

Expressões regular em linguagem funcional: Módulo de busca usando expressões regulares implementado em Haskell

Bruno Gomes

16 de fevereiro de 2020

1 INTRODUCAO

1.1 Objetivos

Objetivo geral: Implementar uma biblioteca de processamento de expressoes regulares para a linguagem Haskell.

Objetivo especifico: - Entender o paradigma funcional e os componentes da linguagem Haskell - Estudar a teoria de automatas, variantes das expressoes regulares e algoritimos de processadores de expressoes regulares - Implentar biblioteca de busca de Strings em Haskell

1.2 Justificativa

Embora o objetivo deste trabalho seja a construçao de uma biblioteca para processamento de expressões regulares, a motivacao para isso é de realizar um estudo sobre o paradigma funcional. Dentre as inumeras linguagens de programcao ja inventadas, as que realmente se destacam e sao utilizadas em grandes aplicacoes sao poucas. Dentre elas, podemos destacar as languages Java, Python, C, C++, C#, Go, Shell Scripting e JavaScript. Embora a diferenca entre essas linguagens serem gritantes para pessoas com um minimo conhecimento sobre elas, todas podem ser classificas como linguagens imperativas. Por mais diferentes que essas linguagens sejam, elas compartilham um paradigma para a resolucao do problema computacional em questao. Linguagens imperativas focam em dizer *como* o problema devera ser resolvido.

Embora ser muito utilizada, a programacao imperativa nao é o unico paradigma computacional existente. Um dos paradigmas que contrasta diretamente é o paradigma de programacao funcional. Em poucas palavras, a grande diferenca entre esses paradigmas é que enquanto a programacao imperativa foca no como resolver o problema, a programacao funcional foca no *que* o programa faz. O ponto de constate sao as palavras como e que, embora seja uma definicao simplista, consegue capturar a essencia desses paradigmas.

Embora o paradigma funcional seja muito menos dissiminado, ele eh de extrema importancia. Grandes descobertas nesse campo de estudo transcende o paradigma

funcional e 'infecta' as linguagens procedurais. Alguns exemplos são as list comprehensions da linguagem Python e as lambda expressions e streams introduzidos na versão 8 da linguagem Java, um dos gigantes do mercado. Sendo assim, esse trabalho tem como objetivo introduzir o paradigma funcional, focado na linguagem Haskell. Serão introduzidos os principais conceitos e ferramentas das linguagens funcionais com um foco empírico. Esses conceitos serão utilizados para implementar uma engine de processamento de texto usando expressões regulares.

Em conclusão, o paradigma funcional faz parte do dia a dia da grande maioria dos profissionais da área, embora muitas vezes não saibam disso. Elementos de sucesso de linguagens funcionais acabam sendo implementados nas grandes linguagens imperativas pois os mesmos se mostraram úteis e atraentes. Sendo assim, desenvolvedores devem conhecer um mínimo sobre o paradigma funcional, não só pois ele está presente em grande parte das linguagens mas porque alguns problemas podem ser resolvidos de maneira mais breve.

2 DESENVOLVIMENTO

2.1 FUNDAMENTAÇÃO TEÓRICA

Como visto nos objetivos, esse trabalho pode ser quebrado em 3 partes: expressões regulares (regex), programação funcional e algoritmos para processamento de regex.

2.1.1 Expressões Regulares

Uma expressão regular (regex) é um string de caracteres que define um padrão de busca. Essas strings são construídas a fim de encontrar/extrair informações em texto que seguem uma estrutura pré conhecida com elementos dinâmicos.

Um exemplo, é possível escrever uma regex para encontrar em um texto todas as datas escritas no formato usual de DD/MM/AAAA, onde DD indica o dia, MM o mês e AAAA o ano. Uma regex capaz de encontrar esse padrão é "[0-9]2/[0-9]2/[0-9]4". Essa regex indica que para um texto ser aceito ele deve contar os dígitos de 0 a 9, repetidos duas vezes (indicado pelo substring [0-9]2), seguido por uma barra, outro grupo de dois dígitos, outra barra e finalmente um grupo de 4 dígitos.

A notação de regexes pode variar dependendo da implementação, porém seguindo o modelo Perl de regexes, a notação "[]" indica um grupo de caracteres que são aceitos para aquela posição de texto, [012] indica o conjunto de números 1,2,3 que deve ser aceito como correto. Em algumas situações, é possível indicar um intervalo de caracteres, como foi feito acima. O grupo "[0-9]" aceita os números 0,1,2,3,4,6,7,8,9 como corretos para fins de procura, é possível também definir um conjunto de letras. O conjunto "[a-z]" aceita as letras de "a" até "z", minúsculas como corretas para o padrão. O caractere "." é definido como um caractere reservado em diversas implementações de regexes diferentes, "." simboliza um grupo onde qualquer caractere é aceito.

É possível enumerar o número de vezes que um caractere, ou grupo de caracteres, podem aparecer em um regex, isso é indicado pela notação "". Ainda usando o exemplo de datas, o substring "2" indica que o caractere que antecede o enumerador será repetido duas vezes. Note que para fins de processamento, um grupo é consi-

derado como um unico character pois somente um unico character ira dar match por vez, logo a regex "[0-9]2" ira dar match em qualquer dois digitos adjacentes. Enumeradores podem indicar um intervalo com comprimento unidade ou um intervalo arbitrario, como o enumerador 2,5, que define um enumerador que aceita de 2 a 5 characters. Eh possivel omitir o numero a direita ou esquerda da virgula para definir um intervalo flexivel, o intervalo 0, aceita 0 ou qualquer numero de ocorrencias do grupo que o precede. Alguns caracteres sao reservados para indicar intervalos especiais em regexes. O character "?" eh equivalente ao intervalo "0,1", ou seja, o grupo anterior pode ocorrer 0 ou 1 vezes, usado para indicar quando uma condicao eh opcional. Outro character especial eh "?" que indica o intervalo "1,", ou seja um ou mais characters. Finalmente, o character "*" eh equivalente ao intervalo "0," que aceita qualquer numero de ocorrencias do grupo a sua esquerda.

Essa eh uma breve introducao sobre expressoes regulares, existem diversos characters especiais que podem ser usados para indicar posicoes no texto. Para uma referencia completa sobre as funcionalidades avancadas de expressoes regulares consulte ([python-re](#)), ([jeffrey](#)).

2.1.2 Programacao Funcional

Os conceitos de programacao funcional surgiram junto do calculo lambda, inventado por Alonzo Church. Embora existam varias linguagens funcionais (ML, Lisp, Rackett, Haskell), todas elas compartilham essa origem. Tal como em programacao imperativa que diferencia linguagens entre procedural e orientada a objetos, existem diferencas entre as linguagens funcionais, mas primeiramente, o que eh programacao funcional. Segundo (**Bird**), "programacao funcional eh: um metodo para construcao de programas que enfatiza funcoes e suas aplicacoes ao invaz de comandos e suas execucoes; programacao funcional faz uso de notacao matematica simples que permite que problemas sejam descritos de maneira clara e concisa. [...]".

Esse trabalho foca na linguagem Haskell como fonte de exemplos para o paradigma funcional. Haskell, como toda linguagem, possui suas proprias particularidades, porem grande partes dos conceitos apresentados sera geral para linguagens funcionais.

Haskell faz parte do conjunto de linguagens conhecidas como "linguagens funcionais puras". Uma linguagem pura permite a *definicao* de simbolos uma unica vez, exemplo, ao definir "let a = 4" em um programa, tentar mudar o valor da variavel "a" resulta em um error. O valor de um simbolo nao pode mudar durante a execucao do programa, de certa forma isso eh equivalente as keywords "final" do java e "const" da linguagem C. Embora nao poder atualizar o valor de uma variavel soe como uma grande desvantagem, isso permite o que pode ser chamado de "funcoes sem efeito-colateral". Isso significa que uma funcao jamais ira alterar o estado do programa fora dela, toda vez que uma funcao for chamada com os mesmos atributos, ela ira retornar o mesmo valor, isso nao eh o caso em linguagens orientadas a objeto, por exemplo, onde metodos alteram o estado do objeto (setters). O conceito de funcoes sem efeito colateral eh uma vantagem pois isso facilita o entendimento do programador sobre a sequencia de eventos do programa, nao existe preocupacao de que uma funcao altere alguma variavel global ou o estado de um objeto. Lipovaca expressa o conceito de linguagem funcional pura de maneira instrutiva, ele diz que: "Embora isso parece limitante quando voce esta vindo do mundo imperativo, nos vimos que isso eh algo realmente legal. Em uma linguagem imperativa voce nao tem

garantia de que uma simples funcao que deveria somente processar alguns numeros nao ira queimar sua casa, sequestrar seu cachorro e riscar seu carro com uma batata enquanto processa aqueles numeros.”.

Como nao podemos alterar o estado de uma variavel em um programa funcional, isso nos forca a procurar diferentes maneiras de resolver problemas computacionais. Um exemplo classico disso eh o problema de calcular o fatorial de um numero. O fatorial de n ($n!$) eh definido como $n! = (1) * (2) * \dots * (n - 1) * (n)$. Em uma linguagem imperativa, esse problema pode ser resolvido com um `for`, uma variavel acumuladora e uma variavel contadora; como nao podemos mutar o valor de definicoes em linguagens funcionais, podemos fazer uso de recursao.

O programa abaixo define uma funcao ”factorial” que recebe um valor do tipo `int` e retorna um valor `int`. A funcao ”factorial” retorna o valor 1 quando seu argumento eh 0 e retorna o valor $n * (n-1)!$ para qualquer outro valor. Essa curiosa sintaxe faz uso de um construtor muito util do Haskell chamado de *pattern matching*. Segundo ??, ”pattern matching consiste de especificar padroes para o qual algum dado deve tomar forma, verificar se o dado acorda com esses padroes”, eh possivel definir varios padroes para uma unica funcao. Padroes sao executados de cima para baixo e por isso padroes mais gerais (dao `match` em um maior numero de situacoes) devem ser posicionados por ultimo.

```
factorial :: Int -> Int
factorial 1 = 1
factorial n = n * factorial (n-1)
```

Tipos de dados diferentes possuem diferentes possibilidades de pattern matching. Em Haskell, valores contidos dentro de colchetes ”[]” define uma lista, e o padrao $(x:xs)$ define que x eh o primeiro elemento da lista e xs o restante da lista. Por exemplo, dado a lista `[1,2,3]` os valores de $(x:xs)$ em um pattern matching tera o valor de $x = 1$ e $xs = [2,3]$. Usando esse construtor, podemos definir a funcao `sum` que soma todos os numeros em uma lista.

```
sum :: [Int] -> Int
sum [] = 0
sum (x:xs) = x + sum xs
```

O codigo acima define uma funcao `sum` que recebe uma lista de inteiros (”[Int]”) e retorna um `int`, a funcao `sum` faz uso de dois pattern matchings. O primeiro padrao ”[]” faz `match` com uma lista vazia, o valor da funcao `sum` para a lista vazia eh definido como o valor 0. Caso a lista nao seja vazia, o segundo padrao eh executado, e para aquele caso a soma de uma lista eh definida como o primeiro elemento mais a soma do restante da lista. Novamente, isso eh uma funcao que faz uso de recursao para resolver o problema computacional. A funcao `sum` chama a ela mesma ate o valor de xs ser a lista vazia, para o qual definimos que isso seria igual a 0. Caso omitissimos o primeiro padrao, a funcao iria chamar a ela mesma indefinidamente. Pattern matching eh muito util em funcoes recursivas pois nos permite definir a condicao de saida da recursao.

2.1.3 Algoritmos de regex

Nessa secao sera abordado um pouco sobre maquinas de estado, como elas podem ser usadas para representar expressoes regulares. Talvez abordar a diferenca entre maquinas deterministicas e nao deterministicas, e como uma maquina nao deterministica pode ser transformada em deterministica.

Ainda estou em duvida qual algoritmo sera utilizado para a construcao da biblioteca, irei pesquisar mais sobre isso.

2.2 METODOLOGIA

Utilizando os algoritmos abordados na secao de algoritmos, sera criado um modulo escrito funcionalmente em Haskell para realizar a pesquisa de regexes em um texto.

Tenho duvida a respeito dessa secao. O que faz parte do escopo de metodologia? Praticas de desenvolvimento aplicadas como testes unitarios?

Acredito que a maior parte do conteudo abordado aqui seja estabelecido apos a escolha do algoritmo(s) de regex a serem implementados.

A definicao do dialeto da regex, junto dos simbolos especiais, flags e etc deve ser feita nessa secao ou na fundamentacao teorica?

3 CRONOGRAMA

Seq.	Fases	Atividades	Projeto de pesquisa
1	1a Fase do projeto de pesquisa	Definicao do tema	Dez/19, Jan/20
1	1a Fase do projeto de pesquisa	Estudo Haskell	Dez/19 - Abr/20
2	1a Fase do projeto de pesquisa	Pesquisa Regex e algoritmos	Fev/20 - Abr/20
3	1a Fase do projeto de pesquisa	Escrever projeto de pesquisa	Jan/20 - Mar/20
4	2a Fase do projeto de pesquisa	Planejar Algoritmo	Abr/20
5	2a Fase do projeto de pesquisa	Implementar modulo	Abr/20 - Jun/20
6	2a Fase do projeto de pesquisa	Escrever Relatorio	Jun/20 - Jul/20
7	2a Fase do projeto de pesquisa	Revisao	Jun/20 - Jul/20

4 CONSIDERACOES FINAIS

O objetivo desse trabalho foi expor alguns conceitos por traz do paradigma funcional e demonstrar como esses conceitos sao uteis atraves da construcao de um modulo de processamento de expressoes regulares. O modulo foi implementado usando a linguagem Haskell, uma linguagem funcional pura.

Primeiramente foi abordado as origens da programacao funcional, tal como as principais diferencas entre o paradigma funcional e imperativo. Ainda, foram introduzidos alguns conceitos exclusivos das linguagens funcionais, tais como Monads. Foi abordado de maneira abrangente o conceito por traz das expressoes regulares (regexes) e alguns casos onde sao uteis, juntamente com parte da sua historia. Por ultimo foi discutido diferentes metodos para se resolver o problema computacional referente a pesquisa de expressoes regulares.

Usando os conceitos apresentados foi demonstrado o processo de construcao do modulo de regex.