

Relatório Científico Final – Iniciação Científica

Processo FAPESP 2016/25959-7

Projeto de algoritmos baseados em florestas de posets  
para o problema de otimização U-curve

**Beneficiário:** Gustavo Estrela de Matos

**Responsável:** Marcelo da Silva Reis

Relatório referente aos trabalhos desenvolvidos entre 1 de maio e 31  
de dezembro de 2017

Laboratório Especial de Toxinologia Aplicada, Instituto Butantan

São Paulo, 7 de Janeiro de 2018

# Conteúdo

<b>1</b>	<b>Resumo do Projeto Proposto</b>	<b>2</b>
<b>2</b>	<b>Atividades Realizadas</b>	<b>2</b>
2.1	Estudo de algoritmos baseados em florestas . . . . .	3
2.2	Modificações do PFS na escolha de raízes . . . . .	4
2.3	Mudificação do PFS no armazenamento de raízes . . . . .	6
2.4	Paralelização do PFS . . . . .	6
2.5	Elaboração do UBB-PFS . . . . .	6
2.6	Elaboração de um algoritmo de aproximação . . . . .	6
2.7	Testes com instâncias reais do problema de seleção de características . . . . .	6
<b>3</b>	<b>Avaliação e disseminação de resultados</b>	<b>6</b>
<b>4</b>	<b>Conclusão</b>	<b>6</b>
	<b>Referências</b>	<b>7</b>

# 1 Resumo do Projeto Proposto

O problema U-curve é uma formulação de um problema de otimização que pode ser utilizado na etapa de seleção de características em Aprendizado de Máquina, com aplicações em desenho de modelos computacionais de sistemas biológicos. Não obstante, soluções propostas até o presente momento para atacar esse problema têm limitações do ponto de vista de consumo de tempo computacional e/ou de memória, o que implica na necessidade do desenvolvimento de novos algoritmos. Nesse sentido, em 2012 foi proposto o algoritmo **Poset-  
-Forest-Search** (PFS) [1], que organiza o espaço de busca em florestas de posets. Esse algoritmo foi implementado e testado, com resultados promissores; todavia, novos melhoramentos são necessários para que o PFS se torne uma alternativa competitiva para resolver o problema U-curve. Neste projeto propomos modificações ao PFS na escolha de caminhos de percorrimento da floresta de busca, e na estrutura de dados utilizada para armazenar este grafo, com o uso de diagramas de decisão binária reduzidos e ordenados (OBDDs) [3]; também propomos a criação de uma versão paralela e escalável do algoritmo PFS. Além disso, propomos a criação de um algoritmo baseado no PFS que tenha características de um algoritmo de aproximação, no qual o critério de aproximação da solução ótima se baseie no teorema da navalha de Ockham. Os algoritmos desenvolvidos serão implementados no arcabouço *featsel* [2] e testados com instâncias artificiais e também reais, com conjuntos de dados de aprendizado de máquina retirados do University of California Irvine (UCI) Machine Learning Repository.

## 2 Atividades Realizadas

Uma implementação do PFS, feita por Reis, está disponível no arcabouço *featsel*. Usamos esta implementação como base para estudar as modificações feitas ao PFS.

## 2.1 Estudo de algoritmos baseados em florestas

O algoritmo **Poset-Forest-Search** (PFS) é um algoritmo ótimo para resolver o problema de otimização U-Curve e serviu de base para a criação da maioria dos algoritmos elaborados neste trabalho. O PFS é uma generalização de um outro algoritmo mais simples, o **U-curve-Branch-and-Bound** (UBB), que é um algoritmo *branch-and-bound* ótimo que decompõe o espaço de busca em uma árvore, e acha o mínimo global do problema fazendo ramificações e podas nesta árvore.

A árvore de busca do UBB permite que a procura pelo mínimo ocorra de maneira parecida com uma busca em profundidade, que percorre cadeias do reticulado Booleano de subconjuntos menores para maiores. Sempre que o custo de um subconjunto  $X_i$  aumenta em comparação ao anterior  $X_j$  no percorrimento, a hipótese de que a função de custo é decomponível em curvas em U garante que a subárvore que começa em  $X_i$  pode ser removida do espaço de busca; chamamos este procedimento de poda. O algoritmo UBB tem, entretanto, uma limitação, pois quando a função de custo do problema é monótona não-crescente, a condição de poda nunca é verdadeira e o espaço de busca inteiro é visitado, o que compromete a escalabilidade do algoritmo.

O PFS enfrenta esta limitação ao fazer percorrimentos de cadeias do espaço de busca em duas direções, de conjuntos menores para maiores (como faz o UBB) e também o contrário. Para fazer isso, este algoritmo decompõe o espaço de busca em duas árvores, uma para cada direção de percorrimento. Com a criação de duas estruturas para representar o mesmo espaço de busca, torna-se necessário a atualização de uma estrutura sempre que a outra sofrer mudanças, e isto implica na utilização de florestas ao invés de árvores para representar o espaço de busca no PFS. Resumidamente, uma iteração do deste algoritmo deve escolher uma direção de percorrimento; fazer o percorrimento com poda (de maneira similar ao UBB); e, finalmente, atualizar a floresta dual a que foi percorrida.

## 2.2 Modificações do PFS na escolha de raízes

Para se fazer um percorrimento no espaço de busca, o algoritmo PFS deve escolher uma direção de percorrimento, isto é, uma das duas florestas que representam o espaço de busca, e também uma raiz da floresta escolhida. Esta escolha é feita de maneira arbitrária, e por conta disso, investigamos como diferentes escolhas podem afetar o desempenho do algoritmo.

Na implementação de Reis, a estrutura de dados utilizada para armazenar raízes é a *map* do C++. Escolhe-se nesta implementação o primeiro ou último elemento da estrutura, o que coincide com a primeira ou última raiz quando estas estão ordenadas lexicograficamente por seus vetores característicos. Em nosso trabalho, experimentamos duas modificações para esta escolha:

- de maneira aleatória e uniformemente provável;
- e de maneira determinística, com a raiz de maior sub-árvore completa. Neste caso, o tamanho da sub-árvore completa é o tamanho deste grafo quando nenhuma poda foi feita.

A justificativa para se fazer uma escolha uniforme entre as raízes é ter um algoritmo que não possui viés na escolha de percorrimentos, assim podemos investigar se o viés da escolha arbitrária feita na implementação de Reis compromete a execução do algoritmo. Para fazer a implementação da escolha aleatória e uniforme fizemos uma pequena modificação ao código de Reis. Dado uma floresta  $\mathcal{F} = \{r_1, r_2, \dots, r_l\}$  sorteia-se um número pseudo-aleatório  $a$  entre 1 e  $l$  e escolhe-se a raiz  $r_a$  para o percorrimento.

A escolha de raiz com maior sub-árvore foi feita com a intuição de que percorrimentos em árvores maiores implicariam em maiores podas e consequentemente menos nós visitados no espaço de busca. A decomposição do espaço de busca em árvore feito por Reis faz com que ordenar as raízes por tamanho decrescente de sub-árvores completas coincida com uma ordenação lexicográfica da direita para a esquerda dos vetores característicos das raízes. Para fazer isto, basta modificar a ordenação feita pela estrutura *map* do C++.

Tabela 1: Comparação entre os algoritmos *PFS* e *PFS\_RANDOM*. O tempo de execução do segundo é maior do que o primeiro enquanto a quantidade de chamadas da função custo é parecida em ambos.

Instância		Tempo de execução médio (s)		Número médio de cálculos de custo	
$ S $	$2^{ S }$	PFS	PFS_RANDOM	PFS	PFS_RANDOM
10	1024	$0.013 \pm 0.003$	$0.014 \pm 0.003$	$590.8 \pm 198.5$	$599.5 \pm 177.5$
11	2048	$0.019 \pm 0.004$	$0.022 \pm 0.007$	$1114.8 \pm 331.3$	$1090.1 \pm 350.3$
12	4096	$0.029 \pm 0.008$	$0.036 \pm 0.013$	$1848.6 \pm 600.8$	$1835.7 \pm 683.0$
13	8192	$0.060 \pm 0.018$	$0.090 \pm 0.039$	$4314.4 \pm 1496.4$	$4201.1 \pm 1580.7$
14	16384	$0.100 \pm 0.041$	$0.191 \pm 0.110$	$7323.4 \pm 3318.9$	$7333.8 \pm 3161.0$
15	32768	$0.180 \pm 0.076$	$0.453 \pm 0.311$	$12958.1 \pm 5654.0$	$12807.5 \pm 5753.7$
16	65536	$0.406 \pm 0.185$	$1.715 \pm 1.400$	$27573.8 \pm 12459.5$	$27036.9 \pm 12687.5$
17	131072	$0.717 \pm 0.397$	$5.416 \pm 5.266$	$48176.2 \pm 26938.3$	$47852.1 \pm 26427.6$
18	262144	$1.325 \pm 0.754$	$15.890 \pm 17.726$	$84417.9 \pm 48587.7$	$84025.0 \pm 48882.4$
19	524288	$2.771 \pm 1.603$	$69.600 \pm 82.342$	$167659.1 \pm 99686.7$	$164612.1 \pm 102018.3$

Chamamos o algoritmo que faz a escolha uniforme das raízes de *PFS\_RANDOM*, e os resultados de tempo de execução e número de chamadas da função de custo são apresentados na tabela 1. Podemos observar que esta modificação ao *PFS* não foi benéfica, pois comparando com a implementação de Reis o tempo de execução médio aumentou, e o número médio de chamadas da função de custo continua parecido. O fato do número de chamadas da função de custo não ter diferenças significativas implica que esta escolha de raiz não trouxe mudanças a dinâmica do algoritmo original, logo o tempo a mais de execução veio do código que faz a escolha da raiz. Isto é feito com o percorrimento da estrutura de *map*, o que adiciona tempo linear (sobre a quantidade de raízes) a cada escolha de raiz.

Chamamos de *PFS\_LEFTMOST* o algoritmo que faz a escolha da raiz com maior sub-árvore completa, e a tabela 2 mostra resultados de tempo de execução e número de chamadas da função de custo desta variante comparado a implementação de Reis. Podemos observar que esta modificação também não foi benéfica pois, além do maior tempo de execução, o número de chamadas da função de custo aumentou. Isto significa que o comportamento do algoritmo foi o contrário a intuição que motivou a modificação, pois mais nós do espaço de busca foram visitados.

Tabela 2: Comparação entre os algoritmos PFS e PFS\_LEFTMOST. O tempo de execução e também o número de chamadas da função custo é maior para o PFS\_LEFTMOST.

Instância		Tempo de execução médio (s)		Número médio de cálculos de custo	
$ S $	$2^{ S }$	PFS	PFS_LEFTMOST	PFS	PFS_LEFTMOST
10	1024	$0.013 \pm 0.002$	$0.023 \pm 0.004$	$606.1 \pm 133.5$	$665.0 \pm 165.8$
11	2048	$0.020 \pm 0.004$	$0.042 \pm 0.010$	$1122.1 \pm 351.2$	$1316.6 \pm 382.2$
12	4096	$0.032 \pm 0.008$	$0.078 \pm 0.024$	$2183.7 \pm 733.2$	$2515.8 \pm 871.3$
13	8192	$0.054 \pm 0.017$	$0.160 \pm 0.061$	$3887.7 \pm 1389.9$	$4716.8 \pm 1777.8$
14	16384	$0.107 \pm 0.034$	$0.345 \pm 0.133$	$7851.2 \pm 2793.0$	$9506.8 \pm 3673.9$
15	32768	$0.196 \pm 0.085$	$0.672 \pm 0.274$	$13780.3 \pm 6049.9$	$17071.6 \pm 7005.1$
16	65536	$0.348 \pm 0.189$	$1.271 \pm 0.661$	$24106.5 \pm 13159.9$	$30055.6 \pm 15363.6$
17	131072	$0.785 \pm 0.361$	$3.137 \pm 1.476$	$52369.0 \pm 24751.2$	$67585.6 \pm 30978.4$
18	262144	$1.445 \pm 0.657$	$6.146 \pm 3.032$	$92095.9 \pm 42566.6$	$120635.7 \pm 58039.0$
19	524288	$3.298 \pm 1.883$	$13.881 \pm 7.595$	$199151.0 \pm 112167.8$	$256078.6 \pm 135958.4$

## 2.3 Mudificação do PFS no armazenamento de raízes

## 2.4 Paralelização do PFS

## 2.5 Elaboração do UBB-PFS

## 2.6 Elaboração de um algoritmo de aproximação

## 2.7 Testes com instâncias reais do problema de seleção de características

# 3 Avaliação e disseminação de resultados

# 4 Conclusão

## Referências

- [1] Reis, Marcelo S. “Minimization of decomposable in U-shaped curves functions defined on poset chains—algorithms and applications.” PhD thesis, Institute of Mathematics and Statistics, University of São Paulo, Brazil, (2012).
- [2] Reis, Marcelo S., Gustavo Estrela, Carlos E. Ferreira and Junior Barrera. “featsel: A framework for benchmarking of feature selection algorithms and cost functions”. *SoftwareX* 6 (2017), pp. 193-197.
- [3] Bryant, Randal E. “Graph-based algorithms for boolean function manipulation.” *IEEE Transactions on Computers*, 100.8 (1986): 677-691.