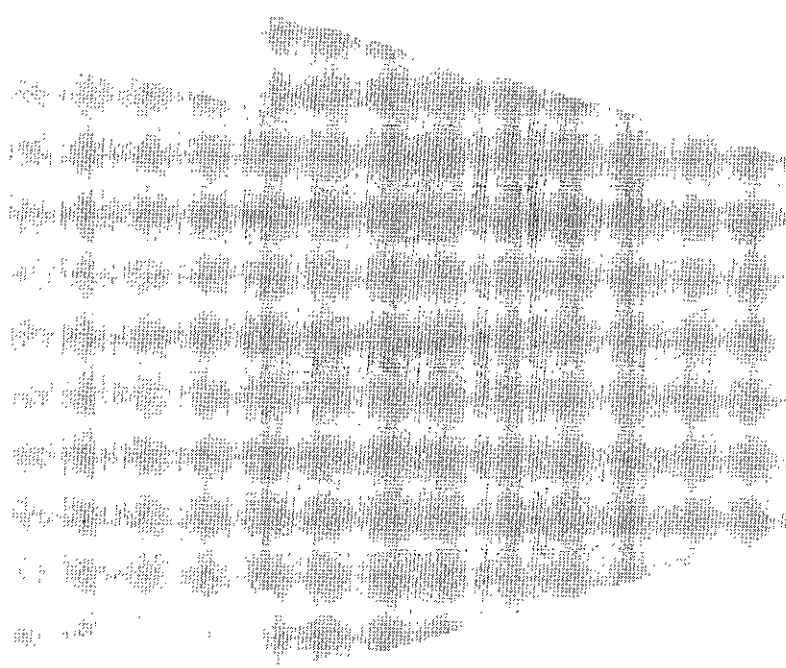Oxford University

# Automated Marking of Exam Papers

## Using Semantic Parsing

Jamie Frost

08

# Contents

# 1. Introduction

## Motivation

'Automated Essay Scoring' has been a large area of research since the 1960s. In such a process, a variety of 'features' are extracted from essays, such as word and sentence length and the structure of sentences, before the data is collaborated to provide a final classification [SHERMIS 03].

'Automated Exam Scoring' is a more objective mode of classification, in which answers are analysed for the presence of concrete facts or statements instead of using any continuous measure. This has been the subject of much research in the Computational Linguistics department at Oxford University, using an online study as the source of data, in which students completed a GCSE Biology paper [PULMAN 05]. The techniques employed are varied, but falls into two main categories. One simulating a human style marker defines the marking scheme via patterns inputted by an administrator, allowing for as many variants of an answer as possible. The clear disadvantage of this method is the hours of work required to painstakingly write these patterns, but this method yields high accuracy. Average accuracy in excess of 95% was obtained.

The latter method adopts a machine learning approach using a set of pre-marked answers for the training process. [PULMAN 06] experimented with a system in which the answers are treated as a 'bag of words' with no semantic structure incorporated. A technique known as 'k nearest neighbour (KNN)' was used, which initially removes stop words from the answer sets (words of common usage such as 'and') and assigns each keyword a TDF-IDF measure ("Term Frequency" and "Document Frequency", concepts in the Information Retrieval domain combined to indicate the 'significance' of a term). Each answer in the training set is then converted to an incidence vector form given its terms (weighting each element with the corresponding term's TDF-IDF measure), as well as the answer to classify. Using a distance measure, the closest $k$ training answers are obtained, and the modal score of these used for classification.

This naive method is subject to a number of problems, as highlighted by Professor Stephen Pulman in this BBC News article:

> *"Professor Srihari asked human examiners to grade 300 answer booklets. Half of the graded scripts were then fed into the computer to "teach" it the grading process. The software identified key words and phrases that were repeatedly associated with high grades. If few of these features are present in an exam script, it generally receives a low grade...*
>
> *Professor Stephen Pulman at the University of Oxford has identified another potential pitfall in Professor Srihari's approach. "You can't just look for keywords, because the student might have the right keywords in the wrong configuration, or they might use keywords equivalents," Professor Pulman says."*

Thus if the answer requirement is a statement such as 'the cat chased the mouse', then an answer of 'the mouse chased the cat' would be accepted despite the clear semantic inequality, due to the identical set of words.
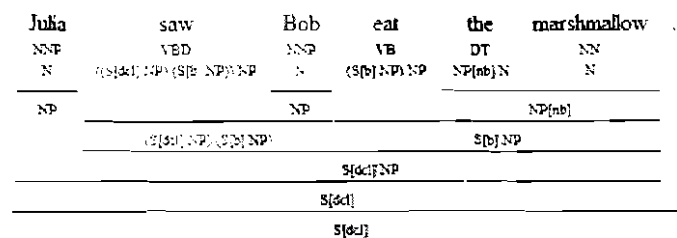
This project aims to extend this method by incorporating the semantic structure of sentences, so that for the above example 'the mouse chased the cat' would be marked as incorrect, whereas 'the mouse was chased by the cat' would be marked as correct.
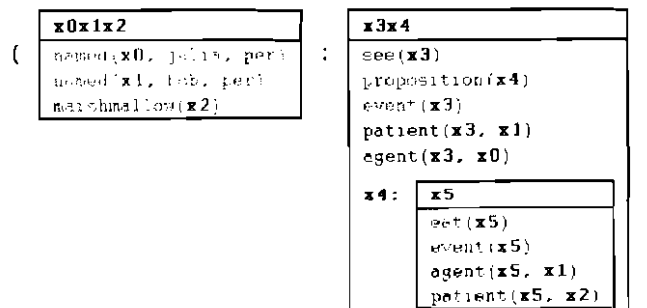
## CAndC Parser & Boxer

The CAndC (Clark and Curran) parser uses statistical methods and 'supertagging' to convert English sentences into a tree representing their structure, as detailed in [CLARK 07]. It was developed by Stephen Clark (a lecturer at Oxford University and member of the Oxford Computational Linguistics Group) and James Curran. The output of the parser is a CCG (Combinatory Categorial Grammar) file, representing this sentence structure.

Alone, this representation is insufficient for machine learning use, given that the semantic interpretation of the sentences is our concern. We therefore use a tool called *Boxer*, which uses Prolog to convert the CCG into a form called DRS (Discourse Representation Structure). This is compatible with first-order logic, and thus can be used to make reasoned logical deductions (with its application extending to other systems such as Question Answering).

Consider the sentence *"Julia saw Bob eat the marshmallow."* Its CCG form is:



The graphical form of its corresponding DRS is:



Verbs and nouns are represented as objects (here x0 to x5) with predicates defining their properties and relations connecting them.

## The Data

All data is from students' answers to GCSE Biology papers (both foundation and higher tiers) from 2001 and 2003. The anonymised data is made available under a confidentiality agreement courtesy of the OCR Exam Board. Maximum marks for answers vary between 1 and 4 marks.
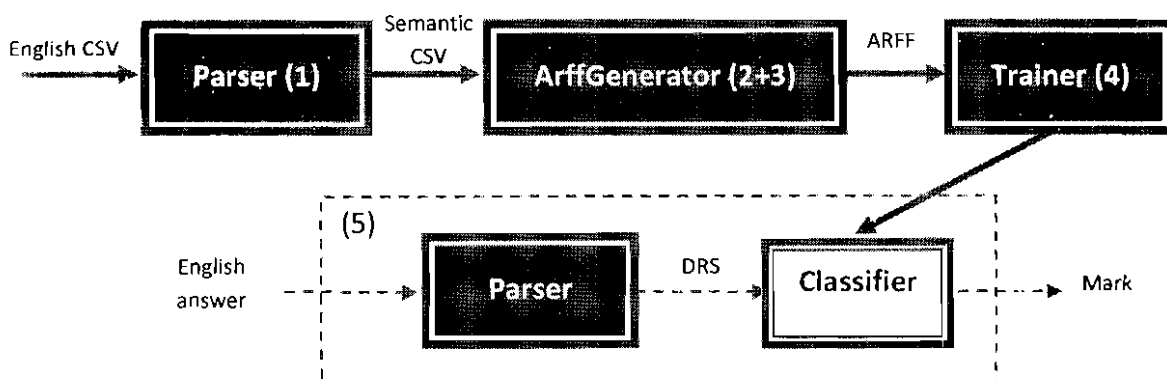
## Weka

Weka is a suite of machine learning tools for Java, developed by the University of Waikato in New Zealand. It provides a variety of classifiers (including Naive Bayesian, Neural Networks and Decision Trees), clustering algorithms and other utilities. Downloads and documentation can be found at http://www.cs.waikato.ac.nz/~ml/weka/.

## 2. System Overview

In the scope of this project, we investigate 3 different methods of exam mark classification: a naive keyword approach in which the answers are treated as a 'bag of words' as described above (but using a more standard classifier than the 'k nearest neighbour' technique), an approach which incorporates the sentence semantics, and a potential optimisation of the latter which uses a new type of classifier. All 3 methods use a 'supervised learning' approach, that is, we train a classifier with pre-marked data, and use this to mark subsequent answers.

Such a classification problem requires a number of stages:

1. *Semantic Conversion:* This takes the answers in English (a CSV file containing the answers and marks), and converts them to a semantic representation, using the *CAndC* parser to identify the sentence structure. This is then fed this into *Boxer*, producing a semantic representation (i.e. a DRS). Clearly this step is skipped in the keyword approach. The conversion process is described in Section 3.

2. *Featurisation:* A 'feature' is a singular property encompassing some fact about an answer. A 'binary feature' takes the value 0, denoting the feature is not present, or 1 if it is. The aim of this stage is to provide some function which maps the answer (either in its raw form for the keyword approach, or the DRS for the semantic approach) to a set of such 'features' which are satisfied. For the keyword approach, the mapping simply produces a set of words the answer contains. For the semantic approach, the DRS is a non-linear structure and thus the mapping is more ambiguous. The latter is dealt with in Section 4.

3. *Data Preparation:* Such features from all answers must be collated to provide a suitable input for a classifier. This data is embodied as a file known as an 'ARFF'. See Section 5.

4. *Training:* A suitable machine learning classifier is chosen, and is 'trained' with this data.

5. *Classification:* Subsequent new answers are parsed (if necessary) and featurised as before. This is fed into the trained classifier to provide an estimated classification (i.e. a mark) for the answer.

# 3. Parsing

Our initial input data is a CSV file containing a list of English sentences, each with a classification. The aim of the parsing process is to produce a new CSV file, instead with DRS expressions representing each of the sentences. There are 2 methods of parsing.

One method uses the online parser found on the CAndC website.[1] However, a more practical solution would incorporate a server which can administer requests from clients and send the parsed expressions back. The CAndC tools include a *SOAP* server and client, which communicate with each other (potentially externally). The advantage of these is that the high overhead cost associated with starting the parser need only occur once in the server's lifetime, relative to once per parse request. Both of these operate on the same host, and we disregard their remote capabilities of these in order to produce a Java 'wrapper' which gives this functionality. The *SOAP* client should not be confused with the remote Java client, the latter of which generates the parse requests.
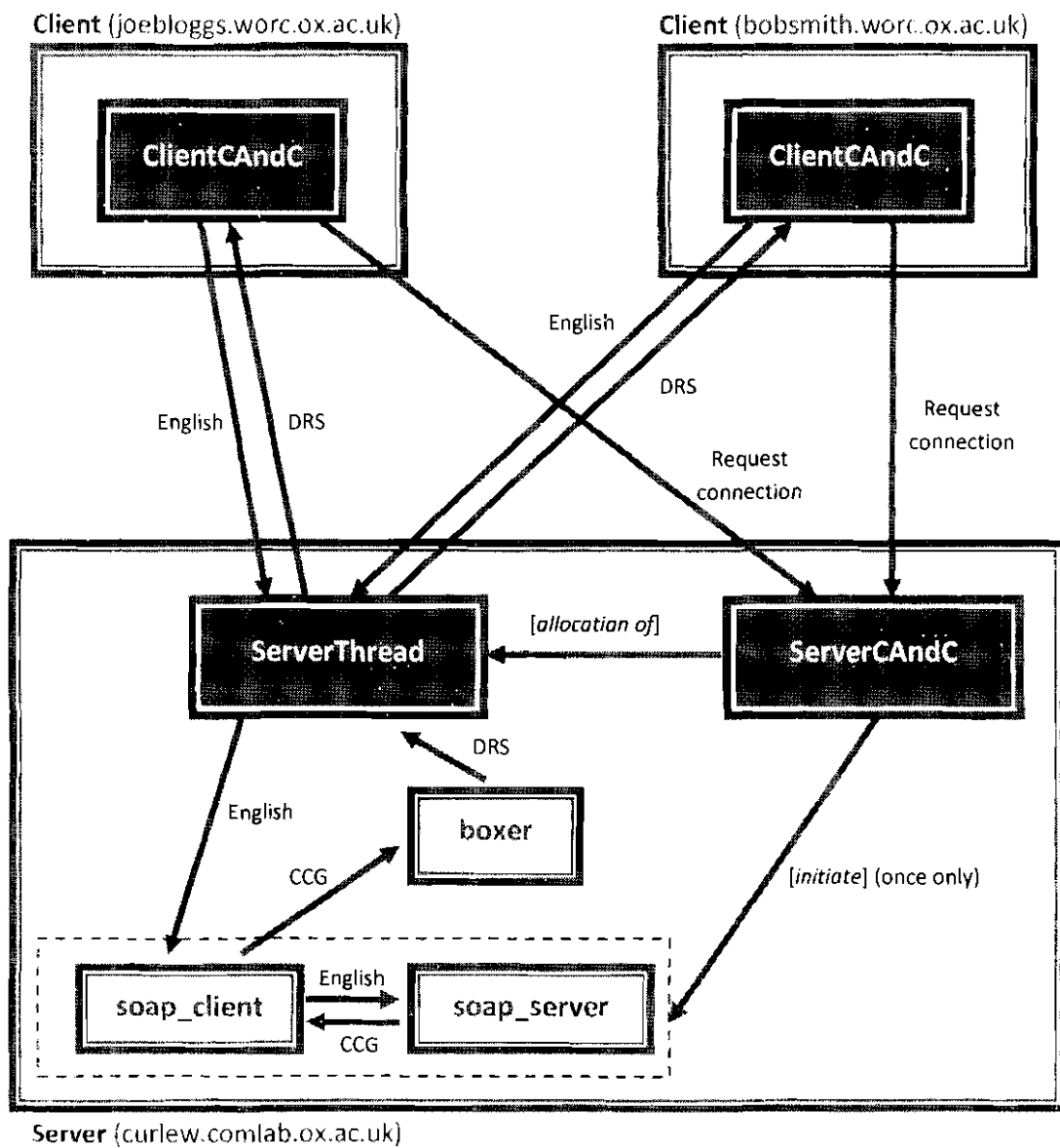
What is required is a server which forwards requests to the *SOAP* client (which in turn communicates with the *SOAP* server). It may need to manage the requests of multiple remote clients, and thus we create a separate manager for each connection, running concurrently. It performs the following:

1. The *SOAP* server is started.
2. We indefinitely wait on a specified port for external clients to make a connection. When an incoming request is received, a socket and a manager representing the connection are created.
3. In this manager, we construct an input and output stream from the socket, and continually read from it until the connection closes. We read an English sentence from the stream, write it to a file (given that the parser reads from a file) and instruct the *SOAP* server to parse the sentence. Finally we use Boxer to convert this to a DRS expression, which is delivered to the output stream.

A client end of the communication is also required. It sends the sentence to the server and delivers the parsed and boxed DRS expression returned to the required resource. As before, we establish a socket using the host name and port of the external machine running the server, and construct an input and output stream to read and write from.

The structure of the entire parsing system can be summarised as follows:

---

[1] This can be found at http://svn.ask.it.usyd.edu.au/trac/candc/wiki/Demo

**Client** (joebloggs.worc.ox.ac.uk)

**ClientCAndC**

**Client** (bobsmith.worc.ox.ac.uk)

**ClientCAndC**

English

DRS

English     DRS

Request
connection

Request
connection

**ServerThread**     [*allocation of*]     **ServerCAndC**

DRS

English

boxer

CCG

[*initiate*] (once only)

English

soap_client     soap_server

CCG

**Server** (curlew.comlab.ox.ac.uk)

7

# 4. Representing Semantic Atomic Facts

As previously discussed, previous attempts at implementing a machine learning system to classify answers have relied on treating them as a 'bag of words', such that marking was based on common keywords. However, this has a major flaw. The ordering and form of words have no bearing, leading to false positives. For example "Jack collects stamps with Jill" and "Jill stamps on Jack as he collects Jill" are likely to have similar classifications, despite their very different meanings and different uses of the word 'stamps'. In other words, there is no capacity for semantic representation. In the latter sentence, although Jill is performing the 'stamping', it could equally be Jack.

The aim therefore is to derive a way of semantically representing sentences that is complete, sound and most importantly, linear. That is, we wish to deconstruct a sentence into a collection of atomic predicates (here described as 'semantic atomic facts') such that these properties are satisfied. This is an example of 'featurisation', that is, we identify a number of possible properties of the input data (in this case, the presence of each of the generated atoms) for the use of machine learning application at a later stage. To establish these concepts of soundness, completeness and linearity, take for example a sentence represented by $\alpha$ and the 'semantic atom facts' we wish to generate from it by a set $\{\beta_1, \dots, \beta_n\}$:

- For soundness, $\forall i \, . \, \alpha \vDash \beta_i$. In other words, all such facts are implied by the original sentence.
- For completeness, $(\wedge|\vee|\neg) \, \beta_i \vDash \alpha$. That is, some disjunctive/conjunctive/negative combination of the atoms implies the original sentence. Allowing the disjunctive connective to yield $\alpha$ may at first seem odd, but the necessity becomes clear when considering sentences with the word 'or'; if for example "The kite is red or green" was split into "Red kite" and "Green kite".
- Linearity simply implies the lack of some recursive or complex data structure, instead keeping the atomic facts 'flat'. Here, atomicity implies the absence of the $\wedge$, $\vee$ and $\neg$ connectives.

We define a suitable function $\varphi$ which produces every possible inference that can be made from a DRS expression $\alpha$, that is, the $\beta_i$ as defined above. Take for example the sentence "When jolly Carol laughs, there is an earthquake". We can make a number of inferences from the sentence, representing looser forms:

1. When Carol laughs, some event occurs.
2. When someone laughs, some event occurs.
3. When Carol performs some action, some event occurs.
4. When jolly Carol laughs, some event occurs.
5. Carol is involved in the sentence.
6. Carol is jolly.
7. An earthquake occurred.
8. When someone laughs, an earthquake occurs.

And so on. To formally define $\varphi$, we consider all the mappings for all possible DRS sub-expressions to sets of 'atoms'. The DRS consists of a number of objects, representing nouns or verbs, named via a specially typed predicate. The atom produced is simply the name of the predicate (note that the term after the colon indicates the type):

$$\varphi(\text{p:noun}(x)) \quad = \quad \{ \text{ p:String } \}$$
$$\varphi(\text{p:verb}(x)) \quad = \quad \{ \text{ p:String } \}$$

For example, $\varphi(\text{cat}(x)) = \{ \text{ "cat" } \}$. In addition, other unary predicates representing adjuncts can modify such an entity; this is an adjective in the case of a noun, and an adverb in the case of a verb. Its corresponding atoms consist of the name of the modifier and the modifier combined with the name of the entity. The *nm(x)* function simply retrieves the assigned 'name' of the entity *x*.

$$\varphi(\text{p:adj}(x)) \quad = \quad \{ \text{ p:String }, \text{ p:String}[nm(x)]\}$$

For convenience, we define $\varphi(x)$ to be all atoms associated with the entity *x*. We similarly take $\varphi$ of the whole DRS expression $\alpha$ as the union of sets produced by applying the function to all expressions contained within it. That is:

$$\varphi(x) = \varphi(p_1(x) ; \ldots ; p_n(x)) \quad = \quad \varphi(p_1(x)) \cup \ldots \cup \varphi(p_n(x))$$
$$\varphi(\alpha) = \varphi(\alpha_1 ; \ldots ; \alpha_n) \quad = \quad \varphi(\alpha_1(x)) \cup \ldots \cup \varphi(\alpha_n(x))$$

Binary predicates (i.e. relations) connect two entities, and typically represent a verb or a prepositional modifier. In mind of providing all inferences, the resulting atoms must incorporate every possible pair of atoms from the atom sets of the constituent entities:

$$\varphi(p(x_1, x_2)) \quad = \quad \{ \text{ p:String}[e_1,e_2] \mid e_1 \in \varphi(x_1) \wedge e_2 \in \varphi(x_2) \}$$

Larger DRS expressions can be combined together using a variety of conjunctives. These must be dealt with appropriately:

$$\varphi( \alpha \vee \beta ) \quad = \quad \varphi(\alpha) \cup \varphi(\beta)$$
$$\varphi( \alpha \wedge \beta ) \quad = \quad \varphi(\alpha) \cup \varphi(\beta)$$
$$\varphi(\alpha =/\text{is } \beta) \quad = \quad \{ \text{ "}e_1 = e_2\text{", "}e_2 = e_1\text{" } \mid e_1 \in \varphi(\alpha) \wedge e_2 \in \varphi(\beta) \}$$
$$\varphi(\alpha \rightarrow \beta) \quad = \quad \{ \text{ "}e_1 \rightarrow e_2\text{" } \mid e_1 \in \varphi(\alpha) \wedge e_2 \in \varphi(\beta) \}$$
$$\varphi( \neg\alpha ) \quad = \quad \varphi(\alpha)$$

Such a function may initially seem disastrously hideous to a logician: negation is ignored, and conjunction and disjunction are treated in a uniform way. The motivation is to have as high coverage as possible of all inferences made by the sentence, with the only requirement being that the original sentence can be reconstructed from some conjunctive/disjunctive/negative combination of a subset of $\varphi(\alpha)$. Note that the equality relation is symmetric, therefore we must supply both orderings in the string representation. We do not treat $\alpha \rightarrow \beta$ as $\neg\alpha \vee \beta$ given we are only interested if the inference was made, not if it was true vacuously when $\alpha$ is false.

To demonstrate this definition, take the sentence "The big dog jumped over the lazy man." This could be represented in DRS by:

$$\text{dog}(x_1) ; \text{big}(x_1) ; \text{man}(x_2) ; \text{lazy}(x_2) ; \text{jump}(x_3) ; \text{agent}(x_3,x_1) ; \text{over}(x_3,x_2)$$

In this example, the jumping is being executed by the dog (i.e. the dog is the agent of the jumping). Using the laws above, the following atoms are produced:

dog
big
big[dog]
man
lazy
lazy[man]
jump

agent[jump,big[dog]]
agent[jump,dog]
agent[jump,big]
over[jump,lazy[man]]
over[jump,lazy]
over[jump,man]

## Relevance to Exam Marking

The above system is based on assumptions relating to a generic marking system. We assume that some number of marks are awarded for the conjunctive/disjunctive combination of these atomic facts. For example, a mark may be awarded for mentioning that the man is lazy ("lazy[man]") and another for identifying a jump occurred over him ("over[jump,man]").

There may initially appear to be a problem with answers either being too specific or not specific enough. Say for example that "over[jump,man]" was required for a mark. If the user was to answer with "jumped over the large man", we obtain "over[jump,large[man]]" and by deconstruction, "over[jump,man]" also. Is a mark awarded for "jumped over the *large* man" if the answer required was "jumped over the man"? This is somewhat of an ambiguous question. In some cases, the extra detail may be unnecessary but not detrimental to the answer; we may consider that our example merits the mark. However, there are some cases where clearly extra detail changes the answer. Take for example the answer "The act of putting magnesium in water causes a reaction." Then "The act of putting magnesium in water causes a *fatal* reaction." is clearly wrong, given that it grossly misrepresents the magnitude of the reaction.

Fortunately, a suitable machine learning method can usually cater for this. Let an answer be $\alpha$ (e.g. red dog) and its looser implication $\beta$ (e.g. dog) such that $\alpha \to \beta$ (i.e. a red dog is a dog):

- For an actual answer $\alpha$ where $\beta$ is not specific enough (e.g. where *red dog* would be awarded a mark but not *dog*): $\alpha \land \neg\beta$ holds.
- For an actual answer $\beta$ where $\alpha$ is an acceptable extension (e.g. if *dog* was accepted, then *red dog* would also be accepted): $\alpha \lor \beta$ holds.
- For an actual answer $\beta$ where $\alpha$ is not a valid extension (such as with the chemistry example above): $\neg\alpha \land \beta$ holds.

These are all simple Boolean expressions which can be 'learned' during the training process. Classifiers' ability to learn any such Boolean expression justifies our function $\varphi$ used for linearisation.

There is however one particular flaw with the definition of $\varphi$, in the way that negation is handled. It is difficult to distinguish between "The dog was red" and "The dog was *not* red", since the negation is ignored. The underlying problem is that the binary value in the vector representation of the answer represents the presence or absence of a clause, not the Boolean interpretation of 'true' or 'false'. We could propagate the $\neg$ through subclauses (resulting in sets like { dog, red, red[dog] } and { dog, $\neg$red, $\neg$red[dog] }), or use 3-valued logic (i.e. 0 for the absence of a feature, 1 if present but negative, and 2 if present and positive), but such attempts generally lowered the accuracy of the system (by up to 10% for some data sets). One might suggest that this is due to students generally avoiding specifying a correct answer before negating it. For example, the question "By which process to plants produce oxygen and food?" is unlikely to be answered with "Not by photosynthesis."

# 5. Data Preparation and Training

Parsing and atomising from the previous stage will produce sets of atoms for each of the answers. In order to produce an input file suitable for a classifier, these must be collated.

In Weka, a training file for a classifier is known as an ARFF (Attribute-Relation File Format) file. The file has two segments: the first is an ordered list of attributes (each with a data range defined), and the second a data table representing the instances (i.e. answers) as vectors, where each element represents a value for the corresponding attribute (the last value being the classification for the instance).

Producing the set of attributes is trivial; we simply take the union of all the sets of atoms. To produce the data table, for each instance in the training set we use the value '1' for an element if the attribute appears in its set of atoms, and '0' otherwise. This is best illustrated by an example; consider these 3 sets of atoms for the 3 answers in the training set:

$$A_1 = \{ \text{cat, dog, mouse} \}, \quad A_2 = \{ \text{dog, horse, rabbit} \}, \quad A_3 = \{ \text{horse, monkey} \}$$

Suppose that $A_1$ received 1 mark, $A_2$ received 0 and $A_3$ received 2. Then the appropriate ARFF file would look like such:

```
@relation animal_example
@attribute "cat" { 0, 1 }
@attribute "dog" { 0, 1 }
@attribute "mouse" { 0, 1 }
@attribute "horse" { 0, 1 }
@attribute "rabbit" { 0, 1 }
@attribute "monkey" { 0, 1 }
@attribute "classification" { 0, 1, 2 }

@data
1, 1, 1, 0, 0, 0,    1
0, 1, 0, 1, 1, 0,    0
0, 0, 0, 1, 0, 1,    2
```

Notice that we must determine the maximum mark in the training set in order to define the range for the "classification" attribute. Upon collating the attributes, it is beneficial to maintain a count of all attributes. Any attribute occurring only once is removed, as it is deemed as having insufficient impact on the classifier. This has the benefit of massively reducing the number of attributes, thus making training and classification considerably faster.

Training a classifier in Weka is near trivial using an ARFF file. We construct an untrained classifier, and an object representing the data in the ARFF file. After providing a reference to this data in the classifier, the training process is finally invoked[2].

---

[2] The trained classifier is required for future use, so we use Java's *FileOutputStream* class to write the classifier to a file. File representations of classifiers are (arbitrarily) given 'model' as a file extension.

There are a number of established classifiers that one may wish to use. Here is a brief description of the options available:

- *Neural Network:* This consists of a number of connected units known as 'neurones', analogous to the functioning of the brain. Each neurone takes the weighted sum of its inputs, and applies a function (usually either a threshold or sigmoid function) to produce its output. The network uses the training set and error minimisation via a process of 'back propagation' to iteratively adjust these weights.
- *Decision Tree:* This repeatedly partitions the training set on attributes, ordered by the 'information gain' (i.e. the extent of diversity in the proposed partitioning) each attribute provides. This results in a tree, with each node splitting into two given the presence of a certain feature, and the 'leaves' designating the resulting classification.
- *Naive Bayesian Classifier:* This uses simple probabilistic theory based on Bayes' Theorem. It calculates the probability of a classification (given the features) by using the probability of each feature given the classification and the prior probability of each classification, both of which are easily calculated. The classification which maximises this probability is selected. It assumes that each of the features are independent.
- *Maximum entropy modelling:* A framework for integrating information from many sources for classification. The idea is that it finds a model that satisfies the given constraints but does not go 'beyond the data', that is, it avoids a model not justified by the empirical evidence (i.e. the training set). Its advantage over the Naive Bayesian classifier is that it does not enforce independence assumptions with regard to the features.

Each of these classifiers has its own advantages and disadvantages. For the purposes of the evaluation process, the decision tree and naive Bayesian classifiers will be used, given their support by the Weka libraries.

Given our 4 stages of semantic parsing, featurisation, data preparation and training, we now have a complete system suitable for testing. However, there are some possible refinements we can make given the specific nature of the classification problem at hand. Two issues raised and explored are that of spelling correction, and a new form of classification that deals with the cumulative nature of exam marking.

# 6. Spelling Correction

Exam answers, typically those of a younger contiguate, must be expected to be filled with both grammatical and spelling errors. While the former are incredibly difficult to correct by autonomous means, it is possible to develop a system which corrects spelling errors with a good degree of accuracy.

One must initially question the importance of this correction in marking. Errors in inconsequential words which form the structure of the sentence will be of little concern to the examiner. However correcting such words would seem vital in order for the parser and boxer to establish a more accurate interpretation of the sentence, and to reduce the number of attributes in training. On the other hand, misspellings of keywords directly encapsulated by the mark scheme may heavily influence the examiner's marking. While the 'examiner's generosity' is possible to approximately model, it raises the issue that 'desirable spelling correction' is a fairly ambiguous and nondeterministic problem. This would suggest that a semi-autonomous approach should be adopted rather than a fully autonomous one. The system should produce a modifiable mapping defining how spelling corrections should be made, before mapping this (whenever the administrator chooses to do so) on the data.

An important observation to make is that a universal set of correct spellings is difficult to obtain. Many words are subject specific and thus will not be found in a standard dictionary. Therefore we must use a combination of the words in the data and words in a dictionary. A simplified version of the proposed correction algorithm can be defined as thus:

correction(*word*) = if correct(*word*) then *word* else
arg max ( { score(x,*word*) | x ∈ suggestedCorrections(*word*) } ∪ { score`(x,*word*) | x ∈ corpus } )

where score(x,y) and score`(x,y) are two different scoring functions for dictionary and corpus comparison. We distinguish between the 2 functions since we wish to take into account the frequency of words in the corpus.

## Dictionary Comparison

Upon obtaining a list of suggestions for an incorrectly spelled word, we obtain a similarity measure between the word *w* and the suggested correction *s*:

score(w,s) = α.FirstLetterBonus(w,s) − β.LevenshteinDistance(w,s)
− LetterDifference(w,s) − SizeDifference(w,s)

It was found $\alpha = 3$ and $\beta = 2$ produced the best results. A number of functions, of the type *(String,String)* ↦ *N*, provide various string metrics:

| Name | Explanation | Motivation | Implementation |
|---|---|---|---|
| Levenshtein Distance | The number of insertion, deletion, or substitution character operations required to change from *w* to *s*. For example, the distance between "kitten" and "sitting" is 3. | This provides an initial well-established comparison measure. | A dynamic programming algorithm using a (n + 1) × (m + 1) matrix, where n and m are the lengths of the two strings. |

| | | | |
|---|---|---|---|
| First Letter Bonus | Returns 1 if the first letter of both words is the same, and 0 otherwise. | Regardless of how badly a word is spelled, typically the first letter is correct. Thus from "recieve", it is more likely that "receive" was intended than "believe", despite both having the same Levenshtein distance. | |
| Letter Difference | Taking each word as a set of letters, it counts the number of letters which are in one word but not the other. | We favour corrections in which the letters used is more consistent. This for example favours a correction from "doen" to "done" rather than the less probable "den". | $\#[\text{letterset}(s) \cap \text{letterset}(w)^c] + \#[\text{letterset}(s)^c \cap \text{letterset}(w)]$ |
| Size Difference | The difference in size between $w$ and $s$. | This gives more preference to a 'letter pair swap' over an addition/deletion, and eliminates candidate corrections too different in size. | $\|\#w - \#s\|$ |

We keep the correction with the highest score, but delay adding this to the mapping given that this may be overwritten by a better match in the corpus.

## Corpus Comparison

When comparing a misspelled word against the corpus (i.e. the exam answers), the assumption above that the words are correctly spelled can no longer be made; therefore the frequency of each term must be taken into account. We can augment the previous metric to provide a new one:

$$\text{score}'(w,s) = \mu \cdot \log \text{frequency}(s) + \text{score}(w,s) - \lambda$$

In essence, we provide extra weighting if the word is commonly occurring, that is, we can be more confident of the word's 'correctness'. The constant $\lambda$ ensures that it is not necessarily the case that score'(w,s) > score(w,s) (which would essentially discard dictionary comparison). We make a number of checks to prevent incorrect mappings:

1. Most obviously, we ignore the pair *(w,s)* if w is equal to s (preventing the identity mapping).
2. We ignore *(w,s)* if the Levenshtein Distance between them is greater than 5. This prevents a mapping to a very dissimilar word if it is the best that exists (i.e. it may be the case there is no mapping available at all).
3. We ignore *(w,s)* if the frequency of w is greater than the frequency of s in the corpus. This prevents mapping cycles, which would cause an infinite loop when attempting to map the corrections on the data.

## Other Considerations

There are a few forms of errors that cannot be easily detected by the above methods. We can correct these accordingly:

1. *Plurality errors:* It is common for the plural form of a word to be misspelled. For example, "babies" is commonly written incorrectly as "babys", and "halves" as "halfs". In the above techniques, it is likely that an alternative word will be suggested, therefore we must compensate accordingly. We hypothetically replace common incorrect suffices with their correct equivalent, such as "ys" for "ies" and "fs" for "ves". If the new word is in the dictionary, we add this to the mapping and ignore any other correction methods.

2. *Words augmented with apostrophes:* Nouns can have "'s" added to the end. The new string is not in the dictionary, despite this extension being perfectly valid. We check if the substring before the apostrophe is in the dictionary, and if so, regard the word as correct.

3. *Word splitting:* It is common for certain pairs of words to be incorrectly merged, for example "aswell", "alot" and "thankyou". To split these, we hypothesise a space insertion at each point in the string; if this forms two words both found in the dictionary, we add this to the mapping. This method can be problematic if not augmented, given that unlikely splits such as "aswell" into "a swell" can be made. This problem can be partially alleviated by favouring splits in which the bigram occurs most frequently in the corpus (i.e. "as well" is likely to occur more frequently than "a swell").

## Mapping the Corrections

The mappings of incorrect words to their (hopefully) correct forms are saved in a file. This allows the user to inspect the suggested mappings and either add their own or modify/correct the ones already present. The mapping is applied simply by searching for occurrences of each member of the mapping key in the data, and replacing it with its corresponding value.

## A Demonstration

Here is a sample of corrections produced when dealing with a question on 'vasoconstriction':

| | | | |
|---|---|---|---|
| alot | a lot | hypothlamus | hypothalamus |
| ammount | amount | insulater | insulator |
| aswell | as well | inturn | in turn |
| bodys | bodies | mantain | maintain |
| capilaries | capillaries | normaly | normally |
| capilary | capillary | prents | parents |
| capileries | capillaries | presonne | person |
| capillarites | capillaries | stoped | stopped |
| capillars | capillaries | suplie | supply |
| capillarys | capillaries | transfered | transferred |
| capilleries | capillaries | traveling | travelling |
| capilliaries | capillaries | vascoconstrict | vasoconstrict |
| capilliery | capillary | vascondriction | vasoconstriction |
| cappillaries | capillaries | vasconstriction | vasoconstriction |
| cappillary | capillary | vascontriction | vasoconstriction |
| cappilleries | capillaries | vasilaconstriction | vasoconstriction |
| colser | colder | vasoconstruction | vasoconstriction |
| dont | don't | vesseles | vessels |
| errect | erect | vessles | vessels |
| everytime | every time | | |

# 7. Cumulative Classification

The approaches used thus far produced promising results despite the small training sets available. However, there is a significant conceptual flaw in the classification approach used, stemming from the fact that classifiers have no abstract notion of 'magnitude' in their outputted values. Marks assigned may as well be replaced with arbitrary letters or symbols, since the classifiers cannot distinguish that a mark of 2 is greater than a mark of 1. This leads to some unfortunate consequences:

- As the maximum mark becomes large, standard classifiers rapidly deteriorate in performance. Consider for example all answers which are awarded 5 out of 10 marks. These marks could have been awarded in $^{10}C_5 = 252$ ways. Clearly, even with a large training set, it would be impossible for classifiers to fully encapsulate 5 mark answers, when one considers that each of these 252 ways are based on multiple features, further increasing the uncertainty.

- We neglect dependency assumptions. If an answer is classified as 2 and satisfies some property $A \wedge B$, then clearly if two answers both classified as 1 have the properties A and B respectively satisfied, then they are strongly correlated. By ignoring these dependencies, we surrender a large quantity of inference that can potentially greatly aid in the training process.

A possible solution we propose is to introduce a notion of 'cumulative classification', with just one fundamental property:

> *If two sets of properties X and Y are satisfied resulting in classifications x and y respectively, then the overall resulting classification for $X \cup Y$ is x + y (given that X is not a subset of Y, and vice versa).*

This is an intuitive principle, and encompasses the way in which an examiner would mark a paper. If an answer gets a mark for mentioning one point and a second for another, then one would expect this to contribute a total of 2 marks to their final mark. Standard machine learning classifiers cannot represent this concept.

The method proposed can be divided into 3 distinct stages:

1. *Segmentation:* We divide the features into sets of features relating to each mark in the mark scheme. These sets may not necessarily be disjoint, neither must they span the whole set of features. We maintain a classifier for each of the marks, with the corresponding sets denoting each one's attributes.

2. *Training:* For each answer, we determine which of the marks are most likely to be satisfied and which aren't, and train each classifier accordingly.

3. *Classification:* We combine the outputs of these classifiers to provide an overall classification for the answer.

In summary, we are delegating the problem to each of the individual marks, combining their results to solve the overall classification task. The following sections deal with the motivation and implementation for each of the three stages.

## Segmentation

The cumulative classifier is a composite of simpler classifiers managing the classification of each individual mark. Each of these classifiers requires a list of attributes, and intuitively one would expect this to incorporate features relating to this particular mark. There are two measures which we must consider to determine the accuracy of this algorithm. The first is precision; the proportion of features which appear in the 'optimum group'. The second is recall; the proportion of features from the 'optimum group' that appear.

Precision is important, when we consider that taking the entire set of features for all groups, which produces 100% recall but low precision, would result in poor classification, given that the system has no way to distinguish between marks. Recall is also important, given that neglecting important features produces a classifier that is not representative of the mark.

Some consideration should be given to the relationship between the sets of features. Should for example the sets be disjoint? This depends greatly on the mark scheme in question. Marks tend to be based on independent points to avoid ambiguity in marking, but it is very possible that two different marks will mention the same item or action. The sets need not provide full coverage of the entire feature set, when we consider that many points contained within some answer may be entirely irrelevant to the mark scheme, and thus not associated with any particular mark.

We construct 3 different methods that attempt to perform this segmentation. The first 2 of these use a 'co-occurrence' matrix, that is, a symmetric $n \times n$ matrix ($n$ is the number of features) where each element $M_{ij} = M_{ji}$ is the number of times feature i occurs with feature j. This can be interpreted as a measure of the connection between 2 features. To compute this, we take a binary matrix $A$ representing the training data, such that each row is an incidence vector representing one answer[3]. M is easily calculated:

$$M = A^T . A$$

This follows from the fact $M_{ij}$ takes the dot product of columns $i$ and $j$ (where such a vector is the incidence vector of a particular feature across all answers), equating to the number of times both features occur across all answers (i.e. if both features are present in a particular answer, their product is 1, contributing to the total count, otherwise their product is 0). We could feasibly use $M^k$ for some k > 1 to improve the connectivity of features. If k = 2 for example, then $M^2_{ij}$ conceptually indicates the extent to which feature $i$ is connected to feature $j$ via some intermediate feature; if for example "cat" occurs with "fur" in an answer and "fur" with "dog", but not "cat" with "dog", then $M^2$ allows us to make such a connection between "cat" and "dog". In practice this has an adverse effect for the type of data in question. We finally map the log function onto all elements of the

---

[3] An initial assumption was that filtering the matrix to rows that received '1' as a classification may provide an optimisation. This is because each row would have present features mostly corresponding to one mark, thus increasing the probability that two features have a high co-occurrence due to them relating to the same mark. However, this reduced the training set to such an extent that this was not a viable approach.

matrix before normalising the columns[4]. This ensures that particular high co-occurrence values do not dominate and force other values towards 0 after normalisation.

## Method 1: Standard K-means Clustering

K-means clustering is a commonly used method that partitions elements into disjoint sets (which span the entire set) using a standard algorithm. The resulting clusters have been partitioned in such a way that the mean distance from each vector to the centroid[5] of its respective cluster is minimised. The process consists of simply clustering on the columns of M into k sets, where k is the maximum mark. Weka provides clustering facilities, thus the implementation involves producing a 'wrapper' which uses the matrix M to construct an input suitable for the clusterer, before returning a list of sets, where each set contains the indices of the features in the cluster.

The algorithm is fast and efficient, and provides moderate accuracy for most questions. However the disjoint set assumption is likely to cause problems for questions where high feature overlap occurs between marks.

## Method 2: Iterative Clique Merging

An alternative approach proposed here generates a number of seeds, where each seed is an island of closely connected features, before merging them into the required number of clusters.

Instead of normalising the columns in M as before, we scale all elements down so that the maximum element is 1. We define a threshold $\alpha$ ($0 < \alpha \leq 1$) such that if $M_{ij} \geq \alpha$, the bond between features i and j is considered to be 'strong', and 'weak' otherwise. This allows us to construct an undirected graph G, where we define the nodes to be the features, and the edges to be the 'strong bonds'.

To demonstrate this, take an example such that the co-occurrence matrix M is:

|   | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| A | 1 | 1 | 1 | 1 | 0 | 0 |
| B | 1 | 1 | 1 | 0 | 0 | 0 |
| C | 1 | 1 | 1 | 1 | 0 | 0 |
| D | 1 | 0 | 1 | 1 | 0 | 0 |
| E | 0 | 0 | 0 | 0 | 1 | 1 |
| F | 0 | 0 | 0 | 0 | 1 | 1 |

with features A to F. The value of $\alpha$ is arbitrary here since all non-zero values are 1. This produces the graph:



---

[4] More specifically, the function is $f(x) = log(x+1)$. This ensures that 0 maps to 0, rather than produce an error.
[5] The centroid of a cluster is the mean of its constituent vectors.

We start by identifying all cliques contained within the graph. A clique is a set of nodes such that there is an edge between every pair of nodes within it. From the graph, we can see that there are 3 cliques: {E, F}, {A, B, C} and {A, C, D}. To identify the cliques, we devise the following algorithm:

> *Let cliques = Ø*
> *For (each row i of M)*
> > *If (cliques is empty) add new clique ( i ) to cliques*
> > *Else (*
> > > *For ( each clique q in cliques)*
> > > > *If all the bonds from i to each of the elements in q are strong,*
> > > > > *add i to q.*
> > > > *Else if all the bonds from i to each of the elements in q are weak*
> > > > *(and additionally, i was not added to any other clique in cliques)*
> > > > > *create a new clique ( i ) and add it to cliques.*
> > > > *Else (i.e. we have a mixture of weak and strong bonds) create*
> > > > > *a new clique:*
> > > > > $q' = ( i ) \cup (j \mid strong(i, j), j \in q)$
> > > > > *and add it to cliques.*

Presuming that *cliques* is not empty, we observe for each of the accumulated cliques which nodes the node $i$ has edges to. Clearly if there are edges to all the nodes, the clique is preserved so we add the node $i$ to the set. Similarly, if $i$ has no edges to any element in the set of cliques (such as row E), we create a new clique (delayed until the end of iteration to ensure this condition is met). Lastly, if there are edges to some but not all the nodes in a clique (such as node D with {A,B,C}), we create a new clique containing the node itself and those nodes which it has an edge to. One might initially think this is a polynomial time algorithm[6] given the nested for loops. However, the inner loop becomes larger as *cliques* expands, leading to an exponential algorithm.

If we follow this algorithm for the above example, on each iteration we obtain: [{A}], [{A,B}], [{A,B,C}], [{A,B,C}, {A,C,D}] (where the split occurs), [{A,B,C}, {A,C,D}, {E}], [{A,B,C}, {A,C,D}, {E,F}].

We desire only $k$ groups however, so the sets above must be merged. Merging is performed by constructing an appropriate algorithm:

> *Let groups = cliques; Discard all g $\in$ groups of size 1*
> *While( |groups| > k)*
> > *Let bestGroup1 = Ø*
> > *Let bestGroup2 = Ø*
> > *For( each group g1 in groups)*
> > > *For(each group g2 in groups)*
> > > > *Take the cross product of g1 and g2 (g1 x g2) and*
> > > > *find the average bond strength over all pairs. If this value*
> > > > *exceeds the best so far, set bestGroup1 = g1, bestGroup2 = g2.*
> > *Add bestGroup1 $\cup$ bestGroup2 to groups.*
> > *Remove bestGroup1 and bestGroup2 from groups.*

---

[6] If this were the case, then $P = NP$ given that a similar problem, the 'maximum clique' problem, appears as one of Richard Karp's original 21 problems shown NP-complete in his seminal 1972 paper "Reducibility Among Combinatorial Problems".

Suppose k = 2. Then just after one iteration of the algorithm, {A,B,C} and {A,C,D} will be merged to {A,B,C,D}, leaving this and {E,F}. Observing the original diagram, such group identification for this particular example appears sensible.

The fundamental flaw in this algorithm is that its unavoidable exponential complexity renders segmentation on hundreds of features (as would usually be expected) infeasible. Even if divide-and-conquer techniques could be used on the merging to produce an $O(n.logn)$ algorithm, there exists tens of thousands of cliques produced from the previous process. Using singleton sets with each of the features as the seeds for the merging algorithm simply does not work as the resulting sets will be disjoint; in which case the standard clustering algorithm would be a superior alternative.

One should also note that a slight alteration in α can drastically alter the output of the algorithm. Clearly tweaking such a parameter dependent on the data set to improve results is undesirable.

### Method 3: Manual specification
A non-autonomous alternative is to manually specify the groups each feature appears in. Clearly manually assigning features to groups removes the problems with the methods listed above, and any loss in accuracy is due only to human error. Given that a data set of 150 answers typically produces 800 features, the task is an arduous one and becomes impractical as the data set increases. A command line interface developed can alleviate this problem if the mark scheme is based on simple keywords (see *Implementation Notes*).

## Training
The segmentation from the previous stage produces k classifiers with assigned attributes ready for training. In order to train multiple networks, we must determine a method of taking one instance and distributing this over the k classifiers.

Suppose an answer in the dataset obtained a mark of j out of k ($0 \leq j \leq k$). Each classifier can only produce a classification of 0 or 1 (indicating whether the mark was obtained or not), so we expect to train j of these classifiers with 1 (where the instance is restricted to the features the feature set corresponding to that classifier contains), and the remaining k-j with 0.

To determine which of these j classifiers are trained with 1, we rank the k groups in order of *cosine similarity* between the cluster's vector (with 1s for attributes that occur in the cluster, and 0s for those that don't) and the instance vector, from highest to lowest. In practical terms, this gives a relative measure of the magnitude of the instance 'matching' each cluster. Once this ranking is achieved, we simply take the top j elements from the list as the clusters that 'received' the mark.

The cosine similarity is a measure of the similarity of two vectors by calculating the angle between them. For two vectors A and B, it is defined to be:

$$\cos \theta = \frac{A \cdot B}{|A||B|}$$

We can ignore the Euclidean length of the vector corresponding to the answer to classify, since this will remain constant as we vary the cluster vector. Thus the measure is simply:

$$\frac{cluster \cdot answer}{|cluster|}$$

# Classification

For a new instance, we simply restrict the answer vector to each of the clusters, classify each instance, and sum over the classifications. This is analogous to a marker determining which marks in the mark scheme were awarded, before totaling their score.

Extending to alternative types of mark scheme requires little change. Take for example a scheme of the form "Award 1 mark for any of the k points, maximum n". The number of clusters remains the same, but in classification we take *min(mark, n)* so that the maximum mark is not exceeded.

# A worked example

Suppose a mark scheme consisted of 2 marks, with the following criteria:

1) Mention of a 'dog'.
2) Mention of the act of 'walking' or 'exercise'.

Suppose also the following training answers were supplied, their classification denoted in brackets:

A. The big dog started barking. (1)
B. The dog buried his bone. (1)
C. Adam enjoys walking and exercise. (1)
D. The dog enjoyed exercising. (2)
E. The big lecturer scared the timid student. (0)

Each of the answers would be semantically analysed and deconstructed into its constituent features. For the sake of simplicity, we extract only 4: 'dog', 'big', 'walk' and 'exercise'. We could therefore generate the following incidence matrix A for the above answers:

|   | dog | big | walk | exercise |
|---|---|---|---|---|
| A | 1 | 1 | 0 | 0 |
| B | 1 | 0 | 0 | 0 |
| C | 0 | 0 | 1 | 1 |
| D | 1 | 0 | 0 | 1 |
| E | 0 | 1 | 0 | 0 |

Computing $M = A^{T} \cdot A$, taking the log of each value, and normalising the columns, we obtain:

|   | dog | big | walk | exercise |
|---|---|---|---|---|
| dog | 0.5 | 0.387 | 0 | 0.279 |
| big | 1 | 0.613 | 0 | 0 |
| walk | 0 | 0 | 0.5 | 0.279 |
| exercise | 1 | 0 | 0.5 | 0.442 |

By inspection of the non-zero values, it is clear that the two groups obtained by clustering are $G_1$ = { dog, big } and $G_2$ = { walk, exercise }. This is consistent with the original mark scheme. We create a standard classifier for each of these groups (say $C_1$ and $C_2$), and now use the rows of the matrix A to

train these. We define an ordered list $L_R$ to be the groups in descending order of their similarity with the row R.

The incidence vectors for the 2 groups are (1, 1, 0, 0) and (0, 0, 1, 1) respectively. Since the Euclidean distance of these are the same, for this example we can ignore these from the cosine measure and simply use the dot product instead.

For the first row A, we have the vector (1, 1, 0, 0). We take its dot product with each of the group incidence vectors:

$$(1, 1, 0, 0) \cdot (1, 1, 0, 0) = 2$$
$$(1, 1, 0, 0) \cdot (0, 0, 1, 1) = 0$$

Since 2 > 0, this gives us $L_A = [G_1, G_2]$. For this first row, 1 out of 2 marks were awarded, so we train $C_1$ (corresponding to $G_1$, the first element of the list) with the row (restricted to the features it contains, i.e. (1, 1) ) and the classification '1', and train $C_2$ with the row (restricted to (0, 0) ) and the classification '0'.

Using this method for each row, we train each classifier 5 times.

| | | | | $C_1$ | | $C_2$ | |
|---|---|---|---|---|---|---|---|
| Row | Row . $G_1$ | Row . $G_2$ | Marks | Training instance | Class | Training instance | Class |
| A | 2 | 0 | 1 | (1, 1) | 1 | (0, 0) | 0 |
| B | 1 | 0 | 1 | (1, 0) | 1 | (0, 0) | 0 |
| C | 0 | 2 | 1 | (0, 0) | 0 | (1, 1) | 1 |
| D | 1 | 1 | 2 | (1, 0) | 1 | (0, 1) | 1 |
| E | 1 | 0 | 0 | (0, 1) | 0 | (0, 0) | 0 |

Now that training is complete, classifying a new answer is a simple process. For example, the answer "The big hamster had some exercise on his wheel" is converted to the vector (0, 1, 0, 1) using the 4 features. We take the vector (0, 1) corresponding to the group $G_1$ and put it through $C_1$. The classifier hopefully should return '0' given that a dog isn't mentioned and the 'big' adjective is irrelevant. Next we take the vector (0, 1) corresponding to the group $G_2$ and put it through $C_2$. The classifier should return '1' given that 'exercise' is mentioned. 0 + 1 = 1, therefore the final classification for this answer is 1.

## Pre-Testing Analysis

Before this system is tested, it is beneficial to establish a hypothetical upper bound on accuracy compared with conventional classifiers. While cumulative classification has some merits, its cumulative nature leads unfortunately to cumulative error.

Consider the case where the maximum mark is 1. Since partitioning the feature set into 1 leaves it unaffected, the algorithm behaves the same as a normal classifier. Suppose that we assign a probability $p$ that the classification for an instance is correct, and for simplification, presume this is the same probability as classifying any mark in a mark scheme correctly. If a mark of $k$ out of $n$ was obtained, we obtain the following probabilistic upper bound for a correct overall classification being received:

$$P(\text{correct}) = \sum_{r=0}^{\min(k,n-k)} \binom{k}{r}\binom{n-k}{r} p^{n-2r}(1-p)^{2r}$$

This upper bound could only theoretically be reached if the segmentation is optimal, and if in training we always identify the correct clusters for which a mark was awarded. This expression seems complex, but is a simple modification of a standard binomial distribution. To understand it, consider firstly the probability of the examiner classifying all marks correctly. This is simply $p^n$. Now consider that the examiner accidently assigns 2 marks incorrectly, such that for one mark it was wrongly identified as incorrect, and another wrongly identified as correct. The total classification would remain the same, since we are not changing the total number of marks. There are $^kC_1 = k$ ways of picking a 1 to switch to a 0, and $^{(n-k)}C_1 = (n-k)$ ways of picking a 0 to switch to a 1. $(n-2)$ marks are assigned correctly with probability p (giving $p^{n-2}$) and 2 marks are assigned incorrectly with probability $(1-p)$. We generalize this to $2r$ mistakes where $r$ is the number of incorrect pairs, and sum over all the possible number of pairs of mistake. There is a maximum of $min(k,n-k)$ mistakes, since we are limited by the number of 0s we can pair with 1s, or 1s we can pair with 0s.

To assess the practical impact of this, take p to be a sensible value of 0.85. Presuming a random distribution of marks out of 4, the average upper bound for the probability of correct classification is 0.55462. Given that errors in clustering and training are likely to lower the accuracy, this upper bound is bad news. However, accuracy is not the only method of evaluation, and a measure known as 'disposition' is introduced later to justify the usefulness of this classification technique.

# 8. Graphical User Interface

Currently all methods and operations discussed previously are only accessible via the command line. When dealing with numerous files across the conversion pipeline, this can be somewhat of a cumbersome process. Therefore, a simple GUI (Graphical User Interface) was developed and used to facilitate such operations.

The GUI consists of 3 tabbed panels. The first allows conversion between various file formats across the overall pipeline. The first button takes a CSV of English answers and parses them to produce a CSV of DRS expressions. The second button uses this and produces an ARFF, and finally the last button uses the ARFF to train a classifier (prompting the user to select between the different classification methods).



A second panel allows the user to test the performance of the system for a given dataset. The user simply selects their desired training type and the folder containing the input data (which must contain a file named 'answers.csv'), and the results are printed.

Finally, a third panel allows individual answers to be classified given a trained classifier.



There are two features of the GUI worth drawing attention to:

- *Drag & Drop file selection box*: The class *FileDragPanel* is simply a component, and allows the user to specify a file or directory quickly. It has 'drag and drop' functionality, allowing the user to drag a file or directory into the box to make the selection. Alternatively (given that the file to specify may not yet exist if it is the output of an operation) the user can click the panel to open a standard browse dialog. The constructor takes a file type (enforced when browsing) and a name, which is displayed at the top of the panel.



- *Loading bar*: The evaluation process can be long, and thus it is desirable to give the user a rough indication of the time remaining. Progress bars are easy to create and update in the Swing framework. However, an issue arises of GUI components becoming 'frozen' until the entire evaluation process is complete, resulting in a jump from 0% to 100% as it finishes. The solution is to run the evaluation process 'in the background'. The abstract class *SwingWorker<Void, Void>* has a method *doInBackground()*; by implementing this, we ensure that updates to the GUI during an operation are administered. The progress is merely an approximation, with 0-20% progress for parsing, 20-40% for generating the ARFF files, 40-80% for training and 80-100% for testing.



The GUI used the Swing toolkit available to Java. The 'Look & Feel' (that is, the interface style) is set to match that of the operation system; as one may have observed, all screenshots were taken while using Windows Vista.

# 9. Evaluation

The form of evaluation used employs a technique known as "10-fold cross-validation", commonly used in Computation Linguistic systems. This splits the training set into 10 segments, trains a classifier on 9 of the segments and classifies the remaining 1 segment. The results are averaged over the 10 possible testing configurations.

Each answer in the test segment is classified, and is compared to the actual classification from the data. A measure of the system's performance is obtained, encompassing two properties:

- *Accuracy:* The percentage of answers that were classified correctly.
- *Disposition:* The average displacement of the estimated marks to the actual marks. This is potentially a more useful measure for higher mark questions, since for a mark of 4 out of 4, an estimation of 3 is superior to an estimation of 0.

The Naive Bayesian Network was chosen as the classifier (and for the cumulative classifier as the sub-classifiers). The decision tree classifier gave on average approximately 5% worse performance, and thus was ignored. Results were obtained for the 3 different approaches: the keyword approach, the semantic approach using the parser, and the cumulative classifier described in Section 7.

## Results

| Q | Max mark | Num answers | Keyword Approach | | Semantic Approach | | Cumulative Classifier | |
|---|---|---|---|---|---|---|---|---|
| | | | Accuracy | Disposition | Acc. | Dis. | Acc. | Dis. |
| 2a | 1 | 200 | 94% | 0.065 | 92.50% | 0.075 | 92.50% | 0.075 |
| 2bi | 1 | 192 | 89% | 0.115 | 84.87% | 0.151 | 84.87% | 0.151 |
| 2bii | 2 | 190 | 78% | 0.225 | 63.16% | 0.374 | 53.68% | 0.489 |
| 2biii | 1 | 200 | 98% | 0.025 | 95% | 0.05 | 95% | 0.05 |
| 2ci | 2 | 189 | 70% | 0.335 | 53.48% | 0.518 | 50.32% | 0.534 |
| 2cii | 2 | 193 | 78% | 0.220 | 79.92% | 0.201 | 72.68% | 0.278 |
| 4a | 2 | 162 | 73% | 0.306 | 74.74% | 0.296 | 69.01% | 0.322 |
| 4bi | 1 | 195 | 92% | 0.080 | 93.34% | 0.067 | 93.34% | 0.067 |
| 4bii | 2 | 192 | 55% | 0.465 | 43.29% | 0.687 | 55.76% | 0.504 |
| 4biii | 1 | 68 | 81% | 0.186 | 78.33% | 0.217 | 78.33% | 0.217 |
| 4ci | 2 | 195 | 71% | 0.305 | 65.13% | 0.379 | 46.92% | 0.612 |
| 4cii | 1 | 193 | 82% | 0.180 | 76.29% | 0.237 | 76.29% | 0.237 |
| 4d | 3 | 163 | 45% | 0.625 | 37.65% | 0.708 | 40.88% | 0.77 |
| 5aii | 2 | 170 | 65% | 0.370 | 69.41% | 0.329 | 65.88% | 0.37 |
| 6a | 3 | 186 | 60% | 0.490 | 53.83% | 0.623 | 56.37% | 0.506 |
| 6b | 1 | 201 | 85% | 0.155 | 91.00% | 0.09 | 91.00% | 0.09 |
| 6c | 4 | 139 | 48.5% | 0.644 | 50.71% | 0.707 | 47.14% | 0.636 |
| 8b | 2 | 186 | 65% | 0.366 | 69.06% | 0.332 | 66.26% | 0.382 |
| 9c | 2 | 182 | 71% | 0.300 | 67.35% | 0.332 | 52.75% | 0.478 |
| 12ci | 2 | 168 | 76% | 0.240 | 74.41% | 0.256 | 66.73% | 0.339 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **12cii** | 2 | 158 | 62% | 0.415 | 70.29% | 0.316 | 58.17% | 0.463 |
| **13bii** | 4 | 119 | 68% | 0.522 | 66.67% | 0.567 | 61.67% | 0.558 |
| **Avg** | | | 73% | 0.309 | 70.47% | 0.340 | 67.07% | 0.370 |

| | Keyword Approach | Semantic Approach | Cumulative Classifier |
|---|---|---|---|
| Mark | Average Accuracy | Average Accuracy | Average Accuracy |
| 1 | 88.71% | 87.33% | 87.33% |
| 2 | 69.45% | 66.39% | 59.83% |
| 3 | 52.50% | 45.74% | 48.63% |
| 4 | 58.25% | 58.69% | 54.41% |



The first noticeable property of the above graph is that the semantic approach was marginally worse than the keyword approach. This is initially surprising, given that one would expect the set of features produced from the keyword approach to be a subset of that for the semantic approach, and thus the performance at least as good. There is a variety of explanations for this decrease in performance:

1. Some variation of performance is expected given that the input data used was different. While the keyword approach used all input data, some badly spelled/ungrammatical answers were rejected by the parser, and thus discarded from the training data.
2. The parser had the effect of transforming some of the words, such as changing plurality or removing tense, and other forms of 'noise'.
3. Similarly, the parser is not 100% accurate, and therefore some answers will have an incorrect semantic representation.

More generally, the results show impressive performance for 1 mark questions, but a considerable decline in accuracy as the mark increased. The cumulative classifier gave relatively poorer performance to the semantic approach for 2 mark questions, but better performance for a few higher mark questions. Accuracy had some correlation with the size of the data set, albeit small.

**Performance against number of answers**

**Performance against maximum mark**

The disposition data did not yield any significant trends across the 2 classifiers, other than being approximately proportional to the accuracy. There was only one dataset in which the accuracy was lower for the cumulative classifier but had a higher disposition. This is likely to be because the modal mark for each dataset tended to be 0; thus if a 4 out of 4 mark answer is classified (incorrectly) by the cumulative technique as 3, and as 0 by a standard classifier, the low frequency of such high-mark answers renders such improvement (i.e. a classification closer to the actual mark) insignificant.

For a given dataset, the wildly variable accuracies for each segment in the 10-fold cross-validation are also indicative of inadequate training data. If a large dataset was used, one would expect these accuracies to be roughly consistent. It should be noted that the parser was not always successful in producing a parse, particularly for long answers. This is unfortunately unavoidable, given the poor quality of grammar and spelling in many of the answers.

## Cumulative Classifier

The poorer than expected results for the cumulative classifier can partially be explained by the often poor datasets. As identified above, the modal mark was often 0, and for higher mark questions, full-mark or near full-mark answers were very infrequent (in some cases appearing only once or twice). This, coupled with the low amount of training data, means that many of the individual mark sub-classifiers predominantly receive a 0 as a classification in the training data.

It is curious that questions *4d* and *6a* had very contrasting performance despite both having the same maximum mark. Inspecting the mark scheme, this is likely to be because there was less overlap of features associated with marks in *6a* (a question on plants, in which the marks were fairly

28

independent concepts) than there was with *4d* (a question on vasoconstriction, in which 2 of the 3 marks made reference to blood). This strongly suggests that the cumulative classifier performs better when the clustering is less ambiguous.

## Spelling Correction

| Question | Accuracy | Change |
|----------|----------|--------|
| 4a | 73.30% | -1.44% |
| 5aii | 68.82% | -0.59% |
| 6a | 52.78% | -1.05% |
| 6b | 88.05% | -2.95% |
| 6c | 50.90% | 0.19% |
| 8b | 68.54% | -0.52% |
| 9c | 66.76% | -0.59% |
| 12ci | 72.46% | -1.95% |
| 12cii | 68.13% | -2.16% |
| 13bii | 67.69% | 1.02% |
| Average: | | -1.00% |

The results show that spelling correction yields no improvement on accuracy. There are a number of reasons why there may have been a slight reduction in performance:

- The correction of words integral to the grammatical structure of the sentence led to less answers being rejected by the parser. This combined with the other alterations to the dataset leads to a minor change in accuracy, either positive or negative.
- The spelling corrections may be too generous – that is, an incoherent answer with key terms misspelled may be penalised by the examiner. This results in some classifications in the training set being too high.

# 10. Conclusion

The results are a mixture of encouragements and discouragements. While the spell correction and cumulative classifier enhancements failed to provide an increase in performance, accuracy for many of the data sets was high. However, there was a disappointing lack of improvement in performance over the 'bag of words' approach, even if this performance was significantly superior to that of the 'k nearest neighbour' approach in [PULMAN 06]. Unless there were fundamental inadequacies in the method of featurisation used, it is clear that semantic information from the sentences bears little or no influence on classification on the type of data used. This is perhaps of little surprise; while theoretically semantics are of paramount value, in the context of factual exam questions its use is limited. Typically one would not expect examinees to write incoherent answers containing the desired 'keywords' in an attempt to fool the system – instead the inclusion of such keywords in an answer is a strong reflection that the student has understood the topic at hand.

Despite this, the very high accuracies in some test segments suggest that accuracies could be significantly improved by datasets of a larger size, and thus one could deem the results 'inconclusive' given the inadequacy of the data.

Ultimately, the problem with a machine learning approach to automated exam marking is the existence of a no-win situation. A small training set leads to poor results, but a large training set is counterproductive for the classification problem at hand; each classifier can only be used for a specific question, and thus unless the total number of exam papers is very large, manually marking a large subset of these is not economical.

The conclusion is that the approach of manually specifying the mark scheme required (via the patterns previously explained in the introduction) is superior to attempting to learn such patterns. The cumulative classifier may aid in identifying possible groups of features for each mark, but given 100% accuracy in the segmentation algorithm is unrealistic, this leads to features occurring in incorrect groups, and thus renders such information limited in use. The cumulative classifier does however provide some 'feedback' where traditional classifiers would not, in that the classifier indicates which of marks in the classification were and were not awarded.

There are number of extensions that may be made to this research. A more accurate segmentation algorithm that doesn't rely on the co-occurrence matrix could be investigated. Relying solely on the co-occurrence matrix is possibly naive, as one cannot assume that two features occurring frequently together necessarily relate to the same mark (and as discussed, limiting the dataset to 1-mark answers is too restrictive). One might also propose that the system may effectively be used in another domain where one large dataset is used to classify instead of multiple smaller ones. Other optimisations could also be investigated, such as 'stemming', which merges classes of words containing variants of the same word (e.g. "fertilise", "fertilised" and "fertilisation") into a single representation.

# 11. Implementation notes

All coding was written in Java (1.6).

## a) Parsing

The *WebsiteParser* class used for 'online parsing' simply inserts the English sentence into the URL (acting as a parameter), obtains the resulting page (using Java's *InputStreamReader* on the URL as a stream), and extracts the DRS expression from it.[7]

*ServerThread* is a class used to manage a particular connection. As an extension of Java's *Thread* class, such managers can be run concurrently to service multiple connections simultaneously.

Instructions to the SOAP server and client can be made via the command-line. We can execute such instructions from Java by constructing a *Process* object. This can be obtained by *Runtime.getRuntime().exec(command)*. The *Process* class usefully has a method *waitFor()*, which causes the method to halt until the scheduler signals completion of the process. One might wonder why inelegant file writing and reading were used in *ServerThread* instead of piping. This is because the command-line functionality is very limited in Java, and thus can only execute basic commands.

The following command-line prompts are used in implementing the server:

- soap_server --server localhost:9000 --candc ../candc/candc-1.00/models/boxer
  is used to start the server on the local machine (using an arbitrary port), specifying the Boxer data as the parsing model to use. Since there is no way of determining whether the parser has fully loaded, we wait a fixed amount of time (20 seconds) before continuing.
- soap_client --input inputFile --output parsedFile --url http://localhost:9000
  This command instructs the SOAP client to communicate with the SOAP server. The parsed output (in CCG form) is outputted to a file. The port must be the same as that specified in initialising the SOAP server.
- boxer --input parsedFile --output finalOutputFile --box true --flat
  This takes the parsed file and uses Boxer to produce the DRS expression. We extract this from the output file and write this to the output stream of the server.

## b) Implementing $\varphi$

Implementing the $\varphi$ function involves 2 stages. The first is to convert the DRS string obtained from parsing into some object oriented form. The second takes this construction to generate the desired atoms. The *DRSreader* class in the 'parser' package is used to execute this first stage, producing a *SemanticModel* object. A *SemanticModel* represents all the variables, predicates and relations (including equality, etc.) within the DRS.

To do this conversion, we establish an abstract class *Expr*. This has 2 unimplemented methods:

---

[7] Such a technique of "web-scraping" is generally discouraged for large amounts of data, and therefore *WebsiteParser* can only be used to obtain parses for individual sentences (i.e. for classification rather than training purposes).

- *int parse(List tokens, Integer index)*: Takes a list of tokens (where each token is a word or character of the DRS string) and the point to which we have currently scanned. This parses the tokens in a way relevant to the part of the DRS tree we are currently in, and builds up the internal structure of this object. The integer returned is the new token position after parsing.
- *void process(SemanticModel model)*: Used for the second phase. It takes the structure of the object accumulated from parsing, and updates the *SemanticModel*.

We use an abstract class as opposed to an interface so that we can define a method *simpleParse*, which parses a single token (containing no information, such as parenthesis) and produces an error if the token is not as expected.

Clearly we must first convert the DRS string to tokenised form. A method *List<Object> tokenize(String str)* in *DRSreader* performs this, using Java's *StringTokenizer* class. To then initialise the parsing process, we just then construct a new *Drs* object (representing the root of the new tree being constructed), and then execute its parse method with these tokens. One might wonder why a *List* of tokens is used and not a more intuitive data structure such as an *Iterator*, where once a token is consumed, it is not seen again. The reason is that sometimes we wish to 'peek' at the current token without consuming it, in scenarios where there are multiple types of children that a node could have.

Now that a tree structure of the DRS has been obtained, we use it to construct *SemanticModel*. This object contains the following items:

- A map of variable names to *SemonticVar*, where a *SemanticVar* represents a variable in the DRS. Each has a reference to identify it (e.g. x1), a name of the object/event it is representing (e.g. dog), and a list of attributes (e.g. [hairy, large]). It has methods to set its name/add an attribute[8], set its cardinality, and to obtain a list of its semantic atoms (as defined earlier).
- A list of *SemanticRel*, representing relations in the DRS. All relations are binary, therefore we have a reference to the left and right variables (we use their name rather than a pointer to the variable itself), and the predicate being applied. It similarly has a method to generate a list of semantic atoms as defined in Section 4.

*SemanticModel* maintains all other relations, that is, equalities and implications, and has methods to generate their atoms appropriately. It is possible to have nested *SemanticModels*, in the case when the DRS expression uses the *'Prop'* (i.e. proposition) construction. When processing this, we construct a new child *SemanticModel*, set its parent to the current model, and redirect the processing of the sub-expression to the child model. Therefore when looking up a *SemanticVar* for a given name (a process required by *SemanticRel*), if it is not contained within the local context, we recursively search its ancestors.

The *KeywordIdentifier* class can be used to display the list of atoms for a DRS expression. While not used in the overall pipeline, it is useful for analytical purposes.

---

[8] We can determine an entity's type by its *'posType'* value. If we have *hairy(x1)* and hairy has the type 'a' (indicating an attribute) we add it to the list. If we have *dog(x1)* and dog has the type 'n' (for noun), we set the variable's name.

## c) Dictionary Methods

The class *SpellUtils* provides access to methods associated with misspellings, such as determining whether a word is incorrect misspelled, and if so, producing a list of suggestions. An external library *"Swabunga Spell Engine"* was used, which provides such methods.

## d) Cumulative Classification

- *Manual Mark Group Specification*

    A *marklist* file simply contains rows of the form *"feature -> [list of groups]"*. If the directory the classifier model is being outputted to contains a file named *'main.marklist'*, then this will override any other segmentation method used. A class *MarkListModifier* provides a command-line interface to reduce the hugely time-consuming burden of manually assigning groups, with commands such as *"CONTAINS word group"* to add any feature containing 'word' to the specified group.

- *Training*

    A class *ClusterUtils* provides methods to perform this ranking and converting vectors to their appropriate form:

    *List<Integer> getClustersBySimilarity(double[] row, List<Set<Integer>> clusters)* takes the row and clusters and produces a ranked list of cluster indices using the cosine measure. *double[] getClusterIncidenceVector(Set<Integer> clusters, int numFeatures)* puts the cluster in a vector form for use in the above method.
    Lastly, *double[] getRowRestrictedToCluster(Set<Integer> cluster, double[] row)* removes values from the vector that correspond to features not in the cluster. This is the 'instance' used in training each classifier.

# 12. Bibliography

[PULMAN 05] S. G. Pulman and J. Z. Sukkarieh (2005): Automatic Short Answer Marking.

[PULMAN 06] Stephen G. Pulman, Nicholas J. Raikes and Jana Z. Sukkarieh (not yet published): Customising an automated marking system.

[SHERMIS 03] Mark D. Shermis, Jill Burstein, (2003): Automated Essay Scoring: A Cross-disciplinary Perspective.

[CLARK 07] James R. Curran, Stephen Clark, and Johan Bos (2007): Linguistically Motivated Large-Scale NLP with C&C and Boxer. Proceedings of the ACL 2007 Demonstrations Session (ACL-07 demo), pp.29-32.

Artificial examiners put to the test: http://news.bbc.co.uk/1/hi/technology/6961088.stm

Data Mining: Practical Machine Learning Tools and Techniques (Second Edition). Ian H. Witten, Eibe Frank (documentotion for the Weko tools).

# 13. Acknowledgements

```java
1: package schemagenerator;
2:
3: import java.io.BufferedReader;
4: import java.io.FileReader;
5: import java.io.Serializable;
6: import java.util.LinkedList;
7: import java.util.List;
8: import java.util.Map;
9: import java.util.Set;
10: import java.util.TreeMap;
11:
12: import weka.core.Matrix;
13:
14: /**
15:  * This class provides access to the information contained
16:  * within an Arff file, such as the number of features, each
17:  * of their counts, and a matrix representing the data.
18:  * @author Jamie
19:  *
20:  */
21: public class ArffReader implements Serializable
22: {
23:     Matrix matrix;
24:     // The classifications as defined by the Arff file.
25:     List<Integer> classifiedAs = new LinkedList<Integer>();
26:     // The total number of times each of the features (represented
27:     // by its index) occurs, that is, the summation of the column.
28:     Map<Integer,Integer> featureCounts = new TreeMap<Integer,Integer>();
29:
30:     String fileName;
31:     int maxMark  0;
32:     List<String> attributes = new LinkedList<String>();
33:     public ArffReader(String fileName)
34:     {
35:         this.fileName = fileName;
36:         read();
37:     }
38:
39:
40:     /**
41:      * Reads the Arff file. Executed by the constructor.
42:      *
43:     public void read()
44:     {
45:         try
46:         {
47:             // Get attributes and read up to "@data"
48:             BufferedReader br = new BufferedReader(new FileReader(fileName));
49:             boolean reading = true;
50:             int numCols=0;
51:             while(br.ready()&&reading)
52:             {
53:                 String line = br.readLine();
54:                 if(attributes.size()!=0&&line.indexOf("@attribute")<0)
55:                     reading = false;
56:                 else if(line.indexOf("@attribute")==0){
57:                     String[] parts = line.split(" ");
58:                     if(parts[1].equals("score")){
59:                         attributes.add(parts[1].replaceAll("\"", ""));
60:                         numCols++;
61:                     }
62:
63:
64:                     br.readLine();  // @data
65:
66:                     // Now read the data.
67:
68:                     int numRows = 0;
69:                     List<List<Double>> vals = new LinkedList<List<Double>>();
70:
71:                     while(br.ready())
72:                     {
73:                         String[] lineParts = br.readLine().split(", ");
74:                         List<Double> row = new LinkedList<Double>();
75:                         numRows++;
76:                         for(int i=0; i<(lineParts.length-1); i++)
77:                         {
78:                             Double val = new Double(lineParts[i]);
79:                             // If necessary, we increment the count for that feature.
80:                             if(val==1.0)
81:                                 if(featureCounts.containsKey(i))
82:                                     featureCounts.put(i, featureCounts.get(i)+1);
83:                                 else featureCounts.put(i, 1);
84:                         }
85:                         row.add(val);
86:                     }
87:                     vals.add(row);
88:
89:                     int mark = new Integer(lineParts[lineParts.length-1]);
90:                     classifiedAs.add(mark);
91:                     if(mark>maxMark)maxMark = mark;
92:                 }
93:
94:                 Matrix mx = new Matrix(numRows,numCols);
95:
96:                 for(int i=0; i<numRows; i++)
97:                 {   for(int j=0; j<numCols; j++)
98:                     {
99:                         mx.setElement(i,j,vals.get(i).get(j));
100:                     }
101:                 }
102:                 br.close();
103:
104:                 matrix = mx;
105:
106:             } catch(Exception e){
107:                 throw new Error(e);
108:             }
109:
110:
111:
112:     /**
113:      * Gets the prescribed mark for an answer with index i.
114:      */
115:     public Integer getActualMark(Integer i)
116:     {   return classifiedAs.get(i);      }
117:
118:     public Integer getFeatureCount(Integer i)
119:     {   return featureCounts.get(i);     }
120:
121:     public Matrix getMatrix()
122:     {   return matrix.x;                 }
123:
124:     /**
```

```
125:    * @deprecated
126:    * The method returns a matrix consisting only of rows with the given mark
127:    * classification. Initially in training, only rows with a classification
128:    * of 1 were taken. The logic was that such answers would contain features
129:    * only relating to predominantly one mark, and thus identifying groups
130:    * would be more accurate. However in practice this did seem to be the
131:    * case, and had the drawback of discarding the majority of the data set.
132:    */
133:   public Matrix getMatrix(int mark)
134:   {
135:
136:       int cols = matrix.numColumns();
137:       int numRowsWithMark = 0;
138:       for(int i=0; i<classifiedAs.size(); i++)
139:           if(classifiedAs.get(i)==mark)numRowsWithMark++;
140:
141:       Matrix matrixForMark = new Matrix(numRowsWithMark,cols);
142:
143:       int currentRow = 0;
144:       for(int i=0; i<matrix.numRows(); i++)
145:       {
146:           if(classifiedAs.get(i).intValue()==mark)
147:           {
148:               for(int j=0; j<cols; j++)
149:                   matrixForMark.addElement(currentRow, j,
150:                       matrix.getElement(i, j));
151:               currentRow++;
152:           }
153:       }
154:       return matrixForMark;
155:   }
156:
157:
158:   public int getMaxMark()
159:   {
160:       return maxMark;
161:   }
162:
163:   public List<String> getAttributes()
164:   {
165:       return attributes;
166:   }
167:
168:
169:   /**
170:    * @deprecated
171:    * This method is no longer used.
172:    * Given a list of groups (i.e. sets of feature indices) and the original
173:    * matrix, it would produce a new matrix, with #groups columns, that
174:    * indicated whether each of the rows in the original matrix 'matched'
175:    * the group. For example, the nth row in the matrix of [1, 0, 1] would
176:    * indicate that the nth row of the original matrix 'matched' groups 1 and
177:    * 3, but not 2. This method was used when group identification was done
178:    * in 2 stages.
179:    */
180:   public Matrix getGroupMatrix(List<Set<Integer>> groups, Matrix mx)
181:   {
182:       Matrix gMx = new Matrix(mx.numRows(), groups.size());
183:       for(int i=0; i<mx.numRows(); i++)
184:
185:           int i;
186:           for(Set<Integer> group : groups)
187:           {
188:               if(matchesGroup(group,mx.getRow(i)))
189:                   gMx.setElement(i, j, 1.0);
190:               else gMx.setElement(i, j, 0.0);
191:               j++;
192:           }
193:
194:
195:       return gMx;
196:   }
197:
198:   /**
199:    * @deprecated This method (referenced above) indicates whether a row
200:    * matches a particular group. This was done by taking the elements of the
201:    * row occurring in the group, counting the number of 1s, and seeing if
202:    * this was greater the root of the size of the group.
203:    */
204:   public static boolean matchesGroup(Set<Integer> s, double[] row)
205:   {
206:       int matching = 0;
207:       for(int i=0; i<row.length; i++)
208:       {
209:           if(row[i]==1.0&&s.contains(i))matching++;
210:       }
211:
212:       if(new Double(matching)>=Math.floor(Math.sqrt(new Double(s.size()))))
213:           return true;
214:       else return false;
215:   }
216:
217: }
```

```java
 1: package schemagenerator;
 2:
 3: import java.util.HashSet;
 4: import java.util.LinkedList;
 5: import java.util.List;
 6: import java.util.Set;
 7:
 8: import weka.clusterers.Clusterer;
 9: import weka.clusterers.SimpleKMeans;
10: import weka.core.Attribute;
11: import weka.core.FastVector;
12: import weka.core.Instance;
13: import weka.core.Instances;
14: import weka.core.Matrix;
15:
16: /**
17:  * Provides a less cumbersome way of accessing Weka's clustering
18:  * facilities, using the incidence matrix as input data.
19:  * @author Jamie
20:  *
21:  */
22: public class ClustererWrapper {
23:
24:     Matrix incidence;
25:     int numClusters;
26:     Clusterer clusterer = new SimpleKMeans();
27:
28:     public ClustererWrapper(Matrix incidence, int numClusters)
29:     {
30:         // We normalise the columns to improve the accuracy
31:         // of the clustering.
32:         this.incidence = MatrixUtils.normaliseColumns(incidence);
33:         this.numClusters = numClusters;
34:         try
35:         {
36:             ((SimpleKMeans) clusterer).setNumClusters(numClusters);
37:         } catch(Exception e){
38:             e.printStackTrace();
39:         }
40:     }
41:
42:     /**
43:      * Converts the Matrix to an 'Instances' object, as required
44:      * by the clustering.
45:      */
46:     public void buildClusterer(List<String> attrs)
47:     {
48:         FastVector fv = new FastVector(attrs.size());
49:         for(String s : attrs)
50:             fv.addElement(new Attribute(s));
51:
52:         // Need matrix in Instances form.
53:         // Each column of the matrix is an instance.
54:         Instances ins = new Instances("somename", fv, attrs.size());
55:         for(int i=0; i<incidence.numColumns(); i++)
56:             ins.add(new Instance(1.0, incidence.getColumn(i)));
57:
58:         try
59:         {
60:             clusterer.buildClusterer(ins);
61:         } catch(Exception e){
62:             e.printStackTrace();
63:         }
64:     }
65:
66:     /**
67:      * Identifies the cluster for each column, and produces a list of sets,
68:      * each set representing a group of columns (in index form).
69:      */
70:     public List<Set<Integer>> getClusters()
71:     {
72:         List<Set<Integer>> clusters = new LinkedList<Set<Integer>>();
73:         for(int i=0; i<numClusters; i++)
74:             clusters.add(new HashSet<Integer>());
75:
76:         for(int i=0; i<incidence.numColumns(); i++)
77:         {
78:             try
79:             {
80:                 int cNum = clusterer.clusterInstance(
81:                     new Instance(1.0, incidence.getColumn(i)));
82:                 clusters.get(cNum).add(i);
83:             } catch(Exception e){
84:                 e.printStackTrace();
85:             }
86:
87:         }
88:
89:         return clusters;
90:     }
91:
92: }
```

```
 1: package schemagenerator;
 2:
 3: import java.util.LinkedList;
 4: import java.util.List;
 5: import java.util.Set;
 6:
 7: /**
 8:  * Provides utilities with regards to clusters (where
 9:  * each cluster represents a group of related features
10:  * intending to represent one mark of the mark scheme).
11:  * @author Jamie
12:  *
13:  */
14: public class ClusterUtils {
15:
16:     // Return an ordered list of cluster indexes,
17:     // such that the first cluster of the list matches
18:     // the row to the greatest extent.
19:     public static List<Integer> getClustersBySimilarity(
20:             double[] row, List<Set<Integer>> clusters)
21:     {
22:
23:         List<Double> cosSims = new LinkedList<Double>();
24:
25:         for (int i = 0; i < clusters.size(); i++) {
26:             // First turn cluster into incidence array
27:             double[] clusterArray = getClusterIncidenceVector(clusters.get(i),
row.length);
28:         }
29:
30:         // Now get vector (cosine) similarity
31:         // Ignore length of row, since constant 'given the row'.
32:         double cosSim = 0.0;
33:         for (int k = 0; k < clusterArray.length; k++)
34:             cosSim += clusterArray[k] * row[k];
35:         cosSim = cosSim / euclideanDistance(clusterArray);
36:         cosSims.add(cosSim);
37:         }
38:
39:         // Generate list clustersBySim
40:
41:         List<Integer> clustersBySim = new LinkedList<Integer>();
42:
43:         double lastMax = 10000;
44:         while (clustersBySim.size() < cosSims.size()) {
45:             int currentMaxCluster = -1;
46:             double currentMax = -1.0;
47:             for (int i = 0; i < cosSims.size(); i++) {
48:                 if (cosSims.get(i) >= currentMax && cosSims.get(i) < lastMax
49:                     && !clustersBySim.contains(i))
50:                 {
51:                     currentMax = cosSims.get(i);
52:                     currentMaxCluster = i;
53:                 }
54:             }
55:             clustersBySim.add(currentMaxCluster);
56:             lastMax = currentMax;
57:         }
58:
59:         // System.out.println("Cos sims: "+cosSims);
60:
61:         return clustersBySim;
62:     }
63:
64:     // Returns a vector representing which features occur with n
65:     // the cluster. e.g. A cluster of
66:     // { 0, 1, 3, 5 } -> (1, 1, 0, 1, 0, 1)
67:     public static double[] getClusterIncidenceVector(Set<Integer> cluster, int
numFeatures) {
68:         double[] clusterArray = new double[numFeatures];
69:         for (int i = 0; i < numFeatures; i++)
70:             clusterArray[i] = 0.0;
71:
72:         for (Integer i : cluster)
73:             clusterArray[i] = 1.0;
74:
75:         return clusterArray;
76:     }
77:
78:
79:     private static double euclideanDistance(double[] vArray) {
80:         double temp = 0.0;
81:         for (int i = 0; i < vArray.length; i++)
82:             temp += Math.pow(vArray[i], 2);
83:         temp = Math.sqrt(temp);
84:         return temp;
85:     }
86:
87:     // This gives a vector containing only the feature values
88:     // corresponding to features contained by the cluster.
89:     // For example, if the cluster is { 0, 1, 4 } and the row
90:     // [1, 0, 0, 0, 1], then we obtain the vector [1, 0, 1]
91:     public static double[] getRowRestrictedToCluster(Set<Integer> cluster,
92:         double[] row)
93:     {
94:         double[] clusterArray = getClusterIncidenceVector(cluster,row.length);
95:         double[] toReturn = new double[getClusterVectorSize(clusterArray)];
96:         int count = 0;
97:         for (int i = 0; i < clusterArray.length; i++)
98:             if (clusterArray[i] == 1.0) {
99:                 toReturn[count] = row[i];
100:                count++;
101:            }
102:        return toReturn;
103:     }
104:
105:     public static int getClusterVectorSize(double[] clusterArray) {
106:         int oneCount = 0;
107:         for (int i = 0; i < clusterArray.length; i++)
108:             if (clusterArray[i] == 1.0)oneCount++;
109:         return oneCount;
110:     }
111: }
```

```java
1: package schemagenerator;
2:
3: import java.io.PrintWriter;
4: import java.util.List;
5: import java.util.Set;
6:
7: import weka.classifiers.Classifier;
8: import weka.classifiers.bayes.NaiveBayes;
9: import weka.core.Attribute;
10: import weka.core.FastVector;
11: import weka.core.Instance;
12: import weka.core.Instances;
13: import weka.core.Matrix;
14:
15: /**
16:  * A classifier which takes into account the cumulative nature of a
17:  * classification. That is, the classifications x and y if 2 independent
18:  * properties X and Y are satisfied correspondingly should result in a
19:  * classification of x+y if both X and Y are satisfied.
20:  *
21:  * @author Jamie
22:  *
23:  */
24: public class CumulativeClassifier extends Classifier {
25:     // The classifiers for each mark.
26:     Classifier[] classifiers;
27:     // The data set relating to each mark.
28:     Instances[] instancesArray;
29:     // The features relating to each mark in the mark scheme.
30:     List<Set<Integer>> clusters;
31:     int totalClusters = 0;
32:     // Note that the maximum mark is not necessarily
33:     // the same as the number of clusters (i.e. if the mark
34:     // is limited to a maximum below the total feasible marks).
35:     int maxMark = 0;
36:     // Outputs the mark scheme to a .markscheme file.
37:     String outputFile = null;
38:     // Optionally inputs a manually constructed .marklist
39:     // file, which overrides the clustering and specifies
40:     // which mark each feature corresponds to.
41:     String inputFile = null;
42:
43:     public CumulativeClassifier() {
44:     }
45:
46:     public CumulativeClassifier(String inputFile, String outputFile) {
47:         this.inputFile = inputFile;
48:         this.outputFile = outputFile;
49:     }
50:
51:     public CumulativeClassifier(String inputFile, String outputFile,
52:             int totalClusters)
53:     {
54:         this(inputFile, outputFile);
55:         this.totalClusters = totalClusters;
56:     }
57:
58:     public void buildClassifier(Instances instances) throws Exception {
59:
60:         // The feature strength matrix is an nxn matrix (where there are n
61:         // features) such that M(i,j) indicates the number of times features
62:         // i and j appear together in the dataset.

63:         Matrix featureStrengthMx = MatrixUtils
64:                 .getIncidenceMatrix( InstancesReader
65:                         .getMatrixFrom instances instances );
66:
67:         maxMark = InstancesReader.getMaxMark(instances);
68:         // if we haven't got the total number of clusters,
69:         if (totalClusters == 0)
70:             totalClusters = maxMark;
71:
72:         List<String> attributes = InstancesReader.getAttributes(instances);
73:
74:         if (inputFile == null)
75:             clusters = GroupIdentifier.identifyGroupsUsingClustering(
76:                     featureStrengthMx, attributes, maxMark);
77:         // else clusters = SchemaUtils.readMarkSchemeFile(inputFile,
78:         // attributes);
79:         else
80:             clusters = SchemaUtils.readMarkListFile(inputFile, attributes,
81:                     totalClusters);
82:
83:         classifiers = new Classifier[totalClusters];
84:         instancesArray = new Instances[totalClusters];
85:
86:         // STEP 0:
87:         // Output mark scheme if necessary.
88:         if (outputFile != null) {
89:             System.out.println("Outputting mark scheme to: "
90:                     + outputFile);
91:             PrintWriter pw = new PrintWriter(outputFile);
92:             pw.println(attributes);
93:             for (Set<Integer> cluster : clusters)
94:                 pw.println(cluster);
95:             pw.close();
96:         }
97:
98:         // STEP 1:
99:         // Initiate the classifiers, and start a fresh instances object
100:        // for each of the classifiers, each with the correct attribute
101:        // information.
102:
103:        for (int i = 0; i < totalClusters; i++) {
104:            classifiers[i] = new NaiveBayes();
105:            // classifiers[i] = new DecisionTable();
106:        }
107:
108:        for (int i = 0; i < maxMark; i++)
109:        {
110:            int sizeOfCluster = ClusterUtils.getClusterVectorSize(ClusterUtils
111:                    .getClusterIncidenceVector(clusters.get(i), attributes
112:                    .size()));
113:            FastVector fv = new FastVector(sizeOfCluster + 1);
114:            for (int j = 0; j < sizeOfCluster; j++)
115:                fv.addElement(new Attribute("attr" + j));
116:
117:            // The classification is a NOMINAL. Either 0 or 1,
118:            // therefore need to make group.
119:            FastVector binary = new FastVector(2);
120:            binary.addElement("0");
121:            binary.addElement("1");
122:            fv.addElement(new Attribute("score", binary));
123:            instancesArray[i] = new Instances("iSet" + i, fv, 100);
124:            instancesArray[i]
```

```java
125:                .setClassIndex(instancesArray[i].numAttributes() - 1);
126:        }
127:
128:        // STEP 2:
129:        // Now populate the instances object for each of the classifiers.
130:
131:        Matrix oMx = InstancesReader.getMatrixFromInstances(instances);
132:
133:        for (int i = 0; i < oMx.numRows(); i++) {
134:
135:            double[] row = oMx.getRow(i);
136:
137:            // Need to find clusters it's most similar to
138:            List<Integer> clustersOrdered = ClusterUtils
139:                .getClustersBySimilarity(row, clusters);
140:            int classification = InstancesReader
141:                .getClassification(instances, i);
142:
143:            // We take the first 'classification' clusters from ordered
144:            // clusters, and train the networks for these clusters.
145:            for (int j = 0; j < clusters.size(); j++) {
146:                // Train cluster
147:
148:                // But first, need to construct instance!
149:                // To do this, we simply map incidence vector for cluster
150:                // on the row.
151:                double[] clusterArray =
152:                    ClusterUtils.getRowRestrictedToCluster(
153:                        clusters.get(clustersOrdered.get(j)), row);
154:                Instance instance = new Instance(clusterArray.length + 1);
155:                instance.setDataset(instancesArray[clustersOrdered.get(j)]);
156:                for (int k = 0; k < clusterArray.length; k++)
157:                    instance.setValue(k, clusterArray[k]);
158:                if (j < classification) {
159:                    instance.setValue(clusterArray.length, "1");
160:                } else {
161:                    instance.setValue(clusterArray.length, "0");
162:                }
163:
164:                instancesArray[clustersOrdered.get(j)].add(instance);
165:
166:            }
167:
168:        }
169:
170:        // STEP 3:
171:        // Finally, build the classifiers using the instances objects.
172:
173:        try {
174:            for (int i = 0; i < totalClusters; i++)
175:                classifiers[i].buildClassifier(instancesArray[i]);
176:        } catch (Exception e) {
177:            e.printStackTrace();
178:        }
179:
180:        // For the user to see, print the mark scheme
181:        SchemaUtils.print(clusters, attributes);
182:    }
183:
184:    public double classifyInstance(Instance instance)
185:            throws java.lang.Exception {
186:        // For each classifier, we need to:
187:        // 1. Restrict the instance only to the attributes it contains.
188:        // 2a. Do some biggery pockery to make the instance suitable for
189:        //     classification.
190:        // 2b. Classify this restricted instance.
191:        // 3a. Add together all classifications.
192:        // 3b. Cap the mark at maxMark.
193:
194:        double total = 0.0;
195:
196:        for (int i = 0; i < classifiers.length; i++) {
197:            double[] restrictedInstance = ClusterUtils
198:                .getRowRestrictedToCluster(clusters.get(i), instance
199:                    .toDoubleArray());
200:            double[] withArbitraryClassification = appendClassification(
201:                restrictedInstance, 2.0);
202:            Instance instanceB = new Instance(1.0,
203:                withArbitraryClassification);
204:            instanceB.setDataset(instancesArray[i]);
205:            try {
206:                double d = classifiers[i].classifyInstance(instanceB);
207:                total += d;
208:            } catch (Exception e) {
209:                throw new Error(e);
210:            }
211:        }
212:
213:        if (total > maxMark) total = maxMark;
214:
215:        return total;
216:    }
217:
218:    private double[] appendClassification(double[] cArray, double toAppend) {
219:        double[] cArrayNew = new double[cArray.length + 1];
220:        for (int i = 0; i < cArray.length; i++)
221:            cArrayNew[i] = cArray[i];
222:        cArrayNew[cArray.length] = toAppend;
223:        return cArrayNew;
224:    }
225:
226: }
```

```java
1: package schemagenerator;
2:
3: import java.util.HashSet;
4: import java.util.LinkedList;
5: import java.util.List;
6: import java.util.Set;
7:
8: import weka.core.Matrix;
9:
10: /**
11:  * This class supplies different methods to identify
12:  * groups of features based on the features' co-occurence.
13:  * The intended output of such methods is a list of sets,
14:  * where each set represents a group of features (in index
15:  * form).
16:  * @author Jamie
17:  */
18:
19: public class GroupIdentifier
20: {
21:     /**
22:      * This method uses standard (simple K-means) clustering.
23:      */
24:     public static List<Set<Integer>> identifyGroupsUsingClustering(Matrix Featu
reStrengthMx, List<String> features, int numClusters)
25:     {
26:         featureStrengthMx = MatrixUtils.normaliseToMax(featureStrengthMx);
27:         ClustererWrapper cw = new ClustererWrapper(featureStrengthMx, numClust
ers);
28:         cw.buildClusterer(features);
29:         return cw.getClusters();
30:     }
31:
32:     // Put in mind matrix is normalised.
33:     // Anything under value is considered a 'weak bond'.
34:     // Anything above is consider 'strong'.
35:     final static double MINIMUM_BOND = 0.4;
36:
37:     /**
38:      * @deprecated
39:      * Produces the groups by identifying cliques (where feature
40:      * nodes have an edge between them if their connection is
41:      * sufficiently strong) before merging them into the desired
42:      * number of groups.
43:      */
44:     public static Set<Set<Integer>> identifyGroupsUsingCliqueMerging(Matrix Fea
tureStrengthMx)
45:     {
46:         featureStrengthMx = MatrixUtils.normaliseToMax(featureStrengthMx);
47:
48:         Set<Set<Integer>> groups = new HashSet<Set<Integer>>();
49:
50:         for(int i=0; i<featureStrengthMx.numRows(); i++)
51:         {
52:             System.out.println("Processing row "+i);
53:             if(groups.isEmpty())
54:             {
55:                 // Simple start case
56:                 Set<Integer> s = new HashSet<Integer>();
57:                 s.add(i);
58:                 groups.add(s);
59:             } else
60:             {
61:                 Set<Set<Integer>> newGroups = new HashSet<Set<Integer>>();
62:                 boolean didntAppend = true;
63:
64:                 for(Set<Integer> group : groups)
65:                 {
66:                     Set<Integer> strongBonds = new HashSet<Integer>();
67:                     Set<Integer> weakBonds = new HashSet<Integer>();
68:
69:                     for(Integer j : group)
70:                     {   if(j != i)
71:                             if(featureStrengthMx.getElement(i,j)>MINIMUM_BOND)s
trongBonds.add(j);
72:                             else weakBonds.add(j);
73:                     }
74:
75:                     // System.out.println("Strong bonds: "+strongBonds);
76:                     // System.out.println("Weak bonds: "+weakBonds);
77:
78:                     // If weakBonds is empty, then the element becomes part of
that set.
79:                     if(weakBonds.isEmpty())
80:                     {
81:                         group.add(i);
82:                         didntAppend = false;
83:                     } else if(strongBonds.isEmpty())
84:                     {
85:                         // If strongBonds is empty, then we need to create a n
ew set
86:                         // containing just that feature. HOWEVER, this only oc
curs
87:                         // if the feature was not added to any set.
88:
89:                     } else {
90:                         // Both strongBonds and weakBonds is non-empty.
91:                         // We simply make a new group with the strongBonds
92:                         // and the element to add. The original group is
93:                         // unaffected.
94:                         strongBonds.add(i);
95:                         newGroups.add(strongBonds);
96:                         didntAppend = false;
97:                     }
98:                 }
99:
100:                if(didntAppend)
101:                {
102:                    Set<Integer> s = new HashSet<Integer>();
103:                    s.add(i);
104:                    newGroups.add(s);
105:                }
106:                groups.addAll(newGroups);
107:            }
108:        }
109:
110:        return groups;
111:    }
112:
113:    public static void main(String[] args)
114:    {
115:        // Test the above group identification
116:    }
```

```
117:        List<String> features = new LinkedList<String>();
118:        features.add("A");
119:        features.add("B");
120:        features.add("C");
121:        features.add("D");
122:        features.add("E");
123:        features.add("F");
124:        features.add("G");
125:        features.add("H");
126:
127:        Matrix mx = new Matrix(4,4);
128:        mx.setElement(0,0, 1.0);
129:        mx.setElement(0,1, 0.0);
130:        mx.setElement(0,2, 1.0);
131:        mx.setElement(0,3, 1.0);
132:        mx.setElement(1,1, 1.0);
133:        mx.setElement(1,2, 1.0);
134:        mx.setElement(1,3, 0.0);
135:        mx.setElement(2,2, 1.0);
136:        mx.setElement(2,3, 1.0);
137:        mx.setElement(3,3, 1.0);
138:        mirrorMx(mx);
139:
140:        Set<Set<Integer>> groups = GroupIdentifier.identifyGroupsUsingCliqueMe
rging(mx);
141:        SchemaUtils.print(groups, features);
142:
143:        System.out.println("");
144:
145:        Matrix mx2 = new Matrix(5,5);
146:        mx2.setElement(0,0, 1.0);
147:        mx2.setElement(0,1, 1.0);
148:        mx2.setElement(0,2, 1.0);
149:        mx2.setElement(0,3, 0.0);
150:        mx2.setElement(0,4, 0.0);
151:        mx2.setElement(1,1, 1.0);
152:        mx2.setElement(1,2, 1.0);
153:        mx2.setElement(1,3, 0.0);
154:        mx2.setElement(1,4, 0.0);
155:        mx2.setElement(2,2, 1.0);
156:        mx2.setElement(2,3, 1.0);
157:        mx2.setElement(2,4, 1.0);
158:        mx2.setElement(3,3, 1.0);
159:        mx2.setElement(3,4, 1.0);
160:        mx2.setElement(4,4, 1.0);
161:        mirrorMx(mx2);
162:
163:        groups = GroupIdentifier.identifyGroupsUsingCliqueMerging(mx2);
164:        SchemaUtils.print(groups, features);
165:        System.out.println("");
166:
167:        Matrix mx3 = new Matrix(6,6);
168:        mx3.setElement(0,0, 1.0);
169:        mx3.setElement(0,1, 1.0);
170:        mx3.setElement(0,2, 1.0);
171:        mx3.setElement(0,3, 1.0);
172:        mx3.setElement(0,4, 0.0);
173:        mx3.setElement(0,5, 0.0);
174:        mx3.setElement(1,1, 1.0);
175:        mx3.setElement(1,2, 1.0);
176:        mx3.setElement(1,3, 1.0);
177:        mx3.setElement(1,4, 0.0);
```

```
178:        mx3.setElement(1,5, 0.0);
179:        mx3.setElement(2,2, 1.0);
180:        mx3.setElement(2,3, 1.0);
181:        mx3.setElement(2,4, 0.0);
182:        mx3.setElement(2,5, 0.0);
183:        mx3.setElement(3,3, 1.0);
184:        mx3.setElement(3,4, 0.0);
185:        mx3.setElement(3,5, 1.0);
186:        mx3.setElement(4,4, 1.0);
187:        mx3.setElement(4,5, 1.0);
188:        mx3.setElement(5,5, 1.0);
189:        mirrorMx(mx3);
190:
191:        groups = GroupIdentifier.identifyGroupsUsingCliqueMerging(mx3);
192:        groups = Merger.mergeGroups(groups, mx3, 2);
193:        SchemaUtils.print(groups, features);
194:
195:        System.out.println();
196:
197:        ArffReader ar = new ArffReader("./files/13b11-0.arff");
198:        Matrix mx4 = MatrixUtils.getIncidenceMatrix(ar.getMatrix());
199:        mx4 = MatrixUtils.normaliseToMax(mx4);
200:
201: //     groups = GroupIdentifier.identifyGroups(mx4);
202: //     SchemaPrinter.print(groups, ar.getAttributes());
203: //     System.out.println("Merging");
204: //     groups = Merger.mergeGroups(groups, mx4, ar.getMaxMark());
205: //     SchemaPrinter.print(groups, ar.getAttributes());
206:
207:        List<Set<Integer>> clusters = identifyGroupsUsingClustering(mx4, ar.ge
tAttributes(), ar.getMaxMark());
208:        SchemaUtils.print(clusters, ar.getAttributes());
209:    }
210:
211:    /**
212:     * Using by the demo to copy a triangular matrix into
213:     * its diagonal reflection so that the resulting matrix
214:     * is symmetric.
215:     */
216:    private static void mirrorMx(Matrix mx)
217:    {
218:        for(int r=0; r<mx.numColumns(); r++)
219:            for(int c=0; c<r; c++)
220:                mx.setElement(r,c,mx.getElement(c,r));
221:    }
222:
223: }
```

```
 1: package schemagenerator;
 2:
 3: import java.util.LinkedList;
 4: import java.util.List;
 5:
 6: import weka.core.Attribute;
 7: import weka.core.Instance;
 8: import weka.core.Instances;
 9: import weka.core.Matrix;
10:
11: /**
12:  * Whereas ArffReader obtains the representative matrix and other
13:  * data from an Arff file, this provides similar extraction from
14:  * an Instances object.
15:  * @author Jamie
16:  *
17:  */
18: public class InstancesReader {
19:     public static Matrix getMatrixFromInstances(Instances instances)
20:     {
21:         Matrix mx = new Matrix(instances.numInstances(),
22:                 instances.numAttributes()-1);
23:
24:         for(int row = 0; row < instances.numInstances(); row++)
25:         {   Instance ins = instances.instance(row);
26:             double[] insArray = ins.toDoubleArray();
27:             for(int col = 0; col < instances.numAttributes()-1; col++)
28:             {
29:                 mx.setElement(row, col, ins.value(col));
30:             }
31:         }
32:         return mx;
33:
34:     }
35:
36:     public static int getMaxMark(Instances instances)
37:     {
38:         double maxMark = 0.0;
39:         for(int row = 0; row < instances.numInstances(); row++)
40:         {
41:             Instance ins = instances.instance(row);
42:             double[] insArray = ins.toDoubleArray();
43:             double classification = insArray[instances.numAttributes()-1];
44:             maxMark = (classification>maxMark) ? classification : maxMark;
45:         }
46:
47:         return (new Double(maxMark)).intValue();
48:     }
49:
50:     public static List<String> getAttributes(Instances instances)
51:     {
52:         List<String> attributes = new LinkedList<String>();
53:         // Don't want the last attribute (which is 'score')
54:         for(int i=0; i<instances.numAttributes()-1; i++)
55:         {
56:             Attribute att = instances.attribute(i);
57:             attributes.add(att.name());
58:         }
59:         return attributes;
60:     }
61:
62:     public static int getClassification(Instances instances, int index
```

```
63:     {
64:         Instance ins = instances.instance(index);
65:         return (new Double(ins.toDoubleArray()[instances.numAttributes
66:                 ].intValue());
67:     }
68: }
```

```java
1: package schemagenerator;
2:
3: import java.io.BufferedReader;
4: import java.io.BufferedReader;
5: import java.io.InputStreamReader;
6: import java.io.PrintWriter;
7: import java.util.LinkedList;
8: import java.util.List;
9: import java.util.Map;
10: import java.util.Set;
11: import java.util.TreeMap;
12: import java.util.TreeSet;
13:
14: /**
15:  * Used to convert a markscheme file into a marklist file.
16:  * A marklist file is a more convenient form for manual
17:  * editing, and has lines of the form:
18:  * "feature [list of groups-it-appears-in]"
19:  * If a file main.marklist is placed in a test directory,
20:  * this will override the mark groups.
21:  * @author Jamie
22:  *
23:  */
24: public class MarkSchemeLister {
25:
26:     static Map<String,Set<Integer>> groupMappings =
27:         new TreeMap<String,Set<Integer>>();
28:
29:     public static void main(String[] args) throws Exception {
30:         if(args.length ==0)
31:         {
32:             System.out.println("USAGE:");
33:             System.out.println("\tMarkSchemeLister in.markscheme out.marklist");
34:             System.out.println("\t\t or");
35:             System.out.println("\tMarkSchemeLister in.marklist out.marklist");
36:             System.exit(0);
37:         }
38:
39:         String inputFile = args[0];
40:         String outputFile = args[1];
41:
42:         // Read in
43:         if(inputFile.contains("markscheme"))
44:         {
45:             BufferedReader br = new BufferedReader(new FileReader(inputFile));
46:
47:             String attrs = br.readLine();
48:             int maxMark = 0;
49:             while(br.ready()){
50:                 maxMark++;
51:                 br.readLine();
52:             }
53:             br.close();
54:             System.out.println("NOTE: There are "+maxMark+" groups.");
55:
56:             // Now put attributes in list.
57:             attrs = attrs.substring(1,attrs.length()-1);
58:             List<String> attributes = new LinkedList<String>();
59:             String[] attrs2 = attrs.split(", ");
60:             for(int i=0; i<attrs2.length; i++)attributes.add(attrs2[i]);
61:
62:             // put in map
63:             // Put in map, with all initialised to empty lists (i.e. no group)
64:             for(String attr : attributes)
65:                 groupMappings.put(attr, new TreeSet<Integer>());
66:         } else {
67:             // Read marklist file
68:             groupMappings = SchemaUtils.getGroupMappings(inputFile);
69:         }
70:
71:
72:         // Provide a command line interface to modify features' groups.
73:         System.out.println("Commands:");
74:         System.out.println("\tASSIGN attr groups");
75:         System.out.println("\tCONTAINS word groups");
76:         System.out.println("\tQUIT");
77:         System.out.println("*******************************");
78:
79:         BufferedReader br = new BufferedReader(
80:             new InputStreamReader(System.in));
81:         boolean quit = false;
82:
83:         while(!quit)
84:         {
85:             System.out.print("> ");
86:             String line = br.readLine();
87:
88:             String[] lineParts = line.split(" ");
89:             if(lineParts[0].equals("ASSIGN"))
90:             {
91:                 String[] groups = lineParts[2].split(",");
92:                 for(int j=0; j<groups.length; j++)
93:                     groupMappings.get(lineParts[1]).add(
94:                         new Integer(groups[j]));
95:                 output(outputFile);
96:             }
97:             if(lineParts[0].equals("CONTAINS"))
98:             {
99:                 String[] groups = lineParts[2].split(",");
100:                 for(String attr : groupMappings.keySet())
101:                     if(attr.contains(lineParts[1])){
102:                         for(int j=0; j<groups.length; j++)
103:                             groupMappings.get(attr).add(new Integer(groups[j]));
104:                 output(outputFile);
105:                 }
106:             if(lineParts[0].equals("QUIT"))quit = true;
107:         }
108:         br.close();
109:     }
110:
111:
112:     public static void output(String fileName) throws Exception {
113:         PrintWriter pw = new PrintWriter(fileName);
114:         for(String attr : groupMappings.keySet())
115:         {
116:             pw.print(attr+"\t");
117:             pw.print(groupMappings.get(attr));
118:             pw.println();
119:         }
120:         pw.close();
121:     }
122: }
```

```java
1: package schemagenerator;
2:
3: import java.io.BufferedReader;
4: import java.io.FileReader;
5: import java.io.InputStreamReader;
6: import java.io.PrintWriter;
7: import java.util.HashSet;
8: import java.util.LinkedList;
9: import java.util.List;
10: import java.util.Set;
11:
12: /**
13:  * @deprecated
14:  * @author Jamie
15:  *
16:  */
17: public class MarkSchemeModifier {
18:         public static void main(String[] args) throws Exception {
19:                 if(args.length= 0)
20:                 {
21:                         System.out.println("USAGE:");
22:                         System.out.println("\tMarkSchemeModifier in.markscheme
out.markscheme");
23:                         System.exit(0);
24:                 }
25:
26:                 String inputFile = args[0];
27:                 String outputFile = args[1];
28:
29:                 BufferedReader br = new BufferedReader(new FileReader(inputFil
e));
30:                 String attrs = br.readLine();
31:                 int maxMark = 0;
32:                 while(br.ready()){
33:                         maxMark++;
34:                         br.readLine();
35:                 }
36:                 br.close();
37:                 System.out.println("NOTE: There are "+maxMark+" groups.");
38:
39:                 // Now put attributes in list.
40:                 attrs = attrs.substring(1,attrs.length()-1);
41:                 List<String> attributes = new LinkedList<String>();
42:                 String[] attrs2 = attrs.split(", ");
43:                 for(int i=0; i<attrs2.length; i++)attributes.add(attrs2[i]);
44:
45:                 System.out.println("For each attribute enter a group from 1 to
"+maxMark+" to place it in.");
46:                 System.out.println("For multiple groups, separate by spaces."
);
47:                 System.out.println("------------------------------------------"
);
48:
49:
50:                 // Now get new groups
51:                 Set<Integer>[] newGroups = new Set[maxMark];
52:                 for(int i=0; i<maxMark; i++)newGroups[i] = new HashSet<Integer
>();
53:
54:
55:                 br = new BufferedReader(new InputStreamReader(System.in));
56:                 for(int i=0; i<attributes.size(); i++
57:                         System.out.println(attributes.get(i));
58:                         System.out.print("> ");
59:                         String groupSelection = br.readLine();
60:                         if(!groupSelection.equals(""))
61:                         {
62:                                 String[] groups = groupSelection.split(" ");
63:                                 for(int j=0; j<groups.length; j++)newGroups[ne
w Integer(groups[j])-1].add(i);
64:                         }
65:
66:                 }
67:
68:
69:                 System.out.println("-------------------------------------------")
;
70:                 System.out.println("Done, writing to file.");
71:
72:                 PrintWriter pw = new PrintWriter(outputFile);
73:                 pw.println(attributes);
74:                 System.out.println(attributes);
75:                 for(int i 0; i<maxMark; i++){ pw.println(newGroups[i]); System
.out.println(newGroups[i]); }
76:                 pw.close();
77:         }
78: }
```

```java
 1: package schemagenerator;
 2:
 3: import weka.core.Instance;
 4: import weka.core.Instances;
 5: import weka.core.Matrix;
 6:
 7: /**
 8:  * Utilities relating to matrix manipulation.
 9:  * @author Jamie
10:  */
11: public class MatrixUtils {
12:     public static Matrix normaliseToMax(Matrix mx)
13:     {
14:         // We take the log of all the values.
15:         // This ensures very high values do not scale
16:         // everything else down.
17:         // Then we scale everything so that the maximum is 1.
18:
19:         Matrix copy = (Matrix) mx.clone();
20:
21:         // Find max
22:         double currentMax = 0.0;
23:         for(int i=0; i<mx.numRows(); i++)
24:         {   for(int j=0; j<mx.numRows(); j++)
25:             {
26:                 // Set element to the log of it.
27:                 // We plus 1 since log(0) is undefined.
28:                 // Therefore 0 -> 0 as wanted.
29:                 copy.setElement(i,j,Math.log(copy.getElement(i,j)+1.0));
30:                 if(copy.getElement(i,j)>currentMax)
31:                     currentMax = copy.getElement(i,j);
32:             }
33:         }
34:
35:         // Divide everything by max
36:         // Diagonals are set to 1.0, such that a feature will always match
37:         // with itself.
38:         for(int i=0; i<mx.numRows(); i++)
39:             for(int j=0; j<mx.numRows(); j++)
40:                 if(i!=j)copy.setElement(i,j,
41:                     copy.getElement(i,j) / currentMax);
42:                 else copy.setElement(i,j, 1.0);
43:
44:         return copy;
45:     }
46:
47:     // Obtains (A^T).A
48:     public static Matrix getIncidenceMatrix(Matrix mx) {
49:         return mx.transpose().multiply(mx);
50:     }
51:
52:     public static Matrix normaliseColumns(Matrix mx)
53:     {
54:         Matrix mx2 = new Matrix(mx.numRows(), mx.numColumns());
55:         // Normalises columns
56:         for(int i=0; i<mx.numColumns(); i++)
57:         {
58:             double squaredSum = 0.0;
59:             for(int j=0; j<mx.numRows(); j++)
60:                 squaredSum+= Math.pow(mx.getElement(j, i), 2.0);
61:             squaredSum = Math.sqrt(squaredSum);
62:             for(int j=0; j<mx.numRows(); j++)
```

```java
63:                 mx2.setElement(j, i, mx.getElement(j, i) / squaredSum);
64:             }
65:
66:         return mx2;
67:     }
68: }
```

```java
 1: package schemagenerator;
 2:
 3: import java.util.HashSet;
 4: import java.util.Set;
 5:
 6: import weka.core.Matrix;
 7:
 8: /**
 9:  * A class to merge groups (connected components) into n groups
10:  * by using the feature strengths between their nodes.
11:  * @deprecated
12:  * @author Jamie
13:  *
14:  */
15: public class Merger
16: {
17:     public static Set<Set<Integer>> mergeGroups(Set<Set<Integer>> groups,
18:         Matrix featureStrengthMx, int n)
19:     {
20:         // We only want n groups, therefore we need
21:         // to merge groups in the current list until
22:         // there is a sufficiently small number.
23:
24:         // To pick which pair of groups to merge, we
25:         // find the maximum inter-group bond strength,
26:         // by taking the average of all the element pairs
27:         // in the cross product.
28:
29:         featureStrengthMx = MatrixUtils.normaliseToMax(featureStrengthMx);
30:
31:         // First, we remove any groups that only have one element. i.e.
32:         // it has no connection with anything.
33:         Set<Set<Integer>> groupsCopy =
34:             (Set<Set<Integer>>) ((HashSet<Set<Integer>>) groups).clone();
35:         for(Set<Integer> group : groups)
36:             if(group.size() == 1)groupsCopy.remove(group);
37:         groups = groupsCopy;
38:
39:         if(featureStrengthMx.numRows()!=featureStrengthMx.numColumns())
40:             throw new Error("Must have same number of rows as columns.");
41:
42:         System.out.println("Merging into "+n+" groups.");
43:
44:         while(groups.size()>n)
45:         {
46:             System.out.println("Current group size: "+groups.size());
47:             Set<Integer> bestGroup1 = null;
48:             Set<Integer> bestGroup2 = null;
49:             double bestPairScore = 0.0;
50:
51:             for(Set<Integer> group1 : groups)
52:                 for(Set<Integer> group2 : groups)
53:                 {   if(!group1.equals(group2))
54:
55:                         int linkCount = 0;
56:                         double currentPairScore = 0.0;
57:                         for(Integer i : group1)
58:                         {   for(Integer j : group2)
59:
60:                                 currentPairScore+=
61:                                     featureStrengthMx.getElement(i, j);
62:                                 linkCount++;
63:                             }
64:                         }
65:
66:                         double avgScore = currentPairScore/(double)
67:                             new Double(linkCount);
68:                         if(avgScore>bestPairScore)
69:                         {
70:                             bestPairScore = avgScore;
71:                             bestGroup1 = group1;
72:                             bestGroup2 = group2;
73:                         }
74:                     }
75:                 }
76:
77:
78:             // System.out.println("Best pair score: "+bestPairScore);
79:             // System.out.println("Best group 1: "+bestGroup1);
80:             // System.out.println("Best group 2: "+bestGroup2);
81:
82:             // Merge 2 groups:
83:             // 1. Add element of second to first.
84:
85:             for(Integer i : bestGroup2)bestGroup1.add(i);
86:             // 2. Delete second from list.
87:             // System.out.println("Removing "+bestGroup2+" from "+groups);
88:             // boolean removed = groups.remove(bestGroup2);
89:             // System.out.println(removed);
90:             Set<Set<Integer>> newGroups = new HashSet<Set<Integer>>();
91:             for(Set<Integer> group : groups)
92:                 if(!group.equals(bestGroup2))newGroups.add(group);
93:             groups = newGroups;
94:         }
95:         }
96:
97:         return groups;
98:
99:     }
100: }
```

```java
1: package schemagenerator;
2:
3: import java.io.BufferedReader;
4: import java.io.FileReader;
5: import java.util.HashSet;
6: import java.util.LinkedList;
7: import java.util.List;
8: import java.util.Map;
9: import java.util.Set;
10: import java.util.TreeMap;
11: import java.util.TreeSet;
12:
13: /**
14:  * Provides utilities with regards to reading mark scheme/list
15:  * files.
16:  * @author Jamie
17:  *
18:  */
19: public class SchemaUtils {
20:     public static void print(Set<Set<Integer>> groups,
21:         List<String> featureNames)
22:     {
23:         for(Set<Integer> group : groups)
24:         {
25:             System.out.print(("{ "));
26:             // Print each group
27:             int c = 0;
28:             for(Integer i : group)
29:             {
30:                 if(c>0)System.out.print(", ");
31:                 System.out.print(featureNames.get(i));
32:                 c++;
33:             }
34:             System.out.println(" }");
35:         }
36:     }
37:
38:     public static void print(List<Set<Integer>> groups,
39:         List<String> featureNames)
40:     {
41:         for(Set<Integer> group : groups)
42:         {
43:             System.out.print(("{ "));
44:             // Print each group
45:             int c = 0;
46:             for(Integer i : group)
47:             {
48:                 if(c>0)System.out.print(", ");
49:                 System.out.print(featureNames.get(i));
50:                 c++;
51:             }
52:             System.out.println(" }");
53:         }
54:     }
55:     public static List<Set<Integer>> readMarkSchemeFile(String inputFile,
56:         List<String> attributesToCheck)
57:     {
58:         List<Set<Integer>> clusters = new LinkedList<Set<Integer>>();
59:
60:         try
61:         {
62:             BufferedReader br = new BufferedReader(new FileReader(inputFile));
63:             int i=0;
64:             while(br.ready())
65:             {
66:                 String line = br.readLine();
67:                 if(i==0){
68:                     List<String> attributes = new LinkedList<String>();
69:                     // Reading attributes
70:                     line = line.substring(1,line.length()-1);
71:                     String[] lineParts = line.split(", ");
72:                     for(int j=0; j<lineParts.length; j++)
73:                         attributes.add(lineParts[j]);
74:                     if(!attributes.equals(attributesToCheck))
75:                         throw new Error("This mark scheme file is not " +
76:                             "compatible with the training model. " +
77:                             "Ensure that their attributes are the same.");
78:
79:                 } else {
80:                     // Remove [ ]
81:                     line = line.substring(1,line.length()-1);
82:                     String[] lineParts = line.split(", ");
83:                     Set<Integer> cluster = new HashSet<Integer>();
84:                     for(int j=0; j<lineParts.length; j++)
85:                         cluster.add(new Integer(lineParts[j]));
86:                     clusters.add(cluster);
87:                 }
88:
89:                 i++;
90:             }
91:
92:             return clusters;
93:         } catch(Exception e){
94:             throw new Error("Problem reading mark scheme file");
95:         }
96:     }
97: }
98:
99:     public static List<Set<Integer>> readMarkListFile(String inputFile,
100:         List<String> attributes, int numClusters)
101:         List<Set<Integer>> clusters = new LinkedList<Set<Integer>>();
102:         for(int i=0; i<numClusters; i++)clusters.add(new HashSet<Integer>());
103:         Map<String,Set<Integer>> groupMappings =
104:             SchemaUtils.getGroupMappings(inputFile);
105:
106:         // We have to combine 'attributes' and 'groupMappings' to obtain
107:         // 'clusters'. In the marklist file, groups numbers go from 1 to n,
108:         // whereas in lists we reference 0 to n-1, so we deduct 1
109:         // accordingly.
110:         for(int i=0; i<attributes.size(); i++)
111:         {
112:             Set<Integer> groups = groupMappings.get(attributes.get(i));
113:             for(Integer group : groups)clusters.get(group-1).add(i);
114:         }
115:
116:         return clusters;
117:     }
118:
119:     /**
120:      * Produces a map of feature names to a list of groups
121:      * each corresponds to. It takes a marklist file as
122:      * its input.
123:      */
124:     public static Map<String,Set<Integer>> getGroupMappings(String inputFile)
```

```
125:
126:            Map<String,Set<Integer>> groupMappings =
127:                    new TreeMap<String,Set<Integer>>();
128:        try{
129:
130:            BufferedReader br = new BufferedReader(new FileReader(inputFile));
131:            while(br.ready())
132:            {
133:                String[] line = br.readLine().split("\t");
134:                String attribute = line[0];
135:                String[] groups =
136:                    line[1].substring(1,line[1].length()-1).split(", ");
137:                groupMappings.put(attribute, new TreeSet<Integer>());
138:                if(!line[1].equals("[]"))
139:                    for(int i = 0; i<groups.length; i++)
140:                        groupMappings.get(attribute).add(
141:                            new Integer(groups[i]));
142:
143:            }
144:        } catch(Exception e){
145:            e.printStackTrace();
146:        }
147:        return groupMappings;
148:    }
149: }
```

```java
 1: package schemagenerator;
 2:
 3: import java.util.HashSet;
 4: import java.util.Set;
 5:
 6: /**
 7:  * @deprecated
 8:  * @author Jamie
 9:  *
10:  */
11: public class SetUtils {
12:     public static Set<Set<Integer>> removeSubsets(Set<Set<Integer>> groups)
13:     {
14:         // Clean up: merge groups when one is a subset of the other.
15:
16:         Set<Set<Integer>> groupsTemp = new HashSet<Set<Integer>>();
17:         for (Set<Integer> s : groups)groupsTemp.add(s); // Make copy of groups
18:
19:         for (Set<Integer> s : groups) {
20:             for (Set<Integer> t : groups) {
21:                 if (isSubset(s, t)&&s!=t)
22:                     groupsTemp.remove(s);
23:             }
24:         }
25:         return groupsTemp;
26:
27:     }
28:
29:     public static boolean isSubset(Set<Integer> s, Set<Integer> t) {
30:         for (Integer i : s)
31:             if (!t.contains(i))
32:                 return false;
33:         return true;
34:     }
35:
36:
```

```java
1: package semantic;
2:
3: import java.util.HashMap;
4: import java.util.LinkedList;
5: import java.util.List;
6: import java.util.Map;
7: import java.util.Set;
8: import java.util.TreeSet;
9:
10: /***
11:  * This class is a representation of the semantic DRS
12:  * expression, in object oriented form.
13:  * Its main method linearises the DRS into a set of
14:  * semantic atoms.
15:  * @author Jamie
16:  *
17:  */
18: public class SemanticModel {
19:         // A list of variables in the DRS.
20:         Map<String,SemanticVar> svars = new HashMap<String,SemanticVar>();
21:         // A list of relations in the DRS.
22:         List<SemanticRel> srels = new LinkedList<SemanticRel>();
23:         // The LHS of equalities within the DRS.
24:         // The ith element corresponds to the ith element of eqRefR.
25:         List<String> eqRefL = new LinkedList<String>();
26:         List<String> eqRefR = new LinkedList<String>();
27:         // Whether each of the above are equalities or inequalities.
28:         List<Boolean> eqSense = new LinkedList<Boolean>();
29:
30:         // The current sense while navigating the DRS tree.
31:         Boolean sense = true;
32:         // The current parent while navigating the DRS tree.
33:         // null if at the root.
34:         SemanticModel parent;
35:
36:         // The set of generated semantic atoms from the model.
37:         Set<String> atoms = new TreeSet<String>();
38:
39:         // Methods to add atoms to the growing set.
40:         public void addAtom(String atom){        atoms.add(atom);        }
41:         public void addAtoms(Set<String> atomsToAdd){    atoms.addAll(atomsToAd
d);    }
42:
43:         // Gets a variable given its name.
44:         public SemanticVar getVar(String ref)
45:         {
46:                 SemanticVar sv = svars.get(ref);
47:                 // The var might be located in the parent model if not local.
48:                 if(sv!=null) return sv;
49:                 else return getParent().getVar(ref);
50:         }
51:
52:         // Declares a variable.
53:         public void declareVar(String ref)
54:         {
55:                 svars.put(ref, new SemanticVar(ref));
56:         }
57:
58:         // Adds a predicate, where ref is the name of the variable,
59:         // symbol is the unary relation being applied to it, and
60:         // posType is the type of relation.
61:         public void addPred(String ref, String symbol, String posType
```

```java
62:         {
63:                 if(svars.containsKey(ref)
64:                 {
65:                         svars.get(ref).addPred(symbol,posType);
66:                 else if(parent!=null)
67:                         if(parent.getVar(ref)!=null)parent.getVar(ref).addPred
(symbol,posType);
68:                         else declareVar(ref);
69:                 }
70:                 else declareVar(ref);
71:         }
72:
73:         // Adds a binary relation, of the form "symbol(ref1,ref2)".
74:         public void addRel(String ref1, String ref2, String symbol)
75:         {
76:                 srels.add(new SemanticRel(ref1,ref2,symbol));
77:         }
78:
79:         // Adds an equality, of the form ref1=ref2 or ref1!=ref2.
80:         public void addEq(String ref1, String ref2,Boolean sense){
81:                 eqRefL.add(ref1);
82:                 eqRefR.add(ref2);
83:                 eqSense.add(sense);
84:         }
85:
86:         // Sets the cardinality of a variable.
87:         public void setCard(String ref, int card, String cardType)
88:         {        svars.get(ref).setCard(card);    }
89:
90:         public void changeSense(){        sense = !sense; }
91:         public boolean getSense(){        return sense;    }
92:
93:         public void setName(String ref, String symbol, String ncType)
94:         {
95:                 svars.get(ref).setName(symbol);
96:         }
97:
98:         // Used for debugging purposes.
99:         // Prints the variables and relations from
100:         // the tree.
101:         public void print()
102:         {
103:                 for(SemanticVar svar : svars.values())
104:                 {
105:                         svar.print();
106:                 }
107:
108:                 System.out.println("*** Relations ****");
109:                 for(SemanticRel srel : srels)
110:                         System.out.println(srel.toString());
111:         }
112:
113:         // Using the information accumulated from the tree,
114:         // a list of semantic atoms is generated.
115:         public Set<String> getAtoms()
116:         {
117:                 for(SemanticVar svar : svars.values())
118:                 {
119:                         svar.getAtoms(this);
120:
121:                 for( SemanticRel srel : srels)
122:                 {
```

```
123:                              srel.getAtoms(this);
124:                  }
125:
126:                      Deal with equalities.
127:                  for int i 0; i<eqRefL.size(); i++
128:                  {
129:                          SemanticVar l = getVar(eqRefL.get(i));
130:                          SemanticVar r = getVar(eqRefR.get(i));
131:                          String senseStringStart = "";
132:                          String senseStringEnd = "";
133:                          if(!eqSense.get(i))
134:                          {
135:                                  senseStringStart = "![";
136:                                  senseStringEnd = "]";
137:                          }
138:                          for(String atomL : l.getAtoms(this))
139:                          {       for(String atomR : r.getAtoms(this))
140:                                  {       // - is commutative.
141:                                          atoms.add(senseStringStart + atomL+" =
"+atomR + senseStringEnd);
142:                                          atoms.add(senseStringStart + atomR+" =
"+atomL + senseStringEnd);
143:                                  }
144:                          }
145:                          i++;
146:                  }
147:
148:                  return atoms;
149:          }
150:
151:          public SemanticModel getParent() {
152:                  return parent;
153:          }
154:
155:          public void setParent(SemanticModel parent) {
156:                  this.parent = parent;
157:          }
158: }
```

```java
 1: package semantic;
 2:
 3: import java.util.HashSet;
 4: import java.util.Set;
 5:
 6: /**
 7:  * Represents a relation in the semantic model.
 8:  * @author Jamie
 9:  *
10:  */
11: public class SemanticRel {
12:
13:         // Represents:
14:         // name:type(ref1,ref2)
15:         String ref1;
16:         String ref2;
17:         String name;
18:         String type = "";
19:
20:         // Relations which we ignore.
21:         String[] rejectableKeywords = { "topic" };
22:
23:         // We may need to produce the list of atoms
24:         // numerous times, so they are cached for
25:         // subsequent use.
26:         Set<String> cachedAtoms;
27:
28:         public SemanticRel(String ref1, String ref2, String symbol)
29:         {
30:                 this.ref1 = ref1;
31:                 this.ref2 = ref2;
32:                 this.name = symbol;
33:         }
34:
35:         public String toString()
36:         {
37:                 return name+"["+ref1+","+ref2+"]";
38:         }
39:
40:         // The rule for getting atoms is as such:
41:         // * We get the sets of atoms A and B for the
42:         //    left and right variables in the relations.
43:         // * We take the set:
44:         //    { name(a,b) | a<-A, b<-B }
45:         public Set<String> getAtoms(SemanticModel model,
46:         {
47:                 if(cachedAtoms! null)return cachedAtoms;
48:                 Set<String> set = new HashSet<String>();
49:
50:                 String senseString = "";
51:                 if(!model.getSense() )senseString = "1";
52:
53:                 if(!shouldReject())
54:                 {
55:                         Set<String> setL = model.getVar(ref1).getAtoms(model);
56:                         Set<String> setR = model.getVar(ref2).getAtoms(model);
57:
58:                         for(String atomL : setL)
59:                                 for(String atomR : setR)
60:                                         set.add(senseString+name+"["+atomL+","
+atomR+"]");
61:
62:                                 model.addAtoms(set);
63:
64:
65:                 cachedAtoms = set;
66:                 return set;
67:         }
68:
69:         private boolean shouldReject()
70:         {
71:                 for(int i = 0; i<rejectableKeywords.length; i++)
72:                         if(rejectableKeywords[i].equals(name))return true;
73:
74:                 return false;
75:         }
76:
77: }
```

```java
1:  package semantic;
2:
3:  import java.util.HashSet;
4:  import java.util.LinkedList;
5:  import java.util.List;
6:  import java.util.Set;
7:
8:  public class SemanticVar {
9:          // Of the form:
10:         //   [attr1 ^ ... ^ attrN][name]
11:         String ref;
12:         String name;
13:         List<String> attrs = new LinkedList<String>();
14:
15:         // Generic variable names which are ignored.
16:         String[] rejectableKeywords = { "event", "proposition" };
17:
18:         String type = "";
19:         int card = 1; // need to change to lower and upper bound
20:
21:         SemanticModel prop;
22:         public void setProp(SemanticModel model){      prop = model;    }
23:
24:         // We may need to produce the list of atoms
25:         // numerous times, so they are cached for
26:         // subsequent use.
27:         Set<String> cachedAtoms;
28:
29:         public SemanticVar(String ref){ this.ref = ref; }
30:
31:         // We the DRS contains propositions such as
32:         // John(x1) and large(x1), we need to distinguish
33:         // between the name of the variable, and attributes
34:         // relating to it. We can determine this by the
35:         // posType.
36:         public void addPred(String symbol,String posType)
37:         {
38:                 if(!shouldReject(symbol))
39:                 {
40:                         if(posType.equals("n")||posType.equals("v"))
41:                         {
42:                                 name = symbol;
43:                         } else
44:                                 attrs.add(symbol);
45:                 } else {
46:                         type = symbol;
47:                 }
48:         }
49:
50:
51:         // Sets the cardinality of the variable.
52:         public void setCard(int card){   this.card = card;      }
53:
54:         // For debug purposes.
55:         public void print()
56:         {
57:                 System.out.println("**** "+ref+" ****");
58:                 System.out.println(name);
59:                 for(String attr : attrs)
60:                         System.out.println(attr+"["+name+"]");
61:         }
62:
63:
64:         // Gets a list of atoms represented by
65:         // the variable.
66:         public Set<String> getAtoms(SemanticModel model)
67:         {
68:                 if(cachedAtoms!=null) return cachedAtoms;
69:
70:                 String senseString = "";
71:                 if(model.getSense) senseString = "!";
72:
73:                 Set<String> set = new HashSet<String>();
74:                 if(shouldReject(name))return set;
75:
76:                 if(name==null&&prop==null)
77:                 {
78:                         return set;
79:                 }
80:
81:                 if(prop==null)
82:                 {
83:                         // Get variable name itself
84:                         set.add(senseString+name);
85:                         // Get cardinality if present
86:                         if(card!=1)set.add(senseString+"|"+name+"|="+card);
87:                         // Get var with attributes
88:                         for(String attr : attrs)
89:                         {
90:                                 set.add(senseString+attr+"["+name+"]"); // add attr with subject
91:                                 set.add(senseString+attr);     // add attr
92:                         }
93:                 } else {
94:                         // We have a proposition
95:                         Set<String> childAtoms = prop.getAtoms();
96:                         // Add original child atoms to this level
97:                         set.addAll(childAtoms);
98:                         // Add attributed child atoms if any attributes
99:                         for(String attr : attrs)
100:                                for(String atom : childAtoms)
101:                                        if(!shouldReject(atom))
102:                                        {
103:                                                set.add(senseString+attr+"["+atom+"]");
104:                                                set.add(senseString+attr);
105:                                        }
106:                 }
107:
108:                 model.addAtoms(set);
109:                 cachedAtoms = set;
110:                 return set;
111:        }
112:
113:        public String getName() {
114:                return name;
115:        }
116:
117:        public void setName(String name) {
118:                this.name = name;
119:        }
120:
121:        public String getType()
```

```
122:              return type;
123:
124:
125:       public void setType(String type) {
126:              this.type = type;
127:
128:
129:       public int getCard() {
130:              return card;
131:       }
132:
133:
134:       private boolean shouldReject(String symbol)
135:       {
136:              for(int i=0; i<rejectableKeywords.length; i++)
137:                     if(rejectableKeywords[i].equals(symbol))return true;
138:
139:              return false;
140:       }
141:
142: }
```

```java
1: package spelling;
2:
3: import java.io.BufferedReader;
4: import java.io.FileReader;
5: import java.io.PrintWriter;
6: import java.util.HashMap;
7: import java.util.Map;
8:
9: /**
10:  * Takes a CSV containing the English sentences and
11:  * "spell correction mapping file" (where each line
12:  * contains a spelling correction), and produces a
13:  * new CSV with the corrections made.
14:  * @author Jamie
15:  *
16:  */
17: public class MapCorrections {
18:     final static String inputFile = "./files/13bii-answers.csv";
19:     final static String outputFile = "./files/13bii-answers2.csv";
20:     final static String mappingFile = "./files/13bii-spelling-corrections.
txt";
21:     Map<String,String> corrections = new HashMap<String,String>();
22:
23:     public static void main(String[] args) {
24:
25:         MapCorrections mc = new MapCorrections();
26:         if(args.length==0)mc.map(mappingFile,inputFile,outputFile);
27:         else mc.map(args[0],args[1],args[2]);
28:     }
29:
30:     public void map(String mappingFile, String inputFile, String outputFil
e)
31:     {
32:         try
33:         {
34:             // Read in spelling corrections
35:             BufferedReader br = new BufferedReader(new FileReader(
mappingFile));
36:             while(br.ready())
37:             {
38:                 String[] parts = br.readLine().split("\t");
39:                 if(parts.length!=2)System.out.println("Problem
with line "+parts[0]);
40:                 corrections.put(parts[0], parts[1]);
41:             }
42:             br.close();
43:
44:             // Read in input file, and for each line, make the cor
rections
45:             // and write to output file.
46:             br = new BufferedReader(new FileReader(inputFile));
47:             PrintWriter pw = new PrintWriter(outputFile);
48:
49:             while(br.ready())pw.println(transformLine(br.readLine(
)));
50:
51:             br.close();
52:             pw.close();
53:
54:         } catch(Exception e){
55:             e.printStackTrace();
56:         }
57:     }
58:
59:     public String transformLine(String line)
60:     {
61:         for(String key : corrections.keySet())
62:             if(line.toLowerCase().indexOf(key)>=0)
63:             {
64:                 int pos = line.toLowerCase().indexOf(key);
65:
66:                 // Need to check it's not part of some other w
ord
67:                 if(pos>0&&pos<line.length()-key.length())
68:                 {
69:                     if(!Character.isLetter(line.charAt(pos
-1))&&!Character.isLetter(line.charAt(pos+key.length())))
70:                     {   // Not part of a word, so repl
ace!
71:                         line = line.substring(0,pos)+c
orrections.get(key)+line.substring(pos-key.length());
72:                     }
73:                 } else if(pos==0){
74:                     if(!Character.isLetter(line.charAt(pos
+key.length())))
75:                     {
76:                         line = corrections.get(key)+li
ne.substring(pos+key.length());
77:                     }
78:                 } else {
79:                     if(!Character.isLetter(line.charAt(pos
-1)))
80:                     {
81:                         // Not part of a word, so repl
ace!
82:                         line = line.substring(0,pos)+c
orrections.get(key);
83:                     }
84:                 }
85:             }
86:         }
87:
88:         return line;
89:     }
90: }
```

```java
1: package spelling;
2: import java.io.BufferedReader;
3: import java.io.FileReader;
4: import java.io.PrintWriter;
5: import java.util.Map;
6: import java.util.TreeMap;
7:
8: import com.swabunga.spell.engine.Word;
9: import com.swabunga.spell.event.StringWordTokenizer;
10:
11: /**
12:  * Takes a dictionary as well as the original English
13:  * answers, and produces a "spell correction mapping
14:  * file" which indicates what words should be replaced
15:  * and what to replace them with.
16:  * Outputting this correction file rather than modifying
17:  * the CSV immediately allows manual alterations to be
18:  * made.
19:  * @author Jamie
20:  *
21:  */
22: public class SpellCorrectionGenerator
23: {
24:         // The dictionary.
25:         final static String dicFile = "C:/Users/Jamie/Desktop/Programming/exam
marking/dictionary/eng_com.dic";
26:         // The CSV containing potentially misspelled English answers.
27:         final static String fileToCheck = "./files/13bii-answers.csv";
28:         // The file to output.
29:         final static String mappingFile = "./files/13bii-spelling-corrections.
txt";
30:
31:         public static void main(String[] args) throws Exception {
32:
33:                 SpellCorrectionGenerator scg = new SpellCorrectionGenerator();
34:                 if(args.length   )scg.generate(dicFile,fileToCheck,mappingFile
);
35:                 else scg.generate(dicFile,args[0],args[1]);
36:         }
37:
38:         public void generate(String dicFile, String fileToCheck, String mappin
gFile) throws Exception
39:         {
40:                 // A utility class to check spellings.
41:                 SpellUtils spellUtils = new SpellUtils(dicFile);
42:                 BufferedReader br = new BufferedReader(new FileReader(fileToCh
eck));
43:
44:                 // We keep a count of each of the words in the data-set.
45:                 // This allows us to detect correctly spelled words which
46:                 // are not in the dictionary.
47:                 Map<String, Integer> wordCounts = compareWords(fileToCheck);
48:
49:                 // Remove common words
50:                 wordCounts.remove("the");
51:                 wordCounts.remove("has");
52:
53:                 // Want to use wordCounts to produce a map for misspelled word
s.
54:                 Map<String, String> spellingCorrections = new TreeMap<String,
String>  ;
55:
56:         for  String val : wordCounts.keySet
57:
58:                         // Correct plurality errors such as "babys -> babies".
59:                         boolean correctedPlurality = spellUtils.correctPlural
ityErrors(
60:                                 spellingCorrections, val);
61:
62:                         if (!correctedPlurality) {
63:                                 // Is the word in the dictionary?
64:                                 boolean correct = spellUtils.isCorrect(val);
65:                                 // Is the word before the 's if applicable c
orrect?
66:                                 boolean correctA = spellUtils.isApostropheSA
dCorrect(val);
67:
68:                                 if (correct || correctA || val.length()<=2)
rds alone.
69:                                         // Leave short or correctly spelled wo
LEAVE");
70:                                         System.out.println("** " + val + " ->
71:                                 } else {
y.
72:                                         // We have a word not in the dictionar
73:
74:                                         // If it's commonly occurring in the d
ata-set,
75:                                         // leave it alone.
76:                                         if (wordCounts.get(val) > 5) {
77:                                                 // Do nothing.
78:                                         } else {
79:                                                 // SpellUtils.split(val) is no
longer used,
80:                                                 // and simply returns null.
81:                                                 String splitWord = spellUtils.
split(val);
82:                                                 if(splitWord!=null)
83:
84:                                                         spellingCorrections.pu
t(val, splitWord);
85:                                                 else
86:                                                         double bestScore = Dou
ble.NEGATIVE_INFINITY;
87:                                                         String bestWord = "";
88:
89:                                                         // Find the word in th
e data-set which
90:                                                         // has the strongest r
elationship with
91:                                                         // the misspelled word
.
92:                                                         for (String val2 : wor
dCounts.keySet())  {
93:                                                                 ce = StringUtils
94:                                                                 .getLevenshteinDistance(val, val2);
95:                                                                 StringUtils
96:                                                                 .getLetterDifference val, val2 ;
97:                                                                 rBonus = StringUtils
```

```
 98:
.getFirstLetterBonus(val, val2);
  99:
StringUtils.getSizeDifference(val,
 100:
val2);
 101:
 102:
(Math.log(wordCounts.get(val2)    Math
 103:
.log(2))
 104:
- (3 * firstLetterBonus)
 105:
2
 106:
 107:
* editDistance
 108:
- letterDiff    sizeDiff;
stScore
 109:
&& !val2.equals(val)
 110:
&& editDistance <= 5 // We limit edit
 111:
                                        // distance
 112:
&& wordCounts.get(val2) > wordCounts
 113:
                .get(val) // Can only map to a
 114:
                                // more popular word.
 115:
                                // Avoids cycles.
 116:
 117:
ore = score;
 118:
rd = val2;
 119:
 120:
 121:
 122:                            // Additionally find a
 word in the dictionary
 123:                            // which has the closes
 relationship with the
 124:                            // word. This may over
 ride a selection from
 125:                            // the data-set.
 126:
llUtils.getSuggestions(val)) {
 127:
 128:
("Word" val3).toString();
 129:
ce   StringUtils
 130:
.getLevenshteinDistance(val, val2);
 131:
  StringUtils
 132:
```

```
        .getLetterSetDifference(val, val2);
 133:
int sizeDiff   rBonus   StringUtils
 134:
.getFirstLetterBonus(val, val2);
 135:
double score   StringUtils.getSizeDifference(val,
 136:
val2);
 137:
 138:
1   (2 * editDistance)
 139:
- (3 * firstLetterBonus) - letterDiff
 140:
- sizeDiff;
 141:
if (score > be   stScore) {
 142:
ore = score;
 143:
rd = val2;
 144:
 145:                                    }
 146:
 147:                                        spellingCorrections.pu
t(val, bestWord);
 148:                                    }
 149:                                }
 150:                            }
 151:                        }
 152:                    }
 153:
 154:            // Adds common spelling mistakes.
 155:            spellUtils.addCommonMistakes(spellingCorrections);
 156:
 157:            // Write correction mapping to a file.
 158:            PrintWriter pw = new PrintWriter(mappingFile);
 159:
 160:            for (String val : spellingCorrections.keySet()) {
 161:                System.out.print("** " + val);
 162:                String lastVal = null;
 163:
 164:                // A word may map to a word which maps to another.
 165:                String temp   spellingCorrections.get(val);
 166:                while (temp != null) {
 167:                    lastVal = temp;
 168:                    System.out.print(" -> " + temp);
 169:                    // Need to ignore identity mappings
 170:                    if(!temp.equals(spellingCorrections.get(temp)
)temp = spellingCorrections.get(temp);
 171:                    else temp = null;
 172:                }
 173:
 174:                if lastVal! null)pw.print(val "\t"+lastVal);
 175:                System.out.println();
 176:            }
 177:            pw.close ;
 178:
 179:
 180:
 181:        * get ...'s for all the words in the data set.
```

```
int sizeDiff

double score

if (score > be

        bestSc

        bestWo

            }
    }

// Additionally find a
// which has the cloes

// word. This may over

/ the data-set.
for (Object val3 : spe

    String val2 =

    int editDistan

    int letterDiff
```

```
182:               *
183:           public static Map<String, Integer> collateWords(String fileName ...
184:               Map<String, Integer> wordCounts = new TreeMap<String, Integer>
...;
185:
186:               try {
187:                   BufferedReader br = new BufferedReader(new FileReader(
fileToCheck));
188:                   while (br.ready()) {
189:
190:                       String line = br.readLine();
191:                       if(line.indexOf("\t")>0)
192:                           .
193:                           String sentence = line.split("\t")[1];
194:                           StringWordTokenizer swt = new StringWo
rdTokenizer(sentence);
195:                           while (swt.hasMoreWords()) {
196:                               String word = swt.nextWord().t
oLowerCase();
197:
198:                               if (wordCounts.get(word) == nu
ll)
199:                                   wordCounts.put(word, 0
);
200:                                   wordCounts.put(word, wordCount
s.get(word) + 1);
201:
202:                               }
203:                           }
204:                       }
205:
206:               } catch (Exception e) {
207:                   e.printStackTrace();
208:               }
209:
210:               System.out.println("*** Words ***");
211:               for(String s : wordCounts.keySet())
212:                   System.out.println(s+" ("+wordCounts.get(s)+")");
213:               System.out.println("*************");
214:
215:               return wordCounts;
216:           }
217:
218: }
```

```java
1: package spelling;
2: import java.io.File;
3: import java.io.IOException;
4: import java.util.List;
5: import java.util.Map;
6:
7: import com.swabunga.spell.engine.SpellDictionary;
8: import com.swabunga.spell.engine.SpellDictionaryHashMap;
9: import com.swabunga.spell.event.SpellCheckEvent;
10: import com.swabunga.spell.event.SpellCheckListener;
11: import com.swabunga.spell.event.SpellChecker;
12: import com.swabunga.spell.event.StringWordTokenizer;
13:
14: /**
15:  * A wrapper to make the methods of a dictionary
16:  * library used more accessible. The necessity is
17:  * mainly because the library has no direct method
18:  * of determining a word's correctness, and alerts
19:  * a SpellCheckListener when a misspelling is
20:  * detected.
21:  * @author Jamie
22:  *
23:  */
24: public class SpellUtils {
25:
26:         boolean flag;
27:         List suggestions;
28:         SpellDictionary dictionary;
29:         SpellChecker spellChecker;
30:
31:         public SpellUtils(String fileName) throws IOException
32:         {
33:                 this(new SpellDictionaryHashMap(new File(fileName)));
34:         }
35:
36:         public SpellUtils(SpellDictionary dictionary)
37:         {
38:                 this.dictionary = dictionary;
39:                 spellChecker = new SpellChecker(dictionary);
40:             spellChecker.addSpellCheckListener(new SuggestionListener());
41:         }
42:
43:         public class SuggestionListener implements SpellCheckListener
44:                 public void spellingError(SpellCheckEvent event) {
45:                         flag = false;
46:                         suggestions = event.getSuggestions();
47:                 }
48:
49:         }
50:
51:         public boolean isCorrect(String word)
52:         {
53:                 flag=true;
54:                 spellChecker.checkSpelling(new StringWordTokenizer(word));
55:                 return flag;
56:         }
57:
58:         public List getSuggestions(String word)
59:         {
60:                 spellChecker.checkSpelling(new StringWordTokenizer(word));
61:                 return suggestions;
```

```java
63:         }
64:
65:         /**
66:          * Makes corrections with regards to incorrect suffixes
67:          * for plurals. For example "babys" and "halfs".
68:          */
69:         public boolean correctPluralityErrors(Map<String,String> spellCorrect
ions, String val)
70:         {
71:                 String transformedPlural = "";
72:                 if(val.endsWith("ys"))transformedPlural = val.substring(0,val.
length()-2) + "ies";
73:                 if(val.endsWith("lfs"))transformedPlural = val.substring(0,val
.length()-3) + "lves";
74:
75:                 if(!transformedPlural.equals(""))
76:                 {
77:                         if(isCorrect(transformedPlural))
78:                         {
79:                                 spellCorrections.put(val, transformedPlural);
80:                                 return true;
81:                         } else {
82:                                 return false;
83:                         }
84:                 } else {
85:                         return false;
86:                 }
87:         }
88:
89:         public void addCommonMistakes(Map<String,String> spellCorrections)
90:         {
91:                 spellCorrections.put("aswell","as well");
92:                 spellCorrections.put("dont", "don't");
93:                 spellCorrections.put("everytime","every time");
94:                 spellCorrections.put("colour","colour");        // colour not
in dictionary!
95:                 spellCorrections.put("color", "colour");
96:                 spellCorrections.put("childs", "children");
97:                 spellCorrections.put("i", "I");
98:         }
99:
100:         public boolean isApostrophedWordCorrect(String word)
101:         {
102:                 // Note, if isCorrect for the whole word returns true,
103:                 // then the output of this method does not matter.
104:
105:                 int aIndex = word.indexOf("'");
106:                 // If it doesn't contain apostrophe, test not applicable.
107:                 if(aIndex== -1) return false;
108:
109:                 if(word.endsWith("'s"))
110:                 {
111:                         String subword = word.substring(0,aIndex);
112:                         if(isCorrect(subword)) return true;
113:                         else return false;
114:                 } else {
115:                         return false;
116:                 }
117:         }
118:
119:         /**
120:          * @deprecated
```

```
121:            * Used to detect if two words had been merged
122:            * together. split("aswell") would return "as well"
123:            * and split "alot") return "a lot". It was found
124:            * the method led to poor results as it split words
125:            * often when it shouldn't have.
126:            *
127:           public String split String word)
128:           {
129:                   for(int i 2; i<,word.length()-1; i++)
130: //                 {
131: //                         String wordL = word.substring(0,i);
132: //                         String wordR   word.substring(i,word.length());
133: //                         if(isCorrect(wordL)&&isCorrect(wordR))return wordL+" "
+wordR;
134: //                 }
135:                   return null;
136:           }
137:
138: }
```

```java
  1: package spelling;
  2: import java.util.HashSet;
  3: import java.util.Set;
  4:
  5: /**
  6:  * Provides comparison methods for strings.
  7:  * @author Jamie
  8:  *
  9:  */
 10: public class StringUtils {
 11:
 12:     public static int getSizeDifference(String s, String t)
 13:     {
 14:             return Math.abs(s.length() - t.length());
 15:     }
 16:
 17:     /**
 18:      * Treats the words as a 'bag of letters', and
 19:      * for 2 sets of letters A and B, computes:
 20:      * #((A^¬B)v(B^¬A))
 21:      */
 22:     public static int getLetterSetDifference(String s, String t)
 23:     {
 24:             Set<Character> sLetters = new HashSet<Character>();
 25:             Set<Character> tLetters = new HashSet<Character>();
 26:             for(int i=0; i<s.length(); i++)sLetters.add(s.charAt(i));
 27:             for(int i=0; i<t.length(); i++)tLetters.add(t.charAt(i));
 28:
 29:             int diff = 0;
 30:             for(Character c : sLetters)if(!tLetters.contains(c))diff++;
 31:             for(Character c : tLetters)if(!sLetters.contains(c))diff++;
 32:
 33:             return diff;
 34:     }
 35:
 36:     /**
 37:      * As a heuristic, we class two words as particularly
 38:      * similar if their first letter is the same.
 39:      * Therefore "receive" and "recieve" would be more
 40:      * similar than "believe" and "relieve".
 41:      */
 42:     public static int getFirstLetterBonus(String s, String t)
 43:     {
 44:             if(s.charAt(0)==t.charAt(0))return 1;
 45:             else return 0;
 46:     }
 47:
 48:     /**
 49:      * String edit distance.
 50:      * Unknown author.
 51:      */
 52:     public static int getLevenshteinDistance (String s, String t, )
 53:             if (s == null || t == null) {
 54:                 throw new IllegalArgumentException("Strings must not be nu
11");
 55:             }
 56:
 57:             int n = s.length(); // length of s
 58:             int m = t.length(); // length of t
 59:
 60:             if (n == 0) {
 61:                 return m;

 62:             } else if (m == 0) {
 63:                 return n;
 64:             }
 65:
 66:             int p[] = new int[n+1]; //'previous' cost array, horizontall
 67:             int d[] = new int[n+1]; // cost array, horizontally
 68:             int _d[]; //placeholder to assist in swapping p and d
 69:
 70:             // indexes into strings s and t
 71:             int i; // iterates through s
 72:             int j; // iterates through t
 73:
 74:             char t_j; // jth character of t
 75:
 76:             int cost; // cost
 77:
 78:             for (i = 0; i<=n; i++) {
 79:                 p[i] = i;
 80:             }
 81:
 82:             for (j = 1; j<=m; j++) {
 83:                 t_j = t.charAt(j-1);
 84:                 d[0] = j;
 85:
 86:                 for (i = 1; i<=n; i++) {
 87:                     cost = s.charAt(i-1)==t_j ? 0 : 1;
 88:                     // minimum of cell to the left+1, to the top+1, diagon
ally left and up +cost
 89:                     d[i] = Math.min(Math.min(d[i-1]+1, p[i]+1),  p[i-1]+co
st);
 90:                 }
 91:
 92:                 // copy current distance counts to 'previous row' distanc
e counts
 93:
 94:                 _d = p;
 95:                 p = d;
 96:                 d = _d;
 97:             }
 98:
 99:             // our last action in the above loop was to switch d and p,
so p now
100:             // actually has the most recent cost counts
101:             return p[n];
102:     }
103: }
```

```java
 1: package server;
 2: import java.io.*;
 3: import java.net.*;
 4: import java.util.Set;
 5: import java.util.TreeSet;
 6:
 7: import parser.DRSReader;
 8:
 9: public class ClientCAndC {
10:         final static int port = 49157; // port to which to connect
11:         final static String computer = "curlew.comlab.ox.ac.uk";
12:         final static String endString = ".";
13:
14:         PrintWriter toServer;
15:         BufferedReader fromServer;
16:
17:         public static void main(String[] args) throws IOException {
18:                 BufferedReader in = new BufferedReader(new InputStreamReader(System.in));
19:
20:
21:                 // initialize socket
22:                 System.out.println("** Attempting to connect. **");
23:                 Socket t = new Socket(computer, port);
24:                 System.out.println("** Client connected.");
25:                 System.out.println("** Please input a sentence to parse and box.");
26:
27:                 // Associate writer with socket
28:                 PrintWriter pw = new PrintWriter(t.getOutputStream(), true);
29:
30:                 // Also associate reader with socket for receiving back response.
31:                 BufferedReader br = new BufferedReader(new InputStreamReader(t
32:                                 .getInputStream()));
33:
34:                 // Repeatedly input string from terminal, and write to socket, until
35:                 // endString is input
36:                 String st;
37:                 System.out.print("> ");
38:                 st = in.readLine();
39:                 while (!st.equals(endString)) {
40:                         pw.println(st);
41:
42:                         String response = br.readLine();
43:                         // System.out.println(response);
44:
45:                         Set<String> set = DRSReader.atomize(response);
46:                         for(String s : set)System.out.println(s);
47:
48:                         System.out.print("> ");
49:                         st = in.readLine();
50:                 }
51:                 ;
52:
53:                 System.out.println("Connection closed");
54:         }
55:
56:         public void connect() throws IOException
57:         {
58:                 // Initialize socket
59:                 Socket t = new Socket(computer, port);
60:
61:                 // Associate writer with socket
62:                 toServer = new PrintWriter(t.getOutputStream(), true);
63:
64:                 // Also associate reader with socket for receiving back response.
65:                 fromServer = new BufferedReader(new InputStreamReader(t
66:                                 .getInputStream()));
67:
68:                 System.out.println("*-- Connected to server");
69:         }
70:
71:         public String getDRS(String sentence) throws IOException
72:         {
73:                 toServer.println(sentence);
74:                 String response = fromServer.readLine();
75:                 System.out.println("*-- Got response from client.");
76:                 return response;
77:         }
78:
79:         public void close() throws IOException
80:         {
81:                 // Giving the sentence '.' tells the server thread to shut down.
82:                 toServer.println(".");
83:                 // But we still need to close the connection.
84:                 toServer.close();
85:                 fromServer.close();
86:
87:                 System.out.println("*-- Closed connection to server");
88:         }
89:
90: }
```

```java
1: package server;
2:
3: import java.io.BufferedReader;
4: import java.io.FileReader;
5: import java.io.InputStreamReader;
6: import java.io.PrintWriter;
7: import java.net.ServerSocket;
8: import java.net.Socket;
9:
10: /**
11:  * Listens to client requests to parse English sentences,
12:  * and sends back the parsed expression in DRS form.
13:  * Can operate over a remote connection.
14:  * @author Jamie
15:  *
16:  */
17: public class ServerCAndC {
18:
19:         final static int port = 49157; // port to which to listen
20:         final static String INPUT_FILE = "input.txt";
21:         final static String INTERMEDIATE_FILE = "output.txt";
22:         final static String OUTPUT_FILE = "outputBOXED.txt";
23:
24:         // Starts the soap server, which receives parse commands.
25:         final static String soapServerCommand = "../candc/candc-1.00/bin/soap_
server --server localhost:9000 --candc ../candc/candc-1.00/models/boxer";
26:         // Starts the soap client, which sends parse commands.
27:         final static String soapClientCommand = "../candc/candc-1.00/bin/soap_
client --input "+INPUT_FILE+" --output "+INTERMEDIATE_FILE+" --url http://localhost:90
00";
28:         // A command to 'box' the CCG to produce a DRS.
29:         final static String boxerCommand = "../candc/candc-1.00/bin/boxer --in
put "+INTERMEDIATE_FILE+" --output "+OUTPUT_FILE+" --box true --flat";
30:
31:         public static void main(String[] args)
32:                 Runtime r = Runtime.getRuntime();
33:                 try
34:                 {
35:
36:                         ServerCAndC cle = new ServerCAndC();
37:                         execute(soapServerCommand,false,20000); // Start the S
OAP server
38:                         startServerListener();
39:
40:                 } catch(Exception e){
41:                         e.printStackTrace();
42:                 }
43:         }
44:
45:         public static void execute(String command,boolean wait,long delay)
46:         {
47:                 try
48:                 {
49:                         // Execute the command via command line prompt.
50:                         Process p = Runtime.getRuntime().exec(command);
51:                         // Some commands we can detect when they're finished.
52:                         // For starting the soap server, we don't know when it
's
53:                         // ready, so instead we wait a suitable amount of time
54:                         if(wait)
55:                         {
```

```java
56:                                 p.waitFor();
57:                         } else {
58:                                 System.out.println("Sleeping for "+delay+"ms t
o allow server to start.");
59:                                 Thread.sleep(delay);
60:                         }
61:                 } catch(Exception e){
62:                         e.printStackTrace();
63:                 }
64:         }
65:
66:         /**
67:          * Reads the file produced from boxing, and extracts
68:          * the DRS string from it.
69:          */
70:         public static String extractDRS()
71:         {
72:                 try
73:                 {
74:                         BufferedReader br = new BufferedReader(new FileReader(
"./outputBOXED.txt"));
75:                         int i=0;
76:                         String line = "";
77:                         while(i<=13&&br.ready())
78:                         {
79:                                 line = br.readLine();
80:                                 i++;
81:                         }
82:
83:                         if(i==14)
84:                         {
85:                                 int lastBracketPos = line.lastIndexOf(')');
86:                                 line = line.substring(0,lastBracketPos).trim();
87:                                 return line;
88:                         } else {
89:                                 // Case where the C&C parser failed
90:                                 return "";
91:                         }
92:
93:                 } catch(Exception e){
94:                         throw new Error(e);
95:                 }
96:         }
97:
98:         /**
99:          * Listens to incoming connection requests.
100:          */
101:         public static void startServerListener()
102:         {
103:                 try
104:                 {
105:                         System.out.println("*** Listening for incoming connect
ions. ***");
106:
107:                         // Initialise server socket
108:                         ServerSocket listener = new ServerSocket(port);
109:
110:                         // We want the program to always accept connections, s
o we loop continuously.
111:                         while(true){
112:
```

```
113:                    // Wait for, and then accept, connections
114:                    Socket s = listener.accept();
115:                            // Create a new thread using this socket.
116:                    ServerThread st = new ServerThread( );
117:                    t.start();      // Start the thread.
118:
119:                }
120:
121:            } catch(Exception e){
122:                    e.printStackTrace();
123:            }
124:
125:        }
126:
127:        /**
128:         * The CAndC Parser reads the input string from a file.
129:         * The method puts this string in a file to produce
130:         * the suitable input for the parser.
131:         */
132:        public static void writeInputToFile(String str)
133:        {
134:            try
135:            {
136:                    PrintWriter pw = new PrintWriter(INPUT_FILE);
137:                    pw.println(str);
138:                    pw.close();
139:            } catch(Exception e){
140:                    e.printStackTrace();
141:            }
142:
143:        }
144: }
```

```java
1: package server;
2:
3: import java.io.BufferedReader;
4: import java.io.InputStreamReader;
5: import java.io.PrintWriter;
6: import java.net.Socket;
7:
8: /**
9:  * A ServerThread in instantiated by ServerCAndC, and
10:  * represents a service for a particular client.
11:  * Several ServerThreads may be running if multiple
12:  * clients are being concurrently served
13:  * @author Jamie
14:  *
15:  */
16: public class ServerThread extends Thread {
17:
18:         Socket s;
19:
20:         public ServerThread(Socket s){  this.s = s;       }
21:
22:         public void run()
23:         {
24:             try
25:             {
26:                 System.out.println("*** Accepted incoming connection. ***");
27:
28:
29:                 // Associate reader with socket
30:                     BufferedReader br =
31:                             new BufferedReader(new InputStreamReader(s.getIn
putStream()));
32:
33:                     PrintWriter pw = new PrintWriter(s.getOutputStream(),true)
;
34:
35:
36:                     // Read just 1 line
37:                     String st = br.readLine();
38:
39:                     // Now read until we terminate.
40:                     while(!st.equals("."))
41:                     {
42:                             System.out.println("Parsing: "+st);
43:
44:                             ServerCAndC.writeInputToFile(st);
45:                             ServerCAndC.execute(ServerCAndC.swapClientComm
and,true,);       // parse the input using CandC
46:                             ServerCAndC.execute(ServerCAndC.boxerCommand,t
rue,);            // box the output CCG
47:                             String drs = ServerCAndC.extractDRS();
48:
49:                             System.out.println("DRS: "+drs);
50:
51:                             // Prints the DRS string to the socket.
52:                             pw.println(drs);
53:
54:                             System.out.println("*** Request successfully p
rocessed. ***");
55:
56:                             st = br.readLine();
57:                             if(st==null)break;
58:                     }
59:
60:
61:                     br.close();
62:                     pw.close();
63:
64:                     System.out.println("*** Connection terminated ***");
65:
66:             } catch(Exception e){
67:                 e.printStackTrace();
68:
69:
70:             }
71: }
```

```java
 1: package pipeline;
 2:
 3: import java.io.BufferedReader;
 4: import java.io.FileReader;
 5: import java.io.PrintWriter;
 6:
 7: /**
 8:  * Converts the pl file into a neater CSV.
 9:  * Rows of the CV have the form:
10:  * "mark it answer"
11:  * @author Jamie
12:  *
13:  */
14: public class AnswerConverter {
15:
16:     static final String inputFile = "./files/answers/2/4(d).pl";
17:     static final String outputFile = "./files/answers/2/4d-answers.csv";
18:
19:     public static void main(String[] args) {
20:         try
21:         {
22:             BufferedReader br = new BufferedReader(new FileReader(inputFile));
23:             PrintWriter pw = new PrintWriter(outputFile);
24:
25:             while(br.ready())
26:             {
27:                 String line = br.readLine();
28:
29:                 if(line.indexOf("%") ==0||line.indexOf("noanswer")> 0||line.ind
exOf("answer")!=0)
30:                 {
31:                     // Reject
32:                 } else {
33:                     int openBrac = line.indexOf("[");
34:                     int nextComma = line.indexOf (",",openBrac);
35:                     char score = getScore(line,nextComma+1);
36:                     int openQuote = line.indexOf("\"",nextComma);
37:                     int closeQuote = line.indexOf("\"",openQuote+1);
38:                     String answer;
39:                     if(closeQuote>=0)answer = line.substring(openQuote+1,c
loseQuote);
40:                     else answer   = line.substring(openQuote+1);
41:
42:                     answer = answer.replaceAll("\t","");
43:                     answer = answer.replaceAll(","," ,");
44:                             // make sure there's space before comma
45:                     answer = answer.replaceAll("\\."," \\.");
46:                             // make sure there's space before full stop
47:                     answer = answer.replaceAll("  "," ");
48:                             // remove duplicate spaces
49:                     answer = answer.trim();
50:
51:                     System.out.println(score+"\t"+answer);
52:                     pw.println(score+"\t"+answer);
53:                 }
54:             }
55:
56:             br.close();
57:             pw.close();
58:
59:         } catch(Exception e)
60:             e.printStackTrace();
```

```java
61:         }
62:     }
63:
64:
65:     private static char getScore(String line, int startPos)
66:     {
67:         if(Character.isDigit(new Character(line.charAt(startPos))))
68:             return line.charAt(startPos);
69:         else return getScore(line,startPos+1);
70:     }
71: }
```

```java
1: package pipeline;
2: import java.io.BufferedReader;
3: import java.io.FileReader;
4: import java.io.PrintStream;
5: import java.io.PrintWriter;
6: import java.util.HashMap;
7: import java.util.LinkedList;
8: import java.util.List;
9: import java.util.Map;
10: import java.util.Set;
11: import java.util.TreeSet;
12:
13: import parser.DRSreader;
14:
15:
16: /**
17:  * Generates an arff file.
18:  * The input file is a CSV with tab separated fields "score    DRS string".
19:  * It extracts all the semantic atoms from each DRS, collates them,
20:  * and then produces a binary table with the atoms as the columns.
21:  * @author Jamie
22:  *
23:  */
24: public class ArffGenerator {
25:
26:     static String inputFile   "./files/8b-answers-processed.csv";
27:     static String outputFile = "./files/8b.arff";
28:
29:     // If this is set to true, atoms which occur only 1 time
30:     // in the corpus are removed.
31:     final static boolean CONDENSE = true;
32:     public static void main(String[] args)
33:     {
34:         System.out.println("Usage: ArffGenerator INPUTFILE.csv "+
35:                 "OUTPUTFILE.arff");
36:         generateFile(args[0],args[1],1,0,false,System.out);
37:     }
38:
39:     public static void generateFile(String inputFile,String outputFile,
40:             int numsegments, int segnumber, boolean exclude,
41:             PrintStream ps)
42:     {
43:
44:         ps.println("EXECUTING ARFF GENERATOR");
45:         ps.println("Input file: "+inputFile);
46:         ps.println("Output file: "+outputFile);
47:
48:         List<Integer> scores = new LinkedList<Integer>();
49:         List<Set<String>> atomSets = new LinkedList<Set<String>>();
50:         Map<String,Integer> allAtoms = new HashMap<String,Integer>();
51:         Set<String> finalAtomList = new TreeSet<String>();
52:
53:         try
54:         {
55:             BufferedReader br = new BufferedReader(
56:                 new FileReader(inputFile));
57:
58:             int maxScore = 0;
59:             int numLines=0;
60:
61:             while(br.ready())
62:             {

63:                 String line = br.readLine();
64:                 if(line.indexOf("%")!=0)
65:                 {
66:                     String[] lineParts = line.split("\t");
67:
68:                     int theScore = new Integer(lineParts[0]);
69:                     if(theScore>maxScore)maxScore = theScore;
70:                     scores.add(theScore);
71:                     Set<String> set = DRSreader.atomize(lineParts[1]);
72:                     atomSets.add(set);
73:
74:                     ps.println("Successfully parsed line "+numLines+".");
75:                 }
76:                 numLines++;
77:             }
78:
79:             br.close();
80:
81:             ps.println("Num lines: "+numLines);
82:
83:             // Work out segments
84:             int startPos = 0;
85:             int segmentLength = 0;
86:             int k=0;
87:
88:             while(k<=segnumber)
89:             {
90:                 startPos += segmentLength;
91:                 segmentLength = numLines / numsegments; // takes divisor;
92:                 numsegments--;
93:                 numLines -= segmentLength;
94:                 k++;
95:             }
96:
97:             if(!exclude){
98:                 startPos = -1;
99:                 segmentLength   0;
100:             }
101:
102:             // Want a set with all the atoms (but we want
103:             // to ignore the testing segment)
104:             for(int i=0; i<atomSets.size(); i++)
105:             {
106:                 if(!(i>=startPos&&i<startPos+segmentLength))
107:                 {
108:                     Set<String> atoms   atomSets.get(i);
109:                     for(String atom : atoms)
110:                     {
111:                         if(allAtoms.get(atom)==null)allAtoms.put(atom, 1);
112:                         else allAtoms.put(atom, allAtoms.get(atom)+1);
113:                     }
114:                 }
115:             }
116:
117:             // IF an atom has a count of 1, we only include it if
118:             // we're not condensing the list.
119:             for(String atom : allAtoms.keySet())
120:                 if(allAtoms.get(atom)>1)finalAtomList.add(atom);
121:                 else if(!CONDENSE)finalAtomList.add(atom);
122:
123:
124:             PrintWriter pw   new PrintWriter(outputFile);
```

```java
125:
 26:            if   x       w
127:            {
 28:                pw.print n "% This is part of a test suite.";
129:                pw.print ln("% Took all lines except " startPos " to "-
  3 :                       (startPos+segmentlength-1 ".";
131:                pw.println();
 32:            .

133:
 34:            pw.print ln("@relation myexample");

 35:
 36:            for(String atom : finalAtomList)
137:                pw.print ln("@attribute \""+atom+"\" { 0, 1 }",;
 38:
139:            String scoreString   "@attribute score { 0";
 40:            for(int i=1; i<=maxScore; i++)scoreString+= ", "+i;
141:            scoreString+= " }";
142:            pw.println(scoreString);

 43:
144:            pw.println();
145:            pw.print ln("@data");
146:
147:            int i = 0;
 48:            for(Set<String> atomSet : atomSets)
149:            {
 50:                if(!(i>=startPos&&i<startPos+segmentlength))
151:                {
152:                    for(String atom : finalAtomList)
 53:                    {
154:                        // Check to see if the atomSet has this atom.
 55:                        if(atomSet.contains(atom))pw.print("1");
156:                        else pw.print("0");
 57:                        pw.print(", ");
158:                    }
159:
 60:                    pw.println(scores.get(i));
161:                }
 62:                i++;
 6 :            }
164:
 65:            pw.close();
 66:            ps.println("Successfully generated file " outputFile);
 6 :
 68:        } catch(Exception e){
 69:            e.printStackTrace();
170:            throw new Error("Stopped");
    :        }
 72:    }
   :
174:
   :
```

```java
1: package pipeline;
2:
3: import java.io.BufferedReader;
4: import java.io.FileReader;
5: import java.io.PrintStream;
6: import java.io.PrintWriter;
7:
8: import server.ClientCAndC;
9:
10: /**
11:  * Takes a CSV file of English answers and converts
12:  * them into DRS semantic form using an external
13:  * parser.
14:  * @author Jamie
15:  *
16:  */
17: public class FileParser {
18:
19:     static String inputFile = "./files/8b-answers.csv";
20:     static String outputFile = "./files/8b-answers-processed.csv";
21:
22:     public static void main(String[] args) {
23:         processSentences(args[0], args[1], System.out);
24:     }
25:
26:     /**
27:      * Processes a single sentence.
28:      */
29:     public static String processSentence(String sentence)
30:     {
31:         ClientCAndC client = new ClientCAndC();
32:         try
33:         {
34:             client.connect();
35:             String drsString = client.getDRS(sentence);
36:             return drsString;
37:
38:         } catch(Exception e){
39:             throw new Error("Failed to connect to server.");
40:         }
41:     }
42:
43:     public static void processSentences(String inputFile,
44:         String outputFile, PrintStream ps)
45:     {
46:         try {
47:             ps.println("*** Processing sentences. ***");
48:
49:             ClientCAndC client = new ClientCAndC();
50:             client.connect();
51:
52:             BufferedReader br = new BufferedReader(new FileReader(inputFile));
53:             PrintWriter pw = new PrintWriter(outputFile);
54:
55:             while (br.ready()) {
56:                 ps.println("Reading line");
57:                 String line = br.readLine();
58:
59:                 int textPos = line.indexOf("\t") + 1;
60:                 String score = line.substring(0, textPos - 1);
61:
62:                 String sentence = line.substring(textPos);
63:                 if(!sentence.equals(""))
64:                 {
65:                     String drsString = client.getDRS(sentence);
66:                     if(!drsString.equals(""))
67:                     {
68:                         ps.println("!!! PARSED OK");
69:                         pw.println(score + "\t" +
70:                             drsString.replace("/"," or "));
71:                     } else {
72:                         ps.println("XXX FAILED");
73:                     }
74:                 }
75:             }
76:         }
77:
78:         ps.println("*** Finished processing sentences. ***");
79:
80:         br.close();
81:         pw.close();
82:
83:         client.close();
84:
85:     } catch (Exception e) {
86:         e.printStackTrace();
87:     }
88:
89:     }
90:
91: }
```

```java
1: package pipeline;
2: import java.io.BufferedReader;
3: import java.io.FileReader;
4: import java.util.HashMap;
5: import java.util.Map;
6: import java.util.Set;
7: import java.util.TreeSet;
8:
9: import parser.DRSreader;
10:
11:
12: /**
13:  * This identifies the most popular atoms in the dataset.
14:  * Used for testing purposes only, and is not used in the
15:  * overall pipeline.
16:  * @author Jamie
17:  *
18:  */
19: public class KeywordIdentifier {
20:
21:     static String inputFile = "./files/5aii-answers-processed2.csv";
22:     public static void main(String[] args)
23:     {
24:         if(args.length> 1)generateFile(args[0]);
25:         else generateFile(inputFile);
26:     }
27:
28:     public static void generateFile(String inputFile)
29:     {
30:         Map<String,Integer> atomCounts = new HashMap<String,Integer>();
31:
32:         try
33:         {
34:             BufferedReader br = new BufferedReader(new FileReader(inputFile));
35:             int numLines=0;
36:             while(br.ready())
37:             {
38:                 String line = br.readLine();
39:                 if(line.indexOf("%")! 0)
40:                 {
41:                     String[] lineParts = line.split("\t");
42:                     Set<String> set = DRSreader.atomize(lineParts[1]);
43:                     for(String atom : set)
44:                     {
45:                         if(atomCounts.get(atom) ==null)atomCounts.put(atom, 1);
46:                         else atomCounts.put(atom, atomCounts.get(atom)+1);
47:
48:                     }
49:                     numLines++;
50:                 }
51:
52:             br.close();
53:
54:             // want sorted list by the atoms' counts.
55:
56:             Set<Integer> counts = new TreeSet<Integer>(atomCounts.values());
57:                 for(Integer i : counts)
58:                     for(String val : atomCounts.keySet())
59:                         if(atomCounts.get(val).equals(i))
60:                             System.out.println(i": " + val);
61:
62:         }catch(Exception e){
63:                 e.printStackTrace();
64:             }
65:         }
66:     }
67: }
```

```
1:  package pipeline;
2:  import java.io.File;
3:  import java.io.FileOutputStream;
4:  import java.io.FileReader;
5:  import java.io.ObjectOutputStream;
6:  import java.io.PrintStream;
7:
8:  import schemagenerator.CumulativeClassifier;
9:  import weka.classifiers.Classifier;
10: import weka.classifiers.bayes.NaiveBayes;
11: import weka.classifiers.rules.DecisionTable;
12: import weka.core.Instances;
13:
14: /**
15:  * Takes an arff file to train a decision tree with, and an instance
16:  * to classify, initially in DRS form. Produces the number of marks
17:  * this instance has been awarded.
18:  * @author Jamie
19:  *
20:  */
21: public class Trainer
22: {
23:     public static void train(String trainingFile, String trainedModel,
24:         String learningType, PrintStream ps)
25:     {
26:         try
27:         {
28:             ps.println("** No trained model found. Creating new one.");
29:
30:             FileReader reader   new FileReader(trainingFile);
31:             Instances instances - new Instances(reader);
32:             instances.setClassIndex(instances.numAttributes()-1);
33:
34:             Classifier dt - new NaiveBayes(); // default
35:             if(learningType.charAt(0)--'d')dt = new DecisionTable();
36:             if(learningType.charAt(0)=='j'){
37:                 // Learning type is either 'j' (a standard mark scheme)
38:                 // or of the form 'j-k', where k is the total possible marks
39:                 // obtainable before it is capped the maximum mark.
40:                 // (e.g. "Any of the k following, maximum n")
41:
42:                 // We also need to see if there is a 'main.marklist' file
43:                 // in the same folder as the trainingFile. If so, use it.
44:                 String markListFile = getFolder(trainingFile)+
45:                         "/main.marklist";
46:                 if(!(new File(markListFile)).exists())markListFile   null;
47:
48:                 if(markListFile!=null)
49:                     ps.println("Using the manually created mark scheme.");
50:
51:                 if(learningType.length) - 1 dt = new CumulativeClassifier(
52:                     markListFile,
53:                     changeExtension(trainingFile,"markscheme"));
54:                 else dt = new CumulativeClassifier(
55:                     markListFile,
56:                     changeExtension(trainingFile,"markscheme"),
57:                     new Integer(learningType.split("-")[1]));
58:             }
59:             if(learningType.charAt(0)  'n')dt = new NaiveBayes();
60:
61:             dt.buildClassifier(instances);
62:
63:                     ps.println("########################");
64:                     ps.println("EXECUTING TRAINING");
65:                     ps.println("Using learning type: "+learningType);
66:                     ps.println("** Saving trained model to '"+trainedModel+"'.");
67:
68:                     // Write trained model to file
69:                     ObjectOutputStream oos = new ObjectOutputStream(
70:                         new FileOutputStream(trainedModel));
71:                     oos.writeObject(dt);
72:                     oos.flush();
73:                     oos.close();
74:
75:             } catch (Exception e )
76:                 e.printStackTrace();
77:             }
78:         }
79:
80:     /**
81:      * Changes the extension of a filename.
82:      */
83:     private static String changeExtension(String fileName, String newExt)
84:     {
85:         int pos = fileName.lastIndexOf(".");
86:         return fileName.substring(0,pos) + "." + newExt;
87:     }
88:
89:     private static String getFolder(String fileName)
90:     {
91:         int pos = fileName.lastIndexOf("\\");
92:         pos = Math.max(pos, fileName.indexOf("/"));
93:         return fileName.substring(0,pos);
94:     }
95:
96:     static final String trainingFile = "./files/testing/5aii-0.arff";
97:     static final String trainedModel = "./files/trained/5aii-0.model";
98:     public static void main(String[] args) {
99:         train(trainingFile,trainedModel,"d",System.out);
100:    }
101: }
```

```java
 1: package pipeline;
 2:
 3: import java.io.BufferedReader;
 4: import java.io.FileReader;
 5: import java.io.InputStreamReader;
 6: import java.io.PrintWriter;
 7: import java.net.URL;
 8: import java.net.URLEncoder;
 9:
10: /**
11:  * An alternative way of parsing English sentences
12:  * which uses an online version of the C&C parser
13:  * and boxer. This is acceptable for single parses,
14:  * but it would be considered antisocial behaviour
15:  * to parse and entire data set.
16:  * @author Jamie
17:  *
18:  */
19: public class WebsiteParser
20: {
21:     // URL to use.
22:     final static String url = "http://svn.ask.it.usyd.edu.au/"+
23:         "demo/demo3.cgi?sentence=";
24:     final static String url2 = "&printer=prolog&model=a&resolve=yes"+
25:         "&submit=Parse+and+Box%21";
26:
27:     // Used in the main method for parsing an entire file.
28:     static String inputFile = "./files/8b-answers.csv";
29:     static String outputFile = "./files/8b-answers-processed.csv";
30:     final static int exampleLimit = 1000; // do all of them
31:     final static long delay = 20000; // 20seconds
32:
33:     public static String getDrsFromWeb(String sentence)
34:     {
35:         try
36:         {
37:             String fullUrl = url+URLEncoder.encode(sentence)+url2;
38:             System.out.println("Reading from: "+fullUrl);
39:             URL webUrl = new URL(fullUrl);
40:             BufferedReader urlReader = new BufferedReader(
41:                 new InputStreamReader(webUrl.openStream()));
42:
43:             int i = 0;
44:             String drsString = "";
45:             while(urlReader.ready()&&i<=13)
46:             {
47:                 String line = urlReader.readLine();
48:                 if(i == 13)
49:                 {
50:                     int lastBracketPos = line.lastIndexOf(')');
51:                     drsString = line.substring(0,lastBracketPos).trim();
52:                 }
53:                 i++;
54:             }
55:
56:             urlReader.close();
57:
58:             if(drsString.equals(""))throw new Error("Problem obtaining "+
59:                 "DRS string.");
60:             System.out.println("Obtained the DRS: "+drsString);
61:
62:             return drsString;
63:         } catch(Exception e){
64:             e.printStackTrace();
65:             return null;
66:         }
67:
68:     }
69:
70:
71:     /**
72:      * @deprecated
73:      * Use of the online parser for an entire data set
74:      * is strongly discouraged. Use FileParser instead.
75:      */
76:     public static void main(String[] args) {
77:
78:         try
79:         {
80:             BufferedReader br = new BufferedReader(new FileReader(inputFile));
81:             PrintWriter pw = new PrintWriter(outputFile);
82:             int exampleCount = 0;
83:
84:             while(br.ready()&&exampleCount<exampleLimit)
85:             {
86:                 String line = br.readLine();
87:
88:                 int textPos = line.indexOf("\t")+1;
89:                 String score = line.substring(0,textPos-1);
90:
91:                 String drsString = getDrsFromWeb(line.substring(textPos));
92:                 System.out.println("Obtained DRS: "+drsString);
93:                 pw.println(score+"\t"+drsString);
94:
95:                 exampleCount++;
96:
97:                 System.out.println("** Sleeping for 20 seconds");
98:                 Thread.sleep(delay);
99:             }
100:
101:             br.close();
102:             pw.close();
103:
104:         } catch(Exception e){
105:             e.printStackTrace();
106:         }
107:     }
108: }
```

```
: package    ;

: import    ...
-: import    ...
-: import    ...
: import    ...
: import    ...
-: import    ...
9: import    ...
: import    ... Map;
: import    ...;

: import    ...
: import    javax.swing...;
: import    javax...;
: import    javax.swing.JTree ;
: import    ...;
: import    javax.swing...;
: import    javax.swing.JTextArea;
: import    javax.swing.JScrollPane ;
2': import    javax.swing.border.Border;
: import    javax...;
23:
2: import    ...;
25:
26:  **
27:  * The panel ...
28:  * of the system.
29:  * @author ...
30:  *
31:  *
32: public class EvaluatorPane  extends  JPane  implements PropertyChange listener

: final String INSTRUCTION  = "Create a new folder, containing either "
:        "of the following:\n" 
:        "* answers.csv - a file of the source answers in English.\n"
:        "* answers-processed.csv - a file of the parsed sentences.\n"
:        "   This skips the parsing stage.";

:      // Used to map learning type selection
:      // to radio registration ...
: final Map<Integer,String> typeMap = new HashMap<Integer,String> ;
: final ...
44:      // The choice of learning types.
: ...
46:  ...  private  ...;
: ...
: final ... testChooser = new ...this,
:        "Test directory","DIRECTORY" ;

: public  EvaluatorPanel()

:      super();
:      ... new DirChooser... ;

:      ...  new  JPanel new  ... ;
:      ... = new JPanel ;
:      // TextArea instructions  = new JTextArea,30;
:           ...
: startButton  new JButton "Start" ;
: ...  new  ...;
60: ...
```

```
:        ...
:        ... , "d" ;
:        ... , "n" ;
:        ... , "j" ;
:        ...          "Decision Tree",
:               "Naive Bayesian Network", "Cumulative Classifier"  ;
: ...  new  JComboBox ... ;
: ...
: ...

: ... ;

: ...  new  ... ;
: ...

: JPanel statusPanel  = new JPanel  ... "Status" ;
: JPanel status... =   new  JPanel ;
: statusPanel...

: statusArea   new JTextArea(4,  ;
: statusArea...
: JScrollPane scrollPane   new JScrollPane(statusArea);
: statusArea...Editable  false ;
: statusListe.add scrollPane ;

: JPanel buttonPanel   new JPanel ;
: buttonPanel.add testFolder ;

: add(... Panel,BorderLayout.NORTH);
: ... statusPanel ,BorderLayout.CENTER ;
: all(buttonPanel, BorderLayout.SOUTH);

public void propertyChange(PropertyChangeEvent ... ) 
    if ("progress" == ... get ...
        int progress = (Integer) evt.getNewValue();
        ... progressbar setValue to progress;

    public class Evaluator...  implements ActionListener

        public void actionPerformed(ActionEvent ...

            // Erase textarea
            statusArea.setText "" ;

            // Get learning type
            final ...type ... style   listServer
                   learningList.getSelectedIndex();

            final String folderName   testFolder.getFileName   ;
            if (folderName == null ...  ... .setText
                   "ERROR: No folder selected" ;
            else
                start... .setEnabled false ;
                ...
                Task task   new Task Evaluator(...,  ... type) ;
                ... Evaluator  .this ;
```

```
public class .... extends ...

        public .... 
            this.... = ...
            this.learningType    learningType

        public .... doInBackground

                if learningType.equals "j"

                        String s
                                (String)JOptionPane....with j...    s
                                "OPTIONAL: If the mark scheme is of"-
                                " the form 'Any of the k following,"-
                                " maximum n', enter k.";
                            if(s!=""&&s!=null)learningType    "j-"...
                        }

                Tester.test(folderName, learningType,
                        new GUIOutput(statusArea), ....
                return null;
            }

        public void done()
                ...kit.getDefaultToolkit().beep();
                startButton.setEnabled(true);
                setCursor null ;   turn off the wait ...
```

```java
: package    ;

: import                   ;
: import                   ;
: import                   ;
: import                   ;
: import                   ;
: import                   ;
: import                   ;
: import                   ;
: import                   ;
: import                   ;
: import                   ;
: import                   ;
: import                   ;
: import                   ;
: import                   ;

: import                   ;
: import                   ;
22: import                  ;
: import                   ;
24: import                  ;
25: import                  ;
26: import                  ;
27: import                  ;
28: import                  ;
29: import                  ;

   /**
    * A square panel that represents a file or directory
    * graphically. The user can specify a file by one of
34: * two methods:
    *   . click on a panel and select the file interactively
    *     in the normal fashion.
    *   . drag a file icon from into the panel
    *
    *
41: public class FileDragPanel extends      Panel implements DropTargetListener {

        static String             = "./icons/document.gif";
44:     static String             = "./icons/document_greyed.gif";

        final static               = new               ;

                    ;
            JPanel filePanel   ;
                   filePanel ;
            JButton fileType;
                   type;
                   label    true;
                   file;
                        panel;

        public                   return        ;

        public                                        {

          super  ;
          this.               ;
```

```java
:           this.               ;
:                new          this,               ;

:                      new       ( "" ;
:                                            ;
:                      new Font("Tahoma",         ;
:                                        ;

:    = new          new          ;

:             new         new                  ;
:                      new Font       Type ;
:                                        ;
:                                        ;
:       add     new          new                  ;
:                     ;
:          d ;

:              = "type"        new       ("type ;
: Font      2 = new Font "Tahoma",      Bold, 4 ;
:                                ;

        // set up our text area on the edge dropp...
        // this class will handle drop events
dt = new DropTarget(this, this );


public void dragEnter  (DropTargetDragEvent dtde,
public void dragExit   (DropTargetEvent dte)  {
public void dragOver   (DropTargetDragEvent dtde,  {
public void dropActionChanged (DropTargetDragEvent dtde)

public void drop (DropTargetDropEvent  dtde
    try
        // Get the dropped object and try to figure out what it is
        Transferable tr = dtde.getTransferable();
        DataFlavor[] flavors = tr.getTransferDataFlavors();
        for  int      = 0;           getLength;
            // Check for file lists specifically
            if (flavors[i].isFlavorJavaFileListType()
                    // Great! Accept copy drops
                    dtde.acceptDrop(DnDConstants.ACTION_COPY_OR_MOVE);
                    setMessage("Successful file list drop.");

                    javaUtilList list
                                                            list  ;
                    dataToPanel mkdirList.get(0),"" ;

                    // If we made it this far, everything worked.
                    dtde.dropComplete true ;
                    return;

                // Is it another Java object?
            else if  flavors[i].isFlavorSerialized ==dObjectType
                    dtde.acceptDrop DnDConstants.ACTION_COPY_OR_MOVE ;
                        setMessage "Copy:     text dropped: " + ;
                    // Object o = tr.getTransferData(dataFlavors[i]);
                    setMessage "          "       ;
                    dtde.dropComplete true;
                    return;

                // How about an input stream.
```

```
        else if                                              ;

                                    true ;
            return;

                                 dragger
         ;
            catch
                            ;
                           ;


    public void

            this.

    private static

            int    = Math.max fileName.lastIndexOf("\\" ,
                     fileName.lastIndexOf("/"  ;
            String fileNameOnly = fileName.substring       ;
            if(fileNameOnly.length()<=17) return fileNameOnly;
            // Otherwise we need to shorten it.
            // First find dot
            int p     = fileNameOnly.lastIndexOf(".");
            // Remove middle.
            return           fileNameOnly.substring(0, 8)+"..."+
                     fileNameOnly.substring(p-5  ;


    public class              implements

                            new                ;
            final              ;
            public         final

                    this.       = type;
                            = new
                    if           ("DIRECTORY"
                             //   choose DIRECTORIES ONLY ;
                                      ;


    public void              MouseEvent

                    // Handle open button action.
            int                                  ;

            if                                          .
                                               File ;
                                                   ;
            else
```

```
    public void
    public void
    public void
    public void

public class              extends

        String       ;
        public My
                this.

        public
            return                "DIRECTORY"   . "Directory" :
                 "         ;

        public boolean

                    if                   return true;

                                                    ;
                    if extens     !=null

                            if                        return true;
                            else return false;

                    return false;

        public String getExtension         f) {
                    String    ext =   null;
                    String s =
                    int    = s.lastIndexOf('.');

                    if
                        ext

                    return ext;
```

```
: package  ;
:
: import        ;
:
: import        ;
:
:
:
:
:
:
: public class    extends
:
:
:                   ;
:         public
:                   super         ; this...
:
:         public void                        .append "\n" ;
:
:         public void    (String                        .append     ;
:
:         public void    int                            .append    -"" ;
:
:         public void                                    .append "\n" ;
:
:
```

```java
: package    ;
:
: import    ...          .    .   .   .   ;
: import    ..       .      .    .   ;
: import    ...  .        .       .  .  ..;
: import    ...       .    .  .  .  .  ;
:
:   ..
:   .        .      : :.
:   .    .   .       ..
:   .
:   .
: public class  ...  .  :
    : public static void  ..    .   .  .   ..    ..
    :             ..  .   ..    new   ...  - "Marker" ;
    :             .. . ...    . ..  .   ,  .  ;
    :             ..  .. .. .  .   true ;
    :
    :         try
    :             ..Ya...  ..  .   ..  .F..
    :                     ..Manager. .. . . .... .. .. .. . ... ..  ..  ;
    :           catch  .Except.. ..   .  .
    :
24:
25:         ..  ..    new  .. ..  ;
26:         ...  ...  "Pipeline",
27:             new  .     ("./icons/palmpilot.gif"),
28:             new  ..    .. ...  ;
29:         ...  ... "Evaluator",
30:             new  .   ..  ("./icons/cog_small.gif"),
31:             new  Eva.. .. .. ..  ;
32:         ...  ... "Marker",
33:             new  ...  ("./icons/construct.gif"),
34:             new  Marker.. ..  ;
35:
36:         ..   .. . .. .. ..  .add .. ;
37:
38:         .. ..  .. ..  ;
39:
40:
```

```java
package    ;

import          ;
import          ;
import          ;
import          ;
import          ;
import          ;
import          ;
import          ;
import          ;
import          ;
import          ;
import          ;
import          ;

import javax.swing.       ;
import javax.swing.       ;
import javax.swing.       ;
import javax.swing.       ;
import                    ;
import javax.swing.border.         dBorder;

import parser.ORBReader;
import pipeline.WebsiteParser;
import weka.class       e     e        ;
import weka.core.Instances;
import weka.       .         ;

/**
 * A panel           one by one and
 *                classifications.
 * @author     Jones
 *
 */
public class Marker    . extends

                             = new FileOpenPanel this,"Arff","arff" ;
                             = new FileOpenPanel this,"Model","model" ;

         public Marker    ()

                       new B       ay       ;

                       new       new FlowLayout    ;
                       new          (    ),2/    /    ;

                       new        new F  wLay       ;
                       new JLabel "Sentence:   " ;
          new      new Te   a        ;
                                      ;
                       new       "Submit" ;
                             new                          ;
                                  new             step     ;
                                    s            ;
                       new             = "Input" ;
                                          ;
```

```java
                                                         ;

                          new      new             ;
                                                                    ;

                                                       ;

                                       new      new              ;
                               new       " " ;
                       new       "Tahoma",              ;

                                                        ;

                       new              "Mark" ;

                                                        ;
}

     //
public class Submit Action         implements A         proces

     public void actionPerformed A          e

          // arff and model             ored i  d.
          if (a           .getFileName     = null return;
          if (mo   FileL gn.FileName    =null return;

               String t  Classify = FileMarker           s  Sentence
                          answer  .getText());
               String     Classify = WebsiteParser.   D   s  rWeb
                          answer.getText());
               Set<String>  set = ORBReader.           .         ;
               List<Double> cols = new LinkedList Double ;

          try

               B           e   new File        r
                    new FileProfile   tu this.e.getFileName   ;
               System. .printn("** Reading ARFF.");

               boolean reading   true;
               int i=0;
               while(   .ready)//          )

                         St          =        i  e ;
                         if                 ("@attribute")
                              reading   false;
                         else if l          x("@attribute")
                              String    parts = line.sp    "  ";
                              if(parts     equals("score",
                                   if                parts       ep      "\"",""
                                             parts 01    ;
                                   else                  ;

                         i  ;

                         line     ;

                    }//      ( )      eof    ;
```

```
        int

        double         = new double            ;
        for int    ,                   ;



                    new

                            "** Found trained model '".
                                          "'." ;


        new                         new                  ;
                                                          ;



                            new                          ;
                            new Fil-Pa                         ;
                            new                       ;




        double                                              ;
                              "Classified as: "      ;

                       .settext .new                  ""  ;



} catch              x){
            .printStackTrace ;
}
```

```
: package   . ;

:: import
: import
:: import
: import
:: import
:: import

: import
: import
: import
:: import
: import
: import
:: import

:: import
:: import
.: import

     * this panel allows the user to perform each
     * stage of the parsing/training pipeline, and
     * converts between the 4 file formats used:
     * 1. Source CSV   - English sentences.
     * 2. Semantic CSV - parsed sentences.
     * 3. Arff   - classifier input file.
     * 4. Model   A trained classifier.
     * @author Jamie
     *
: public class  PipelinePanel  extends  Panel {

          FileDragPanel              = new FileDragPanel(this,"Source CSV","csv" ;
                                  = new FileDragPanel(this,"Semantic CSV","csv
";
                            new FileDragPanel(this,"Arff","arff" ;
                            = new FileDragPanel(this,"Model","model" ;

          public PipelinePanel(

               setLayout( new  BorderLayout(null) ;

               Border                = new           ("Status" ;
                                   = new  JPanel  ;
                                   = new

                                   = new          ;
                                   new          ( new
                                                    false ;

                                  = ,         ,    .CENTER ;

                                  new    Panel ;
                                   = new   Panel ;
                            .Y new
```

```
                                              = new        ">" ;
                                                          new

                                                     ;
                                                   ;

                                   =  new        ">" ;
                                               = new

                                              = new       ">" ;
                                                = new

                            =,BorderLayout.NORTH ;

          public class             implements Action

               public void
                         String
                         String
                         if(           =null              = null
                                          "ERROR: Files not specified
.";
                               return;

                                          new                    ;

          public class            implements ActionListener

               public void actionPerformed(ActionEvent
                         String sourceFileName
                         String outputFileName   arff          ;
                         if(outputFileName= null            =null
                                          Text "ERROR: Files not specified
.";
                               return;

                                 new
                                        , false, new

          public class              implements Act
                                new
                  public =
                               ("Decision Tree", "d" ;
                               ("Naive Bayesian Network", "n" ;
                               ("Cumulative Classifier", "j" ;

                  public void
```

```
:          . . . .           . .-Na . .- .- . .,  .- = .- .  . . . ;
:          . . . .           . .. .  . . .- .- .-  .  .- .  .- .  . ; .- .  .   ;
:          if ,              . . . ===null    . .     . -. . ===null
                   . . . .- .- . .  . . "ERROR: Files not specified
." ;
:                            .
:                            return;
:                            .
:
:          . . . .           . . .- . .- .  . . . . .- .,.-. . . . . .- .- .  .
y. ;
.- . :          . .- . .    . .- .- . .- . . .      . .- . . .- .  . . . .
. . :                              . . .       . . .-  .this,
:                                  "Please choose a learning type",
:                                  "Learning Type".
. . :                              . .      . .-.    . . .-. .    .
. . :                              null,
. . . :                            . .- . .   .  . . es,
. . . :                            "Naive Bayesian Network" ;
. . . :                                                                   |
. . :          if(.=. null) . . . . . . . . . . . . . . . . . . .  . . . |
.=, . .- . . .- . . .- .- . .   "", new . . . . . . . . . . . . . . . . ;
. . . :
:
```