**Math Foundations of ML, Fall 2022**        **Instructor: Justin Romberg**

**Homework 3**

**Due: Friday, October 7, 2022 at 5:00 pm EST**

**As stated in the syllabus, unauthorized use of previous semester course materials is strictly prohibited in this course.**

1. The file `hw3p1_data.mat` contains two variables: `udata` and `ydata`. We will use this data to estimate a function $f : \mathbb{R}^2 \to \mathbb{R}$. The columns of `udata` contain sample locations, of which there are $M = 100$. The entries of $\boldsymbol{y}$ are the corresponding responses. We want to estimate $f$ such that

$$f(\boldsymbol{u}_m) \approx y_m, \quad m = 1, \ldots, M, \quad \text{where} \quad \boldsymbol{u}_m = \begin{bmatrix} s_m \\ t_m \end{bmatrix}.$$

   We will restrict $f$ to be a second-order polynomial on $[0,1] \times [0,1]$:

$$f(s,t) = \alpha_1 s^2 + \alpha_2 t^2 + \alpha_3 st + \alpha_4 s + \alpha_5 t + \alpha_6, \tag{1}$$

   which means that $f$ lies in a six dimensional subspace of $L_2([0,1]^2)$.

   (a) Explain how to compute the $100 \times 6$ matrix $\boldsymbol{A}$ so that $\boldsymbol{y} \approx \boldsymbol{A\alpha}$, where $\boldsymbol{y}$ contains the 100 response values in `ydata`. Write the code to compute $\boldsymbol{A}$ and turn it in.
   *Solution.*

$$\boldsymbol{A} = \begin{bmatrix} s_1^2 & t_1^2 & s_1 t_1 & s_1 & t_1 & 1 \\ s_2^2 & t_2^2 & s_2 t_2 & s_2 & t_2 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ s_{100}^2 & t_{100}^2 & s_{100} t_{100} & s_{100} & t_{100} & 1 \end{bmatrix}$$

   See script "Problem_1a.py" for the code.

   (b) Solve

$$\underset{\boldsymbol{\alpha} \in \mathbb{R}^6}{\text{minimize}} \ \|\boldsymbol{y} - \boldsymbol{A\alpha}\|_2^2.$$

   Turn in your code and the numerical value of your solution $\hat{\boldsymbol{\alpha}} \in \mathbb{R}^6$.
   *Solution.*

$$\boldsymbol{\alpha} = (\boldsymbol{A}^{\mathrm{T}} \boldsymbol{A})^{-1} \boldsymbol{A}^{\mathrm{T}} \boldsymbol{y}$$

   See script "Problem_1b.py" for the code and the console output below for the value of $\hat{\boldsymbol{\alpha}}$.

```
alpha:
[[-0.63387231]
 [ 1.22817315]
 [ 0.19346972]
 [ 0.97008806]
```

```
[ 0.22554988]
[ 1.23638021]]
```

(c) Make a contour plot of the corresponding

$$\hat{f}(s,t) = \hat{\alpha}_1 s^2 + \hat{\alpha}_2 t^2 + \hat{\alpha}_3 st + \hat{\alpha}_4 s + \hat{\alpha}_5 t + \hat{\alpha}_6.$$

Include 50 contour lines, just so we have a very clear picture of what this function looks like.

*Solution.*
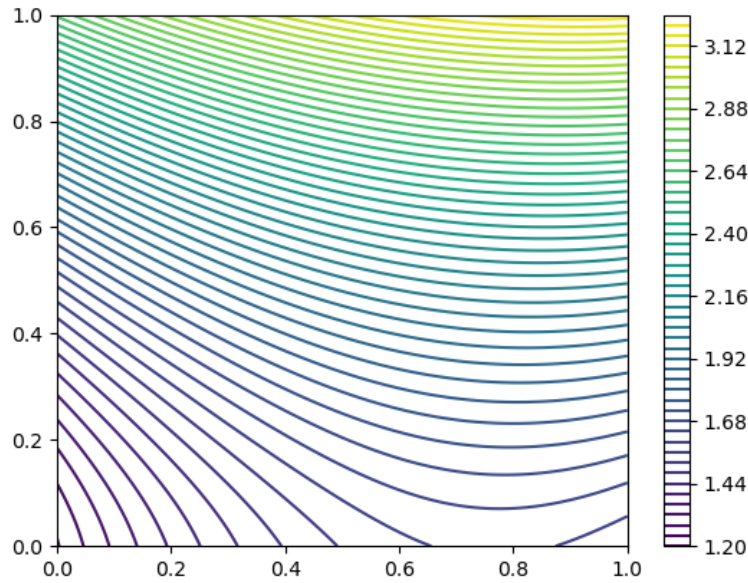See script "Problem_1c.py" for the code and Figure 1 for the contour plot.

Figure 1: Contour plot of $\hat{f}$.

2. Consider the space $\mathcal{P}_2$ of second-order polynomials on $[0,1]^2$ specified by $\boldsymbol{\alpha} \in \mathbb{R}^6$ as in (1) above.

   (a) At every point $(s,t)$, the gradient $\nabla f(s,t)$ of a function $f \in \mathcal{P}_2$ is a vector in $\mathbb{R}^2$. As every $f \in \mathcal{P}_2$ is specified by a vector $\boldsymbol{\alpha} \in \mathbb{R}^6$, we can think of the gradient at $(s,t)$ as a mapping from $\mathbb{R}^6$ to $\mathbb{R}^2$. Show that this mapping is linear, which means, for a specified $(s,t)$, there is a $2 \times 6$ matrix $\boldsymbol{G}_{s,t} \in \mathbb{R}^{2 \times 6}$ such that

   $$\nabla f(s,t) = \boldsymbol{G}_{s,t} \boldsymbol{\alpha}$$

   *Solution.*
   Since

   $$f(s,t) = \alpha_1 s^2 + \alpha_2 t^2 + \alpha_3 st + \alpha_4 s + \alpha_5 t + \alpha_6$$

2

we have

$$\nabla f(s,t) = \begin{bmatrix} 2\alpha_1 s + \alpha_3 t + \alpha_4 \\ 2\alpha_2 t + \alpha_3 s + \alpha_5 \end{bmatrix} = \begin{bmatrix} 2s & 0 & t & 1 & 0 & 0 \\ 0 & 2t & s & 0 & 1 & 0 \end{bmatrix} \boldsymbol{\alpha} = \boldsymbol{G}_{s,t}\boldsymbol{\alpha}$$

where $\boldsymbol{G}_{s,t} = \begin{bmatrix} 2s & 0 & t & 1 & 0 & 0 \\ 0 & 2t & s & 0 & 1 & 0 \end{bmatrix}$ and $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \alpha_3, \alpha_4, \alpha_5, \alpha_6]^\top$.

(b) Find the $6 \times 6$ matrix $\boldsymbol{H}_{s,t} \in \mathbb{R}^{6\times 6}$ such that[1]

$$\|\nabla f(s,t)\|_2^2 = \boldsymbol{\alpha}^{\mathrm{T}} \boldsymbol{H}_{s,t}\boldsymbol{\alpha}.$$

What kinds of functions $f$ are in the null space of $\boldsymbol{H}_{s,t}$ for all $s$ and $t$? Why?

*Solution.*

$$\begin{aligned}
\|\nabla f(s,t)\|_2^2 &= (\nabla f(s,t))^{\mathrm{T}}(\nabla f(s,t)) \\
&= (\boldsymbol{G}_{s,t}\boldsymbol{\alpha})^{\mathrm{T}}(\boldsymbol{G}_{s,t}\boldsymbol{\alpha}) \\
&= \boldsymbol{\alpha}^{\mathrm{T}}\boldsymbol{G}_{s,t}^{\mathrm{T}}\boldsymbol{G}_{s,t}\boldsymbol{\alpha} \\
&= \boldsymbol{\alpha}^{\mathrm{T}}\boldsymbol{H}_{s,t}\boldsymbol{\alpha}
\end{aligned}$$

where $\boldsymbol{H}_{s,t} = \boldsymbol{G}_{s,t}^{\mathrm{T}}\boldsymbol{G}_{s,t} = \begin{bmatrix} 4s^2 & 0 & 2st & 2s & 0 & 0 \\ 0 & 4t^2 & 2st & 0 & 2t & 0 \\ 2st & 2st & t^2+s^2 & t & s & 0 \\ 2s & 0 & t & 1 & 0 & 0 \\ 0 & 2t & s & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}.$

Since every $f \in \mathcal{P}_2$ is specified by a vector $\boldsymbol{\alpha} \in \mathbb{R}^6$, we need to find all possible $\boldsymbol{\alpha}$ such that

$$\boldsymbol{\alpha} \in \mathrm{Null}(\boldsymbol{H}_{s,t}) \quad \text{for all } (s,t) \in [0,1]^2.$$

Fix an arbitrary pair $(s,t) \in [0,1]^2$. If $\boldsymbol{\alpha} \in \mathrm{Null}(\boldsymbol{H}_{s,t})$, then we have $\boldsymbol{H}_{s,t}\boldsymbol{\alpha} = 0$, which implies that $\|\nabla f(s,t)\|_2^2 = \boldsymbol{\alpha}^{\mathrm{T}}\boldsymbol{H}_{s,t}\boldsymbol{\alpha} = 0$, i.e., $\nabla f(s,t) = 0$. Therefore, the $\boldsymbol{\alpha}$ we are seeking must satisfy:

$$\nabla f(s,t) = 0 \quad \text{for all } (s,t) \in [0,1]^2.$$

The $f \in \mathcal{P}_2$ that satisfies the above condition should be a constant function, i.e., $\boldsymbol{\alpha} = [0,0,0,0,0,\alpha_6]^\top$ for any $\alpha_6 \in \mathbb{R}$.

(c) Compute the matrix

$$\boldsymbol{Q} = \int_0^1 \int_0^1 \boldsymbol{H}_{s,t}\, ds\, dt.$$

(This is done simply by integrating each entry individually.)

*Solution.*

$$\boldsymbol{Q} = \begin{bmatrix} 4/3 & 0 & 1/2 & 1 & 0 & 0 \\ 0 & 4/3 & 1/2 & 0 & 1 & 0 \\ 1/2 & 1/2 & 2/3 & 1/2 & 1/2 & 0 \\ 1 & 0 & 1/2 & 1 & 0 & 0 \\ 0 & 1 & 1/2 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

---

[1]Hint: $\|\nabla f(s,t)\|_2^2 = \|\boldsymbol{G}_{s,t}\boldsymbol{\alpha}\|_2^2 = \langle \boldsymbol{G}_{s,t}\boldsymbol{\alpha}, \boldsymbol{G}_{s,t}\boldsymbol{\alpha}\rangle = \cdots$.

(d) Describe how to set up and solve the optimization program

$$\underset{\boldsymbol{f}\in\mathcal{P}_2}{\text{minimize}} \sum_{m=1}^{M}(y_m - f(s_m, t_m))^2 + \delta \int_0^1 \int_0^1 \|\nabla f(s,t)\|_2^2\, ds dt.$$

What is the regularizer above penalizing? What kinds of solutions do we expect for large $\delta$?

*Solution.*

$$\underset{\boldsymbol{f}\in\mathcal{P}_2}{\text{minimize}} \sum_{m=1}^{M}(y_m - f(s_m, t_m))^2 + \delta \int_0^1 \int_0^1 \|\nabla f(s,t)\|_2^2\, ds dt$$

$$= \underset{\boldsymbol{\alpha}\in\mathbb{R}^6}{\text{minimize}} \sum_{m=1}^{M}(y_m - [s_m^2, t_m^2, s_m t_m, s_m, t_m, 1]^\top\boldsymbol{\alpha})^2 + \delta \int_0^1 \int_0^1 \boldsymbol{\alpha}^\mathrm{T}\boldsymbol{H}_{s,t}\boldsymbol{\alpha}\, ds dt$$

$$= \underset{\boldsymbol{\alpha}\in\mathbb{R}^6}{\text{minimize}} \ (\boldsymbol{y} - \boldsymbol{A}\boldsymbol{\alpha})^\mathrm{T}(\boldsymbol{y} - \boldsymbol{A}\boldsymbol{\alpha}) + \delta\boldsymbol{\alpha}^\mathrm{T}\boldsymbol{Q}\boldsymbol{\alpha} \quad (\boldsymbol{A} \text{ is defined in Problem 1(a)})$$

Taking the gradient with respect to $\boldsymbol{\alpha}$ and setting it equal to zero (note: $\boldsymbol{A}^\mathrm{T}\boldsymbol{A} + \delta\boldsymbol{Q}$ is symmetric) leads to

$$-2\boldsymbol{A}^\mathrm{T}\boldsymbol{y} + 2(\boldsymbol{A}^\mathrm{T}\boldsymbol{A} + \delta\boldsymbol{Q})\boldsymbol{\alpha} = \boldsymbol{0}$$

Thus, we have $\widehat{\boldsymbol{\alpha}} = (\boldsymbol{A}^\mathrm{T}\boldsymbol{A} + \delta\boldsymbol{Q})^{-1}\boldsymbol{A}^\mathrm{T}\boldsymbol{y}$. The regularizer is penalizing the integrated squared magnitude of the function's gradient over the region of interest. In other words, functions with more oscillation are penalized. For large $\delta$, we should expect the resultant function to be a nearly constant function.

(e) Apply your answer to part (d) to the data set from Problem 1. Play around with the value of $\delta$, and produce estimates for three different $\delta$ that are interesting. Discuss why you think those values are indeed "interesting".

*Solution.*

See script "Problem_2e.py" for the code and Figure 2 for the plots, and console output below for different values of $\delta$. I chose a small $\delta$, a medium $\delta$, and a large $\delta$ that put $\boldsymbol{A}^\mathrm{T}\boldsymbol{A}$ and $\delta\boldsymbol{Q}$ at roughly the same magnitude. Observe that for small $\delta$, we get roughly the pure least squares solution. For large $\delta$, we get a nearly constant function. For the medium $\delta$, we get a smoother, more linear function. It seems as though regularization affects the returned function in a fairly smooth manner.

```
Results for delta=0.0001
alpha:
[[-0.63385213]
 [ 1.22810731]
 [ 0.19346862]
 [ 0.9700667 ]
 [ 0.22560036]
 [ 1.23638087]]
Results for delta=53.39007960128817
```

4

```
alpha:
[[-0.04726227]
 [ 0.06574148]
 [-0.01368813]
 [ 0.13965968]
 [ 0.15446664]
 [ 1.92434391]]
Results for delta=10000.0
alpha:
[[-2.87066019e-04]
 [ 3.86242113e-04]
 [-1.09711413e-04]
 [ 8.94402034e-04]
 [ 9.87040696e-04]
 [ 2.07142530e+00]]
```
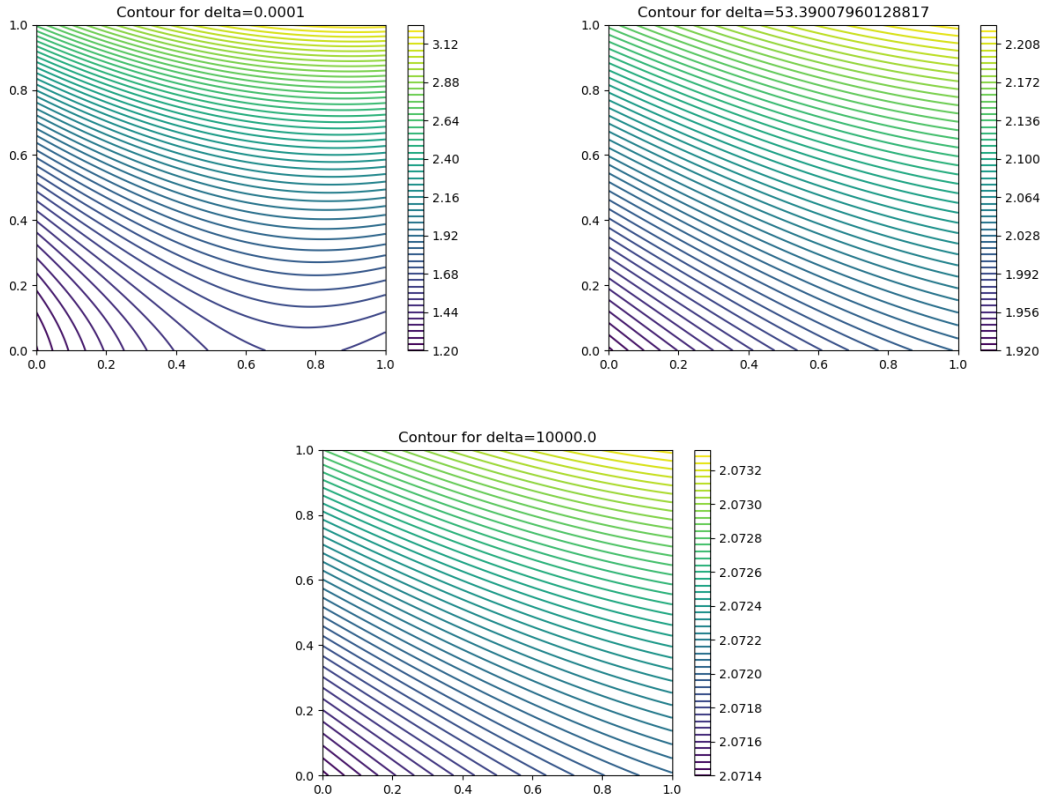


Figure 2: Plots for Problem 2(e)

3. Let $\boldsymbol{A}$ be an $M \times N$ matrix with $\operatorname{rank}(\boldsymbol{A}) < N$. We have seen in this case that the least-squares problem

$$\underset{\boldsymbol{x} \in \mathbb{R}^N}{\operatorname{minimize}} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_2^2 \tag{2}$$

5

has an infinite number of solutions. We have also seen, however, that the regularized least squares problem

$$\underset{\boldsymbol{x}\in\mathbb{R}^N}{\operatorname{minimize}} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_2^2 + \delta\|\boldsymbol{x}\|_2^2 \tag{3}$$

has a unique solution for every $\delta > 0$. In this problem, we will show that as $\delta \to 0$, the regularized solution goes to the minimum norm solution of

$$\underset{\boldsymbol{x}\in\mathbb{R}^N}{\operatorname{minimize}} \|\boldsymbol{x}\|_2^2 \quad \text{subject to} \quad \boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}\boldsymbol{x} = \boldsymbol{A}^{\mathrm{T}}\boldsymbol{y}. \tag{4}$$

(a) Start by showing that if $\boldsymbol{x}_1 \in \operatorname{Row}(\boldsymbol{A})$ and $\boldsymbol{x}_2 \in \operatorname{Row}(\boldsymbol{A})$ then $\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}\boldsymbol{x}_1 \neq \boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}\boldsymbol{x}_2$ unless $\boldsymbol{x}_1 = \boldsymbol{x}_2$.

*Solution.*

We know from thethe Technical Details section in "Notes-mfml-regls.pdf" that $\operatorname{Col}(\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}) = \operatorname{Row}(\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}) = \operatorname{Row}(\boldsymbol{A})$, and so $\operatorname{Null}(\boldsymbol{A})$ is orthogonal to $\operatorname{Row}(\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A})$. This means that for all $\boldsymbol{v} \in \operatorname{Row}(\boldsymbol{A})$, $\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}\boldsymbol{v} = \boldsymbol{0}$ if and only if $\boldsymbol{v} = \boldsymbol{0}$. If $\boldsymbol{x}_1, \boldsymbol{x}_2 \in \operatorname{Row}(\boldsymbol{A})$ then $\boldsymbol{x}_1 - \boldsymbol{x}_2$ is also in $\operatorname{Row}(\boldsymbol{A})$, so $\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}(\boldsymbol{x}_1 - \boldsymbol{x}_2) = 0$ (meaning $\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}\boldsymbol{x}_1 = \boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}\boldsymbol{x}_2$) if and only if $\boldsymbol{x}_1 = \boldsymbol{x}_2$.

(b) Use part (a) to argue that the solution to (4) is always unique.

*Solution.*

We know from "Notes-mfml-regls.pdf" that any solution to (4) must be in $\operatorname{Row}(\boldsymbol{A})$. And we also know from part (a) that there is exactly one point $\boldsymbol{x}^\star \in \operatorname{Row}(\boldsymbol{A})$ that satisfies $\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}\boldsymbol{x}^\star = \boldsymbol{A}^{\mathrm{T}}\boldsymbol{y}$. Thus this point $\boldsymbol{x}^\star$ must be the unique solution to (4).

(c) In fact, something stronger than what we showed in part (a) is true.

There exists a constant $C > 0$ such that

$$\|\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}(\boldsymbol{x}_1 - \boldsymbol{x}_2)\|_2 \geq C\|\boldsymbol{x}_1 - \boldsymbol{x}_2\|_2 \quad \text{for all} \quad \boldsymbol{x}_1, \boldsymbol{x}_2 \in \operatorname{Row}(\boldsymbol{A}).$$

(This follows very easily from work we do later in the course, so we will defer its proof for now.) Use this fact to show that the solution of (3) goes to the solution of (4) as $\delta \to 0$. In particular, if $\boldsymbol{x}^\star$ is the (always unique) minimizer of (4), and $\hat{\boldsymbol{x}}_n$ is the (always unique) minimizer of (3) with[2] $\delta = 1/n$, show that

$$\lim_{n \to \infty} \hat{\boldsymbol{x}}_n = \boldsymbol{x}^\star,$$

i.e. $\lim_{n \to \infty} \|\boldsymbol{x}^\star - \hat{\boldsymbol{x}}_n\|_2 = 0$.

*Solution.*

Let $\delta_n := \frac{1}{n}$. Note that $\boldsymbol{x}^\star \in \operatorname{Row}(\boldsymbol{A})$ and $\hat{\boldsymbol{x}}_n \in \operatorname{Row}(\boldsymbol{A})$ for all $n \geq 1$. We also know that

$$\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}\boldsymbol{x}^\star = \boldsymbol{A}^{\mathrm{T}}\boldsymbol{y},$$
$$(\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A} + \delta_n\mathbf{I})\hat{\boldsymbol{x}}_n = \boldsymbol{A}^{\mathrm{T}}\boldsymbol{y}.$$

Hence, $\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}(\boldsymbol{x}^\star - \hat{\boldsymbol{x}}_n) = \delta_n\hat{\boldsymbol{x}}_n$ and

$$\|\boldsymbol{x}^\star - \hat{\boldsymbol{x}}_n\|_2 \leq \frac{1}{C}\|\boldsymbol{A}^{\mathrm{T}}\boldsymbol{A}(\boldsymbol{x}^\star - \hat{\boldsymbol{x}}_n)\|_2 = \frac{\delta_n}{C}\|\hat{\boldsymbol{x}}_n\|_2. \tag{5}$$

---

[2]There is nothing special about taking $\delta = 1/n$ ... your argument should work for any sequence of $\delta$s that goes to zero.

Now we show that $\delta_n \|\widehat{\boldsymbol{x}}_n\|_2$ can be upper bounded by $\delta_n \|\boldsymbol{x}^\star\|_2$ for all $n \geq 1$. Note that since $\boldsymbol{x}^\star$ is a minimizer of (2), we have

$$\|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}^\star\|_2^2 \;\leq\; \|\boldsymbol{y} - \boldsymbol{A}\widehat{\boldsymbol{x}}_n\|_2^2$$

for all $n \geq 1$. Moreover, we know that

$$\|\boldsymbol{y} - \boldsymbol{A}\widehat{\boldsymbol{x}}_n\|_2^2 + \delta_n\|\widehat{\boldsymbol{x}}_n\|_2^2 \;\leq\; \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}^\star\|_2^2 + \delta_n\|\boldsymbol{x}^\star\|_2^2.$$

Thus, we have

$$\delta_n\|\widehat{\boldsymbol{x}}_n\|_2^2 \leq \delta_n\|\boldsymbol{x}^\star\|_2^2 - \left(\|\boldsymbol{y} - \boldsymbol{A}\widehat{\boldsymbol{x}}_n\|_2^2 - \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}^\star\|_2^2\right) \leq \delta_n\|\boldsymbol{x}^\star\|_2^2. \qquad (6)$$

Combining (5) and (6) yields

$$\|\boldsymbol{x}^\star - \widehat{\boldsymbol{x}}_n\|_2 \;\leq\; \frac{\delta_n}{C}\|\boldsymbol{x}^\star\|_2,$$

and thus $\lim_{n\to\infty} \|\boldsymbol{x}^\star - \widehat{\boldsymbol{x}}_n\|_2 = 0$.