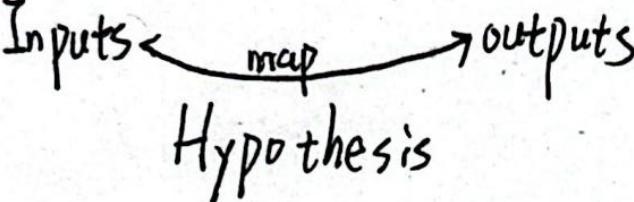


before we came into next section,

ML - consist of 3 parts:

1. Inputs $\xrightarrow{\text{map}}$ outputs
Hypothesis
2. Loss function : specifies how well Hypothesis perform
3. Optimization method : minimize the sum of losses

(linear to non-linear

$$h_{\theta}(x) = \theta^T \phi(x) \quad \theta \in \mathbb{R}^{d \times k} \quad \phi: \mathbb{R}^n \rightarrow \mathbb{R}^d$$

我发现他把单独的 θ^T 叫作 linear classifier

老师给出

$$\theta^T \sigma(W^T x) \neq \tilde{\theta}^T x$$

$W^T \in n \times d$, $\sigma: \mathbb{R}^d \rightarrow \mathbb{R}^d$, σ is a nonlinear function

How to create ϕ ?

1. manual engineering i.e. W be random Gaussian, σ be cosine
→ random Fourier Features
2. "learn itself"

Now we want to optimize both θ and W

哈哈

原话: Cease pretending any requirement on depth

译: 不线性就是 Deep learning!

Beyond "Just not"

So "deep learning" just means "ML using neural network"

A neural network is a particular type of hypothesis class

multiple, parameterized, differentiable functions (a.k.a "layers") composed together from inputs to outputs

定义:

$$h_{\theta}(x) = W_2^T \sigma(W_1^T x), \quad \theta = \{W_1, W_2\}$$

σ is e.g. sigmoid, ReLU

batch form:

$$h_{\theta}(X) = \sigma(XW_1)W_2 \quad X \in \mathbb{R}^{m \times n}$$

Theorem:

Two-layer-network (one-hidden-layer-network)

Is Universal function Approximation

可以无限拟合任何平滑函数 $f: \mathbb{R} \rightarrow \mathbb{R}$

数学表达: given closed region DCR

for any $f: \mathbb{R} \rightarrow \mathbb{R}$ and $\epsilon > 0$

$$\max_{x \in D} |f(x) - \hat{f}(x)| \leq \epsilon, \text{ arbitrarily closely!!}$$

课上给了感性的证明哈哈

注意名词:

layers, activations, neurons

$Z_1 = X, Z_{i+1} = \sigma_i(Z_i W_i), i=1, \dots, L, L\text{-layer NN}$

~~h_θ(x)~~ $Z_{L+1} = h_{\theta}(x), Z_i \in \mathbb{R}^{m \times n_i}, W_i \in \mathbb{R}^{n_i \times n_{i+1}}$

手动链式求导：

$$\frac{\partial \text{ce}(\sigma(xW_1)W_2, y)}{\partial W_2} = \frac{\partial \dots}{\partial \sigma(xW_1)W_2} \cdot \frac{\partial \sigma(xW_1)W_2}{\partial W_2}$$

$$= (S - I_y) \cdot \sigma(xW_1)$$

\downarrow
 $m \times k$, 每行是 hot-encoded
 $m \times k$, 每行是梯度 (梯享)

embarrassingly

$$= [\sigma(xW_1)]^T (S - I_y)$$

$d \times m$ $m \times k$

$$\frac{\partial \text{ce}(\sigma(xW_1)W_2, y)}{\partial W_1} = \frac{\partial \dots}{\partial \sigma(xW_1)W_2} \cdot \frac{\partial \sigma(xW_1)W_2}{\partial \sigma(xW_1)} \cdot \frac{\partial \sigma(xW_1)}{\partial xW_1} \cdot \frac{\partial xW_1}{\partial W_1}$$

$$= (S - I_y) \cdot W_2 \cdot \sigma'(xW_1) \cdot X$$

notice $W_1 \in \mathbb{R}^{n \times d}$ \downarrow just a scalar

embarrassingly

$$= X^T ((S - I_y) \cdot W_2^T \circ \sigma'(xW_1))$$

element-wise multiplication

哈哈，Tim Dettmers 说：

the biggest take-away from all of this is you just want to invest all your time learning automatic differentiation and don't do it anymore.

假定你对 G_i 熟悉了, $G_i = \frac{\partial \ell(Z_{l+1}, y)}{\partial z_i}$

$$G_{i+1} = G_{i+1} \cdot \frac{\partial Z_{i+1}}{\partial z_i}$$

$$= G_{i+1} \cdot \frac{\partial \delta(Z_i W_i)}{\partial Z_i W_i} \cdot \frac{\partial Z_i W_i}{\partial z_i}$$

$$= G_{i+1} \cdot b'(z_i W_i) \cdot W_i \quad 1)$$

$$\text{注: } Z_{i+1} = b_i(Z_i W_i)$$

再

定义: $G_i = \frac{\partial \ell(Z_{l+1}, y)}{\partial z_i} \quad z_i \in \mathbb{R}^{m \times n_i}$

$$= \nabla_{z_i} \ell(Z_{l+1}, y) \quad 2)$$

embarrassingly, according 1) and 2)

$$G_{i+1} = [G_{i+1} \circ b'(z_i W_i)] \cdot (W_i)^T$$

另外, with respect to W_i not Z_i

$$\nabla_{W_i} \ell(Z_{l+1}, y) = \underline{G_{i+1} \cdot \frac{\partial Z_{i+1}}{\partial W_i}}$$

$$G_{i+1} \cdot \frac{\partial b(z_i W_i)}{\partial W_i}$$

$$= G_{i+1} \cdot b'(z_i W_i) \cdot z_i$$

embarrassingly

$$= (z_i)^T (G_{i+1} \circ b'(z_i W_i))$$

可引出 ~~gradient~~ 计算方法:

Backpropagation

I. Forward Pass

$$\text{Init: } Z_1 = X$$

$$\text{Iter: } Z_{i+1} = \delta(Z_i \cdot W_i)$$

II. Backward Pass

$$\left\{ \begin{array}{l} \text{Init: } G_{L+1} = \nabla_{Z_{L+1}} \ell(Z_{L+1}, y) = S - Iy \\ \text{Iter: } G_i = (G_{i+1} \circ \delta'_i(Z_i W_i)) \cdot W_i^T, i = L, \dots, 1 \end{array} \right.$$

and

$$\nabla_{W_i} \ell(Z_{L+1}, y) = Z_i^T (G_{i+1} \circ \delta'_i(Z_i W_i))$$

just the chain rule, huh? also need stored

calculated in forward pass and stored

which indicates trade-off: more efficiently gradient,
more memory needed

最后提了下，所有这些 hacky 的操作，与 vector Jacobian product 有关

鉴于笔记访问的方便，我们仍将继续 CMU 课程作笔录

几 k (labeled) m (training numbers) 记住它们！

hypothesis function $h(x) = \begin{bmatrix} h_1(x) \\ \vdots \\ h_k(x) \end{bmatrix}$ $h: \mathbb{R}^n \rightarrow \mathbb{R}^k$

linear form:

$$h_{\theta}(x) = \theta^T x \quad \theta \text{ 是参数矩阵} \\ (k \times n) \cdot (n \times 1) \quad x \text{ 是一个 column vector}$$

code neater notation:

$$x \in \mathbb{R}^{m \times n} = \begin{bmatrix} -x^{(1)^T} - \\ \vdots \\ -x^{(m)^T} - \end{bmatrix} \quad \text{转置是因为 } X^{(i)} \text{ 是 column vector}$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(k)} \end{bmatrix}$$

$$h_{\theta}(\bar{X}) = \begin{bmatrix} -h_{\theta}(x^{(1)})^T - \\ \vdots \\ -h_{\theta}(x^{(m)})^T - \end{bmatrix}$$

$$= \begin{bmatrix} [\theta^T x^{(i)}]^T \end{bmatrix}$$

$$= \begin{bmatrix} x^{(i)^T} \cdot \theta \end{bmatrix} = X \cdot \theta \quad (m \times k)$$

Loss function

① simplest one: classification error

$$l_{\text{err}}(h(x), y) = \begin{cases} 0 & \text{if } \arg\max_{\text{index}} h(x) = y \\ 1 & \text{otherwise} \end{cases}$$

不好优化，not differentiable

② Softmax or Cross-Entropy loss

转换成更“像”概率的样子：

→ make them all positive and sum to one

→ exponentiate and normalize

$$\rightarrow \frac{\exp(h_i(x))}{\sum_{j=1}^k \exp(h_j(x))} = p(\text{label} = i) = z_i$$

$$\Leftrightarrow z \equiv \text{normalize}(\exp(h(x))) \equiv \text{Softmax}(h(x)) \in \mathbb{R}^k$$

实际的计算并不会直接到上面的表达式。

$$l_{\text{ce}}(h(x), y) = -\log \text{Softmax}(h_y(x)) = -h_y(x) + \log \sum_{i=1}^k \exp(h_i(x))$$

(y在这里像是某种
类)

Think ML-algorithm as a method for solving
associated optimization
problem

~~ML~~ → Core:

$$\underset{\theta}{\text{Minimize}} \left\{ \text{Average} \left\{ f(h_{\theta}(x^{(i)}), y^{(i)}) \right\} \right\}$$

\downarrow loss-func \downarrow hypothesis

Any Supervised algorithm is to solve this form problem

i.e. softmax regression:

$$\underset{\theta}{\text{minimize}} \frac{1}{m} \sum_{i=1}^m \text{ce}(\theta^T x^{(i)}, y^{(i)}) , \text{ which is also minimize } f(\theta)$$

手段：gradient descent (∇)，powers all deep learning!

Though $f(\theta)$ itself is a scalar, the gradient is $n \times k$ matrix
same as $\theta \in \mathbb{R}^n$

$\nabla_{\theta} f(\theta)$, its result belongs to $\mathbb{R}^{n \times k}$

whose direction (sadly $\mathbb{R}^{n \times k}$) most increases $f(\theta)$

So in iteration, we go opposite:

$$\theta := \theta - \alpha \nabla_{\theta} f(\theta)$$

↓
learning rate

如果 m 增长到太大, $f(\theta)$ 中 $\frac{1}{m} \sum_1^m$ 会造成计算负担

→ Stochastic GD (SGD):

sample a minibatch of data $X \in \mathbb{R}^{B \times n}$

轻松的, 对应: $\theta := \theta - \frac{\alpha}{B} \sum_1^B \nabla_{\theta} l(h(x^{(i)}), y^{(i)})$

manually calculate gradient

i.e. the softmax form:

consider $\nabla_h l_{ce}(h, y)$

单独看 $\frac{\partial l_{ce}(h, y)}{\partial h_i} = \frac{\partial}{\partial h_i} (-hy + \log \sum_{j=1}^k \exp h_j)$

$$= -1 \{ i=y \} + \dots$$

~~最好~~ $= -1 \{ i=y \} + P(\text{label} = i)$

$$\nabla_h l_{ce}(h, y) = z - ey$$

for $\nabla_\theta l_{ce}(\theta^T x, y)$, the matrix differential calculus can be
cumbersome

我们粗暴的 treat it as scalar, and check numerically

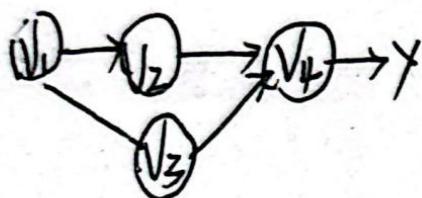
$$\begin{aligned}\frac{\partial l_{ce}(\theta^T x, y)}{\partial \theta} &= \frac{\partial l_{ce}(\theta^T x, y)}{\partial \theta^T x} \cdot \frac{\partial \theta^T x}{\partial \theta} \\ &= (z - ey) \cdot x \\ &= \text{embarrassingly } = x(z - ey)^T\end{aligned}$$

for batch form:

$$\nabla_\theta l_{ce}(X\theta, y) \xrightarrow{\text{观察得}} X^T(z - I_y)$$

Chen's class

- ① forward mode automatic differentiation
- ② DAG 因果图, computational graph
- ③ Reverse mode 非常有趣



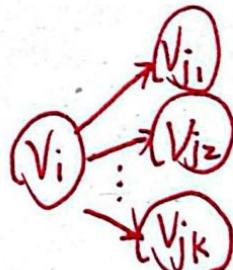
学迷糊了。:-)

$$\text{defi. } \bar{V}_i = \frac{\partial y}{\partial V_i}$$

$$\frac{\partial f(V_2, V_3)}{\partial V_1} = \frac{\partial f(V_2, V_3)}{\partial V_2} \cdot \frac{\partial V_2}{\partial V_1} + \frac{\partial f(V_2, V_3)}{\partial V_3} \cdot \frac{\partial V_3}{\partial V_1}$$

→ Reverse mode 的例子: $\bar{V}_i = \bar{V}_2 \cdot \frac{\partial V_2}{\partial V_i} + \bar{V}_3 \cdot \frac{\partial V_3}{\partial V_i}$

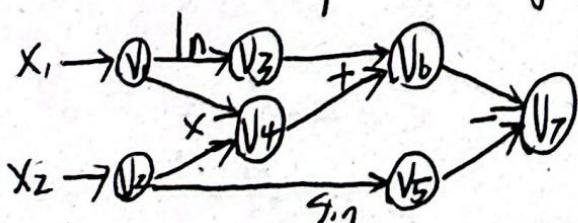
$$\begin{aligned} \text{defi. } \bar{V}_{i \rightarrow j} &\equiv \cancel{\frac{\partial V_j}{\partial V_i}} \\ &\equiv \bar{V}_j \cdot \frac{\partial V_j}{\partial V_i} \end{aligned}$$



这里对应一个 node 有 multiple output 的情况!

$$\bar{V}_i = \sum_{j \in \text{Next}(i)} \bar{V}_{i \rightarrow j}$$

给你画个 computation graph 示意, $y = \ln(x_1) + x_1 x_2 - \sin x_2$



each node is an operation and can represent an immediate value

添加一页：

Numerical differentiation

一般形式： $\frac{\partial f(\theta)}{\partial \theta_i} = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + e_i \epsilon) - f(\theta)}{\epsilon} + O(\epsilon)$

高精度： $\frac{\partial f(\theta)}{\partial \theta_i} = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + e_i \epsilon) - f(\theta - e_i \epsilon)}{2\epsilon} + O(\epsilon^2)$

当然了，这是手算。

which is used for numerical checking

在前面 DAG 的示意图中

我们示例 $\frac{\partial y}{\partial x_1}$ given $x_1=2, x_2=5$

Forward AD trace: → 它的缺点是 input $\in \mathbb{R}^n$, 就要做 n 次

Forward AD
不效率！

$$v_1 = 1 \quad v_2 = 0$$

$$v_3 = \frac{1}{v_1} \cdot v_1 = \frac{1}{2}$$

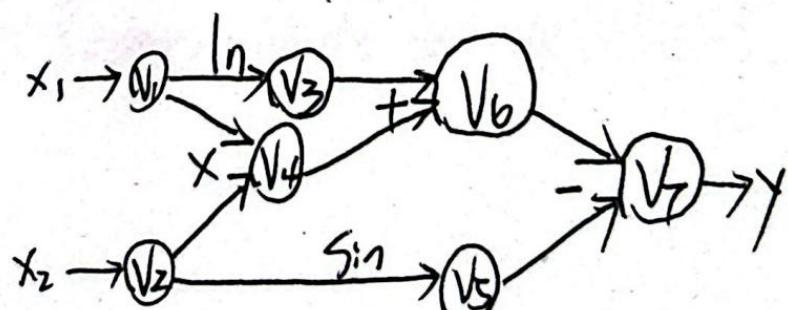
$$\begin{aligned} v_4 &= v_1 \cdot v_2 + v_1 \cdot v_2 \\ &= 1.5 + 2 \cdot 0 \\ &= 5 \end{aligned}$$

$$v_5 = \cos v_2 \cdot v_2 = 0$$

$$v_6 = v_3 + v_4 = 5.5$$

$$v_7 = v_6 - v_5 = 5.5$$

非常糟糕，产生了非常棒可依赖！



as long as i know the former
input node

e.g. Reserve mode AD

还是那个图

$$\frac{\partial y}{\partial V_7} = 1$$

$$\frac{\partial y}{\partial V_6} = \frac{\partial y}{\partial V_7} \cdot \frac{\partial V_7}{\partial V_6} = 1 \cdot 1 = 1$$

$$\frac{\partial y}{\partial V_5} = \frac{\partial y}{\partial V_7} \cdot \frac{\partial V_7}{\partial V_5} = 1 \cdot (-1) = -1$$

$$\frac{\partial y}{\partial V_4} = \frac{\partial y}{\partial V_7} \cdot \frac{\partial V_7}{\partial V_6} \cdot \frac{\partial V_6}{\partial V_4} = 1 \cdot 1 \cdot 1 = 1$$

$$\frac{\partial y}{\partial V_3} = \frac{\partial y}{\partial V_6} \cdot \frac{\partial V_6}{\partial V_3} = 1 \cdot \cancel{1} = 1$$

$$\left\{ \frac{\partial y}{\partial V_2} = \frac{\partial y}{\partial V_5} \cdot \frac{\partial V_5}{\partial V_2} + \frac{\partial y}{\partial V_4} \cdot \frac{\partial V_4}{\partial V_2} = \dots \right.$$

$$\left. \frac{\partial y}{\partial V_1} = \frac{\partial y}{\partial V_3} \cdot \frac{\partial V_3}{\partial V_1} + \frac{\partial y}{\partial V_4} \cdot \frac{\partial V_4}{\partial V_1} = 1 \cdot \cancel{1} + 1 \cdot \cancel{1} = 5.5 \right.$$

$$\Rightarrow \left(\begin{array}{c} \frac{\partial y}{\partial x_2} \\ \frac{\partial y}{\partial x_1} \end{array} \right) = \nabla f, \text{ 所以 Reverse Mode 做一次就够了。}$$

Backprop Vs Reverse mode AD

↓
by extending computational graph

G

- 好处 {
- ① the compation process is another computational graph
可以做梯度的梯度 for free
可拓展性很棒!!!
 - ② output is a computational graph
不需要 run forward
which make process asymmetric

最后提了 - T Reverse mode AD on data structures

"tuple value"? 没怎么听懂

我把伪代码写一遍，可能会清晰些：

```
def grad(out):
```

node-to-be-grad = {最后一个点}

for i in reverse-todo-order(out): → 这个逆序实现经典问题

$$\bar{V}_i = \sum_j \bar{V}_{i \rightarrow j} = \text{sum}(\text{node-to-be-grad}[i]) \rightarrow \text{先把当前 } \bar{V}_i \text{ 算出来}$$

for k ∈ input(i):

$$\bar{V}_{k \rightarrow i} = \bar{V}_i \cdot \frac{\partial V_i}{\partial V_k}$$

~~append~~
node-to-be-grad[k].append($\bar{V}_{k \rightarrow i}$) → 把 i 的前驱结点
算出来，加入维护队列

数学直觉上，每轮到一个 i 时，它的 \bar{V}_i 是全部被计算完毕的

它的那个 R-mode Computational Graph

就是照上面的流程画出来的。

不过是用了 stored memory，弄的就很花。

Fully connected networks

来： $b \in \mathbb{R}^{n_{i+1}}$

$b \cdot \text{reshape}((1, n_{i+1})) \cdot \text{broadcast_to}((m, n_{i+1}))$

newton's method:

directly going to the optimized point.

momentum:

$U_{t+1} = \beta U_t + (1-\beta) \nabla_{\theta} f(\theta_t)$ - β is momentum averaging parameter

$$\theta_{t+1} = \theta_t - \alpha U_{t+1}$$

所有过去的导数都会起作用。

④

momentum 变种：

一开始 warm-up 可能很慢

"Unbiasing momentum"

$$\theta_{t+1} = \theta_t - \frac{\alpha U_{t+1}}{1-\beta^{t+1}}$$

第一步野一些，可能收敛快一些

Nesterov Momentum:

$$U_{t+1} = \beta U_t + (1-\beta) \nabla_{\theta} f(\underbrace{\theta_t - \alpha U_t}_{\downarrow \text{做前瞻预测}})$$

Adam: good optimizer in practice

$$U_{t+1} = \beta_1 U_t + (1-\beta_1) \nabla_{\theta} f(\theta_t)$$

$$V_{t+1} = \beta_2 V_t + (1-\beta_2) (\nabla_{\theta} f(\theta_t))^2$$

$$\theta_{t+1} = \theta_t - \frac{\alpha U_{t+1}}{\sqrt{V_{t+1}} + \epsilon}$$

每一次步长都差不多

几点 tips:

1. Initialization matters

2. Weights doesn't move that much

最后是有趣的 variance 跟随输入现象:

independent random variables: $x \sim N(0, I)$, $w_i \sim N(0, \frac{1}{n} I)$

$$\text{Var}(w^T x) = 1$$

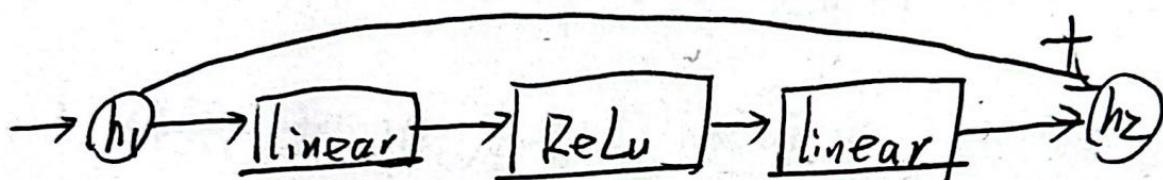
Kaiming normal initialization: $w_i \sim N(0, \frac{2}{n} I)$ 没听懂.

declarative programming 先宣称为计算，再跑
↓ tensorflow.

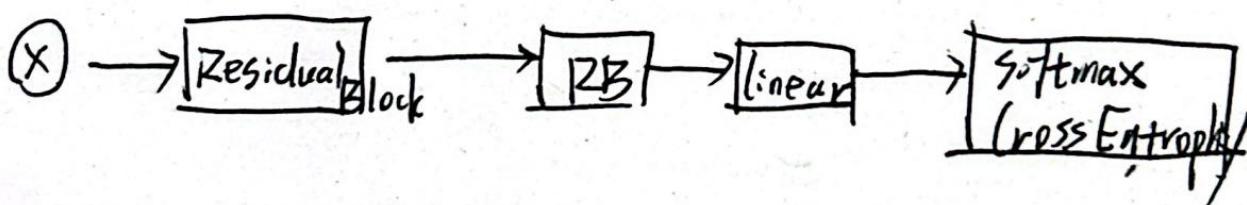
pytorch: imperative
立刻计算

High level modular library components.

0. Residual Block:



1. Multi-layer Residual Net



Deep learning is modular in nature

nn.module 这是把上面 2 个东西凑在一起

things to consider:

1. tensor in tensor out
2. get list of trainable parameters
3. way to initialize the parameters

当然了 Softmax(Cross Entropy) 是 Tensor in . scalar out

Optimizer: 还是 SGD. Momentum 那些

再说 regularization, l2 regularization, 完全没有听懂.

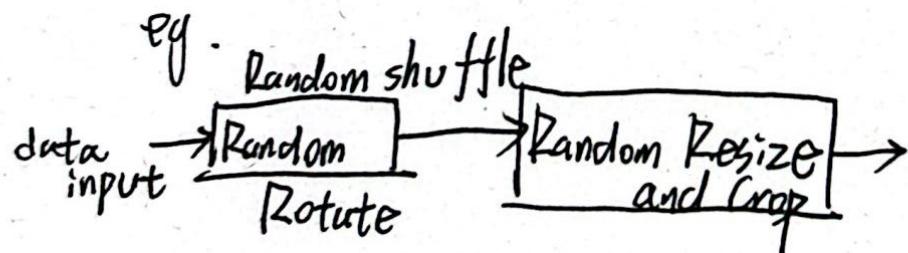
l1 正则化: $\text{loss_func} += \sum |w_i|$

会使很多不重要的 w_i 直接变 0

l2 正则化: $\text{loss_func} += \sum w_i^2$

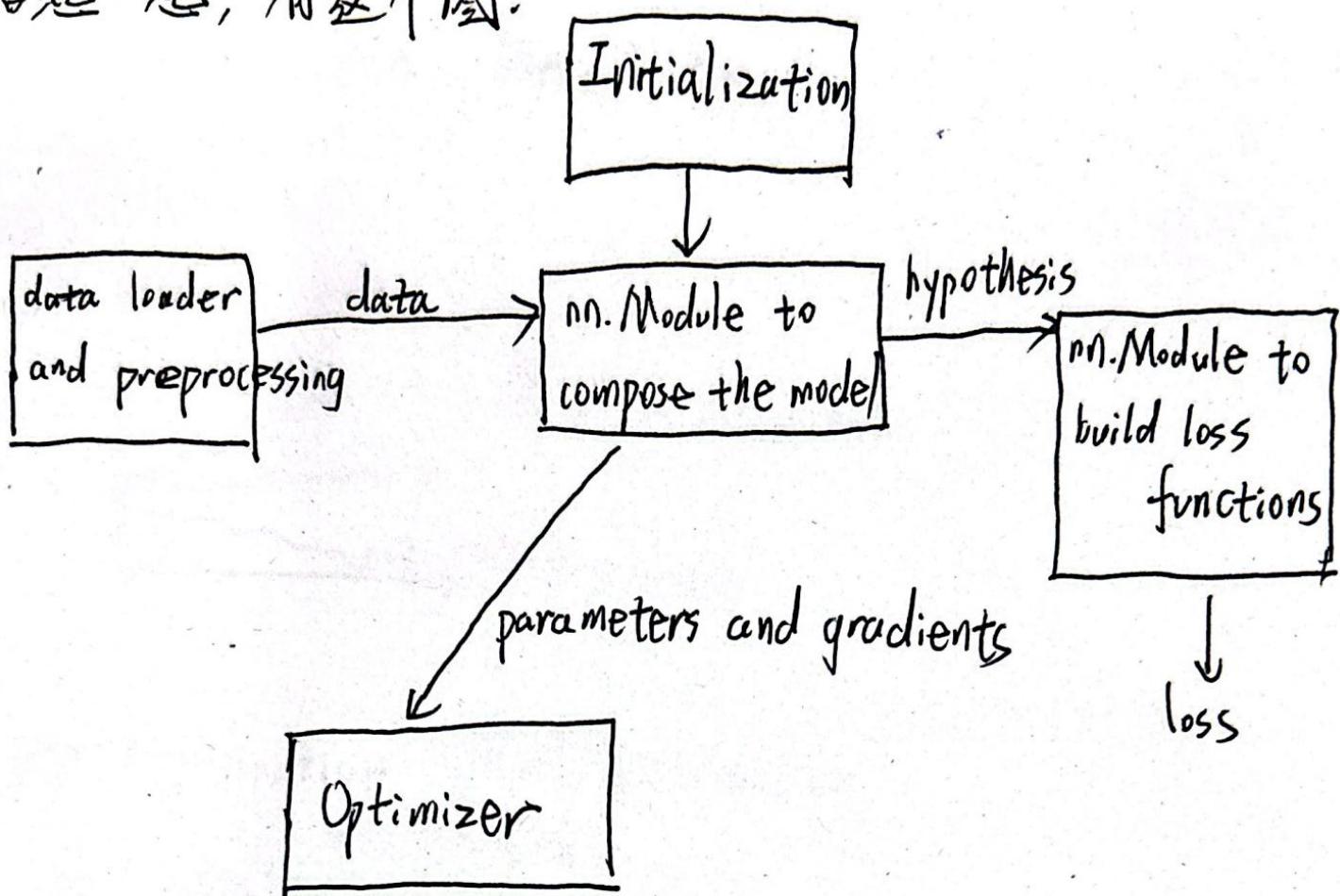
作用: 把 w_i 尽量拉回零点, w_i 过大往往意味着过拟合.

Data load and augmentation



对训练准确也有影响.

合起来，有这个图：



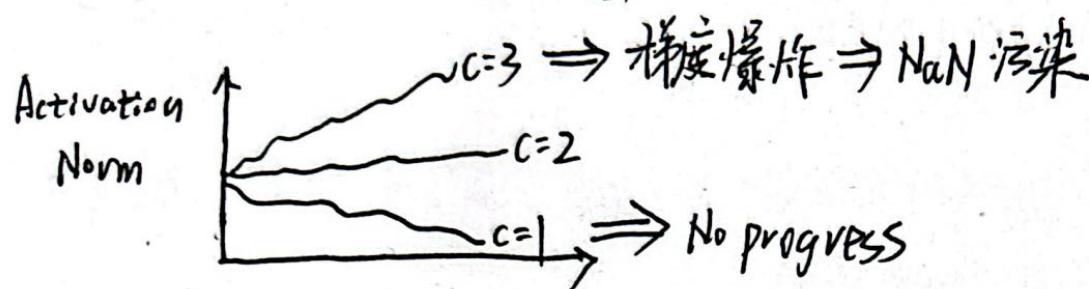
tricks:

Normalization Regularization

一.

权重:

Init: $W_i \sim N(0, \frac{c}{n})$



Layer normalization: 按行来 normalize.

$$\hat{z}_{i+1} = \beta_i (W_i^T z_i + b_i) \xrightarrow{\text{行向量哈哈}}$$

$$z_{i+1} = \hat{z}_{i+1} - E[\hat{z}_{i+1}]$$

$$\sqrt{\text{var}(\hat{z}_{i+1}) + \epsilon} \xrightarrow{\text{极小, 避免 div 0}}$$

Batch normalization: 按列来 normalize
(归一化不知道)

odd的地方是让不同行的样本产生依赖。

二.

~~DNN~~ DNN 是 over-parameterized 的, 因为它的表达力太强, 很容易过拟合, 导致泛化能力不足。

Regularization: Limiting the complexity of function class.
限制表达力, 防止过拟合。

a whole lot of designing DL Sys is to Implicitly regularizing the complexity of functions so we can generalize better.

比如，每次训练，其实 weight 变化不大。

b2 Regularization (or weight decay)

$$\text{minize } \frac{1}{m} \sum_{i=1}^m \ell(h(x_i), y_i) + \frac{\lambda}{2} \sum_{i=1}^L \|W_i\|_F^2$$

梯度下降：

$$W_{i+1} = W_i - \alpha \nabla_{W_i} \ell(\dots) - \lambda W_i = (1 - \alpha \lambda) W_i - \alpha \nabla_{W_i} \ell(\dots)$$

每次都 shrink $(1 - \alpha \lambda)$

Dropout - another common regularization strategy

$$\hat{z}_{i+1} = \beta_i (W_i^T z_i + b_i)$$

→ 除为了保持平均数不变

$$(z_{i+1})_j = \begin{cases} \text{itself_row}_{(1-p)} & \text{with prob } 1-p \\ 0 & \text{with prob } p \end{cases}$$

dropout as a Stochastic approximation

使得DNN不依赖特定的行(神经元), 被迫
学习更鲁棒的特性

差支一毫:

表面看是 Batch Norm reducing internal covariant shift, 训练加速了, 前者真的是后者的原因吗?

并不!

事实是 Batch Norm makes the optimization landscape smoother!

Conventional Operator

定义：略

一个例子，在 image processing

① Gaussian blur：变糊了

② "Grad"：勾画轮廓

$$\left(Z * \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} \right)^2 + \left(Z * \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix} \right)^2$$

在传统图像处理是 pre-fixed
而卷积掩码是可学习的！

Mutichannel - convolution

卷积核 $W_g \in \mathbb{R}^{\text{Cout} \times \text{Cin} \times k \times k}$
 \downarrow
 k is kernel size.

(进一步简化：把 Cout, Cin 也分组)

practical Convolution tricks

1. for kernel size k , pad input with $(k-1)/2$ zeros
in all sides
(when k is odd)

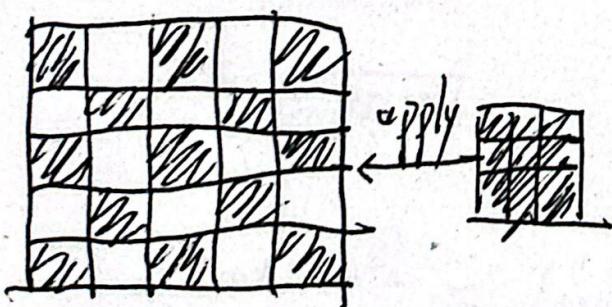
2. Strided Convolution

窗口步长每次移动>1, 像积出来更小.

3. Max pooling

从卷积核中选一个最大的作表.

4. Dilations



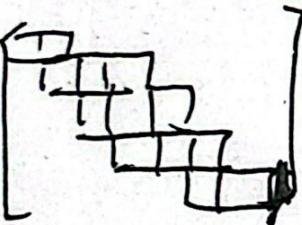
Differentiating Convolution

工程：store all intermediate product before add them up

- 一个 toy 例子

$$z = x * w$$

且 $Z = W^T \cdot x$

$w:$ 

W^T 居然形式上一模一样！

除了 $\underline{[W_1 | W_2 | W_3]}$

变成 $\underline{[W_3 | W_2 | W_1]}$

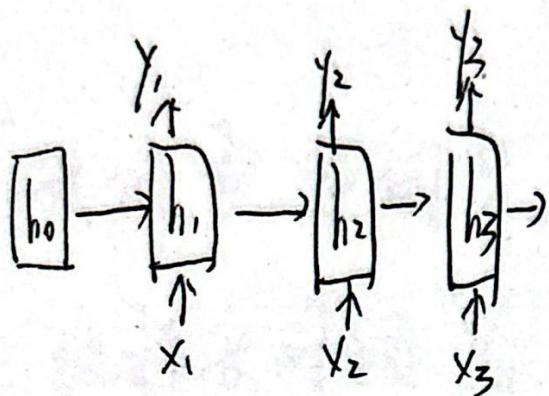
$$\text{有: } \nabla \cdot \frac{\partial \text{Conv}(x, W)}{\partial x} = \text{Conv}(\nabla, \text{flip}(W))$$

如果反过来: $z = \begin{bmatrix} 0 & x_1 & x_2 \\ x_1 & x_2 & x_3 \\ \vdots & \ddots & \vdots \\ x_4 & x_5 & 0 \end{bmatrix} [w]$

这种叫“im2col”阵

Sequence Modeling Recurrent NN

RNN:



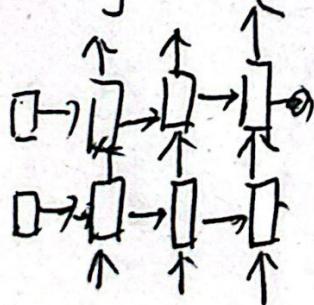
$$h_t = f(W_{hh} \cdot h_{t-1} + W_{hx} \cdot x_t + b_h) \quad y_t = g(W_{hy} \cdot h_t)$$

会忘了 BP TT back propagation through time

先把时间范围一段一段
微分 loss.

(怎么递归微分的?)

Stacking RNNs



就是 h 层加深了。
hidden unit

这个有些缺点：

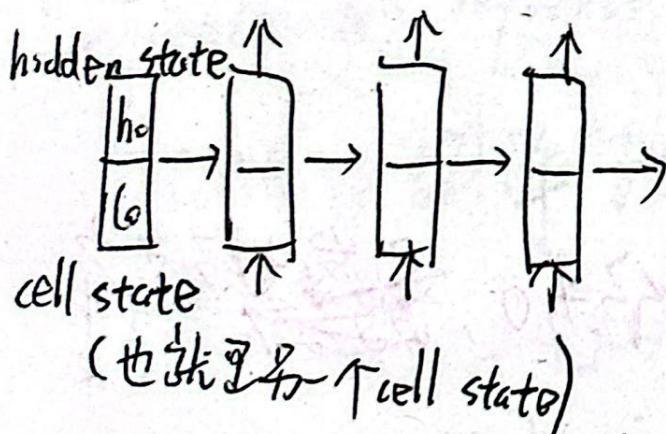
varience一大：梯度爆炸 norm猛增。

Varience太小：Vanishing activation

我们想要的是“long stable range” dependencies
这样长时间 X_i 们都起作用。

ReLU can grow unboundly - why not bounded activation
实际上不起作用。

Long short Term Memory



$$\begin{bmatrix} i_t \\ f_t \\ g_t \\ o_t \end{bmatrix} = \begin{pmatrix} \text{Sigmoid} \\ \text{Sigmoid} \\ \tanh \\ \text{Sigmoid} \end{pmatrix} \left(W_{hh} \cdot h_{t-1} + W_{hx} \cdot x_t + b_n \right)$$

$$\begin{bmatrix} i \\ f \\ g \\ o \end{bmatrix}$$

$$c_t = c_{t-1} \circ f_t + i_t \circ g_t$$

$$h_t = \tanh(c_t) \circ o_t$$

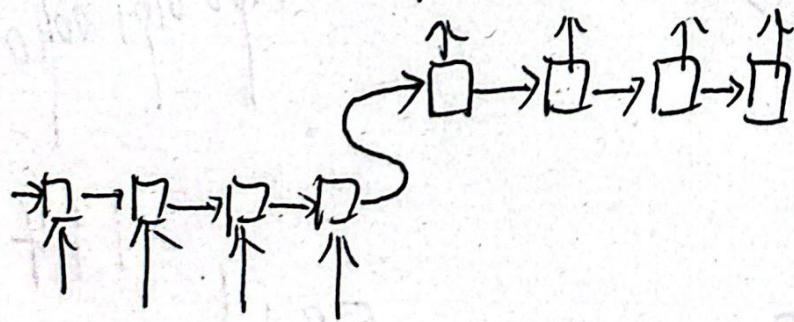
为什么他们能 help with vanishing?

key point:

1. scaling down c_{t-1} , and add a term to i_t .

Sequence-to-Sequence Model

可以作翻译了!



they leverage RNNs in Module!!!

Bidirectional RNNs



微分很难

like translation, we don't only need the current-before things. 这个整个都可以用，

像翻译

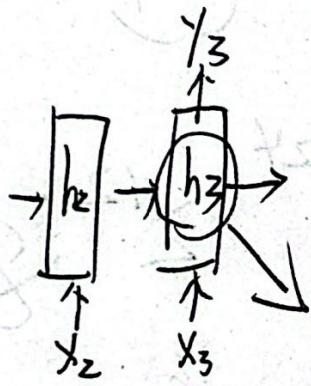
{ Transformers and Attention

{ two approaches to time series modeling

{ 1.1 RNN's "latent state" approach

隐藏层

e.g.



包含了所有前面的信息

RNNs can capture infinite history

Pros: good for long time series prediction

Cons:

Long "compute path" leads to hard to learn

§1.2 "direct prediction" approach

just need a function that can accept
different-sized inputs

Pros: map from past to current
with shorter path

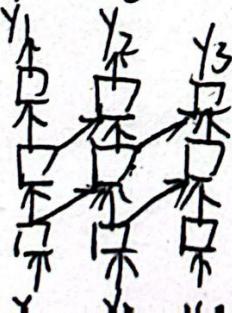
Cons: 每次都要把所有 x_i 輸一遍，
对比之 T-RNN Y, 需要輸入 $x_{current}$

The difference between §1.1. & 1.2 這是 RNNs 和

Transformers

何區別
主要區別

當然了 "direct prediction" 也有 ~~其~~ architecture,
eg. TCNs



"→" 是卷积，
所以这里 CNNs

* 名词: receptive field

显然，卷积层越深，CNNs receptive field
才能大一点

p.s. 为了让卷积层更深，可以增加

显然卷积层的深度跟 kernel size
大小是呈正比例的 receptive field

p.s. 池化 pool 是英文会话 - hha-

22. Self-attention and Transformer

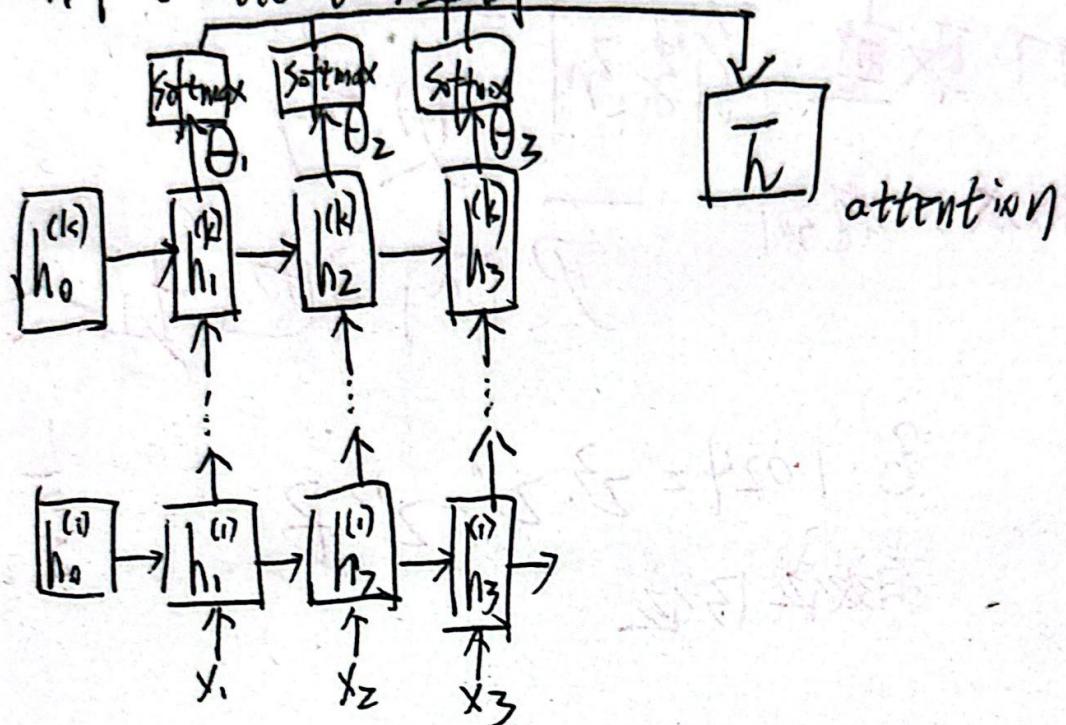
才开一章：

“Attention” in deep learning

前言：

no matter how good these fancy mechanisms skit,
the ultimate state will last less depend on origin
the Influence of first points fade — the ~~exploding~~
or vanishing gradient

RNN 中的 attention 举例



§2.1 Self-attention

Habuta, Detters says 他完全没看出
 K, Q, V 分别为 keys, queries
 values
 何道理

$$\text{SelfAttention} = \text{softmax}\left(\frac{KQ^T}{\sqrt{d}}\right)V$$

where $K, Q, V \in \mathbb{R}^{T \times d}$

$$\text{softmax}\left(\frac{k Q^T}{\sqrt{d}}\right) V$$

理解：

$$K \cdot Q^T \xrightarrow{\text{放大}} \begin{bmatrix} k_i^T q_j \end{bmatrix}$$

a way to measure the similarity between
k_i and q_j

$$\underbrace{\text{softmax}\left(\frac{k Q^T}{\sqrt{d}}\right)}_{T \times T} \underbrace{V}_{T \times d}$$

注意到没有，K在前，Q在后却转置，

它们的项 k_i, q_j 执行向量乘的时候是完整的。

唯独 V 在后面没有转置，粗暴的拆开时间该 V_i，按列来处理

Properties of self-attention

1. permutations

无关性 (K, Q, V 顺序改变, 输出不变)

do same permutation ops to K, Q, V

⇒ output act the same permutation

2. allows influence of K, V, Q over ALL TIMES

without increase parameter count!

(想想TCNs)

~~K, V, Q 是共享参数~~

K, V, Q 不是 Parameter

Nice!
! { Self-attention is just an ops to mix all-time
with no favoritism (because permutationally invariant)
△ △
without increasing parameter count

3. Cons: computation cost is ~~O(T^2 + TD)~~

空间复杂度

$$O(dT^2 + T^2 + T^2 + T \times d \times T) = O(T^2 d)$$

Transformer Block

$$Z^{(i)} \rightarrow Z^{(i+1)}$$

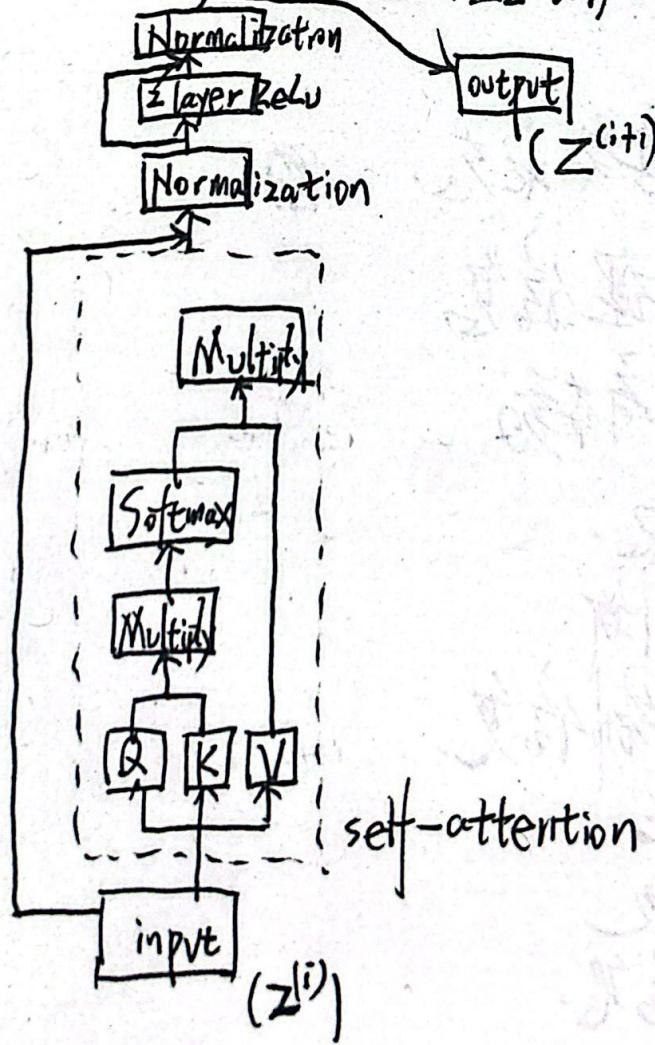
$$\tilde{Z}_1 = \text{SelfAttention}(Z^{(i)}, W_k, Z^{(i)}, W_q, Z^{(i)}, W_z)$$

$$= \text{Softmax} \left(\frac{Z^{(i)} \cdot W_k \cdot W_q^T \cdot Z^{(i)T}}{\sqrt{d}} \right) Z^{(i)} \cdot W_z$$

$$\tilde{Z}_2 = \text{LayerNorm}(Z^{(i)} + \tilde{Z}_1)$$

$$Z^{(i+1)} = \text{LayerNorm}(\tilde{Z}_2 + \text{ReLU}(\tilde{Z}_2 \cdot W_1) \cdot W_2)$$

画图如下：



Transformer:

- Pros:
- full receptive field within a single layer
 - don't add parameter-count

Cons: shuffle 不常用

output depends all inputs

没办法做 time series

解决办法：

Masked self-attention

$$\text{Softmax}\left(\frac{KQ^T}{\sqrt{d}} - M\right) \cdot V, \quad M = \begin{bmatrix} \infty & \\ 0 & \infty \end{bmatrix}$$

会变成

$$\begin{bmatrix} \infty & \\ \text{正端} & \end{bmatrix} \cdot [V]$$

解决 order Invariance 问题

$$\begin{bmatrix} \rightarrow x_i^1 \\ \rightarrow x_i^2 \\ \rightarrow x_i^T \end{bmatrix} + \begin{bmatrix} s_{i,1}(w_1 \cdot 1) & s_{i,1}(w_2 \cdot 1) & \cdots & s_{i,1}(w_d \cdot 1) \\ \vdots & \vdots & & \vdots \\ s_{i,T}(w_1 \cdot T) & s_{i,T}(w_2 \cdot T) & \cdots & s_{i,T}(w_d \cdot T) \end{bmatrix}$$

↓
positional encoding

至于为什么张这样，没讲

最后介绍了 Transformers beyond time series

Transformer, a powerful paradigm used everywhere!