

### Final Exam Solutions

You may use any books, notes, or pdf/text files on your laptop (tablets, cell phones and other devices are prohibited). **You are not permitted to connect to the internet or communicate with anyone in anyway during the exam. Please disconnect wireless, bluetooth etc. now. Any (even seemingly minor) violations will result in a failing grade for the course.**

All sub-problems have roughly equal weight. Some are easy. Others, not so much.

Total	/ 100
Q1	/ 20
Q2	/ 15
Q3	/ 25
Q4	/ 15
Q5	/ 25

**1- Short questions (20 pts). (2.5 points (a) and (b), (5) points (c)-(d))**

Specify whether the functions in parts (a)-(d) are strictly convex, convex, or nonconvex, and give a justification (brief justification in parts a and b)

(a)  $f(x) = \frac{x}{1+\frac{1}{x}}$  for  $x > 0$ .

(b)

$$f(x) = \begin{cases} 0 & |x| \leq 1 \\ |x - 1| & \text{otherwise} \end{cases}$$

(c)  $-\text{trace}(\mathbf{X}) \cdot \log \det \left( \frac{\mathbf{X}}{\text{trace}(\mathbf{X})} \right)$  on  $\{\mathbf{X} \mid \mathbf{X} \succ \mathbf{0}\}$ .

(d)

$$f(\mathbf{x}) = \sup_{t \in [a, b]} \log \left( \frac{p(t)}{q(t)} \right) \quad \text{where} \quad p(t) = e^{x_1 \sin(t)} + e^{x_2 \sin(2t)} + \dots + e^{x_n \sin(nt)}$$

$$\text{and} \quad q(t) = e^{x_1 \sin(t) + x_2 \sin(2t) + \dots + x_n \sin(nt)}$$

Here  $a$  and  $b$  are real constants, with  $a < b$ .

(e) What is the proximal mapping associated with the regularizer

$$\mathcal{R}(z) = \sum_{i=1}^n w_i |z_i|,$$

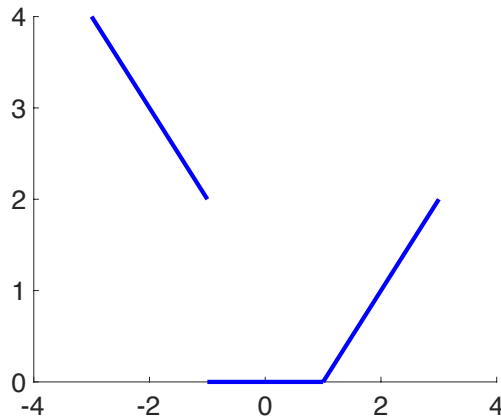
where  $w_i > 0$  are known non-negative weights.

**Solution**

- (a) This function is strictly convex as

$$f''(x) = \frac{2}{(x+1)^3} > 0.$$

- (b) Nonconvex, see drawing below.



- (c) Strictly convex. Note that  $-\log\det(\mathbf{X})$  is convex. Hence its perspective function  $f(\mathbf{X}, t) = -t \cdot \log\det\left(\frac{\mathbf{X}}{t}\right)$  is convex over  $(\mathbf{X}, t)$ . Now compose the latter with a linear function of the form  $\mathbf{X} \rightarrow (\mathbf{X}, \text{trace}(\mathbf{X}))$ .
- (d) Strictly convex. Note that  $f(\mathbf{x})$  can be simplified into the form

$$f(\mathbf{x}) = \sup_{t \in [a, b]} \log\left(\frac{p(t)}{q(t)}\right) = \sup_{t \in [a, b]} \left( \log\left(\sum_{k=1}^n e^{x_k \sin(kt)}\right) - \left(\sum_{k=1}^n x_k \sin(kt)\right) \right)$$

Now note that for a fixed  $t$  the first term is convex (as log-sum-exp is convex) and the second term is linear. So their difference is convex for a fixed  $t$ . Supremum of convex functions preserves convexity so  $f$  is convex.

- (e) The proximal mapping is defined as

$$\mathcal{S}_{\mathcal{R}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|_{\ell_2}^2 + \mathcal{R}(\mathbf{z}).$$

For the weighted  $\ell_1$  norm defined above we have

$$\mathcal{S}_{\mathcal{R}}(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathbb{R}^n} \sum_{i=1}^n \left( \frac{1}{2} |z_i - x_i|^2 + w_i |z_i| \right),$$

Since the loss above is separable we have

$$[\mathcal{S}_{\mathcal{R}}(\mathbf{x})]_i = \arg \min_{z_i \in \mathbb{R}} \frac{1}{2} |z_i - x_i|^2 + w_i |z_i|$$

However, in class we learned that this value is given by soft thresholding at level  $w_i$ . Therefore,

$$[\mathcal{S}_{\mathcal{R}}(\mathbf{x})]_i = \arg \min_{z_i \in \mathbb{R}} \frac{1}{2} |z_i - x_i|^2 + w_i |z_i| = \mathcal{S}_{w_i}(x_i),$$

where

$$\mathcal{S}_{\mu}(x) = \begin{cases} x - \mu & \text{if } x \geq \mu \\ 0 & \text{if } |x| \leq \mu \\ x + \mu & \text{if } x \leq -\mu \end{cases}$$

Thus

$$\mathcal{S}_{\mathcal{R}}(\mathbf{x}) = \begin{bmatrix} \mathcal{S}_{w_1}(x_1) \\ \mathcal{S}_{w_2}(x_2) \\ \vdots \\ \mathcal{S}_{w_n}(x_n) \end{bmatrix}.$$

■

**2- SDP. (15 points)**

We wish to solve the following optimization problem

$$\text{minimize} \quad \frac{(\mathbf{c}^T \mathbf{x})^2}{\mathbf{d}^T \mathbf{x}} \quad \text{subject to} \quad \mathbf{Ax} + \mathbf{b} \geq 0,$$

where we assume that  $\mathbf{d}^T \mathbf{x} > 0$  whenever  $\mathbf{Ax} + \mathbf{b} \geq 0$ .

- (a) Is this optimization problem convex? If yes, justify why. If no, cast it into an equivalent convex optimization problem. (5 points)
- (b) How would you solve this problem using an SDP solver. Justify your answer. (10 points)

**Solution**

- (a) Yes, it is convex. Note that the function  $f(z, t) = z^2/t$  is convex over  $(z, t)$  with  $t > 0$ . Thus, its pre-composition with a linear function of the form  $\mathbf{x} \rightarrow (\mathbf{c}^T \mathbf{x}, \mathbf{d}^T \mathbf{x})$  remains convex.
- (b) Using the epigraph trick we can rewrite the above problem in the form

$$\begin{aligned} &\text{minimize} \quad t \\ &\text{subject to} \quad \mathbf{Ax} + \mathbf{b} \geq 0, \\ &\quad \quad \quad \frac{(\mathbf{c}^T \mathbf{x})^2}{\mathbf{d}^T \mathbf{x}} \leq t. \end{aligned}$$

Using the Shur trick from class the latter is equivalent to

$$\begin{aligned} &\text{minimize} \quad t \\ &\text{subject to} \quad \begin{bmatrix} \text{diag}(\mathbf{Ax} + \mathbf{b}) & 0 & 0 \\ 0 & t & \mathbf{c}^T \mathbf{x} \\ 0 & \mathbf{c}^T \mathbf{x} & \mathbf{d}^T \mathbf{x} \end{bmatrix} \succeq \mathbf{0}. \end{aligned}$$

■

### 3- Its all dual. (25 pts)

Consider the following optimization problem

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } g(\mathbf{x}) \leq 0.$$

where

$$f(\mathbf{x}) = f(x_1, x_2) = e^{-x_1} \quad \text{and} \quad g(\mathbf{x}) = g(x_1, x_2) = \frac{x_1^2}{x_2} \quad \text{with domain } \mathcal{X} = \{(x_1, x_2) | x_2 > 0\}.$$

Please note that the optimization variables are  $x_1$  and  $x_2$ .

- Is the function  $g$  convex over  $\mathcal{X}$ ?
- What is the primal optimal value?
- Derive the dual optimization problem. Note that  $\mathcal{X}$  is an implicit constraint. When creating the dual do not make it explicit.
- What is the optimal value of the dual optimization problem. Does strong duality hold? Why or why not?
- Consider the perturbed problem

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } g(\mathbf{x}) \leq u,$$

with  $u > 0$ . What is the duality gap?

### Solution

- Yes, quadratic over linear is convex when  $x_2 > 0$ . If you would like a proof note that

$$\nabla^2 g(\mathbf{x}) = \begin{bmatrix} \frac{2}{x_2} & -\frac{2x_1}{x_2^2} \\ -\frac{2x_1}{x_2^2} & \frac{2x_1^2}{x_2^3} \end{bmatrix} = \frac{2}{x_2^3} \begin{bmatrix} x_2 & -x_1 \\ -x_1 & x_1^2 \end{bmatrix} \begin{bmatrix} x_2 \\ -x_1 \end{bmatrix}^T \geq 0.$$

- Since  $x_1^2/x_2 \leq 0$  and  $x_2 > 0$  then all feasible solutions are  $(x_1, x_2)$  with  $x_1 = 0$  and  $x_2 > 0$ . The objective only depends on  $x_1$  and so the optimal value is simply  $e^0 = 1$ .
- The Lagrangian is given by

$$\mathcal{L}(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x}) = \begin{cases} +\infty & \text{if } x_2 \leq 0 \\ e^{-x_1} + \lambda \frac{x_1^2}{x_2} & \text{if } x_2 > 0 \end{cases}$$

The Lagrangian obtains it's minimum value of 0 when  $x_2 = x_1^3$  and  $x_1 \rightarrow \infty$ . Thus,

$$g(\lambda) = 0.$$

Therefore, the dual problem is given by

$$\max_{\lambda \geq 0} g(\lambda) = \max_{\lambda \geq 0} 0.$$

- (d) The primal optimal value is equal to 1 and dual optimal value is equal to zero. Therefore, strong duality does not hold. We expect this to happen as Slater's condition does not hold in this case.
- (e) A short answer that would have been sufficient is that Slater's condition holds and therefore we automatically have strong duality with zero duality gap.

A longer answer is that in this case the primal optimization problem is of the form

$$\text{minimize } e^{-x_1} \quad \text{subject to } x_1^2 \leq x_2 u.$$

Now note that for  $x_1 = t$  with  $t > 0$  and  $x_2 = t^2/u$  the objective value is equal to  $e^{-t}$  and the constraint is satisfied. By making  $t$  arbitrary large the primal objective can be made arbitrary small so that the primal optimal value is equal to zero.

The Lagrangian is given by

$$\mathcal{L}(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda (g(\mathbf{x}) - u) = \begin{cases} +\infty & \text{if } x_2 \leq 0 \\ e^{-x_1} + \lambda \left( \frac{x_1^2}{x_2} - u \right) & \text{if } x_2 > 0 \end{cases}$$

The Lagrangian converges to  $-\lambda u$  when  $x_2 = x_1^3$  and  $x_1 \rightarrow \infty$ . Thus,

$$g(\lambda) = -\lambda u.$$

So dual problem is

$$\max_{\lambda \geq 0} g(\lambda) = \max_{\lambda \geq 0} -\lambda u.$$

Therefore, the dual optima value is equal to 0. Thus, strong duality does indeed hold and the duality gap is equal to 0.

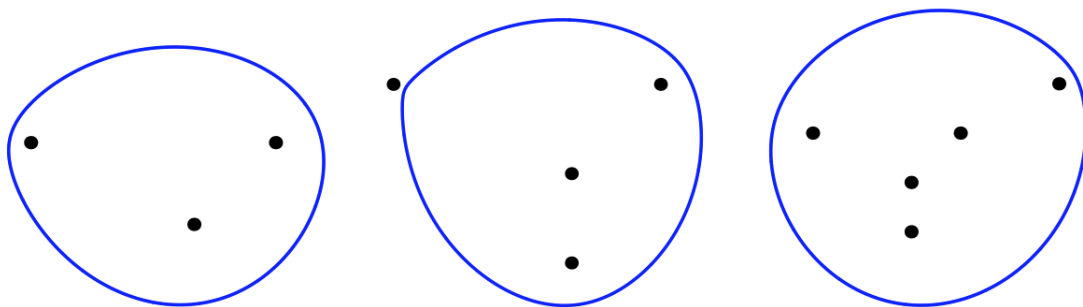
■

**4- Fixed sum of distances from  $k$  given points. (15 pts)**

Fix a positive real number  $d$  and fix  $k$  distinct points  $(u_1, v_1), (u_2, v_2), \dots, (u_k, v_k)$  in  $\mathbb{R}^2$ . The generalized  $k$ -ellipse with foci  $(u_i, v_i)$  and radius  $d$  is the following curve in the plane:

$$\left\{ (x, y) \in \mathbb{R}^2 : \sum_{i=1}^k \sqrt{(x - u_i)^2 + (y - v_i)^2} = d \right\}$$

Figure 1 below depicts such generalized ellipses.



**Figure 1:** A 3-ellipse, a 4-ellipse, and a 5-ellipse, each with its foci.

- (a) Is the set  $\mathcal{E}_k$  defined below convex? Justify your answer. (5 points)

$$\mathcal{E}_k := \left\{ (x, y) \in \mathbb{R}^2 : \sum_{i=1}^k \sqrt{(x - u_i)^2 + (y - v_i)^2} \leq d \right\}$$

- (b) How would you find a point at the intersection of two such generalized ellipses using an SDP solver. Justify your answer. (10 points)

**Solution**

- (a) Yes, it is convex. Consider the vector  $\mathbf{z} = \begin{bmatrix} x \\ y \end{bmatrix}$  and vectors  $\mathbf{p}_i = \begin{bmatrix} u_i \\ v_i \end{bmatrix}$ . Then the set  $\mathcal{E}^{(k)}$  can be rewritten in the form

$$\mathcal{E}_k := \left\{ \mathbf{z} \in \mathbb{R}^2 : f(\mathbf{z}) \leq d \right\} \quad \text{with} \quad f(\mathbf{z}) = \sum_{i=1}^k \|\mathbf{z} - \mathbf{p}_i\|_{\ell_2}.$$

Note that  $f(\mathbf{z})$  is a sum of convex functions and hence convex. Therefore, its sub-level sets including  $\mathcal{E}_k$  is also convex.



(b) Let us consider the two  $k$ -ellipses

$$\mathcal{E}_k := \left\{ \mathbf{z} \in \mathbb{R}^2 : f(\mathbf{z}) \leq d \right\} \quad \text{with} \quad f(\mathbf{z}) = \sum_{i=1}^k \|\mathbf{z} - \mathbf{p}_i\|_{\ell_2}.$$

and

$$\tilde{\mathcal{E}}_k := \left\{ \mathbf{z} \in \mathbb{R}^2 : \tilde{f}(\mathbf{z}) \leq \tilde{d} \right\} \quad \text{with} \quad \tilde{f}(\mathbf{z}) = \sum_{i=1}^k \|\mathbf{z} - \tilde{\mathbf{p}}_i\|_{\ell_2}.$$

To find a point at the intersection (or show non-exists) it suffices to solve the following optimization problem

$$\min_{\mathbf{z} \in \mathbb{R}^2} 0 \quad \text{subject to} \quad f(\mathbf{z}) \leq d \quad \text{and} \quad \tilde{f}(\mathbf{z}) \leq \tilde{d}.$$

The optimal value will be zero if a point exists at the intersection and  $+\infty$  if such a point does not exist. To be able to solve this via SDP we only need to cast the constraints  $f(\mathbf{z}) \leq d$  and  $\tilde{f}(\mathbf{z}) \leq \tilde{d}$  as an LMI.

To this aim note that using the epigraph trick we have

$$\sum_{i=1}^k \|\mathbf{z} - \mathbf{p}_i\|_{\ell_2} \leq d \quad \Leftrightarrow \quad \begin{cases} \|\mathbf{z} - \mathbf{p}_i\|_{\ell_2} \leq t_i & \text{for } i = 1, 2, \dots, k \\ \sum_{i=1}^k t_i \leq d. \end{cases}$$

Using Shur's trick the latter can be written in the form

$$\sum_{i=1}^k t_i \leq d \quad \text{and} \quad \begin{bmatrix} t_i \mathbf{I}_{2 \times 2} & \mathbf{z} - \mathbf{p}_i \\ (\mathbf{z} - \mathbf{p}_i)^T & t_i \end{bmatrix} \geq \mathbf{0} \quad \text{for } i = 1, 2, \dots, k.$$

Putting everything together we arrive at

$$\begin{aligned} & \min_{\mathbf{z} \in \mathbb{R}^2, \mathbf{s}, \mathbf{t} \in \mathbb{R}^k} 0 \\ & \text{subject to} \quad \mathbf{1}^T \mathbf{t} \leq d \quad \text{and} \quad \mathbf{1}^T \mathbf{s} \leq \tilde{d}, \\ & \quad \begin{bmatrix} t_i \mathbf{I}_{2 \times 2} & \mathbf{z} - \mathbf{p}_i \\ (\mathbf{z} - \mathbf{p}_i)^T & t_i \end{bmatrix} \geq \mathbf{0} \quad \text{for } i = 1, 2, \dots, k \\ & \quad \begin{bmatrix} s_i \mathbf{I}_{2 \times 2} & \mathbf{z} - \tilde{\mathbf{p}}_i \\ (\mathbf{z} - \tilde{\mathbf{p}}_i)^T & s_i \end{bmatrix} \geq \mathbf{0} \quad \text{for } i = 1, 2, \dots, k \end{aligned}$$

which is an SDP.

■

### 5- Structured parameter estimation from nonlinear measurements. (25 pts)

We wish to recover an unknown signal  $\mathbf{x}$  from nonlinear measurements of the form

$$y_r = |\langle \mathbf{a}_r, \mathbf{x} \rangle|^2 + w_r.$$

Here,  $\mathbf{a}_r \in \mathbb{R}^n$  are known measurement vectors,  $y_r$  are the measurements and  $w_r$  are the noise in the measurements. We assume that the noise  $w_r$  are distributed i.i.d. according to a Poisson noise model. Specifically the measurements are generated i.i.d. according to the following distribution

$$\text{prob}(y_r = k) = \frac{e^{-\mu_r} \mu_r^k}{k!} \quad \text{with} \quad \mu_r = (\langle \mathbf{a}_r, \mathbf{x} \rangle)^2,$$

with  $\mathbf{x} \in \mathbb{R}^n$  denoting the unknown signal.

- (a) Write down the objective function  $\mathcal{L}(\mathbf{z})$  with variable  $\mathbf{z}$  that you would minimize to find the signal  $\mathbf{x}$  (Note that the objective function does not have to be convex). (5 points)
- (b) In practice, due to various practical considerations (e.g. that taking measurements are time consuming) we do not have a lot of measurements in the sense that  $m \ll n$  so that just minimizing the objective function in part (a) may not be enough as there may be many global optima for the objective in part (a). Luckily however we have a lot of side information that we can utilize. In particular, in our problem we know that an orthonormal rotation of the signal vector  $\mathbf{x}$  is sparse. Specifically,  $\mathbf{U}^T \mathbf{x}$  has only a few non-zero entries with  $\mathbf{U} \in \mathbb{R}^{n \times n}$  a known orthonormal matrix (i.e. obeying  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ ). Based on our side information we wish to solve one of the two following optimization problems. (Below  $\tau, \lambda \geq 0$  are known tuning parameters)

$$\min_{\mathbf{z} \in \mathbb{R}^p} \mathcal{L}(\mathbf{z}) \quad \text{subject to} \quad \|\mathbf{U}^T \mathbf{z}\|_{\ell_1} \leq \tau, \quad (0.1)$$

$$\min_{\mathbf{z} \in \mathbb{R}^p} \mathcal{L}(\mathbf{z}) + \lambda \|\mathbf{U}^T \mathbf{z}\|_{\ell_1}. \quad (0.2)$$

What are the iterative updates when solving (0.1) via the Frank-Wolfe algorithm? (5 points)

- (c) What are the iterative updates when solving (0.2) via ISTA? (Hint: use the case of  $\mathbf{U} = \mathbf{I}$  as a starting point) (5 points)
- (d) What changes would you make to the iterations of part (c) to make the algorithm converge faster. Please write down the corresponding iterations. (5 points)
- (e) How would you solve (0.2) using a stochastic algorithm (i.e. a stochastic version of part (c)). Please write down the iterations. (5 points)

### Solution

- (a) We need to minimize the negative log-likelihood. Following page 354 of the book we arrive at the objective

$$\mathcal{L}(\mathbf{z}) = \sum_{r=1}^m |\langle \mathbf{a}_r, \mathbf{z} \rangle|^2 - y_r \log(|\langle \mathbf{a}_r, \mathbf{z} \rangle|^2)$$

- (b) Note that

$$\nabla \mathcal{L}(\mathbf{z}) = 2 \left( \sum_{r=1}^m \left( 1 - \frac{y_r}{|\langle \mathbf{a}_r, \mathbf{z} \rangle|^2} \right) \mathbf{a}_r \mathbf{a}_r^T \right) \mathbf{z}.$$

Also define

$$\mathcal{K} = \{ \mathbf{z} \in \mathbb{R}^n \mid \| \mathbf{U}^T \mathbf{z} \|_{\ell_1} \leq \tau \}.$$

Note that the Frank-Wolfe iterations are given by

$$\begin{aligned} \mathbf{u}_t &\in \arg \min_{\mathbf{u} \in \mathcal{K}} \mathbf{u}^T \nabla \mathcal{L}(\mathbf{z}_t), \\ \mathbf{z}_{t+1} &= \frac{t-1}{t+1} \mathbf{z}_t + \frac{2}{t+1} \mathbf{u}_t. \end{aligned}$$

To calculate  $\mathbf{u}_t$  note that

$$\begin{aligned} \mathbf{u}_t &\in \arg \min_{\mathbf{u} \in \mathcal{K}} \mathbf{u}^T \nabla f(\mathbf{z}_t) = \arg \min_{\mathbf{u} \in \mathcal{K}} \mathbf{u}^T \nabla \mathcal{L}(\mathbf{z}_t) \\ &= \arg \min_{\mathbf{u}: \| \mathbf{U}^T \mathbf{u} \|_{\ell_1} \leq \tau} \mathbf{u}^T \nabla \mathcal{L}(\mathbf{z}_t) \\ &= \arg \min_{\mathbf{u}: \| \mathbf{U}^T \mathbf{u} \|_{\ell_1} \leq \tau} \mathbf{u}^T \mathbf{U} \mathbf{U}^T \nabla \mathcal{L}(\mathbf{z}_t) \\ &= \arg \min_{\mathbf{u}: \| \mathbf{U}^T \mathbf{u} \|_{\ell_1} \leq \tau} \langle \mathbf{U}^T \mathbf{u}, \mathbf{U}^T \nabla \mathcal{L}(\mathbf{z}_t) \rangle \\ &= \mathbf{U} \cdot \arg \min_{\mathbf{w}: \| \mathbf{w} \|_{\ell_1} \leq \tau} \langle \mathbf{w}, \mathbf{U}^T \nabla \mathcal{L}(\mathbf{z}_t) \rangle \\ &= -\tau \operatorname{sgn}([\mathbf{U}^T \nabla \mathcal{L}(\mathbf{z}_t)]_i) (\mathbf{U} \mathbf{e}_{i_{\max}}) \\ &= -\tau \operatorname{sgn}([\mathbf{U}^T \nabla \mathcal{L}(\mathbf{z}_t)]_i) \mathbf{U}_{i_{\max}} \end{aligned}$$

where

$$i_{\max} = \arg \max_i |[\mathbf{U}^T \nabla \mathcal{L}(\mathbf{z}_t)]_i|.$$

Thus the iterations are equal to

$$\begin{aligned} \mathbf{u}_t &= -\tau \operatorname{sgn}([\mathbf{U}^T \nabla \mathcal{L}(\mathbf{z}_t)]_i) \mathbf{U}_{i_{\max}}, \\ \mathbf{z}_{t+1} &= \frac{t-1}{t+1} \mathbf{z}_t + \frac{2}{t+1} \mathbf{u}_t. \end{aligned}$$

(c) We rewrite (0.2) as

$$\min_{\mathbf{z} \in \mathbb{R}^n} \frac{1}{\lambda} \mathcal{L}(\mathbf{z}) + \mathcal{R}(\mathbf{z}) \quad \text{with} \quad \mathcal{R}(\mathbf{z}) = \|\mathbf{U}^T \mathbf{z}\|_{\ell_1}.$$

The iterations of ISTA are of the form

$$\mathbf{z}_{t+1} = \text{prox}_{\mu \mathcal{R}}(\mathbf{z}_t - \mu \nabla \mathcal{L}(\mathbf{z}_t)),$$

where

$$\text{prox}_{\mu \mathcal{R}}(\mathbf{x}) = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|_{\ell_2}^2 + \mu \mathcal{R}(\mathbf{z}).$$

Note for  $\mathcal{R}(\mathbf{z}) = \|\mathbf{U}^T \mathbf{z}\|_{\ell_1}$  we have

$$\begin{aligned} \text{prox}_{\mu \mathcal{R}}(\mathbf{x}) &= \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|_{\ell_2}^2 + \mu \|\mathbf{U}^T \mathbf{z}\|_{\ell_1}, \\ &= \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{U}^T (\mathbf{z} - \mathbf{x})\|_{\ell_2}^2 + \mu \|\mathbf{U}^T \mathbf{z}\|_{\ell_1}, \\ &= \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{U}^T \mathbf{z} - \mathbf{U}^T \mathbf{x}\|_{\ell_2}^2 + \mu \|\mathbf{U}^T \mathbf{z}\|_{\ell_1}, \\ &= \mathbf{U} \cdot \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w} - \mathbf{U}^T \mathbf{x}\|_{\ell_2}^2 + \mu \|\mathbf{w}\|_{\ell_1}, \\ &= \mathbf{U} \cdot \mathcal{S}_{\mu}(\mathbf{U}^T \mathbf{x}), \end{aligned}$$

where  $\mathcal{S}_{\mu}$  is the soft-thresholding operator and the prox of  $\mu \|\cdot\|_{\ell_1}$  as discussed in class. Thus the ISTA iterations are

$$\mathbf{z}_{t+1} = \mathbf{U} \cdot \mathcal{S}_{\mu} \left( \mathbf{U}^T \left( \mathbf{z}_t - \frac{\mu}{\lambda} \nabla \mathcal{L}(\mathbf{z}_t) \right) \right).$$

(d) In this case we would run FISTA instead of ISTA. The iterations of FISTA are of the form

$$\begin{aligned} \mathbf{v}_t &= \mathbf{z}_t - \beta (\mathbf{z}_t - \mathbf{z}_{t-1}) \\ \mathbf{z}_{t+1} &= \text{prox}_{\mu \mathcal{R}} \left( \mathbf{v}_t - \frac{\mu}{\lambda} \nabla \mathcal{L}(\mathbf{v}_t) \right). \end{aligned}$$

More specifically, it takes the form

$$\begin{aligned} \mathbf{v}_t &= \mathbf{z}_t - \beta (\mathbf{z}_t - \mathbf{z}_{t-1}), \\ \mathbf{z}_{t+1} &= \mathbf{U} \cdot \mathcal{S}_{\mu} \left( \mathbf{U}^T \left( \mathbf{v}_t - \frac{\mu}{\lambda} \nabla \mathcal{L}(\mathbf{v}_t) \right) \right). \end{aligned}$$

(e) Note that the objective is of the form

$$\sum_{r=1}^m \left( |\langle \mathbf{a}_r, \mathbf{z} \rangle|^2 - y_r \log(|\langle \mathbf{a}_r, \mathbf{z} \rangle|^2) + \frac{\lambda}{m} \|\mathbf{U}^T \mathbf{z}\|_{\ell_1} \right)$$

Thus we pick a random  $\phi_t \in \{1, 2, \dots, m\}$  and run the iterations

$$\mathbf{z}_{t+1} = \mathbf{U} \cdot \mathcal{S}_\mu \left( \mathbf{U}^T \left( \mathbf{z}_t - 2 \frac{\mu m}{\lambda} \left( 1 - \frac{y_r}{|\langle \mathbf{a}_{\phi_t}, \mathbf{z} \rangle|^2} \right) (\mathbf{a}_{\phi_t}^T \mathbf{z}_t) \mathbf{a}_{\phi_t} \right) \right).$$

■