

# Machine Learning 1 - Homework 5

Monday, October 10, 2016

Deadline: Wednesday, October 19, 2016, 23:59

## 1 PCA

Suppose we have a dataset of  $N$  vectors  $\{\mathbf{x}_n\}$  of dimension  $D$ . We can write the entire dataset as a  $D$  by  $N$  matrix  $\mathbf{X}$  (column  $n$  is  $x_n$ ). We may wish to perform PCA on this data in the original data space, or in *kernel*-space using kernel-PCA. In the latter case, the data are projected into *feature* space  $\phi$ , such that  $\phi_n = \phi(\mathbf{x}_n)$  is  $M$ -dimensional feature space representation of  $x_n$ . Consider the procedure for PCA (which can be generalized to kernel-PCA):

**Step 1** Center  $\mathbf{X}$ , producing a center data matrix  $\hat{\mathbf{X}}$ .

**Step 2** Compute sample covariance  $\mathbf{S}$  of the centered dataset.

**Step 3** Solve the eigen-value problem  $\mathbf{S} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ , where  $\mathbf{U}$  is a column matrix of eigen-vectors and  $\mathbf{\Lambda}$  is a diagonal matrix of eigen-values  $\lambda_k$ , ie  $\mathbf{\Lambda}_{kl} = \lambda_k\delta_{kl}$ , where  $\delta_{kl} = 1$  iff  $k = l$ .

**Step 4** Pick eigen-vectors with largest eigen-values  $\{\mathbf{u}_1, \dots, \mathbf{u}_K\}$ .

**Step 5** Project data onto  $K$ -dimensional manifold.

Answer the following questions:

- (a) Provide an expression for  $\hat{\mathbf{x}}_n$ .
- (b) Prove that the average of  $\hat{\mathbf{x}}_n$  (over  $N$  data vectors) is the 0 vector.
- (c) Provide an expression for  $\mathbf{S}$  in terms of  $\hat{\mathbf{X}}$ .
- (d) What is the dimensionality of  $\mathbf{S}$ ?
- (e) What is the expression for the linear projection  $\mathbf{L}$  that maps data vectors  $\hat{\mathbf{x}}_n$  onto a  $K$ -dimensional sub-space,  $y_n = \mathbf{L}\hat{\mathbf{x}}_n$ , such that it has zero mean and identity covariance. Prove that the average over  $N$  of  $y_n$  is 0. Prove that the covariance of  $y_n$  is the identity. What is this operation called?

## 2 Mixture Models

Consider a data distribution whose underlying generating process is a mixture of Poisson distributions, but we do not know the parameters of the mixture model. In this question you are asked to derive the update equations for the general Poisson mixture model.

The Poisson distribution is:

$$P(x|\lambda) = \frac{1}{x!} \lambda^x \exp(-\lambda)$$

where  $x = 0, 1, 2, \dots$  (non-negative integers),  $\lambda > 0$  is the ‘rate’ of the data; the expected value of  $x$  is  $\lambda$ . A mixture representation assumes the following:

$$P(x_n) = \sum_{k=1}^K \pi_k P(x_n|\lambda_k)$$

where  $P(x_n|\lambda_k)$  is a Poisson distribution with rate  $\lambda_k$  and  $x_n$  is a single data observation. To answer the following questions assume we are given a dataset  $\{x_1, x_2, \dots, x_N\}$ . Make sure that the constraint  $\sum_k \pi_k = 1$  is satisfied (i.e. think of the log-likelihood or log-joint as  $f$  (an objective to maximize) and  $\sum_k \pi_k - 1 = 0$  as  $g = 0$  (a constraint that must hold)).

- (a) Write down the likelihood (as usual) for the data set in terms of  $\{x_1, x_2, \dots, x_N\}$ ,  $\{\pi_k\}$ ,  $\{\lambda_k\}$ .
- (b) Write down the log-likelihood (as usual) for the data set in terms of  $\{x_1, x_2, \dots, x_N\}$ ,  $\{\pi_k\}$ ,  $\{\lambda_k\}$ .
- (c) Find the expression for the responsibilities  $r_{nk}$ .
- (d) Find the expression for  $\lambda_k$  that maximizes the log-likelihood.
- (e) Find the expression for  $\pi_k$  that maximizes the log-likelihood.
- (f) Now assume priors for  $\pi_k$  and  $\lambda_k$ .  $p(\lambda_k|a, b) = \mathcal{G}(\lambda_k|a, b)$  (a Gamma prior) and  $p(\pi_1, \dots, \pi_k) = \mathcal{D}(\pi_1, \dots, \pi_k|\alpha/K, \dots, \alpha/K)$  (a Dirichlet distribution). These distributions are defined in the appendix of Bishop. Write down the log-joint distribution  $\log p(\mathbf{x}_1, \dots, \mathbf{x}_N, \{\pi_k\}, \{\lambda_k\}|a, b, \alpha, K)$ .
- (g) Find the expression for  $\lambda_k$  that maximizes the log-joint.
- (h) Find the expression for  $\pi_k$  that maximizes the log-joint.

- (i) Write down an iterative algorithm using the above update equations (similar to the ones derived in class for the Mixture of Gaussians); include initialization and convergence check steps.