

第6章 网络优化与正则化

优化问题: 梯度消失 + 参数多

优化问题: 拟合能力过强

本节

机器学习理论

VC 维

Rademacher 复杂度

L_1 和 L_2 正则化在 NN 中的作用有限, 而一些经验做法 提前停止 和 SGD 更有效。

拟合能力的规律。

6.1 网络优化的难点

⇒ 结构多样性

⇒ 高维变量的非凸化

特点

平坦低部

6.2 优化算法

~~1.~~ 1. 小批量梯度下降

2. 随机梯度下降

3. 批量梯度下降

共同问题 \Rightarrow

1. 如何初始化参数

2. 如何处理数据

3. 如何选择合适的学习率

6.2.1 小批量梯度下降

$$\theta_t = \theta_{t-1} - \alpha \frac{\partial R(\theta)}{\partial \theta}$$

- 学习率衰减

- 梯度方向优化

6.2.2 学习率衰减

 \Rightarrow Ada Grad 算法

原因：每个参数的维度与收敛速度都不相同，因此根据每个参数的收敛情况分别设置学习率。

方案：借鉴 L2 正则化思想，每次迭代时自适应地调整每个参数的学习率。

在第 t 次迭代时，先计算每个梯度参数平方的累计值。

$$G_t = \sum_i g_i \odot g_i$$

g_i 是第 i 次迭代时的梯度

$$\Delta \theta_t = - \frac{\alpha}{\sqrt{G_t + \epsilon}} \odot \vec{g}_t$$

α 初始学习率

ϵ 保持数值稳定性 $e^{-7} \sim e^{-10}$

优点: 某些参数的梯度累积较大, 其学习率相对较小.

--- 较小, --- 较大.

整体上随着迭代次数增加, 学习率逐渐减小.

缺点: 经过一定次数的迭代仍然没有找到最优解, 学习率已经很小, 无法找了.

⇒ RMSprop 算法 (改进 Ada Grad)

G_t 每次迭代梯度 \vec{g}_t 的平方指数衰减移动平均.

$$G_t = \beta G_{t-1} + (1-\beta) \vec{g}_t \odot \vec{g}_t$$

β decay rate

$$\Delta \theta_t = - \frac{\alpha}{\sqrt{G_t + \epsilon}} \odot \vec{g}_t$$

与 Ada Grad 的区别

G_t 由累积变为指数衰减移动平均.

优点: 一定程度上缓解了学习率过早衰减的缺点.

⇒ AdaDelta 算法 (改进 Ada Grad)

- 梯度的平方指数衰减 移动平均

- 每次参数更新差的 $\Delta\theta$ 的平方指数移动平均

$$\Delta X_{t+1}^2 = \beta_1 \Delta X_t^2 + (1-\beta_1) \Delta\theta_t \odot \Delta\theta_t$$

β_1 衰减率

$\Delta\theta_t$ 求和, 故只计算到 ΔX_{t+1}

更新公式为

$$\Delta\theta_t = - \frac{\sqrt{\Delta X_{t+1}^2}}{\sqrt{G_t + \epsilon}} \cdot \vec{g}_t$$

优点: 一定程度上平抑了学习率的波动

6.2.3 梯度方向优化

⇒ 动量法

方案: 用之前求出的动量来替代真正的梯度.

每次迭代的梯度看作加速度

$\Delta\theta_t$ 负梯度的移动加权平均.

$$\Delta\theta_t = \rho \Delta\theta_{t-1} - \alpha \vec{g}_t$$

ρ 动量因子, 一般取 0.9

α 学习率

⇒ Nesterov 加速梯度 (改进动量法)

动量法中两个部分

$$\left\{ \begin{array}{l} \Delta \theta_{t-1} \\ \bar{g}_t \end{array} \right.$$

先根据 $\Delta \theta_{t-1}$ 更新得到参数 $\hat{\theta}$

$$\hat{\theta} = \theta_{t-1} + \rho \Delta \theta_{t-1}$$

再用 \bar{g}_t 更新

$$\theta_t = \hat{\theta} - \alpha \bar{g}_t$$

更新方向改为 $\hat{\theta}$ 上的梯度

$$\Delta \theta_t = \rho \theta_{t-1} - \alpha \bar{g}_t (\theta_{t-1} + \rho \Delta \theta_{t-1})$$

⇒ Adam 算法

~~G_t 与 RMS prop 类似, 计算 \bar{g}_t 的平方指数加权平均.~~

~~M_t 与动量法类似, 计算 \bar{g}_t 的指数加权平均.~~

~~$$G_t = \beta_1 G_{t-1} + (1 - \beta_1) \bar{g}_t$$~~

M_t 与动量法类似, 计算 \bar{g}_t 的指数加权平均

G_t 与 RMS prop 类似, 计算 \bar{g}_t 的平方指数加权平均.

$$M_t = \beta_1 M_{t-1} + (1 - \beta_1) \bar{g}_t$$

$$\beta_1 = 0.9$$

$$G_t = \beta_2 G_{t-1} + (1 - \beta_2) \bar{g}_t \odot \bar{g}_t$$

$$\beta_2 = 0.99$$

No.

Date.

当 $M_0 = 0$, $G_0 = 0$ 时, $\beta_1, \beta_2 \rightarrow 1$ 时偏差较大.

$$\hat{M}_t = \frac{M_t}{1 - \beta_1}$$

$$\hat{G}_t = \frac{G_t}{1 - \beta_2}$$

参数更新差值为

$$\Delta \theta_t = - \frac{2}{\sqrt{G_t + \epsilon}} \hat{M}_t$$

⇒ 梯度截度

问题: 梯度爆炸

⇒⇒ 按值截断

大于 a 时, 设为 a ;

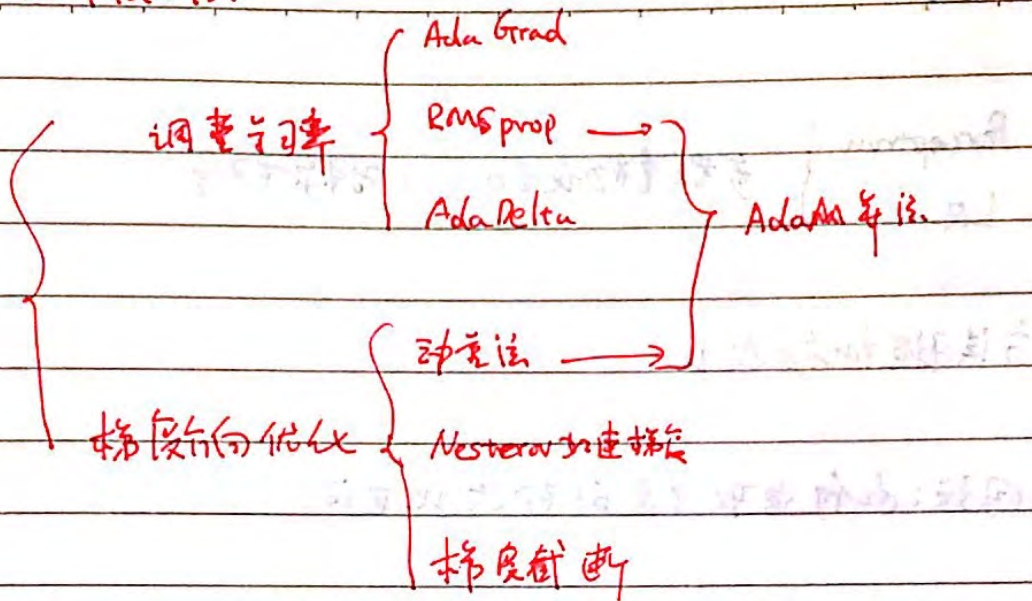
小于 b 时, 设为 b ;

$$\bar{g}_t = \max(\min(\tilde{g}_t, b), a)$$

⇒⇒ 按模截断

$$\bar{g}_t = \frac{b}{\|g_t\|} \tilde{g}_t, \text{ 当 } \|g_t\|^2 > b \text{ 时.}$$

6.2.4 优化算法小结.



No.

Date.

6.3 参数初始化

Perceptron } 参数初始化为0 \Rightarrow 对称权重
LR }

方法: 随机初始化

问题: 如何选取合适的初始化区间.

6.3.1 Gaussian 初始化

固定均值 + 固定方差

输入连接为 n_{in} 时, $N(0, \sqrt{\frac{1}{n_{in}}})$

6.3.2 均匀分布初始化

给定一个区间 $[-r, r]$, 门内采用均匀分布来初始化参数

r 超参数, 按神经网络层数量自适应调整

\Rightarrow Xavier 初始化方法

$$\text{logistic} \quad r = \sqrt{\frac{6}{n^u + n^o}}$$

6.4 数据预处理

问题: 特征值范围差异

⇒ 缩放归一化

将每个特征归一到 $[0, 1]$ 或 $[-1, 1]$ 之间.

$$x^{(1)} = \frac{x - \min}{\max - \min}$$

⇒ 标准归一化

处理成符合标准正态分布, 每个特征.

先计算 μ 和 σ^2

$$\mu = \frac{1}{n} \sum x^{(i)}$$

$$\sigma^2 = \frac{1}{n} \sum (x^{(i)} - \mu)^2$$

新的特征值.

$$\hat{x}^{(i)} = \frac{x^{(i)} - \mu}{\sigma}$$

⇒ 白化

白化后, 每个特征之间相关性较低, 并且所有特征具有相同的方差.

方法: PCA

6.5 逐层归一化

内部协变量偏移

6.5.1 批量归一化 (BN) 对中间层的单例神经元

$$\tilde{a}^{(l)} = f(\tilde{z}^{(l)}) = f(w\tilde{a}^{(l)} + b)$$

逐层归一化在中间层上操作，效率更高 改用标准归一化，归一化净输入 $\tilde{z}^{(l)}$ $\tilde{a}^{(l)}$ 分布性质不如 $\tilde{z}^{(l)}$ 稳定

$$\tilde{z}^{(l)} = \frac{\tilde{z}^{(l)} - E(\tilde{z}^{(l)})}{\sqrt{\text{Var}(\tilde{z}^{(l)}) + \epsilon}}$$

如果用 sigmoid 时，自个取值区间网也是接近线性变换的区间，减弱了网络的非线性性质。

方法：附加的缩放和平移变换改变取值区间。

$$\tilde{z}^{(l)} = \frac{\tilde{z}^{(l)} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \circ \gamma + \beta$$

$$= \text{BN}_{\gamma, \beta}(\tilde{z}^{(l)})$$

批量归一化操作可以看作一个特殊的神经网络，加在每一层非线性激活函数之前。

$$\tilde{a}^{(l)} = f(\text{BN}_{\gamma, \beta}(\tilde{z}^{(l)})) = f(\text{BN}_{\gamma, \beta}(w\tilde{a}^{(l-1)}))$$

6.5.2 层归一化

~~但是~~ BN 要求批量不能太小，无法应用于循环神经网络。

层归一化是对中间层的所有神经元进行归一化

$$\mu^{(l)} = \frac{1}{n} \sum z_i^{(l)}$$

$$\sigma^{(l)} = \frac{1}{n} \sum (z_i^{(l)} - \mu^{(l)})^2$$

层归一化定义为

$$\hat{z}^{(l)} = \frac{z^{(l)} - \mu^{(l)}}{\sqrt{\sigma^{(l)} + \epsilon}} \odot \gamma + \beta$$

$$\hat{z}^{(l)} = LN_{\gamma, \beta}(z^{(l)})$$

⇒ 循环神经网络中的层归一化

$$\vec{z}_t = U \vec{h}_{t-1} + W \vec{x}_t$$

$$\vec{h}_t = f(LN_{\gamma, \beta}(\vec{z}_t))$$

缓解梯度爆炸或消失

⇒ 对比 BN 和 LN

对 k 个样本的小批量集合 $z^{(l)} = \begin{bmatrix} z^{(l,1)} \\ \vdots \\ z^{(l,k)} \end{bmatrix}$ ，层归一化是对矩阵每列进行归一化，而 BN 是对每行进行归一化。

$$\begin{bmatrix} z^{(l,1)} \\ \vdots \\ z^{(l,k)} \end{bmatrix}$$

一般而言，BN 是更好的选择。

Batch 较小时，可以选择 LN。

6.5.3 权重归一化 (WN)

对神经网络的连接权重进行归一化, 再通过参数化方法, 将连接权重分解为长度和方向两种参数

$$\vec{a}^{(l)} = f(W\vec{a}^{(l-1)} + \vec{b})$$

$$w_{iz} = \frac{g_i}{\|\vec{v}_i\|} \vec{v}_i$$

w_{iz} : 表示权重 W 的第 i 行

g_i 为标量

\vec{v}_i 与 $\vec{a}^{(l-1)}$ 维数相同

优点: 归一化可错较小

6.5.4 局部响应归一化 (LRN)

通常用于 CNN

$$\hat{Y}^p = Y^p / (k + \alpha \sum (Y^j)^2)^\beta$$

$$\equiv \text{LRN}_{n, \alpha, \beta}(Y^p)$$

低抑制

6.6 超参数优化

常见超参数有

{ 网络结构
 优化参数
 正则化系数

6.6.1 网格搜索

尝试所有的超参数组合。

6.6.2 随机搜索

随机组合 \Rightarrow 选最好的。 -

一般比网格搜索更加有效。

6.6.3 贝叶斯优化

时序模型优化 SMBO

$$E[Z(x, H)] = \int_{-\infty}^{\infty} \max(y^* - y, 0) P(y | \bar{x}, H) dy$$

No.

Date.

6.6.4 动态资源分配

最优解问题

逐次减半

多臂赌博机

神经网络搜索

6.7 网络正则化

6.7.1 L_1 和 L_2 正则化正则化 \Rightarrow 限制模型复杂度问题: 过度参数下, L_1 和 L_2 正则化效果一般?

方案: 数据增强, 提前停止, 丢弃法, 早停法

$$\theta^* = \arg \min_{\theta} \frac{1}{n} \sum L(y^w, f(\bar{x}^w, \theta)) + \lambda \phi(\theta)$$

等价于

$$\theta^* = \arg \min_{\theta} \frac{1}{n} \sum L(y^w, f(\bar{x}, \theta))$$

subject to $\phi(\theta) < 1$ L_1 范数在零点不可导, 常用下式来近似

$$L_1(\theta) = \sum \sqrt{\theta_i^2 + \epsilon}$$

弹性网络正则化

同时加 λ_1 L_1 和 L_2 正则化

$$\theta^* = \arg \min_{\theta} L(y^w, f(\bar{x}^w, \theta)) + \lambda_1 L_1(\theta) + \lambda_2 L_2(\theta)$$

6.7.2 权重衰减

$$\theta_t = (1-\omega) \theta_{t-1} - \alpha \bar{g}_t$$

在标准的SGD中, 权重衰减与L2正则化效果相同.

在Adam中并不等价.

6.7.3 提前停止

验证集,

用验证集上的误差来代替期望误差.

6.7.4 丢弃法

d(x) 丢弃函数

$$\tilde{y} = f(w d(x) + b)$$

$$d(x) = \begin{cases} \bar{m} \odot \bar{x} & \text{当训练阶段} \\ p \bar{x} & \text{当测试阶段} \end{cases}$$

$\bar{m} \in \{0, 1\}^d$ 丢弃掩码

在测试时, 将每个神经元的输出乘以 p , 相当于不同的网络做了平均.

p 一般等于0.5

⇒ 集成学习的简单解释

每次Dropout, 相当于从原网络中采样得到一个子网络.

n 个神经元 $\Rightarrow 2^n$ 个子网络.

集成了指数级个不同的旧全模型.

⇒ 贝叶斯学习的解释

假设参数 θ 为随机向量, 且先验分布为 $g(\theta)$, 模型为 $f(\bar{x}, \theta)$

贝叶斯方法的预测为

$$E_{g(\theta)}[y] = \int_{\theta} f(\bar{x}, \theta) g(\theta) d\theta$$

$$\approx \frac{1}{m} \sum_m f(\bar{x}, \theta_m)$$

$f(\bar{x}, \theta_m)$ 第 m 次应用 ~~Drop~~ Dropout 后的网络.

⇒ 循环神经网络上的丢弃法.

问题:

不能对每个时刻的隐状态进行随机丢弃, 这样会损害神经网络在时间维度上的记忆能力.

变分丢弃法:

对参数矩阵的每个元素进行随机丢弃, 并且在所有时刻都使用相同的掩码.

6.7.5 数据增强 (图像)

旋转 (Rotation)

翻转 (Flip)

缩放 (Zoom In/out)

平移 (Shift)

加噪声 (Noise)

6.7.6 标签平滑

在输出标签中添加噪声来避免模型过拟合。

硬目标: one-hot 向量

软目标: 假设样本以 ϵ 的概率为某类

$$\tilde{y} = [\frac{\epsilon}{k-1}, \dots, 1-\epsilon, \frac{\epsilon}{k-1}]^T$$