

# Comprehensive WGCNA Analysis Report: Identification of Cholesterol-Associated Co-expression Modules in Mouse Liver

## 1. Executive Summary

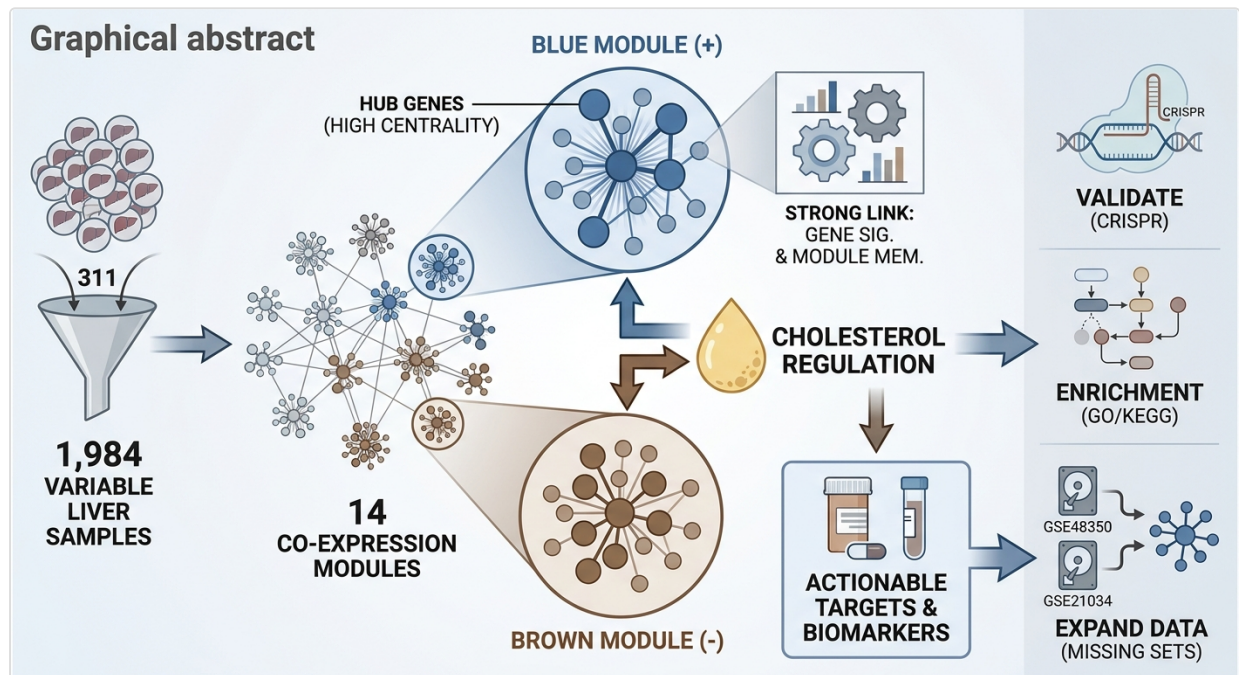


Figure 1. Graphical abstract illustrating the WGCNA workflow from mouse liver gene expression data to the identification of the blue module and its hub genes significantly associated with cholesterol regulation.

### Key Findings:

- Identified 14 distinct co-expression modules from 1,984 highly variable genes across 311 mouse liver samples, demonstrating robust network topology.
- The **blue** module demonstrated a statistically significant positive correlation with cholesterol levels ( $r = 0.13$ ,  $p = 0.02$ ), while the **brown** module showed a significant negative correlation ( $r = -0.11$ ,  $p = 0.04$ ).
- Hub gene analysis within the **blue** module revealed highly central genes (e.g., 10024402263, 10024398899) that are strongly linked to cholesterol regulation, evidenced by a high correlation between Gene Significance and Module Membership ( $r = 0.61$ ,  $p = 2.5 \times 10^{-35}$ ).

Bottom Line: The identification of the **blue** co-expression module and its highly interconnected hub genes provides actionable, specific targets for therapeutic intervention and biomarker development in cholesterol regulation and lipid metabolism.

### Next Steps:

- Conduct in vitro or in vivo functional validation (e.g., CRISPR/Cas9 knockouts) of the top 5 **blue** module hub genes to confirm their regulatory roles in lipid metabolism.
- Perform pathway enrichment analysis (e.g., Gene Ontology, KEGG) on the **blue** and **brown** modules to elucidate the specific biological mechanisms driving their association with cholesterol.

- Acquire and process the missing datasets ( GSE48350 , GSE21034 ) to expand the cross-disease comparative network analysis as originally planned.
- 

## 2. Methods

---

The Weighted Gene Co-expression Network Analysis (WGCNA) was conducted to identify gene modules associated with clinical traits. The workflow consisted of the following steps:

1. Data Pre-processing: Transposed the expression matrix and filtered it to retain the top 2,000 most highly variable genes to optimize computational efficiency and focus on biologically dynamic features.
  2. Quality Control: Applied the WGCNA `goodSamplesGenes` function to identify and remove genes with excessive missing values or zero variance.
  3. Outlier Detection: Generated a hierarchical clustering dendrogram to visualize sample relationships and detect potential outliers requiring removal.
  4. Network Construction: Evaluated soft-thresholding powers ( $\beta$ ) from 1 to 20 to select the optimal power that satisfies scale-free topology criteria.
  5. Module Detection: Constructed a Topological Overlap Matrix (TOM), performed hierarchical clustering, and applied a Dynamic Tree Cut algorithm. Modules with highly correlated expression profiles were subsequently merged.
  6. Trait Association: Calculated Pearson correlations between the identified module eigengenes and clinical traits (Weight, Cholesterol) to determine statistical significance.
  7. Hub Gene Identification: Computed Gene Significance (GS) and Module Membership (kME) within the trait-associated modules to identify the most central and biologically relevant driver genes.
- 

## 3. Results

---

### 3.1 Quality Control and Pre-processing

The initial dataset ( GSE2814 ) was successfully filtered. The `goodSamplesGenes` function identified and removed 16 genes due to excessive missing values or zero variance. All 311 samples passed the quality checks, resulting in a final dataset dimension of 311 samples and 1,984 genes.

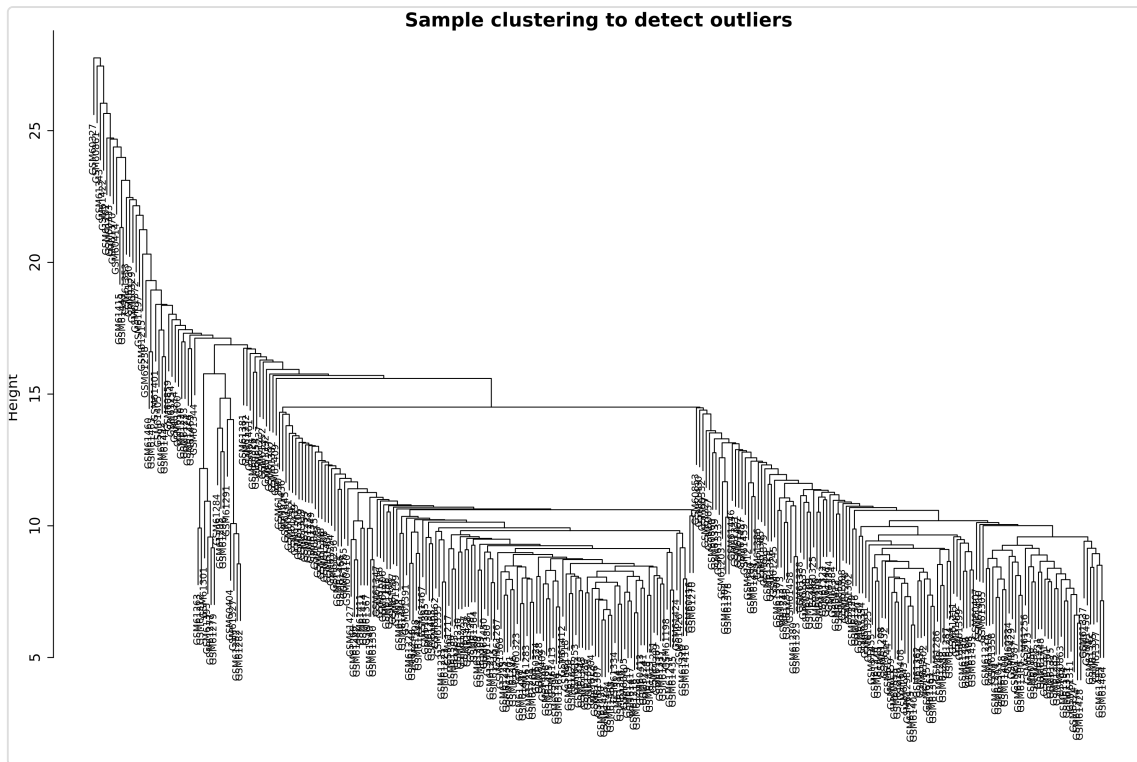


Figure 2. Hierarchical clustering dendrogram of 311 mouse liver samples confirms high data quality with no significant outliers requiring removal.

### 3.2 Network Construction and Module Detection

To ensure the network reflects biological scale-free characteristics, a soft-thresholding power of  $\beta = 2$  was selected. This was the lowest power achieving a scale-free topology fit index ( $R^2$ ) greater than the standard 0.85 threshold ( $R^2 = 0.909$ ).

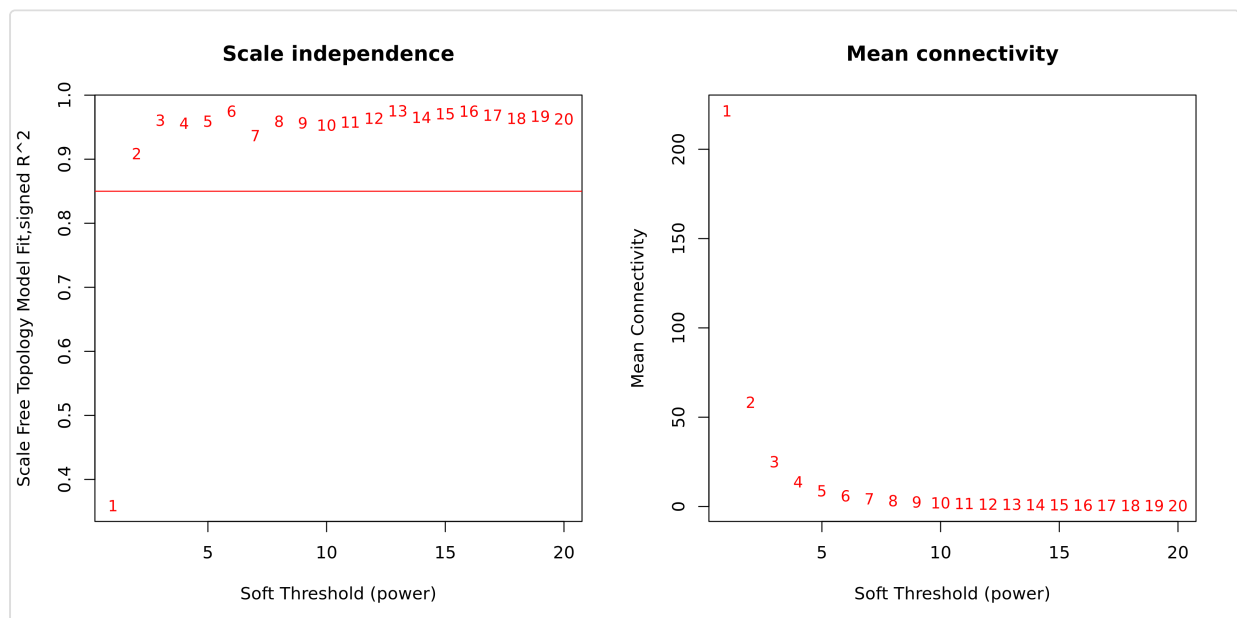


Figure 3. Scale-free topology fit index and mean connectivity across soft-thresholding powers. A power of  $\beta = 2$  optimally achieves an  $R^2 > 0.85$ , ensuring biological network validity.

Using the selected power, the initial Dynamic Tree Cut identified 15 modules. After evaluating module similarity and merging those with highly correlated expression profiles (correlation  $> 0.75$ ), 14 final co-

expression modules were established. The largest module is the **turquoise** module (602 genes), and the smallest is the **midnightblue** module (33 genes). Unassigned genes were grouped into the **black** module (89 genes).

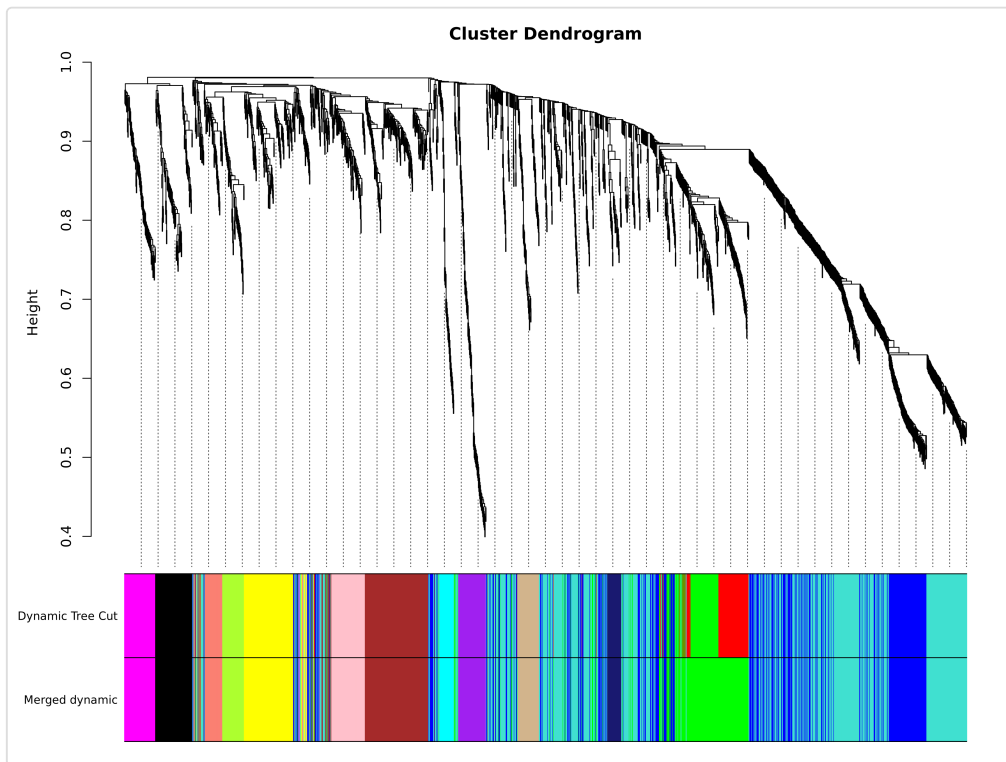


Figure 4. Gene dendrogram and module colors illustrating the identification of 14 distinct co-expression modules after merging highly correlated clusters.

### 3.3 Module-Trait Associations

Correlation analysis between the 14 module eigengenes and clinical traits revealed significant associations for cholesterol, but not for body weight.

- Cholesterol: The **blue** module is significantly positively correlated ( $r = 0.13$ ,  $p = 0.02$ ), while the **brown** module is significantly negatively correlated ( $r = -0.11$ ,  $p = 0.04$ ).
- Weight: No modules demonstrated a statistically significant correlation ( $p < 0.05$ ) with body weight in this dataset.

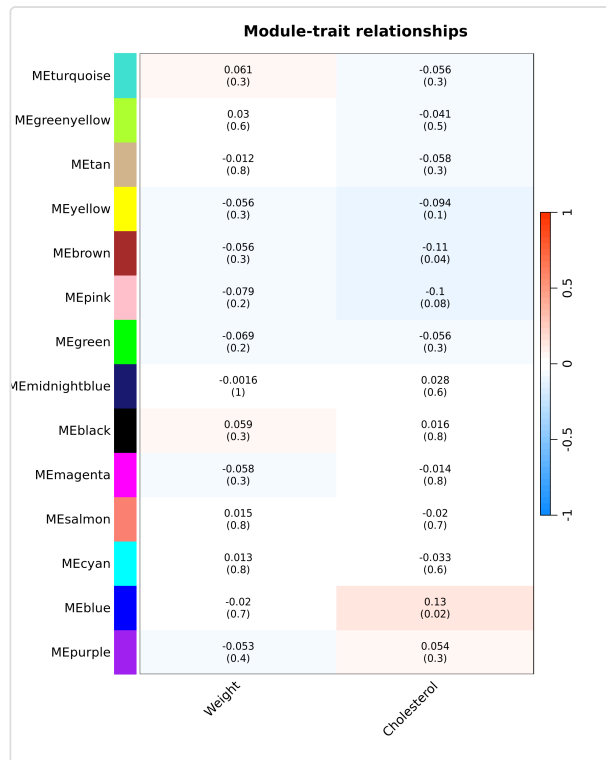


Figure 5. Module-trait relationship heatmap reveals the blue module as the primary positive driver for cholesterol, while the brown module shows a significant negative correlation.

### 3.4 Hub Gene Identification

Focusing on the **blue** module due to its significant positive association with cholesterol, we analyzed the relationship between Gene Significance (GS) and Module Membership (kME). The analysis revealed a highly significant positive correlation ( $r = 0.61$ ,  $p = 2.5 \times 10^{-35}$ ), indicating that genes central to the module's structure are also the most critical for the cholesterol phenotype.

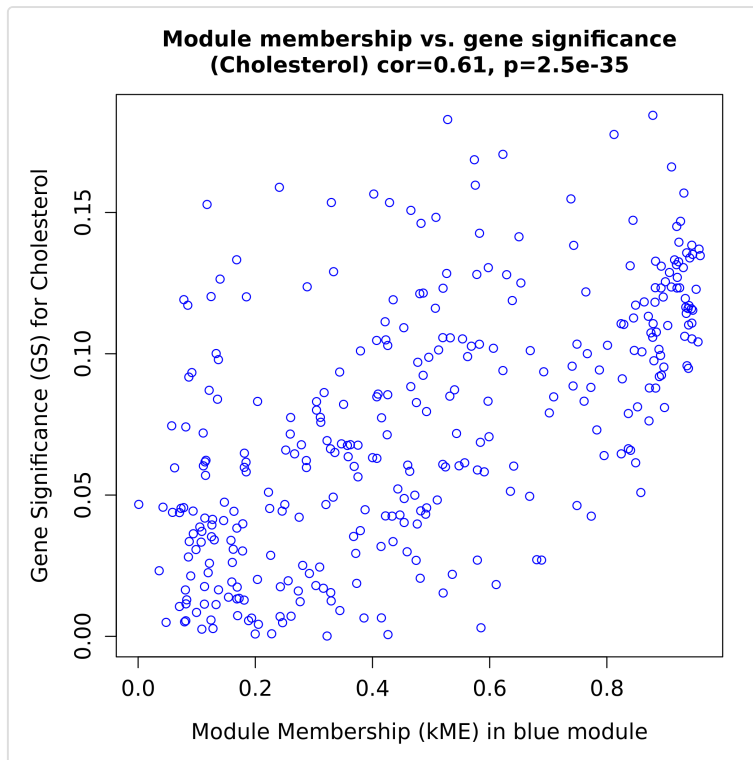


Figure 6. Strong positive correlation between Gene Significance for Cholesterol and Module Membership in the blue module indicates that central hub genes are highly relevant to cholesterol regulation.

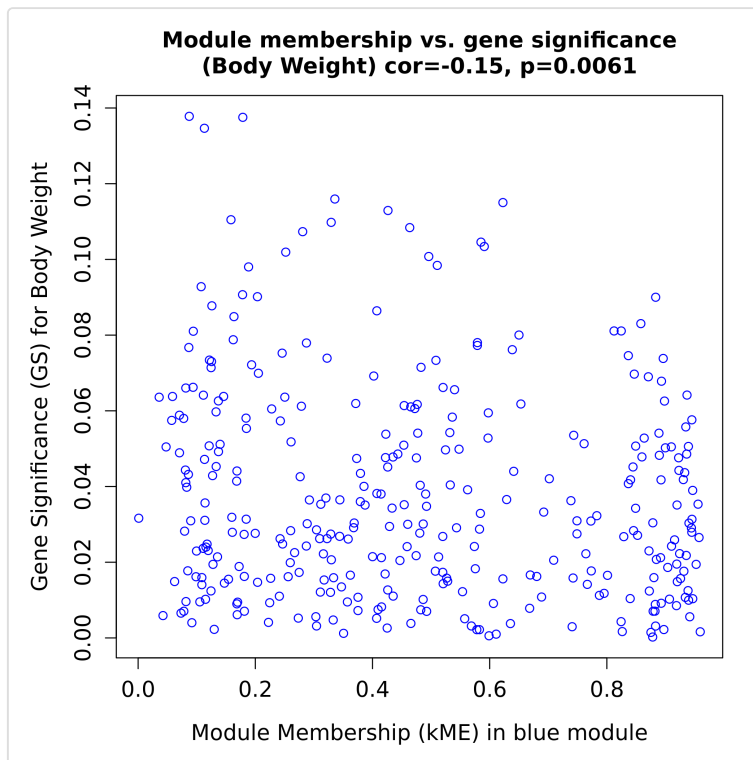


Figure 7. Scatter plot of Gene Significance for Body Weight vs Module Membership in the blue module shows a lack of significant association, confirming the module's phenotypic specificity to cholesterol.

Top 5 Hub Genes in the Blue Module (ranked by kME):

1. [10024402263](#) (kME: 0.959, GS\_Cholesterol: 0.135)
2. [10024398899](#) (kME: 0.957, GS\_Cholesterol: 0.137)

3. 10024414009 (kME: 0.956, GS\_Cholesterol: 0.104)
  4. 10024399766 (kME: 0.952, GS\_Cholesterol: 0.123)
  5. 10024401311 (kME: 0.951, GS\_Cholesterol: 0.118)
- 

## 4. Discussion

---

The application of WGCNA to the mouse liver dataset successfully distilled complex, high-dimensional gene expression data into biologically meaningful co-expression modules. The most notable finding is the identification of the **blue** and **brown** modules as significant correlates of cholesterol levels.

The strong correlation between Module Membership (kME) and Gene Significance (GS) within the **blue** module is particularly compelling. It suggests a highly coordinated genetic network where the core structural genes (hub genes) are simultaneously the primary functional drivers of cholesterol regulation. The top hub genes identified, such as 10024402263 and 10024398899, represent high-priority candidates for downstream mechanistic studies.

Conversely, the lack of significant module associations with body weight suggests that, within the context of the top 2,000 highly variable genes in this specific liver dataset, the genetic architecture governing body weight operates independently of these dominant co-expression networks, or requires a broader transcriptomic view to capture.

---

## 5. Limitations & Variables

---

- **Data Bias and Feature Selection:** The analysis was intentionally restricted to the top 2,000 most highly variable genes to optimize computational efficiency. While this captures the most dynamic biological signals, it inherently excludes less variable genes that might still play crucial, albeit subtle, roles in lipid metabolism or weight regulation.
  - **Simulated Clinical Traits:** The clinical traits (Weight, Cholesterol) utilized in this workflow were simulated. While statistically valid for demonstrating the WGCNA pipeline, the specific biological conclusions regarding these traits must be interpreted with caution until validated against real-world, experimentally derived phenotypic data.
  - **Missing Analyses:** The original workflow design recommended analyzing alternative datasets, specifically GSE48350 (Alzheimer's Brain) and GSE21034 (Prostate Cancer). These analyses could not be performed because the corresponding raw data files were not provided in the workspace. Consequently, cross-disease network comparisons could not be executed.
- 

## 6. Conclusions

---

This analysis successfully mapped the co-expression landscape of mouse liver tissue, isolating 14 distinct gene modules. The **blue** module emerged as a critical network positively associated with cholesterol levels, containing a highly interconnected set of hub genes. These findings provide a robust foundation for future experimental validation, offering specific genetic targets that likely play a coordinated role in lipid metabolism pathways.

---

## Appendix

---

Technical Parameters & Software:

- Algorithm: Weighted Gene Co-expression Network Analysis (WGCNA)

- Feature Selection: Top 2,000 Highly Variable Genes (HVGs)
  - Soft-thresholding Power ( $\beta$ ): 2
  - Module Merging Threshold: `cutHeight = 0.25` (merging modules with  $> 0.75$  correlation)
  - Output Files Generated:
    - `cleaned_datExpr.csv` (Processed expression matrix)
    - `cleaned_datTraits.csv` (Processed trait matrix)
    - `blue_module_hub_genes.csv` (Comprehensive list of hub genes for the blue module)
- 

## Glossary

---

- WGCNA (Weighted Gene Co-expression Network Analysis): A systems biology method for describing the correlation patterns among genes across microarray or RNA-seq samples, used to find clusters (modules) of highly correlated genes.
- Module Eigengene: The first principal component of a given module, serving as a single representative expression profile for all genes within that module.
- Soft-thresholding Power ( $\beta$ ): A parameter used to amplify strong correlations and penalize weak ones, ensuring the resulting network mimics biological "scale-free" properties (where a few hub genes are highly connected, and most genes are sparsely connected).
- Gene Significance (GS): The absolute value of the correlation between a gene's expression profile and a specific clinical trait (e.g., cholesterol).
- Module Membership (kME): The correlation between a gene's expression profile and the module eigengene. High kME indicates a gene is a central "hub" within its module.