

# Sistemas com dados versionados

Thiago Rafael Becker

Dezembro de 2011

## Resumo

A semântica de um sistema com dados versionados é estudada a partir de um mapeamento de um monóide de transformação para um grafo direto reflexivo. O monóide de transformação é composto pelas operações adição de símbolo e remoção de símbolo sobre o conjunto gerado pelo fecho de Kleene  $\Sigma_{\perp}^*$ . O funtor de mapeamento deste grupo para um grafo reflexivo direto apresenta a estrutura e a semântica de um VCS.

## Introdução

Os sistemas com dados versionados (VCS) são sistemas que armazenam um histórico de modificações dos dados que eles contém. Estes sistemas tem diversas funções tais como gerenciamento de documentos, versões de sistemas e configurações destes sistemas para fins de manter histórico de modificações e *rollback* para correção de erros.

O objeto de estudo deste artigo é o cerne da semântica destes sistema, a transição de estados. Para tal, inicialmente vamos definir como são os objetos manipulados por um destes sistemas como um monóide sobre um alfabeto de símbolos,  $\Sigma_{\perp}^*$ . Com essa definição, vamos definir duas operações de adição e remoção de caracteres destes objetos, e provar que estes formam um monóide de transformação,  $T_{\Sigma_{\perp}^*}$ . Por fim vamos definir um funtor **Commit** que mapeia da categoria livremente gerada por  $T_{\Sigma_{\perp}^*}$  para **RGr**, assim fornecendo a estrutura de um grafo para estes sistemas.

## Objetos maniulados

Um SCM gerencia estados de objetos. O conteúdo destes objetos é uma sequência de dados pertencentes ao alfabeto  $\Sigma_{\perp} = \Sigma \cup \perp$ . Todos os estados de um objeto são definidos pelo conjunto gerado pelo fecho de Kleene sobre o alfabeto  $\Sigma_{\perp}$ ,  $\Sigma_{\perp}^* = \langle \Sigma_{\perp}, \circ, \varepsilon \rangle$ .  $\perp$  é um elemento especial que indica indefinição das funções parciais sobre o conjunto  $\Sigma_{\perp}^*$ .

## Algumas definições

**Definição 1** Uma palavra em  $\Sigma_{\perp}^*$  é uma cadeia de símbolos  $\sigma = \{\sigma_1 \sigma_2 \dots \sigma_n\}$  de tamanho  $n$  em  $\Sigma$  ou  $\perp$ .

**Definição 2** O comprimento de uma cadeia de símbolos  $\sigma$  de  $\Sigma_{\perp}^*$  é dado por  $len(\sigma)$ ,  $len(\varepsilon) = len(\perp) = 0$ .

## Adição e remoção de caracteres

Dois funcionais são definidos sobre  $\Sigma_{\perp}^*$ : *adição de caracter* e *remoção de caracter*. A adição de caracter

$$\text{add\_char}_c^k : \mathbb{N} \rightarrow \Sigma \rightarrow (\Sigma_{\perp}^* \rightarrow \Sigma_{\perp}^*) \quad (1)$$

definida como

$$\text{add\_char}_c^k(\sigma) = \begin{cases} \perp & \text{se } len(\sigma) < k, \\ \perp & \text{se } \sigma = \perp, \\ \{\sigma_1 \dots \sigma_{k-1} c \sigma_k \dots \sigma_n\} & \end{cases} \quad (2)$$

insere o símbolo  $c$  na posição  $k$ , e desloca todos os caracteres a direita de  $k$  uma posição à direita, ou retorna  $\perp$  se  $k > \text{len}(\sigma)$ . A remoção de caractere

$$\mathbf{rem\_char}_c^k : \mathbb{N} \rightarrow \Sigma \rightarrow (\Sigma_{\perp}^* \rightarrow \Sigma_{\perp}^*) \quad (3)$$

definida como

$$\mathbf{rem\_char}_c^k(\sigma) = \begin{cases} \perp & \text{se } \text{len}(\sigma) < k, \\ \perp & \text{se } \sigma_k \neq c, \\ \perp & \text{se } \sigma = \perp, \\ \{\sigma_1 \dots \sigma_{k-1} \varepsilon \sigma_{k+1} \dots \sigma_n\} = \{\sigma_1 \dots \sigma_{k-1} \sigma_{k+1} \dots \sigma_n\} & \end{cases} \quad (4)$$

remove o símbolo  $c$  da posição  $k$ , e move todos os caracteres à sua direita uma posição para a esquerda, ou retorna  $\perp$  se

- $k > \text{len}(\sigma)$
- $\sigma_k \neq c$ .

Além disso, a identidade  $\text{id} : \Sigma_{\perp}^* \rightarrow \Sigma_{\perp}^*$  é definida tal que  $\text{id}(\sigma) = \sigma$ .

**Observação 2.1** Note que as operações  $\mathbf{add\_char}_c^k$  e  $\mathbf{rem\_char}_c^k$  não são inversas em  $\Sigma_{\perp}^*$  se  $\mathbf{add\_char}_c^k(\sigma) = \perp$  ou  $\mathbf{rem\_char}_c^k(\sigma) = \perp$  para  $\perp \neq \sigma \in \Sigma_{\perp}^*$ , uma vez que  $(\mathbf{add\_char}_c^k \circ \mathbf{rem\_char}_c^k)(\sigma) = \perp$ . Para todos os outros valores, as operações são inversas,  $(\mathbf{add\_char}_c^k \circ \mathbf{rem\_char}_c^k)(\sigma) = \sigma$  e  $(\mathbf{rem\_char}_c^k \circ \mathbf{add\_char}_c^k)(\sigma) = \sigma$ .

**Observação 2.2** Exemplos de indefinição. Dado o alfabeto  $\Sigma = \{a, b\}$  e a palavra  $\sigma = a^3b^2$  ( $\text{len}(\sigma) = 5$ , as seguintes operações são indefinidas:  $\mathbf{add\_char}_7^a(\sigma)$  e  $\mathbf{rem\_char}_7^a(\sigma)$  (pois  $7 > \text{len}(\sigma)$ ),  $(\mathbf{rem\_char}_3^a \circ \mathbf{rem\_char}_3^a)(\sigma)$  (pois após a primeira operação a posição 3 de  $\sigma'$  é ocupada pelo símbolo  $b$ ).

As funções geradas pelos funcionais  $\mathbf{add\_char}_c^k$  e  $\mathbf{rem\_char}_c^k$  ( $\text{ger}\{\mathbf{add\_char}_c^k\}$ ,  $\text{ger}\{\mathbf{rem\_char}_c^k\}$ ) formam um monóide de transformação  $T_{\Sigma_{\perp}^*} = \langle \text{ger}\{\mathbf{add\_char}_c^k\} \cup \text{ger}\{\mathbf{rem\_char}_c^k\}, \circ, \text{id} \rangle$ . Para provarmos as propriedades deste monóide, definimos o mapeamento  $\Sigma_{\perp}^* \rightarrow \mathbb{N}_{\perp}$

$$\Gamma = \{\perp_{\Sigma_{\perp}} \mapsto \perp_{\mathbb{N} \cup \perp}, \varepsilon_{\Sigma_{\perp}} \mapsto 0_{\mathbb{N} \cup \perp}, c_{\Sigma} \mapsto n_{\mathbb{N}}\} \quad (5)$$

tal que o mapeamento é uma bijeção. O mapeamento das operações  $\mathbf{add\_char}_c^k$  é

$$\Gamma(\mathbf{add\_char}_c^k(\sigma)) = \Gamma(\sigma)p_k^{\Gamma(c)}, p_k \in \mathbb{P} \quad (6)$$

e  $\mathbf{rem\_char}_c^k$  é

$$\Gamma(\mathbf{rem\_char}_c^k(\sigma)) = \Gamma(\sigma)p_k^{-\Gamma(c)}, p_k \in \mathbb{P}, \quad (7)$$

e caso exista  $n > 0$  tal que  $\Gamma(\mathbf{rem\_char}_c^k(\sigma))p_k^{-n} \in \mathbb{N}$ , então  $\mathbf{rem\_char}_c^k(\sigma) = \perp$ . Além disso, se  $\Gamma(\mathbf{rem\_char}_c^k(\sigma)) \notin \mathbb{N}$ , então  $\mathbf{rem\_char}_c^k(\sigma) = \perp$ .

O mapeamento de  $\Gamma(\sigma)$ ,  $\sigma \in \Sigma_{\perp}^*$  é tal que

$$\Gamma(\sigma) = \prod_{i \geq 1} p_i^{\Gamma(\sigma_i)}, p_i \in \mathbb{P}, \sigma_i \in \Sigma_{\perp}^* \quad (8)$$

Este mapeamento é uma bijeção nas classes de equivalência de  $\mathbb{N}$ , tal que  $\Gamma^{-1}(n) = \Gamma^{-1}(n')$ ,  $n \neq n'$ ,  $n, n' \in \mathbb{N}_{\perp}$ , sendo  $\Gamma^{-1}$  o mapeamento inverso a  $\Gamma$ , isto é, o mapeamento de um número para  $\Sigma_{\perp}^*$ . Este mapeamento permite que existam inúmeros números primos entre  $p_k$  e  $p_{k+1}$  usados para representar  $\sigma$ .

Com isso é possível verificar que  $\Gamma(\mathbf{add\_char}_c^k)$  e  $\Gamma(\mathbf{rem\_char}_c^k)$  formam um monóide de transformação sobre  $\mathbb{N} \cup \{\perp\}$ , provando a existencia de  $T_{\Sigma_{\perp}^*}$ .

## Representação como grafo

Com as informações acima, é possível criar um funtor que implementa a operação *commit*, **Commit**. O funtor é tal que

$$\mathbf{Commit}(\Sigma_{\perp}^*) = \sigma_1, \sigma_1 \subset \Sigma_{\perp}^* \quad (9)$$

$$\mathbf{Commit}(ger\{\mathbf{add\_char}_c^k\} \cup ger\{\mathbf{rem\_char}_c^k\} \cup \{id\}) = f : \sigma_1 \rightarrow \sigma_1', \sigma_1, \sigma_1' \in \Sigma_{\perp}^* \quad (10)$$

em que  $\sigma_1$  é um conjunto unário. **Commit** constrói a estrutura do repositório, e é inversível. Intuitivamente, este funtor toma um conjunto  $\sigma_1$  e uma composição de morfismos em  $ger\{\mathbf{add\_char}_c^k\} \cup ger\{\mathbf{rem\_char}_c^k\} \cup \{id\}$  que transforma o estado do repositório.

O grafo que representa um dado repositório é um subgrafo *Repo* do grafo direto reflexivo

$$\langle \mathbf{Commit}(\Sigma_{\perp}^*), \mathbf{Commit}(ger\{\mathbf{add\_char}_c^k\} \cup ger\{\mathbf{rem\_char}_c^k\} \cup \{id\}) \rangle,$$

e é simples verificar as propriedades compositiva e associativa das setas  $\mathbf{Commit}(ger\{\mathbf{add\_char}_c^k\} \cup ger\{\mathbf{rem\_char}_c^k\} \cup \{id\})$  pelo mapeamento  $\Gamma$  acima. Desta forma, a categoria **Repo**

$$\langle \mathbf{Commit}(\Sigma_{\perp}^*), \mathbf{Commit}(ger\{\mathbf{add\_char}_c^k\} \cup ger\{\mathbf{rem\_char}_c^k\} \cup \{id\}), \delta_0, \delta_1, id, \circ \rangle$$

é a categoria livremente gerada pelo grafo *Repo*. O objeto inicial de toda a categoria **Repo** é a palavra vazia  $\varepsilon$ , que representa um repositório vazio.

## Conclusões e desenvolvimentos futuros

O funtor **Commit** da categoria livremente gerada de  $T_{\Sigma_{\perp}^*}$  para **Repo** representa o cerne da semântica dos sistemas com dados versionados. O trabalho é incompleto, entretanto: não trata de sistemas que possuam *merging* e *branching*, nem de metadados comumente presentes nestes sistemas.

É possível perceber que a categoria **Repo** apresenta uma estrutura semelhante ao de um reticulado. Analogamente aos conjuntos parcialmente ordenados vistos como categorias, *branching* poderia ser visto como um produto de dois objetos, e *merging* como o coproduto de dois objetos, mas isto são coisas a se estudar.

Metadados podem ser estudados como endomorfismos etiquetados, porém a construção dos funtores para a criação e movimento destes metadados pode ser complexa. Isso talvez possa ser possível através de uma gramática de grafos sobre *Repo*.

Além disso, o estudo pode se estender para o estudos de todos os sistemas versionados e seus funtores, para verificar que tipo de estrutura tem estes.