

Optical Character Recognition of printed text

Munkhardt, S. & Reisle, C. R.

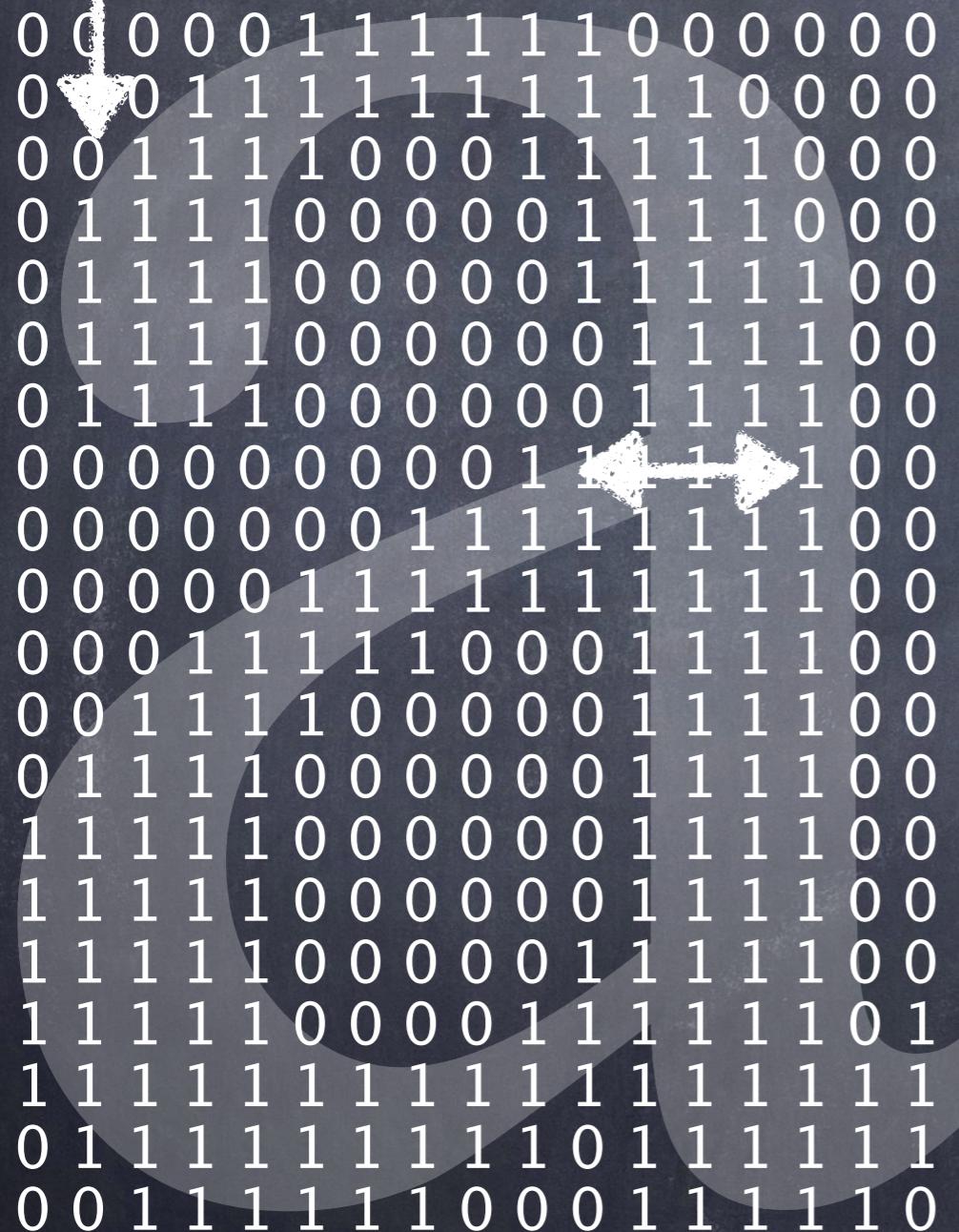
Outline

Project Goal: take jpg file input, extract
text, output to plain text file

- ① Pre-Processing
- ② Classification
- ③ Attributes
- ④ Model
- ⑤ Test data
- ⑥ Results
- ⑦ Next steps

Pre-processing

Simple Clustering



A 2D binary matrix representing a grayscale image. The matrix consists of 25 rows and 16 columns of binary values (0s and 1s). A cluster of 16 adjacent pixels is highlighted with a thick white border. Inside this cluster, the value '1' appears 12 times and '0' appears 4 times. An arrow points to the right from the first column of the cluster, and another arrow points left from the last column, indicating the extent of the cluster across the columns.

0	0	0	0	1	1	1	1	1	0	0	0	0	0	0	0
0	0	1	1	1	1	1	1	1	1	1	0	0	0	0	0
0	0	1	1	1	1	0	0	0	1	1	1	1	0	0	0
0	1	1	1	1	0	0	0	0	0	1	1	1	0	0	0
0	1	1	1	1	0	0	0	0	0	1	1	1	1	0	0
0	1	1	1	1	0	0	0	0	0	0	1	1	1	1	0
0	1	1	1	1	0	0	0	0	0	0	1	1	1	1	0
0	1	1	1	1	0	0	0	0	0	0	1	1	1	1	0
0	1	1	1	1	0	0	0	0	0	0	1	1	1	1	0
0	0	0	0	0	0	0	0	0	1	1	1	1	0	0	0
0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	0
0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	0
0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	0
0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	0
0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	0
0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1

- makes a pass over the columns of binary matrix, first labelling adjacent pixels with the same set number
- on the second pass over the matrix, the sets of adjacent pixels are merged
- on the third pass, the regions are extracted based on the pixels belonging to each set/component

Classification

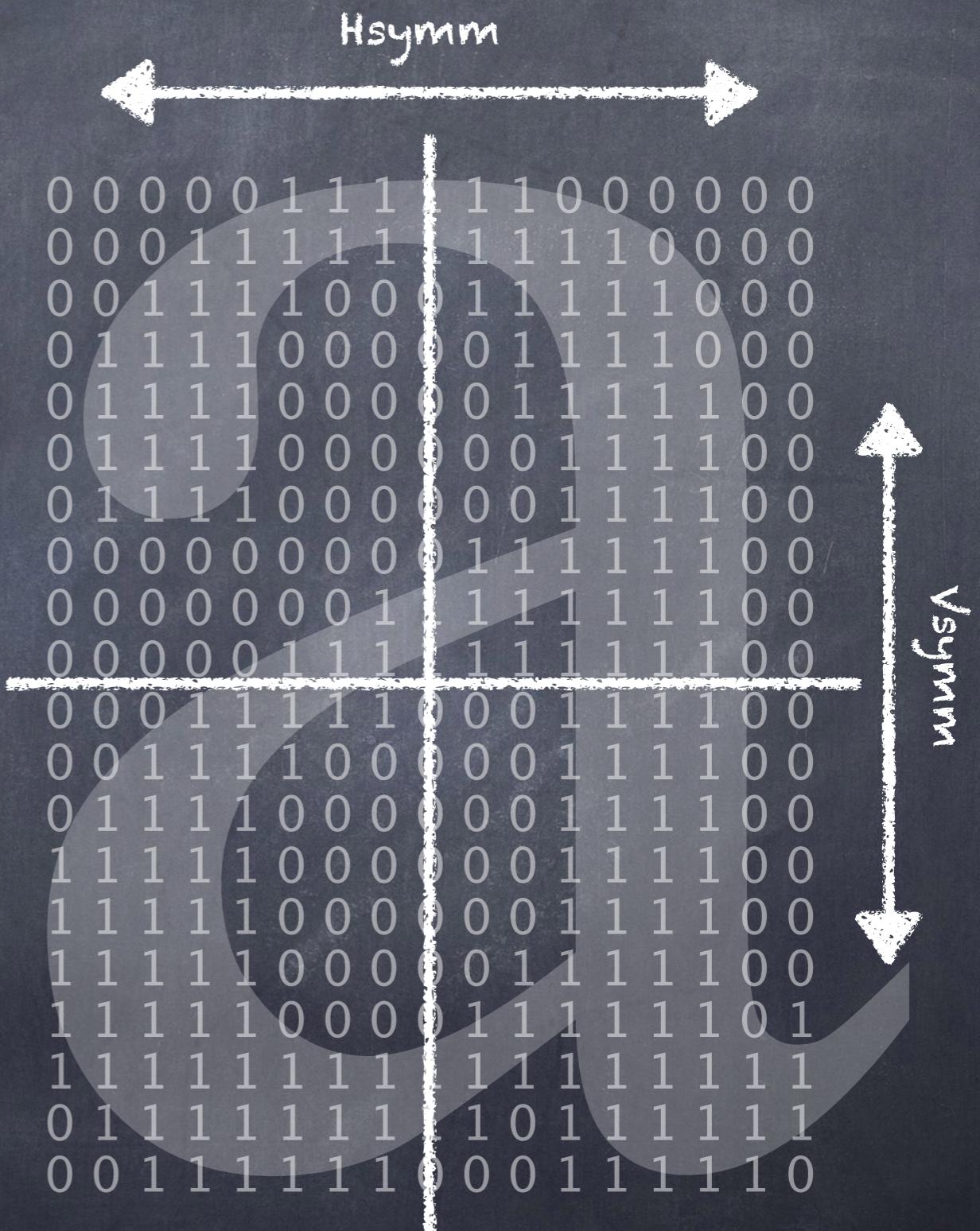
Computing Attribute Vectors

Feature	Function
$H\{10, 30, 50, 80, 90\}$	Sum of black pixels in the row/col of pixels at X% of the total height/width
$V\{10, 30, 50, 80, 90\}$	
H-symm	Measure of the horizontal symmetry
V-symm	Measure of the vertical symmetry
CC	number of closed components
H/W ratio	height to width ratio of the component
C	number of merged components
$Q\{1, 2, 3, 4\}$	ratio of black pixels that fall into the top left, top right, bottom left and bottom right quadrants, respectively, of the image
$Ih\{30, 50, 80\}$	number of times the pixels transition from black to white and white to black at the specified row or column position
$Iv\{30, 50, 80\}$	

Attributes...

00000011111000000
00011111111000000
00111100011110000
011110000001111000
011110000001111000
011110000000111100
011110000000111100
011110000000111100
00000000001111000
00000000011111000
00000011111110000
001111000001111000
011110000001111000
111110000001111000
111110000001111000
111110000001111000
111110000001111000
111110000001111000
111110000001111000
111110000001111000
111110000001111000
111111111111111111
011111111111111111
001111111111111110

C C



Attributes...

v10

V30

V50

V80

vgo

A 16x16 grid of binary digits (0s and 1s). The grid is divided into four quadrants by large, semi-transparent circles. The top-left quadrant is yellow and contains the label "Q1". The top-right quadrant is light blue and contains the label "Q2". The bottom-left quadrant is purple and contains the label "Q3". The bottom-right quadrant is dark blue and contains the label "Q4". The binary values in the grid are as follows:

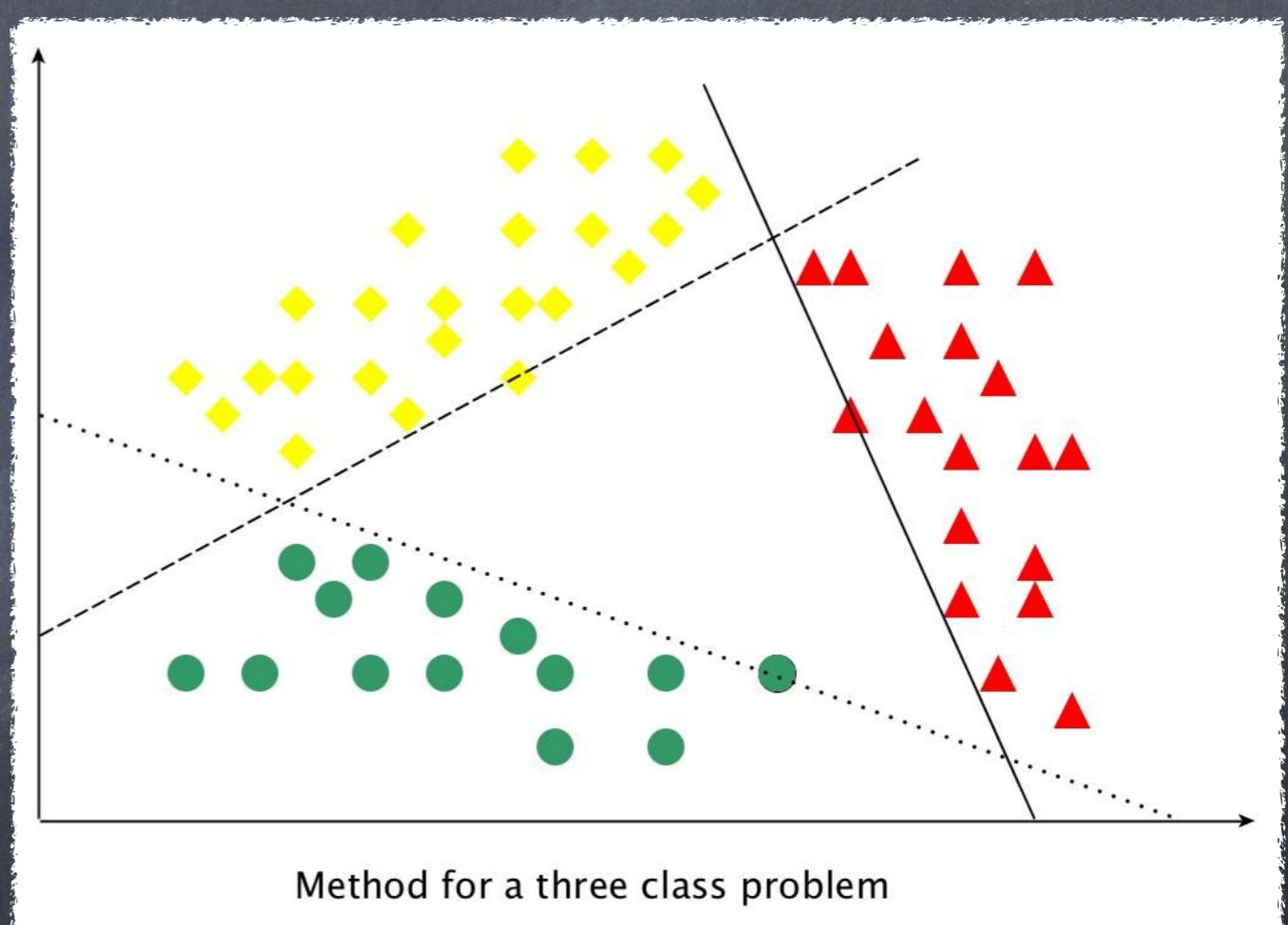
Row	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
1	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1
2	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1
3	0	0	1	1	1	1	1	1	0	0	0	0	0	0	0	0
4	0	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0
5	0	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0
6	0	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0
7	0	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1
10	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1
11	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1
12	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1
13	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1
14	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1
15	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1
16	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1

Classification

Mathematical Models

Support Vector Machine

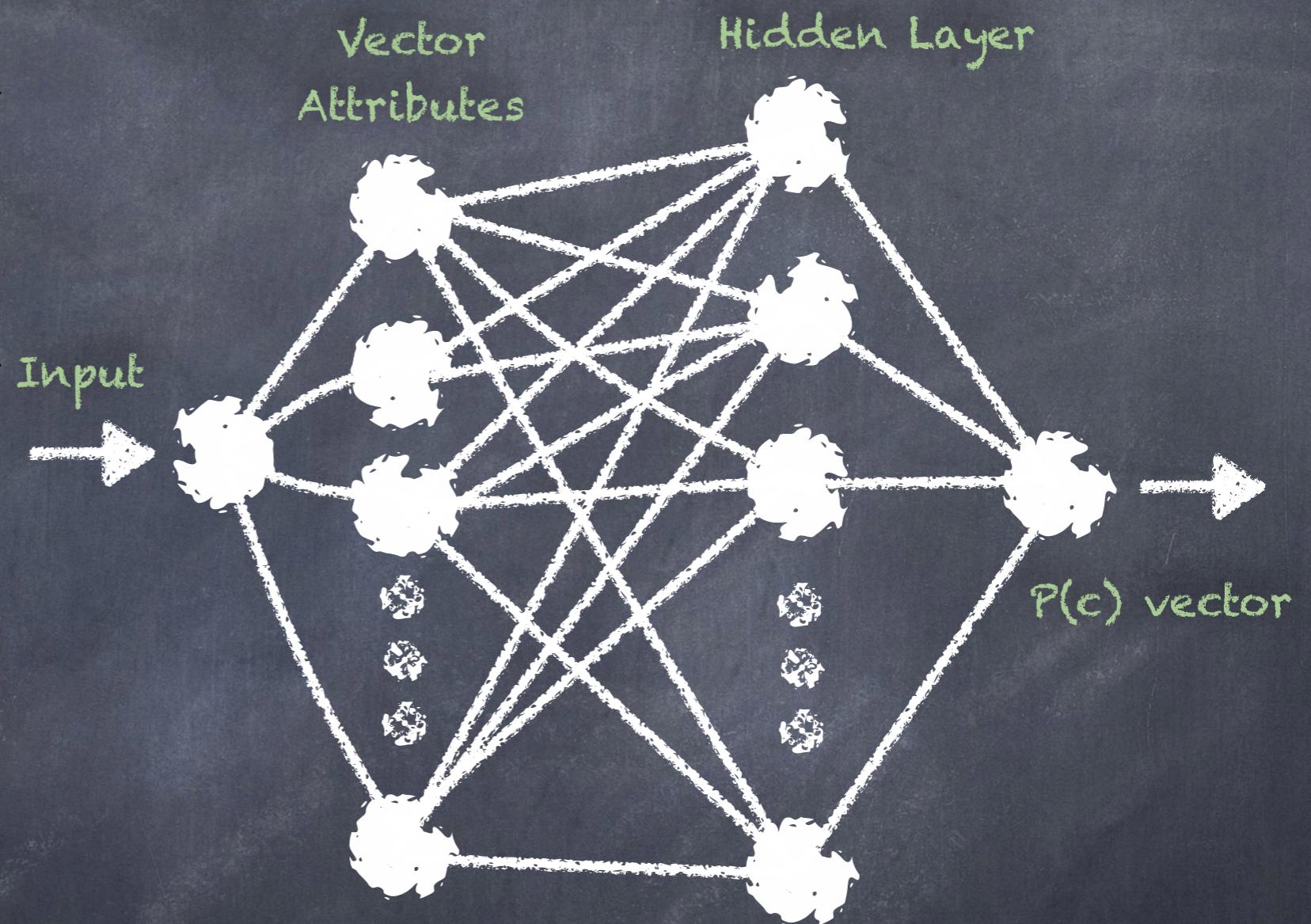
- We have multiple (82) classes of characters
- We need 82 SVM classifiers or hyperplanes to separate characters out from each other
- Example: Classifier for 0, call it SVMO, will separate all the samples of 0 from the other characters



Artificial Neural Network

- An ANN is an abstract copy of biological neural networks

- An ANN is a function that maps input vectors from our input space (a binary matrix) to an output space (class probability vector)



- The function has several layers consisting of neurons (called perceptrons)
- For each input our particular ANN outputs a vector with the probability of an output belonging to a given class

Classification

Model Evaluation

==== Evaluation on training set ====

	<u>ANN</u>	<u>SVM</u>
Correctly Classified Instances	411 (100 %)	410 (97.7567 %)
Incorrectly Classified Instances	0 (0 %)	1 (0.2433 %)
Kappa statistic	1	0.9975
Mean absolute error	0.0033	0.0021
Root mean squared error	0.0175	0.0137
Relative absolute error	13.7823 %	8.8617 %
Root relative squared error	15.9544 %	12.5025 %
Total Number of Instances	411	411

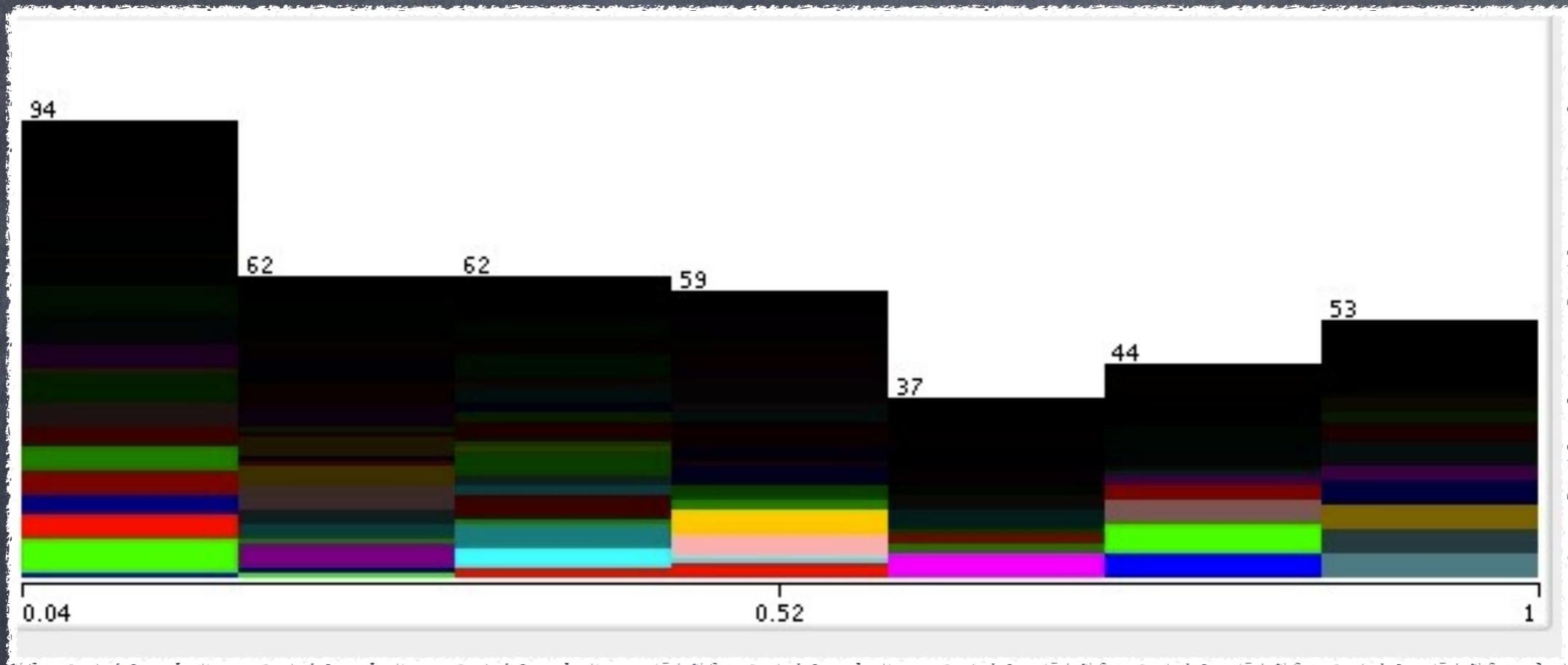
==== Stratified 10-fold cross-validation ====

	<u>ANN</u>	<u>SVM</u>
Correctly Classified Instances	399 (97.0803 %)	399 (97.0803 %)
Incorrectly Classified Instances	12 (2.9197 %)	12 (2.9197 %)
Kappa statistic	0.9704	0.9704
Mean absolute error	0.0031	0.0028
Root mean squared error	0.025	0.0275
Relative absolute error	12.7235 %	11.7249 %
Root relative squared error	22.7414 %	25.0182 %
Total Number of Instances	411	411

Classification

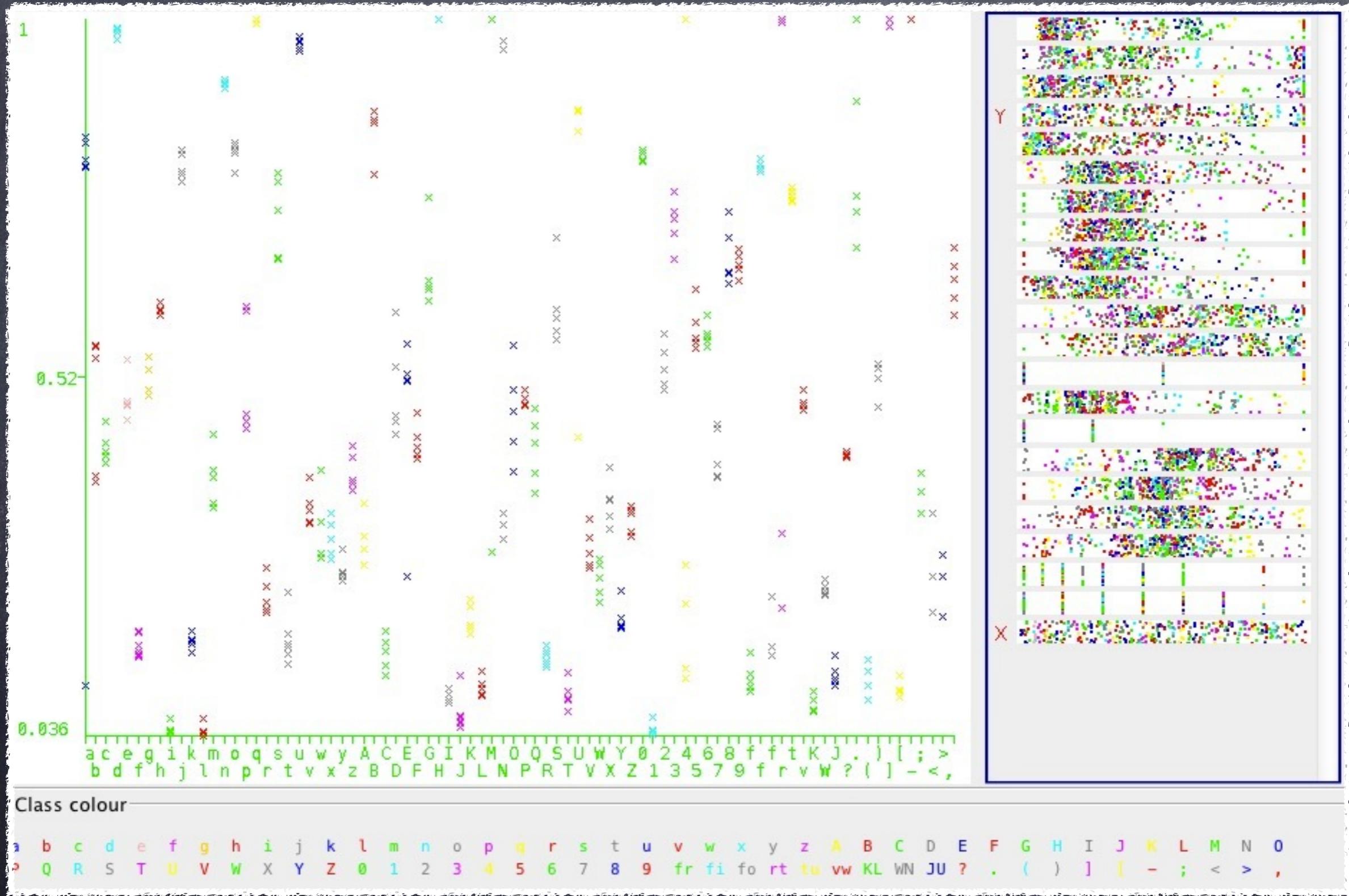
Attribute Evaluation

Attribute: V80

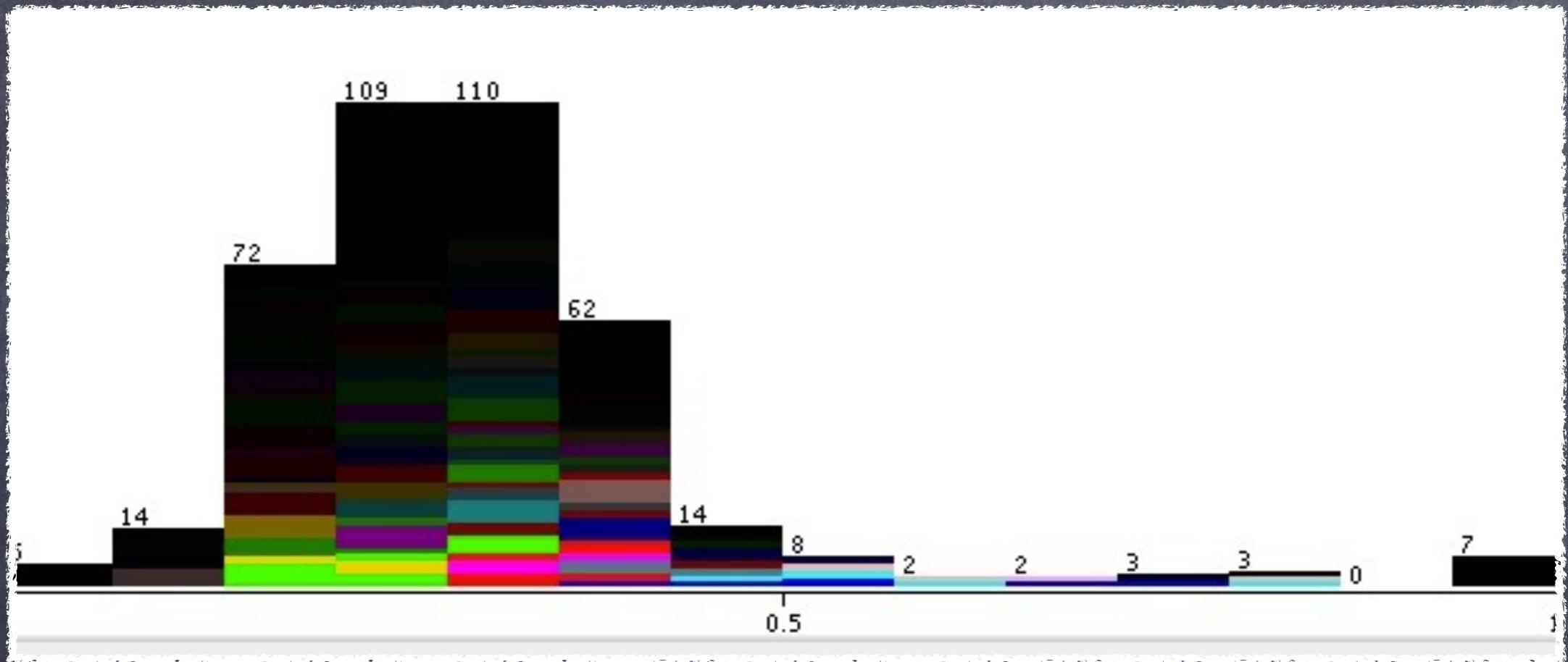


- The V80 values are fairly evenly distributed
- This indicates that this attribute on its own is good at differentiating between different characters

Attribute: Class vs Var

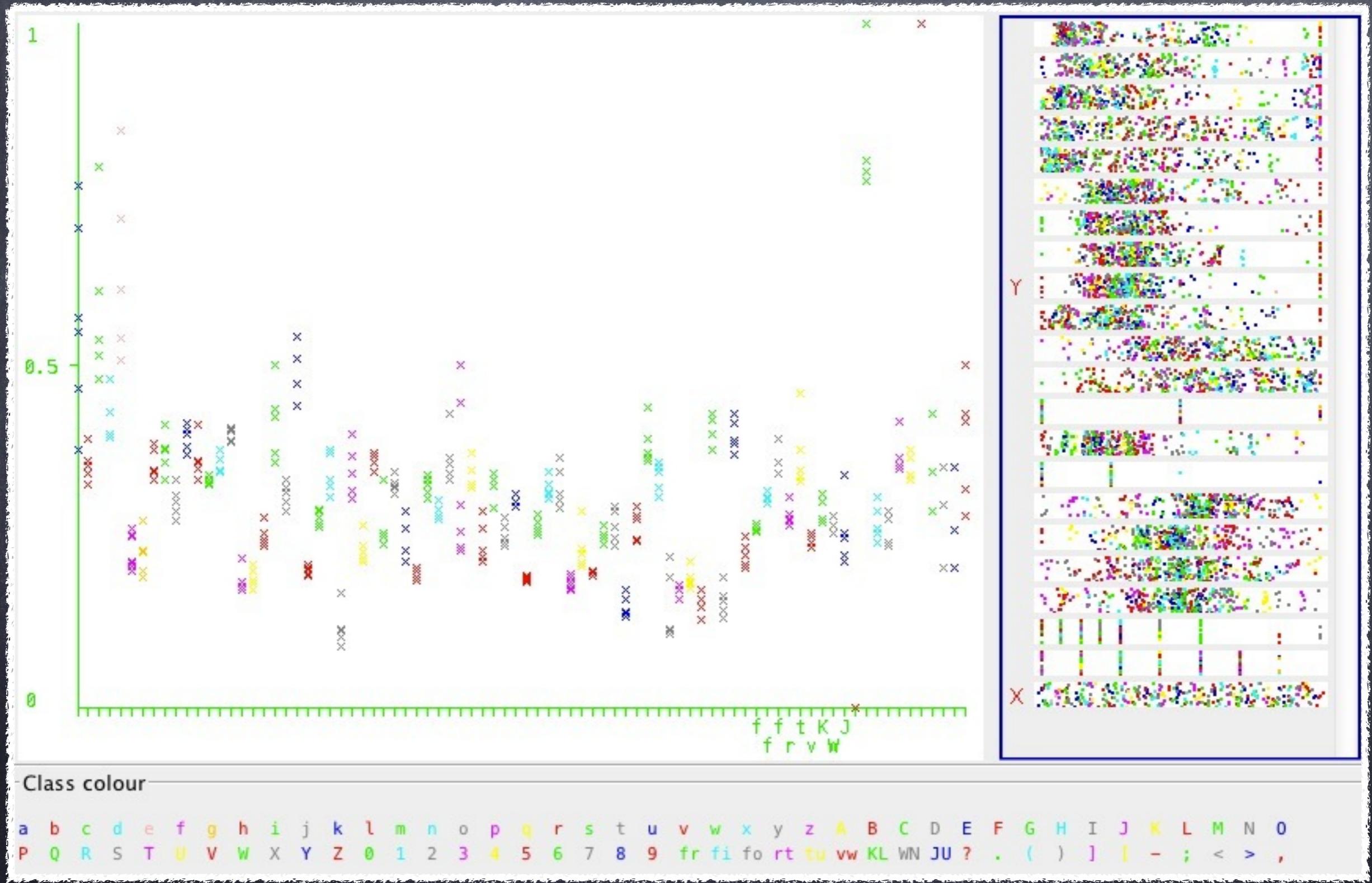


Attribute: H₂O



- The majority of the values tend to be around 0.25
- This indicates that this attribute on its own, is not particularly good at differentiating characters

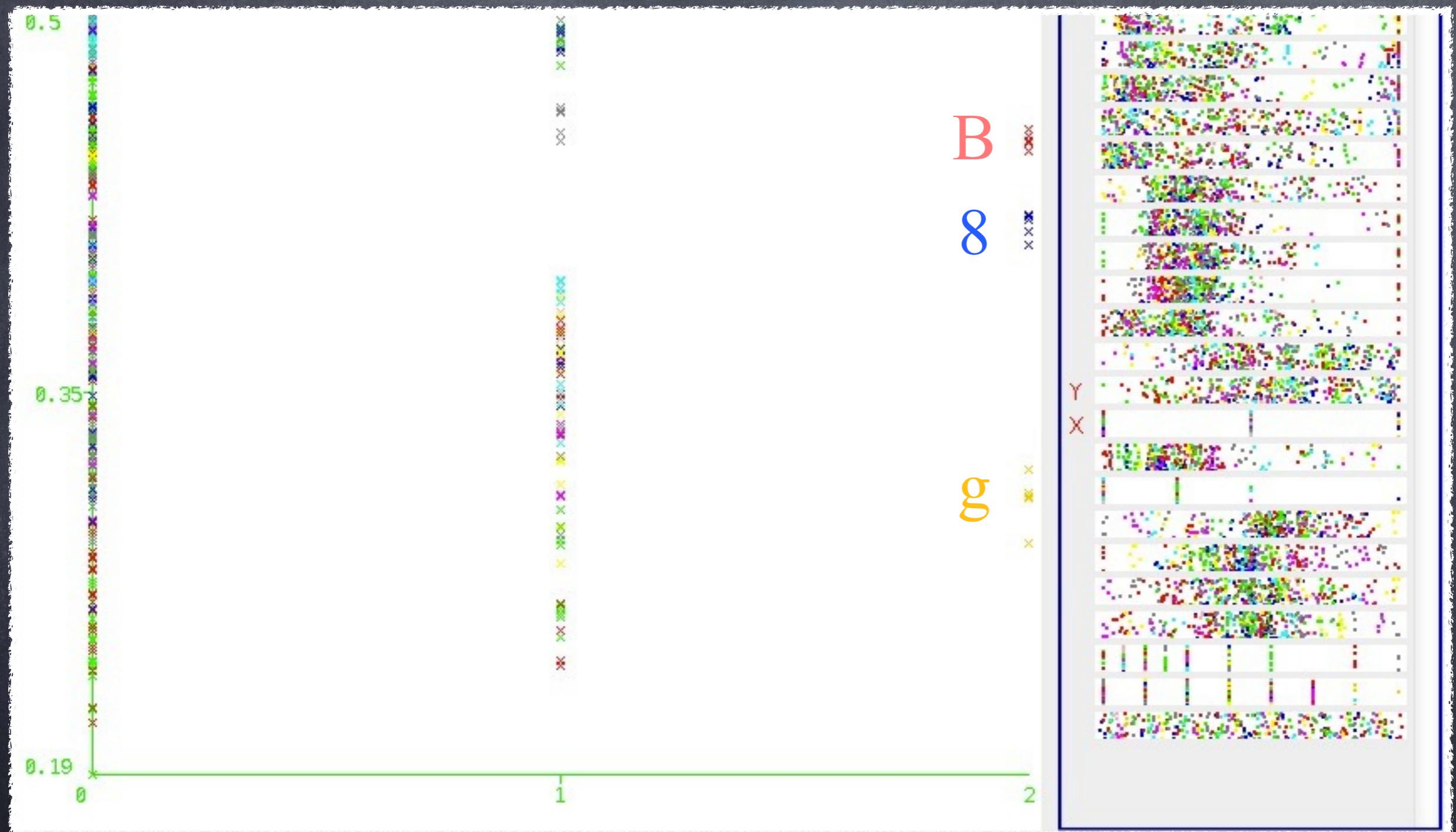
Attribute: CLASS vs H80



Attributes ...

- Of course, this is why we do not have only one point of comparison
- There are a total of 22 attributes that are compared between each test character and the training characters
- One important observation is that although most of the letters have similar H80 values, any outliers, in this case ":" and "-", are effectively isolated from the rest of the data by this attribute

CC vs Vsymm coloured by class value



Examples

Example 1

the quick brown fox jumps
over the lazy dog.

0123456789 THE QUICK
BROWN FOX JUMPS
OVER THE LAZY DOG ;
[] () ?

the quick brown fox jumps
over the lazy dog.

0123456789 THE QUICK
BROWN FOX JUMPS
OVER THE LAZY DOG ;
[] () ?

SVM

ANN

the quick brown fox jumps
over the lazy dog.

0123456789 THE QUICK
BROWN FOX JUMPS
OVER THE LAZY DOG ;

[] () ?

the quick brown fox jumps
over the lazy dog.

0123456789 THE QUICK
BROWN FOX JUMPS
OVER THE LAZY DOG ;

[] () ?

Example 2

c. Repeat the analysis for part (b) using the same
Compare the F-measure results for both classifier
results consistent with what you expect from the

c. Repeat the analysis for part (b) using the same
Compare the F-measure results for both classifier
results consistent with what you expect from the

SVM

c. Repeat the analysis for part (b) using the same

Compare the F-measure results for both classifier]

results consistent with what you expect from the

ANN

c. Repeat the analysis for part (b) using the same

Compare the F-measure results for both classifier]

results consistent with what you expect from the

Next Steps ...

- Pre-processing: Scanning algorithm used to isolated putative letters is not robust in terms of noise, would like to explore other options
- Pre-processing: image straightening step
- post-processing: determining spacing on edge cases (i.e. no separate words, or all separate characters)
- other: find frequent pairs of letters for training on pairs of letters that may end up clustered together

References

1. Neves, E. et al. (1997). IEEE. A Multi-Font Character Recognition Based on its Fundamental Features by Artificial Neural Networks
2. Shrivastava, V. & Sharma, N. (2012). SIPIJ 3(5). Artificial Neural Network Based Optical Character Recognition.
3. Hall, M. et al. (2009); The WEKA Data Mining Software: An Update; SIGKDD Explorations, Volume 11, Issue 1.
4. Tautu E.D. & Leon, F. (2012). Bul. Inst. Polit. Optical Character Recognition using support vector machines.
5. Sadri, J. et al. (2003). MVIP. Application of Support Vector Machines for Recognition of Handwritten Arabic/Persian Digits.
6. Rashnodi, O. (2011). International journal of computer applications 29(12). Persian Handwritten Digit Recognition using Support Vector Machines.