

Here is a detailed review of the R code in `bladder.Rmd`, focusing on bugs, performance, and opportunities for AI tool integration:

1. Bugs and Code Issues

a. Installation Code Within Chunks

- **Problem:** Repeated installation of packages inside the RMarkdown file (`install.packages` and `BiocManager::install`).
- **Impact:** This can significantly slow down the run and lead to issues if the Rmd is re-executed multiple times. Also, package installations should not be in analysis scripts.
- **Fix:** Move installation instructions to a separate setup script or instruct users to install before running the notebook.

b. Mixed Use of Data Objects (`lihc_res` VS `tcga_data`)

- **Problem:** In Step 1, the code checks for `lihc_res` to determine whether to re-download data, but later mainly uses `tcga_data`. If the workspace is loaded, `tcga_data` may not be present.
- **Impact:** Potential error if `tcga_data` is missing after loading workspace.
- **Fix:** Ensure that both objects are saved/loaded, or check for `tcga_data` directly.

c. Use of Deprecated or Non-Standard Functions

- **Problem:** Use of `values(tcga_data)`
`\(gene_name` (line 136) - `values()` is not a standard accessor for `SummarizedExperiment`. Should use `rowData(tcga_data)``
`gene_name` .
- **Impact:** May error on some systems or package versions.
- **Fix:** Use `rowData(tcga_data)$gene_name` .

d. Hard-coded File Paths

- **Problem:** `setwd('C:/Adam/R/')` is hard-coded.
- **Impact:** Will fail on other systems.
- **Fix:** Use relative paths or project root detection (e.g., with `here::here()`).

e. Inconsistent Factor Levels

- **Problem:** Mangling spaces in the definition variable to underscores (line 239), then using `"Solid_Tissue_Normal"` as the reference, but the data may not match this exactly if column names or values change.
- **Impact:** Factor level mismatch can cause errors.
- **Fix:** Consider using factors as-is, or ensure all downstream code uses the same formatting.

f. Potential Memory Issues and Inefficient Data Handling

- Creating a full adjacency or TOM for large datasets without chunking can cause crashes (`blockwiseModules` mitigates this, but commented-out code could be dangerous if uncommented).

g. Deprecated/Conflicting Packages

- **Problem:** `gProfileR` is deprecated; use `gprofiler2` instead. Also, both are loaded.
- **Impact:** Namespace conflicts, unexpected results.
- **Fix:** Remove deprecated package loads.

h. Plotting and Output

- Some plotting functions (e.g., `TCGAVisualize_volcano`) are used without checks for package loading or existence, and may error if not available.

2. Performance Bottlenecks and Speed Improvements

a. Parallelization

- **Current:** Only 2 threads are used in WGCNA (`enableWGCNAThreads(2)`).
- **Improvement:** Increase to the physical core count (e.g., 16 as in the comment) for significant speedup. Consider using `BiocParallel` for parallel operations in other steps.

b. Data Caching

- **Current:** Only basic workspace caching is implemented.

- **Improvement:** Save intermediate outputs (e.g., after GDCdownload, after normalization) to avoid repeating expensive steps.

c. Efficient Filtering

- **Current:** Filtering genes with `filterByExpr` is good, but additional filtering (by logFC, p-value) could further reduce computation before WGCNA.
- **Improvement:** Apply stricter pre-filtering to reduce dimensionality.

d. Use of Fast I/O and Data.table

- Loading and transforming large tables (e.g., clinical data, gene matrices) could benefit from `data.table` for faster read/write and manipulation.

e. Avoid Redundant Calculations

- Avoid recomputing objects (e.g., making new factor columns multiple times).

f. Plot Generation

- For very large heatmaps or PCA plots, consider plotting only top variable genes for visualization.

3. AI Tools to Speed Up and Optimize Code

a. Code Profiling and Optimization

- **Posit Cloud (formerly RStudio Cloud):** Allows profiling and optimization of R code interactively.
- **AI-Powered Linter Plugins:** Use Copilot for R, or Quarto AI assistant for suggestions on function vectorization and bottlenecks.

b. Parallelization/Distributed Computing

- **AI/ML-based Parallelization:** Use the `future` and `furrr` packages, which leverage AI-based scheduling and chunking for parallel tasks.

- **Google Vertex AI Notebooks / AWS SageMaker Studio Lab:** For running large analyses on cloud GPUs/CPUs.

c. Data Preprocessing and Cleaning

- **AI-assisted Data Cleaning:** Use `DataRobot` or `Trifacta` to automatically suggest or apply data cleaning and filtering steps.

d. Efficient Visualization

- **Plotly's AI Visualization Tools:** For interactive, faster plotting of large datasets.

e. Automated Hyperparameter Tuning

- **AutoML Packages** (for downstream ML): `caret` and `mlr` can leverage grid/random search with parallel backends.
-

4. Additional Suggestions

- **Use a Project Structure:** Adopt `renv` or `packrat` for reproducible environments.
 - **Notebook Speed:** Avoid plotting all genes or all modules at once; let the user choose subsets.
 - **Documentation:** Add comments on expected run time and memory requirements.
 - **Error Handling:** Add `tryCatch` blocks around expensive or failure-prone operations (e.g., downloading, WGCNA).
 - **Session Info:** Print session info for reproducibility.
-

Summary Table

Area	Issue/Bug/Opportunity	Suggestion/Tool
Package Installs	In Rmd, slows re-runs	Move to setup script
Data Loading	Mixed object names, missing checks	Check/load <code>tcga_data</code> directly
Parallelism	Only 2 threads in WGCNA	Use all physical cores
Deprecated	Use of old packages	Remove <code>gProfileR</code> , use <code>gprofiler2</code>
Hard Paths	<code>setwd('C:/Adam/R/')</code>	Use <code>here::here()</code> , relative paths
Visualization	Plotting all genes	Plot top N genes only
AI Tools	None used	Copilot for R, Posit Cloud, DataRobot

If you'd like, I can rewrite a specific section of the code for improved performance, or provide example code for parallelization or AI integration!