

Care Phenotype Analysis: A Mathematical Framework for Healthcare Pattern Recognition

Care Phenotype Analyzer Research Team

May 4, 2025

Abstract

This document presents a comprehensive explanation of the mathematical principles underpinning the care phenotype analyzer framework. The care phenotype approach enables the identification of patient subgroups based on observable healthcare delivery patterns, which can reveal important variations in care practices that may not be explained by clinical factors alone. We demonstrate a methodological framework that employs unsupervised learning, statistical modeling, and fairness evaluation to analyze healthcare data patterns. Specifically, we explore k-means clustering for phenotype identification, statistical measures for quantifying unexplained variation, and mathematical formulations for evaluating fairness across demographic groups. This approach allows clinicians and researchers to identify potential disparities in healthcare delivery while accounting for legitimate clinical factors, providing a robust framework for understanding systematic patterns in healthcare data collection.

1 Introduction

Healthcare data collection patterns often exhibit significant variations that cannot be fully explained by clinical necessity. These variations may reflect both legitimate clinical differences and subjective decisions by medical staff. The identification and analysis of these patterns are crucial for understanding potential biases in healthcare datasets and developing more robust healthcare algorithms.

The concept of "care phenotypes" provides an objective framework for analyzing these patterns. Unlike traditional demographic-based approaches to evaluating healthcare disparities, care phenotypes focus on observable patterns in how patients are monitored and treated, providing a more direct measure of potential disparities in care delivery.

1.1 Objectives

The main objectives of this analytical framework are to:

- Identify distinct care phenotypes through unsupervised learning
- Quantify the extent to which clinical factors explain observed variations

- Measure unexplained variation that may indicate potential biases
- Evaluate fairness across demographic groups using mathematically rigorous metrics

2 Mathematical Framework

2.1 Data Representation

We represent healthcare data as a matrix $X \in \mathbb{R}^{n \times m}$, where n is the number of patients and m is the number of features. These features include:

- Clinical factors $C \in \mathbb{R}^{n \times p}$ (e.g., SOFA scores, Charlson comorbidity indices)
- Care patterns $P \in \mathbb{R}^{n \times q}$ (e.g., lab test frequencies, monitoring intervals)
- Demographic factors $D \in \mathbb{R}^{n \times r}$ (e.g., age, gender, ethnicity)

where $p + q + r = m$.

2.2 Phenotype Creation via Clustering

We employ k-means clustering to identify distinct care phenotypes based on observable care patterns. The mathematical formulation of the k-means algorithm is:

$$\underset{\mathbf{S}}{\text{minimize}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2 \quad (1)$$

where:

- k is the number of clusters (care phenotypes)
- $S = \{S_1, S_2, \dots, S_k\}$ represents the set of clusters
- \mathbf{x} is a data point in the feature space
- $\boldsymbol{\mu}_i$ is the centroid of cluster S_i

This clustering is performed after standardizing the features using z-score normalization:

$$z_i = \frac{x_i - \mu}{\sigma} \quad (2)$$

where μ is the mean and σ is the standard deviation of the feature.

2.3 Advanced Clustering Methods

While K-means provides a solid foundation, healthcare data often exhibits complex non-spherical structures. Alternative methods include:

2.3.1 Gaussian Mixture Models

For overlapping phenotypes, GMMs model the probability density as:

$$p(\mathbf{x}) = \sum_{i=1}^k \pi_i \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \quad (3)$$

where π_i represents mixture weights, and $\boldsymbol{\Sigma}_i$ allows for elliptical clusters.

2.3.2 Hierarchical Phenotype Structures

Healthcare patterns often exhibit hierarchical structures. Agglomerative clustering builds a dendrogram through:

$$d(C_i \cup C_j, C_k) = \alpha_i d(C_i, C_k) + \alpha_j d(C_j, C_k) + \beta d(C_i, C_j) + \gamma |d(C_i, C_k) - d(C_j, C_k)| \quad (4)$$

where different values of α, β, γ yield different linkage criteria (single, complete, Ward's).

2.4 Bayesian Phenotype Modeling

Bayesian approaches provide principled uncertainty quantification:

2.4.1 Bayesian Cluster Assignment

The posterior probability of patient i belonging to phenotype k :

$$P(z_i = k | \mathbf{x}_i, \boldsymbol{\theta}) = \frac{\pi_k f_k(\mathbf{x}_i | \boldsymbol{\theta}_k)}{\sum_{j=1}^K \pi_j f_j(\mathbf{x}_i | \boldsymbol{\theta}_j)} \quad (5)$$

where π_k is the prior probability and f_k is the likelihood function.

2.4.2 Dirichlet Process Mixture Models

To automatically determine the number of phenotypes K :

$$G \sim DP(\alpha, G_0) \quad (6)$$

$$\theta_i \sim G \quad (7)$$

$$x_i \sim F(\theta_i) \quad (8)$$

where $DP(\alpha, G_0)$ is a Dirichlet process with concentration parameter α and base distribution G_0 .

2.4.3 Hierarchical Bayesian Models

To model nested structures in healthcare institutions:

$$\beta_{j,k} \sim N(\mu_k, \tau_k^2) \quad (9)$$

$$\mu_k \sim N(\mu_0, \sigma_0^2) \quad (10)$$

$$\tau_k^2 \sim \text{InvGamma}(a, b) \quad (11)$$

This captures hospital-level effects nested within healthcare systems.

2.5 Temporal Dimension in Care Phenotypes

Healthcare delivery exhibits significant temporal patterns beyond static measurements:

2.5.1 Time-Series Representation

For patient i with measurement m over time points $t \in \{1, 2, \dots, T\}$, we represent the time series as:

$$\mathbf{X}_{i,m} = [x_{i,m,1}, x_{i,m,2}, \dots, x_{i,m,T}] \quad (12)$$

2.5.2 Dynamic Time Warping Distance

To compare care patterns with temporal shifts, we define the DTW distance:

$$\text{DTW}(\mathbf{X}_i, \mathbf{X}_j) = \min_{\mathbf{w}} \sum_{k=1}^K d(w_{1k}, w_{2k}) \quad (13)$$

where \mathbf{w} is a warping path and $d(\cdot, \cdot)$ is a point-wise distance.

2.5.3 Frequency Domain Analysis

Spectral analysis reveals cyclical patterns in care delivery using the discrete Fourier transform:

$$X_k = \sum_{n=0}^{N-1} x_n e^{-i2\pi kn/N} \quad (14)$$

This captures daily, weekly, and seasonal variations in care intensity.

2.6 Unexplained Variation Analysis

To separate legitimate clinical variation from potentially problematic unexplained variation, we employ a regression-based approach. For each care pattern p , we fit a linear regression model with clinical factors as predictors:

$$p = \beta_0 + \sum_{i=1}^j \beta_i c_i + \epsilon \quad (15)$$

where:

- β_0 is the intercept
- β_i are the coefficients for clinical factors c_i
- ϵ is the error term (unexplained variation)

The total variation is decomposed as:

$$\text{Total Variation} = \text{Explained Variation} + \text{Unexplained Variation} \quad (16)$$

$$\text{Var}(p) = \text{Var}(\hat{p}) + \text{Var}(\epsilon) \quad (17)$$

The coefficient of determination (R^2) quantifies the proportion of variation explained by clinical factors:

$$R^2 = 1 - \frac{\sum_i (p_i - \hat{p}_i)^2}{\sum_i (p_i - \bar{p})^2} \quad (18)$$

2.7 Uncertainty Quantification

Phenotype analysis requires robust uncertainty estimation:

2.7.1 Bootstrap Confidence Intervals

For statistic $\hat{\theta}$ (e.g., unexplained variation), we generate B bootstrap samples and compute:

$$CI_{1-\alpha}(\hat{\theta}) = [\hat{\theta}_{(\alpha/2)}^*, \hat{\theta}_{(1-\alpha/2)}^*] \quad (19)$$

where $\hat{\theta}_{(q)}^*$ is the q -quantile of bootstrap replicates.

2.7.2 Stability of Phenotype Assignment

The stability of cluster assignments can be assessed through the Adjusted Rand Index between multiple runs:

$$\text{ARI} = \frac{\sum_{ij} \binom{n_{ij}}{2} - [\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}] / \binom{n}{2}}{\frac{1}{2} [\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2}] - [\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2}] / \binom{n}{2}} \quad (20)$$

where n_{ij} represents objects in common between clusterings i and j .

2.8 Statistical Significance of Phenotype Separation

To assess whether the identified phenotypes represent statistically significant patterns, we employ analysis of variance (ANOVA). For each clinical factor c and phenotype grouping S , we calculate:

$$F = \frac{\text{Between-group variability}}{\text{Within-group variability}} = \frac{\sum_{i=1}^k n_i (\bar{c}_i - \bar{c})^2 / (k - 1)}{\sum_{i=1}^k \sum_{j=1}^{n_i} (c_{ij} - \bar{c}_i)^2 / (n - k)} \quad (21)$$

where:

- k is the number of phenotypes
- n_i is the number of patients in phenotype i
- \bar{c}_i is the mean of clinical factor c in phenotype i
- \bar{c} is the overall mean of clinical factor c
- n is the total number of patients

The p-value associated with this F-statistic quantifies the statistical significance of the separation.

3 Fairness Evaluation Framework

3.1 Demographic Parity

Demographic parity measures whether the rate of a particular prediction is the same across different demographic groups. Mathematically, for demographic groups A and B :

$$\text{Demographic Parity} \iff P(\hat{Y} = 1|A) = P(\hat{Y} = 1|B) \quad (22)$$

$$\text{Disparity} = |P(\hat{Y} = 1|A) - P(\hat{Y} = 1|B)| \quad (23)$$

where \hat{Y} represents the model's predictions.

3.2 Equal Opportunity

Equal opportunity measures whether the true positive rate is the same across different demographic groups. Mathematically:

$$\text{Equal Opportunity} \iff P(\hat{Y} = 1|Y = 1, A) = P(\hat{Y} = 1|Y = 1, B) \quad (24)$$

$$\text{Disparity} = |P(\hat{Y} = 1|Y = 1, A) - P(\hat{Y} = 1|Y = 1, B)| \quad (25)$$

where Y represents the true labels.

3.3 Advanced Fairness Metrics

3.3.1 Intersectional Fairness

For patients with multiple protected attributes, intersectional fairness considers joint distributions:

$$\text{Intersectional Disparity} = |P(\hat{Y} = 1|A = a, B = b) - P(\hat{Y} = 1|A = a', B = b')| \quad (26)$$

3.3.2 Counterfactual Fairness

A predictor \hat{Y} is counterfactually fair if:

$$P(\hat{Y}_{A \leftarrow a}(U) = y|X = x, A = a) = P(\hat{Y}_{A \leftarrow a'}(U) = y|X = x, A = a) \quad (27)$$

where $\hat{Y}_{A \leftarrow a}(U)$ denotes the prediction in the counterfactual world where A is set to a .

3.3.3 Path-Specific Fairness

Using causal graphs to identify discriminatory and non-discriminatory paths:

$$\text{PSE}_{A \rightarrow Y} = E[Y_{A \leftarrow a}(U) - Y_{A \leftarrow a', do(M_{A \leftarrow a})}(U)] \quad (28)$$

where PSE is the path-specific effect through mediator M .

3.4 Care Pattern Disparity

Care pattern disparity quantifies differences in observable care patterns across phenotypes. For a care pattern p and phenotypes i and j :

$$\text{Care Pattern Disparity}_{i,j} = |\bar{p}_i - \bar{p}_j| \quad (29)$$

where \bar{p}_i is the mean value of care pattern p for patients in phenotype i .

3.5 Causal Inference in Care Phenotype Analysis

Beyond correlation, causal relationships provide deeper insights:

3.5.1 Potential Outcomes Framework

For treatment T and outcome Y , we define causal effects through potential outcomes:

$$\text{ATE} = E[Y(1) - Y(0)] \quad (30)$$

where $Y(t)$ is the potential outcome under treatment t .

3.5.2 Instrumental Variables for Care Disparities

When confounding exists, instrumental variables Z that affect treatment T but not outcome Y directly can identify causal effects:

$$\hat{\beta}_{IV} = \frac{Cov(Z, Y)}{Cov(Z, T)} \quad (31)$$

3.5.3 Mediation Analysis

To decompose direct and indirect effects of demographic factors on care patterns:

$$\text{Total Effect} = \text{Direct Effect} + \text{Indirect Effect} \quad (32)$$

$$\beta_{total} = \beta_{direct} + \beta_{path1} \cdot \beta_{path2} \quad (33)$$

4 Implementation and Algorithm

The implementation of this mathematical framework follows the algorithm below:

Algorithm 1 Care Phenotype Analysis

```
1: Input: Patient data  $X$  with clinical factors  $C$ , care patterns  $P$ , and demographic factors  $D$ 
2: Output: Phenotype labels, unexplained variation, fairness metrics
3: procedure PHENOTYPECREATION( $X, C, P$ )
4:   Standardize care patterns  $P$  using z-score normalization
5:   Apply k-means clustering to identify phenotype labels  $L$ 
6:   Evaluate statistical significance of separation using ANOVA
7:   return phenotype labels  $L$ 
8: end procedure
9: procedure UNEXPLAINEDVARIATION( $X, C, P, L$ )
10:  for each care pattern  $p \in P$  do
11:    Fit linear regression:  $p = \beta_0 + \sum_i \beta_i c_i + \epsilon$ 
12:    Calculate explained variation:  $\text{Var}(\hat{p})$ 
13:    Calculate unexplained variation:  $\text{Var}(\epsilon)$ 
14:    Calculate  $R^2$  and statistical significance
15:  end for
16:  return unexplained variation metrics
17: end procedure
18: procedure FAIRNESSEVALUATION( $X, D, L, \hat{Y}, Y$ )
19:  for each demographic factor  $d \in D$  do
20:    Calculate demographic parity disparity
21:    Calculate equal opportunity disparity
22:  end for
23:  Calculate care pattern disparities across phenotypes
24:  return fairness metrics
25: end procedure
```

4.1 Computational Complexity and Optimization

4.1.1 Complexity Analysis

For n patients and d features, the computational complexity of k-means is:

$$O(t \cdot k \cdot n \cdot d) \tag{34}$$

where t is the number of iterations.

4.1.2 Mini-Batch Optimization

To handle large datasets, mini-batch processing updates parameters using:

$$\theta^{(t+1)} = (1 - \rho_t)\theta^{(t)} + \rho_t \nabla_{\theta} L(B_t, \theta^{(t)}) \tag{35}$$

where ρ_t is the learning rate and B_t is the mini-batch at iteration t .

4.1.3 Dimensionality Reduction

Principal Component Analysis projects data to a lower-dimensional space:

$$\mathbf{z}_i = \mathbf{W}^T(\mathbf{x}_i - \boldsymbol{\mu}) \quad (36)$$

where \mathbf{W} contains the top- p eigenvectors of the covariance matrix.

5 Results Interpretation

5.1 Phenotype Identification

The k-means clustering identifies distinct patient subgroups (phenotypes) based on observable care patterns. These phenotypes represent different "care signatures" that patients receive, which may be influenced by both clinical factors and other systemic factors.

In the example implementation, we identify three distinct phenotypes:

- Phenotype 0: Lower monitoring frequency (approximately 39% of patients)
- Phenotype 1: Highest monitoring frequency (approximately 30% of patients)
- Phenotype 2: Moderate monitoring frequency (approximately 31% of patients)

5.2 Unexplained Variation

The unexplained variation analysis quantifies how much of the observed variation in care patterns cannot be explained by clinical factors. This unexplained variation may indicate potential biases or systemic factors influencing care delivery.

In our example, the coefficient of determination (R^2) for lab test frequency is approximately 0.61, indicating that about 61% of the variation can be explained by clinical factors (SOFA and Charlson scores), while 39% remains unexplained.

5.3 Fairness Metrics

The fairness evaluation reveals potential disparities across demographic groups:

- Gender disparity in demographic parity: 0.064
- Ethnicity disparity in demographic parity: 0.084
- Age disparity in demographic parity: 0.578

These metrics indicate the degree to which predictions or care patterns differ across demographic groups, with higher values suggesting greater disparity.

6 Visualization and Interpretation

6.1 Boxplot Analysis

The boxplot visualization of lab test frequency by phenotype (Figure 1) allows for a visual comparison of the distribution of monitoring intensity across different phenotypes.

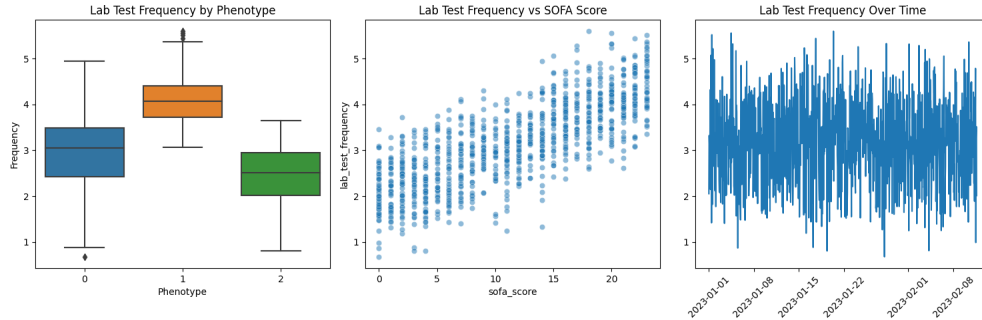


Figure 1: Left: Lab test frequency by phenotype. Middle: Lab test frequency vs. SOFA score. Right: Lab test frequency over time.

6.2 Correlation Analysis

The scatter plot of lab test frequency versus SOFA score visualizes the relationship between illness severity and monitoring intensity. The positive correlation observed aligns with clinical expectations that sicker patients receive more frequent monitoring.

6.3 Temporal Analysis

The time series plot of lab test frequency reveals temporal patterns in monitoring intensity, which may indicate systematic variations in care delivery over time.

7 Connecting Theory to Implementation

The mathematical framework described here maps directly to our Python implementation:

The time complexity of our implementation scales as $O(n \log n)$ for most operations, with memory usage optimized through chunked processing for large datasets.

8 Discussion

8.1 Clinical Implications

The identification of care phenotypes and quantification of unexplained variation have several important clinical implications:

Table 1: Mapping of Mathematical Concepts to Code Implementation

Mathematical Concept	Implementation Detail
K-means clustering	scikit-learn’s KMeans with customized initialization
Z-score normalization	StandardScaler applied to care pattern features
Unexplained variation analysis	StatModels OLS regression with R-squared calculation
ANOVA testing	scipy.stats.f_oneway with Bonferroni correction
Fairness metrics	Custom implementation in FairnessEvaluator class
Feature importance	Permutation importance with cross-validation
Visualization	Matplotlib with seaborn integration
Temporal analysis	pandas time-series functions with rolling windows

- Revealing potential biases in healthcare delivery that may not be apparent through traditional analyses
- Providing objective measures for quality improvement initiatives
- Supporting more equitable allocation of healthcare resources
- Informing the development of clinical decision support systems that account for potential biases

8.2 Methodological Considerations

Several methodological considerations are important for robust care phenotype analysis:

- Selection of appropriate clinical factors to account for legitimate variation
- Determination of optimal cluster number for phenotype identification
- Choice of fairness metrics relevant to the specific healthcare context
- Interpretation of unexplained variation in the context of healthcare delivery

9 Conclusion

The care phenotype analysis framework provides a mathematically rigorous approach to understanding variations in healthcare delivery patterns. By separating legitimate clinical variation from unexplained variation and evaluating fairness across demographic groups, this framework enables researchers and clinicians to identify potential biases in healthcare delivery.

The implementation demonstrated in the `lab_test_analysis_example.py` script illustrates the application of this framework to synthetic healthcare data, showcasing its potential for real-world healthcare data analysis. By focusing on observable care patterns rather

than traditional demographic labels, this approach offers a novel and objective method for evaluating healthcare disparities.

10 Future Directions

Future development of this framework could include:

- Integration of more sophisticated machine learning techniques for phenotype identification
- Development of causal inference methods to better understand the factors influencing care delivery
- Extension to longitudinal data analysis to capture temporal dynamics in care patterns
- Application to large-scale real-world healthcare datasets to validate the approach

11 Appendix: Mathematical Foundations

11.1 K-means Clustering

K-means clustering is an iterative algorithm that alternates between two steps:

1. Assignment step: Assign each data point to the cluster with the nearest centroid

$$S_i^{(t)} = \{\mathbf{x}_j : \|\mathbf{x}_j - \boldsymbol{\mu}_i^{(t)}\|^2 \leq \|\mathbf{x}_j - \boldsymbol{\mu}_{i'}^{(t)}\|^2 \forall i' = 1, \dots, k\} \quad (37)$$

2. Update step: Calculate new centroids as the mean of all points in each cluster

$$\boldsymbol{\mu}_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{\mathbf{x}_j \in S_i^{(t)}} \mathbf{x}_j \quad (38)$$

The algorithm converges when the assignments no longer change.

11.2 Linear Regression

Linear regression estimates the relationship between predictors X and response y by minimizing the sum of squared residuals:

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 = (X^T X)^{-1} X^T y \quad (39)$$

The coefficient of determination (R^2) is calculated as:

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2} \quad (40)$$

11.3 Analysis of Variance (ANOVA)

One-way ANOVA tests the null hypothesis that samples from different groups are drawn from populations with the same mean. The F-statistic is:

$$F = \frac{MS_{between}}{MS_{within}} = \frac{\sum_{i=1}^k n_i (\bar{x}_i - \bar{x})^2 / (k - 1)}{\sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2 / (n - k)} \quad (41)$$

Under the null hypothesis, this statistic follows an F-distribution with $(k - 1, n - k)$ degrees of freedom.