



To make Medium work, we log user data.
By using Medium, you agree to our
Privacy Policy, including cookie policy.

[Read in app](#)[Get started](#)

Published in Analytics Vidhya · [Follow](#)

You have **2** free member-only stories left this month.

[Sign up for Medium and get an extra one](#)



Andrew Matteson · [Follow](#)

Feb 1, 2021 · 7 min read ★



Parsing PGN Chess Games With Python

I have a couple of projects in mind for analysis of chess game positions. They all require collecting a large number of chess games into structured data. [This Week in Chess](#) provides records of thousands of games each week. Each week consists of a file in [Portable Game Notation](#) (PGN), a plain text format for recording chess games.

I want to analyze hundreds of thousands of games, so step one is to gather the plain text files and parse them using the grammar of PGN files.

If you just want the code to parse PGN, you can jump over to [the GitHub repo](#). Otherwise, read on.





To make Medium work, we log user data.
By using Medium, you agree to our
Privacy Policy, including cookie policy.

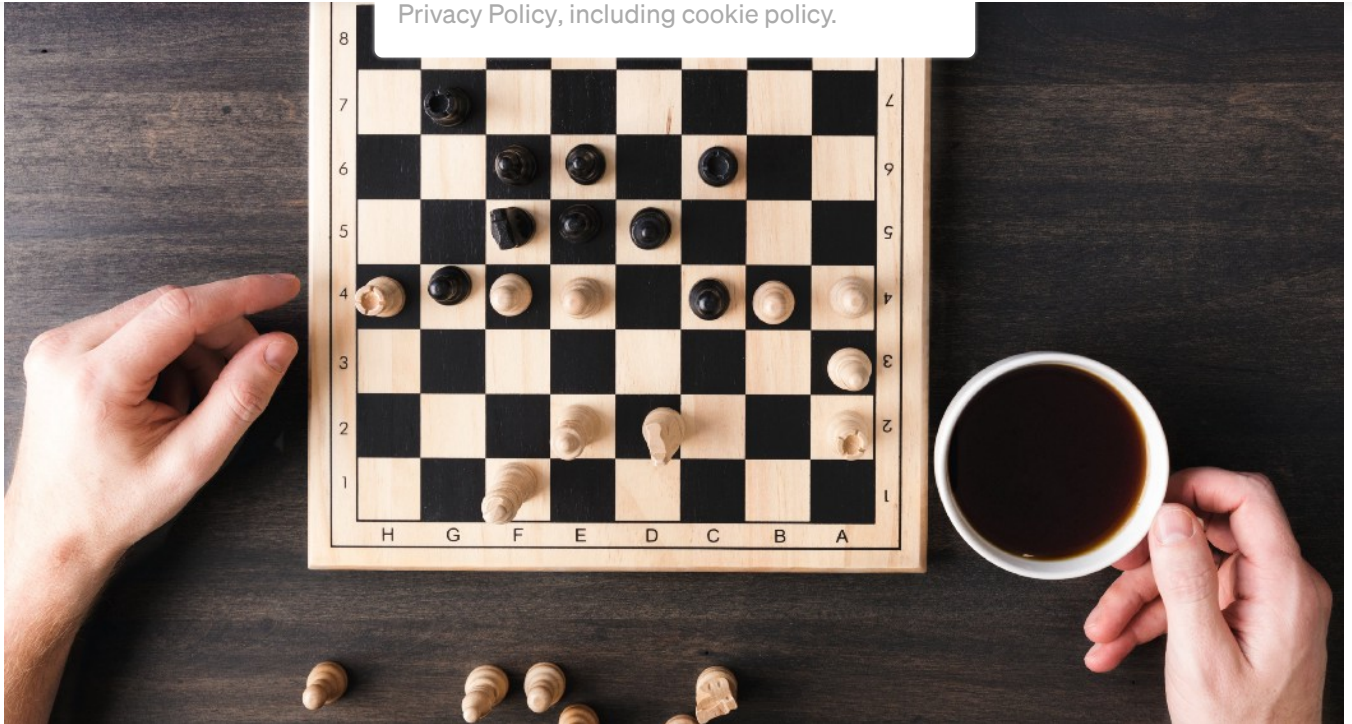
[in app](#)[Get started](#)

Photo by Sarah Pflug, used under CCO

The Grammar of PGN Files

Lots of text data, PGN included, obeys a “grammar”. The grammar of PGN includes annotations that are pairs of a tag and a string, the moves of the game, and the outcome. The important insight is that there are rules for how that text looks. For instance, this 2021 game between Magnus Carlsen and Jorden Van Foreest:

```
[Event "83rd Tata Steel Masters"]  
[Site "Wijk aan Zee NED"]  
[Date "2021.01.19"]  
[Round "4.1"]  
[White "Van Foreest, Jorden"]  
[Black "Carlsen, M"]  
[Result "1/2-1/2"]  
[WhiteTitle "GM"]  
[BlackTitle "GM"]  
[WhiteElo "2671"]  
[BlackElo "2862"]  
[ECO "C78"]  
[Opening "Ruy Lopez"]  
[Variation "Archangelsk (counterthrust) variation"]  
[WhiteFideId "1039784"]  
[BlackFideId "1503014"]  
[EventDate "2021.01.16"]
```





To make Medium work, we log user data.

By using Medium, you agree to our

Privacy Policy, including cookie policy.

in app

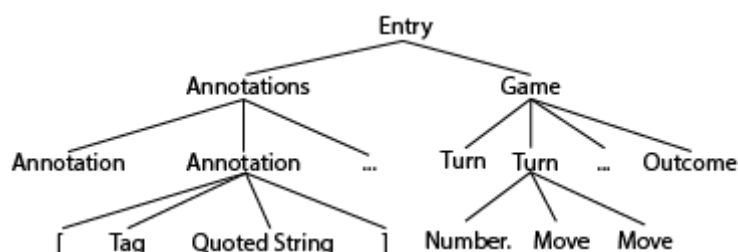
Get started

29. Kf1 Rxd3 30. Rxd3 31. Ng1 Rd7 32. Nf3 Kg7 33. h3 Bf6 34. Kg3 Rb7 35. Kg2 Re7 36. Ra5 Rc7 37. Rd5 Ra7 38. Rb5 Be7 39. Nd4 Rd7 40. Nf3 Rd6 41. Rb7 Kf6 42. Ra7 h6 43. Nh4 Bd8 44. Rh7 Rd2 45. Rxh6 Kg7 46. Rxd6 Kh7 47. Nf3 Rxf2+ 48. Kxf2 Kxg6 49. Kg2 Kh5 50. Nd4 f4 51. Ne6 Bg5 52. Nxg5 Kxg5 53. Kf3 Kh4 54. Kxf4 Kxh3 1/2-1/2

The annotations appear in brackets, starting with the tag followed by a quoted string. After the annotations, the moves are listed by move number followed by white's move, then black's move (in algebraic notation). The last text declares the outcome, 1-0 for white wins, 0-1 for black, and 1/2-1/2 for a draw.

We're going to take a file of PGN entries and parse them. In computer science, a parser-combinator is a higher order function that takes in multiple parsers and returns a single parser. The topic is pretty deep, but the gist is that parser-combinators allow us to write out a parse tree of a grammar and iteratively build complex parsers by combining simpler ones. For this project, I use parsita, an open source python library for parsing written by one of my colleagues.

Writing out the parse tree from the file above is straightforward. An "entry" in the file is composed of two blocks: an "annotations" block and a "game" block. The annotations block is a new line delimited list of the sequence [, TAG , " , STRING , " ,] . The game is a list of NUMBER. , MOVE , MOVE , followed by an OUTCOME . There's some trickiness about the last move because white might make one more move than black, but we'll deal with that later. It's best to visualize our grammar as a tree:



PGN Parse Tree (image by author)

Parsing Simple Text





To make Medium work, we log user data.
By using Medium, you agree to our
Privacy Policy, including cookie policy.

[in app](#)
[Get started](#)

We'll walk through each line of this code:

```
# Parse the " token
quote = lit(r'\"')

# Parse a sequence of unicode characters excluding spaces
tag = reg(r'[\u0021-\u0021\u0023-\u005A\u005E-\u007E]+')

# parse a sequence of unicode characters including spaces
string = reg(r'[\u0020-\u0021\u0023-\u005A\u005E-\u0010FFFF]+')

# An annotation is combines a tag, and a quoted string delimited
# by other characters. We only care about the tag and string.
annotation = '[' >> tag << ' ' & (quote >> string << quote) << ']'

# Function to take a list of parsed annotations and convert to a
# dictionary
def formatannotations(annotations):
    return {ant[0]: ant[1] for ant in annotations}

# The annotation block is a list of 0 or more annotations.
annotations = repsep(annotation, '\n') > formatannotations
```

Parsita does the heavy lifting with literal parsers and regular expression parsers. The simplest parsing operation is to match a single character, the literal. I show two ways of doing this: `lit(r'\"')` matches a quote; alternatively, parsita will accept a Python string into a parser and treat it as a literal. The brackets are parsed using the second technique.

The power of parsita comes from the five parsers we use in the above:

- `A & B` : sequential parser. Match the parser for `A` and the parser for `B`. Return both in a list.
- `A >> B` : discard left parser. Match the same as `&`. Return only `B`.
- `A << B` : discard right parser. Vice-versa from above.





To make Medium work, we log user data.

By using Medium, you agree to our

Privacy Policy, including cookie policy.

in app

Get started

- repsep (A) : repeated matches.

| return a list of

So, expanding on the line defining the annotation:

```
annotation = '[' >> tag << ' ' & (quote >> string << quote) << '']'
```

An annotation is in sequence:

- a bracket (discard),
- a string without spaces,
- a space (discard),
- a quote (discard),
- a string,
- a quote (discard),
- a bracket (discard).

```
# Call the parser --> results in Success or Failure.
# .or_die() either gives the value of the success or raises an error

print(annotation.parse('[parsing "is cool"]').or_die())
> ['parsing', 'is cool']

# Annotations block converts the text to a dict
text = '''[parsing "is cool"
[second "line"]'''

print(annotations.parse(text).or_die())
> {'parsing': 'is cool', 'second': 'line'}
```

Parsing More Complex Text





To make Medium work, we log user data.

By using Medium, you agree to our

Privacy Policy, including cookie policy.

in app

Get started

- $A|B$: the alternative

ls, then try parser B .

- $\text{opt}(A)$: the optional parser. Tries to match parser A , if it succeeds return the value of A in a length one list. If it fails, return an empty list.

We use the alternative parser to simplify the definition of a move:

```
nullmove = lit('--') # Illegal move rarely used in annotations
longcastle = reg(r'O-O-O[+#]?')
castle = reg(r'O-O[+#]?')
regularmove = reg(r'[a-h1-8NBRQKx\+#=]+') # Matches more than just
chess moves
move = regularmove | longcastle | castle | nullmove
```

We match “regular” moves with a regular expression parser. Algebraic notation of a move references 1 or 2 squares with a-h, their file, and 1–8, their rank. Pieces have one letter designations. There are additional notations for capture (\times), check ($+$), checkmate ($\#$), and promotion ($=$). Other, non-regular, moves are long castle, castle, and null move (rarely encountered, but appears for various reasons).

So, a move tries to match a regular move, if it fails, tries to match long castle, then castle, then finally the null move.

```
def handleoptional(optionalmove):
    if len(optionalmove) > 0:
        return optionalmove[0]
    else:
        return None
```

```
def formatgame(game):
    return {
        'moves': game[0],
        'outcome': game[1]
    }
```

```
whitespace = lit(' ') | lit('\n')
```

```
movenumber = (reg(r'[0-9]+') << '.' << whitespace) > int
```





To make Medium work, we log user data.

By using Medium, you agree to our

Privacy Policy, including cookie policy.

in app

Get started

```
draw = lit('1/2-1/2')
white = lit('1-0')
black = lit('0-1')
outcome = draw | white | black
```

```
game = (rep(turn) & outcome) > formatgame
```

We need the optional parser because a game may terminate after white moves. When this occurs, in PGN, the file has a final move number and then a single move listed. So, the last “turn” of our grammar is optional. We parse this into a more useful form for Python by indexing into the returned optional list if it’s not empty, and otherwise returning `None`.

Putting it All Together

All that remains is to download a PGN file, and parse it using the code outlined above.

```
from parsita import *
from parsita.util import constant
import json

# Conversion functions
def formatannotations(annotations):
    return {ant[0]: ant[1] for ant in annotations}

def formatgame(game):
    return {
        'moves': game[0],
        'outcome': game[1]
    }

def formatentry(entry):
    return {'annotations': entry[0], 'game': entry[1]}

def handleoptional(optionalmove):
    if len(optionalmove) > 0:
        return optionalmove[0]
    else:
        return None

# PGN Grammar
quote = lit(r'\"')
tag = reg(r'[\u0021-\u0021\u0023-\u005A\u005E-\u007E]+')
```





To make Medium work, we log user data.
By using Medium, you agree to our
Privacy Policy, including cookie policy.

in app

Get started

Now, I'm off to find fun things to do with thousands of chess games.

Sign up for Analytics Vidhya News Bytes

By Analytics Vidhya

Latest news from Analytics Vidhya on our Hackathons and some of our best articles! [Take a look.](#)

Your email

Get this newsletter

By signing up, you will create a Medium account if you don't already have one. Review our [Privacy Policy](#) for more information about our privacy practices.

