) = 1/(1 + e^{-x}) \quad \sigma'(x) = 1/(1 + e^{-x}) \cdot (-e^{-x})$$

$$= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}}$$

$$= \sigma(x) \cdot (1 - \sigma(x))$$

$$S_i^L = S_i^{L-1} + W_i (\sum_{j=1}^n \sigma(W_{ij} S_j^{L-1}))$$

Let $m_{ij} = S_j^{L-1}$

$$S_i^L = S_i^{L-1} + \frac{W_i}{3} (\sigma(W_{i1} S_1^{L-1}) + \sigma(W_{i2} S_2^{L-1}) + \sigma(W_{i3} S_3^{L-1}))$$

$$P(X|Y) = P(Y|X) P(X) / P(Y)$$

$$N_i \sim \mathcal{N}(0, \sigma^2 I) : \sigma^2 I = E[N_i N_i^T]$$

$$\text{Var}(y) = E[y^2] - E[y]^2$$

$$\text{Var}(x) = \sigma(x)^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

STD