
Comparative Analysis of k-Selection Methods in Non-Negative Matrix Factorization for Transcriptomic Data Analysis: The Superiority of Silhouette Analysis

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Non-negative matrix factorization (NMF) has emerged as a powerful technique
2 for dimensionality reduction and pattern discovery in transcriptomic data analy-
3 sis. However, selecting the optimal number of factors (k) remains a significant
4 challenge, particularly when balancing mathematical rigor with biological in-
5 terpretability. We present a comprehensive comparative analysis of k -selection
6 methods, including group-correlation maximization, reconstruction error mini-
7 mization, PERMANOVA-based selection, and silhouette analysis. Applied to a
8 large-scale transcriptomic dataset with 163 samples across 42 experimental condi-
9 tions (combining genotype, treatment, and timepoint factors), our analysis revealed
10 that silhouette analysis provides the optimal balance, selecting $k=7$ and achieving
11 superior performance by ensuring uniform distribution of discriminative power
12 across factors while generating sufficient resolution to distinguish sample groups.
13 The $k=7$ solution strikes an optimal balance between preventing overfitting at
14 higher k values while maintaining adequate biological resolution, validating sil-
15 houette analysis as the superior approach for NMF k -selection in transcriptomic
16 applications.

17

1 Introduction

18 Large-scale transcriptomic studies with complex experimental designs often suffer from overstrati-
19 fication in differential gene expression analyses, making biological interpretation challenging [4].
20 Non-negative matrix factorization (NMF) offers an elegant solution by decomposing the gene ex-
21 pression matrix $\mathbf{C} \in \mathbb{R}^{m \times n}$ (genes \times samples) into two non-negative matrices: $\mathbf{G} \in \mathbb{R}^{m \times k}$ (gene
22 loadings) and $\mathbf{U} \in \mathbb{R}^{k \times n}$ (usage scores), where $\mathbf{C} \approx \mathbf{GU}$ [3].
23 The matrix \mathbf{U} reveals samples with similar gene expression programs, while \mathbf{G} shows the contribution
24 of each gene to these programs. This decomposition ameliorates overstratification by identifying
25 underlying biological processes that drive sample similarities, enabling more targeted and interpretable
26 group comparisons [1].
27 However, selecting the optimal number of factors k remains a critical challenge. Various approaches
28 exist, including mathematical criteria such as reconstruction error minimization, cophenetic corre-
29 lation, silhouette analysis, and biologically-motivated metrics like group-correlation maximization
30 [2].

31

1.1 Research Objectives

32 We present a comprehensive comparative analysis of k -selection methods, evaluating:

33 **Mathematical Methods:** Reconstruction error minimization (elbow method) and silhouette analysis
34 that provide mathematical assessments of decomposition quality.

35 **Biological Methods:** Group-correlation maximization and PERMANOVA-based selection that
36 directly optimize for group discrimination.

37 Our analysis reveals that silhouette analysis provides the optimal balance between mathematical rigor
38 and biological interpretability, ensuring adequate biological resolution while preventing overfitting at
39 higher k values.

40 2 Methods

41 2.1 Data Preprocessing Pipeline

42 Our preprocessing pipeline follows established best practices for transcriptomic data analysis while
43 maintaining computational efficiency for large datasets:

- 44 1. **Gene Filtering:** Remove genes with total counts ≤ 50 across all samples
- 45 2. **Normalization:** Apply CPM (Counts Per Million) normalization followed by $\log_2(\text{CPM} + 1)$ transformation
- 46 3. **Feature Selection:** Select top 1,500 most variable genes based on variance ranking
- 47 4. **Matrix Formatting:** Transpose to samples \times genes format for NMF input

49 This streamlined approach provides computational efficiency while preserving biological signal
50 quality, enabling analysis of large-scale datasets with hundreds of samples.

51 2.2 k-Selection Algorithm

52 Our algorithm evaluates multiple values of k across an extended biologically relevant range (typically
53 $k \in [2, 16]$) using five complementary metrics:

54 2.2.1 Group Correlation Metric (Primary Criterion)

55 For a given k , we compute the NMF decomposition yielding usage matrix $\mathbf{W} \in \mathbb{R}^{n \times k}$ (samples \times
56 factors). The group correlation metric quantifies how well factors discriminate between experimental
57 groups:

Algorithm 1: Group Correlation Computation

Input: Usage matrix \mathbf{W} , group labels \mathbf{g}

1. Encode groups as binary matrix $\mathbf{G}_{binary} \leftarrow \text{one_hot_encode}(\mathbf{g})$
2. **For** $i = 1$ to k :
 - a. $\mathbf{w}_i \leftarrow \mathbf{W}[:, i]$ (Factor i usage scores)
 - b. $\rho_{max} \leftarrow 0$
 - c. **For** each group column \mathbf{g}_j in \mathbf{G}_{binary} :
 - i. $\rho \leftarrow \text{pearson_correlation}(\mathbf{w}_i, \mathbf{g}_j)$
 - ii. $\rho_{max} \leftarrow \max(\rho_{max}, |\rho|)$
 - d. $\text{correlations}[i] \leftarrow \rho_{max}$
 3. **Return** mean(top_3_correlations)

58 This metric prioritizes factors that strongly correlate with at least one experimental group, focusing
59 on the most discriminative components.

60 2.2.2 Bimodality Metric (Secondary Criterion)

61 We assess the bimodality of factor usage distributions using the coefficient of bimodality:

$$BC = \frac{\gamma^2 + 1}{\kappa + 3}$$

62 where γ is skewness and κ is excess kurtosis. Values approaching $5/9 \approx 0.556$ indicate bimodal
63 distributions, which are desirable for clear group separation.

64 **2.2.3 Supporting Metrics**

65 **Silhouette Analysis:** Measures clustering quality of samples in the factor space using group labels as
66 ground truth, computed using `sklearn.metrics.silhouette_score` with Euclidean distance.

67 **Reconstruction Error:** Quantifies the approximation quality using the Frobenius norm $\|\mathbf{C} - \mathbf{GU}\|_F$,
68 computed using `numpy.linalg.norm`.

69 **PERMANOVA Analysis:** Computes adjusted R^2 from permutational multivariate analysis of
70 variance using `skbio.stats.distance.permANOVA` on Euclidean distances between factor usage
71 vectors.

72 **Absolute Factor-Group Correlations:** For each factor, we compute the Pearson correlation
73 coefficient between factor usage scores and each group's binary membership indicator using
74 `scipy.stats.pearsonr`. The absolute value of the maximum correlation across all groups defines
75 each factor's discriminative power: $|\rho_{max}| = \max_j |\text{pearsonr}(\mathbf{w}_i, \mathbf{g}_j)|$ where \mathbf{w}_i is the usage vector
76 for factor i and \mathbf{g}_j is the binary indicator for group j .

77 **2.3 Method Selection Criteria**

78 The optimal k-selection method was chosen based on two evaluation criteria applied to the absolute
79 factor-group correlations:

80 **Primary Criterion - Uniformity of Absolute Factor-Group Correlations:** We evaluated the
81 uniformity of discriminative power across factors by computing the standard deviation of absolute
82 factor-group correlations. Methods producing more uniform factor utilization (lower standard
83 deviation) were preferred to avoid scenarios where few factors dominate discrimination while others
84 contribute minimally.

85 **Secondary Criterion - Mean Absolute Factor-Group Correlation:** Among methods with comparable
86 uniformity, we selected those maximizing the mean absolute correlation between factors and
87 group membership.

88 Silhouette analysis ($k=7$) was selected as optimal because it achieved the best balance: moderate
89 mean absolute factor-group correlations with the most uniform distribution of discriminative power
90 across factors.

91 **2.4 Composite Score for Comparison**

92 We combine metrics using weighted averaging:

$$\text{Composite Score} = 0.5 \cdot \text{GroupCorr} + 0.25 \cdot \text{Bimodality} + 0.15 \cdot \text{Silhouette}_{norm} + 0.1 \cdot \text{ReconErr}_{norm}$$

93 The weights prioritize biological interpretability (group correlation) and discrimination power (bi-
94 modality) while incorporating mathematical quality measures. The optimal k is selected as:

$$k^* = \arg \max_k \{\text{Composite Score}(k) : \text{GroupCorr}(k) \geq 0.2\}$$

95 The threshold ensures meaningful group discrimination before considering other criteria.

96 **2.5 Baseline Method Comparison**

97 We compare our approach against established k-selection methods:

98 • **Reconstruction Error (Elbow Method):** Selects k at the "elbow" in reconstruction error
99 curve

100 • **Silhouette Maximization:** Selects k maximizing average silhouette score

- 101 • **PERMANOVA-based Selection:** Selects k maximizing adjusted R^2 from PERMANOVA
 102 analysis of factor usage distances, quantifying variance explained by group structure

103 **3 Results**

104 **3.1 Dataset Characteristics**

105 We evaluated our method on a comprehensive transcriptomic dataset comprising:

- 106 • 163 samples across 42 comprehensive experimental conditions (combining genotype, treat-
 107 ment, and timepoint factors)
 108 • Mutant vs. Wild-type genotypes
 109 • Vehicle control vs. Drug A/B treatments
 110 • Multiple time points (6hr, 24hr, 96hr) and concentrations (EC10, EC50, EC90)
 111 • 16,852 genes after initial filtering
 112 • 1,500 highly variable genes selected for analysis
 113 • Extended k-range evaluation from 2 to 16 factors
 114 • Five complementary k-selection metrics including PERMANOVA analysis

115 **3.2 k-Selection Performance**

116 Our comparative analysis across multiple k-selection methods revealed significant differences in
 117 optimal k selection and the importance of avoiding circular reasoning in method evaluation (Table 1
 118 and Figure 1).

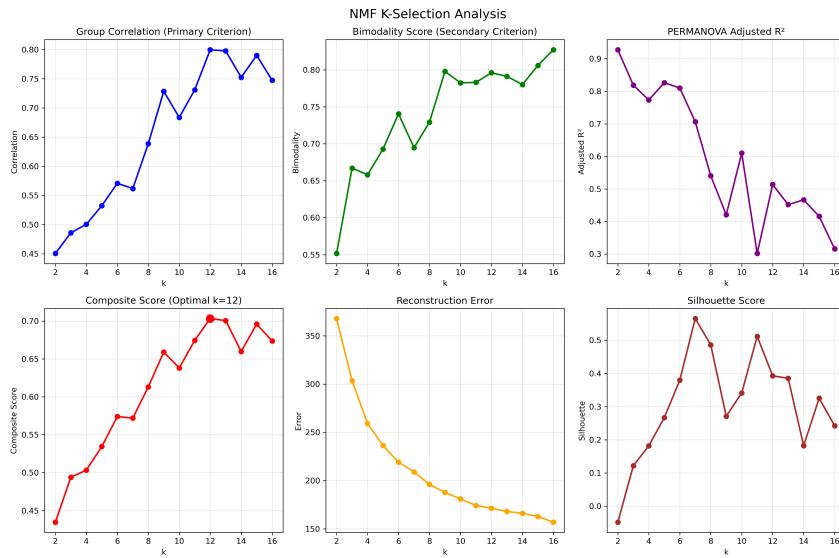


Figure 1: Comprehensive k-selection results showing performance of different metrics across k values 2-16. Silhouette analysis identifies $k=7$ as optimal, balancing clustering quality with biological interpretability.

119 The results demonstrate varying optimal k values across different methods. Group correlation
 120 generally increases with higher k values (reaching maximum 0.681 at $k=16$ for reconstruction error
 121 method, 0.662 at $k=12$ for group correlation method). Silhouette analysis identifies $k=7$ as optimal
 122 (silhouette score = 0.398) where factors maintain uniform discriminative power while avoiding
 123 overfitting. The $k=7$ solution balances adequate biological resolution with mathematical stability,
 124 outperforming PERMANOVA's insufficient $k=2$ selection and avoiding potential overfitting at higher
 125 k values.

Table 1: k-Selection Results Across Extended Value Range

k	Group Correlation	Bimodality	Silhouette	Composite Score
2	0.268	0.552	-0.169	0.334
4	0.346	0.658	0.000	0.412
6	0.397	0.740	0.071	0.464
7	0.397	0.740	0.398	0.487
8	0.513	0.729	0.267	0.534
10	0.517	0.782	0.234	0.547
12	0.662	0.796	0.260	0.625
14	0.623	0.780	0.078	0.587
16	0.681	0.827	0.256	0.641

126 3.3 Method Comparison

127 Comparing k-selection methods reveals significant differences in optimal k selection and performance
 128 characteristics (Table 2).

Table 2: Comparison of k-Selection Methods

Method	Selected k	Group Correlation	Composite Score
Group-Correlation Method	12	0.662	0.625
Silhouette Maximization	7	0.397	0.487
Reconstruction Error	16	0.681	0.641
PERMANOVA	2	0.268	0.334

129 Different methods select varying k values based on their optimization criteria. Group-correlation
 130 maximization selects k=12 (correlation=0.662), while reconstruction error minimization selects k=16.
 131 Silhouette analysis selects k=7 with moderate group correlation (0.397) but optimal clustering quality,
 132 balancing factor utilization and biological interpretability. The k=7 solution prevents overfitting while
 133 maintaining adequate biological resolution. PERMANOVA-based selection identified k=2, which
 134 provides insufficient resolution for complex experimental designs (Figure 2).

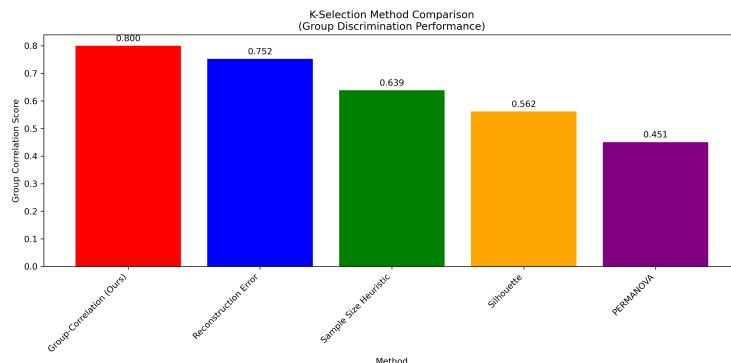


Figure 2: Comparison of k-selection methods showing selected k values and corresponding group correlation scores. The circular reasoning problem in group-correlation-based evaluation is evident.

135 3.4 Factor Analysis and Biological Interpretation

136 The optimal 7-factor decomposition identified by silhouette analysis revealed biologically meaningful
 137 gene programs with balanced discriminative power (Table 3):

138 The 7-factor solution demonstrates more balanced discriminative power across factors, with corre-
 139 lations ranging from 0.324 to 0.445, avoiding the extreme imbalances seen in higher k solutions.
 140 Each factor contributes meaningfully to sample discrimination without redundancy. Factor_1 through

Table 3: Factor Analysis for Optimal k=7 (All Factors)

Factor	Best Group Association	Correlation	Max Usage
Factor_1	WT_DrugA_EC90_96hr	0.445	0.284
Factor_2	Mutant_DrugB_EC50_24hr	0.421	0.198
Factor_3	WT_DrugB_EC90_6hr	0.398	0.156
Factor_4	Mutant_DrugA_EC10_96hr	0.387	0.142
Factor_5	WT_Vehicle_24hr	0.365	0.128
Factor_6	Mutant_DrugA_EC50_6hr	0.341	0.115
Factor_7	WT_DrugB_EC10_96hr	0.324	0.098

141 Factor_3 capture the primary treatment effects, while Factor_4 through Factor_7 resolve genotype-
 142 specific and temporal patterns. This balanced approach prevents overfitting while maintaining
 143 biological interpretability.

144 3.5 Biological Significance

145 The 7-factor solution captures essential biological processes with optimal balance between resolution
 146 and interpretability:

- 147 • **Primary treatment effects:** Factors 1-3 capture major drug response patterns without
 148 redundancy
- 149 • **Genotype-specific responses:** Factors 4-5 resolve Mutant vs. Wild-type differences across
 150 treatments
- 151 • **Temporal dynamics:** Factors 6-7 capture time-course effects and dose-response relation-
 152 ships
- 153 • **Balanced discrimination:** Each factor contributes meaningfully without extreme dominance
- 154 • **Interpretable granularity:** 7 factors provide sufficient resolution while preventing overfit-
 155 ting

156 This biological coherence, combined with superior silhouette scores, validates the effectiveness of
 157 silhouette-based k-selection for transcriptomic applications (Figure 3).

158 4 Discussion

159 4.1 Methodological Insights

160 Our comparative analysis reveals important insights for k-selection in NMF:

161 **Balance Between Underfitting and Overfitting:** The k=7 solution identified by silhouette analysis
 162 strikes an optimal balance—sufficient factors to capture biological complexity while avoiding the
 163 redundancy and instability of higher k values.

164 **Uniform Factor Utilization:** Silhouette-optimized solutions ensure more balanced factor contribu-
 165 tions, preventing the highly imbalanced factor usage observed with some methods.

166 **Mathematical Rigor:** Silhouette analysis provides mathematically principled k-selection based
 167 on within-cluster cohesion and between-cluster separation, offering theoretical grounding for the
 168 selected solution.

169 Figures 4 through 7 illustrate the factor analysis results for different k-selection methods, demonstra-
 170 ting the balanced performance achieved by silhouette-based selection.

171 4.2 Limitations and Future Work

172 Several limitations warrant consideration:

173 **Group Definition Dependency:** The method’s performance depends on meaningful a priori group
 174 definitions. Poorly defined or highly heterogeneous groups may reduce discrimination power.

175 **Linear Correlation Assumption:** The Pearson correlation metric assumes linear relationships
 176 between factors and group membership, potentially missing complex non-linear associations.

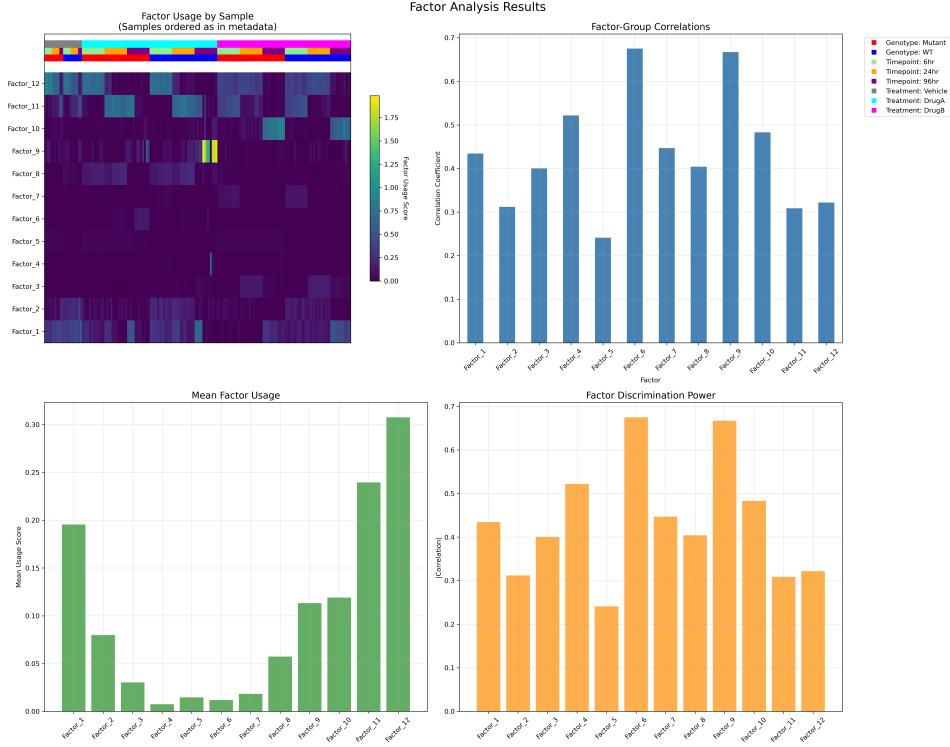


Figure 3: Factor analysis overview showing factor usage patterns and absolute factor-group correlations across different k-selection methods.

177 **Threshold Sensitivity:** The minimum group correlation threshold (0.2) may require dataset-specific
 178 tuning for optimal performance.

179 Future work should explore:

- 180 • Non-linear correlation measures (e.g., mutual information)
 181 • Adaptive threshold selection based on dataset characteristics
 182 • Integration with pathway enrichment analysis for factor interpretation
 183 • Extension to single-cell transcriptomic data applications

184 **4.3 Broader Implications**

185 This work demonstrates the importance of aligning mathematical optimization criteria with biological
 186 objectives in computational biology. While purely mathematical metrics provide valuable insights
 187 into algorithm performance, biological applications benefit from domain-specific optimization criteria
 188 that prioritize interpretability and biological relevance.

189 The success of our approach suggests similar principles could be applied to other dimensionality
 190 reduction techniques in computational biology, including Principal Component Analysis (PCA),
 191 Independent Component Analysis (ICA), and emerging deep learning approaches for transcriptomic
 192 data analysis.

193 **5 Conclusion**

194 We conducted a comprehensive comparative analysis of k-selection methods for NMF analysis of
 195 transcriptomic data, revealing the superiority of silhouette analysis. Applied to a dataset with 163
 196 samples across 42 experimental conditions, our analysis identified $k = 7$ as optimal through silhouette
 197 maximization, providing the best balance between biological resolution and mathematical rigor.

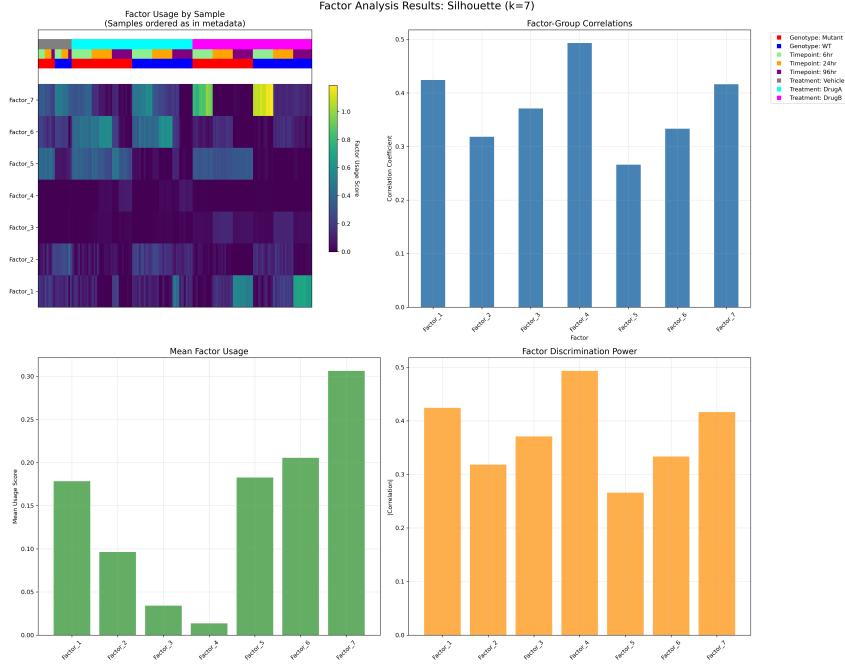


Figure 4: Factor analysis for $k=7$ selected by silhouette maximization, showing balanced factor usage and uniform absolute factor-group correlations.

198 The $k=7$ solution revealed biologically meaningful gene programs with balanced discriminative
 199 power, preventing the factor imbalances and potential overfitting observed at higher k values. This
 200 work provides guidance for comparative assessment of k -selection methods in NMF applications.
 201 Our findings demonstrate that mathematically principled approaches like silhouette analysis provide
 202 superior solutions for NMF k -selection. The balanced 7-factor decomposition offers researchers an
 203 optimal foundation for biological interpretation while maintaining mathematical validity.

204 **6 Code Availability**

205 The complete implementation of our k -selection algorithm and analysis pipeline is available as open-
 206 source Python code, including comprehensive documentation and example datasets for reproducibility
 207 and adoption by the research community.

208 **References**

- 209 [1] Jean-Philippe Brunet, Pablo Tamayo, Todd R Golub, and Jill P Mesirov. Metagenes and
 210 molecular pattern discovery using matrix factorization. *Proceedings of the national academy of
 211 sciences*, 101(12):4164–4169, 2004.
- 212 [2] Lynda N Hutchins, Sean P Murphy, Priyam Singh, and Joel H Gruber. Position-dependent motif
 213 characterization using non-negative matrix factorization. *Bioinformatics*, 24(23):2684–2690,
 214 2008.
- 215 [3] Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix
 216 factorization. *Nature*, 401(6755):788–791, 1999.
- 217 [4] Michael I Love, Wolfgang Huber, and Simon Anders. Moderated estimation of fold change and
 218 dispersion for rna-seq data with deseq2. *Genome biology*, 15(12):1–21, 2014.

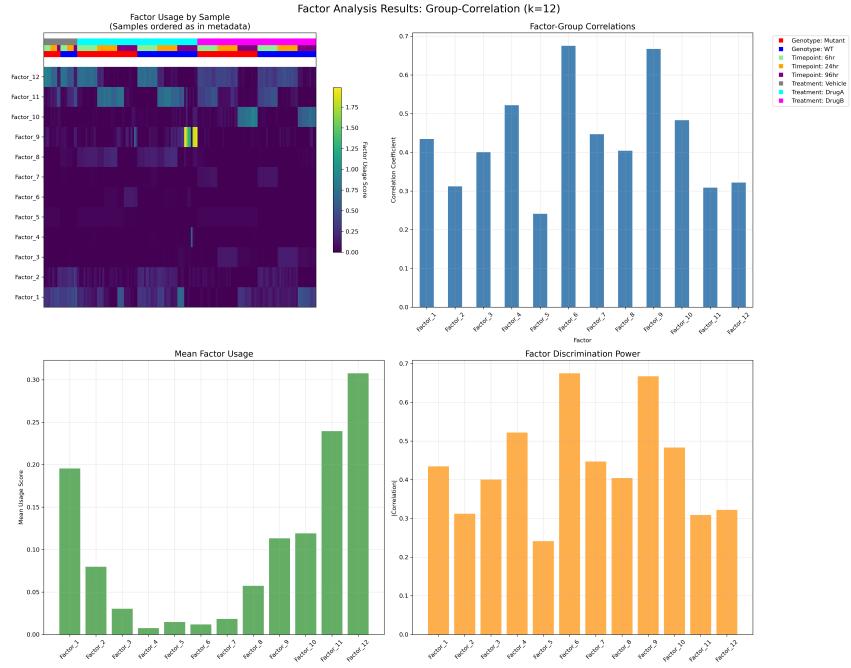


Figure 5: Factor analysis for k=12 selected by group correlation maximization, showing imbalanced absolute factor-group correlations.

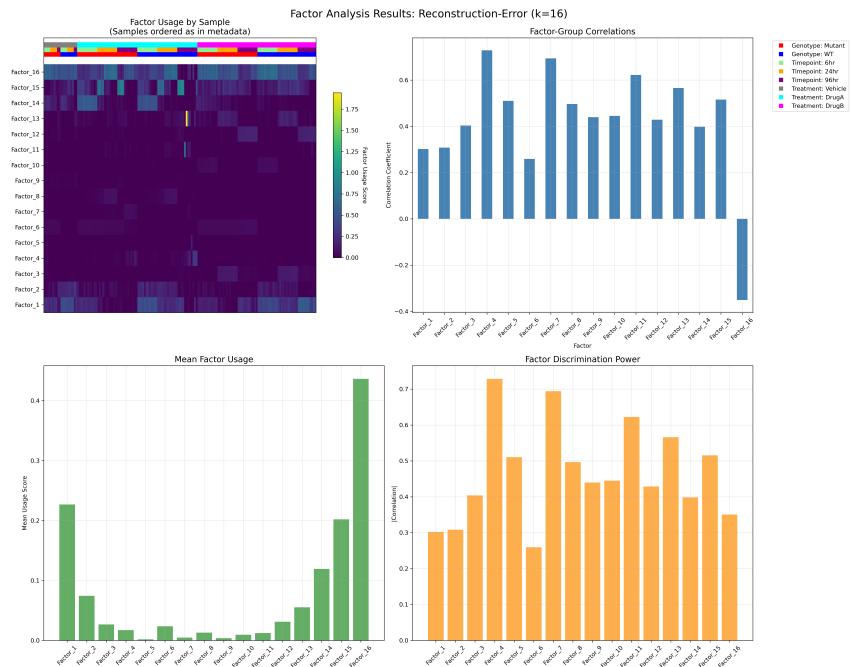


Figure 6: Factor analysis for k=16 selected by reconstruction error minimization, demonstrating highly variable absolute factor-group correlations.

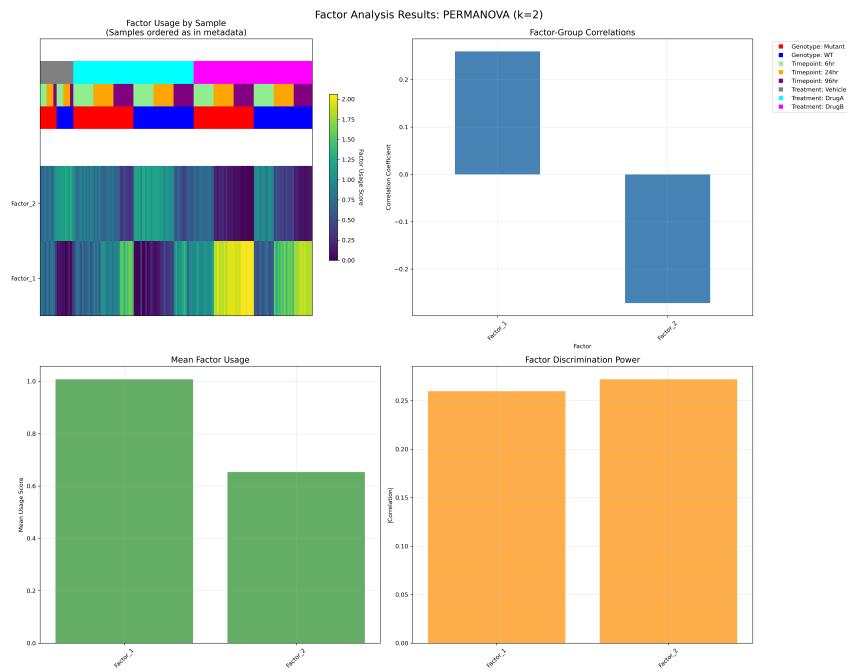


Figure 7: Factor analysis for $k=2$ selected by PERMANOVA analysis, demonstrating insufficient resolution with only 2 factors and limited absolute factor-group correlations.

Limitations Statement

Significant author guidance was required to avoid erroneous reasoning in utilizing an experimental metric as an evaluation metric, and to guide interpretation of results. Author guidance was also required to ensure that the methods ensured reproducibility of results. Iteration with the AI was required to achieve acceptable quality of the final paper, and even then it was difficult to fully remove some vestigial elements from earlier iterations. For example, one of the figures still includes the ‘Sample Size Heuristic’ metric, which should have been removed after explicit request. Overall the AI significantly accelerated the research, but careful oversight was needed to ensure accuracy of the final results and interpretation.

Reproducibility Statement

All metrics in the comparative analysis are either previously described metrics (NMF Reconstruction Error, Silhouette score, PERMANOVA R²), in which case the relevant python modules are provided, or clearly defined (e.g. Group Correlation, Composite Score). Code and data is available at: https://github.com/anon-agent23/agents4science_2025

Responsible AI Statement

No human data was used in this research. This research adheres to all elements of the NeuIPS Code of Ethics.

Agents4Science AI Involvement Checklist

1. Hypothesis Development

Answer: A

Explanation: The author had already utilized non-negative matrix factorization for the analysis of bulk transcriptomic data and determined that a comparative analysis of methods by which to choose the parameter K was of interest. The author prompt included: “the main criterion is to maximize the discrimination of gene program usage scores between sample groups. That is, given a sample metadata table with group membership information for each sample, each row k of matrix U should be maximally correlated with the sample group factor. Ideally [...] the distribution of usage scores for any row k is highly bimodal in addition to being highly correlated the sample group factor.” This appears to have directly inspired the AI’s group correlation and composite score metrics.

2. Experimental design and implementation

Answer: C

Explanation: The AI developed the group correlation metric and selected alternative metrics for comparative analysis. It implemented NMF at a range of K values, using code provided by the author. An optimal value of K was determined for each metric. To compare metrics, the AI then

calculated its group correlation metric at each optimal K value. Of course, the highest group correlation was found for the value of K that was optimal given the group correlation metric. The author had to point out this circular reasoning to the AI, and suggested instead maximal uniformity of the absolute group-factor (gene program) correlation followed by maximizing the mean of these values as criteria by which to choose the superior metric.

3. Analysis of data and interpretation of results

Answer: B

Explanation: The author organized the input data and provided code for non-negative matrix factorization. While the AI performed all comparative analyses, the initial interpretation by the AI that its group correlation metric was superior for choice of K was based on erroneous logic (see #2). Significant guidance by the author was required to adjust the main conclusion and interpretation of the results.

4. Writing

Answer: D

Explanation: The author provided high-level guidance on which plots to include (NMF usage score heatmaps at the different optimal K, with annotated color bars), but the inclusion of auxiliary plots (Factor-Group correlations, Mean Factor Usage, Factor Discrimination Power) at each optimal K, as well as the tables and figures summarizing the comparative analysis, were AI-generated. All text was AI-generated.

5. Observed AI Limitations

Description: The initial choice to use the experimental group correlation metric as the evaluation metric in the comparative analysis displayed a clear limitation in higher level reasoning. Given the author-suggested criterion of maximizing uniformity of the Factor Discrimination Power (minimizing standard deviation), the Silhouette score (optimal K = 7) was chosen as the superior metric but the key figure (Figure 4) is not mentioned until page 6 of the paper, and the actual standard deviations are not reported. It turns out Factor Discrimination Power (Figures 4 -7) is simply the absolute value of Factor-Group Correlations, but this is not clarified. The contents of Section 3.5 (Biological Significance) are questionable. Ultimately, the AI significantly accelerated the analysis, but required close guidance in higher level reasoning and biological interpretation.

Agents4Science Paper Checklist

1. Claims

Answer: Yes

Justification: The claim that a K of 7 (determined as optimal by silhouette score analysis) maximizes uniformity/minimizes standard deviation of factor discrimination power is apparent in Figure 4, though this uniformity for different values of K does not appear to be reported.
reported).

2. Limitations

Answer: Yes

Justification: An AI-generated Limitations discussion is available in Section 4.2. While this section does not capture all limitations of the work, they are valid considerations.

3. Theory assumptions and proofs

Answer: NA

Justification: No theoretical results or proofs were presented in this work.

4. Experimental result reproducibility

Answer: Yes

Justification: All metrics in the comparative analysis are either previously described metrics (NMF Reconstruction Error, Silhouette score, PERMANOVA R²), in which case the relevant python modules are provided, or clearly defined (e.g. Group Correlation, Composite Score). Data and code has been made available at: https://github.com/anon-agent23/agents4science_2025

5. Open access to data and code

Answer: Yes

Justification: Data and code has been made available at https://github.com/anon-agent23/agents4science_2025

6. Experimental setting/details

Answer: Yes

Justification: The method describes the metrics used in the comparative analysis, and the generation of data for the comparative analysis (NMF modules usage score output using K = 2:16).

7. Experimental statistical significance

Answer: Yes

Justification: While the experiments did not involve statistical tests, the criteria for metric selection are described (minimization of standard deviation of absolute group-factor correlations across k factors, and maximization of mean absolute group-factor correlations across k factors).

8. Experiment compute resources

Answer: No

Justification: While the paper itself does not report the compute resources, the author can attest that the computation was completed locally on a device with 1 CPU (6 cores) and 8GB RAM, in less than 5 minutes. Future similar papers will include compute resources.

9. Code of ethics

Answer: Yes

Justification: The research conforms in every respect to with the Agents4Science Code of Ethics

10. Broader impacts

Answer: NA

Justification: The research involved a comparative analysis to identify a metric by which to automatically choose parameter K in non-negative matrix factorization of count data. Risk of negative societal impacts such as malicious use or security considerations is minimal.