



SAPIENZA
UNIVERSITÀ DI ROMA

Non-Sentential Utterances in Dialogue: Experiments in Classification and Interpretation

Department of Computer, Control and Management Engineering
Master of Science in Engineering in Computer Science

Candidate

Paolo Dragone
1370640

Thesis Advisor

Prof. Roberto Navigli

Co-Advisor

Dr. Pierre Lison
(University of Oslo)

Academic Year 2014/2015

Non-Sentential Utterances in Dialogue: Experiments in Classification and Interpretation

Master thesis. Sapienza – University of Rome

© 2015 Paolo Dragone. All rights reserved

This thesis has been typeset by L^AT_EX and the Sapthesis class.

Version: October 14, 2015

Author's email: dragone.paolo@gmail.com

Abstract

Non-sentential utterances (NSUs) are utterances that lack a complete sentential form but whose meaning can be inferred from the dialogue context, such as “OK”, “where?”, “probably at his apartment”. The interpretation of non-sentential utterances is an important problem in computational linguistics since they constitute a frequent phenomena in dialogue and they are intrinsically context-dependent. The interpretation of NSUs is the task of retrieving their full semantic content from their form and the dialogue context.

NSUs also come in a wide variety of forms and functions and classifying them in the right category is a prerequisite to their interpretation. The first half of this thesis is devoted to the NSU classification task. Our work builds upon Fernández et al. (2007) which present a series of machine-learning experiments on the classification of NSUs. We extended their approach with a combination of new features and semi-supervised learning techniques. The empirical results presented in this thesis show a modest but significant improvement over the state-of-the-art classification performance.

The consecutive, yet independent, problem is how to infer an appropriate semantic representation of such NSUs on the basis of the dialogue context. Fernández (2006) formalizes this task in terms of “resolution rules” built on top of the Type Theory with Records (TTR), a theoretical framework for dialogue context modeling (Ginzburg, 2012). We argue that logic-based formalisms, such as TTR, have a number of shortcomings when dealing with conversational data, which often include partially observable knowledge and non-deterministic phenomena. An alternative to address these issues is to rely on probabilistic modeling of the dialogue context. Our work is focused on the reimplementation of the resolution rules from Fernández (2006) with a probabilistic account of the dialogue state. The probabilistic rules formalism (Lison, 2014) is particularly suited for this task because, similarly to the framework developed by Ginzburg (2012) and Fernández (2006), it involves the specification of *update rules* on the variables of the dialogue state to capture the dynamics of the conversation. However, the probabilistic rules can also encode probabilistic knowledge, thereby providing a principled account of ambiguities in the NSU resolution process. In the second part of this thesis, we present our *proof-of-concept* framework for NSU resolution using probabilistic rules.

Acknowledgments

The past two years have been utterly intense, but beautiful nonetheless. In all my wandering around the Earth, first in Melbourne then in Oslo, I have learned an incredible amount of things and met as many interesting people. I have studied hard, explored the world, suffered homesickness and experienced many wonderful moments. It has been an exiting time and now, at the finishing line, I am walking towards the end with a smile on my face. For those opportunities I am grateful to my alma mater Sapienza which, despite all the difficulties, gave me also many positive memories.

This thesis is the product of the period I spent at University of Oslo under the supervision of Dr. Pierre Lison. Pierre has not only been a great supervisor but he has also been my mentor. He has taught me an uncountable number of things. He introduced me to the research on dialogue and proposed the topic of non-sentential utterances in the first place. My collaboration with Pierre has also led to my very first scientific publication at the SemDial 2015 in Göteborg. I can safely say that writing this thesis would not have been possible without his guidance. For these and many other reasons I want to thank him and wish him all the best. I also want to thank all the Language Technology Group at University of Oslo for their hospitality and for the many (or perhaps too few) happy days we spent together.

A thankful note to all my friends too, who support me everyday of my life, even (and especially) when I am not around. In particular, a big thank you to all my friends and colleagues from the course in Engineering in Computer Science with whom I walked side by side for all these years. However, a special thanks goes to my estimate colleague Alessandro Ronca who has practically been my moral supporter as well as my personal reviewer.

Last but not least, I want to express my gratitude to my parents and all my family who have always encouraged me to go beyond my own limits (even if sometimes I went perhaps too far away) and never stopped believing in me.

Rome, 22th of October 2015

Paolo Dragone

Contents

1	Introduction	1
1.1	Motivation	2
1.2	Contribution	2
1.3	Outline	3
2	Background	5
2.1	Non-Sentential Utterances	5
2.1.1	A taxonomy of NSUs	6
2.1.2	The NSU corpus	10
2.1.3	Interpretation of NSUs	11
2.2	A formal model of dialogue	11
2.2.1	Type Theory with Records	11
2.2.2	The dialogue context	14
2.2.3	Update rules	15
2.3	Probabilistic modeling of dialogue	17
2.3.1	Bayesian Networks	17
2.3.2	Probabilistic rules	18
2.4	Summary	20
3	Classification of Non-Sentential Utterances	21
3.1	The data	21
3.2	Machine learning algorithms	23
3.2.1	Classification: Decision Trees	23
3.2.2	Classification: Support Vector Machines	24
3.2.3	Optimization: Coordinate ascent	24
3.3	The baseline feature set	26
3.4	Feature engineering	28
3.5	Semi-Supervised Learning	32
3.5.1	Unlabeled data extraction	32
3.5.2	Semi-supervised learning techniques	33
3.6	Evaluation	36
3.6.1	Metrics	36
3.6.2	Empirical results	38
3.7	Summary	45

4	Resolution of Non-Sentential Utterances	47
4.1	The resolution task	48
4.2	Theoretical foundation	49
4.2.1	Partial Parallelism	49
4.2.2	Propositional lexemes	50
4.2.3	Focus Establishing Constituents	50
4.2.4	Understanding and acceptance	51
4.2.5	Sluicing	52
4.3	Dialogue context design	52
4.3.1	Semantics	52
4.3.2	Dialogue acts	53
4.3.3	Variables of the dialogue context	53
4.4	NSU resolution rules	55
4.4.1	Acknowledgments	55
4.4.2	Affirmative Answers	56
4.4.3	Rejections	57
4.4.4	Propositional Modifiers	58
4.4.5	Check Questions	60
4.4.6	Short Answers	61
4.4.7	Sluices	61
4.4.8	Clarification Ellipsis	63
4.5	Implementation and use case example	64
4.5.1	Dialogue system architecture	65
4.5.2	Example	66
4.6	Summary	71
5	Conclusion	73
5.1	Contributions	73
5.2	Future developments	75
	Appendices	79
A	Context update rules	79

List of acronyms

AL	Active Learning
BNC	British National Corpus
DGB	Dialogue Gameboard
FEC	Focus-Establishing Constituent
MaxQUD	Maximal QUD element
NLG	Natural Language Generation
NLU	Natural Language Understanding
NSU	Non-Sentential Utterance
ParPar	Partial Parallelism
QUD	Question Under Discussion
SMO	Sequential Minimal Optimization
SVM	Support Vector Machine
TSVM	Transductive Support Vector Machine
TTR	Type Theory with Records

Chapter 1

Introduction

In dialogue, utterances do not always take the form of complete sentences. Utterances may sometimes lack some constituents – subject, verb or complements – because they can be understood from the previous utterances or other contextual information. These fragmentary utterances are often called *non-sentential utterances* (NSUs). The following are some examples from the British National Corpus:

- (1.1) A: How do you actually feel about that?
B: **Not too happy.**
[BNC: JK8 168–169]¹
- (1.2) A: They wouldn't do it, no.
B: **Why?**
[BNC: H5H 202–203]
- (1.3) A: So will the tape last for the whole two hours?
B: **Yes, apparently.**
[BNC: J9A 76–77]
- (1.4) A: Right disk number four?
B: **Three.**
[BNC: HDH 377–378]

We can understand without effort the meaning of the NSUs in the short dialogues above, even though they do not have the form of full sentences. We can easily make sense of them by extrapolating their meaning from the surrounding context, which for the above examples is given by the preceding utterance. Other possible contextual factors that affect the intended meaning of the NSUs are, for instance, the history of the dialogue, the shared environment of the conversational participants, their common knowledge and so on. From a computational linguistic perspective, making sense of this kind of utterances is a difficult problem because it involves the formalization of a robust theory of dialogue context.

¹This notation indicates the file name and the line numbers of the portion of dialogues in the British National Corpus.

Moreover, NSUs are a large variety of phenomena that need to be treated in different ways. Fernández and Ginzburg (2002) identify 15 different types of NSUs. One of the problems that must be addressed to make sense of NSUs is determining their type. One possible way is to classify NSUs using machine learning, as previously experimented by Fernández et al. (2007).

To interpret a given NSU, one also has to *resolve* its meaning i.e. construct an high-level semantic representation of the NSU by extracting the relevant information from the dialogue context. To select the right resolution procedure for the given NSU, one needs first to determine its type. That is why the two task are connected. However, they can still be formalized and employed independently.

1.1 Motivation

Non-sentential utterances are interesting in many ways. First of all, they are very frequent in dialogue. According to Fernández and Ginzburg (2002) and related works, the frequency of NSUs in the dialogue transcripts of the British National Corpus is about 10% of the total number of utterances. However, this number may vary greatly if one takes into account a larger variety of phenomena or different dialogue domains e.g. Schlangen (2003) estimates the frequency of NSUs to be 20% of the total number of utterances.

Despite their ubiquity, the semantic content of NSUs is often difficult to extract automatically. Non-sentential utterances are indeed intrinsically dependent on the dialogue context. It is impossible to make sense of them without accessing to the surrounding context. Their high context-dependency makes their interpretation a difficult problem from both a theoretical and computational point of view.

NSUs form a wide range of linguistic phenomena that need to be considered in the formulation of a theory of dialogue context. Only few previous works tackled this problem directly and the majority of them take place in theoretical semantics of dialogue without addressing the possible applications. This means that the interpretation of NSUs is still an understudied problem, making them possibly an even more interesting subject.

1.2 Contribution

Our work follows two parallel paths. On one hand we address the problem of the classification of NSUs by extending the work of Fernández et al. (2007). On the other hand we propose a novel approach to the resolution of NSUs using *probabilistic rules* (Lison, 2015).

The classification task is needed to select the resolution procedure but it is nonetheless an independent problem and it can arise in many different situations. Our contribution to this problem is a small but significant improvement over the accuracy of the previous works as well as the exploration of one way to tackle the scarcity of labeled data.

Our work on the resolution of NSUs takes inspiration from Fernández (2006) and Ginzburg (2012) which provide the theoretical background for our study. Their framework is however purely logic-based therefore it can have some drawbacks in dealing with raw conversational data which often contains hidden or partially observable variables. To this end a probabilistic account of the dialogue state is preferable. In our work we implemented a new approach to NSU resolution based on the probabilistic rules formalism of Lison (2015). Probabilistic rules are similar, in some way, to the rules formalized by Ginzburg (2012), as both express updates on the dialogue state given a set of conditions. However, probabilistic rules can also take into account probabilistic knowledge, thereby making them more suited to deal with the uncertainty often associated with conversational data. Our work does not aim to provide a full theory of NSU resolution but rather be a *proof-of-concept* for the resolution of NSUs via the probabilistic rules formalism. Nevertheless we detail a large set of NSU resolution rules based on the probabilistic rules formalism and provide an actual implementation of a dialogue system for NSU resolution using the OpenDial toolkit (Lison and Kennington, 2015), which can be the baseline reference for future developments.

Our work for this thesis has produced the following publications:

- Paolo Dragone and Pierre Lison. Non-Sentential Utterances in Dialogue: Experiments in classification and interpretation. In: *Proceedings of the 19th workshop on the Semantics and Pragmatics of Dialogue, SEMDIAL 2015 – goDIAL*, p. 170. Göteborg, 2015.
- Paolo Dragone and Pierre Lison. An Active Learning Approach to the Classification of Non-Sentential Utterances. In: *Proceedings of the second Italian Conference on Computational Linguistics, CLiC-IT 2015*, in press. Trento, 2015.

1.3 Outline

Chapter 2

This chapter discusses the background knowledge needed for the development of the following chapters. In particular the chapter describes the concept of non-sentential utterance and the task of interpretation of NSUs with an emphasis on the previous works. Secondly the chapter contains an overview on the formal representation of the dialogue context from the theory of Ginzburg (2012). We discuss briefly the Type Theory with Records, the semantic representation of utterances and the update rules on the dialogue context. Finally, we introduce the probabilistic approach to the definition of the dialogue context from Lison (2014). We discuss the basics of Bayesian Networks (the dialogue context representation) and the probabilistic rules formalism.

Chapter 3

This chapter describes the task of the classification of non-sentential utterances. It provides details on our approach, starting from the replication of the work from Fernández et al. (2007) which we use as baseline. We then discuss the extended feature set we used and the semi-supervised learning techniques we employed in our experiments. Lastly we discuss the empirical results we obtained.

Chapter 4

This chapter describes the problem of resolving non-sentential utterances and our approach to address it through probabilistic rules. First we formalize the NSU resolution task and describe the theoretical notions needed to address it. We then describe our dialogue context design as a Bayesian network and our formulation for the resolution rules as probabilistic rules. In the end we describe our implementation and an extended example of its application to a real-world scenario.

Chapter 5

This is the conclusive chapter of this thesis which summarizes the work and describes possible future works.

Chapter 2

Background

2.1 Non-Sentential Utterances

From a linguistic perspective, Non-Sentential Utterances – also known as *fragments* – has been historically an umbrella term for many elliptical phenomena that often take place in dialogue. In order to give a definition of Non-Sentential Utterances ourselves, we shall start by quoting the definition given by Fernández (2006):

“In a broad sense, non-sentential utterances are utterances that do not have the form of a full sentence according to most traditional grammars, but that nevertheless convey a complete sentential meaning, usually a proposition or a question.”

This is indeed a very general definition, whereas a perhaps simpler approach is taken by Ginzburg (2012) which defines NSUs as “utterances without an overt predicate”. The minimal clausal structure of a sentence in English (as in many other languages) is composed of at least a noun phrase and a verb phrase. However, in dialogue the clausal structure is often truncated in favor of shorter sentences that can be understood by inferring their meaning from the surrounding context. We are interested in those utterances that, despite the lack of a complete clausal structure, convey a well-defined meaning given the dialogue context.

The context of an NSU can comprise any variable in the dialogue context but it usually suffice to consider only the *antecedent* of the NSU. The “antecedent” of an NSU is the utterance in the dialogue history that can be used to infer its underspecified semantic content. For instance, the NSU in (2.1) can be interpreted as “Paul went to his apartment” by extracting its semantic content from the antecedent. Generally, it is possible to understand the meaning of an NSU by looking at its antecedent.

- (2.1) A: Where did Paul go?
B: To his apartment.

It is often the case that an NSU and its antecedent present a certain grade of *parallelism*. Usually the meaning of an NSU is associated to a certain aspect of the antecedent. As described in Ginzburg (2012), the parallelism between an NSU and its antecedent can be of syntactic, semantic or phonological nature. The NSU in (2.1) presents syntactic parallelism – the use of “his” is syntactically constrained by the fact that Paul is a male individual – as well as semantic – the content of an NSU is a location as constrained by the *where* interrogative. This parallelism is one of the properties of NSUs that can be exploited in their interpretation (more details in Chapter 4). Even though it is often the case, the antecedent of an NSU is not always the preceding utterance, especially in multi-party dialogues.

2.1.1 A taxonomy of NSUs

As we briefly mentioned in Chapter 1, non-sentential utterances come in a large variety of forms. We can categorize NSUs on the basis of their form and their intended meaning. For instance NSUs can be affirmative or negative answers to polar questions, requests for clarification or corrections.

In order to classify the NSUs, we use a taxonomy defined by Fernández and Ginzburg (2002). This is a wide-coverage taxonomy resulting from a corpus study on a portion of the British National Corpus (Burnard, 2000). Table 2.1 contains a summary of the taxonomy with an additional categorization of the classes by their function, as defined by Fernández (2006) then refined by Ginzburg (2012).

Other taxonomies of NSUs are available from previous works by e.g. Schlangen (2003), but we opted for the one from Fernández and Ginzburg (2002) because it has been used in an extensive machine learning experiment by Fernández et al. (2007) and it is also used in the theory of Ginzburg (2012), which is our reference for the resolution part of our investigation. A detailed comparison of this taxonomy and other ones is given by Fernández (2006), which also details the corpus study on the BNC that led to the definition of this taxonomy.

Follows a brief description of all the classes with some examples. Fernández (2006) provides more details about the rationale of each class.

Plain Acknowledgment

Acknowledgments are used to signal understanding or acceptance of the preceding utterance, usually using words or sounds like *yeah*, *right*, *mhm*.

(2.2) A: I shall be getting a copy of this tape.

B: **Right.**

[BNC: J42 71–72]

Function	NSU class
<i>Positive Feedback</i>	Plain Acknowledgment Repeated Acknowledgment
<i>Metacommunicative queries</i>	Clarification Ellipsis Check Question Sluice Filler
<i>Answers</i>	Short Answer Affirmative Answer Rejection Repeated Affirmative Answer Helpful Rejection Propositional Modifier
<i>Extension Moves</i>	Factual Modifier Bare Modifier Phrase Conjunct fragment

Table 2.1: Overview of the classes in the taxonomy, further categorized by their function.

Repeated Acknowledgment

This is another type of acknowledgement that make use of repetition or reformulation of some constituent of the antecedent to show understanding.

- (2.3) A: Oh so if you press enter it'll come down one line.
 B: **Enter.**
 [BNC: G4K 102–103]

Clarification Ellipsis

These are NSUs that are used to request a clarification of some aspect of the antecedent that was not fully understood.

- (2.4) A: I would try F ten.
 B: **Just press F ten?**
 [BNC: G4K 72–73]

Check Question

Check Questions are used to request an explicit feedback of understanding or acceptance, usually uttered by the same speaker as the antecedent.

- (2.5) A: So (*pause*) I'm allowed to record you.
Okay?
B: Yes.
[BNC: KSR 5–6]

Sluice

Sluices are used for requesting additional information related to or underspecified into the antecedent.

- (2.6) A: They wouldn't do it, no.
B: **Why?**
[BNC: H5H 202–203]

Filler

These are fragments used to complete a previous unfinished utterance.

- (2.7) A: [...] would include satellites like erm
B: **Northallerton.**
[BNC: H5D 78–79]

Short Answer

The NSUs that are typically answers to *wh*-questions.

- (2.8) A: What's plus three times plus three?
B: **Nine.**
[BNC: J91 172–173]

Plain Affirmative Answer and Plain Rejection

A type of NSUs used to answer polar questions using *yes*-words and *no*-words.

- (2.9) A: Have you settled in?
B: **Yes, thank you.**
[BNC: JSN 36–37]
- (2.10) A: (*pause*) Right, are we ready?
B: **No, not yet.**
[BNC: JK8 137–138]

Repeated Affirmative Answer

NSUs used to give an affirmative answer by repeating or reformulating part of the query.

- (2.11) A: You were the first blind person to be employed in the County Council?
B: **In the County Council, yes.**
[BNC: HDM 19–20]

Helpful Rejection

Helpful Rejections are used to correct some piece of information from the antecedent.

- (2.12) A: Right disk number four?
B: **Three.**
[BNC: H61 10–11]

Propositional and Factual Modifiers

Used to add modal or attitudinal information to the previous utterance. They are usually expressed (respectively) by modal adverbs and exclamatory factual (or factive) adjectives.

- (2.13) A: Oh you could hear it?
B: **Occasionally yeah.**
[BNC: J8D 14–15]
- (2.14) A: You'd be there six o'clock gone mate.
B: **Wonderful.**
[BNC: J40 164–165]

Bare Modifier Phrase

Modifiers that behave like non-sentential adjunct modifying a contextual utterance.

- (2.15) A: [...] then across from there to there.
B: **From side to side.**
[BNC: HDH 377–378]

Conjunct

A Conjunct is a modifier that extends a previous utterance through a conjunction.

- (2.16) A: I'll write a letter to Chris
B: **And other people.**
[BNC: G4K 19–20]

NSU Class	Total	%
Plain Acknowledgment (Ack)	599	46.1
Short Answer (ShortAns)	188	14.5
Affirmative Answer (AffAns)	105	8.0
Repeated Acknowledgment (RepAck)	86	6.6
Clarification Ellipsis (CE)	82	6.3
Rejection (Reject)	49	3.7
Factual Modifier (FactMod)	27	2.0
Repeated Affirmative Answer (RepAffAns)	26	2.0
Helpful Rejection (HelpReject)	24	1.8
Check Question (CheckQu)	22	1.7
Sluice	21	1.6
Filler	18	1.4
Bare Modifier Phrase (BareModPh)	15	1.1
Propositional Modifier (PropMod)	11	0.8
Conjunct (Conj)	10	0.7
Total	1283	100.0

Table 2.2: The distribution of the classes in the NSU corpus.

2.1.2 The NSU corpus

The taxonomy presented in the previous section is the result of a corpus study on a portion of the dialogue transcripts in the British National Corpus, first started by Fernández and Ginzburg (2002), then refined by Fernández (2006). The dialogue transcripts used in the corpus study contain both two-party and multi-party conversations. The transcripts cover a wide variety of dialogue domains including free conversation, interviews, seminars and more. Fernández (2006) also describes the annotation procedure and a reliability test. The reliability test was carried out on a subset of the annotated instances comparing the manual annotation of three annotators. The test showed a good agreement between the annotators with a *kappa*-score of 0.76. From this test it is also clear that humans can reliably distinguish between the NSU classes in the taxonomy. Fernández (2006) provides more details about the complete analysis of the corpus.

In total about 14 000 sentences from 54 files were examined by the annotators, resulting in a corpus of 1 299 NSUs, about 9% of the total of the sentences examined. Of the extracted NSUs, 1 283 were successfully categorized according to the defined taxonomy making up a coverage of 98.9%. Table 2.2 shows the distribution of the classes in the corpus.

The annotated instances were also tagged with a reference to the antecedent of the NSU. About 87.5% of annotated NSUs have their immediately preceding utterance as antecedent. Fernández (2006) describes a study of the distance between NSUs and their antecedents, with a comparison between two-party and multi-party dialogues.

2.1.3 Interpretation of NSUs

Due to their incomplete form, non-sentential utterances do not have an exact meaning by themselves. They need to be “interpreted” i.e. their intended meaning must be inferred from the dialogue context. One way to interpret NSUs is developed by Fernández (2006), in turn based on Schlangen (2003), and it is formed by two consecutive steps, namely the *classification* and the *resolution* of the NSUs. The first step for the interpretation of an NSU is its classification i.e. finding its class according to the taxonomy described in Section 2.1.1. As demonstrated in Fernández et al. (2007), we can infer the class of an NSU using machine learning, i.e. we can train a classifier on the corpus detailed in Section 2.1.2 and use it to classify unseen NSU instances. The type of an NSU is used to determine the right *resolution* procedure to use. The resolution of an NSU is the task of recovering the full clausal meaning from their incomplete form on the basis of contextual information. Fernández (2006) describes a resolution procedure in terms of *rules* that, given some preconditions on the antecedent and other elements of the dialogue states, builds the semantic representation of the NSU. This approach to the resolution of NSUs has been the basis of several implementations of dialogue systems handling the resolution of NSUs such as Ginzburg et al. (2007) and Purver (2006).

Extending the interpretation problem to raw conversational data we need also a way to “detect” an NSU i.e. decide whether an utterance should be considered as an NSU in the first place. Since this is not our direct concern, we employ in our experiments a simple set of heuristics to distinguish between NSU and non-NSU utterances (see Section 3.5.1).

2.2 A formal model of dialogue

As theoretical base of our work we rely on the theory of dialogue context brought up by Ginzburg (2012), which presents a grammatical framework expressly developed for dialogue. The claim of Ginzburg (2012) is that the rules that encode the dynamics of the dialogue have to be built into the grammar itself. The grammatical framework is formulated using *Type Theory with Records* (Cooper, 2005). Type Theory with Records (TTR) is a logical formalism developed to cope with semantics of natural language. TTR is used to build a semantic ontology of abstract entities and events as well as to formalize the *dialogue gameboard* i.e. a formal representation of the dialogue context and its rules. The evolution of the conversation is formalized by means of *update rules* on the dialogue context. Ginzburg (2012) also accounts for NSUs and provides a set of dedicated rules.

2.2.1 Type Theory with Records

We will now briefly introduce the basic notions of the Type Theory with Records (TTR), with just enough detail needed by to understand the following sections, referring to Ginzburg (2012) for a complete description.

In TTR, objects can be of different types. The statement $x : T$ is a *typing judgment*, indicating that the object x is of type T . If x is of type T , x is said to be a *witness* of T . Types can either be *basic* (atomic) such as IND¹ or *complex* i.e. dependent on other objects or types such as $\text{drive}(x, y)$. Types also include constructs such as lists, sets and so on. Other useful constructs are records and record types. A record contains a set of assignments between labels and values whereas a record type contains a set of judgments between labels and types:

$$r : \begin{bmatrix} l_1 = v_1 \\ l_2 = v_2 \\ \dots \\ l_n = v_n \end{bmatrix} \qquad \rho : \begin{bmatrix} l_1 : T_1 \\ l_2 : T_2 \\ \dots \\ l_n : T_n \end{bmatrix}$$

The record r is of record type ρ if and only if $v_1 : T_1 \wedge v_2 : T_2 \wedge \dots \wedge v_n : T_n$. Typing judgment can be used to indicate the record r being of record type ρ as $r : \rho$.

TTR also provides *function types* of the form $T_1 \rightarrow T_2$ which maps records of type T_1 to records of type T_2 . Functional application is indicated as $(x : T_1).T_2$.

Utterance representation

At the basis of the grammatical framework of Ginzburg (2012) lies the notion of *proposition*. Propositions are entities used to represent facts, events and situations as well as to characterize the communicative process. In TTR propositions are records of the type:

$$\text{Prop} = \begin{bmatrix} \text{sit} & : & \text{Record} \\ \text{sit-type} & : & \text{RecType} \end{bmatrix}$$

A simple example of proposition may be the following:

Paul drives a car.

$$\begin{bmatrix} \text{sit} = r_1 \\ \text{sit-type} = \begin{bmatrix} x : \text{IND} \\ p_1 : \text{named}(x, \text{Paul}) \\ y : \text{IND} \\ p_2 : \text{car}(y) \\ c : \text{drive}(x, y) \end{bmatrix} \end{bmatrix}$$

¹The type IND stands for a generic “individual”.

On the other hand, questions are represented as propositional abstracts i.e. functions from the *question domain* Δ_q to propositions. Following the definition of Fernández (2006):

$$\text{Question} = \Delta_q \rightarrow \text{Prop}$$

The question domain Δ_q is a record type containing the *wh*-restrictors of the question q .² The *wh*-restrictors are record types that characterize the necessary information needed to resolve a *wh*-question e.g. for a *where* interrogative the answer must be a place instead for a *when* interrogative it must be a time. Clearly, the right *wh*-restrictor depends on the *wh*-interrogative used. Consider the following example of a *wh*-question:

Who drives?

$$\left(r : \begin{bmatrix} x & : & \text{IND} \\ \text{rest} & : & \text{person}(x) \end{bmatrix} \right) \cdot \begin{bmatrix} \text{sit} & = & r_1 \\ \text{sit-type} & = & [c : \text{drive}(r.x)] \end{bmatrix}$$

Here the question domain of the *who* interrogative is an individual x that is a person.

Polar questions, i.e. bare yes/no-questions, are represented as propositional abstract as the *wh*-questions, with the difference that their question domain is an empty record type. An example of polar question:

Does Paul drive?

$$\left(r : [] \right) \cdot \begin{bmatrix} \text{sit} & = & r_1 \\ \text{sit-type} & = & \begin{bmatrix} x & : & \text{IND} \\ p_1 & : & \text{named}(x, \text{Paul}) \\ c & : & \text{drive}(x) \end{bmatrix} \end{bmatrix}$$

A special type of propositions are used to represent the content of conversational moves which need to take into account the relation that stands between the speaker, the addressee and the content of the move. Those are called *illocutionary propositions* (of type *IllocProp*) and the relation that they contain is called *illocutionary relation*³. Illocutionary relations indicates the *function* of a proposition, such as “Assert”, “Ask”, “Greet”. For a proposition p , the illocutionary proposition that holds p as its content can be indicated as $R(\text{spkr}, \text{addr}, p)$, where R is the illocutionary relation, spkr and addr refer respectively to the speaker and the addressee⁴. Examples of illocutionary propositions are:

Assert($\text{spkr} : \text{IND}$, $\text{addr} : \text{IND}$, $p : \text{Prop}$)

Ask($\text{spkr} : \text{IND}$, $\text{addr} : \text{IND}$, $q : \text{Question}$)

²Ginzburg (2012) extends this field to be a list of record types to take into account situations with multiple question domains.

³Also called *illocutionary act* or *dialogue act*.

⁴For brevity only the semantic content of the illocutionary proposition is shown here.

2.2.2 The dialogue context

In Ginzburg (2012), the dialogue context – also known as the Dialogue Gameboard (DGB) – is a formal representation that describes the current state of the dialogue. It includes a wide range of variables needed to handle different aspects of the dialogue. However, we concentrated on the most basic ones:

- Facts, a set of known facts;
- LatestMove, the latest move made in the dialogue;
- QUD, a partially ordered set of questions under discussion.

The DGB can be represented in TTR as a record in the following way:

$$\left[\begin{array}{ll} \text{Facts} & : \text{Set(Prop)} \\ \text{LatestMove} & : \text{IllocProp} \\ \text{QUD} & : \text{poset(Question)} \end{array} \right]$$

The elements in the DGB represent the *common ground* of the conversation, shared between all the participants. In this representation we abstracted away several details that would be included in the actual DGB presented by Ginzburg (2012) such as the fields to track who is holding the turn, the current time and so on. We now detail the basic variables of the DGB.

Facts

Facts is a set of known facts, shared by all the conversational participants. The elements of Facts are propositions, which are assumed to be sufficient to encode the knowledge of the participants within the context of the dialogue. The Facts encode all the records that are *accepted* by all participants, i.e. facts that will not raise issues in the future development of the conversation. A complementary problem that we marginally address is the understanding – or *grounding* – of a sentence. Ginzburg (2012) develops a comprehensive theory of grounding but we do not include it in our work.

LatestMove

Dialogue utterances are made of coherent responses to the preceding utterances, that is why it is important to keep track of the history of the dialogue. In a two-party dialogue it is usually the case that the current utterance is a response to previous one, instead in a multi-party dialogue can be useful to keep track of a larger window of the dialogue history. Ginzburg (2012) keeps track of the history of the dialogue within the variable Moves while a reference to the latest (illocutionary) proposition is recorded in the field LatestMove.

QUD

QUD is a set of *questions under discussion*. In a general sense, a “question under discussion” represents an issue being raised in the conversation which drives the future discussion. Despite the name, QUDs may arise from both questions and propositions.

Ginzburg (2012) defines QUD as a partially ordered set (poset). Its ordering determines the priority of the issues to be resolved. Of particular importance is the first element in the set according to the defined ordering which is taken as the topic of discussion of the subsequent utterances until it is *resolved*. Such element is referred to as MaxQUD.

The formalization of the ordering is a rather complex matter in a generic theory of context that needs to account for the beliefs of the participants and it is especially problematic when dealing with multi-party dialogues. The usage of QUD is of particular importance in our case because the MaxQUD is used as the antecedent in the interpretation of NSUs.

2.2.3 Update rules

The dynamics of the DGB are defined by a set of *update rules* – also called *conversational rules* – which are applied on the DGB throughout the course of the conversation. Update rules are formalized as a set of effects on the parameters of the DGB given that certain preconditions hold. An update rule can be represented in the following way:

$$\left[\begin{array}{l} \text{pre} \quad : \quad [\dots] \\ \text{effects} \quad : \quad [\dots] \end{array} \right]$$

where both *pre* and *effects* are subsets of the parameters of the DGB and they respectively represent the necessary conditions for the application of the rule and the values of the involved variables right after the application of the rule.

Ginzburg (2012) defines all sorts of rules needed to handle a great variety of conversational protocols. Rules that are particularly interesting with respect to our work are those that handle queries and assertions as well as the ones that describe the dynamics of QUD and Facts.

The following rule describes how QUD is incremented when a question is posed:

$$\left[\begin{array}{l} \text{pre} \quad : \quad \left[\begin{array}{l} \text{q} \quad : \quad \text{Question} \\ \text{LatestMove} = \text{Ask}(\text{spkr}, \text{addr}, \text{q}) \quad : \quad \text{IllocProp} \end{array} \right] \\ \text{effects} \quad : \quad \left[\text{qud} = \langle \text{q}, \text{pre.qud} \rangle \quad : \quad \text{poset}(\text{Question}) \right] \end{array} \right]$$

As argued above, issues are also raised by assertions, as realized by the following rule:

$$\left[\begin{array}{l} \text{pre} \\ \text{effects} \end{array} : \left[\begin{array}{l} p : Prop \\ \text{LatestMove} = \text{Assert}(\text{spkr}, \text{addr}, p) : \text{IllocProp} \\ \text{qud} = \langle p?, \text{pre.qud} \rangle : \text{poset}(\text{Question}) \end{array} \right] \right]$$

The act of answering to a question is nothing else than asserting a proposition that resolves such a question. As a consequence the other speaker can either raise another issue related to the previous one or accept the fact that the issue has been resolved. The acceptance move is realized in the following way:

$$\left[\begin{array}{l} \text{pre} \\ \text{effects} \end{array} : \left[\begin{array}{l} p : Prop \\ \text{LatestMove} = \text{Assert}(\text{spkr}, \text{addr}, p) : \text{IllocProp} \\ \text{qud} = \langle p?, \text{pre.qud} \rangle : \text{poset}(\text{Question}) \\ \text{spkr} = \text{pre.addr} : \text{Ind} \\ \text{addr} = \text{pre.spkr} : \text{Ind} \\ \text{LatestMove} = \text{Accept}(\text{spkr}, \text{addr}, p) : \text{IllocProp} \end{array} \right] \right]$$

The speaker can also query the addressee with a *Check* move in order to ask for an explicit acknowledgment (*Confirm*) to a question-resolving assertion⁵. Acceptance and confirmation lead to an update of Facts and to a “downdate” of the QUD i.e. the removal of the resolved questions in QUD:

$$\left[\begin{array}{l} \text{pre} \\ \text{effects} \end{array} : \left[\begin{array}{l} p : Prop \\ \text{LatestMove} = \text{Accept}(\text{spkr}, \text{addr}, p) \vee \\ \quad \text{Confirm}(\text{spkr}, \text{addr}, p) : \text{IllocProp} \\ \text{qud} = \langle p?, \text{pre.qud} \rangle : \text{poset}(\text{Question}) \\ \text{facts} = \text{pre.facts} \cup \{ p \} : \text{Set}(\text{Prop}) \\ \text{qud} = \text{NonResolve}(\text{pre.qud}, \text{facts}) : \text{poset}(\text{Question}) \end{array} \right] \right]$$

While QUD represents the unresolved issues that have been introduced in the dialogue, Facts contains all the issues that have been resolved instead. That is why their update rules are closely related. The function *NonResolve* in the above rule checks for any resolved issues by the just updated facts and leave the unresolved ones into QUD.

⁵The rules for the *Check* and *Confirm* moves are omitted for brevity.

2.3 Probabilistic modeling of dialogue

In the previous section we detailed a logic-based model of dialogue from Ginzburg (2012). Another possible approach to dialogue modeling relies on probabilistic models to encode the variables and the dynamics of the dialogue context. Arguably this approach can be considered more robust to the intrinsic randomness present in dialogue. This is partially the reason why we explored this strategy as well as other advantages that will be discussed in Chapter 4.

We based our work on the *probabilistic rules* formalism developed by Lison (2012). This formalism is particularly suited for our purpose because of their commonalities with the update rules described in Section 2.2.3. The probabilistic rules formalism is based on the representation of the dialogue state as a Bayesian network. In this section we briefly describe how Bayesian networks are structured, then we detail the probabilistic rules formalism that we employ in Chapter 4 to model the resolution of the NSUs.

2.3.1 Bayesian Networks

Bayesian networks are probabilistic graphical models⁶ representing a set of random variables (nodes) and their conditional dependency relations (edges). A Bayesian network is a directed acyclic graph i.e. a direct graph that does not contain cycles (two random variables cannot be mutually dependent). Given the random variables X_1, \dots, X_n in a Bayesian network, we are interested in the *joint probability distribution* $P(X_1, \dots, X_n)$ of those variables. In general, the size of the joint distribution is exponential in the number n of variables therefore it is difficult to estimate when n grows. In the case of Bayesian networks we can exploit the *conditional independence* to reduce the complexity of the joint distribution. Given three random variables X , Y and Z , X and Y are said to be conditionally independent given Z if and only if (for all combinations of values) $P(X, Y|Z) = P(X|Z)P(Y|Z)$. We can define for a variable X_i in X_1, \dots, X_n the set *parents*(X_i) such that if there is a direct edge from X_j to X_i then $X_j \in \text{parents}(X_i)$. Given a *topological ordering*⁷ of the variables (nodes) of the Bayesian network, a variable X_i is conditionally independent from all its predecessor that are not in *parents*(X_i) therefore the joint probability distribution can be defined as follows:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{parents}(X_i))$$

For each variable X_i , $P(X_i | \text{parents}(X_i))$ is the *conditional probability distribution* (CPD) of X_i . The CPDs together with the directed graph fully determine the joint distribution of the Bayesian network.

⁶A type of probabilistic models represented by graphs.

⁷A topological ordering is an ordering of the nodes such that for every two nodes u and v connected by a directed edge from u to v , u appears before v in such ordering. A topological ordering can only be defined on directed acyclic graphs.

The network can be used for inference by querying the distribution of a subset of variables, usually given some evidence. Given a subset of variables $\mathbf{Q} \subset \mathbf{X}$ and an assignment of values \mathbf{e} of the evidence variables, the query is the posterior distribution $P(\mathbf{Q}|\mathbf{e})$. To compute the posterior distribution one needs an inference algorithm. Such algorithm can be exact – such as the *variable elimination* algorithm (Zhang and Poole, 1996) – or approximate – such as the *loopy belief propagation* algorithm (Murphy et al., 1999).

The distributions of the single variables can be learned from observed data using *maximum likelihood estimation* or *Bayesian learning*.

2.3.2 Probabilistic rules

The probabilistic rules formalism is a domain-independent dialogue modeling framework. Probabilistic rules are expressed as ***if ... then ... else ...*** constructs mapping logical conditions on the state variables to effects encoded by either probability distributions or utility functions. The former are called *probability rules* while the latter are *utility rules*. While we make use of both types of rules in our work, here we concentrate only on the probability rules which are the ones used for the resolution of the NSUs.

Let c_1, \dots, c_n be a sequence of logical conditions and $P(E_1), \dots, P(E_n)$ a sequence of categorical probability distributions⁸ over mutually exclusive effects. A probability rule r is defined as follows:

$$\begin{aligned}
 r : & \\
 & \forall \mathbf{x} \\
 & \mathbf{if } c_1 \mathbf{ then} \\
 & \quad \left\{ \begin{array}{l} P(E_1 = e_{1,1}) = p_{1,1} \\ \dots \\ P(E_1 = e_{1,m_1}) = p_{1,m_1} \end{array} \right. \\
 & \mathbf{else if } c_2 \mathbf{ then} \\
 & \quad \left\{ \begin{array}{l} P(E_2 = e_{2,1}) = p_{2,1} \\ \dots \\ P(E_2 = e_{2,m_2}) = p_{2,m_2} \end{array} \right. \\
 & \mathbf{else} \\
 & \quad \left\{ \begin{array}{l} P(E_n = e_{n,1}) = p_{n,1} \\ \dots \\ P(E_n = e_{n,m_n}) = p_{n,m_n} \end{array} \right.
 \end{aligned}$$

⁸A categorical distribution is a probability distribution of an event having a finite set of outcomes with defined probability.

The random variable E_i encodes a range of possible effects $e_{i,1}, \dots, e_{i,m_i}$, each one with a corresponding probability $p_{i,j}$. The conditions and effects of a rule may include underspecified variables, denoted with \mathbf{x} , which are universally quantified on the top of the rule. The effects are duplicated for every possible assignments (grounding) of the underspecified variables.

Each pair of condition and probability distribution over the effects $\langle c_i, P(E_i) \rangle$ is a branch br_i of the rule. Overall, the rule is a sequence of branches br_1, \dots, br_n . The rule is “executed” by running sequentially through the branches. Only the first condition satisfied triggers the corresponding probabilistic effect, the subsequent branches are ignored (as in programming languages).

The dialogue state is represented as a Bayesian network containing a set of nodes (random variables). At each state update, rules are instantiated as nodes in the network. For each rule, the input edges of the node come from the condition variables whereas the output edges go towards the effect variables. The probability distribution of the rule is extracted by executing it. The probability distribution of the effect variables are then retrieved by probabilistic inference. Lison (2014) details the rules and update procedure.

The probabilistic rules are useful in at least three ways:

- They are expressly designed for dialogue modeling. They combine the expressivity of both probabilistic inference and first order logic. This is an advantage in dialogue modeling where one has to describe objects that relate to each other in the dialogue domain and, at the same time, handle uncertain knowledge of the state variables.
- They can cope with the scarcity of training data of most dialogue domains by exploiting the *internal structure* of the dialogue models. By using logical formulae to encode the conditions for a possible outcome, it is possible to group the values of the variables into *partitions*, reducing the number of parameters needed to infer the outcome distribution and therefore the amount of data needed to learn the distribution.
- The state update is handled with probabilistic inference therefore they can operate under uncertain settings which is often needed in dialogue modeling where variables are best represented as *belief states*, continuously updated by observed evidence.

The probabilistic rules formalism has also been implemented into a framework called OpenDial (Lison and Kennington, 2015). OpenDial is a Java toolkit for developing spoken dialogue systems using the probabilistic rules formalism. Using an XML-based language it is possible to define in OpenDial the probabilistic rules to handle the evolution of the dialogue state in a domain-independent way. OpenDial can either work on existent transcripts or as an interactive user interface. OpenDial can also learn parameters from small amounts of data using either supervised or reinforcement learning.

2.4 Summary

In this chapter we discussed the background knowledge needed for describing our work on non-sentential utterances. We first described the notion of non-sentential utterances and the problem of interpreting them. We showed how those utterances can be categorized with a taxonomy from Fernández and Ginzburg (2002). We described how the interpretation of non-sentential utterances can be addressed by first classifying them using the aforementioned taxonomy and then applying some kind of “resolution” procedure to extract their meaning from the dialogue context. In Chapter 3 we will address the NSU classification problem on the basis of the experiments from Fernández et al. (2007). In Chapter 4 instead we will address the NSU resolution task. Fernández (2006) describes a set of NSU resolution rules rooted in a TTR representation of the dialogue context. Section 2.2 briefly described the TTR notions we employed as well as the dialogue context theory based on TTR from Ginzburg (2012).

At last we described the probabilistic modeling of dialogue from Lison (2014) based on the probabilistic rules formalism. As we mentioned in Chapter 1 this formalism is the framework for our formulation of the NSU resolution rules on the basis of the one developed by Fernández (2006). We described in Section 2.3 the basic notion of Bayesian networks which is the representation of the dialogue state employed by the probabilistic rules formalism. Finally, in Section 2.3.2 we explained the probabilistic rules formalism itself and its advantages.

Chapter 3

Classification of Non-Sentential Utterances

Non-sentential utterances are pervasive dialogue phenomena. Any dialogue processing application (e.g. parsing or machine translation of dialogues, or interactive dialogue systems) has to take into account the presence of NSUs and deal with them. As described in Section 2.1, the NSUs come in a great variety of forms that must be treated differently from one another. To this end, the most basic (and perhaps useful) task is classifying them. In our work we employ the taxonomy and the corpus described in Sections 2.1.1 and 2.1.2. As demonstrated by Fernández et al. (2007), we can use machine learning techniques to automatically classify a given NSU, using the annotated corpus as training data. Fernández et al. (2007) is our main theoretical reference and (to our knowledge) the state-of-the-art in performance for the task of classification of NSUs. We first replicated the approach of the aforementioned work and used it as a benchmark for our experiments. Secondly, we tried to improve the classification performances, starting from an expansion of the feature set, then employing semi-supervised learning to address the scarcity of labeled data.

3.1 The data

The corpus from Fernández et al. (2007) contains 1 283 annotated NSU instances, each one identified by the name of the containing BNC file and their sentence number, a sequential number to uniquely identify a sentence in a dialogue transcript. The instances are also tagged with the sentence number of their antecedent which makes up the context for the classification. The raw utterances can be retrieved from the BNC using this information.

For the classification task, we make the same simplifying restriction on the corpus made by Fernández et al. (2007), that is to consider only the NSUs whose antecedent is their preceding sentence. This assumption facilitates the feature extraction procedure without reducing significantly the size of the dataset (about 12% of the total). The resulting distribution of the NSUs after the restriction is showed in Table 3.1.

NSU class	Total
Plain Acknowledgment (Ack)	582
Short Answer (ShortAns)	105
Affirmative Answer (AffAns)	100
Repeated Acknowledgment (RepAck)	80
Clarification Ellipsis (CE)	66
Rejection (Reject)	48
Repeated Affirmative Answer (RepAffAns)	25
Factual Modifier (FactMod)	23
Sluice	20
Helpful Rejection (HelpReject)	18
Filler	16
Check Question (CheckQu)	15
Bare Modifier Phrase (BareModPh)	10
Propositional Modifier (PropMod)	10
Conjunct (Conj)	5
Total	1123

Table 3.1: Distribution of the classes in the corpus after the simplifying restriction.

As one can see from Table 3.1, the distribution of the instances is quite skewed, largely in favor of some classes than others. Moreover very frequent classes are usually the easiest to classify, leaving the most difficult ones with few instances as examples for the classifiers. Although the scarcity of the training material and the imbalance of the classes are the two major problems for the classification task, we propose a set of methods to address them, as described in the following sections.

The British National Corpus

The British National Corpus (Burnard, 2000) – BNC for short – is a collection of spoken and written material, containing about 100 million words of (British) English texts from a large variety of sources. Among the others, it contains a vast selection of dialogue transcripts covering a wide range of domains. Each dialogue transcript in the BNC is contained in an XML file along with many details about the dialogue settings. The dialogues are structured following the CLAWS tagging system (Garside, 1993) which segmented the utterances both at word and sentence level. The word units contains both the raw text, the corresponding lemma (headword) and the POS-tag according to the *C5* tagset (Leech et al., 1994). Each sentence is identified by an unique ID number within the file. Sentences can also contain information about the pauses and the unclarities. The sentences are sorted in their order of appearance and include additional information about temporal alignment in case of overlapping.

3.2 Machine learning algorithms

We employ two different supervised learning algorithms: decision trees and support vector machines. The former are used mainly as a comparison with Fernández et al. (2007) which employ this algorithm as well. For parameters tuning we implemented a coordinate ascent algorithm. As a framework for our experiments we rely on the Weka toolkit (Hall et al., 2009), a Java library containing the implementation of many machine learning algorithms as well as a general-purpose machine learning API.

3.2.1 Classification: Decision Trees

We employ the *C4.5* algorithm (Quinlan, 1993) for decision tree learning. Weka contains an implementation of this algorithm called J48. The goal of decision tree learning is to create a predictive model from the training data. The construction of the decision tree is performed by splitting the training set into subsets according to the values of an attribute. This process is then repeated recursively on each subset. The construction algorithm is usually an informed search using some kind of heuristics to drive the choice of the splitting attribute. In the case of *C4.5* the metric used for the attribute choice is the expected *information gain*. The information gain is based on the concept of *entropy*.

In information theory, the entropy (Shannon, 1948) is the expected value of information carried by a message (or an event in general). It is also a measure of the “unpredictability” of an event. The more unpredictable an event is, the more information it provides when it occurs. Formally, the entropy of a random variable X is

$$H(X) = - \sum_i P(x_i) \log P(x_i)$$

where $P(x_i)$ is the probability of the i -th value of the variable X . A derived notion is the *conditional entropy* of a random variable Y knowing the value of another variable X :

$$\begin{aligned} H(Y|X) &= \sum_i P(x_i) H(Y|X=x_i) \\ &= \sum_i P(x_i) \sum_j P(y_j|x_i) \log P(y_j|x_i) \end{aligned}$$

where x_i are the values of the variable X and y_j are the value of the variable Y .

For the decision tree construction, the *information gain* of an attribute A is the reduction of the entropy of the class C gained by knowing the value of A :

$$IG(A, C) = H(C) - H(C|A)$$

The attribute with the highest information gain is used as splitting attribute.

3.2.2 Classification: Support Vector Machines

The Support Vector Machines (SVMs) (Boser et al., 1992) is one of the most studied and reliable family of learning algorithms. An SVM is a binary classifier that uses a representation of the instances as points in a m -dimensional space, where m is the number of attributes. Assuming that the instances of the two classes are linearly separable¹, the goal of SVMs is to find an hyperplane that separates the classes with the maximum margin. The task of finding the best hyperplane that separates the classes is defined as an optimization problem. SVMs can also be formulated to have “soft margins” i.e. allowing some points of a class to lay in the opposite side of the hyperplane in order to find a better solution. The SVM algorithm we use regularizes the model through a single parameter C .

The SVMs can also be used with non-linear (i.e. non linearly separable) data using the so called *kernel method*. A kernel function maps the points from the input space into an high-dimensional space where they might be linearly separable. A popular kernel function is the (Gaussian) Radial Basis Function (RBF) which maps the input space into an infinite-dimensional space. Its popularity is partially due to the simplicity of its model which involves only one parameter γ .

Even though SVMs are defined as binary classifiers, they can be extended to a multi-class scenario by e.g. training multiple binary classifiers using a one-vs-all or a one-vs-one classification strategy (Duan and Keerthi, 2005).

The Weka toolkit contains an implementation of SVMs that uses the Sequential Minimal Optimization (SMO) algorithm (Platt et al., 1999). In all our experiments we use the SMO algorithm with an RBF kernel.

3.2.3 Optimization: Coordinate ascent

The parameter tuning of all our experiments is carried out automatically through a simple *coordinate ascent*² optimization algorithm. Coordinate ascent is based on the idea of maximizing a multivariable function $f(\mathbf{X})$ along one direction at a time, as apposed to e.g. gradient descent which follows the direction given by the gradient of the function. Our implementation detects the ascent direction by lookup of the function value. The Algorithm 1 contains a procedure to maximize a function f along the direction k while Algorithm 2 performs the coordinate ascent. The step-size values decay at a rate given by the coefficient α . The minimum step-sizes determine the stopping conditions for the *maximize* function, instead the *coordinateAscent* algorithm stops as soon as the found values do not change between two iterations therefore the maximum is found. The latter algorithm can be easily modified to account for a stopping condition given by a maximum number of iterations.

¹An hyperplane can be drawn in the space such that the instances of one class are all in one side of the hyperplane and the instances of the other class are in the other side.

²Also known as coordinate descent which is the minimization counterpart (distinguished only by changing the sign of the function).

We implemented this algorithm ourselves because, for technical reasons, it was easier than rely on a third-party API. Its simplicity is one of its advantages but it is more prone to be stuck on local maximums than more sophisticated techniques such as gradient ascent.

For all our experiments we use the described algorithm to find the parameters that yield the maximum accuracy of the classifiers (using 10-fold cross-validation). For the SMO algorithm we optimize the parameters C and γ , whereas for the J48 algorithm we optimize the parameters C (confidence threshold for pruning) and M (minimum number of instances per leaf).

Algorithm 1: *maximize*($f, k, \mathbf{X}, \delta_k, m_k$)

Input: Function f to be maximized; index k of the parameter to maximize; vector \mathbf{X} of the current parameter values; initial step-size value δ_k for the k -th parameter; minimum step-size value m_k for the k -th parameter

Output: The value of the k -th parameter that maximizes the function f along the corresponding direction

```

1  $y_{max} \leftarrow f(\mathbf{X});$ 
2 while  $|\delta_k| \geq m_k$  do
3    $\hat{\mathbf{X}} \leftarrow \mathbf{X};$ 
4    $\hat{\mathbf{X}}[k] \leftarrow \delta_k \times \hat{\mathbf{X}}[k];$ 
5    $\hat{y} \leftarrow f(\hat{\mathbf{X}});$ 
6   if  $\hat{y} > y_{max}$  then
7      $y_{max} \leftarrow \hat{y};$ 
8      $\mathbf{X}[k] \leftarrow \hat{\mathbf{X}}[k];$ 
9   else
10     $\delta_k \leftarrow -\delta_k;$ 
11  end
12   $\delta_k \leftarrow \alpha \times \delta_k;$ 
13 end
14 return  $\mathbf{X}[k];$ 

```

Algorithm 2: *coordinateAscent*($f, n, \mathbf{X}, \Delta, \mathbf{m}$)

Input: Function f to be maximized; number n of parameters of the function; vector \mathbf{X} of initial parameter values; vector Δ of initial step-size values; vector \mathbf{m} of minimum step-sizes;

Output: The vector \mathbf{X} that maximizes the function f

```

1 initialize  $\mathbf{X}_{last}$  to random values (different from  $\mathbf{X}$ );
2 while  $\mathbf{X} \neq \mathbf{X}_{last}$  do
3    $\mathbf{X}_{last} \leftarrow \mathbf{X};$ 
4   for  $k \in 1 \dots n$  do
5      $\delta_k \leftarrow \Delta[k];$ 
6      $m_k \leftarrow \mathbf{m}[k];$ 
7      $\mathbf{X}[k] \leftarrow \text{maximize}(f, k, \mathbf{X}, \delta_k, m_k);$ 
8   end
9 end
10 return  $\mathbf{X};$ 

```

3.3 The baseline feature set

Our baseline is set to be the replicated approach of the classification experiments carried out by Fernández et al. (2007), which is our reference work for our study. It contains two experiments, one with a restricted set of classes (leaving out Acknowledgments and Check Questions) and a second taking into account all classes. We are interested in the latter although the former is useful to understand the problem and to analyze the results of our classifier. The aforementioned paper also contains an analysis of the results and the feature contribution which proved useful in the replication of the experiments. For our baseline we use only the features they describe. The feature set is composed of 9 features exploiting a series of syntactic and lexical properties of the NSUs and their antecedents. The features can be categorized as: *NSU features*, *Antecedent features*, *Similarity features*. Table 3.2 contains an overview of the feature set.

NSU features

Different NSU classes are often distinguished by their form. The following is a group of features exploiting their syntactic and lexical properties.

nsu_cont

Denotes the “content” of the NSU i.e. whether it is a question or a proposition. This is useful to distinguish between question denoting classes, such as Clarification Ellipsis and Sluices, and the rest.

wh_nsu

Denotes whether the NSU contains a wh-word, namely: *what*, *which*, *who*, *where*, *when*, *how*. This can help for instance to distinguish instances of Sluices and Clarification Ellipsis knowing that the former are *wh*-questions while the latter are not.

aff_neg

Denotes the presence of a yes-word, a no-word or an ack-word in the NSU. Yes-words are for instance: *yes*, *yep*, *aye*; no-words are for instance: *no*, *not*, *nah*; ack-words are: *right*, *aha*, *mhm*. This is particularly needed to distinguish between Affirmative Answers, Rejections and Acknowledgments.

lex

Indicates the presence of lexical items at the beginning of the NSU. This feature is intended to indicate the presence of modifiers. A modal adverb (e.g. *absolutely*, *clearly*, *probably*) at the beginning of the utterance usually denotes a Propositional Modifier. The same applies for Factual Modifiers, which are usually denoted by factual adjectives (e.g. *good*, *amazing*, *terrible*, *brilliant*) and Conjuncts which are usually denoted by conjunctions. Bare Modifier Phrases are a wider class of NSUs which do not have a precise lexical conformation but they are usually started by lexical patterns containing a Prepositional Phrase (PP) or an Adverbial Phrase (AdvP).

Feature	Description	Values
<code>nsu_cont</code>	the content of the NSU (a question or a proposition)	<code>p,q</code>
<code>wh_nsu</code>	presence of a <i>wh</i> -word in the NSU	<code>yes,no</code>
<code>aff_neg</code>	presence of a <i>yes/no</i> -word in the NSU	<code>yes,no,e(empty)</code>
<code>lex</code>	presence of different lexical items at the beginning of the NSU	<code>p_mod,f_mod,mod,conj,e</code>
<code>ant_mood</code>	mood of the antecedent utterance	<code>decl,n_decl</code>
<code>wh_ant</code>	presence of a <i>wh</i> -word in the antecedent	<code>yes,no</code>
<code>finished</code>	whether the antecedent is (un)finished	<code>fin,unf</code>
<code>repeat</code>	number of common words in the NSU and the antecedent	0-3
<code>parallel</code>	length of the common tag sequence in the NSU and the antecedent	0-3

Table 3.2: An overview of the baseline feature set.

Antecedent features

As for the NSUs, antecedents also show different syntactic and lexical properties that can be used as features for the classification task. This is a group of features exploiting those properties.

`ant_mood`

As defined by Rodríguez and Schlangen (2004), this feature was thought to distinguish between declarative and non-declarative antecedent sentences. This feature is useful to indicate the presence of an answer NSU, if the antecedent is a question, or a modifier, if the antecedent is not a question.

`wh_ant`

As the corresponding NSU feature, this indicates the presence of a *wh*-word in the antecedent. Usually Short Answers are answers to *wh*-questions while Affirmative Answers and Rejections are answers to polar questions i.e. *yes/no*-questions without a *wh*-interrogative.

`finished`

This feature encodes a truncated antecedent sentence as well as the presence of uncertainties at the end of it. Truncated sentences lack a closing full stop, question mark or exclamation mark. Uncertainties are given by the presence of pauses or unclear words or else a last word being “non-closing”, e.g. conjunctions or articles.

Similarity features

As discussed in Section 2.1, some classes show some kind of parallelism between the NSU and its antecedent. The parallelism of certain classes can be partially captured by similarity measures. The following is a group of features encoding the similarity at the word and POS-level between the NSUs and their antecedents.

`repeat`

This feature counts the content words that the NSU and the antecedent have in common (a maximum value of 3 is taken as a simplification). A value greater than 0 is usually a sign of Repeated Acknowledgment or Repeated Affirmative Answers.

`parallel`

This feature encodes whether there is a common sequence of POS tags between the NSU and the antecedent and denotes its length. This feature can help classify Repeated Acknowledgments, Repeated Affirmative Answers and Helpful Rejections.

3.4 Feature engineering

The first and most straightforward method we use to address the classification problem is to find more features to describe the NSU instances. We present here the combination of features that we employ as our final approach. The extended feature set is composed of all the baseline features plus 23 new linguistic features, summing up to a total of 32 features. Our features can be clustered into five groups: *POS-level features*, *Phrase-level features*, *Dependency features*, *Turn-taking features* and *Similarity features*. Table 3.3 shows an overview of the additional features we use in the extended feature set.

POS-level features

Shallow syntactic properties of the NSUs that make use of the pieces of information already present in the BNC such as POS tags and other markers.

`pos_{1,2,3,4}`

A feature for each one of the first four POS-tags in the NSU. If an NSU is shorter than four words the value `None` is assigned to each missing POS tag. Many NSU classes share (shallow) syntactic patterns among their instances, especially at the beginning of the NSU phrase. Those features aim to capture those patterns in a shallow way through the POS tags.

`ending_punct`

A feature to encode the final punctuation mark of the antecedent if any.

`has_pause`

Marks the presence of a pause in the antecedent.

`has_unclear`

Marks the presence of an unclear passage in the antecedent.

Feature	Description	Values
pos_{1,2,3,4}	POS tags of the first four words in the NSU	C5 tag-set
ending_punct	ending punctuation in the antecedent if any	.,?,!,e
has_pause	presence of a pause in the antecedent	yes,no
has_unclear	presence of an “unclear” marker in the antecedent	yes,no
ant_sq	presence of a SQ tag in the antecedent	yes,no
ant_sbarq	presence of a SBARQ tag in the antecedent	yes,no
ant_sinv	presence of a SINV tag in the antecedent	yes,no
nsu_first_clause	first clause-level syntactic tag in the NSU	S,SQ,...
nsu_first_phrase	first phrase-level syntactic tag in the NSU	NP,ADVP,...
nsu_first_word	first word-level syntactic tag in the NSU	NN,RB,...
neg_correct	presence of a negation followed by a correction	yes,no
ant_neg	presence of a <i>neg</i> dependency in the antecedent	yes,no
wh_inter	presence of a <i>wh</i> -interrogative fragment in the antecedent	yes,no
same_who	whether the NSU and its antecedent have been uttered by the same speaker	same,diff,unk
repeat_last	number of repeated words between the NSU and the last part of the antecedent	numeric
abs_len	number of words in the NSU	numeric
cont_len	number of content-words in the NSU	numeric
local_all	the local alignment (at character-level) of the NSU and the antecedent	numeric
lcs	longest common subsequence (at word-level) between the NSU and the antecedent	numeric
lcs_pos	longest common subsequence (at pos-level) between the NSU and the antecedent	numeric

Table 3.3: An overview of the additional features comprised in the extended feature set.

Phrase-level features

Occurrence of certain syntactic structures in the NSU and the antecedent. These features were extracted through the use of the Stanford PCFG parser (Klein and Manning, 2003) on the utterances. Refer to Marcus et al. (1993) for more information about the tag set used for the English grammar.

`ant_{sq,sbarq,sinv}`

Those features indicate the presence of the syntactic tags `SQ`, `SBARQ` and `SINV` in the antecedent. Those tags indicate a question formulated in various ways even when there is no explicit question mark at the end. Useful to recognize e.g. Short Answers.

`nsu_first_clause`

Marks the first clause-level tag (`S`, `SQ`, `SBAR`, ...) in the NSU.

`nsu_first_phrase`

Marks the first phrase-level tag (`NP`, `VP`, `ADJP`, ...) in the NSU.

`nsu_first_word`

Marks the first word-level tag (`NN`, `RB`, `UH`, ...) in the NSU.

`neg_correct`

Presence of a negation word (*no*, *nope*, ...), followed by a comma and a correction. For instance:

- (3.1) A: Or, or were they different in your childhood?
B: No, always the same.
[BNC: HDH 158–159]

This pattern is useful to describe some of the Helpful Rejections such as (3.1).

Dependency features

Presence of certain dependency patterns in the antecedent. These features were extracted through the use of the Stanford Dependency Parser (Chen and Manning, 2014) on the utterances. For more details about the dependency relations tag set please refer to De Marneffe et al. (2014).

`ant_neg`

Signals the presence of a *neg* dependency relation in the antecedent. The *neg* dependency arises from an adverbial negation in the sentence (*not*, *don't*, *never*, ...). This feature helps to capture situations such as the following:

- (3.2) A: You're not getting any funny fits from that at all, June?
B: Er no.
[BNC: H4P 36–37]

Since the question in the antecedent is negative, the NSU in (3.2) is actually an Affirmative Answer, even though it contains a negative word. This feature, in combination with the `aff_neg` feature, addresses this situation.

`wh_inter`

Whether the antecedent contains a *wh-interrogative* fragment such as the one in the following example:

- (3.3) A: And you know what the voltage is
B: Yeah, two forty.
[BNC: GYR 174–175]

The feature looks for a *dobj* dependency with a *wh*-word then for an *nsubj* dependency with the dependent element of the previous dependency, for instance in (3.3) we have `dobj(is-7, what-4)` and `nsubj(is-7, voltage-6)`. This feature tries to mitigate the absence of a question as antecedent for Short Answers such as (3.3).

Turn-taking features

Features indicating certain patterns in the turn-taking of the dialogue.

`same_who`

Denotes whether the NSU and the antecedent were uttered by the same speaker. Sometimes dialogues do not provide the speaker information so an additional value `unk` is added for this cases. This feature is particularly important to capture Check Questions which are almost always uttered by the same speaker.

Similarity features

Additional numeric features and similarity measures between the NSU and its antecedent.

`repeat_last`

This measures the number of words in common between the NSU and the last portion of the antecedent. Often happens that Repeated Affirmative Answers and Repeated Acknowledgments contain the last words in the antecedent.

`abs_len`

The total number of words in the NSU.

`cont_len`

The number of content-words in the NSU.

`local_all`

A feature that denotes the local alignment at the character-level between the NSU and the antecedent, computed using the Smith–Waterman algorithm (Smith and Waterman, 1981).

`lcs`

A feature to express the longest common subsequence at the word-level between the NSU and its antecedent, computed using a modified version of the Needleman–Wunsch algorithm (Needleman and Wunsch, 1970), tailored to account for words instead of characters.

`lcs_pos`

The longest common subsequence at the POS-level between the NSU and its antecedent, computed with the same algorithm of above but using the list of POS tags instead of the list of words.

3.5 Semi-Supervised Learning

The scarcity of labeled data is probably the major problem to face in this classification task. Even though the quality of the data is good enough, it is still difficult for a classifier to learn patterns out of 20 instances or less for some classes (see Table 3.1). However, a large amount of unlabeled data is available in the BNC. There are many classification tasks, such as ours, in which it is hard or costly to label a large amount of instances while instead it is relatively cheap to extract unlabeled ones. The empirical question is whether the use of unlabeled data is useful to improve the classification performances. Semi-Supervised Learning techniques deal with this issue. They exploit the combination of a small amount of labeled data and a large amount of unlabeled data to try improve the classification accuracy. Even though it is still a young research field, semi-supervised learning has already found many fields of application (Liang, 2005; Bergsma, 2010).

3.5.1 Unlabeled data extraction

With the use of some heuristics it is possible to extract NSU instances of good quality from the BNC. We use a set of rules to determine whether an utterance in a dialogue transcript of the BNC is a probable NSU. The following is a list of such rules:

- The number of words in the NSU must be less than a given threshold;
- The number of characters in the NSU must be higher than a given threshold;
- The NSU must not contain only pauses, unclear passages and punctuation;
- The NSU must not contain a greeting (e.g. *hi*, *hello*, *good night*);
- The NSU must not contain a verb in any form.

An accuracy test was run over the corpus of NSUs: of the 1 123 NSUs examined, 1033 were detected correctly by this set of rules, for an accuracy of 0.92. The main flaws of the rules were mostly overlong NSUs, such as (3.4), and presence of verbs, such as (3.5).

- (3.4) A: Was it a coal fire?
 B: Coal fire and er scrubbed the cabin out like that, soda water and soft soap.³
 [BNC: H5G 151–152]
- (3.5) A: [...] the resistance the same the current goes up.
 B: Current goes up.⁴
 [BNC: GYR 112–113]

The detection of NSUs using the rules above is not the only problem to face. Perhaps more challenging is the selection of an antecedent for the NSU. As pointed out in Section 2.1, the antecedent of an NSU is not always the preceding utterance. Nevertheless, as proved in the corpus study of Fernández (2006), the percentage of the utterances whose antecedent is not the preceding utterance is rather low. Another result of the aforementioned work is that the case in which the antecedent is an utterance at distance greater than one is far more probable in a multi-party dialogue context. In light of the above considerations, we restrict the instances we extract to only those from two-party dialogues and we always consider the preceding utterance as the antecedent of an NSU. While there has been some previous work towards using machine learning techniques for the detection of the antecedent of NSUs in multi-party dialogue (Schlangen, 2005), we consider sufficient the amount of unlabeled data we can extract following the previous rule.

In order to maximize the quality of the unlabeled data that we extract we also enforce some rules over the antecedent utterance:

- The number of words in the antecedent must be greater than the number of words in the NSU;
- The antecedent must have a complete clausal form i.e. at least a verb phrase and a noun phrase.

Using the whole set of heuristics we extracted in total 3 198 new unlabeled NSU instances from the BNC (checked not to be already in the corpus).

3.5.2 Semi-supervised learning techniques

As previously mentioned, semi-supervised learning techniques are used when labeled data is scarce and unlabeled data is abundant. Every techniques tries to integrate the information yield by the unlabeled instances inside a learning model based on the available labeled data. In this section we give a brief and high-level description of the semi-supervised learning techniques that we have employed, namely: *Self Training*, *Transductive SVM* and *Active Learning*.

³A Repeated Affirmative Answer, but the additional content after the conjunction makes the NSU much longer. It is still a valid NSU since it does not have a full clausal structure.

⁴The NSU is a Repeated Acknowledgment. Repeating the words in the antecedent, it introduces a verb. It is still considered an NSU according to the definition of Fernández (2006).

Self Training

The simplest way to exploit unlabeled data is to automatically predict some unlabeled instances through a classifier built from the available labeled data then add them to the training data for the next step. This is an iterative process, at each step one or more newly labeled instances are added to the training set then the classifier is retrained and more unlabeled instances are predicted.

Various strategies can be used at each step:

- Add one or a few (random) instances at the time;
- Add a few most confident instances;
- Add all the first time, correct the wrong predictions the next times.

The last strategy as well as other variants can be cast as an Expectation-Maximization problem, especially when using a probabilistic learning model.

Transductive SVM

As already described in Section 3.2.2, Support Vector Machines are one of the most studied and reliable family of classification algorithms. Transductive SVM (TSVM) is a variant of the standard SVM algorithm which exploits unlabeled data to help adjust the SVM model. The basic assumption under TSVM is that unlabeled instances from different classes are separated with large margin. Therefore, similarly to the standard SVM, TSVM tries to find the hyperplane that maximizes the unlabeled data margin i.e. considering unlabeled points as labeled ones. To decide whether an unlabeled point should be considered of one class or the other, clustering techniques are used e.g. k -nearest neighbors (the class of the majority of the neighbors or some other variant).

We will not go into mathematical details so we recommend the interested reader to Vapnik (1998), Collobert et al. (2006).

Active Learning

Annotating data is often a very expensive procedure, mostly because one needs to annotate a lot of instances in order to be able to reliably classify unseen ones. An idea to ease this problem is to let the learning algorithm choose which instance could be the most informative (i.e. the most difficult to predict) then annotate it manually. This technique has the advantage of reducing the cost of manual annotation of the instances by making informed guesses over the instances to label and discarding the redundant ones.

This kind of techniques is typically employed to cope with the scarcity of labeled data. In our case, the lack of sufficient training data is especially problematic due to the strong class imbalance between the NSU classes.

The Active Learning (AL) scheme, which is a special case of semi-supervised learning, trains the model over the available labeled data then queries the user for the label of one (or few more) instances then retrain the model and so on until convergence criteria are met, e.g. the wanted number of new instances is reached.

There can be different query strategies, some of them are:

- **Uncertainty Sampling:** queries the least confident instance (according to the probability of the prediction). A variant of that uses entropy to determine the most informative instance.
- **Query-by-committee:** uses many different classifiers to predict unlabeled data then formulates the most informative query as the instance about which they most disagree.
- **Expected Model Change:** selects the instance that would impart the greatest change to the model, according to a decision-theoretic approach.
- **Expected Error Reduction:** Another decision-theoretic approach that aims to minimize the *risk*, that is the expected future error. The instances are selected on the basis of how much the model generalization error is likely to be reduced. A variant of this approach considers only the output variance of the model.

The particular active learning algorithm we employed in our experiments is a pool-based method⁵ with uncertainty sampling (Lewis and Catlett, 1994). The sampling relies on *entropy* as measure of uncertainty. Given a particular (unlabeled) instance with a vector of feature values \mathbf{f} , we use the existing classifier to predict the class C of the instance, and derive the probability distribution $P(C = c_i | \mathbf{f})$ for each possible output class c_i . We can then determine the corresponding entropy of the class C :

$$H(C) = - \sum_i P(C = c_i | \mathbf{f}) \log P(C = c_i | \mathbf{f})$$

As seen in section 3.2.1, entropy indicates the “unpredictability” of a random variable and also how much information it carries. The higher the entropy of the class of an instance the more information we gain by knowing it. The algorithm we employ (Algorithm 3) selects the instances with highest entropy as the most informative ones. As argued in Settles (2010), entropy sampling is especially useful when there are more than two classes, as in our setting. In practice, we applied the JCLAL active learning library⁶ to extract and annotate 100 new instances of NSUs, which were subsequently added to the existing training data.

⁵That involves drawing labeled instances from a “pool” that remains the same over the iterations, as opposed of stream-based ones in which sampling is done over a stream of data.

⁶cf. <https://sourceforge.net/projects/jclal>.

Algorithm 3: *entropySampling*(Γ, \mathbf{U}, k)

Input: The classifier Γ ; the unlabeled data \mathbf{U} ; the sample size k .

Output: The k instances with highest entropy.

```
1  $\mathbf{H} \leftarrow$  vector of the same size of  $\mathbf{U}$ ;  
2 for  $i \in 1..|\mathbf{U}|$  do  
3    $u \leftarrow \mathbf{U}[i]$ ;  
4    $\mathbf{P}_u \leftarrow \text{classProbDist}(\Gamma, u)$ ;  
5    $H_u \leftarrow -\sum_{p \in \mathbf{P}_u} p \log p$ ;  
6    $\mathbf{H}[i] \leftarrow H_u$ ;  
7 end  
8  $\mathbf{U} \leftarrow \text{sort}(\mathbf{U}, \mathbf{H})$ ; // Sort  $\mathbf{U}$  according to  $\mathbf{H}$  (descending)  
9 return firstK( $\mathbf{U}, k$ );
```

3.6 Evaluation

In this section we discuss the evaluation of our experiments and their empirical results. We first discuss the evaluation metrics for the classification task we employed, then we present the evaluation results on each setting.

3.6.1 Metrics

Given the dataset with a total of N instances, the metrics are based on the amount of true positives (TP), true negatives (TN), false positive (FP) and false negatives (FN).

Accuracy

The ratio of the correctly classified instances over the total

$$Acc = \frac{\sum_{c \in C} TP_c + TN_c}{N}$$

where C is the set of the classes and TP_c and TN_c are respectively the true positives and the true negatives of the class $c \in C$.

Precision

The ratio between the true positives and the total instances classified as positives. In a context with multiple classes (more than two) such as ours, the precision must be calculated per class, where the positive instances are the ones classified with the current class whereas the negative instances are the ones classified otherwise. The per class precision is calculated as follows:

$$Prec_c = \frac{TP_c}{TP_c + FP_c}$$

To have a summary value for all the classes we can compute the weighted average precision:

$$Prec_{avg} = \frac{\sum_{c \in C} N_c \cdot Prec_c}{N}$$

Recall

The recall is the ratio between the true positives and the total instances that are actually positives. As for the precision, we can calculate the per class recall:

$$Rec_c = \frac{TP_c}{TP_c + FN_c}$$

And the weighted average recall:

$$Rec_{avg} = \frac{\sum_{c \in C} N_c \cdot Rec_c}{N}$$

F_1 -score

The F_1 -score is the harmonic mean of precision and recall. As for the other two measures, we compute the per class F_1 -score:

$$\begin{aligned} F_{1,c} &= 2 \cdot \frac{Prec_c \cdot Rec_c}{Prec_c + Rec_c} \\ &= \frac{2 \cdot TP_c}{2 \cdot TP_c + FP_c + FN_c} \end{aligned}$$

Then the weighted average F_1 -score:

$$F_{1,avg} = \frac{\sum_{c \in C} N_c F_{1,c}}{N}$$

Student's t -test

Empirical results alone can not assess whether a classifier performs *better* than another. To assess that the performances of one classifier being higher than a second one is not due to the randomness associated with the data manipulation but to a *statistically significant* difference between the classifiers one needs to prove with high confidence that the *null hypothesis* is false. The *null hypothesis* is a statement that is assumed to be true until evidence indicates otherwise. When comparing two learning systems, the *null hypothesis* states that there is no difference between the performances of the two learning systems. To prove that a classifier performs better than another we need to disprove the *null hypothesis* with a high degree of confidence. For this purpose we employ a Student's t -test, a widespread method to compare two sets of data. The t -test can be used to find the probability p of the performance values of the two classifiers being drawn from the same mean.

To run the t -test, we compare the differences δ_i among the performance values of the two classifiers over the n independent samples. We first compute the mean of the differences:

$$\bar{\delta} = \frac{1}{n} \sum_{i=1}^n \delta_i$$

Then we compute the t -statistic:

$$t = \frac{\bar{\delta}}{\sqrt{\frac{1}{n(n-1)} \sum_{i=1}^n (\delta_i - \bar{\delta})^2}}$$

From the t -statistic we can derive the p -value from a *Student's t -distribution* with $n - 1$ degree of freedom. A small p -value means that it is unlikely that the samples show such a t -statistic by chance therefore we can assess that the difference in performance between the two classifiers is statistically significant.

In our case we use a paired t -test on the accuracy values of the 10-fold cross-validation over the dataset (thus $n = 10$). By convention, an acceptable p -value is $p \leq 0.05$. For our experiments we rely on the `t.test` function from the *R* project, a framework for statistical computing (R Core Team, 2015).

3.6.2 Empirical results

Baseline

As in Fernández et al. (2007), we evaluate our system in a 10-fold cross-validation fashion. Weka's J48 algorithm was used as a comparing classifier. Thanks to the analysis of the resulting trees, we managed to imitate quite closely the behavior of their system as well as reaching a very close performance overall. Even though we use the same feature set and the same algorithm the performance parameters turn out to be slightly lower than the ones claimed in Fernández et al. (2007). That might be for a variety of reasons, for instance the way feature were extracted or how the parameters were tuned. Nevertheless the overall performance is matched as well as many of the patterns in the scores. Table 3.4 shows the comparison between the performance parameters of the reference classification (Fernández et al., 2007) and the values of the same parameters achieved by our implementation.

Self-training and TSVM

Both those two techniques did not perform particularly well, sometimes even worsening the classification accuracy. Self-training was implemented and tested in many variants but none were successful. One possible explanation is that the labeled data added at each step to the training data is always biased by the labeled data available in the initial training set. This may lead to adding redundant data that is not actually useful to improve the classification performances. On the other hand, TSVM has been unsuccessful mostly due to computational performances of the implementation and other technical difficulties. It was impracticable to run it on a large amount of unlabeled data so we managed to test it only on few unlabeled instances and therefore no improvement was shown.

NSU Class	Our replica			Reference classification		
	Precision	Recall	F_1 -score	Precision	Recall	F_1 -score
Ack	0.97	0.97	0.97	0.97	0.95	0.96
AffAns	0.89	0.84	0.86	0.83	0.86	0.84
BareModPh	0.63	0.65	0.62	1.00	0.70	0.82
CE	0.87	0.89	0.87	0.92	0.92	0.94
CheckQu	0.85	0.90	0.87	0.83	1.00	0.91
ConjFrag	0.80	0.80	0.80	0.71	1.00	0.83
FactMod	1.00	1.00	1.00	1.00	0.91	0.95
Filler	0.77	0.70	0.71	0.50	0.37	0.43
HelpReject	0.13	0.14	0.14	0.46	0.33	0.39
PropMod	0.92	0.97	0.93	1.00	0.60	0.75
Reject	0.76	0.95	0.83	0.76	1.00	0.86
RepAck	0.74	0.75	0.70	0.84	0.86	0.85
RepAffAns	0.67	0.71	0.68	0.65	0.68	0.67
ShortAns	0.86	0.80	0.81	0.81	0.83	0.82
Sluice	0.67	0.77	0.71	0.95	1.00	0.98
weighted avg.	0.89	0.89	0.88	0.90	0.90	0.89

Table 3.4: Performances comparison between Fernández et al. (2007) and our replica.

Active Learning

Our active learning experiment was carried out using the JCLAL library. For the active learning process we divided the dataset into three parts: training set (50%), development set (25%) and test set (25%). At each iteration, the JCLAL library builds a classifier on the training set and evaluates it over the development set. The same classifier is then used to select an instance from the unlabeled data, as described in Section 3.5.2. The user is then asked to annotate the selected instance. The process iterates in this manner until the stopping criteria is met, that is when the goal of 100 newly annotated instances is reached. Table 3.5 shows the distribution of the instances annotated with Active Learning. From Table 3.5 we can see that the AL algorithm using the entropy measure prefers the instances that belongs to the classes that are most difficult to classify and, in particular, the ones that are ambiguous, such as Clarification Ellipsis and Sluices. This process has been performed once with the extended feature set and the SMO classifier. Secondly, it has been simulated (i.e. using the data obtained in the previous run) using the baseline feature set instead. The Figures 3.1, 3.2, 3.3, 3.4 show the learning curves⁷, respectively for the accuracy, precision, recall and F_1 -score, of both the extended feature set and the baseline feature set.⁸ All the performance measures are clearly improving as new instances become available, for both the extended feature set and the baseline one.

⁷The graph showing how the performances change as the new labeled data extracted with Active Learning are inserted in the training set.

⁸Notice that the images are scaled on the y -axis to make the change visible.

NSU Class	Instances
Helpful Rejection	21
Repeated Acknowledgment	17
Clarification Ellipsis	17
Acknowledgment	11
Propositional Modifier	9
Filler	9
Sluice	3
Repeated Affirmative Answer	3
Factual Modifier	3
Conjunct Fragment	3
Short Answer	2
Check Question	2
tot.	100

Table 3.5: Distribution of the classes of the instances annotated with Active Learning.

In the end, the test set has been used to evaluate the overall performances of the various settings. Table 3.6 and Table 3.7 show the results of the experiments respectively over the development set and the test set. The results on the test set show that the inclusion of the active learning data is only beneficial when combined with the extended feature set.

We also performed an evaluation of the various settings using 10-fold cross-validation over the full dataset. The evaluation results based on the active learning procedure (AL) refer to the performance of the system after the inclusion of all newly annotated instances. The novel data was added to the training set of each fold.

We compare the results of the various settings using the J48 algorithm (Table 3.8) and SMO algorithm (Table 3.9). The use of active learning was successful and, in the end, the use of the SMO classifier with the extended feature set and the inclusion of the AL instances constitutes our final approach.

The results show a significant improvement of the classification performance between the baseline and the final approach. Using a paired t -test with a 95% confidence interval between the baseline and the final results (as detailed in Section 3.6.1), the improvement in classification accuracy is statistically significant with a p -value of 6.9×10^{-3} .

The SVM algorithm does not perform particularly well with the baseline feature set but scales better than the J48 classifier after the inclusion of the additional features. Overall, the results demonstrate that the classification can be improved using a modest amount of additional training data combined with an extended feature set. However, we can observe from Table 3.10 that some NSU classes remain difficult to classify even with the insertion of additional training data. For instance, Helpful Rejections are still the most difficult classes to classify, even with the addition of 21 new instances. One of the problems with

Training set (feature set)	Accuracy	Precision	Recall	F_1 -score
Train-set (baseline)	0.853	0.857	0.853	0.848
Train-set (extended)	0.860	0.871	0.860	0.858
Train-set + AL (baseline)	0.867	0.883	0.867	0.868
Train-set + AL (extended)	0.884	0.899	0.885	0.886

Table 3.6: Performances of the SMO classifier in the various settings on the development set.

Training set (feature set)	Accuracy	Precision	Recall	F_1 -score
Train-set + Dev-set (baseline)	0.906	0.911	0.906	0.903
Train-set + Dev-set (extended)	0.928	0.937	0.929	0.930
Train-set + Dev-set + AL (baseline)	0.898	0.911	0.898	0.898
Train-set + Dev-set + AL (extended)	0.932	0.945	0.932	0.935

Table 3.7: Performances of the SMO classifier in the various settings on the test set.

Helpful Rejections is that they are connected to their antecedents mainly at the semantic level. Consider the following example of Helpful Rejection that is hard to classify:

- (3.6) A: There was one which you said Ernest Morris was born in 1950.
 B: Fifteen. [BNC: J9A 372–373]

It is clear that, for the Helpful Rejection in (3.6), morpho-syntactic and lexical features, such as the ones we employ, are of little use in classifying this utterance. Most of the connection is at the semantic level therefore we would need to use features that exploit semantic patterns. At the same time, the use of this type of features would add several layers of complexity at the feature extraction process. Other examples of difficult classes are the Repeated Affirmative Answers and Repeated Acknowledgments. They are highly ambiguous because they can be misclassified between each other, with their respective non-repeated classes and sometimes with other NSU classes. An example of ambiguous Repeated Acknowledgment can be the following:

- (3.7) A: Selected period.
 B: Selected period, right, Andrew?⁹
 [BNC: JK8 114–115]

The instance in (3.7) contains also a question therefore it is often misclassified with other question denoting NSU classes. It is clear that handling these type of NSU requires to perform a deeper semantic analysis of the connection with their antecedents then design appropriate semantic features. The extraction of additional labeled data is also especially important for both the feature engineering and the learning process of the classifiers. This two approaches may be the starting points of any future work on this task.

⁹In the dialogue, the speaker B is asking the same question to many people in turns.

Training set (feature set)	Accuracy	Precision	Recall	F_1 -score
Train-set (baseline)	0.885	0.888	0.885	0.879
Train-set (extended)	0.889	0.904	0.889	0.889
Train-set + AL (baseline)	0.890	0.896	0.890	0.885
Train-set + AL (extended)	0.896	0.914	0.896	0.897

Table 3.8: Performances of the J48 classifier in the various settings using 10-fold cross-validation.

Training set (feature set)	Accuracy	Precision	Recall	F_1 -score
Train-set (baseline feature set)	0.881	0.884	0.881	0.875
Train-set (extended feature set)	0.899	0.904	0.899	0.896
Train-set + AL (baseline feature set)	0.883	0.893	0.883	0.880
Train-set + AL (extended feature set)	0.907	0.913	0.907	0.905

Table 3.9: Performances of the SMO classifier in the various settings using 10-fold cross-validation.

NSU Class	Baseline			Final approach		
	Precision	Recall	F_1 -score	Precision	Recall	F_1 -score
Ack	0.97	0.97	0.97	0.97	0.98	0.97
AffAns	0.89	0.84	0.86	0.81	0.90	0.85
BareModPh	0.63	0.65	0.62	0.77	0.75	0.75
CE	0.87	0.89	0.87	0.88	0.92	0.89
CheckQu	0.85	0.90	0.87	1.00	1.00	1.00
ConjFrag	0.80	0.80	0.80	1.00	1.00	1.00
FactMod	1.00	1.00	1.00	1.00	1.00	1.00
Filler	0.77	0.70	0.71	0.82	0.83	0.78
HelpReject	0.13	0.14	0.14	0.31	0.43	0.33
PropMod	0.92	0.97	0.93	0.92	1.00	0.95
Reject	0.76	0.95	0.83	0.90	0.90	0.89
RepAck	0.74	0.75	0.70	0.77	0.77	0.77
RepAffAns	0.67	0.71	0.68	0.72	0.55	0.58
ShortAns	0.86	0.80	0.81	0.92	0.86	0.89
Sluice	0.67	0.77	0.71	0.80	0.84	0.81

Table 3.10: Per class performances comparison between the baseline (J48, baseline feature set) and the final approach (SMO, extended feature set, AL instances).

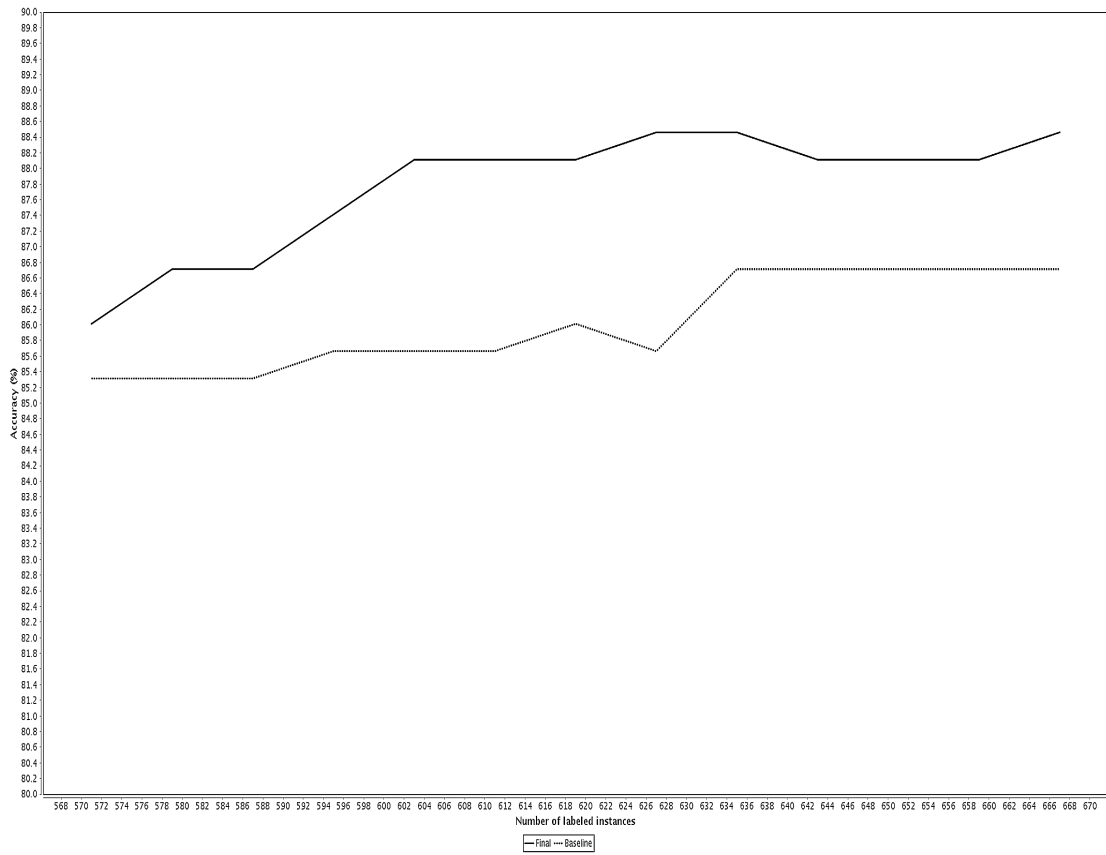


Figure 3.1: Learning curve for the accuracy (output of the JCLAL library).

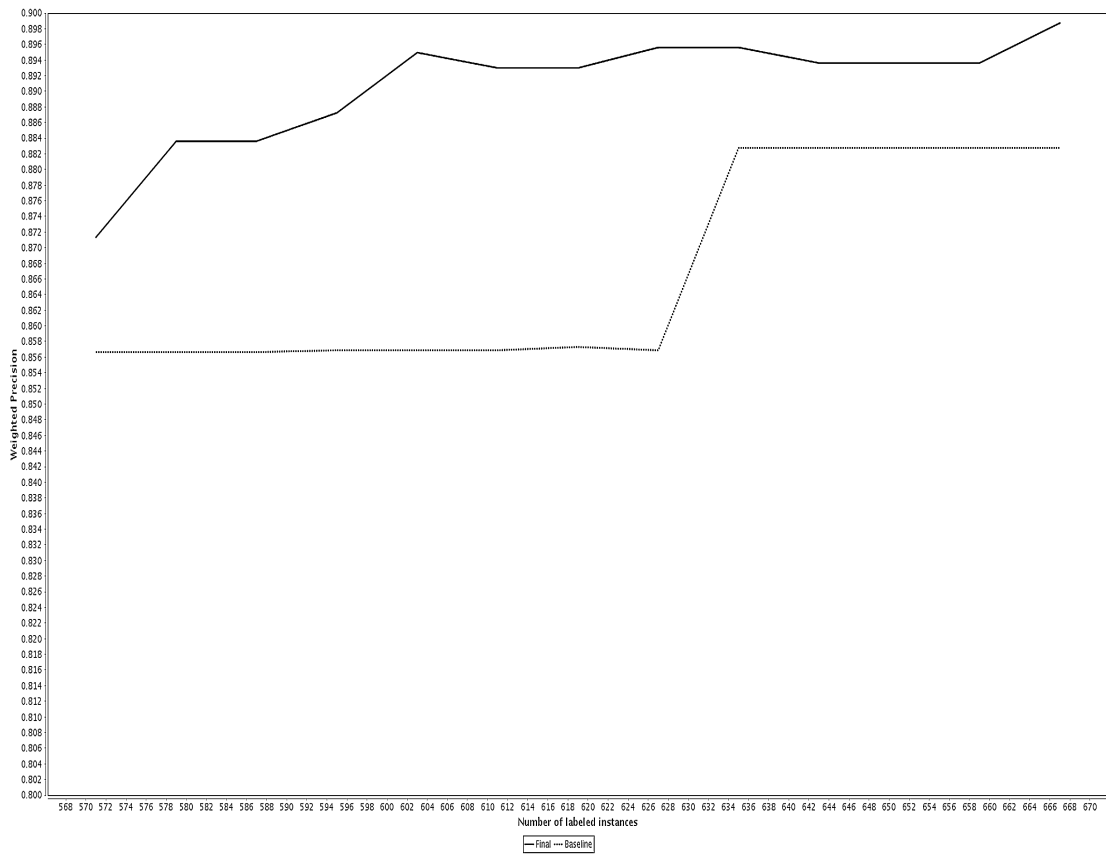


Figure 3.2: Learning curve for the precision (output of the JCLAL library).

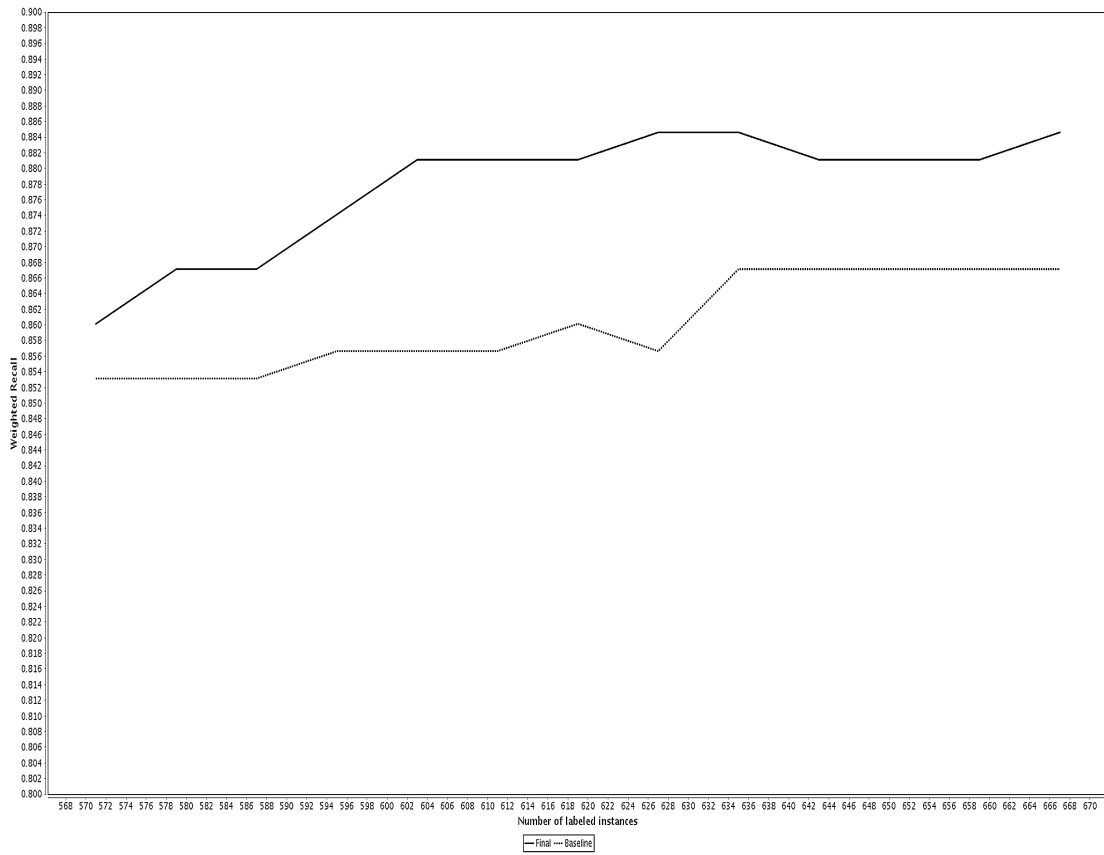


Figure 3.3: Learning curve for the recall (output of the JCLAL library).

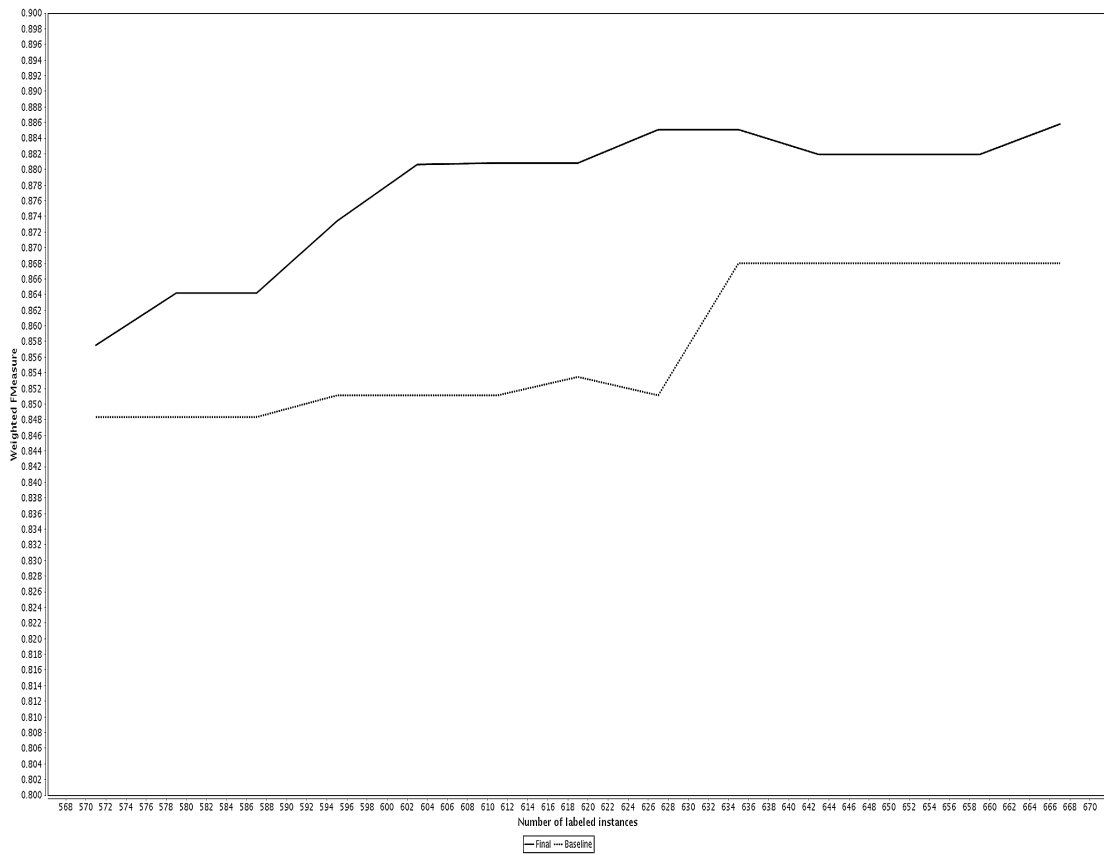


Figure 3.4: Learning curve for the F_1 -score (output of the JCLAL library).

3.7 Summary

This chapter presented the task of classifying non-sentential utterances and our approach to address this problem. This task is formulated as a machine learning problem and we follow and extend the work of Fernández et al. (2007). We use their corpus as a gold-standard and a replica of their approach as a baseline. The data, the machine learning algorithm used and the feature set of the baseline were discussed respectively in Section 3.1, 3.2, 3.3. The two main problems we faced in our work have been the scarcity of the labeled data and the imbalance in the distribution of the classes. To address these problems we extended the baseline approach in two ways: using a larger feature set (detailed in Section 3.4) and employing semi-supervised learning techniques to exploit the abundance of unlabeled data. We described in Section 3.5 the semi-supervised learning techniques that we employed, namely: Self Training, Transductive SVM and Active Learning. Section 3.6 shows the empirical results we got from our experiments. While the extended feature set alone did not make an improvement on the performances of the classifiers, its use in combination with Active Learning made a modest but significant difference.

Chapter 4

Resolution of Non-Sentential Utterances

As introduced in Chapter 2, the resolution of an NSU is the task of reconstructing its meaning from the dialogue context. Fernández (2006) proposes a set of rules to resolve NSUs based on TTR, the logical framework from Cooper (2004) further developed then in Ginzburg (2012). One limitation of logical frameworks such as TTR is their inability to directly represent (and reason over) uncertain knowledge. Moreover, many dialogue domains contain variables that are only partially observable. We have to take into account a certain degree of stochastic behavior when modeling dialogue since we still have an imperfect understanding of its dynamics. The stochastic component is especially important in dealing with NSUs since they do not have a precise meaning by themselves and, as argued in Ginzburg (2012), they are in principle highly ambiguous.

For this reason we propose a new approach to the resolution of NSUs that takes probabilistic account of the variables involved and the procedures used. We employ the probabilistic rules formalism of Lison (2015) (detailed in Section 2.3) to encode the NSU resolution procedures as probabilistic rules. Probabilistic rules are similar, to a certain extent, to the update rules developed by Ginzburg (2012) (described in Section 2.2). For this reason probabilistic rules are particularly suited for our purpose since, in this way, we could reuse many theoretical aspects from Ginzburg (2012) and Fernández (2006). We reinterpreted the variables in the dialogue state as random variables and straightforwardly “converted” the resolution rules into probabilistic rules.

In the next sections we explain how we represented the variables in the dialogue state and how we translated the rules from Fernández (2006) into probabilistic rules. First we describe the theoretical aspects from Ginzburg (2012) we employ in our approach. We present then the design of our dialogue context and the rules to resolve the NSUs.

We surely take a much simpler approach than Ginzburg (2012) in the modeling of the dialogue state, abstracting intentionally from many details that would add complexity to the modeling. Indeed there are a number of issues that arise in the resolution of NSUs

that need to be treated with proper lexical and semantic resources that we did not include. However, in the end our goal for this work is not to formulate a complete theory of NSU resolution but rather to provide a *proof-of-concept* implementation for the resolution of the NSUs in the dialogue context with the probabilistic rules formalism.

This framework has also been implemented and tested with the OpenDial toolkit. We developed a dialogue system able to update the dialogue state probabilistically with update rules similar to the ones from Ginzburg (2012). The system is also able to resolve toy examples in an interactive way. A detailed example of the behavior of the system is given in Section 4.5. More high-level details about the rules for the state update (complementary to the rules for the NSU resolution) can be found in the Appendix A¹.

4.1 The resolution task

The resolution of an NSU is the task of extracting its meaning from the dialogue context. More precisely, let u_a and nsu_a represent respectively the word word sequence making up the NSU and its type according to the taxonomy presented in Section 2.1.1. We also assume MaxQUD to be a high-level semantic representation of the antecedent, as mentioned in Section 2.2.2. Through a resolution procedure, we want to extract a_a i.e. the high-level semantic representation of the NSU. The right resolution procedure is selected on the basis of the type of the NSU. In our case the value of nsu_a is retrieved using the classifier developed in Chapter 3 which takes as input the raw NSU and the antecedent. Figure 4.1 shows a schema of the task just defined. Indeed this is the simplest way to define the task. The resolution procedure may also be dependent of other variables in the dialogue state such as the Facts. In principle, the resolution task is defined independently from the actual semantic representation of the utterances. It is also defined independently from the rules used to update the variables in the dialogue state such as QUD and Facts. In practice, define a set of rules that are generic enough to handle every possible case and behave independently from the state update rules is a difficult task and still an open research problem.

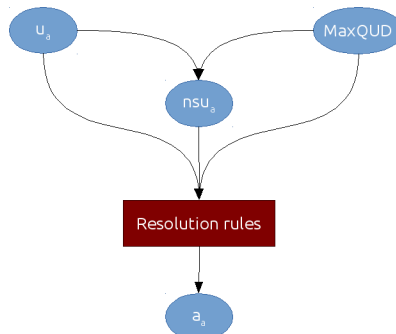


Figure 4.1: The basic schema for the NSU resolution task.

¹For more technical details about the implementation and examples of interaction visit: <https://github.com/paolodragone>

4.2 Theoretical foundation

As previously stated, we rely on Fernández (2006) and Ginzburg (2012) for the theoretical notions needed to represent the dialogue state and to develop the NSU resolution rules. In Section 2.2 we detailed the basic concepts of TTR, the utterance representation and the update rules for the dialogue state. In this section we describe the notions needed for the resolution of the NSUs. In particular we describe how we can exploit the *parallelism* between the NSU and its antecedent that we mentioned in Section 2.1. We discuss here the concepts that Ginzburg (2012) defines to address the resolution of NSUs then we will describe how we adapt those concepts to our needs in the next section.

4.2.1 Partial Parallelism

Instances of NSU classes such as Acknowledgment and Affirmative Answers are related to their antecedent as a whole, that is to understand their meaning one has to consider not a specific aspect of the antecedent but the entire sentence. On the other hand, there are NSU classes, such as Short Answers and Sluices, that show a more fine grained *parallelism* between their instances and their antecedents i.e. they may refer in particular to certain aspects of the antecedent. In the theory of Ginzburg (2012), this concept is named Partial Parallelism². Partial Parallelism is one way to categorize NSU classes according to the relation with their antecedents. NSU classes are categorized as +/-ParPar in order to find the right way to treat them. An NSU class categorized as +ParPar involves the access to one or more sub-utterances from its antecedent. On the contrary, -ParPar NSU classes do not need to know the internal structure of their antecedents to be resolved. Table 4.1 shows how NSU classes are categorized in this way.

-ParPar	+ParPar
Plain Acknowledgment	Short Answer
Plain Affirmative Answer	Repeated Acknowledgment
Plain Rejection	Clarification Ellipsis
Factual Modifier	Repeated Affirmative Answer
Check Question	Sluice
Propositional Modifier	Helpful Rejection
	Filler
	Bare Modifier Phrase
	Conjunct

Table 4.1: An overview of the NSU classes divided according to Partial Parallelism.

²Fernández (2006) previously addressed this concept as Sentential Antecedent (SA).

4.2.2 Propositional lexemes

-ParPar NSU classes are realized (mainly) by *propositional lexemes*, i.e. words that can stand alone and form a proposition with full contextual meaning. Among those classes there are the Plain Affirmative Answers, Plain Rejections and Propositional Modifiers which respectively are realized by the words *yes*, *no* and adverbials such as *probably* and *possibly*.

Those classes of NSU arise from polar questions such as (4.1).

- (4.1) A: Will you go to the party on Saturday?
B: **Yes.** / **No.** / **Probably.**

The semantic content of these stand-alone lexemes can be modeled as a function R of the content of the antecedent polar question.

For Plain Affirmative Answers, R is the *Identity* (Id) relation, i.e. the function that returns the argument itself. This means that the positive answer “yes” to a polar question is equivalent to the assertion of a proposition with the same content as the polar question.

For Plain Rejections, R is the relation Neg . Neg indicates the negation of a proposition p although it is sensitive to the polarity of p , meaning that, when p is positive, $Neg(p)$ is the negation of p (denoted with \bar{p}) whereas, when p is negative, $Neg(p)$ is p itself. This rule is needed to account for the asymmetry in the meaning of negative answers to negative questions. A negative answer to a negative question does not equate a positive one, as exemplified in (4.2) (rephrased from Ginzburg (2012)).

- (4.2) A: Did Paul not leave?
B: No. (= Paul did not leave.)

For Propositional Modifiers, R is a relation $PropRel$ which applies different modalities on the basis of the lexical meaning of the word used as modifier, e.g. “probably” would have a different modality than “clearly”.

4.2.3 Focus Establishing Constituents

To account for the partial parallelism between NSUs and their antecedents stemming out from the instances of the classes of the +ParPar group we need to keep track of the *focal* sub-utterances of the antecedents i.e. of the elements of QUD. For this reason we employ the notion of *focus establishing constituents* (FEC) from the theory of Ginzburg (2012)³. The FECs are relevant constituents in the elements of QUD that may be used to resolve NSUs. Consider the following example:

- (4.3) A: **A friend** is coming to the party.
B: **Who?**

³The concept was previously formalized by Fernández (2006) as *topical constituents*.

The noun phrase “A friend” in the first sentence of (4.3) is the one which the following Sluice is referring to. Roughly the Sluice can be resolved in such a manner: “*Who is your friend that is coming to the party?*”. It is clear that the aforementioned sub-utterance has to be contextually available to allow the resolution of the subsequent Sluice. In this we follow Ginzburg (2012), who defines a set of rules to follow to make FECs contextually available. In particular we are interested in the following ones:

- The FEC associated with a *wh*-interrogative is the *wh*-phrase⁴ itself:

(4.4) A: **Who** is organizing the party?
 B: Paul.

- The FEC associated with a polar interrogative or declarative utterance can be any (quantified) noun phrases:

(4.5) A: **A friend** is organizing a party and **many people** are coming.
 B: Who?

- The FEC associated with a clarification request is the sub-utterance that has to be clarified i.e. any sub-utterance in the antecedent:

(4.6) A: Is Paul organizing a party?
 B: Paul? / Organizing? / A party?

4.2.4 Understanding and acceptance

The classes of Plain Acknowledgments and Check Questions are used to handle understanding and acceptance in the conversation. Plain Acknowledgments are used to send a direct feedback of understanding or acceptance of the previous utterances. Understanding involves grasping successfully the content of an utterance while acceptance is a sign of shared belief which therefore updates the Facts with the accepted utterance and removes the corresponding issue from the QUD. As argued in Fernández (2006), understanding does not always imply acceptance, and Plain Acknowledgments are ambiguous in this distinction. Despite this difference, we assume that Plain Acknowledgments are used to show acceptance, therefore the use of a Plain Acknowledgments also downdates the QUD. On the other hand, understanding is assumed to be shown by any utterance that is not a Clarification Ellipsis.

Check Questions are used in conversation to request an explicit feedback about the understanding/acceptance of the previous utterance.

⁴As in Ginzburg (2012) we consider only unary *wh*-interrogatives. Refer to Fernández (2006) for an account of utterances with multiple *wh*-interrogatives

4.2.5 Sluicing

Sluices can take a wide range of meanings depending on the particular situation. To formalize the meaning of Sluices, Fernández et al. (2007) distinguish four types of Sluices that convey different meanings: *Direct Sluices*, *Reprise Sluices*, *Clarification Sluices*, *Wh-anaphor*.

The aforementioned paper describes a machine learning experiment to automatically classify Sluices according to these types. Ginzburg (2012) describes several different treatments for every group of Sluices.

In our work we do not distinguish between those type of Sluices but we confine ourselves for simplicity to direct Sluices only. Direct Sluices, such as the one in (4.7), are used to query the other speaker for additional information about some aspect of the antecedent.

- (4.7) A: Can I have some toast please?
 B: Which sort?
 [BNC: KCH 104–105]

4.3 Dialogue context design

As mentioned before, the dialogue context is represented as a Bayesian network containing a set of random variables representing the current information state. The values of those random variables can represent virtually anything, from the raw utterances to their semantic representation. The variables in the dialogue context are inspired by Ginzburg (2012). In order to make the transition from the rules of Fernández (2006) to probabilistic rules as direct as possible, we mimic the basic dynamic of the DGB detailed in Section 2.2. For our semantics we do not employ TTR because it would add unnecessary complexity to our formalization. In this section we first describe the semantics we adopt and then we discuss the random variables that compose the dialogue context.

4.3.1 Semantics

The semantic content of the utterance is represented by logical predicates, individuals and variables. Predicates are labeled as words or camel-case phrases and can present zero or more arguments. Individuals are labeled with uppercase abbreviations such as IND for generic individuals or E for events. Variables are labeled with an uppercase X. Both individuals and variables are uniquely identified by a numeric subscript.

Predicates represent the high-level semantic meaning of the constituents of the utterances. Intuitively, predicates without variables as argument can represent propositions such as (4.8). As discussed in Section 2.2.1, polar questions and *wh*-questions can be seen as functions from/to record types. Polar questions take as argument the empty record type.

Following this schema, in our formalism polar questions are denoted by predicates with no variables, whereas *wh*-questions are denoted by predicates containing one or more variables, as exemplified by (4.9).

(4.8) Paul is a friend of yours. friend(addr, Paul)

(4.9) Is Paul a friend of yours? friend(addr, Paul)
 Who is your friend? friend(addr, X₁)

Retrieving the semantic representation from the raw utterances is a Natural Language Understanding (NLU) task, a completely different task with respect to the resolution of NSUs. We do not attempt to generate predicates from raw utterances instead we make use of simple handcrafted predicates in our examples, abstracting the necessity of NLU to retrieve the meaning of all the utterances that are not NSUs. We try to keep the problem of NSU resolution generic separating it as much as possible from the NLU task.

4.3.2 Dialogue acts

As seen in Section 2.2.1, to represent the “purpose” of an utterance, we need to use an *illocutionary relation*, also known as dialogue act. The set of dialogue acts we employ in our dialogue context is a small subset of the ones defined by Ginzburg (2012):

- *Assert*, denoting the act of asserting a proposition;
- *Ask*, denoting the act of posing a question;
- *Ground*, denoting the act of understanding what being previously said;
- *Accept*, denoting the act of accepting what being previously said.

Assertions are applied to propositions and they are implicitly considered truthful unless they violate some predicates in the Facts. Asking a query is the act of posing questions and they are piled up in the QUD until they are resolved by an answer. The act of answering to a query corresponds, in the case of a *wh*-interrogative, to finding the valid arguments to the variables of the question. In case of a polar question, the answer is derived simply by its truth status, denoted by the presence of the same predicate in the Facts. In our formalization we use “Ground” to represent the act understanding. Acceptance is the act of resolving an issue, which involves updating Facts and downdating QUD.

4.3.3 Variables of the dialogue context

For our formalization, as in TTR, we assume the availability of various data structures such as variables, lists, sets and complex types. The probabilistic rules formalism provides those structures out of the box as possible values for the random variables. Random variables are denoted with the notation “**var**₁”. Array element accesses are indicated with the square brackets notation, such as “**array**[0]”. Sets are denoted with the classical mathematical

notation “ $\{e_1, \dots, e_n\}$ ”. Complex type accesses are denoted with the dot notation, such as “`complexVar.var1`”. The classical operations on sets are available such as the union and the intersection. Array concatenation is denoted with the $+$ symbol. Now we describe the variables used in our formalization of the dialogue context.

`ua, ub, aa, ab`

As a convention, raw utterances and dialogue acts are indicated respectively with the letters u and a and a subscript denotes the speaker. We record in separate variables only the last utterance and dialogue act of each speaker.

`nsua`

A random variable that contains the distribution over the NSU classes returned by the classifier for the latest recorded utterance. It uses `max-qud` to refer to the antecedent therefore the probabilistic inference framework takes care of finding the most probable antecedent for the current NSU. Besides the values of the NSU classes, a distinct value `NoNsu` is used to account for input utterances that are not NSUs. To determine whether an utterance is an NSU or not we used the same detection rules explained in Section 3.5.1.

`new-fec`

The set of FECs introduced by the NLU of the last recorded utterance. It is also a buffer variable used in the NSU resolution to encode the focal constituents of the newly resolved NSU. It is used also to hold FECs of the utterance that is being inserted in the `qud`.

`facts`

A set of predicates representing the common knowledge of the users. The predicates in `facts` contain only individuals as arguments (i.e. no variables) and they are implicitly considered truthful.

`qud`

As defined in Section 2.2, the QUD is a partially ordered set containing the issues currently under discussion. Its ordering determines the “priority” of the issues to be resolved. Here instead `qud` is represented as a vector and the `max-qud` variable denotes the index of the MaxQUD element (see below). Each element in `qud` has a number of sub-fields:

- `utt`: The raw utterance associated to the current question under discussion;
- `q`: The semantic representation of the utterance;
- `fec`: An array of topical sub-utterances used in the resolution of the NSUs.

The `qud` is incremented by adding elements in the tail (growing numbers) and decremented in a random-access fashion, usually by removing the MaxQUD element (which could not be the last element) after its resolution. We denote with `qudsize` the size of `qud`.

`max-qud`

Despite being represented as the maximal element of QUD in Ginzburg (2012), here `max-qud` actually denotes the index of such element which therefore is retrieved in this way: `qud[max-qud]`. In Ginzburg (2012), MaxQUD is given from the partial ordering imposed on QUD. This ordering is often similar but not limited to the behavior of a stack. At our disposal we have the full power of probabilistic modeling which enables us to encode `max-qud` as a random variable with a prior that gives more probability to the highest element in `qud`. The function used to give this prior to `max-qud` is

$$P(\text{max-qud} = i) = e^{i - \text{qud}_{\text{size}}}$$

where $i < \text{qud}_{\text{size}}$ is an index in `qud`. In this way, the prior most probable MaxQUD is the last element inserted in the QUD but the probability can be modified by other contextual elements by probabilistic inference on the dialogue state.

4.4 NSU resolution rules

Here we present the probabilistic rules that handle the resolution of NSUs. For each rule we also present an example of usage. Since they are a (almost) direct translation of the deterministic rules from Fernández (2006), most of them have deterministic effects (i.e. a single effect with probability 1). Nonetheless the updates are handled probabilistically by the probabilistic rules framework through probabilistic inference over the Bayesian network representing the dialogue state. We show an example of probabilistic update in Section 4.4.1, which is valid for every other resolution rule.

4.4.1 Acknowledgments

The only requirement for Acknowledgment resolution is to have at least one issue under discussion to be accepted. As explained in Section 4.2.4, we assume that an explicit Acknowledgment is a sign of acceptance of the latest issue under discussion. For Repeated Acknowledgments, Fernández (2006) requires to have co-referentiality between the repeated constituent in the NSU and the relative constituent in the FEC of MaxQUD. We decided to drop this requirement assuming that the co-reference is always present when the classifier assigns the class *RepAck* to the current NSU. This assumption does not affect the system given that the effect on the state variables is the same for both *Acks* and *RepAcks*. The rule for Acknowledgments is the following⁵:

ack :

if $((\text{nsu}_a = \text{Ack} \vee \text{nsu}_a = \text{RepAck}) \wedge \text{max-qud} > 0)$ **then**

$$\left\{ P(\mathbf{a}_a \leftarrow \text{Accept}()) = 1 \right.$$

⁵The symbol \leftarrow indicates the assignment of the right-hand side value to the left-hand side variable.

Consider the following example.

(4.10) B: I am going to the party.
A: OK.

The dialogue context of (4.10) is:

$$\begin{aligned} \text{max-qud} &= 1 \\ \text{qud}[\text{max-qud}].\mathbf{q} &= \text{goingToParty}(\text{IND}_1) \\ \text{nsu}_a &= \text{Ack} \end{aligned}$$

After the application of the rule:

$$\mathbf{a}_a = \text{Accept}()$$

Notice that this may be an oversimplification since often the values of the variables in the dialogue state are not determined with full probability but rather the variables encode a probability distribution over a set of values. For instance, it is often the case that the classifier will retrieve the type of the NSU in a probability distribution with one value with large probability and other few values with smaller probability scores. In this case we could have a situation resembling the following:

$$\text{nsu}_a = \begin{cases} \text{Ack} & \text{with probability } 0.75 \\ \text{AffAns} & \text{with probability } 0.2 \\ \text{CheckQu} & \text{with probability } 0.05 \end{cases}$$

The case above would result in the following distribution of \mathbf{a}_a ⁶:

$$\mathbf{a}_a = \begin{cases} \text{Accept}() & \text{with probability } 0.75 \\ \text{None} & \text{with probability } 0.25 \end{cases}$$

Since the dialogue state is a Bayesian network, the update rules will return a distribution of values that is dependent on both the distribution assigned by the rule (in this case only one value with full probability) and the distributions of the variables the rule depend on. These considerations can be of course extended to all the other classes so in the following sections we will only point out the most relevant use cases.

4.4.2 Affirmative Answers

The context for an Affirmative Answer contains a polar question $q(\mathbf{y})$ as MaxQUD. As for the Acknowledgments, we drop the requirement of co-referentiability between the repeated constituent of the *RepAffAns* and the same constituent in the FECs of the MaxQUD element.

⁶The actual distribution would not necessarily assign None as alternative value because other rules may be triggered by the other values of nsu_a .

An Affirmative Answer to a polar question corresponds to asserting the same semantic content (predicate) of the question. The following is the rule to handle Affirmative Answers^{7,8}:

affAns :

$\forall q, \mathbf{y}$

if $((\mathbf{nsu}_a = \mathit{AffAns} \vee \mathbf{nsu}_a = \mathit{RepAffAns}) \wedge \mathbf{qud}[\mathbf{max-qud}] . \mathbf{q} = q(\mathbf{y}))$ **then**

$$\left\{ P \left(\begin{array}{l} \mathbf{a}_a \leftarrow \mathit{Assert}(q(\mathbf{y})), \\ \mathbf{new-fec} \leftarrow \mathbf{qud}[\mathbf{max-qud}] . \mathbf{fec} \end{array} \right) = 1 \right.$$

An example of application of the *affAns* rule can be:

(4.11) B: Are you going to the party?
 A: Yes.

The context of (4.11) is the following:

$$\begin{array}{ll} \mathbf{max-qud} & = 1 \\ \mathbf{qud}[\mathbf{max-qud}] . \mathbf{q} & = \mathit{goingToParty}(\mathbf{IND}_2) \\ \mathbf{qud}[\mathbf{max-qud}] . \mathbf{fec} & = \{\} \\ \mathbf{nsu}_a & = \mathit{AffAns} \end{array}$$

After the application of the rule:

$$\begin{array}{ll} \mathbf{a}_a & = \mathit{Assert}(\mathit{goingToParty}(\mathbf{IND}_2)) \\ \mathbf{new-fec} & = \{\} \end{array}$$

4.4.3 Rejections

As for the Affirmative Answers, the context of Rejections is a polar question $q(\mathbf{y})$, but, as explained in Section 4.2.2, we need to distinguish the cases in which q is positive or negative. We will define the following function *Neg* indicating the negation of a proposition p (or equivalently a question):

$$\mathit{Neg}(p) = \begin{cases} \bar{p} & \text{if } p \text{ is positive} \\ p & \text{if } p \text{ is negative} \end{cases}$$

where \bar{p} is the negative of p . As an extension of the above notation, we indicate a proposition that is explicitly negative as \bar{p} .

⁷As a convention, quantified variables and quantified individuals in the rule definitions are indicated respectively as x and y . Vectors of variables or individuals are indicated respectively as \mathbf{x} and \mathbf{y} .

⁸As in this case, a probabilistic effect might contain several assignments. Hereafter, for readability, we write the sequence of assignments in a vertical notation.

Rejections are handled by the following rule:

reject :

$$\forall q, \mathbf{y}$$

if ($\text{nsu}_a = \text{Reject} \wedge \text{qud}[\text{max-qud}] . \mathbf{q} = q(\mathbf{y})$) **then**

$$\left\{ P \left(\begin{array}{l} \mathbf{a}_a \leftarrow \text{Assert}(\text{Neg}(q)(\mathbf{y})), \\ \text{new-fec} \leftarrow \text{qud}[\text{max-qud}] . \text{fec} \end{array} \right) = 1 \right.$$

else if ($\text{nsu}_a = \text{Reject} \wedge \text{qud}[\text{max-qud}] . \mathbf{q} = \bar{q}(\mathbf{y})$) **then**

$$\left\{ P \left(\begin{array}{l} \mathbf{a}_a \leftarrow \text{Assert}(\bar{q}(\mathbf{y})), \\ \text{new-fec} \leftarrow \text{qud}[\text{max-qud}] . \text{fec} \end{array} \right) = 1 \right.$$

The following is an example for the above rule:

(4.12) B: Are you going to the party?

A: No.

The context of (4.12) is the following:

$$\begin{array}{ll} \text{max-qud} & = 1 \\ \text{qud}[\text{max-qud}] . \mathbf{q} & = \text{goingToParty}(\text{IND}_2) \\ \text{qud}[\text{max-qud}] . \text{fec} & = \{\} \\ \text{nsu}_a & = \text{Reject} \end{array}$$

After the application of the rule:

$$\begin{array}{ll} \mathbf{a}_a & = \text{Assert}(\text{Neg}(\text{goingToParty})(\text{IND}_2)) \\ \text{new-fec} & = \{\} \end{array}$$

4.4.4 Propositional Modifiers

As Affirmative Answers and Rejections, Propositional Modifiers are triggered by polar questions. As seen in Section 4.2.2, their resolution corresponds to asserting the predicate of the polar question modified by a certain modality given by the lexical meaning of the NSU itself.

We define the function $\text{PropRel}_M(p)$ that modifies the meaning of a proposition p (or equivalently a question) with the modality M . The modality is given by the lexical meaning of the word used in the NSU, here indicated for simplicity as the word itself (contained in the variable \mathbf{u}_a).

The rule for Propositional Modifiers states:

$$\begin{aligned}
& \text{propMod} : \\
& \forall q, \mathbf{y} \\
& \text{if } (\text{nsu}_a = \text{PropMod} \wedge \text{qud}[\text{max-qud}] . \mathbf{q} = q(\mathbf{y})) \text{ then} \\
& \left\{ P \left(\begin{array}{l} \mathbf{a}_a \leftarrow \text{Assert}(\text{PropRel}_{\mathbf{u}_a}(q)(\mathbf{y})), \\ \text{new-fec} \leftarrow \text{qud}[\text{max-qud}] . \text{fec} \end{array} \right) = 1 \right.
\end{aligned}$$

Here is an example of application of the above rule:

- (4.13) B: Are you going to the party?
A: Probably.

The dialogue state of (4.13) before the application of the rule is the following:

$$\begin{aligned}
\text{max-qud} & = 1 \\
\text{qud}[\text{max-qud}] . \mathbf{q} & = \text{goingToParty}(\text{IND}_2) \\
\text{qud}[\text{max-qud}] . \text{fec} & = \{\} \\
\text{nsu}_a & = \text{PropMod}
\end{aligned}$$

After the application of the rule we would have:

$$\begin{aligned}
\mathbf{a}_a & = \text{Assert}(\text{PropRel}_{\text{probably}}(\text{goingToParty})(\text{IND}_2)) \\
\text{new-fec} & = \{\}
\end{aligned}$$

Conversely to Affirmative Answers and Rejections, the Propositional Modifiers need to take into account the lexical meaning of the modifier used to update the dialogue state accordingly. This requires a set of lexicalized update rules to properly react to each possible modality of the modified proposition. However, these rules will only take place at the level of action selection and context update therefore it is still possible to resolve this kind of NSUs in a general way, as previously explained in Section 4.2.2.

An example of lexicalized rule for updating the context in the presence of a modified proposition can be the following.

$$\begin{aligned}
& \text{factsIncrement}_{\text{PropRel}} : \\
& \forall p, \mathbf{y} \\
& \text{if } (\mathbf{a}_b = \text{Accept}(\text{PropRel}_{\text{probably}}(p)(\mathbf{y}))) \text{ then} \\
& \left\{ P(\text{facts} \leftarrow \text{facts} \cup \{p(\mathbf{y})\} \cup \text{new-fec}) = 0.75 \right. \\
& \text{else if } (\mathbf{a}_b = \text{Accept}(\text{PropRel}_{\text{unlikely}}(p)(\mathbf{y}))) \text{ then} \\
& \left. \left\{ P(\text{facts} \leftarrow \text{facts} \cup \{p(\mathbf{y})\} \cup \text{new-fec}) = 0.25 \right. \right.
\end{aligned}$$

This rule handles the update of the **facts** variable when the addressee decides to accept a proposition modified by a “probably” relation or by an “unlikely” relation. The latter is realized by updating **facts** with a high probability while the former updates **facts** with a low probability.

While the above rule has handcrafted probabilities, they can in principle be learned from actual data. Of course this would be only possible given a corpus containing NSU instances annotated with the dialogue acts and state updates at each step. When an instance of Propositional Modifier is encountered, the probabilities of the effects are updated according to the relative state update move.

4.4.5 Check Questions

As defined in Section 4.2.4, Check Questions are used to ask for understanding/acceptance of the latest issue being raised. In practice this means asking the latest asserted proposition as a polar question. The following is the rule to handle this type of NSUs:

$$\begin{aligned}
 & \textit{checkQu} : \\
 & \forall p, \mathbf{y} \\
 & \text{if } (\textit{nsu}_a = \textit{CheckQu} \wedge \textit{qud}[\textit{max-qud}] . \mathbf{q} = p(\mathbf{y})) \text{ then} \\
 & \left\{ P \left(\begin{array}{l} \mathbf{a}_a \leftarrow \text{Ask}(p(\mathbf{y})), \\ \textit{new-fec} \leftarrow \textit{qud}[\textit{max-qud}] . \textit{fec} \end{array} \right) = 1 \right.
 \end{aligned}$$

An example of application of the previous rule is the following:

- (4.14) A: I am going to the party.
 A: OK?

The dialogue context of (4.14) is:

$$\begin{aligned}
 \textit{max-qud} & = 1 \\
 \textit{qud}[\textit{max-qud}] . \mathbf{q} & = \text{goingToParty}(\text{IND}_1) \\
 \textit{qud}[\textit{max-qud}] . \textit{fec} & = \{\} \\
 \textit{nsu}_a & = \textit{CheckQu}
 \end{aligned}$$

After the application of the rule:

$$\begin{aligned}
 \mathbf{a}_a & = \text{Ask}(\text{goingToParty}(\text{IND}_1)) \\
 \textit{new-fec} & = \{\}
 \end{aligned}$$

4.4.6 Short Answers

The antecedent of Short Answers is assumed to be a *wh*-question. As stated previously, in this work we limit ourselves to unary *wh*-interrogatives i.e. questions with only one unknown variable x . Short answers are resolved by applying them to the MaxQUD *wh*-interrogative then asserting the resulting proposition. The application of the Short Answer is done by substituting every occurrences of the variable x with u_a (or equivalently a high-level representation of it). The following is the rule for Short Answers:

shortAns :

$$\begin{aligned} & \forall q, x, \mathbf{y}, p_i, \mathbf{y}_i \\ & \mathbf{if} (\mathbf{nsu}_a = \mathit{ShortAns} \wedge \mathbf{qud}[\mathbf{max-qud}].\mathbf{q} = q(x, \mathbf{y}) \wedge \\ & \quad \{p_1(x, \mathbf{y}_1), \dots, p_n(x, \mathbf{y}_n)\} \subseteq \mathbf{qud}[\mathbf{max-qud}].\mathbf{fec}) \mathbf{then} \\ & \quad \left\{ \left(\begin{array}{l} \mathbf{a}_a \leftarrow \mathit{Assert}(q(\mathbf{u}_a, \mathbf{y})), \\ \mathbf{new-fec} \leftarrow \{p_1(\mathbf{u}_a, \mathbf{y}_1), \dots, p_n(\mathbf{u}_a, \mathbf{y}_n)\} \end{array} \right) = 1 \right. \end{aligned}$$

An example of use of this rule is:

- (4.15) B: Who is your friend organizing the party?
A: Paul.

To make the example more appropriate we added a constituent that will be resolved by the rule together with the dialogue act. The context of (4.15) is:

$$\begin{aligned} \mathbf{max-qud} & = 1 \\ \mathbf{qud}[\mathbf{max-qud}].\mathbf{q} & = \mathit{organizingTheParty}(X_1) \\ \mathbf{qud}[\mathbf{max-qud}].\mathbf{fec} & = \{\mathit{friend}(\mathit{IND}_2, X_1)\} \\ \mathbf{nsu}_a & = \mathit{ShortAns} \end{aligned}$$

After the application of the rule:

$$\begin{aligned} \mathbf{a}_a & = \mathit{Assert}(\mathit{organizingTheParty}(\mathit{Paul})) \\ \mathbf{new-fec} & = \{\mathit{friend}(\mathit{IND}_2, \mathit{Paul})\} \end{aligned}$$

4.4.7 Sluices

As argued in Section 4.2.5, we limit ourselves to the treatment of direct Sluices. Even for this type of Sluices only, the resolution rules are many since they have to account for the lexical meaning of each *wh*-word. Furthermore, the meaning of *wh*-words can be modified in many ways, e.g. “how many”, “how long”, “who else”, “what about”. This would require an extensive treatment for this kind of NSUs that we do not attempt to elaborate.

Nevertheless we will show some rules to treat simple direct Sluices like the ones in (4.16) and (4.17).

(4.16) B: A friend is coming to the party.
A: Who?

(4.17) B: Paul is throwing a party.
A: When?

The context of the Sluices is a MaxQUD with at least one variable (raised by either a *wh*-question or a proposition with some undefined reference). The Sluice is used to request some kind of information regarding one of the FECs of the antecedent. The requested information as well as the context generated by the resolution depend on the lexical meaning of the *wh*-word. For instance, the following rule treats the Sluice “*who?*”.

sluice_{who} :

$\forall q, x, \mathbf{y}, p_i, \mathbf{y}_i$

if ($\text{nsu}_a = \text{Sluice} \wedge \text{“who”} \in \mathbf{u}_a \wedge \text{qud}[\text{max-qud}].\mathbf{q} = q(x, \mathbf{y}) \wedge$

$\{p_1(x, \mathbf{y}_1), \dots, p_n(x, \mathbf{y}_n)\} \subseteq \text{qud}[\text{max-qud}].\text{fec}$) **then**

$$\left\{ P \left(\begin{array}{l} \mathbf{a}_a \leftarrow \text{Ask}(\text{named}(x, \hat{x})), \\ \text{new-fec} \leftarrow \{p_1(x, \mathbf{y}_1), \dots, p_n(x, \mathbf{y}_n)\} \cup \{\text{person}(x)\} \end{array} \right) = 1 \right.$$

where \hat{x} is a newly created variable. Such a Sluice asks about the identity (here simplified by the name) of a *person* which is referred to in the antecedent.

We can use as example the transcript (4.16). The context before the application of the rule is:

$$\begin{array}{ll} \text{max-qud} & = 1 \\ \text{qud}[\text{max-qud}].\mathbf{q} & = \text{comingToParty}(X_1) \\ \text{qud}[\text{max-qud}].\text{fec} & = \{\text{friend}(\text{IND}_2, X_1)\} \\ \text{nsu}_a & = \text{Sluice} \end{array}$$

After the application of the rule we have:

$$\begin{array}{ll} \mathbf{a}_a & = \text{Ask}(\text{named}(X_1)) \\ \text{new-fec} & = \{\text{friend}(\text{IND}_2, X_1), \text{person}(X_1)\} \end{array}$$

An interesting result of this rule in combination with the probabilistic inference employed in OpenDial is how the ambiguity in the FECs is handled. As argued in Ginzburg (2012), the antecedent of a Sluice can contain more than one potential FEC as exemplified by the following transcript.

(4.18) B: **A friend of mine** is organizing a party for **his girlfriend**.
A: Who? (= Who is your friend? / Who is his girlfriend?)

The resolution of this kind of ambiguities is automatically handled in a probabilistic fashion. In (4.18) the representation of the antecedent (MaxQUD) would be the following:

$$\begin{aligned} \text{qud}[\text{max-qud}].\text{q} &= \text{organizingPartyFor}(X_1, X_2) \\ \text{qud}[\text{max-qud}].\text{fec} &= \{\text{friend}(\text{IND}_1, X_1), \text{girlfriend}(X_2, X_1)\} \end{aligned}$$

In this case, the above rule would be applied to both variables X_1 and X_2 , resulting in two possible assignments of \mathbf{a}_a and new-fec with 0.5 probability:

$$\begin{aligned} \mathbf{a}_a &= \begin{cases} \text{Ask}(\text{named}(X_1, X_3)) & \text{with probability 0.5} \\ \text{Ask}(\text{named}(X_2, X_3)) & \text{with probability 0.5} \end{cases} \\ \text{new-fec} &= \begin{cases} \{\text{friend}(\text{IND}_1, X_1), \text{person}(X_1)\} & \text{with probability 0.5} \\ \{\text{girlfriend}(X_2, X_1), \text{person}(X_2)\} & \text{with probability 0.5} \end{cases} \end{aligned}$$

Without any prior, the probabilistic inference over the dialogue state would assign equal probability to each possible assignment of \mathbf{a}_a . A more sophisticated approach may use some notion of *saliency* as a prior to adjust the probabilities of each focal constituent. For example one could adjust the probability according to whether the constituent is a subject or an object in the antecedent. In the previous example the friend would have had more probability mass (e.g. 0.8) and the girlfriend less probability mass (e.g. 0.2). It is also possible that the prior saliency could depend on other variables such as the Facts or other contextual factors. Moreover, the parameters of the saliency function could also be learned from data.

4.4.8 Clarification Ellipsis

Clarification Ellipsis are a kind of *clarification requests*. To resolve clarification requests and their elliptical variants, Ginzburg (2012) includes a general theory of grounding and clarification requests. This theory would add a non-trivial amount of complexity to our formalization so we shall assume in our formalization that the latest utterance always grounded unless a clarification request comes afterwards. Therefore we resolve Clarification Ellipsis on the MaxQUD element without adding any other structure to the dialogue state. We also consider the Clarification Ellipsis to have only a clausal confirmation reading, leaving aside their other possible readings which would require a more elaborate approach (more details in Ginzburg (2012)). The clausal confirmation reading can be exemplified by (4.19).

- (4.19) A: Is Paul coming to the party?
 B: Paul? (= Are you asking if **Paul** is coming to the party?)

The clausal confirmation reading can be interpreted as asking a polar question about the constituent brought about by the Clarification Ellipsis.

CE_{conf} :

$$\forall p, x, \mathbf{y}$$

if ($nsu_a = CE \wedge u_a = "x?" \wedge$
 $(qud[\text{max-qud}] . q = p(x, \mathbf{y}) \vee p(x, \mathbf{y}) \in qud[\text{max-qud}] . fec)$) **then**
 $\left\{ P \left(a_a \leftarrow Ask(p(x, \mathbf{y})) \right) = 1 \right.$

As an example we can show the application of the rule on (4.19). The context is:

$$\begin{aligned} \text{max-qud} &= 1 \\ \text{qud}[\text{max-qud}] . q &= \text{comingToParty}(\text{IND}_1) \\ \text{qud}[\text{max-qud}] . fec &= \{\text{named}(\text{IND}_1, \text{Paul})\} \\ nsu_a &= CE \end{aligned}$$

The result of the application of the rule is:

$$\begin{aligned} a_a &= Ask(\text{named}(\text{IND}_1, \text{Paul})) \\ \text{new-fec} &= \{\} \end{aligned}$$

4.5 Implementation and use case example

In this section we exemplify some usages of the rules on a real-world conversation. It shows some example of behavior of the rules over a selected transcript from the COMMUNICATOR dataset. However, we want to point out that this section is not intended to give an empirical evaluation of the rules which is very far from being trivial since such an evaluation would require (at least) the availability of a fully annotated dataset of transcripts with the dialogue acts and the context updates at each step.

The COMMUNICATOR (Walker et al., 2001) dataset is a set of transcripts of interactions between a dialogue system and human testers. The COMMUNICATOR dataset contains transcripts of conversations for booking flight tickets. The interactions are mainly “machine-driven” meaning that the system drives the conversation, it asks questions and the user only answer to those questions. In this scenario it is possible to find many answers NSUs such as Short Answers, Affirmative Answers and Rejections. To test the resolution rules over this transcript we integrated the rules within a dialogue system developed with the OpenDial toolkit.

Next we will briefly talk about the architecture of our dialogue system and then we elaborate the step-by-step description of the example of interaction with the system using the chosen transcript from the COMMUNICATOR dataset.

4.5.1 Dialogue system architecture

As defined in Lison (2014), “a *spoken dialogue system* is a computational agent that can converse with humans through everyday spoken language”. These systems have a complex structure formed of many different parts, however they are usually formed by the following major components:

- *Natural language understanding* (NLU), maps the textual utterances into a high-level semantic representation;
- *Dialogue management*, updates the dialogue state and plans the actions to perform;
- *Natural language generation* (NLG), generates the linguistic realization of the planned actions or dialogue acts;

The resolution of NSUs is closely related to the NLU task and to the dialogue management in presence of such utterances. For its proper operation, our implementation indeed includes also shallow NLU and NLG modules as well as a very simple action selection procedure. Figure 4.2 shows the work-flow of the system. The system takes as input a user utterance u_u . The NLU module generates the semantic representation of the utterance a_u , the NSU resolution takes place at this stage and involves the recovery of the right semantic form for the incomplete utterance using the information available in the dialogue context. The classifier, right before the resolution, generates the content of the variable nsu_u . The action selection module decides what to do with the user utterance, producing a semantic representation of the action to perform a_m (or dialogue act of the sentence to be uttered). The NLG module transforms a_m into its linguistic form u_m . Throughout the process the context is updated by rules triggered by a_u , a_m and u_m .

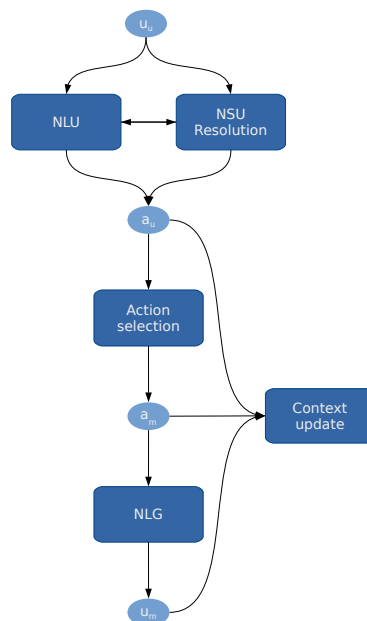


Figure 4.2: Schema of the basic work-flow of our dialogue system.

4.5.2 Example

In this section we present an extended example of the usage of some of the rules on the example transcript (4.20) from the COMMUNICATOR dataset. The transcript (4.20) encodes a conversation between a dialogue system and a user with the purpose of booking a ticket for a flight.

- (4.20) COMMUNICATOR 1693:04
1. M: [...] What are your travel plans?
 2. U: I'd like to travel from Columbus Ohio to Phoenix Arizona on Thursday October 5th.
 3. M: What time do you want to leave Columbus?
 4. U: Before 6 P.M.
 5. M: Do you have a preferred airline?
 6. U: No.
 7. M: I have you going from Columbus Ohio to Phoenix Arizona on October 5th. Is that correct?
 8. U: Yes.
 9. M: Will you return to Columbus from Phoenix?
 10. U: No.
 11. M: Is Phoenix your final destination?
 12. U: Yes.
 13. M: Hold on while I check availability.

As previously remarked, this type of conversations is driven by the system which asks the questions to fill the pieces of information it needs to properly retrieve the available flight options. The user answers to the questions posed by system, often using NSUs such as Short Answers, Affirmative Answers and Rejections. It is clear that this example covers only few of the available NSU classes but still it suffice at explaining the overall behavior of our implementation. More elaborate examples would have been complicated to understand and to read therefore we address the interested reader to the actual implementation⁹ to have a more detailed look at the behavior of the system.

We want to stress that the predicates we use are simple and handcrafted. A much more sophisticated NLU module would be required to automatize the process of extracting those predicates from the sentences which is a completely different problem from the one we addressed in this work.

The following is a step-by-step explanation of the systems behavior through the updates of the dialogue state. The process is intrinsically probabilistic, for every move there are many updates with a certain probability but, for readability, we avoid to show each possible probabilistic update and we show the most probable one instead.

⁹cf. <https://github.com/paolodragone>

1. The system greets and asks for the the travel plans of the user, updating the QUD with the new question:

$u_m = [\dots]$ What are your travel plans?
 $a_m = \text{Ask}(\text{travelPlans}(x_1, x_2, x_3))$
 $\text{qud}[1].q = \text{travelPlans}(x_1, x_2, x_3)$
 $\text{max-qud} = 1$

2. The user asserts its travel plans resolving the current issue under discussion:

$u_u =$ I'd like to travel from Columbus Ohio to Phoenix Arizona on
 Thursday October 5th.
 $a_u = \text{Assert}(\text{travelPlans}(C_1, C_2, D_1))$
 $\text{new-fec} = \{\text{city}(C_1, \text{Columbus}), \text{city}(C_2, \text{Phoenix}), \text{date}(D_1, \text{05-10})\}$

The system then accepts the assertion and resolves the issue:

$\text{facts} = \{\text{travelPlans}(C_1, C_2, D_1),$
 $\text{city}(C_1, \text{Columbus}), \text{city}(C_2, \text{Phoenix}), \text{date}(D_1, \text{05-10})\}$
 $\text{max-qud} = 0$

3. The system asks for time of departure, a new issue arises:

$u_u =$ What time do you want to leave Columbus?
 $a_u = \text{Ask}(\text{departTime}(x_1))$
 $\text{qud}[1].q = \text{departTime}(x_1)$
 $\text{max-qud} = 1$

4. The user resolves the issue with a Short Answer. The meaning of this NSU is inferred from the MaxQUD: the departure time must be before the given hour. The Short Answer resolution rule gives the following result:

$u_u =$ Before 6 P.M.
 $\text{nsu}_u = \text{ShortAns}$
 $a_u = \text{Assert}(\text{departTime}(T_1))$
 $\text{new-fec} = \{\text{before}(T_1, T_2), \text{time}(T_2, \text{18:00})\}$

The system again acknowledges the answer of the user inserting his assertion in the set of Facts and downdates the QUD with the resolved issue:

$\mathbf{facts} = \{\text{travelPlans}(C_1, C_2, D_1),$
 $\quad \text{city}(C_1, \text{Columbus}), \text{city}(C_2, \text{Phoenix}), \text{date}(D_1, \text{05-10}),$
 $\quad \text{departTime}(T_1), \text{before}(T_1, T_2), \text{time}(T_2, \text{18:00})\}$
 $\text{max-qud} = 0$

5. The system asks again a question about the preferred airline of the user. This time it is a polar question (notice the absence of a variable in the question predicate).

$\mathbf{u}_m = \text{Do you have a preferred airline?}$
 $\mathbf{a}_m = \text{Ask}(\text{havePreferredAirline}(\text{user}))$
 $\text{qud}[1].\mathbf{q} = \text{havePreferredAirline}(\text{user})$
 $\text{max-qud} = 1$

6. The user gives a negative answer to the previous question using an NSU. Again the meaning is inferred from the MaxQUD: the user does not have a preferred airline. From Section 4.4.3 we recall that the resolution rule in this case applies the *Neg* function to the predicate to indicate its negative form.

$\mathbf{u}_u = \text{No.}$
 $\mathbf{nsu}_u = \text{Reject}$
 $\mathbf{a}_u = \text{Assert}(\text{Neg}(\text{havePreferredAirline})(\text{user}))$

Again the answer of the user is inserted in the Facts:

$\mathbf{facts} = \{\text{travelPlans}(C_1, C_2, D_1),$
 $\quad \text{city}(C_1, \text{Columbus}), \text{city}(C_2, \text{Phoenix}), \text{date}(D_1, \text{05-10}),$
 $\quad \text{departTime}(T_1), \text{before}(T_1, T_2), \text{time}(T_2, \text{18:00}),$
 $\quad \text{Neg}(\text{havePreferredAirline})(\text{user})\}$
 $\text{max-qud} = 0$

7. The system summarizes the pieces of information gained so far and queries the user for their correctness. A check question like this can be represented as asking the first proposition as a polar question in the following way:

$\mathbf{u}_m = \text{I have you going ... Is that correct?}$
 $\mathbf{a}_m = \text{Ask}(\text{travelPlans}(C_1, C_2, D_1))$
 $\text{qud}[1].\mathbf{q} = \text{travelPlans}(C_1, C_2, D_1)$
 $\text{qud}[1].\mathbf{fec} = \{\text{city}(C_1, \text{Columbus}), \text{city}(C_2, \text{Phoenix}), \text{date}(D_1, \text{05-10})\}$
 $\text{max-qud} = 1$

8. The user confirms the information stored by the system with an Affirmative Answer. We recall that an Affirmative Answer is equivalent to stating the polar question in MaxQUD as it is:

$u_u = \text{Yes.}$
 $nsu_u = \text{AffAns}$
 $a_u = \text{Assert}(\text{travelPlans}(C_1, C_2, D_1))$
 $\text{new-fec} = \{\text{city}(C_1, \text{Columbus}), \text{city}(C_2, \text{Phoenix}), \text{date}(D_1, \text{05-10})\}$

The system acknowledges the answer by downdating the QUD but leaving Facts as it is since no additional information was included.

$\text{max-qud} = 0$

9. The system asks a new question about a possible return flight.

$u_m = \text{Will you return to Columbus from Phoenix?}$
 $a_m = \text{Ask}(\text{return}(C_2, C_1))$
 $\text{qud}[1].q = \text{return}(C_2, C_1)$
 $\text{qud}[1].fec = \{\text{city}(C_1, \text{Columbus}), \text{city}(C_2, \text{Phoenix})\}$
 $\text{max-qud} = 1$

10. The user says that he will not return from its destination.

$u_u = \text{No.}$
 $a_u = \text{Assert}(\text{Neg}(\text{return})(C_2, C_1))$
 $\text{facts} = \{\text{travelPlans}(C_1, C_2, D_1),$
 $\quad \text{city}(C_1, \text{Columbus}), \text{city}(C_2, \text{Phoenix}), \text{date}(D_1, \text{05-10}),$
 $\quad \text{departTime}(T_1), \text{before}(T_1, T_2), \text{time}(T_2, \text{18:00}),$
 $\quad \text{Neg}(\text{havePreferredAirline})(\text{user}),$
 $\quad \text{Neg}(\text{return})(C_2, C_1)\}$
 $\text{max-qud} = 0$

11. The system asks one last question about a its final destination.

$u_m = \text{Is Phoenix your final destination?}$
 $a_m = \text{Ask}(\text{finalDest}(C_2))$
 $\text{qud}[1].q = \text{finalDest}(C_2)$
 $\text{qud}[1].fec = \{\text{city}(C_2, \text{Phoenix})\}$
 $\text{max-qud} = 1$

12. The user confirms that Phoenix is his last destination.

$u_u = \text{Yes}$.

$a_u = \text{Assert}(\text{lastDest}(C_2))$

$\text{facts} = \{\text{travelPlans}(C_1, C_2, D_1),$
 $\text{city}(C_1, \text{Columbus}), \text{city}(C_2, \text{Phoenix}), \text{date}(D_1, \text{05-10}),$
 $\text{departTime}(T_1), \text{before}(T_1, T_2), \text{time}(T_2, \text{18:00}),$
 $\text{Neg}(\text{havePreferredAirline})(\text{user}),$
 $\text{Neg}(\text{willReturn})(\text{user}),$
 $\text{lastDest}(C_2)\}$

$\text{max-qud} = 0$

13. After gathering up all the information needed, contained in the Facts, the system starts checking the availability of the flights for the user.

A possible continuation would be that the system finds some alternative flights to show to the user and ask him which one he or she prefers. The understanding in the dialogues is not always perfect though. It often happens that the system asks for repetition or even that the user resets or interrupts the conversation. This makes the dialogue transcripts from the COMMUNICATOR dataset very unpredictable and a good starting point to evaluate the benefits and limitations of probabilistic approaches to NSU resolution.

This example gives enough detail to understand the basic behavior of the system in this particular dialogue domain. It is of course limited in many ways and extensions could be sought in different directions. Testing the system on a different domain covering other types of NSUs should be perhaps the first way to check the validity of the rules. However, a complete evaluation would require the manual annotation of several dialogue transcripts, which is not an easy task.

Besides, the behavior of the rules rely on a fixed semantics that is highly domain-dependent. To make the system scalable to different domains one needs to integrate other grammatical and lexical resources to make the semantics as generic as possible. It is clear though that a complete domain independent framework would be difficult to achieve.

Our system is of course still a prototype of a complete dialogue system for NSU interpretation. There are many improvements needed to achieve a working system. For instance the inclusion of rules to handle the NSU classes not covered by our work. A more general theory of grounding is also needed to properly account for the clarification requests. Extending every other assumption we made to simplify the development is another important goal for possible future works. Furthermore, to enhance the capabilities of the system, the integration of other natural language understanding modules is needed as well. Anaphora Resolution (Mitkov, 2014) and Named Entity Recognition (Nadeau and Sekine, 2007) are just a few of the problems concerning the correct interpretation of NSUs.

4.6 Summary

In this chapter we presented how to resolve the semantic meaning of NSUs. The resolution is a task that aims at extracting the meaning of an NSU given the dialogue context. To do so, we relied on the previous work of Fernández (2006) which presented a series of rules to resolve the meaning of the NSUs from a TTR-encoded dialogue context. As previously argued throughout this thesis, the use of a purely logic-based formalism, such as TTR, has some disadvantages in dealing with partially observable inputs and stochastic events when compared to a probabilistic approach. We showed how to reformulate the rules from Fernández (2006) using the probabilistic rules formalism (Lison, 2014) in order to include a probabilistic account of the dialogue state. We made use of a portion of the dialogue context theory from Ginzburg (2012) to encode the basic elements of the dialogue state needed for the resolution of the NSUs (see Section 2.2). We presented in Section 4.4 the probabilistic rules for the resolution of the NSUs. In Section 4.5 we also described a step-by-step example of the usage of the rules. The framework presented in this chapter has also been implemented and tested with OpenDial (Lison and Kennington, 2015).

Chapter 5

Conclusion

This chapter concludes the present thesis by summarizing the contributions of our work. The chapter also points out some of the possible developments that can be pursued in future works.

5.1 Contributions

In this thesis we described our research work regarding non-sentential utterances. NSUs are utterances that do not have a complete sentential form but convey a full meaning. However, they require to be “interpreted” i.e. their meaning must be extracted from the context of the dialogue. Our experiments concerned two separate aspects of the interpretation of the NSUs, namely:

- The classification of NSUs given their context;
- The resolution of the semantic content of the NSUs from the dialogue context;

In Chapter 2 we presented the background knowledge needed to the development of our work. We discussed in Section 2.1 the concept of non-sentential utterance, referring to Fernández (2006) as our theoretical basis. From the aforementioned work we employ the same taxonomy and corpus of NSUs. We then explain our methodology for the interpretation of NSUs, also based on the theory from Fernández (2006). To interpret an NSU, we first classify it according to the aforementioned taxonomy using machine learning then we “resolve” its meaning from the dialogue context through a resolution procedure dependent on its type.

Fernández (2006) develops a set of resolution procedures based on Type Theory with Records (Cooper, 2004; Ginzburg, 2012). In Section 2.2 we briefly describe the aspects of TTR and the theory of dialogue context from Ginzburg (2012) that we needed in our work.

In this thesis we argue that a purely logical framework such as TTR may have disadvantages in dealing with the uncertain nature of the NSUs. In our view a proper alternative is a probabilistic approach to the resolution of NSUs. To this end we employ the probabilistic rules formalism and the theory from Lison (2014) as probabilistic representation of the dialogue state and the NSU resolution procedures. We detail the basic aspects of the theory from Lison (2014) in Section 2.3. In Lison (2014) the dialogue state is represented as a Bayesian network and its dynamics are described by probabilistic rules. Probabilistic rules are *if ... then ... else ...* constructs that map logical conditions to probabilistic effects.

The focus of Chapter 3 is on the classification of NSUs, which is the task of inferring the type of a given NSU from its context. The context of a NSU is formed by its “antecedent”, the preceding utterance that holds its hidden meaning. Our work on the classification of NSUs is based on Fernández et al. (2007). We replicated their approach and set it as our baseline, as explained in Section 3.3. In Section 3.4 we describe the new features we use to extend the baseline feature set. The extended feature set alone was not enough to achieve an improvement of the classification performances. The major problems in this respect were the scarcity of labeled data and the class imbalance. To address those problems we employed semi-supervised learning techniques that we detail in Section 3.5. Our experiments show that the combination of the extended feature set and new training instances labeled with Active Learning led to a significant improvement of the classification accuracy. Nevertheless we argue that further analysis and testing need a larger amount of labeled data to be carried out properly. In Dragone and Lison (2015a) we present our findings in the classification of NSUs using Active Learning.

In Chapter 4 we detail the resolution of NSUs. The NSU resolution is the task of extracting the meaning of a given NSU from the dialogue context. We explain the process that we employ for the NSU resolution in Section 4.1. In Section 4.2 we describe the theoretical concepts we need from Ginzburg (2012). The description of our theory starts in Section 4.3 where we explained the design of the dialogue context. To model our dialogue context we take inspiration from Ginzburg (2012), however we reinterpret its constructs as random variables. The random variables in the dialogue state interact with each other through the probabilistic rules. The resolution rules that we developed are explained in Section 4.4. Finally, in Section 4.5, we show a detailed example of the use of the resolution rules over a transcript from the COMMUNICATOR dataset (Walker et al., 2001).

Our approach to the resolution of NSUs is intended to be a proof-of-concept for this task. We showed how we could reuse many concepts from the theory of Fernández (2006) and Ginzburg (2012) and “translate” the resolution rules based on TTR into probabilistic rules. The works on classification and interpretation carried out in this thesis were also presented in Dragone and Lison (2015b).

5.2 Future developments

In this section we list a series of ideas for possible future works that come out directly from our findings and from the assumptions we made.

Improve the NSU classification performances

In our work on the classification of NSUs we experimented many different approaches in seeking an improvement of the classification accuracy. However, there are many other paths that we did not explore or aspects in our approaches that can be improved. Here we discuss a few of the possible extension of our work.

There are classes of NSU that are intrinsically difficult to predict such as Helpful Rejections and other pairs of classes that are difficult to discriminate such as Repeated Acknowledgments and Repeated Affirmative Answers. A common issue in trying to predict those classes is that the parallelism with their antecedent is almost entirely at the semantic level. This requires deeper understanding of the phenomena and the use of features that exploit semantic relations in the NSU instances. We did not use any semantic feature since it would have added a non-trivial amount of complexity to our feature extraction algorithms. The deeper understanding of “difficult” classes and the use of such features may be a good starting point to any feature work on this topic.

Using additional features does not avoid the problem of class imbalance in the dataset. Many techniques could be experimented to try to mitigate this issue. An example may be an over-sampling technique such as *SMOTE*¹ (Chawla et al., 2002). The aforementioned work shows that the combination of *SMOTE* and majority class under-sampling leads to better classification performances on certain domains.

Perhaps the most difficult issue to overcome is the scarcity of labeled data. Our work shows that additional training data is indeed useful to improve the classification performances but we still lack enough data to run proper evaluations. We did not use the instances labeled with Active Learning as test data. Additional data for the gold standard should be composed of high-quality, manually annotated instances extracted within a corpus study that closely follows the original one from Fernández and Ginzburg (2002).

Incorporate additional elliptical phenomena

The corpus study from Fernández and Ginzburg (2002) is focused on data extracted from the British National Corpus therefore confined to only certain kind of dialogue domains. As argued by Raghu et al. (2015), there are many elliptical phenomena that do not fit well in the taxonomy from Fernández and Ginzburg (2002). An interesting follow up work on non-sentential utterances might try to find new elliptical phenomena in different

¹Synthetic Minority Over-sampling Technique.

dialogue domains and try to extend the taxonomy. Another point that may be considered is that some classes include large variety of forms and functions such as Short Answers and Helpful Rejections. A possible development could be to increase the granularity of such classes to try to capture more subtle differences.

Extend our NSU resolution approach

Our study on the resolution of NSUs pioneers a rule-based approach that involves a probabilistic dynamics of the dialogue state. There are still many issues to address:

- Increase the coverage of the rules to all the classes that were not covered by our work: Factual Modifiers, Helpful Rejections and so on.
- Develop a proper mechanism of rule adaptation in the presence of lexical modifiers e.g. for Sluices such as “For how long?”.
- Include grammatical and lexical resources to extract more complex meanings. A simple Short Answer that would not be covered by our rules is:

(5.1) A: Who is coming tomorrow?
 B: Nobody.

- Properly evaluate the rules on testing data from different dialogue domains.

The last point, perhaps the most important one, would require the development of a corpus of dialogue transcripts annotated with each semantic move and state update. In turn this would require to develop a generic representation of the semantic content of the utterances which is a non-trivial task by itself. A possibility could be to use TTR as semantic representation and reformulate the rules accordingly. Using TTR as basic semantic formalism it could be an interesting challenge to develop probabilistic rules to address a larger set of linguistic phenomena besides NSUs.

Combine different NSU resolution approaches

For our work on the NSU resolution we develop a rule-based approach based on a probabilistic representation of the dialogue state. It bares similarities with statistical approaches for the resolution such as Raghu et al. (2015). Their work is concentrated on follow-up NSU questions such as:

(5.2) A: How much for this model?
 B: ...
 A: For this other one?

Their approach is based on the combination of keywords from the follow-up question and the original one. From the combination of keywords they build possible meaningful “completions” of the NSU e.g. a completion for the NSU at the third line in (5.2) would

be “How much for this other model?”. After generating each possible completions they rank them according to some score and pick the best one.

Differently from ours, their approach does not use a high-level semantic representation of the utterances. It would be interesting to try to combine their statistical approach to our probabilistic rule-based one.

Compare our approach with other existing systems

Another useful comparison, and perhaps integration, should be made with the systems originally developed on the theory of Fernández (2006), namely SHARDS (Ginzburg et al., 2007) and one of its extensions CLARIE (Purver, 2006). The former is a system for ellipsis resolution that can handle Short Answer, Sluices and Affirmative Answers. The latter is a dialogue system developed to deal with *clarification requests* and, among them, Clarification Ellipsis, implementing the theory of Ginzburg and Cooper (2004) on top of the GoDiS dialogue system (Larsson et al., 2000). Both are based on the HPSG² framework from Ginzburg and Sag (2000), which is substantially different from our current design. Our framework lacks a grammar and other lexical resources that are indeed needed to build a functional system. It would be interesting to further develop our approach taking advantage from aspects of those systems and perhaps even integrate them into our architecture based on probabilistic rules.

²Head-driven Phrase Structure Grammar.

Appendix A

Context update rules

In this appendix we provide an overview of the probabilistic rules used for updating the context that have been implemented in the dialogue system for the testing of the resolution rules. As described in Section 4.5, the dialogue system that we implemented is focused on conversation of the human-machine kind therefore we will be using the notations a_m and u_m to refer respectively to the dialogue act performed by the system and the corresponding raw utterance. Following the aforementioned interaction model, the dialogue context is only representing the pieces of information known by the system. The context update rules are needed in order to make the dialogue context evolve along with the user acts and relative system reactions. In particular we need the rules for updating the QUD and the Facts variables. These rules are inspired by, but not limited to, Ginzburg (2012). Section 2.2 gives the background knowledge for the rules from Ginzburg (2012).

The rules shown in this appendix are not meant to give an extensive look on the system architecture but rather an high-level insight on the behavior of the system. These rules may differ from the actual implementation due to technicalities but still they fit for the purpose of the explanation. We will not talk about other modules of the system that we used in the implementation (i.e. NLU, NLG and action selection) because they have been implemented merely for toy examples of interaction on simple domains and they are not directly concerning the interpretation of NSUs. Follows a description of each context update rule.

QUD increment

NSUs are mostly reactionary utterances to previously raised issues. We aim at using the NSU resolution rules to interpret the content of the user NSUs. For this reason we concentrate on “machine-driven” conversations such as the one used in Section 4.5. We assume that in this type of dialogues issues are raised only by the system while the user limits to answer. In this scenario is common to find NSUs uttered by the users, as it often happens in the COMMUNICATOR dataset.

Given this setting, we update the QUD only when the system raises a new issue and we downdate it only when the system accepts a user assertion which resolves the maximal element. We also take into account that asking a question and asserting a proposition may have different probabilities to update the QUD. In the rules below we encode the asking with full probability and the asserting with a probability of 0.75, although, as remarked many times throughout this thesis, those probabilities may actually be estimated on real data.

qud-increment :

$\forall q, \mathbf{x}$

if ($\mathbf{a}_m = \text{Ask}(q(\mathbf{x}))$) **then**

$$\left\{ P \left(\begin{array}{l} \text{qud}[\text{qud.size} + 1].q \leftarrow q(\mathbf{x}), \\ \text{qud}[\text{qud.size} + 1].utt \leftarrow \mathbf{u}_m, \\ \text{qud.size} \leftarrow \text{qud.size} + 1 \end{array} \right) = 1 \right.$$

else if ($\mathbf{a}_m = \text{Assert}(q(\mathbf{x}))$) **then**

$$\left\{ P \left(\begin{array}{l} \text{qud}[\text{qud.size} + 1].q \leftarrow q(\mathbf{x}), \\ \text{qud}[\text{qud.size} + 1].utt \leftarrow \mathbf{u}_m, \\ \text{qud.size} \leftarrow \text{qud.size} + 1 \end{array} \right) = 0.75 \right.$$

The rule below handles the update of the FEC of the newly added QUD element. The rule adds to the FEC of the new element of QUD only the **new-fec** predicates sharing at least a variable with the proposition predicate.

fec-update :

$\forall p, \mathbf{x}, p', \mathbf{x}', x$

if ($(\mathbf{a}_m = \text{Ask}(p(\mathbf{x})) \vee \mathbf{a}_m = \text{Assert}(p(\mathbf{x}))) \wedge p'(\mathbf{x}') \in \text{new-fec} \wedge x \in \mathbf{x} \wedge x \in \mathbf{x}'$) **then**

$$\left\{ P \left(\text{qud}[\text{qud.size} + 1].\text{fec} \leftarrow \text{qud}[\text{qud.size} + 1].\text{fec} \cup \{p'(\mathbf{x}')\} \right) = 1 \right.$$

QUD downdate

The QUD is downdated when the system accepts a proposition asserted by the user. The system responds with an Accept act which will remove the MaxQUD from the QUD.

qud-downdate :

if ($\mathbf{a}_m = \text{Accept}(p)$) **then**

$$\left\{ P \left(\begin{array}{l} \text{qud}[\text{max-qud}].q \leftarrow \text{None}, \\ \text{qud}[\text{max-qud}].utt \leftarrow \text{None}, \\ \text{qud}[\text{max-qud}].fec \leftarrow \text{None}, \\ \text{qud.size} \leftarrow \text{qud.size} - 1 \end{array} \right) = 1 \right.$$

Max-qud update

As soon as the QUD array is updated, the MaxQUD is updated too. As explained in Section 4.3.3, the `max-qud` variable is defined as the index of the MaxQUD inside the QUD array and its stack-like behavior is determined by an exponentially decreasing probability with maximum on the last inserted element.

max-qud-update :

$\forall i$
if ($i > 0 \wedge i \leq \text{qud.size} \wedge \text{qud}[i].q \neq \text{None}$) **then**
 $\left\{ P(\text{max-qud} \leftarrow i) = e^{i-\text{qud.size}}$

Facts increment

As mentioned above, the dialogue context encodes the knowledge of the system and so are the Facts. The Facts variable contains only predicates accepted by the system. In the rule below we incorporate the ones presented in Section 4.4.4 for Propositional Modifiers and the ones for handling Rejections and Affirmative Answers. Again we point out that while the probabilities are handcrafted for simplicity here, they can be actually learned from data.

facts-increment :

$\forall p, \mathbf{y}$
if ($\mathbf{a}_m = \text{Accept}(\text{PropRel}_{\text{probably}}(p)(\mathbf{y}))$) **then**
 $\left\{ P(\text{facts} \leftarrow \text{facts} \cup \{p(\mathbf{y})\} \cup \text{new-fec}) = 0.75$
else if ($\mathbf{a}_m = \text{Accept}(\text{PropRel}_{\text{unlikely}}(p)(\mathbf{y}))$) **then**
 $\left\{ P(\text{facts} \leftarrow \text{facts} \cup \{p(\mathbf{y})\} \cup \text{new-fec}) = 0.25$
...
else if ($\mathbf{a}_m = \text{Accept}(\text{Neg}(p)(\mathbf{y}))$) **then**
 $\left\{ P(\text{facts} \leftarrow \text{facts} \cup \{\text{Neg}(p)(\mathbf{y})\} \cup \text{new-fec}) = 1$
else if ($\mathbf{a}_m = \text{Accept}(p(\mathbf{y}))$) **then**
 $\left\{ P(\text{facts} \leftarrow \text{facts} \cup \{p(\mathbf{y})\} \cup \text{new-fec}) = 1$

Bibliography

- Bergsma, S. (2010). *Large-scale semi-supervised learning for natural language processing*. University of Alberta.
- Boser, B. E., I. M. Guyon, and V. N. Vapnik (1992). “A training algorithm for optimal margin classifiers”. In: *Proceedings of the fifth annual workshop on Computational learning theory*. ACM, pp. 144–152.
- Burnard, L. (2000). *Reference guide for the British National Corpus (world edition)*.
- Chawla, N. V. et al. (2002). “SMOTE: synthetic minority over-sampling technique”. In: *Journal of artificial intelligence research*, pp. 321–357.
- Chen, D. and C. D. Manning (2014). “A fast and accurate dependency parser using neural networks”. In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Vol. 1, pp. 740–750.
- Collobert, R. et al. (2006). “Large scale transductive SVMs”. In: *The Journal of Machine Learning Research* 7, pp. 1687–1712.
- Cooper, R. (2004). “A type theoretic approach to information state update in issue based dialogue management”. In: *Catalog '04*, invited paper.
- Cooper, R. (2005). “Records and record types in semantic theory”. In: *Journal of Logic and Computation* 15.2, pp. 99–112.
- De Marneffe, M.-C. et al. (2014). “Universal Stanford Dependencies: A cross-linguistic typology”. In: *Proceedings of LREC*, pp. 4585–4592.
- Dragone, P. and P. Lison (2015a). “An Active Learning Approach to the Classification of Non-Sentential Utterances”. In: *Proceedings of the second Italian Conference on Computational Linguistics*. Trento, Italy, in press.
- Dragone, P. and P. Lison (2015b). “Non-Sentential Utterances in Dialogue: Experiments in classification and interpretation”. In: *Proceedings of the 19th Workshop on the Semantics and Pragmatics of Dialogue*. Göteborg, Sweden, p. 170.
- Duan, K.-B. and S. S. Keerthi (2005). “Which is the best multiclass SVM method? An empirical study”. In: *Multiple Classifier Systems*, pp. 278–285.

- Fernández, R. and J. Ginzburg (2002). “Non-sentential utterances in dialogue: A corpus-based study”. In: *Proceedings of the 3rd SIGdial workshop on Discourse and dialogue-Volume 2*. Association for Computational Linguistics, pp. 15–26.
- Fernández, R., J. Ginzburg, and S. Lappin (2007). “Classifying non-sentential utterances in dialogue: A machine learning approach”. In: *Computational Linguistics* 33.3, pp. 397–427.
- Fernández, R. R. (2006). “Non-sentential utterances in dialogue: Classification, resolution and use”. PhD thesis. King’s College London.
- Garside, R. (1993). “The large-scale production of syntactically analysed corpora”. In: *Literary and Linguistic Computing* 8.1, pp. 39–46.
- Ginzburg, J. (2012). *The interactive stance*. Oxford University Press.
- Ginzburg, J. and R. Cooper (2004). “Clarification, ellipsis, and the nature of contextual updates in dialogue”. In: *Linguistics and philosophy* 27.3, pp. 297–365.
- Ginzburg, J. and I. Sag (2000). *Interrogative investigations*. Stanford: CSLI publications.
- Ginzburg, J. et al. (2007). “SHARDS: Fragment resolution in dialogue”. In: *Computing Meaning*, pp. 125–144.
- Hall, M. et al. (2009). “The WEKA data mining software: an update”. In: *ACM SIGKDD explorations newsletter* 11.1, pp. 10–18.
- Klein, D. and C. D. Manning (2003). “Accurate unlexicalized parsing”. In: *Proceedings of the 41st Annual Meeting on Association for Computational Linguistics-Volume 1*. Association for Computational Linguistics, pp. 423–430.
- Larsson, S. et al. (2000). “GoDiS: an accommodating dialogue system”. In: *Proceedings of the 2000 ANLP/NAACL Workshop on Conversational systems-Volume 3*. Association for Computational Linguistics, pp. 7–10.
- Leech, G., R. Garside, and M. Bryant (1994). “CLAWS4: the tagging of the British National Corpus”. In: *Proceedings of the 15th conference on Computational linguistics-Volume 1*. Association for Computational Linguistics, pp. 622–628.
- Lewis, D. D. and J. Catlett (1994). “Heterogeneous uncertainty sampling for supervised learning”. In: *Proceedings of the eleventh international conference on machine learning*, pp. 148–156.
- Liang, P. (2005). “Semi-supervised learning for natural language”. PhD thesis. Massachusetts Institute of Technology.
- Lison, P. (2015). “A hybrid approach to dialogue management based on probabilistic rules”. In: *Computer Speech & Language* 34.1, pp. 232–255.

- Lison, P. (2012). “Declarative Design of Spoken Dialogue Systems with Probabilistic Rules”. In: *Proceedings of the 16th Workshop on the Semantics and Pragmatics of Dialogue (SemDial 2012)*.
- Lison, P. (2014). “Structured Probabilistic Modelling for Dialogue Management”. PhD thesis. University of Oslo.
- Lison, P. and C. Kennington (2015). “Developing Spoken Dialogue Systems with the OpenDial toolkit”. In: *Proceedings of the 19th Workshop on the Semantics and Pragmatics of Dialogue*. Göteborg, Sweden.
- Marcus, M. P., M. A. Marcinkiewicz, and B. Santorini (1993). “Building a large annotated corpus of English: The Penn Treebank”. In: *Computational linguistics* 19.2, pp. 313–330.
- Mitkov, R. (2014). *Anaphora resolution*. Routledge.
- Murphy, K. P., Y. Weiss, and M. I. Jordan (1999). “Loopy belief propagation for approximate inference: An empirical study”. In: *Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence*. Morgan Kaufmann Publishers Inc., pp. 467–475.
- Nadeau, D. and S. Sekine (2007). “A survey of named entity recognition and classification”. In: *Linguisticae Investigationes* 30.1, pp. 3–26.
- Needleman, S. B. and C. D. Wunsch (1970). “A general method applicable to the search for similarities in the amino acid sequence of two proteins”. In: *Journal of molecular biology* 48.3, pp. 443–453.
- Platt, J. et al. (1999). “Fast training of support vector machines using sequential minimal optimization”. In: *Advances in kernel methods—support vector learning* 3.
- Purver, M. (2006). “CLARIE: Handling Clarification Requests in Dialogue System”. In: *Research on Language and Computation* 4.2-3, pp. 259–288.
- Quinlan, J. R. (1993). *C4.5: Programs for Machine Learning*. Vol. 1. Morgan Kaufmann.
- R Core Team (2015). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria. URL: <http://www.R-project.org/>.
- Raghu, D. et al. (2015). “A Statistical Approach for Non-Sentential Utterance Resolution for Interactive QA System”. In: *Proceedings of the SIGDIAL 2015 Conference*, pp. 335–343.
- Rodríguez, K. J. and D. Schlangen (2004). “Form, intonation and function of clarification requests in German task-oriented spoken dialogues”. In: *Proceedings of Catalog (the 8th workshop on the semantics and pragmatics of dialogue; SemDial04)*.
- Schlangen, D. (2003). “A coherence-based approach to the interpretation of non-sentential utterances in dialogue”. PhD thesis. University of Edinburgh. College of Science and Engineering. School of Informatics.

- Schlangen, D. (2005). “Towards finding and fixing fragments: Using ML to identify non-sentential utterances and their antecedents in multi-party dialogue”. In: *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*. Association for Computational Linguistics, pp. 247–254.
- Settles, B. (2010). “Active learning literature survey”. In: *University of Wisconsin, Madison* 52.55-66, p. 11.
- Shannon, C. (1948). “A mathematical theory of communication”. In: *Bell System Technical Journal, The* 27.3, pp. 379–423.
- Smith, T. F. and M. S. Waterman (1981). “Identification of common molecular subsequences”. In: *Journal of molecular biology* 147.1, pp. 195–197.
- Vapnik, V. N. (1998). *Statistical Learning Theory*. Wiley-Interscience.
- Walker, M. et al. (2001). “DARPA Communicator dialog travel planning systems: The June 2000 data collection”. In: *EUROSPEECH*, pp. 1371–1374.
- Zhang, N. L. and D. Poole (1996). “Exploiting causal independence in Bayesian network inference”. In: *Journal of Artificial Intelligence Research*, pp. 301–328.