

Lecture #5: Optimization

$$\left. \begin{array}{l} \text{GD: } \theta^{n+1} = \theta^n - \eta \nabla_{\theta} \mathcal{L}(\theta^n) \\ \text{Newton: } \theta^{n+1} = \theta^n - \eta H_n^{-1} \nabla_{\theta} \mathcal{L}(\theta^n) \end{array} \right\}$$

Remarks:

1.) Choosing η is often an "art" (GD).

2.) $\theta \in \mathbb{R}^d$, it may not be smart to use the same learning rate for all of them.

e.g. fitting a univariate Gaussian, $\theta := \{\mu, \sigma^2\}$

3.) Exact Hessians are often very expensive to compute/store/invert. \rightarrow Quasi-Newton

\rightarrow 4.) Scalability to big data

Stochastic gradient descent:

In many ML applications the loss function

factorize across data-points:

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \underbrace{\mathcal{L}_i(\theta)} \quad (\text{see eg. linear regression, NNs})$$

$$\text{GD: } \theta^{n+1} = \theta^n - \eta \nabla_{\theta} \mathcal{L}(\theta^n)$$

$$= \theta^n - \eta \underbrace{\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L_i(\theta^n)}_{\text{exact gradient}}$$

SGD : $\theta^{n+1} = \theta^n - \eta \nabla_{\theta} L_i(\theta^n) \leftarrow$ Robins-Monroe

Mini-batch SGD : $\theta^{n+1} = \theta^n - \eta \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L_i(\theta^n), \quad \underline{m \ll n}$

Epoch : # iterations needed to perform
 $\frac{n}{m} \rightarrow$ on full sweep through our data-set (SGD)

Gradient descent variants :

GD with momentum :
$$\begin{aligned} u^{n+1} &= \gamma u^n + \eta \nabla_{\theta} L(\theta^n) \\ \theta^{n+1} &= \theta^n - u^{n+1} \end{aligned} \left. \vphantom{\begin{aligned} u^{n+1} &= \gamma u^n + \eta \nabla_{\theta} L(\theta^n) \\ \theta^{n+1} &= \theta^n - u^{n+1} \end{aligned}} \right\} \begin{array}{l} \eta : 10^{-3} \\ \gamma = 0 \rightarrow \text{GD} \\ \gamma = 0.9 \end{array}$$

Nesterov Accelerated Gradient :
$$\begin{aligned} u^{n+1} &= \gamma u^n + \eta \nabla_{\theta} L(\theta^n - \gamma u^n) \\ \theta^{n+1} &= \theta^n - u^{n+1} \end{aligned}$$

Adaptive Learning rate methods :

RMS Prop : $\underbrace{\mathbb{E}[g^2]}_n :=$ average of the square gradi
variance of the gradients at iteration n .

$$g_n := \nabla_{\theta} L(\theta^n)$$

$$\rightarrow \mathbb{E}[g^2]_{n+1} = \gamma \mathbb{E}[g^2]_n + (1-\gamma) g_n^2, \quad \gamma \sim 0.9$$

 dx

$$\theta^{n+1} = \theta^n - \frac{\eta}{\sqrt{\mathbb{E}[g^2]_{n+1} + \epsilon}} \cdot g_n \quad \left(\begin{array}{l} \text{standard SGD with} \\ \text{adaptive learning} \\ \text{rate} \end{array} \right)$$

Adam: (adaptive moment estimation)

$$m^{n+1} = b_1 m^n + (1 - b_1) g^n \quad : \text{Estimate of the 1st-moment of the gradient}$$

$$v^{n+1} = b_2 v^n + (1 - b_2) g_n^2 \quad : \text{Estimate of the 2nd-moment}$$

$$\hat{m}^n = \frac{m^n}{1 - b_1^n}, \quad \hat{v}^n = \frac{v^n}{1 - b_2^n}$$

$$\theta^{n+1} = \theta^n - \frac{\eta}{\sqrt{\hat{v}^{n+1} + \epsilon}} \hat{m}^{n+1}$$

Parameters: θ

State of optimizer: e.g. $\{\theta^n, \tilde{u}^n\}$ momentum, $\{\theta^n, \hat{m}^n, \hat{v}^n\}$ Adam