

Recall: When we analyzed GD and SGD for strongly convex objectives, the convergence rate depended on the condition number $\kappa = L/\mu$.

For example, the noise ball for SGD w/ constant step size was determined by

$$\begin{aligned} \mathbb{E}[f(\omega_t) - f^*] &\leq (1-\alpha\mu)^T(f(\omega_0) - f^*) + \frac{\alpha\sigma^2 L}{2\mu} \\ &= (1-\alpha\mu)^T(f(\omega_0) - f^*) + \frac{\alpha\sigma^2 \kappa}{2} \end{aligned}$$

and the convergence rate for descent w/ constant step size was

$$\begin{aligned} f(\omega_T) - f^* &\leq \exp\left(-\frac{\mu T}{2}\right)(f(\omega_0) - f^*) \\ &= \exp\left(-\frac{T}{\kappa}\right)(f(\omega_0) - f^*) \end{aligned}$$

Takeaway: When κ high, convergence slow!

How can we speed up when κ high?

Let's consider the simplest possible setting with a high condition number: a two-dimensional quadratic.

Specifically,

$$f(w) = f(w_1, w_2) = \frac{L}{2} w_1^2 + \frac{\mu}{2} w_2^2 = \frac{L}{2} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}^\top \begin{bmatrix} L & 0 \\ 0 & \mu \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

Then

$$\nabla f(w) = \begin{bmatrix} Lw_1 \\ \mu w_2 \end{bmatrix}, \quad \nabla^2 f(w) = \begin{bmatrix} L & 0 \\ 0 & \mu \end{bmatrix}$$

Here, the optimum value occurs at $w^* = 0$, and the Hessian is constan.

Eigenvalues of $\nabla^2 f(w)$ are L, μ where $L \geq \mu > 0$. Thus f is strongly convex w/ strong convexity constant μ and its second derivative is bounded by L .

Gradient descent w/ constant learning rate on this problem has update rule

$$\begin{aligned} (w_{t+1})_1 &= (w_t)_1 - \alpha L (w_t)_1 \\ (w_{t+1})_2 &= (w_t)_2 - \alpha \mu (w_t)_2 \\ \Rightarrow (w_K)_1 &= (1 - \alpha L)^K (w_0)_1 \\ (w_K)_2 &= (1 - \alpha \mu)^K (w_0)_2 \end{aligned}$$

and so

$$f(w_K) = \frac{L}{2} (1 - \alpha_L)^{2K} ((w_0)_1)^2 + \frac{\mu}{2} (1 - \alpha_\mu)^{2K} ((w_0)_2)^2$$

Asymptotically, this will be dominated by whichever term grows slower. That is

$$f(\omega_K) = O\left(\max(|1-\alpha L|, |1-\alpha \mu|)^{2K}\right)$$

Looking at the inner maximum function of α , a curious effect emerges.

So while we would want to pick a smaller learning rate for the L term, and a larger rate for the μ term, we're forced to pick something in the middle.

Aside
$$\alpha = \frac{2}{L+\mu} \Rightarrow \max(|1-\alpha L|, |1-\alpha \mu|) = \frac{L-\mu}{L+\mu} = \frac{\kappa-1}{\kappa+1} = 1 - \frac{2}{\kappa+1}$$

That is, we'd like to set the step size larger for dimensions with less curvature, and smaller for dimensions with more curvature.

Can't do w/ GD or SGD b/c only one parameter.

Today's fix?

Momentum

Motivation: try to tell the difference b/t more and less curved directions using information already available in GD.

Idea: In 1-D case, if gradients are reversing sign, step size is too large, b/c we're overshooting the optimum.

Conversely, if the gradients are staying in the same direction, then the step size is too small.

Polyak Momentum Step: Adds an extra momentum term to gradient descent.

$$\omega_{t+1} = \omega_t - \alpha \nabla f(\omega_t) + \beta(\omega_t - \omega_{t-1})$$

Intuition: If current gradient is in the same direction as the previous step, move a little further in the same direction. If it's in the opposite direction, move a little less far.

Analysis of Polyak Momentum

Let's observe the case of one dimensional quadratics.

$$f(\omega) = \frac{\lambda}{2} \omega^2$$

Using this, Polyak update step looks like

$$\begin{aligned} w_{t+1} &= w_t - \alpha \lambda w_t + \beta (w_t - w_{t-1}) \\ &= (1 + \beta - \alpha \lambda) w_t - \beta w_{t-1}. \end{aligned}$$

Write in terms of matrix multiplication as

$$\begin{bmatrix} w_{t+1} \\ w_t \end{bmatrix} = \begin{bmatrix} 1 + \beta - \alpha \lambda & -\beta \\ 1 & 0 \end{bmatrix} \begin{bmatrix} w_t \\ w_{t-1} \end{bmatrix}$$

This is called the companion matrix of the process.

Call this matrix M .

By induction, we get

$$\begin{bmatrix} w_{t+1} \\ w_t \end{bmatrix} = M \begin{bmatrix} w_t \\ w_{t-1} \end{bmatrix} = M^t \begin{bmatrix} w_1 \\ w_0 \end{bmatrix}$$

Assume M diagonalizable. (basis of eigenvectors $\Rightarrow M = Q D Q^{-1}$)

then

$$\begin{bmatrix} w_{t+1} \\ w_t \end{bmatrix} = (Q D Q^{-1})^t \begin{bmatrix} w_1 \\ w_0 \end{bmatrix} = Q D Q^{-1} \begin{bmatrix} w_1 \\ w_0 \end{bmatrix}$$

Taking norms,

Induced norm: $\|A\| = \max_{\|x\|=1} \frac{\|Ax\|}{\|x\|}$

$$\begin{aligned} \left\| \begin{bmatrix} w_{t+1} \\ w_t \end{bmatrix} \right\| &= \left\| Q D^t Q^{-1} \begin{bmatrix} w_1 \\ w_0 \end{bmatrix} \right\| \leq \|Q\| \|D^t\| \|Q^{-1}\| \left\| \begin{bmatrix} w_1 \\ w_0 \end{bmatrix} \right\| \\ &= \|Q\| \|Q^{-1}\| \left\| \begin{bmatrix} w_1 \\ w_0 \end{bmatrix} \right\| \cdot \max_{i \text{ eigenvals}} \chi_i^t \end{aligned}$$

So, this algorithm will converge at a linear rate determined by the maximum absolute value of an eigenvalue of M .

To understand how Polyak momentum converges, we need to understand the eigenvalues of M .

To do this, observe characteristic polynomial of M is

$$|\lambda I - M| = \left| \begin{bmatrix} \lambda - (1 + \beta - \alpha\lambda) & \beta \\ -1 & \lambda \end{bmatrix} \right|$$

$$= \lambda^2 - (1 + \beta - \alpha\lambda)\lambda + \beta$$

So, the eigenvalues of M are

$$\lambda = \frac{(1 + \beta - \alpha\lambda) \pm \sqrt{(1 + \beta - \alpha\lambda)^2 - 4\beta}}{2}$$

Now we note that if

$$(1 + \beta - \alpha\lambda)^2 - 4\beta < 0 \iff -1 \leq \frac{1 + \beta - \alpha\lambda}{2\sqrt{\beta}} \leq 1$$

then eigenvalues are complex conjugates!

That is,

$$\lambda_1^* = \lambda_2 \quad \text{and} \quad |\lambda_1|^2 = |\lambda_2|^2 = \lambda_1^* \lambda_1 = \lambda_1 \lambda_2$$

But $\lambda_1 \lambda_2$ is the product of eigenvalues of M
So

$$\lambda_1 \lambda_2 = \beta \xleftarrow{\det(M)}$$

So

$$|\lambda_1| = |\lambda_2| = \sqrt{\beta}$$

This means momentum converges at a rate of $\sqrt{\beta}$ in this case!

Or, equivalently w_t^2 will approach zero at a linear rate of β^t for ANY step size (as long as condition above satisfied). i.e $\|w_t\|^2 = O(\beta^t)$

What does this mean for multidimensional quadratics?

$$f(w) = w^T A w \xleftarrow{\text{A symmetric, PSD assumed}}$$
$$\nabla f(w) = Aw$$

$$w_{t+1} = w_t - \alpha A w_t + \beta (w_t - w_{t-1})$$

If u is an eigenvector of A , $Au = \lambda u \rightarrow u^T A^T = \lambda u^T$

$$\begin{aligned} u^T w_{t+1} &= u^T w_t - \alpha u^T A w_t + \beta (u^T w_t - u^T w_{t-1}) \\ &= u^T w_t - \alpha \lambda u^T w_t + \beta (u^T w_t - u^T w_{t-1}) \end{aligned}$$

Importantly, multidimensional quadratics will converge at the same rate β^t , EVEN for directions of different curvature!

This will happen as long as

$$-1 \leq \frac{1 + \beta - \alpha\lambda}{2\sqrt{\beta}} \leq 1.$$

Since $\mu \leq \lambda \leq L$, this will hold iff

$$-1 = \frac{1 + \beta - \alpha L}{2\sqrt{\beta}} \quad \text{and} \quad \frac{1 + \beta - \alpha\mu}{2\sqrt{\beta}} = 1$$

And through solving for α, β we get

$$\alpha = \frac{2 + 2\beta}{L + \mu} \quad \text{and} \quad \sqrt{\beta} = 1 - \frac{2}{\sqrt{\kappa} + 1}$$

$\curvearrowleft \sqrt{\kappa} !$

that is,

$$\|\omega_T\| = O\left(\left(1 - \frac{2}{\sqrt{\kappa} + 1}\right)^T\right)$$

$$= O\left(\exp\left(-\frac{2T}{\sqrt{\kappa} + 1}\right)\right)$$

Momentum has $O(\delta)$ extra memory and computational cost.

How to set momentum?

- Typically, just set $\beta = 0.9$
- Hyperparameter optimization.

Nesterov Momentum Step: Slightly different than Polyak momentum — GUARANTEED to work for convex functions.

$$\begin{aligned}v_{t+1} &= w_t - \alpha \nabla f(w_t) \\w_{t+1} &= v_{t+1} + \beta (v_{t+1} - v_t)\end{aligned}$$

Main difference: Separate momentum state from the point we are calculating gradient at.

Momentum with SGD

$$\begin{aligned}v_{t+1} &= \beta v_t - \alpha \nabla f_{i_t}(w_t) \\w_{t+1} &= w_t + v_{t+1}\end{aligned}$$