

Handling conjunctions in named entities

Pawel Mazur^{*,**} and Robert Dale^{**}

^{*}Institute of Applied Informatics, Wrocław University of Technology /

^{**}Centre for Language Technology, Macquarie University

Introduction

Named entity recognition consists of identifying ‘mentions’ — strings in a text that correspond to named entities — and then classifying each such mention as corresponding to a specific type of named entity, with typical categories being Company, Person and Location. The full range of named entity categories to be identified is usually application dependent.

Introduced for the first time as a separately evaluated task at the Sixth Message Understanding Conference in 1995 (see, for example, [Grishman and Sundheim, 1995; 1996]), named entity recognition has attracted a considerable amount of research effort. Initially handled with hand crafted rules (as, for example, in many of the participating systems in MUC-6 and MUC-7) and later by means of statistical approaches (see, for example, [Sang, 2002; Sang and Meulder, 2003]), the state-of-the-art provides high performance for named entity identification and classification both for specific domains and for language- and domain-independent systems.

However, our experience with existing software tells us that there are still some categories of named entities that remain problematic. In particular, very little work has explored the problem of the potential ambiguity of conjunctions appearing in named entity strings. Consider, for example, a string like the following:

- (1) Seshasayee Paper and Boards Limited

Even if the name has never been seen before, it is fairly obvious to a human that this is a single company name. But emulating this interpretation computationally requires resources such as an appropriate domain lexicon or relevant semantic knowledge, and in the absence of such resources, the string could just as easily be interpreted as two separate names, one being *Seshasayee Paper* and the other being *Boards Limited*. Determining the correct interpretation is clearly important for any application that relies on named entity extraction. We are interested in how

such interpretations can be arrived at relatively cheaply, and particularly without recourse to expensive-to-construct resources, so as to allow for rapid development in new domains.

The significance of this kind of ambiguity depends, of course, on the extent to which the phenomenon of conjunctions in named entities is widespread. Our current work focuses on a corpus of 13000 company announcements released through the Australian Stock Exchange: these are documents provided by companies in order to meet both continuous and periodic disclosure requirements, in which we want to track mentions of companies and individuals across time.

To get some idea of the impact of conjunction ambiguity, we selected 45 documents from this corpus at random. These documents yielded a total of 545 candidate named entity strings, of which 31 contained conjunctions. This informal sampling suggests that conjunctions appear, on average, in around 5.7% of candidate named entity strings; however, in some documents in our sample, the frequency is as high as 23%.

The documents in our corpus have some features that are not necessarily typical for other corpora. In particular, texts in this domain frequently have some of the characteristics of legal documents, where many sometimes apparently arbitrary elements are given initial capitals. Therefore, we might expect some specific domains, such as those dealing with accountancy and law, to have a higher density of names involving conjunctions. However, for comparison, the proportion of candidate named entity strings containing conjunctions in the MUC-7 evaluation data is 4.5%. These frequencies are sufficient to suggest that the seeking of an appropriate means of handling conjunctions is a worthwhile and important pursuit.

1. Problem Description

Our interest is in **candidate named entity strings** that contain conjunctions. For each such candidate string, we want to know whether it consists of a single named entity mention, or whether it consists of a sequence of two or more mentions of distinct named entities.¹

An analysis of the candidate named entity strings appearing in our corpus reveals four distinct uses of the conjunction, as exemplified in the following examples:

- (2) Oil and Gas Ltd
- (3) Agfa and Fuji
- (4) John and Mary Smith
- (5) Mssrs Cochrane and Coventry

In example (2), we have a single named entity that happens to contain an internal conjunction; in example (3), we have a conjunction of two distinct named entities; and in examples (4) and (5), we have conjunctions that, from a linguistic perspective, contain a form of ellipsis, so that one conjunct is incomplete on its own, but can be completed using information provided in the other conjunct.

Following from the above, we distinguish four categories of candidate named entity strings containing conjunctions.

Name Internal Conjunction (NI)

This category covers those cases where the candidate named entity string contains one named entity mention, where the conjunction is part of the name. Here are some examples from our corpus:

- Publishing and Broadcasting Limited
- J B Were & Son
- Hancock and Gore
- Acceptance and Transfer Form
- Fixing and Planning Phase

Name External Conjunction (NE)

This category covers those cases where the conjunction serves to separate two distinct named entity mentions. Some examples from our corpus:

- Italy and Central Europe
- Hardware & Operating Systems
- Mr Danny Fisher and Mr Don Wilson
- American Express and Visa International

Right-Copy Separator (RC)

This category of conjunction separates two named entity mentions, where the first is incomplete in itself but can be completed by copying information from the right-hand conjunct. This is perhaps most common in conjunctions of proper names, as in *John and Mary Smith*, but appears in other contexts as well. Some examples from our corpus:

- State and Federal Government
- Eastern and Western Australia
- General & Miscellaneous Equipment

Left-Copy Separator (LC)

This is similar to the previous category, but instead of copying information from the right-hand conjunct, in order to complete the constituent named entities we need to copy information from the left conjunct. Examples in our corpus:

- Gas Supply and Demand,
- Financial Statements and Reports,
- Hospital Equipment & Systems,
- J H Blair Company Secretary & Corporate Counsel.

Conceptually, we might view the last two categories as subtypes of a more general category we could call **Copying Separator**. However, it makes sense to keep the two categories separate, since the process of reconstructing the unelided conjuncts is different in each case.

2. Previous Work

The fact that conjunctions in named entity mentions are problematic has been noted before. Of course, the processing of conjunctions has long been a focus of interest in linguistics; in particular, Categorical Grammar (see, for example, Steedman, 1985) provides a sophisticated treatment of the syntax of conjunctions. Linguistic analyses tend to focus on conjunctions involving common nouns and adjectives rather than proper names, but as is clear from the examples presented above, it is not so easy to draw a clear line between these two categories. Ultimately, a full-blown and complete account of the use of conjunctions in candidate named entity strings is likely to make significant use of sophisticated syntactic analyses; however, for present purposes we are more concerned with relatively shallow approaches that can be quickly applied to large bodies of text. In this section, we review existing work in this area.

2.1 Rau

Rau (1991) focuses on the problem of identification of unknown names in financial documents. She addresses the problem of conjunctions in candidate named entity strings using the following heuristics:

- **Syntactic agreement:** being the subject of a singular verb is a strong indicator that the candidate named entity string corresponds to one entity; similarly, parallel sentence structures as in *IBM, GE and HP each ...* and *IBM, GE and HP all...* indicate that we have a list of names, rather than a single complex name.

- **Punctuation:** if there are many commas in a candidate named entity string, it is assumed that the string contains a list of named entity mentions.

Rau's work is limited to company names only, and the solution is based on analysis of a large corpus of financial news. For this data, her approach delivers an accuracy of over 95%. It remains to be seen how well the method would extend to other domains and entity types.

2.2 Coates-Stephens

Coates-Stephens (1992) presents FUNES, a system for identification of proper names in general. The definition of proper name he adopts is very broad, and relies on the notion of **key words** as indicators of the presence of named entities. Coates-Stephens' description of the system makes it clear that the handling of conjunctions is one of the most significant limitations of the system. FUNES is able to handle the following cases correctly:

- 'Copy' conjunctions when a plural key word, indicating ellipsis, is present: *the Korean and Vietnam wars; Lakes Huron and Michigan; Presidents Mitterand and Bush.*
- Cases where the candidate named entity string containing the conjunction appears in an appositive: *the authors of the report, Tim Jenkinson and Colin Meyer.*

FUNES is effective in many cases, but has its limitations; for example, as Coates-Stephens points out, this approach is likely to assign different categories to *The Competition and Service Bill* and *The Competition and Service Bills*, with the first being seen (correctly) as single named entity, but the second analyzed as a conjunction of two mentions: in some contexts this will be incorrect.

2.3 McDonald

McDonald (1996) distinguishes name-internal and name-external features of candidate named entities. Name-internal features involve the use of gazetteers, allowing identification of names of cities and countries, company-name-related keywords (somewhat akin to Coates-Stephens' work) like *Church* or *Bank*, company designators (*Inc.*, *GmbH*, etc), person designators (*Jr.*, *Sr.*, *Mr.*, *Dr.*), and so on. After delimiting candidate named entity mentions using capitalisation of tokens and punctuation, these features are used as name-internal evidence to the subsequent classifier. The classification of the mentions into categories like PERSON or COMPANY is based on both name-internal and name-external features. McDonald's

system considers three forms of conjunction: comma, *&* and *and*; the interpretation of conjunctions is implemented by means of a set of rules that analyse the lexical form of the conjunction and the type of conjuncts in order to determine the way in which the conjunction is being used. For example, one rule indicates that candidate strings of the form *X & Company* involve, in our terms, name-internal conjunctions, with the candidate named entity string simultaneously classified as being of type Organization. The system is also capable of handling more complex expressions involving conjunctions, so that forms like *Goldman, Sachs & Co.* are recognised as one named entity.

The approach relies heavily on hand written rules that make use of name-internal and name-external features of candidate named entity strings, and is limited to the names of people and companies in the news domain. McDonald claims accuracy of nearly 100% on a selected sublanguage for “Who’s News” articles from the *Wall Street Journal*, noting that applying this approach to more diverse set of texts would require a new implementation.

2.4 Mikheev et al.

Mikheev et al. (1998) suggest the strategy of examining the preceding document context to identify candidate conjuncts that should be considered as separate named entities. The authors indicate that this approach is used in their entry to the MUC-7 competition, but no data is reported on the accuracy or impact of this kind of heuristic; in our experience, there are many cases where there are no antecedent mentions that can be used in this way.

Furthermore, in the MUC-7 data, strings like *John and Mary Smith* were considered as one named entity, whereas, for many information extraction applications, it is important to recognize that this string represents two distinct entities. This has the consequence that the MUC-7 data cannot be used for evaluation of our approach without further annotation.

2.5 Downey et al

Most recently, work that is highly related to the focus of the present article is presented in Downey et al (2007). Their interest is in the recognition of complex names in web text, where complex names are those that contain lowercase words (typically conjunctions, determiners, and prepositions) within the name; this is common, for example, in book and movie titles. Downey et al report that, in their data, complex names account for 16% of all names.

Their approach makes significant use of previously-seen name strings, relying on the web as a source of previous mentions; so, for example, a candidate named

entity string like *Procter and Gamble* will be considered a single named entity on the basis that, on the basis of corpus data, the name constituent *Procter* is highly predictive of a subsequent *and Gamble*.

3. Our Approach

To recap, we characterise the task of conjunction disambiguation as being one of determining which of the four uses of conjunction is being used in any given candidate named entity string, where we have named the four uses as **Name Internal Conjunction**, **Name External Conjunction**, **Right-Copy Separator** and **Left-Copy Separator**.

Our approach to determining the proper conjunction category in a candidate named entity string is to use a machine-learned classifier. We are particularly interested in seeing how far we can address the task using only limited knowledge sources: in the work described here, we restrict ourselves to very limited gazetteers that contain the most frequent proper nouns that appear in our corpus, and to the use of so-called ‘name-internal’ properties (i.e., characteristics of the candidate string itself, rather than of its surrounding context).

Using only limited gazetteers maximises portability; considering only name internal properties will make it easier to see the impact of subsequently adding contextual information.

In the remainder of this paper, we first describe in detail in Section 4 the experiments we have carried out. Section 4.1 describes the data used; Section 4.2 describes the name-internal text features used as attributes for classification; Section 4.3 describes the data encoding used to encode the features into a feature vector; and Section 4.4 discusses the machine learning algorithms we used.

Section 5 goes on to provide a discussion of the evaluation scheme we adopted, and an overview of the results achieved in the experiments. Section 6 presents details of what went wrong by analysing misclassified examples from our data set. Finally, in Section 7, we present a discussion of possible directions in which the approach described here could be further developed. A particular approach to extending a standard named entity recognizer with the conjunction disambiguation functionality implemented as a separate module is presented in Mazur and Dale (2006).

4. Experimental setup

4.1 Corpus and Data Preparation

The focus of our project is a data set from the Australian Stock Exchange (ASX). This data set consists of a large number of company announcements: for a variety of regulatory reasons, listed companies provide around 100000 documents to the ASX each year, and the ASX subsequently makes these available to users via the web. For more information about the documents in this corpus, and a discussion of our general approach to processing them, see Dale et al (2004).

The corpus used for the research described here consisted of a 13460-document sub-corpus drawn from a much larger corpus of company announcements from the ASX. The documents range in length from 8 to 1000 lines of text.

Evaluation data was prepared as follows. For our purposes, we define a candidate named entity string to be any sequence of words with initial capitals and one embedded conjunction. We also allowed these strings to contain the lowercased preposition *of* and the determiners *a*, *an*, and *the*. Candidate named entity strings from sentences written completely in uppercase or with every word being init-capped (i.e., strings in ‘title case’) were ignored. Using a Perl script, we extracted 10925 candidate named entity string instances from our corpus, corresponding to 6437 unique forms. From the set of unique forms, we randomly selected 600 examples for our test data set. In a small number of cases, problems arising from typographic features such as ASCII formatted tables caused us to manually correct some individual strings. An example of the need for such correction is demonstrated by the candidate extracted string *Name of Entity Hancock & Gore Limited*, where it turns out that *Name of Entity* is a label in a list, and *Hancock & Gore Limited*, being a company name, is the value of that label; however, in our data, the text extraction process has caused the separating formatting to be lost, resulting in the two strings being concatenated. In this case we remove *Name of Entity* from the string extracted by our Perl script, on the assumption that a smarter text extraction technique would be able to interpret the layout more accurately.

The resulting set of strings was then annotated using a set of small gazetteers listing common person names, company names without conjunctions, locations and other named entity elements that are frequent in our corpus and related to our tagset, which is described in the next section.²

Table 1. Example distributions in categories

Name Internal	Name External	Right-Copy	Left-Copy	Sum
185	350	39	26	600
30.8%	58.3%	6.5%	4.3%	100%

The categories of the conjunctions in the candidate named entity strings were assigned by a human annotator. Table 1 presents the distribution of evaluation instances across the four conjunction categories introduced above.

4.2 The Tag Set

We developed a 16-tag tag set, presented in Table 2, to annotate the tokens in our corpus of candidate named entity strings. Most of the tags, such as LOC, ORG, GIVENNAME, ALPHANUM, DIR, and PERSDESIG, are the same as those used by many other named entity recognizers; some, however, are specific to our needs. The SON tag is used to annotate tokens whose surface form is either *Son* or *Sons*: these occur relatively often in company names (as, for example, in *A Davies & Sons Pty Ltd*), and are a strong indicator of the Name Internal Conjunction category. The OF and DET tags are used to mark the preposition *of* and the determiners *the*, *a* and *an*, irrespective of casing. Finally, INITCAPPED is used to annotate any tokens that do not belong to the other categories, or which are ambiguous between those categories.

Table 2. The tagset used for text annotation.

No	Tag	Meaning
1	Loc	The name of a location
2	Org	The name of an organization
3	GivenName	A person's given name
4	FamilyName	A person's family name
5	Initial	An initial in the range A–Z
6	CompPos	A position within a company
7	Abbrev	Abbreviation
8	PersDesig	A person designator
9	CompDesig	A company designator
10	Son	<i>Son(s)</i>
11	Dir	A compass direction
12	AlphaNum	An alphanumeric expression
13	Month	The name of a month
14	Of	Preposition <i>of</i>
15	Det	Determiners <i>the</i> , <i>a</i> , <i>an</i>
16	InitCapped	Unrecognized initcapped token

We also recognize multi-word elements where there is no ambiguity (for example, in the case of unambiguous person, location and company names). For example, although the company name *Australia and New Zealand Banking Group Limited* is not in our gazetteer, *New Zealand* as a country name is, and so this string is recognized

Table 3. The popularity of tags in annotated data

Tag	Occurrences	Percentage
InitCapped	925	42.24
Loc	245	11.19
Org	175	7.99
FamilyName	164	7.49
CompDesig	138	6.30
Initial	108	4.93
CompPos	99	4.52
GivenName	89	4.06
Of	76	3.47
Abbrev	73	3.33
PersDesig	39	1.78
Det	31	1.42
Dir	12	0.55
Son	7	0.32
Month	6	0.27
AlphaNum	3	0.14

as a sequence of tokens whose types are marked as LOC AND LOC ORG COMPDESIG; here the second LOC tag corresponds to the pair of tokens *New Zealand*.

We refer to the sequence of tags assigned to a particular string as a **pattern**. A pattern also indicates the conjunction type present in the string, as determined through the human annotation; so, for the example above, the complete pattern is $\langle \text{LOC AND LOC ORG COMPDESIG, INTERNAL} \rangle$.

Table 3 presents the number of tags of each type used to annotate our data set; in total there were 2190 tags assigned over the 600 candidate named entity strings, for an average of 3.65 tags per instance.

Notably, a significant number of the tokens are tagged as simply being of type INITCAPPED; this is a consequence of our deliberate use of small gazetteers, and is likely to be the case in any domain where new names are constantly being introduced.

4.3 Encoding

For the purposes of machine learning, we encode each pattern in the following way. We create an attribute for each of the 16 tag types for each of the left and right sides of a conjunction, for a total of 32 attributes. The attributes are of integer type with values $\{0,1\}$, thus signalling either the presence or absence of a token of that type anywhere within either conjunct.

We also introduce an additional binary attribute, CONJFORM, for encoding the lexical form of a conjunction in the string: 0 denotes *∅*; 1 denotes *and*.

With each data instance there is associated a categorical CONJTYPE attribute with the values {INTERNAL, EXTERNAL, RIGHT-COPY, LEFT-COPY}; this is used to encode the actual category of the conjunction in the string.

4.4 Baseline and Algorithms

It is quite common to determine a baseline using the 0-R algorithm, which simply predicts the majority class (see Witten and Frank, 2005). On our data set, with this approach we get a baseline accuracy of 58.33%. However, we have found that with the 1-R algorithm, described by Holte (1993), we obtain a better-performing model based simply on the lexical form of the conjunction:

- IF ConjForm='&' THEN PredCat ← Internal
- IF ConjForm='and' THEN PredCat ← External.

This very simple rule provides a baseline of 69.83%.

The experiments were conducted using the WEKA toolkit (described by Witten and Frank, 2005). This provides implementations of several machine learning algorithms, along with the data structures and code needed to perform data input and output, data filtering, and the evaluation and presentation of results.

After some initial exploration using a variety of algorithms for supervised machine learning available in WEKA, we chose the following: the Multilayer Perceptron, two lazy algorithms (IBk and K*), and three tree algorithms: Random Tree, Logistic Model Trees and J4.8. We also include here the results for Naive Bayes and SMO given the popularity of these methods in the field. We provide here only short descriptions of these algorithms, and point to appropriate references for further information:

- LMT (Logistic Model Trees) are decision trees that contain logistic regression functions at the leaves (Landwehr et al, 2005).
- J4.8 is an implementation of the decision tree algorithm C4.5 revision 8, the last public version of the C4.5 decision tree learner (Quinlan, 1993).
- Random Tree is an algorithm for constructing a decision tree that considers K random features at each node (we experimented with different values for K); it performs no pruning.
- IBk is a K -nearest neighbours classifier ($K=1$); see Aha et al (1991).
- K* is an instance-based classifier: i.e, the class of a test instance is based upon the class of those training instances similar to it, as determined by some similarity function. It differs from other instance-based learners in that it uses an entropy-based distance function; see (Cleary and Trigg, 1995).

- Multilayer Perceptron is a multilayer neural network using back propagation for training (Rojas, 1996).
- Naive Bayes is a specialized form of Bayesian network with two assumptions: predictive attributes are conditionally independent given the class; and no hidden or latent attributes influence the prediction process (John and Langley, 1995).
- SMO (Sequential Minimal Optimization) is an algorithm for the fast training of Support Vector Machines (Platt, 1999). We used the quadratic exponent of the polynomial kernel.

5. Results

5.1 Evaluation Scheme

For evaluation, we used the k -fold method with $k=10$, so that our data set of 600 examples was divided into ten folds by random selection of instances from the original data set. Then, for each of the folds, the classification models were built on the remaining 540 examples and tested on the held-out fold. The sum of correctly classified examples for all folds is the final result. There are of course some side effects of this evaluation approach, which we mention in Section 6.2.5; however, it still makes more sense to use this approach for our small data set of 600 examples, than artificially dividing this set into even smaller training and test data sets.

5.2 Classification Results

Table 4. Results for k -fold evaluation

Algorithm	Correctly classified (out of 600)
IBk	84.00% (504)
Random Tree	83.83% (503)
K*	83.50% (501)
SMO	82.33% (494)
Multilayered Perceptron	82.17% (493)
LMT	81.17% (487)
J4.8	79.50% (477)
Naive Bayes	70.67% (424)
Baseline	69.83% (419)

Table 4 presents the results achieved in the experiments. All algorithms scored above the baseline, though Naive Bayes, with the worst result, was very close to the baseline.

Table 5. Detailed accuracy by category of a conjunction for results of the IBk classifier

Category	Precision	Recall	F-measure
Name Internal	0.814	0.876	0.844
Name External	0.872	0.897	0.885
Right-Copy	0.615	0.410	0.492
Left-Copy	0.800	0.462	0.585
weighted mean	0.834	0.840	0.833

Table 6. Confusion matrix for IBk

Name Internal	Name External	Right-Copy	Left-Copy	→ classified as ↓
162	28	6	3	Name Internal
18	314	17	11	Name External
4	6	16	0	Right-Copy
1	2	0	12	Left-Copy

The best classifier turned out to be IBk, the K -nearest neighbours algorithm. The precision, recall and F-measure for this case are presented in Table 5. Table 6 provides a confusion matrix with the desired and actual classification of examples. The best results are for Name Internal and Name External conjunctions. The low results for Right- and Left-Copy Separator conjunction types are mainly because of low recall for these categories: 0.410 and 0.462, respectively. This is most likely caused by the fact that there are very few examples of these categories: 6.5% and 4.3%, respectively (see Table 1).

We used the χ^2 test for equality of distributions and a significance level of 90% to check whether the difference between the result of IBk and other algorithms is statistically significant; on this basis, we find that only the difference between the IBk algorithm and the Random Tree algorithm is no greater than chance.

It is interesting to note that the relatively simple Random Tree algorithm scored so highly. We tried different values for its parameter K , the number of randomly chosen attributes to be considered at each node. The result presented in the table is for $K=22$; for the default $K=1$, the algorithm correctly classified 490 examples.

6. Analysis

6.1 Conjunction Category Indicators

A statistical analysis of the data reveals some strong conjunction category indicators. For the Name External category these are:

- a MONTH tag in the left conjunct (as in *September and December*);

- a COMPDESIG or ABBREV tag in the left conjunct (as in *Alliance Technology Pty Ltd and Suco International* or *NLD and BRL Hardy*); but there are exceptions: *JP Morgan Investment Management Australia Ltd and Associates*; *Association of Mining & Exploration Companies* and *ASX Settlement and Transfer Corporation*, which are all Name Internal;
- a MONTH or PERSDESIG tag in the right hand conjunct (as in *February and March* or *Mr R L Hanwright & Mrs M J Hanwright*);
- a GIVENNAME, DIR or ABBREV tag in the right hand conjunct, although there are exceptions: *BHP Executive Director and Chief Operating Officer Ron McNeilly* and *SMDS and ATM WANS* (both are of the Right-Copy Separator type).

The presence of a SON tag in the right conjunct is a strong indicator of a Name Internal conjunction.

6.2 Error Analysis

Our experiment demonstrates that with supervised machine learning over a simple set of features, we can reach a classification error rate of 16–18%. In this section, we provide some discussion of the classification errors made by the best-performing learner, the IBk algorithm.

6.2.1 *InitCapped*

Of the 96 misclassified examples, 38 (39.58%) have a pattern consisting entirely of INITCAPPED tags. In such cases, classification ends up being determined on the basis of the CONJFORM attribute: if the value is *&*, then the conjunction is classified as being Name Internal, and if its value is *and*, the conjunction is classified as being Name External. Consequently, the following examples are misclassified: *Victorian Casino and Gaming Authority*, *Coal Handling and Preparation Plan*, *Gas Supply and Demand Study*, and *Explanatory Memorandum & Proxy Form*.

At the same time, there were 96 INITCAPPED-only patterns that were classified correctly; this means that out of all 134 INITCAPPED-only patterns, 71.64% were classified correctly, which is quite consistent with the previously discussed baseline.

There were also another 16 misclassified instances consisting mainly of INITCAPPED tags along with some other tags; examples of these are:

- (6) Australian Labor Party and Independent Members
(LOC INITCAPPED ORG AND INITCAPPED INITCAPPED)
- (7) Association of Mining & Exploration Companies

⟨COMPDESIG OF INITCAPPED & INITCAPPED INITCAPPED⟩

- (8) Securities and Exchange Commission
 ⟨INITCAPPED AND INITCAPPED ORG⟩

6.2.2 Long Patterns

Two misclassified instances were represented by relatively long patterns: for example, *Fellow of the Australian Institute of Geoscientists and The Australasian Institute of Mining*, represented by the 12-tag pattern ⟨COMPPOS OF DET LOC ORG OF INITCAPPED AND DET LOC ORG OF INITCAPPED⟩.

6.2.3 Other Interesting Cases

There were two cases of misclassified strings with patterns containing as substrings other, rather common, patterns. The additional tag in the extended pattern turned out to be insufficient information for a classifier to classify the extended pattern properly. As a result, the examples with extended patterns were classified as being of the same type as the embedded pattern. One example is the string *WD & HO Wills Holdings Limited*, being the name of a company (so the conjunction is of Name-Internal type) and having the pattern ⟨INITIAL INITIAL & INITIAL INITIAL FAMILYNAME COMPDESIG⟩; this was incorrectly classified as containing a Right-Copy Separator conjunction. This is because the conjunction is Right-Copy in the common embedded pattern ⟨INITIAL INITIAL & INITIAL INITIAL FAMILYNAME⟩.

A similar case is the string *Vancouver and Toronto Stock Exchanges*, which has the pattern ⟨LOC AND LOC ORG⟩; since the pattern ⟨LOC AND LOC⟩ has category Name External, the model classifies this example in the same way, effectively ignoring the ORG tag, whose presence might be taken as an indicator of the Right-Copy conjunction type. This is a case where only a slightly more sophisticated linguistic analysis might have helped: since *Exchanges* is plural, this might serve to indicate a Right-Copy conjunction. However, there are also possible counter examples to such a heuristic: for example, *Acceptance and Transfer Forms* is a plural entity of type Name Internal.

The string *Wayne Jones and Topsfield Pty Ltd*, which in reality involves a Name External conjunction, was classified as Name Internal. We would note here that, in the absence of additional contextual information, conjunctions of person names and company names are often ambiguous even for humans.

Another related highly ambiguous type of example corresponds to the pattern ⟨FAMILYNAME AND FAMILYNAME⟩, which can either be a conjunction of two person names or just one company name. One would usually expect this to be of the Name External type, but in our domain this is in almost all cases of the Name Internal type. Our mechanism misclassified one instance of this pattern.

There are also some cases which we expected to be handled easily, because there exist some obvious rules that can be used to classify these examples, but which turned out to be problematic for the classifier. Consider the following examples:

- (9) D J Carmichael Pty Limited and Kirke Securities Ltd
- (10) Kookaburra Resources Ltd and KOB Energy Inc

Both of these were classified as Name Internal, although they contain company designators in both conjuncts and the form of conjunction is *and*, which suggests the Name External category. Similarly, the string *Department of Transport and Department of Main Roads* (with the pattern `<ORG OF INITCAPPED AND ORG OF INITCAPPED INITCAPPED>`) was classified as Name Internal instead of Name External.

6.2.4 *Of tag in Right- and Left-Copy Separator Categories*

The confusion matrix presented in Table 6 shows that there were 17 examples classified as being of the Name External category when they should have been Right-Copy, and 11 examples classified as Name External when they should have been Left-Copy. For some of these examples the reason is one of those already discussed above; however there is also a significant number of examples (three and two, respectively) which have a particular structure based on the preposition *of*. Examples (11)–(13) present cases where a Right-Copy separator was classified as Name External, and examples (14) and (15) show cases where a Left-Copy separator was classified as Name External.

- (11) President and Chief Executive Officer of the Company
`<POS AND POS OF DETER DESIG>`
- (12) Retirement & Appointment of Director
`<INITCAPPED & INITCAPPED OF POS>`
- (13) Accounts and Report of the Directors
`<INITCAPPED AND INITCAPPED OF DETER POS>`
- (14) Report of the Directors and Auditors
`<INITCAPPED OF DETER POS AND POS>`
- (15) Details of the Offer and TABs
`<INITCAPPED OF DETER INITCAPPED & ABB>`

Our results suggest that *Of*-based structures are difficult to analyse correctly, as has also been noted by Downey et al (2007).

6.2.5 Other Observations

We also note here the impact of the k -fold evaluation approach. Since a new model is built for each fold, it turns out that the IBk classifier assigned category Name Internal to instances of the pattern $\langle \text{INITCAPPED AND INITCAPPED ORG} \rangle$ in one case, but assigned Right-Copy in another case. Consequently, both *Federal and State Government* (Right-Copy), being in one fold, and *Securities and Exchange Commission* (Name Internal), being in another fold, were misclassified. This phenomenon also makes the error analysis more difficult.

Finally, there is a group of around 24 examples for which it is difficult to provide a clear explanation for misclassification along the lines of the cases above; in these cases, the major issue is the classifier's ability to generalize the rules, which is not necessarily due to a deficiency in the algorithm, but perhaps due to the simple tagset we use or the relatively small size of our data set.

7. Conclusions

We have presented the problem of conjunction disambiguation in named entities and defined four categories of conjunction in candidate named entity strings. We defined the problem as one of classification and showed that it can be handled well using supervised machine learning algorithms and a limited set of name-internal features.

Given the similarity in results for most of the different machine-learned classifiers we used, we conclude that any further improvement in results lies in a richer feature selection rather than in the choice of classifier. This conclusion is also supported by the fact that some examples are difficult for a human to classify without wider context or domain knowledge.

A number of issues arise in the work reported here as candidates for future work.

We have restricted ourselves to candidate named entity strings which contain a single conjunction; however, as we have already noted, there are of course cases where multiple conjunctions appear. One category consists of examples like *Audited Balance Sheet and Profit and Loss Account*, where again the kinds of syntactic ambiguity involved suggest that a more syntactically-driven approach would be worth consideration. Another category consists of candidate named entity strings that contain commas as well as lexicalised conjunctions.

A rudimentary analysis of frequently occurring n -grams in our corpus makes it clear that some strings containing conjunctions appear frequently. For example, in our corpus there are 296 occurrences of the string *Quarter Activities and Cash-flow Report*,³ making it the most frequent 5-gram. Moreover, there are another

34 occurrences of this string with the conjunction *&* in place of *and*, and another six strings with the variant spelling *Cash Flow*. In any real application context, it would make sense to filter out these common cases via table lookup before applying a machine learning process to classify the remaining conjunctions. This kind of preprocessing could identify frequent strings containing either Name Internal or Name External conjunctions. Another form of preprocessing could involve the analysis of abbreviations: for example, in the string *ASX Settlement and Transfer Corporation (ASTC)*, the abbreviation *ASTC* could be used to decide that the conjunction in the preceding string has the category Name Internal.

More generally, there are three directions in which we might move in order to further improve performance.

First, we can always use larger gazetteers to reduce the number of tokens that can only be tagged as *INITCAPPED*. This, of course, has a cost consequence; in current work, we are exploring how performance on this task improves as larger numbers of frequent name elements from the corpus are incorporated into the gazetteers.

Another consequence of extending gazetteers is the problem of the same token being in two or more gazetteers, for example *LOCATION* and *FAMILYNAME*. A naive approach would be to assign these tokens the catch-all *INITCAPPED* tag, but since this is what we want to avoid, we could also assign all the ambiguous tags and indicate this fact in the feature vector. This would of course require a redesign of the feature vector.

Second, we can make more sophisticated use of the name internal properties of the candidate string. This includes, as noted above with regard to the *Exchanges* example, taking account of the syntactic number of the constituent tokens. Armed with a part of speech tagger, we could also attempt heuristic chunking of the candidate strings which might assist in determining conjunction type; and a resource like WordNet might be used to identify terms with shared superordinates, as in the *Paper and Board* example mentioned in the Introduction to this article.

Third, we can extend the learning process to take account of contextual features. In our data, there are cases where the local context cannot be easily determined, such as when the candidate named entity strings appear as the content of table cells; but in many cases local syntactic information such as the number of an associated verb can serve to distinguish the type of conjunction being used. However, as demonstrated here, it is already possible to achieve a high level of accuracy without recourse to name external features.

There are also aspects of our evaluation process which are open to improvement. First, we aim to increase the size of the data set used to see if the present results are robust; and second, the techniques here need to be evaluated in the context of a complete named entity recognition process, so that we can determine the

overall difference in performance when operating with and without a conjunction disambiguation capability.

Notes

1. In the present paper, we restrict ourselves to those situations where there are at most two named entity mentions in the candidate string; this is by far the largest category in the data we have considered.
2. This is part of our strategy for fast deployment in a new domain, where a seed lexicon is constructed from the most frequent words that contain initial capitals.
3. This appears frequently as a substring of longer expressions like *First Quarter Activities and Cashflow Report*, *Second Quarter Activities and Cashflow Report*, and so on.

References

- Aha, D. W., D. Kibler, and M. K. Albert. 1991. Instance-based learning algorithms. *Machine Learning*, 6(1):37–66.
- Cleary, J. G. and L. E. Trigg. 1995. K*: An Instance-based Learner Using an Entropic Distance Measure. In *Proceedings of the 12th International Conference on Machine Learning*, pages 108–114. Morgan Kaufmann.
- Coates-Stephens, S. 1992. The analysis and acquisition of proper names for the understanding of free text. *Computers and the Humanities* V26, 441–456.
- Dale, R., R. Calvo, and M. Tilbrook. 2004. Key Element Summarisation: Extracting Information from Company Announcements. In *Proceedings of the 17th Australian Joint Conference on AI*, 7th–10th Dec. 2004, Australia.
- Downey, D., M. Broadhead, and O. Etzioni. 2007. Locating Complex Named Entities in Web Text. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence (IJCAI)*, 8th–12th Jan 2007, Hyderabad, India.
- Grishman, Ralf and B. Sundheim. 1996. Message Understanding Conference–6: A Brief History. In COLING 1996 Volume 1: *The 16th International Conference on Computational Linguistics*, Los Altos, Ca. Morgan Kaufmann.
- Holte, R. C. 1993. Very simple classification rules perform well on most commonly used datasets. *Machine Learning*, 11:63–91.
- John, G. H. and P. Langley. 1995. Estimating Continuous Distributions in Bayesian Classifiers. In *Proceedings of the 11th Conference on Uncertainty in Artificial Intelligence*, pages 338–345, San Mateo. Morgan Kaufmann.
- Landwehr, N., M. Hall, and E. Frank. 2005. Logistic Model Trees. *Machine Learning*, 59(1/2):161–205.
- Mazur P., and R. Dale. 2006. Named entity extraction with conjunction disambiguation. In *Proc. of 5th LREC*, Italy, 24th–26th May, pages 1752–1755.
- McDonald D. D. 1996. Internal and external evidence in the identification and semantic categorization of proper names. In: B. Boguraev and J. Pustejovsky, editors, *Corpus processing for lexical acquisition*, pages 21–39.

- Mikheev A., C. Grover, and M. Moens. 1998. Description of the LTG System Used for MUC-7. In *Proc. of MUC-7 Conf.*
- Platt J. C.. 1999. Fast training of support vector machines using sequential minimal optimization. In *Advances in Kernel Methods: Support Vector Learning*, pages 185–208, Cambridge, MA, USA. MIT Press.
- Quinlan J. R.. 1993. *C4.5: programs for machine learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.
- Rau L. F.. 1991. Extracting company names from text. In *Proceedings of the Seventh Conference on Artificial Intelligence Applications*, IEEE 189–194
- Rojas R. 1996. *Neural networks: a systematic introduction*. Springer-Verlag New York, Inc., New York, NY, USA.
- Sang E.. 2002. Introduction to the CoNLL–2002 Shared Task: Language-Independent Named Entity Recognition. In D. Roth and A. van den Bosch, editors, *Proceedings of the 6th Conference on Natural Language Learning*, pages 155–158. Taipei, Taiwan.
- Sang E. and F. D. Meulder. 2003. Introduction to the CoNLL–2003 Shared Task: Language-Independent Named Entity Recognition. In W. Daelemans and M. Osborne, editors, *Proceedings of the 7th Conference on Natural Language Learning*, pages 142–147. Edmonton, Canada.
- Steedman M.. 1985. Dependency and Coordination in the Grammar of Dutch and English. *Language*, 61:523–568.
- Witten I. H. and E. Frank. 2005. *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann, San Francisco.

Summary

Although the literature contains reports of very high accuracy figures for the recognition of named entities in text, there are still some named entity phenomena that remain problematic for existing text processing systems. One of these is the ambiguity of conjunctions in candidate named entity strings, an all-too-prevalent problem in corporate and legal documents. In this paper, we distinguish four uses of the conjunction in these strings, and explore the use of a supervised machine learning approach to conjunction disambiguation trained on a very limited set of ‘name internal’ features that avoids the need for expensive lexical or semantic resources. We achieve 84% correctly classified examples using k -fold evaluation on a data set of 600 instances. We argue that further improvements are likely to require the use of wider domain knowledge and name external features.

Keywords: named entity recognition, conjunctions, machine learning

Authors' addresses:

Pawel Mazur
Institute of Applied Informatics
Wrocław University of Technology
Wyb. Wyspiańskiego 27
50–370 Wrocław, Poland
Centre for Language Technology
Macquarie University
NSW 2109, Sydney
Australia

Robert Dale
Centre for Language Technology
Macquarie University
NSW 2109, Sydney
Australia

Robert.Dale@mq.edu.au

Pawel.Mazur@pwr.wroc.pl

Copyright of *Linguisticae Investigationes* is the property of John Benjamins Publishing Co. and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.