

### Problem 1 - Hard Thresholding

Given the following,

$$\text{prox}_h(x) := \arg \min_{z_i \in \mathbb{R}} \{ \lambda 1\{z_i \neq 0\} + \frac{1}{2} (z_i - x_i)^2 \}, \forall i = 1, \dots, n$$

derive the solution formula for the 1-dimensional optimization problem, and then obtain the formula for the proximal mapping,  $\text{prox}_{\lambda \|\cdot\|_0}(x)$ .

Solution:

$\text{prox}_{\lambda \|\cdot\|_1}$  - soft thresholding

$$= (|x| - \lambda)_+ \cdot \text{sgn}(x)$$

$\text{prox}_{\lambda \|\cdot\|_0}$  - look at  $\arg \min_{z \in \mathbb{R}} \{ \lambda 1\{z \neq 0\} + \frac{1}{2} (z - x)^2 \}$

- derive analytical form using cases  $F_x(z)$

$$F_x(z) = \begin{cases} \lambda(1) + \frac{1}{2} (z - x)^2, & z \neq 0 \\ \lambda(0) + \frac{1}{2} (z - x)^2, & z = 0 \end{cases} = \begin{cases} \lambda + \frac{1}{2} (z - x)^2, & z \neq 0 \\ \frac{1}{2} x^2, & z = 0 \end{cases}$$

↳ constant

- find minimizer for each part

$$F_x(z) = \begin{cases} \text{min val is } \lambda, & z = x \\ \text{min val is } \frac{x^2}{2}, & z = 0 \end{cases}$$

} compare 2 min values  
↳ go with what is lower

⇒ If  $\frac{x^2}{2} \leq \lambda \Leftrightarrow |x| \leq \sqrt{2\lambda}$   
 then  $z^* = 0$

If  $\lambda < \frac{x^2}{2} \Leftrightarrow |x| > \sqrt{2\lambda}$   
 then  $z^* = x$

soft thresholding

→ this is hard thresholding  
 ↳ discontinuous  
 ↳ may lead to local min

→ could pick any point. This is sub-differential

$$\text{prox}_{\lambda \|\cdot\|_0}(x) = (1\{|x| > \sqrt{2\lambda}\}) \cdot x$$

## Simulation I - Proximal Gradient Descent for Compressed Sensing

### Gradient Descent Implementation

The first objective of this assignment was to apply gradient descent to solve the least-mean-square problem defined as:

$$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{2} \|Ax - b\|_2^2, \text{ where } x \in \mathbb{R}^n, A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m.$$

The gradient for the definition above is given as  $\nabla f(x) = A^T(Ax - b)$ .

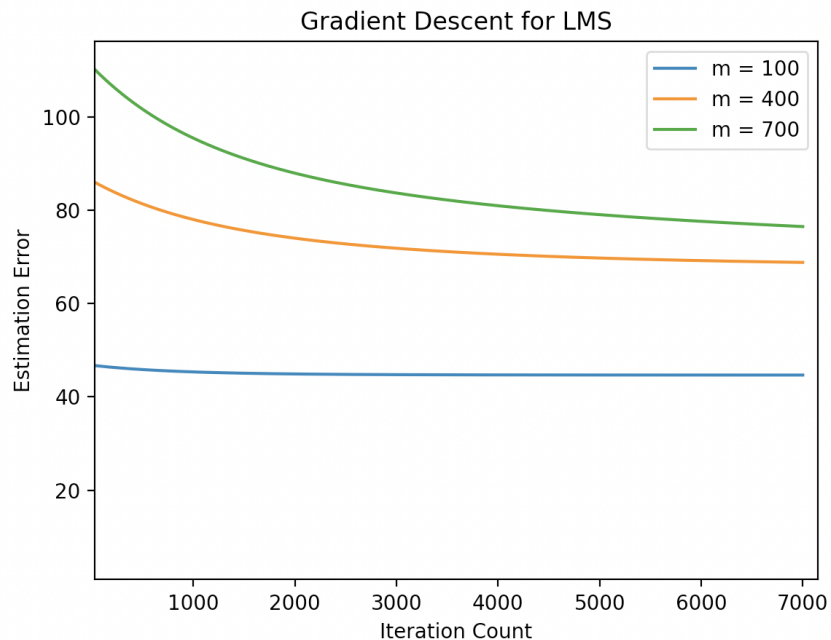
The underlying signal ( $z$ ), and the initialization of  $x_0$  were generated randomly using the `numpy.random.randn` function. The underlying signal was made sparse by randomly replacing 95% of its values with zeros. The initial noise variance was chosen to be 0.1. The variables  $A$  and  $b$  were generated in the same way as in HW2.

The gradient algorithm implementation is iterative and takes the following form:

$$x_{k+1} = x_k - \eta \nabla f(x_k), \text{ where } \eta \text{ is the learning rate.}$$

The results were obtained using the following parameters:  $n = 1000$ , learning rate =  $1e-3$ , variance = 0.1, iteration number  $k = 7000$ . The algorithm was run three times setting the value of  $m$  to 100, 400, 700.

The results of the gradient descent implementation estimation error are shown below.



We can see that as  $m$  increases, the longer it takes for the error to converge. This is logical because a larger  $m$  means that more samples are provided and thus, more gradients need to be updated and obtained.

### Proximal Gradient Descent Implementation

The second objective of this assignment was to apply proximal gradient descent to solve the compressed sensing problem defined as:

$$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1, \text{ where } x \in \mathbb{R}^n, A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m.$$

The proximal operator of  $\lambda \|x\|_1$  is given as:

$$\text{prox}_{\lambda \|\cdot\|_1}(x) = (|x| - \lambda)_+ * \text{sign}(x), \text{ where } (u)_+ := \max\{u, 0\}$$

The underlying signal ( $z$ ), and the initialization of  $x_0$  were generated randomly using the `numpy.random.randn` function. The underlying signal was made sparse by randomly replacing 95% of its values with zeros. The initial noise variance was chosen to be 0.1. The variables  $A$  and  $b$  were generated in the same way as in HW2.

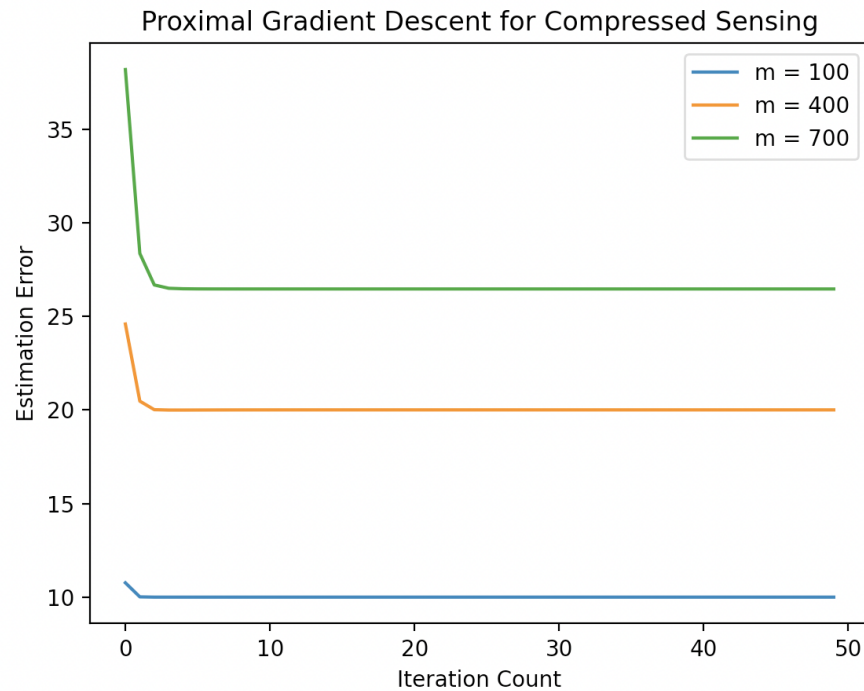
The proximal gradient algorithm implementation is iterative and takes the following update rule:

$$x_{k+1} = \text{prox}_{\lambda \|\cdot\|_1}(x_k - \eta \nabla f(x_k)), \text{ where } \eta \text{ is the learning rate.}$$

The  $\lambda$  was finely tuned by experimentally running the model with the values of  $\lambda = [10, 1, 0.1, 0.01, 0.001, 0.0001]$  and seeing what the lowest error was. The final  $\lambda$  was determined to be 0.1.

The results were obtained using the following parameters:  $n = 1000$ , learning rate =  $1e-3$ ,  $\lambda = 0.1$ , variance = 0.1, iteration number  $k = 50$ . The algorithm was run three times setting the value of  $m$  to 100, 400, 700.

The results of the proximal gradient descent implementation estimation error are shown below.



We see that as  $m$  increases, the slope of the estimation error decreases at a steeper rate and although it is a little tough to see, it seems that it might take the larger  $m$  values a few more iterations to converge.

When comparing LMS and CS, CS is clearly a better choice. It takes far fewer iterations to converge and the convergence happens a lot faster, for all  $m$  values, than in the LMS model.