

Introduction to Machine Learning

Spam Filter Coursework Report

Dylan Cope and Rowan Kypreos

1 Simple Naive Bayes

1.1 Constructing the Model

To prepare the data for a Naive Bayes model boolean features needed to be extracted from the emails. First the emails were preprocessed by splitting the text by whitespace and filtering out any strings that aren't composed purely of alphabetical characters. Next a list of words needed to be determined for a feature extraction process that checked the presence of a word in a email. In constructing this model a process for extracting feature words from the training set needs be provided. The primary reason for including the extraction process here is to ensure the feature words are determined purely from the training data, if the words were chosen from the entire available set of emails any scoring of the classifier would be invalid as it would be impossible to tell if the model had overfit.

The Naive Bayes classifier for is defined by two estimator vectors $\hat{\theta}_s$ and $\hat{\theta}_h$ corresponding to spam and ham. Each element in the vectors then corresponds to a probability associated with the presence of a feature word. To tune the classifier's estimators the training emails E are transformed to a matrix X where each column corresponds to a feature word and each row corresponds to an email. The value of each entry in X is either a 1 or a 0 depending on if the word for the column is in the email for the row. The rows of X are then grouped by class and summated down the columns. The resultant bit vectors are then divided by the number of instances in that class to give the final estimator vectors with elements of the form, $\hat{\theta}_i = \frac{x_i}{N}$ where x_i are the observations and N is the number of trials.

So far what been described is a simple Naive Bayes model with no smoothing. The implemented spam filter uses additive smoothing to allow the assignment of non-zero probabilities to words which do not occur in the sample. This changes the class estimators to be of the form shown in equation 1, additionally the logarithm is taken to avoid arithmetic underflow as we expect some probabilities to be very small,

$$\hat{\theta}_i = \frac{x_i + \alpha}{N + \alpha d} \quad \therefore \log \hat{\theta}_i = \log(x_i + \alpha) - \log(N + \alpha d) \quad (1)$$

where α is the smoothing parameter and d is the number of features.

Finally with the tuned estimators a decision rule is determined for classifying future data. Given an observation \underline{x} and class priors p and $(1 - p)$ for ham and spam respectively we can define a classification function C ,

$$C(\underline{x}) = \begin{cases} \text{ham} & \text{if } pP(\underline{x}|\hat{\theta}_h) > (1 - p)P(\underline{x}|\hat{\theta}_s) \\ \text{spam} & \text{otherwise} \end{cases}$$

Through assuming the observations are probabilistically independent the values $P(\underline{x}|\hat{\theta})$ can be computed by multiplying together the elements $\hat{\theta}$ that correspond to true values in \underline{x} .

1.2 Cross-Validation Scoring

Overfitting can occur by learning model parameters on the same data that's used for testing. A common solution to this is to partition the data into training and testing splits, however when evaluating across a models hyperparameter space there is still a risk of overfitting as the parameters will be tweaked to optimal performance.

This leads to the idea of splitting the data three ways; testing, training and validation sets. The model's parameters can be tuned on the training set, the hyperparameters can be optimised on the validation set, then the final model can be assessed on the testing set. However splitting the data three ways means we drastically reduce the number of samples that can be used for learning the model, and the results can depend on the particular random choices made to split the data.

This is solved with cross-validation (CV), a test set is still held out for final evaluation, but the validation set is no longer needed when doing CV. In the basic approach, called k -fold CV, the training set is split into k smaller sets. A model is then trained using $k - 1$ of the folds and validating on the remaining fold. This process is then repeated k times with differing splits on the folds, and the procedure returns a list of k mean accuracy scores from each trained model. From this list an overall mean score can be determined with a margin of error.

1.3 Comparing to a Baseline

With the methodology for constructing and scoring classifiers determined, a baseline needed to be established to ensure the trained models are statistically significantly better at classifying the data. A classifier that makes predictions by randomly assigning classes obtains an accuracy of 0.5, while this is an acceptable baseline making assessments it doesn't reflect the class proportionality of the data. Instead we consider a classifier that always predicts the majority class, as the data is 80% ham and 20% spam, this classifier always predicts ham and has an accuracy of 0.8.

In section 1.1 the method for constructing a model was details. The construction required a given feature word extraction process and hyperparameters for smoothing α and class priors p . Finding optimal hyperparameters for the classifier will be discussed in section 2.2, for now we follow Laplace's rule of sucession, i.e. $\alpha = 1$, and we let the class prior have no influence, i.e. $p = 0.5$.

A simple feature words extraction process was be provided to the constructor whereby n random words from the training set were taken as feature words.

Figure 1: Number of features against performance

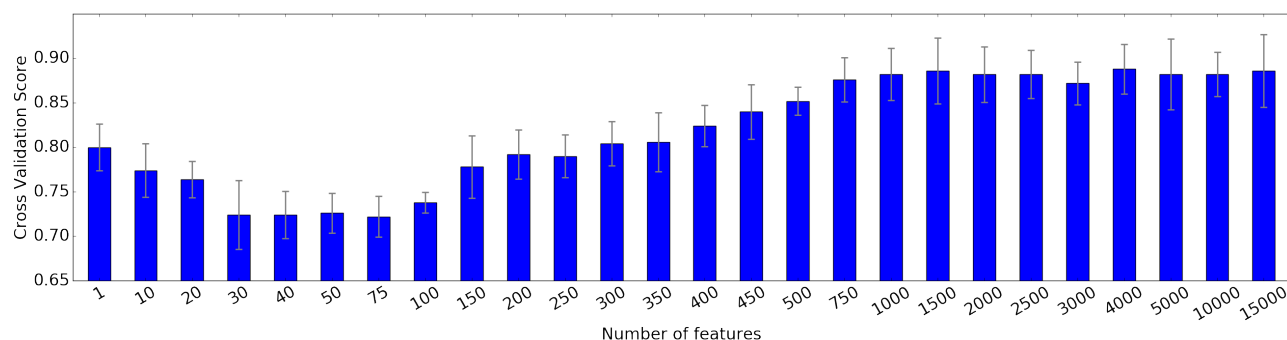


Figure 1 shows the mean 5-fold cross-validation scores and standard deviations for varying n . The peak performance is achieved with $n = 1000$. However its worth observing other values have margins of error that fall within the same score, also values between the shown bars were tested - the displayed bars are chosen for clarity. In choosing $n = 1000$ the resultant classifier had a mean scores of 89% with a standard deviation of 3%, making it better than the majority class baseline by 9%.

2 Improving the Model

2.1 Smarter Preprocessing

Several additional preprocessing techniques were implemented, the following paragraphs describe these techniques and give the associated accuracy values derived from solely applying that technique in conjunction with the previously detailed classifier of section 1.3.

An initial flaw to observe about the simple preprocessing procedure described in section 1.1 is that the filtering of non purely-alphabetical substrings delimited by whitespace discriminates against words that appear next to punctuation, such as “hello,” or “hello.”. This is counter-acted by tokenizing the input string through matching to the regular expression “ $\backslash b \backslash w + \backslash b$ ”, where $\backslash b$ matches to a word-boundary and $\backslash w$ matches to a alphabetical character. This change results in an improved score of 91% accuracy.

Another approach was to extract the content of the emails from the formatting of the emails. The emails were parsed to separate the headings and payload of the email, information from the subject header and the recipients fields was striped and concatenated to the content in the payload. Applying this technique in isolation gave an accuracy of 93%.

The third technique was numeric substitution, whereby numerical or alpha-numerical substrings would be replaced with the words “num” or “alphanumeric”. The motivation for this was that numbers would often appear in emails, but the actual values in question wouldn’t actually be semantically relevant. The hope of this substitution was to allow the system to better generalise. This gave an accuracy of 93%.

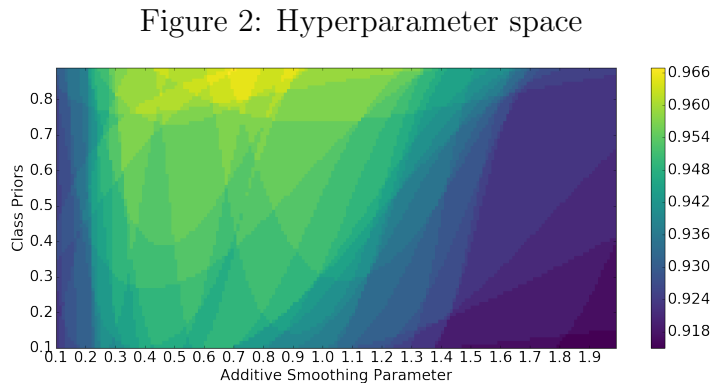
Similarly to numerical substitution was symbol processing - in most cases symbols were simply filtered out, however for symbols that provide meaning in their given contexts would be replaced with appropriate words, for example “\$50” would be replaced by “money 50” and “cat & dog ” would be replaced by “cat and dog”. Substitutions would also be made when symbols appeared in sucession, such as “free!!!” being replaced with “free multibang ”. Symbolic substitutions applied in isolation resulted in an accuracy of 92%.

The final technique was to process letter cases, any all-caps words such as “OFFER” were mapped to “allcaps offer”, and everything else was converted to lowercase. This gave an accuracy of 91%.

The use of all of these methods in conjunction resulted in a mean 5-fold cross-validation score of 94% with a 3% standard deviation - a 5% improvement from the simple preprocessing and a 14% improvement from the baseline.

2.2 Optimising Classifier Hyperparameters

In section 1.1 the hyperparameters α and p of the Naive Bayes classifier were described. Figure 2 visualises the three dimensional mapping $(\alpha, p) \mapsto s$ where s is the mean 5-fold cross-validation score derived from the classifier constructed from the given parameters. The maxima on this surface is the point (0.7, 0.85), where the constructed classifier has an accuracy of 97% with a standard deviation of 2%.



2.3 Comparing to Another Classifier

3 Extending the Model

3.1 Calibration of Naive Bayes Probabilities

The Naive Bayes classifier assumes that within each class, each of the feature values is independent of one another. Though this assumption is often incorrect, the classifier can still yield accurate results. However, the dependence of features within a class can lead to inaccurate probability estimates.

A classifier can be regarded as well calibrated if the probability of a prediction can give an indication of confidence of that prediction. That is, if the Naive Bayes classifier is well calibrated and we select only the instances which have a probability of around 0.7, then around 70% of that sample should belong to the positive class.

Figure 3 below shows the uncalibrated probabilities for the spam class. Each data point indicates the fraction of class which is within ± 0.1 of the class probability value. The dashed lines indicates the ideal trend line for perfectly calibrated probabilities, where at this point the probability exactly reflects the prediction.

Figure 3: Before Calibration

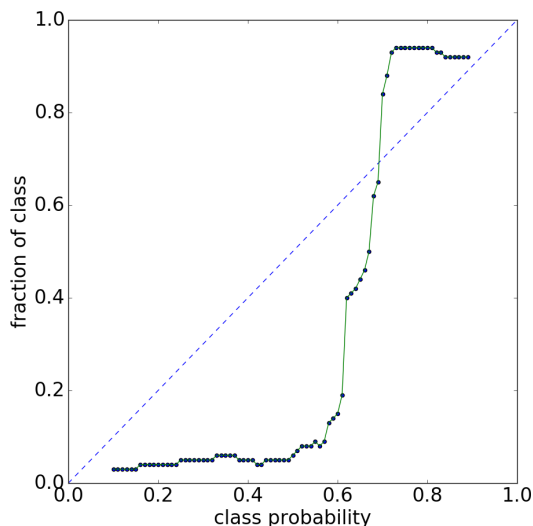
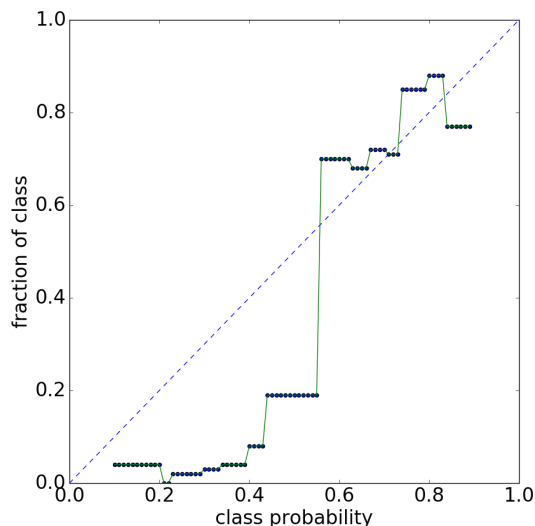


Figure 4: After Calibration



It can be observed that as the class probability tends from 0.0 to 0.5 that the reliability of the prediction also decreases. As the probability tends from 0.5 to 1.0 the classifier appears to tend towards perfect calibration before beginning to under-predict the proportion of the positive class.

To provide a quantification of the level of calibration the Brier score was evaluated. The Brier score metric measures the mean square difference between the predicted probability assigned to the possible outcomes for item i and the actual outcome o_i . In short, the lower a Brier score for a set of predictions, the better calibrated a classifier is.

We first establish a type of baseline for the Brier score to compare our results with. It can be observed that if the classifier always predicted the probability of spam as 50% then the Brier score loss would be 0.25 irrespective of the prediction as $(0.50 - 1)^2 = (0.50 - 0)^2 = 0.25$.

Evaluating the Brier score on the uncalibrated classifier yielded a value of 0.17. This value indicates that the confidence taken from a prediction is better than random, though can still be improved.

Isotonic regression was utilised to achieve a better confidence of a given prediction. Isotonic regression finds a non-decreasing approximation of a function while minimising the mean squared error on the training data. After fitting the model using Isotonic calibration the Brier loss score decreased to 0.15, a small improvement. The marginal improvement can be observed in figure 4, it can be seen that many of the points with a class probability of around 0.5 have move closer to the perfect calibration line, as have the majority of the points with class probabilities of above 0.5. It can be seen that some of the points have been moved further from the line of perfect calibration(points with class probability of curly equals 0.9). This could be explained by the points where the trend becomes non-monotonic.

3.2 Weighted Naive Bayes

Probability values were assigned to every word by dividing the total number of times the word appears across all the emails by the number of emails. More precisely for emails $e \in E$, let $f(e, w)$ detect the presence of a word in an email and $P(w | E)$ be the probability of a word appearing in an email.

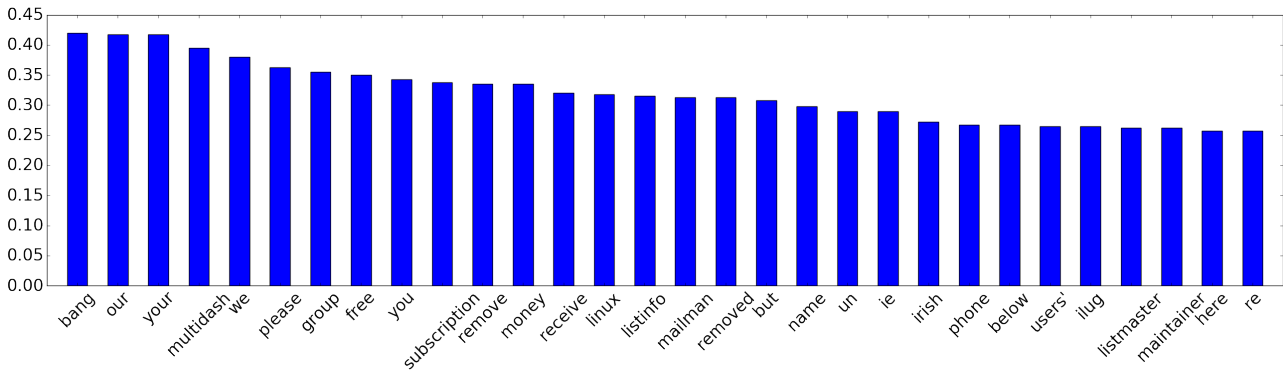
$$f(e, w) = \begin{cases} 1 & \text{if } w \in e \\ 0 & \text{otherwise} \end{cases} \quad \therefore P(w | E) = \frac{1}{|E|} \sum_{e \in E} f(e, w)$$

For the sets S and H of spam and ham emails a metric d was devised to measure how well a given word differentiates between the classes:

$$d(w) = \max(|P(w | S)|, |P(w | H)|)$$

Figure 5 plots this metric for the top thirty words across all the emails.

Figure 5: Top 30 most differentiating words



This mapping d provides a weighting to the Naive Bayes by multiplication with associated elements of the estimator vectors, following from equation 1 in 1.1,

$$\hat{\theta}_i = d(w_i) \frac{x_i + \alpha}{N + \alpha d}$$

$$\log \hat{\theta}_i = \log d(w_i) + \log(x_i + \alpha) - \log(N + \alpha d)$$

Therefore the weighted model is simply implemented by adding together the already existing estimator vectors and the logarithm of the weightings vector. In the actual implementation however the expression $d(w_i)$ was replaced with $d(w_i) + 1$ to avoid taking the logarithm of zero. The result of assessing this classifier by the mean 5-fold cross validation score is an accuracy of 98% with a standard deviation of 1%.