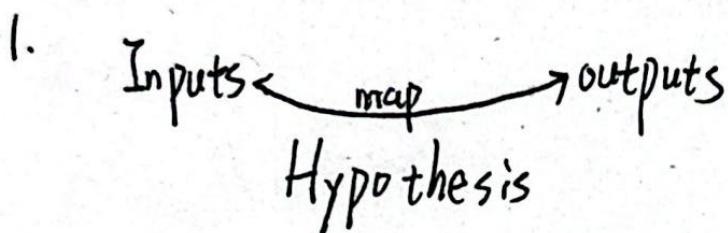


before we came into next section,

ML - consist of 3 parts:



2. Loss function: specifies how well Hypothesis perform

3. Optimization method: minimize the sum of losses



CS 扫描全能王

3亿人都在用的扫描App

(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

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

$W^T \in \mathbb{R}^{n \times d}$, $\sigma: \mathbb{R}^d \rightarrow \mathbb{R}^d$, 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: Please 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)



CS 扫描全能王

3亿人都在用的扫描App

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

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

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



CS 扫描全能王

3亿人都在用的扫描App

手动链式求导：

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

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

\downarrow
 $m \times k$, 每行是 hot-encoded
 $m \times k$, 每行是预测概率

embarrassingly

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

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

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

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

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

embarrassingly

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

\downarrow
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.



CS 扫描全能王

3亿人都在用的扫描App

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

$$\begin{aligned} G_i &= G_{i+1} \cdot \frac{\partial Z_{i+1}}{\partial z_i} \\ &= G_{i+1} \cdot \frac{\partial \delta(Z_i W_i)}{\partial Z_i W_i} \cdot \frac{\partial Z_i W_i}{\partial z_i} \quad \text{注: } Z_{i+1} = \delta_i(Z_i W_i) \\ &= G_{i+1} \circ b'(Z_i W_i) \cdot W_i \end{aligned}$$

等等

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

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

embarrassingly, according 1) and 2)

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

另外, with respect to W_i not Z_i

$$\begin{aligned} \nabla_{W_i} \ell(Z_{l+1}, y) &= G_{i+1} \cdot \frac{\partial Z_{i+1}}{\partial W_i} \\ &= G_{i+1} \cdot \frac{\partial \delta(Z_i W_i)}{\partial W_i} \\ &= G_{i+1} \cdot b'(Z_i W_i) \cdot Z_i \\ &\text{embarrassingly} \\ &= (Z_i)^T (G_{i+1} \circ b'(Z_i W_i)) \end{aligned}$$



CS 扫描全能王

3亿人都在用的扫描App

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

最后提了一下，所有这些 hacky 的操作，与 vector Jacobian product 有关



CS 扫描全能王

3亿人都在用的扫描App

鉴于笔记访问的方便，我们仍将 CMU 课程作笔记

n 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:

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

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

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

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

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



CS 扫描全能王

3亿人都在用的扫描App

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

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

$$\Leftrightarrow z \equiv \text{normalize}(\exp(h(x))) \equiv \text{Softmax}(h(x)) \quad (k, 1)$$

实际的计算并不会直接用上面的表达式。

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

(y在这里像是某种
参数)



CS 扫描全能王

3亿人都在用的扫描App

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

~~ML~~ → Core:

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

\swarrow loss-func \searrow hypothesis

Any Supervised algorithm is to solve this form problem

i.e. softmax regression:

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



CS 扫描全能王

3亿人都在用的扫描App

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

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

$\nabla_{\theta} f(\theta)$, 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: $\theta := \theta - \frac{\alpha}{B} \sum_i^B \nabla_{\theta} l(h(x^{(i)}), y^{(i)})$



CS 扫描全能王

3亿人都在用的扫描App

manually calculate gradient

i.e. the softmax form:

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

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

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

类似

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

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

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

$$\begin{aligned}\frac{\partial l_{ce}(\theta^T x, y)}{\partial \theta} &= \frac{\partial l_{ce}(\theta^T x, y)}{\partial \theta^T x} \cdot \frac{\partial \theta^T x}{\partial \theta} \\ &= (z - ey) \cdot x \\ &= \underset{k \times 1}{(z - ey)} \cdot \underset{n \times 1}{x} \\ &= \text{embarrassingly } = x(z - ey)^T\end{aligned}$$

for batch form:

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

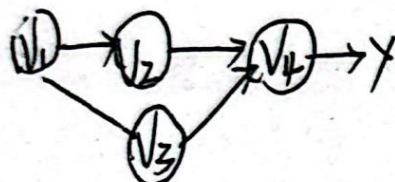


CS 扫描全能王

3亿人都在用的扫描App

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: $\bar{V}_i = \bar{V}_2 \cdot \frac{\partial V_2}{\partial V_i} + \bar{V}_3 \cdot \frac{\partial V_3}{\partial V_i}$

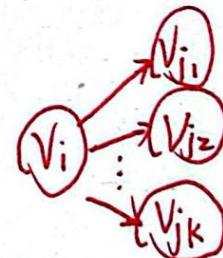
$$\text{defi. } \bar{V}_{i \rightarrow j} = \cancel{\frac{\partial V_j}{\partial V_i}}$$

$$= \bar{V}_j \cdot \frac{\partial V_j}{\partial V_i}$$

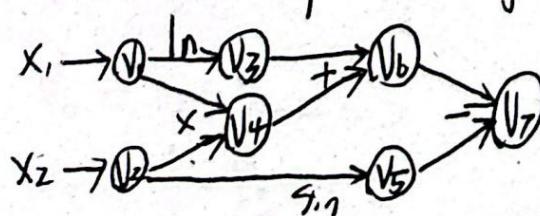
有简洁的形式:

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

这就是对应一个 node 有 multiple output 的情况!



给你画一个 computation graph 示意, $y = \ln(x_1) + x_1 x_2 + g_{\ln x}$



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, 这只是手算.

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: → 它的特点是 input $\in \mathbb{R}^n$, 就要做 n 次

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

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

$$v_4 = v_1 \cdot v_2 + v_1 \cdot v_2$$

$$= 1.5 + 2.0$$

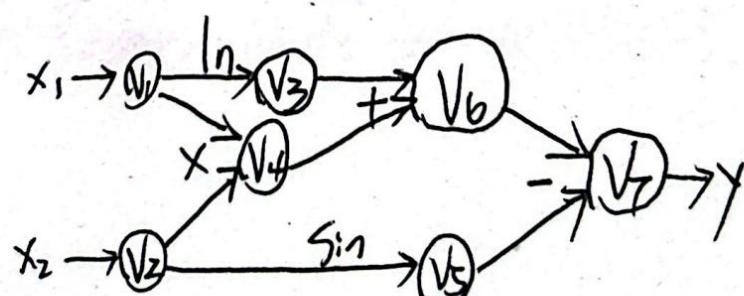
$$= 5$$

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

$$v_6 = v_3 + v_4 = 5.5$$

$$v_7 = v_6 - v_5 = 5.5$$

非常糟糕, 产生了非常棘手依赖!



as long as i know the former
input node



CS 扫描全能王

3亿人都在用的扫描App

e.g. Reserve mode AD

还是那个图

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

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

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

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

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

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

$$\left\{ \begin{array}{l} \frac{\partial y}{\partial V_1} = \frac{\partial y}{\partial V_3} \cdot \frac{\partial V_3}{\partial V_1} + \frac{\partial y}{\partial V_4} \cdot \frac{\partial V_4}{\partial V_1} = 1 \cdot \cancel{1} + 1 \cdot V_2 = 5.5 \\ \end{array} \right.$$

$$\Rightarrow \left(\begin{array}{c} \frac{\partial y}{\partial V_2} \\ \frac{\partial y}{\partial V_1} \end{array} \right) = \nabla f, \text{ 所以 Reverse Mode 做一次就够了。}$$



CS 扫描全能王

3亿人都在用的扫描App

Backprop Vs Reverse mode AD



by extending computational graph

6.

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

最后提了 -T Reverse mode AD on data structures

"tuple value"? 没怎么听懂



CS 扫描全能王

3亿人都在用的扫描App

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

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

for k ∈ input(i):

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

~~opp~~ 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，牵的线很花。



CS 扫描全能王

3亿人都在用的扫描App

Fully connected networks

来: $b \in \mathbb{R}^{n_{\text{it}}}$

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

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 稳一些, 可能收敛快一些



CS 扫描全能王

3亿人都在用的扫描App

Nesterov Momentum:

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

Adam: good optimizer in practice

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

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

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

每次步长都差不多

几条 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: $x \sim N(0, 1)$, $w_i \sim N(0, \frac{1}{n})$

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

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



CS 扫描全能王

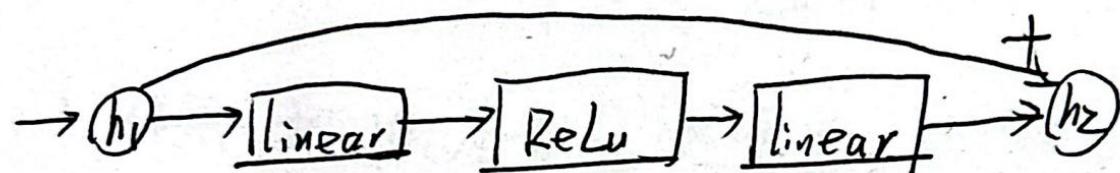
3亿人都在用的扫描App

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

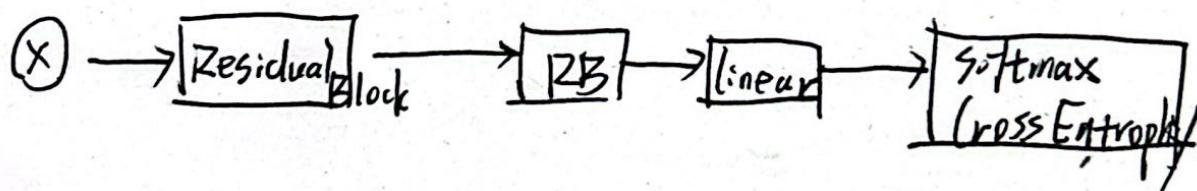
pytorch: imperative
立刻计算

High level modular library components.

0. Residual Block:



1. Multi-layer Residual Net



Deep learning is modular in nature

nn.module 就是把上面 2 个东西连在一起

things to consider:

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



CS 扫描全能王

3亿人都在用的扫描App

当然了 Softmax Cross Entropy 是 Tensor in, scalar out

Optimizer: 会讲 SGD, Momentum 那些

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

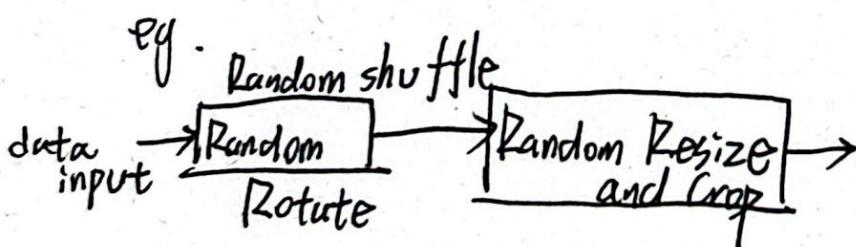
(1) 正则化: loss-func += $\sum |w_i|$

会使很多不重要的 w_i 直接变 0

(2) 正则化: 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



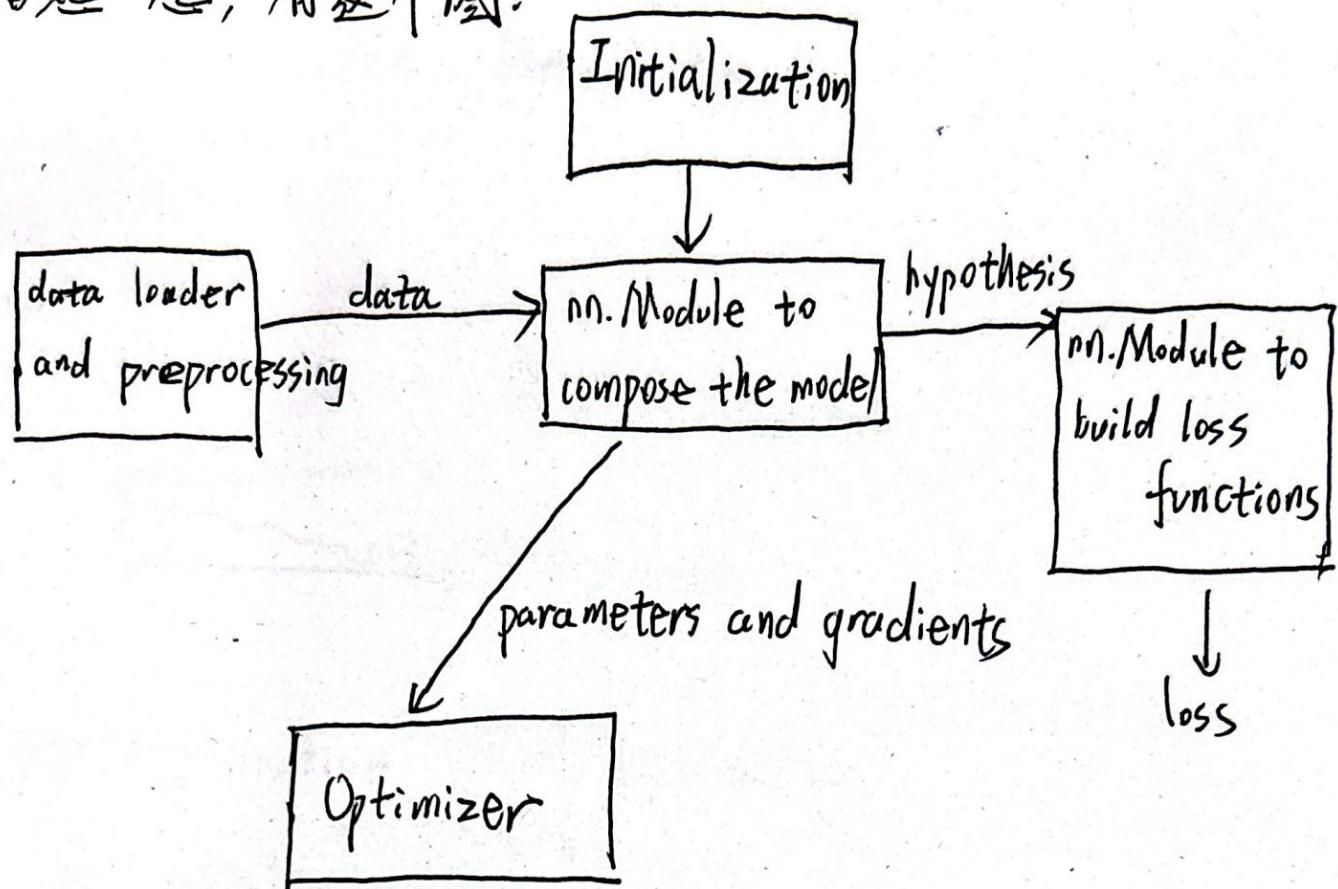
对训练准确也有影响.



CS 扫描全能王

3亿人都在用的扫描App

合起来，有这个图：



CS 扫描全能王

3亿人都在用的扫描App

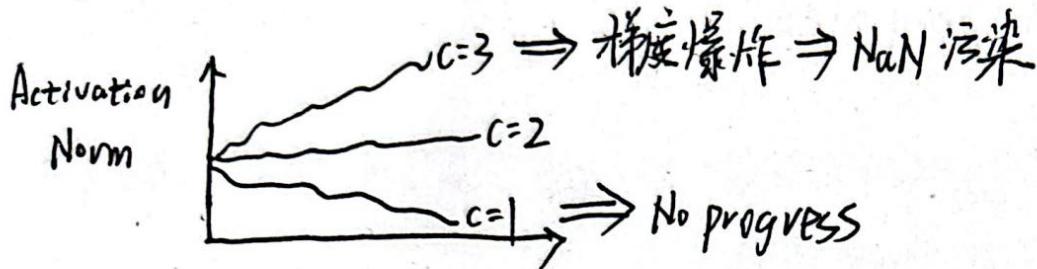
tricks:

Normalization Regularization

一.

观察：

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



Layer normalization: 按行来 normalize.

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

$$z_{i+1} = \frac{\hat{z}_{i+1} - E[\hat{z}_{i+1}]}{\sqrt{\text{var}(\hat{z}_{i+1}) + \epsilon}} \xrightarrow{\text{极小, 避免 div 0}}$$

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，预防过拟合。



CS 扫描全能王

3亿人都在用的扫描App

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。

bz 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：

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

每次都 shrink $(1 - \alpha \lambda)$



CS 扫描全能王

3亿人都在用的扫描App