

Comp 6321 - Machine Learning - Assignment 3

Federico O'Reilly Regueiro

November 10th, 2016

Question 1: Midterm preparation question

Propose an adequate learning algorithm for each instance.

1.a 1000 samples, 6-dimensional continuous space, classify ~ 100 examples.

This could be a candidate for knn, despite the dimensionality approaching higher orders. We can still somewhat escape the curse of dimensionality, since for non-parametric methods, in order to have an best-case error rate e we require at least n samples where $n \sim \left(\frac{c}{e}\right)^{\frac{d+4}{4}}$, c is the number of classes, and d is the number of dimensions¹. Having 1000 samples in a 6-dimensional space would still allow for an error $e = 0.0632$ in the best case². An appropriate value for k would need to be found via cross-validation.

1.b Classifier for children in special-ed, justified to the board before it's implemented.

One of the easiest classification algorithms to explain in layman's terms is decision trees; since the method should be justified to the board, this would probably be an adequate choice. Furthermore, given that the stakes for such a classification are very high, an ensemble approach such as random forests / bagging could increase the classifier's performance by diminishing the tree's inherent variance and tendency to overfit.

1.c Binary classification, train with very large data-set of products / customer preferences. Input - 1 million bits - other clients' preferences. Frequent updates.

A recommender system of this nature could use naive bayes in a similar way to the document classification example presented in class. In this case, the features of each product are the clients who have shown interest in the product. The training would rely on the trends in clients' preferences accross all products. NB could work well given the size of the dataset and the need for frequent updates. However, a drawback to this approach is that the recommender assumes feature-independence while relying on the underlying relation between customer preferences, which in turn implies feature-dependency³. A similar problem arises in document classification, where the presence of words in a document cannot really be considered independent, yet NB still performs well for said task.

¹See the non-parametric methods pdf used in class.

²Said otherwise, a possible success rate of 93.68%

³Eg, If customer A and customer B have both shown interest in a large set of products, they are likely to show similar preferences for new products.

1.d 40 attributes, discrete and continuous, some have noise; only about 50 labeled observations.

With few examples and a fair amount of features, the curse of dimensionality haunts this classification. The presence of noise and the need for some sort of reduction in dimensionality might be well served by logistic regression with L1 regularization. K-fold cross-validation, with a relatively small k given the dataset size, would be necessary in order to find the appropriate rate for regularization.

Question 2: Properties of entropy

2.a Compute the following for (X, Y) :

$$p(0, 0) = 1/3, p(0, 1) = 1/3, p(1, 0) = 0, p(1, 1) = 1/3.$$

$$\text{i } H[x] = -\sum_x p(x) \log_2(p(x)) = -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) = .9182$$

$$\text{ii } H[y] = -\sum_y p(y) \log_2(p(y)) = -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) = .9182$$

$$\text{iii } H[y|x] = -\sum_x p(x) H[Y|X=x] = -\frac{2}{3} \left(\frac{1}{2} \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right) = \frac{2}{3}$$

$$\text{iv } H[x|y] = -\sum_y p(y) H[X|Y=y] = -\frac{2}{3} \left(\frac{1}{2} \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right) = \frac{2}{3}$$

$$\text{v } H[x, y] = -\sum_x \sum_y p(x, y) \log_2(p(x, y)) = 3 \left(-\frac{1}{3} \log_2\left(\frac{1}{3}\right) \right) = 1.5849$$

$$\text{vi } I[x, y] = \sum_x \sum_y p(x, y) \log_2\left(\frac{p(x, y)}{p(x)p(y)}\right) = H[x] - H[x|y] = 0.2516$$

2.b Prove maximum entropy in a discrete distribution happens in U

We wish to find:

$$\arg \max_{p_n} \sum_{n=1}^N p_n \log(p_n)$$

With constraints:

$$1 - \sum_{n=1}^N p_n = 0$$

$$p_i \geq 0, \forall i \in \{1, 2, \dots, N\}$$

We use Lagrange for maximization with constraints with a lagrangian multiplier only for the first constraint⁴:

$$\mathcal{L}(p_1, p_2, \dots, p_n, \lambda) = \sum_{n=1}^N p_n \log(p_n) - \lambda \left(1 - \sum_{n=1}^N p_n\right)$$

⁴The second series of constraints are satisfied by the solution using only $1 - \sum_n p_n = 0$.

And by setting the gradient of the Lagrangian function to 0, $\nabla_{p_1, p_2, \dots, p_N, \lambda} \mathcal{L}(p_1, p_2, \dots, p_N, \lambda) = 0$, we are left with a system of equations:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial p_1} \sum_{n=1}^N p_n \log(p_n) - \lambda(1 - \sum_{n=1}^N p_n) &= 0 \\ \frac{\partial \mathcal{L}}{\partial p_2} \sum_{n=1}^N p_n \log(p_n) - \lambda(1 - \sum_{n=1}^N p_n) &= 0 \\ &\vdots \\ \frac{\partial \mathcal{L}}{\partial p_N} \sum_{n=1}^N p_n \log(p_n) - \lambda(1 - \sum_{n=1}^N p_n) &= 0 \\ \frac{\partial \mathcal{L}}{\partial \lambda} \sum_{n=1}^N p_n \log(p_n) - \lambda(1 - \sum_{n=1}^N p_n) &= 0 \end{aligned}$$

Which in turn yields:

$$\begin{aligned} \log(p_1) + 1 - \lambda p_1 &= 0 \\ \log(p_2) + 1 - \lambda p_2 &= 0 \\ &\vdots \\ \log(p_N) + 1 - \lambda p_N &= 0 \\ 1 - \sum_{n=1}^N p_n &= 0 \end{aligned} \tag{1}$$

From which:

$$\lambda = \frac{\log(p_1) + 1}{p_1} = \frac{\log(p_2) + 1}{p_2} = \dots = \frac{\log(p_N) + 1}{p_N} \tag{2}$$

it is clear from equations 1 and 2 that $p_1 = p_2 = \dots p_N = \frac{1}{N}$, which is precisely a discrete uniform distribution.

2.c Show that T_1 wins

The notes show two possible tests for a decision tree. T_1 , where the left child has $[20+, 10-]$ possible outcomes in its sub-trees and the right node has $[10+, 0-]$. T_2 , on the other hand, yields: *left* = $[15+, 7-]$; *right* = $[15+, 3-]$.

The best choice should yield the maximum mutual information or information gain $I[p, T_n], n \in \{1, 2\}$. So for T_1 :

$$\begin{aligned} H[p] &= -\frac{1}{4} \log_2 \left(\frac{1}{4} \right) - \frac{3}{4} \log_2 \left(\frac{3}{4} \right) = 0.8112 \\ H[p|T_1 = t] &= -\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right) = 0.9182 \\ H[p|T_1 = f] &= 0 \\ H[p|T_1] &= p(T_1 = t)H[p|T_1 = t] + p(T_1 = f)H[p|T_1 = f] \\ &= 0.6887 \\ I[p, T_1] &= H[p] - H[p|T_1] = 0.1225 \end{aligned}$$

Whereas for T_2 we have:

$$\begin{aligned}
H[p|T_2 = t] &= -\frac{15}{22}\log_2\left(\frac{15}{22}\right) - \frac{7}{22}\log_2\left(\frac{7}{22}\right) = 0.9024 \\
H[p|T_2 = f] &= -\frac{15}{18}\log_2\left(\frac{15}{18}\right) - \frac{3}{18}\log_2\left(\frac{3}{18}\right) = 0.65002 \\
H[p|T_2] &= p(T_2 = t)H[p|T_2 = t] + p(T_2 = f)H[p|T_2 = f] \\
&= \frac{22}{40}0.9024 + \frac{18}{40}0.65002 = 0.7888 \\
I[p, T_2] &= H[p] - H[p|T_2] = 0.02245
\end{aligned}$$

From which we can see that we gain much more information from knowing the result of T_1 than by knowing the result of T_2 .

Question 3: Kernels

Suppose $k_1(\mathbf{x}, \mathbf{z})$ and $k_2(\mathbf{x}, \mathbf{z})$ are valid kernels over $\mathbb{R}^n \times \mathbb{R}^n$. Prove or disprove that the following are valid kernels.

Use Mercer's theorem regarding the Gram matrix⁵ or the fact that a kernel can be expressed as $k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z})$.

preliminaries

From Mercer, we know for each $k_1(\mathbf{x}, \mathbf{z})$ and $k_2(\mathbf{x}, \mathbf{z})$ we have corresponding kernel matrices \mathbf{M}_1 and \mathbf{M}_2 which are symmetric and positive semi-definite.

For both \mathbf{M}_1 and \mathbf{M}_2 :

Symmetry:

$$\mathbf{M}_i = \mathbf{M}_i^T \quad (3)$$

Positive semidefiniteness:

$$\mathbf{x}^T \mathbf{M}_i \mathbf{x} \geq 0 \quad (4)$$

$$|\mathbf{M}_i| \geq 0 \quad (5)$$

3.a $k(\mathbf{x}, \mathbf{z}) = ak_1(\mathbf{x}, \mathbf{z}) + bk_2(\mathbf{x}, \mathbf{z})$, $a, b > 0$; $a, b \in \mathbb{R}$

Firstly, we establish that if $k(\mathbf{x}, \mathbf{z})$ is a valid kernel, then $ak(\mathbf{x}, \mathbf{z})$ is also a valid kernel $\forall a > 0$; $a \in \mathbb{R}$:

We know that for a square matrix \mathbf{A} of size $n \times n$, $|a\mathbf{A}| = a^n |\mathbf{A}|$. And, since $a \geq 0$, we know that $a^n \geq 0$. Thus equation 5 holds for both of our summands. Additionally, since the scalar multiplication of a symmetric matrix yields another symmetric matrix, both summands are symmetric and therefore valid kernels.

Now, let us say:

$$ak_1(\mathbf{x}, \mathbf{z}) = k'_1(\mathbf{x}, \mathbf{z})$$

and

$$bk_2(\mathbf{x}, \mathbf{z}) = k'_2(\mathbf{x}, \mathbf{z})$$

are both valid kernels with kernel matrices \mathbf{M}'_1 and \mathbf{M}'_2 . The addition of two symmetric matrices yields a symmetric matrix, so we need to check for positive semi-definiteness.

Since both \mathbf{M}'_1 and \mathbf{M}'_2 are symmetric we can write:

⁵Equivalently known as the kernel matrix.

$$\begin{aligned} \mathbf{M}'_1 &= \mathbf{U}^T \mathbf{\Lambda}_U \mathbf{U} \\ \mathbf{M}'_2 &= \mathbf{V}^T \mathbf{\Lambda}_V \mathbf{V} \end{aligned}$$

and using equation 4:

$$\begin{aligned} (\mathbf{x}^T \mathbf{U}^T \mathbf{\Lambda}_U \mathbf{U} \mathbf{x} + \mathbf{x}^T \mathbf{V}^T \mathbf{\Lambda}_V \mathbf{V} \mathbf{x}) &\geq 0 \\ \mathbf{x}^T (\mathbf{U}^T \mathbf{\Lambda}_U \mathbf{U} + \mathbf{V}^T \mathbf{\Lambda}_V \mathbf{V}) \mathbf{x} &\geq 0 \\ \mathbf{x}^T (\mathbf{M}'_1 + \mathbf{M}'_2) \mathbf{x} &\geq 0 \end{aligned}$$

Which proves that $k(\mathbf{x}, \mathbf{z}) = ak_1(\mathbf{x}, \mathbf{z}) + bk_2(\mathbf{x}, \mathbf{z})$, $a, b > 0$; $a, b \in \mathbb{R}$ is a valid kernel.

3.b $k(\mathbf{x}, \mathbf{z}) = ak_1(\mathbf{x}, \mathbf{z}) - bk_2(\mathbf{x}, \mathbf{z})$, $a, b > 0$; $a, b \in \mathbb{R}$

Suppose:

$$a = 1, b = 1, \mathbf{M}_1 = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \mathbf{M}_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix},$$

Both \mathbf{M}_1 and \mathbf{M}_2 symmetric, positive semi-definite matrices. Yet $\mathbf{M}' = a\mathbf{M}_1 - b\mathbf{M}_2$ would yield:

$$\mathbf{M}_1 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

The eigenvalues of which are $\lambda_1 = -1, \lambda_2 = 1$, making \mathbf{M}' a non-positive semi-definite matrix and thus $k(\mathbf{x}, \mathbf{z})$ is not a valid kernel.

3.c $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z})k_2(\mathbf{x}, \mathbf{z})$

The kernel matrix \mathbf{M}' of the product of two matrices $k_1(\mathbf{x}, \mathbf{z}), k_2(\mathbf{x}, \mathbf{z})$ is equivalent to the element-wise multiplication of the respective two kernel matrices $\mathbf{M}' = \mathbf{M}_1 \odot \mathbf{M}_2$. This is also known as the Hadamard product or the Schur product. The Schur product theorem states that said product of two positive semi-definite matrices is also positive semi-definite. It is very easy to show that symmetry is preserved for the Hadamard product of two symmetric matrices. Thus $k(\mathbf{x}, \mathbf{z}) = k_1(\mathbf{x}, \mathbf{z})k_2(\mathbf{x}, \mathbf{z})$ is a valid kernel.

3.d $k(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})f(\mathbf{z})$, where $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Here we rely on the fact that a kernel can be expressed as $k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z})$ where $\phi(\mathbf{x})$ maps \mathbf{x} onto an n-dimensional space.

If $n = 1$ and $\phi = f$, $f(\mathbf{x})f(\mathbf{z})$ constitutes a valid kernel since it can be expressed as $k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z})$.

3.e $k(\mathbf{x}, \mathbf{z}) = p(\mathbf{x})p(\mathbf{z})$, where p pdf.

The same rationale as question 3.d applies here, $k(\mathbf{x}, \mathbf{z}) = p(\mathbf{x})p(\mathbf{z})$ is a valid kernel.

Question 4: Nearest neighbour vs decision trees, do boundaries coincide?

Boundaries do not necessarily coincide for these two classification strategies. In fact, in typical usage, they would tend to be non-coincidental but in some rare or contrived cases⁶ the boundaries might equate.

⁶Eg A dataset consisting of two points or the usage of decision functions of an arbitrary number of features, etc.

Decision tree boundaries are typically composed of hyper-planes that are orthogonal to the features f_d chosen for each decision; boundaries pass through the midpoint between points neighboring on a projection along the axis of f_d . Thus each segment of a decision-tree boundary will generally have one out of n directions for an n -dimensional space.

Conversely, boundaries for nearest-neighbours form a Voronoi tessellation, where each boundary segment corresponds to a hyper-plane running orthogonal to the line between a given point and its nearest neighbors while passing through the midpoint of such a line (thus the ensemble of said hyperplanes has a wide gammut of directions within the space).

For an example, see figures 1a and 1b.

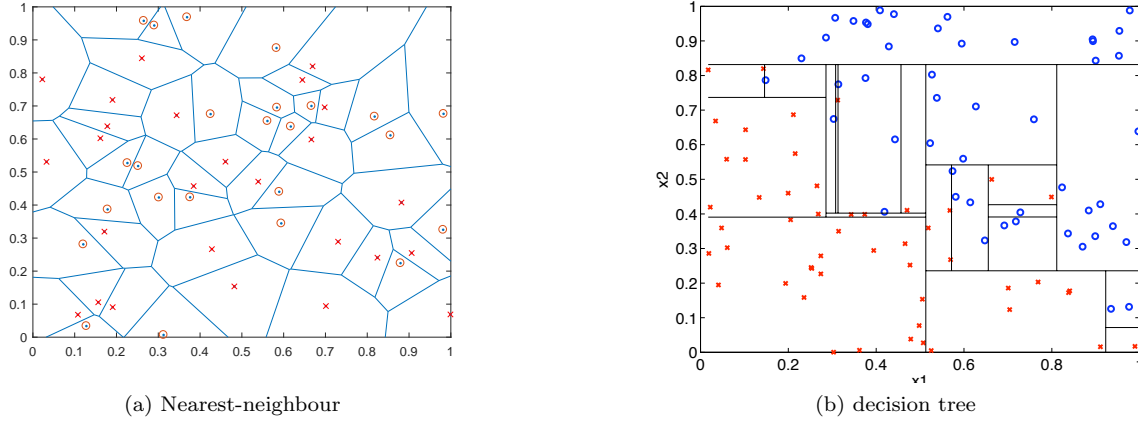


Figure 1: A Voronoi tessellation has boundary segments in many different directions, perpendicular to the lines between any two nearest-neighbors whereas decision-tree boundary segments are typically perpendicular to any one of a given set of features or feature combinations

Question 5: Bayes rate

For the following univariate case where $P(\omega_i) = \frac{1}{c}$ and

$$P(x|\omega_i) = \begin{cases} 1 & 0 \leq x \leq \frac{cr}{c-1} \\ 1 & i \leq x \leq i+1 - \frac{cr}{c-1} \\ 0 & otherwise \end{cases}$$

5.a Show that $P^* = r$

The minimal multi-class classification error rate P^* is given by:

$$P^* = 1 - \int \arg \max_i P(\omega_i|x) P(x) dx$$

We note that $P(\omega_i|x)P(x) = P(x, \omega_i) = P(x|\omega_i)P(\omega_i)$, thus:

$$P(\omega_i|x)P(x) = \begin{cases} \frac{1}{c} & 0 \leq x \leq \frac{cr}{c-1} \\ \frac{1}{c} & i \leq x \leq i+1 - \frac{cr}{c-1} \\ 0 & otherwise \end{cases}$$

Given the class density and probability, we can see that for any region with overlapping densities, the choice of any i will maximize. Additionally, we see that the constraints imposed by existing densities demand that $0 \leq r \leq \frac{c-1}{c}$. This in turn implies that densities overlap only in $[0, \frac{cr}{c-1}]$ thus:

$$\begin{aligned}
P^* &= 1 - \int P(\omega_1|x)P(x)dx \\
&= 1 - \frac{1}{c} \int_0^{\frac{cr}{c-1}} 1dx - \sum_{i=1}^c \frac{1}{c} \int_i^{i+1-\frac{cr}{c-1}} 1dx \\
&= 1 - \frac{1}{c} \frac{cr}{c-1} - 1 - \frac{cr}{c-1} \\
&= \frac{cr-r}{c-1} \\
&= r
\end{aligned}$$

5.b Show the nearest-neighbor rate $P = P^*$

From the piece-wise densities and the class prior, we can derive:

$$\begin{aligned}
P(\omega_i|x) &= \begin{cases} \frac{1}{c} & 0 \leq x \leq \frac{cr}{c-1} \\ 1 & i \leq x \leq i+1 - \frac{cr}{c-1} \\ 0 & otherwise \end{cases} \\
LNN &= \int \left[1 - \sum_{i=1}^c P^2(\omega_i|x) \right] p(x)dx \\
&= \int \left[1 - \sum_{i=1}^c \left(\frac{P(x|\omega_i)P(\omega_i)}{p(x)} \right)^2 \right] p(x)dx \\
&= \int p(x) - \sum_{i=1}^c \frac{P(x|\omega_i)^2 P(\omega_i)^2}{p(x)} dx \\
&= \int p(x) - \sum_{i=1}^c \frac{P(x|\omega_i)P(\omega_i)(P(\omega_i|x)p(x))}{p(x)} dx \\
&= \int p(x) - \sum_{i=1}^c P(x|\omega_i)P(\omega_i)p(\omega_i|x) dx \\
&= 1 - \frac{1}{c} \sum_{i=1}^c \int_0^{\frac{cr}{c-1}} \frac{1}{c} dx - \sum_{i=1}^c \frac{1}{c} \int_i^{i+1-\frac{cr}{c-1}} 1dx \\
&= 1 - \frac{1}{c} \frac{cr}{c-1} - 1 - \frac{cr}{c-1} \\
&= \frac{cr-r}{c-1} \\
&= r
\end{aligned}$$

Question 6: Implementation

In the interest of comparing methods, I have chosen to do both adaboost and knn. Both implementations rely on a single driver script, (please see attached file `A3_q6_driver.m`) which is printed at the end of this sub-section for the reader's convenience. In order to run the code created for this assignment, one must call `A3_q6_driver.m` from within Matlab while including the function files that were created for each method in the working directory or executable path. Adaboost relies on `ada_boost.m`, `stump.m` and `calculate_error.m`, while KNN simply relies on `knn.m`. The function files are included as attachments and shall be detailed and printed in sub-sections 6.a and 6.b for the reader's convenience.

The driver script deals with loading the data, creating the partitions for the k-fold CV, instantiating the maximum number of iterations or maximum number of neighbors for each one of the methods and plotting results. Each classifier, relies on function scripts written for the classifier.

`A3_q6_driver.m` (omitting header and final print instructions):

```
5 X = load('wpbcx.dat');
6 y = load('wpbcy.dat');
7
8 num_folds = 10;
9
10 folds_info = cvpartition(length(y), 'kfold', num_folds);
11 folds_idx = randperm(length(y));
12
13 iters = 50;
14 max_k = 50;
15
16 errs_ada = zeros(iters, num_folds, 2);
17 errs_knn = zeros(max_k, num_folds, 2);
18
19 for fold = 1:num_folds
20     disp(sprintf('Performing %d-fold CV, fold: %d', num_folds, fold));
21     idxs_prev = 1:sum(folds_info.TestSize(1:(fold-1)));
22     if ~isempty(idxs_prev)
23         offset = idxs_prev(end);
24     else
25         offset = 0;
26     end
27     idxs_xcl = (1:folds_info.TestSize(fold))+offset;
28     idx_after_skip = length(y)-(sum(folds_info.TestSize((fold+1):end))-1);
29     idxs_next = idx_after_skip:length(y);
30     X_train = X(folds_idx([idxs_prev, idxs_next]), :);
31     X_test = X(folds_idx(idxs_xcl), :);
32     y_train = y(folds_idx([idxs_prev, idxs_next]));
33     y_test = y(folds_idx(idxs_xcl));
34
35     h = zeros(iters, 3);
36     alphas = zeros(iters,1);
37
38     m = length(y_train);
39     W = ones(m, 1)/m;
40     %%% Adaboost loop %%%
41     for i = 1:iters
42         % keep the user informed on longer runs
43         if mod(i, 100) == 0
44             disp(sprintf('\tworking on iter %d', i));
```



```

45     end
46     [h(i,:), alphas(i), W] = ada_boost(X_train, y_train, W);
47     errs_ada(i, fold, 1) = calculate_error(X_train, y_train, ...
48                                         h(1:i, :), alphas(1:i));
49     errs_ada(i, fold, 2) = calculate_error(X_test, y_test, ...
50                                         h(1:i, :), alphas(1:i));
51 end
52 %%% KNN loop %%%
53 for l = 1:m
54     % remove point i from the training set
55     X_train_temp = X_train;
56     y_train_temp = y_train;
57     X_train_temp(l,:) = [];
58     y_train_temp(l) = [];
59     yi_hats = knn(X_train_temp, y_train_temp, X_train(l, :), max_k);
60     for k = 1:max_k
61         if yi_hats(k) ~= y_train(l)
62             errs_knn(k, fold, 1) = errs_knn(k, fold, 1) + 1/m;
63         end
64     end
65 end
66 n = length(y_test);
67 for l = 1:n
68     yi_hats = knn(X_train, y_train, X_test(l, :), max_k);
69     for k = 1:max_k
70         if yi_hats(k) ~= y_train(l)
71             errs_knn(k, fold, 2) = errs_knn(k, fold, 2) + 1/n;
72         end
73     end
74 end
75 end

```

6.a Adaboost

The implementation relies on the above-mentioned driver script (common to both adaboost and knn) and three function files `ada_boost.m`, `stump.m` and `calculate_error.m`.

The driver calls `ada_boost.m` for as many iterations as indicated in line 13, while retaining a copy of the weak classifiers h , the classifier weights α and the observation weights w after each iteration. Each time it runs, `ada_boost.m` calls `stump.m` with the appropriate observation-weights and the later returns a weak classifier h_i , for which `ada_boost.m` calculates a weight α_i . After each iteration i of training with the training set, the driver script uses the accrued $h_i, h_{i-1}, \dots, h_1, \alpha_i, \alpha_{i-1}, \dots, \alpha_1$ pairs to call `calculate_error.m` both with the training and testing data in order to store training and testing errors (see lines 46 to 50 in the driver script).

Note that most of the effort required to get adaboost working was put into, not so much `ada_boost.m` but into `stump.m`, which returns checks along all dimensions, between any two points of differing classes in order to minimize the error of the stump returned. In order to do so, the stump that yields the error farthest from 0.5 is chosen and polarity is decided accordingly: polarity is positive if the error is less than 0.5 and negative if the error is larger than 0.5⁷.

`ada_boost.m` (omitted header)

```

5 function [h, alpha, W] = ada_boost(X, y, W)

```

⁷That is to say, if we have an extremely large error, the stump is very informative, we just have to take the opposite values for classification.

```

6 %ADA_BOOST takes an observation matrix X (each observation as a row), labels
7 %y and a vector W of weights for each observation. It returns a series
8 %h [thr,dim,polarity] tuples per stump, their associated weights and a a
9 %two-column error matrix.
10 % for the error matrix, row is associated with each iteration and training
11 % and test errors populate the first and second columns respectively.
12
13 [ Threshold, Dim, polarity, err ] = stump(X, y, W);
14 h = [Threshold, Dim, polarity];
15
16 y_hat = X(:,h(2)) > h(1);
17 if h(3) == -1
18     y_hat = ~y_hat;
19 end
20
21 alpha = 0.5 * log((1-err)/err);
22 % a vector containing 1 for correctly classified entries given h(round)
23 % and -1 for misclassified entries.
24 classifiedRightOrNot = ((2*y)-1).*((2*y_hat)-1);
25 W = W .* exp(-alpha*(classifiedRightOrNot));
26 W = W./sum(W);
27 end

```

stump.m (omitted header)

```

5 function [ Threshold, Dim, polarity, err ] = stump( X, y, W )
6 %STUMP Returns the best decision stump as a threshold-dimension pair for
7 %observation contained in the X matrix given class label y and weights W
8 %vectors
9 % X should contain row entries for observations, where columns are
10 % features. y is a vector containing binary classifications and W are the
11 % weights associated with each observation in x (eg weights assigned by a
12 % boosting algorithm. The threshold is the point at which a decision
13 % boundary perpendicular to the selected dimension Dim should pass.
14 % Polarity indicates negative values are to the left/down (+) or right/up(+)
15 [m,d] = size(X);
16 % Normalize W, in case it hasn't been done
17 W = W/sum(W);
18 Obs = zeros(m,d);
19 IDXs = zeros(m,d);
20 % Order each feature, keep a matrix of row indices per feature and then
21 % keep row classes and errors.
22 [Obs, IDXs] = sort(X);
23 Ys = (y(IDXs)*2)-1;
24 E = NaN*(ones(m,d));
25 % we can only have a split between neighbours of different classes, don't
26 % bother checking between neighbours of the same class
27 diff_neighbors = [2*ones(1,d);0.5*abs(diff(Ys))];
28
29 for n = 1:d
30     for i = 1:m
31         if diff_neighbors(i,n) == 1
32             C_hat = X(:,n) >= Obs(i,n);
33             err = sum((C_hat ~= y).*W);
34             E(i,n) = err;
35         end
36     end
37 end

```

```

36         end
37     end
38
39     % get the error farthest from 0.5, the most informative split
40     [mx,idx] = max(abs(0.5-E));
41
42     % we need to know along which direction it happened
43     [~, Dim] = max(mx);
44
45     % recover the true index
46     idx = idx(Dim);
47
48     err = E(idx,Dim);
49
50     % polarity 1 means that positive class observations are contained at
51     % higher values for that feature, polarity -1 means positive observations
52     % live at lower values than the threshold
53     if err > 0.5
54         polarity = -1;
55         err = 1 - err;
56     else
57         polarity = 1;
58     end
59
60     % Now calculate where the split happens, the midpoint between two adjacent
61     % points on that feature dimension, repeat the first point so that decision
62     % at the lower bound is just the smallest value of the feature space
63     cut_point = Obs(idx,Dim);
64     feature_vals = unique(Obs(:,Dim));
65     the_real_idx = cut_point == feature_vals;
66     feature_vals = [(feature_vals(1) - feature_vals(2)); feature_vals];
67     midpoints = feature_vals + 0.5*(diff(feature_vals);0);
68     Threshold = midpoints(the_real_idx);
69 end

```

calculate_error.m (omitted header)

```

5 function [ err, y_hat ] = calculate_error( X, y, h, alphas )
6 %CALCULATE_ERROR calculates the error given a dataset X, its target labels
7 %y and the ensemble adaboost predictor h, alpha.
8 % X are the observations given in row form, y the labels given as a
9 % column vector, h the matrix of [threshold, dim, polarity] weak
10 % classifiers and alphas is the column vector of weights associated with
11 % each weak classifier tuple.
12
13 y_hat = zeros(length(y), 1);
14
15 m = length(y);
16 for i = 1:m
17     for k = 1:length(alphas)
18         threshold = h(k, 1);
19         dim = h(k, 2);
20         pol = h(k, 3);
21         contrib_k = (2*(pol*X(i, dim) >= pol*threshold)-1)*alphas(k);
22         y_hat(i) = y_hat(i) + contrib_k;
23     end

```

```

24     y_hat(i) = y_hat(i) >= 0;
25 end
26
27 err = sum((y_hat > 0) ~= y)/m;
28 end

```

Surprisingly, adaboost's did not improve significantly (it actually degraded) for the test set throughout iterations. Errors on the training set, however, decreased exponentially as expected. Performance considerations, however, are fully addressed in subsection 6.c.

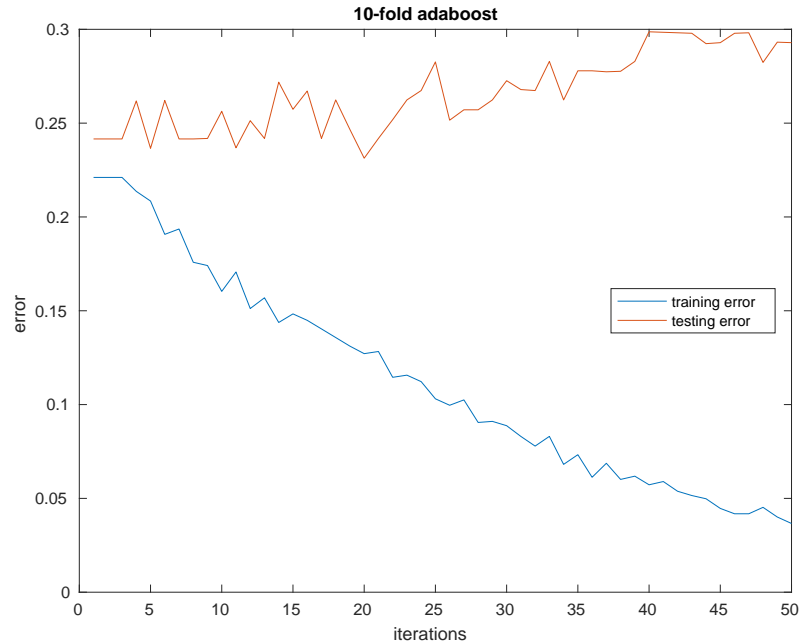


Figure 2: Training and testing errors for a 10-fold adaboost classification of the Wisconsin dataset used for Assignment 2

6.b KNN

KNN was significantly simpler to implement. The difficulty lying more in establishing what it means to have training and testing on a method that does not really *learn* the data. We chose to use the training set as neighbours to both the training and testing instances, making sure to remove the point we sought to classify from the training set during training instances. We refer the reader to the driver script, lines 54 through 73.

knn.m (omitted header)

```

5 function [yi_hats, IDXs] = knn(X, y, xi, max_k)
6 %KNN returns the class estimates for xi, given xi's 1:max_k nearest neighbours.
7 %It also returns the indexes of said neighbours while requiring a
8 %row-entry matrix of observations X and a corresponding column of labels y.
9     m = length(y);
10     dist = zeros(m);
11     for i = 1:m
12         dist(i) = norm(xi - X(i,:));

```

```

13     end
14     [~, idx] = sort(dist);
15     for k = 1:max_k
16         IDXs = idx(1:k);
17         yi_hats(k) = round(sum(y(IDXs))/k);
18     end
19 end

```

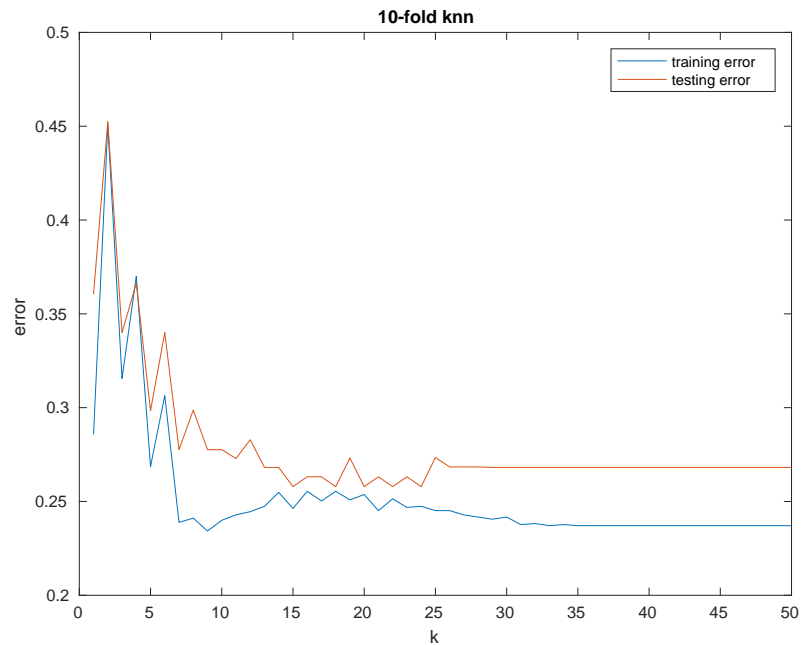


Figure 3: Training and testing errors for a 10-fold KNN of the Wisconsin dataset used for Assignment 2

6.c Results

Despite significant differences in performances for training data, both methods yielded comparable best rates for testing data.

Moreover, a comparison between both classifiers and the use of class priors as a means of classification yielded interesting results. from `A3.q6_driver.m` :

```

99 % how well did we do?
100 disp('Here''s how well we did:')
101 disp(sprintf('best prediction given on test data by adaboost:\n\t%d',...
102             min(mean(errs_ada(:,2), 2)) ));
103 disp(sprintf('best prediction given on test data by knn:\n\t%d',...
104             min(mean(errs_knn(:,2), 2)) ));
105 disp(sprintf('empirical ratio of class 1 to class 0:\n\t%d',...
106             sum(y)/length(y) ));

```

yields the following output.

```

best prediction given on test data by adaboost:
min(mean(errs_ada(:,:,2), 2))
    2.313158e-01
best prediction given on test data by knn:
min(mean(errs_knn(:,:,2), 2))
    2.578947e-01
empirical ratio of class 1 to class 0:
sum(y)/length(y)
    2.371134e-01

```

I refer the reader to the analysis I performed on the same data for assignment 2, where the conclusion was that significant class-overlap across most dimensions yielded classification only slightly more effective than using class priors.

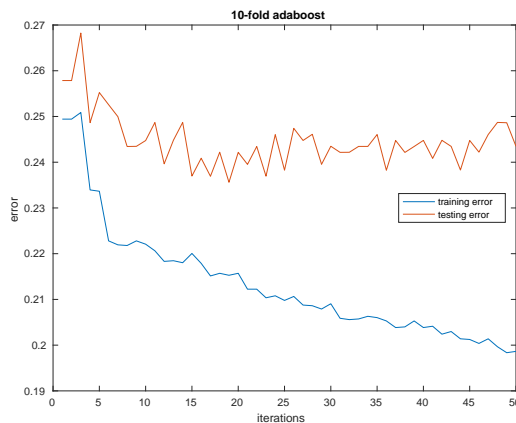
Curious to see the contribution of this particular dataset to classifier accuracy, I downloaded another dataset with less class overlap from UCI: (<https://archive.ics.uci.edu/ml/machine-learning-databases/pima-indians-diabetes/>) which is included as an attachment. For this dataset, the performance of the classifiers is quite different from what we observe in the previous example. See figures 4a and 4b.

We note that in order to load this dataset, the reader need only substitute lines 5 and 6 from the driver script with the following code:

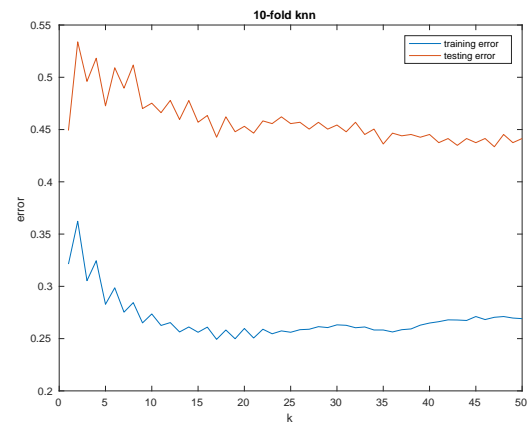
```

pima = load('pima-indians-diabetes.data');
X = pima(:,1:end-1);
y = pima(:,end);

```



(a) adaboost - pima dataset



(b) knn - pima dataset

Figure 4: Adaboost and KNN behave quite differently with the UCI pima-indians-diabetes dataset than they do with the wisconsin dataset.

And the output of the script shows that Adaboost performs significantly better with this dataset however KNN performs significantly worse.

Here's how well we did:

```

best prediction given on test data by adaboost:
    2.356118e-01

```

best prediction given on test data by knn:

4.335441e-01

empirical ratio of class 1 to class 0:

3.489583e-01

This last comparison merely serves as an informal observation regarding the suitability of different methods for different kinds of problems and different datasets.