Joy Ming and Alisa Nguyen (15 March 2013)

1 Problem 1

a. The probability that all of the M dimensions of x-y are between $-\epsilon$ and ϵ is $\rho = (2\epsilon)^M$. For each dimension i of χ , the probability that $|x_i - y_i| \le \epsilon$ is equivalent to

$$P(|x_i - y_i| \le \epsilon) =$$

$$P(-\epsilon \le x_i - y_i \le \epsilon) =$$

$$P(-\epsilon - x_i \le -y_i \le \epsilon - x_i) =$$

$$P(\epsilon + x_i \ge y_i \ge x_i - \epsilon) =$$

$$P(x_i - \epsilon \le y_i \le \epsilon + x_i) =$$

This distribution function is equivalent to $\int_{\epsilon+x_i}^{x_i-\epsilon} f(x)dx$, where f(x) is the PDF of y_i , which we know to have a uniform distribution, so $f(x) = \frac{1}{b-a} = 1$. Thus, we get:

$$\int_{\epsilon+x_i}^{x_i-\epsilon} 1 dx = \epsilon + x_i - (x_i - \epsilon) = 2\epsilon$$

Because we want to know the probability that all of the M dimensions of x-y are between $-\epsilon$ and ϵ , we simply take $\prod_{i=1}^{M} P(|x_i - y_i| \le \epsilon) = (2\epsilon)^M$.

- b. The probability of $\max_{m}|x_m y_m| \le \epsilon$ is at most ρ because as shown in (a), ρ does not depend on x_i and thus holds for all x_i . In addition, logically, if x is the center point, the average distance from it to any other point y is at most $\frac{1}{2}$ for any one dimension. As x moves farther and farther away from the center, the average distance increases so that it becomes at most 1 in any one dimension. So, if x is not in the center $\max_{m}|x_m y_m|$ grows and is less likely to be less than ϵ , decreasing that probability so that it is less than ρ .
- c. We will show that $||x-y|| \ge max_m|x_m-y_m|$.

$$||x - y|| = \sqrt{\sum_{m=1}^{M} (x_m - y_m)^2}$$

$$\sqrt{\sum_{m=1}^{M} (x_m - y_m)^2} \ge \max_{m} |x_m - y_m|$$

$$\sum_{m=1}^{M} (x_m - y_m)^2 \ge (\max_{m} |x_m - y_m|)^2$$

This is true because the left side of the inequality includes the right side in its sum. ||x - y|| is the total Euclidean distance between two points whereas $max_m|x_m - y_m|$ is only the distance between one dimension of two points. The left side must be larger.

If x is any point in χ , and y is a point in χ drawn randomly from a uniform distribution on χ , then the probability that $||x-y|| \le \epsilon$ is also at most p because ||x-y|| is greater than or equal to $max_m|x_m-y_m|$, making it less likely to be less than ϵ and thus giving it a probability lower than ρ of being less than ϵ .

d. Lowerbound on number N of points needed to guarantee that the nearest neighbor of point x will be within a radius ϵ of it is $log \delta / log (1 - (2\epsilon)^M)$.

For the nearest neighbor not to be within a radius ϵ , none of the neighbors can be within a radius ϵ .

The probability that any one neighbor is not within a radius ϵ of x is $1-(2\epsilon)^M$, so the probability that all the nighbors are not within a radius ϵ of x is equivalent to $(1-(2\epsilon)^M)^N$, where N is the number of neighbors. So, the probability that at least one neighbor is within a radius ϵ is 1 - that quantity. Since we want to guarantee with probability at least $1-\delta$ that the nearest neighbor will be within a radius ϵ of it, we can solve for a lower bound for N by setting the two equations equal to each other.

$$1 - \delta = 1 - (1 - (2\epsilon)^M)^N$$
$$1 - 1 + (1 - (2\epsilon)^M)^N = \delta$$
$$(1 - (2\epsilon)^M)^N = \delta$$
$$Nlog(1 - (2\epsilon)^M) = log\delta$$
$$N = log\delta/log(1 - (2\epsilon)^M)$$

e. We can conclude that the effectiveness of the hierarchical agglomerative clustering algorithm in high dimensional spaces is ineffective as the dimension M grows because N would also grow too large and HAC would require too many N points to actually be effective. As M increases, the denominator of the lower bound for N decreases, thus leading to an increase in N overall. In addition, as covered in class, when the size of the dataset gets larger, the probability that two points from different clusters are closer to each other in terms of distance than two points from separate clusters converges to 1/2.

$\mathbf{2}$ Problem 2

a Given a prior distribution $Pr(\theta)$ and likelihood $Pr(D|\theta)$, the predictive distribution Pr(x|D) for a new datum.

(a) ML:
$$Pr(x|D) = Pr(x|\theta) = \arg\max_{\theta} (\ln(Pr(D|\theta)))$$

$$\begin{split} \text{(a)} \ \ \text{ML:} \ ⪻(x|D) = Pr(x|\theta) = \underset{\theta}{\arg\max}(\ln(Pr(D|\theta))) \\ \text{(b)} \ \ \text{MAP:} \ ⪻(x|D) = Pr(x|D) = Pr(x|\theta) = \underset{\theta}{\arg\max}(\ln(P(D|\theta)P(\theta))) \end{split}$$

(c) FB:
$$Pr(x|D) = \int \theta P(\theta|D) d\theta$$

- b MAP can be considered "more Bayesian" than ML because it takes into account the distribution of θ instead of assuming same weight or uniformity.
- c One advantage the MAP method enjoys over the ML method
- d The Beta distribution is the conjugate prior of the Bernoulli.
- e Under the ML approach

3 Problem 3

- a The K-means clustering objective is to minimize the sum of squared distances between prototype and data.
- b PCA relates to K-means

Problem 4 4