

## Introduction to Machine Learning

Instructor: Lara Dolecek

TA: Zehui (Alex) Chen

## 1. Matrix calculus review

(a) Gradient of differentiable function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ :

$$\nabla f(x) = \left[ \frac{\partial}{\partial x_1} f(x), \frac{\partial}{\partial x_2} f(x), \dots, \frac{\partial}{\partial x_n} f(x) \right]^T.$$

•  $\nabla_w(w^T b)$ 

$$\frac{\partial w^T b}{\partial w_i} = \frac{\partial \sum_j w_j b_j}{\partial w_i} = b_i \quad \nabla_w(w^T b) = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} = b$$

•  $\nabla_w(\|w\|^2)$ 

$$\frac{\partial \|w\|^2}{\partial w_i} = \frac{\partial \sum_j w_j^2}{\partial w_i} = 2w_i \quad \nabla_w \|w\|^2 = \begin{bmatrix} 2w_1 \\ 2w_2 \\ \vdots \\ 2w_n \end{bmatrix} = 2w$$

•  $\nabla_w(w^T A w)$ 

$$\begin{aligned} \frac{\partial w^T A w}{\partial w_i} &= \frac{\partial \sum_j \sum_k w_j A_{jk} w_k}{\partial w_i} \\ &= \frac{\partial A_{ii} w_i^2}{\partial w_i} + \frac{\partial \sum_{k \neq i} A_{ik} w_k}{\partial w_i} + \frac{\partial \sum_{j \neq i} A_{ji} w_j}{\partial w_i} \\ &= 2w_i A_{ii} + A(i, :) w + A(:, i)^T w \\ &= A(i, :) w + A(:, i)^T w \end{aligned}$$

*Handwritten notes:*  
 $A(i, :)$  is the  $i$ -th row of  $A$  excluding  $A_{ii}$ .  
 $A(:, i)^T$  is the  $i$ -th column of  $A$  excluding  $A_{ii}$ .

•  $\nabla_w(w^T X^T X w)$ 

$$\begin{aligned} \nabla_w w^T X^T X w &= X^T X w + (X^T X)^T w \\ &= 2X^T X w \end{aligned}$$

$$\nabla_w(w^T A w) = \begin{bmatrix} A(1, :) w + A(:, 1)^T w \\ A(2, :) w + A(:, 2)^T w \\ \vdots \\ A(n, :) w + A(:, n)^T w \end{bmatrix} = A w + A^T w$$

(b) Jacobian/derivative matrix of differentiable function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ :

$$f_1, \dots, f_m: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$f(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \\ \vdots \end{bmatrix}: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$J = \begin{bmatrix} \nabla f_1(x)^T \\ \nabla f_2(x)^T \\ \vdots \\ \nabla f_m(x)^T \end{bmatrix}, J_{ij} = \frac{\partial f_i}{\partial x_j}$$

$m \times n$

$$\Delta f(x) \doteq \bigcup_{m \times n} \Delta x_{n \times 1}$$

- $Ax$ 

$$A = \begin{bmatrix} a_1^T \\ a_2^T \\ \vdots \end{bmatrix} \quad Ax = \begin{bmatrix} a_1^T x \\ a_2^T x \\ \vdots \end{bmatrix} \quad A = \begin{bmatrix} A_{11} & A_{12} & A_{13} & \dots \\ A_{21} & \dots & \dots & \dots \end{bmatrix} \quad a_i = \begin{bmatrix} A_{i1} \\ A_{i2} \\ A_{i3} \\ \vdots \end{bmatrix}$$

$$J = \begin{bmatrix} \nabla f_1(x)^T \\ \nabla f_2(x)^T \\ \vdots \\ \nabla f_n(x)^T \end{bmatrix} = \begin{bmatrix} a_1^T \\ a_2^T \\ \vdots \\ a_n^T \end{bmatrix} = A$$

- Example: transformation from polar  $(\underline{r}, \underline{\theta})$  to Cartesian coordinates  $(\underline{x}, \underline{y})$ :  

$$x = r \cos(\theta), y = r \sin(\theta).$$

$$\begin{matrix} \underline{f}_1 & & \underline{f}_2 \end{matrix}$$

$$\begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix} = \begin{bmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{bmatrix} \begin{bmatrix} \Delta r \\ \Delta \theta \end{bmatrix} \quad J = \begin{bmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{bmatrix}$$

J

(c) Hessian matrix for twice differentiable function  $f: \underline{\mathbb{R}^n} \rightarrow \underline{\mathbb{R}}$ :

$$\nabla^2 f(x)_{ij} = \frac{\partial^2}{\partial x_i \partial x_j} f(x).$$

The Hessian matrix is also the derivative matrix  $J$  of the gradient  $\nabla f(x)$ .

- Affine function  $f(x) = a^T x + b$ .

$$\nabla f(x) = a$$

$$\nabla_x^2 f(x) = 0_{n \times n}$$

$x \in \mathbb{R}^n$

- Least squares cost:  $\|Ax - b\|^2$ .

$$\nabla f(x) = 2A^T Ax - 2A^T b$$

$$\nabla_x^2 f(x) = 2A^T A$$

- Example:  $4x_1^2 + 4x_1x_2 + x_2^2 + 10x_1 + 9x_2$

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 8x_1 + 4x_2 + 10 \\ 4x_1 + 2x_2 + 9 \end{bmatrix}$$

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_2} \end{bmatrix} = \begin{bmatrix} 8 & 4 \\ 4 & 2 \end{bmatrix}$$

2. We now try to provide a probabilistic interpretation of the linear regression problem. Consider a model where each of the  $N$  samples is independently drawn according to a normal distribution

$$P(y_n | x_n, w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_n - w^T x_n)^2}{2\sigma^2}\right) \sim N(w^T x_n, \sigma^2)$$

In this model, each  $y_n$  is drawn from a normal distribution with mean  $w^T x_n$  and variance  $\sigma^2$ . The  $\sigma$  are **known**. Write the log likelihood of this model as a function of  $w$ . Show that finding the maximum likelihood estimate of  $w$  leads to the same answer as solving a linear regression problem.

LS Problem :  $\arg\min_w \sum_{i=1}^N (y_i - w^T x_i)^2$

Maximum Likelihood Estimation for  $w$ . give observation

$$\{x_1, y_1\} \{x_2, y_2\} \dots$$

$$\arg\max_w P(y_1, \dots, y_N | x_1, \dots, x_N; w) \quad \text{MLE of } w \rightarrow$$

$$= \arg\max_w \prod_{i=1}^N P(y_i | x_i; w)$$

$$= \arg\max_w \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - w^T x_i)^2}{2\sigma^2}\right)$$

$$= \arg\max_w \sum_{i=1}^N -\frac{(y_i - w^T x_i)^2}{2\sigma^2} \quad + \text{constant}$$

$$= \arg\min_w \sum_{i=1}^N (y_i - w^T x_i)^2$$

3. We now try to provide a probabilistic interpretation of the weighted linear regression. Consider a model where each of the  $N$  samples is independently drawn according to a normal distribution

$$P(y_n|x_n, w) = \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{(y_n - w^T x_n)^2}{2\sigma_n^2}\right) \sim N(w^T x_n, \sigma_n^2)$$

In this model, each  $y_n$  is drawn from a normal distribution with mean  $w^T x_n$  and variance  $\sigma_n^2$ . The  $\sigma_n^2$  are **known**. Write the log likelihood of this model as a function of  $w$ . Show that finding the maximum likelihood estimate of  $w$  leads to the same answer as solving a weighted linear regression. How do  $\sigma_n^2$  relate to  $\alpha_n$ ?

Weighted LS Problem

$$\arg\min_w \sum_{i=1}^N \alpha_i (y_i - w^T x_i)^2$$

$$\arg\max_w P(y_1, \dots, y_N | x_1, \dots, x_N; w)$$

$$= \arg\max_w \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(y_i - w^T x_i)^2}{2\sigma_i^2}\right)$$

$$= \arg\max_w \sum_{i=1}^N -\frac{(y_i - w^T x_i)^2}{2\sigma_i^2} + \text{constant}$$

$$= \arg\min_w \sum_{i=1}^N \frac{1}{2\sigma_i^2} (y_i - w^T x_i)^2$$

$$\alpha_i = \frac{1}{2\sigma_i^2}$$