

KI Project C

Martijn Dekker 6013368
Egor Dmitriev 6100120
Cody Bloemhard 6231888
L^AT_EX

January 21, 2019

1 Exercise 1

1.a ex1-a

number of input features + what they represent

Images are converted into 2 dimensional arrays with size 28x28 for digits and 60x74 for faces. In these arrays every element represents one pixel of the image.

the possible values and what these represent

For every element in the array there are two options. For digits a 0 indicates a white pixel and a 1 indicates a gray/black pixel. For faces a 0 indicates no edge and a 1 indicates an edge.

the output labels and what they represent

For digits there are 10 possible labels. These are the numbers 0 to 9, and they represent the numbers 0 to 9. For faces there are 2 possible labels: 1 and 0, because it's either a face or it's not. Here a 1 represents a face, and a 0 represents it's not a face

the frequency or probability distributions over the output labels

1.b ex1-b

Most frequent counts for every possible label how often it appears and then uses the most common label to classify an input.

Naive bayes uses the log-joint distribution, so it also looks at how often a label is given, but unlike most frequent, it does so for every feature and normalizes the results it found.

1.c ex1-c

Most frequent does use supervised learning because it directly looks at the labels given to the training data, and uses the most common label to classify all future inputs.

Naive bayes does also use supervised learning because it looks at the labels by counting how often a level A is given to a feature B.

2 Exercise 2

2.a ex2-a

| Classifier | data | Validation | testing | k |
|---------------|--------|------------|---------|---|
| Most frequent | digits | 14% | 14% | x |
| | faces | 56% | 53% | x |
| Naive Bayes | digits | 69% | 55% | 2 |
| | faces | 77% | 75% | 2 |

What is important to know is that for digits there are ten possible labels and for faces there are only two possible labels. This means that if you randomly guess you will have a higher chance of guessing correctly with the faces data set. This is what happens for the most frequent classifier. The classifier found the most common label in the training set and apparently this label also appears with an above average frequency in the validation and testing set.

Therefore the scores for digits is 14% instead of 10%, which would be the case if all labels were equally common in the test and validation set. And the same is true for faces, if both labels were equally common in the test and validation set, then the scores would be 50%, instead of 56% and 53%.

The naive bayes classifier scores higher because it really looks at the image. It tries to find good features that really help distinguish between the different labels. The score for faces is higher for naive bayes because with the digits the different images will overlap, even though they are not necessarily the same label. This is because multiple digits have, for example, a curve in the top. And with faces it is easier to find a feature that distinguishes well, because there are only two possible options.

2.b ex2-b

Imagine we have 5 possible labels, these labels do not occur with the same frequency, but rather one label is quite rare. Then we can have a big training set where this label does not occur. And so the classifier will assign a probability of 0 to it. Then when we come across an input in our test data with this label, the classifier can not possibly assign the correct label. In order to prevent this, and make sure that every label has a chance to be assigned, we use laplace smoothing. This entails that we will always add a standard value to the frequency of a label. The probability used to be: frequency of a label divided by the total number of elements in the training data. But with laplace smoothing it becomes: $(\text{frequency of a label} + 1) / (\text{number of elements} + \text{number of labels})$.

2.c ex2-c

| Classifier | data | Validation | testing | k |
|-------------|--------|------------|---------|------|
| Naive Bayes | digits | 74% | 65% | 0.1 |
| | faces | 85% | 82% | 0.05 |

The autotune option finds the best value for k among these [0.001, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 20, 50]. And then the classifier uses the best k value

that was found for the validation and test set. This option helps to improve the score for the naive bayes classifier and will in general always give better scores. It is better for both data sets, but there could be a dataset where a k of 2 would give the best scores and there the autotune would give worse scores. What does stay the same is the presence of a difference between accuracy on validation and testing set. Autotune will always choose a k that performs best on the validation set, but this is not necessarily the best k for the test set.

2.d ex2-d

For naive bayes: first it gets the classifier by looking at the training set, then the validation set can be used to get a good k value and make sure we do not get a classifier that does really well on the training set but poorly on other sets. And as last it is tested on the test set to get an indication of how the classifier will perform on real life data sets. For most frequent: it gets the classifier by looking at the training set. It looks which label is the most common and then uses that as classifier. It does nothing important with the validation set. It merely looks what score it gets on there. Then it uses the test set to get an indication of how the classifier will perform on real life data sets.

2.e ex2-e

The test set is used to get an indication of the expected performance in the real world. This is because the data in the test set is not used for training. This means it is not overfit on the test set and you can see how well the classifier generalizes. This is important, for it will receive data it has not been trained on either in the real world. The validation set can be use to fine tune the classifier and rule out any form of overfitting. This happens with the naive bayes classifier when the autotune option enable. The validation test is used to get the best k value from the training data.

3 Exercise 3

3.a ex3-a

Because there are a lot of features which vary between zero and one multiplying them together can yield very small numbers. These numbers in some cases cannot be represented using float and double types.

A workaround for this is representing probabilities a log probabilities. This way their values are between $-\infty$ and 0.

Also because of this we no longer use product, but use summation instead. This works since sum of logarithms is equivalent to the log of the product of all the probabilities.

3.b ex3-b

Conditional probability can also be written as $P(A\hat{B}) = P(A|B)P(B)$. With this we can see that joint probability yields the same result as posterior probability multiplied with probability of $P(B)$.

Using log in joint probabilities achieves same result as described earlier because

sum of logarithms of the probabilities is equivalent to the log of the product of all the probabilities.

TODO: I could be wrong here at my first point. Doublecheck

3.c ex3-c

| | 0.55 | 0.60 | 0.65 | 0.70 | 0.75 | 0.80 |
|---------------------|-------|-------|-------|-------|-------|-------|
| accuracy test | 85.0% | 85.0% | 85.0% | 85.0% | 85.0% | 85.0% |
| accuracy validation | 81.0% | 82.0% | 82.0% | 82.0% | 82.0% | 83.0% |

Thresholding does seem to improve the accuracy, but not by much. The fact that accuracy improves suggests that the classifier did have some false positives.//

TODO: this can't be right. Reproduce `dataClassifier.py -c naiveBayes -d faces -k 0.1 --threshold=0.45,0.55`

3.d ex3-d

The naive bayes classifier would be more biased towards classifying numbers as even numbers. Normally one would need to balance the dataset first, but this would require retraining. Since naive bayes works using frequency counting we can also correct the amount of even number occurrences in the frequency table by reducing them by 50/dataset out.

4 Exercise 4

4.a ex4-a

| Training iterations → | 1 | 2 | 3 | 4 | 5 | 6 | 10 | 50 |
|-----------------------|-----|-----|-----|-----|-----|-----|-----|-----|
| Validation accuracy | 63% | 60% | 55% | 55% | 56% | 56% | 56% | 56% |
| Testing accuracy | 57% | 57% | 48% | 54% | 54% | 54% | 54% | 54% |

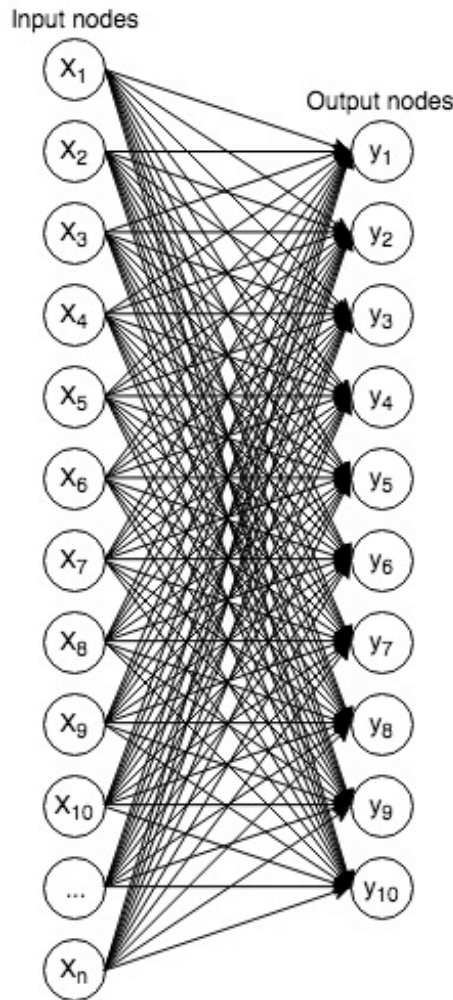
The perceptron goes through the examples same way as they are ordered in `trainingData`. The element that is on position 0 in `trainingData` will be processed first. Whether or not the perceptron is better than Naive Bayes depends heavily on the options one might change when training on the data. As follows from the data given in section 2.1 Naive Bayes gets a 69% score for validation and a 55% score for testing with a K of 2. If you compare this to the score of the perceptron with 5 iterations then the Naive Bayes classifier is obviously better because both the validation and testing score is higher. However, when you compare it with the perceptron with 1 or 2 training iterations, the perceptron scores are better. The Naive Bayes validation accuracy might be higher but the testing accuracy is the important score because this will give the best indication for future scores. The testing accuracy is higher with 57% versus 55%. But once you use autotune to get a better score for Naive Bayes, it is significantly better than the best score of the perceptron with a validation accuracy of 74% and a testing accuracy of 65% for a K of 0.1. So in most cases the Naive Bayes classifier is better than the perceptron.

4.b ex4-b

Based on the scores collected in part 4.1 it is safe to stop training after 5 iterations. This is because the scores you get when training even more are not better than after 5 iterations, in fact, they are exactly the same. However, when faced with a different training set it is possible that the perceptron's score keep improving after 5 iterations.

4.c ex4-c

Figure 1:



The activation function for the perceptron is $y_j = \sum_{i=1}^n X_i * w_i^j$. Here every connection has a weight w_i^j where i indicates the input node and j indicates the output node.

5 Exercise 5

5.a ex5-a

Sequence weights (a) are more representable of the perceptron because these represent features that are weighted the most in classifying the given label. Because these features are mostly used to determine whether an input corresponds to the label it is fair to conclude that perceptron sees weights 'a' as the ones that represent the labels the most.

5.b ex5-b

Perceptron is capable of correctly classifying separable data. In this case digit writings can not be considered separable data. This is because a lot of digits look very similar. They also depend on handwriting and because of this a lot of different digits can look alike when written by different people.

6 Exercise 6

6..1 ex6-a

By limiting weight change extreme outliers can not influence the weights as much. By not limiting the weight change large errors may cause weights to overshoot the minima which causes the weights converge slower. Because there is a large number of features, they may vary a lot in values which may cause such a slowdown. // TODO improve answer (about many features)

6..2 ex6-b

Autotune trains the model with different values for constant C and picks the best value which achieves the best accuracy.

Autotune has used three different values for C 0.002, 0.004, 0.008. Between those values of 0.004 for C has given the highest accuracy on the validation set. MIRA has given 68.0 for validation and 48.0 with k=0.1.

MIRA has better performance than perceptron but performs worse than best naive bayes.

MIRA performs better than perceptron because it limits the learning rate and tries to put a margin between the boundaries which gives it an advantage at classifying unseen data points.

MIRA itself is still perceptron based and does not handle inseparable data. Classify digits with pixels as features is an example of inseparable data and therefore causes MIRA perform worse than naive bayes.

6..3 ex6-c

The first thing to notice when comparing most relevant weights from perceptron and mira is that more dense regions on the perceptron are even denser in mira and less dense regions are less dense at MIRA. This is mostly noticable with digits like 7 and 5.

The only difference between the two algorithms is the learning rate. MIRA adds a variable learning rate which tries to add a margin between the

boundaries. This is responsible for making regions more or less dense than they are on the perceptron.

7 Exercise 7

7.a ex7-a

No naive bayes classifier is not a suitable instrument for classifying digits because it assumes the values of the features are independent of the values of any other features given a class variable. This is not true because digits aren't made of independent pixels. They are made of lines / curve segments which together form a digit. Therefore it is possible that values of features are dependent on other features.

7.b ex7-b

Given command shows all the features which are more likely to belong to label 4 in comparison to label 2. -1 and -2 take as argument the labels that need to be compared. Since features belong to pixels, these features are drawn again as pixels/characters. The likeliness is represented by the weights. The label that has a higher weight for a feature in comparison with another, then it is more likely to represent that label.

7.c ex7-c

In the following list we describe a few features which would have a positive impact at classifying digits. Some of these are non binary or take up multiple binary inputs.

Hole count Some digits have holes. Some have multiple. For example if a digit has 2 holes, then it is quite certain that it is an 8 since only 8 has 2 holes. We can count holes by partitioning image into multiple groups of empty/black pixels which are in direct contact. If we reduce this number by once, since background doesn't count as a hole. Then we can map it into 3 binary features (has no holes, 1 hole, 2 or more). This feature is implemented in our code and has increased the performance.

Curvature/Convexity defects Almost all digits are concave except zero. One can fit a polygon around the borders of white pixels and reduce it to have fewer polygons. This way one can determine the count of convexity defects in the shape using gift wrapping algorithm. Then this can be stored in multiple binary features (Ex. 0, 1, 2, 3 convex defects). This way a number can be determined quite accurately. 0 has 0 defects, 1 has 1 defect 2 has 2 defects, 3 has 3 defects, 4 has 3/1 defects and 9 has 1 defect.

Color cross rate One can count how many times pixel color changes (white -> black) in a row or column and add 3 binary features for each row/column representing 0, 1 or 2 color changes. This is another way to determine if the input has holes. Since 9 would have highest values in top rows meanwhile 6 would have those in lower rows.

Density One could split the shape into multiple smaller squares and calculate the density/amount of white pixels in the given region. Then represent each square in as a binary feature. One could use 0 als "not dense" or "less than 50while 1 gives the opposite. Some digits have high density in their corners. These corners have different positions for every digit.

7.d ex7-d

We have implemented the first described feature. We have done this by taking each pixel not in closed list and expanding it in all directions using depth first search. Neighbours that are out of bounds or are a white pixel are not added. This has increased the accuracy on the test set from 78% to 84% which is quite an improvement and it makes sense since this is a string indicator which type of digit it is. Digit 8 can be easily classified since its the only that has 2 holes. Also numbers 1, 2, 3 cant be mistaken for other numbers since they dont have any holes.

8 Exercise 8

8.a ex8-a

8.b ex8-b

8.c ex8-c

9 Exercise 9