

CS 182 Lecture 4: Optimization

$I(\theta)$

6/9

$$\theta^* \leftarrow \arg \min_{\theta} \left(- \sum_i \log p_{\theta}(y_i | x_i) \right) \quad v_i = - \frac{dI(\theta)}{d\theta_i} \quad (\text{gradient vector})$$

Gradient descent - picks steepest direction, not "best" direction

"Nice" loss surface - all roads lead to Rome.

Not for log-reg. - guaranteed to be convex! Strengthens gradient descent.

Three features of loss landscape $\begin{cases} \rightarrow \text{local optimum} \\ \rightarrow \text{the plateau} \\ \rightarrow \text{saddle} \end{cases}$ all are harmful for GD.

Local optima becomes less of an issue as # of NN params increase...

We need learning rates large enough to get past plateaus.

Hessian: $\nabla^2 \rightarrow$ gradients"

Most critical points in neural net loss landscapes are saddle points.
 \rightarrow unlikely to have same sign for all diagonal entries in Hessian matrix

$\nabla_{\theta} L(\theta_i) \rightarrow$ gradient Newton's Method: able to find better direction for gradient descent

$$\nabla_{\theta}^2 L(\theta_i) \rightarrow \text{Hessian} \quad \theta^* \leftarrow \theta_0 - (\nabla_{\theta}^2 L(\theta_0))^{-1} \nabla_{\theta} L(\theta_0) \quad \text{Hessian inverse}$$

Not a viable way, computationally expensive.

Momentum: Try to tackle oscillation / "slowness" in GD

$$\theta_{k+1} = \theta_k - \alpha g_k, \quad g_k = \nabla_{\theta} L(\theta_k) + \mu g_{k-1} \quad \begin{cases} \text{"blend in"} \\ \text{previous direction} \end{cases}$$

RMSProp: "root-mean-squared" normalize magnitude of gradient along each dimension

$$s_k \leftarrow \beta s_{k-1} + (1-\beta) (\nabla_{\theta} L(\theta_k))^2 \quad \text{"forgetting" property. Better for DL}$$

$$\theta_{k+1} = \theta_k - \frac{\nabla_{\theta} L(\theta_k)}{\sqrt{s_k}} \quad \text{This is rough est. of magnitude. non-convex probs}$$

AdaGrad: estimate per-dimension cumulative magnitude.

$$s_k \leftarrow s_{k-1} + (\nabla_{\theta} L(\theta_k))^2 \quad \text{Good for convex problems.}$$

helps "blow up" values for small k , speed it up

Adams combines momentum and RMSProp.

$$m_k = (1-\beta_1) \nabla_{\theta} L(\theta_k) + \beta_1 m_{k-1} \quad \text{"first moment estimate" (momentum-like)}$$

$$v_k = (1-\beta_2) (\nabla_{\theta} L(\theta_k))^2 + \beta_2 v_{k-1} \quad \text{"second moment estimate"}$$

$$\theta_{k+1} = \theta_k - \frac{m_k}{\sqrt{v_k} + \epsilon} \quad \text{Typically: } \epsilon = 10^{-8}, \alpha = 0.001, \beta_1 = 0.9, \beta_2 = 0.999$$

Stochastic GD:

1. Sample $B \subset \mathcal{D}$:

$$2. \text{Estimate } g_k \leftarrow \nabla_{\theta} \frac{1}{B} \sum_{i=1}^B \log p(y_i | x_i; \theta) \approx \nabla_{\theta} L(\theta)$$

$$3. \theta_{k+1} \leftarrow \theta_k - \alpha g_k \quad (\text{or momentum, Adam, etc.})$$

Each iteration samples diff. B .

Epoch: complete

"look-then" dataset.

Preferably decay learning rate over time w/ SGD momentum.

on aside: Nesterovs accelerated gradient

Tuning SGD:

B (batch size) larger, better

α (learning rate) decay over time

μ (momentum)

0.9, 0.9 are good

β_1, β_2 (keep defaults)

can turn on validation loss.

start with one shuffle for all batches, saves RAM