| Title | WW |
|---|---|
| Author(s) | , |
| Citation | |
| Issue Date | 2007-03 |
| Type | Thesis or Dissertation |
| Text version | author |
| URL | http://hdl.handle.net/10119/3598 |
| Rights | |
| Description | Supervisor: , , |

cannot understand meaning of technical terms like a medical doctor. The dictionary
has five categories about the five kinds of cancers mentioned above. The number of
media information is the largest in the dictionary. We classified 50 Web pages with
this dictionary based on a liner discrimination model.

We used *SVM* (Support Vector Machine) for machine learning and made classifiers
for the five categories. As a result of our experiments, we got the following results.
We evaluated our method by recall and precision. We got high recall and precision
about stomach cancer and leukemia (C1, C3, C4). We think the reason of the good
results is that technical terms in the dictionary represent features of the categories
well. On the other hand, about uterine cancer, we got bad results. We think the
reason is that number of Web pages about uterine cancer is not enough to make a
good classifier. Personal information about every cancer shows bad results, too. We
think the reason is that personal pages have a various kinds of words and it is difficult
to extract features of the categories.

From these reasons, we propose a method which raises accuracy of classification by
adding new technical terms to the dictionary automatically. This method raises recall
and precision by extracting technical terms (cancer name, medicine name etc.) from
misclassified Web pages and adding them to the dictionary. By using this method, we
got a good prospect about the difficult classification of personal information. In our
study, we investigated cancer information on Web pages one by one with a medical
doctor. We made clear the usefulness of automatic classification of cancer information.

Our experiments showed good results about the classification of authorized infor-
mation, media information and commercial information. It was difficult to classify
personal information. We could improve the accuracy of the classification by adding
new technical terms to the dictionary. We think our study is the first step for support-
ing cancer patients and their families who want to get reliable information on the Web.

cannot recognize technical terms about cancer. The dictionary was made from sentences of Web site "National Cancer Center" by hand.

We made a Naive Bayesian classifier that classifies cancer information automatically according to CII, and experimented classification test. The accuracy of the classification is about 80% and we figure out that this classifier has enough performance to classify cancer information. However, from the discussion about the classification result, we found out that it is difficult to classify C4 category with only language features. We guess reasons of this problem as follows: (1) C4 includes information which maliciously leads patients to purchase goods. Some parts of such pages have information for sales and the other parts explain cancer. These confuse the estimation for the correct category. (2) Pages of individuals or some companies refer Web pages of C1 to explain cancer diseases. Considering these, we use web morphological features such as image size and html size of web pages in addition to language features to classify web pages correctly. Accordingly, the classification accuracy is improved.