# Interacting Meaningfully with Machine Learning Systems: Three Experiments[1]

*Simone Stumpf, Vidya Rajaram, Lida Li, Weng-Keen Wong,*
*Margaret Burnett, Thomas Dietterich, Erin Sullivan, Jonathan Herlocker*

Oregon State University
School of Electrical Engineering and Computer Science
Corvallis, OR 97331 USA
1-541-737-3617
{stumpf,rajaramv,lili,wong,burnett,tgd,sullivae,herlock}@eecs.oregonstate.edu

## ABSTRACT

Although machine learning is becoming commonly used in today's software, there has been little research into how end users might interact with machine learning systems, beyond communicating simple "right/wrong" judgments. If the users themselves could somehow work hand-in-hand with machine learning systems, the accuracy of learning systems could be improved and the users' understanding and trust of the system could improve as well. We conducted three experiments to begin to understand the potential for rich interactions between users and machine learning systems. The first experiment was a think-aloud study, aiming to see how willing users were to interact with and about machine learning reasoning, and to help us understand what kinds of feedback users might give to machine learning systems. Specifically, users were shown explanations of machine learning predictions and asked to provide feedback to improve the predictions. The results were that users' feedback was rich, complex, and widely varied, ranging from suggestions for reweighting of features to proposals for new features, feature combinations, relational features, and wholesale changes to the learning algorithm. We then investigated the viability of introducing such feedback into machine learning systems: specifically, how to incorporate some of these types of user feedback into machine learning systems, and impact on the accuracy of the system. Taken together, the results of our experiments show that supporting rich interactions between users and machine learning systems is feasible for both user and machine. This shows the potential of rich human-computer collaboration via on-the-spot interactions as a promising direction for machine learning systems to work more intelligently, hand-in-hand with the user.

## 1. INTRODUCTION

A new style of human-computer interaction is emerging, in which some reasoning and intelligence reside in the computer itself. These intelligent systems and user interfaces attempt to adapt to their users' needs, to incorporate knowledge of the individual's preferences, and to assist in making appropriate decisions based on the user's data history. These approaches use artificial intelligence to support the system's part of the reasoning. One increasingly common approach being brought to intelligent systems and user interfaces is machine learning, in which the system learns new behaviors by examining usage data.

Traditionally, machine learning systems have been designed and implemented off-line by experts and then deployed. Recently however, it has become feasible to allow the systems to continue to adapt to end users by learning from their behavior after deployment. Interactive email spam filters, such as in Apple's Mail system, are prime examples.

Although machine learning is often reasonably reliable, it is rarely completely correct. One factor is that statistical methods require many training instances before they can reliably learn user behavior. Sometimes correctness is not critical. For example, a spam filter that successfully collects 90% of dangerous, virus-infested spam leaves the user in a far better situation than having no spam filter at all. But sometimes correctness is important. For example, recommender systems that recommend substandard suppliers or incorrect parts, language translators that translate incorrectly, decision support systems that lead the user to overlook important factors, and even email classifier algorithms that misfile important messages could

---

[1] An early presentation of Experiment #1 appeared in Stumpf et al, 2007.

cause significant losses to their users and raise significant liability issues for businesses. Further, too much inaccuracy in "intelligent" systems erode users' trust.

When accuracy matters, allowing the user to help could make a crucial difference. Therefore, approaches have begun to emerge in which the user and the system *interact with each other*, not just to accomplish the goal of the moment, but also *to improve the system's accuracy* in its services to the user over the longer term.

This direction is still in its infancy. The norm for the few machine learning systems that communicate with users at all is to allow the user to indicate only that a prediction was wrong or to specify what the correct prediction should have been. This is just a glimpse of the rich knowledge users have about the correct prediction. We began to consider whether end users might be able to provide rich guidance to machine learning systems. We wondered whether enabling them to provide this guidance could substantially improve the speed and accuracy of these systems, especially early on in training when machine learning does not possess a lot of knowledge.

Many questions arise from this possibility. Will users be interested in providing rich feedback to machine learning systems? If so, what kind of advice will they give? Will their feedback be usable by machine learning algorithms? If so, *how* would the algorithms need to be changed so that they could make use of this feedback? In this paper, we begin an exploration of these questions. The overall premise we explore is that, if the machine learning system could explain its reasoning more fully to the user, perhaps the user would, in return, specify *why* the prediction was wrong and provide other, rich forms of feedback that could improve the accuracy of machine learning.

There are implications for both directions of the communication involved in this premise. First, the system's explanations of why it has made a prediction must be usable and useful to the user. Second, the user's explanations of what was wrong (or right) about the system's reasoning must be usable and useful to the system. Both directions of communication must be viable for production/processing by both the system and the user.

To investigate possibilities for both directions of communication, we conducted three studies. The first was a formative think-aloud study with email users, in which machine learning algorithms sorted email messages into folders and explained their reasoning using three different explanation paradigms: Rule-based, Keyword-based, and Similarity-based. The participants were asked to provide feedback to improve the predictions. No restrictions were placed upon the form or content of participants' feedback. We then conducted two follow-up studies, in which we changed our machine learning algorithms to make use of some of what the users advised in the first study, evaluating the results. In both of these studies, we conducted off-line experiments to investigate *how* some of the user feedback could be incorporated into machine learning methods and to evaluate the *effectiveness* of doing so.

Thus, our research questions were:

RQ 1. *Is it possible* for machine learning systems to explain themselves such that (a) users can *understand* the system's reasoning, and (b) users are *willing* to provide the system rich, informative feedback with potential to improve the system's accuracy?

RQ 2. *What types* of feedback will users give? That is, how might we categorize the nature of their feedback from the perspective of machine learning, and what sources of background knowledge underlie the users' feedback?

RQ 3. Can these types of user feedback be assimilated by existing learning algorithms? If so, exactly *how* could some of these types of user feedback be incorporated into machine learning algorithms, and does doing so actually *improve* the performance of algorithms?

## 2. RELATED WORK

The most closely related work is research into users' perceptions and opinions of the algorithms used by intelligent interfaces. For example, in Pazzani's experiment (Pazzani 2000), users were asked which email learning system they trusted more to classify the email correctly as junk or not junk, given the choice between rules, signed-weighted keywords, and a new approach that used general descriptions employing keywords. This study did not include the understandability of machine-generated explanations or consider incorporating users' feedback into the system. A recent field study of eight non-technical users interacting with an intelligent system that predicted a person's interruptibility using sensor data (Tullio et al. 2007) is also relevant to our work. In this study, participants' early models of how the system worked, which varied greatly among the participants, were remarkably persistent even given counterevidence. Still, some of their misconceptions did give way to feedback contradicting their perceptions. This study did have a brief form of machine-generated feedback in one of the treatments, namely listing the top factors that impacted the system's decision. (Presence of this explanation had only

moderate impact on participants' mental models.)  There was no provision of user feedback being brought back into the system.

A first step in any effort to obtain user feedback about a learning system is to ensure that explanations by the learning system are understandable and useful to users.  Previous work  (Herlocker et al. 2000; Myers et al. 2006; Crawford et al. 2002a, 2002b) has shown that explanations that answer why certain outcomes happened, based on user actions, can contribute positively to system use.  Similarly, it has been shown that highlighting the relationship between user actions and ensuing predictions can influence user preference (Billsus et al. 2005). There is also previous research on the characteristics of explanations that help users choose between predictions. For example, showing contrasting features in recommendations can play a role in user trust (Herlocker et al. 2000; Pu and Chen 2006). One way that this relationship can be expressed is by making use of various ways of reasoning, such as analogical reasoning (Lieberman and Kumar 2005).

Different methods for gathering user feedback have also been investigated, along a spectrum of formality and richness. An obvious way to gather user feedback is to allow interactions in natural language (Blythe 2005). Semi-formal types of feedback that have been shown to be preferred by users make use of editing feature-value pairs (McCarthy et al. 2005; Chklovski et al. 2005). Other approaches allow the user to edit models produced by a learning algorithm using a formal description language (Oblinger et al. 2006). However, so far there has been a lack of research that integrates an investigation into the understanding of machine learning systems' explanations with an analysis of the *content* of the rich feedback users give when they have an unconstrained opportunity to do so.

The process of incorporating user feedback into a machine learning algorithm relies on both the type of the user feedback and the learning algorithm itself. One general approach is to treat user feedback as hard constraints to the algorithm. For instance, the constraints can enforce qualitative monotonicities (Altendorf et al. 2005), clamp labels in a Conditional Random Field (CRF) (Culotta et al. 2006) or fix parameters in a graph model (Huang and Mitchell 2006). For a support vector machine, which uses quadratic programming to find the maximum margin hyper-plane, it is natural to add the user-proposed constraints into the constraints of the optimization problem (Fung et al. 2002).  Other research uses the feedback to select features for the machine learning algorithm (Liu et al. 2004; Fails and Olsen 2003).  In Ware et al. (2001), the authors let the user directly build a decision tree for the data set with the help of visualization techniques. Our investigation considers two approaches to incorporate user feedback into machine learning and the effects on accuracy.

Our research takes place in the domain of email classification.  Although various supervised learning algorithms have been developed to automatically classify email messages into a set of categories or folders defined by users (e.g., Brutlag and Meek 2000; Cohen 1996; Shen et al. 2006), the reported accuracy of these algorithms indicates that email classification is a very challenging problem. The challenges stem from numerous factors, such as imbalanced categories, incomplete information in the email messages, and the fact that the categories (folders) set up by the users are often idiosyncratic and non-orthogonal. These challenges further motivated our interest in studying rich user feedback and its potential in improving machine learning algorithms.

## 3. MAIN EXPERIMENT: EXPLAINING SYSTEM BEHAVIOR AND GETTING FEEDBACK FROM USERS

Our main experiment investigated the first two of our research questions:

RQ 1. *Is it possible* for machine learning systems to explain themselves such that (a)  users can *understand* the system's reasoning, and (b) users are *willing* to provide the system rich, informative feedback with potential to improve the system's accuracy?

RQ 2. *What types* of feedback will users give? That is, how might we categorize the nature of their feedback from the perspective of machine learning, and what sources of background knowledge underlie the users' feedback?

To maximize external validity, it was important to base our experiment on real-world data. To allow as thorough investigation of users' potential as possible, it was also important to allow participants to express feedback freely. Thus, our first two design principles were:

*(Principle 1) Real-world email data*: 122 messages from a user's email (farmer-d), which had sufficient content for human and machine classification, were drawn from the publicly available Enron dataset (Klimt and Yang 2004). (Our data will be provided upon request.) The original Enron user had categorized these email messages into four email folders: Personal (36 messages), Resumé (27 messages), Bankrupt (23 messages) and Enron News (36 messages).

*(Principle 2) Rich collecting of result data:* We employed a qualitative "think-aloud" design in order to extract the richest

possible data from the participants. We observed and videotaped their activities and comments throughout the experiment, as well as collecting their work products.

As the first step of our procedure, learning algorithms classified each email message. Then, three explanations of each result were generated: a Rule-based, a Keyword-based, and a Similarity-based explanation. The application of the classification algorithms and the generation of explanations, described in the next section, were all done off-line prior to the experiment. We used the outcomes of these two steps to make possible the participants' interactions with our low-fidelity prototype.

Low-fidelity prototypes are important for experiments aiming to encourage participant feedback, because they avoid the impression of a "finished" product (Rettig 1994). Thus, we used printouts of emails instead of an on-line display. This set-up also allowed for flexibility and ease of feedback. Using pens, printouts, and a big table to support spatial arrangements (Figure 1), participants could move papers around to compare them, scratch things out, draw circles, or write on them in any way they chose (Figure 2).

The participants were 13 graduate and undergraduate students (7 females, 6 males). All had previous experience using computers but did not have computer science backgrounds. All were native English speakers. The experiment followed a within-subject design, in which each participant experienced all three explanation paradigms. We counterbalanced learning effects in our design by randomizing the order of explanation paradigms that each participant experienced.

The experiment was conducted one participant at a time with a facilitator interacting with the participant and an observer taking additional notes. First, the participant was familiarized with thinking aloud. Next, he or she looked through 40 sample pre-classified email messages to become familiar with the folders and to develop an intuition for how new email messages should be categorized; this sample was kept the same across participants. At this point, the main task began, divided into three 15-minute blocks (one per explanation paradigm) of processing the remainder of email messages from the dataset.

For the main task, we randomized assignments of emails to explanation paradigms to avoid exclusive association of an email with just one paradigm. For each message, the facilitator handed a new printout to the participant, who decided whether the predicted folder classification was correct, reclassified the message if needed, and gave feedback to improve the classification if needed. The participants were told that an evolving "virtual email assistant" had been implemented, and that we wanted their help "in order to get these predictions working as well as they possibly can."

After each paradigm's 15-minute block, participants provided subjective self-evaluations of mental effort, time pressure, overall effort, performance success, and frustration level, based on standard NASA TLX questions (Hart and Staveland 1988). They also took a comprehension test for that paradigm. Finally, at the end of the study, they compared and ranked all three explanation paradigms in terms of overall preference, ease of understanding, and ease of feedback.

## 4. EXPLANATIONS OF THE LEARNING ALGORITHMS
We generated the explanations under the following additional three design principles:

*(Principle 3) Common algorithms*: We focused on standard implementations of machine learning algorithms found in Weka



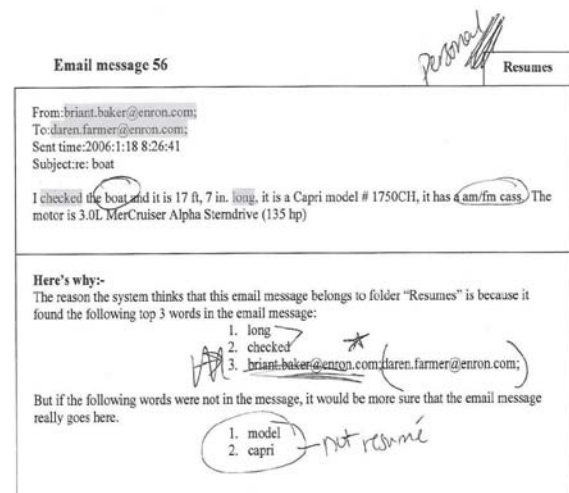Figure 1: Lo-fi prototype set-up with pens, printouts, table.



Figure 2: Example of participant feedback.

(Witten and Frank 2005) that were viable for (a) generating explanations and (b) good performance in the email domain.

*(Principle 4) Simplified but faithful explanations*: It does not seem reasonable to provide end users a complete explanation of a statistical learning algorithm. Instead, we sought to develop explanations that would be informal, yet accurate enough to engender useful mental models of this reasoning, analogous to "naïve physics" descriptions of qualitative physics.

*(Principle 5) Concrete explanations:* The explanations were required to be in terms of specific features that were visible in the current email message. In machine learning terms, each *feature* is an attribute that helps describe an email. In a bag-of-words approach, which is characteristic of our algorithms, the words that appear in emails are considered to be individual features.

## 4.1 Learning Algorithms and Training

We chose two learning algorithms: the Ripper rule-learning algorithm (Cohen 1996) and the Naïve Bayes algorithm (Mitchell 1997). These algorithms have been widely applied for email classification (e.g., Cohen 1996; Dalvi et al. 2004; Shen et al. 2006). To obtain a prediction for each of the 122 email messages, we performed a stratified 5-fold cross-validation.

Prior to training, each email message was preprocessed to remove headers and common "stop" words. The remaining words were stemmed by Porter's method to remove word endings (Porter 1980). Each email message was then represented as a Boolean vector with a Boolean feature for each observed email sender (the From field), one Boolean feature for each observed set of email recipients (the union of the From, To, CC, and BCC fields, in essence identifying the "team" of people to which the message relates (Shen et al. 2006)), and one Boolean feature for each distinct word observed in the Subject and Body fields.

Ripper learns a set of classification rules. The rules are ordered by class but unordered within class. Hence, to make a prediction, Ripper first applied the rules for the least frequent class (Bankrupt in our dataset). If one of these rules matched the email message, it was classified as Bankrupt. Otherwise, Ripper moved on to the rules for the next most frequent class (Resume), and so on. There were no rules for the most frequent class (Enron News); it was treated as the default if none of the rules for the other classes matched.

Naïve Bayes estimates the probability of each folder given the email message. Let the random variable $Y \in \{y_1, y_2, \ldots y_m\}$ represent a folder that the email message can be filed under. The variable $Y$ can take on $m$ possible values, corresponding to the $m$ folders. Additionally, let the vector $X = (x_1, x_2, \ldots x_n)$ be an email message where each component $x_i$ represents a feature of the email. More precisely, each $x_i$ is a Boolean variable indicating the presence or absence of the features described above. The goal of the Naïve Bayes algorithm is to calculate the posterior probability $P(Y = y_k \mid X)$. The Naïve Bayes algorithm simplifies the calculation of this posterior probability by making the assumption that probability of the each feature $x_i$ is conditionally independent given the folder Y. The mathematical details of the Naïve Bayes algorithm are shown in Figure A1 in Appendix A.

The overall accuracy of the predictions was not particularly high: 60% for Naïve Bayes and 75% for Ripper when used to classify the entire set of emails. We would have preferred higher accuracy and equal accuracy between algorithms. Still, high accuracy was not required to answer our experiment's research questions, and our analysis takes accuracy differences into account.

## 4.2 Generating Explanations

The Rule-based explanations (Figure 3) were generated by highlighting the rule, learned by the Ripper algorithm, that made the classification. That rule was listed above all other possible rules in the explanation. The Keyword-based and Similarity-based explanations were both generated from the learned Naïve Bayes classifier algorithm.

Consistent with our fifth design principle ("visible words only"), the Keyword-based explanations (Figure 4) were generated by listing up to five words *present in the email message* having the largest positive weights as well as up to five words *present in the email message* having the most negative weights. (As we will discuss later, users found this latter set of "negative" words counter-intuitive. They are the words whose presence in the message *reduces* the certainty of the classifier in the sense that the classifier would be more confident if these words did *not* appear.)

The Similarity-based explanations (Figure 5) were generated by computing the training email message that, if deleted from

the training set, would have most decreased the score. This was typically the training example most similar to the email message being classified. This example was then displayed and up to five words with the highest weights that appeared in both the example and the target email message were highlighted.

---

Resume

From: toni.graham@enron.com
To: daren.farmer@enron.com
Subject: re: job posting

Daren, is this position budgeted and who does it report to?
Thanks,
Toni Graham

---

The reason the system thinks that this email message belongs to folder "Resume" is because the highest priority rule that fits this email message was:

- **Put the email in folder "Resume" if:**
  **It's from toni.graham@enron.com.**

The other rules in the system are:

...

- Put the email in folder "Personal" if:
  The message does not contain the word "Enron" and
  The message does not contain the word "process" and
  The message does not contain the word "term" and
  The message does not contain the word "link".

- Put the email in folder "Enron News" if:
  No other rule applies.

Figure 3: (Top): Email.
(Bottom): Rule-based explanation excerpt.

---

Personal

From: buylow@houston.rr.com
To: j..farmer@enron.com
Subject: life in general

Good god -- where do you find time for all of that? You should w...

By the way, what is your new address? I may want to come by ... your work sounds better than anything on TV.

You will make a good trader. Good relationships and flexible pri... a few zillion other intangibles you will run into. It beats the hell o... other things.

I'll let you be for now, but do keep those stories coming we love...

---

The reason the system thinks that this email message belongs to folder "Personal" is because it found the following top 5 words in the email message:

1. ill
2. love
3. better
4. things
5. god

But if the following words were not in the message, it would be more sure the email message really goes here.

1. keep
2. find
3. trader
4. book
5. general

Figure 4: (Top): Excerpt from email.
(Bottom): Keyword-based explanation,
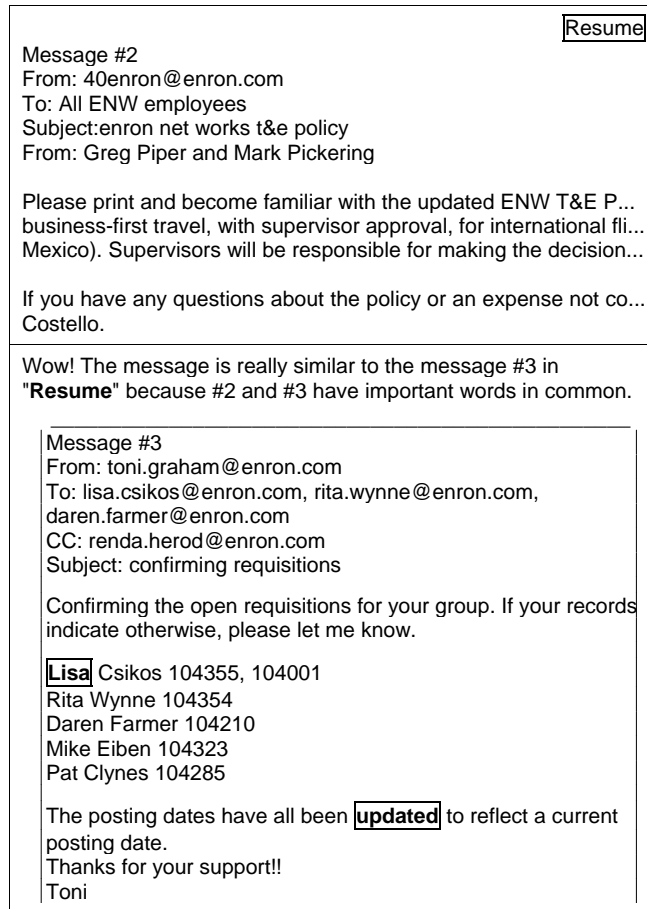supplementing the highlights in the email.

Message #2
From: 40enron@enron.com
To: All ENW employees
Subject:enron net works t&e policy
From: Greg Piper and Mark Pickering

Please print and become familiar with the updated ENW T&E P...
business-first travel, with supervisor approval, for international fli...
Mexico). Supervisors will be responsible for making the decision...

If you have any questions about the policy or an expense not co...
Costello.

Wow! The message is really similar to the message #3 in
"**Resume**" because #2 and #3 have important words in common.

Message #3
From: toni.graham@enron.com
To: lisa.csikos@enron.com, rita.wynne@enron.com,
daren.farmer@enron.com
CC: renda.herod@enron.com
Subject: confirming requisitions

Confirming the open requisitions for your group. If your records
indicate otherwise, please let me know.

**Lisa** Csikos 104355, 104001
Rita Wynne 104354
Daren Farmer 104210
Mike Eiben 104323
Pat Clynes 104285

The posting dates have all been **updated** to reflect a current
posting date.
Thanks for your support!!
Toni

Figure 5: (Top): Excerpt from email.
(Bottom): Its Similarity-based explanation.

# 5. METHODOLOGY FOR QUALITATIVE ANALYSIS

The think-aloud data and questionnaire comments were analyzed using two coding schemes. In the first coding scheme, we coded all user utterances with the goal of determining the reaction of users to the explanations. In the second coding scheme, we performed a more detailed analysis of only the utterances that constituted negative comments about the explanations or suggested changes to the learning algorithms. The goal of the second coding scheme was to classify the users' feedback with respect to requirements for machine learning algorithms and the background knowledge required for the algorithms to be able to incorporate the feedback. The main codes, along with a description and example are shown in the first three columns of Table 1. The second coding scheme is discussed in Section 7.

For both schemes, to ensure consistency in interpretation of the codes and when to use them, two researchers independently coded a small subset. They then iterated on this subset, further refining the codes and developing norms about how to apply them. For the first coding scheme, the total agreement value was 81% for the first subset at the end of these iterations, which indicates high coding reliability. For the second coding scheme, the agreement was 82% after the iterations. At this point, the schemes were deemed robust enough, and the remaining data were then coded.

For both coding schemes, we calculated the inter-rater agreement as the percentage of intersection of the codes divided by the union of all codes applied. For example, if one researcher gave the codes {Breakdown, Suggest Change} for one email and another researcher gave the codes as {Emotion, Suggest Change} for the same email, then the agreement was calculated as 1/3 (33%) as follows:

$$\frac{|\{Breakdown, SuggestCha nge\} \cap \{Emotion, SuggestCha nge\}|}{|\{Breakdown, SuggestCha nge\} \cup \{Emotion, SuggestCha nge\}|}$$

# 6. RESULTS: EXPLAINING TO USERS

Analyzing the video transcripts and questionnaires using the coding scheme just described produced the counts shown in the final column of Table 1. (We will not discuss further the codes making up less than 1% of the total.)

| Code | Description | Example from data | Count (% of total) |
|---|---|---|---|
| Breakdown | Expressing confusion or lack of understanding with the explanation of the algorithm. | I don't understand why there is a second email. | 41 (8%) |
| Understand | Explicitly showing evidence of understanding the explanation of the algorithm. | I see why it used "Houston" as negative | 85 (17%) |
| Emotion | Expressing emotions. | It's funny to me. | 15 (3%) |
| Trust | Stating that he or she trusted the system. | I would probably trust it if I was doing email. | 1 (<1%) |
| Expectation | Expressing an expectation for the system to behave in a certain way. | I hope that eventually the intelligent assistant would learn to give more reasons. | 2 (<1%) |
| Suggest change | Correcting the explanations or otherwise suggesting changes to the system's reasoning. | Different words could have been found in common, like "Agreement," "Ken Lay." | 161 (32%) |
| Negative comment | Making negative comments about the explanation (without suggesting an improvement). | …arbitrary words: "energy" especially bad. | 100 (20%) |
| Positive comment | Making positive comments about the explanation. | The Resume rules are good. | 94 (19%) |

Table 1: The first coding scheme.

## 6.1 Explaining to Users: Understandability

### 6.1.1 Which Paradigms Did They Understand?

According to the participants' responses to the questionnaires, the Rule-based explanation paradigm was the most understandable (Table 2). This was corroborated by their verbal remarks: Rule-based explanations generated three times as many remarks indicating understanding and less than a tenth the remarks indicating breakdowns as either Keyword-based or Similarity-based explanations.

Differentiating between shallow and deep understanding reveals further insights. "Shallow understanding" in this context means that participants were simply able to make the classification decision the same way as the explanation paradigms. To gather data needed to assess this aspect, the questionnaires included an email without a classification or explanation, and asked the participants to categorize it to a folder based on what the paradigm would predict. Nearly all participants categorized it correctly in the case of Rule-based explanations and Keyword-based explanations, but only 4 of the 13 participants categorized the email correctly in the Similarity-based case.

"Deep understanding" implies understanding the *reasoning behind* the classification decision of the explanation paradigms. The questionnaires included an email with a classification but without the explanation, and participants were asked *why* the paradigm would classify an email the way it did. For the Rule-based explanation paradigm, a majority of participants answered by giving a rule, and some even managed to reconstruct a close version of the actual rule that was applied. For Keyword-based explanations, nearly all participants answered with keywords, even managing to identify correctly some of the keywords used in the actual example. However, only three participants answered even close to correctly for the Similarity-based case.

The combined evidence, from the participants' opinions and their inability to explain or reproduce the Similarity logic, is thus quite strong that the Similarity-based explanations had a serious understandability problem.

| Explanation | Rank 1 | Rank 2 | Rank 3 |
|---|---|---|---|
| Rule-based | 9 | 2 | 2 |
| Keyword-based | 3 | 6 | 4 |
| Similarity-based | 1 | 5 | 7 |

Table 2: Participants' rankings from the written questionnaires. (Rank 1 is "understood the most.")

### 6.1.2 What Factors Affected Understanding?

We investigated the factors that contributed to understanding via the Understand, Breakdown, and Negative Comments

codes from the video transcripts and written questionnaire comments. Three factors stood out in affecting understanding of the system's behavior: understanding of the general idea of the algorithm, the Keyword-based explanations' negative keyword list, and appropriateness of word choices.

Regarding understanding of the algorithm, some participants expressed understanding of the algorithm by describing the essential strategy, as in the following two quotes. This enabled them to predict system behavior.

*P6 (on Rule-based):* "I understand why it would just default to Enron, since that's what the rule is."

*P1 (on Similarity-based):* "I guess it went in here because it was similar to another email I had already put in that folder."

In the case of the Keyword-based paradigm, some problems in understanding were caused by the negative keyword list. Nobody had anything positive to say about the inclusion of negative keywords in the explanation:

*P6 (on Keyword-based):* "So what does this mean (referring to 2nd set of words)?"

*P8 (on Keyword-based):* "I guess I really don't understand what it's doing here. If those words weren't in the message?"

Finally, appropriateness of the word choices seemed to have an effect on understanding, especially if they were felt by participants to be common words or topically unrelated:

*P1 (on Similarity-based):* "'Day', 'soon', and 'listed' are incredibly arbitrary keywords."

### 6.1.3 Discussion: Understanding.

In addition to the clear evidence of understandability problems for Similarity-based explanations, we would like to point out three results of particular interest.

First, although Rule-based explanations were consistently understandable to more than half the participants and, at least for this group of participants, seemed to "win" over the other two paradigms, note that about one-third of the participants preferred one of the other explanation paradigms. This implies that machine learning systems may need to support *multiple* explanation paradigms in order to effectively reach all of their users.

Second, Keyword-based explanations seemed to be reasonably understandable except for the negative keyword list, which our results suggest was a problem. There are several potential remedies. One possibility is that the negative keyword list could be explained in some different way to give users a better understanding of how the algorithm works. For example, instead of drawing attention to words with negative weights that are present in the email, the explanation could make use of the strongest negative weights associated with words that are *absent* from emails, since their absence increases the confidence of the learning algorithm. Another possibility is that the negative keyword list should be omitted from the explanation altogether.

Third, the *topical* appropriateness of word choices seemed particularly critical to participants' ability to predict and understand system behavior. This knowledge is too complex to be learned from only 122 email messages, but it could be possible in larger document collections; we will return to this point in Section 7.2.

## 6.2 Explaining to Users: Preferred Paradigms and Why

Following the understanding trends, participants' rankings favored the Rule-based explanations over the other two (Table 3). Still, nearly 50% of the participants chose a paradigm other than Rule-based as their favorite, so the Rule-based paradigm did not receive a clear mandate.

We expected preference to closely follow understanding trends, but analysis of the Positive Comments, the positive Emotion codes, and the questionnaire responses provided additional useful insights into factors that seemed to affect participants' preferences in positive ways. These remarks fell into four categories, three of which (approval of reasoning soundness, clear communication of reasoning, and perceived accuracy) tie at least somewhat to understandability.

Participants' approval of soundness of reasoning was remarked upon often. Also, clear *communication* of reasoning, which is distinctly different from the mere presence of sound reasoning, mattered to a number of our participants. For example, Participant 1's comment below is fairly representative of several about the reasoning itself, whereas Participant 10's comment exemplifies several comments specifically about communication of the reasoning:

*P1 (on Keyword-based):* "The reasons seem like perfectly good reasons…this is a good reason why it shouldn't be in Personal."

*P10 (on Similarity based)*: "I like this one because it shows relationship between other messages in the same folder rather than just

spitting a bunch of rules with no reason behind it."

High accuracy, as perceived by the participants, was remarked upon often. (We will return to the influence of *actual* accuracy shortly). For example:

P11 (on Rule-based): "I think this is a really good filter. Put in Resume if it's from toni.graham"

P2 (on Similarity-based): "Similarity was my favorite - seemed the most accurate, and took email addresses into account."

The fourth category was unexpected: Several participants appreciated Similarity-based explanations' less technical style of expression, a characteristic we inadvertently introduced in our wording that emphasizes similarity ("Wow!"). This introduction of informality in the form of slang produced a number of Positive Comments for that explanation paradigm, pointing out possible benefits from relaxing the language style used in explanations. For example:

P1 (on Similarity-based): "I also appreciate how the computer is excited about its decision... It's funny to me ... Told you, in conversational form, why it was similar."

P10 (on Similarity-based): "This is funny... (laughs) ... This seems more personable. Seems like a narration rather than just straight rules. It's almost like a conversation."

| Explanation | Rank 1 | Rank 2 | Rank 3 |
|---|---|---|---|
| Rule-based | 7 | 4 | 2 |
| Keyword-based | 3 | 4 | 6 |
| Similarity-based | 3 | 5 | 5 |

Table 3: Participants' rankings from the written questionnaires. (Rank 1 is "preferred the most.")

## 6.3 Accuracy

As we have mentioned, the algorithms performed at different accuracy rates, with Ripper outperforming Naïve Bayes. (We define "accurate" as being in agreement with the original Enron user who owned the email.) It surprised us that Ripper was so much more accurate than Naïve Bayes. This suggests that Ripper may be a better choice than Naïve Bayes for accuracy in this type of situation. As to our experiment, accuracy rate was not statistically predictive via linear regression of any of the ratings provided by participants, did not result in differences in participants' willingness to provide feedback, and did not affect their accuracy in doing so.

When the participants disagreed with the machine (28% of the time), participants were usually right (22%), but not always (6%). Also, both the machine and the participants disagreed with the original user 22% of the time, suggesting some knowledge possessed only by the original user and perhaps even some misfiling by the original user. Ultimately, the participant corrections brought the accuracy rates for all paradigms to almost identical levels: 71-72%, a surprising result suggesting that the performance of the algorithms *did not matter to accuracy*: the participants brought them all to the same accuracy level.

As the preceding paragraph also points out, the participants were not perfect oracles. The error rate is consistent with earlier findings regarding end-user programmers' accuracy in serving as oracles when debugging, which have reported error rates of 5-20% (e.g., Phalgune et al. 2005). This range of error rates has been robust across studies, and suggests a similar level of "noise" that users' judgments would introduce into the learning algorithm's data.

There was an odd relationship between actual accuracy and participants' speed (efficiency). With Keyword-based and Similarity-based explanations, there was no significant relationship between accuracy and participants' speed, but in the case of Rule-based explanations, the predictive effect was negative! More specifically, an interaction test of the number of processed emails predicted by accuracy, participant, and the interaction of accuracy and participant, showed an interaction effect of participant and accuracy predicting the number processed with Rule-based explanations (linear regression, $p$=0.0314, $F[3,9]$=2.515, $R^2$=0.456). As expected, there was also a main effect of participant predicting the performance speed in the case of Rule-based explanation paradigm (linear regression, $p$=0.0299, $F[3,9]$=2.515, $R^2$=0.456), which simply says that some participants were faster than others in this paradigm. This says that the number of emails processed depended, in the Rule-based paradigm, on the participant and the accuracy rate he/she experienced, in a negative direction: the higher the accuracy rate for some participants, the *lower* the number of emails processed.

The odd relationship between accuracy and participant efficiency suggests the possibility of a point of diminishing return. That is, the more accurate the algorithms (and their explanations), the harder it may be for a user to spot the flaw when the

algorithm makes a mistake. Thus, there may be a threshold accuracy point beyond which users may view the cost of guiding the algorithm further as exceeding the benefit of doing so.

## 7. RESULTS: THE USERS EXPLAIN BACK

What did participants think machine learning algorithms should change? For all three paradigms, we coded participants' feedback (Negative Comment and Suggest Change) along two dimensions. The rows of Table 4 identify the type of change and the columns identify the knowledge needed to handle the change.

| | KB-English | KB-commonsense | KB-domain | KB-other | Total | % |
|---|---|---|---|---|---|---|
| 1. Select different features (words) | 70 | 64 | 25 | 16 | 175 | 53% |
| 2. Adjust weight | 11 | 11 | 4 | 13 | 39 | 12% |
| 3. Parse or extract in a different way | 7 | 17 | 10 | 0 | 34 | 10% |
| 4. Employ feature combinations | 9 | 5 | 2 | 1 | 17 | 5% |
| 5. Relational features | 0 | 9 | 5 | 0 | 14 | 4% |
| 6. Other | 3 | 12 | 4 | 33 | 52 | 16% |
| Total | 100 | 118 | 50 | 63 | 331 | |
| % | 30% | 36% | 15% | 19% | | |

Table 4: Types of participants' changes (in rows) that required various background knowledge (in columns).

## 7.1 Participants' Suggestions by Type

The rows of Table 4 categorize the participants' suggestions into six types. The types of feedback were not independent of the explanation paradigms: some paradigms seemed to encourage participants to think along the lines of particular types of feedback. We will point out these influences along the way whenever a paradigm represents at least 50% of a category.

The first type, selecting different features, was the most widespread type of feedback, for all three explanation paradigms. Words that appear in an email are considered to be individual features by the machine learning algorithms. Each feature has a weight that influences its contribution to determining a particular classification. More than half of all feedback referred to either adding a new feature for the algorithm to consider or removing a feature from consideration, such as:

*P13 (on Rule-based)*: "It should put email in 'Enron News' if it has the keywords 'changes' and 'policy'. I put down some keywords that I noticed."

The second type was comprised of suggestions in which participants agreed that a particular feature was worthy of note, but wanted to change its weight or importance. Participants' reactions to Keyword-based explanations generated 69% of the feedback of this type, perhaps because of the feature-focused nature of the Keyword-based paradigm. Some participants' suggestions for changing feature weight or importance involved adjusting the weight on features in a general sense, such as P8's below. Other participants, such as P7 and P1, flipped a weight from negative to positive (or vice versa), or focused on the frequency of the word occurrence in the email, akin to term weighting in information retrieval. Finally, participants such as P4 made other adjustments to ordering of relative importance.

*P8 (on Keyword-based)*: "The second set of words should be given more importance."

*P7 (on Keyword-based)*: "Keyword 'include' is not good for the second set. It should be in the first set of words."

*P1 (on Rule-based)*: "'Bankruptcy' is here over and over again, and that seems to be the obvious word to have in here."

*P4 (on Keyword-based)*: "I think that 'payroll' should come before 'year'."

The third type of feedback concerned parsing the text or extracting features from the text in a different way. Some participants suggested a different form of text parsing, such as P1 below. In the simplest case, this could be achieved by an improved stemming procedure (Porter 1980). In some cases, however, the suggested extraction operates on a structure such as a URL, such as P6. In this case, either the system would already need to know about URLs or else the user would need to define them (perhaps by giving examples). Participants such as P13 also suggested using the structure of the email to extract features, such as the "From" and "Subject" field. Finally, some participants such as P6 suggested new kinds of informative cues.

>   *P1 (on Similarity-based):* "Different forms of the same word must be looked at."
>
>   *P6 (on Similarity-based):* "I think it would be good to recognize a URL."
>
>   *P13 (on Rule-based):* "Yea, I mean it has 'job' in the subject line (for sorting into Resumé folder)"
>
>   *P6 (on Keyword-based):* "I think that it should look for typos in the punctuation for indicators toward Personal."

Feature combinations were the fourth type of user feedback. Participants pointed out that two or more features taken together could improve the prediction, especially when they were working with Similarity-based explanations, which generated 63% of the suggestions of this type:

>   *P6 (on Similarity-based):* "I think it would be better if it recognized a last and a first name together."
>
>   *P12 (on Keyword-based):* "I would think like 'authorize signature' or 'w-2 form'."

The fifth type of participant feedback suggested incorporating the use of relational features. In these cases, the participants used relationships between messages (threading) or organizational roles (chairman of Enron) to define a new feature. For example:

>   *P6 (on Rule-based):* "I think maybe it should use the response and automatically put it in the folder with the message that was responded to."
>
>   *P8 (on Keyword-based):* "This message should be in 'EnronNews' since it is from the chairman of the company."

The remaining feedback by users did not fall into the types above and we coded these as "Other". Most of this kind of feedback concerned changes to the learning algorithm itself. These included suggestions such as adding logical NOT to the rule language, eliminating the default rule in Ripper, requiring an equal number of positive and negative keywords in Keyword-based explanations.

There were also cases in which the real problem lay with the way the explanation was constructed, rather than with the learning algorithm:

>   *P13 (on Similarity-based):* "Having 'points' being the only keyword, I mean that kind of sucks."

Although the algorithm actually used all of the keywords in the messages, the explanation only highlighted the shared words with the top weights. This suggests that an automated method for assimilating user feedback would need a component that could diagnose whether the perceived problem is due to approximations in the explanation or design decisions in the learning algorithm.

## 7.2 Participants' Suggestions by Knowledge Source

Table 4's columns categorize the participants' suggestions into four knowledge sources: knowledge of English, commonsense knowledge, domain-dependent knowledge, and other.

Almost a third (30%) of the participations' suggestions relied on knowledge of English (KB-English). We coded feedback in this category if the necessary knowledge could be learned from analysis of large document collections or obtained from other online resources (e.g., Wordnet (Miller 1995) or named-entity recognizers (Zhou and Su 2002)). For example:

>   *P8 (on Rule-based):* "Does the computer know the difference between 'resumé' and 'resume'? You might have email where you're talking about 'resume' but not in a job-hiring sense."
>
>   *P5 (on Keyword-based):* "Last names would be better indicators."
>
>   *P1 (on Similarity-based):* "'day', 'soon' and 'listed' are incredibly arbitrary keywords."

Some knowledge might need to be manually encoded, but it could then be reused across many different organizations and applications (KB-Commonsense). For example, participants indicated that there are "families" of words that are work- or business-related, and also suggested topic-related words:

>   *P4 (on Rule-based):* "'Policy' would probably be a good word that would be used a lot during business talk."
>
>   *P1 (on Keyword-based):* "'Qualifications' would seem like a really good Resume word, I wonder why that's not down here."

Some knowledge is domain-dependent (KB-Domain). Some participant suggestions relied on knowledge specific to Enron. For example, in the following example, the machine learning system would need to know that Ken Lay was the CEO of Enron and that as such he carries special importance to Enron employees, would need to know the implications of being an

Enron employee, and so on. Such knowledge would need to be encoded separately for each organization.

*P11 (on Similarity-based):* "Different words could have been found in common like 'Agreement', 'Ken Lay'."

All remaining feedback was coded KB-Other. This usually occurred when a participant's comment was not specific enough to suggest the underlying knowledge source.

## 8. INCORPORATING PARTICIPANTS' SUGGESTIONS INTO MACHINE LEARNING
Given the participants' feedback, we conducted two follow-up studies to investigate the following research question:

RQ 3. Can these types of user feedback be assimilated by existing learning algorithms? If so, exactly *how* could some of these types of user feedback be incorporated into machine learning algorithms, and does doing so actually *improve* the performance of algorithms?

Currently most learning systems take only a simple form of feedback into account: whether the prediction was right or wrong, or what the correct prediction should have been. *How* rich user feedback could be incorporated into existing algorithms has received little attention previously. Further, the effects of incorporating rich user feedback on accuracy require investigation. For example, one issue is that our participants were not perfect, and sometimes their feedback was erroneous; introducing such errors into the system's reasoning might make the algorithm perform worse rather than better.

First, consider how the user feedback reported in Experiment #1 might be incorporated statically into the algorithms, i.e., how machine learning algorithms might be changed permanently to do as the users suggested in those comments without complex *in situ* reasoning advice from users. Consider Ripper and Naïve Bayes, the widely used algorithms that we used in Experiment #1. For these two algorithms, Type codes 1, 2, and 4 (features, weights, and feature combinations) as supported by KB-English are conducive of direct assimilation. These types of user feedback accounted for 27% of the 331 suggestions. Supporting these three type codes by KB-Commonsense rather than solely by KB-English may also be possible but is more challenging, due to the difficulty of obtaining the relevant common sense topic models. Type codes 1, 2, and 4 that relied on KB-Commonsense accounted for an additional 24% of the suggestions. The remaining type codes appear to require substantial additional machine learning research in feature extraction, relational learning, and user interfaces for providing application-specific feature specifications. Thus, a little over half of the participants' suggestions appear potentially viable for incorporation into these two particular algorithms.

Now consider the second issue, whether incorporating this feedback would improve accuracy. To make consideration of this issue tractable, we narrowed our focus to two of the feedback types, namely weight and keywords changes, without regard to knowledge source. Thus, we assumed that either the knowledge sources would be available, or else that the users themselves would serve as the knowledge source and would enter the changes interactively during reasoning. The changes necessary for these two feedback types could be implemented in the algorithms of Experiment #1, thereby covering 65% of the user feedback given. This approach allowed empirical exploration of the last part of our research question, namely whether incorporating rich user feedback would actually improve the algorithm.

We conducted two off-line experiments that investigated this question under two approaches to incorporating this feedback, namely, a constraint-based approach and a co-training approach. For the constraint-based approach, automatically modifying the scoring functions for feature selection and pruning in Ripper is substantially more complex than for Naïve Bayes, and hence we evaluated the constraint-based approach via Naïve Bayes only in Experiment #2. Experiment #3, the co-training approach, was evaluated with both Ripper and Naïve Bayes. Participant-by-participant details of the evaluations are given in Appendix B.

### 8.1 Experiment #2: Using User Feedback as Constraints
One possible approach is to incorporate rich user feedback of type codes 1 and 2 into existing algorithms as user constraints (e.g., "there should be a positive weight on feature X") (Altendorf et al. 2005). From the constraints, the corresponding weights can be automatically modified in the Naïve Bayes classifier.

From the user's point of view, this approach is appealing for three reasons. First, with the constraint-based approach, the system can heed the user's feedback right away. It does not require a large set of corrections for the feedback to begin making a difference. The reason the system can be so responsive is that, instead of merely converting the user feedback into an additional training example, the constraint-based approach can employ the rich user feedback directly in the actual machine learning algorithm by translating it directly into constraints. Second, the constraint-based approach is fairly easily understandable and explainable: because the user's feedback is translated into constraints, the constraints the system is trying

to heed can be explained much as the user originally expressed them. Third, constraints offer flexibility since they can be varied according to their "hardness." Some constraints that the user specifies may be hard constraints that the system *must* satisfy. Some constraints can be softer and should be considered by the system if possible but are not absolutely necessary.

We therefore implemented the following three types of constraints:

Constraint type 1 (hard):   If the participant reduced the weight or removed the word proposed in the explanation, the proposed word was considered vague or unimportant. The word was removed from the feature set.

Constraint type 2 (soft):   If the participant increased the weight of a word, regardless of whether it was proposed in the explanation or not, the proposed word was assumed to be a strong indicator as to the email's folder. We incorporated this form of feedback by adding a constraint that forces the weight on the word to be positive, which made the word more important for the user-assigned folder than for other folders. In our implementation, we also increased the effect of this constraint by requiring the weight of the word to hold above some amount $\delta \geq 0$. Figure A2 in Appendix A illustrates the mathematical details of this second constraint type.

Constraint type 3 (soft):   For the third constraint type, we assumed that words that had their weights increased by the user were more significant indicators of the user-assigned folder than other words appearing in email messages in that folder. In order to describe this constraint at a high level, we need to quantify the *predictiveness* of a word for the user-assigned folder. This predictiveness is defined as the probability of the email being assigned to the user-specified folder given that the word is present in the email message. We require that there should be a gap between the predictiveness of a user-selected word and the predictiveness of the most informative words for that folder that were not selected by the user. Thus, for each user-selected word, we added a set of 10 constraints to force the user-selected words to exceed by some amount $\theta \geq 0$ the predictiveness of the top 10 most informative words not selected by the user. The formal notation for constraints of type 3 are shown in Figure A3 in Appendix A.

The $\delta$ and $\theta$ parameters above control the hardness of the constraint and can take on any value between 0 and 1. The lower the value of these parameters, the softer the constraint. We varied these values in order to investigate the influence of the constraint parameter values on the accuracy of the classification. To keep the number of runs tractable, we chose increment steps of 0.2. $\delta$ took on the values 0, 0.2, 0.4 or 0.6. Likewise, we varied $\theta$ to be 0, 0.2 or 0.4. We chose not to investigate values over 0.6 as there was a danger of making the constraints too aggressive. (With $\theta$ values over 0.6, the risk was too high of over-generalizing spurious constraints over all the data, resulting in a problem in machine learning known as overfitting.)

Constraints of type 1 were incorporated into Naïve Bayes by simply excluding the feature from the classifier. The remaining two constraints types were incorporated into the parameters of the Naïve Bayes classifier. The standard method for setting the parameters for Naïve Bayes is to calculate the parameters from the training data using a process known as maximum likelihood estimation. We were able to incorporate the constraints by converting the maximum likelihood estimation procedure into a constrained optimization problem as shown in Figure A4 in Appendix A. In our experiment, this optimization was performed by the non-linear programming solver Lindo (Lindo 2007), a software package for resolving optimization problems.

To evaluate the improvements that this approach could achieve, we compared against a baseline. This baseline is the traditional way user feedback is incorporated into machine learning, a simple Naïve Bayes algorithm that takes only the corrected folder assignments provided by users into account ("this should be in Personal"), but not the suggestions' rich content ("because it's from Mom").

To train and evaluate these constraint additions versus the baseline versions of both the keyword and similarity variants of Naïve Bayes, we divided the emails from the Enron farmer-d dataset into three sets: the *training set*, the *feedback set* and the *evaluation set*. The training set contained the same 87 emails that were used to train the machine learning classifier for the main experiment. Emails in the feedback set were the ones that had been presented to participants in the main experiment. Since participants were constrained by time, some participants were faster than others and were able to provide feedback on more emails in the feedback set. See Appendix B1 for the exact number of emails that the participants were able to give feedback on. The final set, known as the evaluation set, consisted of emails that had not been previously seen by either the
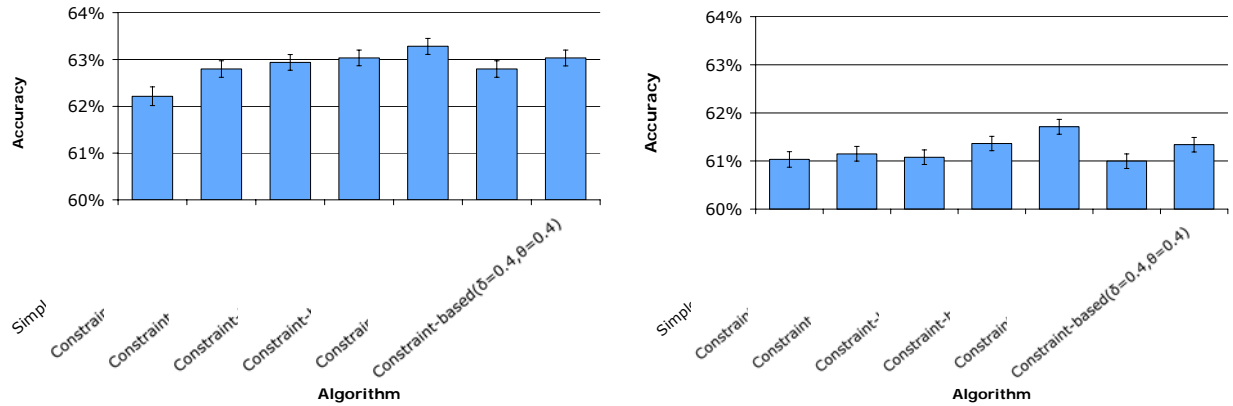
machine learning classifiers or the participants. Due to the limited number of emails in the Bankrupt and Resumes folders, the evaluation set was quite small. As a result, we had a total of 20 email messages for the feedback and evaluation sets. All email messages that a participant provided feedback on were assigned to that participant's feedback set while the remaining email messages out of the 20 were assigned to the evaluation set.

In order to evaluate the algorithms, we could have used a holdout test set approach in which we trained the classifier on the combined training set and a participant's feedback set and then tested it on the evaluation set, one participant at a time. Doing so would have resulted in accuracy numbers on a sample size of 13, the number of participants in the study. This is too small a sample size for statistically significant outcomes. Alternately, we could have used a standard machine learning approach called k-fold cross validation. With k-fold cross validation, we would have divided a participant's feedback set into $k$ subsets, chosen $k$-1 of the subsets to be combined with the original training set to form the cross validation training set, and then appended the remaining subset to the original evaluation set to serve as the cross validation evaluation set. The feedback from the email messages in the $k$-1 subsets would have been applied to the constraint-based classifier. Since there are $k$ possible ways to divide the feedback set into $k$-1 subsets and 1 remaining subset, we would repeat this process $k$ times (corresponding to the $k$ "folds" in cross validation). We would eventually obtain accuracy numbers for each fold and average the accuracy across the folds. However, the average number of emails that a participant was able to review during the feedback stage was approximately 5. The extremely limited size of the feedback set would have resulted in a special case of cross validation known as leave-one-out-cross validation where $k$ is equal to the number of emails in the feedback set. We did not use leave-one-out cross validation due mainly to the computational complexity of the constraint-based approach. Incorporating the feedback from all $k$-1 emails in the feedback set results in a constrained optimization problem that can be computationally expensive to solve with a large number of constraints. Another reason we did not use leave-one-out cross validation is the fact that it would have resulted in a single email in the evaluation set and even though we were averaging the results over the $k$ folds, and we would prefer a much larger evaluation set over which we average our results.

Instead, we adopted an approach that we will call leave-all-but-one-out cross validation to evaluate the constraint-based and baseline classifiers. In this approach, we created a training set by taking one email message from a participant's feedback set and combining it with the original training set. The feedback for the email message selected from the feedback set was applied to the constraint-based algorithm. Then, an evaluation set was created by combining the original evaluation set with the remaining messages from the feedback set without applying the feedback. If there were $f$ messages in the feedback set, each message served as a "training set message", resulting in $f$ rounds of leave-all-but-one-out cross validation. In Appendix B3.1 and B4.1, for each participant, we report the accuracy averaged over the $f$ rounds.

Figure 6 summarizes how the constraint-based approach using the keyword-based paradigm performed against the standard, baseline online algorithm for user feedback. Figure 7 shows a similar graph of how the constraint-based approach using a similarity-based paradigm fares against the baseline algorithm. (Full details of the evaluation are given in Appendix B). As the graphs show, improvements varied according to the constraint parameters applied, where setting δ=0.6 and θ=0 achieved highest accuracy in both paradigms. Thus, although the constraint-based approach achieved some gains in accuracy, the gains were not substantial, resulting in an increase of only approximately 1%.

Such a small gain suggests that this approach is not a feasible way to incorporate user feedback into machine learning systems. Users would be unlikely to regard use of their time to provide feedback to be a good investment for such a small improvement. Of course, this is just one experiment with data from one group of users; follow-up studies are needed to establish the generality of this algorithm's viability, or lack thereof, for incorporating user feedback in the form of constraints.

Figures 6 (left: Keyword) and 7 (right: Similarity): Percentage of emails the constraint-based algorithms sorted correctly, after incorporating participants' rich feedback. (Bars indicate standard error). These accuracy rates are only slightly higher than the algorithm achieved without the incorporation of participants' rich feedback.

## 8.2 Experiment #3: Using User Feedback in Co-training

Since the constraint-based approach resulted in only a little improvement, we decided to explore more complex machine learning approaches that could leverage user feedback more aggressively. The approach we selected is similar in spirit to the co-training algorithm (Blum and Mitchell 1998) which is a commonly used algorithm in a sub-area of machine learning known as semi-supervised learning (Chapelle et al. 2006).

Semi-supervised learning is used when labeled data is scarce but unlabeled data is abundant. For example, consider the email sorting situation in which the user has assigned a small number of emails to folders while the majority of the emails in her inbox remain unsorted. The limited number of emails that have been assigned to folders is considered to be *labeled* data while the remaining unsorted emails are considered to be *unlabeled* data. Even though the majority of emails are unlabeled, they can still provide useful information for classification because they may form natural groupings based on the presence or absence of certain words. Semi-supervised learning aims to improve the performance of a learning algorithm trained on the small amount of labeled data by leveraging the structure of the unlabeled data.

The co-training (Blum and Mitchell 1998; Kiritchenko and Matwin 2001; Balcan et al. 2005) approach to semi-supervised learning employs two classifiers that work on the same data but have two different "views" of the data through independent sets of features. This approach operates under the assumption that the two classifiers produce the same classification even though they have different views. In our email example, the set of features for the first classifier consist of email addresses that appeared in the email header fields and in the body of an email message. The feature set for the second classifier includes words appearing in the subject line and the email body. These two feature sets are from the same emails but they are two independent sets. Initially, the two classifiers are trained on a labeled training set. Then, in the second phase of training, the classifiers compare which of the two can more confidently classify an email message to a folder. The most confidently classified email message along with its folder assignment is then added to the training set for the next round of training. The standard co-training algorithm is described in pseudocode in Figure A5 in Appendix A.

We adapted the notion of co-training by regarding the *user* as if he or she were one of the classifiers in co-training, and a standard machine learning classifier such as Naïve Bayes or Ripper as the other classifier. We call this modified version of co-training *user co-training*. In order to treat the user as the second classifier, we developed a *user feedback classifier* that represents the user and treats the important words selected in the user feedback as a set of features for the specific folder to which the user assigns the email. Thus, associated with each folder $f$ is a vector of important words $v_f$ obtained by taking the union of all the important words selected by the user in the email messages placed into folder $f$. For the purposes of our experiment, this came from the union of each single participant's emails processed in this way during Experiment #1. The user feedback classifier uses this information to classify an email message by finding the folder with the highest number of words in $v_f$ that appear in the message.

Note that unlike the machine learning classifier, the user feedback classifier is, by definition, trained only on the feedback set. Due to the limited number of user-selected words in Experiment #1, computing the classification confidence of the user feedback classifier resulted in too many folders having the exact same assignment probability. This situation would also arise

in the real world, especially for during a user's early weeks of using a new email client. In order to avoid these tie-breaker conditions, we weighted the confidence of the user feedback classifier by the confidence of the machine learning classifier.

Pseudocode for the user co-training algorithm is shown in Figure A6 in Appendix A. The algorithm's variable $Confidence_f$ was computed as the posterior probability of Naive Bayes. In this experiment, we required $Score_m$ to be positive. Due to the fact that the results did not change much for k > 10, we set the value of k to be 10.

In user co-training, the user has an equal influence in classification and training as the machine learning algorithm. We used Naïve Bayes as the standard machine learning classifier. We also included Ripper as a machine learning classifier in the co-training experiment, but the results did not show any substantial improvement. We therefore do not discuss it further here.

As in Experiment #2, we divided the emails from the Enron farmer-d dataset into a training set, a feedback set and an evaluation set. This experiment required a large amount of data without user feedback; hence we needed each folder to have a large amount of data that the participants had not seen previously. Since the Bankrupt and Resume folders in the evaluation set did not have a substantial amount of unseen emails remaining, we needed to restrict the training, feedback and evaluation datasets to only include emails from two folders (Personal and Enron News). As a result, the training set contained 50 messages, the feedback set contained 11 messages, and the evaluation set contained 371 messages (88 emails in Enron News, 283 in Personal).

To evaluate the user co-training approach we compared it against two standard algorithms. The first standard algorithm is the simple Naïve Bayes algorithm that takes the corrected folder assignments provided by users into account, but not the user's algorithmic suggestions, i.e., the same as the baseline of Experiment #2. The second standard algorithm is the standard co-training algorithm, which likewise makes use of the corrected folder assignments but not the user's algorithmic suggestions. To test the performance of the standard co-training algorithm in this experiment, we combined the training set and the feedback set into a data set with user feedback in order to train two classifiers employing separate feature sets. The set of features for the first classifier were email addresses that appeared in the From, To, Cc and Bcc fields and in the body of an email message. The feature set for the second classifier used all words appearing in the subject line and the email body. In our experiment, we set the confidence parameter θ to 0.95. We experimented with θ values between 0.9 and 0.99 in increments of 0.1 and found little difference in the results. We also set the number of iterations (i) equal to 10; we noticed that the results did not change significantly after 10 iterations.

Figures 8 and 9 compare the simple Naïve Bayes, standard co-training and user co-training algorithms. Standard co-training did not improve on accuracy; in fact, it led to decreased performance in the case of similarity-based. We hypothesize that the two independent feature sets do not, by themselves, produce a complete "view" of the data, thereby resulting in impoverished classifiers that can mislead each other in standard co-training.

However, the user co-training approach we developed to represent the user's feedback as the co-trainer substantially increased accuracy for keyword-based and similarity-based user feedback when compared to the standard online training approach. User co-training was significantly more accurate than its closest runner-up, standard online training (McNemar, df=1, p<0.0001).



Figures 8 (left: Keyword) and 9 (right: Similarity): Percentage of emails the co-training-based algorithms sorted correctly. (Bars indicate standard error). User co-training approach, incorporating participants' rich feedback, performs usually better than algorithms without participants' rich feedback.

To understand just how these accuracy improvements came about, we investigated the user co-training approach more

deeply. (Full results of the evaluation are given in Appendix B). Figure 10 and 11 show the accuracy for the baseline online training approach compared with the user co-training approach for individual participants, taking into account keyword-based and similarity-based user feedback respectively. In Figure 10, the results from participants 1, 4, 10 and 11 are particularly noteworthy. For these participants, the standard online algorithm's accuracy was low (below 50%) yet substantial improvements were made through applying the user co-training approach. Similarly, participants 2, 3 and 13 benefited for similarity-based feedback (in Figure 11). There were a few cases, i.e., participant 9 for keyword-based and participants 7 and 8 for similarity-based, in which user co-training's performance was lower than the baseline approach, but for the remaining participants, user co-training was superior.

## 9. DISCUSSION OF EXPERIMENT #2 AND #3 RESULTS

The user's willingness to invest time providing information to an intelligent system can be viewed in terms of a cost-benefit-risk judgment by users, drawing on Attention Investment theory (Blackwell 2002). According to this theory, in weighing up whether to provide feedback, users' perceptions of costs, benefits, and risks each play key roles.

One of the charges of a good user interface is to provide the user with the information needed to make assessments of costs, benefits, and risks that are reasonably close to actual costs, benefits, and risks. For example, some costs can be controlled through careful design of the intelligent interfaces to explain the system's reasoning and to accept feedback. Risk can arise due to the intelligent system's inability to make profitable use of the user's suggestions, because the feedback in some circumstances actually makes accuracy worse. Benefits can be gained by the system's ability to make effective use of the user's suggestions to improve its accuracy.

Results of Experiments #2 and #3 provide some evidence as to the amount of potential costs, risks and benefits to users in providing feedback to machine learning systems. Experiment #2's results suggest that the constraint-based approach may not provide enough benefit through accuracy improvement to offset the user cost to give feedback in the first place. The users' selection of which and how many features to feedback may increase the risk and reduce the benefit. This is due to two reasons. First, participants were inclined to propose the most significant words present in the folder (for example, "enron" for the Enron News folder) and this resulted in redundant constraints that had already been learned by the system. Giving feedback of this kind gave no additional benefits. Second, the number of features proposed by participants made up only a fraction of the total number of features considered by the system and therefore were usually not resulting in enough change to correct decisions made by the classifier. This particular problem could be overcome by increasing the hardness of the constraints to ensure that the system better heeds the constraints. However, we found that picking the right hardness of constraints is very difficult. If constraints are too weak they will have almost no impact on the classifier. On the other hand, if the constraints are too hard then it could also lower the performance through overfitting. Overall, this approach only resulted in slight accuracy improvements. All these considerations of benefit against cost may make the constraint-based approach unsuitable from an Attention Investment perspective.

The user co-training approach appeared to return more benefit for the user's costs. However, although the user co-training approach produced very promising results, it is not without risk. There are situations in which performance decreases by taking user feedback into account, possibly due to introduction of human errors into the learned information. Further investigation is needed to evaluate how well this approach works in real-world situations.

Figures 10 (keyword: top) and 11 (similarity: bottom): Percentage of emails the user co-training algorithm sorted correctly for individual participants compared to the baseline approach. For both, co-training's improvement was highest in cases when improvement was needed most, i.e., when the baseline approach performed badly.

## 10. CONCLUDING REMARKS

The notion of humans interacting with a machine learning system, not just to override outcomes but rather to "rewire" the algorithms—changing the reasoning itself—is new. There are three components to this notion: (1) an intelligent system's ability to explain its reasoning to the user, (2) the user's reasoning corrections reflected in critiques and adjustments, and (3) the system's ability to make use of this rich user feedback in ways that are profitable for the system and ultimately for the user.

We have described three experiments that provide insights into all three of these components. Regarding the first component, an intelligent interface's ability to explain its reasoning to the user, we considered three explanation paradigms: Rule-based, Keyword-based, and Similarity-based. All were at least understandable enough that the participants were willing and reasonably effective at critiquing them. Rule-based explanations were the most understandable. Keyword-based were next, but the negative keywords list interfered. Similarity-based had serious understandability problems. But more important than which was the "winning" explanation paradigm is that there was no one winner: the lack of our participants' consensus suggests that multiple paradigms of explanations may need to be supported. We also reported data identifying the factors that won participant approval, which included perception of reasoning soundness, clear communication of reasoning, and informal wording ("Wow!").

Regarding the second component, the user's reasoning corrections made through critiques and adjustments, participants made a wide variety of reasoning suggestions, including reweighting features, feature combinations, relational features, and even wholesale changes to the algorithms. These suggestions were grounded in a variety of knowledge sources, such as knowledge of English, common-sense knowledge, and knowledge of the domain. Participants were more accurate than the machine—but they were not perfect, and occasionally made mistakes. This explicitly demonstrates the likelihood of rich user feedback introducing errors into the reasoning.

Regarding the third component, we were able to implement a subset of the participants' suggestions via two different algorithms, and empirically evaluated the resulting improvements in accuracy rates. The user co-training algorithm was the most effective, achieving substantial improvements in accuracy by taking into account the user suggestions. It had the highest accuracy on average, showing that although in some cases it may be outperformed by other approaches, in general it led to improvements in accuracy. In situations in which the standard online learning algorithm performed very badly, the user co-training approach led to the greatest improvement. These situations may be exactly the situations in which there is not enough training data present for standard machine learning approaches to succeed, and our results suggest that rich user feedback may provide the key to compensate for this lack of data.

Traditionally, machine learning approaches have been addressed primarily as an artificial intelligence problem, with users (when they are considered at all) as relatively passive spectators. We believe the design of intelligent systems needs to be viewed as a full-fledged HCI problem. Meaningful interaction involves both parties inseparably: approaches must be viable and productive for *both* the user and the system. This paper takes this view, considering both user and machine learning issues as inseparable parts of an intertwined whole. Overall, our results show evidence that intelligent interfaces can explain their reasoning and behavior to users, and that users in turn can provide rich, informative feedback to the intelligent system, and finally that machine learning algorithms can make use of this feedback to improve their performance. This suggests rich human-computer collaboration as a promising direction for machine learning systems to work more intelligently, hand-in-hand with the user.

## ACKNOWLEDGMENTS

## REFERENCES

Altendorf, E., Restificar, E., and Dietterich, T. 2005. Learning from sparse data by exploiting monotonicity constraints. Proc. Uncertainty in Artificial Intelligence.

Balcan,M.F., Blum, A., and Yang, K. 2005. Co-training and expansion: Towards bridging theory and practice. Proc. NIPS.

Billsus, D., Hilbert, D. and Maynes-Aminzade, D. 2005. Improving proactive information systems. Proc. IUI, 159-166.

Blackwell, A. First Steps in Programming: A Rationale for Attention Investment Models. 2002. Proc. Human Centric Computing Languages and Environments.

Blythe, J. 2005. Task learning by instruction in Tailor. Proc. IUI, 191-198.

Blum, A., and Mitchell, T. 1998. Combining labeled and unlabeled data with co-training. Proc. COLT.

Brutlag, J., Meek, C. 2000. Challenges of the email domain for text classification. Proc. ICML, 103-110.

Chapelle, O., Scholkopf, B., and Zien, A.. 2006. Semi-Supervised Learning. MIT Press, Cambridge, MA.

Chklovski, T., Ratnakar, V. and Gill, Y. 2005. User interfaces with semi-formal representations: A study of designing argumentation structures. Proc. IUI, 130-136.

Cohen, W. 1996. Learning rules that classify e-mail. Proc. AAAI Spring Symp. Information Access.

Crawford, E., Kay, J., and McCreath, E. 2002a. An Intelligent Interface for Sorting Electronic Mail. Proc. IUI, 182-183.

Crawford, E., Kay, J., and McCreath, E. 2002b. IEMS – The Intelligent Email Sorter. Proc. ICML, 83-90.

Culotta, A. Kristjansson, T. McCallum, A. and Viola, P. 2006. Corrective Feedback and Persistent Learning for Information Extraction Artificial Intelligence, 170, 1101-1122.

Dalvi, N., Domingos, P., Sanghai, M. S. and Verma, D. 2004. Adversarial classification. Proc. Intl. Conf. Knowledge Discovery and Data Mining, 99-108.

Fails, J. A. and Olsen, D. R. 2003. Interactive machine learning. Proc. IUI, 39-45.

Fung, G., Mangasarian, O. and Shavlik, J. 2002. Knowledge-based support vector machine classifiers. Proc. NIPS.

Hart, S. and Staveland, L. 1988. Development of a NASA-TLX (Task load index): Results of empirical and theoretical research, in: Hancock, P. and Meshkati, N. (eds.), Human Mental Workload, 139-183.

Herlocker, J., Konstan, J. and Riedl, J. 2000. Explaining collaborative filtering recommendations. Proc. CSCW, 241-250.

Huang, Y. and Mitchell, T. M. 2006. Text clustering with extended user feedback. Proc. SIGIR, 413-420.

Kiritchenko, S. and Matwin, S. 2001. Email classification with co-training. Proc. Centre for Advanced Studies Conference, 8-17

Klimt, B. and Yang, Y. 2004. The Enron corpus: A new dataset for email classification research. Proc. European Conf. Machine Learning 2004, 217-226.

Lieberman, H. and Kumar, A. 2006. Providing expert advice by analogy for on-line help. Proc. Intl. Conf. Intelligent Agent Technology, 26-32.

Lindo Systems, Inc., 2007, http://www.lindo.com.

Liu, B. Li, X. Lee, W. and Yu, P. 2004. Text Classification by Labeling Words. Proc. AAAI.

McCarthy, K., Reilly, J., McGinty, L. and Smyth, B. 2005. Experiments in dynamic critiquing. Proc. IUI, 175-182.

Miller, G. 1995. WordNet: A lexical database for English. Comm. ACM 38(11), 39-41.

Mitchell, T. 1997. Machine Learning. McGraw-Hill, Boston, MA.

Myers, B., Weitzman, D., Ko, A. Chau, D. 2006. Answering why and why not questions in user interfaces. Proc. CHI.

Oblinger, D., Castelli, V. and Bergman, L. 2006. Augmentation-based learning. Proc. IUI, 202-209.

Pazzani, M. J. 2000. Representation of electronic mail filtering profiles: a user study. Proc. IUI 2000, 202-206.

Phalgune, A., Kissinger, C., Burnett, M., Cook, C., Beckwith, L. Ruthruff, J. 2005. Garbage in, garbage out? An empirical look at oracle mistakes by end-user programmers. Proc. Symp. Visual Languages and Human Centric Computing, 45-52.

Porter, M. 1980. An algorithm for suffix stripping. Program, 14(3), 130-137.

Pu, P. and Chen, L. 2006. Trust building with explanation interfaces. Proc. IUI, 93-100.

Rettig, M. 1994. Prototyping for tiny fingers. Comm. ACM 37(4), 21-27.

Shen, J., Li, L., Dietterich, T. Herlocker, J. 2006. A hybrid learning system for recognizing user tasks from desk activities and email messages. Proc. IUI 2006, 86-92.

Stumpf, S., Rajaram, V., Li, L., Burnett, M., Dietterich, T. Sullivan, S., Drummond, R., Herlocker, J. 2007. Toward Harnessing User Feedback For Machine Learning. Intelligent User Interfaces, Proc. IUI 2007, 82-91.

Tullio, J., Dey, A., Chalecki, J., and Fogarty, J. 2007. How it works: A field study of non-technical users interacting with an intelligent system, Proc. CHI, 31-40.

Ware, M., Frank, E., Holmes, G., Hall, M., Witten, I. H. 2001. Interactive machine learning: letting users build classifiers. IJCHS, 55, 281-292.

Witten, I., Frank, E. 2005. Data Mining: Practical Machine Learning Tools and Techniques, 2nd Ed., Morgan Kaufmann, 2005.

Zhou, G., Su, J. 2002. Named entity recognition using an HMM-based chunk tagger. Proc. Assoc. Comp. Linguistics.

# Appendix A: Algorithm Details

## A1. NAÏVE BAYES CLASSIFIER

The goal of the Naïve Bayes classifier is to compute $P(Y = y_k \mid X)$

Using Bayes rule, we get:

$$P(Y = y_k \mid X) = \frac{P(X \mid Y = y_k)P(Y = y_k)}{P(X)}$$

Assuming that probability of the each feature $x_i$ is conditionally independent given the folder Y, the formula above can be rewritten as:

$$P(Y = y_k \mid X) = \frac{P(X \mid Y = y_k)P(Y = y_k)}{P(X)} = \frac{\prod_{j=1}^{n} P(x_j \mid Y = y_k)P(Y = y_k)}{P(X)}$$

The optimal decision is the $y_k$ that maximize the posterior probability $P(Y = y_k \mid X)$. This is equivalent to finding out the maximum value of the ratio:

$$Ratio = \frac{P(Y = y_k \mid X)}{P(Y \neq y_k \mid X)} = \frac{P(X \mid Y = y_k)P(Y = y_k)}{P(X \mid Y \neq y_k)P(Y \neq y_k)} = \frac{\prod_{j=1}^{n} P(x_j \mid Y = y_k)P(Y = y_k)}{\prod_{j=1}^{n} P(x_j \mid Y \neq y_k)P(Y \neq y_k)}$$

The log of the ratio is:

$$\log \frac{P(Y = y_k \mid X)}{P(Y \neq y_k \mid X)} = \log \frac{P(x_1 \mid Y = y_k)}{P(x_1 \mid Y \neq y_k)} + \log \frac{P(x_2 \mid Y = y_k)}{P(x_2 \mid Y \neq y_k)} + \cdots + \log \frac{P(x_n \mid Y = y_k)}{P(x_n \mid Y \neq y_k)} + \log \frac{P(Y = y_k)}{P(Y \neq y_k)}$$

In this way, the Naïve Bayes classifier is transformed to a linear classifier where the input is a vector of Boolean variables and the weight $w_{jk}$ for each word $j$ and each email folder $k$ is the log ratio. Hence, to predict the email folder, it computed a score for folder $k$ as

$$score(k) = \sum_j w_{jk} \cdot x_j$$

where $x_j$ was an indicator variable representing the presence of the $j$th word in the email message. The Naïve Bayes algorithm then predicted the folder with the highest score.

Figure A1: The mathematical details of the Naïve Bayes classifier

## A2. CONSTRAINT REASONING DETAILS

Constraints of type 2 were represented by constraining weights to be positive as follows:

$$w_{ik} > 0$$

$$= \log \frac{P(x_j = 1 | Y = y_k)}{P(x_j = 1 | Y \neq y_k)} > 0$$

$$= \log P(x_j = 1 | Y = y_k) > \log P(x_j = 1 | Y \neq y_k)$$

Where:

$w_{jk}$ is the weight associated with the word $j$ and the index $k$ of the user-assigned folder.

$Y$ is the random variable representing the folder

$y_k$ is the value of the random variable $Y$ ie. the folder that the user assigned the email to

$x_j$ is a Boolean variable indicating the presence or absence of the $j$th word.

If the log of the ratio was greater than 0, then the conditional probability in the numerator was greater than the conditional probability in the denominator. Therefore, if this inequality held, the probability of the word $x_j$ being present in the email given that the user assigned the email to folder $y_k$ was greater than the probability of the word $x_j$ being present in the email given that the user assigned the email to any other folder than $y_k$. Consequently, the presence of the word $x_j$ had more of an effect when classifying the email to folder $y_k$. To increase the effect of this constraint, we required the inequality to hold above some amount $\delta \geq 0$ ie. $\log P(x_j = 1 | Y = y_k) > \log P(x_j = 1 | Y \neq y_k) + \delta$.

Figure A2: Constraint Type 2.

Constraints of type 3 were expressed as an inequality of the following form:

$$P(Y = y_k | x_j = 1) > P(Y = y_k | x_k = 1) + \theta$$

Where:

$Y$ is the random variable for the folder

$y_k$ is the folder that the email was assigned to

$x_j$ is a word selected by the user

$x_k$ is one of the top 10 words appearing in the email message that was a strong predictor of the user-assigned folder for the email message but was not selected by the user.

$\theta$ is a parameter controlling the hardness of the constraint.

For each user-selected word, we added 10 constraints of this form. These constraints were determined by sorting the conditional probabilities $P(Y = y_k | x_k = 1)$ for all possible words $x_k$ and creating a constraint $P(Y = y_k | x_j = 1) > P(Y = y_k | x_k = 1) + \theta$ for each word $x_k$ in the top 10 largest conditional probabilities in this sorted list.

Figure A3: Constraint Type 3

> *Maximize:*
> > The likelihood of the Naïve Bayes model on the all the samples in the training set and feedback set
>
> *Subject To*:
> > For each message in the feedback set,
> > > For all words $x_j$ that the user proposed or increased its weight in the message,
>
> $$P(x_j = 1 \mid Y = y_k) > P(x_j = 1 \mid Y \neq y_k) + \delta \qquad \text{[Constraint type 2]}$$
>
> $$P(Y = y_k \mid x_j) > P(Y = y_k \mid x_k) + \theta \qquad \text{[Constraint type 3]}$$
>
> *Where*:
> > $y_k$ is the user assigned folder
> >
> > $x_k$ is one of the words that user did not select but is in the top 10 of $P(Y = y_k \mid x_k)$

Figure A4: The Constraint-based algorithm.

## A3. CO-TRAINING ALGORITHMS

> Create two classifiers $C_1$ and $C_2$ based on the two independent feature sets.
> Repeat *i* times
> > For each message in the evaluation set
> > Classify the message with the two classifiers
> > If classification confidence of any classifier is > θ
> > > Add the classified message to the training data
> > Rebuild $C_1$ and $C_2$ with the new training data

Figure A5: The standard co-training algorithm. In the Blum and Mitchell 1998 version of co-training, after training the two classifiers on the labeled data, the top *N* most confidently classified negative data points and the top *P* most confidently classified positive data points are added to the training set. A minor difference in our implementation is the fact that we add all data points that are classified above some confidence threshold θ to the training data.

> For each folder *f*, create a vector $v_f$ to store the important words proposed by the user
> For each message *m* in the unlabeled set
> > For each folder *f*,
> > > Compute **Confidence**$_f$ from the machine learning classifier
> > > **FolderScore**$_f$=# of words in $v_f$ appearing in the message $\times$ **Confidence**$_f$
> > Find the folder $f_{max}$ that has largest **FolderScore**$_f$ in all folders.
> > Let **FolderScore**$_{other}$ be the **FolderScore** of the other folder (since there are only 2 folders)
> > Save the message $Score_m$=**FolderScore**$_{fmax}$ – **FolderScore**$_{other}$
> Sort $Score_m$ for all messages in decreasing order
> Select the top *k* messages to add to the training set along with their folder label $f_{max}$

Figure A6: Our user co-training algorithm.

# Appendix B: Detailed Results Comparing Learning Approaches to Take User Feedback into Account

## B1. NUMBER OF MESSAGES WITH FEEDBACK

The number of emails processed by each participant:

| | Participant | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Type of Feedback** | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| Keyword-based | 10 | 5 | 5 | 5 | 3 | 2 | 9 | 7 | 5 | 3 | 7 | 5 | 6 |
| Similarity-based | 9 | 5 | 8 | 7 | 5 | 4 | 7 | 6 | 4 | 2 | 5 | 0 | 4 |
| Rule-based | 8 | 1 | 3 | 4 | 3 | 1 | 2 | 5 | 1 | 1 | 4 | 2 | 9 |

## B2. RESULTS FOR RULE-BASED USER FEEDBACK

### B2.1 Applied to 2 Folders Only

| | Participant | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| **Algorithm** | 1 | 4 | 5 | 7 | 8 | 9 | 10 | 11 | 13 |
| Simple online training | 0.73 | 0.73 | 0.70 | 0.73 | 0.70 | 0.70 | 0.73 | 0.73 | 0.73 |
| Standard co-training | 0.72 | 0.70 | 0.70 | 0.72 | 0.70 | 0.70 | 0.72 | 0.70 | 0.70 |
| User co-training | 0.73 | 0.73 | 0.70 | 0.72 | 0.73 | 0.70 | 0.73 | 0.73 | 0.73 |

## B3. RESULTS FOR KEYWORD-BASED USER FEEDBACK

### B3.1 Applied to All 4 Folders

| | Participant | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Algorithm** | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| Simple online training | 0.65 | 0.65 | 0.55 | 0.53 | 0.54 | 0.73 | 0.56 | 0.54 | 0.63 | 0.64 | 0.67 | 0.71 | 0.69 |
| Constraint-based($\delta=0,\theta=0$) | 0.65 | 0.65 | 0.57 | 0.53 | 0.58 | 0.71 | 0.59 | 0.53 | 0.63 | 0.65 | 0.68 | 0.71 | 0.69 |
| Constraint-based($\delta=0.2,\theta=0$) | 0.65 | 0.65 | 0.57 | 0.55 | 0.58 | 0.71 | 0.59 | 0.53 | 0.63 | 0.65 | 0.67 | 0.71 | 0.69 |
| Constraint-based($\delta=0.4,\theta=0$) | 0.65 | 0.66 | 0.57 | 0.56 | 0.58 | 0.71 | 0.59 | 0.53 | 0.62 | 0.65 | 0.67 | 0.71 | 0.69 |
| Constraint-based($\delta=0.6,\theta=0$) | 0.65 | 0.67 | 0.57 | 0.56 | 0.58 | 0.71 | 0.59 | 0.54 | 0.63 | 0.65 | 0.67 | 0.72 | 0.70 |
| Constraint-based($\delta=0,\theta=0.2$) | 0.65 | 0.65 | 0.57 | 0.53 | 0.58 | 0.71 | 0.59 | 0.53 | 0.63 | 0.65 | 0.68 | 0.71 | 0.69 |
| Constraint-based($\delta=0.4,\theta=0.4$) | 0.65 | 0.66 | 0.57 | 0.56 | 0.58 | 0.71 | 0.59 | 0.53 | 0.62 | 0.65 | 0.67 | 0.71 | 0.69 |

### B3.2 Applied to 2 Folders Only

| | Participant | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Algorithm** | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
| Simple online training | 0.44 | 0.65 | 0.82 | 0.44 | 0.82 | 0.65 | 0.82 | 0.82 | 0.65 | 0.44 | 0.44 | 0.65 | 0.65 |
| Standard co-training | 0.32 | 0.73 | 0.83 | 0.32 | 0.83 | 0.73 | 0.83 | 0.83 | 0.73 | 0.32 | 0.32 | 0.73 | 0.73 |
| User co-training | 0.71 | 0.73 | 0.83 | 0.72 | 0.81 | 0.84 | 0.84 | 0.85 | 0.51 | 0.79 | 0.81 | 0.77 | 0.75 |

# B4. RESULTS FOR SIMILARITY-BASED USER FEEDBACK

## B4.1 Applied to All 4 Folders

| Algorithm | Participant | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 13 |
| Simple online training | 0.57 | 0.57 | 0.59 | 0.66 | 0.69 | 0.55 | 0.60 | 0.68 | 0.53 | 0.57 | 0.70 | 0.62 |
| Constraint-based($\delta=0,\theta=0$) | 0.59 | 0.57 | 0.59 | 0.66 | 0.68 | 0.55 | 0.60 | 0.68 | 0.53 | 0.57 | 0.69 | 0.64 |
| Constraint-based($\delta=0.2,\theta=0$) | 0.59 | 0.57 | 0.60 | 0.66 | 0.68 | 0.55 | 0.60 | 0.68 | 0.53 | 0.57 | 0.69 | 0.62 |
| Constraint-based($\delta=0.4,\theta=0$) | 0.59 | 0.57 | 0.61 | 0.66 | 0.68 | 0.55 | 0.62 | 0.68 | 0.53 | 0.57 | 0.69 | 0.63 |
| Constraint-based($\delta=0.6,\theta=0$) | 0.58 | 0.58 | 0.61 | 0.66 | 0.69 | 0.55 | 0.62 | 0.68 | 0.53 | 0.57 | 0.69 | 0.64 |
| Constraint-based($\delta=0,\theta=0.2$) | 0.58 | 0.57 | 0.59 | 0.66 | 0.68 | 0.55 | 0.60 | 0.68 | 0.53 | 0.57 | 0.69 | 0.63 |
| Constraint-based($\delta=0.4,\theta=0.4$) | 0.58 | 0.57 | 0.61 | 0.66 | 0.68 | 0.55 | 0.62 | 0.68 | 0.53 | 0.57 | 0.69 | 0.63 |

## B4.2 Applied to 2 Folders Only

| Algorithm | Participant | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 13 |
| Simple online training | 0.82 | 0.44 | 0.44 | 0.65 | 0.65 | 0.82 | 0.44 | 0.65 | 0.82 | 0.82 | 0.65 | 0.44 |
| Standard co-training | 0.83 | 0.32 | 0.32 | 0.73 | 0.73 | 0.83 | 0.32 | 0.73 | 0.83 | 0.83 | 0.73 | 0.32 |
| User co-training | 0.84 | 0.79 | 0.75 | 0.81 | 0.62 | 0.81 | 0.30 | 0.33 | 0.85 | 0.84 | 0.57 | 0.67 |