

CS 7641 CSE/ISYE 6740 Homework 3

October 8, 2019

- Submit your answers as an electronic copy on T-square.
- No unapproved extension of deadline is allowed. Zero credit will be assigned for late submissions. Email request for late submission may not be replied.
- For typed answers with LaTeX (recommended) or word processors, extra credits will be given. If you handwrite, try to be clear as much as possible. No credit may be given to unreadable handwriting.
- Explicitly mention your collaborators if any.

1 EM application [25 points].

Consider the following problem. There are P papers submitted to a machine learning conference. Each of R reviewers reads each paper, and gives it a score indicating how good he/she thought that paper was. We let $x^{(pr)}$ denote the score that reviewer r gave to paper p . A high score means the reviewer liked the paper, and represents a recommendation from that reviewer that it be accepted for the conference. A low score means the reviewer did not like the paper.

We imagine that each paper has some “intrinsic” true value that we denote by μ_p , where a large value means it’s a good paper. Each reviewer is trying to estimate, based on reading the paper, what μ_p is; the score reported $x^{(pr)}$ is then reviewer r ’s guess of μ_p .

However, some reviewers are just generally inclined to think all papers are good and tend to give all papers high scores; other reviewers may be particularly nasty and tend to give low scores to everything. (Similarly, different reviewers may have different amounts of variance in the way they review papers, making some reviewers more consistent/reliable than others.) We let ν_r denote the “bias” of reviewer r . A reviewer with bias ν_r is one whose scores generally tend to be ν_r higher than they should be.

All sorts of different random factors influence the reviewing process, and hence we will use a model that incorporates several sources of noise. Specifically, we assume that reviewers’s scores are generated by a random process given as follows:

$$\begin{aligned}y^{(pr)} &\sim \mathcal{N}(\mu_p, \sigma_p^2) \\z^{(pr)} &\sim \mathcal{N}(\nu_r, \tau_r^2) \\x^{(pr)}|y^{(pr)}, z^{(pr)} &\sim \mathcal{N}(y^{(pr)} + z^{(pr)}, \sigma^2).\end{aligned}$$

The variables $y^{(pr)}$ and $z^{(pr)}$ are independent; the variables (x, y, z) for different paper-reviewer pairs are also jointly independent. Also, we only ever observe the $x^{(pr)}$ s; thus, the $y^{(pr)}$ s and $z^{(pr)}$ s are all latent random variables.

We would like to estimate the parameters μ_p , σ_p^2 , ν_r , τ_r^2 . If we obtain good estimates of the papers “intrinsic values” μ_p , these can then be used to make acceptance/rejection decisions for the conference.

We will estimate the parameters by maximizing the marginal likelihood of the data $\{x^{(pr)}; p = 1, \dots, P, r = 1, \dots, R\}$. This problem has latent variables $y^{(pr)}$ s and $z^{(pr)}$ s, and the maximum likelihood problem cannot

be solved in closed form. So, we will use EM. Your task is to derive the EM update equations. For simplicity, you need to treat only $\{\mu_p, \sigma_p^2; p = 1 \dots, P\}$ and $\{\nu_r, \tau_r^2; r = 1 \dots R\}$ as parameters. I.e. treat σ^2 (the conditional variance of $x^{(pr)}$ given $y^{(pr)}$ and $z^{(pr)}$) as a fixed, known constant.

1. In this part, we will derive the E-step:

- (a) The joint distribution $p(y^{(pr)}, z^{(pr)}, x^{(pr)})$ has the form of a multivariate Gaussian density. Find its associated mean vector and covariance matrix in terms of the parameters $\mu_p, \sigma_p^2, \nu_r, \tau_r^2$ and σ^2 . [Hint: Recognize that $x^{(pr)}$ can be written as $x^{(pr)} = y^{(pr)} + z^{(pr)} + \epsilon^{(pr)}$, where $\epsilon^{(pr)} \sim \mathcal{N}(0, \sigma^2)$ is independent Gaussian noise.
 - (b) Derive an expression for $Q_{pr}(\theta'|\theta) = \mathbb{E}[\log p(y^{(pr)}, z^{(pr)}, x^{(pr)})|x^{(pr)}, \theta]$ using the conditional distribution $p(y^{(pr)}, z^{(pr)}|x^{(pr)})$ (E-step) (Hint, you may use the rules for conditioning on subsets of jointly Gaussian random variables.)
2. Derive the M-step updates to the parameters μ_p, σ_p^2, ν_r , and τ_r^2 . [Hint: It may help to express an approximation to the likelihood in terms of an expectation with respect to $(y^{(pr)}, z^{(pr)})$ drawn from a distribution with density $Q_{pr}(y^{(pr)}, z^{(pr)})$.]

2 Basic optimization [15 points]

1. Show that log-sum-exp $f(x) = \log(e^{x_1} + \dots + e^{x_n})$ is convex.
2. Show that the log-determinant function $\log \det X$ is concave.
3. (Jensen's inequality.) Use the definition of a concave function f , to show that

$$f\left(\sum_{i=1}^m \alpha_i x_i\right) \geq \sum_{i=1}^m \alpha_i f(x_i)$$

where $\sum_{i=1}^m \alpha_i = 1$, $\alpha_i \geq 0$.

3 Solve simple optimization problems [20 points]

1. Find the minimizer of a linear function

$$f(x) = c^T x$$

on the set

$$V_p = \{x \in \mathbb{R}^n \mid \sum_{i=1}^n |x_i|^p \leq 1\},$$

where $p, 1 < p < \infty$, is a parameter.

2. Find the minimizer of

$$\min_{x > 0: \sum_{i=1}^n x_i = 1} d^T x + \sum_{i=1}^n x_i \log x_i$$

for some given d .

4 Bayes and KNN classifier [20 points]

In this programming assignment, you are going to apply the Bayes Classifier to handwritten digits classification problem. Here, we use the binary 0/1 loss for binary classification, i.e., you will calculate the miss-classification rate as a performance metric.

To ease your implementation, we selected two categories from USPS dataset in `usps-2cls.mat` and provide a skeleton code containing data splitting part, `SplitData.m`, and the main routine `main.m`. `SplitData.m` is designed to split the dataset into training set and testing set according to a predefined ratio.

1. Your first task is implementing the classifier by assuming the covariance matrix is
 - (a) Full matrix, i.e., $\Sigma_1 = \Sigma_2 = \Sigma$ and Σ is a dense matrix, i.e., all entries can be non-zero;
 - (b) Diagonal matrix, $\Sigma_1 = \Sigma_2 = D$, where D is a diagonal matrix and
 - (c) Spherical (the diagonal has a constant value), $\Sigma_1 = \Sigma_2 = \sigma^2 I$.

The first step you will need to do, is to estimate the mean parameter and the covariance parameters using the training data. You can use maximum likelihood estimate to estimate the parameters in each case. Then report the misclassification rate (error rate) over the training set and over the testing set averaged over the 100 random train/test splits for the each of the three covariance matrix cases by using different value of p . Explain and compare the performance of each classifier.

The file you need to edit is `ModelFull.m`, `ModelDiagonal.m` and `ModelSpherical.m` provided with this homework. After implementing these methods, you should evaluate your algorithm on the given set. Repeat 100 times: split the dataset into two parts randomly, use p portion for training and the other $1 - p$ portion for testing. Let p change from 0.1, 0.2, 0.5, 0.8, 0.9.

2. Now repeat the classification again using K -nearest neighbors, for $K = 5, 10, 15, 30$. Repeat 100 times: split the dataset into two parts randomly, use p portion for training and the other $1 - p$ portion for testing. Let p change from 0.1, 0.2, 0.5, 0.8, 0.9. Report the training error and testing error for each case.

5 Logistic regression [20 points]

Repeat the binary classification task above for the same data `usps-2cls.mat` using logistic regression. Repeat 100 times: split the dataset into two parts randomly, use p portion for training and the other $1 - p$ portion for testing. Let p change from 0.1, 0.2, 0.5, 0.8, 0.9. Report the training error and testing error for each case. You may use the matlab built-in function for logistic regression (see class demo code as an example).