

# Artificial Intelligence - Problem Set 5

Mark Xavier (xaviem01)

May 13, 2019

## 1. Problem 1

- (a) For training, we walk through each example once, determine the classification to determine which centroid to compute, then convert the attributes to vectors and sum them. Each of these summed vectors is then multiplied by one over the number of items in that classification. Then based on input, training is simply  $\approx O(N)$  running time, not taking into account the time for division.

For testing, we take the new instance and find the distance between this instance and that of the centroids, then determine which centroid is closest. There are  $c$  centroids, and there are  $k$  predictive attributes, each an element in the vectors. Then, for  $c$  centroids we have to determine the distance between two  $k$  element arrays, leading to a running time of  $\approx O(ck)$ .

- (b) The approach in cases of missing values is to simply use the values that have been obtained and compare them to corresponding values in other vectors. If a new instance comes in  $z = \langle a_1, \emptyset, a_3 \rangle$ , then in all comparisons we simply grab the first and third element of each vector when we test distance. During the training phase, we can simply use the attributes we have. We could even remove attributes from the sets if they are missing from certain instances as long as we can avoid removing a large set of instances (for instance, if many instances do not include  $a_2$  then we can simply remove  $a_2$  from all instances) - this would lower some of the predictive power of our algorithm since we're removing information, but this would help in the case of missing attributes.
- (c) A probabilistic generative model associated with this process is kind of a mix between the model used for k-means clustering and naive-bayes. We can first choose a category, a, b, c, etc. Choosing this category leads us to different distributions to sample from, that is, category a has a certain mean and standard deviation, b has a different one, and c has a different one. Each of these means and deviations contributes to a gaussian distribution for each category from which to sample vectors from. Then the whole process is to "spin a roulette wheel" as noted in class, choose a category, then randomly sample a vector from the resulting distribution. We repeat this process to generate as many vectors as we want.
- (d) We can get unreasonable and/or unpredictable results from the following set-up:

Predictive Attribute	Classification
$\approx 1$	b
1	a
2	a
$\approx 2$	c

In the case above, we see that category a actually sits in-between categories b and c, and that the distance from a to b and from a to c is 1. Then when the centroid for a is calculated, it is calculated as 1.5, which is 0.5 away from centroid a. However, this is directly on the classification line for b and c, each of which are 0 away from either point in a (we can even make them 0.1 or such away). This means that the centroid distances for b and c will classify both points in a as being in either b or c. Therefore, though category a exists in training and is linearly separable, without enough points to define better centroids for b and c, we classify everything as b or c.

Really this works if the points are shifted slightly to (probably making this more obviously linearly separable). That is the meaning of the approx symbol ( $\approx$ ). In this case we still run into the same problem, classification for a relies on centroid 1.5, but for b and c the centroids are the points themselves, which are closer to the individual points of category a than the centroid calculated for a.

- (e) I take for given that the best place to orient a centroid is in the exact center of it's window (horizontally, across the x-axis). I can provide a minor explanation - according to the generative process defined above, we choose vectors by sampling from a distribution with a given  $\mu$  and  $\sigma$ . The given windows define for us what that mean and std. dev. look like, and in order to be as close as possible to all possible points that could be sampled, the centroid must be placed in the center of the window.

Now, with this given information, we note that if we place the centroids in a central location for each given window, firstly we cannot define centroids for categories a and e as their window of possible values is infinite. We

can assume however, that the centroid is as close as possible to the window that is closed (0 for category a, 6 for category e). We could also place them further out, it makes no difference to us. The real issue here is that the inner windows are not of the same size - category b has a window size of  $\approx 1$ , although technically it is infinite if this is a real number line, however category c has a window size of  $\approx 4$ . This means that, even with an ideal centroid placement, category b's centroid has a much shorter reach to category c's lower window than c's centroid, which will lead to points being incorrectly assigned to b instead of c. We may try to alleviate this by moving c's centroid, however the same issue occurs on the other side with category d, which can bleed into c's upper window. Since no matter where we move the given centroids we get issues labeling category c, we find that there is no set of centroids (one per category) that can correctly classify all possible vectors created from our generative process.

- (f) One method for doing this is to compute exemplars by computing some combinations of the training data. For example, every 2 or 3 vectors encountered in a given category  $c$  can be averaged to find a mean vector, but this may cause issues for vectors very far apart. Another method is to simply determine a halfway point between two vectors, and keep that as a centroid. The latter method holds information about both vectors, and as many as  $N - 1$  of these can be created in linear time (we can obtain many more by finding distances between all vectors, but this running time is  $\approx O(N^2)$ ). Either of these sets of centroids stores some kind of information about the vectors used to compute them, and therefore can be used in place of the centroid(s) calculated by the prompt.
- (g) The exemplars here are simple to define, they are all as close as possible to their respective boundaries (2 exemplars where boundaries are two-sided, only when if they are open). An example set is included in the table below:

X.A	X.C
-0.01	a
0.01	b
0.99	b
1.01	c
4.99	c
5.01	d
5.99	d
6.01	e

## 2. Problem 2

- (a) Naive-Bayes will classify vector  $\langle 2, 2, 2 \rangle$  as being in category b.
- (b) 3-Nearest-Neighbors identifies the vectors at indices 0, 1, and 7 as being the nearest to 2, 2, 2 (categories a, a, and c). This identifies the vector as being in category a.

## 3. Problem 3

- (a) The error function simply sums up the weight times the vector  $\langle x, y, 1 \rangle$  for all misclassified points, which turns out to be  $(1 - 1 - 3) + (2 - 2 - 3) + (3 - 4 - 3) = -10$ , so the cost function is 10.
- (b) Since correctly labeled points are multiplied by 0, we can ignore them. Looking at the 3 mislabeled points above, the first was labeled red but was blue, therefore we get  $\langle 1, 1, 1 \rangle$ . For the other two, they were red but labeled blue, so they get multiplied by -1, given  $\langle -2, -2, -1 \rangle$  and  $\langle -3, -4, -1 \rangle$ . Summing these gives:  $\langle -4, -5, -1 \rangle$ .
- (c) The new weight vector =  $\langle 1.4, 1.5, -2.9 \rangle$ .
- (d) The change is that now the updated weights get closer to categorizing the first point correctly to blue, however it doesn't work yet (the result is 0, which is just under the threshold). So in actual terms, classifications don't change at all.
- (e) The error function computes the same value as in the first part (10) as no vectors have been re-categorized after the weights were updated.