

CS 440/ECE 448 - 3 Hours  
Spring 2015  
Assignment 3: Naive Bayes Classification

*Instructor: Svetlana Lazebnik*

*Harrison Kiang - hkiang2, Jonathan Park - jgpark2, Gabriella Quirini - gquirin2*

## Part 1: Digit Classification

In this section, we explore using the Naive Bayes model for classifying images. The images used were examples of numbers 0-9 written in different styles of handwriting. Most of them were easily classifiable to the human eye, but as we know, accurate image classification is a challenge to implement effectively even on the most powerful computer systems. Both parts of this section involve breaking down a given image into a set of features. These features define the image as a whole, but individually are nothing more than a binary value representing foreground (the area of the digit itself), and background (the rest of the 28x28 pixel space). Part 1.1 explores using single pixels as features. Later on, for extra credit we expanded the images to handle ternary values (foreground, background, middleground) for an overall improved accuracy with interesting results.

### Part 1.1 - Single Pixels as Features

Using this approach, each image is looked at pixel by pixel and broken down into a 2-dimensional list of binary values. In order to achieve this, we created a Digit class. Each Digit holds its own feature list, as well as the ASCII "image" (#/+) representation of itself. Additionally, a 'proper class' attribute holds the classification given to it post-initialization. In order to obtain the feature list, we simply parse the ASCII list and append 1's and 0's to our feature list according to foreground and background characters respectively. We use this method to set up both our training and testing data for use. For the ternary version implementation, we assign 1's to #, 0's to blanks, and 0.8's to +'s. We have experimented with different middle ground values, but found 0.8 to improve the overall accuracy more so than 0.1~0.7 or 0.9.

### Training

Training data refers to the set of images we already know the classifications for. We "train" our algorithm to recognize which class an unknown digit belongs to this way. We first set up a "digit database" of sorts, which is a 2-dimensional list of Digit objects. It's first dimension is of size 10, with each index corresponding to a digit class 0-9. Our training occurs by populating this database with the known Digit objects. Our resulting list is used to calculate feature likelihoods for each digit. The result is 10 28x28 lists whose values correspond to that pixel's likelihood of being in the foreground for the given digit. Laplace smoothing is used to ensure no zero-counts, and the priors are estimated based on the occurrences of training samples for each class. For the k constant of Laplace smoothing, we chose k=1 to get the best accuracy.

### Testing

We can now use our trained data set to classify unknown images using maximum a posteriori classification. This method utilizes the likelihood matrix and priors for each number class to generate a probability that the test image belongs to a given class. Looking at each pixel of the test image individually, we determine whether or not the test image fits a class based on the probability that the current pixel matches the corresponding pixel of that class. The prior term comes into play by acting as an overall likelihood of that digit occurring at all. Omitting this term turns it into a maximum likelihood classification. Using the same training and testing data, MAP

classification yields a 77.1% accuracy rating, while maximum likelihood yields 77.0%. The difference comes from accounting for the prior, and this gap would be larger had the training set been less balanced in terms of occurrences for each digit, and the testing set adjusted to match. However, as seen from the overall accuracy, the removal of prior terms for Maximum Likelihood seems to present little difference than MAP classification.

Evaluation			Confusion Matrix										
Classification Accuracy		Priors											
Class 0	84.44%	0.0958	0	1	2	3	4	5	6	7	8	9	
Class 1	96.30%	0.1126	0	.84	.00	.01	.00	.01	.06	.03	.00	.04	.00
Class 2	77.67%	0.0976	1	.00	.96	.01	.00	.00	.02	.01	.00	.00	.00
Class 3	79.00%	0.0986	2	.01	.03	.78	.04	.01	.00	.06	.01	.05	.02
Class 4	76.64%	0.1070	3	.00	.02	.00	.79	.00	.03	.01	.06	.02	.06
Class 5	67.39%	0.0868	4	.00	.01	.00	.00	.77	.00	.03	.01	.02	.17
Class 6	75.82%	0.1002	5	.02	.02	.01	.13	.03	.67	.01	.01	.02	.07
Class 7	72.64%	0.1100	6	.01	.07	.04	.00	.04	.05	.76	.00	.02	.00
Class 8	60.19%	0.0924	7	.00	.06	.03	.00	.03	.00	.00	.73	.03	.13
Class 9	80.00%	0.099	8	.02	.01	.03	.14	.02	.06	.00	.01	.60	.12
			9	.01	.01	.01	.03	.09	.02	.00	.02	.01	.80
orange = pairs most likely confused													

With the ternary version, MAP classification overall accuracy gives 77.8% and ML gives 78%. It is interesting to note the differences caused by adding the middleground.

Class Number	0	1	2	3	4	5	6	7	8	9
Accuracy %	85.6	95.4	77.7	83.0	73.8	64.1	78.0	70.8	66.9	82.0

Certain classes get a boost in accuracy, whereas other classes tend to lose some. Overall, the accuracy is 0.7~1.0% better.

**Odds Ratios:**

One way of inspecting what our program has “learned” is by generating images based off of the likelihood matrix. Another is by using odds ratios. Given classes  $c_1$ ,  $c_2$ , the odds ratio is a matrix formed by comparing corresponding pixels in the following manner:

$$\text{odds}(F_{ij}=1, c_1, c_2) = P(F_{ij}=1 \mid c_1) / P(F_{ij}=1 \mid c_2)$$

Here we are essentially mapping the areas that indicate whether an unknown belongs to  $c_1$  over  $c_2$ . If the unknown digit's foreground is contained within these “hot spots” (illustrated below), we can say that it is more likely to be a member of  $c_1$  over  $c_2$ .

The odds ratios printed in ASCII was easy to generate, but a bit unintuitive:

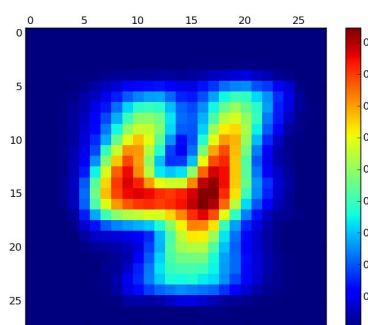
Matrix 2 comparing class 5 with class 3: (5 over 3)

[illegible]

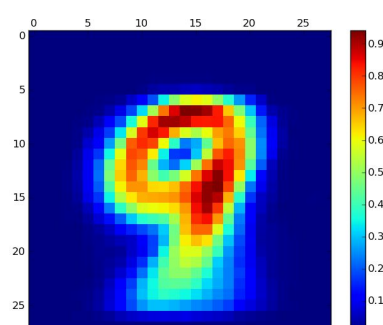
So we went through the extra effort of translating this into generating colored heat maps.

Below are heat maps corresponding to the feature likelihoods and odds ratios of the four most commonly confused digits based on our confusion matrix above.

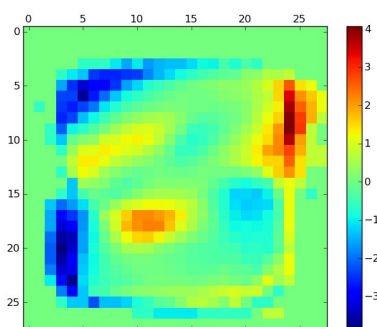
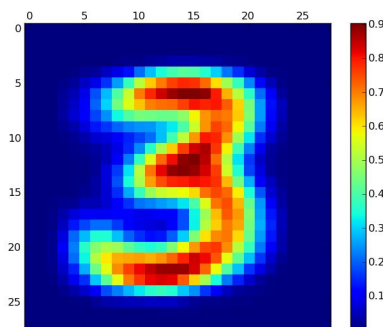
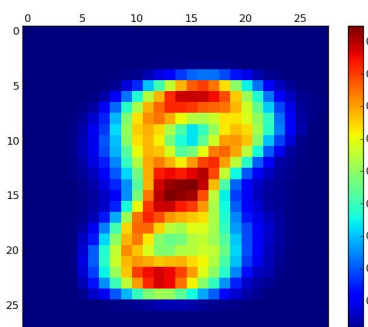
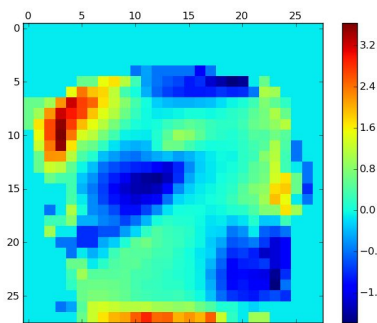
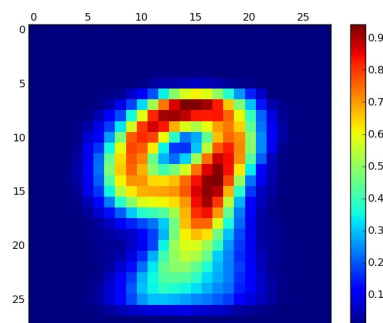
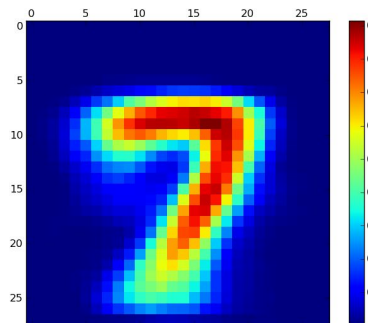
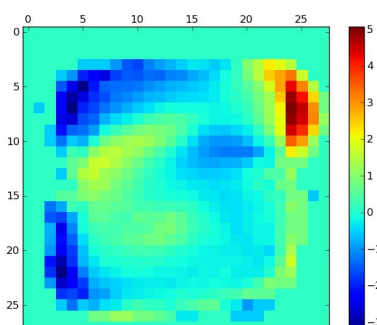
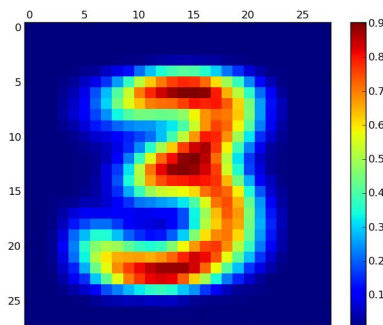
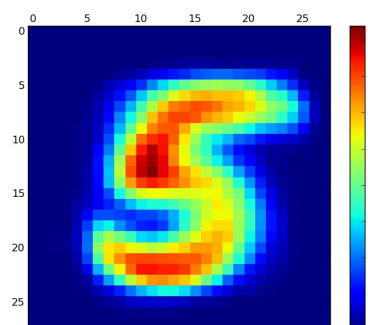
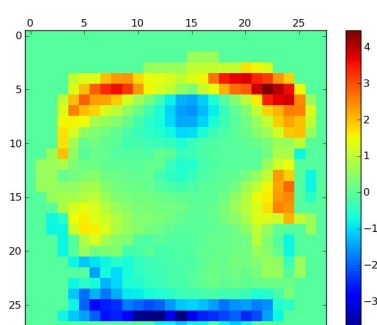
Likelihood Heat Map  $c_1$



Likelihood Heat Map  $c_2$



Log Odds for  $c_1$  over  $c_2$



## Faces:

Our above naive bayes classifier was modified to classify faces. Each face image is either a face (1) or a 'non-face' (0). Each face image is 70x60 and is populated by '#' and ' ' characters, representing edge pixels non-edge pixels, respectively. The methods of training and classification are the same as above.

## Training:

Number of training non-face images : 234

Number of training face images : 217

There are 451 numbers in training set

Prior for non-face class : 0.518847006652

Prior for face class : 0.481152993348

Overall accuracy(MAP): 90.6666666667 %

Overall accuracy (ML): 90.6666666667 %

Confusion Matrix:

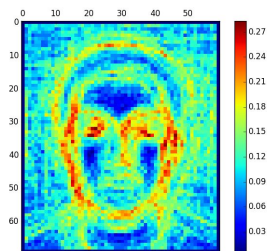
0.8831 0.1169

0.0685 0.9315

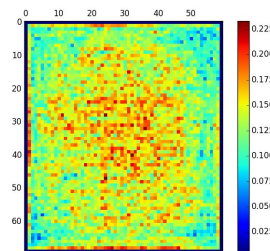
Accuracy for non-face class : 88.31 %

Accuracy for face class : 93.15 %

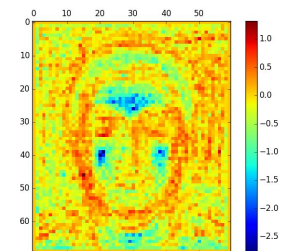
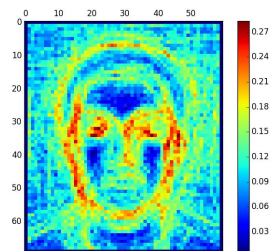
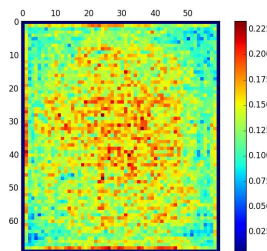
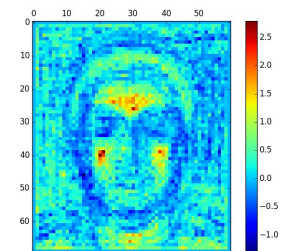
Likelihood Heat Map  $c_1$



Likelihood Heat Map  $c_2$



Log Odds for  $c_1$  over  $c_2$



## Part 2: Text Document Classification

This section aims to apply Naive Bayes method for classifying written documents. The datasets provided are preprocessed into “bag of words” representation, each containing keywords, their labels, and number of occurrences. Labels correspond to the category that word belongs to, and are represented by numbers. In the email dataset, 0 and 1 refer to normal and spam respectively, while numbers 0-7 correspond to the categories of the newsgroup dataset. We begin by creating dictionaries from the training data. Then, we estimate the conditional probability tables over those dictionaries for each class. We use Laplace smoothing to help classify words not found in our dictionary. Overall, the training and testing process is very similar to part 1 above.

Email Evaluation			Confusion Matrix											
Overall Accuracy: 96.53%			<table><tr><td></td><td>0</td><td>1</td></tr><tr><td>0</td><td>.9461</td><td>.0538</td></tr><tr><td>1</td><td>.0153</td><td>.9846</td></tr></table>				0	1	0	.9461	.0538	1	.0153	.9846
	0	1												
0	.9461	.0538												
1	.0153	.9846												
Category Classification Accuracy		Priors												
Spam	94.61%	0.5493												
Normal	98.46%	0.4507												

8 Category Evaluation			Confusion Matrix							
Overall Accuracy: 81.81%										
Category Classification Accuracy		Priors	0	1	2	3	4	5	6	7
sci.space	79.41%	0.1598	0	.79	.00	.00	.00	.20	.00	.00
comp.sys.ibm.hardware	57.58%	0.0701	1	.06	.57	.00	.12	.21	.00	.03
rec.sport.baseball	88.89%	0.1072	2	.00	.00	.88	.00	.08	.00	.02
comp.windows.x	85.71%	0.1364	3	.00	.00	.00	.85	.03	.00	.10
talk.politics.misc	97.87%	0.2034	4	.02	.00	.00	.00	.97	.00	.00
misc.forsale	10.00%	0.0494	5	.20	.10	.00	.00	.30	.10	.30
rec.sport.hockey	93.48%	0.1421	6	.00	.00	.00	.00	.06	.05	.93
comp.graphics	82.76%	0.1313	7	.03	.00	.00	.03	.10	.00	.82
			orange = pairs most likely confused							

### Top 20 Words With Highest Log-Odds Ratios - Email

$c_1 = 0$ over $c_2 = 1$ (normal over spam)	$c_1 = 1$ over $c_2 = 0$ (spam over normal)
linguistic workshop abstract theory syntax grammar chair discourse translation movement computational committee semantic benjamin verb context programme phonology semantics nl	sell nbsp ffa remove capitalfm bulk money floodgate investment bonus aol profit mailing save links reports offshore hundred debt advertise

### Top 20 Words With Highest Log-Odds Ratios - 8 Category

$c_1 = 5$ over $c_2 = 4$ (misc.forsale over talk.politics.misc)	$c_1 = 5$ over $c_2 = 7$ (misc.forsale over comp.graphics)	$c_1 = 1$ over $c_2 = 4$ (comp.sys.ibm.pc.hardware over talk.politics.misc)	$c_1 = 0$ over $c_2 = 4$ (sci.space over talk.politics.misc)
dos wolverine shipping comics spiderman liefeld sabretooth rider bagged cd hobgoblin xforce hulk disks punisher unix marvel ufl panther mutants	wolverine comics hulk spiderman liefeld sabretooth app bagged hobgoblin xforce punisher cable marvel panther mutants pom comic mint bwsmith rider	scsi ide controller bios disk dos isa interface modem floppy motherboard hardware jumper rom adaptec pc vlb card drive cpu	orbit shuttle spacecraft venus lunar mars sky satellites solar hst telescope propulsion launch orbital orbiter jupiter moon saturn kilometers probes



Each dictionary for each category (normal and spam for email, categories 0-7 for newsgroups) was traversed. Each word was added to a text file in a separate line and repeated on that same line according to the frequency. Each dictionary's word's frequency takes into account repeats and is based off of the given training data sets.

[illegible][illegible]

[illegible][illegible]



[illegible]

### Category 5: misc.forsale



### Category 6: rec.sport.hockey



### Statements of Contribution:

Jon: debugging, coding part1: ternary, heatmaps, some report

Gabriella: coding part1, debugging, report

