# Advanced Graph Algorithms and Optimization Graded Homework 1

Jonas Hübotter

April 24th, 2022

## 1 Strongly Convex Functions

Let  $f : \mathbb{R}^n \to \mathbb{R}$  be a  $\mu$ -strongly convex and  $\beta$ -smooth<sup>1</sup> function that is twice continuously differentiable<sup>2</sup> and whose first and second order derivatives are integrable.

#### 1.1 Part A

**Lemma 1.** Let  $h : \mathbb{R}^n \to \mathbb{R}$  be a convex function that is continuously differentiable. Then,

$$(\nabla h(\mathbf{x}) - \nabla h(\mathbf{y}))^{\top} (\mathbf{x} - \mathbf{y}) \ge 0.$$
<sup>(1)</sup>

Proof. We have,

$$\begin{aligned} (\nabla h(\boldsymbol{x}) - \nabla h(\boldsymbol{y}))^\top (\boldsymbol{x} - \boldsymbol{y}) \\ &= \nabla h(\boldsymbol{x})^\top \boldsymbol{x} - \nabla h(\boldsymbol{x})^\top \boldsymbol{y} - \nabla h(\boldsymbol{y})^\top \boldsymbol{x} + \nabla h(\boldsymbol{y})^\top \boldsymbol{y} \\ &= -\nabla h(\boldsymbol{x})^\top (\boldsymbol{y} - \boldsymbol{x}) - \nabla h(\boldsymbol{y})^\top (\boldsymbol{x} - \boldsymbol{y}). \end{aligned}$$

Thus, it suffices to show,

$$\nabla h(\mathbf{x})^{\top}(\mathbf{y}-\mathbf{x}) + \nabla h(\mathbf{y})^{\top}(\mathbf{x}-\mathbf{y}) \leq 0.$$

By the first-order characterization of convexity, we have,<sup>3</sup>

$$h(\mathbf{y}) \ge h(\mathbf{x}) + \nabla h(\mathbf{x})^{\top} (\mathbf{y} - \mathbf{x})$$
 and  
 $h(\mathbf{x}) \ge h(\mathbf{y}) + \nabla h(\mathbf{y})^{\top} (\mathbf{x} - \mathbf{y}).$ 

Rearranging terms, we obtain,

$$\nabla h(\mathbf{x})^{\top}(\mathbf{y}-\mathbf{x}) + \nabla h(\mathbf{y})^{\top}(\mathbf{x}-\mathbf{y}) \le h(\mathbf{y}) - h(\mathbf{x}) + h(\mathbf{x}) - h(\mathbf{y}) = 0.$$

#### 1.2 Part B

We first prove the following lemma.<sup>4</sup>

**Lemma 2.** Let  $h : \mathbb{R}^n \to \mathbb{R}$  be a convex function that is continuously differentiable and  $\beta$ -smooth. Then,

$$h(\boldsymbol{y}) \ge h(\boldsymbol{x}) + \boldsymbol{\nabla} h(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \frac{1}{2\beta} \| \boldsymbol{\nabla} h(\boldsymbol{y}) - \boldsymbol{\nabla} h(\boldsymbol{x}) \|_{2}^{2}.$$
(2)

<sup>4</sup> This is analogous to lemma 3.5 in [1], however, they use a different strategy in their proof.

 <sup>1</sup> I say β-smooth to mean β-gradient Lipschitz as I am more used to this wording.
 <sup>2</sup> We use the notion of Frechét differentiability.

<sup>3</sup> theorem 2.3.7

*Proof.* Let  $\phi_x(z) \doteq h(z) - \nabla h(x)^\top z$ . Note  $\nabla \phi_x(z) = \nabla h(z) - \nabla h(x)$ . We have that  $\phi_x$  is convex,<sup>5</sup>

$$\begin{split} \phi_{\mathbf{x}}(z_{1}) + \nabla \phi_{\mathbf{x}}(z_{1})^{\top}(z_{2} - z_{1}) \\ &= h(z_{1}) - \nabla h(\mathbf{x})^{\top} z_{1} + \nabla h(z_{1})^{\top}(z_{2} - z_{1}) + \nabla h(\mathbf{x})^{\top}(z_{1} - z_{2}) \\ &\leq h(z_{2}) - \nabla h(\mathbf{x})^{\top} z_{2} \\ &= \phi_{\mathbf{x}}(z_{2}). \end{split}$$

We also have that  $\phi_x$  is  $\beta$ -smooth,

$$\begin{split} \|\nabla \phi_{x}(z_{1}) - \nabla \phi_{x}(z_{2})\|_{2} &= \|\nabla h(z_{1}) - \nabla h(x) - \nabla h(z_{2}) + \nabla h(x)\|_{2} \\ &= \|\nabla h(z_{1}) - \nabla h(z_{2})\|_{2} \\ &\leq \beta \|z_{1} - z_{2}\|_{2}. \end{split}$$

Thus,6

 $\phi_x(z) \leq \phi_x(y) + \nabla \phi_x(y)^\top (z-y) + rac{eta}{2} \|z-y\|_2^2$ 

and therefore,

$$\min_{z\in\mathbb{R}^n}\phi_x(z)\leq\min_{z\in\mathbb{R}^n}\phi_x(y)+
abla\phi_x(y)^{ op}(z-y)+rac{eta}{2}\,\|z-y\|_2^2\,.$$

We have  $\min_{z \in \mathbb{R}^n} \phi_x(z) = \phi_x(x)$  as  $\nabla \phi_x(x) = 0$  and  $\phi_x$  is convex. In the lecture,<sup>7</sup> we have seen in an analogous argument that the right-hand side is minimized for  $z = y - 1/\beta \nabla \phi_x(y)$ . The inequality simplifies to,

$$h(\mathbf{x}) - \nabla h(\mathbf{x})^{\top} \mathbf{x} = \min_{\mathbf{z} \in \mathbb{R}^n} \phi_{\mathbf{x}}(\mathbf{z})$$
  
$$\leq \min_{\mathbf{z} \in \mathbb{R}^n} \phi_{\mathbf{x}}(\mathbf{y}) + \nabla \phi_{\mathbf{x}}(\mathbf{y})^{\top} (\mathbf{z} - \mathbf{y}) + \frac{\beta}{2} \|\mathbf{z} - \mathbf{y}\|_2^2$$
  
$$= h(\mathbf{y}) - \nabla h(\mathbf{x})^{\top} \mathbf{y} - \frac{1}{2\beta} \|\nabla \phi_{\mathbf{x}}(\mathbf{y})\|_2^2.$$

By reordering the terms, we obtain,

$$h(\boldsymbol{y}) \geq h(\boldsymbol{x}) + \boldsymbol{\nabla} h(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \frac{1}{2\beta} \|\boldsymbol{\nabla} \phi_{\boldsymbol{x}}(\boldsymbol{y})\|_{2}^{2}$$
  
=  $h(\boldsymbol{x}) + \boldsymbol{\nabla} h(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \frac{1}{2\beta} \|\boldsymbol{\nabla} h(\boldsymbol{y}) - \boldsymbol{\nabla} h(\boldsymbol{x})\|_{2}^{2},$ 

as desired.

**Lemma 3.** Let  $h : \mathbb{R}^n \to \mathbb{R}$  be a convex function that is continuously differentiable and  $\beta$ -smooth. Then,

$$(\boldsymbol{\nabla} h(\boldsymbol{x}) - \boldsymbol{\nabla} h(\boldsymbol{y}))^{\top}(\boldsymbol{x} - \boldsymbol{y}) \geq \frac{1}{\beta} \|\boldsymbol{\nabla} h(\boldsymbol{x}) - \boldsymbol{\nabla} h(\boldsymbol{y})\|_{2}^{2}.$$
 (3)

<sup>5</sup> We show the first-order characterization of convexity.

using the first-order characterization of convexity for h

using that *h* is  $\beta$ -smooth

<sup>6</sup> proposition 3.3.3

7 section 3.3.2

Proof. Recall from section 1.1 that

$$(\boldsymbol{\nabla} h(\boldsymbol{x}) - \boldsymbol{\nabla} h(\boldsymbol{y}))^{\top}(\boldsymbol{x} - \boldsymbol{y}) = -\boldsymbol{\nabla} h(\boldsymbol{x})^{\top}(\boldsymbol{y} - \boldsymbol{x}) - \boldsymbol{\nabla} h(\boldsymbol{y})^{\top}(\boldsymbol{x} - \boldsymbol{y}).$$

Using eq. (2), we obtain,

$$\begin{aligned} &- \nabla h(\mathbf{x})^{\top} (\mathbf{y} - \mathbf{x}) - \nabla h(\mathbf{y})^{\top} (\mathbf{x} - \mathbf{y}) \\ &\geq h(\mathbf{x}) - h(\mathbf{y}) + \frac{1}{2\beta} \| \nabla h(\mathbf{y}) - \nabla h(\mathbf{x}) \|_{2}^{2} \\ &+ h(\mathbf{y}) - h(\mathbf{x}) + \frac{1}{2\beta} \| \nabla h(\mathbf{y}) - \nabla h(\mathbf{x}) \|_{2}^{2} \\ &= \frac{1}{\beta} \| \nabla h(\mathbf{y}) - \nabla h(\mathbf{x}) \|_{2}^{2}. \end{aligned}$$

#### 1.3 Part C

**Lemma 4.** Let  $f : \mathbb{R}^n \to \mathbb{R}$  be a twice continuously differentiable,  $\mu$ -strongly convex,  $\beta$ -smooth function. Then,

- (1)  $h(\mathbf{x}) \doteq f(\mathbf{x}) \frac{\mu}{2} ||\mathbf{x}||_2^2$  is a convex and, if  $\beta \neq \mu$ ,  $(\beta \mu)$ -smooth function; and
- (2)  $(\nabla f(\mathbf{x}) \nabla f(\mathbf{y}))^{\top} (\mathbf{x} \mathbf{y})$  $\geq \frac{\mu\beta}{\beta + \mu} \|\mathbf{x} - \mathbf{y}\|_{2}^{2} + \frac{1}{\beta + \mu} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_{2}^{2}.$

*Proof of (1).* Let us compute the Hessian  $H_h$  of h.

$$\begin{split} H_h(\mathbf{x})(i,j) &= \frac{\partial^2}{\partial \mathbf{x}(i) \, \partial \mathbf{x}(j)} h(\mathbf{x}) \\ &= \frac{\partial^2}{\partial \mathbf{x}(i) \, \partial \mathbf{x}(j)} \left( f(\mathbf{x}) - \frac{\mu}{2} \, \|\mathbf{x}\|_2^2 \right) \\ &= H_f(\mathbf{x})(i,j) - \frac{\mu}{2} \underbrace{\frac{\partial^2}{\partial \mathbf{x}(i) \, \partial \mathbf{x}(j)} \, \|\mathbf{x}\|_2^2}_{=2} \\ &= H_f(\mathbf{x})(i,j) - \mu. \end{split}$$

Thus,  $H_h(x) = H_f(x) - \mu I$  for all  $x \in \mathbb{R}^n$ . In particular, if  $\{\lambda_i\}_i$  are the eigenvalues of  $H_f$ , then  $\{\lambda_i - \mu\}_i$  are the eigenvalues of  $H_h$ .<sup>8</sup>

To show that *h* is convex, it suffices to show that  $H_h$  is positive semi-definite and therefore that  $\lambda_{\min}(H_h(x)) \ge 0$  for all  $x \in \mathbb{R}^{n}$ .9 Using that *f* is  $\mu$ -strongly convex, we have for all  $x \in \mathbb{R}^{n}$ ,

$$\lambda_{\min}(H_h(\mathbf{x})) = \underbrace{\lambda_{\min}(H_f(\mathbf{x}))}_{>\mu} - \mu \ge \mu - \mu = 0.$$

To show that *h* is  $(\beta - \mu)$ -smooth, it suffices to show that  $\lambda_{\max}(\mathbf{H}_h(\mathbf{x})) \leq \beta - \mu$  for all  $\mathbf{x} \in \mathbb{R}^{n}$ .<sup>10</sup> Using that *f* is  $\beta$ -smooth, we

<sup>8</sup> Let  $A \in \mathbb{R}^{n \times n}$  and  $c \in \mathbb{R}$ . Then, for any eigenvalue  $\lambda \in \mathbb{R}$  of A and corresponding eigenvector  $x \in \mathbb{R}^n$ ,

$$(A + cI)x = Ax + cIx$$
$$= \lambda x + cx = (\lambda + c)x.$$

Hence,  $\lambda + c$  is the eigenvalue of A + cI corresponding to the eigenvector x. <sup>9</sup> using theorem 3.2.9 and theorem 3.1.2

<sup>10</sup> using proposition 3.3.2

have for all  $x \in \mathbb{R}^n$ ,

$$\lambda_{\max}(\boldsymbol{H}_h(\boldsymbol{x})) = \underbrace{\lambda_{\max}(\boldsymbol{H}_f(\boldsymbol{x}))}_{\leq \beta} - \mu \leq \beta - \mu. \qquad \Box$$

*Proof of (2).* We consider two cases. First, suppose  $\beta = \mu$ . We have,

$$f(\boldsymbol{y}) \geq f(\boldsymbol{x}) + \boldsymbol{\nabla} f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \frac{\beta}{2} \|\boldsymbol{x} - \boldsymbol{y}\|_{2}^{2}$$
  
$$f(\boldsymbol{y}) \geq f(\boldsymbol{x}) + \boldsymbol{\nabla} f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \frac{1}{2\beta} \|\boldsymbol{\nabla} f(\boldsymbol{y}) - \boldsymbol{\nabla} f(\boldsymbol{x})\|_{2}^{2}.$$

We obtain,

$$\begin{split} (\boldsymbol{\nabla} f(\boldsymbol{x}) - \boldsymbol{\nabla} f(\boldsymbol{y}))^{\top} (\boldsymbol{x} - \boldsymbol{y}) \\ &= -\boldsymbol{\nabla} f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) - \boldsymbol{\nabla} f(\boldsymbol{y})^{\top} (\boldsymbol{x} - \boldsymbol{y}) \\ &\geq f(\boldsymbol{x}) - f(\boldsymbol{y}) + \frac{\beta}{2} \|\boldsymbol{x} - \boldsymbol{y}\|_2^2 + f(\boldsymbol{y}) - f(\boldsymbol{x}) + \frac{1}{2\beta} \|\boldsymbol{\nabla} f(\boldsymbol{y}) - \boldsymbol{\nabla} f(\boldsymbol{x})\|_2^2 \\ &= \frac{\beta}{2} \|\boldsymbol{x} - \boldsymbol{y}\|_2^2 + \frac{1}{2\beta} \|\boldsymbol{\nabla} f(\boldsymbol{y}) - \boldsymbol{\nabla} f(\boldsymbol{x})\|_2^2, \end{split}$$

which is what we wanted to show.

Now, suppose  $\beta \neq \mu$ . Let  $h(\mathbf{x}) \doteq f(\mathbf{x}) - \frac{\mu}{2} \|\mathbf{x}\|_2^2$  be defined as in (1). Using our results from (1), *h* is convex and  $(\beta - \mu)$ -smooth. By eq. (3), we have

$$(\boldsymbol{\nabla} h(\boldsymbol{x}) - \boldsymbol{\nabla} h(\boldsymbol{y}))^{\top}(\boldsymbol{x} - \boldsymbol{y}) \geq rac{1}{eta - \mu} \| \boldsymbol{\nabla} h(\boldsymbol{x}) - \boldsymbol{\nabla} h(\boldsymbol{y}) \|_2^2.$$

Note that  $\nabla h(\mathbf{x}) = \nabla f(\mathbf{x}) - \mu \mathbf{x}$ . This gives us,

$$\begin{aligned} (\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^{\top} (\mathbf{x} - \mathbf{y}) \\ &= (\nabla h(\mathbf{x}) - \nabla h(\mathbf{y}))^{\top} (\mathbf{x} - \mathbf{y}) + \underbrace{(\mu \mathbf{x} - \mu \mathbf{y})^{\top} (\mathbf{x} - \mathbf{y})}_{=\mu \|\mathbf{x} - \mathbf{y}\|_{2}^{2}} \\ &\geq \frac{1}{\beta - \mu} \|\nabla h(\mathbf{x}) - \nabla h(\mathbf{y})\|_{2}^{2} + \mu \|\mathbf{x} - \mathbf{y}\|_{2}^{2} \\ &= \frac{1}{\beta - \mu} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) + \mu(\mathbf{y} - \mathbf{x})\|_{2}^{2} + \mu \|\mathbf{x} - \mathbf{y}\|_{2}^{2} \\ &= \frac{1}{\beta - \mu} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_{2}^{2} - \frac{2\mu}{\beta - \mu} (\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^{\top} (\mathbf{x} - \mathbf{y}) \\ &+ \frac{\beta \mu}{\beta - \mu} \|\mathbf{x} - \mathbf{y}\|_{2}^{2}. \end{aligned}$$

expanding the squared norm

Rearranging the terms, we get,

$$\frac{\beta + \mu}{\beta - \mu} (\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^{\top} (\mathbf{x} - \mathbf{y}) \ge \frac{1}{\beta - \mu} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_{2}^{2} + \frac{\beta \mu}{\beta - \mu} \|\mathbf{x} - \mathbf{y}\|_{2}^{2}.$$

using exercise 19 (A) from the first problem set, where f is  $\beta$ -strongly SSWE have shown for  $\beta$ -smooth f in eq. (2) Finally, multiplying both sides by  $\frac{\beta-\mu}{\beta+\mu} > 0$ , we obtain,

$$\begin{aligned} (\boldsymbol{\nabla} f(\boldsymbol{x}) - \boldsymbol{\nabla} f(\boldsymbol{y}))^{\top}(\boldsymbol{x} - \boldsymbol{y}) &\geq \frac{1}{\beta + \mu} \| \boldsymbol{\nabla} f(\boldsymbol{x}) - \boldsymbol{\nabla} f(\boldsymbol{y}) \|_{2}^{2} \\ &+ \frac{\beta \mu}{\beta + \mu} \| \boldsymbol{x} - \boldsymbol{y} \|_{2}^{2}. \quad \Box \end{aligned}$$

**Lemma 5.** When f is  $\beta$ -smooth and  $\mu$ -strongly convex, we always have  $\mu \leq \beta$ .

*Proof.* Suppose for a contradiction that  $\mu > \beta$ . Recall that for any  $x \in \mathbb{R}^n$ ,  $\lambda_{\min}(H_f(x)) \ge \mu$  and  $\lambda_{\max}(H_f(x)) \le \beta$  as f is  $\mu$ -strongly convex and  $\beta$ -smooth. But this yields,

$$\lambda_{\min}(\boldsymbol{H}_f(\boldsymbol{x})) \geq \mu > \beta \geq \lambda_{\max}(\boldsymbol{H}_f(\boldsymbol{x})).$$

1.4 Part D

**Lemma 6.** Let f be defined as in the beginning of this section. When using a version of gradient descent with  $\mathbf{x}_{i+1} \doteq \mathbf{x}_i - \alpha \nabla f(\mathbf{x}_i)$  for some  $\alpha \in \mathbb{R}$ , then

$$\|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2^2 \le \left(1 - \frac{2\alpha\mu\beta}{\mu+\beta}\right) \|\mathbf{x}_i - \mathbf{x}^*\|_2^2 + \alpha\left(\alpha - \frac{2}{\mu+\beta}\right) \|\nabla f(\mathbf{x}_i)\|_2^2,$$
(4)

where  $\mathbf{x}^* \in \operatorname{arg\,min}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ .

Proof. We have,

$$\begin{aligned} \|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2^2 &= \|\mathbf{x}_i - \mathbf{x}^* - \alpha \nabla f(\mathbf{x}_i)\|_2^2 \\ &= \|\mathbf{x}_i - \mathbf{x}^*\|_2^2 - 2\alpha \nabla f(\mathbf{x}_i)^\top (\mathbf{x}_i - \mathbf{x}^*) + \alpha^2 \|\nabla f(\mathbf{x}_i)\|_2^2. \end{aligned}$$

It suffices to show,

$$2\alpha \nabla f(\mathbf{x}_i)^{\top}(\mathbf{x}_i - \mathbf{x}^*) \geq \frac{2\alpha}{\mu + \beta} \|\nabla f(\mathbf{x}_i)\|_2^2 + \frac{2\alpha\mu\beta}{\mu + \beta} \|\mathbf{x}_i - \mathbf{x}^*\|_2^2.$$

Dividing by  $2\alpha$ , observe that this is precisely what we have proven in part (2) of lemma 4 where  $x \doteq x_i$  and  $y \doteq x^*$ .<sup>11</sup>

<sup>11</sup> Note that  $\nabla f(\mathbf{x}^*) = 0$ .

#### 1.5 Part E

**Lemma 7.** In the setting of lemma 6, we have for  $\alpha \doteq 1/\beta$ ,

$$\|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \le \exp\left(-\frac{\mu}{\beta}k\right) \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2.$$
 (5)

*Proof.* Unraveling the recurrence from eq. (4), we get,

$$\|\mathbf{x}_{k}-\mathbf{x}^{*}\|_{2}^{2} \leq \left(1-\frac{2\alpha\mu\beta}{\mu+\beta}\right)^{k}\|\mathbf{x}_{0}-\mathbf{x}^{*}\|_{2}^{2}+\alpha\left(\alpha-\frac{2}{\mu+\beta}\right)\|\nabla f(\mathbf{x}_{i})\|_{2}^{2}.$$

Plugging in  $\alpha \doteq 1/\beta$ , yields,

$$= \left(1 - \frac{2\mu}{\mu + \beta}\right)^{k} \|\mathbf{x}_{0} - \mathbf{x}^{*}\|_{2}^{2} + \frac{1}{\beta} \left(\frac{1}{\beta} - \frac{2}{\mu + \beta}\right) \|\nabla f(\mathbf{x}_{i})\|_{2}^{2}$$

Using  $\mu \leq \beta$ , we have,

$$\frac{2\mu}{\mu+\beta} \geq \frac{\mu}{\beta}$$
 and  $\frac{2}{\mu+\beta} \geq \frac{1}{\beta}$ .

We obtain,

$$\|m{x}_k - m{x}^*\|_2^2 \le \left(1 - rac{\mu}{eta}
ight)^k \|m{x}_0 - m{x}^*\|_2^2 \le \exp\left(-rac{\mu}{eta}k
ight) \|m{x}_0 - m{x}^*\|_2^2.$$

using that 
$$1 + x \le \exp(x)$$
 for all  $x \in \mathbb{R}$ 

1.6 Part F

We will first show a result for the version of gradient descent we have seen in parts D and E. We will then improve on this result using acceleration. The proof of this statement is not required for our improved version.

**Theorem 8.** Let  $f : \mathbb{R}^n \to \mathbb{R}$  be a  $\mu$ -strongly convex and  $\beta$ -smooth function that is twice continuously differentiable. Then, gradient descent with  $\mathbf{x}_{i+1} \doteq \mathbf{x}_i - 1/\beta \nabla f(\mathbf{x}_i)$  yields an approximate solution  $\mathbf{x}_k$  such that for any  $\epsilon > 0$ ,

$$f(\boldsymbol{x}_k) - f(\boldsymbol{x}^*) \le \epsilon$$

where  $\mathbf{x}^* \in \arg\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$  and the gradient of f is computed at at most  $\kappa \log(\beta \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2/2\epsilon)$  points.<sup>12</sup>

*Proof.* First, note that during each iteration of the given scheme, the gradient of f is evaluated at exactly one point. It therefore suffices to bound the number of iterations until we get an  $\epsilon$ -optimal solution.

As *f* is  $\beta$ -smooth, we have,

$$f(\mathbf{x}_k) \leq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^{\top} (\mathbf{x}_k - \mathbf{x}^*) + \frac{\beta}{2} \|\mathbf{x}_k - \mathbf{x}^*\|_2^2$$

Noting that  $\nabla f(\mathbf{x}^*) = 0$  and rearranging the terms, we obtain,

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \le \frac{\beta}{2} \|\mathbf{x}_k - \mathbf{x}^*\|_2^2.$$

Using eq. (5), we get,

$$\leq rac{eta}{2} \exp\left(-rac{k}{\kappa}
ight) \|oldsymbol{x}_0 - oldsymbol{x}^*\|_2^2 \stackrel{!}{\leq} \epsilon.$$

Solving the inequality for *k*, yields,

$$k \geq \kappa \log \left( rac{eta \| m{x}_0 - m{x}^* \|_2^2}{2 \epsilon} 
ight)$$

as desired.

<sup>12</sup>  $\kappa \doteq \beta/\mu$  is the *condition number* of *f*.

We now show that we can improve the previous result using acceleration to only require order  $\sqrt{\kappa}$  rather than order of  $\kappa$  iterations to converge to an  $\epsilon$ -optimal solution.

**Theorem 9.** Let  $f : \mathbb{R}^n \to \mathbb{R}$  be a  $\mu$ -strongly convex and  $\beta$ -smooth function that is twice continuously differentiable. Let  $\mathbf{x}_0 \in \mathbb{R}$  be any initial guess. Then, the iterative scheme,

$$y_0 \doteq x_0 \tag{6}$$

$$\boldsymbol{y}_{i+1} \doteq \boldsymbol{x}_i - \frac{1}{\beta} \boldsymbol{\nabla} f(\boldsymbol{x}_i) \tag{7}$$

$$\boldsymbol{x}_{i+1} \doteq (1+\theta)\boldsymbol{y}_{i+1} - \theta\boldsymbol{y}_i \quad \text{for } \theta \doteq \frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}, \tag{8}$$

*yields an approximate solution*  $y_k$  *such that for any*  $\epsilon > 0$ *,* 

$$f(\boldsymbol{y}_k) - f(\boldsymbol{x}^*) \le \epsilon$$

where  $\mathbf{x}^* \in \arg\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$  and the gradient of f is computed at at most  $\sqrt{\kappa} \log(\beta \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2/\epsilon)$  points.<sup>13</sup>

Note that the sequence  $\{y_i\}_i$  is similar to the gradient descent scheme that we have examined previously. We choose  $x_i$  as a convex combination of the previous and current best guess. Our approach will be to (1) upper bound  $f(y_i)$  by a function  $\phi_i : \mathbb{R}^n \to \mathbb{R}$ , of which we (2) show that  $\phi_i(x)$  converges to f(x) quickly.

We define  $\phi_i$  iteratively,

$$\phi_0(\mathbf{x}) \doteq f(\mathbf{x}_0) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}_0\|_2^2$$
(9)  
$$\phi_{i+1}(\mathbf{x}) \doteq (1 - \gamma)\phi_i(\mathbf{x}) + \gamma \Big(f(\mathbf{x}_i) + \nabla f(\mathbf{x}_i)^\top (\mathbf{x} - \mathbf{x}_i) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}_i\|_2^2\Big),$$
(10)

as the convex combination of itself and a second-order Taylor approximation of *f* at  $x_i$  where we write  $\gamma \doteq 1/\sqrt{\kappa} = \sqrt{\mu/\beta}$  to simplify notation. It is easy to see that  $\phi_i$  is convex.<sup>14</sup> Our analysis rests on the following two claims, which we will prove later.

**Claim 10** (Upper bound).  $f(y_i) \leq \min_{x \in \mathbb{R}^n} \phi_i(x)$ .

**Claim 11** (Fast convergence).  $\phi_i(\mathbf{x}) \leq f(\mathbf{x}) + (1 - \gamma)^i (\phi_0(\mathbf{x}) - f(\mathbf{x})).$ 

*Proof of theorem 9.* By claim 10,  $f(y_i) \le \phi_i(x^*)$  during all iterations *i*. Therefore,

$$\begin{split} f(\boldsymbol{y}_k) - f(\boldsymbol{x}^*) &\leq \phi_k(\boldsymbol{x}^*) - f(\boldsymbol{x}^*) \\ &\leq (1 - \gamma)^k (\phi_0(\boldsymbol{x}^*) - f(\boldsymbol{x}^*)) \\ &= (1 - \gamma)^k (f(\boldsymbol{x}_0) - f(\boldsymbol{x}^*) + \frac{\mu}{2} \|\boldsymbol{x}_0 - \boldsymbol{x}^*\|_2^2). \end{split}$$

using claim 11

using the definition of  $\phi_0$ , eq. (9)

<sup>13</sup> The proof of this theorem is inspired by the lecture on accelerated gradient descent and section 3.7.1 of [1].

<sup>14</sup> We will later show that  $\phi_i$  is  $\mu$ -strongly convex.

$$\begin{aligned} &\leq (1-\gamma)^k \frac{\mu+\beta}{2} \| \boldsymbol{x}_0 - \boldsymbol{x}^* \|_2^2 \\ &\leq (1-\gamma)^k \beta \| \boldsymbol{x}_0 - \boldsymbol{x}^* \|_2^2 \\ &\leq \exp\left(-\frac{k}{\sqrt{\kappa}}\right) \beta \| \boldsymbol{x}_0 - \boldsymbol{x}^* \|_2^2 \stackrel{!}{\leq} \epsilon \end{aligned}$$

Solving the inequality for *k*, yields,

 $k \geq \sqrt{\kappa} \log \left( rac{eta \| oldsymbol{x}_0 - oldsymbol{x}^* \|_2^2}{\epsilon} 
ight)$ 

as desired.

It remains to prove the two claims.

*Proof of claim 11.* We prove the claim by induction on *i*. In the base case, i = 0, we immediately have,

$$f(\mathbf{x}) + (1 - \gamma)^0 (\phi_0(\mathbf{x}) - f(\mathbf{x})) = \phi_0(\mathbf{x}).$$

Let us now consider any fixed  $i \in \mathbb{N}_0$  and suppose that the statement holds for *i*. We have,

$$\begin{split} \phi_{i+1}(\mathbf{x}) &= (1-\gamma)\phi_i(\mathbf{x}) + \gamma \Big( f(\mathbf{x}_i) + \nabla f(\mathbf{x}_i)^\top (\mathbf{x} - \mathbf{x}_i) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}_i\|_2^2 \Big) \\ &\leq (1-\gamma)^{i+1} (\phi_0(\mathbf{x}) - f(\mathbf{x})) + (1-\gamma)f(\mathbf{x}) \\ &+ \gamma \Big( f(\mathbf{x}_i) + \nabla f(\mathbf{x}_i)^\top (\mathbf{x} - \mathbf{x}_i) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}_i\|_2^2 \Big). \end{split}$$

using the definition of  $\phi_{i+1}$ , eq. (10)

using the induction hypothesis

Finally, observe that

$$f(\mathbf{x}) \ge f(\mathbf{x}_i) + \nabla f(\mathbf{x}_i)^\top (\mathbf{x} - \mathbf{x}_i) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}_i\|_2^2$$

as *f* is  $\mu$ -strongly convex. Noting that  $(1 - \gamma)f(\mathbf{x}) + \gamma f(\mathbf{x}) = f(\mathbf{x})$ , completes the proof.

To prove the final claim, we define  $v_i \doteq \arg \min_{x \in \mathbb{R}^n} \phi_i(x)$  and  $\phi_i^* \doteq \min_{x \in \mathbb{R}^n} \phi_i(x)$ .

Claim 12. 
$$\phi_{i+1}^* \ge (1-\gamma)\phi_i^* + (1-\gamma)\nabla f(x_i)^\top (x_i - y_i) + \gamma f(x_i) - \frac{1}{2\beta} \|\nabla f(x_i)\|_2^2$$
.

*Proof of claim 10.* We prove the claim by induction on *i*. In the base case, i = 0, we have,

$$f(y_0) = f(x_0) \le \min_{x \in \mathbb{R}^n} f(x_0) + \underbrace{\frac{\mu}{2} ||x - x_0||_2^2}_{\ge 0} = \min_{x \in \mathbb{R}^n} \phi_0(x),$$

using that  $y_0 = x_0$ .

using  $f(\mathbf{x}) - f(\mathbf{x}^*) \le \frac{\beta}{2} ||\mathbf{x} - \mathbf{x}^*||_2^2$  as f is  $\beta$ -smooth, see the proof of theorem 24 using  $\mu \le \beta$ 

using that  $1 + x \leq \exp(x)$  for all  $x \in \mathbb{R}$ 

Let us now consider any fixed  $i \in \mathbb{N}_0$  and suppose that the statement holds for *i*. By the  $\beta$ -smoothness of *f*, we have,

$$f(y_{i+1}) \leq f(x_i) + \nabla f(x_i)^{\top} (y_{i+1} - x_i) + \frac{\beta}{2} ||y_{i+1} - x_i||_2^2.$$

By the definition of  $y_{i+1}$ , we have  $y_{i+1} - x_i = -\nabla f(x_i)/\beta$  and the inequality simplifies to,

$$\begin{split} f(\boldsymbol{y}_{i+1}) &\leq f(\boldsymbol{x}_i) - \frac{1}{2\beta} \| \boldsymbol{\nabla} f(\boldsymbol{x}_i) \|_2^2 \\ &= (1-\gamma) f(\boldsymbol{y}_i) - (1-\gamma) f(\boldsymbol{y}_i) + f(\boldsymbol{x}_i) - \frac{1}{2\beta} \| \boldsymbol{\nabla} f(\boldsymbol{x}_i) \|_2^2 \\ &\leq (1-\gamma) \phi_i^* - (1-\gamma) f(\boldsymbol{y}_i) + f(\boldsymbol{x}_i) - \frac{1}{2\beta} \| \boldsymbol{\nabla} f(\boldsymbol{x}_i) \|_2^2 \\ &= (1-\gamma) \phi_i^* + (1-\gamma) (f(\boldsymbol{x}_i) - f(\boldsymbol{y}_i)) + \gamma f(\boldsymbol{x}_i) - \frac{1}{2\beta} \| \boldsymbol{\nabla} f(\boldsymbol{x}_i) \|_2^2 \end{split}$$

using the induction hypothesis

Using the first-order characterization of convexity, we have,

$$f(\mathbf{x}_i) - f(\mathbf{y}_i) \leq \nabla f(\mathbf{x}_i)^\top (\mathbf{x}_i - \mathbf{y}_i)$$

Combining the previous two inequalities, yields,

$$\begin{split} f(\boldsymbol{y}_{i+1}) &\leq (1-\gamma)\boldsymbol{\phi}_i^* + (1-\gamma)\boldsymbol{\nabla}f(\boldsymbol{x}_i)^\top (\boldsymbol{x}_i - \boldsymbol{y}_i) + \gamma f(\boldsymbol{x}_i) \\ &\quad - \frac{1}{2\beta} \|\boldsymbol{\nabla}f(\boldsymbol{x}_i)\|_2^2 \end{split}$$

 $f(y_{i+1}) \leq \phi_{i+1}^*$  follows by claim 12.

**Claim 13.** We use the following simple observations for our proof of claim 12.

(1) 
$$\phi_i(\mathbf{x}) = \phi_i^* + \frac{\mu}{2} \|\mathbf{x} - \mathbf{v}_i\|_2^2$$
.  
(2)  $\mathbf{v}_{i+1} = (1 - \gamma)\mathbf{v}_i + \gamma \left(\mathbf{x}_i - \frac{1}{\mu} \nabla f(\mathbf{x}_i)\right)$ .  
(3)  $\mathbf{v}_i - \mathbf{x}_i = \frac{\mathbf{x}_i - \mathbf{y}_i}{\gamma}$ .

Proof of claim 12. We have,

$$\begin{split} \phi_{i+1}^* + \frac{\mu}{2} \| \mathbf{x}_i - \mathbf{v}_{i+1} \|_2^2 &= \phi_{i+1}(\mathbf{x}_i) \\ &= (1 - \gamma)\phi_i(\mathbf{x}_i) + \gamma f(\mathbf{x}_i) \\ &= (1 - \gamma)\phi_i^* + (1 - \gamma)\frac{\mu}{2} \| \mathbf{x}_i - \mathbf{v}_i \|_2^2 + \gamma f(\mathbf{x}_i). \end{split}$$

using claim 13(1)

using the definition of  $\phi_{i+1}$ , eq. (10)

using claim 13(1)

By rearranging the terms, we get,

$$\phi_{i+1}^* = (1-\gamma)\phi_i^* + (1-\gamma)\frac{\mu}{2} \|\mathbf{x}_i - \mathbf{v}_i\|_2^2 + \gamma f(\mathbf{x}_i) - \frac{\mu}{2} \|\mathbf{x}_i - \mathbf{v}_{i+1}\|_2^2.$$

Using claim 13(2), we have,

$$\|\mathbf{x}_{i} - \mathbf{v}_{i+1}\|_{2}^{2} = \|(1 - \gamma)(\mathbf{x}_{i} - \mathbf{v}_{i}) + \frac{\gamma}{\mu} \nabla f(\mathbf{x}_{i})\|_{2}^{2}$$

$$= (1 - \gamma)^2 \|\mathbf{x}_i - \mathbf{v}_i\|_2^2 - 2(1 - \gamma)\frac{\gamma}{\mu} \nabla f(\mathbf{x}_i)^\top (\mathbf{v}_i - \mathbf{x}_i)$$
  
 
$$+ \frac{\gamma^2}{\mu^2} \|\nabla f(\mathbf{x}_i)\|_2^2 .$$

Combining the two equalities, we obtain,

$$\begin{split} \phi_{i+1}^{*} &= (1-\gamma)\phi_{i}^{*} + \underbrace{\gamma(1-\gamma)\frac{\mu}{2} \|\mathbf{x}_{i} - \mathbf{v}_{i}\|_{2}^{2}}_{\geq 0} \\ &+ \gamma f(\mathbf{x}_{i}) + \gamma(1-\gamma)\nabla f(\mathbf{x}_{i})^{\top}(\mathbf{v}_{i} - \mathbf{x}_{i}) - \frac{1}{2\beta} \|\nabla f(\mathbf{x}_{i})\|_{2}^{2} \\ &\geq (1-\gamma)\phi_{i}^{*} + \gamma f(\mathbf{x}_{i}) + (1-\gamma)\nabla f(\mathbf{x}_{i})^{\top}(\mathbf{x}_{i} - \mathbf{y}_{i}) - \frac{1}{2\beta} \|\nabla f(\mathbf{x}_{i})\|_{2}^{2}, \qquad \text{using claim 13(3)} \end{split}$$

as desired.

We finish by giving formal proofs of the statements in claim 13 even though they are similar to proofs we have seen in class and the weekly problem sets.

*Proof of claim* 13(1). We first show by induction on *i* that  $H_{\phi_i}(\mathbf{x}) = \mu \mathbf{I}$  for all  $\mathbf{x} \in \mathbb{R}^n$  and  $i \ge 0$ . By a simple calculation, we have,

$$abla \phi_0(x) = \mu(x - x_0)$$
 and  $H_{\phi_0}(x) = \mu I_{\phi_0}(x)$ 

Let us consider any fixed  $i \in \mathbb{N}_0$  and suppose that the statement holds for *i*. Following from the definition of  $\phi_{i+1}$ , we have,

$$abla \phi_{i+1}(\mathbf{x}) = (1 - \gamma) \nabla \phi_i(\mathbf{x}) + \gamma (\nabla f(\mathbf{x}_i) + \mu(\mathbf{x} - \mathbf{x}_i)) \quad \text{and}$$
  
 $H_{\phi_{i+1}}(\mathbf{x}) = (1 - \gamma) H_{\phi_i}(\mathbf{x}) + \gamma \mu I.$ 

Using the induction hypothesis, we conclude,

$$H_{\phi_{i+1}}(\mathbf{x}) = (1-\gamma)\mu \mathbf{I} + \gamma\mu \mathbf{I} = \mu \mathbf{I}.$$

In particular, this shows that  $\phi_i$  is  $\mu$ -strongly convex.

Note that the highest-order term in  $\phi_i$  must therefore be of order two. It is easy to see that any quadratic function that satisfies  $H_{\phi_{i+1}}(\mathbf{x}) = \mu \mathbf{I}$  and  $\phi_i^* = \min_{\mathbf{x} \in \mathbb{R}^n} \phi_i(\mathbf{x})$ , can be written as<sup>15</sup>

$$\phi_i(x) = rac{\mu}{2} \|x-z\|_2^2 + \phi_i^*$$

for some  $z \in \mathbb{R}^n$ . We immediately see that  $\phi_i$  is minimized by z, and hence,  $z = v_i$ .

Proof of claim 13(2). Recall,

$$\nabla \phi_{i+1}(\mathbf{x}) = (1 - \gamma) \nabla \phi_i(\mathbf{x}) + \gamma (\nabla f(\mathbf{x}_i) + \mu(\mathbf{x} - \mathbf{x}_i)).$$

<sup>15</sup> Consider an arbitrary quadratic function  $g(\mathbf{x}) = \mathbf{x}^{\top} A \mathbf{x} + \mathbf{x}^{\top} \mathbf{b} + c$  with minimum *m* and  $H_g(\mathbf{x}) = A + A^{\top} = \mu I$ . Thus,  $A = \mu/2I$ . Now, consider the function

$$h(\mathbf{x}) = \frac{\mu}{2} \|\mathbf{x} - \mathbf{z}\|_{2}^{2} + m$$
$$= \frac{\mu}{2} \|\mathbf{x}\|_{2}^{2} - \mu \mathbf{x}^{\top} \mathbf{z} + \frac{\mu}{2} \|\mathbf{z}\|_{2}^{2} + m$$

To get  $h \equiv g$ , we simply need to set  $z = -b/\mu$ . As z is the minimizer of both h and g,  $c = \mu/2 ||z||_2^2 - m$  is uniquely determined using that the minimum of h and g is m.

Using claim 13(1), we get,

$$= (1 - \gamma) \nabla \left( \frac{\mu}{2} \| \mathbf{x} - \mathbf{v}_i \|_2^2 + \phi_i^* \right) + \gamma (\nabla f(\mathbf{x}_i) + \mu(\mathbf{x} - \mathbf{x}_i))$$
  
=  $(1 - \gamma) \mu(\mathbf{x} - \mathbf{v}_i) + \gamma (\nabla f(\mathbf{x}_i) + \mu(\mathbf{x} - \mathbf{x}_i))$   
=  $\mu \mathbf{x} - \mu(1 - \gamma) \mathbf{v}_i - \mu \gamma \mathbf{x}_i + \gamma \nabla f(\mathbf{x}_i) \stackrel{!}{=} 0.$ 

Solving the equation for *x*, yields,

$$\mathbf{x} = (1 - \gamma)\mathbf{v}_i + \gamma \mathbf{x}_i - \frac{\gamma}{\mu} \nabla f(\mathbf{x}_i).$$

As  $\phi_{i+1}$  is convex, x minimizes  $\phi_{i+1}$ , and hence,  $v_{i+1} = x$ .<sup>16</sup>

*Proof of claim* 13(3). We prove the statement by induction on *i*. For i = 0, note that the minimizer  $v_0$  of  $\phi_0$  is *x* and hence,

 $v_0 - x_0 = 0 = x_0 - y_0.$ 

Let us consider any fixed  $i \in \mathbb{N}_0$  and suppose that the statement holds for *i*. We have,

$$\begin{aligned} \boldsymbol{v}_{i+1} - \boldsymbol{x}_{i+1} &= (1-\gamma)\boldsymbol{v}_i + \gamma \boldsymbol{x}_i - \frac{1}{\gamma\beta}\boldsymbol{\nabla}f(\boldsymbol{x}_i) - \boldsymbol{x}_{i+1} \\ &= \frac{1}{\gamma}\boldsymbol{x}_i - \left(\frac{1}{\gamma} - 1\right)\boldsymbol{y}_i - \frac{1}{\gamma\beta}\boldsymbol{\nabla}f(\boldsymbol{x}_i) - \boldsymbol{x}_{i+1} \\ &= \frac{1}{\gamma}\boldsymbol{y}_{i+1} - \left(\frac{1}{\gamma} - 1\right)\boldsymbol{y}_i - \boldsymbol{x}_{i+1} \\ &\stackrel{!}{=} \frac{\boldsymbol{x}_{i+1} - \boldsymbol{y}_{i+1}}{\gamma}. \end{aligned}$$

<sup>16</sup> As  $\phi_{i+1}$  is  $\mu$ -strongly convex, it is strictly convex, and therefore x is its unique minimizer.

using that  $x_0 = y_0$ 

using claim 13(2) and the identity  $\gamma/\mu = 1/\gamma\beta$ 

using the induction hypothesis

using the definition of  $y_{i+1}$ ,  $x_i = y_{i+1} + \frac{1}{\beta} \nabla f(x_i)$ 

Solving the equation for  $x_{i+1}$ , we obtain,

$$\mathbf{x}_{i+1} = (1+ heta)\mathbf{y}_{i+1} - \mathbf{ heta}\mathbf{y}_i \quad ext{for } \mathbf{ heta} = rac{1-\gamma}{\gamma+1} = rac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1},$$

which coincides precisely with our choice of  $x_{i+1}$ .

2 A different kind of smoothness

**Definition 14.** A *norm* on  $\mathbb{R}^n$  is a function  $\|\cdot\| : \mathbb{R}^n \to \mathbb{R}$  such that

(1) for every  $a \in \mathbb{R}$  and  $x \in \mathbb{R}^n$ , ||ax|| = |a| ||x||;

- (2) for every  $x, y \in \mathbb{R}^n$ ,  $||x + y|| \le ||x|| + ||y||$ ; and
- (3) for every  $x \in \mathbb{R}^n$ , ||x|| = 0 implies x = 0.

**Definition 15.** Given the norm  $\|\cdot\|$  on  $\mathbb{R}^n$  its *dual norm*  $\|\cdot\|_*$  is defined as,

$$\|x\|_{*} \doteq \sup\{z^{\top}x \mid z \in \mathbb{R}^{n}, \|z\| = 1\}.$$
(11)

2.1 Part A

#### Lemma 16.

(1) The supremum in the definition of the dual norm is obtained.

(2)  $\|\cdot\|_*$  is a norm on  $\mathbb{R}^n$ .

(3)  $x^{\top}y \leq ||x|| ||y||_{*}$ .

(4) 
$$(\|x\|_*)_* \leq \|x\|^{17}$$

*Proof of (1).* Let  $B \doteq \{z \in \mathbb{R}^n \mid ||z|| = 1\} \subseteq \mathbb{R}^n$  be the unit ball and consider the linear functional,

$$f_x: B \to \mathbb{R}, z \mapsto z^\top x.$$

We want to show that im  $f_x$  has a supremum. By the completeness axiom, it is sufficient to show that im  $f_x$  is nonempty and bounded (as im  $f_x \subseteq \mathbb{R}$ ).

Note that im  $f_x \neq \emptyset$  follows from the simple observation that  $B \neq \emptyset$ .<sup>18</sup>

To show that im  $f_x$  is bounded, recall that the unit ball *B* is bounded. Thus, it suffices to show that  $f_x$  is a bounded operator. For any  $z \in \mathbb{R}^n$ , we have,

$$|f_x(z)| = |z^{ op}x| \le ||z||_2 ||x||_2$$
 ,

using the Cauchy-Schwartz inequality. Now, recall that all norms on  $\mathbb{R}^n$  are equivalent.<sup>19</sup> Using this fact, we obtain,

$$\leq \underbrace{C \|x\|_2}_{\text{const.}} \|z\|$$
,

proving that  $f_x$  is a bounded operator and im  $f_x$  is bounded.

*Proof of* (2). We check the three properties of a norm. We fix arbitrary  $a \in \mathbb{R}$  and  $x, y \in \mathbb{R}^n$ .

(1) 
$$||ax||_{*} = \sup_{\|z\|=1} az^{\top}x = \sup_{\|z\|=1} ||a| z^{\top}x = |a| \sup_{\|z\|=1} z^{\top}x = |a| ||x||_{*}$$
  
(2)  $||x+y||_{*} = \sup_{\|z\|=1} z^{\top}(x+y) = \sup_{\|z\|=1} z^{\top}x + z^{\top}y$   
 $\leq \sup_{\|z\|=1} z^{\top}x + \sup_{\|z\|=1} z^{\top}y = ||x||_{*} + ||y||_{*}.$ 

(3) We prove the contrapositive of positive definiteness. Suppose  $x \neq 0$ . Then, using the unit vector  $z \doteq x/||x||$ ,

$$\|x\|_* = \sup_{\|z\|=1} z^\top x \ge \|x\|_2^2 / \|x\| > 0.$$

*Proof of (*3*)*. Taking the unit vector  $z \doteq x/||x||$ , we get,

$$\|oldsymbol{y}\|_* = \sup_{\|oldsymbol{z}\|=1}oldsymbol{z}^ opoldsymbol{y} \geq rac{oldsymbol{x}^ opoldsymbol{y}}{\|oldsymbol{x}\|}.$$

Rearranging the terms, yields the desired result.

<sup>17</sup> The other direction holds too, but is not shown here.

<sup>18</sup> We have that *B* is nonempty, as we have for any  $x \in \mathbb{R}^n \setminus \{\mathbf{0}\}$  that the unit vector  $x/||x|| \in B$ .

<sup>19</sup> In particular, there exists  $C \in \mathbb{R}$  such that for all  $x \in \mathbb{R}^n$ ,  $||x||_2 \le C ||x||$ .

using that  $\sup z^{\top}x = \sup z^{\top}(-x)$  as ||z|| = 1 implies ||-z|| = 1

using that  $\sup a + b \leq \sup a + \sup b$ 

*Proof of (4)*. First, note that the dual norm can be characterized equivalently as,

$$\|x\|_{*} = \sup_{\|z\|=1} z^{\top} x = \sup_{y \neq 0} \frac{y^{\top} x}{\|y\|},$$
 (12)

by taking the unit vector  $z \doteq y/||y||$ . Using this characterization, we obtain,

$$(\|x\|_*)_* = \sup_{z \neq 0} \frac{z^\top x}{\|z\|_*}$$
$$= \sup_{z \neq 0} \frac{z^\top x}{\sup_{y \neq 0} \frac{y^\top z}{\|y\|}}$$
$$= \sup_{z \neq 0} z^\top x \inf_{y \neq 0} \frac{\|y\|}{y^\top z}$$
$$= \sup_{z \neq 0} \inf_{y \neq 0} \|y\| \frac{z^\top x}{y^\top z}$$
$$\leq \inf_{y \neq 0} \|y\| \sup_{z \neq 0} \frac{z^\top x}{y^\top z}.$$

using the max-min inequality

Observe that when *x* and *y* are not linearly dependent, their fraction can be made arbitrarily large, and hence, in this case the supremum is  $\infty$ . If, on the other hand  $y = \alpha x$  for some  $\alpha \in \mathbb{R}^n$ , then the fraction evaluates to  $1/\alpha$ . This observation yields,

$$= lpha \|\mathbf{x}\| \cdot \frac{1}{lpha} = \|\mathbf{x}\|$$
 ,

where we used absolute homogeneity.

2.2 Part B

#### Definition 17.

- (1) Given any positive definite matrix *M*, the *Mahalanobis norm* is defined as  $||x||_M \doteq \sqrt{x^\top M x}$ .
- (2) The *uniform norm* is defined as  $||\mathbf{x}||_{\infty} \doteq \max_{i} |\mathbf{x}(i)|$ .
- (3) The Manhattan norm is defined as  $\|x\|_1 \doteq \sum_i |x(i)|$ .

#### Lemma 18.

(1) 
$$(\|\cdot\|_M)_* = \|\cdot\|_{M^{-1}}$$
.

(2) 
$$(\|\cdot\|_{\infty})_* = \|\cdot\|_1$$
.

*Proof of* (1). As M is positive definite, it can be factorized uniquely<sup>20</sup> into  $M = LL^{\top}$  where L is lower triangular with positive entries on the diagonal. We write  $M^{1/2} \doteq L$ . Also note that as M is positive definite, it is symmetric. We have for any  $x \in \mathbb{R}^n$ ,

<sup>20</sup> by the Cholesky decomposition

$$(\|x\|_M)_* = \sup_{\|z\|_M=1} z^\top x.$$

We substitute  $y \doteq M^{1/2} z$ ,<sup>21</sup>

$$= \sup_{\|y\|_{2}=1} x^{\top} M^{-1/2} y$$
  
= 
$$\sup_{\|y\|_{2}=1} \left( M^{-1/2} x \right)^{\top} y$$
  
$$\leq \sup_{\|y\|_{2}=1} \left\| M^{-1/2} x \right\|_{2} \underbrace{\|y\|_{2}}_{=1} = \left\| M^{-1/2} x \right\|_{2}$$
  
= 
$$\sqrt{\left( M^{-1/2} x \right)^{\top} M^{-1/2} x} = \sqrt{x^{\top} M^{-1} x} = \|x\|_{M^{-1}}$$

Moreover, for  $y \doteq M^{-1/2} x / \left\| M^{-1/2} x \right\|_2$ , we have,<sup>22</sup>

$$(\|\mathbf{x}\|_{M})_{*} \geq \frac{\|\mathbf{M}^{-1/2}\mathbf{x}\|_{2}^{2}}{\|\mathbf{M}^{-1/2}\mathbf{x}\|_{2}} = \|\mathbf{M}^{-1/2}\mathbf{x}\|_{2} = \|\mathbf{x}\|_{M^{-1}}.$$

Hence,  $(||x||_M)_* = ||x||_{M^{-1}}$ .

**Corollary 19.** In particular, the euclidean norm  $\|\cdot\|_2$  is self-dual.

*Proof.* 
$$(\|\cdot\|_2)_* = (\|\cdot\|_I)_* = \|\cdot\|_I = \|\cdot\|_2.$$

Proof of (2). We have,

$$egin{aligned} & \left(\|\cdot\|_{\infty}
ight)_{*} = \sup_{\|m{z}\|_{\infty}=1}m{z}^{ op}m{x} \ & = \sup_{\max_{i}|m{z}(i)|=1}m{z}^{ op}m{x} \end{aligned}$$

Clearly,

$$oldsymbol{z}(i) \doteq egin{cases} 1 & oldsymbol{x}(i) \geq 0 \ -1 & oldsymbol{x}(i) < 0 \end{cases}$$

is a least upper bound for  $z^{\top}x$ . To see this, suppose for a contradiction that there exists a  $y \in \mathbb{R}^n$  such that  $y^{\top}x > z^{\top}x$ . But then, we must have for at least one coordinate *i* that |y(i)| > 1, contradicting  $||y||_{\infty} = 1$ . We obtain,

$$(\|\cdot\|_{\infty})_{*} = z^{\top} x = \sum_{i} |x(i)| = \|x\|_{1}.$$

<sup>21</sup> We have  $z = M^{-1/2}y$  and

$$\begin{split} \|\boldsymbol{z}\|_{\boldsymbol{M}} &= \sqrt{\boldsymbol{z}^{\top} \boldsymbol{M} \boldsymbol{z}} \\ &= \sqrt{\left(\boldsymbol{M}^{-1/2} \boldsymbol{y}\right)^{\top} \boldsymbol{M} \boldsymbol{M}^{-1/2} \boldsymbol{y}} \\ &= \sqrt{\boldsymbol{y}^{\top} \boldsymbol{y}} = \|\boldsymbol{y}\|_2 \,. \end{split}$$

using the Cauchy-Schwartz inequality

<sup>22</sup> Note that y is normalized to unit length, i.e.,  $\|y\|_2 = 1$ .

#### 2.3 Part C

**Definition 20.** Given a norm  $\|\cdot\| : \mathbb{R}^n \to \mathbb{R}$ , the *dual vector map* is a function  $(\cdot)^{\#} : \mathbb{R}^n \to \mathbb{R}^n$  such that  $\mathbf{x}^{\top} (\mathbf{x})^{\#} = \|\mathbf{x}\|$  and  $\|(\mathbf{x})^{\#}\|_{*} = 1$ .

We will often work with the dual vector map with respect to the dual norm of a given norm  $\|\cdot\|$ . We denote this dual vector map by  $(\cdot)^{\#}_{*}$ . Using the aforementioned properties, we have,

(1) 
$$x^{\dagger}(x)_{*}^{\#} = ||x||_{*}$$
 and  
(2)  $\left( \left\| (x)_{*}^{\#} \right\|_{*} \right)_{*} = \left\| (x)_{*}^{\#} \right\| = 1.$ 

#### Lemma 21.

- (1) The dual vector map for  $\|\cdot\|_M$  is  $(x)^{\#} \doteq Mx/\sqrt{x^{\top}Mx}$  and unique.
- (2) A (non-unique) dual vector map for  $\|\cdot\|_1$  is given by,

$$(\mathbf{x})^{\#}(i) \doteq \begin{cases} 1 & \mathbf{x}(i) \ge 0 \\ -1 & \mathbf{x}(i) < 0. \end{cases}$$
 (13)

(3) A (non-unique) dual vector map for  $\|\cdot\|_{\infty}$  is given by,

$$(\mathbf{x})^{\#}(i) \doteq \begin{cases} 1 & i = j \text{ and } \mathbf{x}(i) \ge 0 \\ -1 & i = j \text{ and } \mathbf{x}(i) < 0 \\ 0 & otherwise, \end{cases}$$
(14)

where  $j \in \arg \max_{i} |\mathbf{x}(j)|$  is arbitrary but fixed.

Note that in our analysis of the dual vector map, we exclude the case x = 0, as any unit vector can be chosen as the dual vector to  $0.^{23}$ 

Proof of (1). We have,

$$\boldsymbol{x}^{\top} (\boldsymbol{x})^{\#} = \frac{\boldsymbol{x}^{\top} \boldsymbol{M} \boldsymbol{x}}{\sqrt{\boldsymbol{x}^{\top} \boldsymbol{M} \boldsymbol{x}}} = \sqrt{\boldsymbol{x}^{\top} \boldsymbol{M} \boldsymbol{x}} = \|\boldsymbol{x}\|_{\boldsymbol{M}} \quad \text{and}$$
$$\left\| (\boldsymbol{x})^{\#} \right\|_{\boldsymbol{M}^{-1}} = \frac{\|\boldsymbol{M} \boldsymbol{x}\|_{\boldsymbol{M}^{-1}}}{\sqrt{\boldsymbol{x}^{\top} \boldsymbol{M} \boldsymbol{x}}} = \frac{\sqrt{(\boldsymbol{M} \boldsymbol{x})^{\top} \boldsymbol{M}^{-1} \boldsymbol{M} \boldsymbol{x}}}{\sqrt{\boldsymbol{x}^{\top} \boldsymbol{M} \boldsymbol{x}}} = \frac{\sqrt{\boldsymbol{x}^{\top} \boldsymbol{M} \boldsymbol{x}}}{\sqrt{\boldsymbol{x}^{\top} \boldsymbol{M} \boldsymbol{x}}} = 1$$

It remains to show that this choice of  $(x)^{\#}$  is unique. Consider the special case where  $M \doteq I.^{24}$  As  $\|\cdot\|_2$  is self-dual, we need that

$$\left\| (x)^{\#} \right\|_{2} = (x)^{\#^{\top}} (x)^{\#} \stackrel{!}{=} 1,$$

implying that  $(x)^{\#}$  must have unit length. Then, to satisfy  $x^{\top}(x)^{\#} \stackrel{!}{=} ||x||_2$ , we must have  $(x)^{\#} = x/||x||_2$ , which corresponds uniquely to our choice of the dual vector map for  $||\cdot||_I$ .

Proof of (2). We have,

$$\mathbf{x}^{ op}(\mathbf{x})^{\#} = \sum_{i} \mathbf{x}(i) \cdot (\mathbf{x})^{\#}(i) = \sum_{i} |\mathbf{x}(i)| = \|\mathbf{x}\|_{1}$$
 and

<sup>23</sup> This would immediately imply that there are infinitely many dual vector maps.

as *M* is positive definite, it is also symmetric

### using that $(\left\|\cdot\right\|_*)_* = \left\|\cdot\right\|$

<sup>&</sup>lt;sup>24</sup> As we have seen,  $\|\cdot\|_I = \|\cdot\|_2$ .

$$\left\|\left(\boldsymbol{x}\right)^{\#}\right\|_{\infty} = \max_{i} \left|\left(\boldsymbol{x}\right)^{\#}\left(i\right)\right| = 1$$

Clearly, our choice of  $(\cdot)^{\#}$  is not unique, as if *x* contains zeros, the coordinates of the dual vector map may be either positive or negative.

*Proof of* (3). Observe that, by definition,  $(x)^{\#}$  has only one non-zero coordinate. This coordinate corresponds precisely to the coordinate of *x* with the largest absolute value. We therefore have,

$$egin{aligned} \mathbf{x}^{ op}\left(\mathbf{x}
ight)^{\#} &= \max_{i} |\mathbf{x}(i)| = \|\mathbf{x}\|_{\infty} \quad ext{ and } \ \left\| (\mathbf{x})^{\#} 
ight\|_{1} &= 1. \end{aligned}$$

Again, our choice of  $(x)^{\#}$  is not unique, as when x has multiple coordinates with maximal absolute value, any one of them can be selected by the dual vector map.

#### 2.4 Part D

**Definition 22.** A differentiable function  $f : \mathbb{R}^n \to \mathbb{R}$  is  $\beta$ -smooth with respect to a norm  $\|\cdot\|$  if for all  $x, y \in \mathbb{R}^n$ ,

$$\|\boldsymbol{\nabla}f(\boldsymbol{x}) - \boldsymbol{\nabla}f(\boldsymbol{y})\|_* \le \beta \|\boldsymbol{x} - \boldsymbol{y}\|.$$
<sup>(15)</sup>

**Lemma 23.** Let  $f : \mathbb{R}^n \to \mathbb{R}$  be differentiable and  $\beta$ -smooth with respect to the norm  $\|\cdot\|$ . Then,

$$f(\boldsymbol{y}) \leq f(\boldsymbol{x}) + \boldsymbol{\nabla} f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \frac{\beta}{2} \|\boldsymbol{y} - \boldsymbol{x}\|^2.$$
 (16)

*Proof.* We fix any  $x, y \in \mathbb{R}^n$ . We define  $g(\theta) \doteq f(x_\theta)$  where we let  $x_\theta \doteq x + \theta(y - x)$ . Note that g(1) - g(0) = f(y) - f(x). We have,

$$\begin{split} f(\boldsymbol{y}) &= f(\boldsymbol{x}) + g(1) - g(0) \\ &= f(\boldsymbol{x}) + \int_0^1 \frac{dg(\theta)}{d\theta} d\theta & \text{by the fundamental theorem of calculus} \\ &= f(\boldsymbol{x}) + \int_0^1 \nabla f(\boldsymbol{x}_\theta)^\top (\boldsymbol{y} - \boldsymbol{x}) d\theta & \text{by the chain rule} \\ &= f(\boldsymbol{x}) + \int_0^1 \nabla f(\boldsymbol{x})^\top (\boldsymbol{y} - \boldsymbol{x}) d\theta & \\ &+ \int_0^1 (\nabla f(\boldsymbol{x}_\theta) - \nabla f(\boldsymbol{x}))^\top (\boldsymbol{y} - \boldsymbol{x}) d\theta & \\ &= f(\boldsymbol{x}) + \nabla f(\boldsymbol{x})^\top (\boldsymbol{y} - \boldsymbol{x}) + \int_0^1 (\nabla f(\boldsymbol{x}_\theta) - \nabla f(\boldsymbol{x}))^\top (\boldsymbol{y} - \boldsymbol{x}) d\theta & \end{split}$$

For the integrand, we obtain,

$$(\nabla f(\boldsymbol{x}_{\theta}) - \nabla f(\boldsymbol{x}))^{\top} (\boldsymbol{y} - \boldsymbol{x}) \leq \|\nabla f(\boldsymbol{x}_{\theta}) - \nabla f(\boldsymbol{x})\|_{*} \|\boldsymbol{y} - \boldsymbol{x}\| \qquad \text{using } \boldsymbol{x}^{\top} \boldsymbol{y} \leq \|\boldsymbol{x}\| \|\boldsymbol{y}\|_{*}$$

$$\leq eta \| oldsymbol{x}_ heta - oldsymbol{x} \| \| oldsymbol{y} - oldsymbol{x} \| \ = heta eta \| oldsymbol{y} - oldsymbol{x} \|^2 \,.$$

We get,

$$f(\boldsymbol{y}) \leq f(\boldsymbol{x}) + \boldsymbol{\nabla} f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \int_0^1 \theta \beta \|\boldsymbol{y} - \boldsymbol{x}\|^2 \ d\theta$$
$$= f(\boldsymbol{x}) + \boldsymbol{\nabla} f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \frac{\beta}{2} \|\boldsymbol{y} - \boldsymbol{x}\|^2. \qquad \Box$$

2.5 Part E

**Theorem 24.** Let  $f : \mathbb{R}^n \to \mathbb{R}$  be continuously differentiable, convex, and  $\beta$ -smooth with respect to the norm  $\|\cdot\|$ . Then, gradient descent with

$$x_{i+1} \doteq x_i - \frac{1}{\beta} \| \nabla f(x_i) \|_* (\nabla f(x_i))_*^{\#}$$
 (17)

*yields an approximate solution*  $x_k$  *such that for any*  $\epsilon > 0$ *,* 

 $f(\mathbf{x}_k) - f(\mathbf{x}^*) \le \epsilon$ 

where  $\mathbf{x}^* \in \operatorname{arg\,min}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ ,  $\nabla f$  and  $(\cdot)^{\#}_*$  are evaluated at most  $\mathcal{O}(\beta \mathbb{R}^2/\epsilon)$  times and at most  $\mathcal{O}(n\beta \mathbb{R}^2/\epsilon)$  additional arithmetic operations are used. Here,

$$R \doteq \max_{\substack{\boldsymbol{x} \in \mathbb{R}^n \\ f(\boldsymbol{x}) \le f(\boldsymbol{x}_0)}} \|\boldsymbol{x} - \boldsymbol{x}^*\|.$$
(18)

*Proof.* We will show that  $k = O(\beta R^2/\epsilon)$  is sufficient. Clearly, by the choice of the update rule, each iteration computes the gradient and dual vector only once. As we work with vectors in *n* dimensions, the addition and scalar multiplications take O(n) time per iteration.

By  $\beta$ -smoothness of *f*, we have,

$$f(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_{i}) + \nabla f(\mathbf{x}_{i})^{\top} (\mathbf{x}_{i+1} - \mathbf{x}) + \frac{\beta}{2} \|\mathbf{x}_{i+1} - \mathbf{x}\|$$
  
=  $f(\mathbf{x}_{i}) - \frac{1}{\beta} \|\nabla f(\mathbf{x}_{i})\|_{*} \underbrace{\nabla f(\mathbf{x}_{i})^{\top} (\nabla f(\mathbf{x}_{i}))_{*}^{\#}}_{=\|\nabla f(\mathbf{x}_{i})\|_{*}}$   
+  $\frac{1}{2\beta} \|\nabla f(\mathbf{x}_{i})\|_{*}^{2} \underbrace{\|(\nabla f(\mathbf{x}_{i}))_{*}^{\#}\|^{2}}_{=1}$   
=  $f(\mathbf{x}_{i}) - \frac{1}{2\beta} \|\nabla f(\mathbf{x}_{i})\|_{*}^{2}$ . (19)

The remainder of the proof is analogous to the proof of gradient descent in  $\|\cdot\|_2$  we have seen in the lecture and the exercises. We define gap<sub>i</sub>  $\doteq f(\mathbf{x}_i) - f(\mathbf{x}^*)$ . We have,

$$\operatorname{gap}_{i} = f(\boldsymbol{x}_{i}) - f(\boldsymbol{x}^{*}) \leq \nabla f(\boldsymbol{x}_{i})^{\top} (\boldsymbol{x}_{i} - \boldsymbol{x}^{*})$$

using the first-order characterization of convexity,  $f(\mathbf{x}^*) \ge f(\mathbf{x}_i) + \nabla f(\mathbf{x}_i)^\top (\mathbf{x}^* - \mathbf{x}_i)$ 

$$\leq \|\nabla f(\mathbf{x}_i)\|_* \|\mathbf{x}_i - \mathbf{x}^*\|$$
  
$$\leq R \|\nabla f(\mathbf{x}_i)\|_*, \qquad (20)$$

where we note that for all  $i, f(x_i) \leq f(x_0)$ . We obtain,

$$gap_{i+1} - gap_i = f(\mathbf{x}_{i+1} - \mathbf{x}_i) \le -\frac{1}{2\beta} \|\nabla f(\mathbf{x}_i)\|_*^2 \le -\frac{1}{2\beta} \left(\frac{gap_i}{R}\right)^2.$$
(21)

using eq. (19) and then rearranging eq. (20)

Claim 25.  $f(x_k) - f(x^*) \le \frac{2\beta R^2}{k+1}$ .

Using this claim, solving

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \le \frac{2\beta R^2}{k+1} \le \epsilon$$

for k, yields  $k = \Omega(\beta R^2/\epsilon)$ . Thus, choosing  $k = \mathcal{O}(\beta R^2/\epsilon)$  is sufficient. 

*Proof of claim 25.* We prove  $1/gap_i \ge i+1/2\beta R^2$  analogously to the proof of exercise 15 on the first problem set by an induction on  $1/gap_i$ .<sup>25</sup> In the base case,

$$\begin{split} \operatorname{gap}_{0} &= f(\mathbf{x}_{0} - f(\mathbf{x}^{*}) \leq \nabla f(\mathbf{x}^{*})^{\top} (\mathbf{x}_{0} - \mathbf{x}^{*}) + \frac{\beta}{2} \|\mathbf{x}_{0} - \mathbf{x}^{*}\|^{2} \\ &= \frac{\beta}{2} \|\mathbf{x}_{0} - \mathbf{x}^{*}\|^{2} \leq 2\beta R^{2}, \end{split}$$

from which we obtain  $1/gap_0 \ge 1/2\beta R^2$ . Let us now consider an arbitrary but fixed  $i \in \mathbb{N}_0$  and suppose the statement holds for *i*. Dividing eq. (21) by  $gap_i \cdot gap_{i+1}$ , yields,

$$\frac{1}{\operatorname{gap}_i} - \frac{1}{\operatorname{gap}_{i+1}} \leq -\frac{1}{2\beta R^2} \cdot \frac{\operatorname{gap}_i}{\operatorname{gap}_{i+1}} \leq -\frac{1}{2\beta R^2} \qquad \qquad \text{using } \operatorname{gap}_i \geq \operatorname{gap}_{i+1}$$

Rearranging and using the induction hypothesis, yields,

$$\frac{1}{\operatorname{gap}_{i+1}} \geq \frac{1}{2\beta R^2} + \frac{1}{\operatorname{gap}_i} \geq \frac{i+2}{2\beta R^2}.$$

2.6 Part F

**Lemma 26.** Let  $f : \mathbb{R}^n \to \mathbb{R}$  be differentiable and convex such that for all  $x, y \in \mathbb{R}^n$ ,

$$f(\boldsymbol{y}) \leq f(\boldsymbol{x}) + \boldsymbol{\nabla} f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \frac{\beta}{2} \|\boldsymbol{y} - \boldsymbol{x}\|^2$$

*Then, f is*  $\beta$ *-smooth with respect to the norm*  $\|\cdot\|$ *, i.e.,* 

$$\left\|\boldsymbol{\nabla}f(\boldsymbol{x})-\boldsymbol{\nabla}f(\boldsymbol{y})\right\|_{*}\leq\beta\left\|\boldsymbol{x}-\boldsymbol{y}\right\|.$$

<sup>25</sup> We assume that  $gap_i > 0$  for all *i*, as otherwise our algorithm has already converged to the optimal solution.

using that f is  $\beta$ -smooth and  $\boldsymbol{\nabla} f(\boldsymbol{x}^*) = 0$ 

*Proof.* We adopt a similar approach to our proof of lemma 2. Let  $\phi_x(z) \doteq f(z) - (f(x) + \nabla f(x)(z-x)^\top)$ . Note that this yields,  $\nabla \phi_x(z) = \nabla f(z) - \nabla f(x)$ . We have that  $\phi_x$  is convex,

$$\begin{split} \phi_{x}(z_{1}) + \nabla \phi_{x}(z_{1})^{\top}(z_{2} - z_{1}) \\ &= f(z_{1}) - f(x) - \nabla f(x)^{\top}(z_{1} - x) \\ &+ \nabla f(z_{1})^{\top}(z_{2} - z_{1}) - \nabla f(x)^{\top}(z_{2} - z_{1}) \\ &\leq f(z_{2}) - f(x) - \nabla f(x)^{\top}(z_{2} - x) \\ &= \phi_{x}(z). \end{split}$$

using the first-order characterization of convexity for f

Using the  $\beta$ -smoothness of f, we have for any  $y \in \mathbb{R}^n$ ,

$$egin{aligned} \phi_{m{x}}(m{z}) &= f(m{x}) - \left(f(m{x}) + m{
abla} f(m{x})(m{z} - m{x})^{ op}
ight) \ &\leq f(m{y}) + m{
abla} f(m{y})^{ op}(m{z} - m{y}) + rac{m{eta}}{2} \, \|m{z} - m{y}\|^2 \ &- \left(f(m{x}) + m{
abla} f(m{x})(m{z} - m{x})^{ op}
ight). \end{aligned}$$

Rearranging to group terms that depend on *z*, we obtain,

$$= f(\boldsymbol{y}) - \left(f(\boldsymbol{x}) + \boldsymbol{\nabla} f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x})\right) \\ + \left(\boldsymbol{\nabla} f(\boldsymbol{y}) - \boldsymbol{\nabla} f(\boldsymbol{x})\right)^{\top} (\boldsymbol{z} - \boldsymbol{y}) + \frac{\beta}{2} \|\boldsymbol{z} - \boldsymbol{y}\|^{2}.$$

As  $\nabla \phi_x(x) = 0$  and  $\phi_x$  is convex,  $\min_{z \in \mathbb{R}^n} \phi_x(z) = \phi_x(x) = 0$ . Taking the minimum of both sides of the previous inequality, we get,

$$\begin{split} 0 &= \min_{z \in \mathbb{R}^n} \phi_x(z) \\ &\leq f(y) - \left( f(x) + \nabla f(x)^\top (y - x) \right) \\ &+ \min_{z \in \mathbb{R}^n} (\nabla f(y) - \nabla f(x))^\top (z - y) + \frac{\beta}{2} \|z - y\|^2 \\ &= f(y) - \left( f(x) + \nabla f(x)^\top (y - x) \right) \\ &+ \min_{\delta \in \mathbb{R}^n} (\nabla f(y) - \nabla f(x))^\top \delta + \frac{\beta}{2} \|\delta\|^2 \,. \end{split}$$

**Claim 27.** For any  $z \in \mathbb{R}^n$ , we have  $\min_{\delta \in \mathbb{R}^n} z^{\top} \delta + \frac{\beta}{2} \|\delta\|^2 = -\frac{1}{2\beta} \|z\|_*^2$ . Using this claim, rearranging the terms of the previous inequality gives,

$$f(\boldsymbol{y}) \ge f(\boldsymbol{x}) + \boldsymbol{\nabla} f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \frac{1}{2\beta} \| \boldsymbol{\nabla} f(\boldsymbol{x}) - \boldsymbol{\nabla} f(\boldsymbol{y}) \|_{*}^{2}.$$
(22)

Now, recall from section 1.1 that

$$(\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^{\top}(\mathbf{x} - \mathbf{y}) = -\nabla f(\mathbf{x})^{\top}(\mathbf{y} - \mathbf{x}) - \nabla f(\mathbf{y})^{\top}(\mathbf{x} - \mathbf{y}).$$

using  $x^{\top}y \leq ||x|| ||y||_*$ 

Using eq. (22), we obtain,

$$\begin{split} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_* \|\mathbf{x} - \mathbf{y}\| \\ &\geq (\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^\top (\mathbf{x} - \mathbf{y}) \\ &= -\nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) - \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \\ &\geq f(\mathbf{x}) - f(\mathbf{y}) + \frac{1}{2\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_*^2 \\ &\quad + f(\mathbf{y}) - f(\mathbf{x}) + \frac{1}{2\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_*^2 \\ &= \frac{1}{\beta} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_*^2. \end{split}$$

Rearranging gives the desired inequality.

*Proof of claim* 27. We will prove both directions separately. To see that  $\min_{\delta \in \mathbb{R}^n} z^{\top} \delta + \frac{\beta}{2} \|\delta\|^2 \leq -\frac{1}{2\beta} \|z\|_*^2$ , we choose  $\delta \doteq -\frac{1}{\beta} \|z\|_* (z)_*^{\#}$ , and obtain,<sup>26</sup>

$$\begin{split} z^{\top}\delta &+ \frac{\beta}{2} \, \|\delta\|^2 = -\frac{1}{\beta} \, \|z\|_* \underbrace{z^{\top} \, (z)_*^{\#}}_{= \|z\|_*} + \frac{1}{2\beta} \, \|z\|_*^2 \underbrace{\left\| (z)_*^{\#} \right\|^2}_{= 1} \\ &= -\frac{1}{2\beta} \, \|z\|_*^2 \, . \end{split}$$

To see that  $\min_{\delta \in \mathbb{R}^n} z^{\top} \delta + \frac{\beta}{2} \|\delta\|^2 \ge -\frac{1}{2\beta} \|z\|_*^2$ , we bound,

$$egin{aligned} &z^ op \delta + rac{eta}{2} \left\| \delta 
ight\|^2 = -(-z)^ op \delta + rac{eta}{2} \left\| \delta 
ight\|^2 \ &\geq - \left\| z 
ight\|_* \left\| \delta 
ight\| + rac{eta}{2} \left\| \delta 
ight\|^2 \ &\geq \min_{\Delta \in \mathbb{R}} - \left\| z 
ight\|_* \Delta + rac{eta}{2} \Delta^2 \ &= \Phi_z(\Delta) \end{aligned}$$

Clearly,  $\Phi_z$  is a quadratic with positive curvature, and hence, convex.

 $rac{d\Phi_{z}(\Delta)}{d\Delta} = - \left\|z\right\|_{*} + eta\Delta \stackrel{!}{=} 0,$ 

 $egin{aligned} z^ op \delta + rac{eta}{2} \left\| \delta 
ight\|^2 \geq -rac{1}{2eta} \left\| z 
ight\|^2_*. \end{aligned}$ 

is solved for  $\Delta = 1/\beta ||z||_*$ , which therefore is a minimizer of  $\Phi_z$ . Substituting for this minimizer in our previous inequality, we obtain, <sup>26</sup> This is similar to our choice of the update rule of gradient descent from the previous section.

using  $x^{\top}y \leq \|x\| \|y\|_*$ 

choosing  $\Delta \doteq \|\delta\|$  and minimizing

2.7 Part G

We have that

**Lemma 28.** Let  $f : \mathbb{R}^n \to \mathbb{R}$  be a twice continuously differentiable function such that for all  $x, y \in \mathbb{R}^n$ ,

$$0 \le \boldsymbol{y}^{\top} \boldsymbol{H}_{f}(\boldsymbol{x}) \boldsymbol{y} \le \beta \left\| \boldsymbol{y} \right\|^{2}.$$
(23)

*Then, f is convex and*  $\beta$ *-smooth with respect to the norm*  $\|\cdot\|$ *.* 

*Proof.* To show that f is convex, it suffices that  $H_f$  is positive semidefinite.<sup>27</sup> This corresponds precisely to the condition that for all  $x, y \in \mathbb{R}^n, y^\top H_f(x)y \ge 0$ . Thus, it only remains to show that f is also  $\beta$ -smooth.

Similarly to our proof of lemma 23, we employ the fundamental theorem of calculus. We fix arbitrary  $x, y \in \mathbb{R}^n$  and let  $g(\theta) \doteq f(x_{\theta})$  for  $x_{\theta} \doteq x + \theta(y - x)$ . Analogously to the mentioned proof, we have,

$$\begin{split} f(\boldsymbol{y}) &= f(\boldsymbol{x}) + g(1) - g(0) \\ &= f(\boldsymbol{x}) + \int_0^1 \frac{dg(\theta)}{d\theta} d\theta \\ &= f(\boldsymbol{x}) + \int_0^1 \boldsymbol{\nabla} f(\boldsymbol{x}_\theta)^\top (\boldsymbol{y} - \boldsymbol{x}) d\theta \\ &= f(\boldsymbol{x}) + \int_0^1 \boldsymbol{\nabla} f(\boldsymbol{x})^\top (\boldsymbol{y} - \boldsymbol{x}) d\theta \\ &\quad + \int_0^1 (\boldsymbol{\nabla} f(\boldsymbol{x}_\theta) - \boldsymbol{\nabla} f(\boldsymbol{x}))^\top (\boldsymbol{y} - \boldsymbol{x}) d\theta \\ &= f(\boldsymbol{x}) + \boldsymbol{\nabla} f(\boldsymbol{x})^\top (\boldsymbol{y} - \boldsymbol{x}) + \int_0^1 (\boldsymbol{\nabla} f(\boldsymbol{x}_\theta) - \boldsymbol{\nabla} f(\boldsymbol{x}))^\top (\boldsymbol{y} - \boldsymbol{x}) d\theta \end{split}$$

 $^{\scriptscriptstyle 27}$  using theorem 3.2.9

by the fundamental theorem of calculus

by the chain rule

Now, let us shift our attention to bounding the integrand. We define  $h(\tau) \doteq \nabla f(\mathbf{x}_{\tau})^{\top} (\mathbf{y} - \mathbf{x})$  where we let  $\mathbf{x}_{\tau} \doteq \mathbf{x} + \tau(\mathbf{x}_{\theta} - \mathbf{x})$ . Note that

$$(\boldsymbol{\nabla} f(\boldsymbol{x}_{\theta}) - \boldsymbol{\nabla} f(\boldsymbol{x}))^{\top} (\boldsymbol{y} - \boldsymbol{x}) = h(1) - h(0).$$

By the chain rule,

$$egin{aligned} & rac{dh( au)}{d au} = (oldsymbol{x}_ heta - oldsymbol{x})^ op oldsymbol{H}_f(oldsymbol{x}_ au)(oldsymbol{y} - oldsymbol{x}) \ &= heta(oldsymbol{y} - oldsymbol{x})^ op oldsymbol{H}_f(oldsymbol{x}_ au)(oldsymbol{y} - oldsymbol{x}). \end{aligned}$$

We obtain the bound,

$$h(1) - h(0) = \int_0^1 \frac{dh(\tau)}{d\tau} d\tau$$
  
=  $\int_0^1 \theta \underbrace{(\boldsymbol{y} - \boldsymbol{x})^\top \boldsymbol{H}_f(\boldsymbol{x}_\tau)(\boldsymbol{y} - \boldsymbol{x})}_{\leq \beta \| \boldsymbol{y} - \boldsymbol{x} \|^2} d\tau$   
=  $\int_0^1 \theta \beta \| \boldsymbol{y} - \boldsymbol{x} \|^2 d\tau$   
=  $\theta \beta \| \boldsymbol{y} - \boldsymbol{x} \|^2$ .

by the fundamental theorem of calculus

using the assumption

Substituting this bound for the integrand, we obtain,

$$f(\boldsymbol{y}) \leq f(\boldsymbol{x}) + \boldsymbol{\nabla} f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \int_0^1 \theta \beta \|\boldsymbol{y} - \boldsymbol{x}\|^2 \ d\theta$$
$$= f(\boldsymbol{x}) + \boldsymbol{\nabla} f(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \frac{\beta}{2} \|\boldsymbol{y} - \boldsymbol{x}\|^2.$$

Using lemma 26, we conclude that *f* is indeed  $\beta$ -smooth.

#### 2.8 Part H

We will consider the function,

$$m: \mathbb{R}^n \to R, \quad \mathbf{x} \mapsto \frac{1}{\lambda} \log\left(\sum_i \exp(\lambda \mathbf{x}(i))\right),$$
 (24)

for some  $\lambda > 0$ . We will see that *m* is a well-behaved approximation to a slight variation of the uniform norm.

#### Lemma 29.

- (1)  $\max_i \mathbf{x}(i) \le m(\mathbf{x}) \le \frac{\log n}{\lambda} + \max_i \mathbf{x}(i).$
- (2) *m* is convex and  $\lambda$ -smooth with respect to  $\|\cdot\|_{\infty}$ .

*Proof of (1).* Fix any  $x \in \mathbb{R}^n$ . We have,

$$m(\mathbf{x}) = \frac{1}{\lambda} \log \left( \sum_{i} \exp(\lambda \mathbf{x}(i)) \right)$$
  
$$\leq \frac{1}{\lambda} \log \left( n \exp(\lambda \max_{i} \mathbf{x}(i)) \right)$$
  
$$= \frac{1}{\lambda} \left( \log n + \lambda \max_{i} \mathbf{x}(i) \right)$$
  
$$= \frac{\log n}{\lambda} + \max_{i} \mathbf{x}(i).$$

For the other direction,

$$m(\mathbf{x}) = \frac{1}{\lambda} \log \left( \sum_{i} \exp(\lambda \mathbf{x}(i)) \right)$$
  

$$\geq \frac{1}{\lambda} \log \left( \exp(\lambda \max_{i} \mathbf{x}(i)) \right)$$
  

$$= \max_{i} \mathbf{x}(i).$$

*Proof of* (2). First, we show that *m* is convex. To begin with, recall Hölder's inequality,

$$\sum_{i} |\mathbf{x}(i)\mathbf{y}(i)| \le \left(\sum_{i} |\mathbf{x}(i)|^{p}\right)^{\frac{1}{p}} \left(\sum_{i} |\mathbf{y}(i)|^{q}\right)^{\frac{1}{q}},\tag{25}$$

for any  $x, y \in \mathbb{R}^n$  and  $\frac{1}{p} + \frac{1}{q} = 1$ . Fix any  $\theta \in [0, 1]$ . Then,

$$\begin{split} m(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) &= \frac{1}{\lambda} \log \left( \sum_{i} e^{\lambda(\theta \mathbf{x} + (1 - \theta) \mathbf{y})} \right) \\ &= \frac{1}{\lambda} \log \left( \sum_{i} e^{\theta \lambda \mathbf{x}} e^{(1 - \theta)\lambda \mathbf{y}(i)} \right) \\ &\leq \frac{1}{\lambda} \log \left( \left( \sum_{i} e^{\lambda \mathbf{x}(i)} \right)^{\theta} \left( \sum_{i} e^{\lambda \mathbf{y}(i)} \right)^{1 - \theta} \right) \end{split}$$

using Hölder's inequality with  $1/p \doteq \theta$ and  $1/q = 1 - \theta$ 

$$=\theta \frac{1}{\lambda} \log \left( \sum_{i} e^{\lambda \mathbf{x}(i)} \right) + (1-\theta) \frac{1}{\lambda} \log \left( \sum_{i} e^{\lambda \mathbf{y}(i)} \right)$$
$$=\theta m(\mathbf{x}) + (1-\theta) m(\mathbf{y}).$$

To prove smoothness of *m*, we first compute its Hessian and then apply lemma 28. For the Hessian of *m*, we have for any fixed  $x \in \mathbb{R}^n$ ,

$$H_m(\mathbf{x})(i,j) = \frac{\partial^2}{\partial \mathbf{x}(i) \partial \mathbf{x}(j)} m(\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}(i)} \left( \frac{1}{\lambda} \frac{\partial}{\partial \mathbf{x}(j)} \log\left(\sum_k e^{\lambda \mathbf{x}(k)}\right) \right).$$

We compute,

$$\frac{\partial}{\partial \mathbf{x}(j)} \log \left( \sum_{k} e^{\lambda \mathbf{x}(k)} \right) = \frac{\frac{\partial}{\partial \mathbf{x}(j)} \sum_{k} e^{\lambda \mathbf{x}(k)}}{\sum_{k} e^{\lambda \mathbf{x}(k)}} = \frac{\lambda e^{\lambda \mathbf{x}(j)}}{\sum_{k} e^{\lambda \mathbf{x}(k)}}.$$

We write  $\Phi \doteq \sum_{k} e^{\lambda x(k)}$  and  $\Phi_{-i} \doteq \Phi - e^{\lambda x(i)}$ . Then,

$$\begin{split} \boldsymbol{H}_{m}(\boldsymbol{x})(i,j) &= \frac{\partial}{\partial \boldsymbol{x}(i)} \frac{e^{\lambda \boldsymbol{x}(j)}}{\Phi} \\ &= \frac{\left(\frac{\partial}{\partial \boldsymbol{x}(i)}e^{\lambda \boldsymbol{x}(j)}\right) \cdot \Phi - e^{\lambda \boldsymbol{x}(j)} \cdot \frac{\partial}{\partial \boldsymbol{x}(i)}\Phi}{\Phi^{2}} \\ &= \frac{1}{\Phi^{2}} \begin{cases} \lambda e^{\lambda \boldsymbol{x}(i)}\Phi - \lambda e^{2\lambda \boldsymbol{x}(i)} & i = j \\ -\lambda e^{\lambda(\boldsymbol{x}(i) + \boldsymbol{x}(j))} & i \neq j \end{cases} \\ &= \frac{1}{\Phi^{2}} \begin{cases} \lambda e^{\lambda \boldsymbol{x}(i)}\Phi_{-i} & i = j \\ -\lambda e^{\lambda(\boldsymbol{x}(i) + \boldsymbol{x}(j))} & i \neq j. \end{cases} \end{split}$$

Fixing any  $y \in \mathbb{R}^n$ , we have,

$$y^{\top} H_m(\mathbf{x}) y = \sum_{i,j} H_m(\mathbf{x})(i,j) \cdot y(i) \cdot y(j)$$
  
$$\leq \lambda \sum_i y(i)^2 e^{\lambda \mathbf{x}(i)} \cdot \frac{\Phi_{-i}}{\Phi^2}$$
  
$$\leq \lambda \|y\|_{\infty}^2 \frac{1}{\Phi} \underbrace{\sum_i e^{\lambda \mathbf{x}(i)}}_{=\Phi}$$
  
$$= \lambda \|y\|_{\infty}^2.$$

#### 2.9 Part I

We consider the flow problem on a weighted and undirected graph G = (V, E, c) with incidence matrix *B* and  $\mathbf{U} \doteq \text{diag}_{e \in E} \mathbf{c}(e)$  for capacities  $\mathbf{c} \in \mathbb{R}_{\geq 0}^{|E|}$ :

$$\min_{\substack{f \in \mathbb{R}^{|E|} \\ B^{\top}f = b}} \left\| U^{-1} f \right\|_{\infty}$$

using the chain rule in each step

using the quotient rule

using that the off-diagonal entries of the Hessian are negative

using  $\frac{\Phi_{-1}}{\Phi} < 1$  and  $y(i) \leq \|y\|_{\infty}$ 

for demands  $\boldsymbol{b} \in \mathbb{R}^{|V|}$ . The flow problem can be characterized equivalently as,

$$\min_{\substack{\boldsymbol{d}\in\mathbb{R}^{|E|}\\\boldsymbol{B}^{\top}\boldsymbol{d}=\boldsymbol{0}}}\left\|\boldsymbol{U}^{-1}(\tilde{f}_{0}+\boldsymbol{d})\right\|_{\infty}$$

for any feasible flow  $\tilde{f}_0$ , i.e.,  $B^{\top} \tilde{f}_0 = b$ . We can also characterize the problem as,

$$\min_{oldsymbol{x}\in\mathbb{R}^{|E|}}\|f_0+oldsymbol{P}oldsymbol{x}\|_\infty$$

where  $\hat{P} \in \mathbb{R}^{|E| \times |E|}$  is a projection matrix such that for all  $x \in \mathbb{R}^{|E|}$  we have  $B^{\top}\hat{P}x = 0$  and for every circulation<sup>28</sup> *d* there exists an  $x \in \mathbb{R}^{|E|}$  so that  $\hat{P}x = d$ . We let  $P \doteq U^{-1}\hat{P}U$  and  $f_0 \doteq (I - P)U^{-1}\tilde{f}_0$ . We write,

$$OPT \doteq \min_{\boldsymbol{x} \in \mathbb{R}^{|E|}} \|f_0 + \boldsymbol{P}\boldsymbol{x}\|_{\infty}.$$

Because the uniform norm is not smooth, we will use a smooth approximation similar to the function m we have seen in the previous section to approximately solve the optimization problem using gradient descent. As a smooth approximation, we use,

$$s: \mathbb{R}^n \to R, \quad x \mapsto \frac{1}{\lambda} \log \left( \frac{\sum_{e \in E} \exp(\lambda x(e)) + \exp(-\lambda x(e))}{2|E|} \right),$$

for some  $\lambda > 0$ , which is convex,  $\mathcal{O}(\lambda)$ -smooth with respect to  $\|\cdot\|_{\infty}$ , and satisfies,

$$\|\mathbf{x}\|_{\infty} \le s(\mathbf{x}) \le 2\frac{\log|E|}{\lambda} + \|\mathbf{x}\|_{\infty}.$$
 (26)

We will therefore optimize the function  $g(x) \doteq s(f_0 + Px)$ . Note that, as *g* is the composition of two convex functions, it is also convex.

**Lemma 30.** g is  $\mathcal{O}(\lambda \|P\|_{\infty \to \infty}^2)$ -smooth with respect to  $\|\cdot\|_{\infty}^2$ .

*Proof.* By  $\mathcal{O}(\lambda)$ -smoothness of *s*, we have for any  $x, y \in \mathbb{R}^n$ ,

$$s(\boldsymbol{y}) \leq s(\boldsymbol{x}) + \boldsymbol{\nabla} s(\boldsymbol{x})^{\top} (\boldsymbol{y} - \boldsymbol{x}) + \frac{\mathcal{O}(\lambda)}{2} \|\boldsymbol{y} - \boldsymbol{x}\|_{\infty}^{2}$$

Let us now fix any  $x', y' \in \mathbb{R}^n$ . We substitute  $x \doteq f_0 + Px'$  and  $y \doteq f_0 + Py'$ . Note that by the chain rule,  $\nabla g(x') = P^\top \nabla s(f_0 + Px')$ .<sup>30</sup> We obtain,

$$\begin{split} g(\boldsymbol{y}') &= s(f_0 + \boldsymbol{P} \boldsymbol{y}') \\ &\leq s(f_0 + \boldsymbol{P} \boldsymbol{x}') + \boldsymbol{\nabla} s(f_0 + \boldsymbol{P} \boldsymbol{x}')^\top \boldsymbol{P}(\boldsymbol{y}' - \boldsymbol{x}') + \frac{\mathcal{O}(\lambda)}{2} \left\| \boldsymbol{P}(\boldsymbol{y}' - \boldsymbol{x}') \right\|_{\infty}^2 \end{split}$$

<sup>28</sup> A *circulation* is a vector  $d \in \mathbb{R}^{|E|}$  such that  $B^{\top}d = 0$ .

<sup>29</sup> Here,  $\|A\|_{\alpha \to \beta} \doteq \max_{\|X\|_{\alpha}=1} \|AX\|_{\beta}$ is the *operator norm* of *A* induced by norms  $\|\cdot\|_{\alpha}$  on the input space and  $\|\cdot\|_{\beta}$ on the output space. We have,

$$\left\|Ax\right\|_{\beta} \leq \left\|A\right\|_{\alpha \to \beta} \left\|x\right\|_{\alpha}.$$
 (27)

<sup>30</sup> The chain rule says that for a function  $g(\mathbf{x}) \doteq s(h(\mathbf{x})),$ 

$$Dg(\mathbf{x}) = Df(h(\mathbf{x}))Dh(\mathbf{x}).$$

Moreover,  $\nabla g(\mathbf{x}) = (Dg(\mathbf{x}))^{\top}$ . In this case,  $h(\mathbf{x}) = f_0 + P\mathbf{x}$ , so  $Dh(\mathbf{x}) = P$ .

$$= g(\mathbf{x}') + (\underbrace{\mathbf{P}^{\top} \nabla s(f_0 + \mathbf{P}\mathbf{x}')}_{= \nabla g(\mathbf{x}')})^{\top} (\mathbf{y}' - \mathbf{x}') + \frac{\mathcal{O}(\lambda)}{2} \|\mathbf{P}(\mathbf{y}' - \mathbf{x}')\|_{\infty}^2$$
  
$$\leq g(\mathbf{x}') + \nabla g(\mathbf{x}')^{\top} (\mathbf{y}' - \mathbf{x}') + \frac{\mathcal{O}(\lambda)}{2} \|\mathbf{P}\|_{\infty \to \infty}^2 \|\mathbf{y}' - \mathbf{x}'\|_{\infty}^2.$$

Applying lemma 26, completes the proof.

using eq. (27)

We denote by  $X^*$  the set of vectors  $x^*$  that are minimizers of g, i.e., for which we have  $g(x^*) = \min_{x \in \mathbb{R}^n} g(x)$ .

**Lemma 31.** For any  $\epsilon > 0$  and  $\hat{\mathbf{x}} \in \mathbb{R}^n$  such that  $g(\hat{\mathbf{x}}) \leq g(\mathbf{x}^*) + \frac{\epsilon}{2}OPT$ , we have,

$$\|f_0 + P\hat{x}\|_{\infty} \le (1+\epsilon)OPT$$

for some  $\lambda = \Theta(\log |E|/\epsilon OPT)$ .

Proof. We know,

$$\|f_0 + P\hat{x}\|_{\infty} \leq s(f_0 + P\hat{x}) = g(\hat{x}) \leq g(x^*) + \frac{\epsilon}{2}OPT.$$

Thus, it suffices to show,  $g(\mathbf{x}^{\star}) \leq (1 + \epsilon/2)OPT$ . We have,

$$g(\mathbf{x}^{\star}) = \min_{\mathbf{x} \in \mathbb{R}^{n}} g(\mathbf{x})$$
  
=  $\min_{\mathbf{x} \in \mathbb{R}^{n}} s(f_{0} + P\mathbf{x})$   
 $\leq 2 \frac{\log |E|}{\lambda} + \min_{\mathbf{x} \in \mathbb{R}^{n}} ||f_{0} + P\mathbf{x}||_{\infty}$   
=  $2 \frac{\log |E|}{\lambda} + OPT$   
 $\stackrel{!}{\leq} \left(1 + \frac{\epsilon}{2}\right) OPT.$ 

Solving for  $\lambda$ , we obtain,

$$\lambda \geq \frac{4\log|E|}{\epsilon OPT}.$$

Hence, choosing  $\lambda = \Theta(\log |E|/\epsilon OPT)$  is sufficient.

#### 2.10 Part J

It can be shown that,

$$R \doteq \max_{\substack{\boldsymbol{x} \in \mathbb{R}^n \\ g(\boldsymbol{x}) \le g(\boldsymbol{x}_0)}} \|\boldsymbol{x} - \boldsymbol{x}^*\|_{\infty} = \max_{\substack{\boldsymbol{x} \in \mathbb{R}^n \\ g(\boldsymbol{x}) \le g(\boldsymbol{x}_0)}} \min_{\substack{\boldsymbol{x}^* \in X^*}} \|\boldsymbol{x} - \boldsymbol{x}^*\|_{\infty}.$$
 (28)

#### Lemma 32.

(1)  $\|f_0\|_{\infty} \leq (1 + \|P\|_{\infty \to \infty})OPT.$ 

(2) For any y such that  $g(y) \leq g(\mathbf{0})$ , we have,

$$g(\boldsymbol{y}) \leq (1 + \|\boldsymbol{P}\|_{\infty \to \infty})OPT + 2\frac{\log|E|}{\lambda}.$$

(3)  $R = \mathcal{O}((1 + \|\boldsymbol{P}\|_{\infty \to \infty})OPT + \log|E|/\lambda)$  when  $\boldsymbol{x}_0 \doteq \boldsymbol{0}$ .

*Proof of (1).* Consider an optimal circulation  $d^* \in \mathbb{R}^{|E|}$ . As the discussed optimization problems are all equivalent, we have,

$$\left\| \boldsymbol{U}^{-1}(\boldsymbol{f}_0 + \boldsymbol{d}^*) \right\|_{\infty} = OPT.$$

We have,

$$\begin{split} \|f_0\|_{\infty} &= \left\| (I-P)U^{-1}\tilde{f}_0 \right\|_{\infty} \\ &= \left\| U^{-1}\tilde{f}_0 + U^{-1}\hat{P}d^* - PU^{-1}\tilde{f}_0 - U^{-1}\hat{P}d^* \right\|_{\infty} \\ &\leq \left\| U^{-1}\tilde{f}_0 + U^{-1}\hat{P}d^* \right\|_{\infty} + \left\| PU^{-1}\tilde{f}_0 + U^{-1}\hat{P}d^* \right\|_{\infty} \\ &= \left\| U^{-1}(\tilde{f}_0 + \hat{P}d^*) \right\|_{\infty} + \left\| PU^{-1}(\tilde{f}_0 + d^*) \right\|_{\infty}. \end{split}$$

using the triangle inequality

using 
$$\hat{P} = UPU^{-1}$$

Recall that  $d^*$  was the result of the projection  $\hat{P}x^*$  for some  $x^*$ . Therefore, due to the idempotency of projections,  $\hat{P}d^* = \hat{P}^2x^* = \hat{P}x^* = d^*$ , and we obtain,

$$= \left\| \mathbf{U}^{-1}(\tilde{f}_{0} + d^{*}) \right\|_{\infty} + \left\| \mathbf{P}\mathbf{U}^{-1}(\tilde{f}_{0} + d^{*}) \right\|_{\infty}$$
 using  $\hat{P} = \mathbf{U}\mathbf{P}\mathbf{U}^{-1}$   

$$\le \left\| \mathbf{U}^{-1}(\tilde{f}_{0} + d^{*}) \right\|_{\infty} + \left\| \mathbf{P} \right\|_{\infty \to \infty} \left\| \mathbf{U}^{-1}(\tilde{f}_{0} + d^{*}) \right\|_{\infty}$$
 using eq. (27)  

$$= (1 + \left\| \mathbf{P} \right\|_{\infty \to \infty}) \left\| \mathbf{U}^{-1}(\tilde{f}_{0} + d^{*}) \right\|_{\infty}$$
  

$$= (1 + \left\| \mathbf{P} \right\|_{\infty \to \infty}) OPT.$$

*Proof of (2).* We have,

$$g(\boldsymbol{y}) \leq g(\boldsymbol{0}) = s(f_0) \leq 2 \frac{\log |E|}{\lambda} + \|f_0\|_{\infty}$$
$$\leq 2 \frac{\log |E|}{\lambda} + (1 + \|\boldsymbol{P}\|_{\infty \to \infty}) OPT. \qquad \Box$$

*Proof of* (3). We have,

$$R = \max_{\substack{\boldsymbol{x} \in \mathbb{R}^n \\ g(\boldsymbol{x}) \le g(\boldsymbol{0})}} \min_{\boldsymbol{x}^{\star} \in X^{\star}} \|\boldsymbol{x} - \boldsymbol{x}^{\star}\|_{\infty}.$$

Observe that given any  $x^* \in X^*$ , we have for  $y \doteq Px^* + x - Px$  that  $Py = Px^*$  and therefore  $y \in X^*$ . This gives a feasible solution for the minimum,

$$\leq \max_{\substack{\boldsymbol{x} \in \mathbb{R}^n \\ g(\boldsymbol{x}) \leq g(\boldsymbol{0})}} \|\boldsymbol{x} - \boldsymbol{y}\|_{\infty}$$

$$= \max_{\substack{x \in \mathbb{R}^n \\ g(x) \le g(0)}} \|Px - Px^*\|_{\infty}$$

$$= \max_{\substack{x \in \mathbb{R}^n \\ g(x) \le g(0)}} \|f_0 + Px - f_0 - Px^*\|_{\infty}$$

$$\leq \max_{\substack{x \in \mathbb{R}^n \\ g(x) \le g(0)}} \|f_0 + Px\|_{\infty} + \|f_0 + Px^*\|_{\infty}$$

$$\leq \max_{\substack{x \in \mathbb{R}^n \\ g(x) \le g(0)}} g(x) + \underbrace{g(x^*)}_{\le g(x)}$$

$$\leq 2 \max_{\substack{x \in \mathbb{R}^n \\ g(x) \le g(0)}} g(x)$$

$$\leq 2(1 + \|P\|_{\infty \to \infty})OPT + 4\frac{\log |E|}{\lambda}.$$

2.11 Part K

We choose  $\hat{P} \doteq I - UBL^+B^\top$  with the Laplacian  $L \doteq B^\top UB^{.31}$  Then,

$$\left\| \boldsymbol{U}^{-1} \hat{\boldsymbol{P}} \boldsymbol{U} \right\|_{\infty \to \infty} = \left\| \boldsymbol{P} \right\|_{\infty \to \infty} \le 1 + 8 rac{\log |E|}{\Phi^2},$$

where  $\Phi$  is the expansion of the graph.

**Theorem 33.** Gradient descent with respect to g and  $\|\cdot\|_{\infty}$  yields a  $(1 + \epsilon)$  approximate solution to OPT in time  $\tilde{\mathcal{O}}(|E|/\epsilon^2 \Phi^8)$  under the assumption that solving a Laplacian linear system exactly is as expensive as finding a  $1/|V|^{100}$ -approximate solution.

*Proof.* First, recall that *g* is continuously differentiable, convex, and  $\mathcal{O}\left(\lambda \|\boldsymbol{P}\|_{\infty \to \infty}^2\right)$ -smooth with respect to  $\|\cdot\|_{\infty}$ . Using our analysis from theorem 24, gradient descent with respect to *g* and  $\|\cdot\|_{\infty}$  will evaluate  $\nabla g$  and  $(\cdot)_*^{\#}$  at most  $\mathcal{O}\left(\beta R^2/\epsilon OPT\right)$  times and use at most  $\mathcal{O}\left(|E|\beta R^2/\epsilon OPT\right)$  additional arithmetic operations to find an  $\hat{\boldsymbol{x}}$  such that  $g(\hat{\boldsymbol{x}}) - g(\boldsymbol{x}^*) \leq \frac{\epsilon}{2} OPT$ . By lemma 31, we know that this yields a  $(1 + \epsilon)$  approximation of *OPT* for some  $\lambda = \Theta(\log |E|/\epsilon OPT)$ .

We have for  $\beta$ ,

$$\begin{split} \beta &= \mathcal{O}\Big(\lambda \, \|\boldsymbol{P}\|_{\infty \to \infty}^2 \Big) = \mathcal{O}\left(\frac{\log |\boldsymbol{E}|}{\epsilon OPT} \left(1 + 8\frac{\log |\boldsymbol{E}|}{\Phi^2}\right)^2\right) \\ &= \mathcal{O}\left(\frac{\log |\boldsymbol{E}|}{\epsilon OPT} \cdot \frac{\log^2 |\boldsymbol{E}|}{\Phi^4}\right) \\ &= \mathcal{O}\left(\frac{\log^3 |\boldsymbol{E}|}{\epsilon OPT\Phi^4}\right) \end{split}$$

and for R,

$$R = \mathcal{O}\left((1 + \|\boldsymbol{P}\|_{\infty \to \infty})OPT + \frac{\log|E|}{\lambda}\right)$$

using the triangle inequality

<sup>31</sup> Note that the definition of the incidence matrix used here is the transpose of the incidence matrix as defined in the lecture notes.

$$= \mathcal{O}\left(\left(2 + 8\frac{\log|E|}{\Phi^2} + \epsilon\right)OPT\right)$$
$$= \mathcal{O}\left(\frac{\log|E|}{\Phi^2}OPT\right).$$

Combining these bounds, we get,

$$\mathcal{O}\left(\frac{\beta R^2}{\epsilon OPT}\right) = \mathcal{O}\left(\frac{\log^5 |E|}{\epsilon^2 \Phi^8}\right) = \tilde{\mathcal{O}}\left(\frac{1}{\epsilon^2 \Phi^8}\right)$$

It therefore remains to show that each iteration of gradient descent takes  $\tilde{O}(|E|)$  time.

For a matrix *A*, let *T*(*A*) be the maximum time to compute *Ax* and  $A^{\top}x$  for any vector *x*. We use the following claim:

**Claim 34.** 
$$T(P) = O(|E|)$$
.

By the chain rule, we have  $\nabla g(\mathbf{x}) = \mathbf{P}^\top \nabla s(f_0 + \mathbf{P}\mathbf{x})$ . Thus,  $\nabla g$  can be computed in time  $T(\mathbf{P})$  plus the time to compute  $\nabla s$ . It is not hard to show that  $\nabla s$  can be computed in time  $\mathcal{O}(|E|)$ .<sup>32</sup> Finally, observe that  $(\cdot)^{\#}_*$  corresponds to the dual vector map of the Manhattan norm, which we stated in eq. (13). Clearly, this mapping can be computed in  $\mathcal{O}(|E|)$  time. Therefore, each iteration of gradient descent can be computed in  $\tilde{\mathcal{O}}(|E|)$  time.

*Proof of claim* 34. We have,

$$P = U^{-1} \hat{P} U = I - BL^+ B^\top U$$
 and  
 $P^\top = I - U^\top B^\top L^+ B.$ 

Trivially, *Ix*, *Ux*, and  $U^{\top}x$  can be computed in  $\mathcal{O}(|E|)$  time. As by definition of the incidence matrix *B*, nnz(*B*) =  $\mathcal{O}(|E|)$ , *Bx* and  $B^{\top}x$  can also be computed in  $\mathcal{O}(|E|)$  time.

It follows from the definition of the incidence matrix that  $1 \in \ker B$ and  $1 \in \ker B^{\top}$ . Therefore, we have for any  $y \in \mathbb{R}^{|V|}$  that  $By \perp 1$  and for any  $x \in \mathbb{R}^{|E|}$  that  $B^{\top}x \perp 1$ .<sup>33</sup> Therefore,  $B^{\top}Ux \perp 1$  and  $By \perp 1$ for any x and y.<sup>34</sup>

Using the result of Kyng and Sachdeva,<sup>35</sup> we can find an  $\epsilon$ -approximate solution  $\tilde{z}$  to Lz = d in time  $\mathcal{O}(|E|\log^3 |V|\log(1/\epsilon))$ , where  $d \perp 1$ . Using our assumption that finding  $\tilde{z}$  for  $\epsilon = 1/|V|^{100}$  is as expensive as finding z exactly, we conclude that  $L^+B^\top Ux$  and  $L^+By$  can be computed in

$$\mathcal{O}\left(|E|\log^3|V|\log|V|^{100}\right) = \mathcal{O}\left(|E|\log^4|V|\right) = \tilde{\mathcal{O}}(|E|)$$

time.

 $^{32}$  The argument is similar to our computation of the first-order derivative of the function *m* in the proof of lemma 29(2).

using  $(L^+)^\top = L^+$ 

<sup>33</sup>  $(By)^{\top}\mathbf{1} = y^{\top}B^{\top}\mathbf{1} = y^{\top}\mathbf{0} = \mathbf{0}$ ; the other case is symmetric

<sup>34</sup>  $B^{\top} Ux$  and By can be interpreted as demand vectors.

<sup>&</sup>lt;sup>35</sup> Corollary 10.2.5 in the lecture notes

# References

 Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends*® *in Machine Learning*, 8(3-4):231–357, 2015.