

Convex Optimisation

Matthew Tam
(University of Melbourne)
matthew.tam@unimelb.edu.au

ACE Network

S1 2021

* Assignment 4

- due June 3rd @ 5pm

Continued from last time

Recall that

$$D = \left\{ (x_1, \dots, x_m) \in \mathbb{R}^n \times \dots \times \mathbb{R}^n : x_1 = \dots = x_m \right\} \stackrel{\cong}{=} \mathbb{R}^{nm}$$

Last time, we needed to compute

$$\text{prox}_D^* = P_D.$$

Proposition 3.5.4

Let $\mathbf{x} = (x_1, \dots, x_m) \in \mathbb{R}^n \times \dots \times \mathbb{R}^n$. Then

$$\mathbf{P} = (P_1, \dots, P_m) = P_D(\mathbf{x}) \text{ where}$$

$$P_j = \frac{1}{m} \sum_{j=1}^m x_j.$$

Proof.

We prove the result by using Corollary 2.9.6.

We first note that D is nonempty, closed and convex. Next note that $p \in D$ by definition. Finally, let $c = (c, \dots, c) \in D$ and observe

$$\langle x - p, c - p \rangle$$

$$= \sum_{i=1}^m \langle x_i - p, c - p \rangle$$

$$= \sum_{i=1}^m \langle x_i, c - p \rangle - m \langle p, c - p \rangle$$

$$= \sum_{i=1}^m \langle x_i, c - p \rangle - m \left\langle \sum_{j=1}^m x_j, c - p \right\rangle$$

$$= \sum_{i=1}^m \langle x_i, c - p \rangle - \sum_{j=1}^m \langle x_j, c - p \rangle$$

$$= 0,$$

which completes the proof. □

3.6 Stochastic Gradient Descent

Consider the problem

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{where } f(x) := \frac{1}{m} \sum_{i=1}^m f_i(x)$$

when m is large and $f_i: \mathbb{R}^n \rightarrow \mathbb{R}$ are convex and differentiable.

Under these assumptions, f is also convex and differentiable, so we could apply gradient descent (Section 3.2).

That is,

$$x_{k+1} = x_k - \lambda \boxed{\nabla f(x_k)} \stackrel{\text{approximate}}{\sim} x_k - \frac{\lambda}{m} \sum_{i=1}^m \nabla f_i(x_k).$$

If m is large, the sum on the RHS might take too long to compute.

We require the stepsize sequence (λ_n)
to satisfy

$$\sum_{i=0}^{\infty} \lambda_k^2 < +\infty, \text{ and } \left[\lambda_k = \frac{1}{k+1} \right]$$

$$\sum_{i=0}^{\infty} \lambda_n = +\infty.$$

Given an initial point $x_0 \in \mathbb{R}^n$, generate
sequences (i_k) and (x_n) according to

$$\begin{cases} i_k \sim \text{unif}\{1, \dots, m\} \\ x_{n+1} = x_n - \lambda_n \nabla f_{i_n}(x_n) \end{cases} \quad \forall k \in \mathbb{N}. \quad (3.27)$$

Here the notation ' $i_n \sim \text{unif}\{1, \dots, m\}$ '
means in the k th iteration, the integer
 i_k is chosen at random from $\{1, \dots, m\}$.

with outcome having equal probability.

Definition 3.6.1

Let Z be a random variable with a finite number of outcomes z_1, z_2, \dots, z_m occurring with probabilities p_1, p_2, \dots, p_m .

The expectation of Z is defined as

$$\mu := E[Z] := \sum_{i=1}^m p_i z_i.$$

Note that expectation is linear: if Z_1 and Z_2 are random variables and $\lambda \in \mathbb{R}$, then

$$E[Z_1 + \lambda Z_2] = E[Z_1] + \lambda E[Z_2].$$

Now, since $i_k \sim \text{unif}\{1, \dots, m\}$ and x_k doesn't depend on i_k , we have.

$$\mathbb{E}[f_{i_k}(x_k)] = \sum_{i=1}^m \frac{1}{m} \cdot f_i(x_k) = \frac{1}{m} \sum_{i=1}^m f_i(x_k)$$

↓
Solu.

$$= f(x_k)$$

$$\mathbb{E}[f_{i_k}(x^*)] = \frac{1}{m} \sum_{i=1}^m f_i(x^*) = f(x^*)$$

$$\mathbb{E}[\nabla f_{i_k}(x_k)] = \frac{1}{m} \sum_{i=1}^m \nabla f_i(x_k) = \nabla f(x_k)$$

⋮

Theorem 3.6.2

Let $f_i: \mathbb{R}^n \rightarrow \mathbb{R}$ be convex and differentiable for all $i \in \{1, \dots, m\}$ and set $f(x) := \frac{1}{m} \sum_{i=1}^m f_i(x)$.

Let $x^* \in \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} f(x)$. Given $x_0 \in \mathbb{R}^n$,

let (x_k) be the sequence generated

by (3.27). Further suppose $\beta > 0$

such that $\|\nabla f_i(x)\| \leq \sigma$ for all $x \in \mathbb{R}^n$ and all $i \in \{1, \dots, m\}$. Then

$$\boxed{\lambda_j \rightarrow 0 \text{ as } j \rightarrow \infty}$$

$$\begin{aligned} & \mathbb{E}[f(\bar{x}_k)] - f(x^*) \\ & \leq \|x_0 - x^*\|^2 + \sigma^2 \left(\sum_{j=0}^k \lambda_j \right)^2 \\ & \quad \overbrace{\left(\sum_{j=0}^k \lambda_j \right)}^{2 \rightarrow +\infty \text{ as } k \rightarrow \infty} \quad \text{if } k \in \mathbb{N}. \end{aligned}$$

where $\bar{x}_k := \sum_{j=0}^k \left(\frac{\lambda_j}{\sum_{i=0}^k \lambda_i} x_i \right)$.

In particular, $E\{f(\bar{x}_k)\} \rightarrow f(x^*)$

as $k \rightarrow \infty$.

Proof

Since f is convex and differentiable, Proposition 2.2.3(a) implies

$$f(x^*) \geq f(x_k) + \langle \nabla f(x_k), x^* - x_k \rangle.$$

Expand the norm squared gives

$$\begin{aligned} & \|x_{k+1} - x^*\|^2 \\ &= \| (x_k - \lambda_k \nabla f_{i_k}(x_k)) - x^* \|^2 \\ &= \|x_k - x^*\|^2 - 2\lambda_k \langle \nabla f_{i_k}(x_k), x_k - x^* \rangle \\ &\quad + \lambda_k^2 \|\nabla f_{i_k}(x_k)\|^2 \end{aligned}$$

$$\leq \|x_n - x^*\|^2 + 2\lambda_n \langle \nabla f_{in}(x_n), x^* - x_n \rangle \\ + \lambda_n^2 \sigma^2.$$

Taking the expectation of both sides, we obtain

$$\begin{aligned} & \mathbb{E} \left\{ \|x_{n+1} - x^*\|^2 \right\} \\ & \leq \mathbb{E} \left\{ \|x_n - x^*\|^2 \right. \\ & \quad \left. + 2\lambda_n \langle \nabla f_{in}(x_n), x^* - x_n \rangle \right. \\ & \quad \left. + \lambda_n^2 \sigma^2 \right\} \\ & = \|x_n - x^*\|^2 + 2\lambda_n \mathbb{E} \left\{ \langle \nabla f_{in}(x_n), x^* - x_n \rangle \right\} \\ & \quad + \lambda_n^2 \sigma^2 \end{aligned}$$

$$= \|x_n - x^*\|^2 + 2\lambda_n \mathbb{E}[\nabla f(x_n)^T x^* - x_n]$$

$$+ \lambda_n^2 \sigma^2$$

$$= \|x_n - x^*\|^2 + 2\lambda_n \langle \nabla f(x_n), x^* - x_n \rangle$$

$$+ \lambda_n^2 \sigma^2.$$

$$\leq \|x_n - x^*\|^2 + 2\lambda_n (f(x^*) - f(x_n))$$

$$+ \lambda_n^2 \sigma^2.$$

Rearranging this inequality gives

$$2\lambda_n (f(x_n) - f(x^*))$$

$$\leq \|x_n - x^*\|^2 - \mathbb{E}[\|x_{n+1} - x^*\|^2]$$

$$+ \lambda_n^2 \sigma^2.$$

Now, consider x_u as being determined by the random variable i_1, i_2, \dots, i_{n-1} .
 Taking the expectation of both sides gives

$$2\lambda_k (\mathbb{E}[f(x_u)] - f(x^*)) \leq \mathbb{E}[\|x_u - x^*\|^2] - \mathbb{E}[\|x_{u+1} - x^*\|^2] + \lambda_u^{-2}.$$

Summing this inequality for $j=0, \dots, k$ gives

$$2 \sum_{j=0}^k \lambda_j (\mathbb{E}[f(x_j)] - f(x^*))$$

$$\leq \sum_{j=0}^k \left(\mathbb{E}[||x_j - x^*||^2] - \mathbb{E}[||x_{j+1} - x^*||^2] \right) \\ + \sigma^2 \sum_{j=0}^k \lambda_j^2$$

$$= \mathbb{E}[||x_0 - x^*||^2] - \mathbb{E}[||x_{k+1} - x^*||^2] \\ + \sigma^2 \sum_{j=0}^k \lambda_j^2$$

$$\leq ||x_0 - x^*||^2 + \sigma^2 \sum_{j=0}^k \lambda_j^2.$$

Dividing both sides by
 $2 \sum_{j=0}^k \lambda_j$ gives

$$\begin{aligned}
 & \sum_{j=0}^k \left(\frac{\lambda_j}{\sum_{i=0}^k \lambda_i} E[f(x_i)] \right) \xrightarrow{f(\bar{x}^*)} \\
 & \leq \|x_0 - \bar{x}^*\|^2 + \sigma^2 \sum_{j=0}^k \lambda_j^2 \\
 & \quad 2 \sum_{j=0}^k \lambda_j
 \end{aligned}$$

RHS of what we want

Finally, using convexity of f , we deduce that

$$\begin{aligned}
 & E[f(\bar{x}_n)] \\
 & = E \left[f \left(\sum_{j=0}^k \left(\frac{\lambda_j}{\sum_{i=0}^k \lambda_i} \right) x_i \right) \right]
 \end{aligned}$$

$\parallel \bar{x}_n$ by defn.

Save.

$$\leq \mathbb{E} \left[\sum_{j=0}^k \left(\frac{\lambda_j}{\sum_{i=0}^k \lambda_i} \right) f(x_j) \right]$$

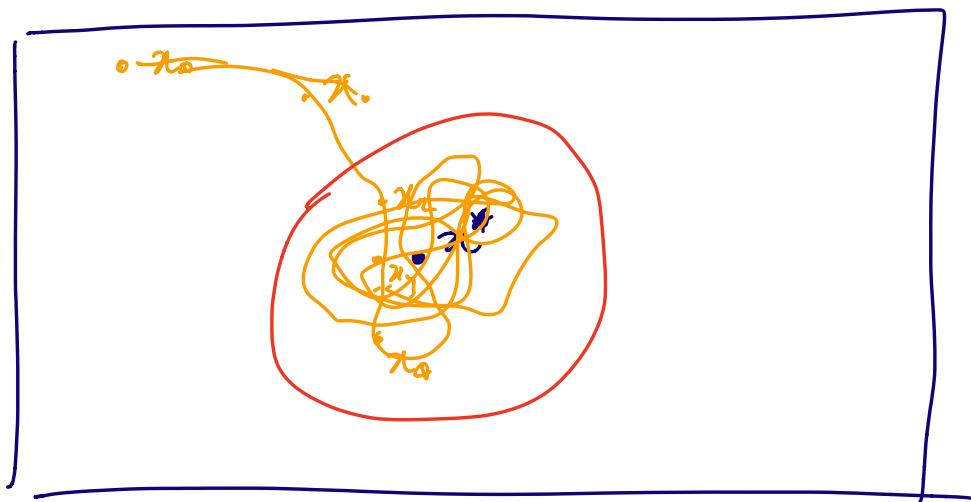
$$= \sum_{j=0}^k \left(\frac{\lambda_j}{\sum_{i=0}^k \lambda_i} \right) \mathbb{E} \{ f(x_j) \}.$$

The result follows by combining
the last two inequalities. \square .

$$\begin{cases} i_k \sim \text{unif}\{1, \dots, n\} \\ x_{k+1} = x_k - \lambda \nabla f_{i_k}(x_k) \end{cases}$$

for some $\lambda > 0$. There exist $\varepsilon > 0$
such that

$$f(x_k) - f(x^*) \leq \varepsilon \quad \forall k \in \mathbb{N}.$$



- variance reduction techniques.