

MASTER'S THESIS

Analysis of internet image search performance

Wang, Xiaoling

Date of Award:
2011

[Link to publication](#)

General rights

Copyright and intellectual property rights for the publications made accessible in HKBU Scholars are retained by the authors and/or other copyright owners. In addition to the restrictions prescribed by the Copyright Ordinance of Hong Kong, all users and readers must also observe the following terms of use:

- Users may download and print one copy of any publication from HKBU Scholars for the purpose of private study or research
- Users cannot further distribute the material or use it for any profit-making activity or commercial gain
- To share publications in HKBU Scholars with others, users are welcome to freely distribute the permanent URL assigned to the publication

Analysis of Internet Image Search Performance

WANG Xiaoling

**A thesis submitted in partial fulfilment of the requirements
for the degree of
Master of Philosophy**

Principal Supervisor: Prof. LEUNG Clement Ho Cheung

Hong Kong Baptist University

July 2011

Abstract

With the rapid advancement of digital cameras and the Internet, a large collection of digital images can be easily created, shared and distributed, using not only computers, but also numerous other portable digital devices. As digital images have become fully ubiquitous in our lives, searching for the relevant image objects has thus become an important activity. It is desirable to be able to find a solution to searching for images effectively and efficiently. However, many raw images are constantly uploaded without meaningful text labeling or with few words based on the interests of the owner, which is not always reliable and informative. The Internet image search engines based on keywords retrieval, such as Google, Yahoo and MSN, tend to return a large number of images which the engines consider to be relevant, and such pool of results generally is very large and may be regarded as effectively inexhaustible. While the images are presented as relevant, it is normally true that many of them are actually irrelevant, and that the distribution of relevant images over the returned results is non-uniform. Therefore, to predict the distribution of relevant images for the Internet image search engines has become a critical and urgent issue.

The prediction of the relevance for individual images is generally difficult since it only takes on binary values and therefore tends to oscillate randomly between relevance and irrelevance with insignificantly noticeable trends. Increasing the range of possible values is necessary to enhance the prediction ability and it is advantageous to accumulate the aggregate relevance for larger groups of images in a sequential manner. Here, we present a partition approach to the number of relevant images. This will involve

appropriately grouping the random binary sequence into non-overlapping groups and converting it into a form which makes them more amenable for prediction.

Here, we present a Regression model and Markov Chain model for predicting Image Search Engines (ISEs) behaviour. The framework of our approach is initially to design a set of benchmark queries, and then the distribution formula or Markov Chain model will be able to fit the experimental observations by using appropriate parameters and so providing a mathematical description of an empirical process. These two models are particularly effective for the prediction of relevant images for image search engines. The experimental results show that they are able to give good and robust predictions of search engine performance. In addition, the results of this research can have a direct bearing on search engine design to provide informative guidance to users on the retrieval of relevant images, and allows the users to optimize their strategy in the recovery and discovery of images.

After developing the Regression model and Markov Chain model, the estimation of recall becomes easier. In addition, an approach called tagged relevant images, which mainly adopts hypergeometric distribution theory, is developed to estimate recall. The experimental results show that the result is slightly inferior to that of the corresponding estimation methods to create a regression model and Markov Chain model.

Keywords: image retrieval, linear regression model, moving average, exponential smoothing, Markov Chain model, recall, tagged relevant images method

Table of Contents

Declaration	i
Abstract	ii
Acknowledgements	iv
Table of Contents	v
List of Figures	ix
List of Tables	xii
Chapter 1 Introduction	1
1.1 Background	1
1.1.1 Performance Behaviour of Internet Image Search Engines	2
1.1.2 Challenges in the Estimation of Recall	4
1.2 Aims and Contributions of the Thesis	6
1.3 Outline of the Thesis	7
Chapter 2 Stochastic Modelling of Cumulative Page Image Relevance	8
2.1 Partition Approach to the Number of Relevant Images	8
2.2 Stochastic Models	11
2.2.1 Regression Model	11
2.2.2 Moving Average	12
2.2.3 Exponential Smoothing	13

2.3	Query Design and Benchmarking	14
2.4	Measures of Forecast Accuracy	14
2.5	Evaluations	15
2.5.1	Precision	15
2.5.2	Recall	16
2.6	Experiment I: Regression Model for Predicting Image Search Engine Behaviour	18
2.6.1	On Selecting the Appropriate Regression Model	18
2.6.2	On Selecting the Appropriate Cumulative Cell Size	20
2.6.3	Accuracy Test of the Regression Model	30
2.7	Experiment II: Time Series Models for Predicting Image Search Engine Behaviour	34
2.7.1	On Selecting the Appropriate Order for Moving Average	34
2.7.2	On Selecting the Appropriate Smoothing Factor for Exponential Smoothing	39
2.7.3	Accuracy Test of Moving Average	45
2.7.4	Accuracy Test of Exponential Smoothing	45
2.8	Comparison of the Performance Behaviour for Image Search Engines	48
2.9	Summary	51
Chapter 3 Markov Chain Model of Cumulative Page Image Relevance		52
3.1	Stochastic Modelling of Cumulative Page Image Relevance	52
3.1.1	One-step Markov Chain Model	53
3.1.2	Two-step Markov Chain Model	55
3.2	Experiment I: Representing Image Search Engine Performance Using the One-step Markov Chain Model	57
3.2.1	Decay Behaviour of the Number of Relevant Images as Cell Index Increases	58
3.2.2	Accuracy Test of the One-step Markov Chain Model	60

3.3	Experiment II: Representing Image Search Engine Performance Using the Two-step Markov Chain Model	71
3.3.1	Behaviour of the Number of Relevant Images as Cell Index Increases	71
3.3.2	Accuracy Test of the Two-step Markov Chain Model	74
3.4	Evaluation of the Performance Behaviour for Image Search Engines	78
3.5	Summary	81
Chapter 4 Estimation of Database Size and Recall		82
4.1	Corresponding Estimation Method for the Regression Model	82
4.2	Corresponding Estimation Method for the Markov Chain Model	86
4.3	Analysis of Recall Characteristics Using the Tagged Relevant Images Method	89
4.3.1	Methodology of the Tagged Relevant Images	89
4.3.2	Experiment	90
4.4	Comparison of the Recall Rate of Different ISEs Using Different Estimation Methods	91
4.5	Summary	92
Chapter 5 Conclusions and Future Works		93
5.1	Our Contributions	93
5.2	Future Works	95
5.2.1	Considering the Number of Overlapping Relevant Images across All Image Search Engines	95
5.2.2	Discovering the Rules of the Oscillation between Relevant and Irrelevant Images Based on Binary Values	96
Bibliography		97
Chapter A		107
Chapter B		112

