

7 优化算法 8 采样

2020年5月31日 星期日 上午9:11

机器学习 = 模型表征 + 模型评估 + 优化算法

1. 有监督的损失函数

1.1 二分类问题 $Y = \{-1, 1\}$.

指示函数: $L_{0-1}(f, y) = 1 \text{ if } f_y \leq 0$

Hinge: $L_{\text{hinge}}(f, y) = \max\{0, 1 - fy\}$

logistic (sigmoid): $L_{\text{logistic}}(f, y) = \log_2(1 + \exp(-fy))$

Cross Entropy: $L_{\text{CE}}(f, y) = -\log_2\left(\frac{1 + fy}{2}\right)$

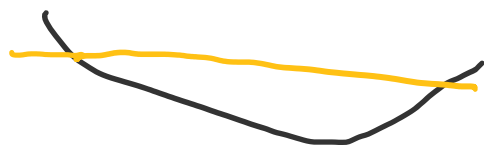
1.2 回归问题 平方损失函数

$L_{\text{square}}(f, y) = (f - y)^2$

2 优化问题

Q: 哪些是凸优化问题, 哪些不是 (机器学习中)

A: 凸函数 $L(\lambda x + (1 - \lambda)y) \leq \lambda L(x) + (1 - \lambda)L(y)$



逻辑回归是凸优化; 支持向量为不是凸优化

3. 经典优化算法

Q: 无约束问题的优化方法有哪些?

A: 梯度法 $\min_{\theta} L(\theta)$

经典优化算法: 直接法和迭代法

迭代法: 一阶泰勒展开 一梯度下降法

二阶泰勒展开 牛顿法

4. 梯度验证

5. 随机梯度下降法

$$L(\theta) = E_{(x, y) \sim p_{\text{data}}} L(f(x, \theta), y)$$

$$\rightarrow \theta = \arg\min_{\theta} L(\theta)$$

$$\text{梯度下降: } L(\theta) = \frac{1}{n} \sum_{i=1}^n L(f(x_i, \theta), y_i) \quad \theta_{t+1} = \theta_t - \alpha \cdot \nabla L(\theta)$$

$$\text{随机梯度下降: } L(\theta) = \frac{1}{m} \sum_{j=1}^m L(f(x_j, \theta), y_j) \quad \nabla L(\theta) = \frac{1}{m} \sum_{j=1}^m \nabla L(f(x_j, \theta), y_j)$$

取 m 一般 2 的幂次, 先对所有数据进行随机排序

6. 随机梯度算法的弊端

Q: SGD 失效原因

A: SGD 相对于大批量样本梯度的量少, 而且波动剧烈甚至不收敛, 局部最优

解决办法: 引入惯性, 动量 (Momentum)

$$\text{SGD: } \theta_{t+1} = \theta_t - \eta g_t \Rightarrow \begin{cases} v_t = \gamma v_{t-1} + \eta g_t \\ \theta_{t+1} = \theta_t - v_t \end{cases} \begin{array}{l} \text{速度更快} \\ \text{更收敛} \end{array}$$

$$\text{— AdaGrad: } \theta_{t+1, i} = \theta_{t, i} - \frac{\eta}{\sqrt{\sum_{s=0}^t g_{s, i}^2 + \epsilon}} g_{t, i} \quad (\text{环境感知})$$

— Adam: 集合 adagrad 和 momentum 的优点.

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta \cdot m_t}{\sqrt{v_t + \epsilon}}$$

β_1, β_2 衰减系数

m_t 一阶矩估计 $E(g_t)$

v_t 二阶矩估计 $E(g_t^2)$

7. 正则化和稀疏性

Q: L_1 正则化使模型具有稀疏性的原因

A:

8. 采样: 从特定的概率分布中抽取样本点的过程

均匀分布随机数