ECE M146                                                                                          Homework 6
Introduction to Machine Learning
Instructor: Lara Dolecek
TA: Zehui (Alex) Chen, Ruiyi (John) Wu

**Please upload your homework to Gradescope by May 21, 11:59 pm.**
**Please submit a single PDF directly on Gradescope**
**You may type your homework or scan your handwritten version. Make sure all**
**the work is discernible.**

1. The pdf for two jointly Gaussian random variables $X$ and $Y$ is of the following form parameterized by the scalars $m_1$, $m_2$, $\sigma_1$, $\sigma_2$ and $\rho_{XY}$:

$$f_{X,Y}(x,y) = \frac{\exp\left\{\frac{-1}{2(1-\rho_{XY}^2)}\left[\left(\frac{x-m_1}{\sigma_1}\right)^2 - 2\rho_{XY}\left(\frac{x-m_1}{\sigma_1}\right)\left(\frac{y-m_2}{\sigma_2}\right) + \left(\frac{y-m_2}{\sigma_2}\right)^2\right]\right\}}{2\pi\sigma_1\sigma_2\sqrt{1-\rho_{XY}^2}}. \quad (1)$$

   The pdf for multivariate jointly Gaussian random variable $Z \in \mathbb{R}^k$ is of the following form parameterized by $\mu \in \mathbb{R}^k$ and $\Sigma \in \mathbb{R}^{k \times k}$.

$$f_Z(z) = \frac{\exp\left\{-\frac{1}{2}(z-\mu)^T\Sigma^{-1}(z-\mu)\right\}}{\sqrt{(2\pi)^k|\Sigma|}}. \quad (2)$$

   Suppose $Z = [X, Y]^T$, i.e., $z = [x, y]^T$.

   (a) Find $\mu$, $\Sigma^{-1}$ and $\Sigma$ in terms of $m_1$, $m_2$, $\sigma_1$, $\sigma_2$ and $\rho_{XY}$.

   (b) Suppose $\rho_{XY} = 0$, what is $\Sigma$ in this case? Can you write $f_{X,Y}(x,y)$ as the product of two single variate Gaussian distributions? Are $X$ and $Y$ independent?

2. The Gaussian Discriminant Analysis (GDA) models the class conditional distribution as multivariate Gaussian, i.e, $P(X|Y) \sim \mathcal{N}(\mu_Y, \Sigma)$. Suppose we want to enforce the **Naive Bayes (NB) assumption**, i.e. $P(X_i|Y, X_j) = P(X_i|Y), \forall j \neq i$, to GDA. Show that all off diagonal elements of $\Sigma$ equal to 0: $\Sigma_{i,j} = 0, \forall i \neq j$ with the **NB assumption**.

3. Consider the classification problem for two classes, $C_0$ and $C_1$. In the generative approach, we model the class-conditional distribution $P(x|C_0)$ and $P(x|C_1)$, as well as the class priors $P(C_0)$ and $P(C_1)$. The posterior probability for class $C_0$ can be written as

$$P(C_0|x) = \frac{P(x|C_0)P(C_0)}{P(x|C_0)P(C_0) + P(x|C_1)P(C_1)}.$$

(a) Show that $P(C_0|x) = \sigma(a)$ where $\sigma(a)$ is the *sigmoid* function defined by

$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$

Find $a$ in terms of $P(x|C_0)$, $P(x|C_1)$, $P(C_0)$ and $P(C_1)$.

(b) In the GDA model, we have the class conditional distribution as follows

$$P(x|C_0) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_0)^T\Sigma^{-1}(x - \mu_0)\right),$$

$$P(x|C_1) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T\Sigma^{-1}(x - \mu_1)\right).$$

Suppose we are able to find the maximum likelihood estimation of $\mu_0, \mu_1, \Sigma, P(C_0)$, and $P(C_1)$. Show that $a = w^Tx + b$ for some $w$ and $b$. Find $w$ and $b$ in terms of $\mu_0, \mu_1, \Sigma, P(C_0)$, and $P(C_1)$. This shows that the decision boundary is linear.

(c) In (b), we modeled the class conditional distribution with same covariance matrix $\Sigma$. Now let us consider two classes that have difference covariance matrix as follows

$$P(x|C_0) = \frac{1}{(2\pi)^{n/2}|\Sigma_0|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_0)^T\Sigma_0^{-1}(x - \mu_0)\right),$$

$$P(x|C_1) = \frac{1}{(2\pi)^{n/2}|\Sigma_1|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T\Sigma_1^{-1}(x - \mu_1)\right).$$

Suppose we are able to find the maximum likelihood estimation of $\mu_0, \mu_1, \Sigma_0, \Sigma_1, P(C_0)$, and $P(C_1)$. Show that $a = x^TAx + w^Tx + b$ for some $A$, $w$ and $b$. Find $w$ and $b$ in terms of $\mu_0, \mu_1, \Sigma_0, \Sigma_1, P(C_0)$, and $P(C_1)$. This shows that the decision boundary is quadratic.

4. We are given a training set $\{(x^{(i)}, y^{(i)}); i = \{1, \cdots, m\}\}$, where $x^{(i)} \in R^n$ 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:

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

In this exercise, we want to maximize $L(\phi, \mu_0, \mu_1, \Sigma)$ with respect to $\phi$, $\mu_0$. The maximization over $\Sigma$ is left for discussion.

(a) Write down the explicit expression for $P(x^{(1)}, \cdots, x^{(m)}, y^{(1)}, \cdots, y^{(m)})$ and $L(\phi, \mu_0, \mu_1, \Sigma)$.

(b) Find the maximum likelihood estimate for $\phi$. How do you know such $\phi$ is the "best" but not the "worst"? Hint: Show that the second derivative of $L(\phi, \mu_0, \mu_1, \Sigma)$ with respect to $\phi$ is negative.

(c) Find the maximum likelihood estimate for $\mu_0$. How do you know such $\mu_0$ is the "best" but not the "worst"? Hint: Show that the Hessian Matrix of $L(\phi, \mu_0, \mu_1, \Sigma)$ with respect to $\mu_0$ is negative definite. You may use the following: if $A$ is positive definite, then $A^{-1}$ is also positive definite. Also $B$ is negative definite if $-B$ is positive definite.

5. In this exercise, you will implement a binary classifier using the Gaussian Discriminant Analysis (GDA) model in MATLAB. The data is given in *data.csv*. The first two columns are the feature values and the last column contains the class labels.

   (a) Visualization. Plot the data from different classes in different colors. Is the data linearly separable?

   (b) In the GDA model, we assume the class label follows a Bernoulli distribution and we model the class conditional distribution as multivariate Gaussian with same covariance matrix $(\Sigma)$ and different means ($\mu_0$ and $\mu_1$). Find the maximum likelihood estimate of the parameters $P(y = 0)$ (parameter for the Bernoulli distribution), $\mu_0, \mu_1$ and $\Sigma$ given this data set.

   (c) Using the result you find in Question 3 and your ML estimate of model parameters, find the decision boundary parameterized by $w^T x + b = 0$. Report $w$, $b$ and plot the decision boundary on the same plot.

   (d) Visualize your results by plotting the contour of the two distributions $P(x, y = 0)$ and $P(x, y = 1)$. For consistency, set 'LevelList' ('level' for python) to logspace(-3,-1,7). Does your decision boundary pass through the points where the two distributions have equal probabilities ? Explain why.