

---

# OCR using Machine Learning Techniques and Hidden Markov Models

---

**Jonathan Yee**  
jyee1@andrew.cmu.edu

**Lingzhang Jiang**  
lingzhaj@andrew.cmu.edu

## Abstract

Optical Character Recognition (OCR) has been the focus of much ML research. A reliable OCR system has a wide range of uses, from recognizing scanned books to license plate recognition. Since a sequence of characters (e.g. an English word) is not generated randomly, we can model the process of writing after a Hidden Markov Model (HMM). The hidden states are the letters the human is thinking of, and the observed states are what he or she actually writes. This HMM exploits the correlation between neighboring letters in the general OCR case to improve accuracy. We have found that using Maximum Likelihood Estimates (MLE) to learn the parameters, and using the Viterbi algorithm to find the most likely set of hidden states, leads to a recognition rate of 69.7%. This is higher than the accuracy rate of using Naive Bayes (62.7%) and Logistic Regression (64.1%).

## 1 Background

We are trying to apply different machine learning techniques to an Optical Character Recognition (OCR) problem, with the objective of obtaining a high level of accuracy comparable on a given data-set comparable to state-of-the-art techniques.

OCR can be analyzed using the noisy channel model, with the source being the human mind, and the channel being the handwriting technique. We want to determine the most likely letter given the handwriting.

For this project, we intend to design a Hidden Markov Model. The hidden states of the HMM will be the letters the human is thinking of, and the observed output is the actual handwriting represented as a pixel vector. This HMM can be graphically represented as a Bayes net. The goal is to find the sequence of hidden states (i.e. letters) that maximize the joint probability of those characters and the pixel vectors. Various OCR techniques (neural networks, naive bayes, logistic regression) can be applied to determine the probabilities that a pixel vector is a certain letter. We also intend to use the Viterbi algorithm to determine the most likely sequence of letters.

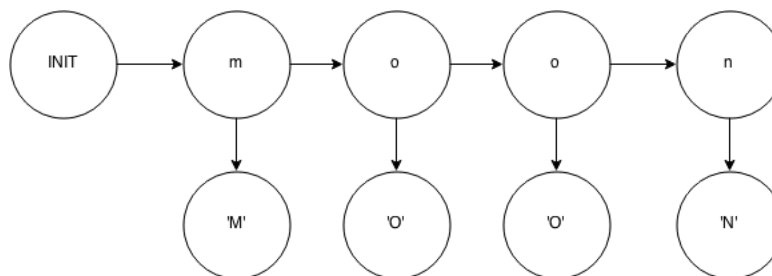


Figure 1: Hidden markov model for the word "moon". Observed states are the written pixels.

## 2 Related Work

The first related work we looked at was a similar study by Ivan Dervisevic[1], which utilized four classifiers(Naive Bayes, Complement Naive Bayes, SMO, and C4.5) to perform OCR and compared the results of these classifiers. SMO is a kind of support vector machine classifier while C4.5 is a decision tree classifier. The dataset that was used contained capital characters from 238 TrueType fonts and the data was 40x40 pixel images of each character. In addition to using the raw unprocessed pixel image data, refined datasets with processed representations of the characters were also used. As our focus is more on the machine learning portion and less on the image processing techniques, we will only be using 2 classifiers as a baseline, leaving out image processing, and focusing on implementing a HMM to optimize classification instead.

The second related work we looked at was an introduction to Hidden Markov Models[2] by Rabiner. Since we have not had any comprehensive experience with HMMs, we thought it would be a good idea to see what has already been accomplished with HMMs. In this paper, Rabiner provides a review of the theoretical aspects of using HMMs modeling, and also shows how they apply to machine recognition of speech. This overview serves to provide a good background on the theory and application of HMMs, as well as the limitations that Rabiner discusses in the final section.

## 3 Methods

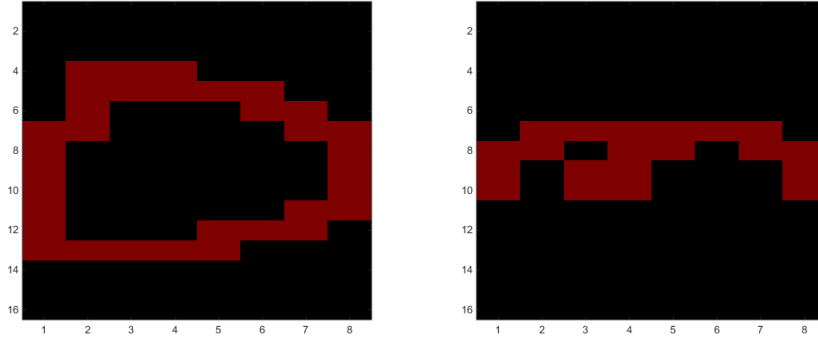
To obtain our baselines, we picked two machine learning methods that were taught in class, Naive Bayes and Logistic Regression. Both of these methods would take as inputs a data array containing the pixel data of a collection of handwritten characters as well as a vector containing the labels for the data array. After training the classifiers on the training set, we ran the trained classifiers on a test data set for which we know the correct labels, and compared the predicted labels with the correct labels to get the accuracy of the classifiers. Also, we did cross-validation to get a confident accuracy rating for our baseline methods.

As an extension of the baselines, we decided to use MLE to learn the parameters of the HMM. First, we learned the parameters by counting the letters and pixels and calculating the transition and emission probabilities using MLE. Then, we used the Viterbi algorithm to calculate the most likely sequence of hidden states given some test sample. This likely sequence of hidden states corresponds to the recognized letter for each pixel vector. We compared the predicted labels with the correct labels to get the accuracy of the classifiers. Finally, we did cross-validation to get a confident accuracy rating for our extension.

## 4 Experiments

### 4.1 Data Set

The dataset we are using contains handwritten words. The dataset was collected by Rob Kassel at MIT Spoken Language Systems Group. A subset of the words were chosen and the images of each letter were rasterized and normalized. Since the first letter of each word was capitalized and the rest were lowercase, the first letter was removed and only the lowercase letters were used. The given tab delimited data file contains a line for each letter, with its label, pixel values, and several additional fields listed in a separate file.



The above are matlab-generated images of the pixel data provided in the data set. The left image is an 'o' and the right image is an 'm'.

The data set can be accessed at <http://ai.stanford.edu/~btaskar/ocr/>

## 5 Results - Baseline

### 5.1 Naive Bayes

For our first baseline we performed Naive Bayes on the raw pixel data. If given some character  $c$  and an array of  $n$  binary pixel values  $\{p_1, \dots, p_n\}$ , the algorithm is based on basic Bayes rule:

$$P(c|p_1, p_2, \dots, p_n) = \frac{P(p_1, p_2, \dots, p_n|c)p(c)}{p(p_1, p_2, \dots, p_n)}$$

If we make the assumptions that individual pixels are conditionally independent given some character then we can expand the conditional probability in the numerator. Using the naive bayes algorithm in the statistics toolbox for Matlab, we ran Naive Bayes over our dataset consisting of pixel data for 52152 characters(26 unique characters, no capitals) with each observation having 128 binary pixel values.

N-fold cross validation splits the data into  $N$  disjoint sets. In each of  $N$  iterations training is done on  $N-1$  sets and testing is done on the remaining 1 set. If we take the mean accuracy we can then average out errors resulting from variance.

Iteration	Accuracy
1	62.53%
2	61.94%
3	62.93%
4	62.86%
5	63.67%
6	62.58%
7	62.59%
8	62.99%
9	61.84%
10	62.84%
Avg	62.7%

Table 1: Output of training and testing MLE-learned HMM over some number of iterations

We performed 10-fold cross validation for our Naive Bayes classifier and were able to obtain a final accuracy of 62.7%.

Figure 2 shows a visualization of the confusion matrix obtained from using a trained Naive Bayes classifier to predict a test set of characters. The y-axis shows the actual characters(1 is 'a', 2 is 'b' etc.), while the x-axis shows the predicted characters. The matrix is normalized as such:

$$I_{x,y} = \text{Classified}_{x,y} / \text{Total}_y$$

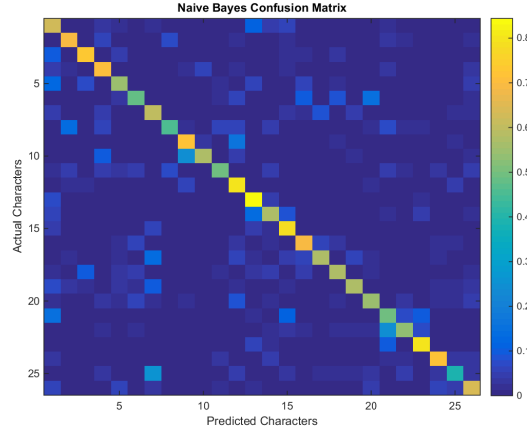


Figure 2: Visualization of the confusion mat for Naive Bayes

The value at each coordinate  $(x,y)$  is equivalent to the number of instances of character  $y$  that has been classified as character  $x$  divided by the total instances of character  $y$  in the test set. A perfectly yellow diagonal with deep blue everywhere else is thus a representation of a perfect classifier.

From the confusion matrix, some of the commonly misclassified letters are 'h', 'k', 'n', and 'y'. Specifically, looking at row 14, 'n' seems to be commonly misclassified as 'm' or 'o'.

## 5.2 Logistic Regression

The other classifier that we chose for our baseline is logistic regression. Logistic regression transforms a given binomial dependent variable, eg. the result of a coin toss, and applies the logistic function to it, effectively transforming it into a continuous variable, as such:

$$F(x) = \frac{1}{1 + e^{-(-\beta_0 + \beta_1 x)}}$$

The logistic function is useful because it is a function mapping from real numbers to an interval between 0 and 1, and hence the output can be treated as a probability. The logistic regression algorithm itself then follows an approach similar to linear regression to train a set of weight  $(\beta_0, \beta_1 \dots)$  that maximizes the likelihood of the data.

In our case, as we have 26 unique characters which cannot be represented by a binomial variable, we chose to use the multinomial logistic regression algorithm(mnrfit) provided in the statistics toolbox in Matlab. The model takes basically the same principles governing logistic regression for a binary dependent variable, except that it assigns one of the categories for the dependent variable as a 'reference category', and calculates the probabilities of an observation being in each of the other categories as opposed to being in the reference category. For the parameters we used the nominal(default) model, as there is no natural ordering among our response variable categories that would better suit an ordinal model.

When running the logistic regression algorithm, we noticed an abnormally long run time. Hence, we eventually reduced the number of examples in the data to 2000 as well as using principal component analysis to reduce the dimensionality of the data. Even so, there were warnings that the model failed to converge. We suspect that this could be a result of the sparseness of the data matrix as there are a large proportion of 0s in the data set.

The above diagram shows the confusion matrix when using a model trained by logistic regression(with 80% variance retained) to classify a test set. Comparing it to Figure 2, it seems like the mistakes that logistic regression are quite different from those made by Naive Bayes. For instance, Naive Bayes was able to classify 'x' and 'z' with fairly high accuracy, but not logistic regression. Furthermore, an interesting fact is that logistic regression seems to get specific characters completely wrong, for instance, classifying 'k' as 'n' in all instances.

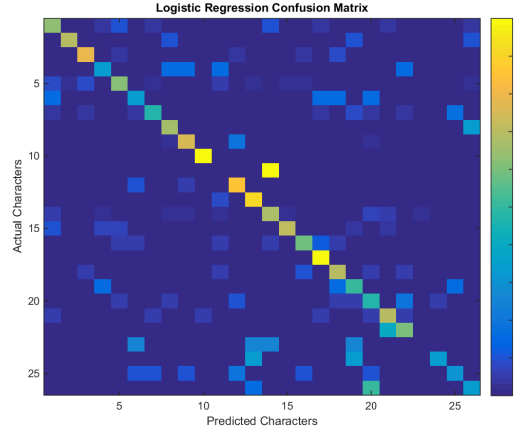


Figure 3: Visualization of the confusion mat for Logistic Regression

Iteration	Accuracy
1	72.00%
2	65.75%
3	64.50%
4	65.25%
5	64.75%
Avg	66.45%

Table 2: Output of training and testing Logistic Regression over some number of iterations

The above table shows the accuracy of a trained logistic regression model in which PCA was performed on a reduced data set with 2000 data points, keeping 70% of the variance, over 5-fold cross validation. There were significant gains over the Naive Bayes model with an average of 66.45% accuracy, but the training took a lot longer(Over 20 minutes even with parallel loops) even on the reduced data set with reduced number of features.

We were also curious about how the amount of dimensionality reduction would affect accuracy, and thus did logistic regression on 2000 data points with different extents of reduction(60%, 70%, 80% variance retained), resulting in overall 25, 34, and 49 features respectively, down from 128.

From Figure 4, it appears that both the variance and accuracy of classification decreased with increasing percentage of variance of the data retained, with 60% variance retained having the best results. This suggests that perhaps logistic regression is not a very good model for our classification problem as it is not able to make use of more data to improve the model's accuracy.

## 6 Results - Intermediate

### 6.1 MLE-Learned HMM Parameters

We can model a sequence of handwritten letters using a Hidden Markov Model (see Figure ??). Each hidden state is a letter  $c$ . Each observed state is a pixel vector  $v$ ,  $v_i \in \{0, 1\}$ , with length 128 (this is a  $16 \times 8$  image).

To obtain a Hidden Markov Model, we need to find the following parameters:

- $\pi_i$  : transition probability from state INIT to  $i$
- $\phi_{i,j}$  : transition probability from state  $i$  to  $j$
- $\theta_i(v)$  : emission probability associated with state  $i$  of a pixel vector  $v$

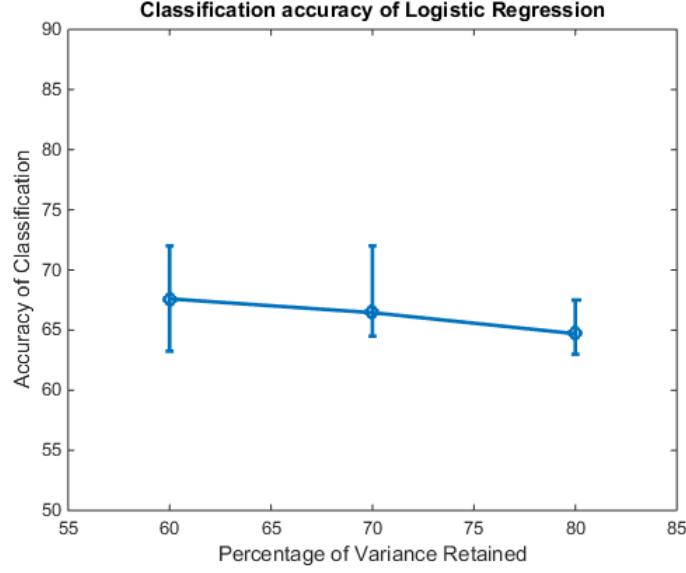


Figure 4: Logistic Regression Classification with different extent of dimensionality reduction

We plan to use Maximum Likelihood Estimates to find  $\phi_{i,j}$ . Let  $c_1, c_2$  be consecutive letters. Then this is equivalent to finding  $P(c_2|c_1)$ . This can be calculated using MLE as follows:

$$P(c_2|c_1) = \frac{\text{count}(c_1, c_2)}{\text{count}(c_1)}$$

We are counting the number of times the bigram  $\langle c_1, c_2 \rangle$  appears, and dividing by the number of times  $\langle c_1 \rangle$  appears.

To find  $\theta_i(v)$ , we find  $P(v|c)$  where  $c$  is the letter that is associated with state  $i$ . Assume Naive Bayes (conditional independence) assumption on  $P(v|c)$ , so we have

$$P(v|c) = \prod_i P(v_i|c)$$

This can be obtained from the training data by looking at all pixel vectors for letter  $c$ , and applying MLE to calculate  $P(v_i = 0|c)$  and  $P(v_i = 1|c)$  using the frequencies of the pixels that are 0 or 1. In other words,

$$P(v_i = 1|c = l) = \frac{\text{count}(c = l, v_i = 1)}{\sum_{j=0}^1 \text{count}(c = l, v_i = j)}$$

A word is a sequence of letters  $Y = \langle c_1, c_2, \dots, c_n \rangle$  with an accompanying sequence of pixel vectors  $\langle v_1, \dots, v_n \rangle$ , where  $n$  is the length of the word. To classify a word, feed the above parameters into the Viterbi algorithm, with the observed states as  $v_1, \dots, v_n$  to obtain the most likely path of hidden states  $\hat{Y} = \langle c_1, \dots, c_n \rangle$ , which are the letters. We can then compare the sequence  $\hat{Y}$  to the correct labels  $Y$  and derive the accuracy of the MLE-learned HMM at OCR.

After 5-fold cross validation, we were able to obtain a mean accuracy of 69.7% (see Table ??), which is more than 5% better than either Naive Bayes or Logistic Regression.

The figure above shows the confusion matrix obtained from using our MLE-learned HMM for classification of a test set. When we compare it to those obtained from Naive Bayes or Logistic Regression, there are some similarities, for instance low accuracy in classifying the character 'y'. Something interesting to note is how much trouble the model has with classifying the character 'o' (row 15), most commonly misclassifying it as 'a', 'c', 'e', 'g', 'n' or 'u', while Naive Bayes does not face the same difficulty. Since the emission probabilities essentially work like Naive Bayes, or a

Iteration	Accuracy
1	69.77%
2	70.08%
3	67.72%
4	70.61%
5	70.11%
Avg	69.7%

Table 3: Output of training and testing MLE-learned HMM over some number of iterations

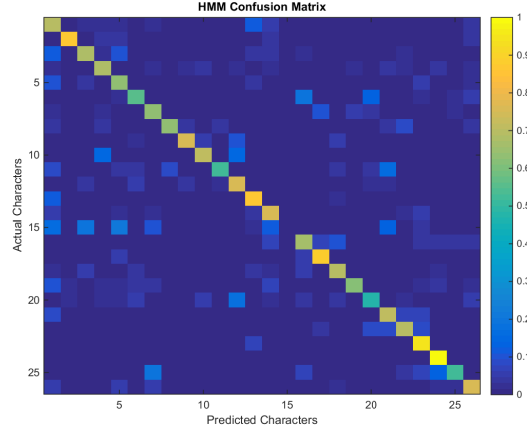


Figure 5: Visualization of the confusion mat for HMM

close approximation of it, in this case the transition probabilities we derived may potentially have decreased the classification accuracy instead, in a scenario where the emission probabilities were not certain to begin with, by increasing the uncertainty, since 'o' is a commonly used vowel and does not have as rigid transition rules as consonants.

## 6.2 MLE-Learned HMM Comparison to Baseline

The mean accuracy of the MLE-learned HMM (69.7%) is more than 5% more accurate than either the Naive Bayes classifier (62.7%) or the Logistic Regression classifier (64.1%). We theorize that this is because the HMM model accounts for the transitions between letters, hence is more likely to predict a character that, in the training data, has appeared after the previous character.

One of the shortcomings thus far might be the Naive Bayes assumption on the pixel vector, where we assume that the value of each pixel is independent of other surrounding pixels, given the same character. This is unfortunately not true, since a writing instrument is likely to have a line-width of more than 1 pixel. Furthermore, writing often happens in sets of continuous strokes. This implies that if one pixel is switched on, it is more likely that the surrounding pixels are switched on than not.

Another shortcoming is the assumption of the hidden states as being modeled as per figure ?? . This assumption lets us use MLE to calculate the transition states. This assumption may result in lower performance, and we will have to wait until we use the Baum-Welch algorithm to estimate the transition and emission probabilities without any assumptions on the hidden states in the learning phase.

## 6.3 Modified Naive Bayes

One thing that caught our attention was the fairly low classification accuracy of the Naive Bayes model. We suspected that one of the main problems with the model is the independence assumption. Basically, when taking the probability  $P(p_1, p_2, \dots, p_n | c)$ , of a certain pixel vector given some

character, we simply assume conditional independence and multiply  $P(p_1|c), P(p_2|c), \dots, P(p_n|c)$ . However, that is not likely to be true because we can imagine that even if we knew what some letter was, how the handwritten version appears is not fixed, and is indeed likely to be dependent on values of neighboring pixels. To put this to the test, we came up with a method of applying some filters to the data first and then classifying the processed data with a trained Naive Bayes model.

We used mainly two filters. The first assigned a unique value to each pattern that a given cluster can assume, with two sizes, 3x3 and 4x4, as below. For convenience we will refer to this as the bitstring filter.

$$\begin{bmatrix} 2^0 & 2^1 & 2^2 \\ 2^3 & 2^4 & 2^5 \\ 2^6 & 2^7 & 2^8 \end{bmatrix} \begin{bmatrix} 2^0 & 2^1 & 2^2 & 2^3 \\ 2^4 & 2^5 & 2^6 & 2^7 \\ 2^8 & 2^9 & 2^{10} & 2^{11} \\ 2^{12} & 2^{13} & 2^{14} & 2^{15} \end{bmatrix}$$

Looking at the 3x3 filter, for instance, we see it projects each pixel and its 8 surrounding pixels onto a bit-string of length 9, depending on the values they assume, hence assigning equal importance to all neighbors.

$$\frac{1}{16} * \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$

$$\frac{1}{256} * \begin{bmatrix} 1 & 4 & 6 & 4 & 1 \\ 4 & 16 & 24 & 16 & 4 \\ 6 & 24 & 36 & 24 & 6 \\ 4 & 16 & 24 & 16 & 4 \\ 1 & 4 & 6 & 4 & 1 \end{bmatrix}$$

The second filter we used was a simple gaussian filter that assigned the most importance to the central pixel and then less importance as the distance from the center increased. We will refer to this as the gaussian filter.

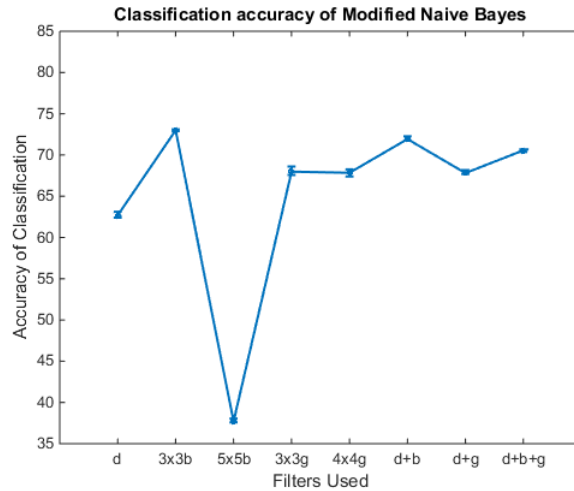


Figure 6: d=default. b=bitstring. g=gaussian. Combinations use 3x3 filters.

Looking at Figure 6, We see that all filters attempted except for the 5x5 bitstring filter helped to improve classification accuracy up from the 63% obtained on the default data. We think that the reason that the 5x5 filter had such bad performance is because the size of the filter is simply too large as compared to the overall size of the character, and hence resulted in too much loss of information. Also, it is noteworthy that the 3x3 bitstring filter had the best performance improvement of close to 10%, surpassing even the HMM model with MLE estimates, and also did better than the modified datasets containing combinations of filtered data and the default data(sets 6-8 on the x-axis).



## 6.4 Support Vector Machine

Support Vector Machine is a powerful machine-learning algorithm that is able to utilize kernels to find a maximum-margin hyperplane that separate the data points into their respective classes. For our project we use the libsvm toolkit created by Chih-Chung Chang and Chih-Jen Lin. The default soft-margin SVM, given some  $n$  training examples, is characterized by the objective function:

$$\operatorname{argmin}_{w,\xi,b} \left\{ \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \right\}$$

subject to the constraint

$$\forall_i y_i (w \cdot x_i - b) \geq 1 - \xi_i$$

where  $\xi_i$  is a slack term that allows a data point to trespass the margin, and hence 'soft-margin', and  $C$  is a term that controls how much to penalize for a given amount that the margin has been trespassed. We use the soft-margin SVM together with the radial basis kernel function for classifying our data set. The radial basis function is given by the equation:

$$k(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2), \gamma > 0$$

Hence, we mainly have 2 terms that can be optimized,  $\gamma$  and  $C$ . We use grid search with 5-fold cross validation to optimize across 5 values of  $C$  and 7 values of  $\gamma$ .

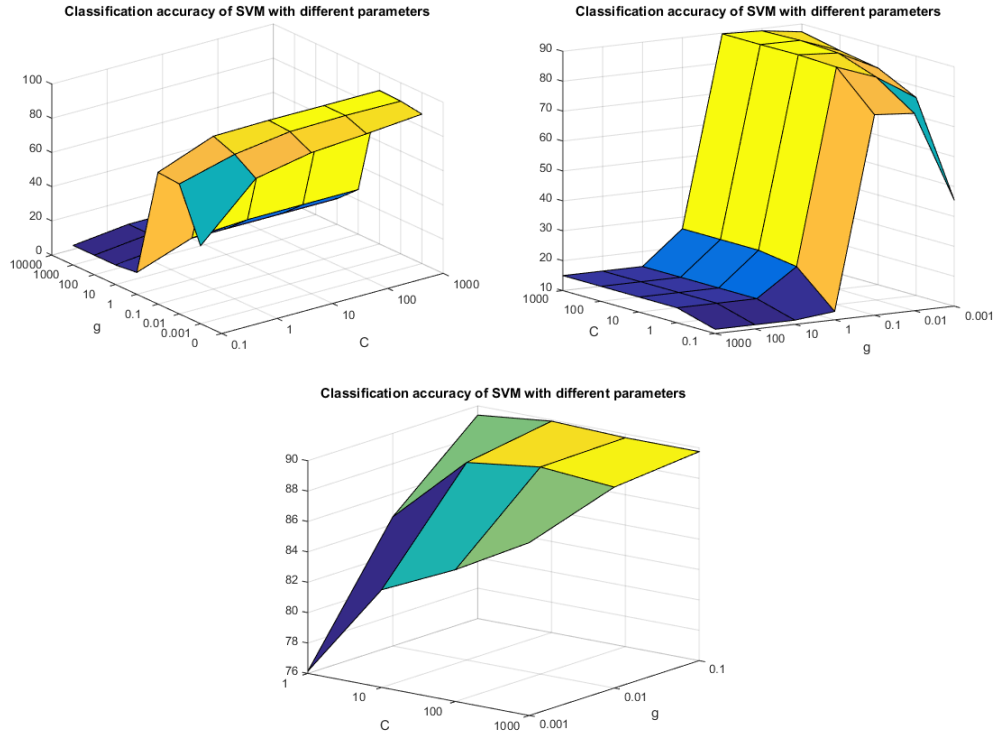


Figure 7: SVM accuracy with different values of  $C$  and  $\gamma$

From the top 2 graphs in Figure 7, we see that values of  $\gamma$  larger than 1 contributed sharp decrease in accuracy regardless of the values of  $C$ , and there seems to be an 'optimal range' of values, shown by the yellow grids, where accuracy hovered around 89%. After removing the sub-optimal range, we obtain the graph at the bottom of Figure 7, which shows our optimal value of  $C = 10$ ,  $g = 0.1$ , which was able to obtain a classification accuracy of almost 90%.

Analyzing the SVM confusion matrix is a lot easier than for the baseline classifiers because of the relatively smaller number of mistakes. Looking at Figure 8, there are a few outstanding pairs of classification mistakes, including ('q', 'g'), ('l', 'i'), ('x', 'y'), ('v', 'u'), and ('y', 'g').

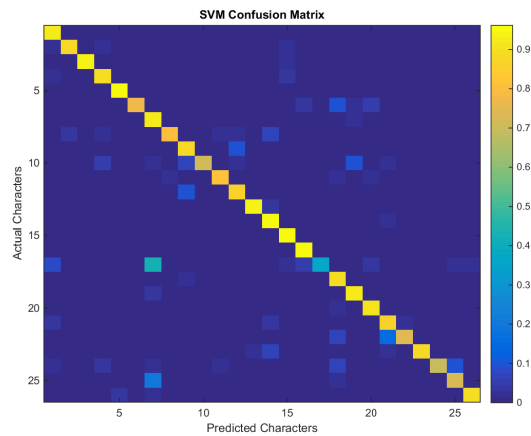


Figure 8: SVM confusion matrix

## 7 Results - Final

SVM and HMM are da bomb.

## 8 Conclusion

$x$  is clearly better than  $y$  and  $f$  this class.

## References

- [1] Ivan Dervisevic(2006) *Machine Learning Methods for Optical Character Recognition*. <http://perun.pmf.uns.ac.rs/radovanovic/dmsem/complete/2006/OCR.pdf>
- [2] Lawrence R. Rabiner (1989) *A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition* <http://www.ece.ucsb.edu/Faculty/Rabiner/ece259/Reprints/tutorial%20on%20hmm%20and%20applications.pdf>