

Optimization Basics

Find \bar{w} to minimize loss; search over space of params

$$\text{loss } \mathcal{L}(\underbrace{(\bar{x}^{(i)}, y^{(i)})}_{\text{training data}}^D, \bar{w}) = \sum_{i=1}^D \mathcal{L}(\bar{x}^{(i)}, y^{(i)}, \bar{w})$$

function of \bar{w}

Stochastic gradient descent: repeatedly pick example i

$$\bar{w} \leftarrow \bar{w} - \underset{\substack{\uparrow \\ \text{step size}}}{\alpha} \frac{\partial}{\partial \bar{w}} \mathcal{L}(i, \bar{w}) \leftarrow \text{loss on } i\text{th example}$$

Step size

Suppose

$$\mathcal{L}(i, \bar{w}) = w^2 \quad \text{one feature}$$

$$\bar{w} = [w]$$

$$w = -1$$

$$\text{if } \alpha = 1: w = -1$$

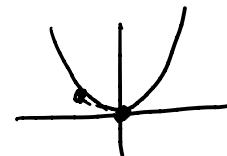
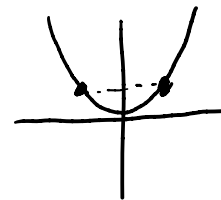
$$-1 - 2(-1)$$

$$\Rightarrow w = 1$$

$$\frac{\partial}{\partial w} \mathcal{L} = 2w$$

Keep $\alpha = 1$: oscillates!

$$\text{if } \alpha = 1/2: w \rightarrow 0$$



Choosing Step Size

How to choose step size?

- Try them out: $1e^0$ $1e^{-1}$ $1e^{-2}$...
- Large \rightarrow small, e.g. $1/t$ for epoch t ($1/\sqrt{t}$...)
(fixed schedule)

or decrease step size when performance stagnates on held-out data

Newton's method: $\bar{w} \leftarrow \bar{w} - \underbrace{\left(\frac{\partial^2}{\partial \bar{w}^2} \mathcal{L} \right)^{-1}}_{\text{inverse Hessian}} \frac{\partial}{\partial \bar{w}} \mathcal{L}$

DL
Adagrad, Adadelta, Adam: "adaptive" methods
approximations to the inverse Hessian
(linear in # feats) very expensive to compute

Regularization: don't really use