

# Homework 1 for Statistical Machine Learning

Instructor: Prof. Jun Zhu

March 11, 2016

You will receive bonus points if you finish the (Bonus) problems.

## 1 Mathematics Basics

### 1.1 Optimization

Use the Lagrange multiplier method to solve the following problem:

$$\begin{aligned} \min_{x_1, x_2} \quad & x_1^2 + x_2^2 - 1 \\ \text{s.t.} \quad & x_1 + x_2 - 1 = 0 \\ & 2x_1 - x_2 \geq 0 \end{aligned}$$

### 1.2 Conjugate Prior

Suppose  $p \sim \text{Beta}(p|\alpha, \beta)$  and  $x|p \sim \text{Bernoulli}(x|p)$ . Show that  $p|x \sim \text{Beta}(p|\alpha + x, \beta + 1 - x)$ , which implies that the Beta distribution can serve as a conjugate prior to the Bernoulli distribution.

Hint:  $\text{Beta}(p|\alpha, \beta) = B(\alpha, \beta)^{-1} p^{\alpha-1} (1-p)^{\beta-1}$ , where  $B(\cdot, \cdot)$  is Euler's Beta function, and  $\text{Bernoulli}(x|p) = p^x (1-p)^{1-x}$ .

### 1.3 Parameter Estimation

There are  $N$  i.i.d. samples  $\mathbf{x}_1, \dots, \mathbf{x}_N$  from a  $K$ -dim multivariate Gaussian distribution  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{K}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}))$ .

1. What is the likelihood of  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  given the observations  $\mathbf{x}_1, \dots, \mathbf{x}_N$ ? (Hint:  $p(\mathbf{x}_1, \dots, \mathbf{x}_N|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ).
2. What is the maximum likelihood estimation (MLE) of  $\boldsymbol{\mu}$ ? Is it unbiased? (Hint: investigate if  $\mathbb{E}f(\mathbf{x}_1, \dots, \mathbf{x}_N) = \boldsymbol{\mu}$ , if  $f(\cdot)$  is your estimator.)
3. Assume the variance  $\boldsymbol{\Sigma}$  is already known, put a multivariate Gaussian prior  $\mathcal{N}(\boldsymbol{\mu}|\boldsymbol{\mu}_0, \lambda^{-1}\boldsymbol{\Sigma})$  on  $\boldsymbol{\mu}$ , what is the maximum *a posteriori* (MAP) estimation of  $\boldsymbol{\mu}$ ? What is the posterior distribution of  $\boldsymbol{\mu}$ ?

4. (Bonus) What is the MLE of  $\Sigma$ ? Is it unbiased? (Hint: you may refer to Appendix C of PRML, or in more detail, Tom Minka’s “Old and New Matrix Algebra Useful for Statistics” for techniques of deriving gradients for matrices. Two useful differentials are  $d \log |X| = \text{tr}(\mathbf{X}^{-1}d\mathbf{X})$  and  $d\mathbf{X}^{-1} = \mathbf{X}^{-1}(d\mathbf{X})\mathbf{X}^{-1}$ .)

## 2 Mixture of Multinomials

### 2.1 MLE for multinomial

Derive the maximum-likelihood estimation for the parameter  $\boldsymbol{\mu} = (\mu_i)_{i=1}^d$  of a multinomial distribution:

$$P(\mathbf{x}|\boldsymbol{\mu}) = \frac{n!}{\prod_i x_i} \prod_i \mu_i^{x_i}, \quad i = 1, \dots, d \quad (1)$$

where  $x_i \in \mathbb{N}$ ,  $\sum_i x_i = n$  and  $0 < \mu_i < 1$ ,  $\sum_i \mu_i = 1$ .

### 2.2 EM for mixture of multinomials

Consider the following mixture-of-multinomials model to analyze a corpus of documents that are represented in the bag-of-words model.

Specifically, assume we have a corpus of  $D$  documents and a vocabulary of  $W$  words from which every word in the corpus is taken. We are interested in counting how many times each word appears in each document, regardless of their positions and orderings. We denote by  $T \in \mathbb{N}^{D \times W}$  the word occurrence matrix where the  $w$ th word appears  $T_{dw}$  times in the  $d$ th document.

According to the mixture-of-multinomials model, each document is generated i.i.d. as follows. We first choose for each document  $d$  a latent “topic”  $c_d$  (analogous to choosing for each data point a component  $z_n$  in the mixture-of-Gaussians) with

$$P(c_d = k) = \pi_k, \quad k = 1, 2, \dots, K; \quad (2)$$

And then given this “topic”  $\boldsymbol{\mu}_k = (\mu_{1k}, \dots, \mu_{Wk})$  which now simply represents a categorical distribution over the entire vocabulary, we generate the word bag of the document from the corresponding multinomial distribution<sup>1</sup>

$$P(d|c_d = k) = \frac{n_d!}{\prod_w T_{dw}!} \prod_w \mu_{wk}^{T_{dw}}, \quad \text{where } n_d = \sum_w T_{dw}. \quad (3)$$

Hence in summary

$$P(d) = \sum_{k=1}^K P(d|c_d = k)P(c_d = k) = \frac{n_d!}{\prod_w T_{dw}!} \sum_{k=1}^K \pi_k \prod_w \mu_{wk}^{T_{dw}}. \quad (4)$$

---

<sup>1</sup>Make sure you understand the difference between a categorical distribution and a multinomial distribution. You may think about a Bernoulli distribution and a binomial distribution for reference.

Given the corpus  $T$ , please design and implement an EM algorithm to learn the parameters  $\{\boldsymbol{\pi}, \boldsymbol{\mu}\}$  of this mixture model and test it on the NIPS dataset (<http://ml.cs.tsinghua.edu.cn/~jianfei/static/nips.tar.gz>).

Set the number of topics  $K$  to be 5, 10, 20, 30 respectively and show the most-frequent words in each topic for each case. Observe the result and try to find the “best”  $K$  value for this dataset and explain why.

(Bonus) Compare the log-likelihood (Eq. 4) of different  $K$ .

(Bonus) Try different initialization strategies, and report the difference of results.

## 3 PCA

### 3.1 Minimum Error Formulation

Complete the proof on the lecture slide which shows that PCA can be equivalently formulated as minimizing the mean-squared-error of a low-dimensional approximation from a subset of orthonormal basis.

### 3.2 MNIST

Implement the PCA algorithm and test it on the MNIST dataset<sup>2</sup>. Plot the principle components, and the reconstructed image using 1, 5, 20, 100 components. Rerun your PCA implementation without centering the dataset (subtracting the sample mean) and report the difference between the results.

---

<sup>2</sup><http://yann.lecun.com/exdb/mnist/>, or <http://www.cs.nyu.edu/~roweis/data.html> for a MATLAB version.