

# A global feature selection algorithm for the model selection step in the identification of cell signaling networks

Gustavo Estrela<sup>1,2,3</sup>, Lulu Wu<sup>1,2</sup>, Vincent Noël<sup>1,3</sup>, Carlos Eduardo Ferreira<sup>2</sup>, Hugo A. Armelin<sup>1,3</sup>, Marco Dimas Gubitoso<sup>2</sup>, Junior Barrera<sup>1,2</sup>, and Marcelo S. Reis<sup>1</sup>

<sup>1</sup>Center of Toxins, Immune-response and Cell Signaling (CeTICS), Instituto Butantan, Brazil

<sup>2</sup>Instituto de Matemática e Estatística, Universidade de São Paulo, Brazil

<sup>3</sup>Laboratório Especial de Ciclo Celular (LECC), Instituto Butantan, Brazil

<sup>4</sup>Instituto de Química, Universidade de São Paulo, Brazil



## Motivation

In the context of Machine Learning, the feature selection problem consists in choosing a subset of features that best explains the classification with minimum redundancy. The space of solutions of this problem induces a Boolean lattice and the cost function commonly describes U-shaped curves on chains of this lattice, what is explained by the increase of estimation error as we include more features. Hence, we can approximate this problem to the U-Curve problem, which is a special case of the feature selection problem where every chain of the search space describes U-shaped curves. Some algorithms in the literature exploit this approximation; still, they show limitations regarding scalability, which might be a problem for the feature selection step in the identification of cell signaling networks. To tackle this issue, we developed the Parallelized U-Curve Search (PUCS).

## Method

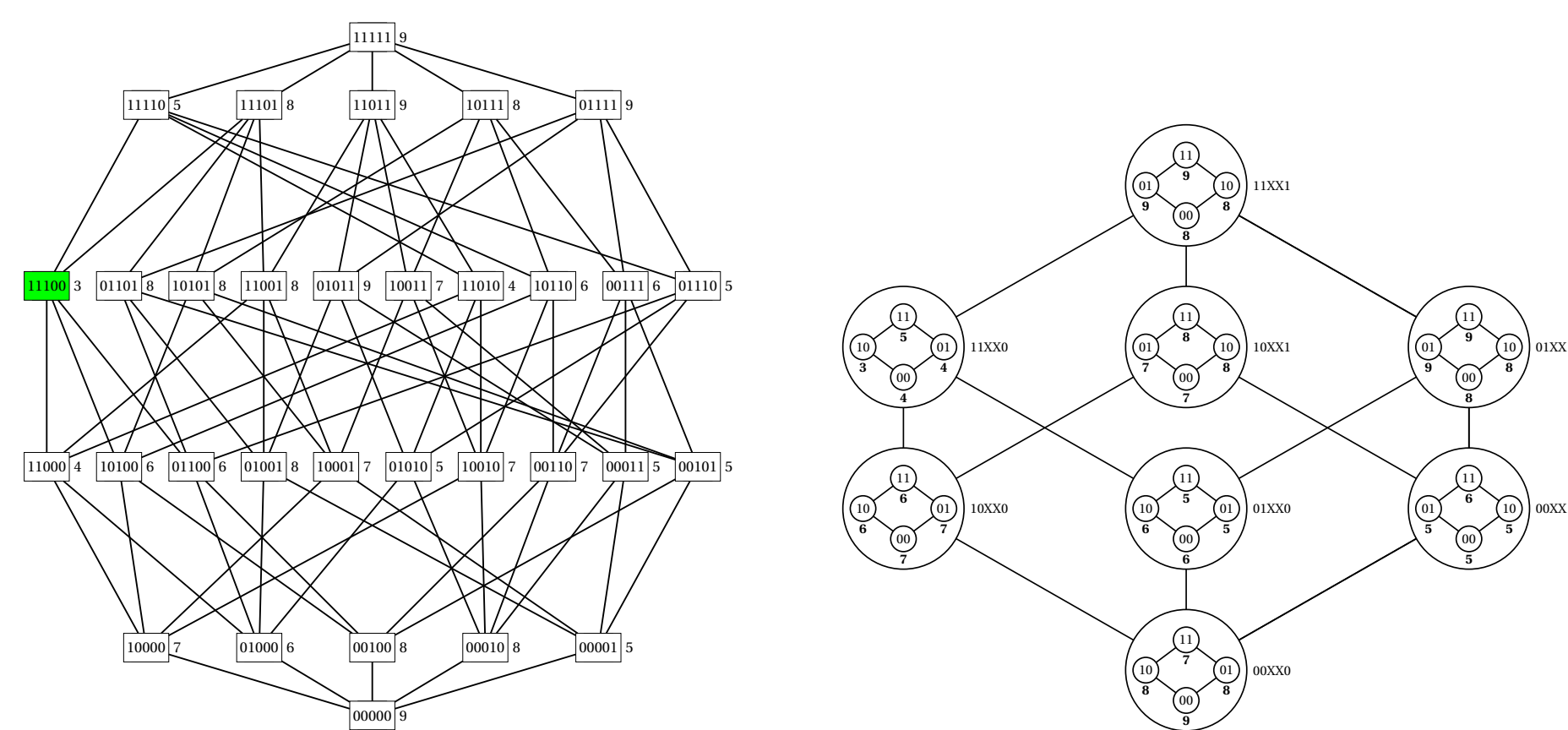


Fig. 1: An instance of the U-curve problem with  $|S| = 5$  (left) and induced search space when the third and fourth elements are regarded as don't care (right).

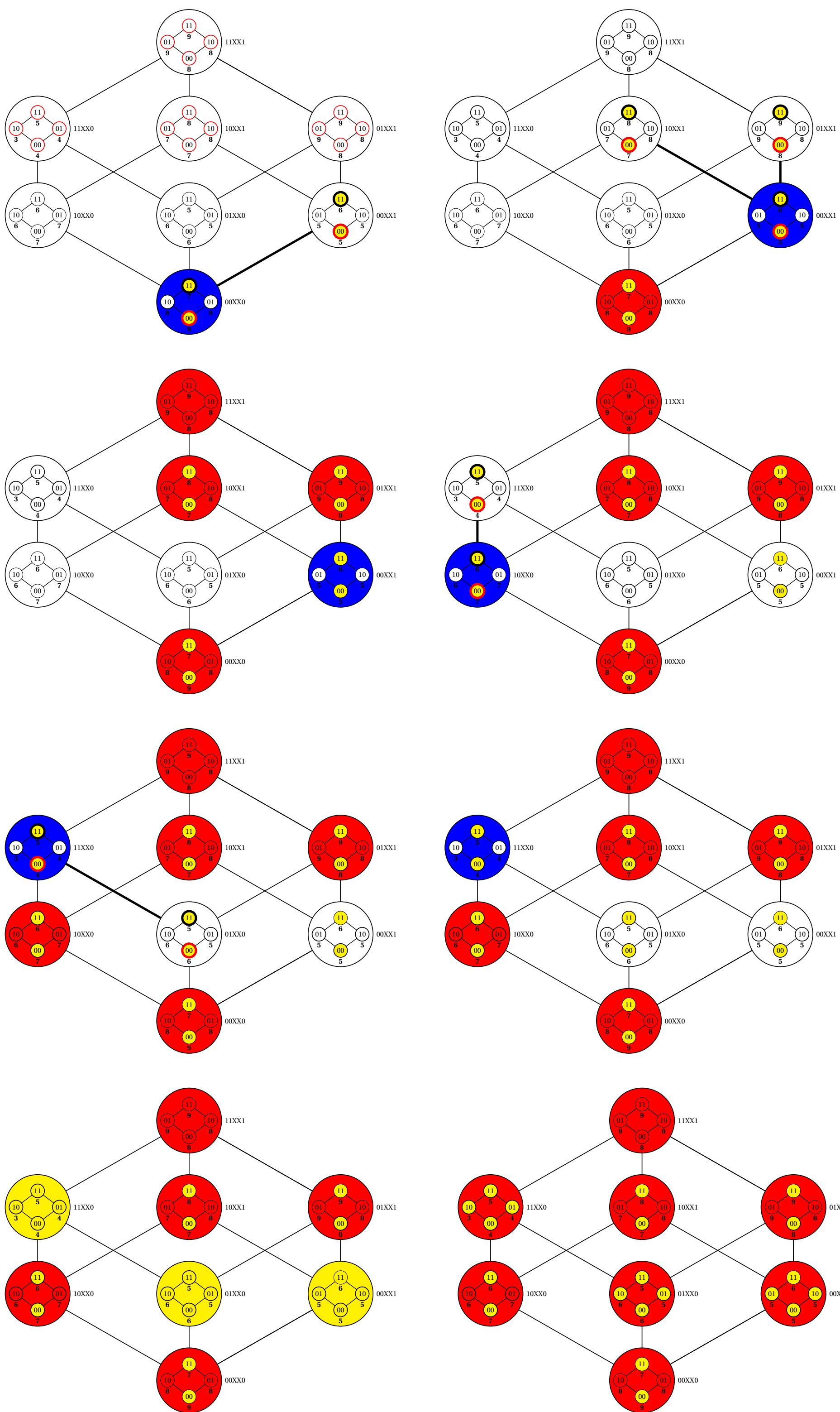
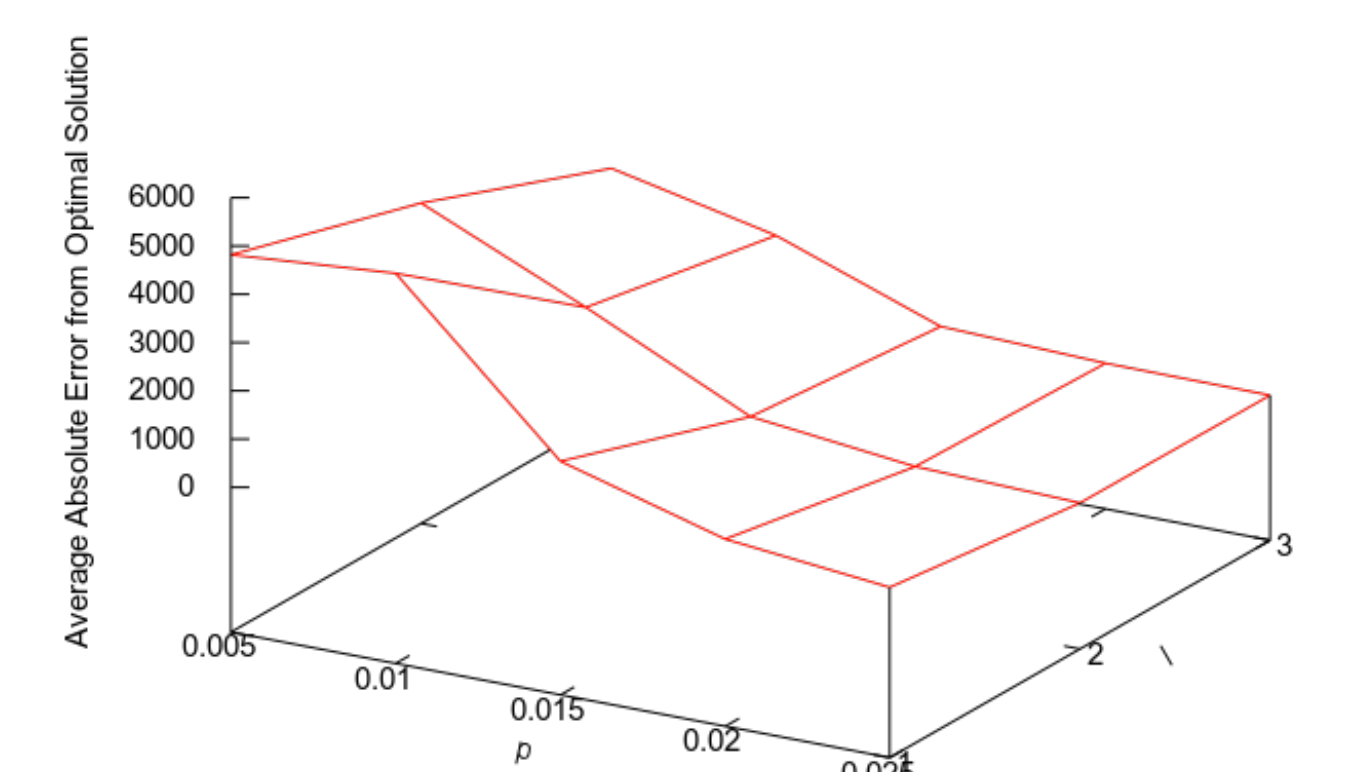


Fig. 2: Dynamics of the PUCS algorithm. [EXPLICAR O QUE SIGNIFICAM AS CORES]

## Acknowledgements

## Results

We implemented the PUCS algorithm in C++ language and used OpenMP to parallelize the code. Using a 64-core, 256 GB RAM server, we were able to confirm our expectations that the algorithm can find solutions as good as it possible (be optimal) as long as we increase the parameters  $p$  and  $l$ . [EXPLICAR O QUE SÃO  $p$  E  $l$ ]



We used the *featsel* framework ([github.com/msreis/featsel](https://github.com/msreis/featsel)) to benchmark PUCS with other algorithms, such as Exhaustive Search (ES), Sequential Forward Selection (SFS) and Backward Feature Selection (BFS).

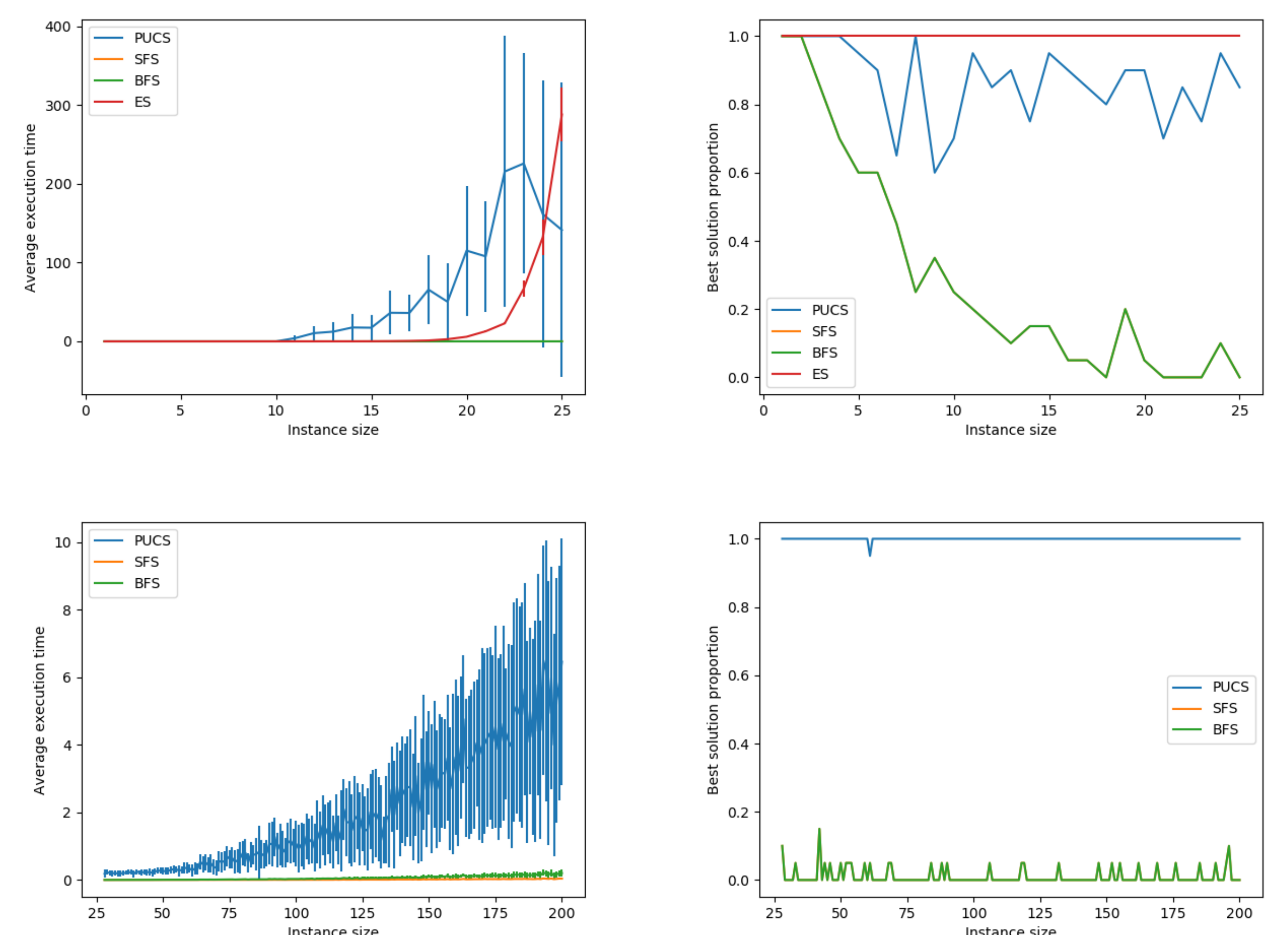


Fig. 3: Average execution time and average number of times each algorithm found the best solution. For these instances, we used SFS to solve recursion base cases of the PUCS algorithm.

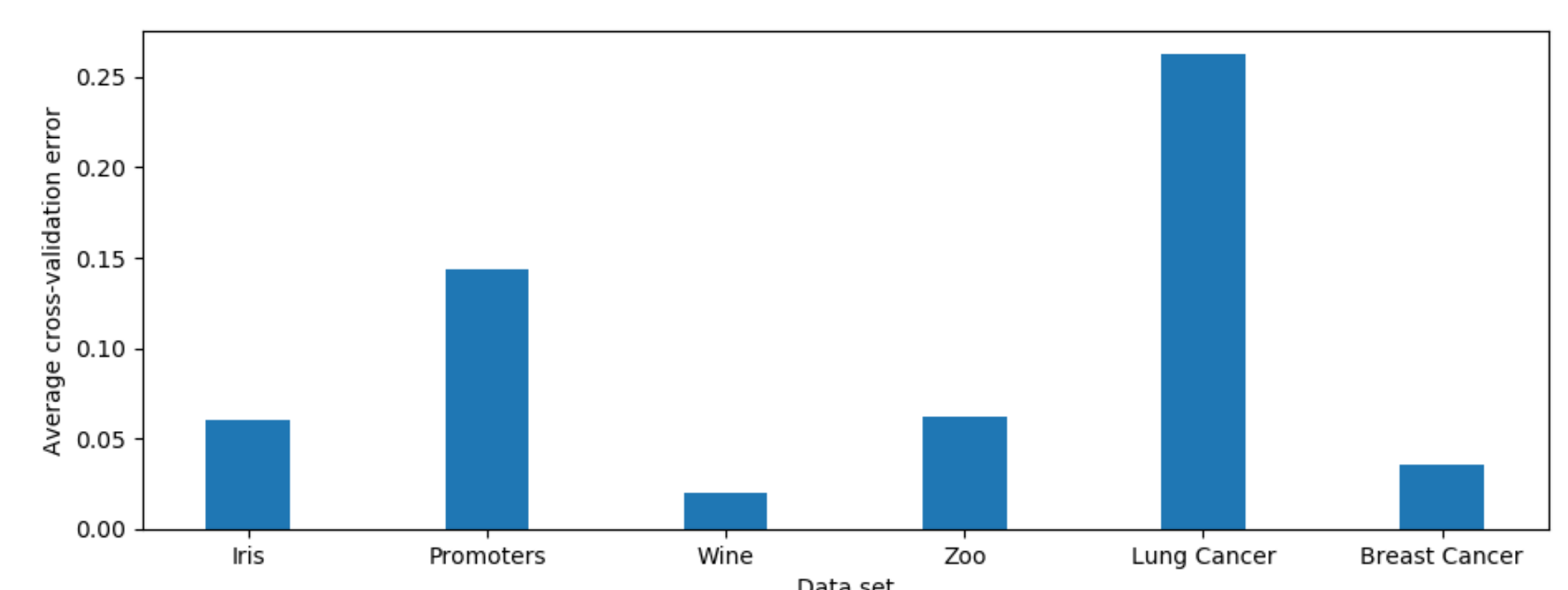


Fig. 4: Average cross-validation error when using the features selected by PUCS to design SVM classifiers for some biological databases available at the UCI Machine Learning Repository.

## Conclusion

Results showed that PUCS has a better time performance than optimal solvers such as ES, and also can find better solutions than heuristics such as SFS and BFS. Future and ongoing work on this research line includes:

- Design of new algorithms, especially ones that generalize the search space for poset forests;
- Application of PUCS on the feature selection step in identification of cell signaling pathways.