The University of Texas at Austin Department of Electrical and Computer Engineering

EE381V: Large Scale Learning — Spring 2013

Assignment 1

Caramanis/Sanghavi

Due: Thursday, Feb. 7, 2013.

(Problems 1 and 2 have been adapted from Jure Leskovec's course on 'Mining Massive Datasets')

1. Locality Sensitive Hashing

The first week of class discussed locality sensitive hashing. For some standard similarity functions, like the Jaccard similarity, we showed that there corresponds a locality sensitive hashing scheme. As it turns out, not all similarity functions have a locality sensitive hashing scheme. This problem explores this issue.

Recall that a locality sensitive hashing scheme is a set F of hash functions that operate on a set S of objects, such that for two objects $x, y \in S$,

$$\Pr_{h \in \mathcal{F}} \left[h(x) = h(y) \right] = \sin(x, y)$$

where $sim(\cdot): S \times S \to [0,1]$ is a pairwise function (the similarity function).

• Let d(x, y) = 1 - sim(x, y). Prove that for $sim(\cdot)$ to have a locality sensitive hashing scheme, d(x, y) should satisfy the triangle inequality.

$$d(x,y) + d(y,z) \ge d(x,z)$$

for all $x, y, z \in S$.

• Consider the following two similarity functions: The so-called Overlap similarity function,

$$\sin_{Over}(A, B) = \frac{|A \bigcap B|}{\min(|A|, |B|)}$$

and the Dice similarity function,

$$\operatorname{sim}_{Dice}(A,B) = \frac{2|A \cap B|}{(|A|+|B|)},$$

where A and B are two sets.

Is there a locality sensitive hashing scheme for either? Prove, or disprove, giving a counterexample.

Prove or disprove (give counterexamples). Let A, B be any two sets.

2. Approximate Near Neighbor Search using LSH

LSH has been used for nearest neighbor search, in numerous applications. This problems explores this.

Given a data set \mathcal{A} , along with a distance function $d(\cdot)$, the Nearest Neighbor problem says the following: given a query point, z, return its nearest neighbors, w.r.t. $d(\cdot)$. The approximate nearest neighbor problem is an approximation in two respects: it only requires us to know the immediate neighborhood of any given point, and also, we need only return approximate nearest neighbors, up to a dilation factor λ . More precisely, the (c, λ) -Approximate Near Neighbor (ANN) problem is as follows: Given a query point, z, for which (we assume) there exists a point $x \in \mathcal{A}$ with $d(x, z) \leq \lambda$, return a point x' from the dataset with $d(x', z) \leq c\lambda$.

We outline an approximate nearest neighbor algorithm, and then, through the parts of this problem, show that with large probability, it outputs a *c*-approximate nearest neighbor, as explained above.

Let \mathcal{A} be a dataset with n points from a metric space with distance measure d(). Let \mathcal{H} be a $(\lambda, c\lambda, p_1, p_2)$ locally sensitive family of hash functions for the distance measure d(). Let $\mathcal{G} = \mathcal{H}^k = \{g = (h_1, \ldots, h_k) | h_i \in \mathcal{H}\}$, where $k = \log_{1/p_2}(n)$ be an amplified family. Choose $L = n^{\rho}$ random members $g_1, \ldots, g_L \in \mathcal{G}$, where $\rho = \frac{\log(1/p_1)}{\log(1/p_2)}$. We then do the following: (a) hash all the data points and the query point z using all g_i 's $(1 \leq i \leq L)$; (b) retrieve at most 3L data points from the buckets $g_j(z)$ $(1 \leq j \leq L)$, and (c) report the closest one as a (c, λ) -ANN.

• Define $W_j = \{x \in \mathcal{A} | g_j(x) = g_j(z)\}$ $(1 \le j \le L)$ as the random set of data points x hashed to the same bucket as the query point z by the hash function g_j . Let $T = \{x \in \mathcal{A} | d(x, z) > c\lambda\}$. Prove:

$$Pr\left[\sum_{j=1}^{L} |T \bigcap W_j| > 3L\right] < \frac{1}{3}.$$

• Let $x^* \in \mathcal{A}$ be a data point such that $d(x^*, z) \leq \lambda$. Prove:

$$Pr\left[g_j(x^*) \neq g_j(z), \forall 1 \le j \le L\right] < \frac{1}{e}.$$

- Find a bound on δ , the probability that the reported point is an actual (c, λ) -ANN.
- 3. This problem tests empirically how nearest-neighbor search using LSH compares to linear search. For this, we provide a dataset of images. Each column in the dataset is a vectorized 20 × 20 image patch. Download the image set and matlab code here: http://http://users.ece.utexas.edu/~cmcaram/LSL2013/lsh.zip¹ and see the ReadMe.txt file for instruction on using the code, and in particular, the functions lsh and lshlookup.

In this problem we use the ℓ_1 distance measure. The LSH function is run with L = 10, k = 24, where L is the number of hash tables generated and k is the length/number of bits of the hash key.

¹This is the same as that provided from the problem's source, Jure Leskovec's course on Mining Massive Data Sets. The dataset and code are adapted from Brown University's Greg Shakhnarovich

• Consider the image patches z_j , of column $100 \times j$, for $j \in \{1, 2, ..., 10\}$. Find the top 3 nearest neighbors for these image patches, (excluding the original patch) using both LSH and linear search.

Compare the average search time for LSH and linear search. (If the bucket contains less than 3 points, rehash till you get enough neighboring points).

• For each z_j $(1 \le j \le 10)$ let $\{x_{ij}\}_{i=1}^3$ denote the approximate near neighbors of z_j found using LSH, and $\{x_{ij}^*\}_{i=1}^3$ to be the actual top 3 near neighbors of z_j found using linear search. Compute the following error measure:

$$error = \frac{1}{10} \sum_{j=1}^{10} \frac{\sum_{i=1}^{3} d(x_{ij}, z_j)}{\sum_{i=1}^{3} d(x_{ij}^*, z_j)}$$

Plot the error value as a function of L (for L = 10, 12, 14, ..., 20, with k = 24). Then plot the error values as a function of k (for k = 16, 18, 20, 22, 24 with L = 10).

• Plot the top 10 near neighbors found using the two methods (using the default L = 10, k = 24) for the image patch in column 100, together with the image patch itself. Use functions reshape() and mat2gray() to convert the matrices to images, then use the functions imshow() and subplot() to display the images. Compare them visually.

4. K-Mean vs. EM for Gaussian Mixture Models

Consider samples being generated from a Gaussian mixture model with the following pdf

$$p(x|z=k) = \frac{1}{2\pi |\Sigma_k|^{1/2}} \exp\left[-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1}(x-\mu_k)\right], \quad k \in \{1,\dots,K\}$$
$$p(x) = \frac{1}{K} \sum_{k=1}^K p(x|z=k)$$

where $x, \mu_k \in \mathbb{R}^2, z \in \{1, ..., K\}, K = 3, \Sigma_k \in \mathbb{R}^{2 \times 2}$ are the covariance matrices. Let $\mu_1 = [5, 5]^T, \mu_2 = [10, 20]^T, \mu_3 = [16, 8]^T$ and

$$\Sigma_1 = \sigma^2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \Sigma_2 = \sigma^2 \begin{bmatrix} 1.2 & 0 \\ 0 & 1.2 \end{bmatrix}, \quad \Sigma_3 = \sigma^2 \begin{bmatrix} .28 & .42 \\ .42 & 1.4 \end{bmatrix}$$

• We want to cluster the points generated from the distribution p(x) using the EM algorithm. Generate a sample set $S = \{x_i\}_{i=1}^n$ of size n = 300 from the distribution p(x), using $\sigma^2 = 2$. Starting with initial estimates of the means as $\hat{\mu}_1^{(0)} = [0, 0]^T$, $\hat{\mu}_2^{(0)} = [10, 15]^T$, $\hat{\mu}_3^{(0)} = [16, 0]^T$ and covariance $\hat{\Sigma}_k^{(0)} = I$, the 2 × 2 identity matrix, find the final estimates of the means $\hat{\mu}_k$ and covariances $\hat{\Sigma}_k$, $k \in \{1, 2, 3\}$ of the three components of p(x) using the EM algorithm. Now using the MAP estimates $\hat{z}_i = \arg \max_k P(z_i = k|S)$, obtained by the EM algorithm, cluster the sample set S in 3 clusters. Plot the clusters using the matlab function *scatter()* using separate colors for each cluster (or plot each cluster in separate graphs).

- Define the set $E_{\pi} = \{x_i : \pi(\hat{z}_i) \neq z_i\}$, for a particular permutation π of the labels of the clusters. Define error set $E = E_{\pi_0}$, where $\pi_0 = \arg \min_{\pi} |E_{\pi}|$. Hence set E contains the sample points that fall in the wrong cluster. The error fraction is given by $e = \frac{|E|}{n}$. Define the probability of error $P_A(E, \sigma^2)$ for a particular clustering algorithm A as the average e over several sample sets S drawn from the same distribution p(x). Now generate several sample sets varying σ^2 between 1 and 30 (also many sample sets for each σ^2) and cluster using both EM and K-mean algorithms starting with same set of initial mean estimates as given in part (a). Plot the probability of error $P_A(E, \sigma^2)$ as a function of σ^2 for both the algorithms. Also plot the average run-time of both the algorithms for different σ^2 .
- Define the algorithm *B* as follows. First run the K-mean algorithm to obtain mean estimates $\hat{\mu}'_k$. Then run the EM algorithm using $\hat{\mu}'_k$ as the initial mean estimates. Plot probability of error $P_B(E, \sigma^2)$ vs. σ^2 for algorithm *B*.
- 5. A few details about spectral clustering.
 - Suppose that $\{u_1, \ldots, u_k\}$ and $\{\hat{u}_1, \ldots, \hat{u}_k\}$ are any two orthonormal bases for the same k-dimensional null-space of a matrix L. Let U and \hat{U} denote the $n \times k$ matrices whose columns are the respective orthonormal bases. Show that there is an orthonormal $k \times k$ matrix Q, for which $\hat{U} = UQ$. Converseley, show that if U is an $n \times k$ orthonormal matrix, and Q a $k \times k$ orthonormal matrix, then $\hat{U} = UQ$ is also orthonormal.
 - Recall that if we have a graph with k connected components, then the Laplacian has a k-dimensional subspace, spanned by k vectors, $\{u_1, \ldots, u_k\}$, where u_i has support only on the elements corresponding to the nodes in the i^{th} connected component. Hence, if we let U be the matrix with these vectors as *columns*, and then let $\{y_1, \ldots, y_n\}$ be the rows of U, then if we normalize each y_i , each point maps to the standard basis element $e_{c(i)}$, with c(i) corresponding to the index of the connected component of i.

Show that if instead of $\{u_1, \ldots, u_k\}$ we have any other orthonormal basis of the nullspace, $\{\hat{u}_1, \ldots, \hat{u}_k\}$, then if we form the matrix \hat{U} in the same way, and let \hat{y}_i be rows of \hat{U} , then again, $x_i = y_i / ||y_i||$ will be one of k distinct, orthonormal vectors.

6. Recall the example we did in class via direct calculation:

$$A = \left(\begin{array}{cc} 0 & 0 \\ 0 & \epsilon \end{array} \right), \qquad \Delta = \left(\begin{array}{cc} 0 & \beta \\ \beta & 0 \end{array} \right).$$

Let E_0 denote the smallest eigenvalue of A and F_0 the smallest eigenvalue of $(A + \Delta)$. Use the sin-theta theorem to bound $d_p(E_0, F_0)$.