



Biblio: Automatic meta-data extraction

Carl Staelin, Michael Elad¹, Darryl Greig²,
Oded Shmueli¹, Marie Vans
HP Laboratories Israel
HPL-2004-190
October 25, 2004*

E-mail: {carl.staelin,darryl.greig,marie.vans}@hp.com, {elad,oshmu}@cs.technion.ac.il

document
understanding,
learning, support
vector machines,
neural networks

Biblio is an adaptive system that automatically extracts meta-data from semi-structured and structured scanned documents. Instead of using hand-coded templates or other methods manually customized for each given document format, it uses example-based machine learning to adapt to customer-defined document and meta-data types.

We provide results from two document corpuses, a set of scanned journal articles and a set of scanned legal documents. The first set is semi-structured, as the different journals use a variety of flexible layouts. The second set is largely free-form text based on poor quality scans of FAX-quality legal documents. We demonstrate accuracy on the semi-structured document set roughly comparable to hand-coded systems, and much worse performance on the legal documents.

* Internal Accession Date Only

¹Department of Computer Science, Technion

² Hewlett-Packard Laboratories Bristol

Biblio: Automatic meta-data extraction

CARL STAELIN MICHAEL ELAD * DARRYL GREIG †
ODED SHMUELI ‡ MARIE VANS

HP Laboratories Israel[§]

Abstract

Biblio is an adaptive system that automatically extracts meta-data from semi-structured and structured scanned documents. Instead of using hand-coded templates or other methods manually customized for each given document format, it uses example-based machine learning to adapt to customer-defined document and meta-data types.

We provide results from two document corpuses, a set of scanned journal articles and a set of scanned legal documents. The first set is semi-structured, as the different journals use a variety of flexible layouts. The second set is largely free-form text based on poor quality scans of FAX-quality legal documents. We demonstrate accuracy on the semi-structured document set roughly comparable to hand-coded systems, and much worse performance on the legal documents.

1 Introduction

Biblio is a system for automatically extracting information from scanned documents. A scanned document is an electronic representation of a collection of pages. The electronic representation typically consists of an image component and an Optical Character Recognized (OCR) component. Electronic documents may be viewed as images, or as text objects, and contain both representations. Documents may be annotated with additional information in an extensible format (XML). For example, documents may contain hyperlinks, tables, logos, line art, or other automatically recognized information. A document may be as small

* *Present address: Department of Computer Science, Technion*

† *Present address: Hewlett Packard Laboratories Bristol*

‡ *Present address: Department of Computer Science, Technion*

§ Haifa, Israel, E-mail: {carl.staelin,darryl.greig,marie.vans}@hp.com, {elad,oshmu}@cs.technion.ac.il

as a receipt, or as large as a thesis. Users annotate documents which Biblio uses to learn how to classify documents and extract user-specified information.

The space of documents contains a range of document structure types from highly structured forms such as invoices or tax returns through partially structured business letters to unstructured text. Biblio was designed to make it easy to add the ability to automatically extract information from the full range of document types. The current version of Biblio works with highly structured and partially structured documents.

The goal is to identify the document type, and then to recognize the relevant meta-data embedded in the document. A chief design goal was maximal flexibility and adaptability, with zero document-type-specific coding of any sort. This is achieved via statistical and machine learning techniques. Another design goal is fast operation during recognition, with the assumption that training might be done during off-hours and if necessary could consume significant computational resources.

Most machine learning methods attempt to learn, or fit a model to, a $\mathbb{R}^n \mapsto \mathbb{R}^m$ function. Unfortunately, it is difficult to map whole documents into a fixed-size \mathbb{R}^n vector in a way that captures sufficient semantic information to make meta-data recognition and extraction possible. The Biblio architecture breaks documents down into a fixed size blocks using sliding windows and layout information to facilitate this requirement.

Processing is done in several phases, with each phase typically analyzing the document in greater detail at a finer resolution. Various alternatives are pruned at each step to minimize the overall required computation.

Biblio operation has two major modes: training and recognition. The training system adaptively modifies the system based on user feedback and input. The recognition system classifies documents it has been trained to recognize. Biblio is intended to be an automated system with as little user interaction as possible. In addition, Biblio does not require the user to understand anything about how Biblio works; all user-interactions are in terms of items that are of interest to the user, such as defining or specifying document and document-specific types.

The user interacts with Biblio indirectly through a document browser. While viewing a document, the user can specify the type of the document (e.g. “bank statement” or “journal article”) and may establish new document types. Each document type has an associated list of *meta-data* types (e.g., “account number”, “author”). The browser allows the user to highlight text on the document and indicate that the highlighted text is of a particular meta-data type. Users can also define new meta-data types for a particular document type. The browser will save the user-specified document type along with the meta-data in the document. Later Biblio can use that information during self-training to improve its recognition capabilities.

During document acquisition, as part of the scanning process, Biblio will annotate documents with the information that it recognizes. Users can over-ride the automatically generated annotations, correct mistakes, or add information using the document browser. Biblio may use the modified document during re-training to improve its performance. Recognition run-time requirements are stringent since it is expected that the user is waiting for the processing to complete.

The training system requires extensive processing and is expected to run without user guidance. It utilizes the sample (user annotated) documents and the document types along with their associated meta-data types to generate the data structures and systems utilized at run-time by the recognition system.

The next section describes prior work in document classification. Section three discusses the machine learning and information retrieval technologies on which Biblio relies, including neural networks and support vector machines. In sections four and five, we give a description of the system and the major software components. Section six delineates our experimental setup for testing Biblio while section seven reports on results of those experiments. We conclude in section eight along with a discussion of future work.

2 Prior Work

A great deal of work on document classification and analysis has been done in the areas of document management systems and document recognition (e.g., page decomposition and optical character recognition). In this article we are concerned with techniques for identifying the type of a new document from some set of known document types. This is intended for a document understanding system that puts no prior conditions on the documents presented to the system. This contrasts with much of the current work (e.g. [18] and the proceedings of ICDAR, DAS), which often focuses on only one or more well defined document types. For example, title pages of books [33], journals [11], business cards [32], business letters [6, 36, 35], or office documents [23]. In each of these articles, the specific semantic or structural characteristics of the particular document type under consideration is exploited to analyze the document contents. There have also been a number of more general systems proposed which deal with multiple document types. Lam [12] presents a general document understanding framework (see also Srihari [26]) which contains different processing elements for different document types. However, it appears that document types are separated by some kind of identity string printed on the document itself and subsequently recognized by the system, rather than structure-based type detection.

In Wenzel [35] documents are represented as directed graphs with vertices given by the compounds obtained from a segmentation algorithm. This technique allows both document type identification and document meta-data classification by using graph isomorphisms to find the best match in existing document databases. Note that this technique is similar to

the constraint solving approach used by Lam [12], which could equally be used to classify document types. Another scheme for processing general document types is presented in Taylor [29], but once again the problem of automatically determining the type of a new document is not broached. In Casey [4], an input document is classified as one of a number of known form types by matching the lines found on the page with the form type database. This procedure relies on the fixed structure and scale of the documents involved.

The work most similar to ours is Wnek [37], which employs automatically generated document templates using inductive learning from annotated example documents.

CiteSeer [13] analyzes PostScript and PDF documents to extract meta-data and citation information. It uses hand-coded parsing and recognition engines to extract both document meta-data such as author, title and date, and citations to other documents.

3 Machine Learning and Information Retrieval

Biblio exploits techniques from both machine learning (neural networks and support vector machines) and information retrieval (term weighting methods in text retrieval) to perform document recognition.

Machine learning can be described abstractly as a black box that predicts a vector of outputs when given vector inputs. Generally this involves a “training” phase, and an operation or prediction phase. During training, the system is exposed to a number of input vectors together with the desired output vector. The system uses these examples to “learn” how to respond to inputs. During operation the system is simply given the input vector and it uses its stored knowledge to predict the output vector.

One key element of Biblio’s design goals are that it can operate independently without a skilled operator tuning the machine learning system(s). Most neural network and support vector machine implementations require the user to tune one or more parameters to obtain the optimal performance. Since this was not possible in our case, we either developed or adopted techniques to automatically control and optimize the machine learning meta-parameters.

3.1 Support Vector Machines

Support vector machines are a kernel-based approach to machine learning [3, 5, 30, 19]. We used the publicly available system, LIBSVM, as the support vector machine engine. LIBSVM includes four kernel functions: linear, polynomial, radial basis function (RBF), and sigmoid. We used the RBF kernel function which is defined as:

$$RBF \quad k(\vec{u}, \vec{v}) = e^{-\gamma|\vec{u}-\vec{v}|^2} \quad (1)$$

During training, Biblio's SVM engine creates input vectors representing the words found in the training documents. Biblio uses the RBF kernel function with autonomous SVM parameter selection [28] to create a set of support vectors for each type of known meta-data. These vectors are then stored along with the parameter settings to be used during the recognition phase.

3.2 Evolutionary Neural Networks

Neural networks are an example of a machine learning technology. Biblio uses Hplinet[27] for both training and operation. Hplinet is an artificial neural network package with an evolutionary mechanism for automatically configuring network architecture based on ideas from Yao [38]. It employs training algorithms that do not require user defined parameters, such as BFGS [2] and LBFGS [14], while using techniques from evolutionary algorithms to search for the optimal network architecture. The networks use a Generalized Feed-Forward architecture. It also uses committees of networks in order to increase the stability of predictions.

Hplinet uses gradient-based training techniques to train networks for a given architecture and training set. Depending on memory availability, the system switches between full quasi-Newton training (BFGS), and a limited memory version (LBFGS). The major difference between these two algorithms is that BFGS stores an approximation of the full inverse Hessian matrix of the network function with respect to the network weights ($nLink$ by $nLink$ elements), whereas LBFGS ignores the off-diagonal elements and stores only the diagonal of this matrix. The LBFGS algorithm pays a penalty in convergence time, but does offer a reasonable alternative in low memory situations.

The network architecture must also be chosen. This means determining the number of hidden nodes and the structure of the links connecting input, output, and hidden nodes. The size and complexity of the architecture is generally a good measure of the complexity of the output function, and it can have a dramatic effect on both the duration of training time and the accuracy of the resulting network.

It is well known that neural networks often converge to local minima in the error surface, and that different initial weight configurations can lead to networks of identical architectures giving quite different solutions on the same training data. In addition, neural networks often have areas of the input space where they perform very well, but they may have areas where they perform poorly because the network has not learned to represent the function accurately in that region of the input space. A more stable overall solution can be obtained by training multiple networks and combining them into *committees*.

Committees are groups of neural networks, often trained on the same data but with different initial weight settings. During runtime each network is given the data and their outputs are combined to form a final output using any of a number of possible methods, such as voting.

In this way, overall accuracy of the system can be increased by factoring out most of the erroneous predictions.

In general, Biblio trains ten networks for each problem. It uses a separate validation data set to rank the performance of the networks, and chooses the best five networks to create a committee. The committee result is created by averaging the outputs of the five committee members. For classification problems, the real-valued network outputs are thresholded to binary values and a majority vote is used instead of a simple numeric average.

3.3 Information Retrieval

Broadly speaking, information retrieval techniques operate on text streams, and they are commonly used for managing and searching large text corpuses such as the internet or a library. They are also used for routing and filing documents such as emails into folders. The most common approach is to build a dictionary of known terms, and to assign each term a unique index. Documents can then be represented as a vector, with the values typically computed as a function of the term frequency within the document. A Boolean model would simply set the value to 1 if the term appears and 0 otherwise. A more powerful and commonly used function is TFIDF [22] which is computed based on the term's frequency within the document multiplied by the inverse of the frequency with which the term appears in all documents.

Biblio uses word dictionaries to determine if text is indicative of a particular type of meta-data. There are two types of dictionaries, meta-data type dictionaries and special dictionaries. Meta-data dictionaries are used to represent the words associated with a given meta-data type. Special dictionaries are used to represent words that may be common across meta-data types, such as names. The function of meta-data dictionaries is to give the probability that the given text stream contains text representing a particular type of meta-data. During training the system is given example text streams from the training documents. The text streams may or may not contain meta-data. Currently the dictionaries are updated manually, however, the dictionaries need to be self-updating without user intervention during training.

There are a number of models and approaches for text classification, such as naive Bayes [21] and SVMs [17], however, we currently use SVM with RBF kernels. The SVM engine also uses a dictionary that contains a unique id for every word it encounters. Each time it discovers a new word, a unique id is generated and the word is stored in the dictionary permanently. During training, the SVM engine builds input vectors using ids from the dictionary to identify words in the document. This allows the SVM to create support vectors that identify the most probable words associated with a particular type of meta-data. During operation, the dictionary is again used by the SVM engine to build the input vectors for use with the SVM model created during training.

4 System Description

There are two modes of operation for Biblio: training and recognition. Biblio uses several analysis engines for both types of processing, however, they are used in different ways. This section describes both systems, the analysis stages and how they differ between the two modes, and the major data structures used in training and recognition.

4.1 Recognition System

For recognition, Biblio is designed to eliminate many possibilities early to reduce the processing requirements.

The recognition system is structured as follows: There are some number of user-visible document types, such as “Business Letter”, “Journal Article”, or “Bank Statement”. Each document type has some meta-data types, such as “Author” or “Account Number”. There are typically many example documents of a single type, which are available to the system as it trains and retrains itself. Within a single document type, Biblio may break documents down into document classes, which are generally documents that are similar to each other within the document type. Each document class has its own set of neural networks and SVM models for detailed processing. The system itself creates the document classes from the pool of documents in the document type.

Recognition is done in stages, with possibilities filtered at each stage. Before the first stage, Biblio has no knowledge about the document. It could belong to any of the document classes, and could contain any type of meta-data. To minimize the computational costs, Biblio tries to eliminate as many candidate document types as quickly and cheaply as possible.

The first stage, compound analysis, is both content and layout-based and is done for each document type. Each compound is flagged with the meta-data types that may be contained in that compound. Compounds perceived to contain no meta-data are discarded and ignored in further processing.

In the second stage, document type analysis, the system evaluates the probability that the document is of a particular type. For each document type, the system uses the information from the compound layout information to compare the document’s structure to documents of that type. Biblio chooses the closest match, or rejects all candidate types. From this point forward, all processing is done with respect to a single candidate document type.

There are three following stages, paragraph analysis, line analysis, and word analysis, which iteratively analyze the compound, paragraph, and line to select or reject paragraphs, lines, and words that contain meta-data. Each paragraph, line, and word is flagged with the meta-data it is perceived to contain, and the flagged pieces are passed on for further processing. At the end, Biblio takes the remaining words and annotates the document with the recognized

document type and meta-data.

4.2 Training System

The training system adaptively modifies the system based on user feedback and input. Training is designed to be run as a batch processing job, preferably during low usage times such as overnight. This is because training may consume resources at an unacceptable rate while the user is actually using the system.

Each document type has several example documents on which to learn. We use neural networks and support vector machines but training may also include actions such as updating word databases or other supporting data structures in response to the user's feedback.

5 Analysis Engines

Biblio uses a variety of analysis engines and techniques to evaluate the document at each analysis stage. Analysis engines create evidence, and are called witnesses. Judges assemble witnesses, collect and collate evidence, make judgements, and propagate verdicts. Judges use special witnesses, called chairmen, to produce final evidence for each meta-data type in the document class, and this evidence is used to produce the verdict. This modular architecture makes it easy to add new analysis engines.

Biblio uses a layered hierarchy of analysis engines to develop and iteratively refine evidence and classification information about each information layer: compound, document, paragraph, line, and word. Each layer serves as a filter; portions that are rejected are not processed further while the remaining portions are passed on to the remaining filters. These layers correspond to the five analysis phases. Each analysis phase has a complete set of analysis engines trained exclusively for that phase and document class.

5.1 Compound Analysis

At the beginning of compound analysis during recognition, each document contains a list of compounds, which are really OCR recognized regions. The compound analysis is run separately for each candidate class, and the output for each class is only passed to the document analysis for that class. Each document class has its own compound analysis neural network.

Compound analysis annotates each compound with the likelihood that each type of meta-data is contained in the compound. It only examines the list of meta-data types that are associated with the particular document class.

Compound analysis uses a neural network to annotate a single compound at a time. Since the number of meta-data types for each document class is fixed (at least until the user modifies

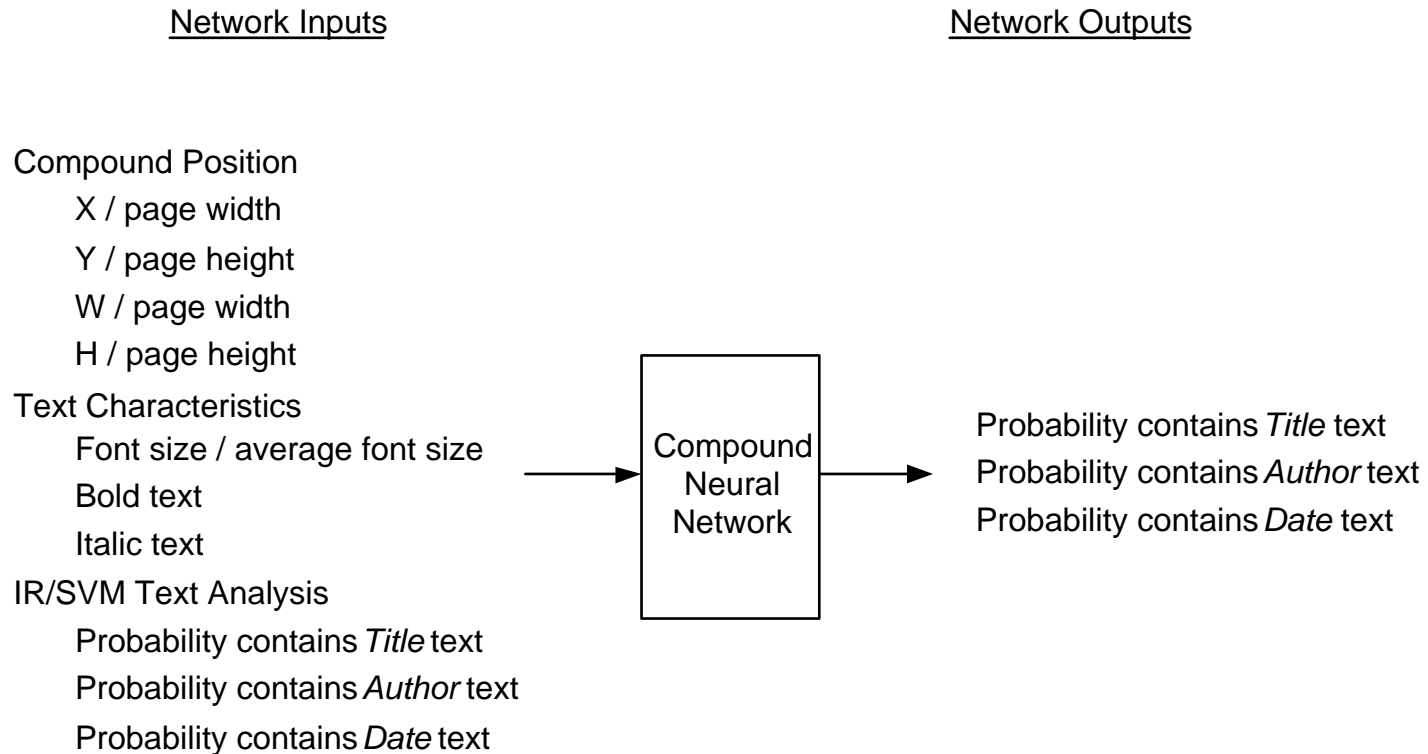


Figure 1: Compound Processing

the list), the neural network can output the probability for each meta-data type directly. Figure 1 is a graphical representation of the network for compound analysis. It shows the inputs and the outputs of the network.

Neural networks require a fixed number of numeric inputs, but document content is not usually fixed size and is usually not described numerically. Biblio translates the compound contents and layout into a fixed number of numeric inputs. The inputs are:

1. the bounding box of the compound,
2. the font size of the text,
3. and the probabilities generated by the support vector machines for each meta-data.
4. the proportion of all words found in the compound that are associated with known meta-types,
5. the proportion of descendent words containing lower case letters,
6. the proportion of descendent words containing upper case letters,

7. the proportion of descendent words containing upper case digits,
8. the proportion of descendent words containing upper case symbols,
9. the proportion of descendent words containing upper case initials,

During training, the system generates the support vector machine output for each compound first. The system then trains a committee of networks to recognize compound types for each document class. It uses the probabilities produced by using the support vectors as well as the layout inputs. It creates sample inputs by feeding the system example compounds with examples containing meta-data and some that do not contain meta-data. The number of training examples is bounded, and training time is reasonable.

5.2 Document Analysis

During recognition, document analysis is done on a per-document class basis. The per-class results are compared to choose the best match. Each document class has its own document analysis neural network.

Document analysis uses a neural network to evaluate the probability that a document belongs to that class. Each meta-data type has a set of inputs to the neural network which include the same inputs as the compound analysis inputs. It uses the compound analysis results to create a pseudo compound based on the union of the compounds containing that meta-data type.

During training, the system trains a committee of networks on each document class. It takes example documents from all classes to train the network to distinguish between document classes. Any time a new class is added, all document analysis networks need to be retrained, which can be very expensive.

In the future, we would like to apply signature filtering [24] that will guarantee only documents from certain document classes will be passed. If we are successful, modifications to radically dissimilar classes should not require retraining of the document network.

5.3 Paragraph and Line Analysis

The paragraph and line analysis is done only for the most likely document type during recognition. It uses the information from the compound analysis to eliminate portions of the document from further processing. Each document class has a neural network for each meta-data type. The neural networks take the previous compound or paragraph analysis results and SVM probabilities as input. For each meta-data type, the networks output the likelihood that the paragraph or line contains that meta-data.

A committee of networks recognizing the paragraphs or lines that contain meta-data are trained for each document class. These networks are trained only with examples from the document class, and training time is generally reasonable. These networks only need to be retrained if there are changes in the target document class.

5.4 Word Analysis

The final step for the recognition process is to examine each word in the remaining lines and identify those that are associated with a particular meta-data type. This is currently done using a neural network that takes token information for three words and outputs information about the middle word.

By giving the neural network information about three words, we give the network some context for the current word, which dramatically improves the network's accuracy. However, these networks are large, requiring more inputs and more hidden units than any other network in the system.

For each document class the system trains, a committee of networks scan the individual words in selected lines. The networks are only trained with examples from lines containing meta-data in the target document class.

Since the system creates a training example for each word, the number of training examples is huge, and the training time is atrocious. This is currently the single largest consumer of CPU time in the training system. These training time problems compel us to look for alternative strategies for extracting the meta-data from each line.

6 Implementation

This section describes the major software classes and their methods during training and operation. Figures 2 and 3 are block diagrams of the major objects and their interactions.

6.1 Portfolios, Hubs, and Evidence

Portfolios contain document hubs, which are objects that contain information about the files on which Biblio is currently working. During training, the portfolio can have as many files as is physically possible to load into memory. In our case, we used a maximum of 50. During processing, the portfolio contains a single document hub, the document it is trying to recognize. The portfolio is mainly responsible for the creation and removal of hubs.

Each hub contains information about the file it represents. The main components include an xml tree representation of the document and evidence objects. The xml representation makes it easy to quickly access any part of the document at any level. In addition, a hash table containing all the words in the document is included for fast word searches. The evidence

Biblio Training Process

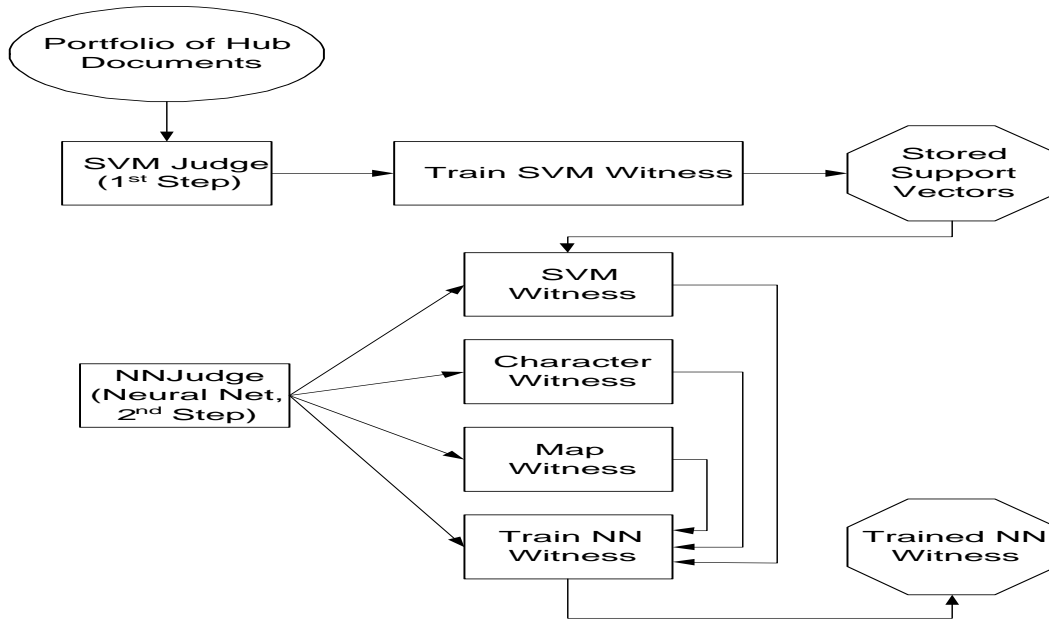


Figure 2: Biblio Training Block Diagram

Biblio Recognition Process

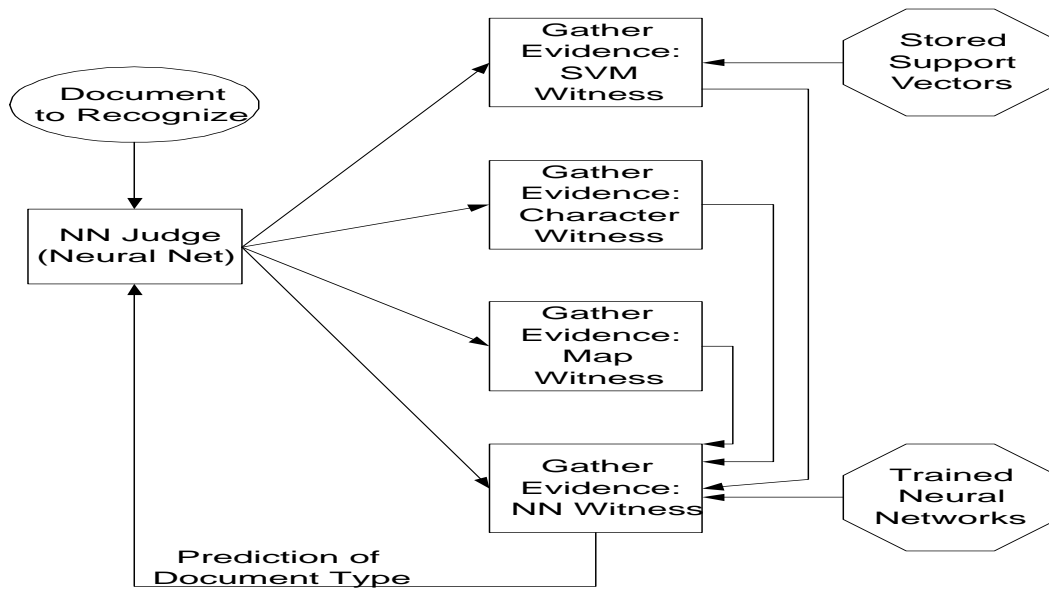


Figure 3: Biblio Operation Block Diagram

objects work as a central repository that is used by witnesses for storing and accessing evidence. During training, the evidence object contains lists of probabilities for each type of known meta-data. Each witness is responsible for determining the probability of existence for each type of meta-data at the current level, i.e. compound, paragraph, line, word. A special witness, the chairman, later collects all the evidence for use in training the neural networks to recognize meta-data. During operation, the evidence object is again used to collect evidence by witnesses and the chairman later collects this evidence. The difference is that during recognition, the chairman uses the evidence as input to a neural net trained to recognize meta-data.

6.2 Judges & Witnesses

Biblio was designed to make it easy to add functionality. This could include alternative machine learning technologies, document processing capabilities, statistical, and mathematical algorithms. This means that whenever a promising technology with application to machine learning or document understanding is invented, it can easily be incorporated into Biblio.

Currently, Biblio has 2 specialized *judges* whose job consists of organizing *witnesses* for specific types of *evidence*. These are the neural-net judge and the support vector judge. Both are derived from a base judge that defines the basic operations such as creating witnesses, training them, initiating prediction, and acting as a repository for evidence returned from witnesses. Another vital task of the judge is to separate data into 3 disjoint sets: training, test, and validation. For our experiments, the test set was given. That is we took one file from the list of files, held it out for test and used the remaining files for training or validation. The judge uses a random procedure that assigns roughly 75% of remaining files to a training set and 25% to the validation set. This design makes it easy to derive any specialized judge that can start up training for specialized witnesses and collect evidence from these witnesses during the prediction phase.

6.2.1 SVM Judge

The SVM judge is responsible for directing the training of the SVM witness. During training, the SVM judge creates the SVM witness and passes it a portfolio of documents, a specification of the type of the documents, and the level at which the witness should train, i.e. compound, paragraph, line, or word. The SVM witness is specialized from a base witness class which contains operations that are generic to all witnesses as well as special operations related to support vector machine processing. Some of the base operations include training and prediction. The main function of this judge is to direct the SVM witness in creating support vector models. These models are stored for use later in the prediction process, a process in which the SVM judge is not currently involved.

THE SVM WITNESS

The SVM witness is responsible for loading dictionary objects. There are two types of

dictionaries used by the SVM witness. The main dictionary object contains all the words from every document ever seen by Biblio during training together with a unique id. When a new document is included during training, the dictionary object generates a unique id and updates the dictionary with the new id/word pairs. This allows the dictionary to grow over time. During training, the SVM witness creates another dictionary consisting of words from the documents in the training set known to be associated with specific meta-data types. For example, for the journal article class of documents, there are files containing words that are specific to the types of meta data found in journal articles like “journal name” and “author”. The main draw back of these files is that there is currently no way to automatically update them, instead they must be edited by hand.

Once both dictionaries are loaded, the SVM witness builds input vectors for the support vector machine, one set of vectors for each type of known meta data. First, the witness creates 5 or 6 support vector entries to record the existence of words from the main dictionary. Each entry in the vector contains the word id and the proportion of the current word contained at current level, e.g. compound, paragraph, line, word. Once the input vectors are created, a target vector is created for each word represented by the vector. The target value indicates whether the word is associated with the particular type of meta-data currently in focus.

Before sending the input/target pairs to the support vector machine to create the final model, a parameter search [28] is initiated to find the best combination of C and gamma values. Support vector machine training is executed using these values and once training is complete, prediction and cross-validation is used to determine how well the system did using those parameter values. Based on these results, new parameter values are generated which are closer to the best values from the previous run and the cycle is repeated. We iterate through the process 5 times, after which the values returned for gamma and C are used for actual training and model creation. Once the model is created, it is written to file where it can be accessed by other instances of the SVM witness during prediction.

6.2.2 Neural Net Judge

In addition to the SVM judge there is also a neural net(NN) judge which is responsible for training and prediction. Similar to the SVM judge, the neural net judge is mainly a supervisor for witnesses. The NN judge creates instances of 4 witnesses: a SVM, a character, a map, and a neural net witness. The character witness is responsible for collecting evidence on upper and lower case letters and the presence of symbols and digits. The map witness collects evidence on bounding boxes, average font size, the relative position of the current compound (paragraph, line, word) within the parent compound, and the relative number of tokens found in the compound. The neural net witness is responsible for collecting evidence from the other three witnesses and creating the input vectors for training the committee of networks.

THE CHARACTER WITNESS

The hub of each file contains a hash table of all the words in the document. Attached to each word is a bit array containing a single bit for each type of attribute the character witness knows about, i.e. upper case, lower case, symbols, digits, and initials. This information can then be used to indicate meta-data with specific formatting characteristics. For example, titles or proper names usually include upper case characters; names often include initials; digits may support the presence of a date or page numbers; URLs and mathematical formulae usually include symbols. This witness first initializes the word bit arrays and fills them in for every word in the document. Later it collects evidence by going through each word in the document and recording the proportion of words containing each attribute in the document. It then stores this information in the evidence object of the hub. This process is the same for both training and recognition.

THE MAP WITNESS

The map witnesses can supply up to seven types of evidence: the bounding box of the type of compound, i.e. compound, paragraph, line, word; the relative position of the compound within its parent, the average number of tokens in the compound, and the average font size in the compound. Boolean flags passed to the witness indicate which evidence needs to be collected. In some cases various types of evidence may not be pertinent to the type of document on which the system is being trained. The xml representation of the document makes calculating the bounding box efficient as the height and width of every component is stored as part of the compound object when it is created. The relative position is also quickly determined because the compound has access to its parent and the parent is aware of all its children. Bounding box and relative position evidence may be important, especially for highly structured documents such as telephone bills which would have the same information displayed in the same relative location on every bill. The average font size is also easy to determine as the font size of each compound is stored in each compound object. We can calculate an average font size for the compound by dividing it by some maximum number, in our case we used fifty. This tells us whether the average font size is larger or smaller than usual. It can be very helpful to know what the average font size is when trying to identify things such as headings and titles. Finally, the average number of tokens is calculated by determining the number of words subordinate to the current compound and dividing that number by some maximum. Similar to the average font size measure, this gives us an idea of whether the number of tokens is small or large.

THE SVM WITNESS

The SVM witness behaves differently depending on the judge that creates it. If the judge is a SVM judge, the witness creates the SVM model as described above. If the judge is a neural net judge, then the SVM witness uses the model it previously created to collect evidence for both the training and recognition processes. As with model creation, the SVM

witness must create input for the support vector machine from the document. It does this in exactly the same way as when building the model. However, during evidence collection, the SVM witness loads the model, calls the support vector machine prediction module, and passes it the input vector. The SVM witness uses this model, along with the previously determined parameters to predict whether each of the words represented by the input vectors contain meta-data. Prediction returns probabilities for each word in the compound. These probabilities are then put into the evidence object of the corresponding hub.

THE NEURAL NET WITNESS

The neural net witness is the chairman. This means it is responsible for collecting all the evidence produced by the other witnesses and using it to either train a committee of neural networks or to use as input for recognition. During training the neural net builds a single training object containing all the evidence from each hub evidence object in the training set. It also builds a similar object using the validation set of hubs. Target arrays are built that associate the true value to each piece of evidence collected. It then creates a committee of evolutionary neural networks, passing in the evidence and target arrays. The number of networks created is configurable. In our experiments we used committees of ten networks. Once the network has found the optimal architecture and configuration, the witness stores the network for use during prediction.

During prediction, the neural network witness collects all the evidence from the hub evidence object and creates a single input vector. The witness then loads the network architecture from the previously stored file and passes it the input vector. Output from the committee consists of probabilities for each type of meta-data. This is stored back into the hub's evidence object and passed back to the judge. The judge will then be able to use the evidence to decide on the type of document based on the probabilities of meta-data found.

7 Experimental Setup

We used a corpus consisting of a few hundred documents. These documents were either articles from various computer science journals or legal documents recording transfer of properties, i.e. deeds. All were scanned into the system and stored in an HP proprietary file format consisting of both images and text recognized by the scanner. Meta data describing the document and identifying key words (i.e. words that are unique to the type of the document trained on) is embedded within the XML.

We trained Biblio on fifty documents at a time, as the memory requirements for larger document sets was prohibitive. During training, the documents are decomposed into regions consisting of major regions, paragraphs, lines, and words. The system considers each type of component separately and determines the support vectors based on known meta data types for each type of component. Once the support vectors are identified the system uses these

vectors to generate probabilities that each component contains a specific type of meta data. The resultant probabilities are fed into the neural network as evidence for training the neural network. The network takes the evidence and builds a parsimonious architecture based on the best results. This means the final network consists of the architecture with the fewest nodes and is close to the absolute best architecture, which may actually have many more nodes. Biblio then saves this network for use during operation.

During operation, or prediction, the neural network acts as the driver for the entire process. Evidence is gathered on a previously unseen document and the trained network uses this evidence to determine the type of the document. The evidence gathering process starts the SVM process. This process takes the stored support vectors for each type of known meta data and determines probabilities that are subsequently fed into the trained network. Once the network has identified the document type and the meta data contained within it, the system annotates the original file by adding the document type and the recognized meta data.

To test the system, we use leave-one-out cross fold validation for our experiments. For each of the fifty documents we hold one document out, train on the remaining documents, and then predict whether the document held for testing contained meta-data appropriate to the type of document we trained on. During prediction we collect confusion matrix data to report the results.

7.1 Documents and Meta-Types

We ran our experiments on two types of documents: journal articles and grant deeds. The journal articles are an example of structured documents while the grant deed documents are an example of unstructured documents. Within the journal article type there are two classes of documents: articles from IEEE Transactions journals and articles from the HP Journal. For grant deeds there are six registered types of meta-data: Date of Recording - Meta0, County - Meta1, Document Number - Meta2, Return Address - Meta3, Grantor - Meta4, and Grantee - Meta5.

For journal articles there are five registered types of meta-data: Title - Meta0, Author - Meta1, Journal Name - Meta2, Pages - Meta3, and Date - Meta4.

8 Results

The results consist of statistics generated from the confusion matrices collected during our experiments on three sets of data.

Type	Total	TrueNeg	FalseNeg	FalsePos	TruePos	PosCount	NegCount
GRDE-Meta0	48629	48496	121	0	12	133	48496
GRDE-Meta1	48629	48368	143	0	118	261	48368
GRDE-Meta2	48629	48477	98	0	54	152	48477
GRDE-Meta3	48629	47051	869	11	698	1567	47062
GRDE-Meta4	48629	47076	1192	14	347	1539	47090
GRDE-Meta5	48629	47106	1152	35	336	1488	47141
JA00-Meta0	19208	18991	28	0	189	217	18991
JA00-Meta1	19208	19048	32	16	112	144	19064
JA00-Meta2	19208	19044	11	0	153	164	19044
JA00-Meta3	19208	19184	15	0	9	24	19184
JA00-Meta4	19208	19161	3	4	40	43	19165
JA01-Meta0	19703	19463	6	0	234	240	19463
JA01-Meta1	19703	19521	27	17	138	165	19538
JA01-Meta2	19703	19638	31	0	34	65	19638
JA01-Meta3	19703	19670	15	0	18	33	19670
JA01-Meta4	19703	19632	33	0	38	71	19632

Table 1: Confusion Matrix

8.1 Confusion Matrix

Table 1 contains the confusion matrix information for all documents. The *Total* column is the total of all the words in the associated document class. For comparison, the total number of actual words containing meta-data (column labelled *PosCount* and the total number of actual words not associated with meta-data (column labelled *NegCount* are included. Obviously, the majority of the words are not associated with meta-data. During testing, the neural networks used a threshold of .50 when deciding whether or not a particular word was meta-data. There are two types of errors, False Negative (prediction of false when the meta-data is actually there) and False Positive (prediction of true when no meta-data exists). Table 1 shows that Biblio does predict incorrectly. However, as Tables 3 and 4 below show, Biblio does not incorrectly label data very often.

8.2 Sensitivity

Table 2 reports the Word and File sensitivity of Biblio. Sensitivity is defined as:

$$Sensitivity = \frac{TruePositiveCount}{TruePositiveCount + FalseNegativeCount}$$

In our case, sensitivity is a measure of how well Biblio can detect real meta-data in the document when it exists. Results were divided into positive and negative predictions with

Type	WordPosCount	WordSens	FilePosCount	FileSens
GRDE-Meta0	133	[0.04 – 0.14]	114	[0.05 – 0.16]
GRDE-Meta1	261	[0.39 – 0.51]	129	[0.37 – 0.54]
GRDE-Meta2	152	[0.28 – 0.43]	145	[0.29 – 0.45]
GRDE-Meta3	1567	[0.42 – 0.47]	137	[0.36 – 0.52]
GRDE-Meta4	1539	[0.20 – 0.25]	123	[0.13 – 0.27]
GRDE-Meta5	1488	[0.20 – 0.25]	128	[0.16 – 0.31]
JA00-Meta0	217	[0.83 – 0.92]	23	[0.72 – 1.00]
JA00-Meta1	144	[0.71 – 0.85]	23	[0.59 – 0.94]
JA00-Meta2	164	[0.89 – 0.97]	23	[0.85 – 1.00]
JA00-Meta3	24	[0.18 – 0.57]	23	[0.17 – 0.57]
JA00-Meta4	43	[0.85 – 1.00]	23	[0.83 – 1.00]
JA01-Meta0	240	[0.96 – 0.99]	32	[0.91 – 1.00]
JA01-Meta1	165	[0.78 – 0.89]	32	[0.68 – 0.95]
JA01-Meta2	65	[0.40 – 0.64]	31	[0.36 – 0.71]
JA01-Meta3	33	[0.38 – 0.72]	32	[0.37 – 0.72]
JA01-Meta4	71	[0.42 – 0.65]	32	[0.37 – 0.71]

Table 2: Word and File Sensitivity

positive predictions shown in the *WordPosCount* and *FilePosCount* columns. These values reflect the actual number of meta-data words, that is, number of true positive plus the number of false negative predictions during testing. *WordPosCount* is the sum of all true positive and false negative predictions for all files. *FilePosCount* is the sum of true positive and false negative predictions only for the files that actually contained the type of meta data associated with the label on the row. This allows us to see the performance of our prediction system both at the file level and at a higher, multi-file level. The file level is indicative of what a typical user would experience, while the higher level allows us to see larger trends for document and meta-data types. Sensitivity measures are based on *WordPosCount* and *FilePosCount*.

One thing should be noted here about the sample populations. We had approximately three times more samples of the GRDE group than each of the Journal Article group samples. The GRDE group consists of three separate runs of fifty documents each, while the Journal Article groups consisted of one run each with approximately fifty documents. We combined the GRDE runs because there is essentially no difference between the three sets of data. However, there are differences between the Journal Article groups, for example, type of journal and layout.

The *WordSens/FileSens* columns are the 95% confidence intervals for the sensitivity measures.

8.2.1 Analysis

From Table 2 we can see a large difference between the GRDE and Journal Article document classes. In general, Biblio performed poorly in recognition of most meta-data types for GRDE, both at the multi-file and the file levels. This is somewhat expected since GRDE falls into the unstructured document category, which we assumed would be harder to recognize. In addition, the GRDE files have numerous OCR errors. If meta-data words contain OCR errors, it would be difficult for Biblio to recognize them as meta-data. It is also possible that a different threshold value would allow better separation.

On the other hand, Biblio did quite well in recognizing meta-data in the journal articles. This is in spite of the fact that Biblio had less data on which to learn than the GRDE category. The main difference in results between the two types of journal articles are with meta-data type 3 (page numbers). We theorize that this is because the first type, the HP Journal articles, are much less structured than the IEEE articles. The format of HP Journals can vary widely. Meta data types 2, 3, and 4 for the HP Journal articles (Journal Name, Pages, Date) are not consistently placed on every page or between different articles. Another factor affecting the results could be related to the number of each meta-data type found in the test file. The table shows significantly better results for meta-data types with many examples. For example, both Meta0 and Meta1 (title and author) appear hundreds of times. The associated confidence intervals for both types of journal articles have upper ranges close to 1.0. The same is true for Meta2 (Journal name) for the IEEE journal articles. This would make sense if we can assume that the number of positive examples of each type of meta-data is representative of what we would find in a “typical” IEEE or HP Journal article, The fact that Biblio did very well for Meta4 on the IEEE journal articles in spite of the relatively small number of samples could be related to the fact that IEEE articles have consistent formats.

8.3 Specificity

The vast majority of the predictions were negative, i.e. words that were not meta-data. This is the value containing the true negative plus the false positive predictions. These are the number of words and words in files containing meta-data, respectively, that did not actually contain any of the associated type of meta-data labelled on the row. Specificity is defined as:

$$Specificity = \frac{TrueNegativeCount}{TrueNegativeCount + FalsePositiveCount}$$

or how well Biblio classifies nonmeta-data. We do not present a table for specificity because all specificities were either 1.00 or 0.99, with 95% confidence intervals between [0.99 – 1.00] both at the file and multi-file levels. In general, this means that our prediction system is very good at identifying words that do not contain meta-data.

Type	FalsePos	TruePos	WordRate
GRDE-Meta0	0	12	[0.000 – 0.000]
GRDE-Meta1	0	118	[0.000 – 0.000]
GRDE-Meta2	0	54	[0.000 – 0.000]
GRDE-Meta3	11	698	[0.006 – 0.024]
GRDE-Meta4	14	347	[0.019 – 0.059]
GRDE-Meta5	35	336	[0.065 – 0.124]
JA00-Meta0	0	189	[0.000 – 0.000]
JA00-Meta1	16	112	[0.068 – 0.182]
JA00-Meta2	0	153	[0.000 – 0.000]
JA00-Meta3	0	9	[0.000 – 0.000]
JA00-Meta4	4	44	[0.006 – 0.176]
JA01-Meta0	0	234	[0.000 – 0.000]
JA01-Meta1	17	138	[0.060 – 0.159]
JA01-Meta2	0	34	[0.000 – 0.000]
JA01-Meta3	0	18	[0.000 – 0.000]
JA01-Meta4	0	38	[0.000 – 0.000]

Table 3: Positive Word Prediction Rate

8.4 False Positive Ratios

Another interesting statistic shows the proportion of meta-data that Biblio predicted as being present that is not really there. Table 3 shows these results. The *WordRate* column reports the 95% confidence interval for the statistic. The ratio is computed using:

$$NegativeWordRate = \frac{FalsePositiveCount}{TruePositiveCount + FalsePositiveCount}$$

This is an important result for potential users of Biblio. Because the majority of the words in any document will most likely not be meta-data, if Biblio reports too many false positives, users will spend time correcting errors. While it is not clear whether this problem is as severe as missing meta-data that does exist, it is an annoyance likely to keep users from using the system. As the table shows, one problem occurs with journal article Meta1 (Author). for both journal article classes. On average, when Biblio predicted Meta1 type for the JA00 class, 12% of the time it was incorrect. For JA01 the rate for Meta1 positive misclassification is 0.11. It is encouraging, however, that for 10 of the 16 metatypes, Biblio made no false positive predictions. The rest of the positive misclassifications occurred at a rate of 9% or less.

Type	FalseNeg	TrueNeg	WordRate
GRDE-Meta0	121	48496	[0.002 – 0.003]
GRDE-Meta1	143	48368	[0.002 – 0.003]
GRDE-Meta2	98	48477	[0.001 – 0.002]
GRDE-Meta3	869	47051	[0.017 – 0.019]
GRDE-Meta4	1195	47076	[0.023 – 0.026]
GRDE-Meta5	1152	47106	[0.023 – 0.025]
JA00-Meta0	28	18991	[0.001 – 0.002]
JA00-Meta1	32	19048	[0.001 – 0.002]
JA00-Meta2	11	19044	[0.000 – 0.001]
JA00-Meta3	15	19184	[0.000 – 0.001]
JA00-Meta4	3	19161	[0.000 – 0.000]
JA01-Meta0	6	19463	[0.000 – 0.001]
JA01-Meta1	27	19521	[0.000 – 0.002]
JA01-Meta2	31	19638	[0.001 – 0.002]
JA01-Meta3	15	19670	[0.000 – 0.001]
JA01-Meta4	33	19632	[0.001 – 0.002]

Table 4: Negative Word Prediction Rate

8.5 False Negative Ratios

The proportion of words that Biblio predicted as nonmeta-data words when they were actually meta-data words is reported in Table 4. The *WordRate* column reports the 95% confidence interval for the statistic and the ratio is computed using:

$$NegativeWordRate = \frac{FalseNegativeCount}{TrueNegativeCount + FalseNegativeCount}$$

The table shows that Biblio does better in terms of predicting false negatives. The highest rate occurs for GRDE-Meta3 and GRDE-Meta4 and the average for these metatypes is around 2%. On the other hand, since the majority of the words in a given document will be negative, the number of false negatives translates into the number of corrections the user would have to make. Because the number of meta types in a given document is small compared to the number of non-meta types, failing to identify even a small percentage of the existing meta data could cause Biblio to incorrectly classify the document as a whole.

9 Conclusions and Future Work

Biblio is a sophisticated document type recognition and meta-data extraction engine that utilizes a combination of machine learning technologies to adapt to user-defined document

types and meta-data fields.

While Biblio is a fully functioning system that represents a major improvement in ease-of-use and adaptability, there are still a number of areas that we should like to improve.

For example, currently all documents within a single user-defined class are treated as a unitary class. However, it is likely that a single class, such as “Journal Article” or “Invoice”, would have sub-classes such as “phone bill” and “gas bill” which are far more regular. However, we would like to automatically discover clusters of similar documents within a user-defined class so as to take advantage of similarities without burdening the user.

Another area that we should like to improve is the ability of the system to recognize when a document is not a member of any known document class. While the system does this today, it is not sufficiently selective so documents can be annotated with spurious meta-data.

We are examining the possibility of using techniques similar to the methods used by Koller and Sahami [10] to identify key words in document clusters. This would involve analyzing the statistical patterns of text containing meta-data versus text that does not contain meta-data to give reasonable probabilities for unseen strings. The Naïve Bayes classification algorithm may be suitable for classifying text as containing / not containing meta-data [20]. The word dictionaries could also be automatically enhanced using the SONIA clustering method.

Additionally, we would like to improve the ability of the system to make use of *tagged* fields, such as fields commonly pre/post-fixed with identifying strings such as “To:” or “Subject:”. This requires analyzing neighboring regions to see if they indicate the presence or absence of meta-data in a given region, and this evidence can then be passed to the region analysis.

Finally, Biblio is currently a single-page analysis system, meaning that each page is analyzed independently. Sometimes documents contain important meta-data on more than just the first page. For example, some journals only put the journal name, volume, and number on subsequent pages of an article. We should like to extend the system so it can analyze a whole document, and extract any relevant meta-data regardless of page.

10 References

- [1] S. BAUMANN AND M. MALBURG, H.-G. HEIN, R. HOCH, T. KIENINGER AND N. KUHN, *Document analysis at DFKI, part 2: Information extraction, German Research Center for Artificial Intelligence (DFKI)*, DFKI Research Report, no. RR-95-03, Kaiserslautern, Germany, March, 1995.
- [2] CHRISTOPHER M. BISHOP, *Neural networks for pattern recognition*, Oxford University Press, Oxford, United Kingdom, 1995.
- [3] CHRISTOPHER J.C. BURGESS, *A tutorial on Support Vector Machines for pattern recognition*, *Data Mining and Knowledge Discovery*, vol.2, no.2, 1998, pp. 121-167.

- [4] R. CASEY AND D. FERGUSON AND K. MOHIUDDIN AND E. WALACH, *Intelligent forms processing system*, *Machine Vision and Applications*, vol.5, 1992, pp. 143–155,
- [5] NELLO CRISTIANINI AND JOHN SHAWE-TAYLOR, *An introduction to Support Vector Machines and other kernel-based learning methods*, *Cambridge University Press*, Cambridge, UK, 2000.
- [6] A. DENGEL AND R. BLEISINGER AND F. FEIN AND R. HOCH AND F. HONES AND M. MALBURG, *OfficeMAID - A system for office mail analysis, interpretation, and delivery*, *Proceedings of the First International Conference on Document Analysis and Recognition*, Oct. 1994, Kaiserslauten, Germany, pp. 253-275.
- [7] MARCO GORI AND FRANCO SCARSELLI, *Are multilayer perceptrons adequate for pattern recognition and verification*, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol.20, no. 11, Nov, 1998, pp. 1121-1132.
- [8] J. J. HULL AND P. HART, *The infinite memory multifunction machine (IM3)*, *Proceedings of the Third IAPR Workshop on Document Analysis Systems*, Nov, 1998, Nagano, Japan, pp. 49-58.
- [9] J. L. HULL AND S. L. TAYLOR, EDS., *Document Analysis Systems (II)*, *World Scientific Publications, Co.*, 1998.
- [10] DAPHNE KOLLER AND MEHRAN SAHAMI, *Hierarchically classifying documents using very few words*, *Proceedings on the Fourteenth Conference on Machine Learning*, 1997,
- [11] M. KRISHNAMOORTHY AND G. NAGY AND S. SETH AND M. VISWANATHAN, *Syntactic segmentation and labeling of digitized pages and technical journals*, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol.15, no.7, July, 1993, pp. 737-747.
- [12] S. W. LAM, A. L. SPITZ AND A. DENGEL, *An adaptive approach to document classification and understanding*, *Proceedings IAPR Workshop on Document Analysis Systems*, World Scientific 1995, Kaiserlauterns, Germany, 1994, pp. 114-134.
- [13] STEVE LAWRENCE AND C. LEE GILES AND KURT BOLLACKER, *Digital libraries and autonomous citation indexing*, *IEEE Computer*, vol.32, no.6, 1999, pp. 67-71.
- [14] D.C. LIU AND J. NOCEDAL, *On the limited memory method for large scale optimization*, *Mathematical Programming B*, vol.45, no.3, 1989, pp. 503–528.
- [15] DANIEL P. LOPRESTI AND JIANYING HU AND RAMANUJAN KASHI, *Document Analysis Systems V*, *Lecture Notes in Computer Science*, vol.2423, August2002, Springer.
- [16] YOELLE MAAREK, *Automatically Organizing Bookmarks per Contents*, *Proceedings WWW5*, 1996.

- [17] LARRY M. MANEVITZ AND MALIK YOUSEF, *One-class SVMs for document classification*, *Journal of Machine Learning Research*, vol.2, December 2001, pp. 139-154.
- [18] L. O'GORMAN AND R. KASTURI, *Document image analysis*, *IEEE Computer Society Press*, 1995.
- [19] J. PLATT, *Fast training of Support Vector Machines using sequential minimal optimization*, *Advances in Kernel Methods - Support Vector Learning*, B. Scholkopf, C. Burges, and A. Smola, Eds., MIT Press, 1998.
- [20] MEHRAN SAHAMI AND SALIM YUSUFALI AND MICHELLE Q. W. BALDONADO, *SONIA: A Service for Organizing Networked Information Autonomously*, *Proceedings of the Third ACM Conference on Digital Libraries*, 1998.
- [21] MEHRAN SAHAMI, *Using machine learning to improve information access*, *Computer Science Department, Stanford University*, thesis, Stanford, CA, 1998.
- [22] , G. SALTON, C. BUCKLEY, *Term Weighting Approaches in Automatic Text Retrieval*, *Information Processing and Management*, vol.24, no.5, 1988, pp. 513-523.
- [23] D. SAVIC, *Automatic classification of office documents: Review of available methods and techniques*, *Records Management Quarterly*, October 1995, pp. 3-18.
- [24] SHMUELI, ODED, STAELIN, CARL, GREIG, DARRYL, AND ELAD, MICHAEL, *Classifying Semi-Structured Documents Using Image Signatures*, *Hewlett-Packard Laboratories, HPL-1999-65*, Dec, 2002, Palo Alto, CA.
- [25] A. L. SPITZ AND A. DENGEL, EDS., *Document Analysis Systems*, World Scientific Publishing, Co., 1995.
- [26] SURGUR N. SRIHARI,STEPHEN W. LAM, VENU GOVINDARAJU,ROHINI K. SRIHARI. AND JONATHAN J. HULL, *DOCUMENT IMAGE UNDERSTANDING: RESEARCH DIRECTIONS*, *Center of Excellence for Document Analysis and Recognition, CEDAR-TR-92-1*, May 1992.
- [27] CARL STAELIN AND DARRYL GREIG, *hplinet: neural networks by design*, *Hewlett-Packard Laboratories, HPL-2002-3??*, Dec,2002, Palo Alto, CA.
- [28] CARL STAELIN, *Parameter selection for support vector machines*, *Hewlett-Packard Laboratories, HPL-2002-3??*, Dec, 2002, Palo Alto, CA.
- [29] S. L. TAYLOR AND M. LIPSHUTZ, *Document understanding system for multiple document representations*, *Document Analysis Systems II*, World Scientific, 1998, pp. 283-300.
- [30] VLADIMIR VAPNIK, *The nature of statistical learning theory*, *Statistics for engineering and information science, 2nd Ed.*, Springer-Verlag, New York, New York, 2000.

- [31] H. WALISCHEWSKI, *Automatic knowledge acquisition for spatial document interpretation, Proceedings of the Fourth International Conference on Document Analysis and Recognition*, August 1997, Ulm, Germany, pp. 243-247.
- [32] T. WATANABE AND X. HUANG, *Automatic acquisition of layout knowledge for understanding business cards, Proceedings of the Fourth International Conference on Document Analysis and Recognition*, Ulm, Germany, August 1997, pp. 216-220.
- [33] S. WEIBEL AND M. OSKINS AND D. VIZINE-GOETZ, *Automated title page cataloging: A feasibility study, Information Processing and Management*, vol. 25, no. 2, 1989, pp. 187-203.
- [34] SHOLOM WEISS AND CASMIR KULIKOWSKI, *Computer systems that learn - Classification and prediction methods from statistics, neural nets, machine learning, and expert systems, Morgan Kaufman*, 1991.
- [35] C. WENZEL AND S. BAUMANN AND T. JAGER, *Advances in document classification by voting of competitive approaches, Document Analysis Systems II*, World Scientific, 1998, pp. 385-405.
- [36] C. WENZEL, *Supporting information extraction from printed document by lexico-semantic pattern matching, Proceedings of the Fourth International Conference on Document Analysis and Recognition*, Aug, 1997, pp. 732-735.
- [37] JANUSZ WNEK, *Machine learning of generalized document templates for data extraction, Document Analysis Systems V, Lecture Notes in Computer Science*, Springer, vol. 2423, August 2002, pp. 457-468.
- [38] XIN YAO AND YONG LIU, *A new evolutionary system for evolving artificial neural networks, IEEE Transactions on Neural Networks*, vol. 8, no. 3, May 1997, pp. 694-713.
- [39] XIN YAO AND YONG LIU, *Making use of population information in evolutionary artificial neural networks, IEEE Transactions on Systems, Man and Cybernetics*, vol. 28, no. 3, June 1998, pp. 417-425.