

Laboratório de Biologia Computacional e Molecular

Centro de Biotecnologia da UFRGS
Universidade Federal do Rio Grande do Sul



R para Ciências da Vida (BCM13065) Aula 6

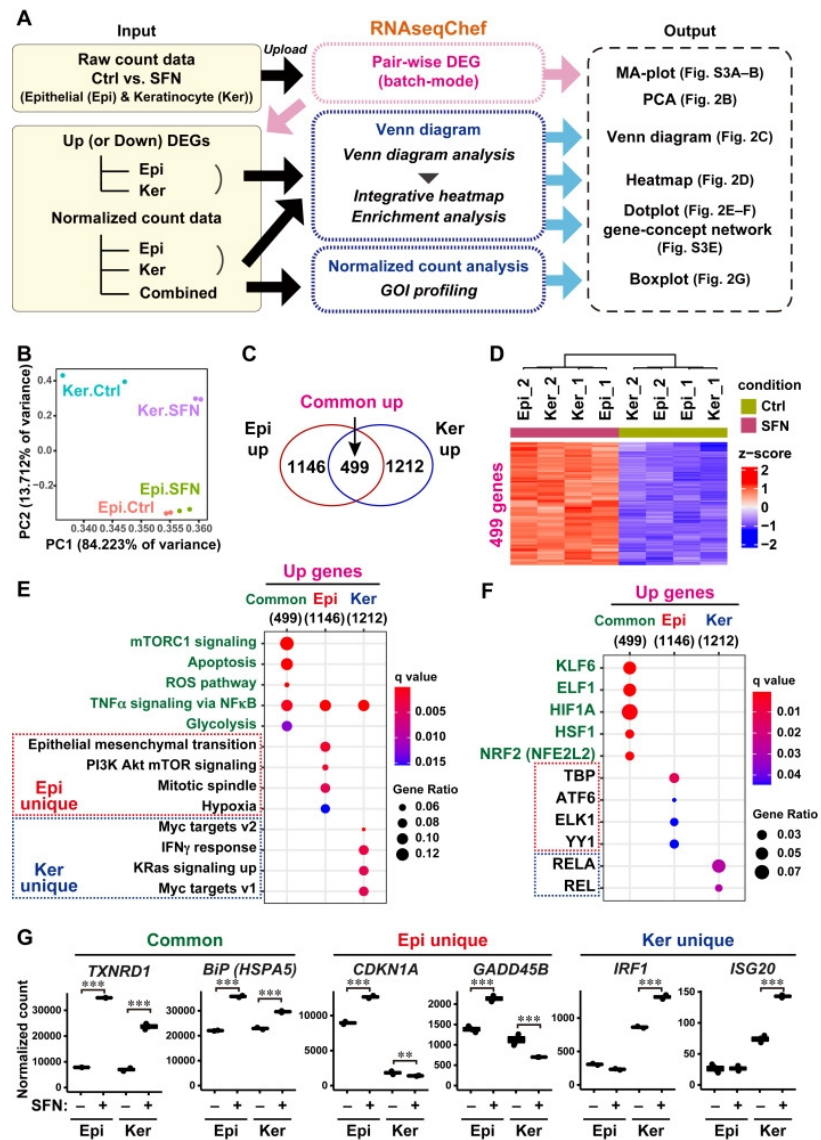
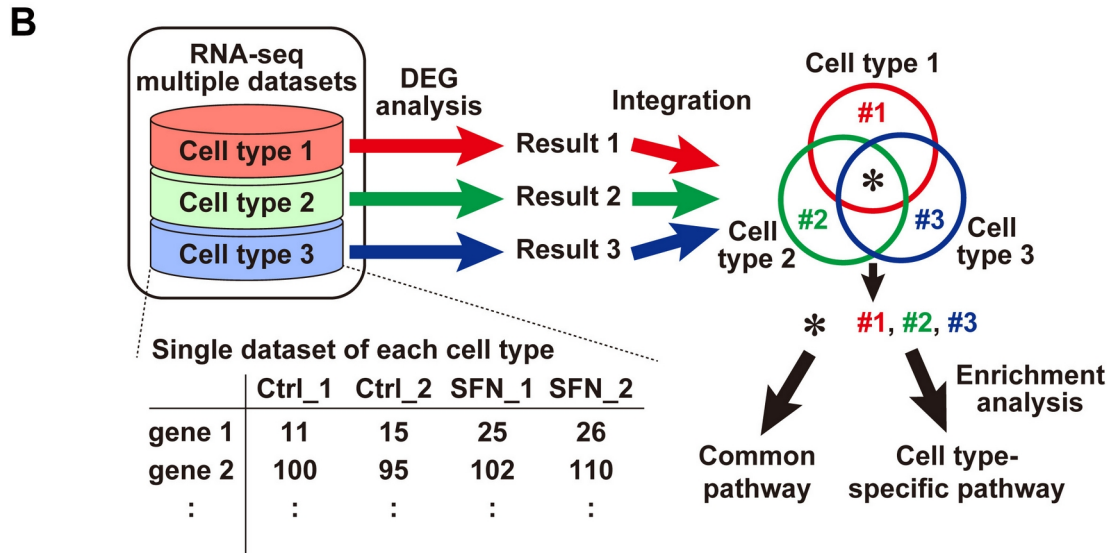
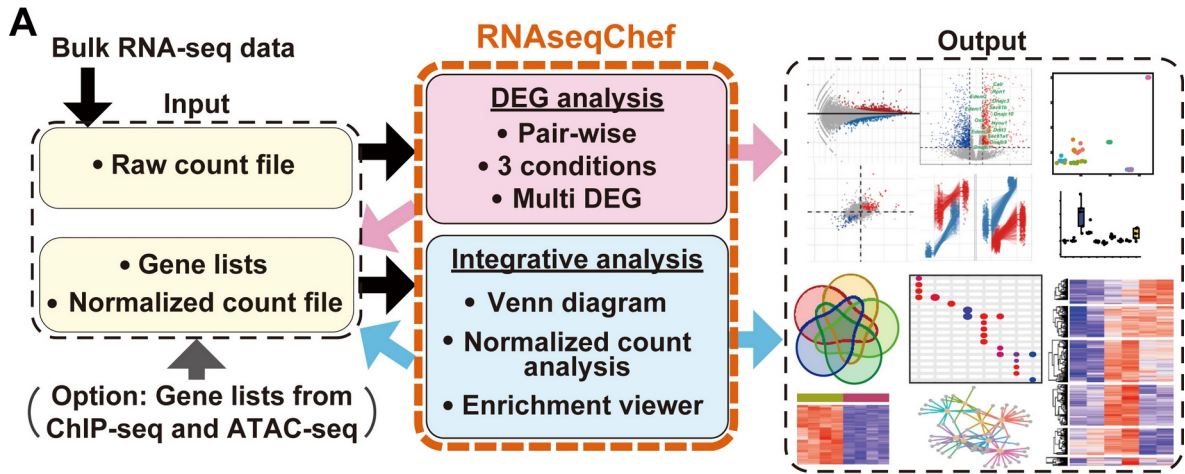
PPGBCM - UFRGS

Diego Bonatto
2024/2

Data mining versus knowledge discovery

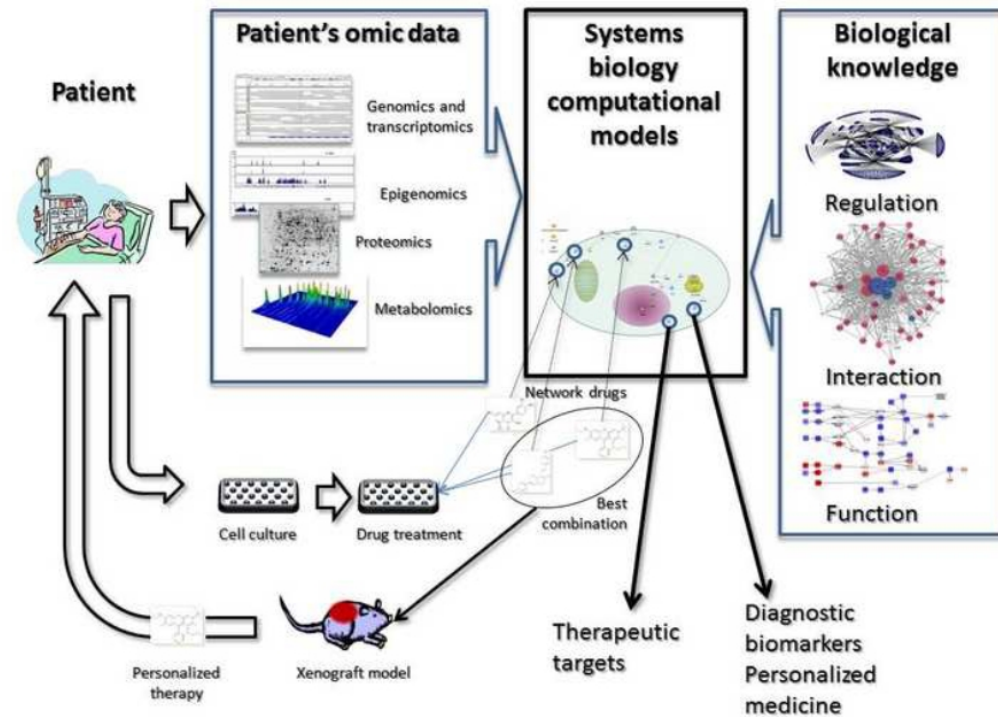
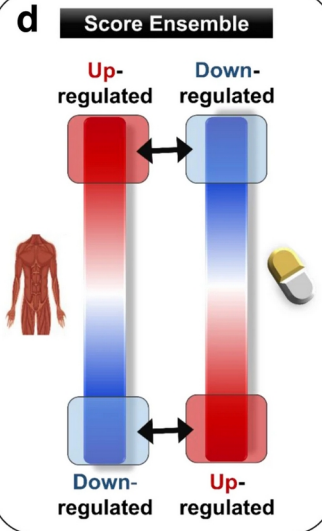
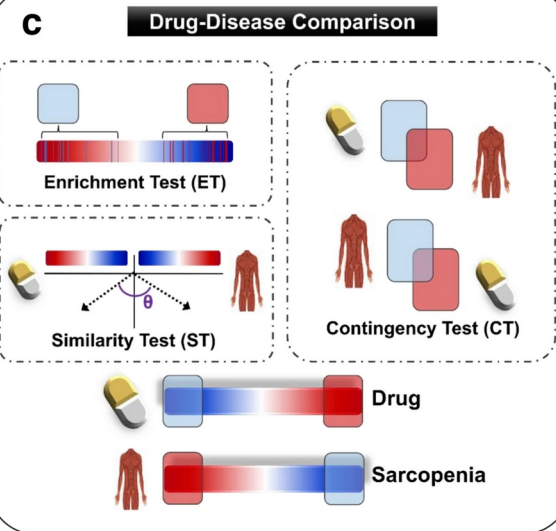
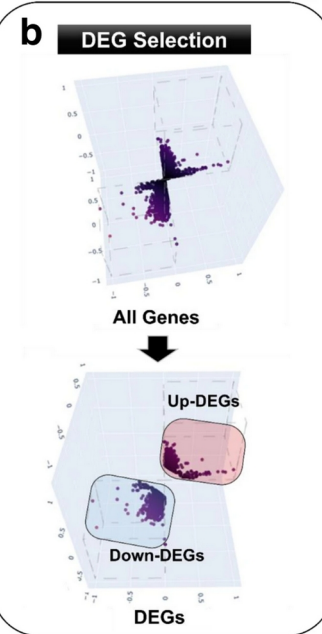
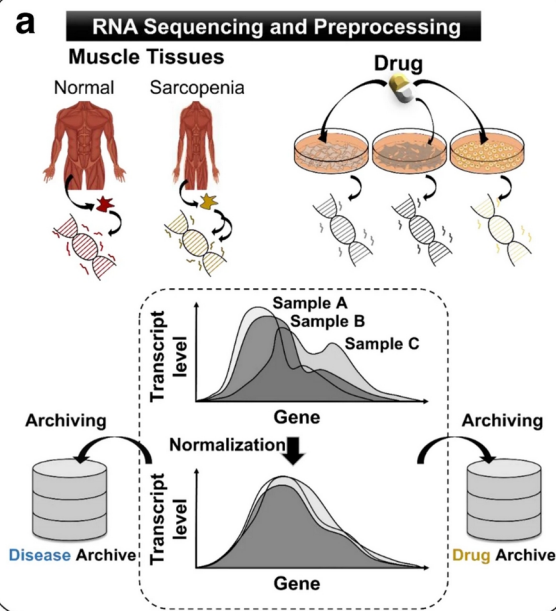
- **O que é data mining?**
 - **Foco:** Extração de padrões e relações a partir de grandes conjuntos de dados.
 - **Técnicas:** Utiliza algoritmos de aprendizado de máquina como árvores de decisão, redes neurais, clustering e associação de regras.
 - **Objetivo:** Encontrar informações relevantes e úteis nos dados.

Exemplo de data mining biológico -10.1016/j.jbc.2023.104810



Data mining versus knowledge discovery

- **O que é knowledge discovery?**
 - **Foco:** Processo mais amplo que inclui todas as etapas para transformar dados brutos em conhecimento acionável.
 - **Etapas:** Seleção de dados, pré-processamento, data mining, interpretação, visualização e geração de conhecimento.
 - **Objetivo:** Não apenas encontrar padrões, mas também entender seu significado e como utilizá-los para tomada de decisões.



10.1016/j.drudis.2013.06.003
10.1038/s12276-024-01189-z

Data mining versus knowledge discovery

- **Em resumo:**

- Data Mining é uma parte do processo de Knowledge Discovery.
- Data Mining se concentra na extração de padrões, enquanto Knowledge Discovery engloba todo o processo de transformar dados em conhecimento útil.

- **Pontos-chave:**

- Data Mining é como minerar um terreno em busca de ouro: você utiliza ferramentas para encontrar algo valioso.
- Knowledge Discovery é como transformar esse ouro em joias: você refina, molda e dá um significado ao que foi encontrado.
- Ambos são importantes: O data mining fornece os dados brutos, enquanto o knowledge discovery transforma esses dados em insights valiosos.

Característica	Data Mining	Knowledge Discovery
Foco	Identificação de padrões e relações em dados de expressão gênica.	Transformação desses padrões em conhecimento biológico útil para a prospecção de drogas.
Técnicas	Algoritmos de clustering, classificação, associação de regras.	Análise estatística, bioinformática, integração de dados de diversas fontes.
Objetivo	Descobrir biomarcadores e genes-alvo.	Descobrir novas drogas, entender mecanismos moleculares e desenvolver diagnósticos.
Exemplo	Identificar um grupo de genes que são consistentemente sobreexpressos em tumores malignos.	Investigar a função desses genes, suas interações com outras proteínas e seu potencial como alvos para novas drogas.
Resultado	Lista de genes ou grupos de genes com padrões de expressão interessantes.	Hipóteses sobre os mecanismos moleculares da doença e novas estratégias terapêuticas.

Por que usar R para data mining e knowledge discovery biológicos?

- O R oferece uma vasta gama de pacotes para realizar diversas tarefas de data mining e knowledge discovery.
- A escolha do pacote ideal dependerá da natureza dos dados, do tipo de análise que se deseja realizar e dos objetivos propostos.

Pacotes do R

Integração de pacotes do R para data mining e knowledge discovery para dados biológicos

- A integração de pacotes R é fundamental para realizar análises complexas e abrangentes em dados biológicos.
- Ao combinar diferentes pacotes, pode-se construir pipelines de análise robustas e eficientes, desde a limpeza dos dados até a interpretação dos resultados.

Tipos de integrações de pacotes do R

- **Sequencial:**

- Descrição: A execução de um pacote após outro, com a saída de um sendo a entrada do próximo.
- Exemplo: Limpeza dos dados com o pacote tidyverse, seguido da análise de expressão gênica com o pacote DESeq2 e, finalmente, a visualização dos resultados com ggplot2.

- **Em Paralelo:**

- Descrição: Múltiplos pacotes são executados simultaneamente em diferentes partes dos dados.
- Exemplo: Utilizar o pacote foreach para paralelizar cálculos intensivos, como a análise de diferentes conjuntos de genes.

Tipos de integrações de pacotes do R

- **Dentro de Funções Personalizadas:**

- Descrição: Criação de funções que combinam a funcionalidade de diversos pacotes.
- Exemplo: Criar uma função que realize a normalização de dados, a seleção de genes e a análise de enriquecimento genético, utilizando pacotes como limma, DESeq2 e clusterProfiler.

- **Utilizando Pipelines:**

- Descrição: Construção de pipelines de análise que definem uma sequência de passos a serem executados, com a possibilidade de parametrização e reutilização.
- Exemplo: Utilizar o pacote magrittr para criar pipelines concisos e legíveis, ou o pacote drake para gerenciar dependências e reprodutibilidade.

Exemplos de Integração de Pacotes

RNA-seq



Redes de interação

