

梯度下降法优化

$$w_{k+1} = w_k - \eta J'(w_k)$$

1. 减小计算量

e.g. 随机梯度下降

2. 优化下降路径

e.g. 牛顿法 (计算量也大)

拟合:

$$J(w) \approx J(w_k) + J'(w_k)(w - w_k) + \frac{1}{2} J''(w_k)(w - w_k)^2 + \dots$$

(泰勒展开)

有

$$J'(w_k) \approx \frac{J(w) - J(w_k)}{w - w_k}$$

$$\Rightarrow J(w) \approx J(w_k) + J'(w_k)(w - w_k)$$

$$\Rightarrow J(w) = J(w_k) + J'(w_k)(w - w_k) + \frac{1}{2} J''(w_k)(w - w_k)^2$$

↓

求得 $J'(w) = J'(w_k) + J''(w_k)(w - w_k)$

使 = 0

推导出

$$w_{k+1} = w_k - \frac{f'(w_k)}{f''(w_k)}$$

补充:

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}$$

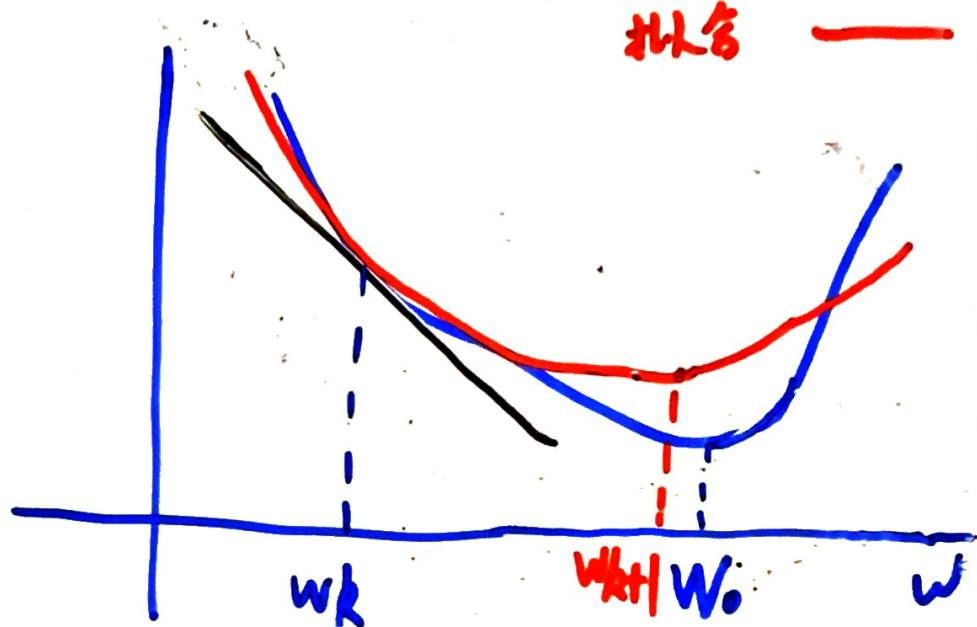
海森矩阵

(函数曲率的信息)

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1 \partial x_1} & \dots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f(x)}{\partial x_n \partial x_n} \end{bmatrix}$$

$J(w)$ 损失函数

$J'(w)$ 拟合
最优
拟合



拟合优于 $J'(w)$, 处于 $J(w)$

与最优两者之间, 平太斜了

2. 优化路径
动量法
(历史)

$$W_{(t+1)i} = W_{(t)i} - \eta \frac{\partial J(W_{(t)i})}{\partial W_i}$$

$$\text{令 } \nabla W_{(t+1)i} = \frac{\partial J(W_{(t)i})}{\partial W_i}$$

$$\text{再令 } V_{(t+1)} = \nabla W_{(t+1)i} + V_{(t)}$$

$$\text{有 } W_{(t+1)i} = W_{(t)i} - \eta V_{(t+1)}$$

使历史越久远的数据指数级减弱

影响：指数加权移动平均法

$$V_{(t+1)} = \underset{0.1}{(1-\beta)} \nabla W_{(t+1)i} + \underset{0.9}{\beta} V_{(t)}$$

Nesterov 算法

(历史 + 未来)

梯度等高线

$$\nabla W_{(t+1)i} = \frac{\partial J(W_{(t)i} + \gamma V_{(t)})}{\partial W_i}$$

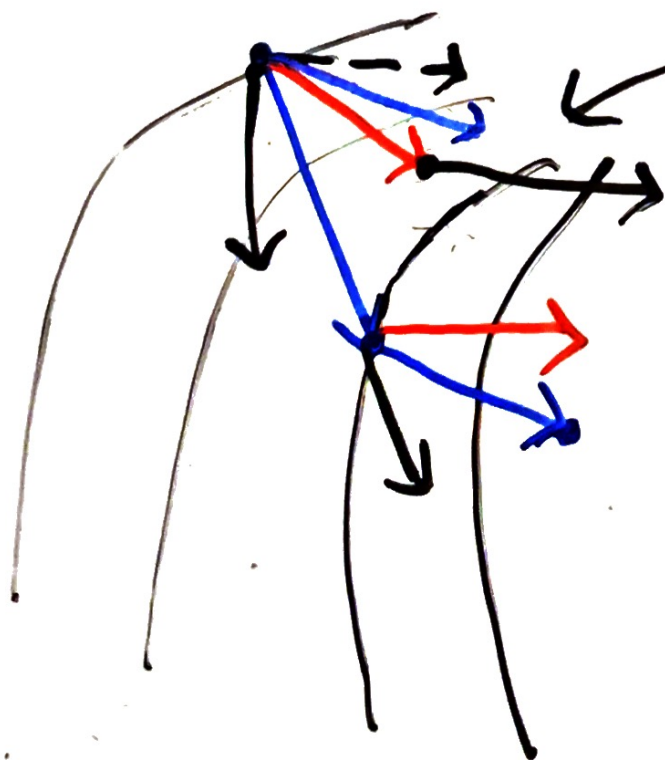
往前走一步

再求“未来点”
的梯度

历史：→

点梯度向量：→

综合：→



2. 优化路径

之前的牛顿法、动量(冲量)法、Nesterov

为优化梯度

现在开始优化 **学习率**

c_i 为某一维度

↳ Ada Grad (自适应)

$$W_{(t+1)i} = W_{(t)i} - \frac{\eta}{\sqrt{S_{(t+1)}} + \epsilon} \cdot \nabla W_{(t+1)i}$$

→ 标小量, 避免分为零

$$S_{(t+1)} = S_{(t)} + \nabla W_{(t+1)i} \cdot \nabla W_{(t+1)i}$$

(依赖历史数据)

梯度越多, 修正越大

补充:

任务的实现主要依赖不同特征维度之间的区别
其数据集称为一个稀疏数据集, 如有毛没毛
1 0

若更注重某一特征的程度不一样, 即不是, 如毛的多少

所以处理稀疏数据多震荡, 可以使用 Ada Grad

锚点问题: "平台区" 移动缓慢, 过了"平台区" 开始时依然缓慢

↳ RMSprop

依然使用指数加权移动平均法

$$S_{(t+1)} = (1-\beta) \nabla W_{(t+1)i}^2 + \beta S_{(t)}$$

RMSprop + 动量法 \Rightarrow Adam

RMSprop + Nesterov \Rightarrow Nadam