

## Lecture 4:

Lecturer: Prof. Yue Qiu &amp; Prof. Ziping Zhao

Scribe: Bing Jiang

**Disclaimer:** These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.

## 1 Least Squares

**problem:** given  $\mathbf{y} \in \mathbb{R}^m, \mathbf{A} \in \mathbb{R}^{m \times n}$ , solve

$$\mathbf{x}_{\text{LS}} = \min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 \quad (1)$$

- called (linear) least squares (LS)
- find an  $\mathbf{x}$  whose residual  $\mathbf{r} = \mathbf{y} - \mathbf{A}\mathbf{x}$  is the smallest in the Euclidean sense

**solution:** suppose that  $\mathbf{A}$  has full-column rank ( $m \geq n$ ). The solution to (LS) is unique and is given by

$$\mathbf{X}_{\text{LS}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} \quad (2)$$

- complexity  $O(mn^2 + n^3)$
- LS solutions to an **overdetermined** system of equations  $\mathbf{y} = \mathbf{A}\mathbf{x}$  ( $m > n$ )
- if  $\mathbf{A}$  is semi-orthogonal, the solution is simplified to  $\mathbf{X}_{\text{LS}} = \mathbf{A}^T \mathbf{y}$
- if  $\mathbf{A}$  is square, the solution is simplified to  $\mathbf{X}_{\text{LS}} = \mathbf{A}^{-1} \mathbf{y}$
- unless specified, in this lecture we will assume  $\mathbf{A}$  to have full column rank without further mentioning

## 2 LS Solution

**Theorem 1** A vector  $\mathbf{X}_{\text{LS}}$  is an optimal solution to the LS problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 \quad (3)$$

if and only if it satisfies

$$\mathbf{A}^T \mathbf{A} \mathbf{X}_{\text{LS}} = \mathbf{A}^T \mathbf{y} \quad (4)$$

- the optimality condition in (4) is true for any  $\mathbf{A}$ , not just full-column rank  $\mathbf{A}$

- suppose that  $\mathbf{A}$  has full column rank
  - (4) is a symmetric PD linear system
  - the Gram matrix  $\mathbf{A}^T \mathbf{A}$  is non-singular
  - the solution to (4) is uniquely given by  $\mathbf{X}_{LS} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$
- (4) is called the **normal equations**
- the same result holds for the complex case, viz.,  $\mathbf{A}^H \mathbf{A} \mathbf{X}_{LS} = \mathbf{A}^H \mathbf{y}$

### 3 Gradient Descent For LS

- consider a general unconstrained optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad (5)$$

where  $f$  is continuously differentiable

- **Gradient Descent:** given a starting point  $x^{(0)}$ , do

$$x^{(k)} = x^{(k-1)} - \mu \nabla f(x^{(k-1)}), k = 1, 2, \dots \quad (6)$$

where  $\mu > 0$  is a step size

- take an optimization course to get more details! It is known that
  - for convex  $f$  and under some appropriate choice of  $\mu$ , gradient descent converges to an optimal solution
  - for non-convex  $f$  and under some appropriate choice of  $\mu$ , gradient descent converges to a stationary point
- gradient descent for LS:

$$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} - 2\mu(\mathbf{A}^T \mathbf{A} \mathbf{x}^{(k-1)} - \mathbf{A}^T \mathbf{y}), k = 0, 1, \dots \quad (7)$$

- complexity for dense  $\mathbf{A}$ 
  - computing  $\mathbf{A}^T \mathbf{A}$  and  $\mathbf{A}^T \mathbf{y}$ :  $O(mn^2)$  and  $O(mn)$ , resp. (same as before)
    - \*  $\mathbf{A}^T \mathbf{A}$  and  $\mathbf{A}^T \mathbf{y}$  are cached for subsequent use in gradient descent
  - complexity of each iteration:  $O(n^2)$
- complexity for sparse  $\mathbf{A}$ 
  - computing  $\mathbf{A}^T \mathbf{y}$ :  $\mathcal{O}(nnz(\mathbf{A}))$
  - complexity of each iteration:  $O(n + nnz(\mathbf{A}))$ 
    - \*  $\mathbf{A}^T \mathbf{A}$  is not necessarily sparse, so we do  $\mathbf{A} \mathbf{x}^{(k-1)}$  and then  $\mathbf{A}^T (\mathbf{A} \mathbf{x}^{(k-1)})$

## 4 Regularized LS

- **Intuition:** replace  $\mathbf{X}_{LS} = \mathbf{A}^\dagger \mathbf{y} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$  by

$$\mathbf{X}_{RLS} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}, \quad (8)$$

for some  $\lambda > 0$ , where the term  $\lambda \mathbf{I}$  is added to improve the system conditioning, thereby attempting to reduce noise sensitivity

- how may we make sense out of such a modification?
- **L2 regularized LS:** find an  $\mathbf{x}$  that solves

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_2^2 \quad (9)$$

for some pre-determined  $\lambda > 0$ .

- the solution is uniquely given by  $\mathbf{X}_{RLS} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}$
- the formulation says that we try to minimize both  $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$  and  $\|\mathbf{x}\|_2^2$ , and  $\lambda$  controls which one should be more emphasized in the minimization.
- **L1 regularized LS:** given  $\lambda > 0$ , solve

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1 \quad (10)$$

- now consider applying Majorization-Minimization to the L2-L1 minimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1 \quad (11)$$

- The Majorization-Minimization method for solving problem updates  $\mathbf{x}$  as

$$\mathbf{x}^{(k+1)} = \mathbf{soft}\left(\frac{1}{c} \mathbf{A}^T (\mathbf{y} - \mathbf{A}\mathbf{x}^{(k)}) + \mathbf{x}^{(k)}, \lambda/c\right),$$

where **soft** is called the soft-thresholding operator and is defined as follows: if  $\mathbf{z} = \mathbf{soft}(\mathbf{x}, \sigma)$ , then  $z_i = \text{sign}(x_i) \max\{|x_i| - \sigma, 0\}$ .