

Classificação de Objetos Astronômicos com Algoritmo Random Forest: um comparativo entre implementações single e multi-threaded

Pedro Montuani

Introdução

- ▶ Humanidade vem produzindo um volume massivo de dados (aumentando a cada ano)
- ▶ Isso se reflete em diversas áreas humanas, incluindo na ciência

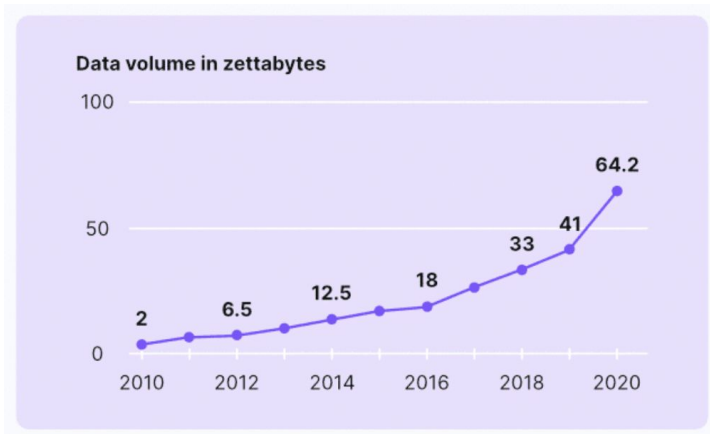


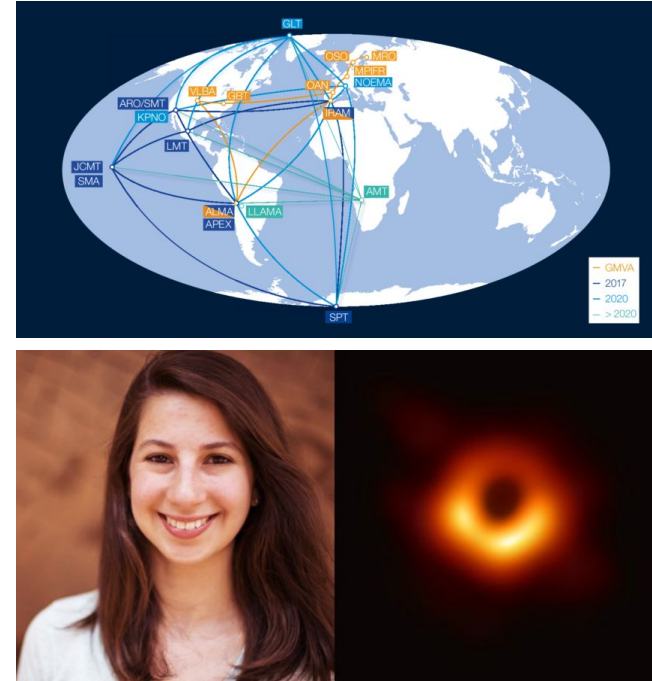
Fig. 01: volume de dados (ZB) produzidos na última década

Introdução

- ▶ Astronomia: instrumentos modernos, novas técnicas e tecnologias -> grande volume de dados
- ▶ Necessidade de metodologias eficientes para lidar com essa quantidade de dados
- ▶ Machine learning, algoritmos de classificação, reconhecimento de padrões, simulações de rotas e órbitas

Introdução

- ▶ Projeto CHIRP (Continuous High-resolution Image Reconstruction using Patch priors)
- ▶ Sincronizou vários telescópios ao redor do planeta para simular uma lente gigante
- ▶ Utilizou algoritmos de classificação para diferenciar buracos negros de outros objetos
- ▶ Utilizou algoritmos de combinação reconstrução de imagem



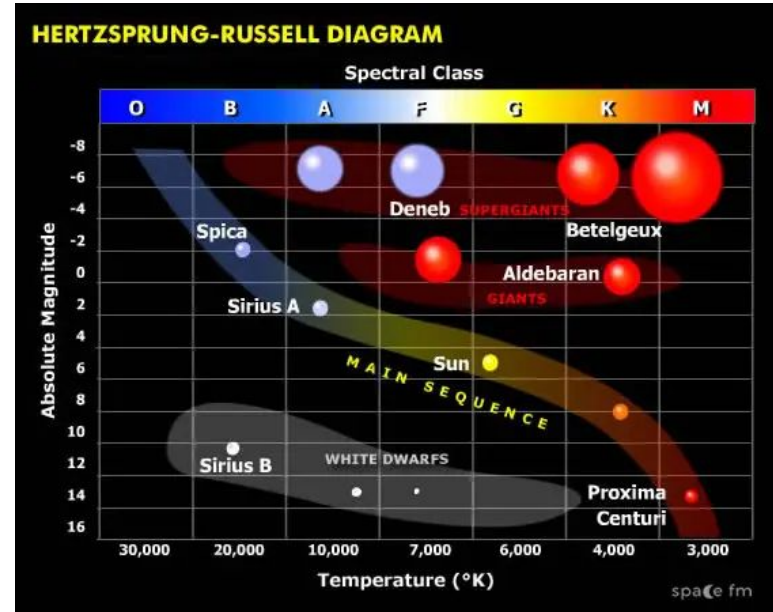
Figs. 02 e 03: Event Horizon Telescope, Katie Bouman e a primeira imagem de um buraco negro

Objetivos

The background features a dark gray area on the left and top, an orange area on the right, and a white horizontal bar at the bottom. A diagonal line separates the dark gray and orange sections.

Objetivos

- ▶ Aplicar implementação de Random Forest em dataset de objetos astronômicos para classificação
- ▶ Implementar uma solução *multi-thread* (OMP)
- ▶ Comparar sklearn, C++ single, e C++ multi-thread

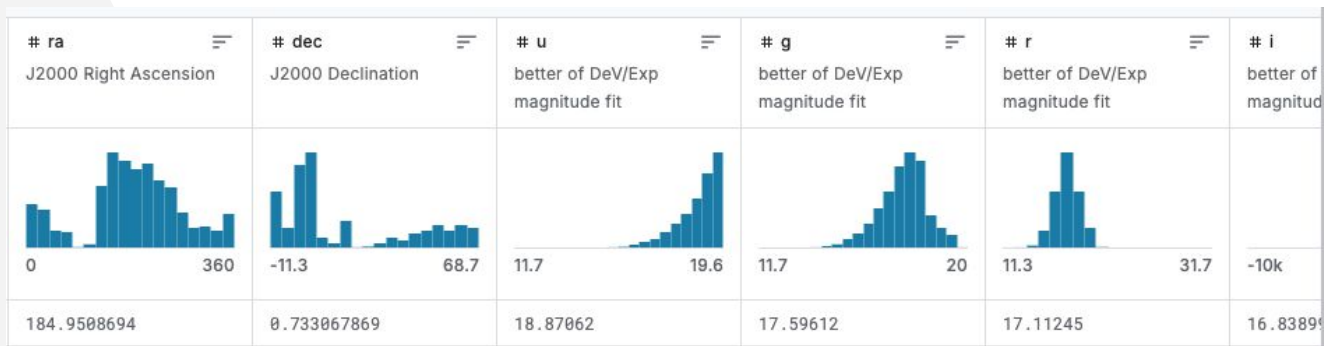


Metodologia

The background features a dark gray upper-left section, a bright orange upper-right section, and a dark gray lower-left section. A white horizontal bar is positioned below the word 'Metodologia'. A diagonal line separates the dark gray and orange areas.

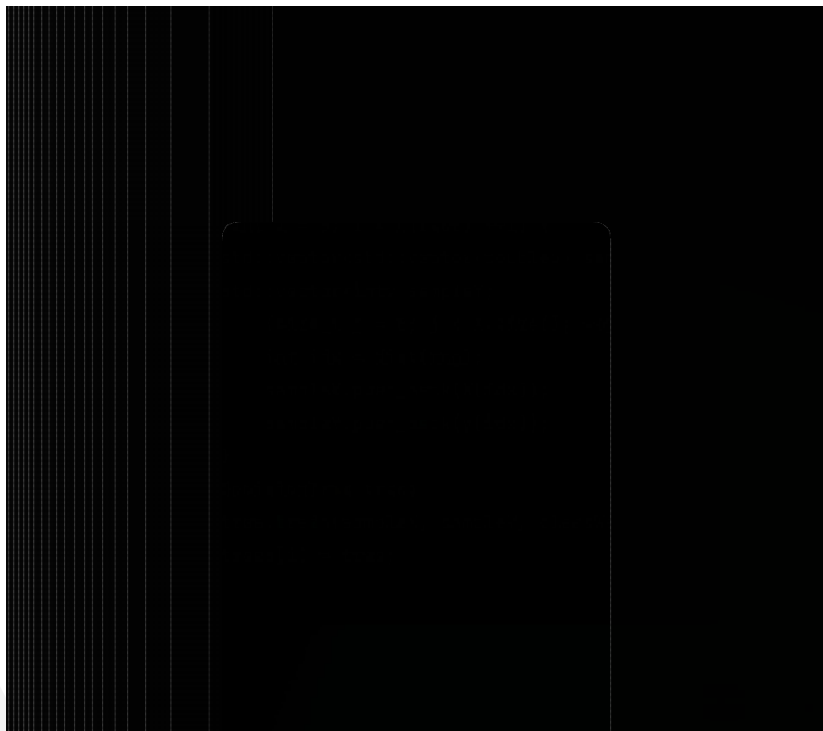
Metodologia

- ▶ Dataset -> Sloan Digital Sky Survey - DR18
 - ▷ CSV com 10.000 entradas de 43 *features*
 - ▷ Objetos classificados em STAR, GALAXY e QSO (*Quasi-Stellar Object* ou Quasar)



- ▶ Implementação
 - ▷ Python
 - ▷ Pré-processamento dos dados
 - ▷ Implementação do Sklearn como referência
 - ▷ Chamada em módulos C++
 - ▷ Avaliação

- ▶ Implementação
 - ▷ C++
 - ▷ Implementação de módulo para Random Forest
 - ▷ Modificação do módulo para paralelização com OMP

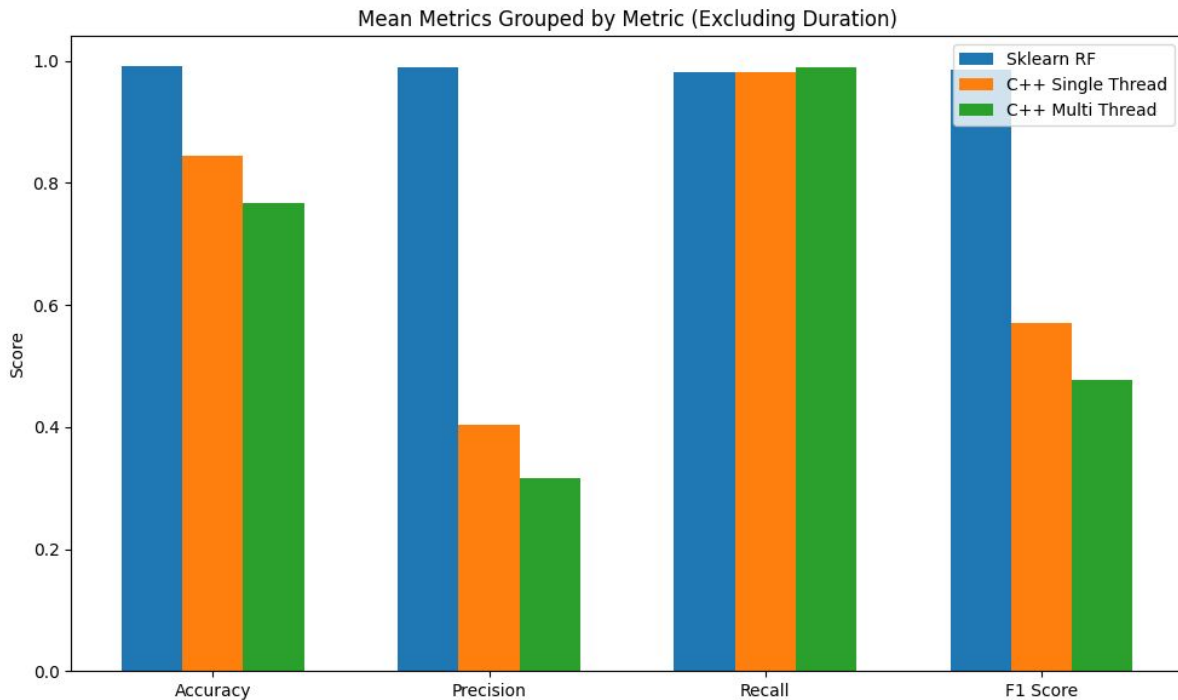


- ▶ Testes
 - ▷ Treinamento com 8000 entradas
 - ▷ Testes com 2000 entradas
 - ▷ 10x execuções -> média aritmética

Resultados

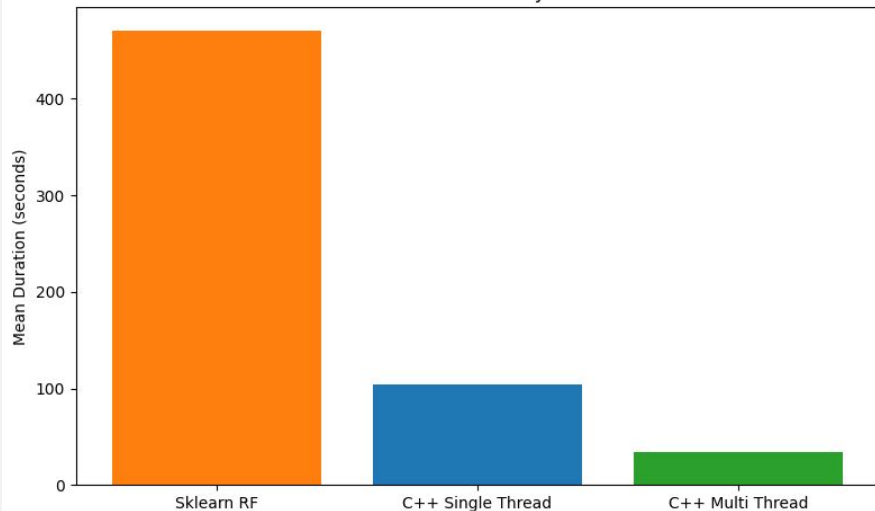


Accuracy, precision, recall, f1-score e tempo de execução

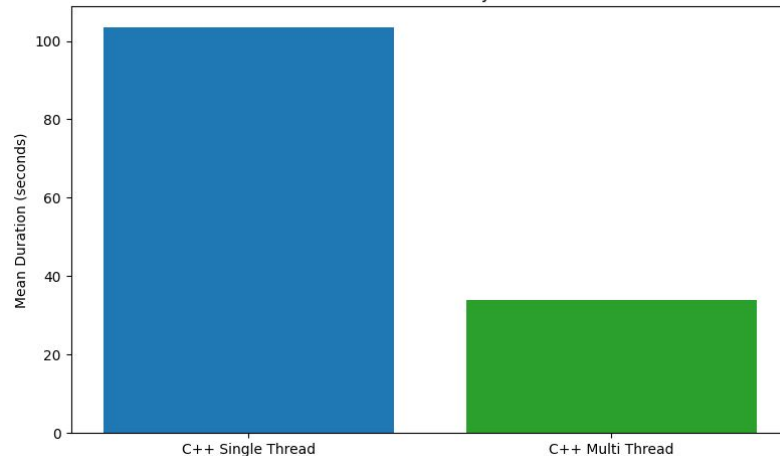


Speedup, eficiência, fração e métrica de Karp-Flatt

Mean Execution Time by Classifier



Mean Execution Time by Classifier



Speedup, eficiência, fração e métrica de Karp-Flatt



Considerações

Metodologia

- ▶ Trouxe um ganho considerável de performance no quesito tempo
- ▶ Em compensação, reduziu a acurácia e precisão
- ▶ Não necessariamente relacionado a paralelização, mas sim a implementação escolhida

- ▶ Baixa eficiência e valor baixo para métrica de Karp-Flatt indicam utilização ineficiente do paralelismo
 - ▷ Núcleos de processamento mal aproveitados
 - ▷ Overhead
- ▶ A paralelização é promissora, e, aplicada corretamente, pode trazer melhorias significativas aos algoritmos

Alguma dúvida?