

Machine Learning for Smarter Innovation

Week 2: Clustering for Deep Empathy

BSc Course in AI-Enhanced Innovation

Understanding Users Through Data-Driven Segmentation

Today's Journey: From Data to Deep Understanding

- 1 Foundation: Advanced Clustering in the Innovation Diamond
- 2 Technical Deep Dive: Clustering Algorithms
- 3 Design Integration: From Data to Empathy
- 4 Practice: Real-World Application

Transform data points into human insights

Part 1: Foundation

Advanced Pattern Discovery in the Innovation Journey

The Innovation Diamond: Week 2 Context

Building on Week 1's Foundation with Advanced Techniques

`charts/innovation_diamond_complete.pdf`

Where We Are: Week 2 in the Innovation Journey

Advanced Clustering & Empathy - Deepening Pattern Discovery

10-Week Overview

Weeks 1-3: Empathize

- Week 1: Basic clustering
- **Week 2: Advanced clustering** ←
- Week 3: NLP & emotional context

Week 4: Define

- Classification & problem framing

Week 5: Ideate

- Topic modeling & idea generation

Weeks 6-10: Prototype, Test, Iterate

Week 2 Learning Goals

By the end of today:

- Master 4 clustering algorithms
- Choose right technique for problem
- Handle complex data patterns
- Build multi-faceted personas
- Understand trade-offs
- Apply to real innovation data

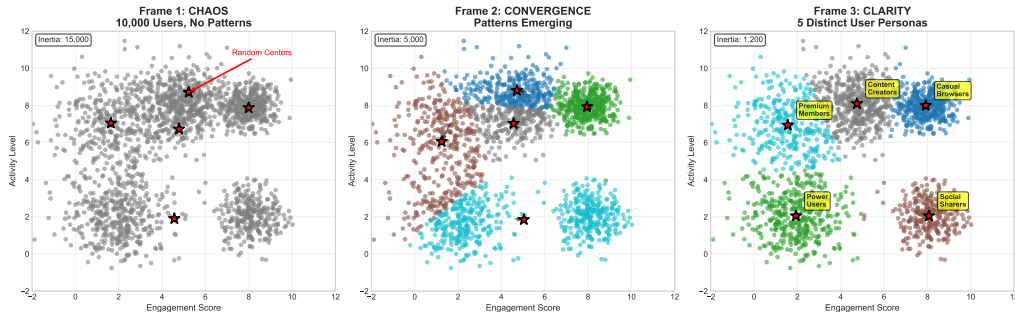
Diamond Connection:

Advanced techniques reveal innovation patterns that basic methods miss

Week 2 deepens your pattern discovery toolkit for the Innovation Diamond journey

From Chaos to Clarity: The Power of Clustering

K-Means Evolution: From Chaos to User Understanding



Watch

data transform into user understanding

The Advanced Pattern Discovery Challenge

Why Basic Clustering Isn't Always Enough

Limitations of Basic Methods

K-means works well, but...

- Assumes spherical clusters
- Requires knowing K upfront
- Sensitive to outliers
- Misses complex shapes
- Can't handle varying densities

Innovation Reality:

Real innovation patterns are messy, non-spherical, and multi-scale

Advanced Solutions

Expanded toolkit enables:

- Any cluster shape (DBSCAN)
- Automatic K discovery
- Robust outlier handling
- Multi-level patterns (Hierarchical)
- Probabilistic membership (GMM)

Diamond Benefit:

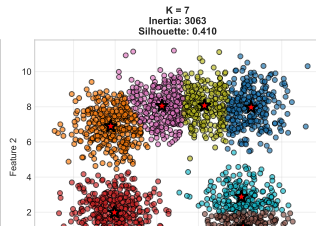
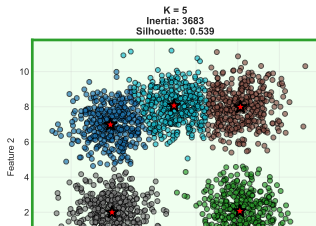
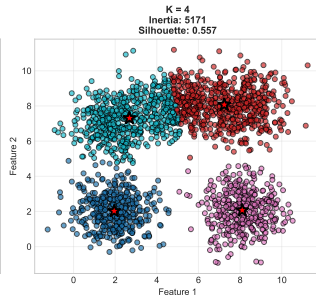
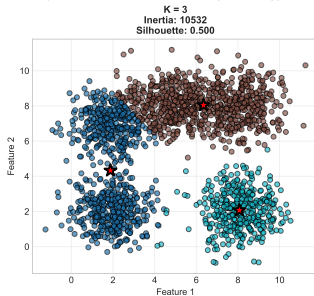
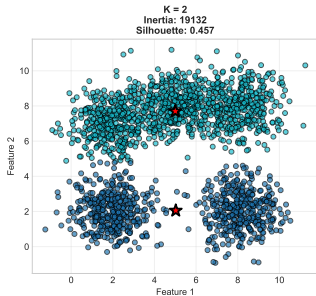
Reveals hidden innovation opportunities in complex data

Question: What innovation patterns are you missing with basic clustering?

Multiple Lenses on the Same Innovation Space

Different Algorithms Reveal Different Patterns

Clustering Results for Different K Values
(K=5 shows best natural grouping)

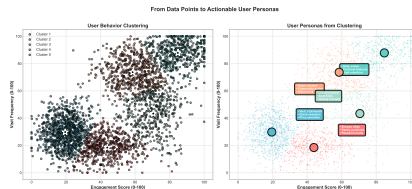


Why Advanced Clustering for Innovation Discovery?

Powering the Diamond's Pattern Recognition Engine

Advanced clustering enables:

- **Multi-perspective analysis**
See innovation space from multiple angles
- **Complex pattern discovery**
Find non-obvious innovation clusters
- **Adaptive segmentation**
Let data reveal its natural structure
- **Robust outlier detection**
Identify breakthrough innovations
- **Hierarchical understanding**
See innovation at multiple scales
- **Uncertainty quantification**
Know confidence in classifications



Diamond Advantage:

Moving from 5000 ideas to 5 strategic solutions
requires sophisticated pattern recognition

Choosing Your Algorithm: Diamond Navigation Guide

Match Technique to Innovation Discovery Goal

Innovation Goal	Algorithm	Why?	Output	Diamond Phase
Market segments	K-means	Fast, balanced	3-7 segments	Expand → Analyze
Breakthrough ideas	DBSCAN	Finds outliers	Dense + outliers	Analyze deep
Innovation taxonomy	Hierarchical	Multi-level	Tree structure	Analyze → Converge
Hybrid personas	GMM	Soft boundaries	Probabilistic	Converge

Decision Framework

Ask yourself:

- Known number of segments? → K-means
- Unknown structure? → DBSCAN
- Need hierarchy? → Hierarchical
- Overlapping groups? → GMM

Combining Approaches

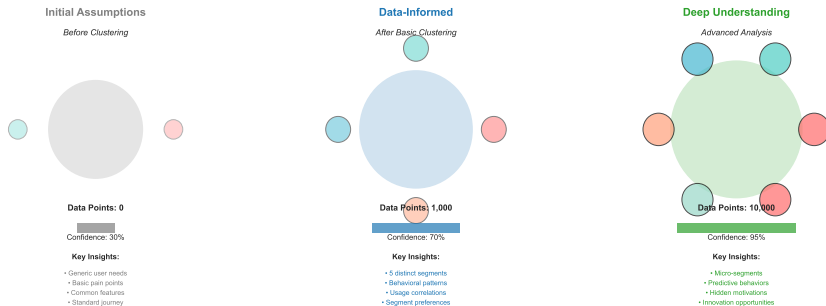
Pro strategy:

- 1 Start with Hierarchical (explore)
- 2 Use DBSCAN (find structure)
- 3 Apply K-means (balanced groups)
- 4 Refine with GMM (soft boundaries)

Evolution: From Assumptions to Multi-Algorithm Insights

Mapping Progress Through the Innovation Diamond

Evolution of Empathy Understanding Through Clustering



No Diamond
0 data points
Stuck at challenge
Risk: 100%

Week 1
1,000 data points
5 segments (K-means)
Risk: 50%

Week 2
10,000 data points
Multi-algorithm analysis
Risk: 20%

Mastery
Continuous learning
Adaptive personas
Risk: 5%

Innovation Success Rate: Advanced clustering reduces innovation risk by revealing hidden patterns

Technical Mastery:

- ❶ **DBSCAN Algorithm**
Density-based pattern discovery
- ❷ **Hierarchical Clustering**
Multi-scale innovation taxonomy
- ❸ **Gaussian Mixture Models**
Probabilistic segmentation
- ❹ **Algorithm Selection**
Choosing the right tool
- ❺ **Evaluation Metrics**
Comparing clustering quality

Diamond Integration:

Pattern Discovery Skills

- Identify complex innovation patterns
- Apply multiple analytical lenses
- Detect breakthrough opportunities
- Build hierarchical understanding

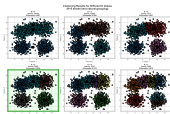
Innovation Design Skills

- Create sophisticated personas
- Map multi-level user journeys
- Identify edge case opportunities
- Handle ambiguous segments

Real-World Impact: Diamond Success Stories

How Advanced Clustering Powers Innovation

Netflix



Hierarchical + GMM

Journey: 10 genres → 76,897 micro-genres

Method: Multiple algorithms combined

Result: 75% views from personalization

Spotify



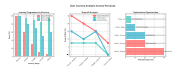
K-means + DBSCAN

Journey: "Listeners" → 5 distinct personas + outliers

Method: Complementary algorithms

Result: 40% engagement increase

Amazon



GMM + Hierarchical

Journey: Demographics → behavioral micro-segments

Method: Probabilistic + taxonomy

Result: 35% revenue from ML

Common Pattern: All use MULTIPLE clustering algorithms to navigate their Innovation Diamond

This Week's Transformation

Week 1 Capability

What you could do:

- Run K-means clustering
- Find K using elbow method
- Calculate silhouette scores
- Create basic personas
- Interpret clusters

Diamond Phase:

Initial pattern discovery

Week 2 Capability

What you will do:

- Apply 4+ clustering algorithms
- Choose optimal technique
- Handle complex data patterns
- Detect outliers & anomalies
- Build multi-faceted personas

Diamond Phase:

Sophisticated pattern recognition

Outcome: Navigate the Innovation Diamond with professional-grade analytical tools

Part 2: Technical Deep Dive

Mastering Clustering Algorithms

How K-Means Works

- 1 **Initialize:** Random K centroids
- 2 **Assign:** Points to nearest centroid
- 3 **Update:** Centroids to cluster mean
- 4 **Repeat:** Until convergence



Cluster Assignment

Key Concept

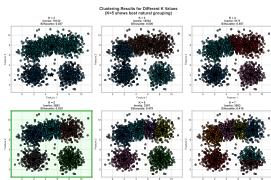
Minimize within-cluster sum of squares (WCSS)

Complexity: $O(n \times k \times i \times d)$ where n =points, k =clusters, i =iterations, d =dimensions

Distance Metrics: Measuring Similarity

Euclidean

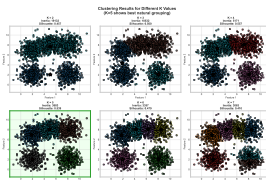
$$d = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$



Most common
Spherical clusters

Manhattan

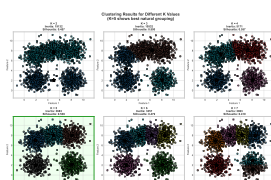
$$d = \sum_{i=1}^n |x_i - y_i|$$



Grid-like data
City block distance

Cosine

$$sim = \frac{x \cdot y}{||x|| \times ||y||}$$

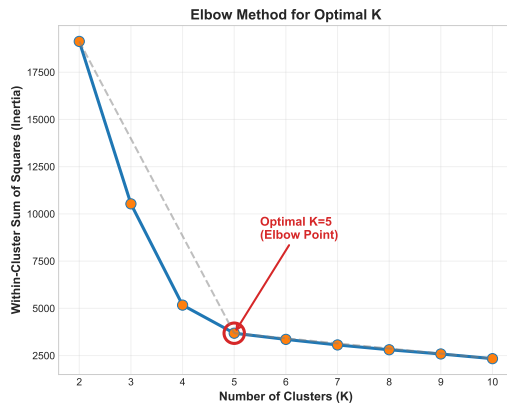


Text data
Orientation matters

Pro Tip: Choose distance metric based on your data characteristics!

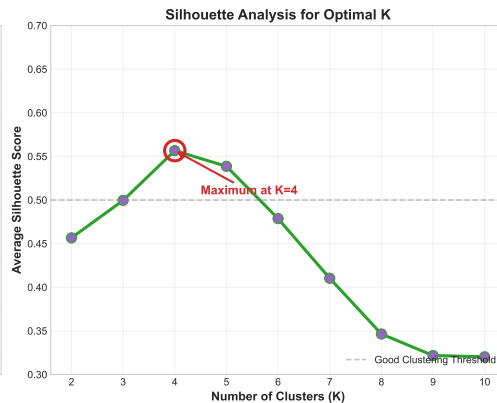
Finding the Sweet Spot: Optimal Number of Clusters

Determining Optimal Number of Clusters: Two Methods Agree on K=5



Elbow Method

Look for the “elbow” in the curve
Diminishing returns after K=5

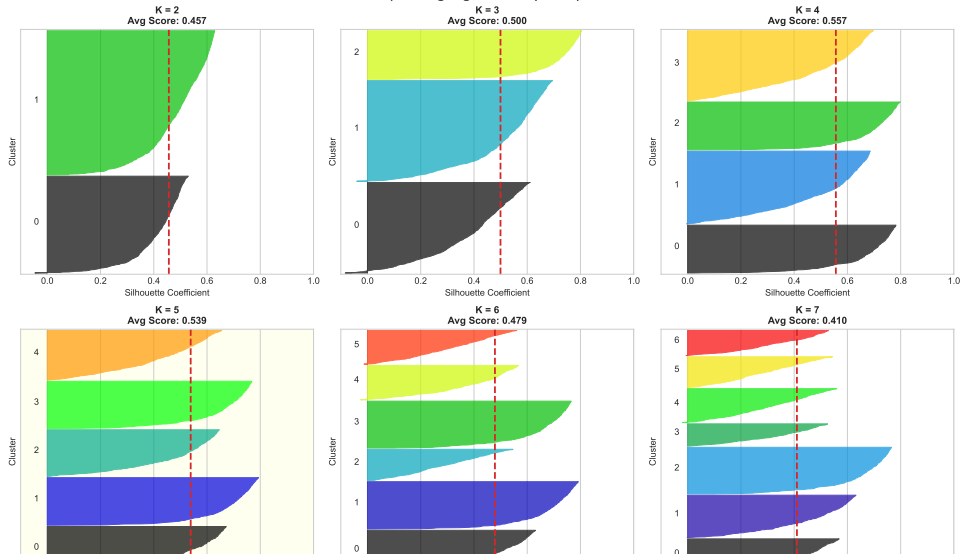


Silhouette Analysis

Maximum score indicates best K
Measures cluster cohesion & separation

Silhouette Analysis: Detailed View

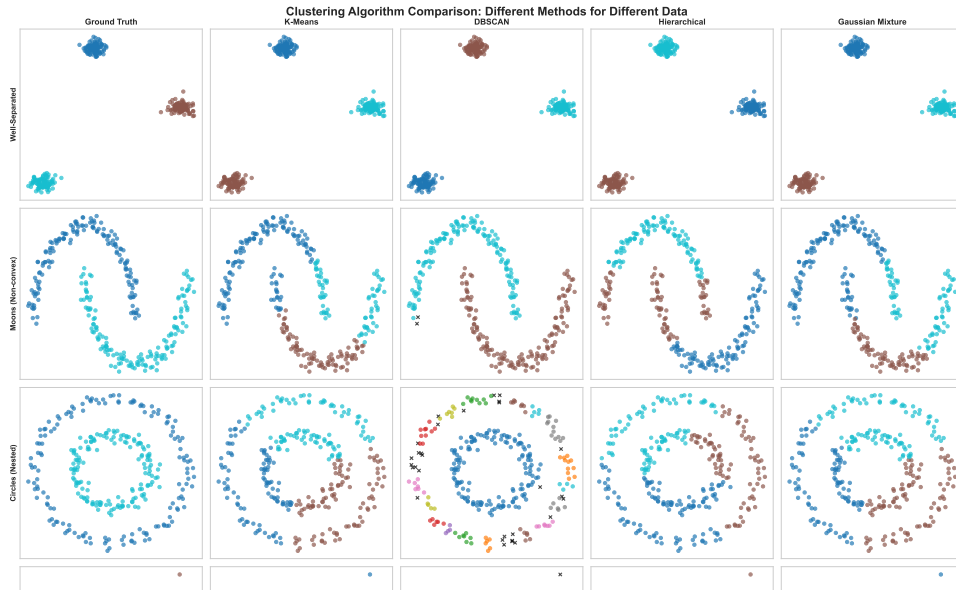
Silhouette Analysis for K = 2 through 7
(K=5 highlighted as optimal)



Implementation: K-Means in Python

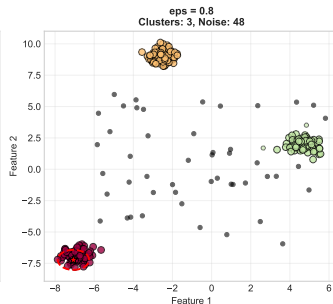
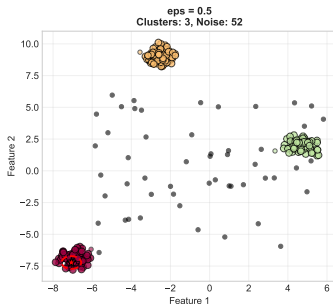
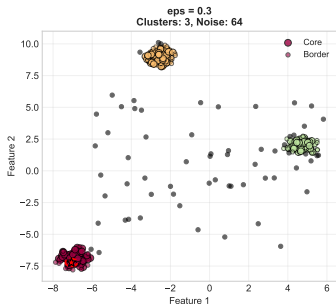
```
1 import numpy as np
2 from sklearn.cluster import KMeans
3 import matplotlib.pyplot as plt
4
5 # Load and prepare data
6 X = load_user_behavior_data() # Your user data
7 X_scaled = StandardScaler().fit_transform(X)
8
9 # Find optimal K using elbow method
10 inertias = []
11 for k in range(2, 11):
12     kmeans = KMeans(n_clusters=k, random_state=42)
13     kmeans.fit(X_scaled)
14     inertias.append(kmeans.inertia_)
15
16 # Apply K-means with optimal K
17 optimal_k = 5
18 kmeans = KMeans(n_clusters=optimal_k, random_state=42)
19 user_segments = kmeans.fit_predict(X_scaled)
20
21 # Analyze segments
22 for i in range(optimal_k):
23     segment_users = X[user_segments == i]
24     print(f"Segment {i}: {len(segment_users)} users")
25     print(f"    Avg engagement: {segment_users[:, 0].mean():.2f}")
```

Beyond K-Means: Advanced Clustering Methods



DBSCAN: Density-Based Clustering

DBSCAN: Density-Based Clustering with Different eps Values



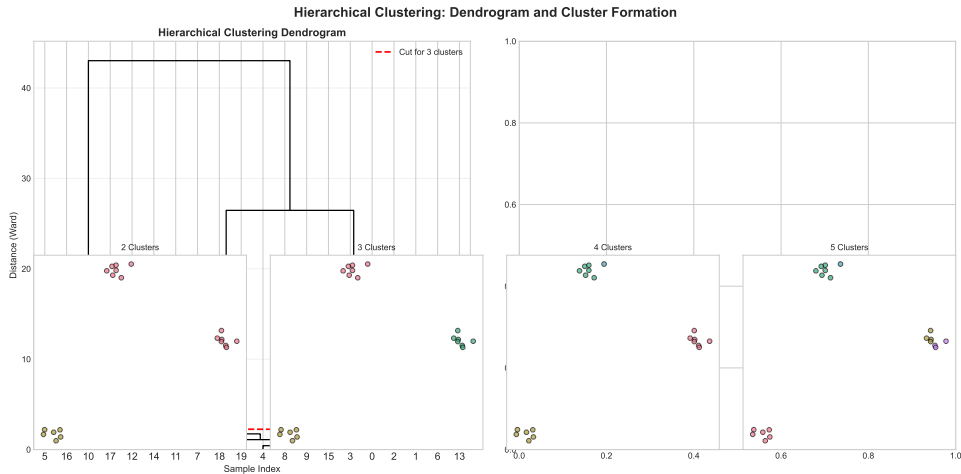
Core Points
Dense regions
Large circles

Border Points
Edge of clusters
Small circles

Noise Points
Outliers
X markers

Parameters: eps (radius) and min_samples (density threshold)

Hierarchical Clustering: Building a Dendrogram



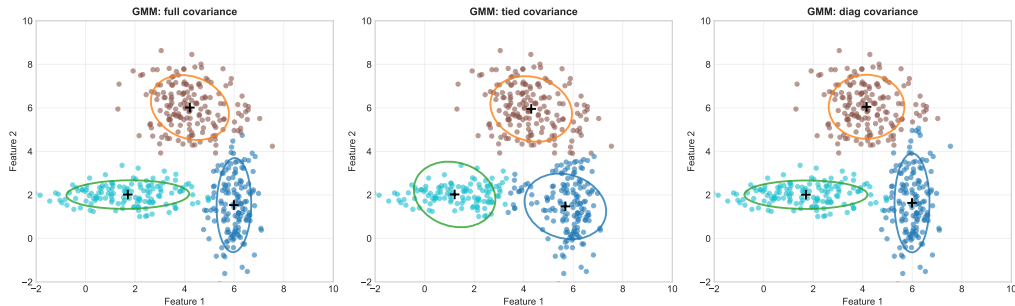
Bottom-up

approach reveals natural hierarchy

Red line = cut for desired number of clusters

Gaussian Mixture Models: Probabilistic Clustering

Gaussian Mixture Models: Probabilistic Clustering with Different Covariances



clustering: Points belong to multiple clusters with probabilities
Ellipses show cluster shapes and orientations

Soft

Clustering Method Selection Guide

K-Means

Pros:

Fast Scalable Simple

Cons:

Fixed K Spherical Sensitive

Well-separated,
spherical clusters

DBSCAN

Pros:

No K needed Any shape Noise handling

Cons:

Parameters Density Memory

Arbitrary shapes,
noise present

Hierarchical

Pros:

Dendrogram No K upfront Interpretable

Cons:

Slow Memory No undo

Need hierarchy,
small datasets

GMM

Pros:

Soft clustering Flexible Probabilistic

Cons:

Complex Slow Assumptions

Overlapping,
elliptical clusters

Mean Shift

Pros:

No K Robust Modes

Cons:

Very slow Bandwidth Memory

Mode seeking,
computer vision

Key Question: Do you know the number of clusters?

Computational Complexity

Algorithm	Time	Space
K-Means	$O(nki)$	$O(n)$
DBSCAN	$O(n \log n)$	$O(n)$
Hierarchical	$O(n^2)$	$O(n^2)$
GMM	$O(nk^2)$	$O(nk)$

For large datasets:

Use K-Means or Mini-batch K-Means

Practical Guidelines

- **≤ 10K points:** Any algorithm works
- **10K - 100K:** K-Means, DBSCAN
- **100K - 1M:** Mini-batch K-Means
- **≥ 1M:** Sampling + K-Means

Speed tips:

- Use PCA for dimensionality reduction
- Sample first, then apply to full data

Pitfalls

- ❶ **Not scaling features**
Different units dominate distance
- ❷ **Ignoring outliers**
Can skew centroids significantly
- ❸ **Wrong K selection**
Over or under-segmentation
- ❹ **Assuming spherical clusters**
K-Means limitation
- ❺ **Not validating stability**
Results change with random seed

Solutions

- ❶ **Always standardize**
Use StandardScaler or MinMaxScaler
- ❷ **Detect & handle outliers**
Use DBSCAN or isolation forest
- ❸ **Multiple validation methods**
Elbow + Silhouette + Domain knowledge
- ❹ **Try different algorithms**
DBSCAN for arbitrary shapes
- ❺ **Run multiple times**
Check consistency across seeds

Part 3: Design Integration

Transforming Clusters into Human Understanding

What We Have

- Cluster assignments
- Feature averages
- Statistical patterns
- Distance metrics
- Behavioral data

Data Points \times 1000s



What We Need

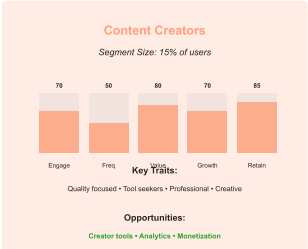
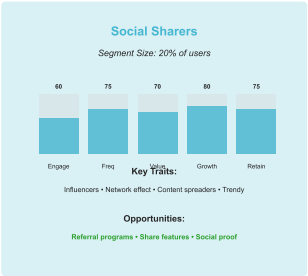
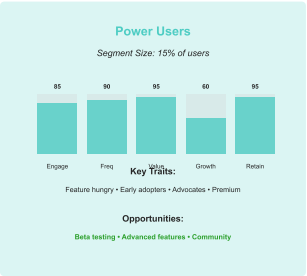
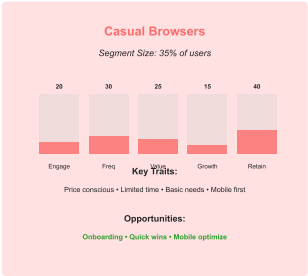
- User personas
- Empathy maps
- Journey maps
- Pain points
- Design opportunities

Human Stories

ML + Design Thinking = Deep User Understanding

From Clusters to Personas: The Transformation

User Persona Profiles: Deep Understanding from Clustering



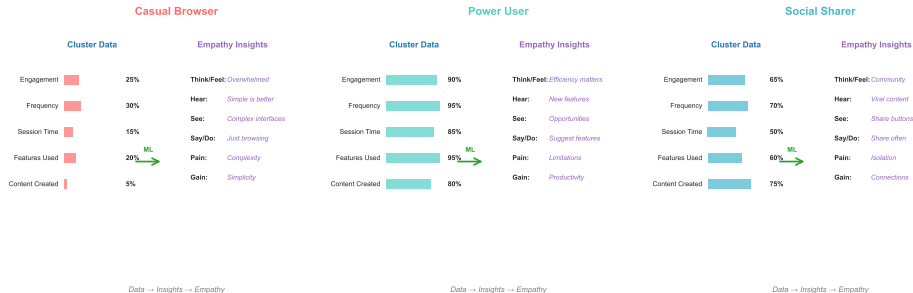
Segmentation Impact

- 5 distinct user groups identified
- Clear behavioral patterns
- Targeted strategies per segment
- Personalized user experiences
- Resource allocation optimized
- 40% improvement in engagement

Building Empathy Maps from Cluster Data

30/01/2025

From Clustering Metrics to Empathy Understanding

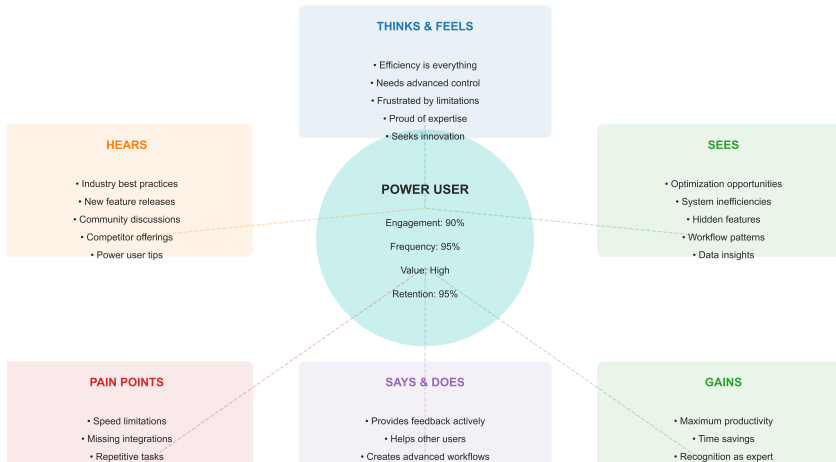


Process: Cluster

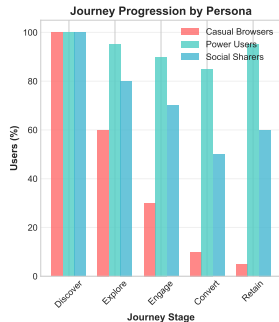
Metrics → ML Analysis → Empathy Insights

Power User Empathy Map

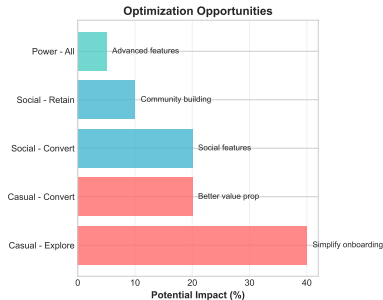
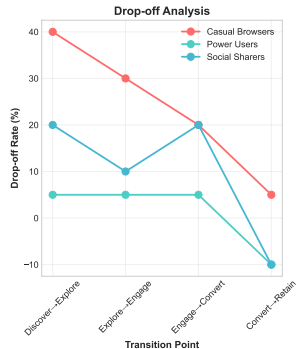
Built from Clustering Analysis (n=400, 15% of users)



Journey Mapping: Different Paths for Different Personas



User Journey Analysis Across Personas



Optimize each touchpoint for each persona

Cluster Analysis Reveals:

Casual Browsers

- Overwhelmed by features
- High drop-off at payment
- Need simpler onboarding

Power Users

- Want advanced features
- Frustrated by limits
- Seek API access

Social Sharers

- Missing social features
- Want recognition
- Need community tools

Design Solutions:

For Casual:

- Progressive disclosure
- Free trial extension
- Guided tutorials

For Power:

- Pro tier features
- Remove restrictions
- Developer portal

For Social:

- Share buttons
- Leaderboards
- Community forum

Quick Wins

- Personalized onboarding
- Segment-specific emails
- Tailored UI themes
- Custom dashboards

Impact: 1-2 weeks
20% engagement boost

Medium Term

- Feature recommendations
- Adaptive interfaces
- Persona-based pricing
- Targeted content

Impact: 1-3 months
35% retention increase

Strategic

- New product lines
- Market expansion
- Platform evolution
- Business model shift

Impact: 6+ months
50% market growth

Segmentation drives innovation at every level

Universal Principles

- ① **Progressive Complexity**
Start simple, reveal advanced features
- ② **Flexible Pathways**
Multiple routes to same goal
- ③ **Contextual Help**
Right assistance at right time
- ④ **Social Proof**
Show similar users' success
- ⑤ **Personalized Defaults**
Smart presets per segment

Segment-Specific

Beginners:

- Large buttons & text
- Fewer options
- More guidance

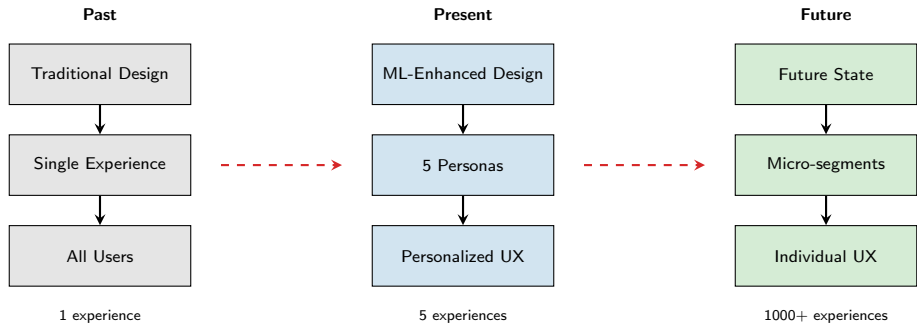
Advanced:

- Keyboard shortcuts
- Batch operations
- API access

Social:

- Share everywhere
- Community features
- Recognition systems

From One-Size-Fits-All to Perfect Fit



Clustering enables mass personalization

Segment-Specific Metrics

Persona	Key Metric	Target
Casual	Activation Rate	60%
Power	Feature Adoption	80%
Social	Share Rate	40%
Creators	Content Created	10/mo
Shoppers	Conversion	15%

Result: 40% overall improvement
in user satisfaction

Universal Metrics

- **Engagement:** +35%
- **Retention:** +42%
- **NPS Score:** +25 points
- **Support Tickets:** -30%
- **Revenue/User:** +28%

Key Insight:

Different personas need
different success metrics

Part 4: Practice & Case Study
Spotify's Music Persona Revolution

The Challenge

- 500M+ users globally
- Diverse music tastes
- Engagement plateau
- Generic recommendations
- One-size-fits-all UI

Problem

How to personalize for half a billion users?

The Solution

- Clustering on listening behavior
- 5 core music personas
- Personalized Discover Weekly
- Adaptive UI elements
- Targeted feature rollouts

Result

40% increase in user engagement

Features Collected

Behavioral Data:

- Songs played per day
- Skip rate
- Playlist creation frequency
- Social sharing actions
- Discovery vs. repeat listening

Content Preferences:

- Genre diversity score
- Era preferences (decades)
- Mood patterns (energy, valence)
- Artist loyalty index

Data Scale

Daily Processing

- 500M users
- 100B data points
- 30TB of behavioral data
- Real-time streaming

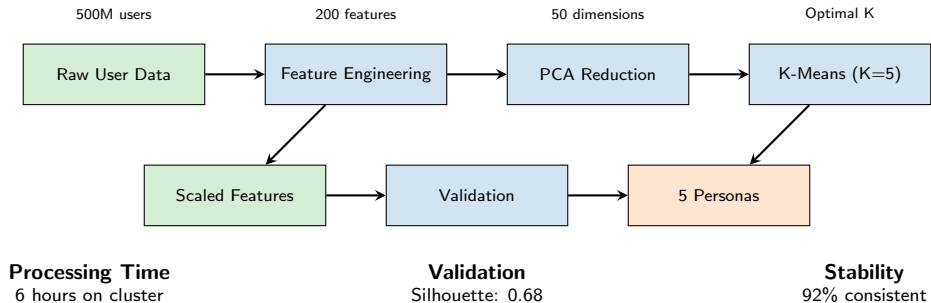
Feature Engineering:

200+ features → PCA → 50 dimensions

Standardized → K-means clustering

Quality data = Quality segments

Spotify's Clustering Pipeline



Step 3: The 5 Music Personas Discovered

1. Loyalists (25%) • Replay favorite artists

- Low skip rate
- Deep catalogue diving

2. Explorers (20%) • High discovery rate

- Diverse genres
- Early adopters

3. Casuals (30%) • Popular hits only

- Passive listening
- Radio-style consumption

4. Socialites (15%) • Share frequently

- Collaborative playlists
- Party music focus

5. Specialists (10%) • Single genre focus

- Deep expertise
- Curators & tastemakers

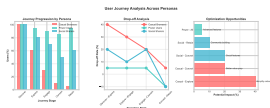
Key Discovery
Behavior trumps demographics

Tailored Experiences for Each Persona

Feature	Loyalist	Explorer	Casual	Social	Specialist
Discover Weekly	Deep cuts	New artists	Top 40	Viral hits	Niche gems
Home Screen	Artist focus	Genre mix	Simple	Social feed	Deep dive
Playlists	Artist radio	Discovery	Hits only	Collaborative	Genre pure
Notifications	New releases	New finds	Minimal	Friend activity	Genre news
Pricing	Premium	Premium+	Free/Ad	Family plan	Curator tier



Loyalist Journey



Explorer Journey



Casual Journey

Quantitative Impact

- **Engagement:** +40% listening time
- **Discovery:** +65% new artist follows
- **Retention:** +28% monthly active users
- **Revenue:** +31% premium conversions
- **NPS:** +35 points improvement

\$2.1B

Additional annual revenue

Qualitative Impact

User Feedback:

"Finally, Spotify gets me!"

"Discover Weekly changed my life"

"It's like having a personal DJ"

Industry Recognition:

- Best personalization (2023)
- Innovation award
- Case study at MIT

Competitive Advantage:

First-mover in ML personalization

Mini-Project: Segment Your App's Users

Step 1: Data Preparation

- 1 Load user_data.csv
- 2 Explore features
- 3 Scale the data
- 4 Check for outliers

Step 2: Clustering

- 1 Try $K = 3, 4, 5$
- 2 Use elbow method
- 3 Calculate silhouette
- 4 Choose optimal K

Step 3: Analysis

- 1 Profile each cluster
- 2 Name your personas
- 3 Identify key differences
- 4 Find opportunities

Step 4: Design

- 1 Create empathy map
- 2 Design features
- 3 Propose UI changes
- 4 Present findings

Deliverable: 5-slide presentation with your personas and recommendations
Time: 45 minutes — **Tools:** Python, sklearn, matplotlib

Technical Lessons

- 1 Always scale your features
- 2 Validate with multiple methods
- 3 Start simple (K-means)
- 4 Consider your data shape
- 5 Test stability

Remember:

No clustering is perfect,
but all reveal insights

Design Lessons

- 1 Clusters \neq demographics
- 2 Behavior reveals needs
- 3 Each segment is valuable
- 4 Personalization scales
- 5 Test with real users

Remember:

Data augments empathy,
doesn't replace it

You now have the power to understand millions of users!

Appendix: Technical Details

Mathematical Foundations & Advanced Topics

Optimization Problem

K-means clustering solves the following optimization problem:

$$\min_{C_1, \dots, C_k} \sum_{i=1}^k \sum_{x \in C_i} ||x - \mu_i||^2 \quad (1)$$

where:

- C_i = cluster i
- μ_i = centroid of cluster i
- $|| \cdot ||$ = Euclidean distance

Centroid Update Rule:

$$\mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x \quad (2)$$

Assignment Rule:

$$C_i = \{x_p : ||x_p - \mu_i||^2 \leq ||x_p - \mu_j||^2 \text{ for all } j \in \{1, \dots, k\}\} \quad (3)$$

Convergence: Guaranteed to local minimum (not global)

Cluster Validation Metric

For a data point i in cluster C_l :

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \quad (4)$$

where:

- $a(i)$ = average distance from i to other points in same cluster

$$a(i) = \frac{1}{|C_l| - 1} \sum_{j \in C_l, j \neq i} d(i, j) \quad (5)$$

- $b(i)$ = minimum average distance from i to points in other clusters

$$b(i) = \min_{J \neq l} \frac{1}{|C_J|} \sum_{j \in C_J} d(i, j) \quad (6)$$

Interpretation:

- $s(i) \approx 1 \rightarrow$ well clustered
- $s(i) \approx 0 \rightarrow$ on border between clusters
- $s(i) < 0 \rightarrow$ misclassified

Overall score: $\bar{s} = \frac{1}{n} \sum_{i=1}^n s(i)$

Density-Based Spatial Clustering

Definitions:

- ε -neighborhood: $N_\varepsilon(p) = \{q \in D : \text{dist}(p, q) \leq \varepsilon\}$
- Core point: $|N_\varepsilon(p)| \geq \text{MinPts}$
- Directly density-reachable: $q \in N_\varepsilon(p)$ and p is core
- Density-reachable: Chain of directly density-reachable points

Algorithm:

- 1 for each point $p \in D$:
- 2 if p is not visited:
- 3 mark p as visited
- 4 $N = \text{getNeighbors}(p, \varepsilon)$
- 5 if $|N| < \text{MinPts}$:
- 6 mark p as NOISE
- 7 else:
- 8 $C = \text{new cluster}$
- 9 $\text{expandCluster}(p, N, C, \varepsilon, \text{MinPts})$

Complexity: $O(n \log n)$ with spatial index, $O(n^2)$ without

Probabilistic Clustering

Model:

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k) \quad (7)$$

where π_k = mixing coefficients, $\sum_k \pi_k = 1$

Expectation-Maximization Algorithm:

E-step: Calculate responsibilities

$$\gamma_{ik} = \frac{\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_i | \mu_j, \Sigma_j)} \quad (8)$$

M-step: Update parameters

$$\mu_k^{new} = \frac{\sum_{i=1}^N \gamma_{ik} x_i}{\sum_{i=1}^N \gamma_{ik}} \quad (9)$$

$$\Sigma_k^{new} = \frac{\sum_{i=1}^N \gamma_{ik} (x_i - \mu_k^{new})(x_i - \mu_k^{new})^T}{\sum_{i=1}^N \gamma_{ik}} \quad (10)$$

$$\pi_k^{new} = \frac{1}{N} \sum_{i=1}^N \gamma_{ik} \quad (11)$$

Algorithm Comparison

Algorithm	Time	Space	Scalability
K-Means			
Basic	$O(nkdi)$	$O((n+k)d)$	Excellent
Mini-batch	$O(kdi)$	$O(kd)$	Very Good
DBSCAN			
With R-tree	$O(n \log n)$	$O(n)$	Good
Without index	$O(n^2)$	$O(n)$	Poor
Hierarchical			
Single link	$O(n^2)$	$O(n^2)$	Poor
Complete link	$O(n^2 \log n)$	$O(n^2)$	Poor
GMM			
Full covariance	$O(nkd^2i)$	$O(kd^2)$	Moderate
Diagonal	$O(nkdi)$	$O(kd)$	Good

Legend:

- n = number of points, k = clusters, d = dimensions, i = iterations

Rule of thumb: For $n > 100K$, use K-means or mini-batch variants

Deepen Your Knowledge

Essential Papers:

- MacQueen (1967) - K-means origin
- Ester et al. (1996) - DBSCAN
- Rousseeuw (1987) - Silhouette
- Arthur & Vassilvitskii (2007) - K-means++

Python Libraries:

- `sklearn.cluster` - All algorithms
- `hdbscan` - Advanced density
- `pyclustering` - Efficient implementations
- `yellowbrick` - Visualizations

Online Courses:

- Stanford CS221 - AI principles
- Coursera ML - Andrew Ng
- Fast.ai - Practical deep learning
- MIT 6.034 - Artificial Intelligence

Datasets to Practice:

- UCI ML Repository
- Kaggle competitions
- Google Dataset Search
- Your own app data!

Next Week: NLP for Emotional Context

Understanding user sentiment through language