

INFO411/INFO911

Data Mining and Knowledge Discovery

Project 3

Instructions:

This task is a real-world data mining problem. You are required to prepare a set of presentation slides which must include (1) the full name and student number of each student in the group, the contribution (in percent) of each group member, (2) a description of the task, (3) your proposed data mining approach and methodology; (4) the strengths and weaknesses of your proposed approach; (5) the performance measures that can evaluate your data mining results; (6) the results a brief discussion and a conclusion.

Below is the recommended structure of your slides:

- Introduction (define the problem and the goal)
- Methods (propose approaches, and discuss their strengths and weaknesses)
- Results (Figures and tables of data analysis)
- Discussion (discovered knowledge from data mining)

Task: Document Classification (Web Spam Detection)

Background:

The term *web spam* refers to the results of activities with an intention to mislead search engines into believing that a particular web page has a high authority value on a particular query, while in fact that particular web page may contain little or no relevant information. Search engines sort the URLs returned in response to a user query on the basis of a score that is usually composed of two parts: a measure of the relevance of the page content with respect to the query (see for example Manning, C., Raghavan, P., & Schütze, H. (2008), “An introduction to information retrieval”, Cambridge University Press) and a measure of the popularity of the page (see e.g. Brin, S., & Page, L. (1998). The anatomy of a large-scale hypertextual Web search engine. *Computer Networks and ISDN Systems*, 30, 107–117.).

In general, spam techniques can be classified into content-based and link-based according to whether their focus is on the former or on the latter measure, respectively (see i.e. Gyöngyi, Z., & Garcia-Molina, H. (2005). “Web spam taxonomy”. *Adversarial Information Retrieval on the Web*).

In a link-based spam, a web page is linked by a large number of links from other web sites. This will increase the popularity of the spam web page as its popularity is due largely to the number of web pages which are linked to it. Thus it is possible to create many web sites, and have each of these web sites linked to a particular spam web page. This is commonly called a “link farm” where links to other web pages can be automatically generated. On the other hand, in content-based spam, a web page is automatically provided with terms that are visually hidden from users and irrelevant to the actual content of the web page. But, these terms are indexable by search engines, so that whatever query that a user issues to a search engine, there is a high possibility that the spam web page will be returned.

We will make use of the benchmark dataset known as the UK2006 benchmark problem. The benchmark is based on a crawl of the .uk domain, carried out in May 2006 and it includes 77.9 million pages and over 3 billion links in about 11,400 hosts. The collection was tagged at the host level by a group of volunteers, who labelled a subset of the hosts as “normal”, “borderline” or “spam.” The dataset and (class/target) labels are available from:

<http://webspam.lip6.fr/wiki/pmwiki.php?n=Main.PhaseIFeatures>
<http://chato.cl/webspam/datasets/uk2006/>

The dataset comes with different feature sets such as text-based features, content based features, link-based features, neighborhood-based features, and direct features.

Definition of the task:

Use the UK2006 spam detection dataset. This is a classification problem. We introduced a number of classification methods in the lectures.

1. Deploy these classification methods.
2. Analyse and compare their results. Use AUC (Area under the ROC curve) as a basis for the comparisons.
3. Which feature set produces the best results?
4. Which combination of feature sets produce the best results?

Requirements:

1. Present a general description of the dataset and present a the general properties of the dataset.
2. Deploy the classification methods to each of the feature sets and present the results. Analyse and compare the results then discuss the strengths and weaknesses of the classification methods used (in the context of this learning problem).
3. Deploy the classification methods to combinations of the feature sets. Present the results and offer a qualitative comparison.
4. Summarize: What new and interesting things did you discover while working on this project?