# A SIMPLE PROOF OF HEAVY BALL CONVERGENCE

OMAR HIJAB

Let $f(w)$ be a scalar function of a point $w$ in euclidean space. A basic problem is to minimize $f(w)$, that is, to find or compute a minimizer $w^*$,

$$f(w) \geq f(w^*), \qquad \text{for every } w.$$

A *descent sequence* is a sequence $w_0$, $w_1$, $w_2$, ... satisfying

$$f(w_0) \geq f(w_1) \geq f(w_2) \geq \dots.$$

In a descent sequence, the point after $w = w_n$ is $w^+ = w_{n+1}$, and the point before $w$ is $w^- = w_{n-1}$. Then $(w^-)^+ = w = (w^+)^-$.

We assume the loss function is quadratic,

$$(1) \qquad f(w) = \frac{1}{2} w \cdot Qw - b \cdot w,$$

where the eigenvalues of the symmetric matrix $Q$ lie strictly between positive constants $m < L$. Then there is a unique global minimizer $w^*$. Let

$$(2) \qquad C^2 = \max_\lambda \frac{(L - m)(L - m)}{(L - \lambda)(\lambda - m)},$$

where the maximum is over the eigenvalues of $Q$.

**Theorem** (Polyak [1, 2, 3, 4]). *Suppose $f(w)$ is quadratic, let $r = m/L$, and set $E(w) = |w - w^*|$. Let $C$ be given by (2). Then the descent sequence $w_{-1} = w_0$, $w_1$, $w_2$, ... given by*

$$(3) \qquad w^+ = w - t\nabla f(w) + s(w - w^-)$$

*with learning rate and the decay rate*

$$t = \frac{1}{L} \cdot \frac{4}{(1 + \sqrt{r})^2}, \qquad s = \left(\frac{1 - \sqrt{r}}{1 + \sqrt{r}}\right)^2,$$

*converges to $w^*$ at the rate*

$$(4) \qquad E(w_n) \leq C \left(\frac{1 - \sqrt{r}}{1 + \sqrt{r}}\right)^n E(w_0), \qquad n = 1, 2, \dots$$

*Proof.* Since $\nabla f(w) = Qw - b$, the sequence satisfies

$$(5) \qquad w_{n+1} = w_n - t(Qw_n - b) + s(w_n - w_{n-1}), \qquad n = 0, 1, 2, \dots.$$

To initialize this recursion, we set $w_{-1} = w_0^- = w_0$. This implies $w_1 = w_0 - t(Qw_0 - b)$.

Let $v$ be an eigenvector of $Q$ with eigenvalue $\lambda$. To solve (5), we assume a solution of the form

$$(6) \qquad w_n = w^* + \rho^n v, \qquad Qv = \lambda v.$$

Inserting this into (5) and using $Qw^* = b$ leads to the quadratic equation

$$\rho^2 = (1 - t\lambda + s)\rho - s$$

with discriminant

$$\Delta = (1 - \lambda t + s)^2 - 4s.$$

Now $\Delta < 0$ exactly when

(7)
$$\frac{(1 - \sqrt{s})^2}{\lambda} < t < \frac{(1 + \sqrt{s})^2}{\lambda}.$$

If we assume

(8)
$$\frac{(1 - \sqrt{s})^2}{m} \leq t \leq \frac{(1 + \sqrt{s})^2}{L},$$

then

(9)
$$\Delta \leq -(1 - s)^2 \frac{(L - \lambda)(\lambda - m)}{mL},$$

for every eigenvalue $\lambda$ of $Q$. When $\Delta < 0$, the roots are conjugate complex numbers $\rho, \bar{\rho}$, where

(10)
$$\rho = x + iy = \frac{(1 - \lambda t + s) + i\sqrt{-(1 - \lambda t + s)^2 + 4s}}{2}.$$

It follows the absolute value of $\rho$ equals

$$|\rho| = \sqrt{x^2 + y^2} = \sqrt{s}.$$

To obtain the fastest convergence, we choose $s$ and $t$ to minimize $|\rho| = \sqrt{s}$, while still satisfying (8). This forces (8) to be an equality,

$$\frac{(1 - \sqrt{s})^2}{m} = t = \frac{(1 + \sqrt{s})^2}{L}.$$

These are two equations in two unknowns $s$, $t$. Solving, we obtain the choices for $s$ and $t$ made above.

Since (5) is a 2-step linear recursion, the general solution depends on two constants $A$, $B$. Let $\lambda_1$, $\lambda_2$, ... be the eigenvalues of $Q$ and let $v_1$, $v_2$, ... be the corresponding orthonormal basis of eigenvectors in the euclidean space. Since (5) is a 2-step vector linear recursion, $A$ and $B$ are vectors, and the general solution depends on constants $A_k$, $B_k$ corresponding to each $\lambda_k$, $k = 1, 2, \ldots$.

If $\rho_k$, $k = 1, 2, \ldots$, are the corresponding roots (10), then (6) is a solution of (5) for each of the roots $\rho = \rho_k$ and $\rho = \bar{\rho}_k$, $k = 1, 2, \ldots$. Therefore the linear combination

(11)
$$w_n = w^* + \sum_k \left( A_k \rho_k^n + B_k \bar{\rho}_k^n \right) v_k, \qquad n = 0, 1, 2, \ldots$$

is the general solution of (5). Inserting $n = 0$ and $n = 1$ into (11), then taking the dot product of the result with $v_k$, we obtain two linear equations for two unknowns $A_k$, $B_k$. Solving for $A_k$, $B_k = \bar{A}_k$, then using (9),

$$|A_k| = |B_k| \leq \frac{1}{2}C \left| (w_0 - w^*) \cdot v_k \right|.$$

By orthonormality of the basis,

$$\begin{aligned}
|w_n - w^*|^2 &= \sum_k |A_k \rho_k^n + B_k \bar{\rho}_k^n|^2 \\
&\leq \sum_k (|A_k| + |B_k|)^2 s^n \\
&\leq C^2 s^n \sum_k |(w_0 - w^*) \cdot v_k|^2 \\
&= C^2 s^n |w_0 - w^*|^2.
\end{aligned}$$

$\square$

Note: Since the proof is dimension-independent, a version of the result should hold in Hilbert space.

## References

[1] Sébastien Bubeck, *Convex Optimization: Algorithms and Complexity*, Now Publishers (2015).
[2] Yurii Nesterov, *Lectures on Convex Optimization*, Springer (2018).
[3] Boris Teodorovich Polyak, *Some methods of speeding up the convergence of iteration methods*, USSR Computational Mathematics and Mathematical Physics, 4(5) 1-17 (1964).
[4] Stephen J. Wright and Benjamin Recht, *Optimization for Data Analysis*, Cambridge University (2022).

TEMPLE UNIVERSITY
*Email address*: hijab@temple.edu