# Gradient Descent with Weighted Inner Product

## Problem

Our goal is to solve the optimization problem

$$\min_x f(x),$$

where $f \colon \mathbb{R}^n \to \mathbb{R}$ is a convex and differentiable function.

Let $P \in \mathbb{R}^{n \times n}$ be a symmetric positive definite matrix. This matrix induces a new (weighted) inner product defined as $\langle x, y \rangle_P = \langle Px, y \rangle$ which, in turn, induces a new gradient operator $\nabla_P f(x)$ with respect to this inner product. (Why does this make sense?)

This motivates us to consider a generalized version of gradient descent using the new gradient:

$$x_{k+1} = x_k - \alpha \nabla_P f(x_k), \tag{1}$$

where $\alpha > 0$ is the step size.

**TODO:**

- Find an explicit form of $\nabla_P f(x)$.

- Think of other ingredients/assumptions you need (such as the Lipschitzness of $\nabla_P f$) and prove the convergence of (1).

- Hint: you should understand well main ingredients of the standard proof of GD in the convex case and adjust them to your setting.

## Base Statements

**Lemma 2.28.** If $f$ is $L$–smooth and $\gamma > 0$, then for all $x, y \in \mathbb{R}^d$,

$$f(x - \gamma \nabla f(x)) - f(x) \leq -\gamma \left( 1 - \frac{\gamma L}{2} \right) \| \nabla f(x) \|^2. \tag{10}$$

If moreover $\inf f > -\infty$, then for all $x \in \mathbb{R}^d$,

$$\frac{1}{2L} \| \nabla f(x) \|^2 \leq f(x) - \inf f.$$

# General Proof of Convergence of Gradient Descent

**Theorem** Consider the Problem (Differentiable Function) and assume that $f$ is convex and $L$-smooth, for some $L > 0$. Let $(x_t)_{t \in \mathbb{N}}$ be the sequence of iterates generated by the (GD) algorithm, with a stepsize satisfying $0 < \gamma \le \frac{1}{L}$. Then, for all $x^* \in \arg \min f$, for all $t \in \mathbb{N}$, we have:

$$f(x_t) - \inf f \le \frac{\|x_0 - x^*\|^2}{2\gamma t}.$$

**Proof** Let $f$ be convex and $L$–smooth. It follows that

$$\|x_{t+1} - x^*\|^2 = \left\| x_t - x^* - \frac{1}{L} \nabla f(x_t) \right\|^2$$

$$= \|x_t - x^*\|^2 - 2 \cdot \frac{1}{L} \langle x_t - x^*, \nabla f(x_t) \rangle + \frac{1}{L^2} \|\nabla f(x_t)\|^2$$

$$\overset{(1)}{\le} \|x_t - x^*\|^2 - \frac{1}{L^2} \|\nabla f(x_t)\|^2. \tag{18}$$

Thus, $\|x_t - x^*\|^2$ is a decreasing sequence in $t$, and consequently

$$\|x_t - x^*\| \le \|x_0 - x^*\|. \tag{19}$$

Calling upon (10) and subtracting $f(x^*)$ from both sides gives

$$f(x_{t+1}) - f(x^*) \le f(x_t) - f(x^*) - \frac{1}{2L} \|\nabla f(x_t)\|^2. \tag{20}$$

Applying convexity we have that

$$f(x_t) - f(x^*) \le \langle \nabla f(x_t), x_t - x^* \rangle$$

$$\le \|\nabla f(x_t)\| \cdot \|x_t - x^*\|$$

$$\overset{(19)}{\le} \|\nabla f(x_t)\| \cdot \|x_0 - x^*\|. \tag{21}$$

Suppose now that $x_0 \ne x^*$, otherwise the proof is finished. Isolating $\|\nabla f(x_t)\|$ in the above and inserting in (20) gives

$$f(x_{t+1}) - f(x^*) \overset{(20)+(21)}{\le} f(x_t) - f(x^*) - \frac{1}{2L} \frac{1}{\|x_0 - x^*\|^2} (f(x_t) - f(x^*))^2 \tag{22}$$

Let $\beta = \frac{1}{2L} \frac{1}{\|x_0 - x^*\|^2}$ and $\delta_t = f(x_t) - f(x^*)$. Since $\delta_{t+1} \le \delta_t$, and by manipulating (22) we have that

$$\delta_{t+1} \le \delta_t - \beta \delta_t^2 \xleftarrow{\times \frac{1}{\delta_t \delta_{t+1}}} \beta \frac{\delta_t}{\delta_{t+1}} \le \frac{1}{\delta_{t+1}} - \frac{1}{\delta_t} \xleftarrow{\delta_{t+1} \le \delta_t} \beta \le \frac{1}{\delta_{t+1}} - \frac{1}{\delta_t}.$$

Summing up both sides over $t = 0, \ldots, T-1$ and using telescopic cancellation we have that

$$T\beta \leq \frac{1}{\delta_T} - \frac{1}{\delta_0} \leq \frac{1}{\delta_T}.$$

Re-arranging the above we have that

$$f(x^T) - f(x^*) = \delta_T \leq \frac{1}{\beta T} = \frac{2L\|x^0 - x^*\|^2}{T}.$$

# Proof of Convergence of Gradient Descent with weighted inner product

From now on, we will use the following notions:

$$\nabla f(x) = P\nabla_P f(x),$$
$$P^{-1}\nabla f(x) = \nabla_P f(x),$$
$$x_{t+1} = x_t - \eta\nabla_P f(x_t) \quad \Leftrightarrow \quad x_{t+1} = x_t - \eta P^{-1}\nabla f(x_t).$$

If you see an inner product written without the subscript $P$, this is done deliberately and refers to the standard Euclidean inner product.

**Proof.** Consider the norm induced by $P$: $\|x\|_P^2 = x^\top P x$. Then the gradient step becomes

$$x_{t+1} = x_t - \eta P^{-1}\nabla f(x_t),$$

which can be written as

$$x_{t+1} - x^* = x_t - x^* - \eta P^{-1}\nabla f(x_t).$$

Taking the squared $P$-norm of both sides:

$$\|x_{t+1} - x^*\|_P^2 = \left\|x_t - x^* - \eta P^{-1}\nabla f(x_t)\right\|_P^2$$
$$= \|x_t - x^*\|_P^2 - 2\eta\langle P^{-1}\nabla f(x_t), x_t - x^*\rangle_P + \eta^2\|P^{-1}\nabla f(x_t)\|_P^2$$
$$= \|x_t - x^*\|_P^2 - 2\eta\langle \nabla f(x_t), x_t - x^*\rangle + \eta^2\nabla f(x_t)^\top P^{-1}\nabla f(x_t),$$

where we used the identity $\langle u, v\rangle_P = u^\top P v$ and the fact that $PP^{-1} = I$.

Now, suppose $f$ is convex and $L_P$-smooth with respect to the $P$-norm.

**Yura:** *What does it mean? And why is the next inequality true?*

Then, from standard smoothness inequality:

$$f(x_{t+1}) \leq f(x_t) + \langle \nabla f(x_t), x_{t+1} - x_t\rangle + \frac{L_P}{2}\|x_{t+1} - x_t\|_P^2. \qquad (?)$$

Substitute $x_{t+1} - x_t = -\eta P^{-1}\nabla f(x_t)$:

$$f(x_{t+1}) \leq f(x_t) - \eta\nabla f(x_t)^\top P^{-1}\nabla f(x_t) + \frac{L_P\eta^2}{2}\nabla f(x_t)^\top P^{-1}\nabla f(x_t)$$
$$= f(x_t) - \left(\eta - \frac{L_P\eta^2}{2}\right)\nabla f(x_t)^\top P^{-1}\nabla f(x_t).$$

3

Choosing $\eta = \frac{1}{L_P}$, we obtain:

$$f(x_{t+1}) \le f(x_t) - \frac{1}{2L_P} \nabla f(x_t)^\top P^{-1} \nabla f(x_t).$$

From convexity, we also have:

$$f(x_t) - f(x^*) \le \langle \nabla f(x_t), x_t - x^* \rangle.$$

Using Cauchy-Schwarz in $P$-norm (It was assumed that Cauchy-Schwarz inequality hold in any weighted inner product space? Link):

**Yura:** *Not clear*

$$\langle \nabla f(x_t), x_t - x^* \rangle \le \|x_t - x^*\|_P \cdot \|P^{-1} \nabla f(x_t)\|_P.$$

Note that:

$$\|P^{-1} \nabla f(x_t)\|_P^2 = \nabla f(x_t)^\top P^{-1} \nabla f(x_t).$$

So we get:

$$f(x_t) - f(x^*) \le \|x_t - x^*\|_P \cdot \sqrt{\nabla f(x_t)^\top P^{-1} \nabla f(x_t)} \le \|x_0 - x^*\|_P \cdot \sqrt{\nabla f(x_t)^\top P^{-1} \nabla f(x_t)}.$$

Solving for $\nabla f(x_t)^\top P^{-1} \nabla f(x_t)$ and plugging into the earlier bound:

$$f(x_{t+1}) - f(x^*) \le f(x_t) - f(x^*) - \frac{1}{2L_P} \cdot \frac{(f(x_t) - f(x^*))^2}{\|x_0 - x^*\|_P^2}.$$

Letting $\delta_t = f(x_t) - f(x^*)$, and $\beta = \frac{1}{2L_P \|x_0 - x^*\|_P^2}$, we obtain:

$$\delta_{t+1} \le \delta_t - \beta \delta_t^2.$$

As in standard analysis, we get ($t = 0, \ldots, T-1$):

$$\delta_{t+1} \le \delta_t - \beta \delta_t^2 \xleftrightarrow{\times \frac{1}{\delta_t \delta_{t+1}}} \beta \frac{\delta_t}{\delta_{t+1}} \le \frac{1}{\delta_{t+1}} - \frac{1}{\delta_t} \xleftrightarrow{\delta_{t+1} \le \delta_t} \beta \le \frac{1}{\delta_{t+1}} - \frac{1}{\delta_t}.$$

$$\beta \le \frac{1}{\delta_{t+1}} - \frac{1}{\delta_t} \quad \Rightarrow \quad T\beta \le \frac{1}{\delta_T} - \frac{1}{\delta_0} \le \frac{1}{\delta_T},$$

which implies:

$$f(x_T) - f(x^*) = \delta_T \le \frac{1}{\beta T} = \frac{2L_P \|x^0 - x^*\|_P^2}{T}.$$

$\square$