SIGCSE: U: Focused Retrieval of University Course Descriptions from Highly Variable Sources

THOMAS EFFLAND*

SUNY, University at Buffalo tom.effland@gmail.com

Abstract

Finding topically relevant content from sparse disparate sources on the Web requires robust techniques. A focused web crawler is a type of crawler that attempts to make predictions about page relevance and traverse the web efficiently to retrieve relevant information. In this work, we design and test a novel framework of focused crawling tailored to extracting semantically relevant information from disparate seed domains with highly variant structure that do not reference each other. We utilize machine learning techniques to predict the normalized link distance of current pages to target pages by employing two separate Random Forest regressors that rank the current page and potential relevance gain of hyper-links. We use a novel reformulation of page relevance as the normalized link distance to efficiently tunnel through irrelevant pages and reach target pages with close to optimal paths for domains with large inter- and intra-site variability. We evaluate this system on a concrete problem: retrieving relevant course description and information pages from many university websites with little training data. We use evaluate the training efficiency of the system using mean regression error and evaluate the retrieval efficiency in practice using the harvest rate metric.

I. PROBLEM & MOTIVATION

The exponential growth of the World Wide Web (WWW) has given rise to an explosion of publicly accessible data in the form of unstructured natural language text and semistructured hyper-text. Often we are interested in semantically similar content that is available in many different web domains because the aggregate of the content can contain richer information than each individual datum. The connections between different data sources are commonly non-existent and a user is forced to search for individual data items separately. However, finding and retrieving this data manually is often unrealistic at scale.

Topical crawlers are the specialization of traditional search engines or crawlers to finding relevant pages across domains: a form of "smart" search. A topical crawler makes a basic assumption of "topical locality" within the web that asserts that relevant pages are typically close together between sites, i.e. there is a short link distance between them which enables traversal across sites to find relevant content. Here we address a different reformulation of the topical crawling problem where the topical locality assumption does not necessarily hold:

> How can we find semantically similar information from separate known sources that do not reference each other?

In this work we are motivated by the example application of mining course information and descriptions from multiple university websites. This specific task is difficult for multiple reasons:

Retrieving content on one site doesn't di-

^{*}Partially supported by NSF Grant DUE-CCLI-0920335

rectly lead to content on another site, i.e. no inter-site topical locality.

- University websites are highly varying in structure from one institution to the next, so reaching the relevant information with a general, but efficient technique is difficult.
- The highly variate structure makes identifying a relevant seed page difficult. Without a relevant seed page, we are forced to start at the root domain page, e.g. www.buffalo.edu.
- Course information in university websites is typically very sparse. Thus we must significantly restrict and navigate the search space in an efficient manner.

In many ways, this problem is akin to the automation of mimicking a user's browsing decision-making process for finding relevant content on many sites by starting at the sites' root domains.

II. BACKGROUND & RELATED WORK

The rapid growth of the World Wide Web presents many challenges to finding specific topical information for traditional web crawlers. Often a crawler is targeted towards a particular topic and is unable to find topical information in an efficient manner. A *focused* crawler attempts to make decisions on which pages to crawl based on a classifier trained to estimate the relevance of a page [6]. Typically, a focused crawler represents a page using the "Vector Space Model" in which a page is represented by a vector of features (often words) so that a traditional classifier may be applied.

There are multiple challenges that arise for focused crawlers. The main issue in the context of our problem is "tunneling", in which a relevant page may only be reachable by traversing irrelevant pages. [5] used reinforcement learning to assign credit to irrelevant pages on the relevant path. [2] address the issue by using a maximum depth counter that resets when the crawler reaches a relevant page. [4] trained a set of Naive-Bayes classifiers that attempt to predict the link-distance the current page is from a target page; This is called "Context-Focused" crawling. [1] evaluates both page content and link structure separately to predict the best links to follow. Both [3] and [7] used Named-Entity Recognizers to extract namedentities as features in addition to words.

III. UNIQUENESS OF APPROACH

In this work we build on previous works for developing efficient focused crawlers and apply the techniques to a crawler that can traverse highly variate, disparate sources where the topical locality assumption from one site to the next does not necessarily hold. Our framework is designed to develop an efficient system that utilizes two machine learning regressors to make relevance predictions for pages and their links with limited training data. This is presented in four parts: the page feature representation (input features), the page relevance representation (target variable), the training phase, and the deployment phase.

I. Page Representation

Web page source code provides semistructured information in the form of HTML. To utilize traditional machine learning techniques, we convert this original data into a feature vector using the vector space model (cite VSM). We use HTML tree parsing to split a page *P* into a content vector \vec{p} and a set of url vectors $U = {\vec{u}_1, ..., \vec{u}_k}$, where *k* is the number of <a> tags on the page. We then quantify these two representational elements separately using traditional information retrieval techniques.

Representing the Page Content

To represent the page content, we take the following steps:

- 1. Convert all the text inside the <body>, <title>, and <a> tags into vectors $\vec{b}, \vec{t}, \vec{a}$ of words by splitting on whitespace.
- 2. Remove special characters and English stop-words using NLTK (citeNLTK).

- 3. Stem the vectors using WordNetLemmatizer (cite NLTK).
- 4. Expand the vectors by adding in bigrams of the terms (cite bigrams).
- 5. Take the TF-IDF (cite tfidf) of the vectors to get the relative frequencies of the terms within the document and the vocabulary.
- 6. For the <body> vector only: use Latent Semantic Analysis (cite LSA) to embed the vector in a lower dimensional semantic space. This greatly reduces the size of the input vector.
- 7. Take the concatenation of these three transformed vectors as the page feature vector $\vec{p} = \langle \vec{b}', \vec{t}', \vec{a}' \rangle^{-1}$.

Thus we have a numerical feature vector that represents the content of a page. We note that this representation is generic with no domain-specific engineered features.

Representing the Link Content

To represent each url $\vec{u} \in U$ we extract and segment its href attribute to get a vector of terms \vec{h} . We also extract the anchor-text within the <a> tag as we did for the page to get the vector \vec{a} . We then form the final representation $\vec{u} = \langle \vec{h}', \vec{a}' \rangle$ using the same sequence of steps described above for the page representation (excluding step 6).

Defining a Relevance Metric

To reach relevant pages from the irrelevant root url of a site, we must identify and traverse many irrelevant pages that lie on a path to the relevant content. We also need to address this issue in a way that is robust to the structural variation of the different sites and differing path lengths. Thus we define the target relevance R of a page P to be the normalized link distance from P to a target page T, with the link distance from the starting page S to T as the normalization factor. Formally,

$$R(P) = 1 - \frac{LinkDist(P,T)}{LinkDist(S,T)}$$
(1)

Intuitively this metric can be thought of as what fraction of the total path length is this page P from our target page T. In practice a page may lie on many relevant paths and in this case we average all of its relevance scores for the final score.





This relevance score may now be used in combination with their corresponding page feature vector and machine learning methods to make predictions about the relevance of a page. We also define the relevance of each link \vec{u} to be the relevance of the page the link leads to. Thus we may use the links to make predictions about which pages to traverse to next.

¹The length of \vec{p} is constant and ensured through the TF-IDF and LSA transformations.

II. Training the System



Figure 2: Flowchart of self-training procedure for the regressors. This iterative approach allows for training of accurate regressors with little manually-gathered training data.

Since the goal of this work is to automate a manual process, we approach the problem of learning accurate regressors with as little manually gathered training data as possible. We utilize a semi-supervised learning technique called self-training (cite selftraining) to iteratively improve our regressors' performance as follows:

- 1. The user manually marks a few ² sample traversals from the start page to a target page as a sequence of urls.
- 2. A training crawler follows these paths and downloads each page on the path as labeled data. The crawler also downloads the neighboring pages (within a link of a page on the path) as unlabeled data.
- 3. The features of the labeled pages and links are extracted along with their corresponding relevance metrics.
- 4. Two random forest regressors (cite RF) ³ are trained on these training data so we can make predictions of page and url relevance as defined in I.

- 5. The regressors are used to make predictions on the unlabeled dataset. The most confident predictions (typically the highest and lowest ranked 25 pages) are then given labels of 1.0 and 0.0, respectively. All unlabeled pages lying on paths to the new target pages are labeled with their calculated relevancies also.
- 6. Iteratively continue steps 2-5 until regressors reach desired accuracy or the unlabeled dataset is used up.

Using this technique, we are able to automatically train accurate regressors with little manually-labeled training data.

III. Deploying the System





After training the regressors, we can deploy them to retrieve pages of interest on other sites.

Beginning at the starting url, we traverse a site by featurizing the urls on a page and using the url regressor to make predictions of the relevance for the page the urls lead to. We then rank the urls by highest predicted relevance and insert them into a priority queue. To choose the next page to crawl, we simply pop from the queue. Thus we use a greedy

²In practice, we used 10 per training site.

³We choose to use Random Forests as our machine learning regressors because they have been shown to be robust to sparse, highly-variable data. (cite robust RF)

approach in traversal by consistently following the link of highest predicted relevance in search of relevant pages. This has the advantage of being able to significantly restrict our search space in a site.

At each page, we also use the page regressor to make a prediction about the page content's relevance. If the predicted relevance is above a threshold (in practice we use .85), then we classify the page as relevant and add it to the set of gathered pages.

To continually grow the training set throughout deployment, we utilize activelearning (cite active learning) to query the user for input on pages that may be confusing the regressors. For example, the url regressor may predict a page has a relevance of .95 while the page regressor predicts a relevance of .25. In this event, correctly classifying the page as relevant or irrelevant may provide discriminative information to the training set. Utilizing activelearning, the user may periodically retrain the system and further increase accuracy.

IV. Results & Contributions

To evaluate the system, we present two metrics. We first evaluate the accuracy of the page and url regressors using absolute regression error. We measure this as a function of the size of the training data set as grown through selftraining. We then evaluate the efficacy of the system in practice by measuring the harvest rate (cite harvest rate) of retrieving relevant course description pages on university websites.

We tested on five university sites: buffalo.edu, illinois.edu, bu.edu, washington.edu, and northwestern.edu ⁴. For each site, we marked ten sample traversal paths from the root url to course descriptions pages. In evaluation, we utilize hold-one-out methods by training the system on all sites except the site we test on.

Training Results

We evaluate the accuracy of the page and url regressors using mean absolute regression error in prediction, i.e. if the true score was .66 and the predicted score was .60, then the absolute regression error would be .06.

Using the 50 sample paths generates a little less than 100 initial training examples among the five sites. We then grow the training dataset approximately 50 pages per iteration using selftraining. From the graph, we see that at about 550 training example the average mean page prediction error among all five sites (represented by the thick green line) reaches .048. A similar plot of the mean url prediction error would show that the average among the sites reaches .052 for the same training set. We omit this plot due to space constraints.



Figure 4: This plot shows mean absolute prediction error for the page regressor as we iteratively increase the training dataset using self training. Each university is represented in its school color and the average among all five in green. The star represents the lowest average of .048 at 542 pages. Here we see that initially the system is quite inaccurate, but is eventually able to make better predictions while training on other sites through self-training.

Testing Results

We evaluate the efficacy of the trained system in retrieving relevant course descriptions pages using the standard metric in focused crawling: harvest rate. Harvest rate is the fraction of

⁴These five sites were chosen for evaluation because they each have a different organizational structure for accessing course descriptions, but still present the data in HTML instead of requiring database querying.

pages retrieved out of the number of pages visited so far. $HarvestRate = \frac{\#Retrieved}{\#Visited}$.

We again test on the previously mentioned five sites, using hold-one-out training methods. In the graph we see that we reach a harvest rate of almost 80% for 4 of 5 sites. This shows the system efficiently navigating to the content rich areas of the site. We also note that illinois.edu is less successful. This is because the system initially followed a pseudorelevant path ⁵ before discovering the correct region of the site. However even in this case, we find that the system is able to find relevant content in less than 200 visits.



Figure 5: This plot shows the harvest rate of the system in deployment using hold-one-out training. From this plot we see that the system is able to achieve a harvest rate of 80% within about the first 50 pages crawled for 4 of the 5 sites. This means that 40 of the first 50 pages were retrieved pages these sites and shows that the system is able to navigate to the content-rich areas of the sites. We note that *illinois.edu* does not have the same success, as the system initially followed a pseudo-relevant branch, however after exhausting this branch, the system does eventually catch on to the course descriptions area of the site.

Conclusions and Impact

In this paper we have presented a novel focused crawling architecture tailored to retrieving semantically similar content from many highly variable disparate sites with no initial assumptions about site organization or topical locality. Using self-training, we require little manually-gathered data, yet are able to train accurate regressors to predict page and url relevance. The novel reformulation of page relevance as normalized link distance is key to addressing the issue of high organizational variation with a general scheme. This is because instead of attempting to classify integer link distances, which leads to binary classification errors, we are able to encode a notion of closeness by reframing the relevance as a scalar.

We evaluated the training stage of our system using mean absolute prediction error and showed that it significantly improved its regressors using self-training and was able to make accurate predictions of page and url relevance. We then evaluated the system in practice by measuring the harvest rate of the trained system. We showed that in 4 of 5 cases, the system reaches a harvest rate of almost 80% within 50 visits and in all 5 cases the system was able to navigate to content-rich areas of the sites within 200 visits.

We note that this system has been designed as a general information retrieval framework and can be used in any scenario where the user wants to automate the process of collecting data from many disparate sources, but only has a list of the source domain names. In the future, this system could be significantly improved by utilizing crowd-sourcing services to generate considerably larger and more informative training datasets and produces extremely accurate regressors. Coupled with domainspecific information extraction tools, it may be used to automatically generate large databases of cross-site information.

References

[1] CHEN, X., AND ZHANG, X. Hawk: A focused crawler with content and link analysis. In *e-Business Engineering*, 2008. ICEBE'08. IEEE International Conference on (2008), IEEE, pp. 677–680.

⁵A pseudo-relevant branch is a path with high initial predicted values, but never leads to a relevant page: a shortcoming of the greedy strategy.

- [2] DEVI, P., AND THAKUR, R. Comprehensive review of web focused crawling.
- [3] DI PIETRO, G., ALIPRANDI, C., DE LUCA, A. E., RAFFAELLI, M., AND SORU, T. Semantic crawling: An approach based on named entity recognition. In Advances in Social Networks Analysis and Mining (ASONAM), 2014 IEEE/ACM International Conference on (2014), IEEE, pp. 695–699.
- [4] DILIGENTI, M., COETZEE, F., LAWRENCE, S., GILES, C. L., GORI, M., ET AL. Focused crawling using context graphs. In *VLDB* (2000), pp. 527–534.
- [5] MCCALLUM, A. K., NIGAM, K., RENNIE, J., AND SEYMORE, K. Automating the construction of internet portals with machine learning. *Information Retrieval 3*, 2 (2000), 127– 163.
- [6] NASRAOUI, O. Web data mining: Exploring hyperlinks, contents, and usage data. ACM SIGKDD Explorations Newsletter 10, 2 (2008), 23–25.
- [7] SAMARAWICKRAMA, S., AND JAYARATNE, L. Focused web crawling using named entity recognition for narrow domains. *IJRET* | *DEC* (2012).