

An Ambiguous Grammar for Mathematics

Richard Warfield

Abstract

Mathematical expressions as they appear in technical documents cannot be confined to a single simple grammatical class. We have developed a framework for more realistic understanding of such expressions than is possible from a traditional programming language parsing perspective. One application for this is for reading expressions from a technical document in preparation for storing them in a content-retrievable form or manipulating them further. In this case a relatively simple character-recognition and sorting technique can often convert 2-dimensional expressions into 1-dimensional strings. Applying the grammar and parser presented here can find one (or more, in cases of ambiguity) parses for such a string. Another application is the reading of “naive” input to a web page requesting mathematical input.

Introduction

The grammar described in the following sections is intended for parsing a one-dimensional representation of mathematical equations as they would appear in publication (as opposed to the way they appear in computer programs). A typical grammar for parsing simple expressions as they appear in, for example, a C program is shown below¹.

$$\begin{aligned} \text{expr} &\rightarrow \text{expr} + \text{term} \mid \text{expr} - \text{term} \mid \text{term} \\ \text{term} &\rightarrow \text{term} * \text{factor} \mid \text{term} / \text{factor} \mid \text{factor} \\ \text{factor} &\rightarrow (\text{expr}) \mid -\text{factor} \mid \mathbf{ident} \mid \mathbf{ident}(\text{expr}) \end{aligned}$$

The language specified by this grammar has several properties not present in the language of "real" math that make it easy to parse with an LALR parser:

- It is context free. For example, because multiplication always requires * to be explicitly present, there is no ambiguity as to whether an expression such as A(X+Y) is a function application or a multiplication (A*X + A*Y).
- Any binary operation can be applied to either side of any term. This would not be true if function applications without parentheses and products without an asterisk were allowed in the language. For example, "sin x" cannot be multiplied on the right by the identifier "y", since "sin xy" clearly means sin(x*y) and not (sin(x)) * y.
- The location of all operands to these unary and binary operators is clear from the syntax. In contrast, the location of the variable of integration in an indefinite integral, for example, can only be found by digging through the operand looking for a differential.

Given that mathematical notation was (unfortunately) not designed with context-free

parsers in mind, parsing this language presents some difficulties that require many complex grammar rules (relative to the short grammar presented above) to overcome, some difficulties that are impossible to overcome without semantic analysis (such as the function application versus multiplication ambiguity), and some which can be overcome easily by making certain assumptions about the language.

The following is a partial list of the structures and operations that are recognized by this parser and limitations and assumptions imposed on them:

<i>Expressions</i>	<i>Limitations</i>
Variables and Function names	Only one-character identifiers are supported, except for a list of functions recognized by the parser (sin, cos, log, etc...)
Binary operators +, -, *, /	"/" operator requires three sets of parentheses in some expressions to produce a correct parse: one around the numerator, one around the denominator, and one around the entire fraction ¹ .
Products with no operator (such as $ab = a*b$)	
Exponentiation	
Function applications with explicit parentheses (such as $f(x)$ or $\tan(x+y)$)	Whenever an expression such as $f(x)$ or $a(b+c)$ is encountered, an ambiguous parse is produced. Without semantic information it is impossible to distinguish between a function application of this form and the product of a variable and expression in parentheses
Function applications with no parentheses (such as $\sin xy$)	Operand must consist of either a product of non-parenthesised variables and numbers, or some other sort of single factor.

¹ The justification for this seemingly severe limitation is that this information is present in fractions as they appear on a two-dimensional page, since the fraction bar delimits the numerator, the denominator, and the fraction itself from an enclosing expression. In general, I have assumed that no additional ambiguity is introduced in the process of converting the two-dimensional expression into a one-dimensional expression.

<i>Expressions</i>	<i>Limitations</i>
Definite and indefinite integrals	The differential must appear in one of two places: either on the far right of the integrand (such as in $\int ax + bx^2 dx$) or on the far right of the numerator of the rightmost factor (such as in $\int \sin x \frac{xdx}{(x+1)}$)
Summations	Associativity of summation with the * operator is ambiguous: $\sum_{i=0}^n x * y$ will produce two parses, one with y inside the sum and one with y outside
Relational operators = , ≤ , ≥ , < , >	

Overview of the Parser and Grammar

This grammar was developed using the context-free parser from chapter 19 of *Paradigms of Artificial Intelligence Programming* by Peter Norvig² together with a scanner written by Richard Fateman. The parser is written in Lisp and accepts a list of grammar rules as input. It is capable of producing multiple parses in the case of ambiguous grammars. There is no distinguished start symbol in the grammars used by this parser; any reduction that uses up all of the input results in a valid parse (duplicate parses are removed, however).

The following sections give a rough description of the grammar. Details specific to the parser and scanner being used are omitted. The grammar is structured loosely the way the computer language expression grammar shown above is. At the top level is the nonterminal *S*, which encapsulates the relational operators:

$S \rightarrow S = expr$
 $S \rightarrow S < expr$
 $S \rightarrow S > expr$
 $S \rightarrow S \leq expr$
 $S \rightarrow S \geq expr$
 $S \rightarrow expr$

The *expr* nonterminal appears in its usual form:

$expr \rightarrow expr + term$
 $expr \rightarrow expr - term$
 $expr \rightarrow term$

When we come to the *term* nonterminal, the grammar departs rather dramatically from the usual expression grammar. As discussed above, certain terms cannot be multiplied on the right by some kinds of factors unless an asterisk or parentheses are used. In other words, if α is the term in question then $\alpha\beta$ is not allowed for some factors β . *term* is divided into three nonterminals based on what β are allowed. One nonterminal, *lara_term* can be multiplied on the left by anything and on the right by anything. This nonterminal includes simple numbers and variables, function applications with parentheses, all parenthesized subexpressions, and products of these entities. The second term type, *larf_term*, can be multiplied on the left by anything on the right by (non-variable, non-number) factors. It includes function applications without parentheses, such as $\sin x$. The third type, *larn_term*, can be multiplied on the left by anything and on the right by nothing. It includes summations and integrations.

The basic *term* nonterminal appears only to attach the other term types into *expr*:

$$term \rightarrow term * any_term$$

$$term \rightarrow term / any_term$$

$$term \rightarrow any_term$$

$$any_term \rightarrow lara_term$$

$$any_term \rightarrow larf_term$$

$$any_term \rightarrow larn_term$$

$$any_term \rightarrow -any_term$$

Function Applications without Parentheses : *larf_term*

This section describes the nonterminal *larf_term* and related nonterminals. *larf_term* represents structures which may appear in a product only if the factor immediately to the right is not a number or variable. Function applications such as $\sin x$ and $\coth 2nx$ fall into this category because an expression like the latter should be interpreted as $\coth(2nx)$ and not as $\coth(2n)*x$ or $\coth(2)*nx$. However, an expression such as $\sin x(a+b)$ is likely to mean $(\sin x)(a+b)$. Other structures may fall into this category as well and the grammar could easily be extended to include them.

$$larf_term \rightarrow larf_term\ larf_factor$$

$$larf_term \rightarrow lara_term\ larf_factor$$

$$larf_term \rightarrow larf_factor$$

$$larf_factor \rightarrow funapp_noparens$$

$$funapp_noparens \rightarrow fun\ funarg_noparens$$

$$funapp_noparens \rightarrow fun \wedge left_delimited_item\ funarg_noparens$$

The first grammar rule allows products of *larf_factors* (for example, $\sin x \cos x$). The

second rule says that when a *lara_term* is multiplied on the right by a *larf_term*, the result is a *larf_term* because it still cannot be multiplied on the right by variables or numbers. The last grammar rule allows for expressions like $\sin^2 x$. *funarg_noparens* rewrites as:

funarg_noparens \rightarrow *var_prod*
funarg_noparens \rightarrow *lara_factor*
funarg_noparens \rightarrow *left_delimited_item*
funarg_noparens \rightarrow $-$ *funarg_noparens*

var_prod \rightarrow *var_prod* *var_power*
var_prod \rightarrow *var_power*

var_power \rightarrow **var** ^ *left_delimited_item*
var_power \rightarrow **number** ^ *left_delimited_item*
var_power \rightarrow **var**
var_power \rightarrow **number**

var_prod represents a product of numbers and variables possibly raised to powers (for example, $7nx^2y$). *left_delimited_item* represents essentially anything that can act as a factor:

left_delimited_item \rightarrow *delimited_nonvar_factor*
left_delimited_item \rightarrow *larf_factor*
left_delimited_item \rightarrow **var**
left_delimited_item \rightarrow **number**
left_delimited_item \rightarrow $-$ *left_delimited_item*

delimited_nonvar_factor includes all non-variable, non-number factors that are both left-delimited and right-delimited. By delimited I mean that the boundary of the factor is marked by some terminal or otherwise doesn't depend on the surrounding expression. For example, anything surrounded by parentheses is both left and right-delimited, as are function applications with parentheses. Function applications without parentheses are not right-delimited because adding additional characters to the right may (or may not) alter the original factor (e.g., adding "y" to the right of "sin x" changes the expression to "sin xy", but adding "cos y" does not change the argument to sin).

delimited_nonvar_factor \rightarrow (*expr*)
delimited_nonvar_factor \rightarrow *funapp_parens*

funapp_parens \rightarrow *fun* (*exprlist*)

Integrations and Summations : *larn_term*

larn_term represents mathematical structures which cannot be multiplied on the right

without some form of explicit separation of factors (parentheses or an asterisk). Integrations and summations fall into this category: $\sum xy$ of course means $\sum (xy)$ rather than $(\sum x)y$, and an expression such as $\int x dx \frac{x}{1-x}$ is presumed to mean

$$\int \frac{x^2}{1-x} dx \text{ rather than } (\int x dx) \left(\frac{x}{1-x} \right).$$

The structure of *larn_term* is similar to that of *larf_term*:

larn_term → *larf_term larn_factor*
larn_term → *lara_term larn_factor*
larn_term → *larn_factor*

larn_factor → *summation*
larn_factor → *integration*

The first two rules say that when either a *larf_term* or *lara_term* is multiplied on the right by a *larn_term*, the result is a *larn_term* since the rightmost factor is an integration or summation. The rule for summation is:

summation → *sigma any_term*

where *simga* is a terminal or sequence of nonterminals representing \sum and possibly the bounds of the summation. Note that this rule does not allow * to appear in the summand unless there are parentheses around the expression. An alternative rule would use *term* instead of *any_term*. This rule would result in an ambiguous parse when * is encountered in the top-level of the summand.

Rules for correctly parsing integration are more difficult to formulate. Such a set of rules would have to extract the variable of integration by finding a differential (such as dx), and ensure that the integrand contains exactly one such differential. Such a set of rules would necessarily be long and complicated, since it would have to keep track of whether or not the variable of integration has been while the integrand is built. An alternative solution is to assume that the differential will always appear in one of two places: either to the far right of the integrand (as in $\int x + x^2 dx$) or, if the integrand consists of a single term, at the right of the numerator of the rightmost fraction (as in

$$\int \sin x \cos x \frac{x^2 dx}{x-1}).$$

I represent this rightmost factor containing the differential by the nonterminal *diffntl_fraction*. Rules to parse integrals, given these assumptions, are given below:

integration → *integral expr differential*
integration → *integral integrand*

integrand → *integrand / any_term*
integrand → *term differential*
integrand → *differential*
integrand → *term diffntl_factor*
integrand → *diffntl_factor*

diffntl_factor → (*integrand*)

where *integral* is some terminal or sequence of terminals used to represent \int and possibly the limits of integration. These rules work by starting with the differential and building the integrand from right to left, and from top to bottom (as in a fraction).

Everything Else: *lara_term*

lara_term represents any term which can be multiplied by any factor on either side. Most single-factor terms fall into this category, including variables, numbers, parenthesized subexpressions, and function applications with parentheses. It also includes these items when raised to exponents, and it includes function applications without parentheses that have been multiplied on the right by anything other than another function application without parentheses:

lara_term → *lara_term lara_factor*
lara_term → *lara_factor*
lara_term → *larf_term delimited_nonvar_factor*

lara_factor → *lara_factor ^ left_delimited_item*
lara_factor → *delimited_factor*

delimited_factor → *delimited_nonvar_factor*
delimited_factor → *var*
delimited_factor → *number*

Conclusion

I have presented a grammar that covers the most common structures used with the operators supported. It could be expanded to cover many other structures. However, I find that the way different types of mathematical expressions interact with each other tends to be somewhat unpredictable. As a result adding support for a new structure tends to require rethinking of the way others are implemented, despite efforts to keep the

grammar modular.

A semantic processing phase following the parsing phase would allow some of the difficult syntactic checks to be offloaded from the parser. It could potentially also allow some of the context-sensitive areas of the language to be processed. For example, the search for the variable integration in an integrand could be performed by allowing any expression to be in an integrand then throwing away parses that have integrands without a unique differential in a legitimate location.

That is only one of the possible directions this method of processing formulas on a page could take. A variety of techniques could doubtless improve the balance between efficiency, accuracy, and flexibility of the procedure.

Appendix A: The Complete Grammar

The following is a complete copy of the full grammar in one piece. Semantic rules for constructing the output are not included. Nonterminal symbols are shown in italics. Terminal symbols, some of which are described below, are shown in bold:

- **var** represents single-character identifiers. It could also be implemented as a non-terminal representing other sorts of variables as well (such as subscripted variables or special symbols).
- **integral** represents an integration sign: \int . This could also be implemented as a nonterminal incorporating more information about the integration (limits, for example).
- **sigma** represents \sum . It could also be implemented as a nonterminal incorporating information about the bounds of the summation.

$S \rightarrow S = expr$

$S \rightarrow S < expr$

$S \rightarrow S > expr$

$S \rightarrow S \leq expr$

$S \rightarrow S \geq expr$

$S \rightarrow expr$

$expr \rightarrow expr + term$

$expr \rightarrow expr - term$

$expr \rightarrow term$

$term \rightarrow term * any_term$

$term \rightarrow term / any_term$

$term \rightarrow any_term$

$any_term \rightarrow lara_term$

$any_term \rightarrow larf_term$

$any_term \rightarrow larn_term$

$any_term \rightarrow -any_term$

$lara_term \rightarrow lara_term\ lara_factor$
 $lara_term \rightarrow lara_factor$
 $lara_term \rightarrow larf_term\ delimited_nonvar_factor$

$lara_factor \rightarrow lara_factor \wedge left_delimited_item$
 $lara_factor \rightarrow delimited_factor$

$delimited_factor \rightarrow delimited_nonvar_factor$
 $delimited_factor \rightarrow \mathbf{var}$
 $delimited_factor \rightarrow \mathbf{number}$

$delimited_nonvar_factor \rightarrow (expr)$
 $delimited_nonvar_factor \rightarrow funapp_parens$
 $delimited_nonvar_factor \rightarrow \text{'sqrt'} (expr)$

$left_delimited_item \rightarrow delimited_nonvar_factor$
 $left_delimited_item \rightarrow larf_factor$
 $left_delimited_item \rightarrow \mathbf{var}$
 $left_delimited_item \rightarrow \mathbf{number}$
 $left_delimited_item \rightarrow -\ left_delimited_item$

$larf_term \rightarrow larf_term\ larf_factor$
 $larf_term \rightarrow lara_term\ larf_factor$
 $larf_term \rightarrow larf_factor$

$larf_factor \rightarrow funapp_noparens$

$funapp_parens \rightarrow fun (exprlist)$

$exprlist \rightarrow exprlist , expr$
 $exprlist \rightarrow expr$

$funapp_noparens \rightarrow fun\ funarg_noparens$
 $funapp_noparens \rightarrow fun \wedge left_delimited_item\ funarg_noparens$

funarg_noparens → *var_prod*
funarg_noparens → *lara_factor*
funarg_noparens → *left_delimited_item*
funarg_noparens → *- funarg_noparens*

var_prod → *var_prod var_power*
var_prod → *var_power*

var_power → **var** ^ *left_delimited_item*
var_power → **number** ^ *left_delimited_item*
var_power → **var**
var_power → **number**

larn_term → *larf_term larn_factor*
larn_term → *lara_term larn_factor*
larn_term → *larn_factor*

larn_factor → *summation*
larn_factor → *integration*

integration → *integral expr differential*
integration → *integral integrand*

integrand → *integrand / any_term*
integrand → *term differential*
integrand → **differential**
integrand → *term diffntl_factor*
integrand → *diffntl_factor*

diffntl_factor → (*integrand*)

summation → **sigma** *any_term*

Appendix B: Test Cases and Sample Output

The following is a group of test inputs designed to demonstrate the capabilities and limitations of the grammar. I have arbitrarily chosen the following mappings from two-dimensional mathematics to the one-dimensional string of characters that is the input to the parser:

$$\int x dx \rightarrow \text{'[:integral:]xdx'}$$

$$\int_a^b x dx \rightarrow \text{'[:integral a, b:]xdx'}$$

$$\sum_a^b x \rightarrow \text{'[:sum a, b:]x'}$$

The output is shown in LISPstyle prefix notation. Indefinite integrals are shown as:

(INTEGRAL <integrand> <variable>)

Definite integrals are shown as:

(INTEGRAL <integrand> <variable> <lower bound> <upper bound>)

Summations are shown as:

(SUM <expression> <initialization statement> <upper bound>)

Here are the test cases. Note that, where possible, the output has been simplified and redundant parses removed. I have used a scanner and simplifier written by Richard J. Fateman together with Norvig's parser:

- "xy" \Rightarrow (* x y)
- "(a+b)(c+d)" \Rightarrow (* (+ c d) (+ a b))
- "abc/xyz" \Rightarrow (/ (* a b c) (* x y z))
- "f(x)" \Rightarrow (f x) AND (* f x)
- "c(a+b)" \Rightarrow (c (+ a b)) AND (* c (+ a b))
- "xsinx" \Rightarrow (* x (sin x))
- "tanxy" \Rightarrow (tan (* x y))

- "xsinxtanxy" $\Rightarrow (* x (\sin x) (\tan (* x y)))$
- "xsinycos-xztan(a+bx)" $\Rightarrow (* (\cos (* -1 x z))$
 x
 $(\tan (+ a (* b x)))$
 $(\sin (* x y)))$
- "sinxycosxy/sin5xycos5xy"
 $\Rightarrow (/ (* (\cos (* x y))$
 $(\sin (* x y)))$
 $(* (\sin (* 5 x y))$
 $(\cos (* 5 x y))))$
- "sin^2x + cos^2x = 1" $\Rightarrow (= (+ (\text{EXPT} (\sin x) 2) (\text{EXPT} (\cos x) 2)) 1)$
- "sin^3xyz" $\Rightarrow (\text{EXPT} (\sin (* x y z)) 3)$
- "[:integral:]xdx" (meaning $\int x dx$)
 $\Rightarrow ((\text{INTEGRATE } x x))$
- "[:integral:]x" $\Rightarrow \text{NIL}$
(error, because there is no differential)
- "[:integral a,b:]x^2y^2tanxdx" (meaning $\int_a^b x^2 y^2 \tan x dx$)
 $\Rightarrow (\text{INTEGRATE } (* (\tan x)$
 $(\text{EXPT } x 2)$
 $(\text{EXPT } y 2))$
 $x a b)$
- "[:integral a,b:][:integral c,d:]xydxdy" (meaning $\int_a^b \int_c^d xy dxdy$)
 $\Rightarrow (\text{INTEGRATE } (\text{INTEGRATE } (* x y) x c d)$
 y
 a
 $b)$
- "[:integral:]lnx(xdx/(1+x))" (meaning $\int \ln x \frac{xdx}{1+x}$)
 $\Rightarrow (\text{INTEGRATE } (/ (* x (\ln x)) (+ 1 x))$
 $x)$
- "[:integral:]dxsinx" $\Rightarrow \text{NIL}$
(an error because dx is not on the far right
of the integrand)

• "[:integral:](1+x+px^2+p^2q^-2cos^2x)/(1-pq)(dx/x)"

(meaning $\int \frac{1+x+px^2+p^2q^{-2}\cos^2x}{1-pq} \frac{dx}{x}$)

=>
 (INTEGRATE (/ (+ 1 x (* p (EXPT x 2))
 (* (EXPT (cos x) 2) (EXPT q -2) (EXPT p 2)))
 (* x (+ 1 (* -1 p q)))))
 x)

• "[:integral:]sintcosx+sinxcostdx"

=> (INTEGRATE (+ (* (sin x) (cos t))
 (* (sin t) (cos x)))
 x)

• "[:integral:](xydx)/z"

=> (INTEGRATE (/ (* x y) z) x)

• "[:integral:]xydx/z" => (/ (INTEGRATE (* x y) x) z)
 AND (INTEGRATE (/ (* x y) z) x)
 (division really requires parentheses to work
 predictably, especially where integrals are involved)

• "[:integral:]sinxy(sinxcosxdx)/x"

=>
 (/ (INTEGRATE (* (|sin| |x|) (|sin| (* |x| |y|)) (|cos| |x|)) |x|) |x|)
 AND (INTEGRATE (/ (* (|sin| |x|) (|sin| (* |x| |y|)) (|cos| |x|)) |x|) |x|))

• "n[:sum i=1, n:]3i^2" (meaning $n \sum_{i=1}^n 3i^2$)

=> ((* n (SUM (* 3 (EXPT i 2)) (= i 1) n)))

• "[:sum i=0,10:][:sum j=0,10:]ij" (meaning $\sum_{i=0}^{10} \sum_{j=0}^{10} ij$)

=>(SUM (SUM (* i j) (= j 0) 10) (= i 0) 10)

• "[:sum i=0, n:]i+1" => (+ 1 (SUM i (= i 0) n))
 (note that summation takes precedence over + and -)

• "[:sum i=0, n:]3i^3+[:sum j=0, m:]2j^2"

=> (+ (SUM (* 2 (EXPT |j| 2)) (= |j| 0) |m|)
 (SUM (* 3 (EXPT |i| 3)) (= |i| 0) |n|))

- 1 Aho, Alfred, Ravi Sethi, and Jefferey D. Ullman, *Compilers: Principles, Techniques, and Tools*, Addison-Wesley, Reading, Massachusetts, 1986
- 2 Norvig, Peter, *Artificial Intelligence Programming: Case Studies in Common Lisp*, Morgan Kaufmann Publishers, San Mateo, California, 1992