

Apr 13, 2020

The goals of this particular article are:

- This post assumes a fair bit of familiarity with parsing techniques, and, for example, does not explain what a context free grammar is.

Parsing is the process by which a compiler turns a *sequence* of tokens into a *tree* representation:

There are many approaches to this task, which roughly fall into one of the broad two categories:

- Using a DSL to specify an abstract grammar of the language
- Hand-writing the parser

Pratt parsing is one of the most frequently used techniques for handwritten parsing.

The pinnacle of syntactic analysis theory is discovering the context free grammar notation (often using BNF concrete syntax) for decoding linear structures into trees:

```
1 Item =  
2     StructItem  
3     | EnumItem  
4     | ...  
5  
6 StructItem =  
7     'struct' Name '{' FieldList '}'  
8  
9     ...
```

I remember being fascinated by this idea, especially by parallels with natural language sentence structure. However, my optimism quickly waned once we got to describing expressions. The natural expression grammar indeed allows one to see what is an expression.

```
1 Expr =  
2     Expr '+' Expr  
3     | Expr '*' Expr  
4     | '(' Expr ')'  
5     | 'number'
```

Although this grammar looks great, it is in fact ambiguous and imprecise, and needs to be rewritten to be amendable to automated parser generation. Specifically, we need to specify precedence and associativity of operators. The fixed grammar looks like this:

```
1 Expr =  
2     Factor  
3     | Expr '+' Factor  
4  
5 Factor =  
6     Atom  
7     | Factor '*' Atom  
8  
9 Atom =  
10    'number'  
11    | '(' Expr ')'
```

To me, the “shape” of expressions feels completely lost in this new formulation. Moreover, it took me three or four *courses* in formal languages before I was able to reliably create this grammar myself.

And that’s why I love Pratt parsing — it is an enhancement of recursive descent parsing algorithm, which uses the natural terminology of precedence and associativity for parsing expressions, instead of grammar obfuscation techniques.

Recursive descent and left-recursion

The simplest technique for hand-writing a parser is recursive descent, which models the grammar as a set of mutually recursive functions. For example, the above item grammar fragment can look like this:

```

1 fn item(p: &mut Parser) {
2     match p.peek() {
3         STRUCT_KEYWORD => struct_item(p),
4         ENUM_KEYWORD   => enum_item(p),
5         ...
6     }
7 }
8
9 fn struct_item(p: &mut Parser) {
10    p.expect(STRUCT_KEYWORD);
11    name(p);
12    p.expect(L_CURLY);
13    field_list(p);
14    p.expect(R_CURLY);
15 }
16
17 ...

```

Traditionally, text-books point out left-recursive grammars as the Achilles heel of this approach, and use this drawback to motivate more advanced LR parsing techniques. An example of problematic grammar can look like this:

```

1 Sum =
2     Sum '+' Int
3     | Int

```

Indeed, if we naively code the `sum` function, it wouldn't be too useful:

```

1 fn sum(p: &mut Parser) {
2     // Try first alternative
3     sum(p); ❶
4     p.expect(PLUS);
5     int(p);
6
7     // If that fails, try the second one
8     ...
9 }

```

❶ At this point we immediately loop and overflow the stack

A theoretical fix to the problem involves rewriting the grammar to eliminate the left recursion. However in practice, for a hand-written parser, a solution is much simpler — breaking away with a pure *recursive* paradigm and using a loop:

```

1 fn sum(p: &mut Parser) {
2     int(p);
3     while p.eat(PLUS) {
4         int(p);
5     }
6 }

```

Pratt parsing, the general shape

Using just loops won't be enough for parsing infix expressions. Instead, Pratt parsing uses *both* loops *and* recursion:

```

1 fn parse_expr() {
2     ...
3     loop {
4         ...
5         parse_expr()
6         ...
7     }
8 }

```

Not only does it send your mind into Möbeus-shaped hamster wheel, it also handles associativity and precedence!

From Precedence to Binding Power

I have a confession to make: I am always confused by “high precedence” and “low precedence”. In $a + b * c$, addition has a lower precedence, but it is at the top of the parse tree...

So instead, I find thinking in terms of binding power more intuitive.

```

1 | expr:   A       +       B       *       C
2 | power:   3       3       5       5

```

The $*$ is stronger, it has more power to hold together B and C, and so the expression is parsed as $A + (B * C)$.

What about associativity though? In $A + B + C$ all operators seem to have the same power, and it is unclear which $+$ to fold first. But this can also be modeled with power, if we make it slightly asymmetric:

```

1 | expr:   A       +       B       +       C
2 | power:  0       3       3.1     3       3.1     0

```

Here, we pumped the right power of $+$ just a little bit, so that it holds the right operand tighter. We also added zeros at both ends, as there are no operators to bind from the sides. Here, the first (and only the first) $+$ holds both of its arguments tighter than the neighbors, so we can reduce it:

```

1 | expr:   (A + B)   +       C
2 | power:  0       3       3.1     0

```

Now we can fold the second plus and get $(A + B) + C$. Or, in terms of the syntax tree, the second $+$ really likes its right operand more than the left one, so it rushes to get hold of C. While he does that, the first $+$ captures both A and B, as they are uncontested.

What Pratt parsing does is that it finds these badass, stronger than neighbors operators, by processing the string left to right. We are almost at a point where we finally start writing some code, but let's first look at the other running example. We will use function composition operator, $.$ (dot) as a *right* associative operator with a high binding power. That is, $f . g . h$ is parsed as $f . (g . h)$, or, in terms of power

```

1 | f       .       g       .       h
2 | 0     8.5     8     8.5     8     0

```

Minimal Pratt Parser

We will be parsing expressions where basic atoms are *single character* numbers and variables, and which uses punctuation for operators.

Let's define a simple tokenizer:

```

1  #[derive(Debug, Clone, Copy, PartialEq, Eq)]
2  enum Token {
3      Atom(char),
4      Op(char),
5      Eof,
6  }
7
8  struct Lexer {
9      tokens: Vec<Token>,
10 }
11
12 impl Lexer {
13     fn new(input: &str) -> Lexer {
14         let mut tokens = input
15             .chars()
16             .filter(|it| !it.is_ascii_whitespace())
17             .map(|c| match c {
18                 '0'..'9' |
19                 'a'..'z' | 'A'..'Z' => Token::Atom(c),
20                 _ => Token::Op(c),
21             })
22             .collect::<Vec<_>>();
23         tokens.reverse();
24         Lexer { tokens }
25     }
26
27     fn next(&mut self) -> Token {
28         self.tokens.pop().unwrap_or(Token::Eof)
29     }
30     fn peek(&mut self) -> Token {
31         self.tokens.last().copied().unwrap_or(Token::Eof)
32     }
33 }

```

To make sure that we got the ~~precedence~~ binding power correctly, we will be transforming infix expressions into a gold-standard (not so popular in Poland, for whatever reason) unambiguous notation — S-expressions:

$1 + 2 * 3 == (+ 1 (* 2 3)).$

```

1 use std::fmt;
2
3 enum S {
4     Atom(char),
5     Cons(char, Vec<S>),
6 }
7
8 impl fmt::Display for S {
9     fn fmt(&self, f: &mut fmt::Formatter<'_>) -> fmt::Result {
10         match self {
11             S::Atom(i) => write!(f, "{}", i),
12             S::Cons(head, rest) => {
13                 write!(f, "({}", head)?;
14                 for s in rest {
15                     write!(f, " {}", s)?
16                 }
17                 write!(f, ")")
18             }
19         }
20     }
21 }

```

And let's start with just this: expressions with atoms and two infix binary operators, + and *:

```

1 fn expr(input: &str) -> S {
2     let mut lexer = Lexer::new(input);
3     expr_bp(&mut lexer)
4 }
5
6 fn expr_bp(lexer: &mut Lexer) -> S {
7     todo!()
8 }
9
10 #[test]
11 fn tests() {
12     let s = expr("1 + 2 * 3");
13     assert_eq!(s.to_string(), "(+ 1 (* 2 3))")
14 }

```

So, the general approach is roughly the one we used to deal with left recursion — start with parsing a first number, and then loop, consuming operators and doing ... something?

```

1 fn expr_bp(lexer: &mut Lexer) -> S {
2     let lhs = match lexer.next() {
3         Token::Atom(it) => S::Atom(it),
4         t => panic!("bad token: {:?}", t),
5     };
6
7     loop {
8         let op = match lexer.next() {
9             Token::Eof => break,
10            Token::Op(op) => op,
11            t => panic!("bad token: {:?}", t),
12        };
13
14        todo!()
15    }
16
17    lhs
18 }
19
20 #[test]
21 fn tests() {
22     let s = expr("1"); ❶
23     assert_eq!(s.to_string(), "1");
24 }

```

❶ Note that we already can parse this simple test!

We want to use this power idea, so let's compute both left and right powers of the operator. We'll use `u8` to represent power, so, for associativity, we'll add 1. And we'll reserve the `0` power for the end of input, so the lowest power operator can have is 1.

```

1 fn expr_bp(lexer: &mut Lexer) -> S {
2     let lhs = match lexer.next() {
3         Token::Atom(it) => S::Atom(it),
4         t => panic!("bad token: {:?}", t),
5     };
6
7     loop {
8         let op = match lexer.peek() {
9             Token::Eof => break,
10            Token::Op(op) => op,
11            t => panic!("bad token: {:?}", t),
12        };
13        let (l_bp, r_bp) = infix_binding_power(op);
14
15        todo!()
16    }
17
18    lhs
19 }
20
21 fn infix_binding_power(op: char) -> (u8, u8) {
22     match op {
23         '+' | '-' => (1, 2),
24         '*' | '/' => (3, 4),
25         _ => panic!("bad op: {:?}", op)
26     }
27 }

```

And now comes the tricky bit, where we introduce recursion into the picture. Let's think about this example (with powers below):

```

1 | a   +   b   *   c   *   d   +   e
2 | 1   2   3   4   3   4   1   2

```

The cursor is at the first `+`, we know that the left `bp` is 1 and the right one is 2. The `lhs` stores `a`. The next operator after `+` is `*`, so we shouldn't add `b` to `a`. The problem is that we haven't yet seen the next operator, we are just past `+`. Can we add a lookahead? Looks like no — we'd have to look past all of `b`, `c` and `d` to find the next operator with lower binding power, which sounds pretty unbounded. But we are onto something! Our current right priority is 2, and, to be able to fold the expression, we need to find the next operator with lower priority. So let's recursively call `expr_bp` starting at `b`, but also tell it to stop as soon as `bp` drops below 2. This necessitates the addition of `min_bp` argument to the main function.

And lo, we have a fully functioning minimal Pratt parser:


```

1 fn expr(input: &str) -> S {
2     let mut lexer = Lexer::new(input);
3     expr_bp(&mut lexer, 0) ⑤
4 }
5
6 fn expr_bp(lexer: &mut Lexer, min_bp: u8) -> S { ①
7     let mut lhs = match lexer.next() {
8         Token::Atom(it) => S::Atom(it),
9         t => panic!("bad token: {:?}", t),
10    };
11
12    loop {
13        let op = match lexer.peek() {
14            Token::Eof => break,
15            Token::Op(op) => op,
16            t => panic!("bad token: {:?}", t),
17        };
18
19        let (l_bp, r_bp) = infix_binding_power(op);
20        if l_bp < min_bp { ②
21            break;
22        }
23
24        lexer.next(); ③
25        let rhs = expr_bp(lexer, r_bp);
26
27        lhs = S::Cons(op, vec![lhs, rhs]); ④
28    }
29
30    lhs
31 }
32
33 fn infix_binding_power(op: char) -> (u8, u8) {
34     match op {
35         '+' | '-' => (1, 2),
36         '*' | '/' => (3, 4),
37         _ => panic!("bad op: {:?}", op),
38     }
39 }
40
41 #[test]
42 fn tests() {
43     let s = expr("1");
44     assert_eq!(s.to_string(), "1");
45
46     let s = expr("1 + 2 * 3");
47     assert_eq!(s.to_string(), "(+ 1 (* 2 3))");
48
49     let s = expr("a + b * c * d + e");
50     assert_eq!(s.to_string(), "(+ (+ a (* (* b c) d)) e)");
51 }

```

① `min_bp` argument is the crucial addition. `expr_bp` now parses expressions with relatively high binding power. As soon as it sees something weaker than `min_bp`, it stops.

② This is the “it stops” point.

③ And here we bump past the operator itself and make the recursive call. Note how we use `l_bp` to check against `min_bp`, and `r_bp` as the new `min_bp` of the recursive call. So, you can think about `min_bp` as the binding power of the operator to the left of the current expressions.

- ④ Finally, after parsing the correct right hand side, we assemble the new current expression.
- ⑤ To start the recursion, we use binding power of zero. Remember, at the beginning the binding power of the operator to the left is the lowest possible, zero, as there's no actual operator there.

So, yup, these 40 lines *are* the Pratt parsing algorithm. They are tricky, but, if you understand them, everything else is straightforward additions.

Bells and Whistles

Now let's add all kinds of weird expressions to show the power and flexibility of the algorithm. First, let's add a high-priority, right associative function composition operator: `::`

```
1 fn infix_binding_power(op: char) -> (u8, u8) {
2     match op {
3         '+' | '-' => (1, 2),
4         '*' | '/' => (3, 4),
5         '.' => (6, 5),
6         _ => panic!("bad op: {:?}", op),
7     }
8 }
```

Yup, it's a single line! Note how the left side of the operator binds tighter, which gives us desired right associativity:

```
1 let s = expr("f . g . h");
2 assert_eq!(s.to_string(), "(. f (. g h))");
3
4 let s = expr("1 + 2 + f . g . h * 3 * 4");
5 assert_eq!(s.to_string(), "(+ (+ 1 2) (* (* (. f (. g h)) 3) 4))");
```

Now, let's add unary `-`, which binds tighter than binary arithmetic operators, but less tight than composition. This requires changes to how we start our loop, as we no longer can assume that the first token is an atom, and need to handle minus as well. But let the types drive us. First, we start with binding powers. As this is a unary operator, it really only have right binding power, so, ahem, let's just code this:

```
1 fn prefix_binding_power(op: char) -> ((), u8) { ❶
2     match op {
3         '+' | '-' => ((), 5),
4         _ => panic!("bad op: {:?}", op),
5     }
6 }
7
8 fn infix_binding_power(op: char) -> (u8, u8) {
9     match op {
10        '+' | '-' => (1, 2),
11        '*' | '/' => (3, 4),
12        '.' => (8, 7), ❷
13        _ => panic!("bad op: {:?}", op),
14    }
15 }
```

- ❶ Here, we return a dummy `()` to make it clear that this is a prefix, and not a postfix operator, and thus can only bind things to the right.

- ② Note, as we want to add unary `-` between `.` and `*`, we need to shift priorities of `.` by two. The general rule is that we use an odd priority as base, and bump it by one for associativity, if the operator is binary. For unary minus it doesn't matter and we could have used either 5 or 6, but sticking to odd is more consistent.

Plugging this into `expr_bp`, we get:

```
1 fn expr_bp(lexer: &mut Lexer, min_bp: u8) -> S {
2     let mut lhs = match lexer.next() {
3         Token::Atom(it) => S::Atom(it),
4         Token::Op(op) => {
5             let (_, r_bp) = prefix_binding_power(op);
6             todo!()
7         }
8         t => panic!("bad token: {:?}", t),
9     };
10    ...
11 }
```

Now, we only have `r_bp` and not `l_bp`, so let's just copy-paste half of the code from the main loop? Remember, we use `r_bp` for recursive calls.

```

1  fn expr_bp(lexer: &mut Lexer, min_bp: u8) -> S {
2      let mut lhs = match lexer.next() {
3          Token::Atom(it) => S::Atom(it),
4          Token::Op(op) => {
5              let ((), r_bp) = prefix_binding_power(op);
6              let rhs = expr_bp(lexer, r_bp);
7              S::Cons(op, vec![rhs])
8          }
9      };
10     t => panic!("bad token: {:?}", t),
11 };
12
13 loop {
14     let op = match lexer.peek() {
15         Token::Eof => break,
16         Token::Op(op) => op,
17         t => panic!("bad token: {:?}", t),
18     };
19
20     let (l_bp, r_bp) = infix_binding_power(op);
21     if l_bp < min_bp {
22         break;
23     }
24
25     lexer.next();
26     let rhs = expr_bp(lexer, r_bp);
27
28     lhs = S::Cons(op, vec![lhs, rhs]);
29 }
30
31 lhs
32 }
33
34 #[test]
35 fn tests() {
36     ...
37
38     let s = expr("--1 * 2");
39     assert_eq!(s.to_string(), "(* (- (- 1)) 2)");
40
41     let s = expr("--f . g");
42     assert_eq!(s.to_string(), "(- (- (. f g)))");
43 }

```

Amusingly, this purely mechanical, type-driven transformation works. You can also reason why it works, of course. The same argument applies; after we've consumed a prefix operator, the operand consists of operators that bind tighter, and we just so conveniently happen to have a function which can parse expressions tighter than the specified power.

Ok, this is getting stupid. If using `()`, `u8` “just worked” for prefix operators, can `u8`, `()` deal with postfix ones? Well, let's add `!` for factorials. It should bind tighter than `-`, because `-(92!)` is obviously more useful than `(-92)!`. So, the familiar drill — new priority function, shifting priority of `.` (this bit *is* annoying in Pratt parsers), copy-pasting the code...

```
1 let (l_bp, ()) = postfix_binding_power(op);  
2 if l_bp < min_bp {  
3     break;  
4 }  
5  
6 let (l_bp, r_bp) = infix_binding_power(op);  
7 if l_bp < min_bp {  
8     break;  
9 }
```

Wait, something's wrong here. After we've parsed the prefix expression, we can see either a postfix or an infix operator. But we bail on unrecognized operators, which is not going to work... So, let's make `postfix_binding_power` to return an option, for the case where the operator is not postfix:

```
1 fn expr_bp(lexer: &mut Lexer, min_bp: u8) -> S {
2     let mut lhs = match lexer.next() {
3         Token::Atom(it) => S::Atom(it),
4         Token::Op(op) => {
5             let ((), r_bp) = prefix_binding_power(op);
6             let rhs = expr_bp(lexer, r_bp);
7             S::Cons(op, vec![rhs])
8         }
9     };
10    t => panic!("bad token: {:?}", t),
11
12    loop {
13        let op = match lexer.peek() {
14            Token::Eof => break,
15            Token::Op(op) => op,
16            t => panic!("bad token: {:?}", t),
17        };
18
19        if let Some((l_bp, ())) = postfix_binding_power(op) {
20            if l_bp < min_bp {
21                break;
22            }
23            lexer.next();
24
25            lhs = S::Cons(op, vec![lhs]);
26            continue;
27        }
28
29        let (l_bp, r_bp) = infix_binding_power(op);
30        if l_bp < min_bp {
31            break;
32        }
33
34        lexer.next();
35        let rhs = expr_bp(lexer, r_bp);
36
37        lhs = S::Cons(op, vec![lhs, rhs]);
38    }
39
40    lhs
41 }
42
43 fn prefix_binding_power(op: char) -> ((), u8) {
44     match op {
45         '+' | '-' => ((), 5),
46         _ => panic!("bad op: {:?}", op),
47     }
48 }
49
50 fn postfix_binding_power(op: char) -> Option<(u8, ())> {
51     let res = match op {
52         '!' => (7, ()),
53         _ => return None,
54     };
55     Some(res)
56 }
57
58 fn infix_binding_power(op: char) -> (u8, u8) {
59     match op {
60         '+' | '-' => (1, 2),
61         '*' | '/' => (3, 4),
62         '.' => (10, 9),
63         _ => panic!("bad op: {:?}", op),
64     }
65 }
66 }
```

```

67 #[test]
68 fn tests() {
69     let s = expr("-9!");
70     assert_eq!(s.to_string(), "(- (! 9))");
71
72     let s = expr("f . g !");
73     assert_eq!(s.to_string(), "(! (. f g))");
74 }

```

Amusingly, both the old and the new tests pass.

Now, we are ready to add a new kind of expression: parenthesised expression. It is actually not that hard, and we could have done it from the start, but it makes sense to handle this here, you'll see in a moment why. Parens are just a primary expressions, and are handled similar to atoms:

```

1 let mut lhs = match lexer.next() {
2     Token::Atom(it) => S::Atom(it),
3     Token::Op('(') => {
4         let lhs = expr_bp(lexer, 0);
5         assert_eq!(lexer.next(), Token::Op(')'));
6         lhs
7     }
8     Token::Op(op) => {
9         let (_, r_bp) = prefix_binding_power(op);
10        let rhs = expr_bp(lexer, r_bp);
11        S::Cons(op, vec![rhs])
12    }
13    t => panic!("bad token: {:?}", t),
14 };

```

Unfortunately, the following test fails:

```

1 let s = expr("(((0)))");
2 assert_eq!(s.to_string(), "0");

```

The panic comes from the loop below — the only termination condition we have is reaching eof, and `)` is definitely not eof. The easiest way to fix that is to change `infix_binding_power` to return `None` on unrecognized operands. That way, it'll become similar to `postfix_binding_power` again!

```
1 fn expr_bp(lexer: &mut Lexer, min_bp: u8) -> S {
2     let mut lhs = match lexer.next() {
3         Token::Atom(it) => S::Atom(it),
4         Token::Op('(') => {
5             let lhs = expr_bp(lexer, 0);
6             assert_eq!(lexer.next(), Token::Op(')'));
7             lhs
8         }
9         Token::Op(op) => {
10             let ((, r_bp)) = prefix_binding_power(op);
11             let rhs = expr_bp(lexer, r_bp);
12             S::Cons(op, vec![rhs])
13         }
14         t => panic!("bad token: {:?}", t),
15     };
16
17     loop {
18         let op = match lexer.peek() {
19             Token::Eof => break,
20             Token::Op(op) => op,
21             t => panic!("bad token: {:?}", t),
22         };
23
24         if let Some((l_bp, ())) = postfix_binding_power(op) {
25             if l_bp < min_bp {
26                 break;
27             }
28             lexer.next();
29
30             lhs = S::Cons(op, vec![lhs]);
31             continue;
32         }
33
34         if let Some((l_bp, r_bp)) = infix_binding_power(op) {
35             if l_bp < min_bp {
36                 break;
37             }
38
39             lexer.next();
40             let rhs = expr_bp(lexer, r_bp);
41
42             lhs = S::Cons(op, vec![lhs, rhs]);
43             continue;
44         }
45
46         break;
47     }
48
49     lhs
50 }
51
52 fn prefix_binding_power(op: char) -> (((), u8)) {
53     match op {
54         '+' | '-' => (((), 5),
55         _ => panic!("bad op: {:?}", op),
56     }
57 }
58
59 fn postfix_binding_power(op: char) -> Option<(u8, ())> {
60     let res = match op {
61         '!' => (7, ()),
62         _ => return None,
63     };
64     Some(res)
65 }
66
```



```
67 fn infix_binding_power(op: char) -> Option<(u8, u8)> {  
68     let res = match op {  
69         '+' | '-' => (1, 2),  
70         '*' | '/' => (3, 4),  
71         '.' => (10, 9),  
72         _ => return None,  
73     };  
74     Some(res)  
75 }
```

And now let's add array indexing operator: `a[i]`. What kind of -fix is it? Around-fix? If it were just `a[]`, it would clearly be postfix. if it were just `[i]`, it would work exactly like parens. And it is the key: the `i` part doesn't really participate in the whole power game, as it is unambiguously delimited. So, let's do this:

```

1  fn expr_bp(lexer: &mut Lexer, min_bp: u8) -> S {
2      let mut lhs = match lexer.next() {
3          Token::Atom(it) => S::Atom(it),
4          Token::Op('(') => {
5              let lhs = expr_bp(lexer, 0);
6              assert_eq!(lexer.next(), Token::Op(')'));
7              lhs
8          }
9          Token::Op(op) => {
10             let ((), r_bp) = prefix_binding_power(op);
11             let rhs = expr_bp(lexer, r_bp);
12             S::Cons(op, vec![rhs])
13         }
14         t => panic!("bad token: {:?}", t),
15     };
16
17     loop {
18         let op = match lexer.peek() {
19             Token::Eof => break,
20             Token::Op(op) => op,
21             t => panic!("bad token: {:?}", t),
22         };
23
24         if let Some((l_bp, ())) = postfix_binding_power(op) {
25             if l_bp < min_bp {
26                 break;
27             }
28             lexer.next();
29
30             lhs = if op == '[' {
31                 let rhs = expr_bp(lexer, 0);
32                 assert_eq!(lexer.next(), Token::Op(']'));
33                 S::Cons(op, vec![lhs, rhs])
34             } else {
35                 S::Cons(op, vec![lhs])
36             };
37             continue;
38         }
39
40         if let Some((l_bp, r_bp)) = infix_binding_power(op) {
41             if l_bp < min_bp {
42                 break;
43             }
44
45             lexer.next();
46             let rhs = expr_bp(lexer, r_bp);
47
48             lhs = S::Cons(op, vec![lhs, rhs]);
49             continue;
50         }
51
52         break;
53     }
54
55     lhs
56 }
57
58 fn prefix_binding_power(op: char) -> ((), u8) {
59     match op {
60         '+' | '-' => ((), 5),
61         _ => panic!("bad op: {:?}", op),
62     }
63 }
64
65 fn postfix_binding_power(op: char) -> Option<(u8, ())> {
66     let res = match op {

```

```

67         '!' | '[' => (7, ()), ❶
68         _ => return None,
69     };
70     Some(res)
71 }
72
73 fn infix_binding_power(op: char) -> Option<(u8, u8)> {
74     let res = match op {
75         '+' | '-' => (1, 2),
76         '*' | '/' => (3, 4),
77         '.' => (10, 9),
78         _ => return None,
79     };
80     Some(res)
81 }
82
83 #[test]
84 fn tests() {
85     ...
86
87     let s = expr("x[0][1]");
88     assert_eq!(s.to_string(), "([ ([ x 0) 1])");
89 }

```

❶ Note that we use the same priority for ! as for [. In general, for the correctness of our algorithm it's pretty important that, when we make decisions, priorities are never equal. Otherwise, we might end up in a situation like the one before tiny adjustment for associativity, where there were two equally-good candidates for reduction. However, we only compare right bp with left bp! So for two postfix operators it's OK to have priorities the same, as they are both right.

Finally, the ultimate boss of all operators, the dreaded ternary:

```
1 | c ? e1 : e2
```

Is this ... all-other-the-place-fix operator? Well, let's change the syntax of ternary slightly:

```
1 | c [ e1 ] e2
```

And let's recall that `a[i]` turned out to be a postfix operator + parenthesis... So, yeah, `?` and `:` are actually a weird pair of parens! And let's handle it as such! Now, what about priority and associativity? What associativity even is in this case?

```
1 | a ? b : c ? d : e
```

To figure it out, we just squash the parens part:

```
1 | a ?: c ?: e
```

This can be parsed as

```
1 | (a ?: c) ?: e
```

or as

```
1 | a ?: (c ?: e)
```

What is more useful? For `?`-chains like this:

```

1 | a ? b :
2 | c ? d :
3 | e

```

the right-associative reading is more useful. Priority-wise, the ternary is low priority. In C, only `=` and `,` have lower priority. While we are at it, let's add C-style right associative `=` as well.

Here's our the most complete and perfect version of a simple Pratt parser:

```
1 use std::{fmt, io::BufRead};
2
3 enum S {
4     Atom(char),
5     Cons(char, Vec<S>),
6 }
7
8 impl fmt::Display for S {
9     fn fmt(&self, f: &mut fmt::Formatter<'_>) -> fmt::Result {
10         match self {
11             S::Atom(i) => write!(f, "{}", i),
12             S::Cons(head, rest) => {
13                 write!(f, "({}", head)?;
14                 for s in rest {
15                     write!(f, " {}", s)?
16                 }
17                 write!(f, ")")
18             }
19         }
20     }
21 }
22
23 #[derive(Debug, Clone, Copy, PartialEq, Eq)]
24 enum Token {
25     Atom(char),
26     Op(char),
27     Eof,
28 }
29
30 struct Lexer {
31     tokens: Vec<Token>,
32 }
33
34 impl Lexer {
35     fn new(input: &str) -> Lexer {
36         let mut tokens = input
37             .chars()
38             .filter(|it| !it.is_ascii_whitespace())
39             .map(|c| match c {
40                 '0'..'9'
41                 | 'a'..'z' | 'A'..'Z' => Token::Atom(c),
42                 _ => Token::Op(c),
43             })
44             .collect::<Vec<_>>();
45         tokens.reverse();
46         Lexer { tokens }
47     }
48
49     fn next(&mut self) -> Token {
50         self.tokens.pop().unwrap_or(Token::Eof)
51     }
52
53     fn peek(&mut self) -> Token {
54         self.tokens.last().copied().unwrap_or(Token::Eof)
55     }
56 }
57
58 fn expr(input: &str) -> S {
59     let mut lexer = Lexer::new(input);
60     expr_bp(&mut lexer, 0)
61 }
62
63 fn expr_bp(lexer: &mut Lexer, min_bp: u8) -> S {
64     let mut lhs = match lexer.next() {
65         Token::Atom(it) => S::Atom(it),
66         Token::Op('(') => {
67             let lhs = expr_bp(lexer, 0);
```

```

67     assert_eq!(lexer.next(), Token::Op(' '));
68     lhs
69 }
70 Token::Op(op) => {
71     let (_, r_bp) = prefix_binding_power(op);
72     let rhs = expr_bp(lexer, r_bp);
73     S::Cons(op, vec![rhs])
74 }
75 t => panic!("bad token: {:?}", t),
76 };
77
78 loop {
79     let op = match lexer.peek() {
80         Token::Eof => break,
81         Token::Op(op) => op,
82         t => panic!("bad token: {:?}", t),
83     };
84
85     if let Some((l_bp, ())) = postfix_binding_power(op) {
86         if l_bp < min_bp {
87             break;
88         }
89         lexer.next();
90
91         lhs = if op == '[' {
92             let rhs = expr_bp(lexer, 0);
93             assert_eq!(lexer.next(), Token::Op(']'));
94             S::Cons(op, vec![lhs, rhs])
95         } else {
96             S::Cons(op, vec![lhs])
97         };
98         continue;
99     }
100
101     if let Some((l_bp, r_bp)) = infix_binding_power(op) {
102         if l_bp < min_bp {
103             break;
104         }
105         lexer.next();
106
107         lhs = if op == '?' {
108             let mhs = expr_bp(lexer, 0);
109             assert_eq!(lexer.next(), Token::Op(':'));
110             let rhs = expr_bp(lexer, r_bp);
111             S::Cons(op, vec![lhs, mhs, rhs])
112         } else {
113             let rhs = expr_bp(lexer, r_bp);
114             S::Cons(op, vec![lhs, rhs])
115         };
116         continue;
117     }
118
119     break;
120 }
121
122 lhs
123 }
124
125 fn prefix_binding_power(op: char) -> (((), u8)) {
126     match op {
127         '+' | '-' => (((), 9)),
128         _ => panic!("bad op: {:?}", op),
129     }
130 }
131
132 fn postfix_binding_power(op: char) -> Option<(((), u8))> {

```

```

132 fn postfix_binding_power(op: char) -> Option<(u8, ())> {
133     let res = match op {
134         '!' => (11, ()),
135         '[' => (11, ()),
136         _ => return None,
137     };
138     Some(res)
139 }
140
141 fn infix_binding_power(op: char) -> Option<(u8, u8)> {
142     let res = match op {
143         '=' => (2, 1),
144         '?' => (4, 3),
145         '+' | '-' => (5, 6),
146         '*' | '/' => (7, 8),
147         '.' => (14, 13),
148         _ => return None,
149     };
150     Some(res)
151 }
152
153 #[test]
154 fn tests() {
155     let s = expr("1");
156     assert_eq!(s.to_string(), "1");
157
158     let s = expr("1 + 2 * 3");
159     assert_eq!(s.to_string(), "(+ 1 (* 2 3))");
160
161     let s = expr("a + b * c * d + e");
162     assert_eq!(s.to_string(), "(+ (+ a (* (* b c) d)) e)");
163
164     let s = expr("f . g . h");
165     assert_eq!(s.to_string(), "(. f (. g h))");
166
167     let s = expr(" 1 + 2 + f . g . h * 3 * 4");
168     assert_eq!(
169         s.to_string(),
170         "(+ (+ 1 2) (* (* (. f (. g h)) 3) 4))",
171     );
172
173     let s = expr("--1 * 2");
174     assert_eq!(s.to_string(), "(* (- (- 1)) 2)");
175
176     let s = expr("--f . g");
177     assert_eq!(s.to_string(), "(- (- (. f g)))");
178
179     let s = expr("-9!");
180     assert_eq!(s.to_string(), "(- (! 9))");
181
182     let s = expr("f . g !");
183     assert_eq!(s.to_string(), "(! (. f g))");
184
185     let s = expr("(((0)))");
186     assert_eq!(s.to_string(), "0");
187
188     let s = expr("x[0][1]");
189     assert_eq!(s.to_string(), "[ ([ x 0) 1]");
190
191     let s = expr(
192         "a ? b :
193         c ? d
194         : e",
195     );
196     assert_eq!(s.to_string(), "(? a b (? c d e))");
197

```

```
198     let s = expr("a = 0 ? b : c = d");
199     assert_eq!(s.to_string(), "(= a (= (? 0 b c) d))")
200 }
201
202 fn main() {
203     for line in std::io::stdin().lock().lines() {
204         let line = line.unwrap();
205         let s = expr(&line);
206         println!("{}", s)
207     }
208 }
```

The code is also available in [this repository](#), Eof :-)

 Fix typo  Subscribe  Get in touch  matklad