

HW1 整理

Topic 1: Why learning can't be too big / lr & o

Simplest Formulation: $\sigma w = y$ (Scalar equation)

Loss: $L(w) = (y - \sigma w)^2$, η is learning rate

$$\begin{aligned} \text{GD: } w_{t+1} &= w_t + 2\eta \sigma (y - \sigma w_t) \\ &= (1 - 2\eta \sigma^2) w_t + 2\eta \sigma y \end{aligned}$$

Important thought!

We know $w^* = y/\sigma$, and want $|w_t - w^*|$ converge!

$$w_{t+1} - \frac{y}{\sigma} = (1 - 2\eta \sigma^2) (w_t - \frac{y}{\sigma})$$

$$w_{t+1} = \frac{y}{\sigma} + (1 - 2\eta \sigma^2) (w_0 - \frac{y}{\sigma})$$

$$\text{We need } |1 - 2\eta \sigma^2| < 1, \quad \therefore \eta < \frac{1}{\sigma^2}$$

If we want: $|w_t - w^*| < \varepsilon |w^*|$, and assume $w_0 = 0$

$$\text{Then } |(1 - 2\eta \sigma^2)^t| < \varepsilon \quad t > \frac{\ln \varepsilon}{\ln(1 - 2\eta \sigma^2)}$$

The above is scalar problem, what if: $\begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_s \end{bmatrix} \begin{bmatrix} w[1] \\ w[2] \end{bmatrix} = \begin{bmatrix} y[1] \\ y[2] \end{bmatrix}$
 (And assume $\sigma_1 > \sigma_s$)

Assume this is $\sum \vec{w} = \vec{y}$, then $L(w) = \|y - \sum w\|_2^2$, GD:

$$\vec{w}_{t+1} = (I - 2\eta \sum) \vec{w}_t + 2\eta \sum \vec{y} \quad \text{similar to convergence analysis above:}$$

$$\left\{ \begin{array}{l} |1 - 2\eta \sigma_1^2| < 1 \\ |1 - 2\eta \sigma_s^2| < 1 \end{array} \right. \quad \therefore \eta < \min \left\{ \frac{1}{\sigma_1^2}, \frac{1}{\sigma_s^2} \right\} = \frac{1}{\sigma_1^2}$$

If we want fastest convergence speed? $|1 - 2\eta \sigma_1^2| = |1 - 2\eta \sigma_s^2|$

$$\eta' = \frac{1}{\sigma_1^2 + \sigma_s^2}$$

If $\sum \in \mathbb{R}^{k \times k}$, diagonal ($\neq \sigma_s$), then $\eta < \min \left\{ \frac{1}{\sigma_k^2} \right\}$,

$$\eta' = \frac{1}{\sigma_1^2 + \sigma_s^2} ? \quad \text{Not exactly, hard to tell}$$

Finally, for trivial $Xw = y$, we can SVD: $\underbrace{U} \Sigma \underbrace{V^T w}_{} = y$

We can: $\sum (V^T w) = U^T y = \tilde{y}$, let $V^T w = \hat{w}$, i.e., $w = V \hat{w}$

$\sum \hat{w} = \tilde{y}$, and it connects to the previous analysis

Topic 2: Why SGD can work?

there are solutions ↑

Suppose: $Xw = y$, $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$, $y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$, and X has full rank, $d > n$, $x_i \in \mathbb{R}^d$

Assume $w_0 = 0$, we know $w^* = X^T(XX^T)^{-1}y$, rewrite $w' = w - w^*$
 \therefore Goal is $Xw' = 0$, $w_0' = -w^*$

Let's Use SVD: $\Sigma = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \vec{o} = [\Sigma_1, \vec{o}]$, then:

$$Xw' = U\Sigma V^T w' = U[\Sigma_1, \vec{o}_{n \times d-n}] V^T w' = [\tilde{X}, \vec{o}] V^T w' = [\tilde{X}, \vec{o}] w''$$

$$\therefore w'' = -V^T X^T (XX^T)^{-1} y = \dots = -[\Sigma_1, V^T y] \xrightarrow{\text{all zeros}} \begin{bmatrix} \vec{o} \\ \vec{o}_{d-n} \end{bmatrix}$$

Now we know: $[\tilde{X}, \vec{o}] w'' = 0$ is Goal, and $(d-n)$ col/rows will not work in this

$\therefore \tilde{X}$ is: $U\Sigma_1$, \tilde{X} is also $[\tilde{x}_i^T]_{i \in [1, n]}$, we have: $\tilde{x}_i^T \vec{o} = 0$

Solving $\tilde{X}\vec{o} = 0$ is exactly calculating the original problem, since $\tilde{X}\vec{o}$ has linear relationship with X and w . Now: $\tilde{X}\vec{o} = 0$, $L(\vec{o}) = (\tilde{X}\vec{o})^T \tilde{X}\vec{o}$

SGD: $\tilde{w}_{t+1} = \tilde{w}_t - \eta \tilde{x}_{it} \tilde{x}_{it}^T \tilde{w}_t$, It $\in [1, n]$, one sample

Try to decompose: $L(\tilde{w}_{t+1}) = L(\tilde{w}_t) + A + B$, $L(\tilde{w}) = (\tilde{X}\tilde{w})^T \tilde{X}\tilde{w}$

Bring in, and $\left| \begin{array}{l} A = 2\tilde{w}_t^T \tilde{X}^T \tilde{X} (\tilde{w}_{t+1} - \tilde{w}_t) \\ B = (\tilde{w}_{t+1} - \tilde{w}_t)^T \tilde{X}^T \tilde{X} (\tilde{w}_{t+1} - \tilde{w}_t) \end{array} \right.$ A. B are in form of $(\tilde{w}_{t+1} - \tilde{w}_t) \& \tilde{w}_t$

Insight: $E(A|\tilde{w}_t)$, bring in A, $\tilde{w}_{t+1} = \tilde{w}_t - \eta \tilde{x}_{it} \tilde{x}_{it}^T \tilde{w}_t$

$$E(A|\tilde{w}_t) = 2\tilde{w}_t^T \tilde{X}^T \tilde{X} E[-\eta \tilde{x}_{it} \tilde{x}_{it}^T \tilde{w}_t | \tilde{w}_t] = -\frac{4}{n} \eta \tilde{w}_t^T \tilde{X}^T \tilde{X} \tilde{X}^T \tilde{X} \tilde{w}_t$$

$$\leq -\frac{4}{n} \eta \min_{\tilde{w}_t^T \tilde{X}^T \tilde{X} \tilde{w}_t} \tilde{w}_t^T \tilde{X}^T \tilde{X} \tilde{w}_t = -\frac{4}{n} \eta \min_{\tilde{w}_t} L(\tilde{w}_t) = -C_1 \eta L(\tilde{w}_t), C_1 > 0$$

Similarly $E(B|\tilde{w}_t) \leq \frac{4}{n} \eta^2 \max_{\tilde{w}_t} p^2 L(\tilde{w}_t) = C_2 \eta^2 L(\tilde{w}_t)$, $C_2 > 0$
 \downarrow largest norm of \tilde{x}_i

$$\therefore E(L(\tilde{w}_{t+1})|\tilde{w}_t) = E(L(\tilde{w}_t)) + E(A|\tilde{w}_t) + E(B|\tilde{w}_t)$$

$$\leq (1 - C_1 \eta + C_2 \eta^2) L(\tilde{w}_t)$$

So, we can set η s.t. $1 - C_1 \eta + C_2 \eta^2 < 1$, if $\eta < \frac{C_1}{C_2}$

To achieve fastest convergence, $\eta^* = \frac{C_1}{2C_2}$

Topic 3. Relationship between Noise and Regulation

Consider $(X_i, y_i)_{i=1}^m$, $X_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$, assume we add noise:

$$\tilde{X}_i = X_i + N_i, \quad N_i \sim \mathcal{N}(0, \sigma^2 I_n)$$

For noisy data matrix, we want $\underset{w}{\operatorname{argmin}} \mathbb{E} [\|\tilde{X}w - y\|_2^2]$

$$\begin{aligned} \mathbb{E} [\|\tilde{X}w - y\|_2^2] &= \mathbb{E} [\sum_{i=1}^m ((X_i + N_i)^T w - y_i)^2] \quad \mathbb{E}: \sigma^2 I_n \\ &= \sum_{i=1}^m \mathbb{E} [(X_i^T w - y_i)^2] - 2 \mathbb{E} [(\underbrace{N_i^T w}_{\text{sample in } t\text{-th step}})(X_i^T w - y_i)] + \mathbb{E} [w^T \underbrace{N_i N_i^T w}_{\sigma^2 I_n}] \\ &= \|Xw - y\|_2^2 + m\sigma^2 \|w\|_2^2 \end{aligned}$$

If in SGD t -th step, we have N_t to \tilde{X}_t , then:

$$\frac{\partial L}{\partial w} = (w^T \tilde{X} - y) \tilde{X} = w(x^2 + 2xN_t + N_t^2) - y(x + N_t)$$

$$W_{t+1} = W_t - \eta \nabla_w L(W_t) = W_t (1 - \eta (x^2 + 2xN_t + N_t^2)) + \eta y (x + N_t)$$

$$\mathbb{E}(W_{t+1}) = \mathbb{E}(W_t) (1 - \eta (x^2 + \boxed{0})) + \eta y x$$

$$* \operatorname{Var}(N_t) = \mathbb{E}(N_t^2) - [\mathbb{E}(N_t)]^2 = \mathbb{E}(N_t^2) = \sigma^2$$

$$\therefore \eta \text{ should satisfy: } |1 - \eta (x^2 + \sigma^2)| < 1, \quad 0 < \eta < \frac{1}{x^2 + \sigma^2}$$

Topic 4. Another perspective of: $(X^T X + \lambda I)^{-1} X^T y = X^T (X X^T + \lambda I)^{-1} y$

We want build r.v. Z , s.t., $Z = W - AY$, but Z, Y are not correlated:

$$\begin{aligned} \operatorname{Cov}(Z, Y) &= \dots = \mathbb{E}[(W - \mu_W)(Y - \mu_Y)^T] - A \cdot \mathbb{E}[(Y - \mu_Y)(Y - \mu_Y)^T] \\ &= \sum w_y - A \cdot \sum r_y \quad \therefore A = \sum w_y \cdot \sum r_y^T \end{aligned}$$

$$\therefore W = Z + AY, \quad \mathbb{E}[W|Y=y] = \mathbb{E}[Z|Y=y] + Ay, \quad \text{bring in } Z = W - AY$$

$$= \dots = \mu_W + \sum w_y \sum r_y^T (y - \mu_y) \quad Z, Y \sim N, \text{ so correlation}$$

Assume: $\mu_W, \mu_Y = 0$, i.e., W, y has zero mean \Rightarrow independent

$$\therefore \mathbb{E}[W|Y=y] = \sum w_y \sum r_y^T y, \quad \text{assume } *: Y = XW + \sqrt{N}N. \text{ then}$$

$$\sum w_y = \mathbb{E}[W(XW + \sqrt{N}N)^T] = X^T \quad | \quad W \sim N(0, I)$$

$$\sum r_y^T = \mathbb{E}[(XW + N)(XW + N)^T] = XX^T + \lambda I$$

$$\therefore \hat{W} = \mathbb{E}[W|Y=y] = X^T (XX^T + \lambda I)^{-1} y \quad \text{Constrained! As in ridge regression, } \|W\|_2$$

* MAP version to explain ridge: prior $W \sim N(0, I)$, view random Y as $Y = X^T W + \sqrt{N}N$, N is i.i.d. $\sim N(0, I)$. I should refer to HWO 整理

* Why mean? A normal r.v. has its density maximized at its mean

HW0, HW2 整理

Topic 1. MAP perspective of Ridge Regression. We can think $W \sim N(0, I)$.

and view Y as being generated using $Y = X^T W + \sqrt{\lambda} N$, where $N \sim N(0, I)$

Vector version: $Y = XW + \sqrt{\lambda} N$

$$\therefore \text{MAP}(w | Y=y) = \underset{w}{\operatorname{argmax}} f(w | Y=y) = \underset{w}{\operatorname{argmax}} \frac{f(w, y)}{f(y)}$$

$$\Rightarrow \underset{w}{\operatorname{argmax}} f(w) \cdot f(y|w) = \underset{w}{\operatorname{argmax}} f(w) \cdot \prod_{i=1}^n f(y_i|w)$$

$$\text{For } N(0, I), \text{ PDF: } f_Z(z) = \frac{e^{-\frac{z^2}{2}}}{\sqrt{2\pi}} \quad \therefore \underset{w}{\operatorname{argmax}} \frac{e^{-\frac{\|w\|^2}{2}}}{\sqrt{2\pi}} \prod_{i=1}^n \frac{e^{-\frac{(y_i - x_i^T w)^2}{2}}}{\sqrt{2\pi}}$$

$$\stackrel{*}{=} \underset{w}{\operatorname{argmin}} \|w\|_2^2 + \frac{1}{\lambda} \|y - Xw\|_2^2 \stackrel{*}{=} \underset{w}{\operatorname{argmin}} \lambda \|w\|_2^2 + \|Xw - y\|_2^2$$

*: negation, $\operatorname{argmax} \rightarrow \operatorname{argmin}$, take log

Topic 2. Originally in optimizer: $\underset{\|w\| < \eta}{\operatorname{argmin}} \langle \nabla_w L(w), \Delta w \rangle$

We set constraint as: $\|\Delta w\| < \eta$

What if we take it as an optimization problem term?

$$\Rightarrow u = \underset{\Delta \theta}{\operatorname{argmin}} g^T \Delta \theta + \frac{1}{\alpha} d(\Delta \theta), \text{ where } g = \nabla f(\theta), d \in \mathbb{R}^{\dim(\theta)} \rightarrow \mathbb{R}^+$$

$$\text{Then: } \nabla_{\Delta \theta} u = g + \frac{2}{\alpha} \Delta \theta = 0 \quad \text{if } d \text{ is the euclidean distance (L2-Norm)} \\ \Delta \theta^* = -\frac{\alpha}{2} g$$

$$\text{let } M = \max_i |\Delta \theta_i|$$

$$\text{But if } d \text{ is } \infty \text{ norm? } u = \underset{\Delta \theta}{\operatorname{argmin}} g^T \Delta \theta + \frac{1}{\alpha} \max_i |\Delta \theta_i|$$

$$\therefore u = \underset{\Delta \theta, M}{\operatorname{argmin}} g^T \Delta \theta + \frac{1}{\alpha} M^2, \quad |\Delta \theta_i| \leq M$$

Since $|\Delta \theta_i| \leq M$, we want $\sum_{i=1}^n g_i^T \Delta \theta_i$ as negatively small as possible, $|\Delta \theta_i| = M$

$$\therefore \Delta \theta_i = -M \cdot \text{sign}(g_i), \quad u = -M \|g\|_1 + \frac{1}{\alpha} M^2 \quad \text{for } \forall i$$

$$\Rightarrow \nabla_M u = -\|g\|_1 + \frac{2}{\alpha} M = 0, \quad M = \frac{\alpha}{2} \|g\|_1$$

$$\therefore u = -\frac{\alpha}{2} \|g\|_1 \cdot \text{sign}(g), \quad \Delta \theta^* = \left\{ \Delta \theta_i = -M \text{sign}(g_i) \right\}_{i=1}^{\dim(\Delta \theta)}$$

HW3 整理

Topic 1. Maximal Update Parameterization

(a). Consider dense layer $X \in \mathbb{R}^{d_1}$, $W \in \mathbb{R}^{d_2 \times d_1}$, W is initialized with i.i.d standard Gaussian entries, and entries of X are sampled i.i.d from unit Gaussian.

Then: $y = Wx$, its squared RMS norm? consider $W_i \in \mathbb{R}^{1 \times d_1}$, $y_i = W_i x$
 $y_i = \sum_{j=1}^{d_1} w_{ij} x_j$, $w_{ij}, x_j \sim N(0, 1)$.

Regard $\{w_{ij} x_j\}$, $\text{Var}(w_{ij} x_j) = \text{Var}(z) = E(z^2) - [E(z)]^2$

$E(z) = E(w_{ij}) E(x_j)$ (since w_{ij} and x_j is individually independent)

$$E(z^2) = E(w_{ij}^2 x_j^2) = E(w_{ij}^2) E(x_j^2) = 1 \cdot 1 = 1$$

$$\therefore \text{Var}(z) = 1 \Rightarrow \text{Var}(y_i) = d_1 \cdot \text{Var}(z) \stackrel{*}{=} d_1$$

* If $S = \sum_{i=1}^n X_i$, then $\text{Var}(S) = \text{Var}(\sum_{i=1}^n X_i) = \sum_{i=1}^n \text{Var}(X_i)$ if X_i i.i.d

$$\therefore \|y\|_{\text{RMS}}^2 = \frac{1}{d_1} \sum_{i=1}^{d_1} y_i^2 = \frac{1}{d_1} \sum_{i=1}^{d_1} [y_i - E(y_i)]^2 = \frac{1}{d_1} \cdot d_1 \cdot \text{Var}(y_i) = d_1$$

$$\therefore E(\|y\|_{\text{RMS}}^2) = O(d_1) \|X\|_{\text{RMS}}^2, \text{ so } W \text{ should be } \frac{1}{\sqrt{d_1}} W_{\text{previous}}$$

$$\Delta \text{Var}(z) = E(x_j^2) E(w_{ij}^2) = 1/d_1, \text{Var}(w_{ij}) = \frac{1}{d_1}, w_{ij} \sim N(0, 1/d_1)$$

$$\Rightarrow \text{If } S = tX \Rightarrow \text{Var}(S) = t^2 \text{Var}(X).$$

These two seems to contradict, but they are both true

(b). Consider using RMS norm in weight update: $W_{t+1} \leftarrow W_t + \eta \text{ sign}(g_i x_i^\top)$
 $y' = \text{sign}(g_i x_i^\top) x_i$, what is $E(\|y'\|_{\text{RMS}}^2)$? $g_i \in \mathbb{R}^{d_2}$, $x_i \in \mathbb{R}^{d_1}$.

Dimension analysis: $\text{sign}(g_i x_i^\top) \in \mathbb{R}^{d_2 \times d_1}$, $x_i \in \mathbb{R}^{d_1}$, $y' \in \mathbb{R}^{d_2}$

$$\therefore y'_j = \underbrace{\text{sign}(g_{ij} x_i)}_{\text{i.i.d.}} x_i, E(y'^2_j) = E\left(\eta^2 \left[\sum_{i=1}^{d_1} \text{sign}(g_{ij} x_i) x_i\right]^2\right) = \eta^2 \left[E\left(\sum_{i=1}^{d_1} |x_i|\right)\right]^2$$

$$\therefore E(\|y'\|_{\text{RMS}}^2) = \frac{1}{d_2} \sum_{j=1}^{d_2} E(y'^2_j) = \frac{1}{d_2} \cdot d_2 \cdot \eta^2 \left[E\left(\sum_{i=1}^{d_1} |x_i|\right)\right]^2 = \eta^2 d_1^2$$

\therefore For $\Delta y, y'$, scale is d_1 !! Hence need to multiply Δy by $\frac{1}{d_1}$

* sign will force all x_i to face toward the same +/- direction with g_{ij} , and square will erase +/- $\Rightarrow 1$!!

Topic 2: Maximal Update Parameterization Research

Consider $\Delta W \in \mathbb{R}^{n \times n_{l-1}}$, $W \in \mathbb{R}^{n \times n_l}$, if

$$\|\Delta W\|_2 = \Theta(\sqrt{\frac{n_l}{n_{l-1}}}) \quad \|\Delta w\|_2 = \Theta(\sqrt{\frac{n_l}{n_{l-1}}}).$$

then suppose $h(x) \in \mathbb{R}^{n_l}$ as the feature of input x at layer l .
 $\Delta h_l(x) \in \mathbb{R}^{n_l}$ denote the change, then:

$$\|h\|_2 = \Theta(\sqrt{n_l}) \quad \|\Delta h\|_2 = \Theta(\sqrt{n_l}), \text{ i.e., their RMS norm is } \Theta(1)$$

\Rightarrow Each entry of hidden vectors moves and remains by $\Theta(1)$

Proof: If condition is satisfied: then $\|W\|_{\text{RMS}} \Rightarrow \text{RMS} = \Theta(1)$.

$$* \quad \|W\|_2 = \sqrt{\frac{n_l}{n_{l-1}}} \quad \|W\|_{\text{RMS}} = \Theta(\sqrt{\frac{n_l}{n_{l-1}}})$$

$$\therefore \|h_l(x)\|_2 = \|W_l h_{l-1}(x)\|_2 \leq \|W_l\|_2 \|h_{l-1}(x)\|_2 = \Theta(\sqrt{\frac{n_l}{n_{l-1}}}) \cdot \Theta(\sqrt{n_{l-1}}) = \Theta(\sqrt{n_l}).$$

\Rightarrow if $h_0(x)$ is initialized properly, above induction holds true

$$\begin{aligned} \|\Delta h_l(x)\|_2 &= \|W_l \Delta h_{l-1}(x) + \Delta W_l h_{l-1}(x) + \Delta W_l \Delta h_{l-1}(x)\|_2 \\ &\leq \Theta(\sqrt{\frac{n_l}{n_{l-1}}}) \Theta(\sqrt{n_{l-1}}) + \Theta(\sqrt{\frac{n_l}{n_{l-1}}}) \Theta(\sqrt{l-1}) + \Theta(\sqrt{\frac{n_l}{n_{l-1}}}) \Theta(\sqrt{l-1}) = \Theta(\sqrt{n_l}) \end{aligned}$$

Topic 3: Policy Gradient. & Reparameterization Gradient Estimator.

Consider $\mathcal{F}(\theta) = \mathbb{E}_{x \sim p_\theta}[f(x)]$, its derivative to θ is a challenge!
 how x is sampled. Sort 'eval' func

Suppose x is k -D multivariate Gaussian r.v., $\theta = \mu$, then.

$$p(x) = (2\pi\sigma^2)^{-\frac{k}{2}} \exp(-\|x - \mu\|_2^2 / (2\sigma^2)).$$

$$\nabla_\mu \mathcal{F}(\mu) = \underbrace{\nabla_\mu \int_K (2\pi\sigma^2)^{-\frac{k}{2}} \exp(-\|x - \mu\|_2^2 / (2\sigma^2)) f(x) dx}_{= \int_K (2\pi\sigma^2)^{-\frac{k}{2}} \exp(-\|x - \mu\|_2^2 / (2\sigma^2)) \frac{x - \mu}{\sigma^2} f(x) dx}$$

$$= \int_K (2\pi\sigma^2)^{-\frac{k}{2}} \exp(-\|x - \mu\|_2^2 / (2\sigma^2)) \frac{x - \mu}{\sigma^2} f(x) dx = \mathbb{E}_{x \sim p_\mu(x)} \left[\frac{x - \mu}{\sigma^2} f(x) \right]$$

For any p_θ : $\nabla_\theta \mathcal{F}(\theta) = \int_K \nabla_\theta p_\theta(x) f(x) dx = \int_K p_\theta(x) \frac{\nabla_\theta p_\theta(x)}{p_\theta(x)} f(x) dx$

Else, decouple: $x = g(z, \theta)$ Distribution unaffected by θ $\rightarrow \mathbb{E}_{x \sim p_\theta} [\nabla_\theta \log p_\theta(x) f(x)]$

Then $\nabla_\theta \mathcal{F}(\theta) = \mathbb{E}_{z \sim p(z)} [\nabla_\theta g(z, \theta)^\top \nabla_x f |_{x=g(z, \theta)}]$ (Chain Rule).

HW4 整理

Topic 1: Newton-Schulz Runtime

In muon, we use NS to Find:

$$G_t = U \Sigma V^T, P(G_t) = U P(\Sigma) V^T, [P^n(G_t) = U P^n(\Sigma) V^T = UV^T] \quad (*)$$

to replace finding SVD of G_t , so as to save time!

(Recall in Muon: $W_{t+1} = W_t - \eta UV^T$, i.e., discard Σ)

How can NS save runtime? Consider:

$$p(W) = \frac{1}{2} (3Id_{out} - WW^T)W$$

Assume 1 operation takes one unit of time. Then:

WW^T : $WER^{dout \times din}$, takes $2dout \cdot din$ *

*: For $A \times B$, $A \in R^{P \times m}$, $B \in R^{m \times n}$, time: $2mpn$

pxn elements, each need m multiplication & $m-1 \approx m$ summation

$(3Id_{out} - WW^T)^T$ takes $2cdout \cdot din$, c is a constant

∴ Total: $\Theta(dout \cdot din)$

But, if $dout >> din$, can we be faster?

$$\text{re-write: } p(W) = \frac{1}{2} (3Id_{out} - WW^T)W = \frac{1}{2} W (3Id_{in} - W^T W)$$

Then total: $\Theta(din \cdot dout)$.

Topic 2. MuP at unit scale.

If wanting to train on fp8, its number of representable values is small.

So for Xavier Init, if $\sim N(0, 1/d)$, can be a bug: rounding issue.

So it makes sense to: $\sim N(0, 1)$

However, to avoid exploding, we can assign a constant float scalar:

$$y = cWx$$

So we can set $c = \sqrt{din}$ to set variance $\frac{1}{din}$ scaled.

Now, $\Delta y = c \Delta W x$, if $\|x\|_{rms} = 1$, and we want $\|y\|_{rms} \leq 1$, then:

$$\|\Delta y\|_{rms} \leq \frac{1}{\sqrt{din}} \cdot \frac{\sqrt{din}}{\sqrt{dout}} \|\Delta W\|_2 \cdot 1 \leq 1 \quad \therefore \|\Delta W\|_2 \leq \sqrt{dout}$$

*: $\|AB\| \leq \|A\| \|B\|$, $\|A\|_{rms} = \sqrt{\frac{din}{dout}} \|A\|_2$, $A \in R^{dout \times din}$

We want to abide by this constraint!

If using SignSGD: $\Delta W = \alpha \cdot \text{sign}(\nabla_w L)$

$$\text{so } \|\Delta W\|_2 \leq \alpha \|\text{sign}(g) \text{sign}(x^T)\|_2 = \alpha \sqrt{din \cdot dout} = \sqrt{dout} \therefore \alpha = \sqrt{din}$$

where $g = \nabla_y L(y)$ ($y = Wx + b$).

But if we use Muon - Style method:

Consider: $U, \Sigma, V^T = SVD(\nabla wL)$, $\Delta W = \alpha UV^T$, we use NS to get UV^T

$$\text{Then } \|W\|_2 \leq \alpha \|U\|_2 \|V^T\|_2 = \alpha \leq \sqrt{\text{dout}} \Rightarrow \alpha \leq \sqrt{\text{dout}}$$

⇒ SignGD, Adam and Muon: global scale of raw gradients doesn't effect final update direction.

Topic3 Feature Dimensions of CNN

Suppose input: $W \times H \times C$, Kernel size, padding, number and stride are

K, P, S, F . Then:

of weights: FK^2C

of bias: F

$$W' = \left\lfloor \frac{W-K+2P}{S} \right\rfloor + 1$$

$$H' = \left\lfloor \frac{H-K+2P}{S} \right\rfloor + 1$$

$$C' = F$$

And a max pooling layer has stride S and size K . Then:

$$W' = \frac{W-K}{S} + 1$$

$$H' = \frac{H-K}{S} + 1$$

$$C' = C$$

Suppose L layers with kernel size of K , and stride is 1. The final output's receptive field size is $L \times (K-1) + 1 = O(LK)$

Topic4: Weights and Gradients in a CNN

Input: $X = [x_{i,j}]$, $i, j \in [1, n]$; $W = [w_{i,j}]$, $i, j \in [1, k]$; Output $Y = [y_{i,j}]$, $i, j \in [1, m]$

Then: $y_{i,j} = \sum_{h=1}^K \sum_{l=1}^k x_{ith-1, j+l-1} w_{h,l}$ (assume no padding, stride is 1, channel is 1)

Denote as $Y = X * w$. We have L and $dY = [dy_{i,j}]$, where:

$dy_{i,j} = \frac{\partial L}{\partial y_{i,j}}$ So we're curious about dW .

$$\text{We have } \frac{\partial L}{\partial w_{h,l}} = \sum_{i,j} \frac{\partial L}{\partial y_{i,j}} \cdot \frac{\partial y_{i,j}}{\partial w_{h,l}} = x_{ith-1, j+l-1}, i, j \in [1, m]$$

$$\therefore dW_{h,l} = \sum_{i=1}^m \sum_{j=1}^m x_{ith-1, j+l-1} dy_{i,j} \Rightarrow dW = X * dY$$

Suppose $x_{i,j}, dy_{i,j}$ are r.v. with 0 mean, σ_x^2, σ_y^2 variance. then:

$$\mathbb{E}(dw_{ij}) = 0, \text{Var}(dw_{ij}) = m^2 \sigma_x^2 \sigma_y^2, \text{ where } m = n - k + 1$$

$$\therefore \text{Std}(dw_{ij}) = \Theta(n). \text{ Finally, for pooling:}$$

If max: only one of input pixel has gradient of magnitude, and the rest are zeros

If average: distribute gradient evenly across each input blocks

HWS 整理

Topic 1: BN.

For x_i : $\frac{\partial L}{\partial x_i} = \sum_{j=1}^n \frac{\partial L}{\partial y_j} \cdot \frac{\partial y_j}{\partial x_i}$, and $y_i = r(x_i - \frac{1}{n} \sum_{j=1}^n x_j) + \beta$

$$\therefore \frac{\partial y_i}{\partial x_i} = \begin{cases} r(1 - \frac{1}{n}), & \text{if } i=j \\ r(-\frac{1}{n}), & \text{if } i \neq j \end{cases} \quad \therefore \frac{\partial L}{\partial x_i} = \left[\sum_{i=1}^n - \frac{r}{n} \frac{\partial L}{\partial y_i} \right] + \left[r \frac{\partial L}{\partial y_i} \right]$$

$$= r \left[\frac{\partial L}{\partial y_i} - \frac{1}{n} \sum_{j=1}^n \frac{\partial L}{\partial y_j} \right]$$

*: CNN: BN is applied to every value in a channel & in a batch

i.e., average over $BS \times H \times W$ entries and scale & bias

Topic 2: Regularization and Dropout

Recall: $L(w) = \|y - Xw\|_2^2$. One way of dropout during SGD on $X \in \mathbb{R}^{d \times n}$: keep each feature $\sim i.i.d$ Bernoulli(p), and zero out if not kept.

It turns out, our goal now is:

$L(\tilde{w}) = \mathbb{E}_{R \sim \text{Bernoulli}(p)} [\|y - (RX)\tilde{w}\|_2^2]$, while $R_{i,j} \sim i.i.d$ Bernoulli(p) and \odot is the element-wise product.

Consider Generalized Ridge-Regression involving:

$$L(w) = \|y - Xw\|_2^2 + \|\tilde{T}w\|_2^2$$

It turns out we can rewrite and get: $L(\tilde{w}) = \|y - P\tilde{X}\tilde{w}\|_2^2 + p(1-p)\|\tilde{T}\tilde{w}\|_2^2$ with \tilde{T} being a diagonal matrix whose j -th diagonal entry [is the norm of the j -th column of the X]

* Long Proof and we ignore

Then we can tell: $\tilde{T} = \sqrt{\frac{1-p}{p}} \tilde{I}$. If we want to finally write to:

$$L(\tilde{w}) = \|y - \tilde{X}\tilde{w}\|_2^2 + \lambda \|\tilde{w}\|_2^2$$

Then: $w = \tilde{T}^{-1} \tilde{w}$, $\tilde{X} = X \tilde{T}^{-1}$, with the definition of $\tilde{T} = \sqrt{\frac{1-p}{p}} \tilde{I}$ we conclude that: dropout is actually BN (even include rescale!).

HU6 整理

Topic 1: Connection Between graph dynamics & GNN.

Let's assume: ① For a network, underlying graph is adjacency matrix A .

② Has n vertices in each layer, corresponding to the n vertices of underlying graph

③ Each vertex has n channels

④ The input to each node in the 0-th layer is a one-hot encoding of own identity. That is, the node i in the graph has input $(0, \dots, \uparrow, 0, \dots)$

⑤ Weight connecting node i in layer k to node j in layer $k+1$: $A_{i,j}$; i -th entry

⑥ At each layer, operation at each node: weighted sum up

(a): j -th node at layer k : A_j^k . Since $G^{t+1} = AG^t$, $G^0 = I$

(b): So the i -th output of j -th node at layer k , equals to **the number of path from i to j who total distance is k** . $\uparrow j$'s neighbor

(c): In this simple setting, update for v_j is: $v_j = \sum_{i \in N(j)} v_i$

(d): If replace 'sum' func with 'max'. Then the i -th output of node j at layer k represents: if there's a path

(e). Analogy of CNN to GNN: (CNN vs GNN).

Semantic Segmentation \leftrightarrow Node-level prediction

Color Jitter \leftrightarrow Jittering node value

Dropout \leftrightarrow In GNN layers, set some node or edge to 0

Blur \leftrightarrow Average with Neighbors ResNet \leftrightarrow $v_i^{t+1} = v_i^t + \text{update}(v_i^t)$

Topic 2: GNN

For GNN node update, a function which satisfies **permutation invariance** is valid. For loss, we usually use cross-entropy loss:

$$\frac{1}{n} \sum_k (y(x) \log \hat{y}(x) + (1-y(x)) \log (1 - \hat{y}(x)))$$

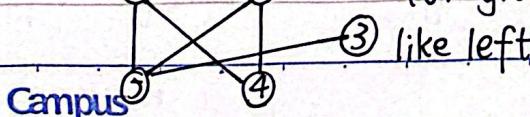
$y(x)=1$ if true else 0; $\hat{y}(x)$ is the prob models predict that $y(x)$ is 1

Example Update: $S_i^l = S_i^{l-1} + W_1 \underbrace{\sum_{j=1}^{n_i} \tanh(W_2 m_{i,j})}_{n_i}$

This is permutation invariant and using connection of residual.

If $W_1 \in \mathbb{R}^{d \times k}$, W_2 k rows, then $W_1 \in \mathbb{R}^{d \times k} + \tanh(W_2 S_i^{l-1})$

For graph: e.g: $S_2^l = S_2^{l-1} + \frac{W_1}{3} [\tanh(W_2 S_1^{l-1}) + \tanh(W_2 S_4^{l-1})]$



$$S_3^l = S_3^{l-1} + \frac{W_1}{1} (\tanh(W_2 S_5^{l-1}))$$

HW7 整理 & Discussion 8 整理

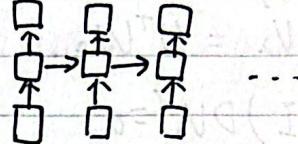
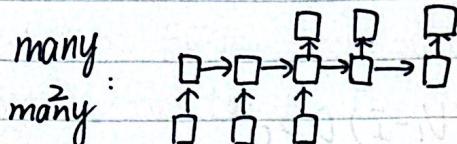
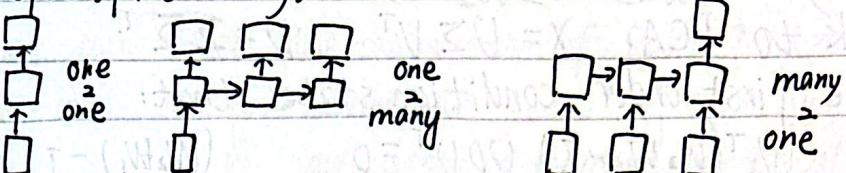
Topic 1: Implementing RNN

In RNN, the update of $h_t \rightarrow h_{t+1}$: $h_t = \sigma(W^h h_{t-1} + W^x x_t + b)$

and the output of h_t : $\hat{y}_t = W^T h_t + b^T$

RNN can be used

for many prediction:



...

Consider a many \rightarrow one at timestep T. Then:

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial h_T} \cdot \frac{\partial h_T}{\partial W} + \frac{\partial L}{\partial h_{T-1}} \cdot \frac{\partial h_{T-1}}{\partial W} + \dots + \frac{\partial L}{\partial h_1} \cdot \frac{\partial h_1}{\partial W}$$

The later terms in the sum often end up small or large \rightarrow Vanishing & Exploding

If nonlinearity is not considered and use MSE:

$$\frac{\partial L}{\partial h_t} = 2(\hat{y}_t - y) W^T (W^h)^{T-t}, \text{ so if the magnitude of the largest}$$

eigenvalue of W^h is much larger than 1, gradient explode; vice versa.

So if $t \ll k$, $\frac{\partial L}{\partial h_t}$ will be very small \Rightarrow hard to learn long-range dependency

Topic 2: PCA & AutoEncoders (k < m)

Assume data $x_i \in \mathbb{R}^m$ and encoder $W_1 \in \mathbb{R}^{k \times m}$ & decoder $W_2 \in \mathbb{R}^{m \times k}$

We have dataset $X \in \mathbb{R}^{m \times n}$, the Loss is:

$$\mathcal{L}(W_1, W_2; X) = \|X - W_2 W_1 X\|_F^2$$

We assume $\sigma_1^2 \geq \dots \geq \sigma_k^2 > 0$ are top-k eigenvalues of XX^T

So set $S = \text{diag}(\sigma_1, \dots, \sigma_k)$ with corresponding eigenvectors are the columns of $U_k \in \mathbb{R}^{m \times k}$. Then, first want to find out first order optimality conditions that \mathcal{L} satisfy.

$$\|A\|_F^2 = \text{tr}(AA^T), \text{ tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB), \text{ tr}(A) = \text{tr}(A^T)$$

$$\nabla_A \text{tr}(AB) = B^T, \quad \nabla_B \text{tr}(AB) = A^T.$$

$$\therefore \nabla_{A^T} \text{tr}(ABBA) = \nabla_{A^T} \text{tr}(ABBC) / C=A + \nabla_{A^T} \text{tr}(CBBA) / C=A$$

$$= (BBA)^T + (ABB)^T$$

KOKUYO

$$\begin{aligned}\nabla_{W_2} \mathcal{L} &= \nabla_{W_2} \text{tr}(XX^T W_2 W_1 X X^T W_1^T W_2^T - 2 W_2 W_1 X X^T) \\ &= 2(W_1 X X^T W_1^T W_2^T)^T - 2(W_1 X X^T)^T \\ &= 2(W_2 W_1 - I) X X^T W_1^T\end{aligned}$$

Similarly : $\nabla_{W_1} \mathcal{L} = 2 W_2^T (W_2 W_1 - I) X X^T$

Back to PCA. $X = U \Sigma V^T$, $D = \Sigma^T$.

The first order condition states that:

$$W_2^T (W_2 W_1 - I) U D U^T = 0 \quad (W_2 W_1 - I) U D U^T W_1^T = 0$$

Let $V_1 = W_1 U$, $V_2 = U^T W_2$, then:

$$(V_2 V_1 - I) D V_1^T = 0 \quad V_2^T (V_2 V_1 - I) D = 0$$

If $W_2 = W_1^T$ are the first k columns of U , then

$$V_1 = \begin{bmatrix} 1 & \dots & 0 \end{bmatrix} \quad V_2 = V_1^T \Rightarrow W_1 = U_k^T, \quad W_2 = U_k$$

Topic 3: SSM as kernel

Consider $y_t = \sum_{i=0}^T K_{t-i} u_i$. Then: $x_1 = Bu_0, x_2 = ABu_0 + Bu_1$,

$$x_3 = A^2Bu_0 + ABu_1 + Bu_2 \dots \quad y_t = C x_t$$

$$\therefore K = (0, CB, CAB, CA^2B, \dots)$$

Topic 4: AutoEncoder

① If train two different autoencoder on the same dataset, the one which produces lower validation loss **not necessarily perform better on downstream task!**

E.g.: Add noise or bottleneck smaller, loss \uparrow , but can produce more useful representation.

② Using representations by a learned autoencoder is most useful when only having a few data samples for downstream task. (Can also be useful if plenty, of course)

③ For training autoencoder with noise or masking, when inferencing for the downstream task, no augmentation for input is needed. (Or: uncommon!)

$$\text{Topic 5: } \mathcal{L}_X(W_1, W_2; X) = \frac{1}{n} \|X - W_2 W_1 X\|_F^2 + \lambda \|W_1\|_F^2 + \lambda \|W_2\|_F^2$$

Represent inductive bias of: Forcing W_1, W_2 to be orthogonal.

Think of $W_1 = U_1 \Sigma_1 V_1^T$, $W_2 = U_2 \Sigma_2 V_2^T$. W_1, W_2 should be pseudo inverse.

Regularization: $\lambda \sum_{i=1}^k (\frac{1}{6} \epsilon_i^2 + \frac{1}{6} \epsilon_i^{-2}) \Rightarrow$ forcing $\epsilon_i \approx 1$

$$\Rightarrow W_1 W_1^T = I, \quad W_2 W_2^T = I$$

HW28 線性代數

Topic 1: SSM as convolution kernel.

Previously, we've known: $\{y\} = K * \{u\}$, i.e., $y_k = \sum_{l=0}^L k_l u_{k-l}$
 $(k \in [1, L], \text{ any } u_{\leq 0} \text{ is set to 0})$ and $K = (D, CB, CAB, CA^2B, \dots)$

Then: how much can we parallelize the computation of $\{y_k\}$? ($u \in \mathbb{R}^n$)

If unroll in a recursive way, we have to compute y_{L-1}, \dots if want y_L .

then critical path is of $O(L)$ length, and depth is $O(\log n)$, so total: $O(L \log n)$

But if using FFT, depth of computation graph: $O(\log L)$; each with depth $O(\log n)$, so total: $O(\log L \log n)$

Also, given A, B, C , how can we compute K efficiently? We can tell that the bottleneck is $A^L \Rightarrow (I, A, \dots, A^{L-1}) = X^L$

We can use divide and conquer for this:

$$X^L = (X^{L/2}, A^{L/2} X^{L/2}), \quad A^L = A^{L/2} A^{L/2}$$

So for each $L = 2^t$, we can compute X^L : t times, each with depth of $O(\log n)$. Therefore, total: $O(\log n \cdot \log L)$

But If we add structure to A ? Suppose A is diagonal, then we can avoid matrix multiplication and degrade complexity to $O(\log L)$

However A can't always be that perfect. Suppose $A = I_n + PP^T$, $P \in \mathbb{R}^n$

$$\text{Now: } A^L = (I_n + PP^T)^L = I_n + LPP^T + \binom{L}{2} P P^T P P^T + \dots + (P P^T)^L$$

$$\text{and } (P P^T)^L = P(P^T P)^{L-1} P^T = P P^T \cdot L^{L-1}, \text{ where } L = \|P\|_2^2$$

Then $A^L = I_n + P P^T \sum_{l=1}^L \left[\binom{L}{l} L^{l-1} \right]$, and computing this can be done in $O(\log L)$ using divide and conquer (No matrix multiplication!)

Topic 2: Self - Supervised Linear Purification

Consider non-trivial purification scenario: $L(W; X, \lambda) = \|X - WX\|_F^2 + \lambda \|W\|_F^2$

Where: $W \in \mathbb{R}^{m \times m}$, $X \in \mathbb{R}^{m \times n}$.

Assume $\sigma_1 > \dots > \sigma_m \geq 0$, $X = U \Sigma V^T$. Then we can actually work out the optimal \hat{W} :

$$\hat{W} = U \cdot \begin{bmatrix} \frac{\sigma_1^2}{\sigma_1^2 + \lambda} & & \\ & \frac{\sigma_2^2}{\sigma_2^2 + \lambda} & \\ & & \ddots & \frac{\sigma_m^2}{\sigma_m^2 + \lambda} \end{bmatrix} U^T$$

KOKUYO

$$\text{Proof: } \mathcal{L}(W; X, \lambda) = \|X - W X\|_F^2 + \lambda \|W\|_F^2 = \sum_{i=1}^m \|X_i - W_i X\|_2^2 + \lambda \|W_i\|_2^2$$

$$= \sum_{i=1}^m \|X_i^T - X^T W_i^T\|_2^2 + \lambda \|W_i^T\|_2^2$$

$$\text{We're familiar with this: } \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \Rightarrow \hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$$

$$\therefore \hat{W}_i^T = (X X^T + \lambda I)^{-1} X X_i^T \Leftrightarrow \hat{W} = X X^T (X X^T + \lambda I)^{-1}$$

$$= \dots = U \Sigma \Sigma^T (\Sigma \Sigma^T + \lambda I)^{-1} U^T$$

Topic 3. Ridge Attention

Firstly: updating simple averaging can be efficient! $m = \frac{1}{n} \sum_{i=1}^n x_i$.

then if x_{n+1} is looped in, $m = \frac{mn + x_{n+1}}{n+1}$, $x_i \in \mathbb{R}^d$, y_i : scalar

Consider Ridge Regression. $A = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$, $y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$, so $w^* = (A^T A + \lambda I)^{-1} A^T y$.

Now want to rewrite $(A^T A + \lambda I)$ and $A^T y$ as sums involving x_i and y_i

$$\textcircled{1} A^T A + \lambda I = \lambda I + \sum_{i=1}^n x_i x_i^T \quad \textcircled{2} A^T y = \sum_{i=1}^n x_i y_i$$

Now suppose we want to do ridge-self-attention*, with a context length of n and d -dim q, k, v . * : n queries is applied to the same pool of (k, v) pairs

k : A matrix; q, v : $\mathbb{R}^{d \times d}$, so actually we want W^* , and inference: $W^* q$

For $W^* q$: $W^* \in \mathbb{R}^{d \times d}$, $q \in \mathbb{R}^{d \times n} \Rightarrow O(nd^2)$.

For computing $W^* = (I + A^T A)^{-1} A^T B$, $A = \begin{bmatrix} k_1^T \\ \vdots \\ k_n^T \end{bmatrix}$ and $B = \begin{bmatrix} v_1^T \\ \vdots \\ v_n^T \end{bmatrix}$

$O(nd^2)$, inverse $O(d^3)$. $\downarrow O(nd^2)$

\therefore Complexity: $O(nd^2) + O(d^3)$; if $n > d$, then $O(nd^2)$

But for inverse: a nice trick: Sherman-Morrison Formula.

$$(M + UV^T)^{-1} = M^{-1} - \frac{1}{1 + V^T M^{-1} U} (M^{-1} U)(V^T M^{-1})$$

if U, V are two d -dimension vector and $M \in \mathbb{R}^{d \times d}$, $O(d^2)$

HW9 整理 & Discussion 9, 10

* q_h of $q_h \sim \mu_i, \sigma^2 I$

Topic 1. Justifying Scaled-Dot Product Attention

Suppose $q_h, k \in \mathbb{R}^d \sim N(\mu, \sigma^2 I)$. Then: ($\mu \in \mathbb{R}^d$)

$$E[q_h^T k] = E\left[\sum_{i=1}^d q_{hi} k_i\right] = \sum_{i=1}^d E[q_{hi}] E[k_i] = \sum_{i=1}^d \mu_i = \|\mu\|_2^2$$

If $\mu=0, \sigma=1$, then $\text{Var}(q_h^T k) = E[(q_h^T k)^2] - E[q_h^T k]^2 \rightarrow d \times \sigma^4$

$$E[(q_h^T k)^2] - E[k^T q_h q_h^T k]^2 = \dots ?$$

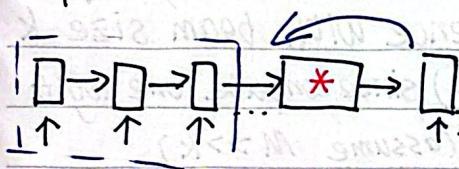
$$\text{So } \text{Var}(q_h^T k) = \text{Var}\left(\sum_{i=1}^d (q_{hi} k_i)\right) = d \text{Var}(q_{hi} k_i)$$

$$= * d(E[q_h]^2 \text{Var}(k_i) + E[k_i]^2 \text{Var}(q_{hi}) + \text{Var}(q_{hi}) \text{Var}(k_i))$$

$$= d\sigma^4 = d \quad *: \text{Var}(XY) = \text{Var}(X)\text{Var}(Y) + \text{Var}(X)E[Y]^2 + \text{Var}(Y)E[X]^2$$

So now, we want $E\left[\frac{q_h^T k}{s}\right] = 0$, and $\text{Var}\left[\frac{q_h^T k}{s}\right] = \frac{d}{s^2} = 1 \Rightarrow s = \sqrt{d}$

Topic 2: Attention Mechanism to Mitigate RNN bottleneck



In RNN, many-to-many, the bottleneck is that
* need to store all information about the input.

But now, we allow the hidden states from decoder as queries and hidden states from encoder as keys and values.

So decoder hidden state can compute similarity with input's hidden state, 'look at them' $\Rightarrow \hat{h}_i = \sum_{j=0}^L \underline{q_j}, \underline{l_j}$ here, attention score.

Topic 3: RoPE

PE is used to ensure that word position is known. Because attention is applied symmetrically to all input vectors. For RoPE: $\text{RoPE}(x_t^{(1)}, x_t^{(2)}, t)$

$$= \begin{pmatrix} \cos tw & -\sin tw \\ \sin tw & \cos tw \end{pmatrix} \begin{pmatrix} x_t^{(1)} \\ x_t^{(2)} \end{pmatrix} \text{ for 2-dim vectors. Moreover:}$$

if $z_t = x_t^{(1)} + j x_t^{(2)}$, then $\text{RoPE}(z_t, t) = e^{jtw} z_t$, hence $\text{RoPE}(x, m)^\top \text{RoPE}(y, n)$

So we can extend to $2n$ -dim vector:

$$= \text{RoPE}(x, m)^\top \text{RoPE}(y, n)$$

$$\text{RoPE}((x^{(1)}, x^{(2)}, \dots, x^{(n)}), t) = \begin{pmatrix} x^{(1)} \cos tw^{(1)} - x^{(2)} \sin tw^{(1)} \\ x^{(2)} \cos tw^{(1)} + x^{(1)} \sin tw^{(1)} \\ \vdots \end{pmatrix} \text{ d:dim of input vector}$$

RoPE doesn't add any learnable parameter! w is fixed! $w^{(1)} = 10000^{-2} \cdot \frac{i-1}{d}$

△ Note that in some model (specific. like GPT), it might be possible for attention layer itself to learn to encode positional information.

KOKUYO

HW10 & HW11 整理

Topic 1: Kernelized Linear Attention

$O(N^2)$ complexity in memory and computation is bad. Let $x \in \mathbb{R}^{N \times F}$ be the sequence. The transformer is a function $T: \mathbb{R}^{N \times F} \rightarrow \mathbb{R}^{N \times F}$:

$$T_L(x) = f_L(A(x)x + x)$$

Formally: $W_Q \in \mathbb{R}^{F \times D}$, $W_K \in \mathbb{R}^{F \times D}$, $W_V \in \mathbb{R}^{F \times M}$, then:

$$\begin{cases} Q = xW_Q, \quad K = xW_K, \quad V = xW_V \\ A_L(x) = V' = \text{softmax}\left(\frac{QK^T}{\sqrt{D}}\right)V \end{cases}$$

$$\text{sim}(q_h, k) = \exp\left(\frac{q_h^T k}{\sqrt{D}}\right)$$

$$\text{But generically: } V_i' = \frac{\sum_{j=1}^N \text{sim}(q_{i \cdot}, k_{j \cdot}) V_j}{\sum_{j=1}^N \text{sim}(q_{i \cdot}, k_{j \cdot})}$$

One potential attention variant is 'polynomial kernel attention'. E.g., degree is 2: $\text{sim}(q_h, k) = (q_h^T k + 1)^2 = \phi(q_h)^T \phi(k)$

Then one benefit of using kernelized attention is that we can represent a kernel using a feature map: $\phi(\cdot)$. In this case:

$$\phi(x) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \dots, \sqrt{2}x_D, x_1^2, \dots, x_D^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_{D-1}x_D]$$

Then, kernel can be: $K(q_h, k) = (\phi(q_h)^T \phi(k))$, then:

$$V_i' = \frac{\sum_{j=1}^N \phi(q_{i \cdot}) \phi(k_{j \cdot})^T V_j}{\sum_{j=1}^N \phi(q_{i \cdot}) \phi(k_{j \cdot})^T} = \frac{\phi(q_{i \cdot}) \sum_{j=1}^N \phi(k_{j \cdot})^T V_j}{\phi(q_{i \cdot}) \sum_{j=1}^N \phi(k_{j \cdot})^T}$$

what's better: j is irrelevant to $\phi(q_{i \cdot})$: \Rightarrow

For quadratic kernel, space complexity is $O(ND + D^3)$.

If $N \gg D^2$, then lots of memory can be saved.

Topic 2: LoRA

A common strategy to do pretrained model adapting is to update only a subset of its parameters and keep the rest frozen.

In this topic, focus on $W \in \mathbb{R}^{m \times l}$. W_0 is the pretrained value.

In LoRA: $W = W_0 + AB$, $A \in \mathbb{R}^{m \times k}$, $B \in \mathbb{R}^{k \times l}$, $k \ll \min(m, l)$

(a) Suppose using LoRA to adapt a new task and observe underfitting

① Try $k \uparrow$ to allow more expressive power

② Change initialization to leverage SVD of either the original

matrix or matrix-valued gradient for a large amount of training data. By aligning initializations to natural gradient directions, we are likely to be able to have more significant impacts faster.

③ Change learning rate or other hyperparameter

(b). If A, B are initialized to all zero \Rightarrow Zero gradient for $A \& B$, No!

But if both Xavier initialized, $AB \neq 0$, and feature of pretrained model is not used.

Topic 3: A brief introduction of Transformer Interpretability.

A core concept is residual stream:

$$X^{l+1} = X^l + f^{(l)}(X^l) \Rightarrow X^{\text{final}} = X^0 + \sum_{l=0}^{L-1} f^{(l)}(X^l)$$

Consider attention head with H heads. In 'concatenation' view, each head's output are stacked together vertically, followed by a W_0^H . So $H = W_0^H \cdot R^H$ $W_0^H \in \mathbb{R}^{d_{\text{model}} \times (d_{\text{head}} \cdot H)}$.

Consider $W_0^h \in \mathbb{R}^{d_{\text{model}} \times d_{\text{head}}}$ head output $r^h \in \mathbb{R}^{d_{\text{model}} \times n_{\text{context}}}$

So H is also : $H = \sum_{h=1}^H W_0^h r^h$, these two are equivalent.

$$H = W_0^H R^H = [W_0^1, W_0^2, \dots, W_0^H] \begin{bmatrix} r^1 \\ r^2 \\ \vdots \\ r^H \end{bmatrix} = \sum_{h=1}^H W_0^h r^h \quad \text{independently.}$$

It means that we can analyze the contributions of each head

The 'concat' view can be efficient, but 'additive' view reveal independency

'Circuit Theory' treats Transformer as a collection of interacting circuits built from attention and MLP.

QK circuit OV circuit

QK circuit : Determining Attention Patterns \downarrow information to move

An attn head can be split into two operations : where to look and what

Pre-softmax score : $S_{ij} = Q_i^T K_j$, $Q_i = W_Q X_i$, $K_j = W_K X_j \Rightarrow S_{ij} = X_i^T (W_Q^T W_K) X_j$

If $W_Q = I$ \Rightarrow 内积; if $W_K = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$, then this head only cares about first dim.

There can be 负极 in W_K which encourage dissimilarity. These mean that

W_K can be highly robust and obtain powerful expressivity. Attn pattern from QK circuit

OV Circuit: $h(X^0) = W_0^h W_V X^0 (A^h)^T = W_0^h \underbrace{X^0 (A^h)^T}_{\text{(post-softmax)}}$

$X^0 (A^h)^T$: attention-weighted average. Then: $X^{\text{final}} = X^0 + \sum_{i=1}^H h_i(X^0)$

$= X^0 + \sum_{i=1}^H W_0^i W_V X^0 (A^i)^T$ We can notice that $h_i(X^0)_t$ (t -th column) $\in \text{Col}(W_V^i)$

Campus \Rightarrow Head can only write to a low-dim subspace of the d_{model} -dim residual stream