

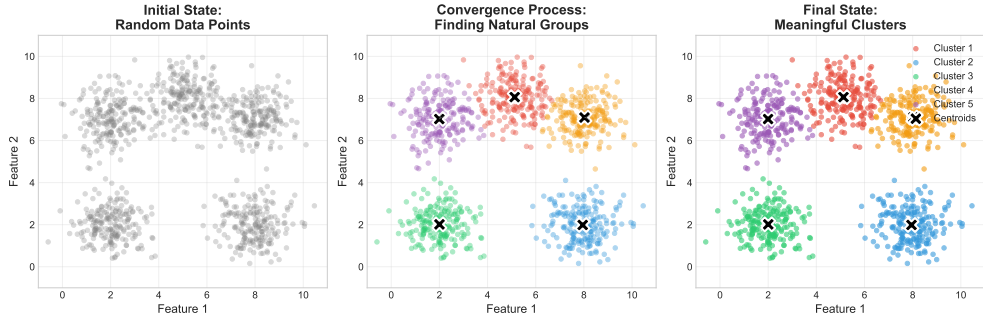
Machine Learning for Smarter Innovation

Week 1: Foundations & Clustering

Augmenting the Empathize Phase with ML

BSc Course in AI-Enhanced Innovation

The Convergence Flow: From Chaos to Clarity



The Convergence Flow: Order from Chaos
Watch 5000 data points self-organize into meaningful clusters

The Innovation Challenge

Why Traditional Design Needs AI Enhancement

Traditional Design Limits

- **Scale:** Can interview 50 users, not 50,000
- **Speed:** Months for insights
- **Bias:** Designer's perspective dominates
- **Patterns:** Miss hidden connections
- **Iteration:** Slow feedback loops

AI-Enhanced Innovation

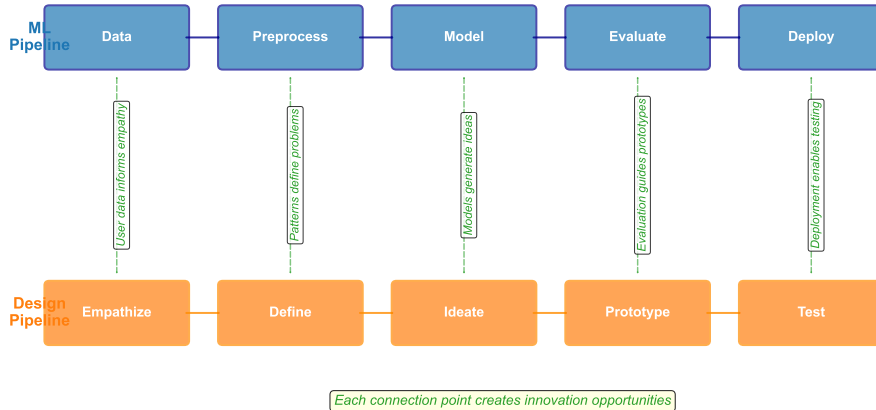
- **Scale:** Analyze millions of data points
- **Speed:** Real-time insights
- **Objectivity:** Data-driven discovery
- **Patterns:** Find non-obvious relationships
- **Iteration:** Continuous learning

The Promise: 100x more insights, 10x faster innovation

The Dual Pipeline

Where ML Meets Design Thinking

The Convergence: ML Meets Design Thinking



The Dual Pipeline (Continued)

Understanding Both Worlds

ML Pipeline

Data → Preprocess → Model → Evaluate → Deploy

- Collect user behavior
- Clean and transform
- Train algorithms
- Validate accuracy
- Scale to production

Design Pipeline

Empathize → Define → Ideate → Prototype → Test

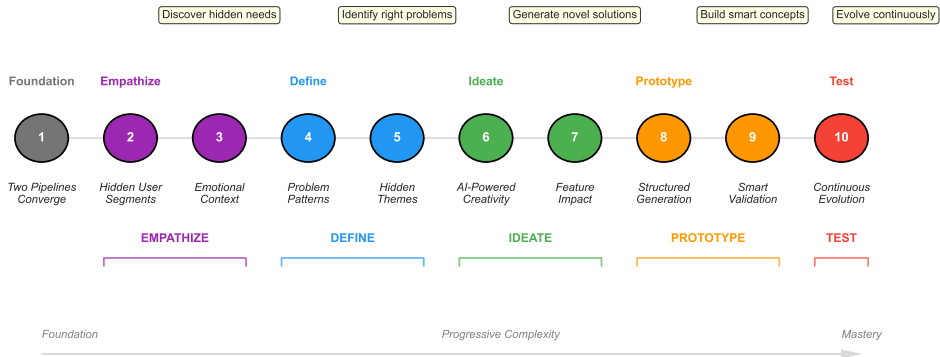
- Understand users
- Frame problems
- Generate solutions
- Build concepts
- Validate with users

Integration = Innovation at Scale

Your Innovation Journey

10 Weeks to AI-Powered Design Mastery

10-Week Innovation Journey



Your Innovation Journey (Continued)

What You'll Master in Each Stage

Stage	Weeks	Innovation Unlocked
Empathize	1-2	Discover hidden user needs at scale
Define	3-4	Identify the right problems to solve
Ideate	5-6	Generate novel solutions with AI
Prototype	7-8	Build smart, adaptive concepts
Test	9-10	Evolve through continuous learning

This Week: Clustering for Deep User Understanding

Week 1: Clustering for Empathy

From Random Data to User Understanding

What We'll Learn:

- How clustering reveals user segments
- K-means algorithm fundamentals
- Finding the optimal number of clusters
- Quality metrics for validation
- Advanced clustering techniques

Design Applications:

- Create data-driven personas
- Map user journeys by segment
- Identify pain points systematically
- Prioritize design efforts
- Scale empathy to thousands

Goal: Transform data points into human insights

What is Clustering?

Finding Natural Groups in Data

From Chaos to Clarity Through Clustering



Clustering Finds:

- Natural groupings
- Similar behaviors
- Hidden segments
- Pattern relationships

Key Insight:

Users who behave similarly likely have similar needs

K-Means: The Workhorse Algorithm

How It Organizes Your Users

The Process:

- 1 Choose K (number of clusters)
- 2 Place K random centroids
- 3 Assign points to nearest centroid
- 4 Move centroids to cluster mean
- 5 Repeat until stable

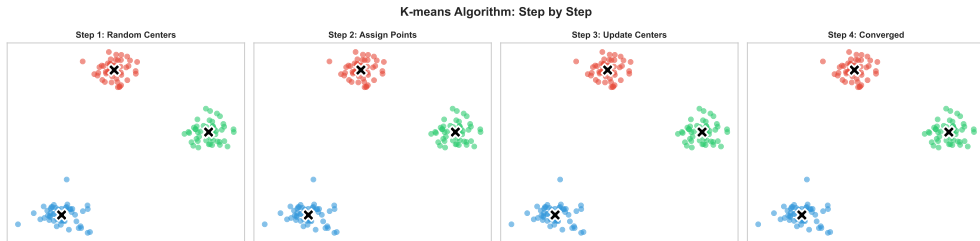
Strengths:

- Fast and scalable
- Easy to understand
- Works well for spherical clusters



K-Means in Action

Step-by-Step Convergence

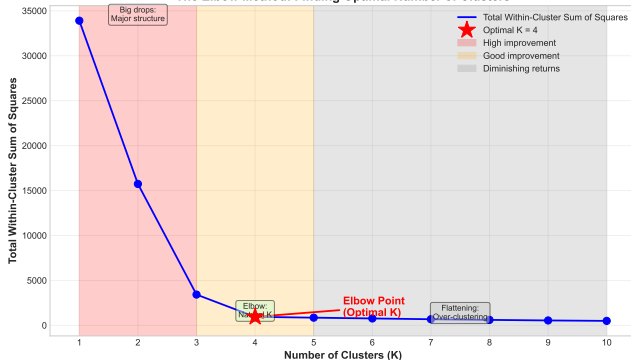


Iteration 1 → Iteration 3 → Iteration 5 → **Converged**

How Many Clusters?

The Elbow Method

The Elbow Method: Finding Optimal Number of Clusters



Finding the Elbow:

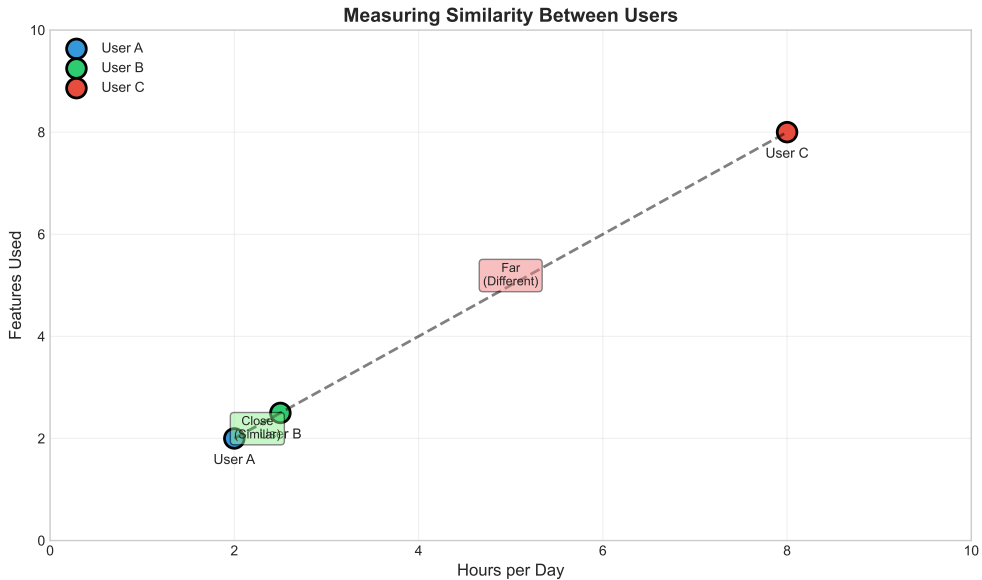
- Plot inertia vs K
- Look for the “elbow”
- Balance between:
 - Too few: Mixed groups
 - Too many: Overfitting

Optimal K = 5

Best trade-off between simplicity and accuracy

Distance Metrics

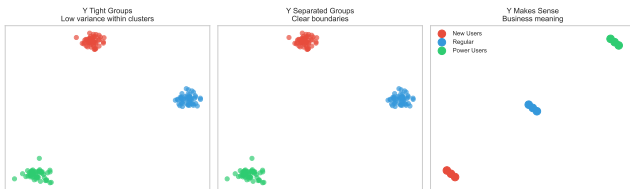
How We Measure Similarity



Cluster Quality Metrics

How Good Are Your Groups?

Three Checks for Good Clusters



Silhouette Score:

- Ranges from -1 to +1
- Higher = better separation
- Our score: **0.73**

What it measures:

- Within-cluster cohesion
- Between-cluster separation
- Overall cluster validity

0.73 = Strong clusters!

Beyond K-Means: DBSCAN

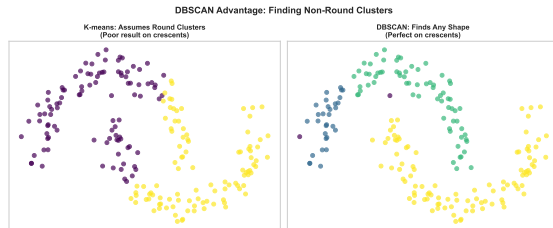
Finding Arbitrary Shaped Clusters

DBSCAN Advantages:

- No need to specify K
- Finds arbitrary shapes
- Identifies outliers
- Handles noise well

Perfect for:

- Non-spherical patterns
- Varying densities
- Outlier detection
- Exploratory analysis

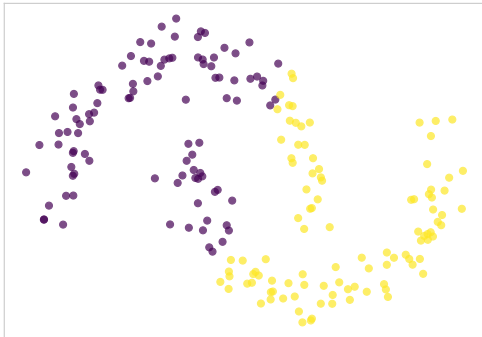


DBSCAN: Complex Patterns

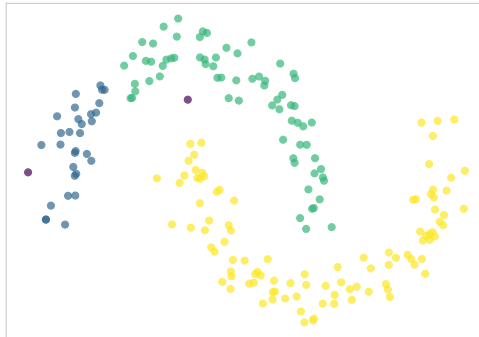
When K-Means Isn't Enough

DBSCAN Advantage: Finding Non-Round Clusters

K-means: Assumes Round Clusters
(Poor result on crescents)



DBSCAN: Finds Any Shape
(Perfect on crescents)

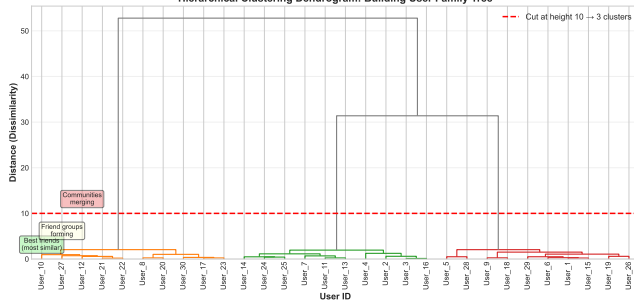


K-Means: Forces spherical shapes — DBSCAN: Finds natural boundaries

Hierarchical Clustering

Building a Tree of Relationships

Hierarchical Clustering Dendrogram: Building User Family Tree



Dendrogram Benefits:

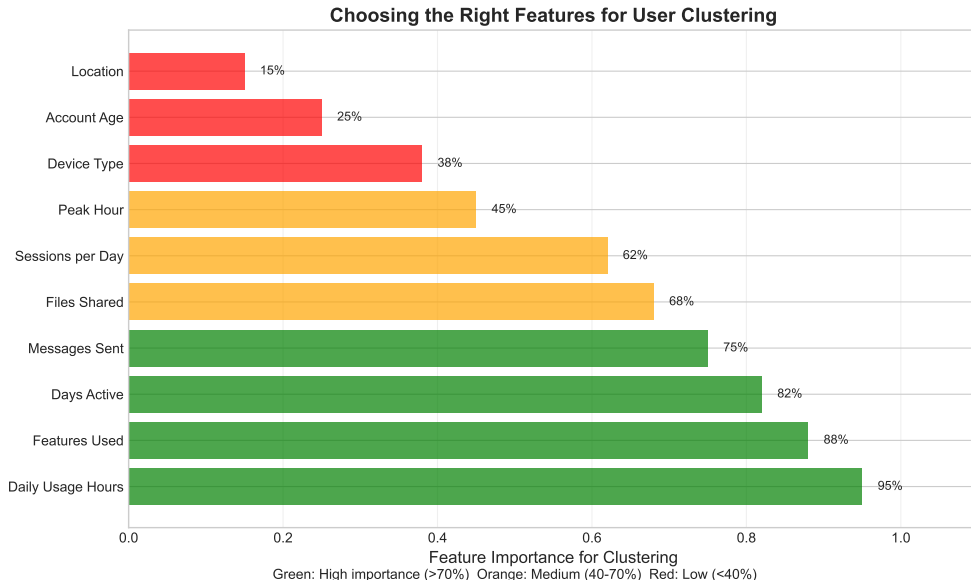
- Shows cluster hierarchy
- Multiple granularities
- Natural relationships
- No preset K needed

Cut the tree at any level:

- High cut = Few clusters
- Low cut = Many clusters
- Choose based on needs

What Drives the Clusters?

Feature Importance Analysis



From Data Points to Human Understanding

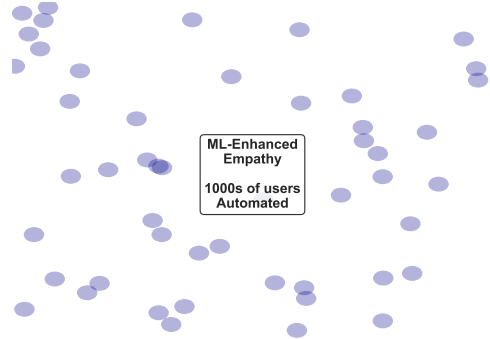
Bridging the Technical-Human Gap

Scaling Empathy with Machine Learning

Traditional Approach



ML-Enhanced Approach



Each cluster represents real human needs

AI-Generated User Personas

Data-Driven Character Development

Data-Driven Persona Cards

Power Paula

Age: 32

Role: Manager

Usage: 7h/day

Regular Rob

Age: 28

Role: Developer

Usage: 4h/day

Casual Carl

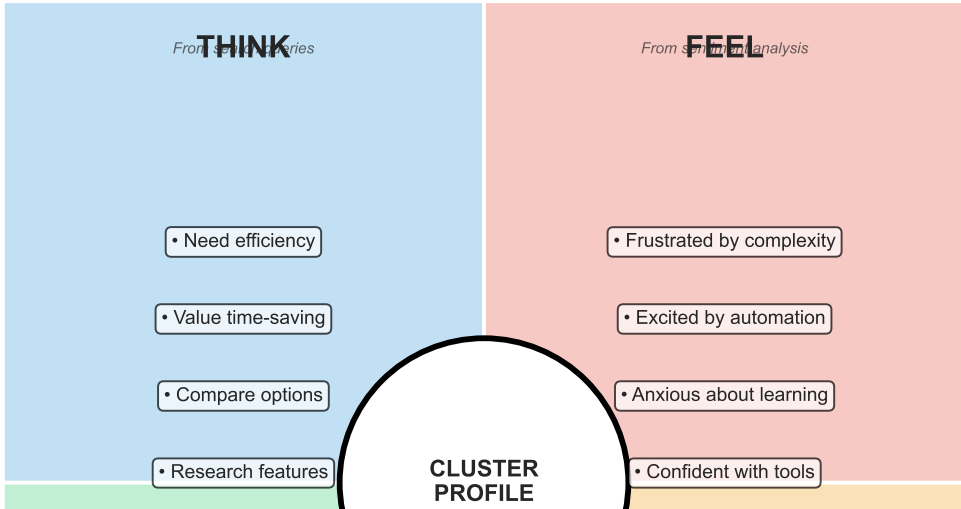
Age: 24

Role: Student

Usage: 1h/day

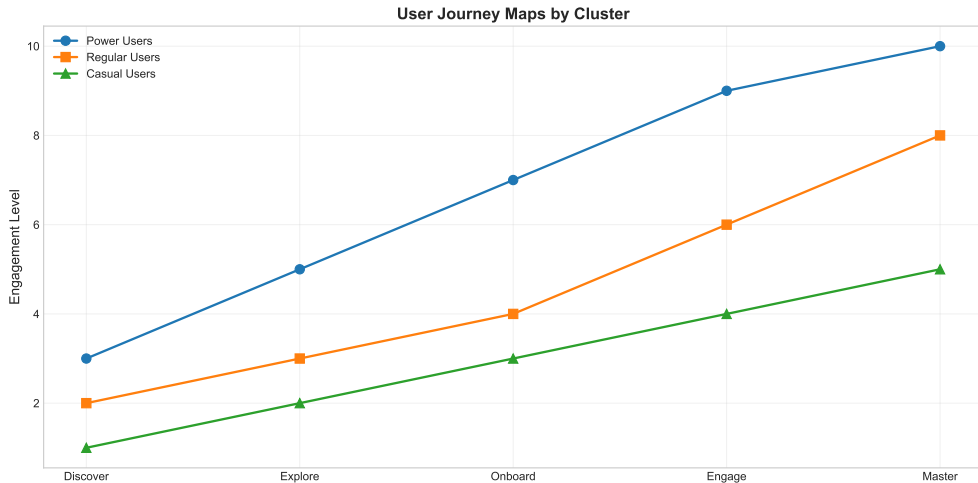
Power Users — Casual Browsers — Price-Conscious — Feature Seekers — New Users

Empathy Map: Data-Driven User Understanding



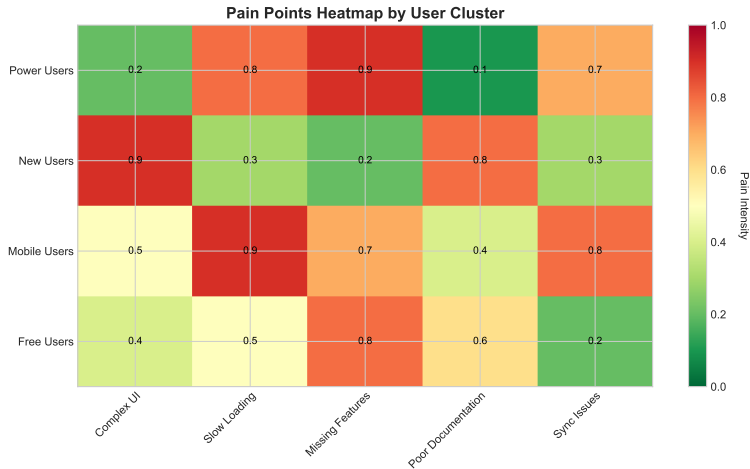
Different Journeys for Different Clusters

Personalized Path Understanding



Pain Points by Cluster

Where Each Segment Struggles



Key Findings:

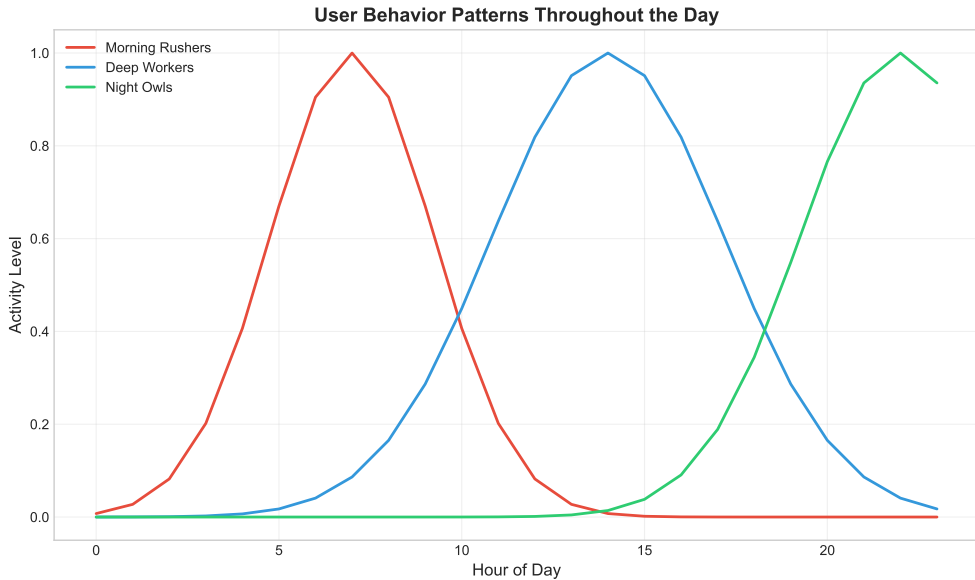
- New users: Onboarding
- Power users: Speed
- Casual: Complexity
- Price-conscious: Value

Design implication:

One solution won't fit all!

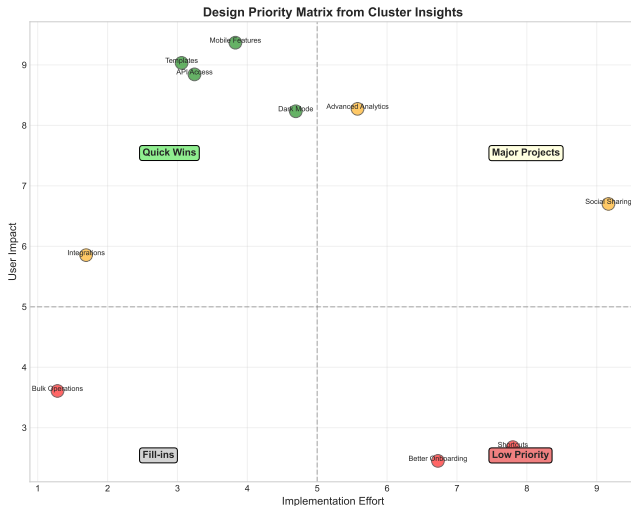
Behavioral Patterns Revealed

What Clusters Tell Us About Usage



Design Priority Matrix

Where to Focus Your Efforts



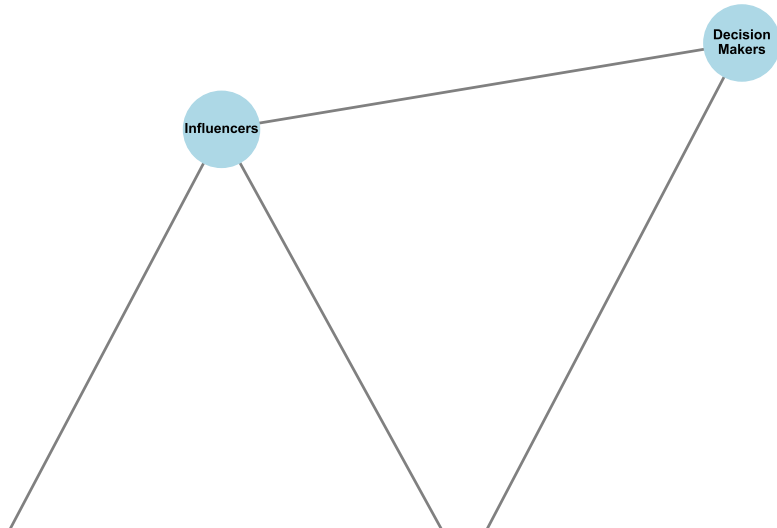
Priority Quadrants:

- **High Impact + High Effort**
Strategic initiatives
- **High Impact + Low Effort**
Quick wins
- **Low Impact + Low Effort**
Fill-ins
- **Low Impact + High Effort**
Avoid

Understanding Stakeholder Connections

Network Analysis of User Relationships

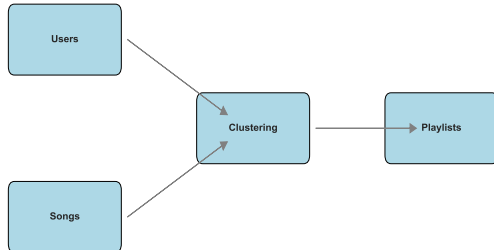
Stakeholder Network from Cluster Analysis



Case Study: Spotify's Clustering Success

Real-World Application

Spotify's Discover Weekly: Clustering in Action



Spotify Uses Clustering For:

- Music taste profiles
- Discover Weekly playlists
- User segmentation
- Recommendation engine

Results:

- Personalized experience
- Increased engagement
- Better retention
- Discovery of new artists

Key Takeaways

What We've Learned

Technical Skills

- K-means clustering algorithm
- Choosing optimal K with elbow method
- Silhouette scores for validation
- DBSCAN for complex shapes
- Hierarchical clustering

Design Applications

- Data-driven personas
- Segment-specific journeys
- Pain point identification
- Priority matrices
- Scaled empathy

Clustering transforms data into actionable user insights

Your Turn: Practice Exercise

Apply What You've Learned

Exercise: Segment Your Users

Scenario: You have data from 1000 app users including:

- Usage frequency
- Feature preferences
- Time spent
- Purchase history

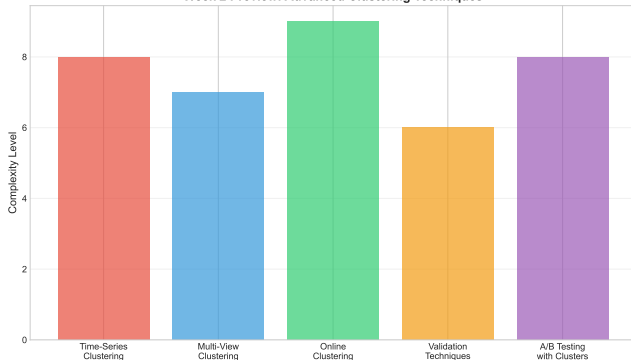
Tasks:

- 1 Choose appropriate features for clustering
- 2 Determine optimal number of clusters
- 3 Interpret what each cluster represents
- 4 Create one persona per cluster
- 5 Identify key pain points for each segment

Next Week: Advanced Clustering

Going Deeper into User Understanding

Week 2 Preview: Advanced Clustering Techniques



Week 2 Topics:

- Density-based clustering
- Gaussian mixture models
- Clustering validation
- Feature engineering
- Real-time clustering

Design Focus:

- Dynamic personas
- Evolving segments
- Predictive empathy
- Micro-segmentation

Technical Resources

Papers:

- MacQueen, J. (1967). K-means
- Ester et al. (1996). DBSCAN
- Rousseeuw (1987). Silhouettes

Tools:

- scikit-learn clustering
- Orange data mining
- KNIME analytics

Design Resources

Books:

- "Design Thinking" - Tim Brown
- "Sprint" - Jake Knapp
- "Lean UX" - Jeff Gothelf

Applications:

- Miro (journey mapping)
- Figma (persona creation)
- Optimal Workshop

Questions? Let's discuss!

Objective Function (Inertia):

$$J = \sum_{i=1}^n \sum_{j=1}^k w_{ij} ||x_i - \mu_j||^2$$

Where:

- n = number of data points
- k = number of clusters
- $w_{ij} = 1$ if x_i belongs to cluster j , 0 otherwise
- μ_j = centroid of cluster j

Update Rules:

- 1 Assignment: $c^{(i)} = \arg \min_j ||x^{(i)} - \mu_j||^2$
- 2 Update: $\mu_j = \frac{1}{|S_j|} \sum_{i \in S_j} x^{(i)}$

Appendix: Distance Metrics

Mathematical Definitions

Euclidean Distance:

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Manhattan Distance:

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

Minkowski Distance:

$$d(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}$$

Cosine Similarity:

$$\cos(\theta) = \frac{x \cdot y}{||x|| \cdot ||y||}$$

Jaccard Distance:

$$J(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

Mahalanobis Distance:

$$d(x, y) = \sqrt{(x - y)^T S^{-1} (x - y)}$$

Silhouette Score for point i :

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

Where:

- $a(i)$ = average distance to points in same cluster
- $b(i)$ = average distance to points in nearest neighbor cluster

Interpretation:

- $s(i) \approx 1$: Well clustered
- $s(i) \approx 0$: On border between clusters
- $s(i) \approx -1$: Misclassified

Overall Score:

$$S = \frac{1}{n} \sum_{i=1}^n s(i)$$

Appendix: PCA for Cluster Visualization

Dimensionality Reduction

User Clusters Visualization (PCA Reduced from 10D to 2D)



PCA Process:

- 1 Standardize data
- 2 Compute covariance matrix
- 3 Find eigenvectors/values
- 4 Select top 2 components
- 5 Transform data

Variance Explained:

- PC1: 45.2%
- PC2: 28.7%
- Total: 73.9%

Key Parameters:

- ϵ (eps): Maximum distance between points
- MinPts: Minimum points to form dense region

Point Classification:

- **Core point:** Has \geq MinPts within ϵ
- **Border point:** Within ϵ of core point
- **Noise point:** Neither core nor border

Algorithm Steps:

- 1 Find all core points
- 2 Form clusters from core points within ϵ
- 3 Assign border points to clusters
- 4 Mark remaining as noise

Appendix: Implementation Guidelines

Practical Considerations

Data Preparation

- Standardize features
- Handle missing values
- Remove outliers (if needed)
- Feature selection/engineering
- Consider scaling methods

Validation Methods

- Silhouette score
- Davies-Bouldin index
- Calinski-Harabasz score
- Visual inspection
- Domain expert review

Algorithm Selection

- K-means: Spherical, similar size
- DBSCAN: Arbitrary shapes
- Hierarchical: Nested structure
- GMM: Overlapping clusters

Common Pitfalls

- Not scaling features
- Wrong distance metric
- Ignoring outliers
- Over-clustering
- Forcing clusters