

## AMBIENTE PARA APLICAÇÃO DAS TÉCNICAS DE MINERAÇÃO DE DADOS - SURVEYS DE GALÁXIAS

Thiago Crestani<sup>1</sup>; Fábio José Rodrigues Pinheiro<sup>2</sup>; Fábio Rafael Herpich<sup>3</sup>; Fernando José Braz<sup>4</sup>; Marcelo Massocco Cendron<sup>5</sup>; Vinícius Barreto Klein<sup>6</sup>; Leila Lisiane Rossi<sup>7</sup>

### INTRODUÇÃO

Como resultado da evolução tecnológica dos tempos atuais, a Astronomia passou a ser uma ciência que se utiliza de volumes de dados gigantescos, obtidos através de projetos de varredura do céu em grande escala. Para otimizar o processamento de informação obtida através deste sistema, surgiu a necessidade de uma nova forma de pesquisa na área de Astrofísica Observacional, a mineração de dados. Com isso tornou-se imprescindível o desenvolvimento de ferramentas de armazenamento, organização, integração, análise e exploração de dados, introduzindo um novo conceito de ciência observacional chamado de Observatório Virtual (VO).

Para cada *survey* são desenvolvidas as ferramentas necessárias para a obtenção da informação nele contida da forma mais eficiente. Porém, não há nenhuma forma de correlação já desenvolvida quando trata-se da comparação entre bancos distintos. Fazendo-se o sistema de *cross matching* entre dois bancos distintos, cada qual com sua função específica, pode-se obter um novo e mais amplo banco de dados, com um maior número de informações para cada objeto. E para que as pesquisas destes bancos de dados aconteçam de forma eficiente, é importante o uso das tecnologias de Banco de Dados e Mineração de Dados resultando em novos *surveys* com detalhes revelados, para o mesmo objeto, por *surveys* distintos. E ainda com a aplicação de algoritmos otimizados é possível executar este processo de *matching* com taxas razoáveis de tempo de resposta. O

---

<sup>1</sup>Aluno do Curso Ciência da Computação, 5º semestre. E-mail: thiagocrestani@gmail.com

<sup>2</sup>Professor Orientador. E-mail: fabio@ifc-videira.edu.br

<sup>3</sup>Co- Orientador. E-mail: fabiorafaelh@gmail.com

<sup>4</sup>Professor Colaborador do Projeto. E-mail: fernando.braz@ifc-araquari.edu.br

<sup>5</sup>Professor Colaborador do Projeto. E-mail: marcelo.cendron@ifc-videira.edu.br

<sup>6</sup>Professor Colaborador do Projeto. E-mail: vinius@ifc-videira.edu.br

<sup>7</sup>Professor Colaborador do Projeto. E-mail: leila.rossi@ifc-videira.edu.br

presente trabalho tem como objetivo descrever as etapas realizadas para a criação do ambiente para a aplicação das técnicas de mineração de dados de galáxias e está organizado conforme segue: O Capítulo 2 apresenta os Procedimentos Metodológicos. No Capítulo 3 são apresentados os Resultados e Discussões. O Capítulo 4 apresenta as Considerações Finais e finalmente no Capítulo 5 são apresentadas as Referências.

## PROCEDIMENTOS METODOLÓGICOS

Os dados utilizados no projeto são provenientes das bases *SDSS* e *WISE*, ambos com vários objetos astronômicos e disponibilizadas para download.

*SDSS* - O *Sloan Digital Sky Survey (SDSS)*, (YORK et.al.2000), (<http://www.sdss.org/>) é um dos maiores projetos da história da astronomia. Em 8 (oito) anos de operação (2000 – 2008), obteve imagens do céu profundo e multicoloridas de mais de um quarto do céu, criando mapas tridimensionais contendo mais de 930000 galáxias e mais de 120000 quasares. Os dados deste período estão disponíveis no *Data Release 7 (DR7)*, (<http://www.sdss.org/dr7/>). Não obstante, o projeto ainda continua com o *Third Sloan Digital Sky Survey (SDSS-III)*, (<http://www.sdss3.org/>), iniciado em Julho de 2008 e tem previsão de durar até 2014.

*WISE* - O *Wide-field Infrared Survey Explorer (WISE)*, (Wright et al. 2010, <http://adsabs.harvard.edu/abs/2010AJ....140.1868W>, <http://wise.ssl.berkeley.edu>) é uma missão astronômica observacional idealizada e executada pela NASA (<http://www.nasa.gov>). O *WISE* é um telescópio espacial que utiliza-se da parte espectral do infravermelho para imagear todo o céu.

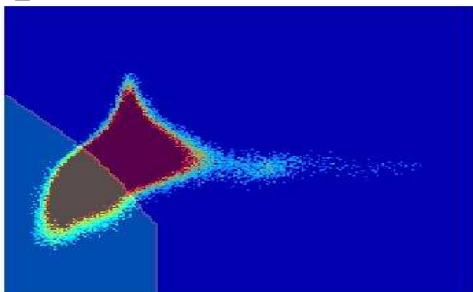
Com o grande volume de dados obtidos nas duas frentes de trabalho (os grandes bancos de dados *SDSS*, *WISE* e o Observatório Municipal Domingos Forlin), o processamento e armazenamento se tornaram limitados. Para suprir essas necessidades, foi criado um ambiente de banco de dados, no qual é possível realizar consultas otimizadas, como por exemplo, de grupos bem definidos de galáxias obtidos a partir da aplicação de técnicas de mineração de dados. As ferramentas da construção dessa infraestrutura foram totalmente desenvolvidas na instituição, proporcionando o contato da pesquisa com a prática aos alunos envolvidos no projeto.

As grandes bases de dados como a usada no projeto dificultam a descoberta de conhecimento quando são usadas as tecnologias tradicionais, pois estas não possuem a capacidade de revelar tendências ou padrões de relacionamentos entre as ocorrências de itens de dados. Uma possível solução é a aplicação da Mineração de Dados, a qual pode ser executada sobre vários tipos de dados como os bancos de dados relacionais, data warehouses, bancos de dados multimídia, bancos de dados espaciais, entre outros. Com isso é possível descobrir padrões nos dados até então desconhecidos. No presente projeto foi aplicada principalmente a técnica de agrupamento permitindo assim a geração de clusters de galáxias conforme as suas características.

## RESULTADOS E DISCUSSÕES

Os principais resultados obtidos no projeto foram a criação do ambiente para aplicação das Técnicas de Mineração de Dados em Pesquisa de *Matching* entre *Surveys* de Galáxias e a aplicação das técnicas de Mineração de Dados *Watershed* (linha divisora de águas) (Figura 1) a qual consiste em inundar a imagem, que nesse caso é um gráfico, e em seguida remover as partes inundadas, sobrando apenas as linhas com os pontos mais altos e *Simple-K-Means* (agrupamento). Além da criação de softwares para a conversão de dados e também gráficos do tipo *WHAN* e BPT que são usados para classificar galáxias.

**Figura 1** - Aplicação da técnica de Watershed em um gráfico *Whan* com os parâmetros *W1\_W4 – U\_R*



Após a obtenção das bases foi iniciado um trabalho de estudo sobre qual ferramenta de banco de dados seria a melhor para armazenar as informações. Escolheu-se o software *PostgreSQL* devido a sua robustez e bom desempenho se tratando de grandes volumes de dados. Então foi desenvolvido um software na

linguagem de programação Java capaz de realizar a conversão dos dados no formato original *Fits*, para o *SQL*. Em seguida foi feita a inserção dos dados na base, através de linha de comando.

Com os dados armazenados foi realizada uma limpeza na base, isto é, foram retirados valores que se tratavam de erros na captura, valores fora dos padrões, que deveriam ser removidos para não interferirem nos resultados finais. Então, através de comandos *SQL* foi feita a remoção desses dados. Com a base limpa efetuou-se a criação de novos campos de dados, esses campos seriam utilizados futuramente para a criação de gráficos, bem como a comparação entre valores. Eles foram gerados a partir de comando *SQL*. Os parâmetros para gerá-los saíram de campos pré-existentes, que foram subtraídos entre si, gerando as combinações.

Após a base estar devidamente pronta para a mineração era necessário a interação dela com uma linguagem de programação capaz de gerar gráficos com os pontos no banco de dados. Realizou-se uma pesquisa bibliográfica e também testes, e constatou-se que a melhor linguagem para essas tarefas seria a linguagem de programação *Python* devido a sua simplicidade a alta produtividade. Então, foi implementado um software que a partir do banco de dados gera gráficos do tipo *WHAN* e *BPT* usados para mostrar características de galáxias. Com a melhoria desse mesmo software, gráficos do tipo *WHAN* podem ser gerados com divisões de classes, permitindo a melhor visualização de diferentes tipos de galáxias.

Tornou-se necessária traçar uma linha divisória nos gráficos do tipo *WHAN* para a visualização de galáxias jovens e antigas. Para este problema foi utilizado a *Watershed*. Neste caso foi implementada uma variação dessa técnica, que consiste em remover os pontos mais altos da gráfico, sobrando apenas uma linha que é traçado pelos pontos mais baixos do gráfico. A técnica foi desenvolvida em duas etapas. Inicialmente os gráficos foram preparados, utilizando a linguagem *python*. Foram criados gráficos em forma de matrizes, que continham informações de quantos pontos apareciam em cada parte do gráfico. Em seguida com o software matemático *matlab* os gráficos foram transformados em imagens, e em seguida submetidos a uma serie de técnicas de tratamento de imagens para a implementação da técnica *Watershed*. Ao final com o próprio *matlab* a técnica foi implantada e então como resultado foi obtida a linha de separação entre galáxias jovens e antigas. (Figura 1)

## CONSIDERAÇÕES FINAIS

Com o uso da informática, a interpretação de dados astronômicos torna-se mais rápida e produz melhores resultados. Com um ambiente específico para a mineração de dados criado e ferramentas adequadas é possível produzir diversos resultados de enorme importância para a comunidade científica. A implantação de técnicas de mineração de dados permite obter várias informações e padrões presentes em grupos de galáxias. Contudo, devido à complexidade do projeto e os resultados obtidos até o momento, considera-se necessária a continuidade do mesmo através da aplicação de outras técnicas de mineração de dados, permitindo assim compará-las e escolher as melhores para a referida base de dados. Espera-se também obter uma melhor separação entre as classes de galáxias, bem como o estudo das propriedades do meio interestelar, bem como a efetiva interferência de gás e poeira sobre a luz produzida pelas estrelas componentes. Estes dados serão disponibilizados no Observatório Virtual, o qual poderá ser acessado tanto pela comunidade acadêmica quanto pelo público externo.

## REFERÊNCIAS

AGRAWAL, R.; IMIELINSKI, T.; SWAMI, A. N. Mining association rules between sets of items in large databases. Proceedings of the ACM International Conference on Management of Data, p. 207–216, 1993.

AGRAWAL, R.; SRIKANT, R. Mining sequential patterns. In: YU, P. S.; CHEN, A. S. P. (Ed.). Eleventh International Conference on Data Engineering. Taipei, Taiwan: IEEE Computer Society Press, 1995. p.3–14. Disponível em: <[citeseer.ist.psu.edu/agrawal95mining.html](http://citeseer.ist.psu.edu/agrawal95mining.html)>.

FAYYAD, U. et al. (Ed.). Advances in knowledge discovery and data mining. Menlo Park, CA, USA: American Association for Artificial Intelligence, 1996. ISBN 0-262-56097-6.7

GANTI, V.; GEHRKE, J.; RAMAKRISHNAN, R. Mining very large databases. Computer, IEEE Computer Society Press, Los Alamitos, CA, USA, v. 32, n. 8, p. 38–45, 1999. ISSN 0018-9162.

SDSS - Disponível em: <<http://www.sdss.org/>> - Visto em 30 de Abril de 2013  
Wright, E. L. et al. The Wide-field Infrared Survey Explorer (WISE): Mission Description and Initial On-orbit Performance. v. 140, p. 1868–1881, dez. 2010.