



End-user feature labeling: Supervised and semi-supervised approaches based on locally-weighted logistic regression

Shubhomoy Das^{a,*}, Travis Moore^a, Weng-Keen Wong^a, Simone Stumpf^b,
Ian Oberst^a, Kevin McIntosh^a, Margaret Burnett^a

^a Oregon State University, OR, USA

^b City University London, UK

ARTICLE INFO

Article history:

Received 16 February 2012

Received in revised form 16 July 2013

Accepted 22 August 2013

Available online 30 August 2013

Keywords:

Feature labeling

Locally-weighted logistic regression

Machine learning

Intelligent interfaces

Semi-supervised learning

ABSTRACT

When intelligent interfaces, such as intelligent desktop assistants, email classifiers, and recommender systems, customize themselves to a particular end user, such customizations can decrease productivity and increase frustration due to inaccurate predictions—especially in early stages when training data is limited. The end user can improve the learning algorithm by tediously labeling a substantial amount of additional training data, but this takes time and is too ad hoc to target a particular area of inaccuracy. To solve this problem, we propose new supervised and semi-supervised learning algorithms based on locally-weighted logistic regression for *feature labeling by end users*, enabling them to point out which features are important for a class, rather than provide new training instances. We first evaluate our algorithms against other feature labeling algorithms under idealized conditions using feature labels generated by an oracle. In addition, another of our contributions is an evaluation of feature labeling algorithms under real-world conditions using feature labels harvested from actual end users in our user study. Our user study is the first statistical user study for feature labeling involving a large number of end users (43 participants), all of whom have no background in machine learning. Our supervised and semi-supervised algorithms were among the best performers when compared to other feature labeling algorithms in the idealized setting and they are also robust to poor quality feature labels provided by ordinary end users in our study. We also perform an analysis to investigate the relative gains of incorporating the different sources of knowledge available in the labeled training set, the feature labels and the unlabeled data. Together, our results strongly suggest that feature labeling by end users is both viable and effective for allowing end users to improve the learning algorithm behind their customized applications.

© 2013 Elsevier B.V. All rights reserved.

1. Introduction

Many applications, powered by machine learning, customize themselves to a particular end user's preferences. Such applications include email classifiers, recommender systems, intelligent desktop assistants, and other intelligent user interfaces. To accomplish this customization, the application must learn from the particular end user—which obviously cannot happen until *after* the system is deployed and training data from that specific end user is obtained.

Customizing to the end user's preferences is challenging, especially when there is limited training data, such as when the application is first deployed. The end user could select additional training instances to label, or the learning algorithm could

* Corresponding author at: 1148 Kelley Engineering Center, Corvallis, OR 97331-5501, USA. Tel.: +1 541 737 3617.

E-mail address: dassh@eecs.oregonstate.edu (S. Das).

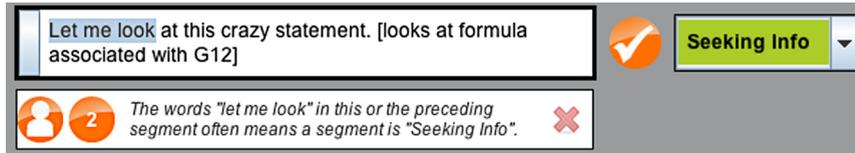


Fig. 1. The user is pointing out that the feature “let me look” is highly indicative of the class “Seeking Info”. (This UI inspired the development of the algorithms we present in this paper.)

ask the user to provide class labels for strategically chosen instances that would most inform the learning algorithm, as is done in traditional active-learning [7,31]. Labeling instances, however, has its drawbacks. First, labeling data instances is a tedious process and a substantial number of instances must often be labeled before a change to the learning algorithm is noticeable to an end user. Second, in a streaming data setting, such as news filtering or email classification, active-learning is not applicable as the system has no control over which data instance arrives next. Finally, if a rare group of instances is incorrectly classified, the learning algorithm cannot be “corrected” until the user labels instances with this rare combination of attributes. Since this group is rare, the cost, in terms of time or effort, to acquire such data instances could be very expensive [1].

To overcome these problems, in this paper¹ we investigate the possibility of end-user *feature labeling* [29,10,33,1], namely allowing end users to label features instead of instances. Here, the term *feature* refers to an attribute of a data instance that is useful for predicting the class label; for example, rather than labeling entire documents, an end user could point out which words (features) in the document are most indicative of certain class labels. Fig. 1 shows this approach in our formative research’s user interface [15], which allowed HCI researchers to point out words that were predictive of that transcript segment’s label. Raghavan et al. [28,29] found that labeling a feature took humans roughly a fifth of the time as a document and the benefits of feature labeling were greatest when the training set sizes were small. However, their work did not evaluate feature labeling in a statistical user study involving a large number of actual end users.

Allowing end users, who are not likely to be educated in machine learning, to use feature labeling introduces new challenges to learning algorithms. End users’ choices of features may be noisy, inconsistent, and might vary greatly in ability to improve the predictive power of the machine learning algorithm. This paper therefore investigates algorithms able to stand up to these challenges.

Our research contributions are as follows. First, we present a new supervised learning algorithm for taking end-user feature labels into account, based on Locally-Weighted Logistic Regression. In order to evaluate our feature labeling algorithm, we perform an empirical comparison on multiple data sets under ideal conditions, using feature labels obtained from an oracle, and under real-world conditions for one particular dataset, using feature labels harvested from actual end users. For the latter study, we present a user study in which ordinary end users, unfamiliar with machine learning, chose the feature labels themselves—with no restrictions as to what they could select as features. Our algorithm was among the best performing feature labeling algorithms in the idealized setting and it was also robust to poor quality feature labels provided by ordinary end users in our study.

Next, we present a semi-supervised version of our feature labeling algorithm which assumes that an unlabeled set of instances is present during training. The semi-supervised setting for feature labeling incorporates knowledge from three sources: a small labeled training set, the feature labels provided by the end user and information from the implicit structure of the unlabeled data. We evaluate our semi-supervised algorithm using both oracle feature labels and end-user feature labels from the user study mentioned above. Our feature labeling algorithm is one of the best performing algorithms with oracle feature labels and the best performer with lower quality feature labels from end users. With our results, we can compare the relative gains using the different sources of knowledge available in the training set, the feature labels, and the unlabeled data. Our analysis shows that incorporating unlabeled data during learning sometimes produces worse performance than just using a purely supervised learning approach, both with and without feature labeling. However, adding the information from feature labels consistently improves performance over not including this information, both in the supervised and semi-supervised settings.

Together, our results strongly suggest that feature labeling by end users is both a viable and an effective solution for allowing end users to improve the learning algorithm behind their customized user interface.

2. Related work

We divide the approaches for feature labeling into supervised and semi-supervised feature labeling algorithms. Supervised feature labeling algorithms require only a training set of labeled instances. On the other hand, semi-supervised feature labeling requires both a labeled training set as well as a pool of unlabeled data, which is assumed to be relatively easy to obtain.

Two of the SVM-based methods presented by Raghavan and Allen [29] involved supervised feature labeling. Their Method 1 scaled features indicated as relevant by the user by a constant a and the rest of the features by a constant d

¹ Early versions of portions of this work appeared in [36,37].

متن کامل مقاله

دریافت فوری ←

ISIArticles

مرجع مقالات تخصصی ایران

- ✓ امکان دانلود نسخه تمام متن مقالات انگلیسی
- ✓ امکان دانلود نسخه ترجمه شده مقالات
- ✓ پذیرش سفارش ترجمه تخصصی
- ✓ امکان جستجو در آرشیو جامعی از صدها موضوع و هزاران مقاله
- ✓ امکان دانلود رایگان ۲ صفحه اول هر مقاله
- ✓ امکان پرداخت اینترنتی با کلیه کارت های عضو شتاب
- ✓ دانلود فوری مقاله پس از پرداخت آنلاین
- ✓ پشتیبانی کامل خرید با بهره مندی از سیستم هوشمند رهگیری سفارشات