

CS280 Fall 2019 Assignment 1

Part A

Basic Neural Networks

Due on October 10, 2019, 23:59 UTC+8

Name:Zhang Chengrui

Student ID:2019233183

1. Hessian in Logistic Regression (10 points)

Let $\sigma(a) = \frac{1}{1+e^{-a}}$ be an activation function, the loss function of LR is

$$f(\mathbf{w}) = - \sum_{i=1}^n [y_i \log(\mu_i) + (1 - y_i) \log(1 - \mu_i)],$$

where $\mu_i = \sigma(\mathbf{w}^T \mathbf{x}_i)$.

- Show that the Hessian of f can be written as $H = X^T S X$, where $S = \text{diag}(\mu_1(1 - \mu_1), \dots, \mu_n(1 - \mu_n))$ and $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T$

Solution:

$$\begin{aligned} \frac{\partial u}{\partial w_a} &= \frac{e^{-w^T x_i}}{(1+e^{-w^T x_i})^2} \cdot \frac{\partial (w^T x_i)}{\partial w_a} \\ &= \frac{e^{-w^T x_i}}{(1+e^{-w^T x_i})^2} \cdot x_i \end{aligned}$$

$$\begin{aligned} \therefore \frac{\partial f(\mathbf{w})}{\partial w_a} &= - \sum_{i=1}^n \left[y_i \cdot \frac{e^{-w^T x_i}}{(1+e^{-w^T x_i})^2} \cdot x_i + (1-y_i) \cdot \frac{e^{-w^T x_i}}{(1+e^{-w^T x_i})^2} \cdot x_i \right] \\ &= - \sum_{i=1}^n [y_i - \mu_i] \cdot x_i \end{aligned}$$

$$\text{and } \frac{\partial^2 f(\mathbf{w})}{\partial w_a \partial w_b} = - \sum_{i=1}^n \frac{\partial [y_i - \mu_i] \cdot x_i}{\partial w_b}$$

$$= \sum_{i=1}^n \frac{\partial \mu_i}{\partial w_b} \cdot x_i = \sum_{i=1}^n x_i \cdot \frac{e^{-w^T x_i}}{(1+e^{-w^T x_i})^2} \cdot x_i = \sum_{i=1}^n x_i \cdot \mu_i (1 - \mu_i) \cdot x_i \quad (1)$$

$$\begin{aligned} \text{while } X^T \cdot S \cdot X &= [\mathbf{x}_1, \dots, \mathbf{x}_n] \cdot \begin{bmatrix} \mu_1(1-\mu_1) & & 0 \\ & \ddots & \\ 0 & & \mu_n(1-\mu_n) \end{bmatrix} \cdot \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix} \\ &= [\mathbf{x}_1 \mu_1(1-\mu_1) + \dots + \mathbf{x}_n \mu_n(1-\mu_n)] \cdot \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix} \\ &= \mathbf{x}_1 \mu_1(1-\mu_1) \mathbf{x}_1 + \dots + \mathbf{x}_n \mu_n(1-\mu_n) \mathbf{x}_n \\ &= \sum_{i=1}^n \mathbf{x}_i \mu_i (1 - \mu_i) \mathbf{x}_i \quad (2) \end{aligned}$$

$$\therefore H = (1) = (2)$$

Therefore, Hessian of f can be written as $H = X^T S X$

2. Linear Regression (5 points)

Linear regression has the form

$$f(x) = E[y|x] = b + \mathbf{w}^T \mathbf{x}.$$

- It is possible to solve for \mathbf{w} and b separately. Show that

$$b = \frac{1}{n} \sum_{i=1}^n y_i - \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^T \mathbf{w} = \bar{y} - \bar{\mathbf{x}}^T \mathbf{w}$$

Solution:

$$f(x) = b + \mathbf{w}^T \mathbf{x} = b + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$$\therefore y_i = b + \mathbf{w}^T \mathbf{x}_i$$

$$\therefore \sum_{i=1}^n y_i = nb + \sum_{i=1}^n \mathbf{w}^T \mathbf{x}_i$$

$$\therefore nb = \sum_{i=1}^n y_i - \sum_{i=1}^n \mathbf{w}^T \mathbf{x}_i$$

$$= \sum_{i=1}^n y_i - \sum_{i=1}^n \mathbf{x}_i^T \mathbf{w} \Rightarrow b = \frac{1}{n} \sum_{i=1}^n y_i - \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^T \mathbf{w} \quad \dots (1)$$

$$\text{suppose } \mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{ia}]^T \quad \mathbf{w} = [w_1, w_2, \dots, w_a]^T$$

$$\therefore \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^T \mathbf{w} = \frac{1}{n} \begin{bmatrix} x_{11} w_1 + x_{12} w_2 + \dots + x_{1a} w_a + \\ x_{21} w_1 + x_{22} w_2 + \dots + x_{2a} w_a + \\ \vdots \\ x_{n1} w_1 + x_{n2} w_2 + \dots + x_{na} w_a \end{bmatrix}$$

$$= \frac{1}{n} [w_1(x_{11} + x_{21} + \dots + x_{n1}) + \dots + w_a(x_{1a} + x_{2a} + \dots + x_{na})]$$

$$= \frac{1}{n} \cdot n \cdot \bar{\mathbf{x}}^T \mathbf{w} = \bar{\mathbf{x}}^T \mathbf{w} \quad \dots (2)$$

put (2) into (1)

$$\therefore b = \frac{1}{n} \sum_{i=1}^n y_i - \bar{\mathbf{x}}^T \mathbf{w}$$

$$= \bar{y} - \bar{\mathbf{x}}^T \mathbf{w}$$

3. Gradient descent for fitting GMM (15 points)

Consider the Gaussian mixture model

$$p(\mathbf{x}|\theta) = \sum_k \pi_k \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)$$

Define the log likelihood as

$$l(\theta) = \sum_{n=1}^N \log p(\mathbf{x}_n|\theta)$$

Denote the posterior responsibility that cluster k has for datapoint n as follows:

$$r_{nk} := p(z_n = k|\mathbf{x}_n, \theta) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)}{\sum_{k'} \pi_{k'} \mathcal{N}(\mathbf{x}_n|\mu_{k'}, \Sigma_{k'})}$$

- Show that the gradient of the log-likelihood wrt μ_k is

$$\frac{d}{d\mu_k} l(\theta) = \sum_n r_{nk} \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)$$

- Derive the gradient of the log-likelihood wrt π_k without considering any constraint on π_k . (bonus: with constraint $\sum_k \pi_k = 1$.)
- Derive the gradient of the log-likelihood wrt Σ_k without considering any constraint on Σ_k . (bonus: with constraint Σ_k be a symmetric positive definite matrix.)

Solution:

$$Q_1: \frac{d l(\theta)}{d\mu_k} = \sum_{n=1}^N \frac{1}{p(\mathbf{x}_n|\theta)} \cdot \frac{d p(\mathbf{x}_n|\theta)}{d\mu_k} \dots (1) \quad \frac{d p(\mathbf{x}_n|\theta)}{d\mu_k} = \sum_{k=1}^K \pi_k \frac{d \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)}{d\mu_k} = \pi_k \frac{d \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)}{d\mu_k} \dots (2)$$

$$\frac{d \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)}{d\mu_k} = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_k|^{\frac{1}{2}}} \cdot \exp\left[-\frac{1}{2}(\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)\right] \cdot \left(-\frac{1}{2}\right) \frac{d(\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)}{d\mu_k} \dots (3)$$

$$\frac{d(\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)}{d\mu_k} = \frac{d\mathbf{Y}^T \Sigma^{-1} \mathbf{Y}}{d\mathbf{Y}} \cdot \frac{d\mathbf{Y}}{d\mu_k} = -\Sigma^{-1} + (\Sigma^{-1})^T \cdot \mathbf{Y} = -2\Sigma^{-1} \mathbf{Y} \quad (\mathbf{Y} = (\mathbf{x}_n - \mu_k)) \dots (4)$$

$$\text{put (2) (3) (4) into (1)} \Rightarrow \frac{d l(\theta)}{d\mu_k} = \sum_n r_{nk} \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)$$

$$Q_2: \frac{d l(\theta)}{d\pi_k} = \sum_{n=1}^N \frac{1}{p(\mathbf{x}_n|\theta)} \cdot \frac{d p(\mathbf{x}_n|\theta)}{d\pi_k} \dots (1) \quad \frac{d p(\mathbf{x}_n|\theta)}{d\pi_k} = \frac{d \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)}{d\pi_k} = \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k) \dots (2)$$

$$\text{put (2) into (1)} \Rightarrow \frac{d l(\theta)}{d\pi_k} = \sum_{n=1}^N \frac{\mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)}{p(\mathbf{x}_n|\theta)} = \sum_{n=1}^N \frac{r_{nk}}{\pi_k}$$

$$Q_3: \frac{d l(\theta)}{d\Sigma_k} = \sum_{n=1}^N \frac{1}{p(\mathbf{x}_n|\theta)} \cdot \frac{d p(\mathbf{x}_n|\theta)}{d\Sigma_k} \dots (1) \quad \text{suppose } A = \exp\left[-\frac{1}{2}(\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \mu_k)\right] \text{ then } \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k) = \frac{A}{(2\pi)^{\frac{D}{2}} |\Sigma_k|^{\frac{1}{2}}}$$

$$\begin{aligned} \frac{d p(\mathbf{x}_n|\theta)}{d\Sigma_k} &= \sum_k \pi_k \cdot \frac{d \left(\frac{A}{(2\pi)^{\frac{D}{2}} |\Sigma_k|^{\frac{1}{2}}} \right)}{d\Sigma_k} \cdot A \\ &= \sum_k \pi_k \cdot \left[\frac{1}{(2\pi)^{\frac{D}{2}}} \cdot \left(-\frac{1}{2}\right) \cdot \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \cdot \left(\Sigma_k^{-1}\right)^T \cdot A \right. \\ &\quad \left. + \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma_k|^{\frac{1}{2}}} \cdot A \cdot \left[-\frac{1}{2} \cdot \left(-\Sigma_k^{-T} \cdot (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T \cdot \Sigma_k^{-T}\right)\right] \right] \\ &= \sum_k \pi_k \cdot A \cdot \left(-\frac{1}{2}\right) \cdot \frac{1}{|\Sigma_k|^{\frac{1}{2}}} \cdot \Sigma_k^{-1} \cdot \left[\mathbf{I} - (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T \cdot \Sigma_k^{-1}\right] \\ &= -\frac{1}{2} \left[\Sigma_k^{-T} - \Sigma_k^{-T} (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T \Sigma_k^{-T}\right] \cdot \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k) \dots (2) \end{aligned}$$

$$\text{put (2) into (1)} \quad \frac{d l(\theta)}{d\Sigma_k} = -\frac{1}{2} \sum_{n=1}^N \left[\Sigma_k^{-T} - \Sigma_k^{-T} (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T \Sigma_k^{-T}\right]$$

$$\text{put (2) into (1)} \quad \frac{d l(\theta)}{d\Sigma_k} = \frac{1}{2} \sum_{n=1}^N \left[\Sigma_k^{-T} - \Sigma_k^{-T} (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T \Sigma_k^{-T}\right]$$