---

## Optimization for Data Science    Final Exam (8 August 2020)    FS20

---

## Candidate

First name:   . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

Last name:   . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

Student ID (Legi) Nr.:   . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

I attest with my signature that I was able to take the
exam under regular conditions and that I have read and
understood the general remarks below.

Signature:   . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

## General remarks and instructions

1. Check your exam documents for completeness (pages numbered from 1 to 13).

2. You can solve the exercises in any order. They are not ordered by difficulty.

3. Immediately inform an assistant in case you are not able to take the exam under regular conditions. Later complaints will not be accepted.

4. Pencils are not allowed. Pencil-written solutions will not be reviewed.

5. No auxiliary material allowed. All electronic devices must be turned off and are not allowed to be on your desk or carried with you to the toilet. We will write the current time on the blackboard every 15 minutes.

6. Attempts to cheat/defraud lead to immediate exclusion from the exam and can have judicial consequences.

7. Provide only one solution to each exercise. Cancel invalid solutions clearly.

8. **All solutions must be understandable and well-founded. Write down the important thoughts in clear sentences and keywords. Unless stated otherwise, no points will be awarded for unfounded or incomprehensible solutions. Please write your solutions in English.**

9. You may use anything that has been introduced and proved in the lecture or in the exercise sessions without reproving it. However, if you need something *different* than what we have shown, you must write a new proof or at least list all necessary changes.

Good luck!

| | achieved points (maximum) |
|---|---:|
| 1 | (15) |
| 2 | (9) |
| 3 | (11) |
| 4 | (18) |
| 5 | (8) |
| 6 | (9) |
| Σ | (70) |

# Useful Facts

In your solutions, you may use the following facts:

**Fact 1.** *Let* $h : \mathbf{dom}(h) \to \mathbb{R}$ *be convex and differentiable, and let* $X \subseteq \mathbf{dom}(h)$ *be a convex set. Then* $\mathbf{x}^*$ *is a minimizer of* $h$ *over* $X$ *if and only if*

$$\nabla h(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) \geq 0 \qquad \forall \mathbf{x} \in X.$$

**Fact 2.** *Let* $G \sim N(0, 1)^{m \times m}$ *for* $m \geq 1$. *Then*

$$\mathbb{P}\left( \|G\| \geq 10\sqrt{m} \right) \leq 0.001 \,.$$

**Fact 3.** *Let* $M, m \in \mathbb{N}$ *and* $G \sim N(0, 1)^{M \times m}$. *Then*

$$\mathbb{P}\left( \left\| \frac{1}{M} G^\top G - \mathrm{Id}_m \right\| \leq 0.001 \right) \geq 1 - 2\exp\left( -\frac{M}{10^6 m} \right) \,.$$

**Assignment 1.** **(15 points)** *Define*

$$\Delta_n := \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} > \mathbf{0}, \mathbf{1}^\top \mathbf{x} = 1\}.$$

*where* $\mathbf{1}$ *is the* $n$-*dimensional all-one vector.*

*Let* $\mathbf{a}$ *be a given* $n$-*dimensional vector and* $f : \mathbb{R}_+^n \to \mathbb{R}$ *be the function given by*

$$f(\mathbf{x}) = \sum_{i=1}^n x_i \log x_i - \sum_{i=1}^n a_i x_i.$$

(a) **(2 points)** *Prove that* $\Delta_n$ *is convex.*

(b) **(2 points)** *Let* $g : \mathbf{dom}(g) \to \mathbb{R}$ *be convex and differentiable. Prove that if there exists* $c \in \mathbb{R}$ *such that* $\nabla g(\mathbf{x}^*) = c\mathbf{1}$, *then* $\mathbf{x}^*$ *is a minimizer of* $g$ *over* $\Delta_n$.

(c) **(4 points)** *Prove that the reverse direction of Part (b) is also true. That is, if* $\mathbf{x}^*$ *is a minimizer of* $g$ *over* $\Delta_n$, *then there exists* $c \in \mathbb{R}$ *such that* $\nabla g(\mathbf{x}^*) = c\mathbf{1}$. *Hint: If two coordinates of* $\nabla g(\mathbf{x}^*)$ *are not the same then you may find a contradiction to one of the useful facts.*

(d) **(4 points)** *Prove that* $f(\mathbf{x})$ *is convex.*

(e) **(3 points)** *Prove that there exists an optimal solution to the problem* $\min\{f(\mathbf{x}) : \mathbf{x} \in \Delta_n\}$ *and compute that optimal solution.*

**Solution:**

(a) For $\mathbf{x}, \mathbf{y} \in \Delta_n$ and $\lambda \in [0, 1]$, we have $\mathbf{z} = \lambda \mathbf{x} + (1 - \lambda)\mathbf{y} > \mathbf{0}$ since $\mathbf{x}, \mathbf{y} > \mathbf{0}$, and

$$\mathbf{1}^\top \mathbf{z} = \lambda \mathbf{1}^\top \mathbf{x} + (1 - \lambda)\mathbf{1}^\top \mathbf{y} = 1.$$

Therefore, $\Delta_n$ is convex.

(b) We have

$$\nabla g(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) = c\mathbf{1}^\top (\mathbf{x} - \mathbf{x}^*) = c\mathbf{1}^\top \mathbf{x} - c\mathbf{1}\mathbf{x}^* = c - c = 0.$$

Hence, $\mathbf{x}^*$ is a minimizer of $g$ over $\Delta_n$, by Fact 1.

(c) Let $\mathbf{y}$ and $\mathbf{z}$ be the vectors which agree with $\mathbf{x}^*$ on all coordinates except for two coordinates $i, j$ where $y_i = x_i + t$, $y_j = x_j - t$, $z_i = x_i - t$ and $z_j = x_j + t$ for some $t > 0$ such that $t < \min\{x_i, x_j\}$. It is easy to see that $\mathbf{y}, \mathbf{z} \in \Delta_n$.

By Fact 1, we have for all $\mathbf{x} \in \Delta_n$,

$$\nabla g(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) \geq 0.$$

Consider $\mathbf{x} = \mathbf{y}$, we have

$$\nabla g(\mathbf{x}^*)^\top (\mathbf{y} - \mathbf{x}^*) = t \nabla g_i(\mathbf{x}^*) - t \nabla g_j(\mathbf{x}^*) \geq 0.$$

Similarly, consider $\mathbf{x} = \mathbf{z}$, we have

$$\nabla g(\mathbf{x}^*)^\top (\mathbf{z} - \mathbf{x}^*) = -t\nabla g_i(\mathbf{x}^*) + t\nabla g_j(\mathbf{x}^*) \geq 0.$$

From the above two equations, we have $\nabla g_i(\mathbf{x}^*) = \nabla g_j(\mathbf{x}^*)$.

As the above holds for any $i, j \in [n]$, it follows that all coordinates of $\nabla g(\mathbf{x}^*)$ have the same value. Hence, $\nabla g(\mathbf{x}^*) = c\mathbf{1}$ for some $c \in \mathbb{R}$.

(d) Firstly, observe that $\mathbf{dom}(f) = \mathbb{R}_+$ is convex. Next, we have

$$\frac{\partial f}{\partial x_i}(\mathbf{x}) = \log x_i + 1 - a_i.$$

Hence, for $i \in [n]$,

$$\frac{\partial^2 f}{\partial x_i \partial x_i}(\mathbf{x}) = \frac{1}{x_i}, \text{ and}$$

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x}) = 0 \text{ for } j \neq i.$$

Therefore, the eigenvalues of $\nabla^2 f(\mathbf{x})$ are $\frac{1}{x_i} > 0$. This implies that $\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$. By the second-order characterization of convexity, we can conclude that $f(\mathbf{x})$ is convex.

(e) Suppose there is an optimal solution $\mathbf{x}^*$. Applying (c), we have for some $c \in \mathbb{R}$ and any $i \in [n]$,

$$\log x_i^* + 1 - a_i = c.$$

Therefore, for any $i \in [n]$,

$$x_i^* = e^{c-1+a_i} = e^{c-1} e^{a_i}.$$

Since $\mathbf{1}^\top \mathbf{x}^* = 1$, we have

$$e^{c-1} = \frac{1}{\sum_{j=1}^n e^{a_j}}.$$

It follows that for any $i \in [n]$,

$$x_i^* = \frac{e^{a_i}}{\sum_{j=1}^n e^{a_j}}.$$

On the other hand, the value of $\mathbf{x}$, where

$$x_i = \frac{e^{a_i}}{\sum_{j=1}^n e^{a_j}},$$

satisfies that for any $i \in [n]$, $\log x_i + 1 - a_i = c$, where $c$ is defined above. Hence, applying part (b), we have that value $\mathbf{x}$ is an optimal solution. Therefore, there exists an optimal solution of $f$ over $\Delta_n$.

**Assignment 2.** **(9 points)** *Let* $f : \mathbb{R}^d \to \mathbb{R}$ *be a convex function, and let* $X$ *be a closed, convex set. Consider the following optimization problem*

$$\begin{aligned} minimize \quad & f(\mathbf{x}) \\ subject\ to \quad & \mathbf{x} \in X. \end{aligned} \tag{1}$$

*To solve this problem, we introduced Projected Gradient Descent with update rule for* $\gamma > 0$,

$$\mathbf{y}_{t+1} := \mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t)$$
$$\mathbf{x}_{t+1} := \Pi_X(\mathbf{y}_{t+1}) := \underset{\mathbf{x} \in X}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{y}_{t+1}\|^2$$

*as an optimization technique to solve this problem.*

*A potential issue is that we may get "stuck", i.e., for some* $t$*, we have* $\mathbf{x}_{t+1} = \mathbf{x}_t$*. In this assignment, you will show that this is actually not an issue, as when it occurs, we have found the optimal solution.*

*More specifically, prove that* $\mathbf{x}^*$ *is an optimal solution to (1) if and only if* $\mathbf{x}^* = \Pi_X(\mathbf{x}^* - \gamma \nabla f(\mathbf{x}^*))$*.*
*Hint: There are actually two constrained optimization problems.*

**Solution:** By Fact 1, $\mathbf{x}^*$ is the optimal solution of Problem (1), if and only if $\nabla f(\mathbf{x}^*)^\top (\mathbf{y} - \mathbf{x}^*) \geq 0$ for all $\mathbf{y} \in X$.

Let $g(\mathbf{x}) = \|\mathbf{x} - (\mathbf{x}^* - \gamma \nabla f(\mathbf{x}^*))\|^2$. Observe that $\Pi_X(\mathbf{x}^* - \gamma \nabla f(\mathbf{x}^*))$ is the minimizer of the problem

$$\min\{g(\mathbf{x}) : \mathbf{x} \in X\}. \tag{2}$$

Again, by Fact 1, a vector $\mathbf{w}$ is the optimal solution to (??), if and only if $\nabla g(\mathbf{w})^\top (\mathbf{y} - \mathbf{w}) \geq 0$ for all $\mathbf{y} \in X$.

Observe that $\nabla g(\mathbf{x}^*)^\top (\mathbf{y} - \mathbf{x}^*) = 2\gamma \nabla f(\mathbf{x}^*)^\top (\mathbf{y} - \mathbf{x}^*)$, for all $\mathbf{y} \in X$. Since $\gamma > 0$, $\nabla g(\mathbf{x}^*)^\top (\mathbf{y} - \mathbf{x}^*) \geq 0$ if and only if $\nabla f(\mathbf{x}^*)^\top (\mathbf{y} - \mathbf{x}^*) \geq 0$.

Putting all the above together, we obtain that the following statements are equivalent:

- $\mathbf{x}^*$ is the optimal solution of Problem (1);

- $\nabla f(\mathbf{x}^*)^\top (\mathbf{y} - \mathbf{x}^*) \geq 0$;

- $\nabla g(\mathbf{x}^*)^\top (\mathbf{y} - \mathbf{x}^*) \geq 0$; and

- $\mathbf{x}^*$ is the optimal solution of Problem (??), that is $\mathbf{x}^* = \Pi_X(\mathbf{x}^* - \gamma \nabla f(\mathbf{x}^*))$.

The result then follows.

**Assignment 3.** **(11 points)** *In class, we learned that the Newton's method with step size 1 can find the minimizer of a quadratic function in one step. In this assignment, we will investigate what happens when we vary the step size.*

*Consider the Newton's method with step size $\gamma$, where the update step is as follows:*

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma \nabla^2 f(\mathbf{x}_t)^{-1} \nabla f(\mathbf{x}_t), \qquad t \geq 0.$$

*Suppose $\mathbf{M} \in \mathbb{R}^{n \times n}$ is positive definite and $\mathbf{x}^* \in \mathbb{R}^n$ is a given vector. Consider the function*

$$f(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^\top \mathbf{M}(\mathbf{x} - \mathbf{x}^*).$$

*Answer the following questions and provide justifications for your answers.*

*(a)* **(2 points)** *Prove that $f(\mathbf{x})$ has only one minimizer and that minimizer is $\mathbf{x}^*$.*

*(b)* **(3 points)** *If we run the Newton's method with some initial vector $\mathbf{x}_0 \in \mathbb{R}^n$, for which values of $\gamma$ does the sequence $(\mathbf{x}_t)_{t \geq 0}$ converge to $\mathbf{x}^*$?*

*(c)* **(3 points)** *For the values of $\gamma$ in Part (b), how many iterations are required to obtain an $\mathbf{x} \in \mathbb{R}^n$ such that $\|\mathbf{x} - \mathbf{x}^*\| < \epsilon$ for some $\epsilon > 0$?*

*(d)* **(3 points)** *For the values of $\gamma$ in Part (b), how many iterations are required to obtain an $\mathbf{x} \in \mathbb{R}^n$ such that $f(\mathbf{x}) - f(\mathbf{x}^*) < \epsilon$ for some $\epsilon > 0$?*

**Solution:**

(a) Since $\mathbf{M}$ is positive definite, by the definition of positive definiteness, we have for $\mathbf{x} \neq \mathbf{x}^*$, $f(\mathbf{x}) > 0$. Further, $f(\mathbf{x}^*) = 0$. Hence, $f$ only has one minimizer $\mathbf{x}^*$.

(b) Calculating the gradient and the Hessian, we obtain

$$\nabla f(\mathbf{x}) = \mathbf{M}(\mathbf{x} - \mathbf{x}^*),$$
$$\nabla^2 f(\mathbf{x}) = \mathbf{M}.$$

Therefore, the update step becomes

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \mathbf{M}^{-1} \mathbf{M}(\mathbf{x} - \mathbf{x}^*)$$
$$= \mathbf{x}_t - \gamma(\mathbf{x} - \mathbf{x}^*).$$

Hence, $\mathbf{x}_{t+1} - \mathbf{x}^* = (\mathbf{x}_t - \mathbf{x}^*)(1 - \gamma)$. Applying this many times, we obtain

$$\mathbf{x}_{t+1} - \mathbf{x}^* = (\mathbf{x}_0 - \mathbf{x}^*)(1 - \gamma)^{t+1}.$$

Since $\mathbf{x}_0$ and $\mathbf{x}^*$ are given, the sequence $\{\mathbf{x}_t\}_{t \geq 0}$ converges to $\mathbf{x}^*$ if and only if $|1 - \gamma| < 1$ or $0 < \gamma < 2$.

(c) We have
$$\|\mathbf{x}_t - \mathbf{x}^*\| = (1-\gamma)^t \|\mathbf{x}_0 - \mathbf{x}^*\|.$$

For $\gamma = 1$, we note that $\mathbf{x}_1 = \mathbf{x}^*$, so we only need one step.

For $\gamma \neq 1$, in order to achieve $\|\mathbf{x}_t - \mathbf{x}^*\| < \epsilon$ for some t, we should have
$$(1-\gamma)^t \|\mathbf{x}_0 - \mathbf{x}^*\| < \epsilon.$$

This is equivalent to
$$t \geq \frac{1}{-\log(1-\gamma)} \log \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|}{\epsilon}.$$

(d)

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) = (1-\gamma)^{2t} (\mathbf{x}_0 - \mathbf{x}^*)^\top \mathbf{M} (\mathbf{x}_0 - \mathbf{x}^*).$$

Again, for $\gamma = 1$, we note that $\mathbf{x}_1 = \mathbf{x}^*$, so we only need one step.

For $\gamma \neq 1$, in order to achieve $f(\mathbf{x}_t) - f(\mathbf{x}^*) < \epsilon$, noting that $f(\mathbf{x}^*) = 0$, we should have
$$\frac{1}{2}(1-\gamma)^{2t} (\mathbf{x}_0 - \mathbf{x}^*)^\top \mathbf{M} (\mathbf{x}_0 - \mathbf{x}^*) \leq \epsilon.$$

This is equivalent to
$$t \geq \frac{1}{-\log((1-\gamma)^2)} \log \frac{(\mathbf{x}_0 - \mathbf{x}^*)^\top \mathbf{M} (\mathbf{x}_0 - \mathbf{x}^*)}{2\epsilon}.$$

**Assignment 4. (18 points)** *Let $n \geq 10^{10} \cdot d$ and suppose that $n/d \in \mathbb{N}$. Let $W \in \mathbb{R}^{n \times d}$ be a matrix with i.i.d. rows $w_1^\top, \ldots, w_n^\top$, where $w_1, \ldots, w_n \sim N(0, \Sigma)$ for some positive definite matrix $\Sigma \in \mathbb{R}^{d \times d}$ such that $\|\Sigma\| \geq 1$.*

  *(a)* **(7 points)** *Show that with probability at least 0.999,*

$$\left\| \frac{1}{n} \sum_{i \in [n]} w_i w_i^\top - \Sigma \right\| \leq 0.001 \cdot \|\Sigma\| .$$

  *(b)* **(5 points)** *Let $X \in \mathbb{R}^{n \times d}$ be a matrix with rows $x_1^\top, \ldots, x_n^\top$ obtained as follows: let $S_1, \ldots, S_d \subset [n]$ be disjoint sets, each of size $n/d$. For all $i \in S_j \subset [n]$, set row $x_i^\top = e_j^\top$ where $e_j$ is the $j$-th vector of the standard basis.*

   *Show that with probability at least 0.999,*

$$\left\| \frac{1}{n} \sum_{i \in [n]} w_i x_i^\top \right\| \leq 0.001 \cdot \sqrt{\|\Sigma\|} ,$$

   *where $w_1, \ldots, w_n \sim N(0, \Sigma)$ as above.*

  *(c)* **(6 points)** *Suppose $d \geq 1000$, $n \geq 10^{10} \cdot d$. Consider a matrix of the form $Y = X + W$, where $X, W$ are defined as above. Assuming Part (a) and Part (b), show that with probability at least 0.99,*

$$\left\| \frac{1}{n} \sum_{i \in [n]} y_i y_i^\top - \Sigma \right\| \leq 0.01 \cdot \|\Sigma\| ,$$

   *where $y_1^\top, \ldots, y_n^\top$ are the rows of $Y$.*

**Solution:**

  (a) Notice that

$$\frac{1}{n} \sum_{i \in [n]} w_i w_i^\top - \Sigma = \Sigma^{1/2} \left( \frac{1}{n} \sum_{i \in [n]} g_i g_i^\top - \mathrm{Id} \right) \Sigma^{1/2} ,$$

   where $g_1, \ldots, g_n \sim N(0, \mathrm{Id})$. Thus by Fact 3 with probability at least 0.999,

$$\left\| \Sigma^{1/2} \left( \frac{1}{n} \sum_{i \in [n]} g_i g_i^\top - \mathrm{Id} \right) \Sigma^{1/2} \right\| \leq \|\Sigma\| \left\| \frac{1}{n} \sum_{i \in [n]} g_i g_i^\top - \mathrm{Id} \right\| \leq 0.001 \cdot \|\Sigma\| .$$

  (b) Notice that

$$\frac{1}{n} \sum_{i \in [n]} w_i x_i^\top = \frac{1}{n} \Sigma^{1/2} \sum_{i \in [n]} g_i x_i^\top ,$$

where $g_1, \ldots, g_n \sim N(0, \mathrm{Id})$. A matrix $\sum_{i \in [n]} g_i x_i^\top$ has iid Gaussian entries with mean 0 and variance $n/d$. Hence by Fact 2 with probability at least 0.999,

$$\left\| \frac{1}{n} \sum_{i \in [n]} w_i x_i^\top \right\| \leq \frac{1}{n} \cdot \sqrt{\|\Sigma\|} \cdot 10\sqrt{d} \cdot \sqrt{\frac{n}{d}} \leq \frac{10}{\sqrt{n}} \cdot \sqrt{\|\Sigma\|} \leq 0.001 \cdot \sqrt{\|\Sigma\|} \, .$$

(c) Opening the product,

$$\frac{1}{n} \sum_{i \in [n]} y_i y_i^\top = \frac{1}{n} \sum_{i \in [n]} \left( x_i x_i^\top + w_i w_i^\top + x_i w_i^\top + w_i x_i^\top \right) \, .$$

So by triangle inequality,

$$\left\| \frac{1}{n} \sum_{i \in [n]} y_i y_i^\top - \Sigma \right\| \leq \left\| \frac{1}{n} \sum_{i \in [n]} w_i w_i^\top - \Sigma \right\| + \left\| \frac{1}{n} \sum_{i \in [n]} x_i x_i^\top \right\| + \left\| \frac{1}{n} \sum_{i \in [n]} x_i w_i^\top \right\| + \left\| \frac{1}{n} \sum_{i \in [n]} w_i x_i^\top \right\| \, .$$

We bound each term separately. Observe that by construction,

$$\frac{1}{n} \sum_{i \in [n]} x_i x_i^\top = \frac{1}{d} \mathrm{Id} \, .$$

Combining with the bounds from Parts (a) and (b) and the fact that $\|\Sigma\| \geq 1$, we get the desired result.

**Assignment 5. (8 points)** *Let $Z \in \mathbb{R}^{n \times d}$ be a matrix with linearly independent columns (known to the data analyst).*

*Consider the following family of an $n$-by-$d$ matrix-valued random variables*

$$\left\{ \mathbf{Y} = Z X + \mathbf{W} \mid X \in \mathbb{R}^{d \times d}, \; \mathbf{W}_{ij} \overset{i.i.d.}{\sim} N(0,1) \right\}.$$

*Given a realization of a random variable $\mathbf{Y} = Z X + \mathbf{W}$ in this family for unknown $X$, the goal is to estimate the matrix $ZX$ in Frobenius norm.*

*Let $\hat{X} : \mathbb{R}^{n \times d} \to \mathbb{R}^{d \times d}$ be an arbitrary estimator for this problem. Show that for every $\epsilon > 0$, there exists a random variable $\mathbf{Y} = ZX + \mathbf{W}$ in the above family such that*

$$\frac{1}{nd} \mathbb{E} \left\| ZX - Z\hat{X}(\mathbf{Y}) \right\|_F^2 \geq (1 - \epsilon) \cdot \frac{d}{n}.$$

*Remark: We have discussed statistical lower bounds for general estimators in the lectures and in the exercises. The lecture notes also discussed lower bounds for the special case of unbiased estimators. You can also get full points if you only prove the above lower bound for unbiased estimators.*

**Solution:** This problem is a special case of linear regression: if we consider $\mathbf{Y}$, $X$ and $\mathbf{W}$ as vectors $y' \in \mathbb{R}^{nd}$, $x' \in \mathbb{R}^{d^2}$ and $w' \in \mathbb{R}^{nd}$, then we can write $y' = Z'x' + w'$, where $Z'$ is an $nd \times d^2$ block-diagonal matrix with blocks equal to $Z$. Notice that $Z'$ has linearly independent columns, hence using the lower bound for linear regression, we get the desired bound.

**Assignment 6. (9 points)** *Consider the following variant of the matrix completion problem in which we see* $n \geq d \log d$ *independent samples of the form*

$$y_i = \langle X_i, \Theta^* + E \rangle + w_i \,,$$

*where* $X_i = e_a e_b^\top$ *for some* $a, b$ *sampled independently and uniformly at random from* $\{1, \ldots, d\}$, $\Theta^* \in \mathbb{R}^{d \times d}$ *satisfies* $\operatorname{rank}(\Theta^*) \leq r$, $\|\Theta^*\|_\infty \leq R$, $w_1, \ldots, w_n \sim N(0, \operatorname{Id}_d)$ *and* $E \in \mathbb{R}^{d \times d}$ *is an arbitrary matrix (possibly dependent on* $\Theta^*, X_1, \ldots, X_n, w_1, \ldots, w_n$*) satisfying* $\|E\|_\infty \leq \sqrt{\frac{R \cdot \log d}{n \cdot d}}$. *Note that the difference with the problem discussed in class is the additional error matrix* $E \in \mathbb{R}^{d \times d}$.

*Let*

$$Y := \frac{1}{n} \sum_{i \in [n]} y_i \cdot d^2 \cdot X_i \,.$$

*Consider the estimator*

$$\hat{\Theta} := \arg\min \left\{ \|\Theta - Y\|_F^2 \;\middle|\; \|\Theta\|_* \leq r, \|\Theta\|_\infty \leq R \right\}$$

*Show that with probability at least 0.99, despite the presence of* $E$,

$$\left\| \hat{\Theta} - \Theta^* \right\|_F^2 \leq O\left( \frac{d^3 \cdot \log d}{n} \cdot R \cdot r \right) \,.$$

*As shown in the lecture notes, you may directly assume the following facts*

**F.1** *It holds that* $\left\| \hat{\Theta} - \Theta^* \right\|_F^2 \leq 10r \left\| Y - \Theta^* \right\|^2$.

**F.2** *With probability at least 0.99, it holds that*

$$\left\| \frac{1}{n} \sum_{i \in [n]} \left( \langle X_i, \Theta^* \rangle + w_i \right) \cdot d^2 \cdot X_i - \Theta^* \right\|^2 \leq O\left( \frac{d^3 \log d}{n} \cdot R \right) \,.$$

**Solution:** We may rewrite

$$Y - \Theta^* = \frac{1}{n} \sum_{i \in [n]} y_i \cdot d^2 \cdot X_i - \Theta^*$$

$$= \frac{1}{n} \sum_{i \in [n]} \left( \left( \langle X_i, \Theta^* \rangle + w_i \right) \cdot d^2 \cdot X_i - \Theta^* \right) + \frac{1}{n} \sum_{i \in [n]} \langle X_i, E \rangle \cdot d^2 \cdot X_i \,.$$

Using the hint, it remains to bound the last term.

$$\left\| \frac{1}{n} \sum_{i \in [n]} \langle X_i, E \rangle \cdot d^2 \cdot X_i \right\| \leq \frac{1}{n} \sum_{i \in [n]} d^2 \cdot |\langle X_i, E \rangle| \leq d^2 \cdot \|E\|_\infty \leq \sqrt{\frac{d^3 \cdot \log d}{n} \cdot R} \,.$$

Putting everything together, the result follows.