

CMPT 409/981: Optimization for Machine Learning

Lecture 7

Sharan Vaswani

September 26, 2024

Recap

Polyak Momentum: Compute the gradient at w_k and then extrapolate:

$$v_k = w_k + \beta_k(w_k - w_{k-1}); w_{k+1} = v_k - \eta \nabla f(w_k).$$

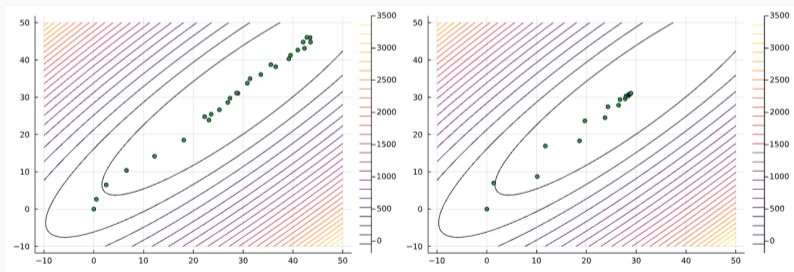
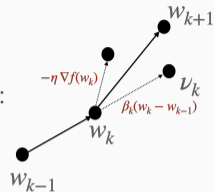


Figure 1: Comparing GD vs HB momentum (with theoretical (η, β)) on a strongly-convex quadratic

Minimizing strongly-convex quadratics with HB momentum

Update: $w_{k+1} = w_k - \eta \nabla f(w_k) + \beta(w_k - w_{k-1})$

Claim: For L -smooth, μ -strongly convex quadratics s.t. $f(w) = \frac{1}{2}w^\top A w - bw + c$ where A is symmetric, positive semi-definite, HB momentum with $\eta = \frac{4}{(\sqrt{L} + \sqrt{\mu})^2}$ and $\beta = \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}\right)^2$

converges as: $\|w_T - w^*\| \leq \sqrt{2} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} + \epsilon_T\right)^T \|w_0 - w^*\|$, where, $\lim_{T \rightarrow \infty} \epsilon_T \rightarrow 0$.

Proof:

$$\begin{aligned} \begin{bmatrix} w_{k+1} - w^* \\ w_k - w^* \end{bmatrix} &= \begin{bmatrix} w_k - w^* - \eta \nabla f(w_k) + \beta(w_k - w_{k-1}) \\ w_k - w^* \end{bmatrix} \\ &= \begin{bmatrix} w_k - w^* - \eta A(w_k - w^*) + \beta(w_k - w^*) - \beta(w_{k-1} - w^*) \\ w_k - w^* \end{bmatrix} \end{aligned}$$

(Since $\nabla f(w) = Aw$, $Aw^* = b$)

$$\implies \begin{bmatrix} w_{k+1} - w^* \\ w_k - w^* \end{bmatrix} = \begin{bmatrix} (1 + \beta)I_d - \eta A & -\beta I_d \\ I_d & 0 \end{bmatrix} \begin{bmatrix} w_k - w^* \\ w_{k-1} - w^* \end{bmatrix}$$

If $\beta = 0$, we can recover the same equation as GD.

Minimizing strongly-convex quadratics with HB momentum

$$\underbrace{\begin{bmatrix} w_{k+1} - w^* \\ w_k - w^* \end{bmatrix}}_{:=\Delta_{k+1} \in \mathbb{R}^{2d}} = \underbrace{\begin{bmatrix} (1 + \beta)I_d - \eta A & -\beta I_d \\ I_d & 0 \end{bmatrix}}_{:=\mathcal{H} \in \mathbb{R}^{2d \times 2d}} \underbrace{\begin{bmatrix} w_k - w^* \\ w_{k-1} - w^* \end{bmatrix}}_{:=\Delta_k \in \mathbb{R}^{2d}} \implies \Delta_{k+1} = \mathcal{H} \Delta_k$$

Recurring from $k = 0$ to $T - 1$, and taking norm,

$$\|\Delta_T\| = \|\mathcal{H}^T \Delta_0\| \leq \|\mathcal{H}^T\| \left\| \begin{bmatrix} w_0 - w^* \\ w_{-1} - w^* \end{bmatrix} \right\| \quad (\text{By definition of the matrix norm})$$

Define $w_{-1} = w_0$ and lower-bounding the LHS,

$$\|w_T - w^*\| \leq \sqrt{2} \|\mathcal{H}^T\| \|w_0 - w^*\|$$

Hence, we have reduced the problem to bounding $\|\mathcal{H}^T\|$.

Minimizing strongly-convex quadratics with HB momentum

Recall that for symmetric matrices, $\|B\|_2 = \rho(B)$. Unfortunately, this relation is not true for general asymmetric matrices, and $\|B\| \geq \rho(B)$.

Gelfand's Formula: For a matrix $B \in \mathbb{R}^{d \times d}$ such that $\rho(B) := \max_{i \in [d]} |\lambda_i|$, then there exists a sequence $\epsilon_k \geq 0$ such that $\lim_{k \rightarrow \infty} \epsilon_k = 0$ and,

$$\|B^k\| \leq (\rho(B) + \epsilon_k)^k.$$

Using this formula with our bound,

$$\|w_T - w^*\| \leq \sqrt{2} (\rho(\mathcal{H}) + \epsilon_T)^T \|w_0 - w^*\|$$

Hence, we have reduced the problem to bounding $\rho(\mathcal{H})$.

Minimizing strongly-convex quadratics with HB momentum

Similar to the GD case, let $A = U\Lambda U^\top$ be the eigen-decomposition of A , then, $(1 + \beta)I_d - \eta A = USU^\top$ where $S_{i,i} = 1 + \beta - \eta\lambda_i$. Hence,

$$\mathcal{H} = \begin{bmatrix} U^\top & 0 \\ 0 & U^\top \end{bmatrix} \underbrace{\begin{bmatrix} (1 + \beta)I_d - \eta\Lambda & -\beta I_d \\ I_d & 0 \end{bmatrix}}_{:=H} \begin{bmatrix} U & 0 \\ 0 & U \end{bmatrix}$$

Since U is orthonormal, $\rho(\mathcal{H}) = \rho(H)$. Hence we have reduced the problem to bounding $\rho(H)$.

Minimizing strongly-convex quadratics with HB momentum

Let P be a permutation matrix such that:

$$P_{i,j} = \begin{cases} 1 & i \text{ is odd, } j = i \\ 1 & i \text{ is even, } j = d + i \\ 0 & \text{otherwise} \end{cases} \quad B = P H P^T = \begin{bmatrix} H_1 & 0 & \dots & 0 \\ 0 & H_2 & \dots & 0 \\ \vdots & \ddots & & \\ 0 & & 0 & H_d \end{bmatrix}$$

where,

$$H_i = \begin{bmatrix} (1 + \beta) - \eta\lambda_i & -\beta \\ 1 & 0 \end{bmatrix}$$

Note that $\rho(H) = \rho(B)$ (a permutation matrix does not change the eigenvalues). Since B is a block diagonal matrix, $\rho(B) = \max_i [\rho(H_i)]$. Hence we have reduced the problem to bounding $\rho(H_i)$.

Minimizing strongly-convex quadratics with HB momentum

For a fixed $i \in [d]$, let us compute the eigenvalues of $H_i \in \mathbb{R}^{2 \times 2}$ by solving the characteristic polynomial: $\det(H_i - uI_2) = 0$ w.r.t u .

$$u^2 - (1 + \beta - \eta\lambda_i)u + \beta = 0 \implies u = \frac{1}{2} \left[(1 + \beta - \eta\lambda_i) \pm \sqrt{(1 + \beta - \eta\lambda_i)^2 - 4\beta} \right]$$

Let us set β such that, $(1 + \beta - \eta\lambda_i)^2 \leq 4\beta$. This ensures that the roots to the above equation are complex conjugates. Hence,

$$1 + \beta - \eta\lambda_i \geq -2\sqrt{\beta} \implies (\sqrt{\beta} + 1) \geq \sqrt{\eta\lambda_i} \implies \beta \geq (1 - \sqrt{\eta\lambda_i})^2$$

If we ensure that $\beta \geq (1 - \sqrt{\eta\lambda_i})^2$

$$\begin{aligned} u &= \frac{1}{2} \left[(1 + \beta - \eta\lambda_i) \pm i\sqrt{4\beta - (1 + \beta - \eta\lambda_i)^2} \right] \\ \implies |u|^2 &= \frac{1}{4} \left[(1 + \beta - \eta\lambda_i)^2 + 4\beta - (1 + \beta - \eta\lambda_i)^2 \right] = \beta \implies |u| = \sqrt{\beta}. \end{aligned}$$

Hence, if $\beta \geq (1 - \sqrt{\eta\lambda_i})^2$, $\rho(H_i) = \sqrt{\beta}$ and $\rho(B) = \max_i [\rho(H_i)] = \sqrt{\beta}$.

Minimizing strongly-convex quadratics with HB momentum

Using the result from the previous slide, if we ensure that for all i , $\beta \geq (1 - \sqrt{\eta\lambda_i})^2$, then, $\rho(B) = \sqrt{\beta}$. Hence, we want that,

$$\beta = \max_i \{(1 - \sqrt{\eta\lambda_i})^2\} \leq \max_{\lambda \in [\mu, L]} \{(1 - \sqrt{\eta\lambda})^2\} = \max\{(1 - \sqrt{\eta\mu})^2, (1 - \sqrt{\eta L})^2\}$$

Similar to GD, we equate the two terms in the max,

$$1 + \eta\mu - 2\sqrt{\eta\mu} = 1 + \eta L - 2\sqrt{\eta L} \implies \eta = \frac{4}{(\sqrt{L} + \sqrt{\mu})^2}.$$

With this value of η , $\rho(\mathcal{H}) = \rho(H) = \rho(B) \leq \sqrt{\beta} = \sqrt{\left(1 - \frac{2\sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}}\right)^2} = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}$.

Putting everything together,

$$\|w_T - w^*\| \leq \sqrt{2} \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} + \epsilon_T \right)^T \|w_0 - w^*\|$$

Questions?

Gradient Descent and Newton's method

For L -smooth, μ -strongly convex functions,

- Gradient Descent (GD) results in an $O(\exp(-T/\kappa))$ rate.
- Nesterov acceleration can speed up the convergence and results in an $\Theta(\exp(-T/\sqrt{\kappa}))$ rate.
- Lower-Bound: Without additional assumptions, no first-order algorithm (one that only relies on gradient information) can attain a dimension-free rate faster than $\Omega(\exp(-T/\sqrt{\kappa}))$.

Next, we will use second-order (Hessian) information to minimize twice differentiable, L -smooth and μ -strongly convex functions and get faster rates under additional assumptions.

Gradient Descent and Newton's method

Recall the GD update: $w_{k+1} = w_k - \eta \nabla f(w_k)$. This can also be written as:

$$w_{k+1} = \arg \min_w \left[\underbrace{f(w_k) + \langle \nabla f(w_k), w - w_k \rangle}_{\text{First-order Taylor series approximation}} + \underbrace{\frac{1}{2\eta} \|w - w_k\|^2}_{\text{Stay close to } w_k} \right]$$

i.e., approximate the function by a first-order Taylor series expansion, and minimize it while staying close (in the Euclidean norm) to the current point.

If f is twice-differentiable, and we approximate it by a second-order Taylor series expansion,

$$w_{k+1} = \arg \min_w \left[\underbrace{f(w_k) + \langle \nabla f(w_k), w - w_k \rangle + \frac{1}{2} (w - w_k)^\top \nabla^2 f(w_k) (w - w_k)}_{\text{Second-order Taylor series approximation}} \right]$$
$$\implies w_{k+1} = w_k - [\nabla^2 f(w_k)]^{-1} [\nabla f(w_k)] \quad (\text{Newton Update})$$

Digression - Preconditioned Gradient Descent

Recall that GD achieves an $O\left(\kappa \log\left(\frac{1}{\epsilon}\right)\right)$ convergence rate, and the condition number $\kappa \geq 1$ is the measure of problem difficulty.

Idea: Reparameterize the space so that the minimum function value remains the same, but condition number in the reparameterized space is smaller enabling GD to converge faster.

Example: $\min_{w \in \mathbb{R}^2} f(w) = \frac{1}{2} w^T A w$ where $A = \begin{bmatrix} L & 0 \\ 0 & \mu \end{bmatrix}$. For the above problem, $w^* = 0$, $f(w^*) = 0$ and $\kappa = \frac{L}{\mu}$.

Let us choose a **preconditioning matrix** $Q \in \mathbb{R}^{2 \times 2}$ such that $w = Qv$, and write the reparameterized function $g(v) := \frac{1}{2} [Qv]^T A [Qv] = \frac{1}{2} v^T Q^T A Q v$.

If we choose $Q = \begin{bmatrix} \frac{1}{\sqrt{L}} & 0 \\ 0 & \frac{1}{\sqrt{\mu}} \end{bmatrix}$, $Q^T A Q = I$, $g(v) = \frac{1}{2} v^T v$. Clearly, $v^* = 0$ and $g(v^*) = 0$ and $w^* = Qv^* = 0$. For this problem, $\kappa = 1$ making it easier to solve using GD.

Digression - Preconditioned Gradient Descent

Formalizing the intuition on the previous slide, define a positive definite, symmetric matrix $Q \in \mathbb{R}^{d \times d}$ such that $w = Qv$ and hence, $v = Q^{-1}w$. Define $g(v) := f(Qv)$.

Q: If $w^* = \arg \min_w f(w)$ and $v^* = \arg \min_v g(v)$, is $f(w^*) = g(v^*)$?

Computing the gradient of $g(v)$, $\nabla g(v) = Q^T \nabla f(Qv)$. Running GD on $g(v)$, we get that,

$$\begin{aligned} v_{k+1} &= v_k - \eta \nabla g(v_k) = v_k - \eta [Q^T \nabla f(Qv_k)] = v_k - \eta [Q^T \nabla f(w_k)] \\ \implies Q^{-1}w_{k+1} &= Q^{-1}w_k - \eta [Q \nabla f(w_k)] \implies w_{k+1} = w_k - \eta [QQ^T \nabla f(w_k)] \end{aligned}$$

Define a positive definite, symmetric P such that $P = QQ^T$. Since Q is symmetric, $Q = P^{\frac{1}{2}}$. Hence, for $w = P^{\frac{1}{2}}v$,

$$w_{k+1} = w_k - \eta [P \nabla f(w_k)] \quad (\text{Preconditioned GD})$$

i.e., compute the gradient, “precondition” it by matrix P and then do the GD step.

Digression - Preconditioned Gradient Descent

Equivalent formulations of preconditioned gradient descent to minimize $f(w)$,

- Reparameterizing the space using a positive definite, symmetric matrix $P^{\frac{1}{2}}$ such that $v = P^{-\frac{1}{2}}w$ and using GD to minimize $g(v) := f(P^{\frac{1}{2}}v)$.
- Use GD with the preconditioned gradient $P\nabla f(w)$.
- The preconditioned GD update at iteration k can be written as:

$$w_{k+1} = \arg \min \left[\underbrace{f(w_k) + \langle \nabla f(w_k), w - w_k \rangle}_{\text{First-order Taylor series approximation}} + \underbrace{\frac{1}{2\eta} \|w - w_k\|_{P^{-1}}^2}_{\text{Stay close to } w_k} \right]$$

i.e., approximate the function by a first-order Taylor series expansion, and minimize it while staying close (in the norm induced by matrix P^{-1}) to the current point.

We can also use a different preconditioner at every iteration, i.e.

$$w_{k+1} = w_k - \eta[P_k \nabla f(w_k)]$$

Digression - Preconditioned Gradient Descent

- But what is the “best” P_k around a specific iterate for a specific problem? For this, consider the Hessian of $g(v) = f(P^{\frac{1}{2}}v)$ and let us choose P such that the resulting $\kappa = 1$.

Recall that $\nabla g(v) = P^{\frac{1}{2}} \nabla f(P^{\frac{1}{2}}v)$ and hence, $\nabla^2 g(v) = P^{\frac{1}{2}} [\nabla^2 f(P^{\frac{1}{2}}v)] (P^{\frac{1}{2}})^{\top}$. If $P = [\nabla^2 f(P^{\frac{1}{2}}v)]^{-1} = [\nabla^2 f(w)]^{-1}$, then,

$$\nabla^2 g(v) = [\nabla^2 f(P^{\frac{1}{2}}v)]^{-\frac{1}{2}} [\nabla^2 f(P^{\frac{1}{2}}v)] [\nabla^2 f(P^{\frac{1}{2}}v)]^{-\frac{1}{2}} = I_d$$

Around iterate w_k , define $P_k := [\nabla^2 f(w_k)]^{-1}$ and using the equivalence to preconditioned gradient descent, the resulting update can be written as:

$$w_{k+1} = w_k - \eta [\nabla^2 f(w_k)]^{-1} \nabla f(w_k)$$

If $\eta = 1$, we have recovered the Newton method! Hence, the Newton method can be thought of as finding the best preconditioner (one that minimizes the condition number) at every iteration of preconditioned GD.

Newton Method

Using the equivalence to preconditioned GD, the Newton method is also equivalent to:

$$w_{k+1} = \arg \min \left[\underbrace{f(w_k) + \langle \nabla f(w_k), w - w_k \rangle}_{\text{First-order Taylor series approximation}} + \underbrace{\frac{1}{2\eta} \|w - w_k\|_{\nabla^2 f(w_k)}^2}_{\text{Stay close to } w_k} \right]$$

i.e., approximate the function by a first-order Taylor series expansion, and minimize it while staying close (in the “local norm” induced by the Hessian at w_k) to the current point.

Example: Consider solving $w^* = \arg \min f(w) := \frac{1}{2} w^T A w - b w + c$. We know that $\nabla f(w) = A w - b = A(w - w^*)$ and $\nabla^2 f(w) = A$. Starting from point w_0 , consider the Newton update with $\eta = 1$,

$$w_1 = w_0 - [A^{-1}] A(w_0 - w^*) = w^*$$

i.e. the Newton method can minimize quadratics in one step. In this case, $P_k = P = A^{-1}$ and hence, $g(v) = f(A^{-\frac{1}{2}} v) = \frac{1}{2} [A^{-\frac{1}{2}} v]^T A [A^{-\frac{1}{2}} v] - b [A^{-\frac{1}{2}} v] + c = \frac{1}{2} v^T v - b A^{-\frac{1}{2}} v + c$. Computing the Hessian of $g(v)$, $\nabla^2 g(v) = I_d$ which has $\kappa = 1$.

Questions?