

CSC411 Fall 2018: Homework 5

Niloufar Afsariardchi

November 15, 2018

1 Problem 1

a) Based on Bayes rule:

$$P(y = k|\mathbf{x}, \mu, \Sigma) = \frac{P(\mathbf{x}|y = k, \mu, \Sigma)P(y = k)}{P(\mathbf{x}|\mu, \sigma)} \quad (1)$$

$$= \frac{P(\mathbf{x}|y = k, \mu, \Sigma)P(y = k)}{\sum_t P(\mathbf{x}|y = t, \mu, \sigma)P(y = t)} \quad (2)$$

$$\log P(y = k|\mathbf{x}, \mu, \Sigma) = \log P(\mathbf{x}|y = k, \mu, \Sigma) + \log P(y = k) \quad (3)$$

$$- \log \sum_t P(\mathbf{x}|y = t, \mu, \sigma)P(y = t) \quad (4)$$

$$\log P(y = k|\mathbf{x}, \mu, \Sigma) = \log P(\mathbf{x}|y = k, \mu, \Sigma) + \log 0.1 \quad (5)$$

$$- \log \sum_t 0.1P(\mathbf{x}|y = t, \mu, \sigma) \quad (6)$$

where the first term of RHS is given in the equation (1) of the handout. The LHS of above equation is conditional log-likelihood or log posterior distribution.

These are the results I got for average conditional log-likelihood:

conditional log-likelihood for training set: -0.125

conditional log-likelihood for test set: -0.197

conditional likelihood for training set: 0.88

conditional likelihood for test set: 0.82

This shows that the test set has the higher inherent uncertainty and perhaps is has more noise than the training set.

b) Here is the accuracy of the algorithm obtained by assigning each datum to the digit with highest conditional log likelihood:

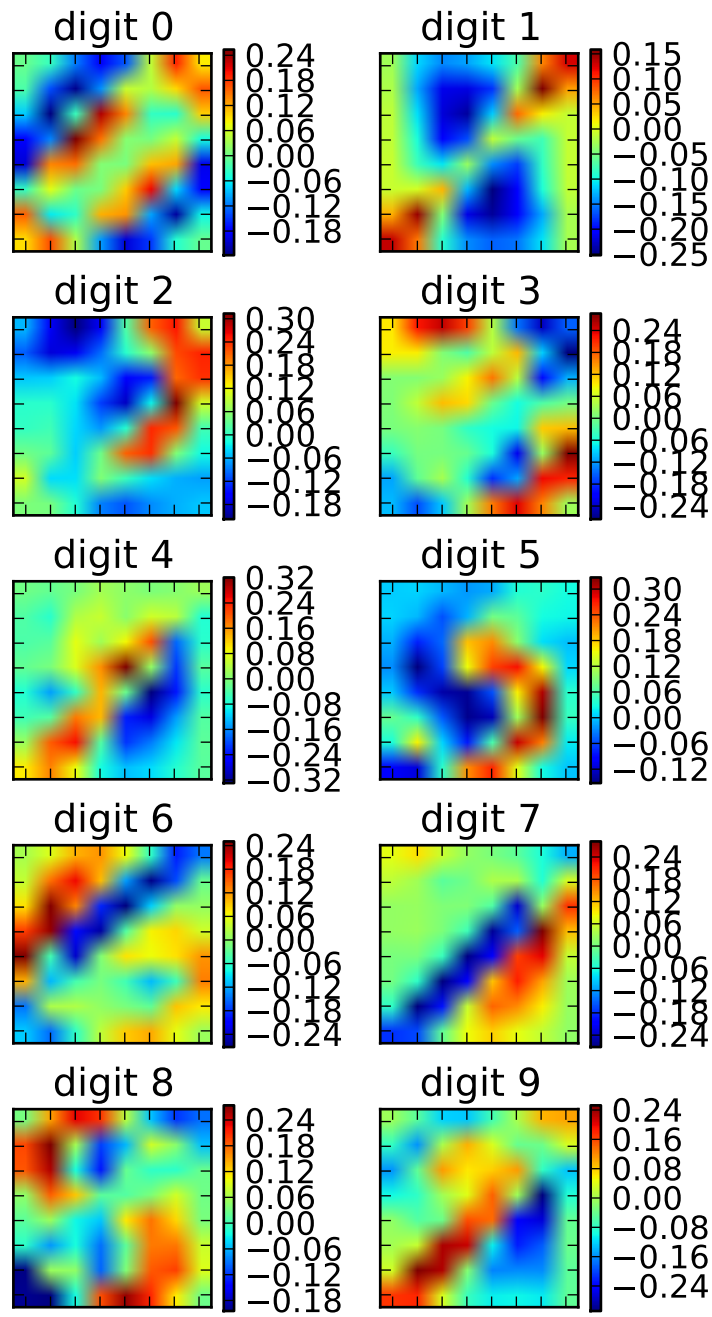
accuracy of training set: 0.98

accuracy of test set: 0.97

c) See attached Figure 1 below. There are some patterns in the eigenvectors, however they are not very clear. It is not expected for eigenvectors to reflect the digits in the exact form either.

2 Problem 2

Continue to next page



Question 2

$$(a) P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

assuming i.i.d x_i

$$\frac{\prod_{i=1}^m P(x^{(i)}|\theta) P(\theta)}{\prod_{i=1}^m P(x^{(i)})}$$

$$\propto \frac{\prod_{i=1}^m \prod_{k=1}^K \theta_k^{x_k^{(i)} - a_k - 1}}{\prod_{i=1}^m \int_{\theta} P(x^{(i)}|\theta') P(\theta') d\theta'}$$

$$= \frac{\prod_{k=1}^K \theta_k^{N_k - a_k - 1}}{\prod_{i=1}^m \int \theta_{x^{(i)}} \text{Dir}(\theta) d\theta}$$

Let's assume $x^{(i)}$ is the index of the category for which $x_k^{(i)} = 1$

$$= \frac{\prod_{k=1}^K \theta_k^{N_k - a_k - 1}}{\prod_{i=1}^m E[\theta_{x^{(i)}}]}$$

$$= \frac{\prod_{k=1}^K \theta_k^{N_k - a_k - 1}}{\prod_{i=1}^m \frac{a_{x^{(i)}}}{\sum_{x'} a_{x'}}$$

$$= \frac{(\sum_{x'} a_{x'})^K \prod_{k=1}^K \theta_k^{N_k - a_k - 1}}{\prod_{k=1}^K a_k^{N_k}}$$

constant

$$= \left(\sum_{x'} a_{x'} \right)^K \prod_{k=1}^K \left(\frac{\theta_k}{a_k} \right)^{N_k - a_k - 1} \theta_k$$

$$\propto \prod_{k=1}^K \left(\frac{\theta_k}{a_k} \right)^{N_k - a_k - 1} \theta_k$$

← which is Dirichlet

$$P(\theta|D) \sim \text{Dirichlet}(a_1 - N_1, a_2 - N_2, \dots, a_K - N_K)$$

Posterior predictive distribution is

$$P(D'|D) = P(x_{k'}=1|D)$$

that the next observation belongs to k -th category.

$$= \int P(\theta|D) P(x_{k'}=1|\theta) d\theta$$

$$\propto \int \left(\prod_{k=1}^K \left(\frac{\theta_k}{a_k} \right)^{N_k - a_k - 1} \right) \cdot \theta_{k'} d\theta$$

We know $\int P(\theta_k|D) d\theta_k = 1$

$$= \int P(\theta_{k'}|D) \theta_{k'} d\theta_{k'}$$

$$= \int \text{Dir}(a_{k'}+N_{k'}) \theta_{k'} d\theta_{k'}$$

$$= E(\theta_{k'})_{P(\theta|D)}$$

$$= \frac{a_{k'} + N_{k'}}{-N + \sum_{k=1}^K a_k}$$

$$(b) \hat{\theta} = \underset{\theta}{\text{argmax}} \log P(\theta|D)$$

$$= \underset{\theta}{\text{argmax}} \log \text{Dirichlet}(a_1+N_1, \dots, a_K+N_K)$$

$$= \underset{\theta}{\text{argmax}} \log \prod_{k=1}^K \theta_k^{N_k - a_k - 1} = \underset{\theta}{\text{argmax}} J$$

$$J = \sum_{k=1}^K (N_k - a_k - 1) \log \theta_k$$

Optimization:

$$\underset{\theta}{\text{maximize}} J(\theta)$$

$$\text{assuming } \sum_{k=1}^K \theta_k = 1$$

$$\text{maximize } Q = J(\theta) - \lambda \left(1 - \sum_{k=1}^K \theta_k\right)$$

$$\textcircled{1} \quad \frac{\partial Q}{\partial \theta_i} = \frac{N_i - a_i - 1}{\theta_i} + \lambda = 0 \Rightarrow \theta_i = \frac{1 + a_i - N_i}{\lambda}$$

$$\textcircled{2} \quad \frac{\partial Q}{\partial \lambda} = 1 - \sum_{k=1}^K \theta_k = 0$$

$$\text{from } \textcircled{2}: \sum_{i=1}^K \frac{1 + a_i - N_i}{\lambda} = 1$$

$$\Rightarrow \lambda = \sum_{i=1}^K 1 + a_i - N_i$$

$$\text{from } \textcircled{1} \Rightarrow \theta_i = \frac{1 + a_i - N_i}{-N + K + \sum_{i=1}^K a_i}$$

which is a little bit different from predictive distribution