



Title	VSOPのruby拡張ライブラリについて
Author(s)	中元, 政一
Citation	2010年度科学技術振興機構ERATO湊離散構造処理系プロジェクト講究録. p.333-335.
Issue Date	2011-06
Doc URL	http://hdl.handle.net/2115/48382
Type	conference presentation
Note	ERATO 湊離散構造処理系プロジェクト春のワークショップ (キックオフシンポジウム). 2010年5月28日 (金) ~ 29日 (土). ERATO湊プロジェクト研究室.
File Information	21.nakamoto_06.pdf



[Instructions for use](#)

VSOPのruby 拡張ライブラリについて

2010年5月28日
ERATO 湊離散構造処理系プロジェクト
統計・マイニング応用グループ 技術員 中元 一

背景と目的

背景

情報大公開プロジェクトにおいて、データマイニングツールであるIGVminerをRubyで開発。このIGVminer内部でVSOP(ZDD)を使用する必要があった

目的

RubyでVSOP(ZDD)を使えるようにすることによって、VSOP(ZDD)をシームレスにruby上で利用可能になる

Rubyを使用した理由

日本で普及しているRubyで仕様を固め、次の段階でPerlやRythonの他言語への移植を考えた。

記述例

VSOP

```
A = a b c
B = 4
C = A * B
print C

D = Lcm("FQ", "tra.txt", 2)
D.restrict(x1 x2).show
```

RUBY拡張

ruby 1.8.6 を使用

```
require "rubygems"
require "vsop"
A = VSOPItemset("a b c")
B = VSOPconstant(4)
C = A * B
C.show

D = VSOPlcm("FQ", "tra.txt", 2)
D.restrict("x1 x2")
```

- vsopオブジェクトを作成および出力をのぞけば、ほぼ同じ記述で処理を行える
- vsopオブジェクトをrubyオブジェクトとして作成しているので、rubyの制御構造(while, ifなど)使用できる
- vsopではITEM名を[a-z][a-zA-Z0-9_]*にする必要があったがruby拡張では意識する必要がない
- rubyの文字列や数値あるいはCSVファイルを相互に変換可能である

VSOPの命令とruby拡張ライブラリのメソッドとの対応表

VSOP	RUBY拡張	VSOP	RUBY拡張	VSOP	RUBY拡張	VSOP	RUBY拡張
Maxval	maxval	+	+	print	show	print /plot	未対応
Minval	minval	/	/	print /bit	bit	print /plot0	未対応
TotalVal	totalval	%	%	print /hex	hex	symbol	未対応
Items	items	+	+	print /map	map	print /maxcover	maxcover
TermsEQ	termsEQ	+(単項演算子)	+(単項演算子)	print /rmap	rmap	print /mincover	mincover
TermsNE	termsNE	-	-	print /case	case	print /nodesize	nodesize
TermsGT	termsGT	-(単項演算子)	-(単項演算子)	print /size	size	print /out	out
TermsGE	termsGE	==	==	print /count	count	print /export	export
TermsLT	termsLT	>	>	print /density	density		
TermsLE	termsLE	>=	>=	print /value	value		
Restrict	restrict	<	<	print /maxcover	maxcover		
Permit	permit	<=	<=	print /maxcost	maxcost		
Permitsym	permitsym	!=	ne!	print /mincover	mincover		
FreqPatA	freqpatA	!	if	print /mincost	mincost		
FreqPatM	freqpatM	import	import	print /decomp	decomp		
FreqPatC	freqpatC	Lcm	lcm	print /nodesize	nodesize		
SymGrp	symgrp			print /export	export		

RUBY拡張固有メソッド

- itemset**
機能: 引数で指定したアイテムから新しいVSOPオブジェクトを作成
書式: s = VSOPitemset("a b c")
- constant**
機能: 引数で指定した整数から新しいVSOPオブジェクトを作成
書式: c = VSOPconstant(10)
- corece**
機能: rubyオブジェクトからVSOPオブジェクトを作成
書式: s = VSOPitemset(rubyobj)
- to_j**
機能: VSOPオブジェクトから数値を生成 (VSOPオブジェクトが数値の場合のみ)
書式: val = c.to_j
- lcm_ub**
機能: パターン長制約を指定できるlcm over add
書式: VSOP.lcm_ub(["FMC"], tra_filename, minimum_support, unper_bound, order, filename)
- csvout**
機能: VSOPオブジェクトを指定したファイルにCSV出力
書式: s.csvout(outfilename)
- hashout**
機能: VSOPオブジェクトをrubyのハッシュオブジェクトへ変換
書式: h.hashout

出力例(重み: itemset)

```
2 a b
3 c d
2 e
```

→ 2 a b + 3 c d + 2 e

応用例

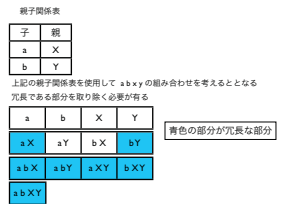
rubyのハッシュオブジェクトに変換する事により
h["x"]とする事で
"x"に対する重みを取得できるようになる

使用例(1)

taxonomyの冗長パターンの削除

前述のIGVminerで使用した記述の抜粋であり、lcmで計算した頻出パターンから親子関係のあるitemのセットを含むパターンを取り除いている

```
pat = VSOPlcm("FQ", "tradata", 4)
taxo.each{|f| f.val}
oyako = VSOPconstant(0)
items = f.val.split
items.each{|item1|
  items.each{|item2|
    if item1 != item2 then
      oyako = oyako + VSOPitemset("#{item1} #{item2}")
    end
  }
}
pat = pat.restrict(oyako).if(0, pat)
```



使用例 (2)

● n-queen問題

NXNのマス上にN個のQUEENを互いにた
れないように配置するパターンを求める

```
require "rbygens"
require "rsop"
n = 8
n = ARGV[0].to_i if ARGV[0]
s_base=VSOPconstant(0)
s_new=VSOPconstant(0)
(1..n).each{|i|
  (1..n).each{|j|
    tmp = s_base
    (1..i).each{|k|
      tmp |= ("x" + i + j + k) if i+k > 0
      tmp |= ("x" + i + j + k) if i+k > 0 && j+k > 0
      tmp |= ("x" + i + j + k) if i+k > 0 && j+k < n
    }
    if i > 1 then
      s_new += (tmp + s_base) if j == 1
    else
      s_new += VSOPzset("x" + i + j)
    end
  }
}
s_base = s_new
s_base = VSOPconstant(0)
s_base.show
s_base.csvout("x.csv")
```

参考文献: Minato, Shin-ichi. Zero-suppressed BDDs and their applications. 2001

今後について

問題点 (開発上の課題)

- ソース整理
- データの前処理
- 異なるデータのZDDでの組み合わせ
- VSOPで重みに小数をつかえるようにする

応用課題

- OLAP (online analytical processing) への応用
- ZDDのstream

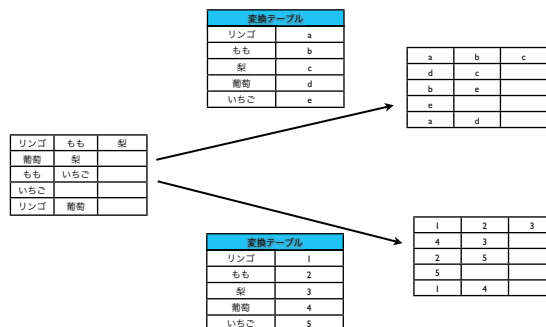
ソースの整理

現状、sed変換により元々あるソースの変換
をおこない実装している。今後の事を考え
ると、この変換処理をなくし、ソースを整理
していく必要が有る

symbol名の変換処理

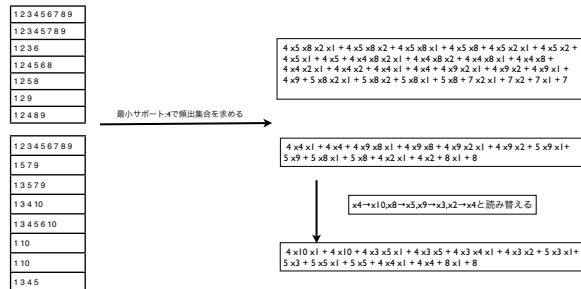
使用者がデータを [a-z][a-zA-Z0-9]*の形へ変換する必要がある
(LCM over zddを使用する場合には[0-9][1-9][0-9]*)

ruby拡張の場合、現在は
LCM over zddを使用するときの変換必要



LCM over ZDDの使用時の 異なるデータのZDDでの組合せ

異なるデータでLCM over ZDDを使用する際、アイテムの読み替えが発生する



ZDDで小数をデフォルトで対応

小数桁数を固定すれば現状でも対応可能だ
が、前処理、後処理が必要となる。また
OLAPでは小数の処理が多くなると考えら
れ、デフォルトで小数に対応した方がよい
と考える。

OLAP(online analytical processing)への応用

多次元データをZDDで作成

key1	key2	key3	val
a	A	X	5
a	A	Y	2
a	A	Z	5
a	B	X	6
a	B	Y	8
a	B	Z	10
b	A	X	3
b	A	Y	14
b	A	Z	13
b	B	X	6
b	B	Y	30
b	B	Z	31

$DB = 5aAX + 2aAY + 5aAZ + 12aA + 6aBX + 8aBY + 10aBZ + 24aB + 11aX + 10aY + 15aZ + 36a + 3bAX + 14bAY + 13bAZ + 30bA + 6bBX + 30bBY + 31bBZ + 67bB + 9bX + 44bY + 44bZ + 97b + 8AX + 16AY + 18AZ + 42A + 12BX + 38BY + 41BZ + 91B + 20X + 54Y + 59Z + 133$

keyによる分析が可能

複数のZDDを作成することで複数valに対応可能
小数桁を固定すれば実数対応可能

ITEMに属性を持たせることで、より詳細な分析が可能
例) ドリルダウン

参考文献: Minato, Shin-ichi, Zero-suppressed BDDs and their applications, 2001

ZDDのstream

ZDDをstream化する事で並列処理することが可能となり、効率よく処理する事が可能。
例えば、ある企業の支店ごとのデータ並列処理し、集める事がファイル出力なしに実行可能になる

参考文献: 渡真一、石原晋也、BDDの規模によらず一定の実記憶の範囲内で動作するストリーム形式BDD処理アルゴリズム、1999

今後についてのまとめと計画

以上の内容を実行していくにはZDD、LCMのアルゴリズムおよび実装方法について、深い理解が必要である。

まずは理解し、その後、整理、および応用へとつなげていく。