

The Best Way to Instil Confidence is by Being Right

An Evaluation of the Effectiveness of Case-Based Explanations in providing User Confidence

Conor Nugent, Pádraig Cunningham, and Dónal Doyle

Department of Computer Science
Trinity College Dublin
Conor.Nugent@cs.tcd.ie
Padraig.Cunningham@cs.tcd.ie
Donal.Doyle@cs.tcd.ie

Abstract. Instilling confidence in the abilities of machine learning systems in end-users is seen as critical to their success in real world problems. One way in which this can be achieved is by providing users with interpretable explanations of the system's predictions. CBR systems have long been understood to have an inherent transparency that has particular advantages for explanations compared with other machine learning techniques. However simply supplying the most similar case is often not enough. In this paper we present a framework for providing interpretable explanations of CBR systems which includes dynamically created discursive texts explaining the feature-value relationships and a measure of confidence of the CBR systems prediction being correct. We also present the results of a preliminary user evaluation we have carried out on the framework. It is clear from this evaluation that being right is important. It appears that caveats and notes of caution when the system is uncertain damage user confidence.

1 Introduction

CBR systems have long been understood to have an inherent transparency that has particular advantages for explanations compared with other machine learning techniques [1]. The realisation that there is a need to make machine learning systems more interpretable and user friendly has brought this fact back into focus in recent years. Research by Cunningham et al. found CBR explanations where the user is simply supplied with the most similar case are more convincing than rule-based explanations in some domains [2].

Recently researchers have begun to look at ways in which this method can be improved upon. The issue with case-based explanations lies in the perceived appropriateness of the presented cases to the validity of the prediction. This is an issue that has received a lot of attention in the CBR community. In CBR explanations, the ability of the user to make meaningful comparisons between the

query and the retrieved explanation case is of critical importance to the success of the explanation [3]. CBR systems are not wholly transparent and much domain knowledge can be contained within the similarity metrics used in the system. It is implicitly assumed in simple CBR explanations systems that the user has this same domain knowledge and so the appropriateness of the explanation case is clear. However, this may not be the case and the relevance of the retrieved case may be lost on novice users. This is an issue that McSherry has addressed in his ProCon System [4]. McSherry has focused on making the relationship between the feature values within a case and its predicted value explicit. Similarly we address this issue in our case-based explanation system for black-box systems [5]. However in our approach we used localised information to ensure that our system captured any non-linear feature interactions that occurred in the feature space.

In other work, Doyle et al. have focused on the observation that the nearest retrieved case in a CBR system may not be the best case to present as an explanation [6]. They use these cases to form *a fortiori* arguments in favour of the CBR systems prediction. They argue that in classification tasks, cases that are between the query case and the decision boundary, provide more convincing explanations. That is, cases that are more marginal on the important criteria are more convincing. With such cases the user is better able to assess whether the classification of the target case is justified.

The primary motivation in providing users of CBR systems with interpretable explanations is to increase their confidence in the system. However, as is pointed out by Cheetham and Price people can quickly lose confidence in a system if it makes predictions which then turn out to be incorrect [7]. To address this issue Cheetham and Price propose using confidence measures so as to alert the user of when a system may be making a mistake.

We have developed an explanation framework for CBR systems which attempts to address the issue of providing user confidence by providing interpretable explanations coupled with a measure of confidence of the systems prediction. We have performed preliminary evaluations on the explanation framework and the results are presented.

The paper is structured as follows. Section 2 outlines how the framework works. Section 3 describes the evaluation we have carried out and presents the results of those evaluations. Finally we end with the conclusions in Section 4.

2 Explanation Framework

We have developed a framework for providing interpretable explanations in CBR systems. The explanations produced by the framework contain a number of elements;

- Cases that form *a fortiori* arguments
- Discursive text describing the effects of differences feature-values between the Query Case and the Explanation Case.

- A measure of confidence in the system’s prediction

The framework expands on earlier work in which we used localised models to help explain the feature-value relationships in regression tasks [5]. The two key aspects of our localised approach are; the generation of a local case-base and the use of a local model. The local model is used to help describe the feature-value relationships and to inform the search for an explanation case. To build a local case-base we simply use a Nearest Neighbour algorithm to create a subset case-base of the original case-base. First we find the Query Case’s nearest neighbours and include them in our new subset case-base until we have at least K cases of each class. This ensures that our local case-base traverses the decision boundary in the area of our Query case. Once we have our localised case-base we then build our local model on it. As a model to use to capture the local information stored in the casebase we have selected to use logistic regression model. Logistic regression models are quite simple yet powerful and allow us to realise all the elements of our explanation framework listed above. In the coming sections we discuss the logistic regression model and how it is used in the generation of explanations.

2.1 Logistic Regression

Logistic regression, like linear regression, produces a set of coefficients from which the relationship of an input variable to the target class variable can be deduced. However unlike linear regression, logistic regression coefficients don’t directly correspond to slope values in the same way. Logistic regression models are restricted to binary tasks and the two possible class values are coded as being either 0 or 1. Because the value predicted by the model, the conditional mean, is no longer an unbounded value as in linear regression but a value between 0 and 1, the data is fitted to a distribution that ensures the outputted value always meets this bounding criteria. To do this the logistic distribution is applied as can be seen below (1).

$$Y(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} \quad (1)$$

Here $Y(x)$ is the conditional mean for a particular value of x while β_0 and β_1 are the model parameters. The distribution produces the conditional mean, a value between 0 and 1, for any given inputted value of x . Importantly, for binary problems the conditional mean is in fact the probability of class 1 given x .

At first glance this model looks quite intimidating and seems to offer no hope of offering an insight into the relationship between x and our class variable. However, the logistic distribution is chosen because it can be easily transformed into another form which has many of the desirable properties of a linear regression model. By applying the logit transform, equation 2, we end up with a simple and interpretable model, the logit (3).

$$g(x) = \ln \frac{Y(x)}{1 - Y(x)} \quad (2)$$

$$g(x) = \beta_0 + \beta_1 x \quad (3)$$

The parameters of the logit model can easily be converted into odds ratios. The odds ratio of an event is the odds of that event occurring over the odds of it not happening. For instance if someone were to state the odds ratio of smokers to non-smokers getting cancer is 2 then this would mean smokers are twice as likely to develop cancer as non-smokers. Alternatively, if we looked at the relationship the other way round, non-smokers to smokers, we would get a odds ratio of 0.5. This means that non-smokers are half as likely to get cancer. In general an odds ratio greater than one for possibility A over possibility B means A makes the event more likely than the alternative while an odds ratio of less than one means it makes it less likely. The logistic regression model makes the calculation of odds ratios quite easy and this is extremely useful and informative. It is this simple relationship between the model coefficients and the odds ratio and their natural interpretation that has made logistic regression such a popular tool. We will discuss in a very general sense how this is done as it will be of use in section 2.3 where we use the logistic regression model to explain the differences in feature-values between the query case and the explanation case.

In order to extract the odds ratio, two steps are taken. First the logit difference is found. Imagine we are interested in the odds ratio of two different events, $x = c$ and $x = d$. the logit difference can be calculated as in equation 4. The logit difference, ld , is simply the difference in the logit function for the two values of x we are interested in. Once this value has been obtained it can then be converted into an odds ratio, see equation 5.

$$LogitDifference(x = c, x = d) = g(c) - g(d) = ld \quad (4)$$

$$OddsRatio(x = c, x = d) = e^{ld} \quad (5)$$

The trick with the logistic regression model is that in many cases it isn't necessary to calculate the logit difference. If the model variables have been properly coded then the desired information usually can be got by simply looking at the model coefficients. As an example consider a hypothetical situation where we have developed a model that relates smoking to the development of cancer. Our hypothetical model might look something like that shown in equation 6.

$$g(Smoker) = 0.3 + 0.69Smoker \quad (6)$$

If we code our smoking variable as being equal to 1 if someone smokes and 0 if they don't then the calculation of the logit difference is simply equal to the *Smoker* coefficient (7).

$$\begin{aligned} g(Smoker = 1) - g(Smoker = 0) &= 0.3 + 0.69(1) - (0.3 + 0.69(0)) \\ &= 0.69 \\ OddsRatio(Smoker) &= e^{0.69} \end{aligned} \quad (7)$$

$$OddsRatio(Smoker) = 2$$

As can be seen above in 6 we need not have bothered calculating the logit difference and instead just used the model coefficient. This is also true for continuous and multi-value nominal variables if they are coded correctly (chapter 4, [8]). Once we have the odds ratio the relationship between input variable and the class variable is clear. We have focused most of our discussion on examples with only a single input variable for simplicity sake but the above observations are also true in multi-variable problems. In the next section we discuss how exactly information derived from the logistic regression model can be used to provide convincing explanations.

2.2 Finding A Fortiori Cases and a Measure of Confidence

Using the local logistic regression model we can generate *a fortiori* arguments dynamically and without any prior domain knowledge. As discussed in section 2.1 Logistic Regression models allow us to generate a probability for a given set of inputs. In the explanation case retrieval process we can then use this to find an explanation case that is nearer the decision boundary and so a more convincing argument. We consider each of the cases in our localised case-base as a candidate case for inclusion in the explanation. By passing each of our candidate explanation cases through our local logistic model we can generate a probability for each. A case that is nearer the decision boundary and of the same class as our CBR system has predicted will have a more marginal probability and so this should be the case we select.

After the Query Case has been classified we can then build our logistic model on our local data. In table 1 we can see the Query Case, its predicted classification and three candidate explanation cases which are in fact the Nearest Neighbours used to classify it. In order to select a case to use in our explanation we first run each of the cases including the Query Case through our local logistic regression model. This gives us the set of probabilities that can also be seen in table 1. We can see that Nearest Neighbour 2 has the lowest probability and so is the case nearest the decision boundary. This is an alternative to the explanation utility framework described by [6] for selecting the case to present to the user to make the most convincing argument. Although Nearest Neighbour 2 had consumed more units of alcohol and weighed less, they were under the limit so it seems reasonable that our Query Case should be too.

We can make this argument more explicit to the end user by explaining the effects of the feature differences between the Query Case and Explanation Case. In the next section we will outline how this can be done using the local logistic regression model. As Cheetham and Price point out being able to provide a measure of prediction confidence is an extremely useful asset in maintaining end-users confidence in a system [7]. Using the localised logistic regression models we have got a probability of a our CBR systems prediction being correct. If this probability is below a certain threshold we can inform the user that confidence is low. How this threshold might be decided upon is discussed in section 3.2.

Table 1. Explanation Case Retrieval Process

Features	Query Case	Nearest Neighbour 1	Nearest Neighbour 2	Nearest Neighbour 3
Weight	88	82	79	76
Duration	120	120	120	120
Gender	Male	Male	Male	Male
Meal	Full	Full	Full	Full
Units	5.2	5.0	7.2	4.6
BAC	Under	Under	Under	Under
Probability	0.98	0.97	0.89	0.96

2.3 Explaining Feature-Value Relationships

Using equations 4 and 5 from section 2.1 we can substitute each of the feature differences into the equations individually and get the odds ratio for each. Using the odds ratio we can then determine the effect of the change. The kind of dialog that can be produced can be seen in Table 2. In this sample explanation we can see the advantage of using our local model to classify the Query Case as this gives us a measure of confidence in the prediction.

Table 2. Sample Explanation

	Query Case	Explanation Case
Weight (kgs)	57.0	79.0
Duration (mins)	240.0	240.0
Gender	Male	Male
Meal	Full	Full
Amount (Units)	12.6	9.6
BAC		Over
The prediction for the individual in the Query Case is: Over the limit		
The confidence that this prediction is correct is: high		
Discursive Text:		
In support of this prediction we have the person presented by the Explanation Case who was also Over the limit. Weight being higher and Amount being bigger have the effect of making the Query individual more likely to be Over the limit than the Explanation individual.		

In our second example (table 3) the confidence measure is low and so the explanation is adjusted so as to include a counter example. The counter example selected is the case of the other classification from the local case-base that is nearest the decision boundary. This is intended to assist the end user in determining whether the prediction might be correct. Again the local logistic regression model is used to explain the differences in the feature values. It is worth noting that if the case-base used to build the local model doesn't adequately represent the problem counter intuitive explanations can be produced. For instance we found that if too few cases were used **duration** could be heavily correlated with **units** and so a larger **duration** value could be seen as evidence in favour of being over the limit.

Table 3. Sample Explanation with Counter Example

	Explanation Case	Query Case	Counter Example
Weight (kgs)	52.0	53.0	73.0
Duration (mins)	270.0	330.0	210.0
Gender	Male	Male	
Meal	Lunch	Lunch	Lunch
Amount (Units)	9.1	10.4	9.0
BAC	Over		Under
The prediction for the individual in the Query Case is: Over the limit			
The confidence that this prediction is correct is: low			
Discursive Text:			
In support of this prediction we have the person represented by the Explanation Case who was also Over the limit. Gender being Female and Amount being bigger have the effect of making the Query individual more likely to be Over the limit than the Explanation individual. However, Weight being heavier and Duration being longer have the effect of making the Query individual less likely to be Over the limit than the Explanation individual			
As there is low confidence in the prediction we also have a counter example of someone who is similar but Under the limit for you to inspect			
Duration being longer has the effect of making the Query individual more likely to be Under the limit than the counter example individual. However, Weight being lighter, Gender being Female and Amount being bigger have the effect of making the Query individual less likely to be Under the limit than the counter example individual			

3 Evaluation

In this section we examine the effectiveness of the explanation framework. There were two principle aspects of the framework which we wished to assess. Firstly the usefulness of the explanations and secondly how effective the framework is at predicting confidence. In order to assess the usefulness of the framework’s explanations we performed a user trial. The effectiveness of the confidence measure was assessed on a number of different data sets using standard machine learning techniques. We will discuss the details of each evaluation in turn.

3.1 User Trials

In designing the user trial there were three principle questions we wished to address; do people find the explanations understandable and useful, do the explanations increase users’ confidence in the case-based system and finally can the explanations alert users to when the system might be in error. The case-base on which the trial was carried out was the blood alcohol case-base [2, 6]. The task involves using information about peoples weight, gender, number of units of alcohol consumed, etc. to predict whether someone’s blood alcohol content (BAC) exceeds the drink driving limit. The full set of features used can be seen in Table 4. We built a simple Nearest Neighbour algorithm on the data set and applied our framework to providing explanations of it’s predictions.

Table 4. The features in the BAC data set

Weight (Kg)	Duration (Time Spent Drinking)
Meal (None, Snack, Lunch, Full)	Amount (In Units)
Gender	BAC (Blood Alcohol Content)

In the trial subjects were given a questionnaire in which they were shown three different forms of explanation; that given by the framework,

- **The Full Framework Explanation:** This is an explanation that includes the selected *a fortiori* explanation case, a discursive text and a measure of confidence as seen in table 2.
- **Case-based Explanation:** In this form of explanation the subject is just shown the selected *a fortiori* case as evidence in favour of the prediction.
- **No Explanation:** The user is just presented with the feature-values of the query and the systems prediction.

The trial subjects were shown four examples of each type of explanation and asked three questions after each example shown;

- **Question One:** Do you think the prediction is correct?

- **Question Two:** How would you rate this Explanation?
- **Question Three:** Did the explanation help you in answering question one?

Below each question the trial subject had five options to select from. In both question one and three the options were; No, Maybe No, Don't Know, Maybe Yes and Yes. In question two the options were; Poor; Fair; Okay, Good and Very Good.

To assess the use of explanations in terms of alerting users to when the system might be in error one of the four examples shown of each explanation type was a mis-classification. Twelve people from a number of different backgrounds took part in the evaluation and the results are discussed in the next section.

User Trial Results: In question one we looked at the frequencies with which users chose each of the five options when the prediction made by the system was correct. These can be seen in figure 1. It is clear that the explanations given by the framework give the users far greater confidence in the system than either of the other two schemes. The trial subjects answered *Yes* 88% for the time with just four answers being anything other than yes. Three people answered *Maybe Yes*, one *Don't Know* and there were no negative answers. We also examined the users

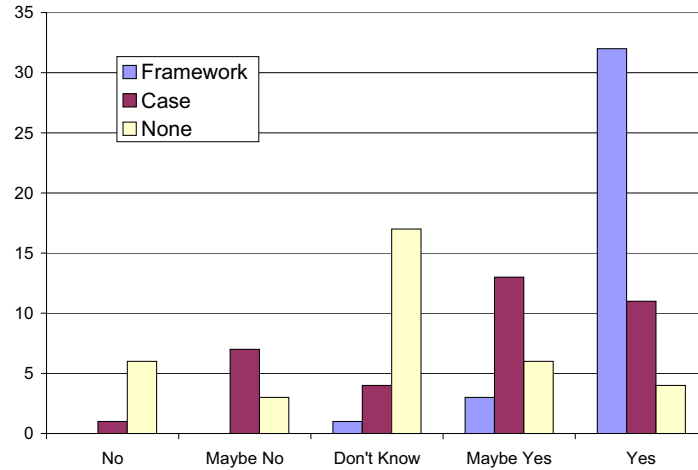


Fig. 1. The distribution of user responses when the system predictions were correct

responses when the system had made an incorrect classification and the results can be seen in figure 2. The graph of frequencies reveals a very different user response pattern. Although no one responded *Yes* in the case of the explanations produced by the framework there is far less certainty in the users responses.

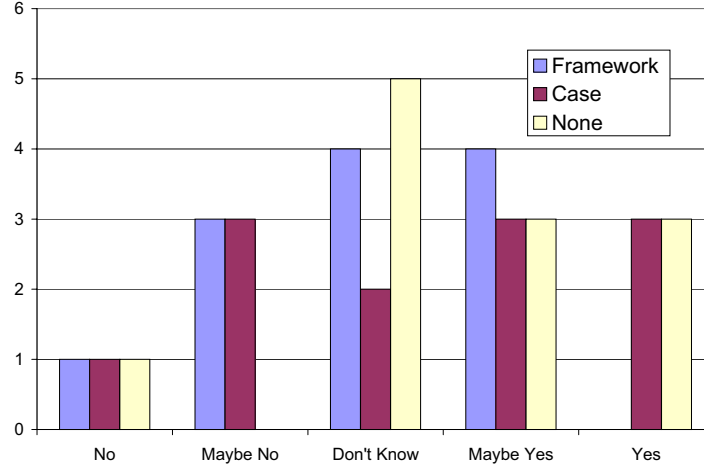


Fig. 2. The distribution of user responses when the system predictions were wrong

In question two we were trying to determine how satisfactory people found the explanations. We coded the trial subjects responses as being a number between one and five. One being *fair* and five being *Very Good*. We then looked at the average value given to each explanation for each scheme. The results are shown in figure 3. Clearly people found the framework explanations to be far more satisfying then the other two schemes and generally the level rating for the framework explanation was quite high.

Finally in question three we were interested in the difference in behaviour when the system was correct and when it was incorrect. We wanted to see how useful users found the explanations in these two situations. We coded the responses as before and we can see that again generally the rating for the framework is quite high (figure 4). However, it dips considerable in the case of the system being incorrect. It is clear that in these circumstances users confidence in the system has been damaged. From comments returned by test subjects the addition of a counter example at times of uncertainty led to confusing explanations. We would like to do a further survey to investigate this matter in greater detail. It can be seen that generally the framework explanations added to users confidence in the system's predictions however user confidence was damaged when the system made errors.

3.2 Confidence Measure Evaluation

We evaluated our confidence measure scheme on two data sets; the BAC set and a Spam data set from the UCI repository. The key aspect in providing any confidence measure is ensuring that when it is confident it is correct while also

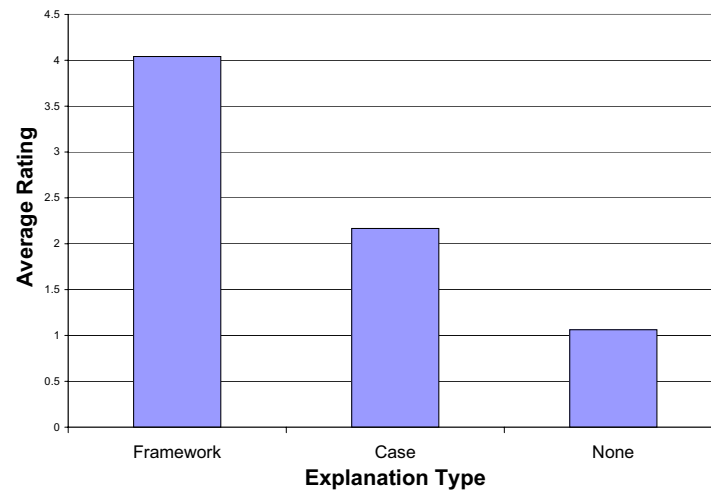


Fig. 3. The Average Rating Scores for Question Two of the Explanations Produced by Each Different Scheme

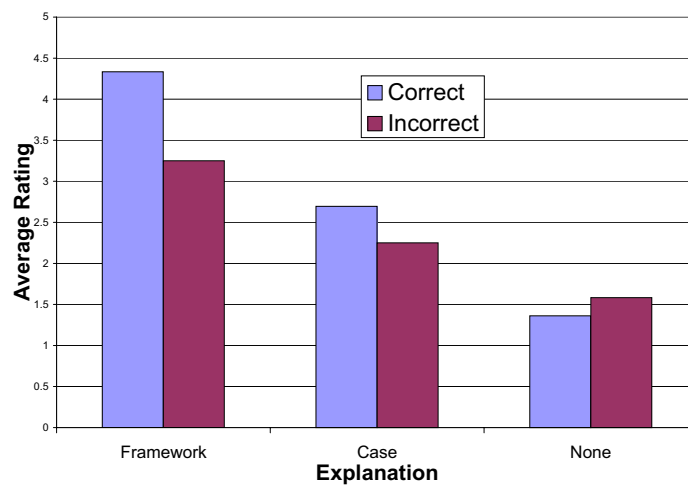


Fig. 4. The Average Rating Scores for Question Three of the Explanations Produced by Each Different Scheme

not being overly pessimistic and saying a lot of correct predictions are incorrect. Constantly supplying users predictions that we are unsure about is bound to damage their confidence in the system. There often is a trade-off between the two and a tolerance level where the level of confidence versus pessimism is acceptable must be chosen. This can make comparing different schemes less than straightforward as one scheme may be better at one level of tolerance and another at a different level. The characteristics of this problem led us to investigate adapting ROC curves to the task [9]. We can characterize our wish for accurate confidence as being our Confident Correct Rate (CCR) as defined in Equation 8. Likewise we can encapsulate our need to minimise pessimism in the Not Confident Correct Rate ($NCCR$) as defined in Equation 9.

$$CCR = \frac{CC}{CC + CI} \quad (8)$$

$$NCCR = \frac{NCC}{NCC + NCI} \quad (9)$$

Where CC is the number of times the measure is confident and the system is correct and CI is the number of times measure is confident and the system is incorrect. Likewise NCC is the number of times the measure is not confident and the system is correct and NCI is the number of times the system is not confident is right to be so. To make the definition of these parameters a little clearer we have displayed them in the form of truth table in Table 5. Our scheme

Table 5. A Truth Table Defining the Equation Parameters

	Incorrect	Correct
Confident	CI	CC
Not Confident	NCI	NCC

for confidence requires one parameter K , the number cases of each type of class value that is required in order to stop the local case-base building process. In our confidence scheme we must chose a level of probability that we must have in a prediction in order to be confident in it. We performed leave-one-out cross-validations on both data sets recording the required statistics while both varying K and the confidence threshold. We then plotted the results of the evaluation on Characteristic Confidence Curves which are very similar to ROC curves as can be seen in figures 5 and 6. For each scheme there is a separate curve and the points on those curves represent different threshold levels for those schemes. Like in ROC curves our ideal solution would lie in the top left hand corner and the solution which is nearest this point optimises the trade-off. However different applications may have restrictions about how often the system can be confident and incorrect. It is quite easy using the characteristic curves to find the scheme that best meets these requirements. It is also possible eliminate certain schemes

as being definitely worse than another (like in ROC curves) if the curve of one scheme lies entirely inside another then it is worse than that scheme.

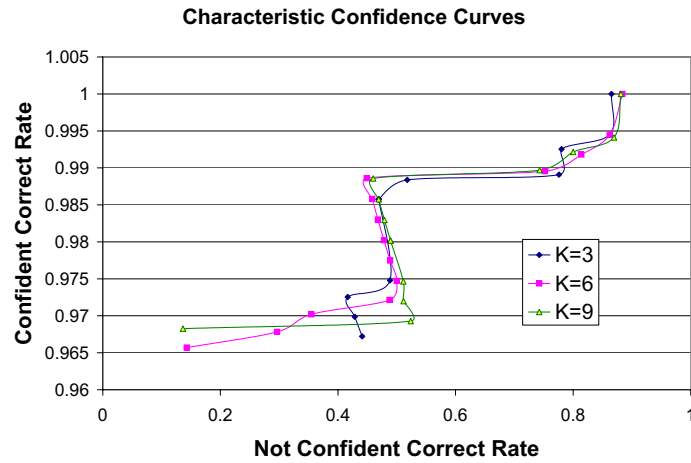


Fig. 5. The Characteristic Confidence Curves for the UCI Spam Data set for a Range of K Values

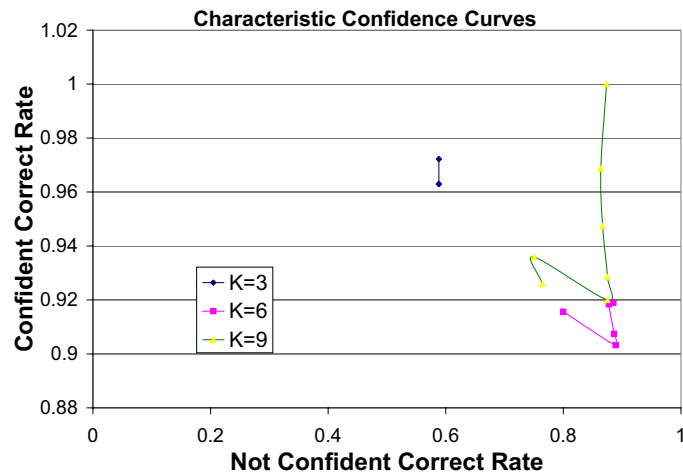


Fig. 6. The Characteristic Confidence Curves for the BAC Data set for a Range of K Values

In figure 5 we can see that the three different schemes are all quite closely aligned but that generally the scheme for $k=6$ out performs the others although at certain points $k=3$ is slightly more favourable. Likewise in figure 6 we can see that our best solution is clearly when $k=3$ as the two points on its curve lie far closer to the upper left hand corner than any others. The reason that there are only two points on the $k=3$ curve is that it very quickly goes from being entirely not confident to reaching the minimum possible threshold value. The minimum threshold value is the probability of 0.5 as any believe below this actually represents a believe in the other class in binary problems. As an example of how accurately we can predict confidence we chose the two points on both graphs that maximised the trade-off. These can be seen in table 6.

Table 6. The Confidence Measures Results

	Spam		BAC	
	Incorrect	Correct	Incorrect	Correct
Confident	13	366	3	78
Not Confident	18	3	7	10

In the case of the Spam data set we are Confident and Correct (CC) 91.5% of the time while being Confident and Incorrect (CI) just 3.25% of the time. Importantly we are not confident when correct less than 1% of the time. In the alcohol data set CC 79% of the time while CI 3% of the time. If the amount of CI predictions is of critical importance then the axis of the graphs can be weighted appropriately.

4 Conclusions

In this paper we have addressed the issue of instilling confidence in the ability of machine learning systems in the users. We have developed an explanation framework which supplies users with interpretable explanations of the systems predictions along with a measure of confidence in that prediction. We have also presented a means by which the trade-off between being overly confident or overly pessimistic can be inspected and different methods compared.

We carried out a preliminary evaluation on the explanation framework and have found that the use of interpretable explanations does indeed increase confidence in the system as can be seen in figure 1. The addition of discursive text explaining the relationship between the presented explanation and the query cases clearly had an effect in evoking this confidence as can be seen in the satisfaction ratings shown in figure 3. However, when the system fails this confidence can be damaged. This can be clearly seen in figure 2 as the users display far less certainty about the system prediction compared with when the system is correct. This is coupled with a drop in satisfaction in the ability of the explanation to inform the user's opinion of whether the system is correct or not (see

figure 4). This could be a result of the extra cognitive load associated with the explanations produced when the level of confidence is low. However users were still unable to reliably perceive that the system was making an error and so their confidence in the system could be lost when the resulting error becomes evident. Clearly notifying the user of uncertainty in the recommendation from the system creates an element of doubt and confidence is damaged. The use of an explanation including a counter example does not seem to make clearer what the correct prediction might be.

This a matter that has only been touched on in our preliminary investigation and it is one which we would like to address further in a more comprehensive study. In the future we would also like to investigate localised logistic regression as a CBR classification technique as well as find improved means by which we can generate local case-bases.

References

1. Leake, D.: Case-Based Reasoning: Experiences, Lessons and Future Directions. AAAI/MIT Press (1996)
2. Cunningham, P., Doyle, D., Loughrey, J.: An evaluation of the usefulness of case-based explanation. In Ashley, K.D., Bridge, D.G., eds.: Case-Based Reasoning Research and Development, 5th International Conference on Case-Based Reasoning (ICCBR 2003). Volume 2689 of Lecture Notes in Computer Science., Springer (2003) 122–130
3. Sormo, F., Cassens, J.: Explanation goals in case-based reasoning. In: 1st Workshop on Case-Based Explanation (Proceedings of the ECCBR 2004 workshops). (2004) 165–174
4. McSherry, D.: Explanation in case-based reasoning: an evidential approach. In: 8th UK Workshop on Case-Based Reasoning. (2003) 47–55
5. Nugent, C., Cunningham, P.: A case-based explanation system for 'black-box' systems. In: 1st Workshop on Case-Based Explanation (Proceedings of the ECCBR 2004 workshops). (2004) 155–164
6. Doyle, D., Cunningham, P., Bridge, D., Rahman, Y.: Explanation oriented retrieval. In Funk, P., Calero, P.A.G., eds.: Advances in Case-Based Reasoning, 7th. European Conference on Case-Based Reasoning (ECCBR 2004). Volume 3155 of Lecture Notes in Computer Science., Springer (2004) 157–168
7. Cheetham, W., Price, J.: Measures of solution accuracy in case-based reasoning systems. In Funk, P., Calero, P.A.G., eds.: Advances in Case-Based Reasoning, 7th. European Conference on Case-Based Reasoning (ECCBR 2004). Volume 3155 of Lecture Notes in Computer Science., Springer (2004) 106–118
8. Hosmer, D., Lemeshow, S.: Applied Logistic Regression. 2nd edn. Wiley (2000)
9. Flach, P., Blockeel, H., Ferri, C., Hernandez-Orallo, J., Struyf, J. In: Decision support for data mining: introduction to ROC analysis and its application. Kluwer Academic Publishers (2003) 81–90