

# Significance Information for Translation: Air Quality Data Integration

Andrew Philpot, Patrick Pantel and Eduard Hovy

Digital Government Research Center  
USC Information Sciences Institute  
4676 Admiralty Way  
Marina del Rey, CA 90292

{philpot,pantel,hovy}@isi.edu

## ABSTRACT

The management of air quality involves local, state, regional, national, and international organizations. At each level, data are collected and used for analysis, assessment, and regulatory enforcement. Effective air quality management requires coordination among multiple organizations and, therefore, requires integration among their respective data sets. This integration remains a complex IT challenge due to the variety of collection, storage, format, and dissemination methods employed by each organization. In most cases today when organizations need to share data, specialized arrangements and significant manual effort are required to create usable mappings between the data sources. More general methods are required to bridge the gap between representations and organization schemes. We present an interactive web-based demo of our preliminary work adapting the statistical alignment and clustering methods from cross-language statistical machine translation. Using the demo, users can discover new relations and test likely candidate relations between two similar data sources from local and California state air quality management agencies.

## Categories and Subject Descriptors

H.2.2 [Database Management]: Heterogeneous Databases.

## General Terms

Algorithms, Experimentation.

## Keywords

Database Integration, Mutual Information.

## 1. APPROACH

Aligning comparable databases is a challenging and important problem. Despite the wealth of information contained in these information sources, lack of metadata, differences in formats, representational emphases, and encoding schemes all combine to render the integration effort human-centered, time-consuming and error-prone. This problem is faced by thousands of large enterprises with numerous data collections, from Government agencies at all levels to the chemical and automotive industries to startup companies that link together and integrate websites.

Recently, automatic alignment approaches have shifted focus to data-driven techniques which can discover relationships inherent in data sets, without use of any particular form of metadata. Inspired by the analogous problem of cross-language Machine

Translation (MT), we are investigating the use of an information theoretic model to perform data-driven alignments.

The key to our approach is to first identify, using an information-theoretic model, the most informative data elements and then match data sources that share these informative elements. For example, in a data set modeling Santa Barbara, California facilities required by statute to submit emissions data to local regulatory agencies, terms such as "lease" (as in oil rig lease), "Vandenberg AFB," "auto body", and especially "Santa Barbara" are very common; while terms like "Ford" or "Oregon" or "Wingerden" occur only rarely. Therefore, a pair of columns both of which contain "Vandenberg AFB" are intuitively not as similar as two columns both of which contain "Wingerden." Our system, which we call Similarity Information for Translation (*SIFT*) automatically detects and exploits situations corresponding to this intuitive case, assigning higher similarity relationships to those columns (or rows, or tables, with some modifications) which contain less common and thus more significant similarities.

In information retrieval, clustering, and related areas, objects (typically documents) are often represented in a feature vector space. Each word or aspect of the domain corresponds to one dimension of a multi-dimensional feature space. A document is then a vector trajectory through this hyperplane, with the value in a given dimension being a statistic (e.g., frequency of that occurrence) of the associated feature. Similarly, in our task, we model each database column using its data elements; the feature statistics are the pointwise mutual-information between a column and a data element (Church and Hanks, 1989). We also represent columns using expanded feature sets suggested by the data domains of the columns. For example, a column of mostly English text phrases can be modeled by extracting the words or three-character "trigrams" from each text field; a column of 10 digit phone numbers can be modeled by the area code, exchange, and suffix, etc. Expanded feature sets can increase the likelihood of a match given two different representations; they can also introduce spurious matches.

Given two columns represented by such feature vectors, their similarity is computed using a vector distance or alignment metric. We use cosine similarity (Salton and McGill, 1983) which has the very nice property of not being very sensitive to 0-frequency features. In other words, the absence of a matching feature does not indicate dissimilarity as strongly as the presence of a matching feature indicates in similarity.

Similar columns in *SIFT* are discovered using a clustering algorithm called CBC (Pantel and Lin, 2002). We also report

