



# REDUCE DATA OVERLOAD

*A geographic search constraint does an excellent job of reducing data overload, freeing limited computer and personnel resources to focus on more relevant information.*

## DON'T FIGHT MASSIVE FILES.

Discover the advantages of geographically faceted data searches.

By Dr. Craig A. Knoblock, Geosemble Technologies ([www.geosemble.com](http://www.geosemble.com)), Manhattan Beach, Calif.

**D**ata overload is a major problem for organizations, and the problem is getting worse across all industries. Recent IDC research indicates that information workers typically spend a staggering 17.8 hours per week searching for and gathering information. That's a cost of more than \$31,000 per worker each year, assuming a \$55,000 average employee salary with 30 percent benefits.

At the enterprise level, larger companies and government organizations are wasting vast resources—as much as \$5.7 million annually for a 1,000-person organization—searching for and re-creating existing information. Tellingly, IDC notes that “Automating repetitive steps and eliminating those that waste time will increase information worker productivity and save an organization millions of dollars.” In this context, let's discuss automating those repetitive steps to save time and resources.

### How Does the World Deal with Data Overload?

There are two main tools for reducing data overload through search: topic filtering and

time filtering. Both do a good job of reducing data overload, as they eliminate information that's irrelevant to your interests. However, most searches for information carry an unexpressed or under-expressed user assumption: “Limit my results only to the area in which I'm interested.”

Whether you're searching near where you are, where you plan to be or within some pre-defined area of responsibility, a geographic search constraint does an excellent job of reducing data overload, freeing limited computer and personnel resources to focus on more relevant information.

In short, with topic and time filtering, an effective geofaceted search capability can be an important contributor to reducing costs and accelerating knowledge in organizations that have areas of geographic responsibility. Given these benefits, it's worth considering some technical approaches for automatically linking textual content to places and compare some best-fit scenarios for the different techniques.

### Geographically Faceted Search and Discovery

The National Academy of Sciences estimates that 80 percent of online content contains geographic information—much of it unassociated with any address or latitude/





## Term Weighting with TF-IDF

Document source: Old Bailey Online t18100110-41

Term Weighting:  None  Term Frequency  Raw Document Frequency  Inverse Document Frequency  TF-IDF

Stopwords greyed out

charles bailey was indicted for feloniously stealing on the 29th of december two dressed deer skins value 20 the property of samuel savage and richard savage richard savage is a leather seller 63 chiswell street my partner name is samuel savage few days previous to the 29th of december looked out seventy skins for an order these skins being of a bad colour directed them to be brimstoned to make them of equal colour pale on the 29th in the afternoon saw them all smooth on a horse few hours afterwards they appeared very much tumbled and one was thrown into the yard and dirtied caused them to be brought in the warehouse and counted there was gone our foreman went to worship street and brought armstrong and vickrey they searched and found this skin in the prisoner breeches and the other was found in the workshop carter is a foreman to samuel and richard savage the seventy skins was with a savage looking them out i took them out of the stove counted them on the horse and on friday counted them three times over there were no more than sixty eight instead of seventy went to worship street brought an armstrong and vickrey with as they waited till the men left work and when they came down they were searched and on the prisoner one skin was found john armstrong went to this gentleman house after the men came down vickrey and i were searching in one minute vickrey called as i received this skin from him it was taken out of the prisoner breeches i have had it ever since john vickrey q you were with armstrong a yes while i was searching another man i saw the prisoner very uneasy and his breeches were unbuttoned i put my hand in and took this skin out he said he could not tell how it came there the property produced identified the prisoner said nothing in his defence called four witnesses who gave him a good character guilty aged 27 confined six months in the house of correction and fined l a second middlesex jury before an recorder

A term's weight increases proportionally to the number of times it appears in a document, but it's offset by the term's frequency in a repository of documents. This helps to control for the fact that some words are generally more common than others.

document. When such a system performs this linking, it also computes a corresponding score, which captures the confidence level that a document is about a given location. In some cases a document may have more than one geographic focus, and the system assigns a score to each location.

Because the term-frequency approach doesn't need to separate out geographic references, the system can use other types of information to perform the linking, such as names of businesses, street names, people who work there, phone numbers and other associated information. In general, just because a geographic location is mentioned in a document, the system wouldn't link it to that location. Rather, there would need to be sufficient evidence in the document that the location was a topic of the document.

Another important advantage of the term-frequency approach is the ability to perform fine-grained linking to locations. This means that instead of merely linking documents to a city or general area, the approach can link documents down to specific buildings, individual businesses or even people associated with locations.

The GeoXray product performs this fine-grained linking by using a gazetteer with specialized "place signatures" for a region and then computing the documents that link to each of the individual locations. The fact that the term-frequency approach determines the overall geographic focus for the docu-

ments makes this fine-grained linking possible. Otherwise, a system would be overwhelmed with a detailed gazetteer if it tried to link each item mentioned in a gazetteer to an individual location.

For example, consider what would happen using the NLP-based approach if you put "McDonald's Restaurant" in the gazetteer—McDonald's has more than 30,000 locations worldwide! This ability to perform fine-grained linking makes it possible to build

ments that mention a specific location can be performed quickly. The primary disadvantage of this approach is that it focuses on disambiguating only the geographic terms in the gazetteer, and it's difficult to accurately compute the overall geographic or topical focus of a document. In addition, because of the complexity of processing natural language syntax and rules, it's computation heavy, and new rules must be produced and processed for each language.

One disadvantage of the term-frequency approach is that instead of preprocessing all documents and being able to determine every

## THE TERM-FREQUENCY APPROACH combines all of the terms in a document to determine the document's geographic focus.

applications where linking down to specific buildings or businesses, such as "the McDonald's on Culver Blvd.," is required.

### Pros and Cons of Each Approach

Both approaches have advantages and disadvantages. The most applicable approach depends on the details of the specific application.

The NLP-based approach works well with a large repository of documents and an application that requires finding any mention of a geographic location within those documents. Because each reference can be assigned at processing time, it also means the documents can be fully processed in advance and finding the docu-

ment that mentions a location, the approach must be given the location of interest first, then it finds the relevant documents. However, in practice, most users know their geographic area of interest, potentially mitigating this disadvantage.

On the positive side, the term-frequency approach combines all of the terms in a document to determine the document's geographic focus. In addition, this approach can handle a much more fine-grained gazetteer and exploit nongeographic terms, thereby improving the ability to accurately link documents to locations and opening the option to geofacet entities as well as places. And because it's a text-matching technique, term frequency scales well and works in any language. [E]

The screenshot shows the GeoXray web application. At the top, there's a search bar with "Time Filter" and "Topic Filter" options. Below that is a map of Syria with a red pin and a callout box. The callout box contains the text: "Syria faces outrage 'smell of death' in Homs - Reuters" and lists related terms: "Baba Amro (Political or Postal Boundaries)", "Baba Amro, Baba Amro, Bab Amro, Bab Amro, Bab Amro, Bab Amro, Bab Umar, Bab Amro". Below the map is a "Mapped Documents" section with a table of results:

Document Title	Source	Date
سوريا: هجمات جديدة في حمص	Twitter	2012-03-06
#Syria		
(03-06-12) Al-House   Homs   (+19) Individual with MR Tortured, Mother Prays Against Assault	Blog	2012-03-06
بريطانيا: المذبحة 2012-3-6 سوريا	Blog	2012-03-06
BBC News - Syrians fleeing Homs accuse troops of atrocities	Blog	2012-03-06
UN says has similar video of Syrian hospital torture - Reuters	News	2012-03-06
Syrian Hospital Torturing Patients, Staffer Says	Blog	2012-03-06
Syria: Homs: Alqaqar: Martyr died under torture 6 March 2012	Blog	2012-03-06

GeoXray allows users to understand the content landscape and drill into areas of interest. For example, by reviewing content about Syria, relevant topical content is pinned to the actual point of interest.