

Using Active Learning with Text Filtering to Generate a Support Vector Machine Training Set

Joseph M. Geyer

23 March 2009

Abstract

The need to understand software systems is an important part of their update and maintenance. If one does not understand a software system, he/she will have difficulty modifying, maintaining, or updating it. This can be costly in terms of both time and money. Reverse software engineering alleviates this by creating models of a system to aid system comprehension. A well known problem in this domain is the concept assignment problem [2]. This is the task of assigning human level concepts or meaning to the code that actuates it. This problem can be extended to the class level where the goal is to assign concepts to classes. Carey and Gannod [6] have automated the classification of the concept classes in a software system by using support vector machines. Support vector machines require a training set to train the learner, however, manually labeling this training set is inefficient. The goal of the proposed research is to present a method and tool to semi-automate the creation of a training set for support vector machine learning in the context of reverse software engineering.

1 Introduction

The documentation of the design of legacy software systems is often neglected as the system ages. It is easy for updates and modifications to be made to a software system without those changes reflected in the design models. Knowing the high level architecture of a software system is critical to its continued efficacy. Chikofsky and Cross [7] state that reverse engineering is the process of examining a software systems to extract design knowledge in order to facilitate maintenance and future updates. There exists an increasing demand for the ability to create models and blueprints of the design of software systems.

Dietterich [8] defines machine learning as “the study of methods for programming computers to learn”. There are many applications for machine learning. Data mining uses machine learning algorithms to find knowledge from data. For example, data mining medical records can lead to medical knowledge [12]. Some other examples are speech recognition, spam classification, autonomous driving, and book recommendation programs. In each of these examples, the actions of the computer are not explicitly programmed. Learning can be thought of as the process of training a function to correctly give an output based on the previous examples it was given. A binary classification learning problem is where each data example is in one of two distinct classes. Thus, the learning function is being trained to take as input the features of the example and to give as output the classification to which the input features is mapped. In our problem, the training

examples are classes in a software system where the input features are object oriented metrics about the class. The output is whether that example should be classified as a concept class or a non-concept class.

Ontology creation from text has its roots in text mining. Often, clustering is used [5, 19, 20, 11, 10] to group words into groups. Hierarchy can be added to these clusters by parsing the text and looking for hierarchical phrasing. For example, the phrase “...*skeleton, bobsleigh, and other winter olympic sports...*” indicates that *skeleton* and *bobsleigh* are types of *sports*. It also indicates that *olympic sport* is a type of *sport* and that *winter olympic sport* is a type of *winter sport*. [5] Key concepts from the text can then be readily accessed from the ontology that was created.

A well known problem in the domain of reverse software engineering is the concept assignment problem [2]. This is the task of assigning human level concepts to the code that actuates it. Carey and Gannod [6] have developed a tool to automate the classification of the concept classes in a software system by using support vector machines. However, this requires a great amount of work to manually label a training set of classes as either concept or non-concept to train the learner. The goal of the research presented in this paper is to present a method to reduce the load on this aspect of training the support vector machine

2 Background and Related Work

2.1 Background

In the domain of reverse software engineering, the concept assignment problem is a classic task for recovering design rationale. Another important technique is that of active learning where the user observes the results of the software tool and indicates whether it was correct or not. Finally, in machine learning supervised learning is used successfully in many cases. We are interested in classification area of machine learning where we can classify an item in one of two different classes.

2.1.1 Machine Learning

Machine learning uses algorithms so that computers can solve problems without being explicitly programmed. The machine changes its behavior to gain improved performance for subsequent iterations of the problem. [12]. Machine learning is a specialized field of artificial intelligence and can be broken up into further categories. Supervised machine learning is where training samples are labeled with the appropriate output. If the output is continuous, it is known as a regression problem. If the output is in the range of a small number of discrete variables, it is a classification problem. For our problem, we will focus on this - classification supervised learning. Unsupervised learning is where no knowledge is given about the output of the training samples. A common solution for this type of problem is to use clustering. Unsupervised learning problems can also be regression or classification. Reinforcement learning is where a sequence of decisions are made with a reward function. An example of this is training a robot to drive a car. There is not a single classification being done for this task, or a single target value to achieve. The goal is a series of acceptable actions for a larger goal.

For our research, we are interested in supervised learning for a binary classification problem. This means that we have a training set with labeled outputs that can be used to train the learner. In a general sense, the problem can be defined in the following way. The variable x is a vector of input variables and y is the output variable and can either be a value from $\{-1, 1\}$. The vector x consists of any number of features

depending on the problem. Thus, the feature vector or input vector with n features is

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

The ordered pair $(x^{(i)}, y^{(i)})$ is called a training examples. The training set is composed of m training examples so that i has a range from 1 to m .

$$trainingset = \begin{bmatrix} (x^{(1)}, y^{(1)}) \\ (x^{(2)}, y^{(2)}) \\ \vdots \\ (x^{(m)}, y^{(m)}) \end{bmatrix}$$

We want to use the training set to learn a function $h : X \rightarrow Y$, called the hypothesis, where X is the space of input variables and Y is the space of output variables. This mapping should accurately classify a sample based on the value of its feature vector. The goal of learning this function is to find the parameter vector θ that can be used to get accurate performance.

One way to understand support vector machines is to visualize the problem geometrically. Each training example in the training set can be placed in a hyperplane of degree n where n is the number of features in the feature vector x . Each example point has output of either $+1$ or -1 and can be identified with some binary identification scheme, like using a circle for $+1$ or a triangle for -1 . If the two sets of points are linearly separable, a hyperplane can be used to separate the samples into the two classifications. But, there are infinitely many separating hyperplanes. The goal then is to find the hyperplane that gives the maximal distance between the closest vectors on either side. These vectors that are closed are called the support vectors. This problem becomes more complicated when the training examples are not linearly separable. However, by projecting these samples into a higher dimensional hyperplane, the points can be linearly separable.

2.1.2 Reverse Engineering and the Concept Assignment Problem

In order to maintain or update a software systems, it is important to understand it. Reverse software engineering takes a software systems and extracts knowledge about it's design to further the understanding of the system. Chikofsky and Cross [7] define reverse software engineering as

- identify the system's components and their interrelationships and
- create representations of the system in another form or at a higher level of abstraction

With the number and complexity of legacy software systems, there is a real demand for reverse software engineering.

The famous concept assignment problem presented by Biggerstaff [2] involves two parts. The first is to identify the real human concepts that are in a software system. For example, a concept might be to

“calculate revenue”. Then the user assigns that concept (calculate revenue) to the code that accomplishes that concept. This was first described on a line-by-line micro level of the code. However the principle still applies to a larger scale when we consider the class level. In Carey and Gannod’s [6] work, they consider classes as either concept or non-concept classes. They assigned each class a vector of object oriented metrics and used machine learning to classify the concept and non-concept classes. They used supervised learning with cross validation to assess the validity of their classifier.

2.1.3 Active Learning

Active learning is where the user can play a part in the acceptance or rejection of the classifications of the learner. This was successfully used in the the work of Bowring et al. [3]

2.2 Related Work

Sartipi did some research on reverse engineering with machine learning. Others also...

2.2.1 Reverse Engineering through Data Mining

TODO: put Sartipi [14, 13] here also Shirabad [16] also

2.2.2 Automatic Training Set Creation

TODO: Tang [17] developed a method for automatically creating a training set for video annotation.

2.2.3 Ontology creation through Clustering

Clustering is a method used for ontology building from text [5, 19, 20, 11, 10] . Early work involved simpler similarity metrics and algorithms while current work is much more sophisticated. The general idea is to build the ontology bottom up. After identifying concept terms, a similarity metric is used to measure the closeness of terms. Depending on the closeness, the concept terms will be clustered together. These clusters can represent a new concept. The process continues with these new clusters, building a hierarchy.

Caraballo [5] developed an algorithm for automatically creating a hypernym-labeled noun hierarchy from text corpus. Using the Wall Street Journal, she identified 50,000 nouns. When two or more nouns are used together in a conjunction or an appositive, the assumption is that they are semantically related. Each noun is given a vector where the values are the number of times each other noun is found in either a conjunction or an appositive in the text. Similarity between nouns is then calculated by comparing the angle between each pair of vectors. The most similar nouns are clustered into a new parent node, and the process continues. Hierarchy is then extracted by looking for the use of the word *other* with nouns in a conjunction clause. For example, *birds, squirrels, and other small mammals* indicates that mammal is a hypernym of birds and squirrels. This algorithm correctly assigned hyponyms to randomly selected hypernyms with a 33% accuracy according to human judges participating in the experiment. One of the top three hypernyms that the algorithm produced for a noun matched the judge’s hypernym with a 60% accuracy.

Instead of a bottom up approach, Yang and Callan [18] proposed an incremental approach where each term is considered and placed in the correct spot in the hierarchy.

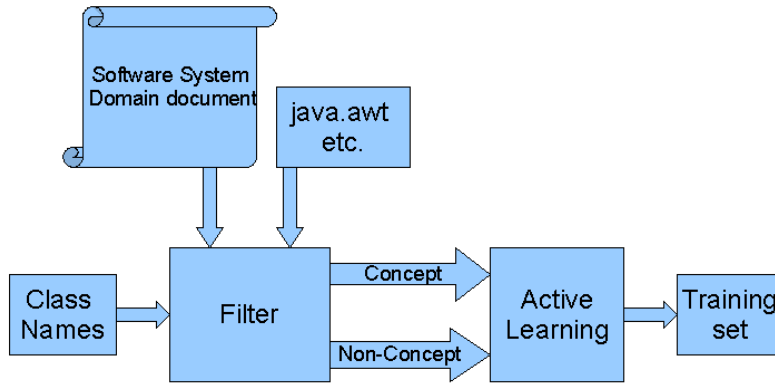


Figure 1: Flowchart Model of Proposed Research

3 Proposed Research

This section discusses the research to be conducted. It starts by defining the research problem and the the proposed solution. Next it examines the feasibility and validity measures of the proposed solution. This section also includes some preliminary results conducted on a UML Modeling software system. Finally, a timetable for completing the research is presented.

3.1 Specific Research Problem

Classification problems with supervised learning algorithms allow the machine to create a learner from the training set that will be applied to the rest of the data set. The user must manually classify a subset of the data so that the learner can automatically classify the rest of the instances in the data set. As an example, if a researcher wants to create a training set that is 10% of a software system with 5000 classes, he/she would have to manually classify 500 classes. This creates a bottleneck in the productivity of the reverse software engineering process. Through this research, we want to reduce the load of manually classifying the training set.

3.2 Proposed Solution and Methodology

The proposed solution will be implemented in an eclipse plugin that utilizes filtering and active learning techniques. This plugin will be an extension on the plugin by Carey et al. [6] The flowchart in figure 1 is shows a model of the proposed solution.

3.2.1 Training Set Selection via Filtering

We start by using a filtering technique. In our problem, every class is either a *concept class* or *other*. The SVM learner created by Carey [6] requires the user to manually classify at least one positive example and at least one negative example. The goal of this step is to automatically filter out probable concept classes and

those classes that are most likely not a concept class. Software systems have a concept domain associated with its high level functioning. Some domain examples are given below.

- Database
- IDE
- Data Mining
- UML Modeling
- Chat Client
- Text Editor
- Online Store
- Chat Client
- Project Manager

In the UML Modeling, some examples of conceptual terms are *inheritance*, *use case*, and *association*. These terms can be taken directly out of a textbook or other document on UML. The same is true of other domains as well. We can extract these terms from any document (user's guides, manuals, textbooks, tutorials, etc.) associated with the concept domain of the system. These keywords together make up the *concept filter*. Likewise there are other terms that are usually associated with the inner workings of a software system. We compiled a list of 98 terms that would filter out classes that are probably not concept classes. These terms are from the java.awt package. Some examples are *JButtonItem*, *Scrollbar*, and *JTextPane*. These non-concept keywords are used in the *non-concept filter*.

The filtering of the classes takes place by examining matches of substrings of the software systems class names with substrings elements of the *concept filter* and the *non-concept filter*. If a match occurs on the *concept filter*, then the class name is placed in a list of probable concept classes. If a match occurs on the *non-concept filter*, then the class name is placed in a list of probable non-concept classes. If a class does not get matched by either filter, then it is not placed in either list. In this way, we have reduced the number of classes to choose from for manual classification.

3.2.2 Training Set Selection via Active Learning

Next we allow the user to accept or reject the automatic filtering. Figure 1 shows the concept and non-concept classes as input into the active learning module. Here the user will be presented with the list of probable concept classes generated by the *concept filter*. The user accepts individual classes as concept by using a checkbox. In accepting a class as a concept class, the user is finalizing the positive examples for the training set for the learner. Similarly, the user will be presented with the list of probable non-concept classes generated by the *non-concept filter*. Again, the user accepts individual classes as non-concept classes by clicking a checkbox. After this process is complete, the training set has been finalized, complete with positive and negative examples. This training set can then be used with the support vector machine to classify all of the classes in the system as either concept or other.

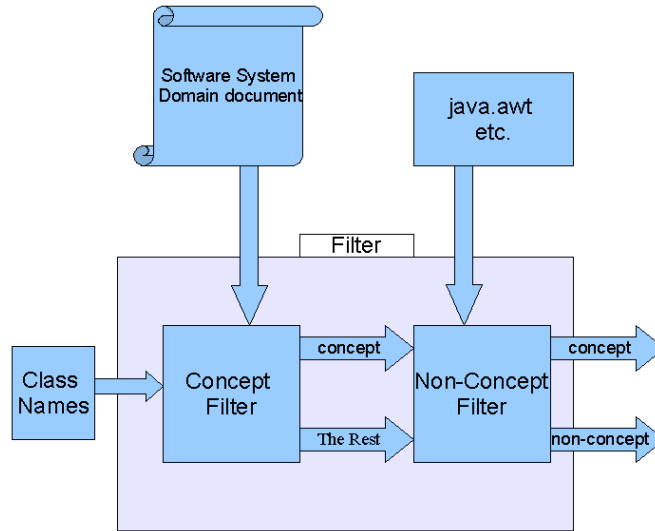


Figure 2: Model of Filter

3.3 Feasibility

3.4 Validation Approach and Measures

TODO: For this section, review some of the statistics used by Carey. Also consider exploring the quality of the training set as a result. For example, maybe the technique for generating a training set falls within the acceptable categories, but is it complete insofar as it describes or defines what a concept class should and should not be?

3.5 Preliminary Results

TODO: After cleaning up some of the statistics, this section will contain some of the initial results from the ArgoUML experiments.

Argo UML is a UML modeling software systems. Since the domain of the systems is UML, I used a glossary of UML terms found online. TODO: Find this source. Also put the document as an appendix to this proposal? I parsed through this document to extract all of the terms.

3.6 Timetable

4 Conclusion

References

- [1] Kristin P. Bennett and Colin Campbell. Support vector machines: hype or hallelujah? *SIGKDD Explor. Newsl.*, 2(2):1–13, 2000.

- [2] Ted J. Biggerstaff, Bharat G. Mitbander, and Dallas E. Webster. Program understanding and the concept assignment problem. *Commun. ACM*, 37(5):72–82, 1994.
- [3] James F. Bowring, James M. Rehg, and Mary Jean Harrold. Active learning for automatic classification of software behavior. In *ISSTA '04: Proceedings of the 2004 ACM SIGSOFT international symposium on Software testing and analysis*, pages 195–205, New York, NY, USA, 2004. ACM.
- [4] Gerardo CanforaHarman and Massimiliano Di Penta. New frontiers of reverse engineering. In *FOSE '07: 2007 Future of Software Engineering*, pages 326–341, Washington, DC, USA, 2007. IEEE Computer Society.
- [5] Sharon A. Caraballo. Automatic construction of a hypernym-labeled noun hierarchy from text. In *Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics*, pages 120–126, Morristown, NJ, USA, 1999. Association for Computational Linguistics.
- [6] M.M. Carey and G.C. Gannod. Recovering concepts from source code with automated concept identification. pages 27–36, June 2007.
- [7] E.J. Chikofsky and II Cross, J.H. Reverse engineering and design recovery: a taxonomy. *Software, IEEE*, 7(1):13–17, Jan 1990.
- [8] T. G. Dietterich. Machine learning. In *Nature Encyclopedia of Cognitive Science*. Macmillan, 2003.
- [9] Y.-G. Gueheneuc, K. Mens, and R. Wuyts. A comparative framework for design recovery tools. pages 10 pp.–134, March 2006.
- [10] He Hu and Da-You Liu. Learning owl ontologies from free texts. volume 2, pages 1233–1237 vol.2, Aug. 2004.
- [11] Hyunjang Kong, Myungwon Hwang, and Pankoo Kim. Design of the automatic ontology building system about the specific domain knowledge. volume 2, pages 4 pp.–1408, Feb. 2006.
- [12] Mitchell. *Machine Learning*. McGraw-Hill Education (ISE Editions), October 1997.
- [13] K. Sartipi. Software architecture recovery based on pattern matching. pages 293–296, Sept. 2003.
- [14] K. Sartipi, K. Kontogiannis, and F. Mavaddat. Architectural design recovery using data mining techniques. pages 129–139, Feb 2000.
- [15] Jochen Seemann and Jürgen Wolff von Gudenberg. Pattern-based design recovery of java software. In *SIGSOFT '98/FSE-6: Proceedings of the 6th ACM SIGSOFT international symposium on Foundations of software engineering*, pages 10–16, New York, NY, USA, 1998. ACM.
- [16] Jelber Sayyad Shirabad, Timothy C. Lethbridge, and Stan Matwin. Supporting maintenance of legacy software with data mining techniques. In *CASCON '00: Proceedings of the 2000 conference of the Centre for Advanced Studies on Collaborative research*, page 11. IBM Press, 2000.

- [17] Jinhui Tang, Yan Song, Xian-Sheng Hua, Tao Mei, and Xiuqing Wu. To construct optimal training set for video annotation. In *MULTIMEDIA '06: Proceedings of the 14th annual ACM international conference on Multimedia*, pages 89–92, New York, NY, USA, 2006. ACM.
- [18] Hui Yang and Jamie Callan. Metric-based ontology learning. In *ONISW '08: Proceeding of the 2nd international workshop on Ontologies and nformation systems for the semantic web*, pages 1–8, New York, NY, USA, 2008. ACM.
- [19] Hui Yang and Jamie Callan. Ontology generation for large email collections. In *dg.o '08: Proceedings of the 2008 international conference on Digital government research*, pages 254–261. Digital Government Society of North America, 2008.
- [20] Jingtao Zhou, Mingwei Wang, Han Zhao, Shusheng Zhang, and Chao Zhang. Concept capture based on column matching and clustering. pages 71–71, Nov. 2005.