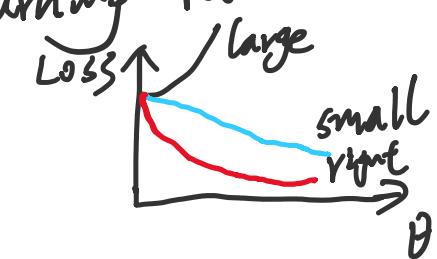


3 Gradient Descent

2020年5月20日 星期三 下午2:42

梯度下降: $\nabla L(\theta) = \begin{bmatrix} \partial L(\theta) / \partial \theta_1 \\ \partial L(\theta) / \partial \theta_2 \end{bmatrix}$

Tips: 1. Learning rate $\theta^i = \theta^{i-1} - \eta \nabla L(\theta^{i-1})$

Visualize Loss: 

Q: 可以自动调整学习率吗?
1. 刚开始 η 比较大
2. 之后 η 逐渐减小

技巧: 自动梯度 adgrad

$$\begin{cases} w^{t+1} \leftarrow w^t - \eta^t g^t \\ w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t \end{cases} \quad \text{adgrad} - \sigma^t: \text{root mean square of } w \text{ 的偏导数的值}$$

example:

$$w^1 \leftarrow w^0 - \frac{\eta^0}{\sigma^0} g^0$$

$$\sigma^0 = \sqrt{(g^0)^2}$$

$$w^2 \leftarrow w^1 - \frac{\eta^1}{\sigma^1} g^1$$

$$\sigma^1 = \sqrt{\frac{1}{2}[(g^0)^2 + (g^1)^2]}$$

\vdots

\vdots

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t$$

$$\sigma^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^i)^2}$$

$$\eta^t = \frac{\eta}{\sqrt{t+1}} \quad \text{decay}$$

best step:

first derivative

second derivative

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$

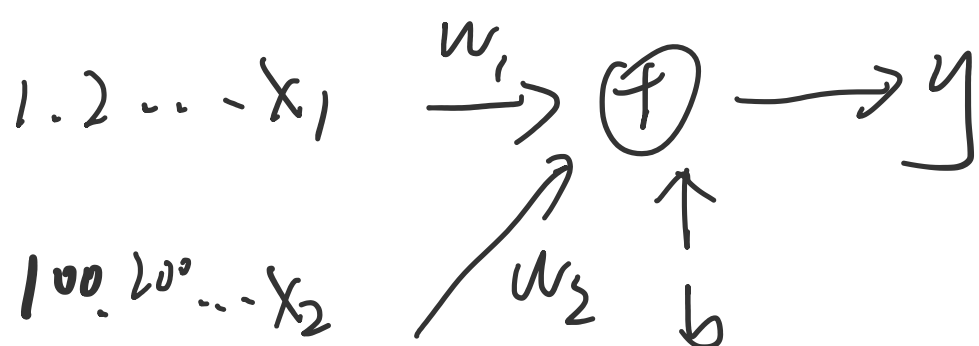
— 强调反差效果

2. Stochastic Gradient Descent 随机梯度下降

pick an example: x^n

$$L^n = (y^n - (b + \sum w_i x_i^n))^2 \quad \theta^i = \theta^{i-1} - \eta \nabla L^n(\theta^{i-1})$$

3. Feature Scaling 特征编码



将数据转换到相同维度

$$x_i^r = \frac{x_i - m_i}{\sigma_i}$$

m_i : 均值

σ_i : 标准差

$$x_i^r = \frac{x_i - \min}{\max - \min}$$

最大最小化

问题:

1. Stuck local minima $\frac{\partial L}{\partial w} = 0$
2. Stuck saddle point
3. Very slow at plateau $\frac{\partial L}{\partial w} \approx 0$

