

OF COURSE! Using Bayesian Inference to Build a More Dynamic Course Search

CS221 Final Project by Michael Dickens and Mihail Eric

December 12, 2013

Introduction

ExploreCourses is a search engine used regularly by the Stanford University community for browsing and finding courses. At the present moment, ExploreCourses can support basic query searches, including somewhat accurate retrieval given a course code and the exact title of a course. Almost any query search that does not fall into one of these two classes of searches will either return unrelated class results or more often no results at all. For this project, we sought to develop an improved course searching program that would be more robust in that it could support more diverse user input queries and would also be more dynamic in that the courses that were returned were more related given a universal metric that we define for assessing relatedness. The metric will be explained in a later section.

To achieve improved robustness, we implemented some basic natural language processing schemes for extracting useful and relevant information from a user input. The information that we were specifically looking for included course titles, course codes, department codes, and instructor names. To find more related courses, we extracted a variety of features that we considered relevant from all the data we could attain about a course and then we created a course-relatedness “graph” that assigns a relatedness score to each pair of classes, given their extracted features. To compute the relatedness score, we utilized a Bayesian inference scheme whereby we calculated the probabilities that two courses are related given that they have a pair of features in common. The Bayesian approach and the features used will be explained in later sections.

Feature Extraction

In order to determine an accurate label for relatedness between two courses, we had to extract a useful collection of features from the data for each course. This required some experimentation in order to find a good balance: too many features and the algorithm runs slowly; too few, and we cannot perform useful inference.

We used the following initial set of features:

- words in the title
- words in the description
- course code name
- course code ones digit
- course code tens digit
- instructors
- minimum units
- maximum units

Later, we took each of the course code features and combined them with each title, instructor, and description feature to create a set of binary features. This roughly quadruples the total number of features in the set. We also tried using word bigrams in the title and description, but this did not add substantial benefit.

A Bayesian Approach to Course-Relatedness

To determine course relatedness, we began simply by taking the total number of matching features between two courses. This approach proved too coarse: it matched courses with many common features that weren't actually all that related to each other.

A sudden insight came when we realized we could do much better by taking a Bayesian approach. Instead of simply counting the number of features in common, interpret each

feature in common as Bayesian evidence that the two courses are related and perform a Bayesian probability update. Similarly, if a feature does not occur in common between two courses, consider this evidence that they are not related.

Then, instead of considering each feature as equally strong evidence, weigh a feature against the prior probability of that feature occurring. We figure out the prior probability of a feature by counting how frequently it occurs in the database.

Thus, to find the probability that two courses are related, we update a prior probability estimate with the evidence given by each feature found in the two courses.

Query Parsing

In order to satisfy a user's input queries, we need to extract useful information from a given query. Using the Python Natural Language Processing Toolkit, we employed the following natural language processing scheme: tokenize query, tag with parts of speech, chunk appropriately using regular expression grammars. Once a query is chunked, useful information can be derived through analysis of the corresponding parse tree. Using this scheme, we are able to support user query searches consisting of more complex phrases

We search for instructors using a simple grammar that searches a string for sequences of proper nouns. We support course code and department code searches by tokenizing a query into unigram and bigram tokens and checking a set of course/department codes for matches. We also support title searches by searching for the input string in a set of all course titles.

*****TODO: Talk about use of NLTK POS-Tagger, chunker, regex grammar.*****

Data

We used the ExploreCourses Java API to acquire information related to the 11,613 courses listed for enrollment for the 20132014 academic year. We populated a SQL database, using the sqlite3 Python library, with the following information for each course: course title, course

code, instructors teaching the course, minimum units of credit received for taking the class, maximum units of credit received for taking the class, and course description.

However, for the purposes of the assessment, we used a database of a reduced subset of approximately 170 random classes taken from the CS and MATH departments. We had to utilize a reduced database because it was too time-consuming to create a comprehensive relatedness graph for all 11,613 courses.

*****TODO: Explain how we find the top match and then find the most related courses.*****

Metric for Assessment

We developed a simple point-based metric to objectively determine the quality of our course searcher as compared to ExploreCourses. To evaluate the performance of our searcher, we generated a random subset of 26 courses taken from our reduced database. We subdivided these 26 courses into five types of searches:

1. Specific course code such as CS109
2. Instructor name such as ‘Mehran Sahami’
3. Course title such as ‘Introduction to Probability for Computer Scientists’
4. Subset of a course title such as ‘Probability for Computer Scientists’
5. More complex queries

The more complex queries included two of each of the following:

1. ‘courses taught by [instructor]’
2. ‘[course 1] and [course 2]’
3. ‘courses in the [department] department’

We used the following point system to evaluate the accuracy of queries:

50 points if desired course is the first hit, 30 points if the second hit, then 25, 20, 15.

For each of the following, give the full points for the first hit, $4/5$ of the points for the second hit, $3/5$ for the third, etc.

- 10 points if same sequence as searched-for course (e.g. Math 51, 52, 53)
- 5 points if in the same department
- 10 points if is the same course in a different department
- 10 points if has same instructor
- 10 points if course descriptions are sufficiently similar, judged subjectively

Searches for instructors are judged more simply: 50 points if the first course is taught by the instructor, 45 for the second course, then 40, 35, 30.

These numbers are somewhat arbitrary. To get a better idea of the reliability of our search algorithm, we should collect data on user satisfaction with real-world queries. Unfortunately, such a metric is not feasible at this time, so for now we will stick with the point-based system described above.

Results

Further Work

**Support searches based on user's history using more standard supervised machine learning classifiers.

**More complex query searches

**Spell-correction (edit distance)

**Play around with other schemas for determining relatedness coefficient

**Incorporate into the actual site.

References