

Machine Learning: Pattern Recognition

Final Exam

Wednesday, October 23, 2013
9:00 - 12:00

Before you start

1. Indicate your name and student number of everything you hand in.
2. On the first page also list the master program you are currently following and your previous education (e.g. bachelor in XX at the U. of YY).
3. This is a closed book exam. You are allowed a single double-sided A4 cheat sheet. No calculators are required.

1 Kernel Ridge Regression

/25

We are given the following dataset: $\{\mathbf{x}_i, y_i\}$, $i = 1..N$, where each $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$. We are also given a collection of feature functions $\{\phi_a(\cdot)\}$, $a = 1..A$. Now consider the following optimization problem for ridge regression,

$$\min_{\mathbf{w}} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2$$

subject to: $\|\mathbf{w}\|^2 \leq B$

- (a) Provide an expression for the Lagrangian. (Check that the sign of the Lagrange multiplier is correct). /5
- (b) Is this a convex optimization problem (why?). Assume you are given the expression for $\mathbf{w}^*(\lambda)$ (which you will compute in the next question). Write the *dual optimization problem* in terms of $\mathbf{w}^*(\lambda)$ (don't forget dual constraints!). Is this problem concave? /6
- (c) Write down all KKT equations and solve for \mathbf{w} . You may assume that λ is fixed and known. /6
- (d) Assume we have a test case \mathbf{x}^* and you are given a kernel $K(\cdot, \cdot)$ (assume you do not know the features ϕ explicitly or that there are infinitely many of them). Provide an expression for the predicted value of y using the above ridge regression model. The expression may only involve kernel evaluations (instead of feature evaluations). /6
- (e) Assume the kernel does not have any free parameters to change. Now consider a situation where you suspect overfitting. What would you do to reduce this overfitting? /2

2 K-means Clustering

/10

We are given data $\{x_i\}$, $i = 1..N$ with $x_i \in \mathbb{R}^d$. Consider minimizing the following K-means cost function,

$$C = \frac{1}{2} \sum_c \sum_{i \in S_c} \|x_i - \mu_c\|_{L_2}^2$$

where S_c is the subset of data-items assigned to cluster c and L_2 means the L_2 norm (the same as used in class).

- (a) Derive the K-means update rule for μ_c by computing the gradient $\partial C / \partial \mu_c$ and equating it to 0. It is implicitly assumed that the assignments are held fixed. /5
- (b) Is the K-means algorithm (which alternates the above updates for $\{\mu_c\}$ with reassigning data-items to clusters) guaranteed to converge eventually? Why? /5

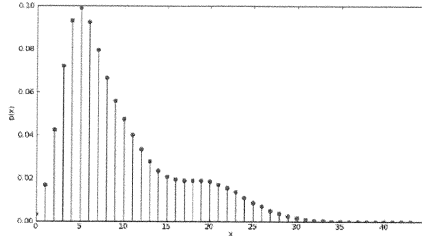


Figure 1: The empirical distribution of data drawn from a mixture of $K = 3$ Poisson distributions.

3 Mixture Models

/25

Consider the data distribution shown in Figure 1. Each vertical line represents the empirical probability distribution of a dataset of discrete data values x (the frequency of a particular value x in the dataset). We are told that the 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 λ . 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 λ_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.

- Write down the log-likelihood (as usual) for the data set in terms of $\{x_1, x_2, \dots, x_N\}$, $\{\pi_k\}$, $\{\lambda_k\}$. /5
- Find the expression for π_k that maximizes the log-likelihood. /5
- Find the expression for λ_k that maximizes the log-likelihood. /5
- Find the expression for the responsibilities r_{nk} . /5
- 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. /5

4 PCA and kernel-PCA

/30

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 \mathbf{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 ϕ , such that $\phi_n = \phi(\mathbf{x}_n)$ is M -dimensional feature space representation of \mathbf{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 λ_k , ie $\Lambda_{kl} = \lambda_k \delta_{kl}$, 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 $\hat{\mathbf{x}}_n$. /2
- (b) Prove that the average of $\hat{\mathbf{x}}_n$ (over N data vectors) is the 0 vector. /4
- (c) Provide an expression for \mathbf{S} in terms of $\hat{\mathbf{X}}$. /2
- (d) What is the dimensionality of \mathbf{S} ? /2
- (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, $\mathbf{y}_n = \mathbf{L}\hat{\mathbf{x}}_n$, such that it has zero mean and identity covariance. Prove that the average over N of \mathbf{y}_n is 0. Prove that the covariance of \mathbf{y}_n is the identity. What is this operation called? /10
- (f) For kernel-PCA, the centering step cannot in general be performed in the feature space.
 - (i) Write the equation for centered feature vector $\hat{\phi}_n$. /2
 - (ii) The eigen-vector problem solved in kernel-PCA uses a kernel matrix $\hat{\mathbf{K}}$ that is centered in kernel-space. Using your result above for the center feature vector, expand the gram matrix entry $\hat{\mathbf{K}}_{nm} = \hat{\phi}_n^T \hat{\phi}_m$ in terms of \mathbf{k} only (ie in terms of the non-centered kernel functions). /6
- (g) In terms of memory requirements, what is the disadvantage of kernel-PCA versus PCA? /2

5 Decision Trees

/10

We are given the following dataset: $\{\mathbf{x}_i, y_i\}$, $i = 1..N$, where x_{ai} is the a 'th attribute ($a = 1..A$) of the i 'th data-case and $y_i = \{-1, +1\}$. Each x_a can take one of 2 discrete values $\{A, B\}$.

- (a) Assume that you have partially trained a decision tree and are considering to add another attribute to the tree at some branch. At this point you may assume you have n (negative) data cases in class -1 and p (positive) data cases in class $+1$. Moreover, after you apply the attribute, you will have n_A negative data cases with attribute value A , n_B negative data cases with attribute value B , p_A positive data cases with attribute value A and p_B positive data cases with attribute value B . Provide the expression for the *information gain* for this attribute.

/3

- (b) At some point you find that in some branch you have used all your attributes, but that there are still some positive and negative data-cases that did not get resolved (or split). Say there are p_L positive and n_L negative data-cases left in this leaf L . Consider now a test case that ends up at this leaf node of this branch in the decision tree. What procedure would you follow to classify this test case as either positive or negative?

/2

- (c) Describe how the random forest algorithm works (based on decision trees). Again, use pseudo-code style to explain your answer.

/5