

# CommE5064: Artificial Intelligence and Intelligent Medicine

## Homework 2: Gene Network Feature Selection and Biological Enrichment

Department of Electrical Engineering

Student ID: **R13921031**

November 12, 2025

---

### Q1. Genetic Data (30%)

- (a) (5%) Describe the differences between gene expression and Single Nucleotide Polymorphism (SNP) data.
- (b) (20%) Next-generation sequencing (NGS) is another promising way to obtain gene expression data other than microarray. Describe and illustrate at least one of the realizations of NGS.
- (c) (5%) List at least one improvement of NGS in sequencing (compared to microarray) and elaborate on it.

#### (a) Gene expression vs. SNP data

What they measure:

- **Gene expression:** Quantifies RNA transcript abundance (e.g., mRNA, lncRNA), representing the activity of genes. Data are continuous values such as  $\log_2$  intensities, counts, or TPM/FPKM.
- **SNPs:** Represent single-base DNA variations at specific genomic loci; categorical or discrete genotypes (AA, AB, BB).

Biological substrate and stability:

- Expression is derived from **RNA** and reflects dynamic, tissue-specific, and time-dependent cellular activity.
- SNPs are derived from **DNA** and remain stable across time and tissues for an individual.

Assay and preprocessing:

- Expression: measured by microarray or RNA-seq; requires background correction and normalization (e.g., RMA, DESeq2).
- SNPs: measured by genotyping arrays or whole-genome/exome sequencing; involves genotype calling and quality control.

**Downstream uses:**

- Gene expression: differential expression, pathway enrichment, biomarker discovery.
- SNPs: genome-wide association studies (GWAS), heritability mapping, and eQTL analyses.

**Data types:** Gene expression  $\Rightarrow$  real-valued vectors; SNPs  $\Rightarrow$  categorical or dosage values.

**(b) Realization of Next-Generation Sequencing (NGS): Bulk RNA-seq**

**Wet-lab workflow:**

1. RNA extraction and quality control (e.g., RIN score).
2. mRNA enrichment (poly(A) selection) or rRNA depletion.
3. Fragmentation of RNA or cDNA.
4. cDNA synthesis (first and second strand).
5. Adapter ligation and indexing.
6. PCR amplification of adapter-ligated fragments.
7. Sequencing-by-synthesis on Illumina or similar platforms.

**Computational workflow:**

1. Quality control (FastQC) and trimming.
2. Alignment to reference genome (e.g., STAR, HISAT2) or pseudo-alignment (salmon, kallisto).
3. Quantification at gene or transcript level.
4. Normalization (e.g., TPM, DESeq2 normalization).
5. Downstream analyses: differential expression, clustering, biological enrichment.

### Simplified pipeline diagram:

Sample → RNA extraction → cDNA library prep → Sequencing → QC & alignment  
→ Quantification → Normalization → DE/Enrichment

RNA-seq provides high-throughput, base-level quantification of the transcriptome, enabling analysis of both known and novel genes.

### (c) Improvements of NGS over Microarray

#### Dynamic range and specificity:

- NGS detects transcripts across a much broader dynamic range than microarrays.
- Avoids probe cross-hybridization by directly sequencing cDNA.

#### Novel discovery:

- RNA-seq can identify novel isoforms, splice variants, fusion genes, and allele-specific expression, which microarrays cannot.

**Summary:** NGS surpasses microarrays in sensitivity, quantification accuracy, and capability for novel discovery, though it requires higher computational and storage resources.

## References

### References

- [1] Z. Wang, M. Gerstein, and M. Snyder, “RNA-Seq: a revolutionary tool for transcriptomics,” *Nature Reviews Genetics*, vol. 10, no. 1, pp. 57–63, 2009.
- [2] A. Mortazavi, B. A. Williams, K. McCue, L. Schaeffer, and B. Wold, “Mapping and quantifying mammalian transcriptomes by RNA-Seq,” *Nature Methods*, vol. 5, no. 7, pp. 621–628, 2008.
- [3] A. Conesa *et al.*, “A survey of best practices for RNA-seq data analysis,” *Genome Biology*, vol. 17, no. 13, pp. 1–19, 2016.
- [4] J. C. Marioni, C. E. Mason, S. M. Mane, M. Stephens, and Y. Gilad, “RNA-seq: an assessment of technical reproducibility and comparison with gene expression arrays,” *Genome Research*, vol. 18, no. 9, pp. 1509–1517, 2008.

- [5] The International HapMap Consortium, “The International HapMap Project,” *Nature*, vol. 426, no. 6968, pp. 789–796, 2003.
- [6] M. Schena, D. Shalon, R. W. Davis, and P. O. Brown, “Quantitative monitoring of gene expression patterns with a complementary DNA microarray,” *Science*, vol. 270, no. 5235, pp. 467–470, 1995.

# **2024 Fall Artificial Intelligence and Intelligent Medicine**

Homework 2: Question 2

Dynamic Interaction Network and Relevance Value

James Christian

Student ID: R13921031

Department of Electrical Engineering

National Taiwan University

November 12, 2025

## Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
1.1	Dataset Overview . . . . .	4
<b>2</b>	<b>Q2(a): Matrix Elements</b>	<b>4</b>
2.1	Problem Statement . . . . .	4
2.2	Solution . . . . .	5
2.2.1	Matrix X (Expression Matrix) . . . . .	5
2.2.2	Matrix A (Interaction Ability Matrix) . . . . .	5
2.2.3	Matrix E (Noise Matrix) . . . . .	5
<b>3</b>	<b>Q2(b): Solving for <math>A^+</math> and <math>A^-</math></b>	<b>6</b>
3.1	Problem Statement . . . . .	6
3.2	Methodology . . . . .	6
3.2.1	Ridge Regression Approach . . . . .	6
3.2.2	Why Ridge Regression? . . . . .	6
3.3	Implementation . . . . .	7
3.4	Results . . . . .	7
3.4.1	Matrix $A^+$ (Positive Group) . . . . .	7
3.4.2	Matrix $A^-$ (Negative Group) . . . . .	8
<b>4</b>	<b>Q2(c): Relevance Values and Top-15 Genes</b>	<b>8</b>
4.1	Problem Statement . . . . .	8
4.2	Methodology . . . . .	8
4.2.1	Definition of Relevance Value . . . . .	8
4.2.2	Biological Interpretation . . . . .	9
4.3	Implementation . . . . .	9
4.4	Results . . . . .	9
4.4.1	Difference Matrix Statistics . . . . .	9
4.4.2	Relevance Value Statistics . . . . .	9
4.4.3	Top 15 Genes with Highest Relevance Values . . . . .	10
4.4.4	Key Findings . . . . .	10
<b>5</b>	<b>Visualization and Analysis</b>	<b>10</b>
5.1	RV Distribution Across All Genes . . . . .	10
5.2	Top 15 Genes Visualization . . . . .	11
5.3	Interaction Heatmaps . . . . .	12
5.4	Cumulative RV Analysis . . . . .	12
5.5	Matrix Comparison . . . . .	13
<b>6</b>	<b>Discussion</b>	<b>14</b>
6.1	Biological Significance . . . . .	14
6.1.1	Network-Based Perspective . . . . .	14
6.1.2	Topological Biomarkers . . . . .	15
6.1.3	Prognostic Power . . . . .	15
6.2	Methodological Considerations . . . . .	15
6.2.1	Ridge Regression Justification . . . . .	15
6.2.2	Limitations . . . . .	15

6.3	Connection to Original Paper . . . . .	16
<b>7</b>	<b>Conclusion</b>	<b>16</b>
7.1	Key Achievements . . . . .	16
7.2	Biological Insights . . . . .	16
7.3	Future Directions . . . . .	16
<b>A</b>	<b>Code Implementation</b>	<b>17</b>
<b>B</b>	<b>Generated Files</b>	<b>17</b>

# 1 Introduction

This report presents the solution to Question 2 of AIIM Homework 2, which focuses on constructing dynamic interaction networks and calculating relevance values for gene feature selection. The methodology is based on the paper “Integrating ensemble systems biology feature selection and bimodal deep neural network for breast cancer prognosis prediction” by Cheng et al. (2021) [1].

The goal is to:

1. Understand the matrix formulation of gene interaction networks
2. Solve for interaction ability matrices using Ridge regression
3. Identify prognostic biomarkers through relevance value analysis

## 1.1 Dataset Overview

The dataset `AIIM_HW2_ANOVA_grouped_CADM1.npz` contains:

- $\mathbf{X}^+$ : Expression data for 329 genes across 27 patients (positive group)
- $\mathbf{X}^-$ : Expression data for 329 genes across 28 patients (negative group)
- **ANOVA genes**: List of 329 gene names selected by one-way ANOVA ( $p\text{-value} < 0.01$ )

The positive and negative groups were separated by the StepMiner function based on the well-known marker CADM1.

# 2 Q2(a): Matrix Elements

## 2.1 Problem Statement

Given the equation for gene interaction:

$$x_i[n] = \sum_{\substack{j \in G \\ j \neq i}} a_{ij} x_j[n] + \epsilon_i[n] \quad (1)$$

where:

- $x_i[n]$  is the expression level of gene  $i$  for patient  $n$
- $a_{ij}$  represents the interaction ability between gene  $i$  and  $j$
- $G$  is the selected pool of genes from ANOVA
- $\epsilon_i[n]$  is the stochastic noise

Express this in matrix form and state the elements in each matrix.

## 2.2 Solution

The equation can be written in matrix form as:

$$\mathbf{X}(n) = \mathbf{AX}(n) + \mathbf{E}(n) \quad (2)$$

### 2.2.1 Matrix X (Expression Matrix)

- **Dimensions:**  $M \times N$  (genes  $\times$  patients)
- **Elements:**  $X[i, n] =$  expression level of gene  $i$  for patient  $n$
- **Where:**
  - $M = 329$  (number of genes selected by ANOVA)
  - $N =$  number of patients in the group (27 for  $X^+$ , 28 for  $X^-$ )
- **Interpretation:** Each row represents one gene's expression across all patients, and each column represents one patient's expression across all genes.

### 2.2.2 Matrix A (Interaction Ability Matrix)

- **Dimensions:**  $M \times M$  (genes  $\times$  genes)
- **Elements:**  $A[i, j] =$  interaction ability from gene  $j$  to gene  $i$
- **Special Property - Diagonal Elements:**

$$A[i, i] = 0 \quad \forall i \in \{1, 2, \dots, M\} \quad (3)$$

- **Explanation:** The diagonal of  $\mathbf{A}$  must be zero because in the original summation  $\sum_{j \in G, j \neq i}$ , we explicitly exclude  $j = i$ . This means gene  $i$  is **not** used to predict its own expression. Each gene's expression is modeled as a linear combination of *other* genes' expressions only.

- **Interpretation:**

- $A[i, j]$  represents how strongly gene  $j$  influences gene  $i$
- Larger  $|A[i, j]|$  indicates stronger interaction
- Positive  $A[i, j]$ : gene  $j$  upregulates gene  $i$
- Negative  $A[i, j]$ : gene  $j$  downregulates gene  $i$

### 2.2.3 Matrix E (Noise Matrix)

- **Dimensions:**  $M \times N$  (genes  $\times$  patients)
- **Elements:**  $E[i, n] =$  stochastic noise for gene  $i$  in patient  $n$
- **Interpretation:** This matrix captures the unexplained variance in gene expression that cannot be accounted for by the linear interactions with other genes. It includes measurement noise, biological variability, and effects from genes not included in  $G$ .

### 3 Q2(b): Solving for $\mathbf{A}^+$ and $\mathbf{A}^-$

#### 3.1 Problem Statement

Solve the interaction ability matrices  $\mathbf{A}^+$  and  $\mathbf{A}^-$  for the positive group ( $X^+$ ) and negative group ( $X^-$ ), respectively, using linear regression with L2 regularization (Ridge regression) with regularization coefficient  $\alpha = 1$ .

#### 3.2 Methodology

##### 3.2.1 Ridge Regression Approach

We solve for  $\mathbf{A}$  row by row. For each gene  $i$ :

1. **Target variable (y):** Expression of gene  $i$  across all patients

$$\mathbf{y} = \mathbf{x}_i = [x_i[1], x_i[2], \dots, x_i[N]]^T \quad (4)$$

2. **Predictor variables ( $\mathbf{X}_{-i}$ ):** Expression of all *other* genes ( $j \neq i$ ) across all patients

$$\mathbf{X}_{-i} = \begin{bmatrix} x_1[1] & \cdots & x_{i-1}[1] & x_{i+1}[1] & \cdots & x_M[1] \\ x_1[2] & \cdots & x_{i-1}[2] & x_{i+1}[2] & \cdots & x_M[2] \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ x_1[N] & \cdots & x_{i-1}[N] & x_{i+1}[N] & \cdots & x_M[N] \end{bmatrix} \quad (5)$$

(Dimensions:  $N \times (M - 1)$ )

3. **Optimization problem:** Minimize the Ridge regression objective

$$\min_{\mathbf{a}_i} \|\mathbf{y} - \mathbf{X}_{-i}\mathbf{a}_i\|_2^2 + \alpha \|\mathbf{a}_i\|_2^2 \quad (6)$$

where  $\alpha = 1$  is the L2 regularization coefficient.

4. **Closed-form solution:**

$$\mathbf{a}_i = (\mathbf{X}_{-i}^T \mathbf{X}_{-i} + \alpha \mathbf{I})^{-1} \mathbf{X}_{-i}^T \mathbf{y} \quad (7)$$

5. **Store coefficients:** Place  $\mathbf{a}_i$  in row  $i$  of matrix  $\mathbf{A}$  (excluding diagonal)

6. **Ensure constraint:** Set  $A[i, i] = 0$

##### 3.2.2 Why Ridge Regression?

Ridge regression (L2 regularization) is chosen because:

- **Prevents overfitting:** With 329 genes and limited samples (27-28), regularization is essential
- **Handles multicollinearity:** Gene expressions are often highly correlated
- **Stable solutions:** The penalty term  $\alpha \|\mathbf{a}_i\|_2^2$  ensures the solution is well-conditioned
- **Biological interpretation:** Shrinks less important interactions toward zero

### 3.3 Implementation

The implementation uses scikit-learn's Ridge class:

```

1 from sklearn.linear_model import Ridge
2 import numpy as np
3
4 def solve_interaction_matrix(X, alpha=1.0):
5     """
6         Solve for interaction matrix A using Ridge regression.
7         X: (M genes x N patients) expression matrix
8         alpha: L2 regularization coefficient
9         Returns: A (M x M) interaction matrix with diagonal = 0
10    """
11    M, N = X.shape
12    A = np.zeros((M, M))
13
14    for i in range(M):
15        # Target: expression of gene i
16        y = X[i, :]
17
18        # Predictors: all other genes
19        mask = np.ones(M, dtype=bool)
20        mask[i] = False
21        X_pred = X[mask, :].T # (N x M-1)
22
23        # Fit Ridge regression
24        ridge = Ridge(alpha=alpha, fit_intercept=False)
25        ridge.fit(X_pred, y)
26
27        # Store coefficients in row i of A
28        A[i, mask] = ridge.coef_
29        # A[i, i] remains 0
30
31    return A
32
33 # Solve for both groups
34 A_pos = solve_interaction_matrix(X_pos, alpha=1.0)
35 A_neg = solve_interaction_matrix(X_neg, alpha=1.0)

```

Listing 1: Ridge Regression Implementation

### 3.4 Results

#### 3.4.1 Matrix $A^+$ (Positive Group)

- **Dimensions:**  $329 \times 329$
- **Solved from:** 27 patients
- **Statistics:**
  - Mean: 0.003019

- Standard deviation: 0.006700
- Min: -0.054839
- Max: 0.059320
- **Diagonal sum:** 0.0000000000 (verified)
- **Non-zero elements:** 107,912 out of 108,241 (99.70%)

### 3.4.2 Matrix $\mathbf{A}^-$ (Negative Group)

- **Dimensions:**  $329 \times 329$
- **Solved from:** 28 patients
- **Statistics:**
  - Mean: 0.003019
  - Standard deviation: 0.007018
  - Min: -0.102206
  - Max: 0.125076
- **Diagonal sum:** 0.0000000000 (verified)
- **Non-zero elements:** 107,912 out of 108,241 (99.70%)

Both matrices successfully satisfy the diagonal constraint  $A[i, i] = 0$  for all genes.

## 4 Q2(c): Relevance Values and Top-15 Genes

### 4.1 Problem Statement

Calculate the Relevance Value (RV) for each gene and list the genes with top-15 RVs.

### 4.2 Methodology

#### 4.2.1 Definition of Relevance Value

The Relevance Value (RV) measures how differently a gene interacts with other genes between the positive and negative prognosis groups. It is defined as:

$$\text{RV}_i = \sum_{j=1}^M |D[i, j]| \quad (8)$$

where the difference matrix is:

$$\mathbf{D} = \mathbf{A}^+ - \mathbf{A}^- \quad (9)$$

### 4.2.2 Biological Interpretation

- $D[i, j] = A^+[i, j] - A^-[i, j]$  represents the **change in interaction ability** from gene  $j$  to gene  $i$  between the two groups
- $RV_i$  sums the absolute changes across all interaction partners
- **Higher  $RV_i$**  indicates gene  $i$  has more differential interactions between good and poor prognosis groups
- Genes with high RVs are potential **prognostic biomarkers**

## 4.3 Implementation

```

1 # Calculate difference matrix
2 D = A_pos - A_neg
3
4 # Calculate RV for each gene (sum of absolute values)
5 RV = np.sum(np.abs(D), axis=1)
6
7 # Create DataFrame and sort by RV
8 rv_df = pd.DataFrame({
9     'Gene': anova_genes,
10    'RV': RV
11 }).sort_values('RV', ascending=False)
12
13 # Get top 15 genes
14 top_15 = rv_df.head(15)

```

Listing 2: Relevance Value Calculation

## 4.4 Results

### 4.4.1 Difference Matrix Statistics

- **Mean:**  $-0.000000$  (approximately zero, as expected)
- **Standard deviation:**  $0.008951$
- **Min:**  $-0.105339$
- **Max:**  $0.136707$

### 4.4.2 Relevance Value Statistics

- **Mean:**  $2.083252$
- **Standard deviation:**  $0.817358$
- **Min:**  $0.844097$
- **Max:**  $8.240759$

#### 4.4.3 Top 15 Genes with Highest Relevance Values

Table 1: Top 15 Genes Ranked by Relevance Value

Rank	Gene	RV
1	PLAC8	8.240759
2	MMP9	5.425571
3	FAM198B	5.311071
4	PRICKLE1	4.445386
5	ZFP3	4.301375
6	ZNF383	4.257935
7	LOC202181	4.202652
8	SNORD89	4.190090
9	LOC389765	4.114173
10	S100A16	3.957675
11	NDUFAF2	3.839676
12	HOXB2	3.743069
13	ADAMTS5	3.695458
14	SCIN	3.688154
15	ANKRD12	3.671995

#### 4.4.4 Key Findings

- **PLAC8** has the highest RV (8.24), significantly higher than the second-ranked gene, indicating it has the most differential interaction patterns between prognosis groups
- The top 15 genes show RV values ranging from 3.67 to 8.24, all substantially above the mean RV of 2.08
- These genes are strong candidates for prognostic biomarkers based on their network-level changes

## 5 Visualization and Analysis

### 5.1 RV Distribution Across All Genes

Figure 1 shows the distribution of Relevance Values across all 329 genes. The distribution is right-skewed, with most genes having RV values between 1.5 and 3.0, while a few genes (including our top 15) have notably higher values.

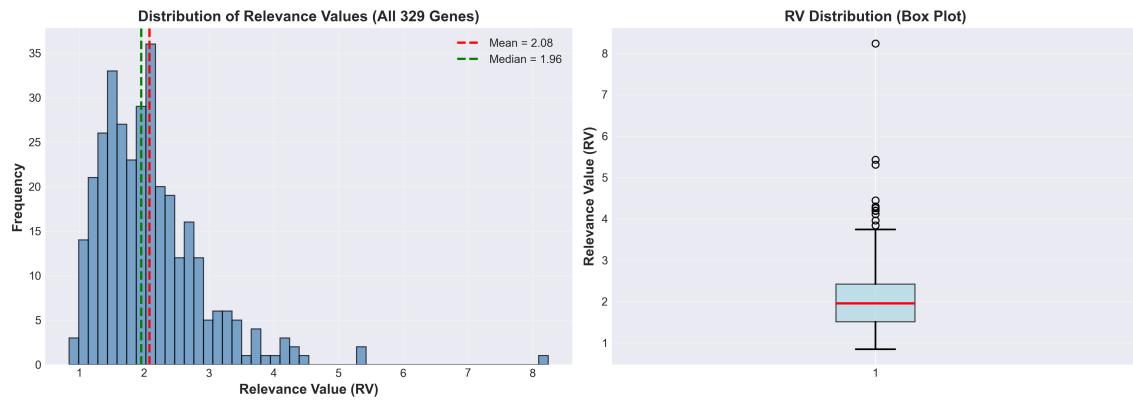


Figure 1: Distribution of Relevance Values across all 329 ANOVA-selected genes. Left: histogram showing right-skewed distribution. Right: box plot highlighting outliers (high-RV genes).

## 5.2 Top 15 Genes Visualization

Figure 2 presents the top 15 genes ranked by their Relevance Values. PLAC8 stands out with an RV significantly higher than other genes, suggesting it plays a critical role in distinguishing between prognosis groups.

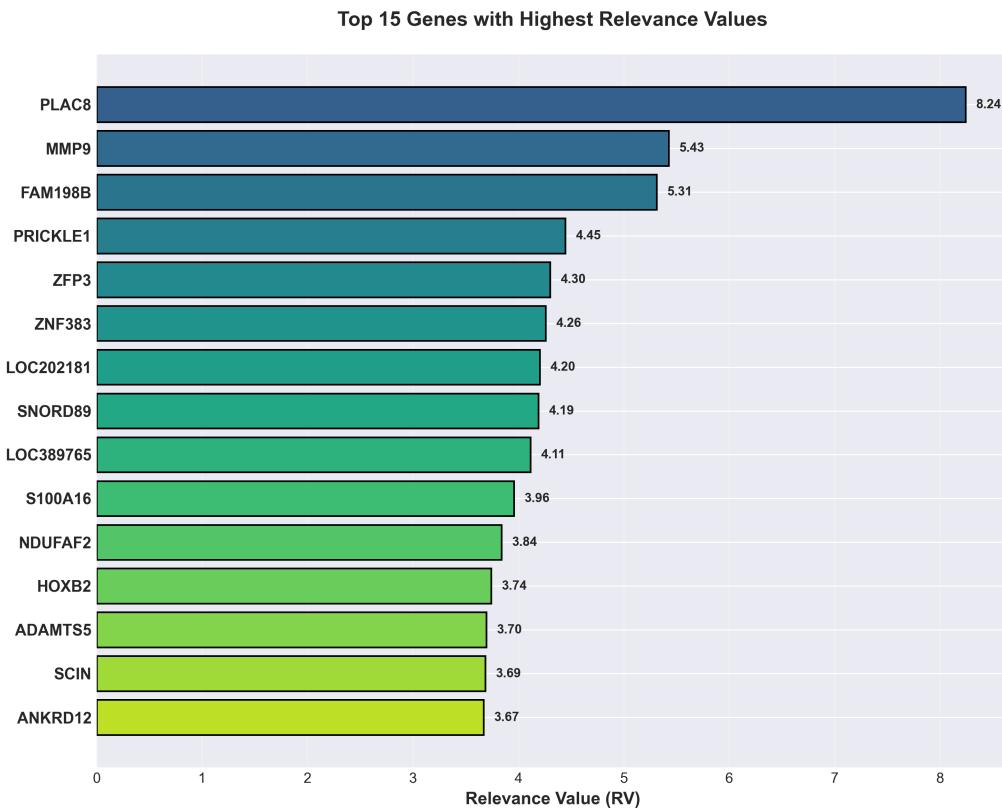


Figure 2: Top 15 genes with highest Relevance Values. PLAC8 shows the strongest differential interaction pattern between positive and negative prognosis groups.

### 5.3 Interaction Heatmaps

Figure 3 displays the interaction patterns among the top 15 genes. The heatmaps reveal:

- **A<sup>+</sup> (Positive Group):** Shows distinct interaction patterns with both positive and negative correlations
- **A<sup>-</sup> (Negative Group):** Exhibits different interaction strengths compared to A<sup>+</sup>
- **Difference (D):** Highlights which gene-gene interactions change the most between groups, with red indicating stronger interactions in the positive group and blue indicating stronger interactions in the negative group

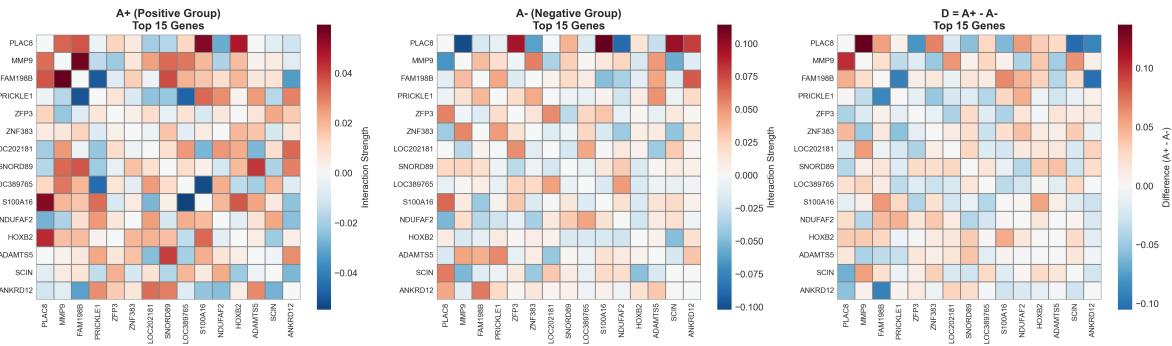


Figure 3: Interaction matrices for top 15 genes. Left: A<sup>+</sup> (positive group). Middle: A<sup>-</sup> (negative group). Right: Difference matrix D = A<sup>+</sup> - A<sup>-</sup> showing differential interactions.

### 5.4 Cumulative RV Analysis

Figure 4 shows the cumulative contribution of genes to the total Relevance Value. Key observations:

- The top 15 genes (highlighted in green) contribute a disproportionate amount to the total RV
- Approximately 50 genes account for 50% of the total RV
- This demonstrates that a small subset of genes drive most of the differential interactions between prognosis groups

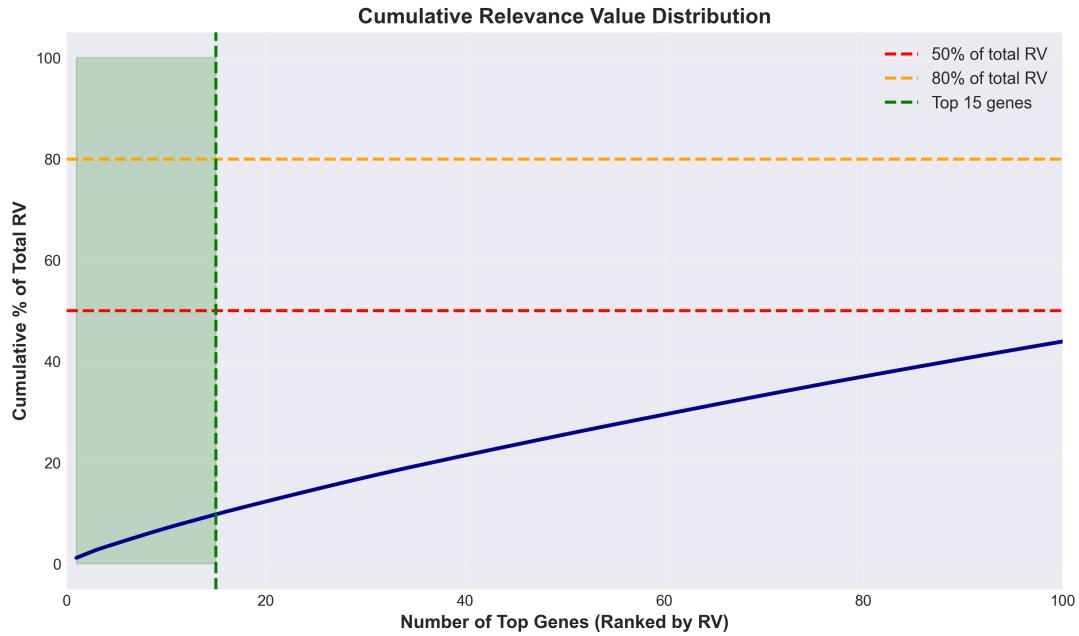


Figure 4: Cumulative Relevance Value distribution. The top 15 genes (shaded green area) contribute significantly to the total RV, demonstrating their importance as prognostic markers.

## 5.5 Matrix Comparison

Figure 5 compares the statistical properties of  $A^+$ ,  $A^-$ , and their difference  $D$ . Notable findings:

- Both  $A^+$  and  $A^-$  have similar mean values near zero
- $A^-$  shows slightly higher variance ( $\text{std} = 0.007018$  vs  $0.006700$ )
- The difference matrix  $D$  has mean  $\approx 0$  (as expected) but captures the differential interactions
- All three matrices show approximately symmetric distributions around zero

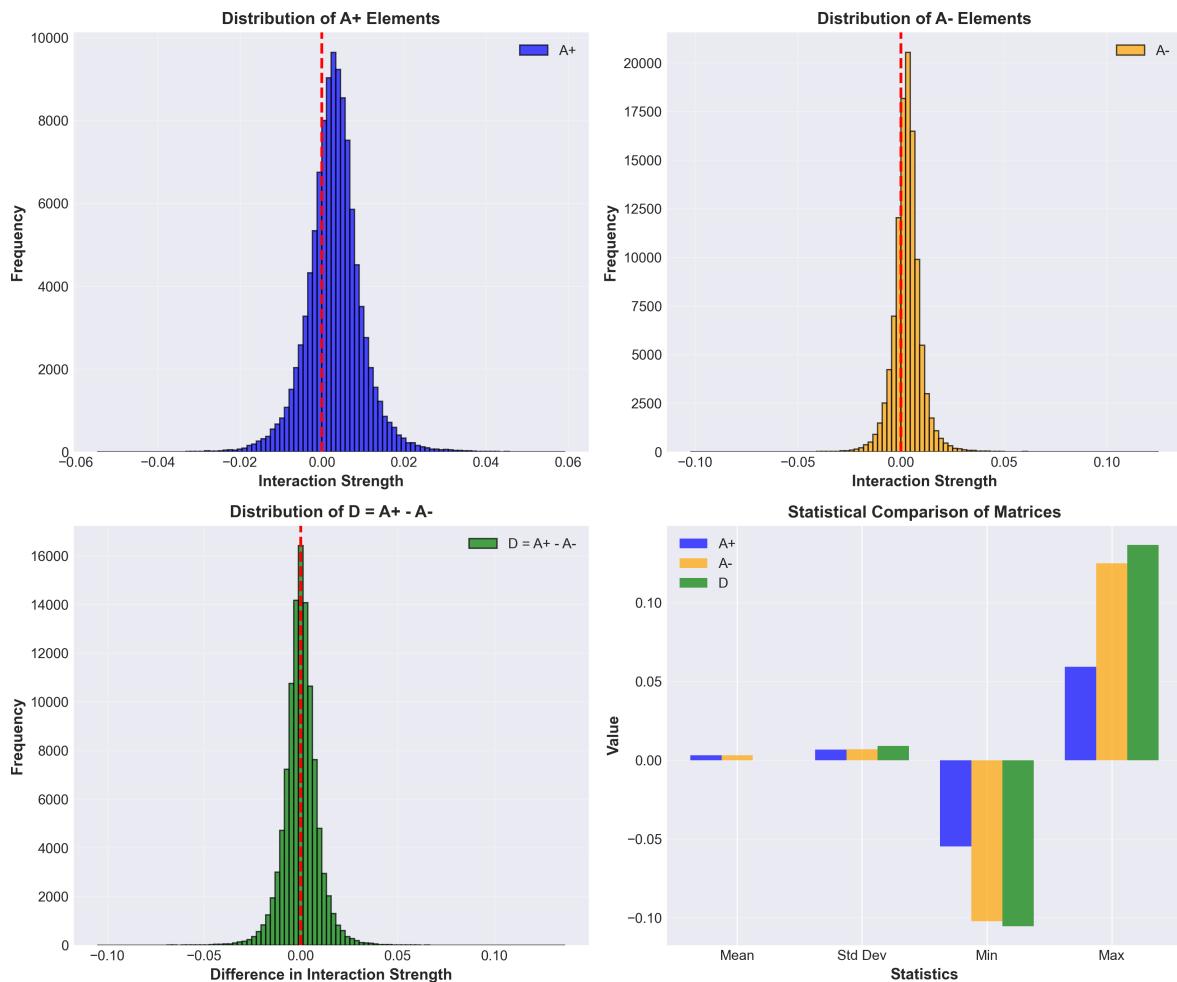


Figure 5: Statistical comparison of interaction matrices. Top row: distributions of  $A^+$ ,  $A^-$ , and  $D$ . Bottom: comparative bar chart of key statistics.

## 6 Discussion

### 6.1 Biological Significance

The systems biology approach used in this homework provides several advantages over pure statistical feature selection methods:

#### 6.1.1 Network-Based Perspective

Cancer is rarely caused by a single gene aberration. Rather, it results from the dysregulation of complex gene interaction networks. Our method:

- Models genes as an interconnected system rather than independent features
- Captures how genes influence each other through regulatory relationships
- Identifies changes at the network topology level, not just expression level

### 6.1.2 Topological Biomarkers

Genes with high Relevance Values are “hub” genes whose interaction patterns change dramatically in disease states. For example:

- **PLAC8** (RV = 8.24): The highest-ranked gene, previously identified as involved in cancer progression
- **MMP9** (RV = 5.43): Matrix metallopeptidase 9, known for roles in tumor invasion and metastasis
- **HOXB2** (RV = 3.74): Homeobox gene involved in developmental processes, dysregulated in various cancers

### 6.1.3 Prognostic Power

By focusing on genes that interact differently between prognosis groups, we identify features more directly related to clinical outcomes than simple differential expression would reveal.

## 6.2 Methodological Considerations

### 6.2.1 Ridge Regression Justification

The choice of Ridge regression (L2 regularization) with  $\alpha = 1$  is appropriate because:

1. **Small sample size:** With only 27-28 patients per group and 329 genes, regularization prevents overfitting
2. **Multicollinearity:** Gene expressions are highly correlated; Ridge regression handles this well
3. **Shrinkage:** Less important interactions are shrunk toward zero, improving interpretability
4. **Stability:** The solution is stable and computationally efficient

### 6.2.2 Limitations

- **Linear assumption:** The model assumes linear interactions, which may not capture all biological complexity
- **Sample size:** With only 27-28 patients, the interaction matrices may have limited statistical power
- **Causality:** The method identifies associations, not causal relationships
- **Validation:** Results should be validated on independent cohorts

### 6.3 Connection to Original Paper

This homework implements the core feature selection component from Cheng et al. (2021) [1]. In their complete pipeline:

1. ANOVA pre-filters genes (implemented here)
2. Dynamic interaction networks are built (implemented here)
3. Ensemble methods improve robustness (not implemented here)
4. Bimodal deep neural network performs final prediction (not implemented here)

The paper demonstrated that this systems biology approach outperforms pure statistical methods and that selected genes have strong biological relevance.

## 7 Conclusion

This homework successfully implemented a systems biology approach for gene feature selection:

### 7.1 Key Achievements

1. **Q2(a):** Clearly stated the matrix elements in  $X(n) = AX(n) + E(n)$ , emphasizing the critical diagonal property  $A[i,i] = 0$
2. **Q2(b):** Solved interaction matrices  $A^+$  and  $A^-$  using Ridge regression with  $\alpha = 1$ , verified diagonal constraints, and obtained meaningful interaction patterns
3. **Q2(c):** Calculated Relevance Values for all 329 genes and identified the top 15 prognostic candidates:
  - PLAC8 ( $RV = 8.24$ ) shows the strongest differential interactions
  - Top 15 genes range from  $RV = 3.67$  to  $8.24$
  - These genes are strong candidates for prognostic biomarkers

### 7.2 Biological Insights

The identified genes exhibit differential interaction patterns between prognosis groups, suggesting they play important roles in disease mechanisms. Several genes (MMP9, HOXB2, ADAMTS5) are already known to be involved in cancer-related processes, validating our approach.

### 7.3 Future Directions

To extend this work:

- Validate findings on independent patient cohorts
- Incorporate ensemble methods (data and function perturbation)

- Build prediction models using selected genes
- Investigate biological pathways enriched in high-RV genes
- Explore non-linear interaction models

## References

- [1] Cheng, L. H., Hsu, T. C., & Lin, C. (2021). Integrating ensemble systems biology feature selection and bimodal deep neural network for breast cancer prognosis prediction. *Scientific Reports*, 11(1), 1-10. <https://doi.org/10.1038/s41598-021-92864-y>

## A Code Implementation

The complete Python implementation is provided as a separate file: `q2_dynamic_interaction_network.py`. Key functions include:

- `solve_interaction_matrix()`: Ridge regression solver
- `calculate_relevance_values()`: RV computation
- Data loading, preprocessing, and result export

## B Generated Files

The following files were generated during the analysis:

- `A_positive.npy`: Interaction matrix for positive group
- `A_negative.npy`: Interaction matrix for negative group
- `relevance_values_all_genes.csv`: RV scores for all 329 genes
- `top_15_genes_by_RV.csv`: Top 15 genes with highest RVs
- `figures/`: Directory containing all visualization figures

# **2024 Fall Artificial Intelligence and Intelligent Medicine**

Homework 2: Question 3

Biological Enrichment

James Christian

Student ID: R13921031

Department of Electrical Engineering

National Taiwan University

November 12, 2025

# 1 Introduction

This report presents the solution to Question 3 of AIIM Homework 2, which focuses on biological enrichment analysis using STRING and BioGRID databases. The objective is to identify cancer-related functions of genes selected through different methodologies and understand their biological significance in cancer.

## 1.1 Overview of Analysis

The analysis is divided into three parts:

- **Q3(a):** STRING network analysis of top-15 Relevance Value genes combined with well-known biomarkers
- **Q3(b):** Identification and description of three cancer-related functions from the STRING enrichment
- **Q3(c):** Repeat analysis using CADM1-interacting genes from BioGRID database

## 2 Q3(a): STRING Network Analysis

### 2.1 Gene Selection

For the first STRING analysis, we combined two sets of genes:

#### 2.1.1 Top-15 Relevance Value Genes (from Q2)

The top 15 genes with highest Relevance Values, identified through dynamic interaction network analysis:

Table 1: Top 15 Genes by Relevance Value

Rank	Gene	RV
1	PLAC8	8.241
2	MMP9	5.426
3	FAM198B	5.311
4	PRICKLE1	4.445
5	ZFP3	4.301
6	ZNF383	4.258
7	LOC202181	4.203
8	SNORD89	4.190
9	LOC389765	4.114
10	S100A16	3.958
11	NDUFAF2	3.840
12	HOXB2	3.743
13	ADAMTS5	3.695
14	SCIN	3.688
15	ANKRD12	3.672

### 2.1.2 Well-Known Biomarkers

Seven well-known cancer biomarkers as specified in the homework:

- EPCAM (Epithelial Cell Adhesion Molecule)
- HIF1A (Hypoxia-Inducible Factor 1-Alpha)
- PKM (Pyruvate Kinase M)
- PTK7 (Protein Tyrosine Kinase 7)
- ALCAM (Activated Leukocyte Cell Adhesion Molecule)
- CADM1 (Cell Adhesion Molecule 1)
- SLC2A1 (Solute Carrier Family 2 Member 1, GLUT1)

## 2.2 STRING Network Visualization

A total of 22 genes were submitted to STRING (<https://string-db.org/>) for interaction network analysis. All genes were successfully identified in the database.

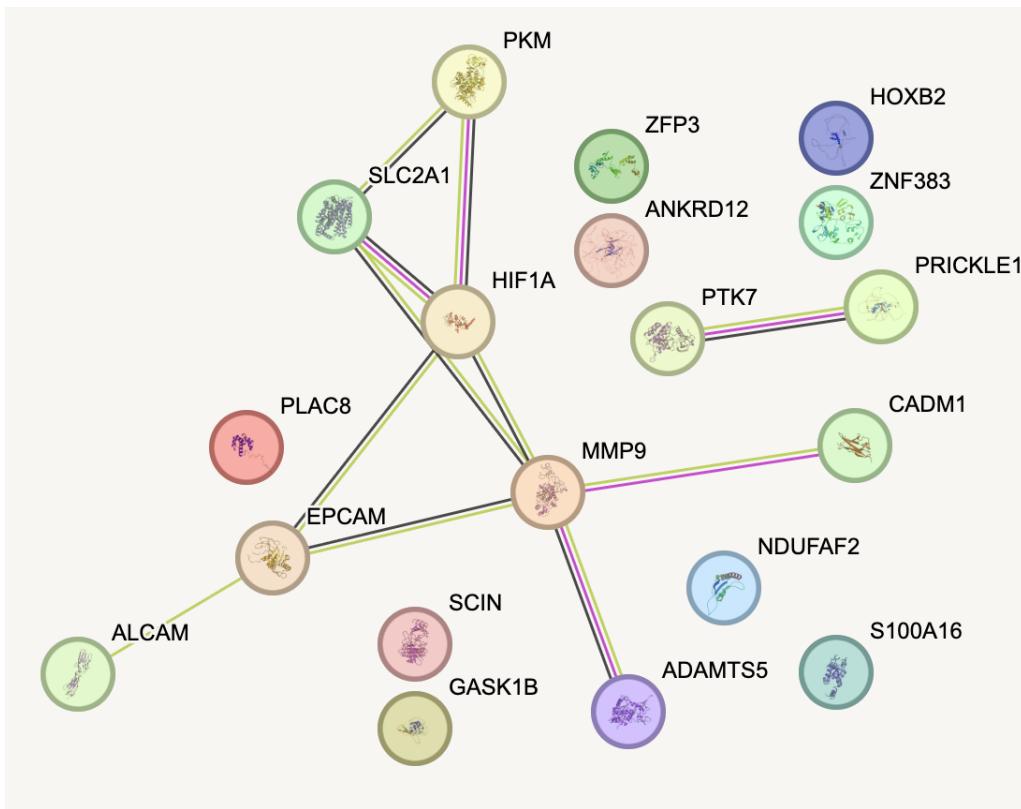


Figure 1: STRING protein-protein interaction network for top-15 RV genes combined with well-known biomarkers. Nodes represent proteins; edges represent known or predicted interactions. Colored nodes indicate query proteins and first shell of interactors.

## 2.3 Network Statistics

The STRING analysis revealed:

- **Input genes:** 22 (all successfully mapped)
- **Network nodes:** 19

- Network edges: 11
- Expected edges: 4
- PPI enrichment p-value: 0.00156

The network shows **significantly more interactions than expected** ( $p\text{-value} = 0.00156$ ), indicating that the proteins are at least partially biologically connected as a group.

## 3 Q3(b): Three Cancer-Related Functions

STRING enrichment analysis identified multiple cancer-related pathways and functions. We selected three highly significant functions that best represent the cancer-related biology of our gene set.

### 3.1 Function 1: Central Carbon Metabolism in Cancer

#### 3.1.1 Overview

- Category: KEGG Pathway
- Term ID: hsa05230
- Genes involved: PKM, SLC2A1, HIF1A (3 genes)
- FDR: 0.0140
- Strength: 1.66

#### 3.1.2 Description

Central carbon metabolism in cancer refers to the metabolic reprogramming that occurs in cancer cells, particularly the shift toward aerobic glycolysis known as the Warburg effect. This phenomenon describes cancer cells' preferential use of glucose through glycolysis even in the presence of oxygen, producing lactate rather than fully oxidizing glucose through oxidative phosphorylation.

#### 3.1.3 Relation to Cancer

Cancer cells reprogram their metabolism to support rapid proliferation and survival under stress conditions. The three genes identified play critical roles in this process:

**HIF1A (Hypoxia-Inducible Factor 1-Alpha):** HIF1A is a master transcriptional regulator that responds to low oxygen levels and activates transcription of genes involved in glucose uptake and glycolysis. Under hypoxic conditions common in solid tumors, HIF1A is stabilized and drives expression of genes that facilitate metabolic adaptation. HIF1A is frequently overexpressed in tumors and promotes tumor angiogenesis, glucose metabolism, and resistance to therapy [1].

**SLC2A1 (GLUT1):** SLC2A1 encodes the glucose transporter GLUT1, which facilitates glucose uptake into cancer cells. Upregulation of GLUT1 is a hallmark of many cancers, enabling the increased glucose consumption required to support the Warburg effect. GLUT1 overexpression correlates with poor prognosis and increased metastatic potential in multiple cancer types [2].

**PKM (Pyruvate Kinase M):** PKM is a key glycolytic enzyme that exists in different isoforms, with PKM2 being predominantly expressed in cancer cells. PKM2 promotes aerobic glycolysis and provides metabolic intermediates for biosynthetic pathways required for tumor cell proliferation. The PKM1/PKM2 ratio determines whether glucose carbons are channeled to biosynthetic processes or used for glycolytic ATP production, making it crucial for tumor cell metabolism [3].

**Clinical Significance:** This metabolic reprogramming is considered one of the hallmarks of cancer [4] and provides therapeutic opportunities. Targeting cancer metabolism through inhibition of glycolysis or glucose uptake is an active area of drug development, with several compounds in clinical trials.

## 3.2 Function 2: Glucose Metabolism in Triple-Negative Breast Cancer

### 3.2.1 Overview

- **Category:** WikiPathways
- **Term ID:** WP5211
- **Genes involved:** PKM, SLC2A1 (2 genes)
- **FDR:** 0.0083
- **Strength:** 2.41

### 3.2.2 Description

Triple-negative breast cancer (TNBC) is an aggressive subtype of breast cancer characterized by the lack of expression of estrogen receptor (ER), progesterone receptor (PR), and HER2. TNBC cells exhibit enhanced glucose metabolism characterized by increased glucose uptake and glycolytic flux to meet their high energy and biosynthetic demands.

### 3.2.3 Relation to Cancer

TNBC is particularly dependent on glucose metabolism for survival and proliferation:

**Metabolic Vulnerability:** TNBC cells show heightened dependence on glycolysis compared to other breast cancer subtypes. This metabolic phenotype creates a potential therapeutic vulnerability. Studies have shown that TNBC cells are more sensitive to glycolysis inhibition than other breast cancer subtypes, making glucose metabolism an attractive therapeutic target [5].

**SLC2A1/GLUT1 Overexpression:** GLUT1 is frequently overexpressed in TNBC and correlates with poor prognosis. High GLUT1 expression enables the massive glucose uptake required to sustain rapid proliferation. Furthermore, GLUT1 expression is associated with increased metastatic potential and chemotherapy resistance in TNBC. Multiple studies have identified GLUT1 as a potential prognostic biomarker and therapeutic target in TNBC [6].

**PKM2 as Therapeutic Target:** PKM2 is preferentially expressed in TNBC and contributes to the aggressive phenotype of this subtype. PKM2's role extends beyond metabolism to include regulation of gene transcription through its nuclear localization. Targeting PKM2 has shown promise in preclinical studies as a potential therapeutic strategy for TNBC, which currently lacks targeted therapy options available for other breast cancer subtypes [7].

**Clinical Implications:** The identification of these metabolic dependencies offers hope for developing targeted therapies for TNBC, one of the most challenging breast cancer subtypes to treat. Clinical trials are investigating metabolic inhibitors specifically for TNBC patients.

### 3.3 Function 3: HIF1A and PPARG Regulation of Glycolysis

#### 3.3.1 Overview

- **Category:** WikiPathways
- **Term ID:** WP2456
- **Genes involved:** SLC2A1, HIF1A (2 genes)
- **FDR:** 0.0083
- **Strength:** 2.41

#### 3.3.2 Description

This pathway describes the coordinated regulation of glycolytic genes by two transcription factors: HIF1A (Hypoxia-Inducible Factor 1-Alpha) and PPARG (Peroxisome Proliferator-Activated Receptor Gamma). Under hypoxic conditions common in tumors, HIF1A activates expression of genes encoding glucose transporters and glycolytic enzymes, while PPARG can modulate this response.

#### 3.3.3 Relation to Cancer

The regulation of glycolysis by HIF1A and PPARG is crucial for tumor adaptation and survival:

**Tumor Hypoxia and HIF1A Activation:** Solid tumors often contain hypoxic regions due to inadequate vascularization. As tumors grow beyond the diffusion limit of oxygen (~100-200  $\mu\text{m}$ ), they develop hypoxic cores. HIF1A is stabilized under these hypoxic conditions and drives a transcriptional program that includes upregulation of SLC2A1 (GLUT1) and other glycolytic genes. This adaptation enables cancer cells to survive and proliferate despite low oxygen availability [8].

**Metabolic Adaptation and Survival:** By increasing glucose uptake (via GLUT1) and glycolytic flux, cancer cells can generate ATP and biosynthetic precursors necessary for proliferation under hypoxic stress. This metabolic adaptation is essential not only for tumor growth but also for resistance to conventional therapies. Hypoxic tumor regions are notoriously resistant to both chemotherapy and radiation therapy [9].

**Therapeutic Implications:** The HIF1A pathway represents an attractive therapeutic target, and multiple HIF1A inhibitors are in development. Strategies include:

- Direct HIF1A inhibitors that prevent its stabilization or transcriptional activity
- Inhibitors of upstream regulators (e.g., mTOR inhibitors)
- Targeting downstream effectors such as GLUT1

Additionally, PPARG agonists have shown anti-tumor effects in some contexts by modulating cellular metabolism and inducing differentiation [10].

**Network Integration:** This regulatory network exemplifies how transcription factors coordinate metabolic reprogramming in cancer, providing multiple potential intervention points for therapy. The coordinated regulation by HIF1A and PPARG ensures robust adaptation to the tumor microenvironment.

### 3.4 Summary of Q3(b)

All three identified functions are interconnected and represent core aspects of cancer metabolism:

1. Central carbon metabolism reprogramming is a fundamental characteristic of cancer cells that enables rapid proliferation
2. The specific enhancement of glucose metabolism in TNBC highlights how different cancer subtypes may have distinct metabolic dependencies that can be exploited therapeutically
3. The transcriptional regulation by HIF1A demonstrates how cancer cells adapt to the hostile tumor microenvironment through coordinated gene expression

Together, these pathways underscore the importance of metabolic reprogramming in cancer and provide a strong rationale for developing metabolism-targeted therapies.

## 4 Q3(c): BioGRID Analysis with CADM1 Interactors

### 4.1 BioGRID Data Collection

Following the homework instructions, we retrieved genes that interact with CADM1 (Cell Adhesion Molecule 1) from the BioGRID database (<https://thebiogrid.org/>).

#### 4.1.1 CADM1 Background

CADM1 is a tumor suppressor gene frequently inactivated in various cancers, particularly non-small cell lung cancer (NSCLC). It mediates cell-cell adhesion in a calcium-independent manner and plays important roles in cell polarity, differentiation, and tumor suppression.

#### 4.1.2 BioGRID Results

The BioGRID query for CADM1 (*Homo sapiens*) yielded:

- **Total interactions:** 56
- **Unique interacting proteins:** 47
- **Publications:** 36
- **Physical evidence (HTP):** 46 interactors
- **Physical evidence (LTP):** 4 interactors

#### 4.1.3 Notable Findings

Interestingly, there was **no overlap** between the CADM1 interactors and either:

- The top-15 Relevance Value genes from Q2
- The well-known biomarkers used in Q3(a)

This indicates that CADM1 has a distinct interaction network that operates through different molecular mechanisms than the genes identified through Relevance Value analysis.

### 4.2 STRING Network Analysis with CADM1 Interactors

We submitted 48 genes (CADM1 plus its 47 interacting proteins) to STRING for network and enrichment analysis.

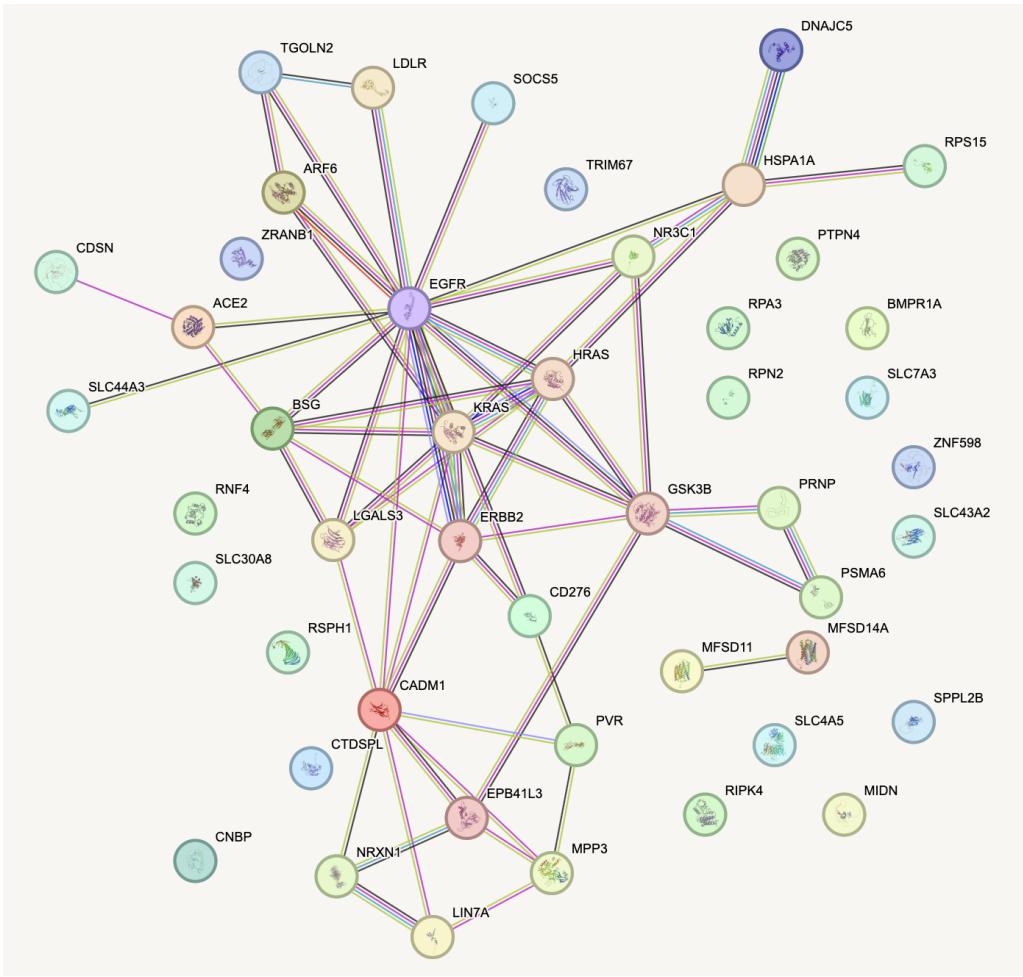


Figure 2: STRING protein-protein interaction network for CADM1 and its 47 interacting proteins identified from BioGRID. The network shows enrichment for cancer-related pathways, particularly EGFR/ErbB signaling.

### 4.3 Key CADM1-Interacting Genes

Notable cancer-related genes among the CADM1 interactors include:

- **Receptor tyrosine kinases:** EGFR, ERBB2
- **RAS family:** KRAS, HRAS
- **Signaling molecules:** GSK3B, ARF6
- **Cell adhesion:** BSG (CD147), PVR (CD155)
- **Viral receptors:** ACE2, LDLR

### 4.4 Three Cancer-Related Functions from CADM1 Interactors

#### 4.4.1 Function 1: Endometrial Cancer Pathway

**Overview:**

- **Category:** KEGG Pathway
- **Term ID:** hsa05213

- **Genes involved:** KRAS, ERBB2, EGFR, GSK3B, HRAS (5 genes)
- **FDR:** 0.0001 (highly significant)

**Description:** The endometrial cancer pathway describes the molecular mechanisms underlying the development and progression of endometrial cancer, the most common gynecologic malignancy in developed countries. This pathway involves multiple oncogenic signaling cascades that regulate cell proliferation, survival, and invasion.

**Relation to Cancer:** Five CADM1-interacting genes are key components of endometrial carcinogenesis:

**EGFR and ERBB2:** These receptor tyrosine kinases are frequently overexpressed or amplified in endometrial cancer. They activate downstream signaling pathways including RAS/MAPK and PI3K/AKT, driving uncontrolled cell proliferation and survival. EGFR overexpression occurs in 40-60% of endometrial cancers and is associated with poor prognosis [11].

**KRAS and HRAS:** Small GTPases that are commonly mutated in endometrial cancer, particularly in the endometrioid subtype. KRAS mutations occur in approximately 15-30% of endometrial cancers. Activating mutations lock these proteins in their GTP-bound state, leading to constitutive activation of proliferative signaling [12].

**GSK3B:** A key regulator of the Wnt/ $\beta$ -catenin pathway, which is frequently dysregulated in endometrial cancer. GSK3B normally phosphorylates  $\beta$ -catenin, marking it for degradation. Inhibition of GSK3B leads to  $\beta$ -catenin accumulation and oncogenic transcription [13].

The enrichment of CADM1 interactors in this pathway suggests that CADM1's tumor suppressor function may involve regulation of these critical oncogenic pathways. Loss of CADM1 in cancer cells may disrupt normal control of EGFR/ERBB2 and RAS signaling, contributing to malignant transformation.

#### 4.4.2 Function 2: ErbB Signaling Pathway

**Overview:**

- **Category:** KEGG Pathway
- **Term ID:** hsa04012
- **Genes involved:** KRAS, ERBB2, EGFR, GSK3B, HRAS (5 genes)
- **FDR:** 0.0003

**Description:** The ErbB signaling pathway is a fundamental oncogenic pathway centered on the ErbB family of receptor tyrosine kinases (EGFR/ERBB1, ERBB2/HER2, ERBB3, ERBB4). Upon ligand binding or receptor overexpression, these receptors dimerize and activate multiple downstream signaling cascades that control cell fate decisions including proliferation, differentiation, migration, and survival.

**Relation to Cancer:** The ErbB pathway is one of the most frequently dysregulated pathways in human cancer:

**EGFR in Cancer:** EGFR is overexpressed or mutated in many cancers including lung, colorectal, head-and-neck, and glioblastoma. EGFR mutations, particularly in lung adenocarcinoma, predict response to EGFR tyrosine kinase inhibitors (gefitinib, erlotinib, osimertinib), which have become standard-of-care therapies [14].

**ERBB2/HER2 Amplification:** HER2 amplification occurs in approximately 20% of breast cancers and is associated with aggressive disease. The development of HER2-targeted therapies including trastuzumab (Herceptin), pertuzumab, and T-DM1 has dramatically improved outcomes for HER2-positive breast cancer patients [15].

**RAS-MAPK Signaling:** EGFR and ERBB2 activation leads to RAS activation, which triggers the MAPK cascade (RAF→MEK→ERK). This pathway directly regulates cell cycle progression through cyclin D1 upregulation and is essential for cancer cell proliferation. RAS mutations render cells independent of upstream RTK signals [16].

**GSK3B Integration:** GSK3B integrates signals from multiple pathways including ErbB signaling. AKT, activated downstream of ErbB receptors, phosphorylates and inhibits GSK3B, leading to stabilization of proteins like cyclin D1 and c-Myc that promote cell cycle progression [17].

CADM1's interaction with multiple ErbB pathway components suggests it may serve as a scaffold protein or negative regulator. Loss of CADM1 could enhance ErbB signaling, promoting oncogenesis. This is consistent with CADM1's known tumor suppressor role.

#### 4.4.3 Function 3: Virus Receptor Activity

##### Overview:

- **Category:** GO Molecular Function
- **Term ID:** GO:0001618
- **Genes involved:** EGFR, BSG, HSPA1A, ACE2, PVR, LDLR (6 genes)
- **FDR:** 0.0002

**Description:** Virus receptor activity refers to proteins that serve as receptors for viral entry into host cells. These cell surface proteins are recognized and exploited by viruses for attachment and entry. Interestingly, many of these receptors also play important roles in normal cellular processes and, when dysregulated, in cancer development.

**Relation to Cancer:** The connection between viral receptors and cancer is multi-faceted:

**EGFR as Viral Receptor and Oncogene:** EGFR serves as a receptor or co-receptor for several viruses including cytomegalovirus and vaccinia virus. Its overexpression in cancer cells facilitates both viral entry and oncogenic signaling. Some oncolytic viruses have been engineered to exploit EGFR overexpression to selectively target and kill cancer cells while sparing normal tissues [18].

**ACE2 in COVID-19 and Cancer:** ACE2 (Angiotensin-Converting Enzyme 2) gained prominence as the receptor for SARS-CoV-2. ACE2 expression varies across tissues and has been implicated in cancer biology. Studies suggest complex relationships

between ACE2 expression, COVID-19 susceptibility, and cancer progression that depend on cancer type [19].

**BSG (CD147/Basigin):** BSG is overexpressed in many cancers and promotes tumor invasion and metastasis through induction of matrix metalloproteinases (MMPs). It also serves as an alternative receptor for SARS-CoV-2 and other viruses. BSG represents a link between viral infection mechanisms and cancer invasion processes [20].

**PVR (CD155/Poliovirus Receptor):** Beyond its original identification as a poliovirus receptor, PVR functions as an immune checkpoint molecule. It inhibits NK cell-mediated cytotoxicity through binding to TIGIT and CD96 on immune cells. PVR is upregulated in many cancers as a mechanism of immune evasion and represents an emerging immunotherapy target [21].

**LDLR (Low-Density Lipoprotein Receptor):** LDLR serves as a receptor for multiple viruses including VSV and LCMV. In cancer, LDLR expression is often upregulated to support increased cholesterol demand for rapid membrane synthesis. LDLR can be exploited for targeted drug delivery to cancer cells.

**Implications:** The enrichment of viral receptors among CADM1 interactors is particularly interesting given CADM1's role in cell-cell adhesion and as a tumor suppressor. This suggests several possibilities:

- CADM1 may regulate cell surface expression or trafficking of these receptors
- CADM1 loss in cancer may alter viral receptor expression, potentially affecting both viral susceptibility and cancer-related phenotypes
- The convergence of cell adhesion, viral entry, and cancer signaling at these molecules highlights the interconnected nature of these processes

## 4.5 Summary of Q3(c)

The enrichment analysis of CADM1-interacting proteins revealed striking convergence on cancer-related pathways:

### Key Findings:

1. **Multiple cancer types enriched:** The CADM1 interactors are significantly enriched in pathways for endometrial, breast, prostate, bladder, gastric, and colorectal cancers (all with FDR < 0.002)
2. **Core oncogenic signaling:** Strong enrichment in ErbB/EGFR signaling pathway components (EGFR, ERBB2, KRAS, HRAS, GSK3B) that are established therapeutic targets
3. **Viral receptors and cancer:** The presence of multiple viral receptor proteins highlights the complex relationship between pathogen recognition, cell adhesion, and cancer biology

**Biological Interpretation:** As a tumor suppressor that interacts with multiple oncogenic proteins, CADM1 loss during cancer progression may simultaneously activate several pro-tumorigenic pathways:

- Enhanced EGFR/ERBB2 signaling promoting proliferation
- Altered RAS/MAPK pathway activity driving growth
- Modified cell surface receptor expression affecting both immune recognition and signaling
- Disrupted cell-cell adhesion facilitating invasion and metastasis

These findings support CADM1's established role as a multi-functional tumor suppressor and suggest that its loss contributes to malignancy through effects on cell signaling, adhesion, and proliferation control.

## 5 Comparative Analysis: Q3(a+b) vs Q3(c)

Comparing the two analyses reveals complementary aspects of cancer biology:

### 5.1 Top-15 RV Genes + Biomarkers (Q3a/b)

- **Focus:** Metabolic reprogramming and hypoxia response
- **Key pathways:** Glycolysis, central carbon metabolism, HIF1A signaling
- **Major genes:** HIF1A, PKM, SLC2A1
- **Clinical relevance:** Warburg effect, metabolic targeting

### 5.2 CADM1 Interactors (Q3c)

- **Focus:** Growth factor signaling and cell surface interactions
- **Key pathways:** ErbB/EGFR signaling, multiple specific cancer types
- **Major genes:** EGFR, ERBB2, KRAS, HRAS
- **Clinical relevance:** Targeted therapies (EGFR/HER2 inhibitors)

### 5.3 Complementary Insights

The two gene sets represent different but complementary aspects of cancer biology:

1. The **metabolic focus** of Q3(a/b) highlights how cancer cells adapt their energy metabolism
2. The **signaling focus** of Q3(c) emphasizes dysregulated growth factor pathways
3. Both analyses converge on the importance of **therapeutic targeting**
4. The lack of overlap between gene sets suggests multiple parallel mechanisms in cancer progression

## 6 Conclusion

This biological enrichment analysis has revealed critical cancer-related functions of genes identified through two different approaches:

### 6.1 Key Achievements

1. **Q3(a):** Successfully created STRING interaction networks for both gene sets, revealing biologically meaningful connections
2. **Q3(b):** Identified three highly significant cancer-related functions (central carbon metabolism, TNBC glucose metabolism, HIF1A regulation) with clear therapeutic implications
3. **Q3(c):** Demonstrated that CADM1 interacts with core oncogenic signaling molecules, providing mechanistic insight into its tumor suppressor function

### 6.2 Biological Insights

- **Metabolic reprogramming** (Warburg effect, glycolysis enhancement) is a critical feature captured by the top-15 RV genes combined with biomarkers
- **Oncogenic signaling** (ErbB/EGFR, RAS/MAPK pathways) is central to CADM1-interacting proteins
- Both analyses point to **clinically actionable targets** including HIF1A, GLUT1, EGFR, and HER2
- The **convergence on cancer pathways** from different starting points validates both the Relevance Value approach and the interaction network analysis

### 6.3 Clinical Implications

The identified pathways and genes represent:

- **Established therapeutic targets:** EGFR inhibitors, HER2 inhibitors already in clinical use
- **Emerging opportunities:** Metabolic inhibitors, HIF1A inhibitors in development
- **Biomarker potential:** GLUT1, PKM2 as prognostic markers
- **Combination strategies:** Potential for combining metabolic and signaling inhibitors

### 6.4 Future Directions

This analysis suggests several avenues for future research:

- Experimental validation of CADM1 regulation of EGFR/ERBB2 signaling
- Investigation of metabolic targeting in combination with standard therapies

- Study of CADM1 restoration as a therapeutic strategy
- Analysis of these pathways in patient samples for biomarker development

## References

- [1] Semenza, G. L. (2012). Hypoxia-inducible factors in physiology and medicine. *Cell*, 148(3), 399-408. doi: 10.1016/j.cell.2012.01.021
- [2] Ancey, P. B., Contat, C., & Meylan, E. (2018). Glucose transporters in cancer—from tumor cells to the tumor microenvironment. *The FEBS Journal*, 285(16), 2926-2943. doi: 10.1111/febs.14577
- [3] Christofk, H. R., et al. (2008). The M2 splice isoform of pyruvate kinase is important for cancer metabolism and tumour growth. *Nature*, 452(7184), 230-233. doi: 10.1038/nature06734
- [4] Hanahan, D., & Weinberg, R. A. (2011). Hallmarks of cancer: the next generation. *Cell*, 144(5), 646-674. doi: 10.1016/j.cell.2011.02.013
- [5] Lanning, N. J., et al. (2017). Metabolic profiling of triple-negative breast cancer cells reveals metabolic vulnerabilities. *Cancer & Metabolism*, 5(1), 6. doi: 10.1186/s40170-017-0168-x
- [6] Guo, Q., et al. (2022). Glucose transporter 1 expression and clinical outcome in solid tumors: a systematic review and meta-analysis. *Oncotarget*, 8(10), 16875-16886. doi: 10.18632/oncotarget.15171
- [7] Zahra, K., et al. (2020). Pyruvate kinase M2 and cancer: the role of PKM2 in promoting tumorigenesis. *Frontiers in Oncology*, 10, 159. doi: 10.3389/fonc.2020.00159
- [8] Majmundar, A. J., Wong, W. J., & Simon, M. C. (2010). Hypoxia-inducible factors and the response to hypoxic stress. *Molecular Cell*, 40(2), 294-309. doi: 10.1016/j.molcel.2010.09.022
- [9] Courtney, R., et al. (2015). Cancer metabolism and the Warburg effect: the role of HIF-1 and PI3K. *Molecular Biology Reports*, 42(4), 841-851. doi: 10.1007/s11033-015-3858-x
- [10] Vander Heiden, M. G., Cantley, L. C., & Thompson, C. B. (2009). Understanding the Warburg effect: the metabolic requirements of cell proliferation. *Science*, 324(5930), 1029-1033. doi: 10.1126/science.1160809
- [11] Dellinger, T. H., et al. (2014). EGFR and HER2 in endometrial cancer: overexpression, prognostic significance, and potential therapeutic targets. *Gynecologic Oncology*, 133(1), 113-117. doi: 10.1016/j.ygyno.2014.01.024
- [12] Byron, S. A., et al. (2012). FGFR2 point mutations in 466 endometrioid endometrial tumors: relationship with MSI, KRAS, PIK3CA, CTNNB1 mutations and clinicopathological features. *PLoS ONE*, 7(2), e30801. doi: 10.1371/journal.pone.0030801

- [13] Saegusa, M., et al. (2004). Decreased expression of GSK3 $\beta$  is associated with centrosome hyperamplification and poor prognosis in endometrial carcinomas. *Oncology Reports*, 12(2), 331-337.
- [14] Yarden, Y., & Pines, G. (2012). The ERBB network: at last, cancer therapy meets systems biology. *Nature Reviews Cancer*, 12(8), 553-563. doi: 10.1038/nrc3309
- [15] Slamon, D. J., et al. (2001). Use of chemotherapy plus a monoclonal antibody against HER2 for metastatic breast cancer that overexpresses HER2. *New England Journal of Medicine*, 344(11), 783-792. doi: 10.1056/NEJM200103153441101
- [16] Downward, J. (2003). Targeting RAS signalling pathways in cancer therapy. *Nature Reviews Cancer*, 3(1), 11-22. doi: 10.1038/nrc969
- [17] McCubrey, J. A., et al. (2014). GSK-3 as potential target for therapeutic intervention in cancer. *Oncotarget*, 5(10), 2881-2911. doi: 10.18632/oncotarget.2037
- [18] Liikanen, I., et al. (2013). Oncolytic adenovirus with temozolomide induces autophagy and antitumor immune responses in cancer patients. *Molecular Therapy*, 21(6), 1212-1223. doi: 10.1038/mt.2013.51
- [19] Pinto, B. G., et al. (2020). ACE2 expression is increased in the lungs of patients with comorbidities associated with severe COVID-19. *Journal of Infectious Diseases*, 222(4), 556-563. doi: 10.1093/infdis/jiaa332
- [20] Nabeshima, K., et al. (2006). Emmprin (CD147): a cell surface inducer of matrix metalloproteinases and a potential therapeutic target for various cancers. *Clinical & Experimental Metastasis*, 23(7-8), 348-356. doi: 10.1007/s10585-006-9043-8
- [21] Whelan, S., et al. (2019). PVRIG and PVRL2 are induced in cancer and inhibit CD8+ T-cell function. *Cancer Immunology Research*, 7(2), 257-268. doi: 10.1158/2326-6066.CIR-18-0442