

EN.553.662: Optimization for Data Science

Homework 4

Ha Manh Bui (CS Department)
hbui13@jhu.edu

Spring 2023

1 Problem 1

Let X be a real-valued random variable modeled with a normal distribution $\mathcal{N}(m, \sigma^2)$. Fixing a positive number ρ , we propose to estimate m and $\sigma^2 > 0$, based on a sample (x_1, \dots, x_N) by maximizing the penalized log-likelihood

$$\sum_{k=1}^N \log \varphi_{m, \sigma^2}(x_k) - N\rho \frac{|m|}{\sigma},$$

where φ_{m, σ^2} is the p.d.f. of $\mathcal{N}(m, \sigma^2)$, i.e.,

$$\varphi_{m, \sigma^2}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-m)^2}{2\sigma^2}\right).$$

We will assume in the following that the observed samples are not equal to the same constant value.

(1) Let μ_1 and μ_2 denote the first- and second-order empirical moments:

$$\mu_1 = \frac{1}{N}(x_1 + \dots + x_N)$$

$$\mu_2 = \frac{1}{N}(x_1^2 + \dots + x_N^2)$$

Prove that the pair (m, σ) is the optimal solution of the penalized likelihood problem if and only if there exists a scalar ξ such that (α, β, ξ) minimizes

$$F(\alpha, \beta, \xi) = \frac{1}{2}\alpha^2 - \mu_1\alpha\beta + \frac{1}{2}\mu_2\beta^2 - \log \beta + \rho\xi$$

subject to the constraints $\beta > 0$, $\xi - \alpha \geq 0$ and $\xi + \alpha \geq 0$, with $\alpha = m/\sigma$ and $\beta = 1/\sigma$.

Proof. Let denote

$$G : (m, \sigma) \mapsto \sum_{k=1}^N \log \varphi_{m, \sigma^2}(x_k) - N\rho \frac{|m|}{\sigma}.$$

We know that maximizing $G(m, \sigma)$ is equivalent to minimizing

$$-G(m, \sigma) = -\sum_{k=1}^N \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_k - m)^2}{2\sigma^2}\right) \right] + N\rho \frac{|m|}{\sigma} \quad (1)$$

$$\begin{aligned} &= \sum_{k=1}^N \left(\frac{(x_k - m)^2}{2\sigma^2} - \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) \right) + N\rho \frac{|m|}{\sigma} \\ &= \frac{1}{2\sigma^2} \sum_{k=1}^N x_k^2 - \frac{m}{\sigma^2} \sum_{k=1}^N x_k + \frac{Nm^2}{2\sigma^2} - N \log\left(\frac{1}{\sqrt{2\pi}} \frac{1}{|\sigma|}\right) + N\rho \frac{|m|}{\sigma}. \end{aligned} \quad (2)$$

Since N is positive and $\log(\frac{1}{\sqrt{2\pi}})$ is a constant, minimizes Equation 1 is equivalent to minimizing

$$\frac{1}{2\sigma^2} \frac{1}{N} \sum_{k=1}^N x_k^2 - \frac{m}{\sigma^2} \frac{1}{N} \sum_{k=1}^N x_k + \frac{m^2}{2\sigma^2} - \log\left(\frac{1}{|\sigma|}\right) + \rho \frac{|m|}{\sigma}. \quad (3)$$

Let consider $\frac{1}{\sigma} > 0$, then $\log(\frac{1}{|\sigma|}) = \log(\frac{1}{\sigma})$ and 3 becomes

$$\frac{1}{2\sigma^2} \frac{1}{N} \sum_{k=1}^N x_k^2 - \frac{m}{\sigma^2} \frac{1}{N} \sum_{k=1}^N x_k + \frac{m^2}{2\sigma^2} - \log\left(\frac{1}{\sigma}\right) + \rho \frac{|m|}{\sigma}.$$

Replace $\alpha = m/\sigma$ and $\beta = 1/\sigma$, we obtain

$$\frac{1}{2}\mu_2\beta^2 - \mu_1\alpha\beta + \frac{1}{2}\alpha^2 - \log\beta + \rho \frac{|m|}{\sigma} \text{ s.t. } \beta > 0. \quad (4)$$

Combining with $\rho > 0$, we obtain the optimal solution of minimizing 4, i.e., maximizing $G(m, \sigma)$ exists if and only if exist a scalar ξ s.t. (α, β, ξ) minimize

$$F(\alpha, \beta, \xi) = \frac{1}{2}\alpha^2 - \mu_1\alpha\beta + \frac{1}{2}\mu_2\beta^2 - \log\beta + \rho\xi \quad (5)$$

subject to the constraints $\beta > 0$, $\xi - \alpha \geq 0$ and $\xi + \alpha \geq 0$, with $\alpha = m/\sigma$ and $\beta = 1/\sigma$. \square

(2) Define $\hat{F}(\alpha, \beta, \xi) = F(\alpha, \beta, \xi)$ if $\beta > 0$ and $+\infty$ otherwise. Prove that F is a closed convex function.

Proof. Calculate gradient of F , we get

$$\nabla_{\alpha, \beta, \xi} F(\alpha, \beta, \xi) = \left(\frac{\partial F}{\partial \alpha}, \frac{\partial F}{\partial \beta}, \frac{\partial F}{\partial \xi} \right) = \left(\alpha - \mu_1\beta, -\mu_1\alpha + \mu_2\beta - \frac{1}{\beta}, \rho \right).$$

Therefore, the Hessian matrix of F is

$$\mathbf{H}_F = \begin{bmatrix} \frac{\partial^2 F}{\partial \alpha^2} & \frac{\partial^2 F}{\partial \alpha \partial \beta} & \frac{\partial^2 F}{\partial \alpha \partial \xi} \\ \frac{\partial^2 F}{\partial \beta \alpha} & \frac{\partial^2 F}{\partial \beta^2} & \frac{\partial^2 F}{\partial \beta \partial \xi} \\ \frac{\partial^2 F}{\partial \xi \alpha} & \frac{\partial^2 F}{\partial \xi \beta} & \frac{\partial^2 F}{\partial \xi^2} \end{bmatrix} = \begin{bmatrix} 1 & -\mu_1 & 0 \\ -\mu_1 & \mu_2 + \frac{1}{\beta^2} & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

Since \mathbf{H}_F is a symmetric matrix, $\det(\mathbf{H}_{F_1}) = 1$,

$$\det(\mathbf{H}_{F_2}) = \begin{vmatrix} \mu_2 + \frac{1}{\beta^2} & 0 \\ 0 & 1 \end{vmatrix} = \mu_2 + \frac{1}{\beta^2} > 0,$$

and

$$\det(\mathbf{H}_{F_3}) = \begin{vmatrix} 1 & -\mu_1 & 0 \\ -\mu_1 & \mu_2 + \frac{1}{\beta^2} & 0 \\ 0 & 0 & 1 \end{vmatrix} = \mu_2 + \frac{1}{\beta^2} - \mu_1^2 \geq 0 \text{ (Jensen's inequality)},$$

we obtain \mathbf{H}_F is C^2 and a positive semi-definite matrix, as a result, F is a convex function.

Continue to consider the Epigraphs of F w.r.t. β

$$\text{epi}(F) = \{(\beta, \hat{F}(\alpha, \beta, \xi)) \in \mathbb{R} \times \mathbb{R} : F(\beta) \leq \hat{F}(\alpha, \beta, \xi)\},$$

where $\beta \in (0, +\infty)$. Therefore, we have $\lim_n \beta_n = \beta$, $\lim_n \hat{F}(\alpha_n, \beta_n, \xi_n) = \hat{F}(\alpha, \beta, \xi)$, $F(\beta_n) \leq \hat{F}(\alpha_n, \beta_n, \xi_n)$, and $F(\beta) \leq \hat{F}(\alpha, \beta, \xi)$, i.e., $\text{epi}(F)$ is closed subset of $\mathbb{R} \times \mathbb{R}$. As a result, we obtain F is a closed convex function. \square

(3) Prove that there exists $\epsilon > 0$ that only depends on μ_1 and μ_2 such that, if $\beta \leq \epsilon$, then (α, β, ξ) cannot be an optimal solution of the problem in Question (1).

Proof. Assume $F(\alpha, \beta, \xi) \geq \frac{1}{2}(\mu_2 - \mu_1^2)\beta^2 - \log \beta$, i.e.,

$$\begin{aligned} & \frac{1}{2}\alpha^2 - \mu_1\alpha\beta + \frac{1}{2}\mu_2\beta^2 - \log \beta + \rho\xi \geq \frac{1}{2}(\mu_2 - \mu_1^2)\beta^2 - \log \beta \\ \Leftrightarrow & \frac{1}{2}\alpha^2 - \mu_1\alpha\beta + \rho\xi \geq -\frac{1}{2}\mu_1^2\beta^2 \\ \Leftrightarrow & \left(\frac{1}{\sqrt{2}}\right)^2 - 2\frac{1}{\sqrt{2}}\frac{1}{\sqrt{2}}\mu_1\alpha\beta + \left(\frac{1}{\sqrt{2}}\mu_1\beta\right)^2 + \rho\xi \\ \Leftrightarrow & \left(\frac{1}{\sqrt{2}}\alpha - \frac{1}{\sqrt{2}}\mu_1\beta\right)^2 + \rho\xi \geq 0 \text{ (always true since } \rho\xi \geq 0 \text{ due to } \xi - \alpha \geq 0 \text{ and } \xi + \alpha \geq 0). \end{aligned}$$

As a consequence, we obtain

$$F(\alpha, \beta, \xi) \geq \frac{1}{2}(\mu_2 - \mu_1^2)\beta^2 - \log \beta.$$

Consider $(\alpha, \beta, \xi) = (0, 1, 0)$, we get

$$F(\alpha, \beta, \xi) = F(0, 1, 0) \geq \frac{1}{2}(\mu_2 - \mu_1^2) \geq \frac{1}{2}(\mu_2 - \mu_1^2) - \log \beta, \forall \beta > 0. \quad (6)$$

Assume there exists $\epsilon > 0$ s.t. $\beta \leq \epsilon$, then

$$F(\alpha, \beta, \xi) = F(0, 1, 0) \geq \frac{1}{2}(\mu_2 - \mu_1^2) \geq \frac{1}{2}(\mu_2 - \mu_1^2) - \log \beta \geq \frac{1}{2}(\mu_2 - \mu_1^2) - \log \epsilon.$$

This means that (α, β, ξ) is no longer the minimizer of the F function. As a result, it can not be an optimal solution to the problem in Question (1). \square

(4) Prove that the minimization problem in Question (1) is then equivalent to minimizing F subject to the constraints $\beta \geq \epsilon$, $\xi - \alpha \geq 0$ and $\xi + \alpha \geq 0$, where ϵ satisfies the conditions of Question (3).

Proof. We have

$$\nabla_{\alpha, \beta, \xi} F(\alpha, \beta, \xi) = \left(\frac{\partial F}{\partial \alpha}, \frac{\partial F}{\partial \beta}, \frac{\partial F}{\partial \xi} \right) = \left(\alpha - \mu_1\beta, -\mu_1\alpha + \mu_2\beta - \frac{1}{\beta}, \rho \right).$$

Since $\nabla_{\alpha, \beta, \xi} F(0, 1, 0) = 0$, we obtain $(\alpha, \beta, \xi) = (0, 1, 0) \in \arg \min F$. On the other hand, since $\beta \geq \epsilon$ and $\epsilon > 0$ in Question (3), combining with the Inequality 6, we obtain

$$F(\alpha, \beta, \xi) = F(0, 1, 0) \geq \frac{1}{2}(\mu_2 - \mu_1^2) \geq \frac{1}{2}(\mu_2 - \mu_1^2) - \log \epsilon \geq \frac{1}{2}(\mu_2 - \mu_1^2) - \log \beta.$$

As a result, the minimization problem in Question (1) is then equivalent to minimizing F subject to the constraints $\beta \geq \epsilon$, $\xi - \alpha \geq 0$ and $\xi + \alpha \geq 0$, where $\epsilon > 0$. \square

(5) Prove that the KKT conditions characterize the solutions of the problem in Question (4), and that they can be reduced to:

$$\begin{cases} \alpha - \mu_1\beta = \lambda_2 - \lambda_1 \\ -\mu_1\alpha + \mu_2\beta - \frac{1}{\beta} = 0 \\ \rho - \lambda_1 - \lambda_2 = 0 \\ \lambda_1(\alpha - \xi) = 0 \\ \lambda_2(\alpha + \xi) = 0 \end{cases}$$

where $\lambda_1, \lambda_2 \geq 0$ are Lagrange multipliers.

Proof. Let the set constraints $\mathcal{C} = \mathcal{E} \cup \mathcal{I}$, where \mathcal{E} is equality and \mathcal{I} is inequality set constraints. Since the constraints in Question (4) include $\beta \geq \epsilon$, $\xi - \alpha \geq 0$ and $\xi + \alpha \geq 0$, where $\epsilon > 0$, we obtain the constraints \mathcal{C} follows

$$\begin{cases} -\beta < 0 \\ \alpha - \xi \leq 0 \\ -\xi - \alpha \leq 0 \end{cases}$$

Let the active constraints $A(\alpha, \beta, \xi) = \{i \in \mathcal{C}, \gamma_i(\alpha, \beta, \xi) = 0\}$ at $(\alpha, \beta, \xi) \in \Omega$, we get

$$\begin{cases} \gamma_1(\alpha, \beta, \xi) = \alpha - \xi = 0 \\ \gamma_2(\alpha, \beta, \xi) = -\xi - \alpha = 0 \end{cases}$$

So, we have $\nabla \gamma_1(\alpha, \beta, \xi) = (1, 0, -1)$ and $\nabla \gamma_2(\alpha, \beta, \xi) = (-1, 0, -1)$, therefore, $(\nabla \gamma_i(\alpha, \beta, \xi), i \in A(\alpha, \beta, \xi))$ are linearly independent, as a consequence, $(\alpha, \beta, \xi)^* \in \arg \min F$ satisfy MF-CQ and $\exists \lambda_i, i \in \mathcal{C}$ s.t. satisfy the KKT condition

$$\begin{cases} \nabla_{\alpha, \beta, \xi} L((\alpha, \beta, \xi)^*, \lambda) = 0 & \text{(I)} \\ \lambda_1 \geq 0, i \in \mathcal{I} & \text{(II)} \\ \lambda_i \gamma_i((\alpha, \beta, \xi)^*) = 0, i \in \mathcal{I} & \text{(III)} \end{cases}$$

Consider condition (I) $\nabla_{\alpha, \beta, \xi} L((\alpha, \beta, \xi)^*, \lambda) = 0$, we have

$$\begin{aligned} L((\alpha, \beta, \xi)^*, \lambda) &= F(\alpha, \beta, \xi) + \sum_{i \in \mathcal{C}} \lambda_i \gamma_i(\alpha, \beta, \xi) \\ &= \frac{1}{2} \alpha^2 - \mu_1 \alpha \beta + \frac{1}{2} \mu_2 \beta^2 - \log \beta + \rho \xi + \lambda_1 (\alpha - \xi) + \lambda_2 (-\xi - \alpha). \end{aligned}$$

Calculate gradient w.r.t. α, β, ξ , we get

$$\nabla_{\alpha, \beta, \xi} L((\alpha, \beta, \xi)^*, \lambda) = \left(\alpha - \mu_1 \beta + \lambda_1 - \lambda_2, -\mu_1 \alpha + \mu_2 \beta - \frac{1}{\beta}, \rho - \lambda_1 - \lambda_2 \right).$$

So, $\nabla_{\alpha, \beta, \xi} L((\alpha, \beta, \xi)^*, \lambda) = 0$ equivalent to

$$\begin{cases} \alpha - \mu_1 \beta = \lambda_2 - \lambda_1 \\ -\mu_1 \alpha + \mu_2 \beta - \frac{1}{\beta} = 0 \\ \rho - \lambda_1 - \lambda_2 = 0 \end{cases} \quad (7)$$

Consider condition (II) $\lambda_1 \geq 0, i \in \mathcal{I}$, we get

$$\begin{cases} \lambda_1 \geq 0 \\ \lambda_2 \geq 0 \end{cases} \quad (8)$$

Consider condition (III) $\lambda_i \gamma_i((\alpha, \beta, \xi)^*) = 0, i \in \mathcal{I}$, we get

$$\begin{cases} \lambda_1 (\alpha - \xi) = 0 \\ \lambda_2 (-\xi - \alpha) = 0 \end{cases} \Leftrightarrow \begin{cases} \lambda_1 (\alpha - \xi) = 0 \\ \lambda_2 (\alpha + \xi) = 0 \end{cases} \quad (9)$$

Combine the result of 7, 8, and 9, we obtain the KKT conditions characterize the solutions of the problem in Question (4), and that they can be reduced to:

$$\begin{cases} \alpha - \mu_1 \beta = \lambda_2 - \lambda_1 \\ -\mu_1 \alpha + \mu_2 \beta - \frac{1}{\beta} = 0 \\ \rho - \lambda_1 - \lambda_2 = 0 \\ \lambda_1 (\alpha - \xi) = 0 \\ \lambda_2 (\alpha + \xi) = 0 \end{cases}$$

where $\lambda_1, \lambda_2 \geq 0$ are Lagrange multipliers. □

(6) Find a necessary and sufficient condition on μ_1, μ_2 and ρ for a solution to this system to satisfy $\alpha = 0$, and provide the optimal β in that case.

Solution. The necessary and sufficient condition on μ_1, μ_2 and ρ for a solution to this system to satisfy $\alpha = 0$ follows

$$\begin{cases} -\mu_1 \beta = \lambda_2 - \lambda_1 \\ \mu_2 \beta - \frac{1}{\beta} = 0 \\ \rho - \lambda_1 - \lambda_2 = 0 \\ -\lambda_1 \xi = 0 \\ \lambda_2 \xi = 0 \\ \lambda_1, \lambda_2, \xi \geq 0 \\ \rho, \beta > 0 \end{cases} \Leftrightarrow \begin{cases} \mu_1 = \frac{\lambda_1 - \lambda_2}{\beta} \\ \mu_2 = \frac{1}{\beta^2} \\ \rho = \lambda_1 + \lambda_2 \end{cases}$$

where $\lambda_1, \lambda_2 \geq 0$. Since the condition $\beta > 0$, we obtain the optimal of β in this case is $\beta = \frac{1}{\sqrt{\mu_2}}$.

(7) Assume that the condition in Question (6) is not satisfied.

(a) Prove that in that case, α has the same sign as μ_1 .

Proof. Let's consider the condition in Question (6) is not satisfied, i.e., $\alpha \neq 0$, then there are two cases $\alpha < 0$ and $\alpha > 0$. From the KKT conditions in Question (4), α must satisfy the conditions that follow

$$\begin{cases} \lambda_1(\alpha - \xi) = 0 \\ \lambda_2(\alpha + \xi) = 0 \\ \lambda_1 + \lambda_2 = \rho > 0 \\ \lambda_1, \lambda_2, \xi \geq 0 \end{cases}$$

Therefore, the solution to these equations are

$$\begin{cases} \alpha = \xi > 0 \\ \lambda_2 = 0 \\ \lambda_1 > 0 \end{cases} \quad \text{or} \quad \begin{cases} \alpha = -\xi < 0 \\ \lambda_1 = 0 \\ \lambda_2 > 0 \end{cases} \quad (10)$$

Replace these solutions to the KKT condition $\alpha - \mu_1\beta = \lambda_2 - \lambda_1$, we obtain

$$\beta = \begin{cases} \frac{\alpha + \lambda_1}{\mu_1} & \text{with } \alpha, \lambda_1 > 0 \\ \frac{\alpha - \lambda_2}{\mu_1} & \text{with } \alpha < 0, \lambda_2 > 0 \end{cases}$$

Assume α and μ_1 has a different sign, then $\beta \leq 0$ for both case and contradict with the condition that $\beta > 0$. As a consequence, we obtain α has the same sign as μ_1 . \square

(b) Completely describe the solution of the system, separating the cases $\mu_1 > 0$ and $\mu_1 < 0$.

Solution. Recall the system from the KKT condition,

$$\begin{cases} \alpha - \mu_1\beta = \lambda_2 - \lambda_1 & \text{(I)} \\ -\mu_1\alpha + \mu_2\beta - \frac{1}{\beta} = 0 & \text{(II)} \\ \rho - \lambda_1 - \lambda_2 = 0 & \text{(III)} \\ \lambda_1(\alpha - \xi) = 0 & \text{(IV)} \\ \lambda_2(\alpha + \xi) = 0 & \text{(V)} \end{cases}$$

where $\lambda_1, \lambda_2 \geq 0$. From (I), we have $\alpha = \lambda_2 - \lambda_1 + \mu_1\beta$. Replace this α to (II), and we get

$$\begin{aligned} & -\mu_1(\lambda_2 - \lambda_1 + \mu_1\beta) + \mu_2\beta - \frac{1}{\beta} = 0 \\ \Leftrightarrow & (\mu_2 - \mu_1^2)\beta + (\lambda_1 - \lambda_2)\mu_1 - \frac{1}{\beta} = 0 \\ \Leftrightarrow & (\mu_2 - \mu_1^2)\beta^2 + (\lambda_1 - \lambda_2)\mu_1\beta - 1 = 0 \quad (\text{due to } \beta > 0). \end{aligned}$$

Calculate Δ from this Equation, we have

$$\Delta = (\lambda_1 - \lambda_2)^2\mu_1^2 + 4(\mu_2 - \mu_1^2).$$

Therefore, we obtain the solution of β includes

$$\beta = \begin{cases} \frac{(\lambda_2 - \lambda_1)\mu_1 - \sqrt{(\lambda_1 - \lambda_2)^2\mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)} \\ \frac{(\lambda_2 - \lambda_1)\mu_1 + \sqrt{(\lambda_1 - \lambda_2)^2\mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)} \end{cases}$$

Let's consider $\mu_1 > 0$, using the result $\alpha = \xi > 0, \lambda_2 = 0, \lambda_1 > 0$ from 10, we obtain

$$\beta = \begin{cases} \frac{-\lambda_1\mu_1 - \sqrt{\lambda_1^2\mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)} \leq 0 \quad (\text{contradiction with } \beta > 0) \\ \frac{-\lambda_1\mu_1 + \sqrt{\lambda_1^2\mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)} \end{cases}$$

Replace this β to (I), we obtain

$$\alpha = -\lambda_1 + \mu_1 \frac{-\lambda_1 \mu_1 + \sqrt{\lambda_1^2 \mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)}.$$

As a consequence, for $\mu_1 > 0$, we have the solution of the system from KKT conditions are

$$\begin{cases} \alpha = -\lambda_1 + \mu_1 \frac{-\lambda_1 \mu_1 + \sqrt{\lambda_1^2 \mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)} \\ \beta = \frac{-\lambda_1 \mu_1 + \sqrt{\lambda_1^2 \mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)} \\ \xi = -\lambda_1 + \mu_1 \frac{-\lambda_1 \mu_1 + \sqrt{\lambda_1^2 \mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)}. \end{cases}$$

Similarly, consider $\mu_1 < 0$, using the result $\alpha = -\xi < 0, \lambda_1 = 0, \lambda_2 > 0$ from 10, we obtain

$$\beta = \begin{cases} \frac{\lambda_2 \mu_1 - \sqrt{\lambda_2^2 \mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)} \leq 0 \text{ (contradiction with } \beta > 0) \\ \frac{\lambda_2 \mu_1 + \sqrt{\lambda_2^2 \mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)} \end{cases}$$

Replace this β to (I), we obtain

$$\alpha = \lambda_2 + \mu_1 \frac{\lambda_2 \mu_1 + \sqrt{\lambda_2^2 \mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)}.$$

As a consequence, for $\mu_1 < 0$, we have the solution of the system from KKT conditions are

$$\begin{cases} \alpha = -\lambda_2 + \mu_1 \frac{\lambda_2 \mu_1 + \sqrt{\lambda_2^2 \mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)} \\ \beta = \frac{\lambda_2 \mu_1 + \sqrt{\lambda_2^2 \mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)} \\ \xi = -\lambda_2 - \mu_1 \frac{\lambda_2 \mu_1 + \sqrt{\lambda_2^2 \mu_1^2 + 4(\mu_2 - \mu_1^2)}}{2(\mu_2 - \mu_1^2)}. \end{cases}$$

(8) Summarize this discussion and prove that one of the two following statements holds (with extra credit if you prove both, i.e., prove that they give the same solution).

(a) The optimal parameters can be computed as follows

- If $\mu_1^2 \leq \rho^2 \mu_2$: $\hat{m} = 0, \hat{\sigma} = \sqrt{\mu_2}$.
- If $\mu_1^2 \geq \rho^2 \mu_2$:

$$\begin{aligned} \hat{m} &= \mu_1 - \text{sign}(\mu_1) \rho \hat{\sigma} \\ \hat{\sigma} &= \frac{2s^2}{\sqrt{\rho^2 \mu_1^2 + 4s^2} - \rho |\mu_1|}, \end{aligned}$$

with $s^2 = \mu_2 - \mu_1^2$.

Proof. Consider $\mu_1^2 \leq \rho^2 \mu_2$, from the result of Question (6), we know $\mu_1 = \frac{\lambda_1 - \lambda_2}{\beta}$, $\mu_2 = \frac{1}{\beta^2}$, and $\rho = \lambda_1 + \lambda_2$ is necessary and sufficient condition for a solution of the system from KKT conditions to satisfy $\alpha = 0$ with the optimal of $\beta = \frac{1}{\sqrt{\mu_2}}$. So, replace these condition into $\mu_1^2 \leq \rho^2 \mu_2$, we obtain

$$\frac{(\lambda_1 - \lambda_2)^2}{\sigma^2} \leq \frac{(\lambda_1 + \lambda_2)^2}{\sigma^2} \text{ (true for } \lambda_1, \lambda_2 \geq 0). \quad (11)$$

As a result, $\mu_1^2 \leq \rho^2 \mu_2$ satisfy $\alpha = 0$. Since $\alpha = \frac{\hat{m}}{\hat{\sigma}}$ and $\beta = \frac{1}{\hat{\sigma}}$, with $\alpha = 0$ and $\beta = \frac{1}{\sqrt{\mu_2}}$ for $\mu_1^2 \leq \rho^2 \mu_2$, we obtain

$$\hat{m} = 0, \hat{\sigma} = \sqrt{\mu_2}.$$

Consider $\mu_1^2 \geq \rho^2 \mu_2$, Inequality 11 and the result above implies $\mu_1^2 \geq \rho^2 \mu_2$ satisfy $\alpha > 0$ and $\alpha < 0$. Since $\alpha = \frac{\hat{m}}{\hat{\sigma}}$ and $\beta = \frac{1}{\hat{\sigma}}$, combining with the result of α, β from Question (7) for both case $\alpha > 0$ and $\alpha < 0$, we obtain

$$\begin{aligned}\hat{m} &= \mu_1 - \text{sign}(\mu_1)\rho\hat{\sigma} \\ \hat{\sigma} &= \frac{2s^2}{\sqrt{\rho^2\mu_1^2 + 4s^2} - \rho|\mu_1|},\end{aligned}$$

with $s^2 = \mu_2 - \mu_1^2$. □

(b) The optimal parameters are

$$\begin{aligned}\hat{m} &= \text{sign}(\mu_1) \max(0, |\mu_1| - \rho\hat{\sigma}) \\ \hat{\sigma} &= \min\left(\sqrt{\mu_2}, \frac{2s^2}{\sqrt{\rho^2\mu_1^2 + 4s^2} - \rho|\mu_1|}\right)\end{aligned}$$

with $s^2 = \mu_2 - \mu_1^2$.

Proof. From Question (6) and (7), we know there are three cases for the solution of the system from KKT conditions, which includes

$$\begin{cases} \hat{m} = 0, \hat{\sigma} = \sqrt{\mu_2}, & \alpha = 0 \\ \hat{m} = \mu_1 - \text{sign}(\mu_1)\rho\hat{\sigma}, \hat{\sigma} = \frac{2s^2}{\sqrt{\rho^2\mu_1^2 + 4s^2} - \rho|\mu_1|}, & \alpha > 0 \text{ and } \alpha < 0 \end{cases} \quad (12)$$

Since α has the same sign with μ_1 if $\alpha \neq 0$ (result from Question (7)), the system 12 equivalents

$$\begin{cases} \hat{m} = 0, \hat{\sigma} = \sqrt{\mu_2}, & \alpha = 0 \\ \hat{m} = \text{sign}(\mu_1)(|\mu_1| - \rho\hat{\sigma}), \hat{\sigma} = \frac{2s^2}{\sqrt{\rho^2\mu_1^2 + 4s^2} - \rho|\mu_1|}, & \alpha > 0 \text{ and } \alpha < 0 \end{cases}$$

Therefore, to maximize the penalized log-likelihood

$$G(m, \sigma) = \sum_{k=1}^N \log \left[\frac{1}{\sqrt{2\pi}\sigma^2} \exp \left(-\frac{(x_k - m)^2}{2\sigma^2} \right) \right] + N\rho \frac{|m|}{\sigma},$$

the optimal parameter must be the maximum of m and minimum of σ , i.e.,

$$\begin{aligned}\hat{m} &= \text{sign}(\mu_1) \max(0, |\mu_1| - \rho\hat{\sigma}) \\ \hat{\sigma} &= \min\left(\sqrt{\mu_2}, \frac{2s^2}{\sqrt{\rho^2\mu_1^2 + 4s^2} - \rho|\mu_1|}\right)\end{aligned}$$

with $s^2 = \mu_2 - \mu_1^2$. □

(9) What happens when $\rho \geq 1$?

Solution. Since $s^2 = \mu_2 - \mu_1^2$ and $s^2 \geq 0$, we get $\mu_1^2 \leq \mu_2$. If $\rho \geq 1$, let's consider two cases of condition for μ_1, μ_2 , we have

$$\begin{cases} \mu_1^2 \leq \rho^2 \mu_2 \\ \mu_1^2 \geq \rho^2 \mu_2 \text{ contradict with } \mu_1^2 \leq \mu_2 \end{cases}$$

Therefore, only $\mu_1^2 \leq \rho^2 \mu_2$ satisfy $\rho \geq 1$. Using the result from Question (8), we obtain the corresponding optimal parameter of the objective function when $\mu_1^2 \leq \rho^2 \mu_2$, i.e., $\rho \geq 1$ is

$$\hat{m} = 0, \quad \hat{\sigma} = \sqrt{\mu_2}.$$

(10) Program the estimator described in Question (8), and apply it to each of the $K = 500$ columns in the file project4_Gaussians.csv (each column has dimension $N = 100$). Let m_k, σ_k denote the mean and standard deviation estimated for column k .

Provide, sorted in increasing order, the indexes of the columns for which $m_k \neq 0$. Also provide a histogram of $(\sigma_1, \dots, \sigma_k)$. Take $\rho = 0.25$.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

def estimate_Gaussians(X, rho):
    Mu, Sigma = [], []
    Mu_1 = np.mean(X, axis = 0)
    Mu_2 = np.mean(np.square(X), axis = 0)
    S_square = Mu_2 - np.square(Mu_1)

    tmp_1 = np.sqrt(Mu_2)
    tmp_2 = (2*S_square)/(np.sqrt(rho**2 * Mu_1**2 + 4*S_square) - rho*np.abs(Mu_1))
    for i in range(len(tmp_1)):
        tmp_s = min(tmp_1[i], tmp_2[i])
        Sigma.append(tmp_s)
        Mu.append(np.sign(Mu_1[i]) * max(0, np.abs(Mu_1[i]) - rho * tmp_s))

    return np.array(Mu), np.array(Sigma)

if __name__ == "__main__":
    X = pd.read_csv('homework4_data/project4_Gaussians.csv')
    X = X.drop(['Unnamed: 0'], axis = 1).to_numpy()
    X = np.squeeze(X)
    Mu, Sigma = estimate_Gaussians(X, 0.25)
    print("The indexes of the columns for which " + r'$m_k \neq 0$' + " is: ")
    for i in range(len(Mu)):
        if Mu[i] != 0:
            print(i, end = " ")
    print()

    plt.hist(Sigma)
    plt.xlabel('Standard deviation ' + r'$\sigma$' + ' value')
    plt.ylabel('Frequency')
    plt.title("Histogram of " + r'$(\sigma_1, \dots, \sigma_k)$')
    plt.savefig("problem1.pdf")
```

Result:

The indexes of the columns for which $m_k \neq 0$ is:

1 23 35 37 61 69 115 159 183 247 252 254 263 275 277 279 304 343 348 355 369 372 380 419 423 424 453
467 475 478 483 489

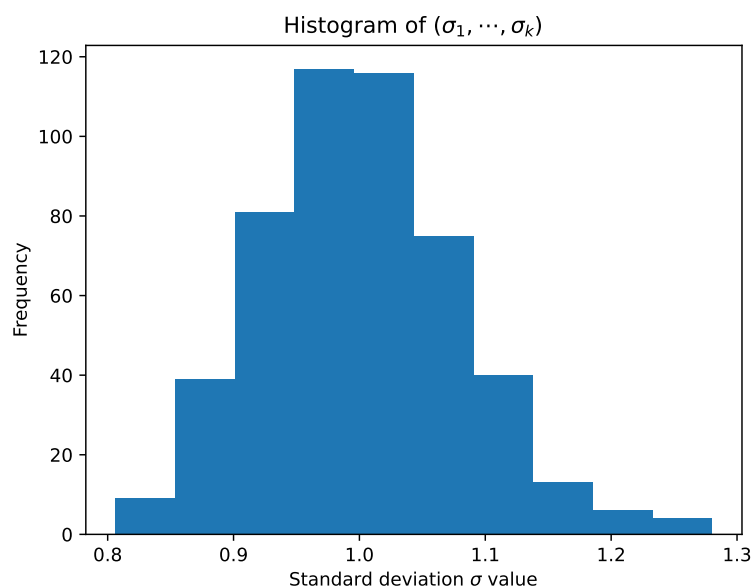


Figure 1: Histogram of $(\sigma_1, \dots, \sigma_k)$.

2 Problem 2

Fix a number $c > 0$ and let Ω_c be the ℓ^1 -norm ball with radius c ,

$$\Omega_c = \{x \in \mathbb{R}^n : |x^{(1)}| + \cdots + |x^{(n)}| \leq c\}.$$

One can prove (see Problem set 3, Problem V) that, for $y \in \mathbb{R}^n$, $y \notin \Omega_c$, $x = \text{proj}_{\Omega_c}(y)$ is such that

$$x^{(j)} = \text{sign}(y^{(j)}) \max(|y^{(j)}| - \lambda, 0)$$

with

$$f_y(\lambda) := \sum_{j=1}^n \max(|y^{(j)}| - \lambda, 0) = c$$

(1) Define

- $\lambda_0 = \max\{|y^{(j)}| : f_y(|y^{(j)}|) \geq c\}$ if there exists j such that $f_y(|y^{(j)}|) \geq c$, and $\lambda_0 = 0$ otherwise,
- $\lambda_1 = \min\{|y^{(j)}| : f_y(|y^{(j)}|) \leq c\}$.

Justify the existence of λ_0, λ_1 when $y \notin \Omega_c$ and prove that the value of λ such that $f_y(\lambda) = c$ is then given by

$$\lambda = \begin{cases} \lambda_0 + \frac{c - f_y(\lambda_0)}{f_y(\lambda_0) - f_y(\lambda_1)}(\lambda_0 - \lambda_1) & \text{if } \lambda_0 < \lambda_1 \\ \lambda_0 & \text{if } \lambda_0 = \lambda_1 \end{cases}$$

Proof. When $y \notin \Omega$, if there exists j such that $f_y(|y^{(j)}|) \geq c$, then we have

$$\{|y^{(j)}| : f_y(|y^{(j)}|) \geq c\} \neq \emptyset$$

therefore, if $\lambda_0 = \max\{|y^{(j)}| : f_y(|y^{(j)}|) \geq c\}$, then λ_0 always exists. Otherwise, with $\lambda_0 = 0$, then λ_0 also exists.

On the other hand, consider $a, b \in \{1, \dots, n\}$, due to $|y^{(a)}|, |y^{(b)}| \geq 0$, then if $|y^{(a)}| \leq |y^{(b)}|$, we always have

$$f_y(|y^{(a)}|) \geq f_y(|y^{(b)}|).$$

Therefore, for $|y^{(j)}| \geq 0, \forall j \in \{1, \dots, n\}$, we obtain

$$f_y(0) \geq |y^{(j)}| \Rightarrow \{|y^{(j)}| : f_y(|y^{(j)}|) \leq c\} \neq \emptyset.$$

As a result, if $\lambda_1 = \min\{|y^{(j)}| : f_y(|y^{(j)}|) \leq c\}$, then λ_1 always exists.

If $\lambda_0 = \lambda_1$, then we have

$$\max\{|y^{(j)}| : f_y(|y^{(j)}|) \geq c\} = \min\{|y^{(j)}| : f_y(|y^{(j)}|) \leq c\} \Leftrightarrow f_y(|y^{(j)}|) = c \Rightarrow \lambda = \lambda_0 = \lambda_1. \quad (13)$$

If $\lambda_0 < \lambda_1$, we have

$$f_y(\lambda_0) - f_y(\lambda_1) = \sum_{j=1}^n \max(|y^{(j)}| - \lambda_0 - |y^{(j)}| + \lambda_1, 0) = \sum_{j=1}^n \max(\lambda_1 - \lambda_0, 0) = n(\lambda_1 - \lambda_0).$$

Replace this equation with the expressions

$$\lambda_0 + \frac{c - f_y(\lambda_0)}{f_y(\lambda_0) - f_y(\lambda_1)}(\lambda_0 - \lambda_1)$$

we get

$$\begin{aligned} \lambda_0 + \frac{c - f_y(\lambda_0)}{f_y(\lambda_0) - f_y(\lambda_1)}(\lambda_0 - \lambda_1) &= \lambda_0 + \frac{1}{n}(f_y(\lambda_0) - c) \\ &= \lambda_0 + \frac{1}{n} \sum_{i=1}^n \max(|y^{(i)}| - \lambda_0 - |y^{(i)}| + \lambda, 0) \\ &= \lambda_0 + \frac{1}{n} \sum_{i=1}^n \max(\lambda - \lambda_0, 0). \end{aligned} \quad (14)$$

Let's consider $\lambda \geq \lambda_0$, from 14, we obtain

$$\lambda_0 + \frac{1}{n} \sum_{i=1}^n \max(\lambda - \lambda_0, 0) = \lambda_0 + \frac{1}{n} n(\lambda - \lambda_0) = \lambda. \quad (15)$$

Otherwise, if $\lambda < \lambda_0$, we get

$$\sum_{i=1}^n \max(\lambda - \lambda_0, 0) = 0 \Rightarrow f_y(\lambda_0) = c \Rightarrow \lambda = \lambda_0 \Rightarrow \text{contradict with } \lambda < \lambda_0.$$

As a consequence, combining the result from 14 and 15, if $\lambda_0 < \lambda_1$, then

$$\lambda = \lambda_0 + \frac{c - f_y(\lambda_0)}{f_y(\lambda_0) - f_y(\lambda_1)} (\lambda_0 - \lambda_1).$$

Combining with the result from 13, we obtain

$$\lambda = \begin{cases} \lambda_0 + \frac{c - f_y(\lambda_0)}{f_y(\lambda_0) - f_y(\lambda_1)} (\lambda_0 - \lambda_1) & \text{if } \lambda_0 < \lambda_1 \\ \lambda_0 & \text{if } \lambda_0 = \lambda_1 \end{cases}$$

□

(2) Write a program that takes y and c as inputs and returns $proj_{\Omega_c}(y)$ using this solution. Apply this program to the vector in the file `project4_vector.csv`, which has dimension 10000. Return, for $c \in \{5, 10, 20, 100\}$, the value of $|y - proj_{\Omega_c}(y)|^2$.

```
import numpy as np
import pandas as pd

def func_f(y, lmda):
    out = np.abs(y) - lmda
    out = out * (out >= 0)
    return np.sum(out)

def find_lambda(y, c):
    list_lmd_0, list_lmd_1 = [], []
    for i in range(len(y)):
        tmp = func_f(y, np.abs(y[i]))
        if tmp >= c:
            list_lmd_0.append(np.abs(y[i]))
        if tmp <= c:
            list_lmd_1.append(np.abs(y[i]))
    lmda_1 = min(list_lmd_1)
    if len(list_lmd_0) != 0:
        lmda_0 = max(list_lmd_0)
    else:
        lmda_0 = 0
    if lmda_0 == lmda_1:
        return lmda_0
    return lmda_0 + ((c - func_f(y, lmda_0)) / (func_f(y, lmda_0) - func_f(y, lmda_1))) * (lmda_0 - lmda_1)

def project(y, c):
    lmda = find_lambda(y, c)
    x = np.abs(y) - lmda
    x = x * (x >= 0)
    x = np.sign(y) * x
    return x

def problem2.2():
    y = pd.read_csv('homework4_data/project4_vector-1.csv')
    y = y.drop(['Unnamed: 0'], axis = 1).to_numpy()
    y = np.squeeze(y)
    list_c = [5, 10, 20, 100]
    for c in list_c:
        out = np.linalg.norm(y - project(y, c))**2
        print("The value of " + r'$|y-proj_{\Omega_c}(y)|^2$' + " for c = " + str(c) + " is: " + str(out))

if __name__ == "__main__":
    problem2.2()
```

Result:

The value of $|y - \text{proj}_{\Omega_c}(y)|^2$ for $c = 5$ is: 9797.346376053942

The value of $|y - \text{proj}_{\Omega_c}(y)|^2$ for $c = 10$ is: 9768.217407278022

The value of $|y - \text{proj}_{\Omega_c}(y)|^2$ for $c = 20$ is: 9713.883606540363

The value of $|y - \text{proj}_{\Omega_c}(y)|^2$ for $c = 100$ is: 9340.632835889344

(3) Let $u \in \mathbb{R}^n$ be given. Define $g_u(x) = u^\top x$. Let $J = \arg \max\{|u^{(j)}|, j = 1, \dots, n\}$. Prove that

$$\arg \min_{\Omega_c} g_u = \left\{ -c \sum_{j \in J} \text{sign}(u^{(j)}) \rho^{(j)} e_j : \sum_{j \in J} \rho^{(j)} = 1, \rho^{(j)} \geq 0, j \in J \right\}$$

where e_1, \dots, e_n form the canonical basis of \mathbb{R}^n .

Proof. We have $\Omega_c = \{x \in \mathbb{R}^n : |x^{(1)}| + \dots + |x^{(n)}| \leq c\}$, with $c > 0$, so, Ω_c being defined by affine inequalities, the KKT conditions apply with the Lagrangian

$$\begin{aligned} L(x, \lambda) &= u^\top x + \lambda(|x^{(1)}| + \dots + |x^{(n)}| - c) \\ &= u^{(1)}x^{(1)} + \dots + u^{(n)}x^{(n)} + \lambda(|x^{(1)}| + \dots + |x^{(n)}| - c). \end{aligned}$$

Calculate the gradient, and we

$$\nabla_x L(x, \lambda) = \left(\frac{\partial L}{\partial x^{(1)}}, \dots, \frac{\partial L}{\partial x^{(n)}} \right) = \left(u^{(1)} + \lambda \frac{x^{(1)}}{|x^{(1)}|}, \dots, u^{(n)} + \lambda \frac{x^{(n)}}{|x^{(n)}|} \right).$$

If $x \in \arg \min_{\Omega_c} g_u$, then x must satisfies the KKT conditions

$$\begin{cases} \nabla_x L(x, \lambda) = 0 & \text{(I)} \\ \lambda \geq 0 & \text{(II)} \\ \lambda(|x^{(1)}| + \dots + |x^{(n)}| - c) = 0 & \text{(III)} \end{cases}$$

From (I), we have $u^{(i)} + \lambda \frac{x^{(i)}}{|x^{(i)}|} = 0, \forall i \in \{1, \dots, n\}$, replace into (III) and we get

$$\lambda \left(\frac{-\lambda x^{(1)}}{u^{(1)}} - \dots - \frac{\lambda x^{(n)}}{u^{(n)}} - c \right) = 0. \quad (16)$$

Combining condition (II) with the result from Question (1), we have $\lambda = \max\{|u^{(j)}| : f_u(|u^{(j)}|) \geq c\}$ if $f_u(|u^{(j)}|) \geq c$, otherwise, $\lambda = 0$. So, if $\lambda > 0$, we have 16 equivalents

$$\begin{aligned} & -\lambda \left(\frac{x^{(1)}}{u^{(1)}} + \dots + \frac{x^{(n)}}{u^{(n)}} \right) - c = 0 \\ \Rightarrow x^{(i)} &= u^{(i)} \left(\frac{-c}{\lambda} - \sum_{j=1, j \neq i}^n \frac{x^{(j)}}{u^{(j)}} \right) \\ \Leftrightarrow x^{(i)} &= -c \left(\frac{u^{(i)}}{\lambda} + \sum_{j=1, j \neq i}^n \frac{u^{(i)}x^{(j)}}{u^{(j)}c} \right), \forall i \in \{1, \dots, n\}. \end{aligned} \quad (17)$$

Replace $J = \arg \max\{|u^{(j)}|, j = 1, \dots, n\}$ into 17, we obtain

$$x = \arg \min_{\Omega_c} g_u = \left\{ -c \sum_{j \in J} \text{sign}(u^{(j)}) \rho^{(j)} e_j : \sum_{j \in J} \rho^{(j)} = 1, \rho^{(j)} \geq 0, j \in J \right\}$$

where e_1, \dots, e_n form the canonical basis of \mathbb{R}^n . □

(4) Let $\alpha_t = 2/(2+t)$. Prove that the iterations

$$\begin{cases} x_{t+1} = (1 - \alpha_t)x_t - c\alpha_t \text{sign}(x_t^{(j_t)} - y^{(j_t)})e_{j_t} \\ j_t \in \arg \max\{|x_t^{(j)} - y^{(j)}|, j = 1, \dots, n\} \end{cases}$$

with $x_0 = 0$, converges to $x = \text{proj}_{\Omega_c} y$.

Proof. Using the result from Question (3), we have the iterations equivalents to

$$\begin{aligned} x_{t+1} &= (1 - \alpha_t)x_t + \alpha_t \arg \min_{\Omega_c} (x \mapsto (x_t - y)^\top x) \\ \Leftrightarrow x_{t+1} &= (1 - \alpha_t)x_t + \alpha_t \arg \min_{\Omega_c} ((x_t - y)^\top (x - x_t)). \end{aligned}$$

Let $\nabla F(x_t) = x_t - y$, then we have $x^* = \text{proj}_{\Omega_c} y$, with $\alpha_t = 2/(2+t)$, we always have

$$(1 - \rho\alpha_t)\alpha_t \leq \alpha_{t+1} \leq \alpha_t$$

for $\rho = 1/2$. Therefore, apply the Theorem of the conditional gradient algorithm, we get

$$F(x_t) - F(x^*) \leq \frac{2LD^2}{2+t}, t \geq 1$$

with $D = \max\{|x - y| : x, y \in \Omega_c\}$. As a consequence, with $x_0 = 0$, we obtain x_t converges to $\text{proj}_{\Omega_c} y$. \square

(5) Program this algorithm, taking as input y and c and using $T = 10^4$ iterations. Using again the vector in `project4_vector.csv`, return, for $c \in \{5, 10, 20, 100\}$, the value of $|y - x_T|^2$ and of $|x_T - \text{proj}_{\Omega_c}(y)|^2$, where x_T is the output of the algorithm, and $\text{proj}_{\Omega_c}(y)$ was computed in Question (2).

```
def conditional_grad(y, c):
    def find_j(x, y):
        j, max_j = 0, 0
        for i in range(len(y)):
            tmp = np.abs(x[i] - y[i])
            if max_j < tmp:
                max_j = tmp
                j = i
        return j

    x = np.zeros(len(y))
    for t in range(10000):
        alpha = 2/(2+t)
        j = find_j(x, y)
        e = np.zeros(len(y))
        e[j] = 1
        x = (1-alpha)*x - c*alpha*np.sign(x[j] - y[j])*e
    return x

def problem2.5():
    y = pd.read_csv('homework4_data/project4_vector-1.csv')
    y = y.drop(['Unnamed: 0'], axis = 1).to_numpy()
    y = np.squeeze(y)
    list_c = [5, 10, 20, 100]
    for c in list_c:
        x_T = conditional_grad(y, c)
        print("The value of " + r'$|y-x_T|^2$' + " for c = " + str(c) + " is: "
              + str(np.linalg.norm(y - x_T)**2))
        print("The value of " + r'$|x_T-\text{proj}_{\{\Omega_c\}}(y)|^2$' + " for c = " + str(c) + " is: "
              + str(np.linalg.norm(x_T - project(y, c))**2))

if __name__ == "__main__":
    problem2.5()
```

Result:

The value of $|y - x_T|^2$ for $c = 5$ is: 9797.346377942433
 The value of $|x_T - \text{proj}_{\Omega_c}(y)|^2$ for $c = 5$ is: 1.8884932352928527e-06
 The value of $|y - x_T|^2$ for $c = 10$ is: 9768.217424264963
 The value of $|x_T - \text{proj}_{\Omega_c}(y)|^2$ for $c = 10$ is: 1.2463814546991488e-05
 The value of $|y - x_T|^2$ for $c = 20$ is: 9713.883736900732
 The value of $|x_T - \text{proj}_{\Omega_c}(y)|^2$ for $c = 20$ is: 0.00010089560776718761
 The value of $|y - x_T|^2$ for $c = 100$ is: 9340.645931864432
 The value of $|x_T - \text{proj}_{\Omega_c}(y)|^2$ for $c = 100$ is: 0.010320629431984777

(6) Let $m > 0$ be an integer and assume that a vector $b \in \mathbb{R}^m$ and an $m \times n$ matrix A is given. Let $F(x) = \frac{1}{2}|Ax - b|^2$, for $x \in \mathbb{R}^n$. Prove that the iterations

$$\begin{cases} x_{t+1} = (1 - \alpha_t)x_t - c\alpha_t \text{sign}(u_t^{(j_t)})e_{j_t} \\ j_t \in \arg \max\{|u_t^{(j)}|, j = 1, \dots, n\} \\ u_t = A^\top(Ax_t - b) \end{cases}$$

with $x_0 = 0$, $\alpha_t = 2/(2+t)$, converges to $x \in \arg \min_{\Omega_c} F$.

Proof. We have

$$F(x) = \frac{1}{2} (|Ax|^2 - 2(Ax)^\top b + |b|^2)$$

Calculate the gradient, and we get:

$$\nabla F(x) = A^\top(Ax - b) = u_t.$$

Using the result from Question (3), we have the iterations equivalent to

$$\begin{aligned} x_{t+1} &= (1 - \alpha_t)x_t + \alpha_t \arg \min_{\Omega_c} (x \mapsto u_t^\top x) \\ \Leftrightarrow x_{t+1} &= (1 - \alpha_t)x_t + \alpha_t \arg \min_{\Omega_c} (\nabla F(x_t)^\top (x - x_t)). \end{aligned}$$

So, let $x^* \in \arg \min_{\Omega_c} F$, with $\alpha_t = 2/(2+t)$, we always have

$$(1 - \rho\alpha_t)\alpha_t \leq \alpha_{t+1} \leq \alpha_t$$

for $\rho = 1/2$. Therefore, apply the Theorem of the conditional gradient algorithm, we get

$$F(x_t) - F(x^*) \leq \frac{2LD^2}{2+t}, t \geq 1$$

with $D = \max\{|x - y| : x, y \in \Omega_c\}$. As a consequence, with $x_0 = 0$, we obtain x_t converges to $\arg \min_{\Omega_c} F$. \square

(7) Program this algorithm, taking as input the matrix A and the vector b . Run this program with $T = 10^5$ iterations, using the data in project4_regression_A.csv and project4_regression_b.csv, for which $m = 500$ and $n = 1000$. For each value of $c \in \{5, 10, 15\}$ provide the values of the residual error, $|Ax_T - b|^2$, and the indexes, listed in increasing order, of the c largest numbers in the set $\{x_T^{(k)}, k = 1, \dots, n\}$.

```

def conditional_grad_2(A, b, c):
    def find_j(u):
        j, max_j = 0, 0
        for i in range(len(u)):
            tmp = np.abs(u[i])
            if max_j < tmp:
                max_j = tmp
                j = i
        return j

    x = np.zeros(A.shape[1])
    for t in range(100000):
        u = np.dot(A.T, np.dot(A, x) - b)
        alpha = 2/(2+t)
        j = find_j(u)
        e = np.zeros(A.shape[1])
        e[j] = 1
        x = (1-alpha)*x - c*alpha*np.sign(u[j])*e
    return x

def problem2.7():
    A = pd.read_csv('homework4_data/project4_regression_A-1.csv')
    A = A.drop(['Unnamed: 0'], axis = 1).to_numpy()
    A = np.squeeze(A)
    b = pd.read_csv('homework4_data/project4_regression_b.csv')
    b = b.drop(['Unnamed: 0'], axis = 1).to_numpy()
    b = np.squeeze(b)
    list_c = [5, 10, 15]
    for c in list_c:
        x_T = conditional_grad_2(A, b, c)
        print("The value of " + r'$|Ax_T - b|^2$' + " for c = " + str(c) + " is: "
              + str(np.linalg.norm(np.dot(A, x_T) - b)**2))
        list_max_x_T = np.sort(x_T)[-c:]
        print("The indexes of the c = " + str(c) + " largest number in the set "
              + r'$\{x_T^{(k)}, k=1, \dots, n\}$' + " is: ")
        for k in range(A.shape[1]):
            if x_T[k] in list_max_x_T:
                print(k, end = " ")
        print()

if __name__ == "__main__":
    problem2.7()

```

Result:

The value of $|Ax_T - b|^2$ for $c = 5$ is: 1361.2936729535782

The indexes of the $c = 5$ largest number in the set $\{x_T^{(k)}, k = 1, \dots, n\}$ is: 104 362 372 619 767

The value of $|Ax_T - b|^2$ for $c = 10$ is: 42.69954220008175

The indexes of the $c = 10$ largest number in the set $\{x_T^{(k)}, k = 1, \dots, n\}$ is: 104 362 372 451 477 538 619 629 767 789

The value of $|Ax_T - b|^2$ for $c = 15$ is: 0.002448971584146125

The indexes of the $c = 15$ largest number in the set $\{x_T^{(k)}, k = 1, \dots, n\}$ is: 104 217 362 372 451 477 520 523 538 619 629 767 789 863 990

(8) Program the projected gradient descent algorithm

$$x_{t+1} = \text{proj}_{\Omega_c}(x_t - \alpha \nabla F(x_t))$$

for the function F in Question (6), with $x_0 = 0$. This algorithm should take as input A , b and α and the vector b . Run this algorithm with the data in project4_regression_A.csv and project4_regression_b.csv and for $c \in \{5, 10, 15\}$, with $T = 10^5$ iterations. For each value of c provide the values of the residual error, $|Ax_T - b|^2$, and the indexes, listed in increasing order, of the c largest numbers in the set $\{x_T^{(k)}, k = 1, \dots, n\}$. Use $\alpha = 10^{-3}$ for $c = 5, 10$ and $\alpha = 10^{-4}$ for $c = 15$.

```

def projected_GD(A, b, c, alpha):
    x = np.zeros(A.shape[1])
    for t in range(100000):
        grad = np.dot(A.T, np.dot(A,x)-b)
        x = project(x-alpha*grad, c)
    return x

def problem2.8():
    A = pd.read_csv('homework4_data/project4_regression_A-1.csv')
    A = A.drop(['Unnamed: 0'], axis = 1).to_numpy()
    A = np.squeeze(A)
    b = pd.read_csv('homework4_data/project4_regression_b.csv')
    b = b.drop(['Unnamed: 0'], axis = 1).to_numpy()
    b = np.squeeze(b)
    list_c = [5, 10, 15]
    list_alpha = [1e-3, 1e-3, 1e-4]
    for i in range(len(list_c)):
        c = list_c[i]
        alpha = list_alpha[i]
        x_T = projected_GD(A, b, c, alpha)
        print("The value of " + r'$|Ax_T-b|^2$' + " for c = " + str(c) + " is: "
              + str(np.linalg.norm(np.dot(A,x_T)-b)**2))
        list_max_x_T = np.sort(x_T)[-c:]
        print("The indexes of the c = " + str(c) + " largest number in the set "
              + r'$\{x_T^{\{k\}}, k=1,\cdots,n\}$' + " is: ")
        for k in range(A.shape[1]):
            if x_T[k] in list_max_x_T:
                print(k, end = " ")
        print()

if __name__ == "__main__":
    problem2.8()

```

Result:

The value of $|Ax_T - b|^2$ for $c = 5$ is: 1361.2936662674565

The indexes of the $c = 5$ largest number in the set $\{x_T^{(k)}, k = 1, \dots, n\}$ is: 104 362 372 619 767

The value of $|Ax_T - b|^2$ for $c = 10$ is: 42.699437957778756

The indexes of the $c = 10$ largest number in the set $\{x_T^{(k)}, k = 1, \dots, n\}$ is: 104 362 372 451 477 538 619 629 767 789

The value of $|Ax_T - b|^2$ for $c = 15$ is:

ValueError: invalid value encountered in double_scalars

return $\text{lmda}_0 + ((c - \text{func.f}(y, \text{lmda}_0))/(\text{func.f}(y, \text{lmda}_0) - \text{func.f}(y, \text{lmda}_1))) * (\text{lmda}_0 - \text{lmda}_1)$