

CS280 Fall 2019 Assignment 1

Part A

Basic Neural Networks

Due on October 6, 2019

Name:

Student ID:

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 $\mu_i = \sigma(\mathbf{w}^\top \mathbf{x}_i)$. (Assume $\mathbf{w}, \mathbf{x}_i \in \mathbb{R}^d, X \in \mathbb{R}^{n \times d}, X_i \in \mathbb{R}^{1 \times d}$)

- Show that the Hessian of f can be written as $H = X^\top S X$, where $S = \text{diag}(\mu_1(1 - \mu_1), \dots, \mu_n(1 - \mu_n))$ and $X = [X_1, \dots, X_n]^\top$

2. Linear Regression (5 points)

Linear regression has the form

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

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

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

3. Gradient descent for fitting GMM (15 points)

Consider the Gaussian mixture model

$$p(\mathbf{x}|\theta) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)$$

where $[\dots, \pi_k, \dots] \sim \text{Multinomial}(\phi)$, $\phi_k \geq 0$, $\sum_{j=1}^K \phi_j = 1$. (Assume $\mathbf{x}, \mu_k \in \mathbb{R}^d$, $\Sigma_k \in \mathbb{R}^{d \times d}$)

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:

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

- (5 points) Show that the gradient of the log-likelihood wrt μ_k is

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

- (5 points) Derive the gradient of the log-likelihood wrt π_k , if without considering the constraint on π_k .

Bonus (2 points): what if with the constraint $\sum_k \pi_k = 1$. (hint: reparameterization using the softmax function)

- (5 points) Derive the gradient of the log-likelihood wrt Σ_k .

Bonus (3 points): what if with the constraint that Σ_k is symmetric positive definitive. (hint: reparameterization using Cholesky Decomposition)