

Automated assessment of the quality of depression websites

Kathleen M Griffiths, PhD, Thanh Tin Tang, B.IT (Hons), David Hawking, PhD & Helen Christensen, PhD

Abstract: 376 words

Text: 4327 words

Dr Kathy Griffiths. Director, Depression & Anxiety Research Unit, Centre for Mental Health Research, The Australian National University, Canberra, Acton, Australia 0200. Kathy.Griffiths@anu.edu.au. (Corresponding author)

Mr Thanh Tang, PhD scholar, Computer Science Department, The Australian National University, Canberra, Australia 0200

Dr David Hawking, Senior Principal Research Scientist, ICT Centre, GPO Box 664, Canberra, Australia, 2601.

Professor Helen Christensen. Director, Centre for Mental Health Research, The Australian National University, Canberra, Acton, Australia 0200.

Abstract

Background: Since health information on the World Wide Web is of variable quality, methods are needed to assist consumers to identify health websites containing evidence-based information. Manual assessment tools may assist consumers to evaluate the quality of sites. However, these tools are poorly validated and often impractical. There is a need to develop better consumer tools, and in particular to explore the potential of automated procedures for evaluating the quality of health information on the web.

Objective: This study (i) describes the development of an automated quality assessment procedure (AQA) designed to automatically rank depression websites according to their evidence-based quality; (ii) evaluates the validity of the AQA relative to human rated evidence-based quality scores; and (iii) compares the validity of Google PageRank™ and the AQA as indicators of evidence-based quality.

Method: The AQA was developed using a quality feedback technique and a set of training websites previously rated manually according to their concordance with statements in the Oxford University Centre for Evidence-Based Mental Health's guidelines for treating depression. The validation phase involved 30 websites compiled from the DMOZ, Yahoo! and LookSmart Depression Directories by randomly selecting six sites from each of the Google PageRank™ bands of 0, 1-2, 3-4, 5-6 and 7-8. Evidence-based ratings from two independent raters (based on concordance with the Oxford guidelines) were then compared with scores derived from the automated AQA and Google algorithms. There was no overlap in the websites used in the training and validation phases of the study.

Results: The correlation between the AQA score and the evidence-based ratings was high and significant ($r=0.85$, $P<.001$). Addition of a quadratic component improved the fit, the combined linear and quadratic model explaining 82 percent of the variance. The correlation between Google PageRank™ and the evidence-based score was lower than that for the AQA. When sites with zero PageRanks were included the association was weak and non-significant ($r=0.23$, $P=.22$). When sites with zero PageRanks were excluded, the correlation was moderate ($r=.61$, $P=.002$).

Conclusions: Depression websites of different evidence-based quality can be differentiated using an automated system developed for the purpose (AQA). If replicable, generalisable to other health conditions and deployed in a consumer-friendly form, the automated procedure described here could represent an important advance for consumers of Internet medical information.

Keywords: quality indicators, depression, Internet, World Wide Web, validity, Information Retrieval.

Introduction

At least 50% of households in the United States, the United Kingdom and Australia are connected to the Internet.[1-3] In addition, many people have access to the Internet outside the home (eg, at work and in public libraries). Recent research has revealed that 79% of American Internet users have searched for health information online.[4] Moreover, there is evidence that online health information can improve health knowledge and health outcomes.[5, 6]

To date, however, few health information websites have been subjected to rigorous assessment of their effectiveness in improving health outcomes. In the absence of such evidence, a minimum requirement for a high quality health information website should be that its content is evidence based; that is, that its content and recommendations are consistent with evidence from a systematic review of the available medical literature. A non-evidence-based site would recommend treatments that are not supported by the evidence as effective, oppose the use of beneficial treatments of demonstrated benefit or fail to mention some effective treatments at all. For example, at the time of writing a depression webpage on the official site of the Church of Scientology claimed that Dianetics was the only ‘proven effective technology of the human mind’ and failed to mention other medical, psychological and alternative treatments of demonstrated effectiveness. Elsewhere on the site, Zoloft and Prozac were described as ‘even more damaging than street drugs’.[7]

Evidence-based health information is routinely disseminated to health professionals with the aim of assisting clinical decision making, improving healthcare and thereby improving health outcomes. Similarly, such evidence, when provided to consumers has the potential to improve health outcomes by assisting consumers to select effective, rather than ineffective self-help techniques and by supporting shared decision making and consumer-provider collaborative care.[8] Unfortunately, reviews of the content of health websites have demonstrated that their quality, when assessed relative to evidence-based standards, is variable.[9-12] People seeking health information may therefore require assistance to filter out lower quality websites.

Various mechanisms have been proposed for assisting consumers to access high quality health websites.[13] These include the use of quality portals (such as OMNI in the United Kingdom and HealthInsite, Australia), adherence to a code of conduct (such as the HON code [14] and the use of consumer tools for assessing a site (eg, DISCERN[15]). However, the criteria for inclusion in a quality portal typically do not include an evidence-based assessment so their validity as guides to quality material is uncertain. Moreover, since such portals require time, effort and training to maintain, it may be difficult to update the database in a timely fashion. Similarly, codes of conduct and consumer tools employ accountability criteria (such as the identification of the author, their affiliations and qualifications) which are typically not validated against an evidence-based standard.[eg,10, 16, 17] One exception is DISCERN, a tool designed to assist users without technical expertise to assess the quality of health information users. We have reported some preliminary evidence that DISCERN may be a valid indicator of the

evidence-based quality of websites when used by consumers.[18] However, the tool may not be valid for all consumers. Moreover, DISCERN requires training and is lengthy, involving 15 items and requiring assessment of all the pages on the target website. In practice, individual consumers may lack the time and motivation to undergo the necessary training or to apply the tool to individual websites.

A potential solution to these problems is to develop assessment tools based on algorithms that automatically evaluate the quality of health information websites. To date there has been little work directed to this possibility. There is some evidence of a relationship between GooglePage Rank™ and evidence based quality from two recent studies.[18, 19] However, the association appears to be only of moderate strength, suggesting that a more valid automatic indicator of website quality may be required.

This paper describes the development of a computer algorithm, the Automatic Quality Assessment procedure (AQA), designed to automatically rank depression websites according to the evidence-based quality of their treatment information. In addition, it describes the results of an evaluation of the validity of the AQA as an indicator of human-rated evidence-based quality of treatment content. It also compares the validity of the AQA and Google PageRank™ as indicators of evidence-based quality.

The study focused on the evaluation of treatment information in depression websites since depression is a primary source of disability burden[20] and it has been reported to be a condition for which users commonly seek information on the Internet.[4] There is also a high degree of unmet need in the treatment of depression.[21]

Method

This section comprises two parts. The first describes the AQA and its development. The second describes the methodology used for the validation of the AQA as an indicator of evidence-based treatment quality. The evidence-based rating scale[11] employed in developing and validating the AQA was based on clinical practice depression guidelines developed by the Centre for Evidence-Based Mental Health in Oxford from a systematic review of the evidence.[22]

1. The Automatic Quality Assessment Procedure (AQA)

Here we present the procedure for calculating AQA scores and note its dependence upon two learned queries and three numerical parameters. We then describe the development phase during which the queries were learned and the parameter values chosen. The development phase employed websites/webpages not in the validation set but for which we had collected human-rated relevance or evidence-based quality measures from our previous studies.[11, 18, 23]

The AQA assumes the availability of search engine software that incorporates a web crawler and has the ability to effectively score the relevance of documents to a query. The current study employed the Panoptic search engine for this purpose. However, we believe that other similar search engines could be substituted with similar results. A set

of computer scripts were written to learn queries, to set values of tuning parameters and to collect and analyse output from the search engine. These scripts are not part of the Panoptic search engine.

The procedure

The AQA procedure comprised six steps as follows:

1. The target websites were downloaded using web crawler software;
2. These downloaded pages were aggregated with a large set of arbitrarily chosen general English language web pages and the resulting collection indexed using the search engine. This was necessary to avoid the extremely biased term frequency distribution of a depression-only collection;
3. A previously learned relevance query (see below) was processed over the collection created in Step 2 using the search engine to produce relevance scores for all documents. The relevance query consisted of many words and phrases, each with a numerical importance weighting. Documents with non-zero scores were not retrieved. For each site to be evaluated, the number of retrieved documents $|R|$ were counted and the mean relevance score (\bar{r}) computed.
4. A previously learned quality query (see below) was processed in the same fashion as in Step 3, yielding $|Q|$ and \bar{q} .
5. Site relevance and site quality scores were computed using Equations (1) and (2). These scores were normalized such that the highest S_r became 1.0 and the highest S_q was also 1.0.
6. An overall site score was computed using Equation (3). Gamma is a scaling parameter designed to make scores comparable with those from the human rating scale.

$$S_r = \alpha * \bar{r} + (1 - \alpha) * |R| \quad (1)$$

$$S_q = \alpha * \bar{q} + (1 - \alpha) * |Q| \quad (2)$$

$$S = \gamma * (\beta * S_q + (1 - \beta) * S_r) \quad (3)$$

The following sections describe how relevance and quality queries were learned and the values of α , β and γ chosen.

Learning relevance and quality queries

Relevance and quality queries were learned using an extension and novel application of the relevance feedback technique from the field of information retrieval. In the relevance feedback approach, a complex query consisting of weighted terms (words and phrases), is automatically generated by comparing the term frequency distributions of sets of relevant and irrelevant documents. Good terms occur frequently in relevant text but seldom otherwise. The resulting query is used by a text retrieval system to derive relevance

scores for documents. We extended this method to learn a 'quality' query from sets of high and low quality webpages.

Relevance query: During development of the relevance query, query terms were selected by computing Term Selection Value (TSV)s[24] for each candidate term, ranking them in descending order and taking all the terms above a cutoff. Numerical weights were applied to the selected terms using the Robertson-Sparck Jones formula.[25]

Using 347 documents previously judged relevant to the topic of depression[23] and 9,000 documents with very low probability of relevance to that topic, we generated a relevance query consisting of the words with the 20 highest TSVs and the two-word phrases with the 20 highest TSVs. The cutoff of 20 was arbitrary but consistent with past information retrieval practice.

Quality query: We generated a 'quality' query in the same fashion, using 110 documents judged to be relevant to depression and of high quality as the 'relevant' set and 3002 documents which were judged either irrelevant or relevant but not of high quality.[23] In this case the number of words (29) and phrases (20) in the query was the minimum number needed to ensure the inclusion of the names of all the evidence-based depression treatments listed in our previously published systematic review of the effectiveness of medical, psychological and alternative treatments for depression.[26]

Choosing parameter values

All the documents from 29 training sites for which we had human evidence-based (Oxford) ratings from previous studies[11, 18] were fetched using the Panoptic crawler and combined with 10,000 documents from the Yahoo! Directory which were not in the depression category, as per Step 2 of the AQA procedure.

In following Steps 3 and 4 of the procedure during training, we computed $|R|$, \bar{r} , $|Q|$ and \bar{q} based on scores obtained using the Okapi BM25[27] relevance scoring mode of the Panoptic search engine. Okapi BM25 takes into account the frequency of occurrence of query term in a document, the discriminating power of each query term, and the length of the document in calculating a relevance score.

The parameter α adjusts the balance between the average document score and the coverage of a site. We then arbitrarily chose $\alpha = 0.75$. The parameter β adjusts the balance between the relevance and quality scores for a site. We stepped through the range of values between 0.0 and 1.0 and chose the value which, when used in Equation 3, maximized the correlation between the computed site scores and the human-assigned quality scores. The best combination found, $\alpha = 0.75$ and $\beta = 0.70$, yielded a correlation of 0.94 on the training data. It is possible that better values could be found with a more exhaustive optimization of parameters.

The parameter γ does not affect the correlation but scales the raw AQA scores to match the range of the human assigned scores. We chose $\gamma = 17.27$ which caused the highest AQA score to be the same as the highest human-rated score.

The values determined in training ($\alpha = 0.75$, $\beta = 0.70$ and $\gamma = 17.27$) were used in the validity testing phase.

2. The validity of the AQA as an indicator of evidence-based treatment quality

Here we describe the methodology used in a comparative validation study of the AQA and the Google PageRank™ procedures. Each of two judges provided evidence-based ratings of 30 new depression websites. These ratings were compared with automated scores derived from the AQA and Google PageRank™ (see *Site Assessment* below).

Selection of sites

The 30 depression information test websites were selected in the following manner.

First, we compiled a master list of all depression websites from the Open Directory (<http://dmoz.org/>), Yahoo (www.yahoo.com) and LookSmart (www.looksmart.com) main and personal and treatment depression subdirectories as of September, 2004. DMOZ, Yahoo and LookSmart are the three major human-compiled search engines on the World Wide Web. The human-compiled directories of many major crawler-based search engines such as Google are derived from Open Directory and the human-compiled content of the Lycos Directory is currently supplied by LookSmart.

Two hundred and eight websites remained after excluding websites that were no longer accessible, websites that were a subdirectory of an already included website, or 'websites' that were links to an individual article.

Using the Google Toolbar, Google PageRank™ scores were recorded for each of these 208 websites. Sites were then pooled into five PageRank™ bands (0, 1-2, 3-4, 5-6 and 7-8) and six websites randomly selected (using a computer generated random number function) from each of the five PageRank™ bands to form an initial set of 30 depression websites. Sites were stratified by PageRank™ prior to sampling to avoid generating a spuriously low correlation due to restricted range effects. A further three sites were excluded because the content was not free (one site), there was no depression information on the primary site (one site) or the site comprised only a single clinical tool for clinicians (one site). These sites were each replaced from the equivalent PageRank™ band using the same computer generated random function.

Site content for each of the 30 websites was printed out in its entirety by systematically following all internal links. Any audio or video material content on a site was accessed online by the evidence-based raters and incorporated into the overall evidence-based score.

Content within a site was included for evaluation if it was free, written in the English language, comprised core informational material and focused on unipolar depression. Since the evidence-based rating scale employed in the current study was based on systematic guidelines for the treatment of major depressive disorder, pages in a site were excluded from evaluation if they focused on bipolar disorder, pre-menstrual syndrome, pre-menstrual dysphoric disorder or seasonal affective disorder. In addition, the following content on the target websites was excluded from the evaluation: news sections, videos of research conferences, book reviews, collections of PubMed abstracts,

poetry, message board and chatroom content. This content was excluded because it often unmanageably large (eg, poetry archives and chatrooms) and peripheral to the core educational material contained on the websites. General clinician assessment instructions and survey databases were excluded as they were not relevant to the study. Non-English content was excluded for practical reasons.

Site assessment

Site characteristics

Test sites were coded independently by the two raters according to their ownership structure (individual vs organisation); whether or not they had an editorial board; whether or not the site was depression specific, was somewhat broader in scope, or comprised a clearinghouse; whether a health professional was involved; and whether the site promoted products or services (see Table 1). Where the two coders disagreed, the final categorisation was assigned by a third rater (KG).

Evidence-based score

Each test site was rated independently by two raters using a 20-item rating scale previously developed by us[11](see Appendix) from statements in the evidence-based, systematically developed clinical practice guidelines for the management of depression in primary care published by the Oxford University's Centre for Evidence-Based Mental Health guidelines.[22] Only statements directly relevant to treatment were incorporated in the 20-item scale. The 30 test sites were rated in a different computer generated random order by the two raters. This rating scale has previously shown high inter-rater reliability.[11, 18] In the current study, inter-rater reliability was also very high ($r=0.93$, $P<.001$) and there was no significant difference between scores for the two raters (mean difference= 0.17, 95% Confidence Interval (CI)= -0.96–0.62, $P=.67$).

Computing AQA scores for the test sites

AQA scores were computed for the test websites by following Steps 1 to 6 of the AQA procedure using the relevance and quality queries and the values of α and β that were derived during training.

Google PageRank™

The Google PageRank™ was recorded for each home page. Google PageRank™ is a measure employed by the Google Search Engine company to evaluate the 'importance' of a webpage. The PageRank™ is based on a computer algorithm that computes iteratively the number and importance of links to a webpage and in turn the number and importance of links to these linking pages.[28] As noted above, we identified the PageRank™ for each test site by downloading the Google Tool bar and recording the integer value (range 0 to 10) on the toolbar for the homepage of the site.

Statistical analysis

The sample size was considered sufficient to justify meaningful parametric analysis of the data. Intercorrelations between variables were computed using Pearson correlation tests. The validity of the automatic measure as an indicator of evidence-based quality was evaluated using hierarchical multiple regression. These analyses were performed using

SPSS 12.0.1.[29] Tests of the significance of differences between dependent correlations were computed using the SISA online calculator.[30]

Results

Site characteristics

The characteristics of the 30 depression test sites are summarized in Table 1. Two-thirds of the sites were depression-specific, a little over one-half were owned by an individual, and a health professional was involved in approximately half of the sites. One-fifth of the sites had an editorial board and a substantial minority promoted products or services or both.

Quality scores

The mean (and standard deviation) of the evidence-based, Google PageRank™ and AQA scores were 5.92 (SD=5.46; n=30), 3.67 (2.59; n=30) and 8.07 (SD=5.22; n=29) respectively. AQA scores were available for 29 test sites only as one website included a robots.txt exclusion, an indicator that the administrator of the site prohibited external crawlers from accessing the website.

The relationship between the AQA score and the evidence-based ratings is shown in Figure 1. The linear correlation between these two measures was high and significant ($r=0.85$, $P<.001$). Visual inspection of the scatterplot in Figure 1 indicated a possible quadratic component to the relationship. A hierarchical multiple regression analysis was therefore performed to determine if adding the square of the AQA score to the linear solution improved the prediction (see Table 2). A substantial 71.4% of the variance in the evidence-based quality score was explained by the automatic quality score alone. Addition of the quadratic component significantly improved the fit ($\Delta R^2=0.10$, $F_{\Delta}(1,26)=14.3$, $P=.001$), the combined linear and quadratic model explaining 82 percent of the variance.

By contrast, the correlation between Google PageRank™ and the evidence based score was small and non-significant ($r=0.23$, $P=.22$, $n=30$; see Figure 2). Excluding the missing case for which AQA could not be computed, this association between Google PageRank™ and the evidence-based score was significantly lower than the association between the AQA score and the evidence-based score ($r(\text{difference})$ ($df=26$) = .64, $t=4.82$, $P=.0001$).

Since it has been argued by some members of the search engine optimisation community that PageRanks of 0 may constitute a special subset of PageRank™ values (see Comment below), the above analyses were recomputed after excluding sites with PageRank™ of 0. The association between the PageRank™ and evidence-based scores for remaining sites was significant ($r=0.61$, $P=.002$, $n=24$), but remained significantly lower than the association between the AQA and evidence-based quality scores ($r(\text{difference})$ ($df=20$) = 0.22, $t=2.61$, $P=.02$).

Discussion

A recent paper in the *Medical Informatics and the Internet in Medicine* journal concluded that ‘quality benchmarking of health-related resources will always depend on a human assessor ...’ [p. 103,31] We have demonstrated that an algorithm based on relevance feedback (and involving no human judgment) is a valid indicator of evidence-based quality of the treatment content of depression sites. To our knowledge, this is the first published study of the validity of a custom designed automated tool for identifying the evidence-based quality of health information. If replicable and generalisable to other health conditions, the current findings may have major practical implications for e-health, consumer empowerment and self-managed healthcare.

Previous researchers have developed search systems designed to identify medical-specific web-based information.[32] However, these systems focus on identifying material that is relevant to the medical domain rather than selecting sites of high content quality. One published study has described a prototype system for rank ordering web-based health information by ‘quality’.[33] However, this paper used accountability criteria (eg, presence of authorship information, detection of an HONcode logo, detection that the page included information about the date it was last updated) rather than content accuracy as a benchmark of quality. There is little or no evidence that these accountability measures singly or together correlate with evidence-based content quality.[eg,10, 16, 17] Moreover, the researchers did not evaluate the content quality of the retrieved pages in order to validate their system of ranking against an evidence-based standard. Finally, in contrast to the procedure described in the current study, the system focused on individual pages rather than on the entire website on a topic. It may be that only by examining all the content of a site is it possible to gain a comprehensive picture of its quality.

The finding in this study that websites can be automatically evaluated for content quality is of considerable practical significance. Suitably adapted, refined and integrated into or used to post-process websites retrieved by a general search engine, this system could assist consumers to identify websites of higher quality. In the shorter term, the system can be used to compile lists of websites for use in a focused search engine for depression, such as that used on the BluePages Depression Information website (<http://bluepages.anu.edu.au>). In addition, the system may prove useful as a screening device for the use by web developers interested in maintaining quality health portals or links of high quality. Once filtered by the automatic quality evaluator, developers could evaluate the remaining sites manually to confirm accuracy of content, and to assess sites according to other dimensions of quality (eg, useability) and according to the needs of the organisation and its users.

Consistent with previous studies[18, 19] we have demonstrated in the current study that Google PageRank™ may be an indicator of evidence-based quality. However, the fact that Google PageRank was unable to provide a meaningful quality assessment for sites with a zero score is a significant impediment to its practical use. Moreover, even with zero PageRanks excluded, the association between PageRank™ and content quality is less strong than the association between the AQA score and content quality. This suggests that relevance feedback may be superior to link structure (employed by Google

PageRank™) as a method for identifying evidence-based quality for a specific health domain. It also provides evidence that a specialised tool such as the AQA is warranted. It might be argued that a combination of relevance feedback and link structure methods could improve the validity of the AQA further. However, the preliminary indications are that this is not the case. Supplementary analyses showed that when the Google PageRank™ was added into the regression model along with the AQA scores, it did not improve the fit of the model. This was the case whether zero page ranks were included (explained 71% of the variance compared to 82% in the original model) or excluded (explained 65% compared to 82%).

The current tool is relevant only for identifying the quality of depression websites. However, the relevance feedback method used in developing the AQA is likely to generalise to other areas of mental health and medicine. To apply the AQA in alternative health domains requires that the procedure is trained to learn new terms and parameter values specific to the new domain. The validity of the technique in these other domains is a matter for empirical investigation.

There are some limitations of the system described here. First, some minor changes to the AQA scaling procedure are necessary before the system is used in practice (particularly for evaluating a single website at a time). Applying a non-linear transformation to the raw AQA scores (rather than linear scaling) might result in a better prediction of the values of the evidence-based scores. However, given that the base (linear) correlation coefficient is already very high (0.85), the scope for improvement is limited. Secondly, the AQA scores could be compromised if publishers use ‘spamming’ methods for optimizing their automatic quality scores. General search engines are faced with a similar problem when website developers attempt to artificially inflate the relevance rankings for their websites. This problem is not peculiar to automated methods for processing websites. It is likely that website developers also use strategies to maximise their scores on manually applied quality evaluation tools while leaving unchanged the substantive content. Public search engines incorporate algorithms for detecting attempts to distort rankings. It is likely that AQA could be refined to do the same. The third limitation of the AQA system is that it is focused on treatment (as indeed is DISCERN) and does not necessarily reflect the quality of other information on a site. In addition, the system may not adequately rate sites which present only one treatment type. The system will be most useful in identifying sites containing high quality, comprehensive treatment information. Another limitation relates to the ‘gold standard’ employed for rating the evidence-based quality of the sites. It might be argued that to be considered valid, the rating system should be validated against ‘health outcomes’ or another scale that has been thus validated. No such scales exist, and given the paucity of efficacy studies of websites, such validation exercises are not currently practical. Finally, one of the evidence-based raters had been involved in the initial selection of the sites which may have led to bias in the findings from the validity studies. However, the pattern of results and conclusions is identical if the findings are recomputed using the ‘blind’ rater’s data only. This is not surprising given the higher inter-rater reliability in this study for the evidence-based scale.

The time has come to acknowledge that consumers do and will continue to use the Internet as a source of health information. We need to provide them with convenient,

effective tools that optimize the usefulness of this process. This study demonstrates that automated methods offer considerable promise in this respect. The task before us now is to refine these methods, and to test the useability, robustness and generalisability of the systems we develop. In the process we need to test alternative strategies for quality filtering, to identify if multiple methods when combined improve the validity of the automatic algorithm, and to evaluate whether the techniques generalise across health domains. We must also construct a user-interface for the procedure, and conduct consumer user and satisfaction studies on the resulting system.

Acknowledgments

This study was funded by the National Health and Medical Research Council Australia Program Grant to the Centre for Mental Health Research (Chief investigators: A F Jorm, H Christensen, B Rodgers, K Dear, and K Anstey). The work described here was carried out independently of the funders. The authors would like to thank the following people for their contributions to the study: Michelle Banfield who systematically identified and printed out the webpages for the test collection, provided evidence-based rating of the test data, set up the database for the study and conducted some preliminary analyses of the results.

CONFLICTS OF INTEREST

DH and TT are employed by CSIRO which owns and has commercialised the Panoptic search technology used in the study. DH is leader of the Panoptic R&D team. The AQA is not currently part of any commercial product and no patents have been applied for in respect of it.

References

1. National Telecommunications and Information Administration, A nation online: Entering the broadband age. 2004, US Department of Commerce: Washington.
2. Office for National Statistics, Individuals accessing the Internet - National Statistics Omnibus Survey. Access to Internet from home. Family Expenditure Survey 2004., Office for National Statistics: London.
3. Australian Bureau of Statistics, Household Use of Information Technology, Australia, 2002 and 2003. 2005, Australian Bureau of Statistics: Canberra.
4. Fox, S. Health information online.
http://www.pewinternet.org/pdfs/PIP_Healthtopics_May05.pdf: Accessed: July 21, 2005.
5. Christensen, H., K. Griffiths, and A. Jorm, Delivering interventions for depression by using the internet: randomised controlled trial. *BMJ*, 2004. doi:10.1136/bmj.37945.566632.EE.
6. Griffiths, K.M., et al., Effect of web-based depression literacy and cognitive-behavioural therapy interventions on stigmatising attitudes to depression. *Br J Psychiatry*, 2004. 185: p. 342-349.
7. Church of Scientology International. www.scientology.org.: Accessed: 30th October, 2005,
8. Hibbard, J.H., Engaging health care consumers to improve the quality of care. *Med Care*, 2003. 41(1): p. I61-I70.
9. Eysenbach, G., et al., Empirical studies assessing the quality of health information for consumers on the World Wide Web. *JAMA*, 2002. 287(20): p. 2691-2700.
10. Griffiths, K.M. and H. Christensen, Quality of web based information on treatment of depression: cross sectional survey. *BMJ*, 2000. 321: p. 1511-1515.
11. Griffiths, K.M. and H. Christensen, The quality and accessibility of Australian depression sites on the World Wide Web. *Med J Aust*, 2002. 176: p. S97-S104.
12. Berland, G.K., et al., Health information on the internet. Accessibility, quality, and readability in English and Spanish. *JAMA*, 2001. 285(20): p. 2612-2621.
13. Risk, A. and C. Petersen, Health information on the internet. Quality issues and international initiatives. *JAMA*, 2002. 287(20): p. 2713-2715.
14. Health On the Net Foundation. HON Code of Conduct (HONcode) for medical and health web sites. <http://www.hon.ch/HONcode/Conduct.html>: Accessed: 16 July, 2005.
15. Charnock, D., et al., DISCERN: an instrument for judging the quality of written consumer health information on treatment choices. *J Epidemiol Community Health*, 1999. 53: p. 105-111.
16. Fallis, D. and M. Frické, Indicators of accuracy of consumer health information on the Internet. *JAMIA*, 2002. 9: p. 73-79.
17. Kunst, H., et al., Accuracy of information on apparently credible websites: Survey of five common health topics. *BMJ*, 2002. 321: p. 581-582.
18. Griffiths, K. and H. Christensen, Website quality indicators for consumers. *JMIR*, In Press.

19. Frické, M., et al., Consumer health information on the web about carpal tunnel syndrome. *Am J Med*, 2005. 118: p. 168-174.
20. Murray, C.J.L. and A. Lopez, eds. *The global burden of disease: A comprehensive assessment of mortality and disability from diseases, injuries and risk factors in 1990 and projected to 2020*. 1996, Harvard University Press: Cambridge, MA.
21. Kohn, R., et al., The treatment gap in mental health care. *Bull World Health Organ*, 2004. 82(11): p. 858-866.
22. Centre for Evidence Based Mental Health. A systematic guide for the management of depression in primary care: treatment. <http://www.psychiatry.ox.ac.uk/cebmh/guidelines/depression/depression1.html>: Accessed: 26th October, 2005.
23. Tang, T.T., et al., Quality and relevance of domain-specific search: a case study in mental health. *Journal of Information Retrieval*, In press.
24. Robertson, S.E., On term selection for query expansion. *Journal of Documentation*, 1990. 46(4): p. 359-364.
25. Robertson, S.E. and K.S. Jones, Relevance weighting of search terms. *Journal of Documentation*, 1976. 27(3): p. 129-146.
26. Jorm, A., et al., *Help for depression: What works (and what doesn't)*. 2001, Canberra: Centre for Mental Health Research.
27. Robertson, S.E., et al., *Okapi at TREC-3*, in *Proceedings of TREC-3. NIST Special publications 500-225*. 1994, NIST: Gaithersburg. p. 109-126.
28. Brin, S. and L. Page, Anatomy of a large-scale hypertextual web search engine. *Proceedings of the 7th International World Wide Web conference*, 1998. 7: p. 107-117.
29. SPSS for Windows, *Release 12.0.1*. 2003, Chicago Illinois: SPSS.
30. SISA online statistical analysis. <http://home.clara.net/sisa/correl.htm>: Accessed: May, 2005.
31. Curro, V., et al., A quality evaluation methodology of health web-pages for non-professionals. *Med Inform Internet Med*, 2004. 29(2): p. 95-107.
32. Suarez, H.H., X. Hao, and I.F. Chang, Searching for information on the Internet using the UMLS and Medical World Search. *Proceedings of the AMIA Annual Fall Symposium*, 1997: p. 824-828.
33. Price, S. and W. Hersh, Filtering Web Pages for Quality Indicators: An Empirical Approach to Finding High Quality Consumer Health Information on the World Wide Web. *Proceedings of the AMIA 1999 Annual Symposium*, 1999: p. 911-915.

Table 1. Characteristics of the websites employed in the test phase of the study

Site characteristic	n (%) of sites
<i>Ownership structure</i>	
Individual	17 (56.7%)
Organisation ^a	12 (40.0%)
Unknown	1 (3.3%)
<i>Editorial board</i>	
Yes	6 (20%)
No	24 (80%)
<i>Scope</i>	
Depression specific	20 (66.7%)
Broadscope	9 (30.0%)
Clearing house	1 (3.3%)
<i>Health professional involved</i>	
Yes	16 (53.3%)
No	14 (46.7%)
<i>Promotion of products/ services</i>	
Yes	11 (36.7%)
No	14 (63.3%)

^acommercial, consumer or other organised group

Table 2. Summary of hierarchical regression analysis for predicting evidence-based quality from automatic quality

Variable	B	SE(B)	β	<i>p</i>
<i>Model 1</i>				
AQA	.85	.10	.85	<i>p</i> <.001
<i>Model 2</i>				
AQA	-.22	.30	-.22	<i>p</i> = .45
AQA ²	.07	.02	1.11	<i>p</i> =.001

Note: Model 1: $R^2 = .71$; Model 2: $\Delta R^2 = .10$ (*p*=.001).

AQA=Automatic Quality Tool Score.

**Mean
evidence-
based score**

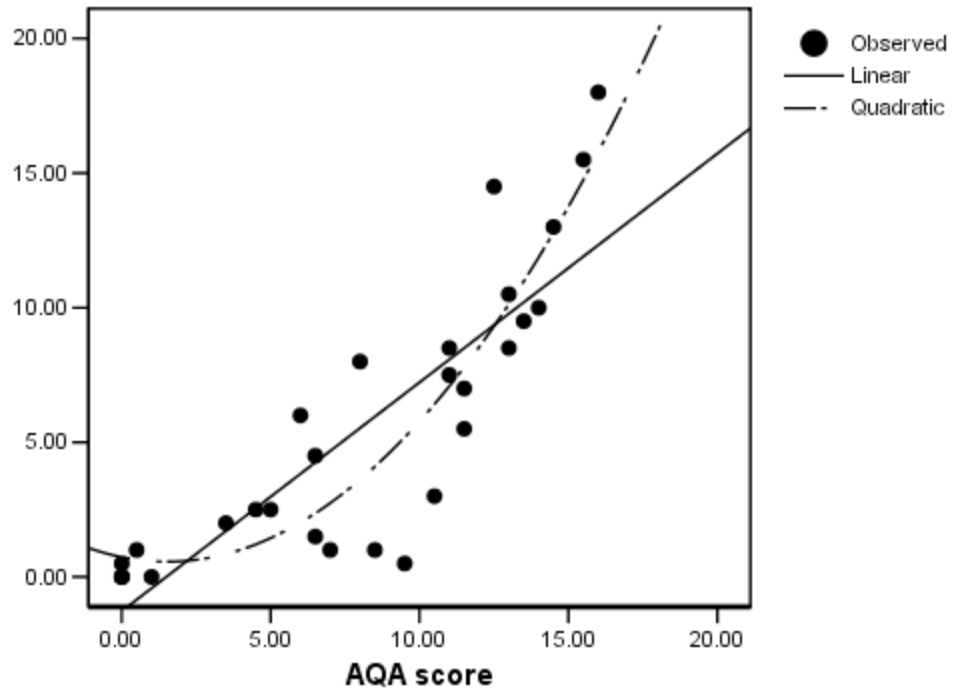


Figure 1. The relationship between the AQA and evidence-based scores.

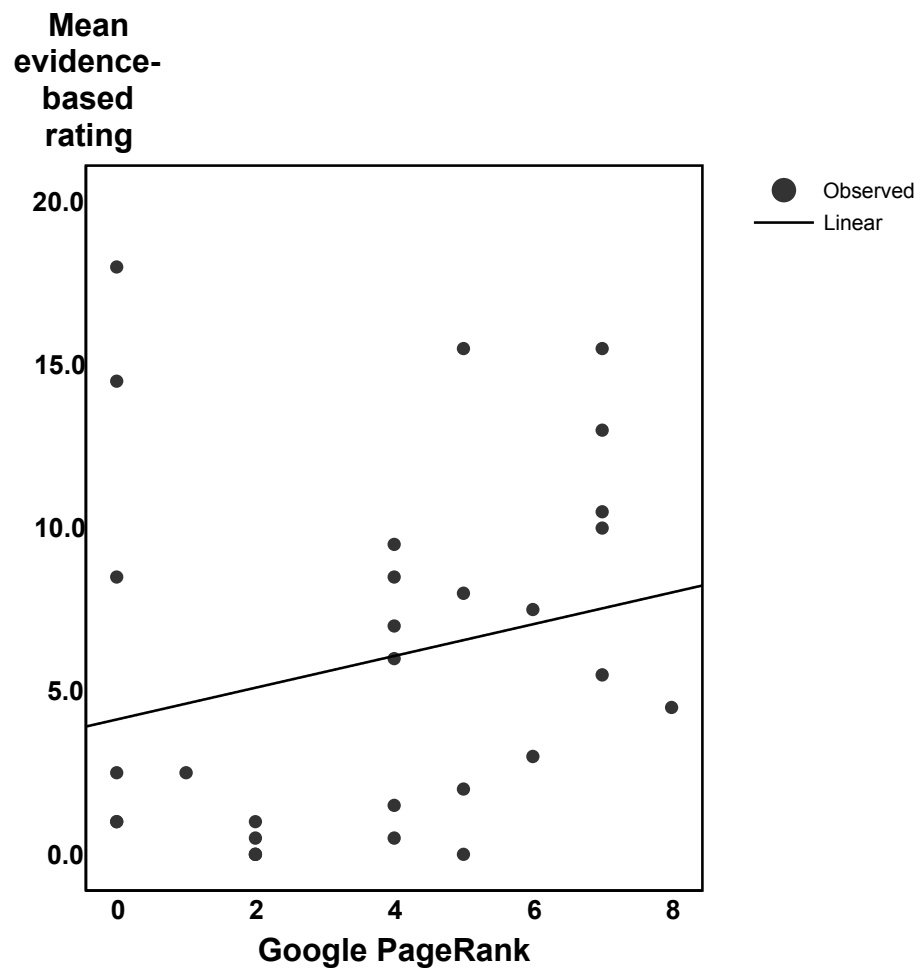


Figure 2. The relationship between the Google PageRank™ and evidence-based scores.

APPENDIX

Evidence-based rating scale

The evidence-based rating scale[11] was developed from statements in the treatment section of *A systematic guide for the management of depression in primary care* published by the Centre for Evidence-based mental health, Oxford[22]

1. Antidepressant medication is an effective treatment for major depressive disorder.
2. Antidepressants are all equally effective.
3. The effectiveness of antidepressants is around 50 to 60%.
4. Full psychosocial recovery can take several months.
5. Drop out rate is same for different antidepressants.
6. The side effect profile varies for different antidepressants.
7. The choice of antidepressant should depend on individual patient factors (eg presence of co-morbid psychiatric or medical conditions, previous response to a particular drug, patient preference regarding the desirability of specific side-effects, concurrent drug therapy, suicidal risk)
8. Antidepressants are not addictive
9. A trial of 6 weeks at full dose is needed before a drug can be considered to have failed and another tried.
10. A second-line drug should probably be from a different class of antidepressant.
11. Once improved continuation treatment at the same dose for at least 4-6 months should be considered.
12. Discontinuation syndrome may occur with abrupt cessation of any antidepressant so antidepressants should not be stopped suddenly. Where possible antidepressants should be withdrawn over a 4 week period, unless there are urgent medical reasons to stop the drug more rapidly. [To score 1, need to make general points that abrupt cessation can cause discontinuation syndrome and that antidepressants should not be stopped suddenly]
13. St John's Wort appears to be as effective as tricyclic antidepressants and causes fewer side effects, but little is known about any long term adverse effects.
14. Cognitive therapy can be an effective treatment for depression
15. Cognitive behaviour therapy is at least as effective as drug treatment in mild-to-moderate depression.
16. Cognitive behaviour therapy may be valuable for people who respond to the concept of Cognitive behaviour therapy, prefer psychological to antidepressant treatment, have not responded to antidepressant therapy. [Score 1 if mention at least one of these]
17. Problem-solving may be effective for depression.
18. [Generic] counselling is probably no more effective than treatment as usual from the GP for depression.
19. Written information (usually based on a cognitive model of depression) can improve mild-to-moderate depression. [Score 1 if cognitive model]
20. Exercise can be effective - alone or as an adjunct to other treatments.

For each item, score 1 if the site information is consistent with the statement. Cumulate item scores across the scale to yield a total evidence-based score for the site.