

MF DNN

Homework 1

24/03/08

Using underline for column vectors  
 using CAPITAL with no subscripts for matrix  
 All others are scalar.

$$\underline{\theta} \in \mathbb{R}^p$$

$$1. (a) \quad \frac{\partial}{\partial \theta_j} \ell_i(\underline{\theta}) = \frac{\partial}{\partial \theta_j} \left( \frac{1}{2} (\underline{x}_i^T \underline{\theta} - y_i)^2 \right) = \frac{\partial}{\partial \theta_j} \left( \frac{1}{2} \left( \sum_{n=1}^p x_{in} \theta_n - y_i \right)^2 \right)$$

$$= x_{ij} (\underline{x}_i^T \underline{\theta} - y_i)$$

$$\underline{\nabla}_{\underline{\theta}} \ell_i(\underline{\theta}) = \sum_{j=1}^p \left( \frac{\partial}{\partial \theta_j} \ell_i(\underline{\theta}) \right) \underline{e}_j$$

$$= \sum_{j=1}^p (x_{ij} (\underline{x}_i^T \underline{\theta} - y_i)) \underline{e}_j$$

$$= (\underline{x}_i^T \underline{\theta} - y_i) \underline{x}_i \quad \blacksquare \quad \text{since } \sum_{j=1}^p x_{ij} \underline{e}_j = \begin{pmatrix} x_{i1} \\ 0 \\ \vdots \end{pmatrix} + \begin{pmatrix} 0 \\ x_{i2} \\ \vdots \end{pmatrix} + \dots + \begin{pmatrix} 0 \\ 0 \\ x_{ip} \end{pmatrix} = \underline{x}_i$$

$$(b) \quad \frac{\partial}{\partial \theta_j} \mathcal{L}(\underline{\theta}) = \frac{\partial}{\partial \theta_j} \left( \frac{1}{2} (\underline{x} \underline{\theta} - \underline{y})^T (\underline{x} \underline{\theta} - \underline{y}) \right) = \frac{\partial}{\partial \theta_j} \left( \frac{1}{2} \left( \sum_{i=1}^N (\underline{x}_{:,i} \theta_i) - \underline{y} \right)^2 \right) \quad \text{using given}$$

$$= \underline{x}_{:,j}^T (\sum_{i=1}^N (\underline{x}_{:,i} \theta_i) - \underline{y})$$

$$= \underline{x}_{:,j}^T (\underline{x} \underline{\theta} - \underline{y}) \quad \text{using given again}$$

$$\underline{y}, \underline{x}_{:,i} \in \mathbb{R}^{N \times 1}$$

$$\underline{x}_{:,i}^T \in \mathbb{R}^{1 \times N}$$

$$\underline{\nabla}_{\underline{\theta}} \mathcal{L}(\underline{\theta}) = \sum_{j=1}^p \underline{e}_j (\underbrace{\underline{x}_{:,j}^T}_{\text{scalar}} (\underline{x} \underline{\theta} - \underline{y}))$$

$$\underline{e}_j \in \mathbb{R}^{p \times 1}$$

$$= \sum_{j=1}^p \left[ (\underline{e}_j \underline{x}_{:,j}^T) (\underline{x} \underline{\theta} - \underline{y}) \right] \quad \text{since matrix multiplication is associative}$$

$$= \left( \sum_{j=1}^p \underline{e}_j \underline{x}_{:,j}^T \right) (\underline{x} \underline{\theta} - \underline{y})$$

$$= \underline{X}^T (\underline{x} \underline{\theta} - \underline{y}) \quad \blacksquare$$

$$\in \mathbb{R}^p$$

since  $\underline{x}_{:,j}^T$  is  $j$ th column of  $X$  as a row vector  
 i.e.  $(x_{1j}, x_{2j}, \dots, x_{Nj})$

$$\underline{e}_j \underline{x}_{:,j}^T = \begin{pmatrix} 1 \\ 0 \\ \vdots \end{pmatrix} (x_{1j}, x_{2j}, \dots, x_{Nj}) = \begin{pmatrix} x_{1j} & x_{2j} & \dots & x_{Nj} \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix} \in \mathbb{R}^{p \times N}$$

$$\Rightarrow \sum_{j=1}^p \underline{e}_j \underline{x}_{:,j}^T = \underline{X}^T$$

$$\underline{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & \dots & \dots & \vdots \\ \vdots & \dots & \dots & \vdots \\ x_{N1} & \dots & \dots & x_{Np} \end{pmatrix} \in \mathbb{R}^{N \times p}$$

$$\begin{aligned}
 2. \quad \theta^{k+1} &= \theta^k - \alpha f'(\theta^k) & f(x) &= x^2/2 \\
 \theta^{k+1} &= \theta^k - \alpha \theta^k & f'(x) &= x \\
 \theta^{k+1} &= \theta^k (1 - \alpha) \\
 \frac{\theta^{k+1}}{\theta^k} &= 1 - \alpha & \text{valid for all } k \text{ since } \theta^0 &\neq 0
 \end{aligned}$$

$$\lim_{k \rightarrow \infty} \left| \frac{\theta^{k+1}}{\theta^k} \right| = \lim_{k \rightarrow \infty} |1 - \alpha| \quad \text{By the ratio test, the sequence } \theta^k \text{ diverges (i.e. } |\theta^k| \rightarrow \infty \text{) if } \lim_{k \rightarrow \infty} |1 - \alpha| > 1.$$

$$\Rightarrow |1 - \alpha| > 1 \Leftrightarrow 1 - \alpha > 1 \text{ or } 1 - \alpha < -1$$

$$\Leftrightarrow 0 > \alpha \text{ or } 2 < \alpha$$

Since  $\alpha$  is positive,  $\alpha > 2$  only relevant condition

i.e.  $\theta$  sequence diverges if  $\alpha > 2$   $\square$

3.  $f(\theta) = \frac{1}{2} \|X\theta - y\|^2$

$\nabla_{\theta} f(\theta) = X^T(X\theta - y)$  from problem 1.

$0 = X^T(X\theta - y)$

$0 = X^T X \theta - X^T y \Rightarrow \theta^* = (X^T X)^{-1} X^T y$  is optimal solution (assuming  $X^T X$  invertible)

NB: Also implies  $X^T X \theta^* = X^T y$

Consider,  $\theta^{k+1} - \theta^* = \theta^k - \alpha \nabla f(\theta^k) - \theta^*$   
 "error vector"  $\nearrow$   
 $= \theta^k - \alpha (X^T(X\theta^k - y)) - \theta^*$   
 $= \theta^k - \alpha X^T X \theta^k + \alpha X^T y - \theta^*$   
 $= \theta^k - \alpha X^T X \theta^k + \alpha X^T X \theta^* - \theta^*$   
 $= \theta^k - \theta^* - \alpha X^T X (\theta^k - \theta^*)$   
 $= (I - \alpha X^T X) (\theta^k - \theta^*)$

Set  $A := I - \alpha X^T X$

$\Rightarrow$  The error vector will grow/diverge if  $\rho(A) > 1$

(aside: to use this, consider  $\underline{x} = \theta^k - \theta^*$  and  $A = I - \alpha X^T X$ .)

Let  $A$  have unique eigenvectors  $\underline{v}_1, \dots, \underline{v}_n$  along with eigenvalues  $\lambda_1, \dots, \lambda_n$  ( $n \leq p$ )

$\underline{x}$  can be expressed in basis of eigenvectors:  $\underline{x} = c_1 \underline{v}_1 + \dots + c_n \underline{v}_n$

In this case,  $A\underline{x} = A(c_1 \underline{v}_1 + \dots + c_n \underline{v}_n)$   
 $= c_1 A \underline{v}_1 + \dots + c_n A \underline{v}_n$   
 $= c_1 \lambda_1 \underline{v}_1 + \dots + c_n \lambda_n \underline{v}_n$

If we apply  $A$  again, i.e.  $AA\underline{x} = \dots = c_1 \lambda_1^2 \underline{v}_1 + \dots + c_n \lambda_n^2 \underline{v}_n$ , it's clear that if  $|\lambda_i| < 1$  then the components of  $\underline{x}$ , and hence  $\|\underline{x}\|$  will shrink and, conversely, if  $|\lambda_i| > 1$ ,  $\|\underline{x}\|$  will grow/diverge. Hence we want to bound the largest eigenvalue (spectral radius) of  $A$ .

Eigenvalues of  $A$  are  $1 - \alpha \lambda_i$  where  $\lambda_i$  are eigenvalues of  $X^T X$ .

$\Rightarrow$  For divergence,  $\exists i$  s.t.  $|1 - \alpha \lambda_i| > 1$

In worst case, only one such  $i$ :  $|1 - \alpha \rho(X^T X)| > 1$

$\nearrow$   
 $\rho(X^T X) := \max_i \{|\lambda_1|, \dots, |\lambda_n|\}$   
 $\Leftrightarrow |1 - \alpha \rho(X^T X)| > 1$  or  $|1 - \alpha \rho(X^T X)| < -1$   
 $\Leftrightarrow 0 > \alpha \rho(X^T X)$  or  $2 < \alpha \rho(X^T X)$   
 N/A  $\quad \quad \quad 2/\rho(X^T X) < \alpha$

$\Rightarrow$  Error vector, by extension GD, diverges if  $\alpha > 2/\rho(X^T X)$   
 (for most  $\theta^0$ . There may be some  $\theta^0$  for instance with no component in the direction of a diverging eigenvector/value but then there are very specific in the context of  $\mathbb{R}^n$  space)

5. (Workings)

$$A \in \mathbb{R}^{(n-r+1) \times n} \quad \begin{matrix} n-r+1 > 1 \\ n > r \end{matrix} \quad \underline{x} \in \mathbb{R}^{n \times 1}$$

$$\underline{Ax} = \underbrace{\begin{pmatrix} k_1 & \dots & k_r & 0 & \dots & 0 \\ 0 & k_1 & \dots & k_r & 0 & \dots & 0 \\ 0 & 0 & k_1 & \dots & k_r & 0 & \dots & 0 \\ \vdots & & & & & & & \\ 0 & \dots & 0 & k_1 & \dots & k_r & 0 & \vdots \\ 0 & \dots & 0 & 0 & k_1 & \dots & k_r \end{pmatrix}}_{n \text{ columns}} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} k_1 x_1 + k_2 x_2 + \dots + k_r x_r \\ k_1 x_2 + k_2 x_3 + \dots + k_r x_{r+1} \\ \vdots \\ k_1 x_{n-r} + k_2 x_{n-r+1} + \dots + k_r x_{n-1} \\ k_1 x_{n-r+1} + k_2 x_{n-r+2} + \dots + k_r x_n \end{pmatrix} \in \mathbb{R}^{n-r+1}$$

↓ in python: `dens`

$$\Rightarrow \underline{Ax}_i = \sum_{j=1}^r k_j x_{i+j-1}$$

$$\begin{matrix} k_0 x_0 + \dots + k_{r-1} x_{r-1} \\ \vdots \\ k_0 x_{n-r} + \dots + k_{r-1} x_{n-1} \end{matrix}$$

$$\left\{ \begin{pmatrix} k_0 \\ k_1 \\ \vdots \\ k_{r-1} \end{pmatrix} * \begin{pmatrix} x_i \\ x_{i+1} \\ \vdots \\ x_{i+r-1} \end{pmatrix} \right\} \text{ i.e. dot product np.dot}$$

$$k \in \mathbb{R}^r$$

$$\underline{A^T V} = \underbrace{\begin{pmatrix} k_1 & 0 & 0 & \dots & 0 & 0 \\ \vdots & k_1 & 0 & & & \\ & k_r & & & & \\ 0 & 0 & k_r & & & \\ & & & k_1 & & \\ & & & & k_1 & \\ 0 & 0 & 0 & \dots & k_r & \vdots \\ & & & & & 0 & k_1 \end{pmatrix}}_{\substack{\cap \\ \mathbb{R}^{n \times (n-r+1)}}} \underbrace{\begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_{n-r+1} \end{pmatrix}}_{\substack{\cap \\ \mathbb{R}^{(n-r+1) \times 1}}} = \begin{pmatrix} k_1 v_1 \\ k_1 v_1 + k_2 v_2 \\ \vdots \\ k_r v_1 + k_{r-1} v_2 + \dots + k_1 v_r \\ k_r v_2 + k_{r-1} v_3 + \dots + k_1 v_{r+1} \\ \vdots \\ k_r v_{n-2r-1} + k_{r-1} v_{n-2r} + \dots + k_1 v_{n-r} \\ k_r v_{n-2r} + k_{r-1} v_{n-2r+1} + \dots + k_2 v_{n-r+1} \\ \vdots \\ k_r v_{n-r} + k_{r-1} v_{n-r+1} \\ k_r v_{n-r+1} \end{pmatrix} \cap \mathbb{R}^n$$

$$\begin{matrix} \text{reverse} \\ \begin{pmatrix} k_r \\ k_{r-1} \\ \vdots \\ k_2 \\ k_1 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ v_1 \end{pmatrix} \dots [i] \\ \vdots \\ \downarrow v_2 \quad \downarrow \text{sliding} \\ \vdots \\ v_{n-r+1} \\ 0 \\ \vdots \\ 0 \end{matrix} \left. \vphantom{\begin{pmatrix} k_r \\ k_{r-1} \\ \vdots \\ k_2 \\ k_1 \end{pmatrix}} \right\} r-1$$