

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 $\boxed{\rho = (2\epsilon)^M}$.
For each dimension i of χ , the probability that $|x_i - y_i| \leq \epsilon$ is equivalent to

$$\begin{aligned} P(|x_i - y_i| \leq \epsilon) &= \\ P(-\epsilon \leq x_i - y_i \leq \epsilon) &= \\ P(-\epsilon - x_i \leq -y_i \leq \epsilon - x_i) &= \\ P(\epsilon + x_i \geq y_i \geq x_i - \epsilon) &= \\ P(x_i - \epsilon \leq y_i \leq \epsilon + x_i) &= \end{aligned}$$

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:

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

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

- b. The probability of $\max_m |x_m - y_m| \leq \epsilon$ 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 $\|x - y\| \geq \max_m |x_m - y_m|$.

$$\begin{aligned} \|x - y\| &= \sqrt{\sum_{m=1}^M (x_m - y_m)^2} \\ \sqrt{\sum_{m=1}^M (x_m - y_m)^2} &\geq \max_m |x_m - y_m| \\ \sum_{m=1}^M (x_m - y_m)^2 &\geq (\max_m |x_m - y_m|)^2 \end{aligned}$$

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 $\|x - y\| \leq \epsilon$ is also at most ρ 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 $\boxed{\log \delta / \log(1 - (2\epsilon)^M)}$.
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 neighbors 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.

$$\begin{aligned}
1 - \delta &= 1 - (1 - (2\epsilon)^M)^N \\
1 - 1 + (1 - (2\epsilon)^M)^N &= \delta \\
(1 - (2\epsilon)^M)^N &= \delta \\
N \log(1 - (2\epsilon)^M) &= \log \delta \\
\boxed{N = \log \delta / \log(1 - (2\epsilon)^M)}
\end{aligned}$$

- 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 $1/2$.

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

where $Dist()$ is a distribution applied to the θ we obtain with the given formula.

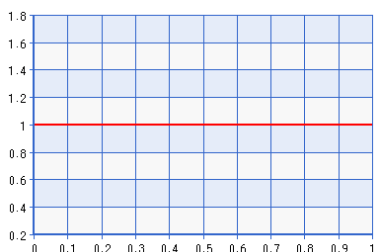
(b) MAP: $Pr(x|D) = Dist(\arg \max_{\theta} (\ln(P(D|\theta)P(\theta))))$,

where $Dist()$ is a distribution applied to the θ we obtain with the given formula.

(c) FB: $Pr(x|D) = \int p(x|\theta)P(\theta|D)d\theta$, where $P(\theta|D) = \frac{p(D|\theta)p(\theta)}{\int_{\theta} p(D|\theta')p(\theta')d\theta'}$

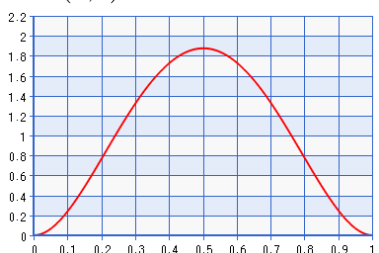
- b. MAP can be considered "more Bayesian" than ML because it takes into account the probability distribution of θ , the prior, instead of assuming same weight or uniformity like the ML method. The MAP method conditions on more information, specifically the prior distribution of theta and incorporates that information when finding the new probability.
- c. One advantage the MAP method enjoys over the ML method is that it accounts for the more likely distribution, as opposed to simply assuming uniformity, as with ML. In addition, MAP tends to be better than ML in terms of not overfitting the distribution to the data because while ML considers all the data uniformly, MAP considers the prior distribution of the parameters and does not assume equal weight. FB, on the other hand, maintains a probability distribution over the set of all possible parameter values. However, because the normalizing factor contains an integration over all parameter values, it can be difficult to compute. This means that it also unnecessarily takes into account the less likely, meaning that FB is less practical than MAP to calculate. Therefore, MAP sits in the "sweet spot" of taking more into account in terms of distribution than ML but less excessively than FB.
- d. Soccer team example based on different Beta distributions as priors for the probability of a win, the different possible distributions are:

- Beta(1,1)



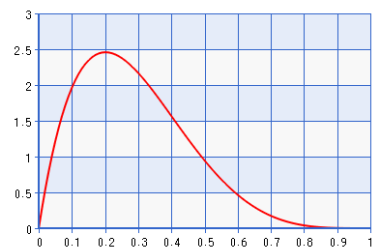
The Beta(1,1) distribution is equated with the Uniform(0,1) distribution, which is a number uniformly chosen between 0 and 1. Beta(1,1) is a uniform prior on θ which says there was just one positive and one negative example, essentially makes the probability of winning and losing equally likely. Often used as a "non-informative prior", this assumes uniformity and is not very descriptive in depicting the distribution of wins and losses otherwise.

- Beta(3,3)



Based on the equation for the mode of a variate with a beta distribution, $x = \frac{\alpha-1}{\alpha+\beta-2} = \frac{3-1}{3+3-2} = \frac{2}{4}$. This means that this distribution says that it is as likely to win as it is to lose as the distribution is symmetrical. However, this is more descriptive than simply a Beta(1,1) distribution, as seen in the plot because it looks more at the distribution over the entire space, with a lower variance as it makes values in the middle more likely.

- Beta(2,5)



In this case, the plot seems to be skewed toward the left, which the equation for the mode of a variate with a beta distribution show $x = \frac{\alpha-1}{\alpha+\beta-2} = \frac{2-1}{2+5-2} = \frac{1}{5}$ that this is true. This is more descriptive than Beta(1,1) in describing a distribution that is not necessarily uniform and it is also more descriptive than Beta(3,3) in describing a distribution that is not necessarily as symmetrical.

- The Beta distribution is the conjugate prior of the Bernoulli, so it is useful when used together with MAP estimation with Bernoulli variables since the posterior and prior will then have the same form. We can use this conjugate relationship to more easily use MAP estimation.
- Looking at the Harvard football team in the 2011-2012 season, we find that our team had 9 wins and 1 loss. In the 2012-2013 season, we won the first 3 games. Using this information, we can predict whether or not Harvard will win the next (fourth) game in the 2012-2013 season with the following approaches:

- Using the ML approach, we basically take the naive probability looking at the number of wins over the total number of games, so the probability they would win is $\frac{3}{3} = 1$. So, we predict that they will win next after the first three games of the 2012-2013 season. The ML approach does not take into account the prior distribution, so their record during the 2011-2012 season is not incorporated into the final prediction.
- Using the MAP approach, we take into account the prior distribution, assuming a Beta distribution of $\text{Beta}(\alpha, \beta)$,

$$\theta_{MAP} = \frac{n_w + \alpha - 1}{n + \alpha + \beta - 2}$$

Using the prior information of 9 wins and 1 loss during the 2011-2012 season, we get that the parameters of the Beta distribution are $\alpha = 9, \beta = 1$. And from the data from the first three games of 3 wins, we get $n_w = 3, n = 3$. So, plugging these into our equation we get

$$\theta_{MAP} = \frac{n_w + \alpha - 1}{n + \alpha + \beta - 2} = \frac{3 + 9 - 1}{3 + 9 + 1 - 2} = \frac{11}{11} = 1$$

. Thus, we predict that they will win the next game.

- Using the FB approach, we integrate over the entire distribution, using

$$p(x = 1|D) = \int_{\theta} p(x, \theta|D) d\theta = \int_{\theta} p(x|\theta) p(\theta|D) d\theta$$

Like in the MAP approach, we will use the prior distribution $p(\theta) = \text{Beta}(9, 1)$ to get $p(\theta|D) = \text{Beta}(12, 1)$ by incorporating the data of the first three games of the 2012-2013 season as well. So,

$$p(x = 1|D) = \int_0^1 \theta p(\theta|D) d\theta = \frac{n_w + \alpha}{n + \alpha + \text{beta}} = \frac{3 + 9}{3 + 9 + 1} = \frac{12}{13}$$

Because this probability is very close to 1, we will predict that our team will win the next game of the 2012-2013 season.

3 Problem 3

- The K -means clustering objective is to minimize the sum of squared distances between prototype and data. The loss function we seek to minimize is

$$L(\{\mu_k\}_{k=1}^K, \{r_n\}_{n=1}^N) = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

The update steps can be derived by performing gradient descent on it by taking the partial derivative of the loss function over all means and responsibilities. We minimize error with respect to $\{r_{nk}\}$ by assigning each example to its closest prototype at every iteration. We minimize error with respect to $\{\mu_k\}$ by assigning each prototype mean μ_k to the mean of the examples with which it is associated.

For each of the examples, we minimize the loss function by taking the partial derivative with respect to the responsibilities and setting it equal to 0:

$$\frac{\partial L}{\partial r_{nk}} = \|x_n - \mu_k\|^2 = 0$$

So, our update rule for responsibilities is $r_{nk} = 1$ for $k = \text{argmin}_{k'} \|x_n - \mu_{k'}\|^2$, and $r_{nk} = 0$ otherwise.

For each of the prototypes, we minimize the loss function by taking the partial derivative with respect to the means of each cluster:

$$\frac{\partial L}{\partial \mu_k} = \sum_{n=1}^N -2r_{nk}|x_n - \mu_k| = 0$$

$$\sum_{n=1}^N r_{nk}\mu_k = \sum_{n=1}^N r_{nk}x_n$$

Using this equation, we get $\mu_k = \frac{\sum_{n=1}^N r_{nk}x_n}{\sum_{n=1}^N r_{nk}}$, the update rule for each of the clusters.

Each iteration of the K-means algorithm updates the prototypes to decrease the error, following gradient descent until the prototypes converge to a local minimum and the loop ends. Because each iteration does decrease the error, K-means is guaranteed to converge in the direction that decreases the error the most, or in the direction of steepest descent.

- b. The principal components analysis (PCA), which looks to compress features using a linear reduction into lower dimension space, has an interpretation that looks to minimize the difference between data points through minimum square loss. One way they are related is that principal components are the continuous solutions to the discrete cluster membership indicators for K -means clustering. As PCA can be represented by a covariance matrix, K -means can be represented by a spectral expansion of the data covariance matrix truncated at $k-1$ terms. The relaxed solution of k -means clustering is given by the PCA principal components, where the PCA subspace spanned by the principal directions is identical to the cluster centroid subspace.

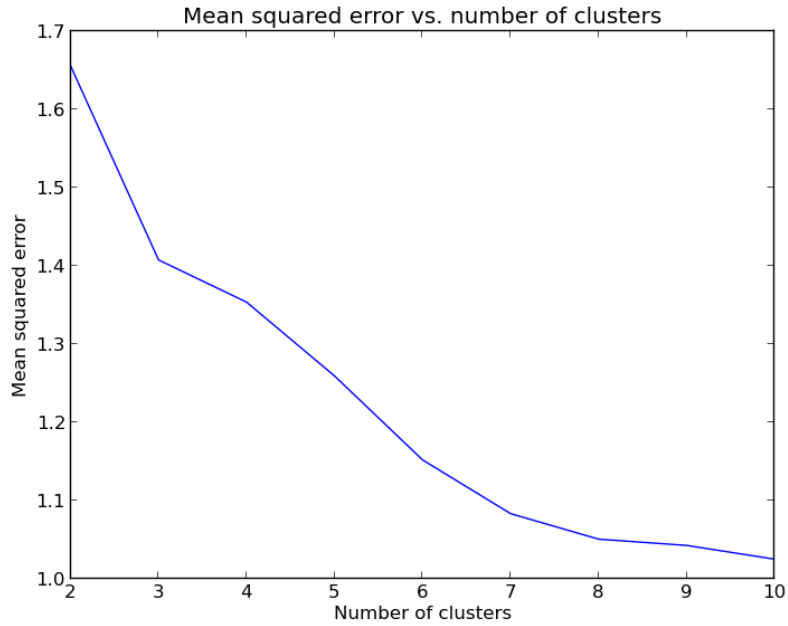
Although there exists a relationship between PCA and K-means clustering, the two algorithms are often used for vastly different situations. PCA is generally used to identify patterns in data and to compress (reduce) the number of dimensions in the data while minimizing the amount of information lost. K-means clustering, on the other hand, is used to identify clusters, or categories of datapoints by minimizing the inter-example distances. Some example situations and hypothetical data sets include:

- PCA is appropriate when dealing with a large dataset that is more easily dealt with when dimensions are decreased. For example, PCA can be used to compress images to minimize the information lost when discarding a dimension of the pixel data. K-means is not appropriate in this case because the goal is to minimize the space of the image file, and while K-means might be able to identify groups of pixels, it is unable to decide which parameter to throw away to compress the file.
- K -means is appropriate when classifying data points because it finds the centroids of clusters that minimize the distance between examples of the same cluster. For example, given multiple characteristics of plants, such as color, shape, etc. as parameters, K-means clustering can be used to classify the species of each plant. PCA is not appropriate in this case because although it can identify patterns in the data, it minimizes the loss by reducing dimensionality, thereby removing necessary data needed to accurately classify.

4 Problem 4

- a. K -means clustering algorithm

- (a) Plot of mean squared error versus K



- (b) If we were to choose the best K for this data based on the plot generated in part (i) we would choose the value of K that minimizes the mean squared error, which in this case is 10. It is important to minimize mean squared error because it is a measure of the difference between estimated values and true values, following in line with objectives mentioned in problem 3.

b. HAC algorithm

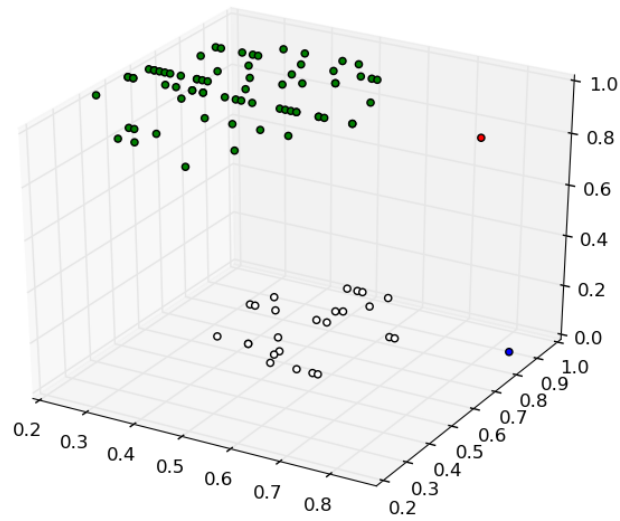
- (a) Comparing clusters formed using *min* distance metric against clusters formed using *max* distance metric

- i. Table showing number of instances in each cluster

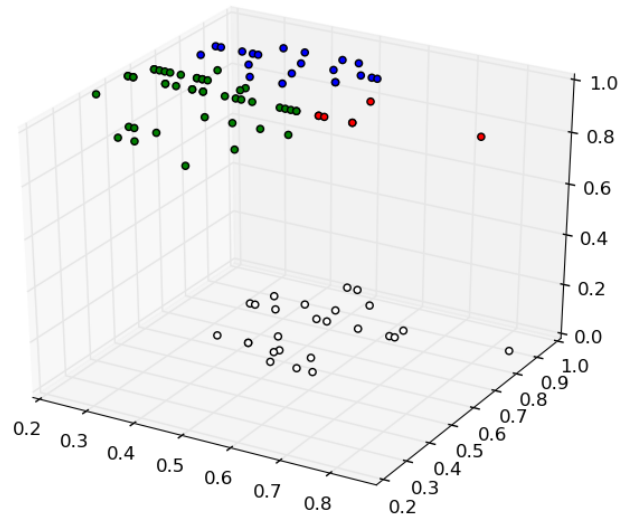
Metric	C1	C2	C3	C4
<i>min</i>	1	1	73	25
<i>max</i>	7	21	46	26

- ii. Scatterplot of the instances in 3-dimensions

- Based on *min*:



- Based on *max*:



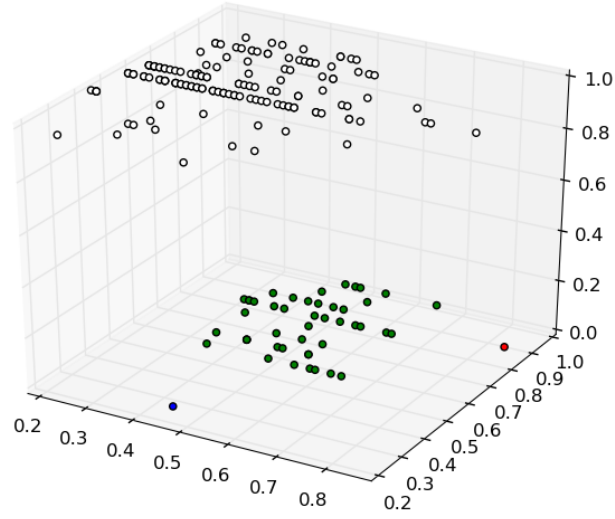
The *min* distance metric considers the distance between the two closest points in two separate clusters, but the *max* distance metric considers the distance between the two farthest points in the clusters. This results in the *min* measure creating more distinct and consolidated clusters than the *max* metric, which makes sense because by minimizing the smallest distance this brings the points closer to each other.

- Comparing clusters formed using *mean* distance metric against clusters formed using *centroid* distance metric
 - Table showing number of instances in each cluster

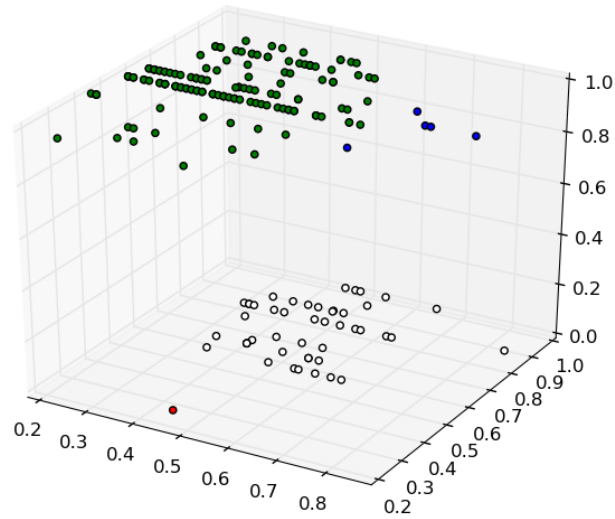
Metric	C1	C2	C3	C4
<i>mean</i>	1	1	46	152
<i>cent</i>	1	5	147	47

ii. Scatterplot of the instances in 3-dimensions

- Based on *mean*:



- Based on *cent*:



The *mean* distance metric looks at the average distances between each of the points in a cluster and *cent* looks at the distances to the center of a cluster. In the clusters based on *mean* the clusters are more distinct in a naive belief of what the clusters would look like whereas the *cent* is a little different and less distinct.

c. Autoclass clustering algorithm

- (a) It takes 10 iterations for the parameters to converge based on checking to see if the percent change in likelihood is greater than $1e-03$.
- (b) Plot of the log likelihood of the data versus number of iterations
- (c) The run-time performance of auto-class seems to be much faster than that of K -means. (???)