## Machine Learning for Smarter Innovation
### Week 2: Clustering for Deep Empathy

BSc Course in AI-Enhanced Innovation

Understanding Users Through Data-Driven Segmentation

**Transform data points into human insights**

# Part 1: Foundation

Understanding Users Through Data Patterns

K-Means Evolution: From Chaos to User Understanding

**Watch** data transform into user understanding

**Traditional Challenges**

- Generic personas based on assumptions
- Missing hidden user segments
- Biased by loudest voices
- Static, outdated profiles
- Limited sample sizes

**ML-Enhanced Solutions**

- Data-driven segment discovery
- Uncover unexpected patterns
- Balanced representation
- Dynamic, evolving insights
- Scale to millions of users

**Question: How can we truly understand ALL our users?**

## Clustering reveals natural user groups

- **Discover** hidden behavioral patterns
- **Segment** users by actual behavior, not demographics
- **Identify** underserved user groups
- **Personalize** experiences at scale
- **Predict** user needs and preferences

> **Key Insight:** Users naturally form groups based on behaviors, needs, and preferences



From Data Points to Actionable User Personas

Evolution of Empathy Understanding Through Clustering

## You Will Master:

1. **K-means Algorithm**
   Understanding the mechanics

2. **Optimal Cluster Selection**
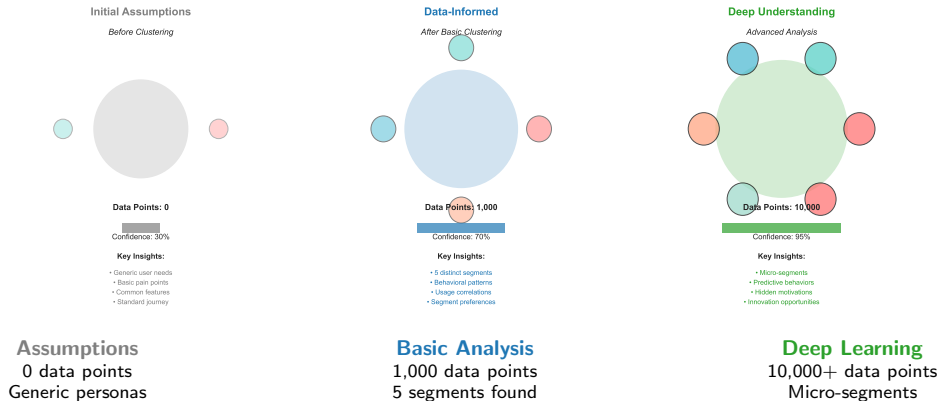   Elbow method & silhouette analysis

3. **Advanced Methods**
   DBSCAN, Hierarchical, GMM

4. **Persona Creation**
   From clusters to empathy maps

5. **Real Implementation**
   Spotify case study

## Key Outcomes:

### Technical Skills

- Implement clustering in Python
- Evaluate cluster quality
- Choose right algorithm

### Design Skills

- Create data-driven personas
- Build empathy maps
- Map user journeys

## Netflix



**2000+ taste groups**
Personalized recommendations
75% of views from algorithms

## Spotify



**5 music personas**
Discover Weekly success
40% engagement increase

## Amazon



**Micro-segments**
Purchase prediction
35% of revenue from ML

**These companies understand users through clustering**

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

**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!

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

```python
import numpy as np
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt

# Load and prepare data
X = load_user_behavior_data()  # Your user data
X_scaled = StandardScaler().fit_transform(X)

# Find optimal K using elbow method
inertias = []
for k in range(2, 11):
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(X_scaled)
    inertias.append(kmeans.inertia_)

# Apply K-means with optimal K
optimal_k = 5
kmeans = KMeans(n_clusters=optimal_k, random_state=42)
user_segments = kmeans.fit_predict(X_scaled)

# Analyze segments
for i in range(optimal_k):
    segment_users = X[user_segments == i]
    print(f"Segment {i}: {len(segment_users)} users")
    print(f"  Avg engagement: {segment_users[:, 0].mean():.2f}")
```

# Beyond K-Means: Advanced Clustering Methods



Clustering Algorithm Comparison: Different Methods for Different Data

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: Dendrogram and Cluster Formation

**Bottom-up approach reveals natural hierarchy**

Red line = cut for desired number of clusters

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

1. **Not scaling features**
   Different units dominate distance

2. **Ignoring outliers**
   Can skew centroids significantly

3. **Wrong K selection**
   Over or under-segmentation

4. **Assuming spherical clusters**
   K-Means limitation

5. **Not validating stability**
   Results change with random seed

## Solutions

1. **Always standardize**
   Use StandardScaler or MinMaxScaler

2. **Detect & handle outliers**
   Use DBSCAN or isolation forest

3. **Multiple validation methods**
   Elbow + Silhouette + Domain knowledge

4. **Try different algorithms**
   DBSCAN for arbitrary shapes

5. **Run multiple times**
   Check consistency across seeds

## Part 3: Design Integration

Transforming Clusters into Human Understanding

**What We Have** $\longrightarrow$ **What We Need**

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

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

**Data Points × 1000s**

**Human Stories**

**ML + Design Thinking = Deep User Understanding**

**User Persona Profiles: Deep Understanding from Clustering**

### Casual Browsers

*Segment Size: 35% of users*

| Engage | Freq | Value | Growth | Retain |
|--------|------|-------|--------|--------|
| 20 | 30 | 25 | 15 | 40 |

**Key Traits:**

Price conscious • Limited time • Basic needs • Mobile first

**Opportunities:**

Onboarding • Quick wins • Mobile optimize

### Power Users

*Segment Size: 15% of users*

| Engage | Freq | Value | Growth | Retain |
|--------|------|-------|--------|--------|
| 85 | 90 | 95 | 60 | 95 |

**Key Traits:**

Feature hungry • Early adopters • Advocates • Premium

**Opportunities:**

Beta testing • Advanced features • Community

### Social Sharers

*Segment Size: 20% of users*

| Engage | Freq | Value | Growth | Retain |
|--------|------|-------|--------|--------|
| 60 | 75 | 70 | 80 | 75 |

**Key Traits:**

Influencers • Network effect • Content spreaders • Trendy

**Opportunities:**

Referral programs • Share features • Social proof

### Content Creators

*Segment Size: 15% of users*

| Engage | Freq | Value | Growth | Retain |
|--------|------|-------|--------|--------|
| 70 | 50 | 80 | 70 | 85 |

**Key Traits:**

Quality focused • Tool seekers • Professional • Creative

**Opportunities:**

Creator tools • Analytics • Monetization

### Window Shoppers

*Segment Size: 25% of users*

| Engage | Freq | Value | Growth | Retain |
|--------|------|-------|--------|--------|
| 40 | 20 | 30 | 50 | 35 |

**Key Traits:**

Researchers • Comparison • Price sensitive • Cautious

**Opportunities:**

Trust signals • Reviews • Free trials

### Segmentation Impact

→ 5 distinct user groups identified

→ Clear behavioral patterns

→ Targeted strategies per segment

→ Personalized user experiences

→ Resource allocation optimized

→ 40% improvement in engagement

**From Clustering Metrics to Empathy Understanding**

### Casual Browser

**Cluster Data**

| | |
|---|---|
| Engagement | 25% |
| Frequency | 30% |
| Session Time | 15% |
| Features Used | 20% |
| Content Created | 5% |

**Empathy Insights**

**Think/Feel:** *Overwhelmed*
**Hear:** *Simple is better*
**See:** *Complex interfaces*
**Say/Do:** *Just browsing*
**Pain:** *Complexity*
**Gain:** *Simplicity*

### Power User

**Cluster Data**

| | |
|---|---|
| Engagement | 90% |
| Frequency | 95% |
| Session Time | 85% |
| Features Used | 95% |
| Content Created | 80% |

**Empathy Insights**

**Think/Feel:** *Efficiency matters*
**Hear:** *New features*
**See:** *Opportunities*
**Say/Do:** *Suggest features*
**Pain:** *Limitations*
**Gain:** *Productivity*

### Social Sharer

**Cluster Data**

| | |
|---|---|
| Engagement | 65% |
| Frequency | 70% |
| Session Time | 50% |
| Features Used | 60% |
| Content Created | 75% |

**Empathy Insights**

**Think/Feel:** *Community*
**Hear:** *Viral content*
**See:** *Share buttons*
**Say/Do:** *Share often*
**Pain:** *Isolation*
**Gain:** *Connections*

*Data → Insights → Empathy*

*Data → Insights → Empathy*

*Data → Insights → Empathy*

**Process:** Cluster

Metrics → ML Analysis → Empathy Insights

# Power User Empathy Map

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

**THINKS & FEELS**

• Efficiency is everything
• Needs advanced control
• Frustrated by limitations
• Proud of expertise
• Seeks innovation

**HEARS**

• Industry best practices
• New feature releases
• Community discussions
• Competitor offerings
• Power user tips

**POWER USER**

Engagement: 90%

Frequency: 95%

Value: High

Retention: 95%

**SEES**

• Optimization opportunities
• System inefficiencies
• Hidden features
• Workflow patterns
• Data insights

**PAIN POINTS**

• Speed limitations
• Missing integrations
• Repetitive tasks

**SAYS & DOES**

• Provides feedback actively
• Helps other users
• Creates advanced workflows

**GAINS**

• Maximum productivity
• Time savings
• Recognition as expert

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

| Past | Present | Future |
|------|---------|--------|
| Traditional Design | ML-Enhanced Design | Future State |
| Single Experience | 5 Personas | Micro-segments |
| All Users | Personalized UX | Individual UX |
| 1 experience | 5 experiences | 1000+ experiences |

**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**



| 500M users | 200 features | 50 dimensions | Optimal K |
|---|---|---|---|
| Raw User Data | Feature Engineering | PCA Reduction | K-Means (K=5) |

| | | |
|---|---|---|
| Scaled Features | Validation | 5 Personas |

**Processing Time**
6 hours on cluster

**Validation**
Silhouette: 0.68

**Stability**
92% consistent

**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 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:

$$\min_{C_1,\ldots,C_k} \sum_{i=1}^{k} \sum_{x \in C_i} ||x - \mu_i||^2 \tag{1}$$

where:

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

**Centroid Update Rule:**

$$\mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x \tag{2}$$

**Assignment Rule:**

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

**Convergence:** Guaranteed to local minimum (not global)

## Silhouette Coefficient: Mathematical Definition

### Cluster Validation Metric

For a data point $i$ in cluster $C_I$:

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

where:

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

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

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

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

**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:**

- $\varepsilon$-neighborhood: $N_\varepsilon(p) = \{q \in D : dist(p, q) \leq \varepsilon\}$
- Core point: $|N_\varepsilon(p)| \geq MinPts$
- 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.         $N = getNeighbors(p, \varepsilon)$
5.         **if** $|N| < MinPts$:
6.             mark $p$ as NOISE
7.         **else**:
8.             $C = $ new cluster
9.             expandCluster($p$, $N$, $C$, $\varepsilon$, $MinPts$)

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

## Probabilistic Clustering

**Model:**

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

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

**Expectation-Maximization Algorithm:**
**E-step:** Calculate responsibilities

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

**M-step:** Update parameters

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

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

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

**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 | $O(nkd^2 i)$ | $O(kd^2)$ | 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