

1. We are given a training set  $\{(x^{(i)}, y^{(i)}); i = \{1, \dots, m\}\}$ , where  $x^{(i)} \in \mathbb{R}^n$  and  $y^{(i)} \in \{0, 1\}$ . We consider the Gaussian Discriminant Analysis (GDA) model, which models  $P(x|y)$  using multivariate Gaussian. Writing out the model, we have:

$$P(y = 1) = \phi = 1 - P(y = 0)$$

$$P(x|y = 0) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0)\right)$$

$$P(x|y = 1) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right)$$

The log-likelihood of the data is given by:

$$L(\phi, \mu_0, \mu_1, \Sigma) = \ln P(x^{(1)}, \dots, x^{(m)}, y^{(1)}, \dots, y^{(m)}) = \ln \prod_{i=1}^m P(x^{(i)}|y^{(i)})P(y^{(i)}).$$

In this exercise, suppose we already find  $\mu_0$  and  $\mu_1$ , we want to maximize  $L(\phi, \mu_0, \mu_1, \Sigma)$  with respect to  $\Sigma$ .

- (a) Write down the explicit expression for  $P(x^{(1)}, \dots, x^{(m)}, y^{(1)}, \dots, y^{(m)})$  and  $L(\phi, \mu_0, \mu_1, \Sigma)$ .
- (b) Differentiate  $L(\phi, \mu_0, \mu_1, \Sigma)$  with respect to  $\Sigma$  and set it to 0. Show that the maximum likelihood result for  $\Sigma$  is:

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T.$$

Hints: You may use the following properties without proof:  $a = \text{Tr}(a)$  for scalar  $a$ ;  $\text{Tr}(A) + \text{Tr}(B) = \text{Tr}(A + B)$ ;  $\frac{\partial \ln |A|}{\partial A} = A^{-T}$ ;  $\frac{\partial \text{Tr}(A^{-1}B)}{\partial A} = -(A^{-1}BA^{-1})^T$ .

2. Here is an example where we would want to regularize clusters. Suppose  $n$  students are seated for taking an endterm exam in an  $\mathbf{R}^2$  Euclidean room. There are  $K$  teaching assistants who must collect the answer scripts once the time is up. The TAs need to figure out good locations to position themselves so that the students can walk to the nearest TA and submit their answers. Once the TAs have all the answer sheets, they must return to the front desk located at  $(0,0)$  while handling the returned answer sheets carefully. To reduce the possibility of mishaps related to handling of the papers, write down an objective which can be used to minimize the total distance that both students and TAs need to walk to bring the papers to the front desk. Assume that everyone can walk by taking the shortest path between two points.

3. We learned that the  $\mu$  that maximize

$$\sum_{i=1}^N (x_i - \mu)^2$$

for  $x_i \in \mathbf{R}$  is given by the mean of  $\{x_1, \dots, x_N\}$ , i.e.,  $\mu^* = \frac{1}{N} \sum_{i=1}^N x_i$ .

Show that the  $\mu$  that maximize

$$\sum_{i=1}^N (x_i - \mu)^0$$

is given by the mode of  $\{x_1, \dots, x_N\}$ . What if  $x_i \in \mathbf{R}^n$ ?