

Abstract Argumentation

ROBERT A. KOWALSKI AND FRANCESCA TONI

Department of Computing, Imperial College of Science, Technology and Medicine 180, Queen's Gate, London SW7 2BZ, UK

Received and accepted 5 September 1996

Abstract. In this paper we explore the thesis that the role of argumentation in practical reasoning in general and legal reasoning in particular is to justify the use of defeasible rules to derive a conclusion in preference to the use of other defeasible rules to derive a conflicting conclusion. The defeasibility of rules is expressed by means of non-provability claims as additional conditions of the rules.

We outline an abstract approach to defeasible reasoning and argumentation which includes many existing formalisms, including default logic, extended logic programming, non-monotonic modal logic and auto-epistemic logic, as special cases. We show, in particular, that the 'admissibility' semantics for all these formalisms has a natural argumentation-theoretic interpretation and proof procedure, which seem to correspond well with informal argumentation.

In the admissibility semantics there is only one way for one argument to attack another, namely by undermining one of its non-provability claims. In this paper, we show how other kinds of attack between arguments, specifically how rebuttal and priority attacks, can be reduced to the undermining of non-provability claims.

Key Words: argumentation, default reasoning, priority

1. Introduction

The purpose of this paper is to outline a formal theory of argumentation, which promises to have relevance for practical reasoning in general and for legal reasoning in particular. We will argue that the role of argumentation is to justify the use of certain defeasible rules deriving a conclusion in preference to the use of other defeasible rules deriving conflicting conclusions.

We explore the thesis that defeasible reasoning with rules of the form

P if Q.

can be understood as 'exact' reasoning with rules of the form

P if Q and S can not be shown.

containing one or more defeasible 'non-provability claims' of the form

S can not be shown, for some sentence *S*.

In many cases, the sentence *S* is the 'contrary', *not P*, of the conclusion, *P*, of the rule. In other cases, *S* may state that the rule itself is 'defeated' by another rule.

With this understanding of defeasibility, argumentation is a dialectic process whereby a proponent presents an 'exact' argument for a conclusion, which is based, however, upon defeasible non-provability claims of the form *S can not be shown*. Such a claim, and the argument it helps to support, can be defeated if an opponent undermines the claim by presenting an argument for *S*. Like the proponent, the opponent can also base her/his argument upon non-provability claims.

The argumentation process can be viewed, therefore, as a game in which the proponent moves first. In the 'credulous' version of the game, which we study in this paper, by moving first, the proponent has the advantage of being able to use his/her previously used claims of non-provability to defeat the opponent's counter-claims. Such credulous argumentation can be used to justify more than one point of view. If the opponent moves first, the opponent can use her/his previously used claims to defeat the proponent. In the 'sceptical' version of the game, which we do not explore in this paper, the proponent can use only undisputed claims of non-provability to defeat the opponent's counter-claims.

A similar view of defeasible reasoning and argumentation has been put forward by several authors, e.g. [Pollock, 1987; Prakken & Sartor, 1996; Gordon, 1993]. Our approach differs, however, in several respects. The most important of these are:

1. Our approach to argumentation focuses on the acceptability of the non-provability claims of an argument, rather than on the acceptability either of the argument or of the conclusion of the argument. This concentrates attention, therefore, on the contentious parts of the argument and avoids being distracted by those parts of the argument that are beyond dispute.
2. We reduce all forms of defeasibility to that of non-provability claims. As a consequence, the only way to defeat an argument is by undermining one of its claims, namely by presenting an argument for a sentence *S* which is claimed to be unprovable. Indirect defeat (also called 'rebutting' [Pollock, 1987; Prakken & Sartor, 1995] or 'reductio-ad-absurdum' [Dung, 1993]), showing that a rule

P if Q.

leads to contradiction, is transformed into undermining defeat, by rewriting the rule in the 'exact' form

P if Q and not P can not be shown.

making explicit in the conditions of the rule that the contrary, *not P*, of its conclusion, *P*, can not be shown.

3. Our semantics is credulous rather than sceptical. A sceptical semantics can be obtained from credulous semantics by regarding a conclusion as justified if and only if it can be derived from any credulously 'acceptable' set of claims.
4. Our methodology deals with priorities between rules by introducing explicit rule names into the language and adding conditions to the rules expressing that the rules are not defeated by other higher priority rules. This approach avoids the need to deal with priorities in the semantics.

5. Our approach does not require the introduction of a new language or of a new semantics. On the contrary, we base our formalisation upon a variant of an abstract approach to defeasible reasoning [Bondarenko et al., 1996] which has been shown to include many existing formalisms, including (extended) logic programming, default logic [Reiter, 1980], non-monotonic modal logic [McDermott, 1982] and autoepistemic logic [Moore, 1985] as special cases. Thus the approach taken in this paper is abstract and can be formalised in any one of these and other formalisms. Although all these formalisms use the ‘stability’ semantics, we argue that the ‘admissibility’ semantics is more adequate.

The paper is organised as follows. In section 2, we illustrate our approach by means of examples, which are used elsewhere in the paper. In section 3, we describe the abstract framework for defeasible reasoning, and present the stability semantics. In section 4, we show how default logic, extended logic programming, non-monotonic modal logic, and auto-epistemic logic can all be considered special cases of the framework. Thus our approach to defeasible reasoning and argumentation can be formulated in any one of these and other formalisms for default reasoning. In section 5, we define the admissibility semantics and its argumentation proof procedure. In section 6, we show how our approach reduces rebutting attacks to undermining attacks. In section 7, we show how the approach reduces priorities to non-provability conditions. In section 8, we conclude.

We assume the reader is already familiar with the general subject of argumentation and its relevance to legal reasoning. In particular, we recommend the article by Prakken and Sartor in this issue [1996] for an overview of related work.

The framework we use in this paper is a variant of one we have developed elsewhere [Bondarenko et al., 1996]. The reader will find formal definitions and results in [Bondarenko et al., 1996]. The main technical contribution of this paper, therefore, is the methodology we present for eliminating rebuttal attacks and priorities from the semantics. The main general contribution is our use of the methodology for argumentation which is abstract and language-independent.

2. Examples^o

EXAMPLE 2.1. Consider the following statement of the principle that, by default, a person accused of a crime should be assumed innocent unless the person can be shown to be not innocent.

- (1) *A person is innocent of a crime*
 if the person is accused of committing the crime
 *and it can not be shown that the person is not innocent of the crime.*¹

The inclusion in the rule of the explicit non-provability condition transforms a defeasible rule

- (1') *A person is innocent of a crime*
 if the person is accused of committing the crime.

^o These and all other examples in this paper can easily be formalised using the techniques described in [Kowalski, 1995].

¹ We use the phrases ‘*S can not be shown*’ and ‘*it can not be shown that S*’ interchangeably.

which might be subject to argument, into an 'exact' rule, which is beyond dispute. Any dispute about the defeasible rule is transformed, instead, into a dispute about the non-provability condition.

In contrast, the 'exact' rule

(2) *A person is not innocent of a crime*

if the person was observed committing the crime.

does not contain any non-provability conditions and consequently does not represent a defeasible rule.

Suppose now that we are given the fact

(3) *John is accused of theft.*

We can justify the conclusion that John is innocent by the following argument.

Proponent: *John is innocent of theft,*

because (although John is accused of theft)

it can not be shown that John is not innocent of theft.

The proponent's claim, *it can not be shown that John is not innocent of theft*, is upheld, because the only way to defeat it, by using rule (2), fails to apply, given the lack of any 'fact' ('exact' rule) recording an observation of his having committed the crime. The claim can be defeated, however, if such information becomes available at a later time. It is because of this that such logic is often said to be 'non-monotonic', in contrast to the case in 'monotonic' logic, where once a conclusion is established it continues to hold no matter what further information is added.

Although the two rules, (1) and (2) above, are both 'exact', the inclusion of the extra condition in (1) can be understood as giving the defeasible rule, (1'), that (1) represents, a lower priority than (2). We shall argue more generally, in section 7, that priorities involving defeasible rules can be represented by explicit non-provability conditions of 'exact' rules.

EXAMPLE 2.2. Consider the two defeasible rules

A person inherits an estate

if he/she is the beneficiary of a valid will.

A person does not inherit an estate

if he/she murders the owner of the estate.

and the facts

John is the beneficiary of Henry's valid will.

John has murdered Henry.

Everything else being equal, in a credulous approach to defeasible reasoning, there are two equally 'acceptable' arguments: one concluding that John inherits Henry's estate because he is the beneficiary of Henry's valid will; the other that he does not inherit the estate because he has murdered Henry. In a sceptical approach (which we do not discuss in this paper, but which we do discuss in [Bondarenko et al., 1996]) neither conclusion holds.

In this paper, we will consider two kinds of ‘acceptable’ arguments, namely ‘stable’ (section 3) and ‘admissible’ (section 5). The two arguments above, one concluding that John inherits Henry’s estate and the other concluding that he does not inherit the estate, are ‘acceptable’ in both senses.

In our analysis of defeasibility and argumentation, we justify the credulous reasoning of this example in the following way:

To say that the first two rules are defeasible is to say that they have implicit conditions to the effect that the contrary of their conclusions can not be shown. We make these conditions explicit by rewriting the rules in the ‘exact’ form

A person inherits an estate
if he/she is the beneficiary of a valid will
and it can not be shown that the person does not inherit the estate.
A person does not inherit an estate
if he/she murders the owner of the estate
and it can not be shown that he/she inherits the estate.

We can now justify the conclusion that John inherits the estate by the following dialogue game.

Proponent: *John inherits the estate,*
because John is the beneficiary of a valid will
and because
it can not be shown that John does not inherit the estate.

Opponent: On the contrary,
it can be shown that John does not inherit the estate
because John murdered the owner of the estate
and because
it can not be shown that John inherits the estate.

Proponent: But, on the contrary,
it can be shown that John inherits the estate,
as I have already argued.

The rules of the game are designed to be liberal about the arguments that are ‘admissible’. An argument is ‘admissible’ if its non-provability claims can be defended against any counter-claims of an opponent. By being the first player, the proponent has the advantage of being able to use his/her earlier claims to defeat the opponent’s counter-claims. This enables the proponent, in this example, to have the last word and to win the argument. If the opponent is allowed to move first, then the opponent’s argument will be equally ‘admissible’.

It is also possible to design a conservative version of the game, which is sceptical about what conclusions can be established.

EXAMPLE 2.3. (Adapted from [Prakken & Sartor, 1995])

Consider the following general principle of European Community law:

*A product can be sold in a country
 if the country is part of the EC
 and the product can be sold in another country
 and the other country is part of the EC
 and it can not be shown that
 the selling of the product in the country endangers public health
 and it can not be shown that
 the selling of the product in the country prejudices the consumer.*

This principle potentially conflicts with the following rule of Italian law:

*A product can not be sold in Italy
 if the product is called 'pasta'
 and the product is not made of durum wheat.*

The potential conflict becomes real, when we are given only the following facts:

*Italy is part of the EC.
 BDR is part of the EC.
 alpha is sold in BDR.
 alpha is called 'pasta'.
 alpha is not made of durum wheat.*

If the two rules above are understood as 'exact' rules, then it is possible to derive the inconsistent conclusion that *alpha* can be sold in Italy and *alpha* can not be sold in Italy. This derivation is based upon the two non-provability claims of the first rule, neither of which can be undermined, using the given rules and facts.

However, if the two rules are understood as defeasible, then in our methodology for representing defeasibility we would transform them into 'exact' rules by adding extra non-provability conditions to both, rewriting them as:

*A product can be sold in a country
 if the country is part of the EC
 and the product can be sold in another country
 and the other country is part of the EC
 and it can not be shown that
 the selling of the product in the country endangers public health
 and it can not be shown that
 the selling of the product in the country prejudices the consumer
 and it can not be shown that
 the product can not be sold in the country.*

*A product can not be sold in Italy
 if the product is called 'pasta'
 and the product is not made of durum wheat
 and it can not be shown that the product can be sold in the country.*

As in example 2.2, it would now be possible to construct two separate ‘admissible’ arguments, each of which undermines the other, rather than one ‘admissible’ argument with an inconsistent conclusion.

However, if, as is the case, we want to give higher priority to European Community law than to Italian law, then we would simply retain the original formulation of the European Community law without an extra non-provability condition, following the lead of example 2.1.

Thus, by adding appropriate non-provability conditions to rules, expressing that the contrary of their conclusions can not be shown, we obtain the same result as Prakken and Sartor [1995], who, in addition to using undermining attacks, also employ indirect rebutting attacks and deal with priorities in the semantics. We will see later in section 7 that the treatment of priorities in this example can be seen as an optimised version of a more general approach that can handle more complex priorities, like those of Prakken and Sartor [1995], which can be defined by means of rules. This approach is illustrated by the following example.

EXAMPLE 2.4. (Adapted from [Hage, 1995])

Consider the principles:

r_1 : Except as provided for by r_2 , all thieves should be punished.

r_2 : Except as provided for by r_3 , thieves who are minors should be rehabilitated and not punished.

r_3 : Any thief who is violent should be punished.

These can be represented by the defeasible rules

r_1 : *A person should be punished
if the person is a thief.*

r_2 : *A person should not be punished
if the person is a thief
and the person is a minor.*

r_3 : *A person should be punished
if the person is a thief
and the person is violent.*

where rule r_3 has higher priority than rule r_2 , and r_2 higher priority than rule r_1 ($r_3 > r_2$ and $r_2 > r_1$).

Given the facts

f_1 : *John is a thief.*

f_2 : *John is a minor.*

rules r_1 and r_2 conflict. Intuitively, the priority $r_2 > r_1$ can be used to derive the conclusion *John should not be punished* of the second rule, in preference to the conclusion of the first. If we are given the additional fact

f_3 : *John is violent.*

then rules r_1 , r_2 and r_3 now conflict. Intuitively, the priority $r_3 > r_2$ can be used to derive the conclusion *John should be punished* of the third rule.

In the previous example 2.3 we used a simple method to represent defeasible rules as 'exact' rules while taking into account the priorities between them. The method consists of adding non-provability conditions, that the contrary of their conclusions can not be shown, to defeasible rules with lower priority, without adding similar conditions to the higher priority rules. This method can not be used in this example, since rule r_2 has higher priority than rule r_1 but lower priority than rule r_3 .

In such an example, we need to distinguish the different ways a conclusion can be established, and add non-provability conditions to defeasible rules with lower priority stating that no higher priority rule applies. In this way, using this more refined methodology, we obtain the following representation of the example.

A person should be punished
if the person is a thief
and it can not be shown that r_1 is defeated for the person.

r_1 is defeated for a person
if the person is a thief
and the person is a minor
and it can not be shown that r_2 is defeated for the person.

A person should not be punished
if the person is a thief
and the person is a minor
and it can not be shown that r_2 is defeated for the person.

r_2 is defeated for a person
if the person is a thief
and the person is violent
and it can not be shown that r_3 is defeated for the person.

A person should be punished
if the person is a thief
and the person is violent
and it can not be shown that r_3 is defeated for the person.

Now, given the facts f_1 and f_2 , there is only one 'admissible' argument for the conclusion *John should not be punished*, supported by the non-provability claim *it can not be shown that r_2 is defeated for John*. However, given the facts f_1 , f_2 and f_3 , there is only one 'admissible' argument for the conclusion *John should be punished*, supported by the non-provability claim *it can not be shown that r_3 is defeated for John*. This argument undermines the conflicting claim *it can not be shown that r_2 is defeated for John*, supporting the conclusion that *John should not be punished*.

3. An Abstract Framework for Representing Defeasibility

All the examples above have been formulated in a form of English which can be represented directly in many existing (and also in many presently uninvestigated) formalisms for default reasoning. Indeed, to formalise such examples, we need only to choose a logic of ‘exact’ (i.e. monotonic, non-defeasible) reasoning (such as classical first-order logic) and extend it, if necessary, by adding a non-provability operator. We then need to define a semantics and proof theory for the extended language, which appropriately interprets the intended meaning of the non-provability operator. More precisely

1. We assume an underlying **monotonic** (or **exact**) **logic** with a logical connective representing implication (to which modus ponens applies) and a universal quantifier (to which universal quantifier elimination applies).
2. We extend the language of the underlying logic, if necessary, by means of a **non-provability operator**, to express statements of the form.

S can not be shown

where *S* names a sentence of the language of the underlying logic. We restrict the occurrence of these non-provability statements to conditions of implications.

3. Defeasible rules can then be represented in the extended language by transforming them into (exact) universally quantified implications with non-provability conditions.
4. Defeasible reasoning from a set *T* (called a **theory**) of sentences in the extended language, is reduced to exact reasoning in the underlying logic, from an extended set of sentences $T \cup N$, where *N* is an ‘acceptable’ set of non-provability claims of the form

S cannot be shown.

Since, in most cases², the underlying language does not contain non-provability formulae, we have to explain how, in such cases, the underlying logic can be used to derive conclusions from the extended theory $T \cup N$. For this purpose, we assume that sentences in *N* are treated as atomic propositions. We also assume that instantiation can be applied to universally quantified implications representing defeasible rules, obtaining instantiated implications with atomic non-provability conditions. Modus ponens can then be applied to such instantiated implications and sentences in *N*, obtaining sentences of the original underlying language.

It remains to specify, given a theory *T*, when a set of non-provability claims *N* is ‘acceptable’. Almost all existing logics for default reasoning solve this problem in the same way, by means of a requirement that *N* be ‘stable’, i.e.

A candidate set *N* of sentences of the form *S can not be shown* is said to be **stable** if and only if for every *S* in the underlying language,

S can not be shown is in *N* if and only if

S is not derivable from $T \cup N$ in the underlying logic.

² As we will see later, non-monotonic modal logic and autoepistemic logic are exceptions.

Different logics for default reasoning formulate this stability requirement in different ways, which are all, nonetheless, equivalent [Bondarenko et al., 1996].³

Despite its being the dominant semantics for default reasoning, stability semantics has a number of drawbacks, the most obvious being that to verify the acceptability of a claim *S* *can not be shown* necessitates determining an entire stable set *N* to show that *S* is not derivable from $T \cup N$. This is computationally infeasible in most interesting cases. There are also cases where a stable set does not exist because of local paradoxical sentences in *T*, but where an intuitively acceptable argument can be constructed from other parts of *T* which are not affected by the paradox.

Consider again the theory *T* consisting of the set of sentences (1), (2) and (3) of example 2.1. Intuitively, the set *N* consisting of the single sentence

(4) *John is not innocent of theft can not be shown*

is acceptable, because there is no way to show, given only *T* and *N* that *John was observed committing theft*.

The conclusion that *John is innocent of theft* then follows from $T \cup N$ in the underlying logic, by instantiating (1), obtaining the variable-free implication

John is innocent of theft
if John is accused of theft
and John is not innocent of theft can not be shown.

and then applying modus ponens, using (3) and (4).

Intuitively, in this example it is unnecessary to consider the set of all sentences unprovable from $T \cup N$ to be convinced that the single sentence in *N* is unprovable. Similarly, even if *T* is extended to a much larger set of sentences *T'*, the same argument presented above continues to hold, provided *T'* does not contain any new sentences which can be used to derive either that *John was observed committing theft* or to derive more directly that *John is not innocent of theft*. Even a paradoxical sentence of the form

P if P can not be shown

in *T'* need not affect the correctness of the argument.

In section 5, we will show that these drawbacks of the stability semantics can be overcome by means of the more liberal 'admissibility' semantics [Dung, 1993], which has an argumentation-theoretic interpretation, as illustrated informally already in examples 2.1 and 2.2. The 'admissibility' semantics is compatible with stability semantics, in the sense that whenever a stable extension exists it is 'admissible'. Moreover, in many cases an 'admissible' set can be extended to a stable set [Bondarenko et al., 1996].

First, we show how a number of existing formalisms for default reasoning can be understood as instances of the abstract framework.

³ In non-monotonic modal logic and autoepistemic logic, non-provability sentences not in *N* may be derivable from $T \cup N$. In such cases, for *N* to be stable, *N* must also be **closed** in the sense that it contains all non-provability sentences derivable from $T \cup N$.

4. Special Instances

4.1. DEFAULT LOGIC

1. The underlying logic is first-order classical logic, augmented with domain-specific inference rules. These inference rules can equivalently be formulated as sentences of the form

$$P \text{ if } Q_1 \text{ and } \dots \text{ and } Q_n.$$

where P, Q_1, \dots, Q_n are all first-order formulae of the underlying language, $n \geq 0$, and *if* is a new implication sign, for which the inference rule of modus ponens applies. Notice that in this formulation there are two logical connectives for implication, the new implication symbol and the material implication of classical logic.

Any variables in these new implications which are not explicitly quantified are implicitly universally quantified, with scope the entire implication. The rule of instantiation applies to such variables.

2. Non-provability is expressed by means of a logical operator, M , where the intended meaning of $M P$, where P is a sentence of first-order logic, is that the contrary of P , namely *not* P , can not be shown. $M P$ can be interpreted equivalently as expressing that P is consistent.
3. Thus defeasible rules are expressed in default logic by translating them into exact sentences of the form

$$P \text{ if } Q_1 \text{ and } \dots \text{ and } Q_n \text{ and } M P_1 \text{ and } \dots \text{ and } M P_m.$$

where $n \geq 0, m > 0$ and apparently free variables are implicitly universally quantified. Statements of the form $M P$ are interpreted as atoms in the underlying first-order language.

4. The standard semantics of default logic [Reiter, 1980] has been shown [Bondarenko et al., 1996] to be a special case of the stable semantics defined in the previous section for abstract frameworks in general.

4.2. EXTENDED LOGIC PROGRAMMING

1. The underlying language consists of implications of the form

$$P \text{ if } Q_1 \text{ and } \dots \text{ and } Q_n.$$

where P, Q_1, \dots, Q_n are literals, i.e. atomic formulas A or the ‘classical’ negation *not* A of an atomic formula, $n \geq 0$ and *if* is an implication sign. Variables in implications are implicitly universally quantified with scope the entire implication. The only inference rules are modus ponens for *if* and instantiation for the implicit universal quantifiers.

2. Non-provability is expressed by means of a logical operator, *naf* (negation as failure), where the intended meaning of *naf* P , where P is a literal, is that P can not be shown. Therefore, *naf not* P expresses that the contrary of P can not be shown, i.e. that P is consistent.

3. Thus defeasible rules are expressed in extended logic programming by translating them into exact sentences of the form

$$P \text{ if } Q_1 \text{ and } \dots \text{ and } Q_n \text{ and } \text{naf } P_1 \text{ and } \dots \text{ and } \text{naf } P_m.$$

where $n \geq 0$, $m > 0$ and apparently free variables are implicitly universally quantified. Statements of the form $\text{naf } P$ are interpreted as atoms in the underlying language.

4. The answer set semantics of extended logic programming [Gelfond & Lifschitz, 1990] has been shown [Bondarenko et al., 1996] to be a special case of the abstract stable semantics, where the theory corresponding to an extended logic program is understood as implicitly containing all implications of the form

$$P \text{ if } Q \text{ and not } Q.$$

where P is any literal and Q is any atomic formula of the underlying language.

Normal logic programs can be understood as extended logic programs without any occurrence of 'classical' negative literals of the form $\text{not } P$. Normal logic programming is the minimal language that is an instance of the abstract framework. The stable model semantics of normal logic programming [Gelfond & Lifschitz, 1988] has been shown [Bondarenko et al., 1996] to be a special case of the abstract stable semantics.

As argued in [Kowalski, 1995; Sartor, 1991; Prakken & Sartor, 1995], extended logic programming is better suited for legal knowledge representation than normal logic programming, due to the presence of both 'classical' negation (not) and non-provability (naf).

4.3. NON-MONOTONIC MODAL LOGIC AND AUTOEPISTEMIC LOGIC

1. In non-monotonic modal logic, the underlying logic can be any modal logic with a modal operator, L , which can be interpreted as representing provability. Clearly, some modal logics are not appropriate for this purpose. However, all modal logics incorporate an inference rule, called necessitation, which sanctions deriving $L P$ from P . One of the most frequently advocated modal logics for this purpose is K45, which contains the axiom schemata

$$(K) \quad L \psi \text{ if } (L (\psi \text{ if } \phi) \text{ and } L \phi).$$

$$(4) \quad L L \phi \text{ if } L \phi.$$

$$(5) \quad L \text{ not } L \phi \text{ if not } L \phi$$

where if is material implication. Non-monotonic modal logic with K45 as underlying logic is equivalent to autoepistemic logic [Shvarts, 1990].

2. The underlying logic is already sufficiently expressive to represent non-provability. Non-provability is represented by means of formulae of the form $\text{not } L P$, expressing that P can not be shown. Therefore, $\text{not } L \text{ not } P$ expresses that the contrary of P can not be shown, i.e. that P is consistent. The first version of non-monotonic modal logic was defined in terms of a modal operator M , standing for consistency. Obviously, the two modal operators, L and M , are interdefinable.

3. Defeasible rules are expressed in non-monotonic modal logic by translating them into exact sentences of the form

$$P \text{ if } L Q_1 \text{ and } \dots \text{ and } L Q_n \text{ and not } L P_1 \text{ and } \dots \text{ and not } L P_m.$$

prefixed by explicit universally quantified variables, where $n \geq 0$ and $m > 0$.

4. The fixed point semantics of non-monotonic modal logics [McDermott, 1982] has been shown [Bondarenko et al., 1996] to be a special case of the abstract stable semantics defined in section 3.

Arbitrary non-monotonic modal logics are more expressive than default logic and extended logic programming, because they allow the expression of complex modal sentences, with nested modal operators. However, Konolige [1994] has shown that in propositional K45 any theory T can be reformulated equivalently as a theory in which all sentences have the form

$$P \text{ if } L Q_1 \text{ and } \dots \text{ and } L Q_n \text{ and not } L P_1 \text{ and } \dots \text{ and not } L P_m.$$

where the Q_i , P_j and P are non-modal, and $n, m \geq 0$.

5. Admissibility Semantics and Argumentation

Like the stability semantics, the admissibility semantics specifies that a conclusion C follows from a theory T if and only if C follows from $T \cup N$ in the underlying logic, where N is an acceptable set of non-provability claims. However, to be 'admissible', N need only contain enough non-provability claims to derive C , supplemented by sufficiently many additional non-provability claims needed to defend itself against any arguments which attack it.

Thus, the admissibility semantics can be understood in argumentation-theoretic terms. Given a theory T , an **argument** is simply a derivation in the underlying logic of some conclusion C from $T \cup N$, where N is some set of non-provability claims. In such a case, we also say both that N **supports the argument** and that the **argument is based on** N . One argument, based on N' , is said to **attack** (or **undermine**) another argument, based on N , if, for some claim S *can not be shown* in N , the first argument derives the conclusion S . In such circumstances, we also say that N' is an (undermining) attack against N and that N' attacks (or undermines) the claim S *can not be shown*.

A candidate set N of non-provability sentences is said to be **admissible** if and only if

$$\text{for every (undermining) attack } N' \text{ against } N, N \text{ attacks } N' - N.^4$$

Thus, an argument is admissible if and only if it is based upon an admissible set of non-provability claims.

Note that the stability semantics can also be understood in argumentation-theoretic terms, i.e. a set of non-provability claims is stable if and only if it does not attack itself and it attacks every non-provability claim it does not contain [Bondarenko et al., 1996].

⁴ In non-monotonic modal logic and some other logics, where non-provability sentences not in N may be derivable from $T \cup N$, the sets N and N' in the definition of admissibility should be closed (see footnote 3).

The definition of admissibility is a semantics in the sense that it specifies what is admissible, without indicating how such admissible sets of claims and arguments are to be constructed. For this latter purpose, we have developed an abstract proof procedure [Toni & Kakas, 1995; Dung et al., 1996], which applies to any instance of the framework.

The proof procedure is initiated by a proponent presenting an argument based upon sufficiently many non-provability claims N_0 to derive a desired conclusion C . The proponent then needs to defend N_0 against any attack N'_0 that an opponent might make. For this purpose, the proponent extends N_0 to a possibly larger set N_1 which, for every such attack N'_0 attacks $N_0 - N_0$. The proponent then similarly needs to defend any new claims in $N_1 - N_0$ against all attacks. The process continues with the proponent generating an increasingly larger set of supporting claims N_0, N_1, \dots, N_n .

The process terminates with a win for the proponent, with $N = N_n$, if N_n successfully counter-attacks all attacks by the opponent on the set of non-provability claims $N_n - N_{n-1}$. It terminates with a win for the opponent, if, for every sequence of moves $N_0, N_1, \dots, N_n, \dots$ by the proponent, there is an n for which the opponent has an attack against $N_n - N_{n-1}$ and the proponent has no defence.

The formal definition of the proof procedure and a formal proof of its correctness and completeness, relative to the admissibility semantics, are given in [Dung et al., 1996].

6. Reduction of Rebuttal Attacks to Undermining Attacks

The abstract admissibility semantics and its dialectic proof procedure reduce all attacks to undermining attacks. In contrast, starting from [Pollock, 1987], a number of argumentation-based frameworks [Dung, 1993; Prakken & Sartor, 1995] allow, not only undermining attacks ('ground attacks' in [Dung, 1993]), but also – **rebutting attacks** ('reductio-ad-absurdum attacks' in [Dung, 1993]), between arguments which derive contradictory conclusions.

Examples of such contradictory conclusions are *John is innocent of theft* and *John is not innocent of theft* in example 2.1 and *John inherits Henry's estate* and *John does not inherit Henry's estate* in example 2.2.

Undermining and rebutting attacks do not have equal status: undercutting attacks are stronger than rebutting attacks. Dung [1993] motivates the difference in status between the two kinds of attack by an example (example 1 in [Dung, 1993]) which has the form of the following variant of our example 2.2.

EXAMPLE 6.1. Consider the following two rules:

- r_1 : *A person inherits an estate*
 if he/she is the beneficiary of a valid will
 and it can not be shown that r_1 is defeated for the person.
- r_2 : *A person does not inherit an estate*
 if he/she has murdered the owner of the estate
 and it can not be shown that r_2 is defeated for the person.

The additional, third rule

*r_1 is defeated for a person
if it can not be shown that r_2 is defeated for that person.*

expresses the priority of the second rule (r_2) over the first (r_1). This rule might be justified, for example, by the general principle that no one should benefit from committing a crime. (In the next section, 7, we will show how to obtain this rule by compiling explicit priorities into non-provability claims.)

Suppose that, as in example 2.2, the facts

*John is the beneficiary of a will.
John has murdered Henry.*

are also given. Consider now the two conflicting arguments

*Arg_1 , based upon the set N_1 consisting of the single claim
it can not be shown that r_1 is defeated for John,
deriving the conclusion John inherits an estate, and*

*Arg_2 , based upon the set N_2 consisting of the single claim
it can not be shown that r_2 is defeated for John,
deriving the conclusion John does not inherit an estate.*

Arg_1/N_1 and Arg_2/N_2 attack each other via a rebutting attack. In addition, Arg_2/N_2 attacks Arg_1/N_1 via an undermining attack. If undermining attacks were not stronger than rebutting attacks, then the two alternative arguments would be equally admissible, and the priority given by the third rule would fail to be taken into account. Instead, if undermining attacks are stronger than rebutting attacks, then only Arg_2 is admissible.

In our argumentation framework only undercutting attacks are considered. In this example 6.1 we obtain the same intuitively correct result, as the more complicated frameworks, namely that only Arg_2 is admissible, since there is no undermining attack against Arg_2/N_2 . On the other hand, Arg_1 is defeated, not by Arg_2 , but by the third rule using the claims N_2 upon which Arg_2 is based.

The viability of our reduction of rebutting attacks to undermining attacks depends upon adopting a specific knowledge representation methodology, which introduces explicit non-provability conditions into rules. We will describe the full methodology, taking priorities into account, in the next section.

7. Reduction of Priorities to Non-Provability Conditions

Several authors, in particular Prakken and Sartor [1995] and Dimopoulos and Kakas [1995], have investigated extended logic programming with priority relations between rules. Other authors, such as Brewka [1989], have studied the problem of assigning priorities to default rules more generally. Here, we present a methodology for dealing with priorities by adding extra non-provability conditions to rules, without changing the

semantics. This methodology has been illustrated already, in condensed form, in examples 2.4 and 6.1.

For every rule which might be defeated by other rules we introduce a new predicate, which can be understood as naming the rule. The new predicate can also be interpreted as introducing a new concept standing for the way in which the rule establishes its conclusion.

For example, given several rules all of which imply that a particular person is a British citizen, the methodology introduces new predicates for each such rule. These new predicates can be thought of as representing new concepts, such as British citizenship by acquisition at birth, British citizenship by descent, or British citizenship by naturalisation, associated with different rules, all of which imply British citizenship. Arguably, the introduction of such predicates simply enables a more exact statement of the meaning of the originally given rules.

Thus each defeasible rule of the form

$$P(X) \text{ if } Q.$$

where X is a tuple of variables, is rewritten as two exact rules

$$\begin{aligned} &P(X) \text{ if holds(rule}(X)). \\ &\text{holds(rule}(X)) \text{ if } Q \text{ and defeated(rule}(X)) \text{ can not be shown.} \end{aligned}$$

where $\text{rule}(X)$ can be understood both as the name of the original rule and as a more precise name for the conclusion of the rule.

Notice that if we do not need to establish separate conclusions of the form $\text{holds(rule}(X))$, then the two rules can be condensed into one

$$P(X) \text{ if } Q \text{ and defeated(rule}(X)) \text{ can not be shown.}$$

as in examples 2.4 and 6.1.

A rule is defeated if there is a conflicting rule which holds and has higher priority:

$$\begin{aligned} &\text{defeated}(R_1(X)) \\ &\text{if } R_2(Y) > R_1(X) \\ &\text{and conflicting}(R_1(X), R_2(Y)) \\ &\text{and holds}(R_2(Y)). \end{aligned}$$

The priority relation can be defined by means of facts or by means of rules such as

$$r_1(X) > r_2(Y) \text{ if } R.$$

Similarly, the predicate *conflicting* can be defined by means of facts or by means of rules such as

$$\begin{aligned} &\text{conflicting}(R_1(X), R_2(X)) \\ &\text{if conclusion}(R_1(X), P(X)) \\ &\text{and conclusion}(R_2(X), \text{not } P(X)) \end{aligned}$$

conflicting($R_1(X, Y), R_2(X, Y)$)
 if *conclusion*($R_1(X, Y), R_3(X) > R_4(Y)$)
 and *conclusion*($R_2(X, Y), R_4(Y) > R_3(X)$).

where *conclusion*(R, P) holds when R is a rule of the form P if Q .

Notice that, given two rules

$r_1(X)$: P if Q .
 $r_2(X)$: not P if R .

with priority

$r_2(X) > r_1(X)$

the definition of defeated can be simplified to

defeated($r_1(X)$) if R and *defeated*($r_2(X)$) can not be shown.

as in examples 2.4 and 6.1.

Our treatment of priority is illustrated by the following example, adapted from [Prakken & Sartor, 1995].

EXAMPLE 7.1.

Consider the following set of defeasible rules, including priority defining rules, before our transformation:

$r_1(X)$: X 's exterior may not be modified if X is a protected building.
 $r_2(X)$: X 's exterior may be modified if X needs restructuring.
 $r_3(X, Y)$: $R_1(X) > R_2(Y)$ if $R_1(X)$ concerns artistic buildings
 and $R_2(Y)$ concerns town planning.
 $T(X, Y)$: $R_1(X) > R_2(Y)$ if $R_1(X)$ is later than $R_2(Y)$.

and the following facts/exact rules:

$r_1(X)$ concerns artistic buildings.
 $r_2(X)$ concerns town planning.
 $r_2(X)$ is later than $r_1(X)$.
 $r_3(X, Y)$ is later than $T(X, Y)$.
 villa is a protected building.
 villa needs restructuring.

After the transformation, and after simplifying the transformed rules by eliminating (using modus ponens) conditions that match the facts, we obtain

villa's exterior may not be modified if holds($r_1(villa)$).
villa's exterior may be modified if holds($r_2(villa)$).
 $r_1(villa) > r_2(villa)$ if *holds*($r_3(villa, villa)$).

$r_2(villa) > r_1(villa)$ if $holds(T(villa, villa))$.

$r_3(X, Y) > T(X, Y)$ if $holds(T((X, Y), (X, Y)))$.

$holds(r_1(villa))$ if $defeated(r_1(villa))$ can not be shown.

$holds(r_2(villa))$ if $defeated(r_2(villa))$ can not be shown.

$holds(r_3(villa, villa))$ if $defeated(r_3(villa, villa))$ can not be shown.

$holds(T(villa, villa))$ if $defeated(T(villa, villa))$ can not be shown.

$holds(T((X, Y), (X, Y)))$ if $defeated(T((X, Y), (X, Y)))$ can not be shown.

Using these rules, the definition of *defeated* can be simplified first to

$defeated(r_1(villa))$ if $holds(r_2(villa))$ and $holds(T(villa, villa))$.

$defeated(r_2(villa))$ if $holds(r_1(villa))$ and $holds(r_3(villa, villa))$.

$defeated(T(villa, villa))$ if $holds(r_3(villa, villa))$
and $holds(T((villa, villa), (villa, villa)))$.

After further simplification, eliminating the predicate *holds*, the original problem reduces to:

villa's exterior may not be modified

if $defeated(r_1(villa))$ can not be shown.

villa's exterior may be modified

if $defeated(r_2(villa))$ can not be shown.

$defeated(r_1(villa))$

if $defeated(T(villa, villa))$ can not be shown

and $defeated(r_2(villa))$ can not be shown.

$defeated(r_2(villa))$

if $defeated(r_3(villa, villa))$ can not be shown

and $defeated(r_1(villa))$ can not be shown.

$defeated(T(villa, villa))$

if $defeated(T((villa, villa), (villa, villa)))$ can not be shown

and $defeated(r_3(villa, villa))$ can not be shown.

We can now justify the conclusion that *villa's exterior may not be modified* by the following dialogue game.

Proponent: *villa's exterior may not be modified*

because $defeated(r_1(villa))$ can not be shown.

Opponent: On the contrary, $defeated(r_1(villa))$ can be shown

because $defeated(T(villa, villa))$ can not be shown

and because $defeated(r_2(villa))$ can not be shown.

Proponent: But, on the contrary, $defeated(r_2(villa))$ can be shown

because *defeated*($r_3(villa, villa)$) *can not be shown*
 and because *defeated*($r_1(villa)$) *can not be shown*,
 as I have already argued.

Alternatively, *defeated*($T(villa, villa)$) can be shown
 because
defeated($T((villa, villa), (villa, villa))$) *can not be shown*
 and because *defeated*($r_3(villa, villa)$) *can not be shown*.

Note that the proponent has two alternative ways to defeat the opponent. The second argument is ‘more conclusive’ than the first because it does not rely on its own claims to defeat the opponent. This argument can also be shown to be justified via a sceptical semantics as defined in [Bondarenko et al., 1996].

The general methodology presented in this section applies also to the case of rules having contradictory conclusions where priorities are not given explicitly, as in example 2.2. In such cases, we treat each rule as having higher priority than every other rule with a contradictory conclusion. Applying the general transformation and simplifying the transformed rules in the manner of the discussion of this section, we obtain the same result as that illustrated in example 2.2.

For example, given the two defeasible rules

P(X) if Q.
not P(X) if R.

with contradictory conclusions, the general transformation introduces rule names, say $r_1(X)$ and $r_2(X)$, for the two rules and replaces the defeasible rules by four exact rules

P(X) if holds($r_1(X)$).
holds($r_1(X)$) *if Q and defeated*($r_1(X)$) *can not be shown.*
not P(X) if holds($r_2(X)$).
holds($r_2(X)$) *if R and defeated*($r_2(X)$) *can not be shown.*

Using the two priorities

$r_1(X) > r_2(X)$.
 $r_2(X) > r_1(X)$.

the definition of *defeated* simplifies to

defeated($r_1(X)$) *if holds*($r_2(X)$).
defeated($r_2(X)$) *if holds*($r_1(X)$).

Further simplifying, by removing some of the *holds* conditions, we obtain

P(X) if Q and defeated($r_1(X)$) *can not be shown.*
holds($r_1(X)$) *if Q and defeated*($r_1(X)$) *can not be shown.*
not P(X) if R and defeated($r_2(X)$) *can not be shown.*

holds($r_2(X)$) if R and $\text{defeated}(r_2(X))$ can not be shown.
 $\text{defeated}(r_1(X))$ if $\text{holds}(r_2(X))$.
 $\text{defeated}(r_2(X))$ if $\text{holds}(r_1(X))$.

But now the conclusions *holds($r_1(X)$)* and *P(X)* are interchangeable, as are *holds($r_2(X)$)* and *not P(X)*. Therefore, rewriting *holds($r_1(X)$)* as *P(X)* and *holds($r_2(X)$)* as *not P(X)*, we obtain

P(X) if Q and $\text{defeated}(r_1(X))$ can not be shown.
not P(X) if R and $\text{defeated}(r_2(X))$ can not be shown.
 $\text{defeated}(r_1(X))$ if not P(X).
 $\text{defeated}(r_2(X))$ if P(X).

Eliminating the predicate *defeated* gives

P(X) if Q and not P(X) can not be shown.
not P(X) if R and P(X) can not be shown.

which is exactly the result of the simplified transformation needed to convert defeasible rules into exact rules, when no priorities are explicitly given. This simplification can be generalised to the general case where there are several rules having conclusion *P(X)* or *not P(X)*.

In all the examples considered up until now, we obtain the same results as those obtained by Prakken and Sartor [1995] (but we employ a credulous semantics while they employ a sceptical one). In general, however, our treatment of priorities gives different results from those of Prakken and Sartor [1995], as illustrated by the following example.

EXAMPLE 7.2. (Taken from [Prakken and Sartor, 1995])

Consider the defeasible rules

r_1 : *A.*
 r_2 : *B if A.*
 r_3 : *not A if B can not be shown.*

with $r_3 > r_1$. Then, according to [Prakken and Sartor, 1995], the derivation consisting of r_1 and r_2 is a justified argument, while the derivation consisting of r_3 and the claim *B can not be shown* is not justified. Therefore, they derive *A* but not *not A*. Instead, our method produces three admissible arguments:

Arg₁ concluding A,
based upon the set of claims N_1 consisting of
 $\text{defeated}(r_1)$ can not be shown and $\text{defeated}(r_2)$ can not be shown,

Arg₂ concluding not A,
based upon the set of claims N_2 consisting of
 B can not be shown and $\text{defeated}(r_3)$ can not be shown,

*Arg₃ concluding not A,
 based upon the set of claims N₃ consisting of
 B can not be shown, defeated(r₃) can not be shown and
 defeated(r₂) can not be shown.*

Note that our approach takes into account the given priority, while their approach does not.

8. Conclusion

We have argued that our abstract framework for defeasible reasoning, with the admissibility semantics and its dialectic proof procedure, is adequate for the formalisation of many aspects of practical reasoning. We have outlined a methodology for transforming inexact, defeasible rules into exact rules with explicit non-provability conditions; and we have argued that this transformation eliminates the need for rebuttal attacks and for dealing with priorities in the semantics. Our transformation is an elaboration of transformations [Kowalski & Sadri, 1990; Toni & Kowalski, 1995] we have developed earlier.

Because our approach is abstract, it can be formalised in any one of the many formalisms, including default logic, extended logic programming, non-monotonic modal logic, and auto-epistemic logic, which are special cases of our framework.

It needs to be noted, however, that our representation of priorities uses metapredicates, such as ‘holds’ and ‘defeated’, in the manner of meta-logic programming. In this respect, our approach bears many resemblances to that of Hage [1995], who uses similar metapredicates. However, his approach differs from ours in his introduction of a new formalism for this purpose, while our approach can be employed with many existing formalisms. Moreover, as we saw in several examples and in one general case, in our simple use of them, these meta-predicates can generally be eliminated in favour of simpler non-provability conditions.

In this paper, we focussed entirely on the credulous admissibility semantics. However, as has been shown elsewhere [Bondarenko et al., 1996], the abstract framework also admits sceptical semantics and proof procedures.

Although we have given several examples to show how our methodology eliminates the need for rebuttal attacks and for dealing with priorities in the semantics, we do not have any formal results proving that the methodology is always adequate. This remains to be done for future work.

Acknowledgements

This research was partially supported by Fujitsu Research Laboratories. The authors are grateful to Phan Minh Dung and Henry Prakken for many helpful discussions, and to the referees and Henry Prakken for helpful suggestions.

References

- A. Bondarenko, P. M. Dung, R. A. Kowalski, F. Toni (1996) An abstract, argumentation-theoretic framework for default reasoning. To appear in *Artificial Intelligence*, Elsevier.
- G. Brewka (1989) Preferred subtheories: an extended logical framework for default reasoning. *Proceedings of the 11th International Joint Conference on Artificial Intelligence*, Detroit, Mi Morgan Kaufmann (N. Sridharan, ed.) 1043–1048
- Y. Dimopoulos, A. C. Kakas (1995) Logic programming without negation as failure. *Proceedings of the International Logic Programming Symposium*, Portland, Oregon MIT Press (J. Lloyd, ed.) 369–384
- P. M. Dung (1993) The acceptability of arguments and its fundamental role in non-monotonic reasoning and logic programming. *Proceedings of the 13th International Joint Conference on Artificial Intelligence*, Chambéry, France Morgan Kaufmann (R. Bajcsy, ed.) 852–857
- Dung, P. M. (1993) An argumentation semantics for logic programming with explicit negation. *Proceedings of the 10th International Conference on Logic Programming*. Paris MIT Press (K. Furukawa, ed.)
- P. M. Dung, R. A. Kowalski, F. Toni (1996) Argumentation-theoretic proof procedures for non-monotonic reasoning. *Proc. 6th Logic Programming Synthesis and Transformation* (J. Gallagher ed.)
- M. Gelfond, V. Lifschitz (1988) The stable model semantics for logic programming. *Proceedings of the 5th International Conference on Logic Programming*, Washington, Seattle MIT Press (K. Bowen and R. A. Kowalski, eds.) 1070–1080
- M. Gelfond, V. Lifschitz (1990) Logic programs with classical negation. *Proceedings of the 7th International Conference on Logic Programming*, Jerusalem MIT Press (D. H. D. Warren and P. Szeredi, eds.) 579–597
- T. F. Gordon (1993) The pleadings game: an exercise in computational dialectics, *Artificial Intelligence and Law* 2(4), Kluwer Academic Publishers 239–292
- J. Hage (1995) Teleological reasoning in reason-based logic. *Proceedings of the 5th International Conference on Artificial Intelligence and Law*, College Park, MD ACM Press, 11–20
- K. Konolige (1994) Autoepistemic logic. *Handbook of Logic in Artificial Intelligence and Logic Programming* 3, Oxford University Press (D. Gabbay, C. Hogger, J. A. Robinson, eds.)
- R. A. Kowalski (1995) Legislation as logic programs. *Informatics and the Foundations of Legal Reasoning*, Kluwer Academic Publishers (Z. Bankowski et al., eds.) 325–356
- Kowalski, R. A., Sadri, F. (1990) Logic programs with exceptions. *Proceedings of the 7th International Conference on Logic Programming*, Jerusalem MIT Press (D. H. D. Warren and P. Szeredi, eds.) 598–613
- D. McDermott (1982) Nonmonotonic logic II: non-monotonic modal theories. *Journal of ACM* 29(1) 33–57
- R. Moore (1985) Semantical considerations on non-monotonic logic. *Artificial Intelligence* 25, Elsevier 75–94
- J. L. Pollock (1987) Defeasible reasoning. *Cognitive Science*, 11 481–518
- H. Prakken, G. Sartor (1995) On the relation between legal language and legal argument: assumptions, applicability and dynamic priorities. *Proceedings of the 5th International Conference on Artificial Intelligence and Law* College Park, MD ACM Press, 1–10
- H. Prakken, G. Sartor (1996) A dialectical model of assessing conflicting arguments in legal reasoning. In this issue.
- R. Reiter (1980) A logic for default reasoning. *Artificial Intelligence* 13, Elsevier 81–132
- G. Sartor (1991) The structure of norm conditions and non-monotonic reasoning in law. *Proceedings of the 3th International Conference on Artificial Intelligence and Law*, Oxford ACM Press, 155–164
- G. Shvarts (1990) Autoepistemic modal logics. *Proc. 3rd Conference on Theoretical Aspects of Rationality and Knowledge*, Pacific Grove, CA Morgan Kaufmann (R. Parikh, ed.) 97–110
- F. Toni, A. C. Kakas (1995) Computing the acceptability semantics. *Proceedings of the 3rd International Workshop on Logic Programming and Non-monotonic Reasoning*, Springer Verlag LNAI 928 (V. Marek, A. Nerode, M. Truszczynski, eds.) 401–415
- F. Toni, R. A. Kowalski (1995) Reduction of abductive logic programs to normal logic programs. *Proceedings of the 12th International Conference on Logic Programming*, Japan MIT Press (Leon Sterling, ed.) 367–381