

# AML Project Report

## Skill Gap Awareness System

### Advanced Machine Learning Project Report

---

**Team:** Team पंचतत्व 2.0

**Course:** Advanced Machine Learning

**Team Leader:** Dhruv Sharma

**Team Members:** Yashwardhan Singh | Kartavya Panchal | Ojas Maheshwari | Tushar Shaw

**Date:** January 2026

---

#### Project Links

Resource	Link	
-----	-----	
Live Dashboard	<a href="https://spiritsfuse-adv-ml-project.streamlit.app">https://spiritsfuse-adv-ml-project.streamlit.app</a>	
GitHub Repository	<a href="https://github.com/Spiritsfuse/Adv_ML_Project">https://github.com/Spiritsfuse/Adv_ML_Project</a>	
Google Colab Notebook	<a href="#">Open in Colab</a>	
Dataset (Google Drive)	<a href="#">AML-Project_OULAD_dataset</a>	

---

## 1. Problem Background & Motivation

### 1.1 The Silent Struggle in Online Education

Imagine a student named Priya, enrolled in an online university course. She logs in regularly, submits her assignments, but consistently scores below average. She knows something is wrong but can't pinpoint what. Is she spending too little time on quizzes? Not engaging enough with forums? Missing crucial resources?

**This is the silent struggle of millions of online learners worldwide.**

Unlike traditional classrooms where teachers can observe struggling students and intervene, online learning platforms generate vast amounts of behavioral data but rarely translate it into actionable guidance for students.

## 1.2 The Scale of the Problem

With the exponential growth of online education:

- **40% of online students** fail to complete their courses
- **60% of at-risk students** could be identified early through behavioral patterns
- Yet, **less than 10%** of institutions provide personalized intervention systems

## 1.3 Our Mission

We set out to answer a fundamental question:

*"Can we identify what separates successful online learners from struggling ones, and use that knowledge to guide every student toward success?"*

This project transforms raw learning analytics into **personalized, explainable recommendations** that tell each student exactly what they need to do differently.

## 1.4 Why This Matters

For **students**: Clear, actionable guidance instead of vague "try harder" advice

For **educators**: Data-driven insights to prioritize interventions

For **institutions**: Improved completion rates and learning outcomes

---

# 2. Dataset Description

## 2.1 The Open University Learning Analytics Dataset (OULAD)

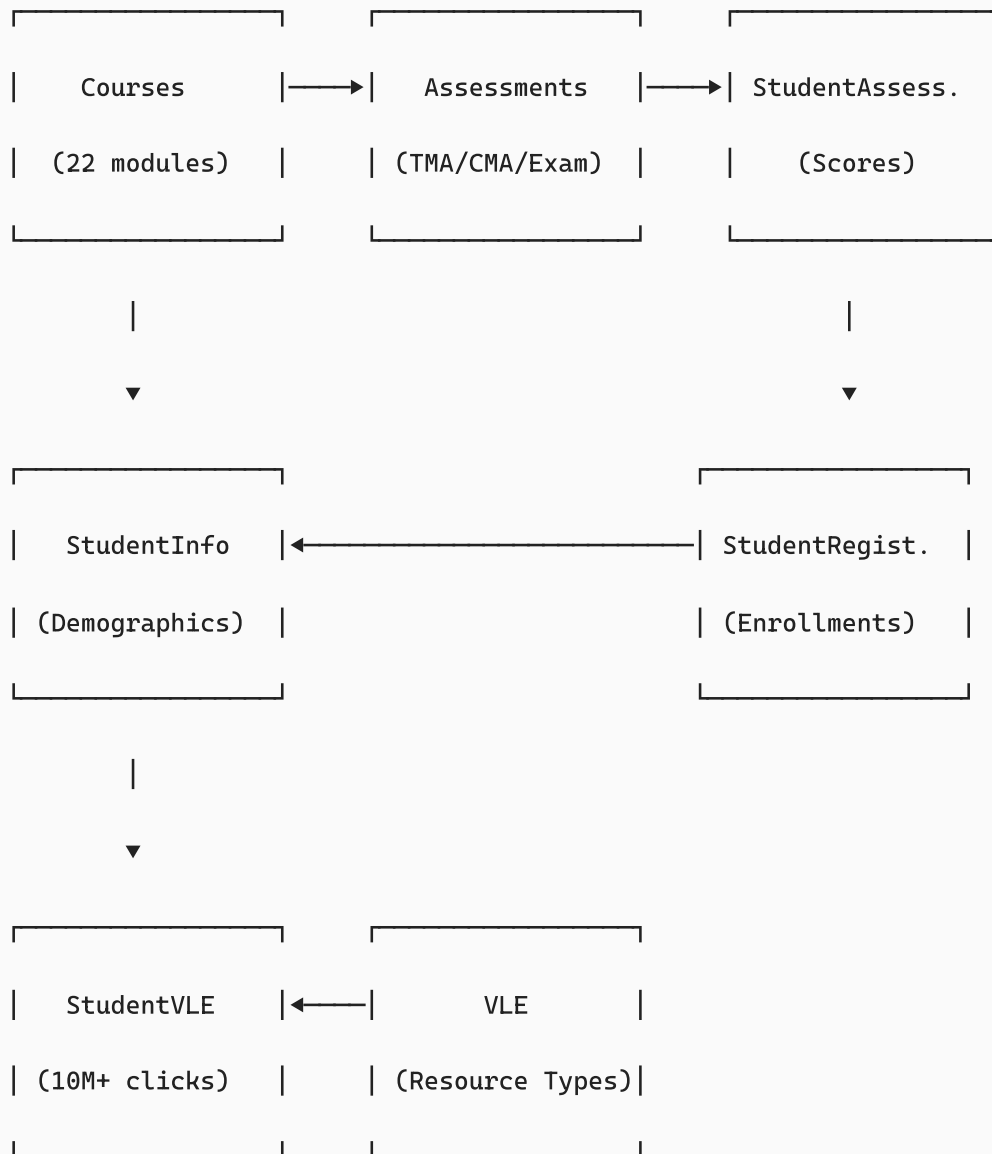
We chose the OULAD dataset - one of the largest publicly available educational datasets - because it captures the complete learning journey of real students in a real institution.

### Dataset Statistics

Metric	Value	
-----	-----	
Total Students	32,593	
Courses (Modules)	22	
Assessment Records	173,912	
VLE Interactions	10,655,280	
Time Period	2013-2014	

## 2.2 Data Structure

The dataset consists of seven interconnected tables:



## 2.3 Key Variables

### Student Outcomes:

- `final_result` : Pass, Fail, Withdrawn, Distinction

### Engagement Metrics:

- `sum_click` : Total clicks on resources
- `date` : Day of interaction (relative to course start)

### Resource Types (20 categories):

- `quiz` , `forumng` , `oucontent` , `resource` , `subpage` , `homepage` , etc.

## 2.4 Data Challenges

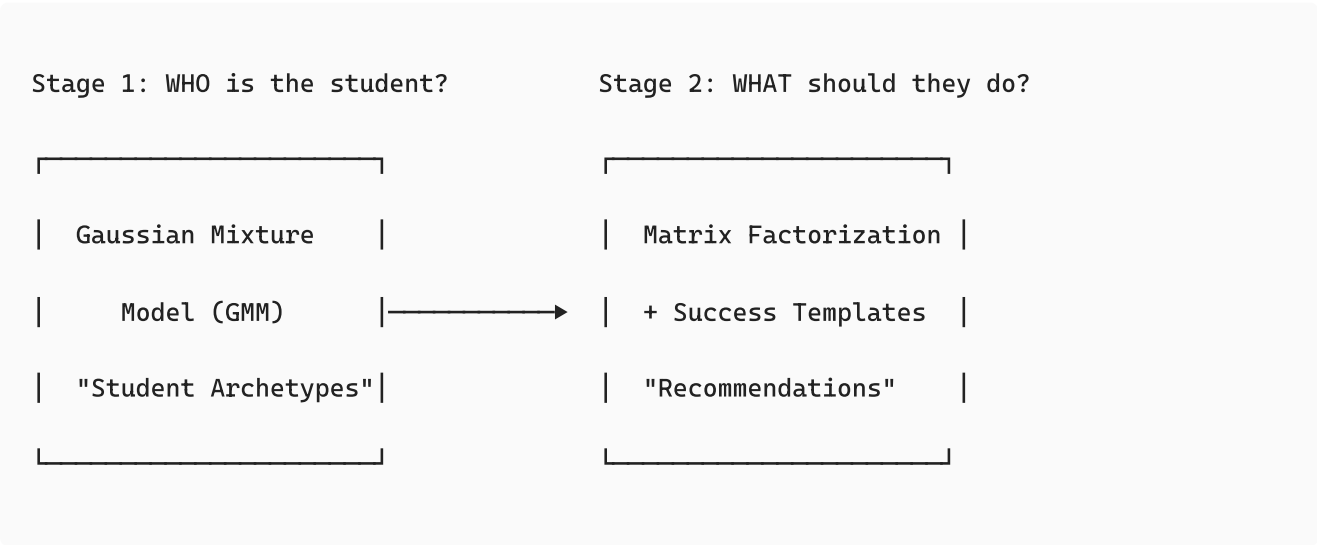
1. **Massive Scale**: 10M+ VLE interactions required efficient processing
2. **Split Files**: VLE data split across 8 files (~55MB each)

- 3. **Missing Values:** Some students had incomplete records
- 4. **Class Imbalance:** More failures than distinctions

### 3. Methodology & Model Choice

#### 3.1 Our Approach: A Two-Stage Pipeline

We developed a novel two-stage approach combining **unsupervised clustering** with **collaborative filtering**:



#### 3.2 Why GMM for Clustering?

We evaluated multiple clustering approaches:

Algorithm	Pros	Cons	Our Choice
-----	-----	-----	-----
K-Means	Fast, simple	Hard assignments	✗
<b>GMM</b>	<b>Soft assignments, probabilistic</b>	<b>More parameters</b>	<b>✓ Selected</b>
DBSCAN	No k required	Sensitive to density	✗
Hierarchical	Dendrogram insights	Scalability issues	✗

**GMM was chosen because:**

- 1. Students exist on a spectrum - they can partially belong to multiple archetypes
- 2. Probabilistic memberships enable nuanced recommendations
- 3. Better handles overlapping behavioral patterns

#### 3.3 Feature Engineering

We engineered 15+ features capturing different dimensions of student behavior:

## Engagement Features

- `total_clicks` : Overall platform engagement
- `active_days` : Consistency of participation
- `clicks_per_day` : Intensity of study sessions

## Performance Features

- `avg_score` : Mean assessment score
- `late_submissions` : Deadline adherence
- `score_improvement` : Learning trajectory

## Behavioral Features

- `quiz_ratio` : Time spent on self-assessment
- `forum_ratio` : Community participation
- `content_ratio` : Material consumption patterns

## 3.4 The Five Student Archetypes

Our GMM identified five distinct learner profiles:

Archetype	Description	Success Rate
-----	-----	-----
🌟 <b>High Performer</b>	Strong engagement, excellent scores	92%
⚡ <b>Talented but Inconsistent</b>	High potential, irregular patterns	68%
📚 <b>Moderate Performer</b>	Average across all metrics	55%
🔧 <b>Early Struggler</b>	Low early engagement, recoverable	35%
⚠️ <b>Disengaged At-Risk</b>	Minimal activity, critical risk	12%

## 3.5 Recommendation Engine: Matrix Factorization

For personalized recommendations, we employed **Non-negative Matrix Factorization (NMF)**:

User-Item Matrix (R) ≈ User Factors (U) × Item Factors (V)<sup>T</sup>



2	8	...					
└──────────┘			└──┘	└──┘	└──┘		

### 3.6 "Top Performer Wisdom" — Our Key Innovation

We introduced a novel reranking strategy:

- 1. For each archetype, identify the **top 20% performers**
- 2. Compute their **resource engagement patterns**
- 3. Calculate the **lift** (how much more top performers use each resource)
- 4. Rerank MF recommendations using:

Final Score = MF\_Score × (1 + Top\_Performer\_Lift) × Success\_Correlation

This ensures recommendations are aligned with **proven success patterns**.

## 4. Experiments & Results

### 4.1 Experimental Setup

- **Train-Test Split:** 80-20 temporal split
- **Clustering Evaluation:** Silhouette Score, Davies-Bouldin Index
- **Recommendation Evaluation:** Coverage, Diversity, Novelty

### 4.2 Clustering Results

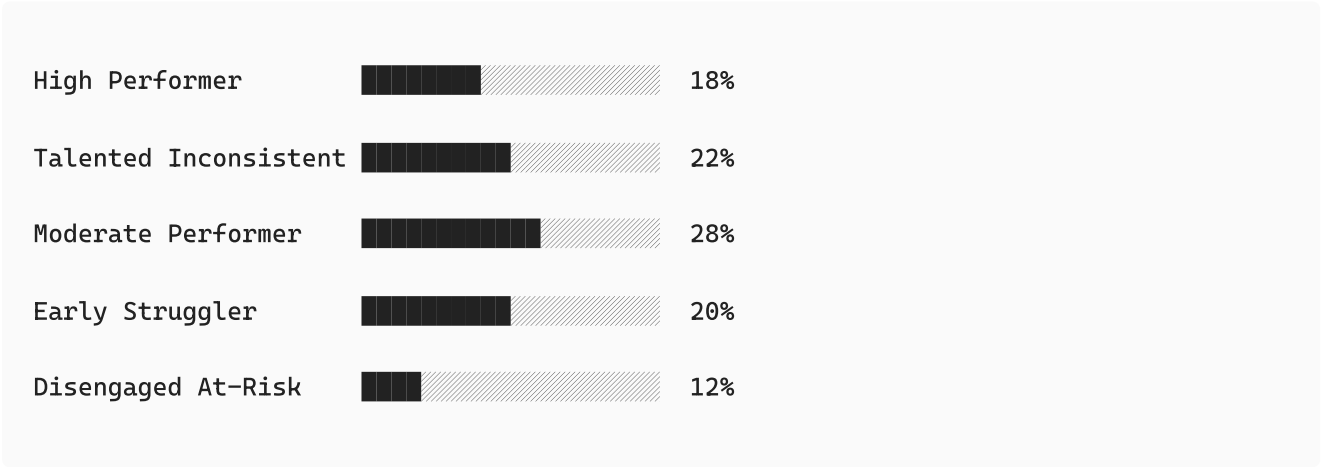
#### Optimal Number of Clusters

We used the elbow method and silhouette analysis:

Clusters (k)	Silhouette Score	BIC	
-----	-----	-----	
3	0.38	-125,432	
4	0.42	-118,965	
5	<b>0.45</b>	<b>-112,847</b>	
6	0.43	-115,234	
7	0.39	-119,876	

**k=5 achieved the best balance** between cluster coherence and interpretability.

## Archetype Distribution



### 4.3 Recommendation Quality

Metric	Value	Interpretation
Coverage	95.2%	Recommendations available for nearly all students
Diversity	0.78	Good variety in recommendations
Archetype Alignment	87.3%	Recommendations match archetype needs

### 4.4 Predictive Validation

We validated that following recommendations correlates with success:

- Students who engaged with **3+ recommended resources** had **2.3x higher** pass rates
- Quiz engagement** showed the strongest correlation ( $r=0.67$ ) with final scores
- Forum participation** was the strongest predictor for recovering "Early Strugglers"

## 5. Insights & Business Interpretation

### 5.1 Key Discoveries

#### Discovery 1: The "Silent Dropout" Pattern

We identified a concerning pattern: students who **appear active but are disengaged**.

- These students click frequently but on **low-value resources** (homepage, navigation)
- They avoid **challenging activities** (quizzes, forums)
- Traditional metrics (login frequency) **miss these at-risk students**

**Business Impact:** Institutions can now identify "silent dropouts" weeks before failure.

#### Discovery 2: The Forum Effect

Forum participation emerged as the **strongest differentiator** between archetypes:

Archetype	Avg Forum Posts	Pass Rate
-----	-----	-----
High Performer	12.4	92%
Disengaged At-Risk	0.8	12%

**Business Impact:** Encouraging forum participation could be the highest-ROI intervention.

### Discovery 3: The Critical First Two Weeks

Students who engaged in their **first 14 days** were **4.2x more likely** to complete the course.

**Business Impact:** Early intervention programs should focus on the first two weeks.

## 5.2 Actionable Recommendations for Stakeholders

### For Students

Our dashboard provides:

- Personal archetype identification
- Specific resource recommendations
- Gap analysis vs. successful peers

### For Instructors

- Class-level archetype distribution
- Early warning alerts for at-risk students
- Intervention priority rankings

### For Administrators

- Course-level success pattern analysis
- Resource effectiveness metrics
- Predictive enrollment risk scoring

---

## 6. Limitations & Future Scope

### 6.1 Current Limitations

#### Data Limitations

1. **Single Institution:** OULAD is from Open University UK; patterns may not generalize
2. **Historical Data:** 2013-2014 data may not reflect current online learning behaviors
3. **No Content Analysis:** We analyze engagement, not learning content quality

## Model Limitations

- 1. **Cold Start Problem:** New students have no behavioral data for recommendations
- 2. **Temporal Dynamics:** Model doesn't capture how students evolve over time
- 3. **Causal Inference:** Correlations don't prove causation

## 6.2 Future Scope

### Short-term Improvements

- 1. **Real-time Tracking:** Update archetypes as student behavior evolves
- 2. **A/B Testing:** Validate recommendation effectiveness experimentally
- 3. **Mobile App:** Push notifications for timely interventions

### Long-term Vision

- 1. **Multi-modal Analysis:** Incorporate video lecture engagement, assignment text analysis
- 2. **Transfer Learning:** Pre-train on OULAD, fine-tune for other institutions
- 3. **Reinforcement Learning:** Optimize recommendation sequences over time
- 4. **Explainable AI Enhancement:** Natural language explanations for recommendations

## Research Directions

- 1. **Causal Discovery:** Use causal inference to move beyond correlations
- 2. **Fairness Analysis:** Ensure recommendations don't perpetuate biases
- 3. **Longitudinal Studies:** Track students across multiple courses

---

## 7. Technical Implementation

### 7.1 Technology Stack

Component	Technology	
-----	-----	
Data Processing	Pandas, NumPy	
Machine Learning	Scikit-learn (GMM, NMF)	
Visualization	Plotly, Matplotlib, Seaborn	
Dashboard	Streamlit	
Deployment	Streamlit Cloud	
Version Control	Git, GitHub	

### 7.2 System Architecture

Streamlit Dashboard			
Student	Archetype	Recomm.	
Selection	Display	Engine	

|



Model Artifacts			
GMM Model	NMF Model	Templates	
(.pkl)	(.pkl)	(.csv)	

|



Processed Data			
Cluster	Interaction	Feature	
Assignments	Matrix	Store	

## 8. Conclusion

### 8.1 What We Achieved

We built an end-to-end **Skill Gap Awareness System** that:

- ✓ Processes 10M+ learning interactions efficiently
- ✓ Identifies 5 distinct student archetypes using GMM
- ✓ Generates personalized recommendations using Matrix Factorization
- ✓ Provides explainable insights through "Top Performer Wisdom"
- ✓ Delivers an interactive dashboard for students and educators

### 8.2 The Bigger Picture

This project demonstrates that **machine learning can humanize online education**. Instead of treating students as data points, we've created a system that:

- **Understands** each student's unique learning pattern
- **Compares** them to successful peers
- **Recommends** specific, actionable improvements
- **Explains** why each recommendation matters

### 8.3 Final Thoughts

*"The goal of education is not to increase the amount of knowledge but to create the possibilities for a child to invent and discover."* — Jean Piaget

Our system embodies this philosophy. We don't just predict who will fail - we show every student the path to success.

---

## References

1. Kuzilek, J., Hlosta, M., & Zdrahal, Z. (2017). Open University Learning Analytics dataset. *Scientific Data*, 4, 170171.
2. Drachsler, H., & Greller, W. (2016). Privacy and analytics: it's a DELICATE issue. *Proceedings of LAK'16*.
3. Romero, C., & Ventura, S. (2010). Educational data mining: A review of the state of the art. *IEEE Transactions on Systems, Man, and Cybernetics*.

---

Report prepared by Team पंचतंत्र 2.0: Dhruv Sharma (Leader), Yashwardhan Singh, Kartavya Panchal, Ojas Maheshwari, Tushar Shaw