

CS182 Note & Review

Optimization: Implicit Regulation of GD; SGD; and Momentum

(x, y) data, need model $f_\theta(x) \rightarrow y$

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n L(y_i - f_\theta(x_i))$$

Assume $L(\vec{w}) = \|\vec{x}\vec{w} - \vec{y}\|_2^2$, then $\vec{w}^* = (\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{y}$

$$\vec{w}_{t+1} = \vec{w}_t - \eta \nabla_{\vec{w}} L(\vec{w}, \vec{x}_{\text{train}}, \vec{y}_{\text{train}}) \quad (\text{GD})$$

$$\vec{w}_{t+1} - \vec{w}_t = -2\eta \vec{x}^T (\vec{x}\vec{w} - \vec{y}) \quad ((\vec{x}\vec{w} - \vec{y})^T (\vec{x}\vec{w} - \vec{y}) \neq 0)$$

Want to rewrite into: $\vec{w}_{t+1} - \vec{w}^* = ? \cdot (\vec{w}_t - \vec{w}^*)$

W_b 衡量 $\vec{w}_t \rightarrow \vec{w}^*$ 距离在下降/上升 / 如何变换

$$\vec{w}_{t+1} - \vec{w}^* = (I - 2\eta \vec{x}^T \vec{x})(\vec{w}_t - \vec{w}^*)$$

∴ We want $(I - 2\eta \vec{x}^T \vec{x})^n$ converge when $n \uparrow$

Consider SVD for $I - 2\eta \vec{x}^T \vec{x}$, then the largest eigenvalue in Σ shouldn't be greater than 1!

$$\therefore |1 - 2\eta \lambda_{\max}(\vec{x}^T \vec{x})| < 1 \quad \eta < \frac{1}{\lambda_{\max}(\vec{x}^T \vec{x})}$$

⇒ 学习率 (or step size) 不可过大！

* 同时也暗含着: singular values & spread of singular values matter
Why? 不妨 $\vec{x} = U \Sigma V^T$ 代入 \vec{w}^* (with ridge considered!)

$$\begin{aligned} \vec{w}^* &= (V \Sigma^T U^T (U \Sigma V^T + \lambda I)^{-1} V \Sigma^T U^T \vec{y}) \\ &= V (\Sigma^T \Sigma + \lambda I)^{-1} V^T \Sigma^T U^T \vec{y} \\ &= V (\Sigma^T \Sigma + \lambda I)^{-1} \Sigma^T U^T \vec{y} \\ &= V \begin{bmatrix} \sigma_1^2 + \lambda & 0 \\ 0 & \sigma_2^2 + \lambda & 0 \\ & \ddots & \ddots \end{bmatrix} \Sigma^T U^T \vec{y} = \sum_{i=1}^n \vec{v}_i \frac{\sigma_i}{\sigma_i^2 + \lambda} \vec{u}_i^T \cdot \vec{y} \end{aligned}$$

可见 $\lambda \ll \sigma_i$, 则等于不正则化; $\lambda \gg \sigma_i$, 全为 0, 优化目标只与 $\|\vec{w}\|_2^2$ 有关了
因此提供了从奇异/特征值解算正则化的方向

带有 Ridge 正则化项的优化称为“岭回归”

\Rightarrow : Regularization \Leftrightarrow Weight Decay

在GD中, $\vec{W}_{t+1} = \vec{W}_t - 2\eta X^T (\vec{X}\vec{w}_t - \vec{y})$. remains in the span of data

GD for ridge:

$$\vec{W}_{t+1} = \underbrace{(1 - 2\eta\lambda)}_{\text{weight decay}} \vec{W}_t - 2\eta X^T (\vec{X}\vec{w}_t - \vec{y})$$

奇异/特征值的威力更直观一些:

$$\vec{w}_{t+1} = \vec{w}_t - 2\eta V \Sigma^T V^T (\vec{X}\vec{w}_t - \vec{y}) \quad \text{左右同乘 } V^T: \begin{cases} V^T \vec{y} = \vec{q} \\ \vec{W} = V\vec{w} \end{cases}$$

$$\vec{W}_{t+1} = \vec{W}_t - 2\eta \Sigma^T (\Sigma \vec{W}_t - \vec{q}), \vec{w}_{t+1,i} = \vec{w}_{t,i} - 2\eta \sigma_i [\sigma_i \cdot \vec{w}_{t,i} - \vec{q}_i]$$

可见 σ_i 越小, i -th 维度越不重要, 更新越慢

这可视为一种 "early stopping", 因为 GD 在奇异值大的方向上对数据拟合, 拟合之后, 奇异值小的几乎不再拟合

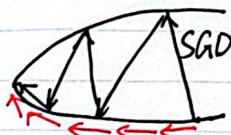
从岭回归的 $\frac{\sigma_i}{\sigma_i + \lambda}$, 若 σ_i 很小, 防止原本的 $\frac{1}{\sigma_i}$ 会很大; 到 GD 中的 σ_i 越小, 对应方向更新速率越小。上述两种都认为是 Implicit Regulation

SGD: GD 对凸函数很好, 但也很 expensive, 因为一次更新要全部数据点都 inference 一遍。

Add some noise to gradient direction! In expectation, it descends:

$$\theta_{t+1} = \theta_t - \eta \frac{1}{B} \sum_{i=1}^B \nabla f_i(\theta_t) \quad B: \text{mini batch, batch size}$$

SGD 能 work, 还有利于跳出 local minima & saddle points, 且 less cost!
但训练中的 oscillation 会难搞:



How to smooth it?

如 batch size \uparrow , low-pass filter ...
以及 Momentum!

Momentum: $\left\{ \begin{array}{l} \vec{W}_{k+1} = \vec{W}_k - \eta \vec{\pi}_{k+1} \\ \vec{\pi}_{k+1} = \beta \vec{\pi}_k + (1-\beta) \nabla f(\vec{w}_k) \end{array} \right.$ 相当于梯度 β damping!

$$\vec{\pi}_{k+1} = \beta \vec{\pi}_k + (1-\beta) \nabla f(\vec{w}_k), \vec{\pi}_0 = 0, \beta \in [0, 1]$$

\Rightarrow Add history of 梯度, 并按 history 时间 decay weight exponentially

$$\vec{\pi}_n = \beta^{n-1} (1-\beta) \nabla f(\vec{w}_0) + \dots = \sum_{i=1}^n \beta^{i-1} (1-\beta) \nabla f(\vec{w}_{n-i})$$

Adam (Who is Adam?)

GD 在 singular value 大的方向上, step 更大。优点在于 Implicit Regularization

缺点在于: 1. learning rate. 非常难选, 因为要兼顾

奇异值大的方向, 梯度大, 反之很小 \Rightarrow normalize the gradients

对梯度向量除以模长, 使之变为单位 vector

$$\vec{W}_{t+1} = \vec{W}_t - \eta \frac{\nabla_{\vec{W}} L(\vec{W})}{\|\nabla_{\vec{W}} L(\vec{W})\|_2 + \epsilon} \rightarrow \text{范数且开正根}$$

但这不是“方向奇异值越大, 变化越快”的最优算法! \Rightarrow 不考虑梯度大小信息!

Instead: $\vec{m}_t \leftarrow \nabla_{\vec{W}} L(\vec{W}) \quad \vec{v}_t \leftarrow (\nabla_{\vec{W}} L(\vec{W})) \downarrow \text{element-wise square}$

$$\vec{W}_{t+1} = \vec{W}_t - \eta \frac{\vec{m}_{t+1}}{\sqrt{\vec{v}_{t+1}} + \epsilon} \rightarrow \epsilon > 0, \text{ avoid division by zero}$$

而 $\frac{\vec{m}_{t+1}}{\sqrt{\vec{v}_{t+1}}} = \text{sgn}(\nabla_{\vec{W}} L(\vec{W}))$, “sign SGD”

但也有 issue: Even we are close, it keeps overshooting / bouncing

\Rightarrow smooth

$$\begin{cases} \vec{m}_{t+1} = \beta \vec{m}_t + (1-\beta) \nabla_{\vec{W}} L(\vec{W}_t) \\ \vec{v}_{t+1} = \delta \vec{v}_t + (1-\delta) [\nabla_{\vec{W}} L(\vec{W}_t)]^2 \\ \vec{W}_{t+1} = \vec{W}_t - \eta \frac{\vec{m}_{t+1}}{\sqrt{\vec{v}_{t+1}} + \epsilon} \end{cases}$$

且 $\begin{cases} \tilde{\vec{m}}_{t+1} = \frac{\vec{m}_{t+1}}{1-\beta^{t+1}} \\ \tilde{\vec{v}}_{t+1} = \frac{\vec{v}_{t+1}}{1-\delta^{t+1}} \end{cases}$, Adam: $\vec{W}_{t+1} = \vec{W}_t - \eta \frac{\tilde{\vec{m}}_{t+1}}{\sqrt{\tilde{\vec{v}}_{t+1}} + \epsilon}$

Adam 常用于 SGD, 故 $\tilde{\vec{m}}, \tilde{\vec{v}}$ 由 mini-batch 计算得到

Normalization (统一化, 非归一化) helps with exploding / vanishing gradients

But: size of the step is largely dictated by the learning rate.

而 AdamW: Adam with weight decay:

$$\vec{W}_{t+1} = (1-\lambda) \vec{W}_t - \eta \frac{\tilde{\vec{m}}_{t+1}}{\sqrt{\tilde{\vec{v}}_{t+1}} + \epsilon}$$

Feature perspective and Taylor Expansion

Network from another perspective: $f(\vec{x}, \vec{\theta})$ △: f 可视为 loss function
 And further regard it as a function of $\vec{\theta}$:

Consider $f: R^m \rightarrow R$ (like regression / binary classification)

Taylor Expansion from current x, θ : change, 且负!! omit.

$$f(x, \theta_0 + \Delta\theta) = f(x, \theta_0) + [\langle \nabla_{\theta} f|_{\theta=\theta_0}, \Delta\theta \rangle] + [\dots \text{higher term}]$$

$\nabla_{\theta} f|_{\theta=\theta_0}$ 为 f 在 θ_0 时对 θ 的梯度， $\langle \cdot, \cdot \rangle$ 为内积

这个近似的核心前提： $\Delta\theta$ 很小！ (converge fast). hold true

因此陷入 dilemma: 欲知大 $\langle \nabla_{\theta} f|_{\theta=\theta_0}, \Delta\theta \rangle$ 大，又 $\Delta\theta$ 需小以靠近 θ_0

⇒ 优化问题: want: $\|\Delta\theta\|_2 \leq \eta$, 可能是 2 范数等。这里假设 $= \infty$

$$\text{i.e., } \|\Delta\theta\|_\infty = \max_j |\Delta\theta[j]|$$

在此优化约束下, 优化: $\underset{\|\Delta\theta\|_\infty \leq \eta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \langle \nabla_{\theta} f|_{x=x_i}, \Delta\theta \rangle$

$$= \underset{\|\Delta\theta\|_\infty \leq \eta}{\operatorname{argmin}} \langle \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} f|_{x=x_i, \theta=\theta_0}, \Delta\theta \rangle \quad \text{内积拆开}$$

$$= \underset{|\Delta\theta[j]| \leq \eta}{\operatorname{argmin}} \sum_j \left[\left(\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} f|_{x=x_i, \theta=\theta_0} \right) [j] \cdot \Delta\theta[j] \right] \Rightarrow \text{notation: } \nabla_{\theta_0} f_{\text{batch}}$$

可见, 可分别优化: $\underset{-\eta \leq \Delta\theta[j] \leq \eta}{\operatorname{argmin}} \nabla_{\theta_0} f_{\text{batch}} [j] \cdot \Delta\theta[j]$

$\nabla_{\theta_0} f_{\text{batch}}$ 与 $\Delta\theta$ 无关, 可见 $\Delta\theta[j]$ 取值实则很简单。

$$\Delta\theta = \begin{cases} +\eta & \text{if } \nabla_{\theta_0} f_{\text{batch}} < 0 \\ -\eta & \text{if } \nabla_{\theta_0} f_{\text{batch}} > 0 \end{cases} \Rightarrow \Delta\theta^* = -1 \operatorname{sign}(\nabla_{\theta_0} f_{\text{batch}})$$

这就是 Adam 算法的核心!

Initialization & Optimizer

标准(归一)化 (standardization) 在 ML / DL 中很重要!

- 不同 feature, 不同 numerical magnitude (else large features dominate)

Consider ReLU : scalar case:

$$\text{elbow: } -\frac{b}{w}, \text{ say, } b \sim N(0, 1) \quad w \sim N(0, 1)$$

则 $-\frac{b}{w} \sim \text{Cauchy Distribution}$

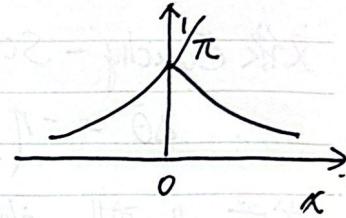
△ 柯西分布定义: x_0 : 位置参数, 峰值中心位置; r : 尺度参数, 分布离散程度

$$f(x) = \frac{1}{\pi r [1 + (\frac{x-x_0}{r})^2]}$$

当 $x_0=0$, $r=1$ 时, 得标准柯西分布 $\Rightarrow f(x) = \frac{1}{\pi (x^2+1)}$

Recall ReLU 中:

$$x' = -\frac{b}{w}$$

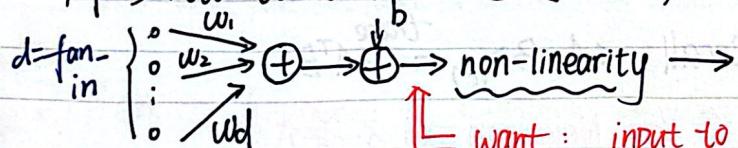


若是 input 分布在蓝区间中, 能充分利用 non-linearity

故若 input 不归一化, 则模型表达能力严重下降

故 input 归一化很重要, 对应地, model 参数初始化也很重要!

How to initialize? 若 0 初始化, gradients are zero, weights never move
所以 how to do? 注意到, 上一层的 activated output 是下一层输入



则考虑: $b: 0$ 初始化*; $w_i \sim N(0, 1/d)$ $\xrightarrow{\text{Xavier}}$ 初始化

$$\begin{aligned} \text{Why? } \text{var} \left[\sum_i w_i h_i \right] &= E \left[\left(\sum w_i h_i \right)^2 \right] - \underbrace{\left[E \left(\sum w_i h_i \right) \right]^2}_{\rightarrow 0} \xrightarrow{\text{initialization}} \\ &= \sum_{i=1}^d E(w_i^2) E(h_i^2) \quad (\text{due to i.i.d.}) \end{aligned}$$

2 倍!

还有 He 初始化: ReLU 导致方差减半 \Rightarrow w 的初始化预先把权重方差放大
 $w_i \sim N(0, 2/d)$.

* bias 初始化不只	{	use small number
0 初始化, 还可以:	b 合并至 w, 视为 $d+1$ 维, 然后全都 Xavier ($\sim N(0, 1/d+1)$)	

当然, 它更受欢迎

* 通过调整不同层的初始化方差和学习率缩放，来解耦“特征学习”与“Lazy training”动态。在一个小模型上找到最优学习率，且可直接应用于一个任意宽的大模型上（小宽度 \rightarrow 大宽度）

注意到：初始化时， $w(b)$ 采用随机 initialization，而在 Adam 中，因为 Taylor 近似视角， $\Delta\theta$ 不能很大

\Rightarrow 意味着 models barely change from random initialization

这被视为一种 **lazy training** (NTK, Neural Tangent Kernel)，且随机初始化不应视为能总是 work \Rightarrow 即神经网络的 power 还未全部使用

在之前的“ $\Delta\theta$ 不能很大”中，用 ∞ 范数 — 若 L_2 范数呢？

$$\|\Delta\vec{\theta}\|_2 \leq \eta \quad \text{欲 } \underset{\|\Delta\vec{\theta}\|_2 \leq \eta}{\operatorname{argmin}} \langle \nabla_{\vec{\theta}} L(\vec{\theta}), \Delta\vec{\theta} \rangle$$

注意到： $\langle x, y \rangle = x^T y = \|x\|_2 \|y\|_2 \cos \varphi$, $\cos \varphi$ 应尽可能 -1

又依 Cauchy - Schwartz : $|\langle \cdot, \cdot \rangle|$ 最大时，两个 vector 应对齐 (align)

$$\therefore \Delta\vec{\theta} = -\eta \frac{\nabla_{\vec{\theta}} L(\vec{\theta})}{\|\nabla_{\vec{\theta}} L(\vec{\theta})\|_2}$$

或者， $\|\Delta\vec{\theta}\|_2$ 的约束以正则化体现： $\underset{\vec{\theta}}{\operatorname{argmin}} \underbrace{\langle \nabla_{\vec{\theta}} L(\vec{\theta}), \Delta\vec{\theta} \rangle}_{g(\vec{\theta})} + \lambda \|\Delta\vec{\theta}\|_2^2$

$$\nabla g(\vec{\theta}) = \nabla_{\vec{\theta}} L(\vec{\theta}) + 2\lambda \Delta\vec{\theta}, \text{ convex}$$

$$\Delta\vec{\theta}^* = -\frac{1}{2\lambda} \nabla_{\vec{\theta}} L(\vec{\theta}) \xrightarrow[\text{step size}]{\text{gradient descent}} \text{只考虑一层}$$

之前的 $\vec{\theta}$ 总是 vector，但真正的 network 不只一层 \Rightarrow matrix!

$$\underset{\|\Delta W\| \leq 1}{\operatorname{argmin}} \langle \nabla_W L(W), \Delta W \rangle, \text{ Recall } \langle A, B \rangle_{TF} = \frac{\text{trace}}{(A^T B)}.$$

设 $\nabla_W L(W) = U \sum V^T$, 则依 Von Neumann 不等式：

$$\Delta W^* = -U V^T$$

这便是 Shampoo 优化器！通过上述三例，we have a recipe:

- Given a norm, make an optimizer & examples for $\|\cdot\|_2$, $\|\cdot\|_\infty$
(下一部分预告)

Additional: RMSnorm: $\|\vec{x}\|_{RMS} = \frac{1}{\sqrt{d_{in}}} \|\vec{x}\|_2$ 诱导范数

Matrix version: Recall: induced Matrix-Norm: $\|A\|_{\alpha \rightarrow \beta} = \max_{\|\vec{x}\|_\alpha=1} \|A\vec{x}\|_\beta$

$$\therefore \text{RMS} \rightarrow \text{RMS Norm: } \|A\|_{RMS \rightarrow RMS} = \max_{\|\vec{x}\|_{RMS}=1} \|A\vec{x}\|_{RMS} = \max_{\|\vec{x}\|_2=\sqrt{d_{in}}} \frac{1}{\sqrt{d_{in}}} \|A\vec{x}\|_2 = \frac{\sqrt{d_{in}}}{\sqrt{d_{in}}} \max_{\|\vec{x}\|_2=1} \|A\vec{x}\|_2$$

* 谱范数 (Spectral): 衡量一个矩阵能对另一个向量产生的最大拉伸程度

\hookrightarrow 其实就是诱导范数 of L_2

$$\Delta: \|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |A_{ij}|^2} = \sqrt{\text{tr}(A^T A)} = \sqrt{\sum_{i=1}^m \min_{j \in [n]} \alpha_i^2}$$

* RMS \rightarrow RMS Norm is rescaling of the spectral norm

No. _____
Date _____

RMS Norm & Maximal update normalization (UP)

Recall, induced matrix norm: for matrix, the def of its

$$\|A\|_{\alpha \rightarrow \beta} = \max_{\|\vec{x}\|_\alpha = 1} \|A\vec{x}\|_\beta \quad \text{spectrum norm (like L2)}$$

$\|\vec{x}\|_\alpha = 1 \rightarrow$ for vector, we know their norm!

RMS \rightarrow RMS Norm \leftrightarrow Motivated by Xavier Initialization:

$$\|A\|_{\text{RMS} \rightarrow \text{RMS}} = \max_{\|\vec{x}\|_{\text{RMS}} = 1} \|A\vec{x}\|_{\text{RMS}}$$

$$\text{And } \|\vec{x}\|_{\text{RMS}} = \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2}, \quad \text{RMSNorm}(\vec{x}) = \frac{\|\vec{x}\|}{\|\vec{x}\|_{\text{RMS}}}$$

$$= \frac{1}{\sqrt{n}} \|\vec{x}\|_2 \quad * \quad \vec{x} \in \mathbb{R}^{\text{dim}}, A \in \mathbb{R}^{\text{dout} \times \text{din}}$$

$$\|A\|_{\text{RMS} \rightarrow \text{RMS}} = \sqrt{\frac{1}{\text{dout}} \max_{\|\vec{x}\|_2 = 1} \|A\vec{x}\|_2} = \sqrt{\frac{1}{\text{dout}} \|A\|_2}$$

Recall: optimizer recipe: $\underset{\Delta w}{\operatorname{argmin}} \langle \nabla_w \mathcal{L}(w), \Delta w \rangle$
 \rightarrow choose norm for $\|\Delta w\| \leq \eta$ \rightarrow choose $\eta \Rightarrow$ get optimizer

In optimizer like Adam, we've tried $\|\cdot\|_0$ and $\|\cdot\|_2$, what about $\|\cdot\|_{\text{RMS}}$?

$$\Delta w \in \mathbb{R}^{\text{dout} \times \text{din}}, \text{ then } \|\Delta w\|_{\text{RMS}} \leq \eta \Rightarrow \|\Delta w\|_2 \leq \sqrt{\frac{1}{\text{dout}}} \cdot \eta$$

Then we can adjust the learning rate dynamically

Different Approach: Commit to SignSGD/Adam:

$$W_{t+1} = W_t - \eta \operatorname{sign}(\nabla_w \mathcal{L}(w))$$

Rethink about step size!! A fixed η is not reasonable

Cause in different layers of different architectures, W are different!

Want η adjust in accordance with W shape \Rightarrow RMS Norm can help

\rightarrow Remains in span of minibatch data

If SGD, $\nabla_w \mathcal{L}(w)$ is low rank, suppose $\nabla_w \mathcal{L}(w) = \sigma \vec{u} \vec{v}^T$ (SVD)

$$\text{Then } \|\nabla_w \mathcal{L}\|_2 = \sigma, \quad \|\nabla_w \mathcal{L}\|_F^2 = \sigma^2, \quad \text{def } \operatorname{sign}(\nabla_w \mathcal{L}(w)) = S$$

$$\therefore S = \operatorname{sign}(\sigma \vec{u} \vec{v}^T) = * \operatorname{sign}(\vec{u}) \cdot \operatorname{sign}(\vec{v}^T) \Rightarrow \operatorname{rank}(S) = 1$$

$$* \operatorname{sign}(mn) = \operatorname{sign}(m) \cdot \operatorname{sign}(n)$$

$$\Rightarrow \|\nabla_w \mathcal{L}\|_2^2 = \sigma^2 \quad \|\nabla_w \mathcal{L}\|_F^2 = \text{tr}(\sigma^2 \vec{u} \vec{v}^T \vec{v} \vec{u}^T) = \sigma^2$$

$$W_{t+1} = W_t - \underbrace{\eta \cdot S}_{\Delta w}, \quad \text{we want } \Delta w \text{ small to ensure Taylor approximation}$$

Now we choose $\|W\|_{\text{RMS}} \leq \delta \leftarrow \text{step size threshold}$

$\therefore \eta \leftarrow \frac{\sqrt{d_{in}}}{d_{out}} \|S\|_2 \leq \delta$, $\|S\|_2 ?$, S has all +1 or -1, and rank is 1!

Intuitively and in fact, $\|S\|_2 = \sqrt{d_{in} d_{out}}$, $\eta \leq \frac{1}{d_{in}} \delta$

\therefore If using RMS norm, learning rate can scale with layer!

\Rightarrow Adam / signSGD lr can scale as $\frac{1}{d_{in}}$ to maintain $\|W\|_{\text{RMS}} \leq \delta$

Essence: μP optimization: maximal update parameterization

$\Rightarrow \mu P$: How to change hyperparameter when NN width changes?

Can we optimize hyperparameter on smaller NN and 'transfer' to larger ones?

E.g. $\frac{1}{\sqrt{n}} (X_1 + \dots + X_n) \xrightarrow{n \rightarrow \infty} N(0, 1)$, $X_i \sim N(0, 1)$. but: $\frac{1}{\sqrt{n}}$ is the right scaling factor!

\therefore So for $\eta \leq \frac{1}{d_{in}} \cdot \delta$, $\eta_{d_{in}}$ is the right scaling

for: $\vec{h}_{l-1} \rightarrow [w_l] \rightarrow \oplus \rightarrow [\text{non-}l_i] \rightarrow \vec{h}_l$, we want $\|\vec{h}_l\|_{\text{RMS}} = 1$

$\Rightarrow \|\vec{h}_l\|_2 = \sqrt{d_{out}}$, and suppose don't care bias & non- l_i :

$\|\vec{h}_l\|_2 = \|w_l\|_2 \|\vec{h}_{l-1}\|_2$, $\|w_l\|_2 = \Theta(\sqrt{\frac{d_{out}}{d_{in}}}) \Rightarrow \|w_l\|_{\text{RMS}} = \Theta(1)$

\Rightarrow fixed RMS norm as we started with!!

\Rightarrow If $\|W\|_{\text{RMS}} = \Theta(1)$, $\|\vec{h}_l\|_{\text{RMS}} = \Theta(1)$, then $\|\vec{h}_l\|_{\text{RMS}} = \Theta(1)$

This 'zero-shot hyperparameter transfer' is good! up!

For more info, please refer to HW3 Topic 2

Muon

Shampoo Optimizer: if $\nabla_w L(w) = U\Sigma V^T$, then:

$$W_{t+1} = W_t - \eta UV^T.$$

Why discard Σ ? It can be confusing!



For example, direction 1's gradient is large, and w can be busy swinging (like blue line), but pay less attention to direction 2.

But SVD can be expensive! Muon: if we can find:

$$G_t = U\Sigma V^T, P(G_t) = U P(\Sigma) V^T, P^n(G_t) = U P^n(\Sigma) V^T = U V^T$$

This method is called Newton-Schulz

$$\begin{cases} B_t = \mu B_{t-1} + \nabla_w L(w) \\ O_t = \text{Newton Schulz}(B_t) \\ W_t = W_{t-1} - \eta O_t \end{cases}$$

CNN & GNN

* Architecture Introduction Order: MLP \rightarrow CNN \rightarrow GNN \rightarrow RNN \rightarrow Transformer

Basic Data Form: Images: $h \times w \times 3$ (RGB)

Overall: CNN

Key Ideas: (Expressivity) convolutional filters called "kernels"

- ① Locality \rightarrow convolution*
- ② Invariance / Equivalence / Symmetries
- ③ Hierarchical Structure \hookrightarrow weight sharing & data augmentation (Train)
 \Rightarrow Depth & Multi-Resolution

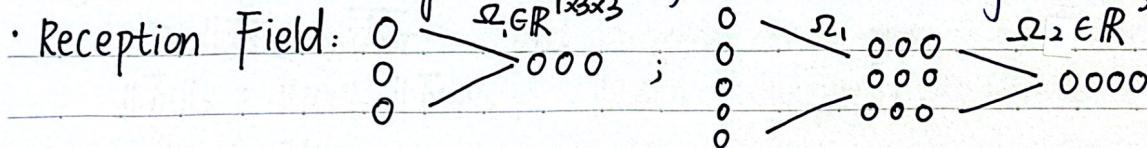
Key Ideas: (Work more Effectively)

- ① Normalization Layers
- ② Dropout
- ③ Residual / Skip Connection

Convolution: In DL, we don't flip !! And can have a bias.

- For kernels, weights are not shared across channels. Elementwise multiplication and sum, followed by a ReLU.
- Also: If wanting more channels at the output \Rightarrow repeat with fresh weights & bias
- Padding: "0", "same", "mirror" or no padding

Zero Padding \Rightarrow Black frame around Image



Suppose. $K \times K$ conv with L layer, one position in the final layer can see the information of range up to: $\frac{K+1}{2} L - 1$

- Down Sampling via Stride i
- Pooling: stride / subscript (取固定位置) maxpooling; mean pooling
- $K \times K$ conv with C_{in} channel input and C_{out} outputs:
 \Rightarrow Weights: $K^2 C_{in} C_{out}$ Bias: C_{out} (#)

- Backprop w.r.t. weights & bias? PyTorch will handle it.

每一步算完后的 $\frac{\partial L}{\partial w_i}$, 最后累加起来 ($\text{Conv}(x, w) = x_1 \otimes w + x_2 \otimes w \dots, \frac{\partial L}{\partial w} = \sum_i \frac{\partial L}{\partial w_i} x_i$)
 \Leftrightarrow (Scale of $\frac{\partial L}{\partial w}$ change with the number of terms in \sum_i)

- Data Augmentation: Robustness

Normalization Layers for Stability and Effectiveness

Idea: RMSNorm: d-dim $\vec{h_i} \xrightarrow{\text{RMSNorm}} \tilde{h_i}$, $\|\tilde{h_i}\|_{\text{rms}} = 1$

$\therefore \tilde{h_i} = \frac{\vec{h_i}}{s}$, where $s = \sqrt{\frac{1}{d} \sum h_i[j]^2 + \epsilon}$, s.t.

$$\Leftrightarrow \tilde{h_i} = \frac{\vec{h_i}}{\sqrt{d} / (\|h_i\|_{\text{rms}} + \epsilon)} = \vec{h_i} / (\|h_i\|_{\text{rms}} + \epsilon) \quad \checkmark \text{ 满足}$$

Batch Norm: Average Over Space & Batch

Input: (bs, channel, H, W): 对每一个channel, 计算这个channel在整个批次中的均值 & 方差

$$m = \frac{1}{|B|} \sum_{i \in B} h_{in}, s = \sqrt{\frac{1}{|B|} \sum_{i \in B} (h_{in} - m)^2}, h_{out}[i, j] = \gamma \left(\frac{h_{in}[i, j] - m}{s + \epsilon} \right) + \delta$$

where γ and δ are learnable, initialized with $\gamma=1, \delta=0$

But for test/inference time: Batchnorm don't have a batch. Need m&s!

How to get? ① 在 train epoch 最后一个 m&s

② Run an extra epoch of training with no grad update

③ Just learn m&s using held-out data

Dropout "Data Augmentation" inside our network during training

Idea: Zero out random activation inside net \Rightarrow Encourage Internal Representation to be redundant (Typical for MLP)

Training do dropout, but Inference? 若仅是放弃(zero out)通常不 work!

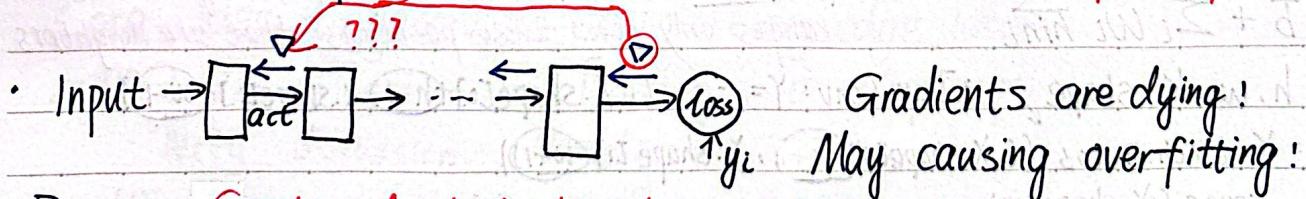
因为会 mess up scale

\Rightarrow Solution: Scale activations by p. (Replace multiplication Noise by its mean)

(Done by PyTorch Dropout: Inference:

\Rightarrow (在 inference 时 scale back)

$$\tilde{h_i}[i] = \begin{cases} 0, & 1-p \\ \frac{h_i[i]}{p}, & p \end{cases}$$

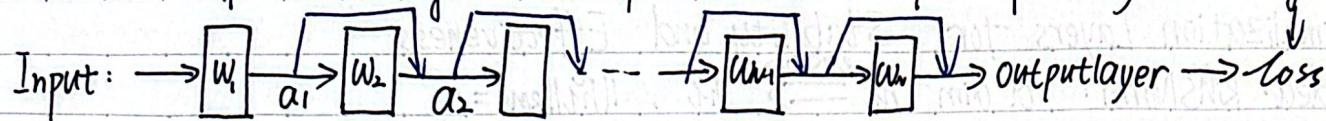


Desire: Create A "checkpoint superhighway" layer

Idea 1: Add links from output of all earlier layers to the input of a single

\Rightarrow Blow up architectures! Not working!

Idea 2: Skip Connecting: Add Input to the Output of a layer! y_i



Inside: ... $\rightarrow w \rightarrow \oplus \rightarrow \text{ReLU} \rightarrow \text{Norm} \rightarrow \dots$ ($w \& b$: Conv style or MLP style).
 $f_i(x)$.

\Rightarrow Neural Network: $f_1(f_2(\dots f_L(x) \dots))$

Now: Resnet Skip Connection $f_{i+1}(x) = f_i(x) + g_i[f_i(x)]$

Norm: Modern Default: Inside the block near the start

\rightarrow Norm $\rightarrow w \rightarrow \oplus \rightarrow \text{ReLU} \rightarrow \dots$ Pre-Norm Convention
 \Rightarrow stable training

But one challenge: if $f_i(x)$ changes the dimension of input?

Simple addition will have dimension mis-match!

\Rightarrow Need Something that changes shape, e.g., Linear Layer

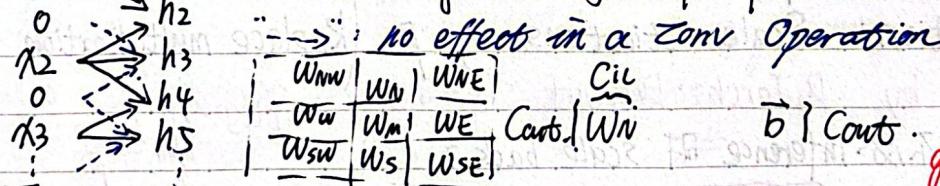
Note: If entire block is just a maxpool downsampling or a subsampling, no skip connection needed for this block since gradients can pass cleanly.

• Δ: Depthwise-Conv: Just treat each channel separately as if it was a one-input-channel & one-output-channel conv. No weight sharing across channels

• Upsampling? For pooling's inverse, we can zero or duplicate pad

\Rightarrow Idea: Fill with zeros and use a learnable Transpose Conv

Tran Conv: $x_1 \rightarrow h_1$, essentially: zero-fill upsample Conv.



3x3 example:
 $h_{out} = \vec{b} + \sum_i w_i \vec{h}_{in i}$, ranges only over those positions that are neighbors

Code: $h, w = k. \text{shape}$, \rightarrow Tran Conv: $Y = \dots ((X. \text{shape}[0] - h + 1) \times X. \text{shape}[1] - w + 1))$

Conv: $Y = \text{torch.zeros}((X. \text{shape}[0] - h + 1), X. \text{shape}[1] - w + 1))$

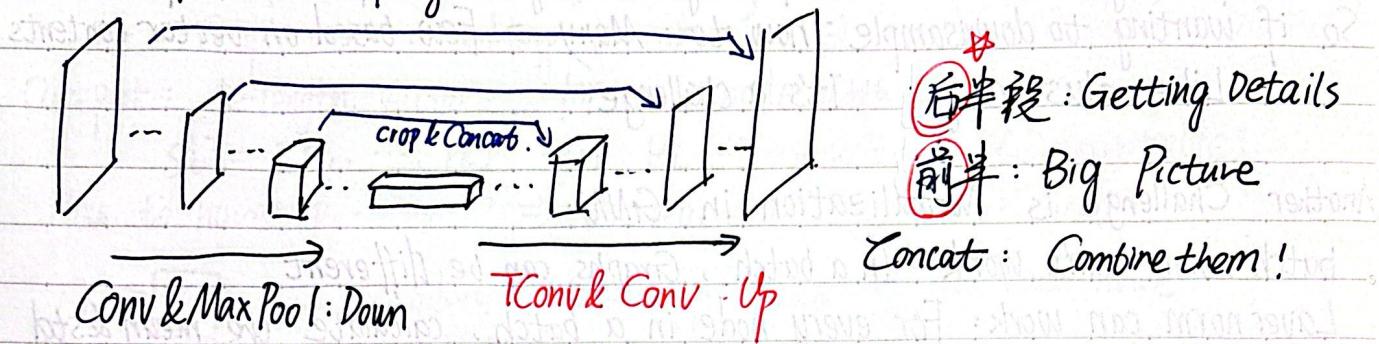
for i in range(X.shape[0]):

 for j in range(X.shape[1]):

 Conv $Y[i, j] = (X[i:i+h, j:j+w] * K). \text{sum}()$

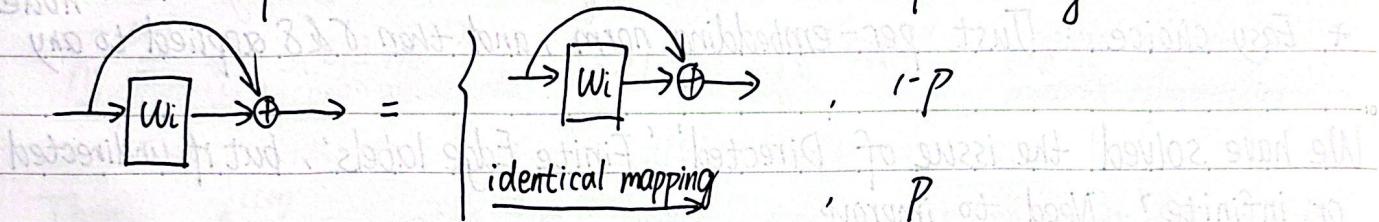
Tran: $Y[i:i+h, j:j+w] += X[i, j] * K$

With Up & Down Sampling \Rightarrow U-net:



Addition: Stochastic Depth Regularization

Idea: Dropout the entire block with some probability



It helps with overfitting and create functional redundancy

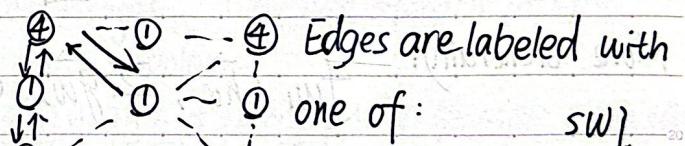
* Unlike dropout, no need to adjust anything during test time!

(Due to identical mapping & Presence of Normalization inside block)

GNN: Graph Neural Networks (Generalizing CNNs)

Now: Image \rightarrow General Graph

Image as a Graph? $\begin{matrix} 4 & 1 & 4 \\ 1 & 1 & 1 \\ 1 & 4 & 1 \end{matrix}$ \Rightarrow



Then, a 3×3 Conv can be viewed as receiving informations from neighbors, each of whom has their own labeled edge to me, multiplying by an edge-label-specific weight matrix, and adding up together with a special self-specific matrix (node's info & bias).

Core idea to generalize CNN \Rightarrow Every edge labeled from a discrete finite set

* What does not generalize? [Pooling, down/up sampling]

Need discussion: Normalization

* Why do we do this? \Rightarrow Grow Receptive Field Faster

Hack it! Add a global node connected to Everything

But it can have side effect: blur away all fine-grained details

② Oversquashing: info loss since squashing so many nodes' info into a finite embedding
 So if wanting to downsample: how to: Many \rightarrow Few based on vector contents
 (Like clustering) It's a challenge.

Another Challenge is Normalization in GNN:

batch norm can't work: In a batch, Graphs can be different

Layer norm can work: For every node in a batch, calculate the mean & std of all channels of this node and normalize

And then, use shared σ and δ (scale and shift) to process every node

* Easy choice: Just per-embedding norm, and then σ & δ applied to any

We have tackled scenario of 'Directed' 'Finite Edge labels', but if undirected or infinite? Need to improve.

For a node to update, which information it will have:

① myself ② degree: # of neighbors ③ set of neighbors

Also, for the info from neighbor set, our processing should be permutation-insensitive

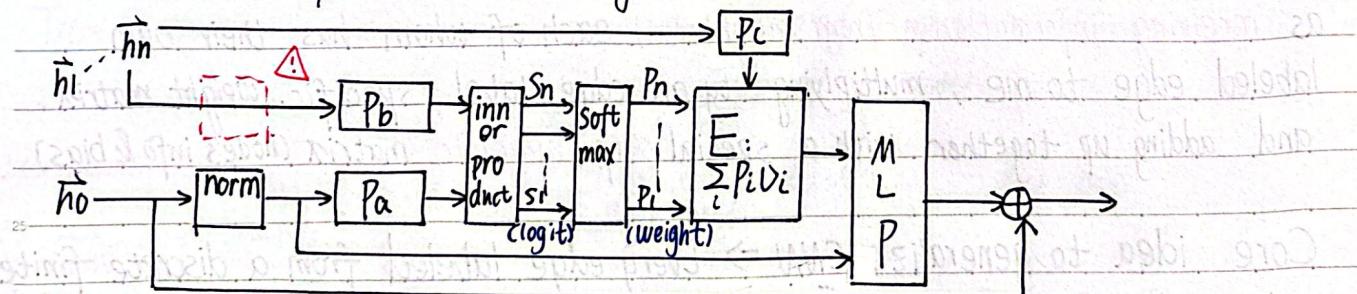
So: Generalized Framework: For node with info h_o , update should be:

$f_{W_1}(h_o, \sum_i g_{W_2}(h_i))$ not necessarily Σ , can be \prod , min/max mean, medium ...

More Generally:

$f_{W_1}(h_o, \sum_i g_{W_2}(h_o, h_i))$ or $f_{W_1}(h_o, \sum_i w_i(h_o, h_i)g_{W_3}(h_i))$

\Rightarrow Natural Example: GNN Building Block



Obviously, this is related to attention mechanism

Back to challenge of Downsampling: Introduce: Diff Pool

Idea: When we downsample, we get a smaller graph, but it's fully connected with probability on each edge

Campus

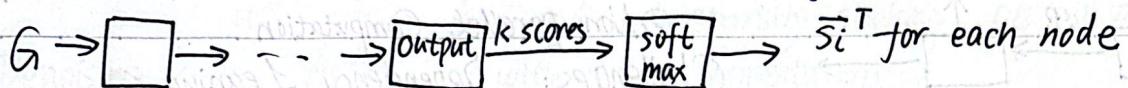
\Rightarrow Make each point belong to every cluster, just with some probability

Input : Graph with embedding vectors

Output : k-length vector for each node with prob init:

$S_i^T = [p_1, \dots, p_k]$ for $p_i \leftarrow \text{use softmax! (Differentiable)}$

How to have $p_1 \dots p_k$? \Rightarrow Use learnable algorithm: GNN



But \vec{S}_i^T aren't a graph. Not yet. What to do?

Need New Vectors (Embeddings) for k clusters.

$$\vec{h}_j^{(l+1)} = \sum_{i \in \text{cluster } j} S_i \cdot h_i^{(l)} \quad \Rightarrow \text{Let } A^{(l)} \in \mathbb{R}^{N \times N} \leftarrow \text{original adjacency matrix (symmetric)}$$

\curvearrowright j-th cluster in layer $l+1$ (after pooling)

Then $\underbrace{A^{(l+1)}}_{\in \mathbb{R}^{k \times k}} = S^T A^{(l)} S \in \mathbb{R}^{N \times k}$. A_{ij} represents how strong the edge is between cluster i and j .

\Rightarrow Now : A new graph with k nodes and a new adjacency matrix



Can help get desirable properties * \leftarrow Auxiliary Loss

\triangle Not to immediately normalize these messages (A_{ij} has strength info)

* Property 1: Have clusters respect graph topology

$L = \|A - SS^T\|_F$ cause encourage 原始图中相互连接的节点，在簇分配上也应相连

Property 2: Want to be close to hard cluster assignment

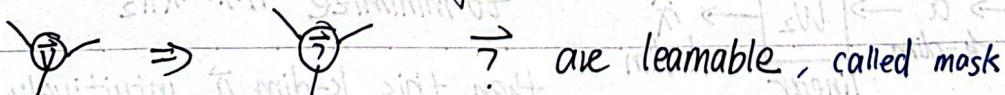
$$L = \frac{1}{N} \sum_i H(\vec{S}_i^T) \quad (\text{want assignment to be more deterministic})$$

\curvearrowright Entropy

Final thing : Held-out Data (Node-Level Problem)*

Approach 1: Excise entire nodes from the graph : but impact connectivity

2: Remove content from the node but leave it in



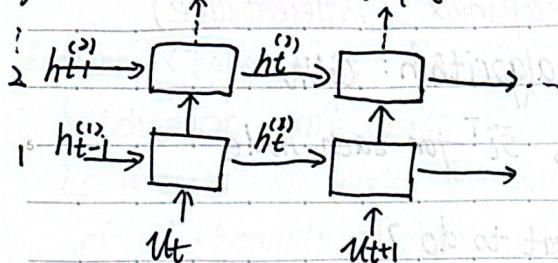
*: Like node classification, tagging, etc

: How to make vali / test dataset

RNN / Self Supervision

Traditional RNN: $h_{t+1} = \sigma(W_h h_t + W_u u_{t+1} + b)$

Layer: output : $y_{t+1} = g_j h_{t+1} + b_j$



Key problem: ① Gradients Exploding/Vanishing

② Non parallel Computation

Challenges: Dependency Leaning

Recap on Kalman Filter: $\hat{x}_{t+1} = A \hat{x}_t + B u_{t+1} + F(\vec{y}_{t+1} - C(A \hat{x}_t + B u_{t+1}) - D u_{t+1})$
 $= \tilde{A} \hat{x}_t + \tilde{B} u_{t+1} + F(\vec{y}_{t+1})$

Where $\tilde{A} = (I - FC)A$ $\tilde{B} = (I - FC)Bu - FD$

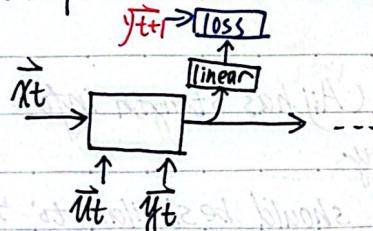
and $F = P C^T (C P C^T + K V)^{-1}$

$P = A P A^T + K u - A P C^T (C P C^T + K V)^{-1} C P A^T$ (P satisfy)

Anyway, essence is: $\vec{y}_{t+1} \rightarrow \vec{x}_{t+1} \rightarrow \hat{x}_{t+1}$, we can regard: $A \vec{B} F$:
 measurable

And \vec{y}_{t+1} : gt! For supervision!

But if we don't have label, how to supervise? Use next y !



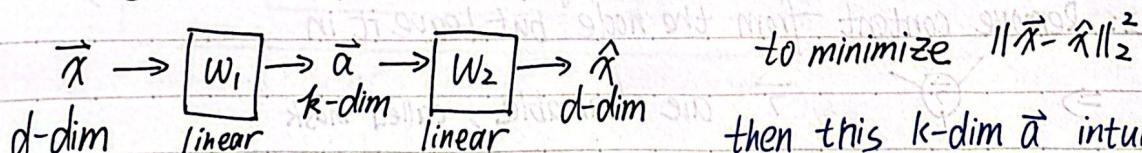
This is self-supervision!

- ① We can learn a partial pattern that can be useful
- ② Might need scaffolding (脚手架) parts of NN
- ③ Generic idea of "next-thing" prediction
in causal sequence modeling

Step back: Consider PCA: Classical Solution:

$X = U \Sigma V^T = \sum \sigma_i \vec{u}_i \vec{v}_i^T$, For oil select top-k high

IF design following network:

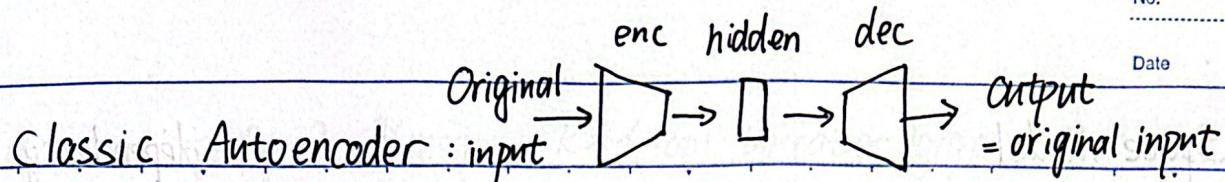


to minimize $\| \vec{x} - \hat{\vec{x}} \|_2^2$

then this k-dim $\vec{\alpha}$ intuitively contains the most rich info to reconstruct

Has a name: auto-encoder

△ Traditional Perception: Decoder is scaffolding



This hidden tensor contains rich feature though in low dimension!

Sparse Autoencoder :

Can make hidden "bottleneck" even bigger than original

with \Rightarrow Auxiliary loss in sparsity-seeking. Eg L1/L2 penalty

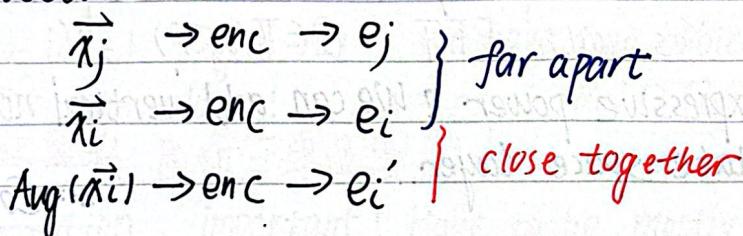
Denoising ones: Above two with data augmentation.

Masked Autoencoder: 'Aug' is masking out some patches and reconstruct them

Besides PCA, another self-supervised ML task: clustering

\Leftrightarrow Contrastive Learning. core idea: an example and its own augmentations should be in the same cluster. Different examples should be in diff cluster.

Basic Architecture :



*: Note that in Teacher-Forcing: input is gt, not last prediction!

But for inferencing, input is the last prediction (auto-regressive)

State Space Model

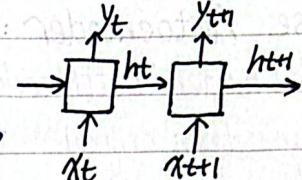
dependency

Problem with traditional RNN: ① Hard to parallelize ② Hard to learn long-range

For ①: Idea: Eliminate 'horizontal' nonlinearity

$$\text{In RNN: } \vec{h}_{t+1} = \underline{\sigma}(\vec{A}\vec{h}_t + \vec{B}\vec{x}_t + \vec{b})$$

Nonlinearity! Have to wait for \vec{h}_t and can't skip



What if: $\vec{h}_{t+1} = \vec{A}\vec{h}_t + \vec{B}\vec{x}_t$?

Then Lets say: $\begin{cases} \vec{h}_t = A\vec{h}_{t-1} + B\vec{x}_t \\ \vec{y}_t = C\vec{h}_t + D\vec{x}_t \end{cases} \Rightarrow \text{Linear-time invariant system}$

ABCD: Learnable

(A, B, C, D not dependent on time)

We can unroll \vec{y}_k !

$$\vec{y}_k = CA^k B \vec{x}_0 + \dots + CB\vec{x}_k + D\vec{x}_t$$

Now \vec{y}_k only depends on previous input! No on hidden states!

At time t, computation needs $O(t)$; But if Max seq length is T, so time range from 1 to t, computation cost is $O(T^2)$

Can we improve? Yes with FFT: $O(T \log T)$

But what about expressive power? We can add vertical nonlinearity, aka, MLP (e.g.) on state-space layer.

So to summarize: Training is actually sort of convolution*

For inference: All past information is in the state. \vec{h}_t

In \vec{y}_k , we have $CA^k B$ matrix. Can we speed up its computation?

For A^k , if A is **diagonal**, then A^k is easy to compute.

Now back to $\vec{h}_t = A\vec{h}_{t-1} + B\vec{x}_t$ and $\vec{y}_t = C\vec{h}_t + D\vec{x}_t$

Let's say: \vec{y}_t : dy-dim; \vec{x}_t : dx-dim; \vec{h}_t : dn dim ("Memory")

$\vec{y}_t = \begin{bmatrix} y_t[1] \\ \vdots \\ y_t[dy] \end{bmatrix}$, think of time-evolution sequence on a **Component** by component basis. such sequences

$y[k]: y_1[k], y_2[k], \dots, y_t[k]$, and there are dy

Regard $x[k]$ in a same way. Def Seq($x[t]$): $x_1[t] x_2[t] \dots x_t[t]$

Then an important formula is: in fact **Convolution**

$$\text{Seq}(y[k]) = \sum_{l=1}^{dx} F_{k,l} * \text{Seq}_l(x[t])$$

Appropriate Filter

Can we simplify? I mean: $k \times l$ can be rather large!

In imaging kernels we use are Depthwise / Channel-Wise.

If A is diagonal, then $\vec{h}_t[l] = \lambda_l \vec{h}_{t-1}[l] + B \vec{x}_{t-1}$ channel wise

So the update of $\vec{h}_t[l]$ is not dependent of $\vec{h}_{t-1}[i]$ ($i \neq l$) $\leq T$

For B , block diagonal, for simulating 'locality'.

Efficiency v.s. Efficacy Tradeoff

Larger \vec{h}_t , more compute, but also more memory/history tracking.

So goal: want to make \vec{h}_t as large as possible, while still being able to run. Now: Introduce: [S4] model

In a LTI (Linear Time Invariance) System, use FFT.

Let $A = \begin{bmatrix} \lambda^1 & 0 \\ 0 & \lambda^N \end{bmatrix}$, then how is λ related to system stability?

if $|\lambda| > 1$: Bad! But if $|\lambda| < 1$: history die out fast!

So for $\lambda \in \mathbb{R}$, $|\lambda| = 1 \Rightarrow \lambda = \pm 1$ we have two choices and far apart ($-1 \leftrightarrow 1$)

Solution: Complex eigenvalues: $\lambda = \exp(-\text{ReLU}(\lambda R) + j\lambda I)$

欲 A 特征值为复数且位于单位圆附近

GD hard

And for initialization, important! Have to be mostly on or very close to unit circle with some inside unit circle. Use Hippo-based initialization.

Also: Freeze A with good initialization and only train the rest also seemed to work well!

(Sampling)

[S6/Mamba]: Think of Discrete-time System as a [discretization] of a continuous time system:

$$\boxed{\frac{d\vec{h}(t)}{dt} = A \vec{h}(t) + B \vec{u}(t)} \quad \Delta: \text{And to convert CT to DT.} \\ \vec{y}(t) = C \cdot \vec{h}(t) \quad \text{we need a sampling interval } \Delta t \\ \text{Scalar case: } x(t) = ax(t) + b u(t)}$$

Assume $u(t)$ is piecewise constant for time Δt . (Δt small: reasonable).

Note that: $a > 0$: explode; $a < 0 \Rightarrow$ decay to 0

$$x(\Delta t) = e^{a\Delta t} \cdot x(0) + \frac{b}{a} (e^{a\Delta t} - 1) u(0)$$

Assume $a < 0$: Then $\Delta t \rightarrow 0 \Rightarrow e^{a\Delta t} \approx 1 \Rightarrow$ Remember!

Δt larg $\Rightarrow e^{a\Delta t} \rightarrow 0 \Rightarrow$ forget

KOKUYO

\hat{T}_{small}

Moreover: $e^{a\Delta T} - 1 \approx a\Delta T$, then $\frac{b}{a}(a\Delta T - 1) \approx b\Delta T$

$\Delta T \rightarrow 0 \Rightarrow b\Delta T \rightarrow 0 \Rightarrow$ input doesn't matter

ΔT large $\Rightarrow b\Delta T$ large \Rightarrow input matters

So via changing ΔT , can change relative importance of past & future.

Idea: Choose ΔT based on input and learn it.

How to generate ΔT ? Softplus: $\text{Softplus}(x) = \ln(1 + e^x)$

$$\text{So: } \begin{cases} \vec{h_{k+1}} = A_k \vec{h_k} + B_k \vec{x_k} \\ \vec{y_{k+1}} = C \vec{h_k} \end{cases} \quad A_k = \begin{bmatrix} e^{\lambda_1 \Delta T} & & \\ & e^{\lambda_2 \Delta T} & \\ & & \ddots \end{bmatrix} \quad B_k = \vec{b_k} = \vec{b} \Delta T$$

Lastly: Initialization:

AI blocks for A : A few methods:

- ① $\lambda_n = -\frac{1}{2} + nL$
- ② $\lambda_n = -\frac{1}{2}$
- ③ $\lambda_n = -(n+1)$

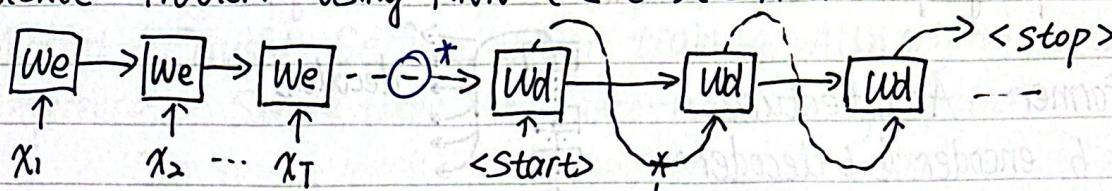
Initialization matrix: A is given as $\begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}$

Initial value: $\vec{h}_0 = \vec{0}$ (zero vector)

Initial value: $\vec{b}_0 = \vec{0}$ (zero vector)

Attention: II Intro to Attn

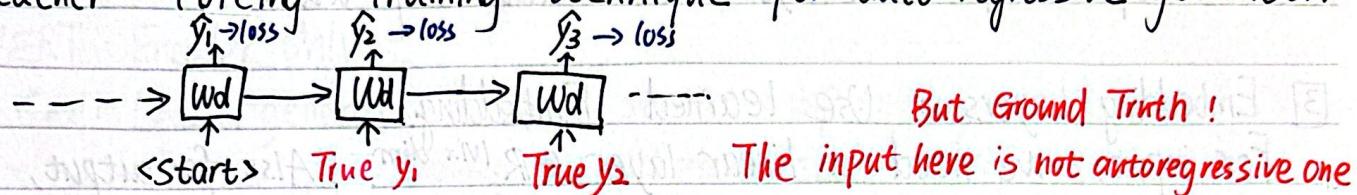
Sequence Problem using RNN (Enc-Dec Architecture)



Loss: Cross - Entropy → With Sampling

* Bottleneck : Need to capture all of the input sequence in this bottleneck

Teacher - Forcing : Training technique for auto-regressive generation



In U-net, we've met bottleneck. In U-net, solution is adding connections.

But for language : cannot set a topology of connects in advance !

Solution: Hash table , but need to be differentiable.

So we take weighted average : $\frac{1}{t} \sum_i \text{sim}(\vec{q}_n, \vec{k}_i) \cdot \vec{v}_i$

$\text{sim}(\vec{q}, \vec{k})$: inner product & softmax

Key question: How to learn them ? W_k, W_q, W_v !

There are many types of attention ...

Cross - Attention : Keys + Values from encoder , Querys from decoder

Self - Attention : K, Q, V from decoder :

$$\begin{cases} \vec{q}_n = \beta_{q_n} + W_q \vec{x}_n \\ \vec{k}_m = \beta_{k_m} + W_k \vec{x}_m \\ \vec{v}_i = \beta_{v_i} + W_v \vec{x}_i \end{cases} \quad \beta_{q_n, k, v} \text{ & } W_{q, k, v} \text{ are learnable}$$

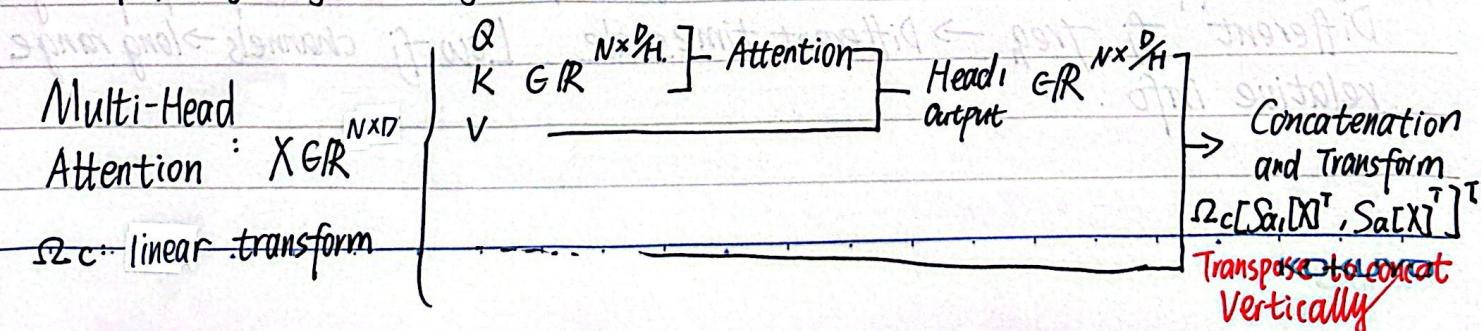
$$a[\vec{x}_m, \vec{x}_n] = \frac{\exp[\langle \vec{k}_m, \vec{q}_n \rangle]}{\sum_{j=1}^N \exp[\langle \vec{k}_j, \vec{q}_n \rangle]}$$

So Higher similarity \Rightarrow Higher attention

Scaled Self - Attention : $\frac{\exp[\langle \vec{k}_m, \vec{q}_n \rangle / \sqrt{D}]}{\sum_{j=1}^N \exp[\langle \vec{k}_j, \vec{q}_n \rangle / \sqrt{D}]}$ where $W_{q, k, v} \in \mathbb{R}^{D \times N}$

\Rightarrow Prevent entries to softmax

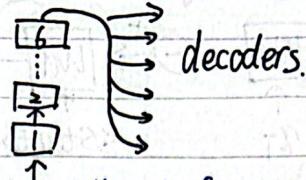
from getting too large



Multi-Query Attention: For each head: keys and values are the same, but change the query for each head. It can save memory!

② Transformer Architecture

Typically 6 encoder & 6 decoder:



And normalization layers are now typically before the attention block.

GPT-3: Decoder Only: Word embeddings are directly fed into a series of transformation layers with masked-self-attention.

③ Embedding Layers: Use learned Embedding.

For input, we need a linear layer $\in \mathbb{R}^{|\mathcal{V}| \times \text{dim}}$; Also for output, we use linear layer to project to $|\mathcal{V}|$ and softmax.

Like before, we may want $W\vec{x} = \vec{y}$, $\|\vec{y}\|/\text{RMS} = 1$, and since \vec{x} are $\in \mathbb{R}^{|\mathcal{V}|}$, one-hot, so we want each column of W to have RMS Norm 1.

④ Positional Encoding:

① In paper, sinusoidal pattern is used, but now no longer used

② NoPE ③ Learned Relative Positional Encoding.

$$\text{score } s_i = \frac{\langle \vec{q}_t, \vec{e}_i \rangle}{\text{not}}, \quad p_{t,i} = \frac{e^{s_i}}{\sum_j e^{s_j}} \Rightarrow \frac{e^{s_i + (bt-i)}}{\sum_j e^{s_j + (bt-i)}} \begin{matrix} \checkmark \text{ learn} \\ \checkmark \text{ these!} \end{matrix}$$

④ RoPE - Modify \vec{k}, \vec{q}_t based on their absolute position so that the score depends on relative position.

$$\langle \vec{k}, \vec{q}_t \rangle = \|\vec{k}\|_2 \|\vec{q}_t\|_2 \cos \theta \quad \text{Rotation: } \begin{bmatrix} \cos \omega t & -\sin \omega t \\ \sin \omega t & \cos \omega t \end{bmatrix}$$

$$\text{So } \langle R_t \vec{q}_t, R_i \vec{k}_i \rangle = \vec{q}_t^T R_t^T R_i \vec{k}_i$$

$\text{Rotate } -\omega t \quad \text{Rotate } \omega i$

$$\text{So design } M = \begin{bmatrix} M_1 & M_2 & \dots & M_d \end{bmatrix} \quad \text{for } q_t, k \in \mathbb{R}^d, M_j = \begin{bmatrix} \cos(2\pi f_j \cdot t) & -\sin(2\pi f_j \cdot t) \\ \sin(2\pi f_j \cdot t) & \cos(2\pi f_j \cdot t) \end{bmatrix}$$

But f_j are different.

Where f_j is fixed and stands for frequency

Different f_j freq \rightarrow Different timescale. Low f_j channels \rightarrow long range relative info.

② Back to GPT-3: Decoder-Only model:

- Focus: Generate the next token in a sequence.
- Masked / Causal Self attention training: Attn only to current and previous
- Inference: Sample token and auto-regressive generation
 - Can reuse much of the computation for attn
- Sampling Strategy: Beam search / Top-k sampling
- Size: 96 Transformer Layers, 96 heads in self-attention layers

③ BERT: Encoder Only

Goal: Learn general info about statistics of language

And use fine-tuning to adapt to solve specific task

Self-Supervision For Pretrain: ① Masked Token Prediction

<cls>: Start token <sep>: token to separate sentences

- Sent A Sent B: Predict if Sent B immediately follows A
→ Predict this on the sep token.

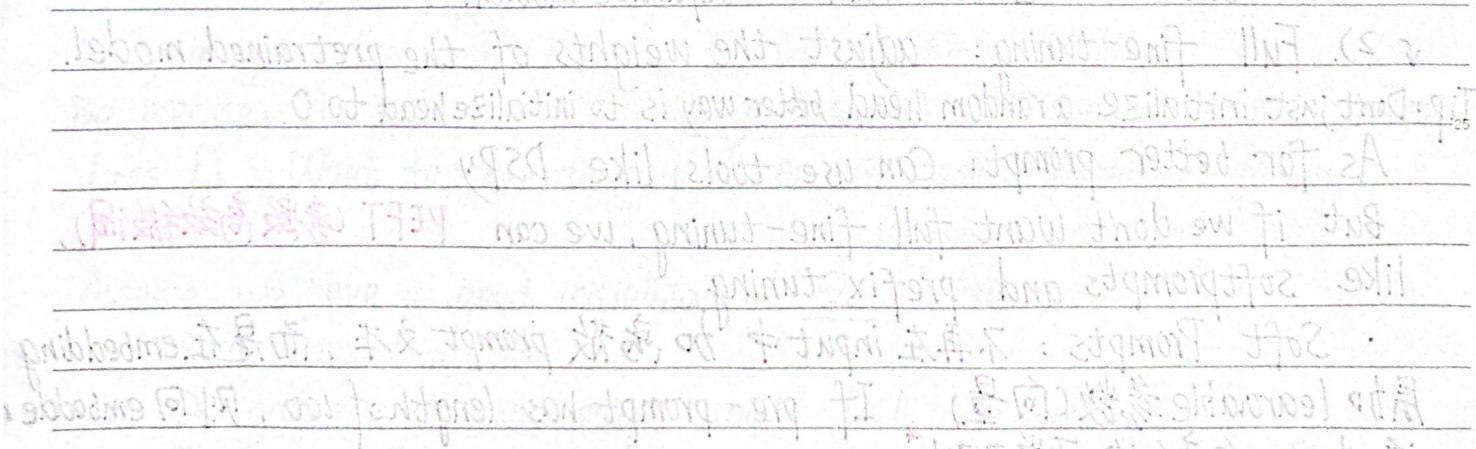
④ MoE: First: Gated Linear Units:

$$GLU(x) = (W_{xtb}) * \text{sigmoid}(V_{xtc})$$

$$\text{SwiGLU}(x) = (W_{xtb}) * \text{Swish}(V_{xtc}) \quad \text{where } \text{Swish}(x) = x \cdot \text{sigmoid}(\beta x)$$

So Mixture of Experts:

Router is designed to Add \rightarrow Router \rightarrow FFN₁ \rightarrow FFN₂ \rightarrow \otimes learnable weights
choose only one or few Experts to activate. Determine which expert to handle only active experts have non-zero weights



ICL / Prompting & Fine-tuning

GPT-style models: Recall Pretraining: Next-token Prediction

△ Data mix matters: 预训练数据组成比例重要!

* Teacher-forcing (Input: GT, not sampled token)

And for inferencing: Autoregressive:

Temp*

→ De-Embedder → Softmax → Sampler → next input token

$$\ast \quad e^{\frac{1}{T}x[i]} / \sum_i e^{\frac{1}{T}x[i]}, T: \text{temp}$$

GPT-3 (2020): With even bigger model, also get in-context learnings.

E.g. Red: Ruby; Blue: Saphire

But now: Can solve a learning problem with a few examples without any finetuning or gradient descent. \Leftrightarrow Pure Prompting.

Challenge: all training data fit in context.

Two core problems: ① How to find better prompt? 更容易被提出来

② How to build model that is more promptable?

↑: Instruction Tuning & SFT

Input Prompt tokens and [do loss on generated response token]

* Masked Loss: No Loss on Question

But take one step back: Fine-tuning or using a pretrained model?

① Pure prompting: no training. Easy to do, but might not work as well

② Linear Probing: treat model as embedder, for a feature extractor.

Do classical ML to train a separate model.

③ Full fine-tuning: adjust the weights of the pretrained model.

Tip: Don't just initialize a random head better way is to initialize head to 0.

As for better prompt: Can use tools like DSPy

But if we don't want full fine-tuning, we can PEFT (参数高效微调), like softprompts and prefix tuning

* Soft Prompts: 不再在 input 中加离散 prompt 文本, 而是在 embedding 层加 learnable 参数(向量). If pre-prompt has length of 100, 则向 embedder 添加 100 个新的可学习层

*: 指的 Vocab 扩充: $W_{vocab}^{D \times V} \xrightarrow{\text{concat}} W_{vocab}^{D \times (V+100)}$

And during training, $W^{D \times V}$ freeze, and $W^{D \times 100}$ (only) are learnable.

Keep GPT model parameters frozen as well.

这样 model 会让这 100 个向量来激发模型完成特定任务, 而不是设计 prompt

• Prefix Tuning: Soft Prompts 之上, 还调整每一层 KV cache.

LoRA:

Idea: Why not do parameter-efficient finetuning beyond soft-prompt?

Soft Prompt & Prefix Tuning 虽高效, 但并未改变权重矩阵。而 LoRA

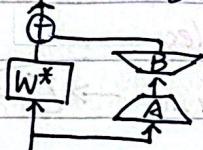
允许以极小成本修改模型深层权重 (e.g. QKV & MLP)

In pretrained model, weight-matrices W^* : We update: $W = W^* + \Delta W$

where ΔW is in rank $r \ll \text{rank}(W^*)$

∴ Let $\Delta W = BA$, $B \in R^{W_{in} \times r}$, $A \in R^{r \times W_{out}}$

Now:



⚠ A: Xavier Init

B: Zero Init (训练开始时, 模型

行为与原来一致, 避免噪声使模型崩塌)

Δ: learning rate of LoRA can be different from the one in pretraining

Meta-Learning: Learning to learn \Rightarrow 让它擅长于“被微调”

Question: If we have a new task, what's the best base model?

① Randomly Init \rightarrow No !! ② General Foundation Model

③ MAML (Model-Agnostic Meta Learning)

Need: { a. Collection of tasks : Data & Loss

b. Approach of fine-tuning

c. Approach of eval

Key insight: Learning Process of SGD is like a RNN

Be precise: A task i has Training Data D_i , Training Loss L_i , Test Loss \tilde{L}_i . Want to train to do well on all tasks from this family.

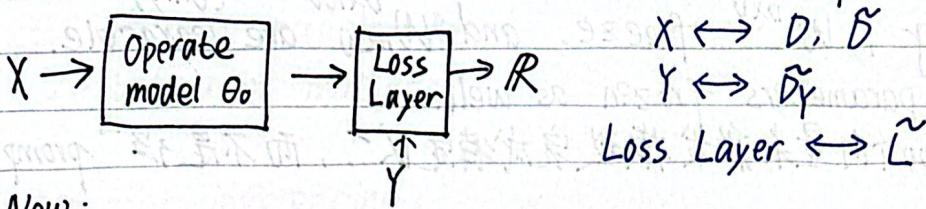
Assume we have a good initialization Θ_0 . How to use it?

1) Start model at Θ_0

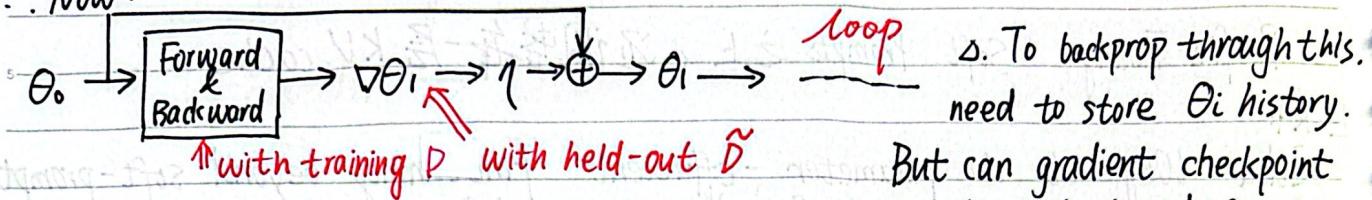
2) Do gradient descent steps with D, L to get Θ_{final}

3) Evaluate Θ_{final} with held-out data \tilde{D} and test loss \tilde{L}

Fit these procedures (1, 2, 3) into our standard form:



Now:



* MAML insight: That's just a deep-model. Can backprop through it

i.e., given a gradient $\nabla \theta_0$ on initial condition

Can take a step in that direction for a better condition

So we can:

- | outer-loop on tasks
- | inner-loop on batches/examples

Task 1: $\theta_{0,0} \xrightarrow{\text{Backprop}} \theta_{1,0} \rightarrow \theta_{1,1} \rightarrow \theta_{1,2} \rightarrow \dots \rightarrow \theta_{1,N_1} \rightarrow \tilde{L} \rightarrow R$

$\downarrow \text{apply}$

2: $\theta_{0,1} \rightarrow \theta_{2,0} \rightarrow \theta_{2,1} \rightarrow \dots \rightarrow \theta_{2,N_2} \rightarrow \tilde{L} \rightarrow R$

3:

$\downarrow \theta_{0,\text{final}}$

Two variants: ① Reptile: Avoid backprop through backprop

$\Rightarrow \text{Gradient } \approx \theta_{\text{final}} - \theta_{\text{initial}} (\theta_0)$

② ANIL / MetaOptnet / R2OZ: Optimize for linear probes

E.g., for regression, have closed form formulas for head. Use closed form solution to compute $W^*(\text{head})$, and it is differentiable, so through it, gradients can back prop to feature extractor Φ .

Δ: Catastrophic Forgetting: When fine-tuning, model forgets how to do what it knew how to do.

Key Practical Solution: During fine-tuning, mix in some pre-training style, like 10%

↳ Can be thought of a kind of overfitting

(VAE&Diffusion).

Generative model & {Post-training and test-time compute}

Generative Model :

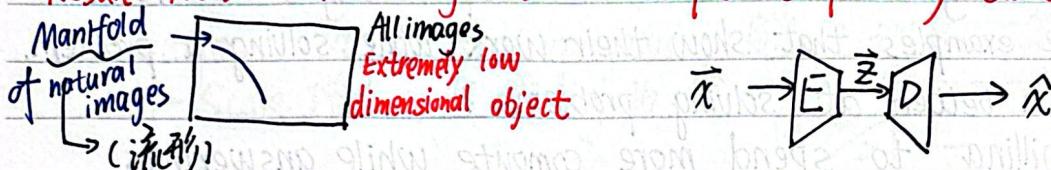
Conditional: Most practically useful

E.g. Prompt in LLM/Images/Video Generative Model

Ideas that don't work:

A) Use a classifier: Image: $\mathbb{R}^{50 \times 50 \times 3}$ \rightarrow Classifier $\rightarrow \mathbb{R}$ "cat score"

Try: Random Uniform Image followed by Gradient Ascent on the cat score.
Result: Noise-Like images that classifier confidently classifies as a cat



B) Use an autoencoder : labels are \vec{x}_i itself, bottleneck in the middle.

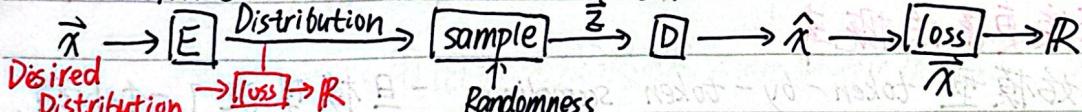
Traditional Perspective : Decoder is scaffolding. 但推理时输入的 \mathbf{z} 通常会得到模糊图像 \Leftrightarrow 训练时 \mathbf{z} 分布与随机采样分布不同

VAE approach: 3 key ingredients:

- 1) make \tilde{z} random during training too.

- 2) Add a loss on distribution of \bar{z}

- 3) Make this work with SGD

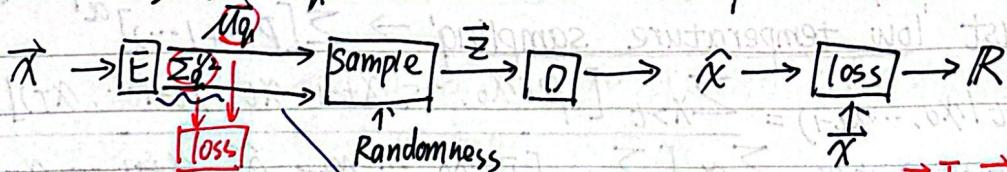


Loss on Distributions: KL Divergence : $KL(Q||P) = \int Q(z) \ln \frac{Q(z)}{P(z)} dz$

Our choice for distribution: $N(0, I_k)$

But directly sample from $z \sim N(\mu, \Sigma)$ cuts gradients flow

$\Rightarrow z = \mu_0 + \sum_{q_n}^{1/2} \nu$, ν is a noise from standard normal distribution.



$$KL(N(\mu_q, \Sigma_Q) || N(0, I_k)) = \frac{1}{2} \text{Tr}(\Sigma_Q) + \mu_q^T \mu_q - k - \log(\det \Sigma_Q).$$

$(\Sigma_Q = \Sigma_Q^{y_2} (\Sigma_Q^{y_2})^\top) \Leftrightarrow$ Enforce semi-definite

- New tricks:
- 1) Using a KL Divergence Regularizer on a distribution
 - 2). Treating Sampling in a way that allows gradients to just pass through it.
 - 3). Accepting the stochasticity of random noise in sampling as just more stochasticity in SGD-style optimization.

Post-training including RLVR: (RL with verifiable rewards)

Recall basic SFT for instruction following: Masked loss (on generated tokens, not questions)

For LLMs, one key advance was "**chain-of-thought**". So it is useful to have examples that show their work while solving a problem.

How to make a better at solving problems?

A) Be willing to spend more compute while answering

B) Train it to be better

"by step"

Test-time Compute:

- 0) Oldest Approach: Pure Promptings. "Think step"

- 1) Repeated generation. Generate N answers, and choose via **majority vote**.
- 2) Sample better.

Δ : Observation on RLVR: empirically show better performance & generalization

但人们注意到，RLVR生成的这些高质量答案，在base model (no RLVR) 中其实也拥有比错误答案更高的概率

难题：若从原始模型 token-by-token sampling，一旦模型犯了一个小错，后面推理会完全崩溃，不知所措 (flounder)

How to solve it:

a): old approach: beam search

b): Alternative: Don't sample from $P(x)$, sample from $(P(x))^\alpha$ ($\alpha > 1$)

$\alpha > 1$: **lower temperature**, 会放大高概率选项，压制低概率项

This is not just '**low temperature sampling**' $\rightarrow \sum [P(x_{t+1} \dots)]^\alpha$

It's: $P_{\text{desired}}(x_t | x_0, \dots, x_{t-1}) = \frac{\sum_{x_{t+1}} [P(x_0, \dots, x_{t-1}, x_t, x_{t+1}, \dots, x_T)]^\alpha}{\sum_{x_{t+1}} [\sum_{x_{t+2}} [P(x_0, \dots, x_{t-1}, \hat{x}_t, x_{t+1}, \dots, x_T)]^\alpha]}$

对整个序列联合概率作规范化

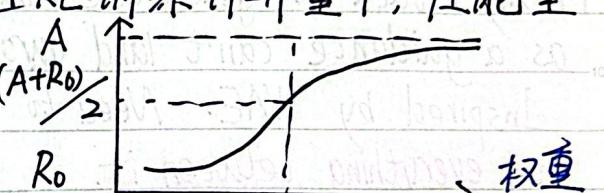
这被称为 power sampling, 能更好地从全局角度筛选高质量的完整推理路径

To sample from $p(x)^\alpha$, in the paper, use MCMC sampling skill.
 以前：模型一句话写完结束，写得烂也认
 现在：模型先写个草稿，然后利用 MCMC 算法，反复地“回退”到句子的某个中间位置，尝试重新生成后缀
 结果：只有当新生成的后缀让整句话的整体概率（经过 α 放大后）变得更高，才会保留修改。

RLVR: Reinforcement Learning with Verifiable Rewards (Turn to 'train it better')

Observation: RL's scaling laws: 随 RL 训练计算量↑, 性能呈 Log-Scale 增长, 且遵训曲线:

$$R_c = R_0 + \frac{A - R_0}{1 + (C_{mid}/C)^B}$$



C_{mid} : 达到总增益 50% 时所需的计算量

Unpacking: $\mathbb{E}_{\theta} \text{Scale}_{RL}(\theta) = \mathbb{E}_{\substack{x \sim D \\ y_i | i=1 \sim G^{\text{old}} \\ \pi_{\text{gen}}(\cdot | x)}} \left[\sum_{g=1}^G \sum_{i=1}^{|\mathcal{Y}_g|} \text{sg}(\min(p_{i,t}, \epsilon)) \hat{A}_i^{\text{norm}} \right]$

目标: 用旧生成策略生成 G 个样本 y_i 归一化 token-level 当前策略生成 y_i 的 token 概率

where $p_{i,t} = \frac{\pi_{\text{train}}^{\theta}(y_i, t)}{\pi_{\text{gen}}^{\text{old}}(y_i, t)}$ 是 Importance Sampling, 因为 π_{gen} 生成数

$\hat{A}_i^{\text{norm}} = \hat{A}_i / \hat{A}_{\text{std}}$ 是优势函数, 通常是 reward - baseline

故 Objective 本质: 用旧 π 生成数据不断更新 π_{train} , 这样便不用

更新一次 π_{train} 后再采样数据

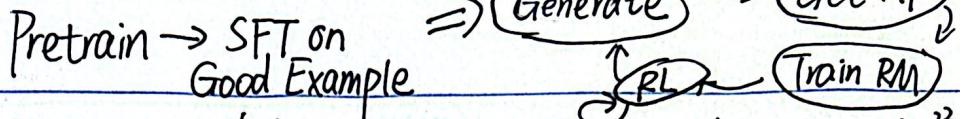
上述公式的重要基石: $\nabla_{\theta} \mathbb{E}_{y \sim p_{\theta}} [F(y)] = \mathbb{E}_{y \sim p_{\theta}} [F(y)] \nabla_{\theta} \log p_{\theta}$

即最大化期望回报, 不需对奖励函数求导, 也可对 $\log p_{\theta}$ 求导并以 $F(y)$ 作为权重. ($\hat{A} \Rightarrow F(y)$; $\log \pi \Rightarrow \log P$)

KL Regularization 也可以帮忙。如有 policy π_{ref} 微调, 希望是在 $r(x)$ 上多得分, 同时不希望 π policy 大变局:

故 Objective: $\underset{\pi}{\operatorname{argmax}} \mathbb{E}_{\substack{y \sim D \\ y \sim \pi(\cdot | x)}} [r(x, y)] - \beta D_{\text{KL}} [\pi || \pi_{\text{ref}}]$

如 GRPO 这个有 KL Regularization



RLHF (human feedback): 三步走: ① SFT ② Reward Model: 收集人类对 model output 的 Binary Comparison (人类不擅长给 real-value)。RM 被训练来预测 哪个回答更好 ③ 使用 RM 为 ref., do reinforcement learning, e.g. $\text{argmax}_{\pi} E[r(x, y)] - \beta DKL[\pi || \pi_{ref}]$

Direct Preference Optimization

DPO: 直接在人类偏好数据上的对数似然优化, 避免 RM 训练。用 π_θ 替代 $L_{DPO}(\pi_\theta; \pi_{ref}) = \bar{E}_{(x, y_w, y_l)} [\log \frac{\pi_\theta(y_w|x)}{\pi_{ref}(y_w|x)} - \log \frac{\pi_\theta(y_l|x)}{\pi_{ref}(y_l|x)}]$
 y_w : Winner 回答; y_l : Loser 回答 Sigmoid

Finally, back-to generative model: diffusion. Recall that: Use a classifier as a guidance can't land on our desired manifold.

Inspired by VAE: Need to see noisy things during training, and why not everything between it.

Encoder: Forward Process: $x_t = \sqrt{1-\beta_t} x_{t-1} + \sqrt{\beta_t} \epsilon_t$ where $\epsilon_t \sim N(0, I)$ and β_t is a variance schedule.

Decoder: Predict the noise. The loss is MSE between ϵ_t & $\hat{\epsilon}_t(x_t)$.

Sampling Approach: ① DDPM (Denoising Diffusion Probabilistic Models)

关于 forward, 其实可以一步写出 x_t & x_0 关系:

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1-\bar{\alpha}_t} \epsilon, \bar{\alpha}_t = \prod_{s=1}^t (1-\beta_s)$$

$$\text{而 } q_h(x_{t-1}|x_t) = \frac{q_h(x_t|x_{t-1}) q_h(x_{t-1})}{q_h(x_t)}$$

但 $q_h(x_{t-1}|x_t, x_0)$ 由于 $q_h(x_t)$ 通过 Bayes 得到

$$q_h(x_{t-1}|x_t, x_0) = N(x_{t-1}; \tilde{\mu}_t(x_t, x_0), \tilde{\sigma}_t^2 I)$$

$$\text{where } \tilde{\mu}_t(x_t, x_0) = \frac{\sqrt{\bar{\alpha}_t}}{1-\bar{\alpha}_t} \beta_t x_0 + \frac{\sqrt{\bar{\alpha}_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} x_t$$

$$\textcircled{2} \text{ DDIM: } x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{x}_0 + \sqrt{1-\bar{\alpha}_{t-1}} \frac{x_t - \sqrt{\bar{\alpha}_t} \hat{x}_0}{\sqrt{1-\bar{\alpha}_t}} \epsilon_t$$

若 $\epsilon_t = 0$:

$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \left(\frac{x_t - \sqrt{1-\bar{\alpha}_t} \epsilon_0(x_t)}{\sqrt{\bar{\alpha}_t}} \right) + \sqrt{1-\bar{\alpha}_{t-1}} \epsilon_0(x_t)$$

$$\text{DDIM} \in \lambda \leftarrow \frac{\sqrt{t}}{\sqrt{t-\Delta t} + \sqrt{\epsilon}}$$

$$x_{t-\Delta t} \leftarrow x_t + \lambda (f_0(x_t, t) - \underline{x})$$

return \underline{x}

图像

△ 假设 $f_0(x_t, t)$ 预测 t 时刻去噪后