

# ADFOCS '21 Summer School: Adaptive Gradient Descent Algorithms

A tentative plan is to cover the first 4 exercises in the first session, and the remaining exercises in the second session.

In the following exercises, we give a proof of the results we used in our analysis of adaptive methods.

**Exercise 1.** Let  $a_1, a_2, \dots, a_n > 0$  be positive scalars. Show that

$$\sqrt{\sum_{i=1}^n a_i} \leq \sum_{i=1}^n \frac{a_i}{\sqrt{\sum_{j=1}^i a_j}} \leq 2\sqrt{\sum_{i=1}^n a_i}$$

*Hint:* Approximate the sum by an integral.

*Proof.* Let  $A_i = \sum_{j=1}^i a_j$ ,  $A_0 = 0$ . We have

$$\begin{aligned} \sum_{i=1}^n \frac{A_i - A_{i-1}}{\sqrt{A_i}} &= \sum_{i=1}^n \frac{(\sqrt{A_i} - \sqrt{A_{i-1}})(\sqrt{A_i} + \sqrt{A_{i-1}})}{\sqrt{A_i}} \\ &\leq 2 \sum_{i=1}^n (\sqrt{A_i} - \sqrt{A_{i-1}}) \\ &= 2(\sqrt{A_n} - \sqrt{A_0}) \end{aligned}$$

$$\begin{aligned} \sum_{i=1}^n \frac{A_i - A_{i-1}}{\sqrt{A_i}} &= \sum_{i=1}^n \frac{(\sqrt{A_i} - \sqrt{A_{i-1}})(\sqrt{A_i} + \sqrt{A_{i-1}})}{\sqrt{A_i}} \\ &\geq \sum_{i=1}^n (\sqrt{A_i} - \sqrt{A_{i-1}}) \\ &= \sqrt{A_n} - \sqrt{A_0} \end{aligned}$$

□

**Exercise 2.** Let  $f$  be a differentiable function that is convex and  $\beta$ -smooth. Show that, for all  $x, y$ , we have

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2\beta} \|\nabla f(y) - \nabla f(x)\|^2$$

*Proof. First proof:* Let

$$z = y - \frac{1}{\beta} (\nabla f(y) - \nabla f(x))$$

We write

$$f(x) - f(y) = f(x) - f(z) + f(z) - f(y)$$

We upper bound the first difference using convexity and the second using smoothness:

$$\begin{aligned} f(x) - f(y) &= f(x) - f(z) + f(z) - f(y) \\ &\leq \langle \nabla f(x), x - z \rangle + \langle \nabla f(y), z - y \rangle + \frac{\beta}{2} \|z - y\|^2 \\ &= \langle \nabla f(x), x - y \rangle - \frac{1}{2\beta} \|\nabla f(x) - \nabla f(y)\|^2 \end{aligned}$$

**Second proof.** Let

$$g(z) = f(z) - (f(x) + \langle \nabla f(x), z - x \rangle)$$

Note that  $g$  is minimized at  $x$ , since  $g$  is convex and  $\nabla g(x) = 0$ . Thus  $\min_{z \in \mathbb{R}^d} g(z) = g(x) = 0$  and thus  $g(z) \geq 0$  for all  $z \in \mathbb{R}^d$ . Since  $f$  is  $\beta$ -smooth,

$$f(z) \leq f(y) + \langle \nabla f(y), z - y \rangle + \frac{\beta}{2} \|z - y\|^2$$

Therefore, for all  $z \in \mathbb{R}^d$ , we have

$$\begin{aligned} 0 &\leq g(z) \\ &= f(z) - (f(x) + \langle \nabla f(x), z - x \rangle) \\ &\leq f(y) + \langle \nabla f(y), z - y \rangle + \frac{\beta}{2} \|z - y\|^2 - (f(x) + \langle \nabla f(x), z - x \rangle) \\ &= f(y) - f(x) + \langle \nabla f(x), x - y \rangle + \langle \nabla f(y) - \nabla f(x), z - y \rangle + \frac{\beta}{2} \|z - y\|^2 \end{aligned}$$

Thus we have

$$\begin{aligned} 0 &\leq f(y) - f(x) + \langle \nabla f(x), x - y \rangle + \min_{z \in \mathbb{R}^d} \left( \langle \nabla f(y) - \nabla f(x), z - y \rangle + \frac{\beta}{2} \|z - y\|^2 \right) \\ &= f(y) - f(x) + \langle \nabla f(x), x - y \rangle - \frac{1}{2\beta} \|\nabla f(y) - \nabla f(x)\|^2 \end{aligned}$$

□

**Adagrad with per-coordinate step sizes** In the next exercise, we consider a version of the Adagrad algorithm from Lecture I that uses per-coordinate step sizes. Let  $\mathbf{D} \in \mathbb{R}^{d \times d}$  be a diagonal matrix with positive diagonal entries. We let  $\mathbf{D}_i$  denote the  $i$ -th diagonal entry of  $\mathbf{D}$ . We also let

$$\|x\|_{\mathbf{D}} := \sqrt{x^\top \mathbf{D} x} = \sqrt{\sum_{i=1}^d \mathbf{D}_i x_i^2}$$

We let  $\|x\|_{\infty}$  denote the  $\ell_{\infty}$ -norm of  $x$ :  $\|x\|_{\infty} = \max_{i \in [d]} |x_i|$ . We let  $\nabla_i f(x)$  denote the  $i$ -th coordinate of  $\nabla f(x)$ .

---

**Algorithm 1** Adagrad algorithm with per-coordinate step sizes.

---

Let  $x_1 \in K$ ,  $R_{\infty} \geq \max_{x, y \in K} \|x - y\|_{\infty}$

For  $t = 1, \dots, T$ :

$$\begin{aligned} \mathbf{D}_{t,i} &= \frac{1}{R_{\infty}} \sqrt{\sum_{s=1}^t (\nabla_i f(x_s))^2} && \forall i \in [d] \\ x_{t+1} &= \arg \min_{u \in K} \left\{ \langle \nabla f(x_t), u - x_t \rangle + \frac{1}{2} \|u - x_t\|_{\mathbf{D}_t}^2 \right\} \end{aligned}$$

Return  $\bar{x}_T = \frac{1}{T} \sum_{t=1}^T x_t$

---

In the unconstrained setting ( $K = \mathbb{R}^d$ ), the update becomes

$$\begin{aligned} x_{t+1} &= x_t - \mathbf{D}_t^{-1} \nabla f(x_t) \\ \Rightarrow x_{t+1,i} &= x_{t,i} - \frac{R_{\infty}}{\underbrace{\sqrt{\sum_{s=1}^t (\nabla_i f(x_s))^2}}_{:= \eta_{t,i}}} \nabla_i f(x_t) \end{aligned}$$

Thus the algorithm can be interpreted as building a diagonal preconditioner based on the gradient information seen so far. This approach is beneficial in practice since the coordinates may be very different (e.g., we can set high learning rates for informative features).

**Exercise 3.** Extend the analysis that we saw in Lecture I to Algorithm 1. You may find the following inequality helpful:

$$\|x - y\|_{\mathbf{D}}^2 \leq \|x - y\|_{\infty}^2 \operatorname{tr}(\mathbf{D})$$

where  $\operatorname{tr}(\mathbf{D}) = \sum_{i=1}^d \mathbf{D}_i$  is the trace of the diagonal matrix  $\mathbf{D}$ .

*Proof.* By convexity, we have

$$f(\bar{x}_T) - f(x^*) \leq \frac{1}{T} \sum_{t=1}^T (f(x_t) - f(x^*)) \leq \frac{1}{T} \sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle$$

Next, we use the optimality condition for  $x_{t+1}$ :

$$x_{t+1} = \arg \min_{u \in K} \left\{ \langle \nabla f(x_t), u - x_t \rangle + \frac{1}{2} \|u - x_t\|_{\mathbf{D}_t}^2 \right\}$$

$$\langle \nabla f(x_t) + \mathbf{D}_t (x_{t+1} - x_t), x_{t+1} - x^* \rangle \leq 0$$

Thus

$$\begin{aligned} \langle \nabla f(x_t), x_{t+1} - x^* \rangle &\leq \langle \mathbf{D}_t (x_t - x_{t+1}), x_{t+1} - x^* \rangle \\ &= \frac{1}{2} \|x_t - x^*\|_{\mathbf{D}_t}^2 - \frac{1}{2} \|x_{t+1} - x^*\|_{\mathbf{D}_t}^2 - \frac{1}{2} \|x_t - x_{t+1}\|_{\mathbf{D}_t}^2 \end{aligned}$$

We fix the mismatch as before:

$$\begin{aligned} \langle \nabla f(x_t), x_t - x^* \rangle &= \langle \nabla f(x_t), x_{t+1} - x^* \rangle + \langle \nabla f(x_t), x_t - x_{t+1} \rangle \\ &\leq \frac{1}{2} \|x_t - x^*\|_{\mathbf{D}_t}^2 - \frac{1}{2} \|x_{t+1} - x^*\|_{\mathbf{D}_t}^2 + \langle \nabla f(x_t), x_t - x_{t+1} \rangle - \frac{1}{2} \|x_t - x_{t+1}\|_{\mathbf{D}_t}^2 \end{aligned}$$

We can extend the Cauchy-Schwartz inequality to work for the norm  $\|\cdot\|_{\mathbf{D}_t}$ :

$$\langle \nabla f(x_t), x_t - x_{t+1} \rangle \leq \|\nabla f(x_t)\|_{\mathbf{D}_t^{-1}} \|x_t - x_{t+1}\|_{\mathbf{D}_t}$$

Thus

$$\begin{aligned} \langle \nabla f(x_t), x_t - x_{t+1} \rangle - \frac{1}{2} \|x_t - x_{t+1}\|_{\mathbf{D}_t}^2 &\leq \|\nabla f(x_t)\|_{\mathbf{D}_t^{-1}} \|x_t - x_{t+1}\|_{\mathbf{D}_t} - \frac{1}{2} \|x_t - x_{t+1}\|_{\mathbf{D}_t}^2 \\ &\leq \frac{1}{2} \|\nabla f(x_t)\|_{\mathbf{D}_t^{-1}}^2 \end{aligned}$$

Hence

$$\langle \nabla f(x_t), x_t - x^* \rangle \leq \frac{1}{2} \|x_t - x^*\|_{\mathbf{D}_t}^2 - \frac{1}{2} \|x_{t+1} - x^*\|_{\mathbf{D}_t}^2 + \frac{1}{2} \|\nabla f(x_t)\|_{\mathbf{D}_t^{-1}}^2$$

Summing up over all iterations,

$$\begin{aligned} \sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle &\leq \sum_{t=1}^T \left( \frac{1}{2} \|x_t - x^*\|_{\mathbf{D}_t}^2 - \frac{1}{2} \|x_{t+1} - x^*\|_{\mathbf{D}_t}^2 \right) + \sum_{t=1}^T \frac{1}{2} \|\nabla f(x_t)\|_{\mathbf{D}_t^{-1}}^2 \\ &= \sum_{t=1}^T \left( \frac{1}{2} \|x_t - x^*\|_{\mathbf{D}_{t-1}}^2 + \frac{1}{2} \|x_t - x^*\|_{\mathbf{D}_t - \mathbf{D}_{t-1}}^2 - \frac{1}{2} \|x_{t+1} - x^*\|_{\mathbf{D}_t}^2 \right) + \sum_{t=1}^T \frac{1}{2} \|\nabla f(x_t)\|_{\mathbf{D}_t^{-1}}^2 \\ &= \underbrace{\sum_{t=1}^T \left( \frac{1}{2} \|x_t - x^*\|_{\mathbf{D}_{t-1}}^2 - \frac{1}{2} \|x_{t+1} - x^*\|_{\mathbf{D}_t}^2 \right)}_{\text{telescopes}} + \sum_{t=1}^T \frac{1}{2} \|x_t - x^*\|_{\mathbf{D}_t - \mathbf{D}_{t-1}}^2 + \sum_{t=1}^T \frac{1}{2} \|\nabla f(x_t)\|_{\mathbf{D}_t^{-1}}^2 \\ &\leq \frac{1}{2} \|x_1 - x^*\|_{\mathbf{D}_0}^2 + \sum_{t=1}^T \frac{1}{2} \|x_t - x^*\|_{\mathbf{D}_t - \mathbf{D}_{t-1}}^2 + \sum_{t=1}^T \frac{1}{2} \|\nabla f(x_t)\|_{\mathbf{D}_t^{-1}}^2 \end{aligned}$$

Now we would like to telescope the second sum as well. Using the inequality  $\|x - y\|_{\mathbf{D}}^2 \leq \|x - y\|_{\infty}^2 \text{tr}(\mathbf{D})$ , we obtain

$$\begin{aligned} \|x_t - x^*\|_{\mathbf{D}_t - \mathbf{D}_{t-1}}^2 &\leq \underbrace{\|x_t - x^*\|_{\infty}^2}_{\leq R_{\infty}^2} \text{tr}(\mathbf{D}_t - \mathbf{D}_{t-1}) \\ &\leq R_{\infty}^2 \text{tr}(\mathbf{D}_t - \mathbf{D}_{t-1}) \\ &= R_{\infty}^2 (\text{tr}(\mathbf{D}_t) - \text{tr}(\mathbf{D}_{t-1})) \end{aligned}$$

This allows us to telescope the sum:

$$\begin{aligned} \sum_{t=1}^T \frac{1}{2} \|x_t - x^*\|_{\mathbf{D}_t - \mathbf{D}_{t-1}}^2 &\leq \sum_{t=1}^T R_{\infty}^2 (\text{tr}(\mathbf{D}_t) - \text{tr}(\mathbf{D}_{t-1})) \\ &= R_{\infty}^2 (\text{tr}(\mathbf{D}_T) - \text{tr}(\mathbf{D}_0)) \end{aligned}$$

Thus

$$\sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle \leq R_{\infty}^2 \text{tr}(\mathbf{D}_T) + \sum_{t=1}^T \frac{1}{2} \|\nabla f(x_t)\|_{\mathbf{D}_t^{-1}}^2$$

Our choice of step sizes gives us:

$$\begin{aligned} \text{tr}(\mathbf{D}_T) &= \sum_{i=1}^d \mathbf{D}_{T,i} \\ &= \frac{1}{R_{\infty}} \sum_{i=1}^d \sqrt{\sum_{t=1}^T (\nabla_i f(x_t))^2} \\ \sum_{t=1}^T \|\nabla f(x_t)\|_{\mathbf{D}_t^{-1}}^2 &= \frac{1}{R_{\infty}} \sum_{i=1}^d \sum_{t=1}^T \frac{(\nabla_i f(x_t))^2}{\sqrt{\sum_{s=1}^t (\nabla_i f(x_s))^2}} \\ &\leq \frac{1}{R_{\infty}} \sum_{i=1}^d 2 \sqrt{\sum_{t=1}^T (\nabla_i f(x_t))^2} \end{aligned}$$

where we applied the inequality in the first exercise.

Thus

$$\sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle \leq 2R_{\infty} \sum_{i=1}^d \sqrt{\sum_{t=1}^T (\nabla_i f(x_t))^2}$$

Since  $\sqrt{x}$  is concave, we can further bound

$$\begin{aligned} \sum_{i=1}^d \sqrt{\sum_{t=1}^T (\nabla_i f(x_t))^2} &= d \cdot \frac{1}{d} \sum_{i=1}^d \sqrt{\sum_{t=1}^T (\nabla_i f(x_t))^2} \\ &\leq d \sqrt{\frac{1}{d} \sum_{i=1}^d \sum_{t=1}^T (\nabla_i f(x_t))^2} \\ &= \sqrt{d} \sqrt{\sum_{t=1}^T \|\nabla f(x_t)\|^2} \end{aligned}$$

Thus we have obtained

$$\sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle \leq 2R_{\infty} \sqrt{d} \sqrt{\sum_{t=1}^T \|\nabla f(x_t)\|^2} \quad (1)$$

We now consider the non-smooth settings separately and complete the analysis.

**Non-smooth setting:** As before, we assume  $\|\nabla f(x)\| \leq G$ , which gives us

$$\sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle \leq 2R_\infty \sqrt{d} G \sqrt{T}$$

and this

$$f(\bar{x}_T) - f(x^*) \leq O\left(\frac{R_\infty \sqrt{d} G}{\sqrt{T}}\right)$$

**Smooth setting:** As in lecture, we apply the second exercise with  $y = x_t$  and  $x = x^*$ . Recall that we are assuming that  $\nabla f(x^*) = 0$  ( $K$  contains a global minimum). As before, we obtain

$$f(x_t) - f(x^*) \leq \langle \nabla f(x_t), x_t - x^* \rangle - \frac{1}{2\beta} \|\nabla f(x_t)\|^2$$

The rest of the argument is the same as in lecture:

$$\begin{aligned} f(\bar{x}_T) - f(x^*) &\leq \frac{1}{T} \sum_{t=1}^T (f(x_t) - f(x^*)) \\ &\leq \frac{1}{T} \sum_{t=1}^T \left( \langle \nabla f(x_t), x_t - x^* \rangle - \frac{1}{2\beta} \|\nabla f(x_t)\|^2 \right) \\ &\stackrel{(1)}{\leq} \frac{1}{T} \left( 2R_\infty \sqrt{d} \sqrt{\sum_{t=1}^T \|\nabla f(x_t)\|^2} - \frac{1}{2\beta} \sum_{t=1}^T \|\nabla f(x_t)\|^2 \right) \\ &\leq \frac{1}{T} \max_{z \geq 0} \left\{ 2R_\infty \sqrt{d} z - \frac{1}{2\beta} z^2 \right\} \\ &= O\left(\frac{\beta d R_\infty^2}{T}\right) \end{aligned}$$

□

**A different adaptive algorithm for constrained optimization** In the following exercises, we will analyze a different algorithm for constrained optimization. The algorithm is based on a variant of gradient descent that uses extrapolation, which we will discuss in more detail in Lecture III. In contrast to the Adagrad+ algorithm that we saw in Lecture I that uses the iterate movement  $\|x_{t+1} - x_t\|^2$  to set the step sizes, this algorithm uses the gradient differences  $\|\nabla f(x_{t+1}) - \nabla f(x_t)\|^2$ .

---

**Algorithm 2** An adaptive version of the Past Extra-Gradient algorithm.

---

Let  $x_0 = z_0 \in K$ ,  $\eta_0 > 0$ .

For  $t = 1, \dots, T$ , update:

$$\begin{aligned} x_t &= \arg \min_{u \in K} \left\{ \langle \nabla f(x_{t-1}), u \rangle + \frac{1}{2\eta_{t-1}} \|u - z_{t-1}\|^2 + \left( \frac{1}{2\eta_t} - \frac{1}{2\eta_{t-1}} \right) \|u - x_t\|^2 \right\} \\ \eta_t &= \frac{R}{\sqrt{\sum_{i=1}^t \|\nabla f(x_i) - \nabla f(x_{i-1})\|^2}} \\ z_t &= \arg \min_{u \in K} \left\{ \langle \nabla f(x_t), u \rangle + \frac{1}{2\eta_{t-1}} \|u - z_{t-1}\|^2 + \left( \frac{1}{2\eta_t} - \frac{1}{2\eta_{t-1}} \right) \|u - x_t\|^2 \right\} \end{aligned}$$

Return  $\bar{x}_T = \frac{1}{T} \sum_{t=1}^T x_t$ .

---

**Exercise 4.** As in lecture, start by upper bounding the function value difference  $f(\bar{x}_T) - f(x^*)$  in terms of the inner products  $\langle \nabla f(x_t), x_t - x^* \rangle$ .

*Proof.* By convexity,

$$\begin{aligned} f(\bar{x}_T) - f(x^*) &\leq \frac{1}{T} \sum_{t=1}^T (f(x_t) - f(x^*)) \\ &\leq \frac{1}{T} \sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle \end{aligned}$$

□

Next, we analyze the inner product  $\langle \nabla f(x_t), x_t - x^* \rangle$ . It is useful to split it as follows:

$$\begin{aligned} \langle \nabla f(x_t), x_t - x^* \rangle &= \langle \nabla f(x_t), z_t - x^* \rangle + \langle \nabla f(x_t), x_t - z_t \rangle \\ &= \langle \nabla f(x_t), z_t - x^* \rangle + \langle \nabla f(x_{t-1}), x_t - z_t \rangle + \langle \nabla f(x_t) - \nabla f(x_{t-1}), x_t - z_t \rangle \end{aligned}$$

The above split is a good one to consider, since it aligns well with the optimality conditions.

**Exercise 5.** Upper bound  $\langle \nabla f(x_t), z_t - x^* \rangle$  and  $\langle \nabla f(x_{t-1}), x_t - z_t \rangle$  using the optimality condition for  $z_t$  and  $x_t$ .

*Proof.* We have

$$\begin{aligned} z_t = \arg \min_{u \in K} &\left\{ \langle \nabla f(x_t), u \rangle + \frac{1}{2\eta_{t-1}} \|u - z_{t-1}\|^2 + \left( \frac{1}{2\eta_t} - \frac{1}{2\eta_{t-1}} \right) \|u - x_t\|^2 \right\} \\ &\left\langle \nabla f(x_t) + \frac{1}{\eta_{t-1}} (z_t - z_{t-1}) + \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) (z_t - x_t), z_t - x^* \right\rangle \leq 0 \end{aligned}$$

$$\begin{aligned} \langle \nabla f(x_t), z_t - x^* \rangle &\leq \frac{1}{\eta_{t-1}} \langle z_{t-1} - z_t, z_t - x^* \rangle + \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \langle x_t - z_t, z_t - x^* \rangle \\ &= \frac{1}{2\eta_{t-1}} \|z_{t-1} - x^*\|^2 - \frac{1}{2\eta_t} \|z_t - x^*\|^2 + \left( \frac{1}{2\eta_t} - \frac{1}{2\eta_{t-1}} \right) \|x_t - x^*\|^2 \\ &\quad - \frac{1}{2\eta_{t-1}} \|z_{t-1} - z_t\|^2 - \left( \frac{1}{2\eta_t} - \frac{1}{2\eta_{t-1}} \right) \|x_t - z_t\|^2 \end{aligned}$$

We have

$$\begin{aligned} x_t = \arg \min_{u \in K} &\left\{ \langle \nabla f(x_{t-1}), u \rangle + \frac{1}{2\eta_{t-1}} \|u - z_{t-1}\|^2 \right\} \\ &\left\langle \nabla f(x_{t-1}) + \frac{1}{\eta_{t-1}} (x_t - z_{t-1}), x_t - z_t \right\rangle \leq 0 \end{aligned}$$

$$\begin{aligned} \langle \nabla f(x_{t-1}), x_t - z_t \rangle &\leq \frac{1}{\eta_{t-1}} \langle z_{t-1} - x_t, x_t - z_t \rangle \\ &= \frac{1}{2\eta_{t-1}} \left( \|z_{t-1} - z_t\|^2 - \|z_{t-1} - x_t\|^2 - \|x_t - z_t\|^2 \right) \end{aligned}$$

$$\begin{aligned} \langle \nabla f(x_t), x_t - x^* \rangle &= \frac{1}{2\eta_{t-1}} \|z_{t-1} - x^*\|^2 - \frac{1}{2\eta_t} \|z_t - x^*\|^2 + \left( \frac{1}{2\eta_t} - \frac{1}{2\eta_{t-1}} \right) \|x_t - x^*\|^2 \\ &\quad - \frac{1}{2\eta_t} \|x_t - z_t\|^2 - \frac{1}{2\eta_{t-1}} \|z_{t-1} - x_t\|^2 \\ &\quad + \langle \nabla f(x_t) - \nabla f(x_{t-1}), x_t - z_t \rangle \end{aligned}$$

□

Next, we consider the third term  $\langle \nabla f(x_t) - \nabla f(x_{t-1}), x_t - z_t \rangle$ . The intuition is that the two gradients  $\nabla f(x_t)$  and  $\nabla f(x_{t-1})$  are close, and thus their difference  $\|\nabla f(x_t) - \nabla f(x_{t-1})\|$  is small. Assuming this intuition is correct, we can try to show that the inner product is small by upper bounding it in terms of  $\|\nabla f(x_t) - \nabla f(x_{t-1})\|$ .

**Exercise 6.** Show that, if  $K = \mathbb{R}^d$ , we have

$$\langle \nabla f(x_t) - \nabla f(x_{t-1}), x_t - z_t \rangle \leq \eta_t \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2$$

*Hint:* Note that  $x_t$  is also the minimizer of  $\|u - x_t\|^2$ .

*Proof.* Since  $x_t$  is the minimizer of both  $\langle \nabla f(x_{t-1}), u \rangle + \frac{1}{2\eta_{t-1}} \|u - z_{t-1}\|^2$  and  $\left(\frac{1}{2\eta_t} - \frac{1}{2\eta_{t-1}}\right) \|u - x_t\|^2$ , it is the minimizer of their sum:

$$\begin{aligned} x_t &= \arg \min_{u \in \mathbb{R}^d} \left\{ \langle \nabla f(x_{t-1}), u \rangle + \frac{1}{2\eta_{t-1}} \|u - z_{t-1}\|^2 + \left(\frac{1}{2\eta_t} - \frac{1}{2\eta_{t-1}}\right) \|u - x_t\|^2 \right\} \\ &= \frac{\eta_t}{\eta_{t-1}} z_{t-1} + \left(1 - \frac{\eta_t}{\eta_{t-1}}\right) x_t - \eta_t \nabla f(x_{t-1}) \end{aligned}$$

By definition,

$$\begin{aligned} z_t &= \arg \min_{u \in \mathbb{R}^d} \left\{ \langle \nabla f(x_t), u \rangle + \frac{1}{2\eta_{t-1}} \|u - z_{t-1}\|^2 + \left(\frac{1}{2\eta_t} - \frac{1}{2\eta_{t-1}}\right) \|u - x_t\|^2 \right\} \\ &= \frac{\eta_t}{\eta_{t-1}} z_{t-1} + \left(1 - \frac{\eta_t}{\eta_{t-1}}\right) x_t - \eta_t \nabla f(x_t) \end{aligned}$$

Thus

$$\|x_t - z_t\| = \eta_t \|\nabla f(x_t) - \nabla f(x_{t-1})\|$$

□

It turns out that we can prove the above result for general  $K$  as well. You may assume this fact without proof, and proceed with the analysis.

By combining the above inequality with the two inequalities that we obtained via the optimality conditions, we obtain an upper bound on  $\langle \nabla f(x_t), x_t - x^* \rangle$ . Next, as in lecture, we sum up these inequalities and telescope as much as possible.

**Exercise 7.** Show that we have

$$\sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle \leq \underbrace{O(R) \sqrt{\sum_{t=1}^T \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2}}_{\text{loss}} - \underbrace{\sum_{t=1}^T \frac{1}{\eta_{t-1}} \left( \|x_t - z_{t-1}\|^2 + \|z_{t-1} - z_{t-1}\|^2 \right)}_{\text{gain}}$$

*Proof.* We showed in the previous exercise that

$$\begin{aligned} \langle \nabla f(x_t), x_t - x^* \rangle &= \frac{1}{2\eta_{t-1}} \|z_{t-1} - x^*\|^2 - \frac{1}{2\eta_t} \|z_t - x^*\|^2 + \left(\frac{1}{2\eta_t} - \frac{1}{2\eta_{t-1}}\right) \underbrace{\|x_t - x^*\|^2}_{\leq R^2} \\ &\quad + \eta_t \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2 - \frac{1}{2\eta_t} \|x_t - z_t\|^2 - \frac{1}{2\eta_{t-1}} \|z_{t-1} - x_t\|^2 \end{aligned}$$

Summing up and telescoping, we obtain

$$\begin{aligned} \sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle &\leq R^2 \frac{1}{2\eta_T} + \sum_{t=1}^T \eta_t \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2 \\ &\quad - \sum_{t=1}^T \left( \frac{1}{2\eta_t} \|x_t - z_t\|^2 + \frac{1}{2\eta_{t-1}} \|z_{t-1} - x_t\|^2 \right) \end{aligned}$$

Now we recall the definition of the step sizes:

$$\eta_t = \frac{R}{\sqrt{\sum_{s=1}^t \|\nabla f(x_s) - \nabla f(x_{s-1})\|^2}}$$

which gives us

$$\begin{aligned}\frac{R^2}{\eta T} &= R \sqrt{\sum_{t=1}^T \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2} \\ \sum_{t=1}^T \eta_t \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2 &= R \sum_{t=1}^T \frac{\|\nabla f(x_t) - \nabla f(x_{t-1})\|^2}{\sqrt{\sum_{s=1}^t \|\nabla f(x_s) - \nabla f(x_{s-1})\|^2}} \\ &\leq 2R \sqrt{\sum_{t=1}^T \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2}\end{aligned}$$

where once again we used the inequality in the first exercise.

Thus we have

$$\sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle \leq \frac{5}{2} R \sqrt{\sum_{t=1}^T \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2} - \sum_{t=1}^T \left( \frac{1}{2\eta_t} \|x_t - z_t\|^2 + \frac{1}{2\eta_{t-1}} \|z_{t-1} - x_t\|^2 \right)$$

□

Using the above exercise, derive the final convergence rate for non-smooth functions.

**Exercise 8.** Suppose that  $f$  is non-smooth and the gradients are bounded:  $\|\nabla f(x_t)\| \leq G$  for all  $t \in [T]$ . Derive the appropriate convergence guarantee for the algorithm.

*Proof.* We have

$$\begin{aligned}\sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle &\leq \frac{5}{2} R \sqrt{\sum_{t=1}^T \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2} \\ &\leq 5RG\sqrt{T}\end{aligned}$$

We further bound

$$\begin{aligned}\|\nabla f(x_t) - \nabla f(x_{t-1})\|^2 &\leq 2\|\nabla f(x_t)\|^2 + 2\|\nabla f(x_{t-1})\|^2 \\ &\leq 4G^2\end{aligned}$$

Thus

$$f(\bar{x}_T) - f(x^*) \leq O\left(\frac{RG}{\sqrt{T}}\right)$$

□

Next, we consider the setting where  $f$  is smooth. Our starting point is the result from the previous exercise:

$$\sum_{t=1}^T \langle \nabla f(x_t), x_t - x^* \rangle \leq \underbrace{O(R) \sqrt{\sum_{t=1}^T \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2}}_{\text{loss}} - \underbrace{\sum_{t=1}^T \frac{1}{\eta_{t-1}} \left( \|x_t - z_{t-1}\|^2 + \|x_{t-1} - z_{t-1}\|^2 \right)}_{\text{gain}}$$

As in the lecture, the goal is to use the gain to offset the loss, and achieve faster convergence.

**Exercise 9.** Use the smoothness of  $f$  to lower bound the terms appearing in the gain in terms of the terms appearing in the loss.

*Proof.* We have

$$\begin{aligned}\|\nabla f(x_t) - \nabla f(x_{t-1})\|^2 &= \|\nabla f(x_t) - \nabla f(z_{t-1}) + \nabla f(z_{t-1}) - \nabla f(x_{t-1})\|^2 \\ &\leq 2\|\nabla f(x_t) - \nabla f(z_{t-1})\|^2 + 2\|\nabla f(z_{t-1}) - \nabla f(x_{t-1})\|^2\end{aligned}$$



By smoothness, we have

$$\begin{aligned}\|x_t - z_{t-1}\|^2 &\geq \frac{1}{\beta^2} \|\nabla f(x_t) - \nabla f(z_{t-1})\|^2 \\ \|x_{t-1} - z_{t-1}\|^2 &\geq \frac{1}{\beta^2} \|\nabla f(z_{t-1}) - \nabla f(x_{t-1})\|^2\end{aligned}$$

Thus

$$\begin{aligned}\|x_t - z_{t-1}\|^2 + \|x_{t-1} - z_{t-1}\|^2 &\geq \frac{1}{\beta^2} \left( \|\nabla f(x_t) - \nabla f(z_{t-1})\|^2 + \|\nabla f(z_{t-1}) - \nabla f(x_{t-1})\|^2 \right) \\ &\geq \frac{1}{2\beta^2} \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2\end{aligned}$$

□

**Exercise 10.** Show that there is a transition point where the gain offsets the loss, and upper bound the net loss. Derive the appropriate convergence guarantee for the algorithm.

*Proof.* We have shown that the net loss is at most

$$O(R) \sqrt{\sum_{t=1}^T \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2} - \sum_{t=1}^T \frac{1}{\eta_{t-1}} \frac{1}{2\beta^2} \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2$$

Let  $\tau$  be the last iteration with  $\frac{1}{\eta_{t-1}} \leq \beta$ . We have

$$\begin{aligned}&O(R) \sqrt{\sum_{t=1}^T \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2} - \sum_{t=1}^T \frac{1}{\eta_{t-1}} \frac{1}{2\beta^2} \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2 \\ &\leq O(R) \sqrt{\sum_{t=1}^{\tau-1} \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2 + \|\nabla f(x_\tau) - \nabla f(x_{\tau-1})\|^2 + \sum_{t=\tau+1}^T \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2} \\ &\quad - \sum_{t=\tau+1}^T \frac{1}{\eta_{t-1}} \frac{1}{2\beta^2} \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2 \\ &\leq O(R) \underbrace{\sqrt{\sum_{t=1}^{\tau-1} \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2}}_{=\frac{R}{\eta_{\tau-1}} \leq R\beta} + O(R) \underbrace{\|\nabla f(x_\tau) - \nabla f(x_{\tau-1})\|}_{\leq \beta \|x_\tau - x_{\tau-1}\| \leq \beta R} \\ &\quad + O(R) \sqrt{\sum_{t=\tau+1}^T \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2} - \sum_{t=\tau+1}^T \underbrace{\frac{1}{\eta_{t-1}} \frac{1}{2\beta^2}}_{\geq \frac{1}{2\beta}} \|\nabla f(x_t) - \nabla f(x_{t-1})\|^2 \\ &\leq O(\beta R^2) + \max_{z \geq 0} \left\{ O(R)z - \frac{1}{2\beta} z^2 \right\} \\ &= O(\beta R^2)\end{aligned}$$

Thus we have

$$f(\bar{x}_T) - f(x^*) \leq O\left(\frac{R^2\beta}{T}\right)$$

□