

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, 4 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 . . . . .	2
2.2	Modificações do PFS na escolha de raízes . . . . .	4
2.3	Mudificação do PFS no armazenamento de raízes . . . . .	4
2.4	Paralelização do PFS . . . . .	4
2.5	Elaboração do UBB-PFS . . . . .	4
2.6	Elaboração de um algoritmo de aproximação . . . . .	4
2.7	Testes com instâncias reais do problema de seleção de características . . . . .	4
<b>3</b>	<b>Avaliação e disseminação de resultados</b>	<b>4</b>
<b>4</b>	<b>Conclusão</b>	<b>4</b>
	<b>Referências</b>	<b>5</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), 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 percurso 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); 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* 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

### 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 percorrimto, 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 percorrimtos 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 percorrimto. 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 percorrimto; fazer o percorrimto 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
- 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., Carlos E. Ferreira, and Junior Barrera. “The U-curve optimization problem: improvements on the original algorithm and time complexity analysis.” arXiv preprint arXiv:1407.6067, (2014).
- [3] Bryant, Randal E. “Graph-based algorithms for boolean function manipulation.” IEEE Transactions on Computers, 100.8 (1986): 677-691.