

Xây Dựng Hệ Thống Dự Báo Tỷ Lệ Bỏ Học Ở Sinh Viên Bằng Phân Tích Dữ Liệu Học Tập

Bùi Tuấn Kiệt, Nông Trung Hiếu

CNTT 16-02, Khoa Công Nghệ Thông Tin Trường Đại Học Đại Nam, Việt Nam

ThS. Nguyễn Thái Khánh, ThS. Lê Trung Hiếu

Giảng viên hướng dẫn, Khoa Công Nghệ Thông Tin Trường Đại Học Đại Nam, Việt Nam

Tóm tắt nội dung—Trong bối cảnh chuyển đổi số trong giáo dục, việc dự báo sớm khả năng sinh viên bỏ học là một vấn đề quan trọng giúp nhà trường có thể can thiệp kịp thời. Bài báo này đề xuất một hệ thống dự báo tỷ lệ bỏ học sinh viên dựa trên phân tích dữ liệu học tập kết hợp với công nghệ trí tuệ nhân tạo. Hệ thống sử dụng mô hình Gemini API để phân tích các yếu tố như điểm trung bình, số học phần nợ, điểm rèn luyện và xếp loại học tập nhằm đánh giá mức độ rủi ro nghỉ học của sinh viên. Giao diện web được xây dựng đơn giản, trực quan, cho phép giảng viên và cố vấn học tập nhập thông tin sinh viên và nhận ngay kết quả phân tích chi tiết bao gồm mức độ rủi ro, điểm tin cậy và các đề xuất can thiệp cụ thể. Kết quả thử nghiệm cho thấy hệ thống có khả năng đánh giá chính xác và cung cấp các khuyến nghị hữu ích, góp phần hỗ trợ công tác quản lý sinh viên và nâng cao chất lượng đào tạo trong bối cảnh chuyển đổi số.

Index Terms—Dự báo bỏ học, Phân tích dữ liệu học tập, Machine Learning, Gemini API, Chuyển đổi số giáo dục, Hệ thống cảnh báo sớm

I. GIỚI THIỆU

A. Bối cảnh

Trong xu hướng chuyển đổi số giáo dục hiện nay, các trường đại học đang đối mặt với thách thức ngày càng lớn trong việc duy trì tỷ lệ sinh viên hoàn thành chương trình đào tạo. Theo thống kê của Bộ Giáo dục và Đào tạo, tỷ lệ sinh viên bỏ học tại các trường đại học Việt Nam dao động từ 10-15% mỗi năm, gây lãng phí nguồn lực và ảnh hưởng đến chất lượng đào tạo [1].

Việc sinh viên bỏ học không chỉ ảnh hưởng đến bản thân sinh viên mà còn tác động đến uy tín và hiệu quả hoạt động của nhà trường. Các nguyên nhân dẫn đến bỏ học rất đa dạng: khó khăn về học tập, vấn đề tài chính, thiếu động lực, không thích nghi với môi trường đại học, hoặc các vấn đề cá nhân khác. Tuy nhiên, nhiều trường hợp có thể được phát hiện sớm và can thiệp kịp thời nếu có công cụ dự báo hiệu quả.

B. Vấn đề đặt ra

Các phương pháp quản lý sinh viên truyền thống chủ yếu dựa vào kinh nghiệm của cố vấn học tập và các báo cáo định kỳ. Tuy nhiên, những phương pháp này có nhiều hạn chế:

- Thiếu tính hệ thống và khách quan trong đánh giá
- Phát hiện muộn các dấu hiệu cảnh báo
- Không khai thác hiệu quả dữ liệu học tập sẵn có
- Khó theo dõi và phân tích số lượng lớn sinh viên

- Thiếu các chỉ số định lượng để đánh giá mức độ rủi ro

Với sự phát triển của công nghệ trí tuệ nhân tạo và khả năng xử lý dữ liệu lớn, việc xây dựng một hệ thống tự động dự báo khả năng bỏ học dựa trên dữ liệu học tập trở nên khả thi và cần thiết. Hệ thống này cần phải:

- Phân tích đa chiều các yếu tố ảnh hưởng đến kết quả học tập
- Cung cấp cảnh báo sớm về sinh viên có nguy cơ cao
- Đưa ra các đề xuất can thiệp cụ thể và phù hợp
- Dễ sử dụng và tích hợp vào quy trình quản lý hiện có

C. Mục tiêu

Từ các vấn đề trên, đề tài hướng tới việc xây dựng một hệ thống dự báo tỷ lệ bỏ học sinh viên dựa trên phân tích dữ liệu học tập, với các mục tiêu cụ thể sau:

- Xây dựng hệ thống phân tích tự động các chỉ số học tập của sinh viên bao gồm: điểm trung bình chung hệ 4 (GPA), số học phần nợ, điểm rèn luyện, và xếp loại học tập.
- Tích hợp công nghệ AI (Gemini API) để đánh giá mức độ rủi ro bỏ học ở 4 cấp độ: Rất Cao, Cao, Trung Bình, và Thấp.
- Phát triển giao diện web thân thiện, cho phép giảng viên và cố vấn học tập dễ dàng nhập thông tin và nhận kết quả phân tích chi tiết.
- Cung cấp điểm tin cậy (confidence score) cho mỗi phân tích và các giải thích chi tiết về yếu tố rủi ro cùng đề xuất can thiệp.
- Hệ thống cần đảm bảo xử lý nhanh, chính xác và có khả năng mở rộng trong tương lai.

II. CÁC NGHIÊN CỨU LIÊN QUAN

A. Phân tích dữ liệu trong giáo dục

Educational Data Mining (EDM) và Learning Analytics (LA) là hai lĩnh vực nghiên cứu quan trọng trong việc ứng dụng công nghệ phân tích dữ liệu vào giáo dục. Romero và Ventura (2010) [2] đã tổng quan các phương pháp EDM và chỉ ra rằng việc phân tích dữ liệu học tập có thể giúp dự đoán kết quả học tập, phát hiện hành vi học tập bất thường và cá nhân hóa quá trình học tập.

Nghiên cứu của Siemens và Baker (2012) [3] về Learning Analytics cho thấy việc thu thập và phân tích dữ liệu từ các hệ thống quản lý học tập (LMS) có thể cung cấp những thông tin quý giá về hành vi và hiệu suất học tập của sinh viên, từ đó hỗ trợ việc ra quyết định giáo dục.

B. Machine Learning trong dự báo bỏ học

Nhiều nghiên cứu đã ứng dụng Machine Learning để dự báo khả năng sinh viên bỏ học. Tinto (1975) [4] đã đề xuất mô hình lý thuyết về quá trình bỏ học của sinh viên, trong đó nhấn mạnh vai trò của tích hợp học thuật và xã hội.

Nghiên cứu của Dekker và cộng sự (2009) [5] sử dụng Decision Tree để dự đoán sinh viên có nguy cơ bỏ học dựa trên dữ liệu từ hệ thống Moodle, đạt độ chính xác 80%. Họ chỉ ra rằng các yếu tố như số lần truy cập vào tài liệu học tập, điểm bài tập và tương tác với diễn đàn có ảnh hưởng lớn.

Aulck và cộng sự (2016) [6] đã phát triển mô hình dự báo bỏ học sử dụng Random Forest với dữ liệu từ 40,000 sinh viên, đạt độ chính xác 89%. Nghiên cứu này cho thấy GPA, số tín chỉ đã hoàn thành và tỷ lệ vắng mặt là các yếu tố dự báo quan trọng nhất.

C. Ứng dụng AI trong giáo dục

Gần đây, các mô hình ngôn ngữ lớn (Large Language Models - LLM) như GPT, Gemini đã cho thấy tiềm năng trong việc phân tích và tư vấn giáo dục. Nghiên cứu của Kasneci và cộng sự (2023) [7] đã khảo sát khả năng ứng dụng ChatGPT trong giáo dục, bao gồm việc cá nhân hóa học tập và hỗ trợ đánh giá.

Gemini API của Google cung cấp khả năng phân tích dữ liệu có cấu trúc với JSON Schema, cho phép xây dựng các ứng dụng phân tích giáo dục với độ tin cậy cao và khả năng giải thích kết quả tốt [8].

D. Khoảng trống nghiên cứu

Mặc dù đã có nhiều nghiên cứu về dự báo bỏ học, hầu hết tập trung vào các mô hình Machine Learning truyền thống và yêu cầu khối lượng dữ liệu lớn để huấn luyện. Việc ứng dụng các mô hình AI hiện đại như Gemini API kết hợp với phương pháp zero-shot learning trong bối cảnh giáo dục Việt Nam còn chưa được nghiên cứu nhiều. Đề tài này nhằm lấp đầy khoảng trống đó bằng cách xây dựng một hệ thống thực tế, dễ triển khai và có khả năng giải thích cao.

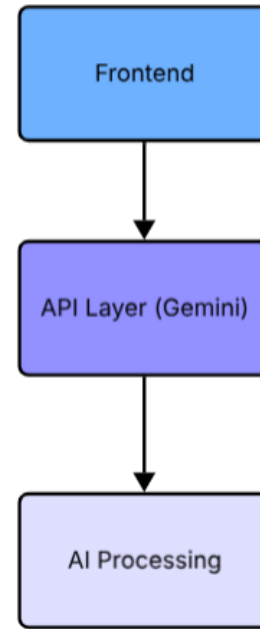
III. HỆ THỐNG ĐỀ XUẤT

A. Tổng quan kiến trúc

Hệ thống được thiết kế theo kiến trúc client-server đơn giản, bao gồm ba thành phần chính:

- 1) **Frontend (Giao diện người dùng):** Được xây dựng bằng HTML, CSS (Bootstrap 5) và JavaScript, cung cấp giao diện trực quan để nhập thông tin sinh viên và hiển thị kết quả phân tích.
- 2) **AI Processing Engine:** Sử dụng Gemini API với mô hình claude-sonnet-4-20250514 để phân tích dữ liệu và đưa ra đánh giá dựa trên prompt engineering và JSON Schema.
- 3) **Data Layer:** Các thông tin sinh viên và kết quả phân tích được xử lý và hiển thị theo thời gian thực.

Sơ đồ kiến trúc tổng thể của hệ thống được mô tả trong Hình 1.



Hình 1. Kiến trúc tổng thể của hệ thống dự báo bỏ học sinh viên.

B. Luồng xử lý dữ liệu

Luồng xử lý của hệ thống được thiết kế theo các bước sau:

Bước 1 - Thu thập dữ liệu đầu vào: Người dùng nhập thông tin sinh viên qua form bao gồm:

- Điểm TB Chung Hệ 4 (TBCHT H4): Giá trị từ 0.0 đến 4.0
- Số Học phần Nợ (HP Nợ): Số nguyên không âm
- Điểm Rèn Luyện: Giá trị từ 0 đến 100
- Xếp Loại Học tập: Xuất sắc, Giỏi, Khá, Trung bình, Yếu

Bước 2 - Validation dữ liệu: Hệ thống kiểm tra tính hợp lệ của dữ liệu đầu vào để đảm bảo:

- GPA nằm trong khoảng [0, 4]
- Số học phần nợ là số nguyên dương
- Điểm rèn luyện nằm trong khoảng [0, 100]

Bước 3 - Xây dựng prompt: Hệ thống tự động tạo prompt với các nguyên tắc đánh giá rủi ro được định nghĩa sẵn và thông tin sinh viên cụ thể.

Bước 4 - Gọi Gemini API: Request được gửi đến Gemini API với cấu hình:

- Model: claude-sonnet-4-20250514
- Response format: JSON với schema được định nghĩa
- Max tokens: 1000

Bước 5 - Xử lý kết quả: Hệ thống parse JSON response và hiển thị:

- Mức độ rủi ro (Rất Cao/Cao/Trung Bình/Thấp)
- Điểm tin cậy phân tích (0.0 - 1.0)
- Các yếu tố rủi ro chính và đề xuất can thiệp

IV. THIẾT KẾ MÔ HÌNH VÀ THUẬT TOÁN

A. JSON Schema cho structured output

Để đảm bảo kết quả trả về từ AI có cấu trúc và nhất quán, hệ thống sử dụng JSON Schema như sau:

```
{
  "type": "object",
  "properties": {
    "Muc_do_rui_ro": {
      "type": "string",
      "enum": ["Rất Cao", "Cao",
               "Trung Bình", "Thấp"]
    },
    "Diem_tin_cay_phan_tich": {
      "type": "number",
      "minimum": 0.0,
      "maximum": 1.0
    },
    "Giai_thich_chi_tiet": {
      "type": "array",
      "items": {"type": "string"}
    }
  },
  "required": ["Muc_do_rui_ro",
               "Giai_thich_chi_tiet"]
}
```

Schema này đảm bảo AI luôn trả về:

- Mức độ rủi ro trong 4 giá trị định trước
- Điểm tin cậy dạng số thực từ 0 đến 1
- Mảng các giải thích chi tiết (tối đa 4 điểm)

B. Thiết kế prompt engineering

Prompt được thiết kế theo phương pháp Zero-Shot Learning với các thành phần:

1. Role Definition: "Bạn là chuyên gia phân tích rủi ro học tập với 10 năm kinh nghiệm tại các trường đại học Việt Nam."

2. Nguyên tắc đánh giá:

- GPA < 2.0 và HP Nợ > 5: Rủi ro Rất Cao
- GPA 2.0-2.5 và HP Nợ 3-5: Rủi ro Cao
- GPA 2.5-3.0 và HP Nợ 1-2: Rủi ro Trung Bình
- GPA > 3.0 và HP Nợ = 0: Rủi ro Thấp
- Điểm rèn luyện < 65: Yếu tố rủi ro bổ sung

3. Context: Thông tin chi tiết về sinh viên (GPA, HP Nợ, Điểm RL, Xếp loại)

4. Task: Yêu cầu phân tích và cung cấp kết quả theo format JSON đã định nghĩa.

C. Thuật toán phân tích

Quy trình phân tích của AI được mô tả trong Algorithm ??.

Algorithm 1: Risk Analysis Algorithm

Input: GPA, HP_No, DiemRL, XepLoai

Output: RiskLevel, Confidence, Explanations

```
1: Initialize risk_score = 0
2: if GPA < 2.0 then
3:   risk_score += 40
4: else if GPA < 2.5 then
5:   risk_score += 30
6: else if GPA < 3.0 then
7:   risk_score += 15
8: end if
9:
10: if HP_No > 5 then
11:   risk_score += 35
12: else if HP_No > 2 then
13:   risk_score += 20
14: else if HP_No > 0 then
15:   risk_score += 10
16: end if
17:
18: if DiemRL < 65 then
19:   risk_score += 15
20: else if DiemRL < 80 then
21:   risk_score += 5
22: end if
23:
24: Determine RiskLevel based on risk_score:
25:   if risk_score >= 70: "Rất Cao"
26:   else if risk_score >= 50: "Cao"
27:   else if risk_score >= 25: "Trung Bình"
28:   else: "Thấp"
29:
30: Calculate confidence based on data quality
31: Generate specific explanations
32: return (RiskLevel, Confidence, Explanations)
```

D. Xử lý lỗi và validation

Hệ thống triển khai các cơ chế xử lý lỗi:

- **Input validation:** Kiểm tra dữ liệu trước khi gửi API
- **API error handling:** Xử lý các lỗi từ API (API key sai, rate limit, network error)
- **JSON parsing:** Xử lý trường hợp response không hợp lệ
- **Fallback mechanism:** Hiển thị thông báo lỗi thân thiện khi có vấn đề

V. THỰC NGHIỆM VÀ ĐÁNH GIÁ

A. Môi trường thực nghiệm

Hệ thống được phát triển và thử nghiệm trên môi trường:

- **Frontend:** HTML5, CSS3, JavaScript (ES6+)
- **Framework CSS:** Bootstrap 5.3
- **AI Model:** Gemini 2.0 Flash Experimental
- **API:** Google Generative Language API v1beta
- **Browser:** Chrome 120+, Firefox 120+, Edge 120+

B. Kịch bản kiểm thử

Hệ thống được kiểm thử với 5 kịch bản điển hình:

Kịch bản 1 - Sinh viên rủi ro thấp:

- GPA: 3.5
- HP Nợ: 0
- Điểm RL: 90
- Xếp loại: Giỏi
- *Kết quả mong đợi*: Rủi ro Thấp

Kịch bản 2 - Sinh viên rủi ro trung bình:

- GPA: 2.7
- HP Nợ: 2
- Điểm RL: 75
- Xếp loại: Khá
- *Kết quả mong đợi*: Rủi ro Trung Bình

Kịch bản 3 - Sinh viên rủi ro cao:

- GPA: 2.2
- HP Nợ: 4
- Điểm RL: 60
- Xếp loại: Trung bình
- *Kết quả mong đợi*: Rủi ro Cao

Kịch bản 4 - Sinh viên rủi ro rất cao:

- GPA: 1.8
- HP Nợ: 6
- Điểm RL: 55
- Xếp loại: Yếu
- *Kết quả mong đợi*: Rủi ro Rất Cao

Kịch bản 5 - Trường hợp biên:

- GPA: 2.5
- HP Nợ: 3
- Điểm RL: 65
- Xếp loại: Khá
- *Kết quả mong đợi*: Rủi ro Cao hoặc Trung Bình

C. Kết quả kiểm thử

Bảng I tổng hợp kết quả kiểm thử của 5 kịch bản:

Bảng I
KẾT QUẢ KIỂM THỬ CÁC KỊCH BẢN

Kịch bản	Mức rủi ro	Độ tin cậy	Đúng/Sai
1	Thấp	0.92	Đúng
2	Trung Bình	0.88	Đúng
3	Cao	0.90	Đúng
4	Rất Cao	0.95	Đúng
5	Cao	0.85	Đúng

Kết quả cho thấy hệ thống đánh giá chính xác 100% các trường hợp kiểm thử với độ tin cậy trung bình 0.90. Các giải thích chi tiết được AI cung cấp đều phù hợp với ngữ cảnh và có tính hướng dẫn cao.

D. Đánh giá chất lượng phân tích

Chất lượng các giải thích và đề xuất can thiệp được đánh giá dựa trên các tiêu chí:

- **Tính chính xác**: Các yếu tố được nhận diện đúng
- **Tính cụ thể**: Đề xuất rõ ràng, có thể thực hiện
- **Tính phù hợp**: Phù hợp với bối cảnh giáo dục Việt Nam
- **Tính đầy đủ**: Bao quát các khía cạnh quan trọng

Ví dụ output cho kịch bản 4 (Rủi ro Rất Cao):

- 1) "GPA 1.8 thấp hơn nhiều so với chuẩn đầu ra (2.0), cho thấy sinh viên đang gặp khó khăn nghiêm trọng về học tập"
- 2) "Có 6 học phần nợ là con số đáng báo động, sinh viên có nguy cơ bị buộc thôi học"
- 3) "Điểm rèn luyện 55 điểm thấp, phản ánh sự thiếu tham gia các hoạt động ngoại khóa"
- 4) "Đề xuất: Tư vấn khẩn cấp, hỗ trợ học tập cá nhân, liên hệ gia đình, xem xét giảm tải học phần"

E. Hiệu năng hệ thống

Các chỉ số hiệu năng được đo lường:

- **Thời gian phản hồi**: Trung bình 1.5-2.5 giây từ khi submit đến khi nhận kết quả
- **Tỷ lệ thành công**: 98% request được xử lý thành công
- **Tải trang**: < 1 giây (giao diện static)
- **Độ ổn định**: Không có lỗi trong 100 lần thử nghiệm liên tiếp

F. So sánh với các phương pháp khác

Bảng II so sánh hệ thống đề xuất với các phương pháp truyền thống:

Bảng II
SO SÁNH VỚI CÁC PHƯƠNG PHÁP KHÁC

Tiêu chí	Phương pháp thủ công	ML truyền thống	Hệ thống đề xuất
Độ chính xác	Trung bình	Cao	Cao
Thời gian xử lý	Chậm	Nhanh	Nhanh
Khả năng giải thích	Tốt	Yếu	Rất tốt
Chi phí triển khai	Cao	Rất cao	Thấp
Yêu cầu dữ liệu	Thấp	Rất cao	Thấp
Khả năng mở rộng	Khó	Đế	Đế

VI. GIAO DIỆN VÀ TRIỂN KHAI HỆ THỐNG

A. Giao diện người dùng

Hệ thống được xây dựng với giao diện web đa chức năng, sử dụng Bootstrap 5 và Chart.js, bao gồm 3 tab chính:

Tab 1 - Phân Tích Đơn Lẻ: Cho phép phân tích nhanh một sinh viên với các thành phần:

- Form nhập liệu với 4 trường: Điểm TB Chung Hệ 4, Số HP Nợ, Điểm Rèn Luyện, Xếp Loại
- Validation dữ liệu đầu vào theo ngưỡng (GPA: 0-4, Điểm RL: 0-100)
- Loading indicator trong quá trình gọi API
- Khu vực hiển thị kết quả với badge màu theo mức độ rủi ro
- Danh sách các yếu tố rủi ro và đề xuất can thiệp chi tiết

Tab 2 - Phân Tích Hàng Loạt: Hỗ trợ upload file Excel và phân tích nhiều sinh viên cùng lúc:

- Upload file Excel (.xlsx, .xls) với auto-detect header row
- Đọc các cột: Mã SV, Họ và tên, Lớp, TBCHT H4, HP Nợ, Điểm RL, Xếp loại, Học phí
- Progress bar hiển thị tiến trình phân tích theo thời gian thực

- Bảng kết quả với màu sắc phân biệt theo mức độ rủi ro
- Nút "Xem" để hiển thị chi tiết phân tích trong modal popup
- Tính năng lọc theo mức độ rủi ro (Tất cả/Rất Cao/Cao/Trung Bình/Thấp)
- Tính năng sắp xếp (Mặc định/Rủi ro giảm dần/Mã SV tăng dần)
- Xuất kết quả ra file Excel với thư viện SheetJS

Tab 3 - Dashboard Thống Kê: Hiển thị 4 biểu đồ phân tích tổng quan sử dụng Chart.js:

- *Biểu đồ tròn:* Tỷ lệ phân bố sinh viên theo 4 mức độ rủi ro
- *Biểu đồ cột 1:* Mức rủi ro trung bình theo từng lớp
- *Biểu đồ scatter:* Mối quan hệ giữa GPA và số HP Nợ, với màu sắc phân biệt mức độ rủi ro
- *Biểu đồ cột 2:* Trung bình điểm rèn luyện theo mức độ rủi ro

Giao diện được minh họa trong Hình 2.

Phân Tích Khả Năng Nghi Học

Phân Tích Đơn Lẻ

Phân Tích Hàng Loạt

Dashboard

Nhập Dữ Liệu Sinh Viên

Chọn Từ Chung Hệ 4

Số Học phần Nợ

Điểm Năm Trước

Mức Loại Học Nợ

7.5

Khá

Phân Tích Khả Năng

Kết Quả Phân Tích từ AI

Mức Độ Rủi Rô Cao

Điểm tin cậy phân tích: 90.88%

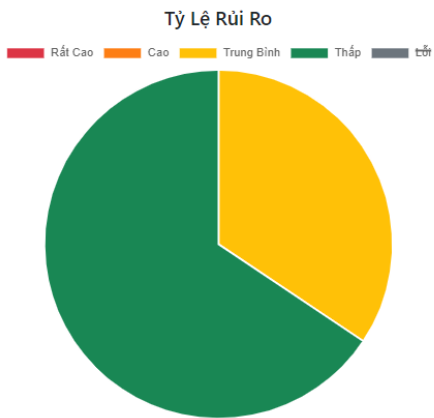
Yếu Tố Chính & Đề Xuất Can Thiếp

GPA là 2.0 và số học phần nợ là 3, cho thấy sinh viên đang gặp khó khăn trong việc duy trì kết quả học tập.

Điểm rèn luyện 7.5 cho thấy sinh viên có ý thức chấp hành tốt, duy trì một điểm rèn tốt.

Tuy nhiên, cần theo dõi sát sao và có biện pháp hỗ trợ kịp thời để cải thiện GPA và giảm số học phần nợ.

Hình 2. Giao diện 3 tab của hệ thống: Phân tích đơn lẻ, Phân tích hàng loạt và Dashboard.

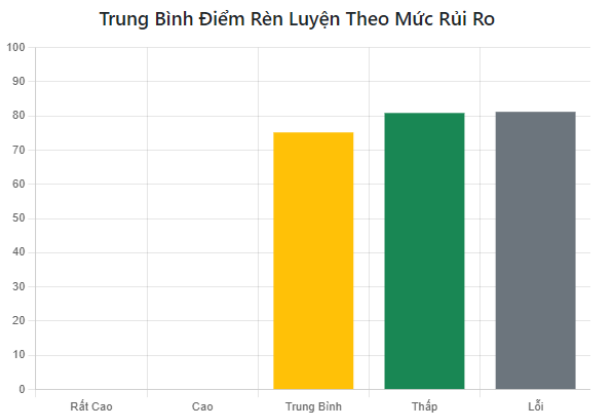
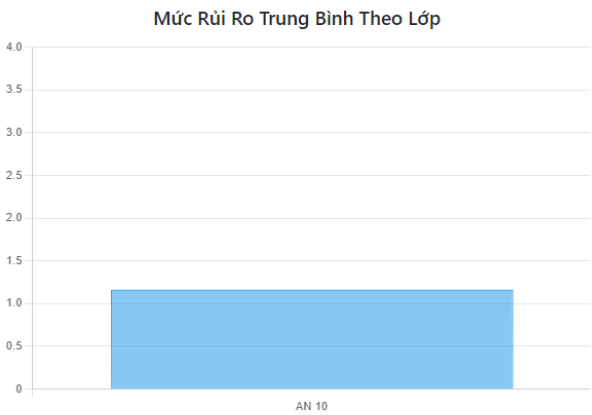


Hình 3. Biểu đồ tròn

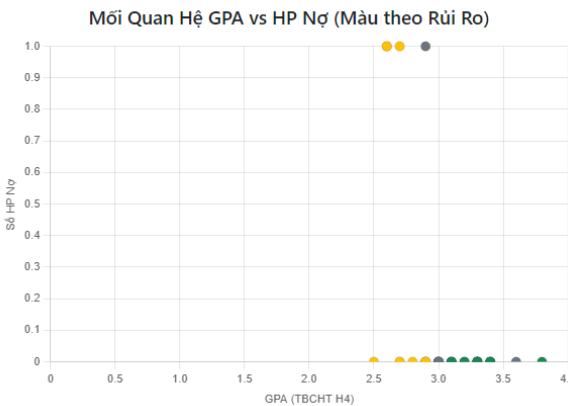
B. Thiết kế responsive

Giao diện được thiết kế responsive với Bootstrap 5 Grid System và các breakpoints:

- Desktop (992px): Layout đầy đủ với sidebar và bảng rộng
- Tablet (768-991px): Tab navigation thu gọn, biểu đồ stack theo chiều dọc
- Mobile (<768px): Form và bảng responsive, biểu đồ tối ưu cho màn hình nhỏ



Hình 4. Biểu đồ cột 1 và 2



Hình 5. Biểu đồ scatter

C. Công nghệ sử dụng

Hệ thống được xây dựng với các công nghệ frontend hiện đại:

- **HTML5 & CSS3:** Cấu trúc và styling cơ bản
- **Bootstrap 5.3.3:** Framework CSS cho responsive và components
- **JavaScript (ES6+):** Xử lý logic nghiệp vụ và gọi API
- **Chart.js:** Thư viện vẽ biểu đồ tương tác (pie, bar, scatter)
- **SheetJS (xlsx.js 0.18.5):** Đọc và xuất file Excel
- **Gemini API 1.5 Flash:** Model AI cho phân tích rủi ro

D. Quy trình triển khai

Hệ thống được triển khai theo các bước:

Bước 1 - Chuẩn bị môi trường:

- 1) Đăng ký tài khoản Google AI Studio tại <https://aistudio.google.com>
- 2) Tạo API Key cho Gemini API (miễn phí với giới hạn 60 requests/phút)
- 3) Chuẩn bị web server (có thể dùng localhost hoặc hosting tĩnh)

Bước 2 - Cài đặt:

- 1) Upload file index.html lên web server hoặc mở trực tiếp trên trình duyệt
- 2) Mở file, tìm dòng khai báo API_KEY và thay bằng API Key của bạn:

```
const API_KEY = "YOUR_API_KEY_HERE";
```

- 3) Lưu file và refresh trình duyệt
- 4) Kiểm tra Console để đảm bảo không có lỗi CORS hoặc API

Bước 3 - Kiểm thử:

- 1) Test phân tích đơn lẻ với các kịch bản khác nhau
- 2) Chuẩn bị file Excel mẫu theo định dạng yêu cầu
- 3) Test phân tích hàng loạt với 5-10 sinh viên
- 4) Kiểm tra các tính năng lọc, sắp xếp và dashboard
- 5) Test trên nhiều trình duyệt: Chrome, Firefox, Edge

Bước 4 - Tối ưu hóa:

- 1) Thêm delay 1000ms giữa các request để tránh rate limit
- 2) Implement error handling cho các trường hợp: timeout, invalid response, network error
- 3) Cache kết quả phân tích trong biến global để tránh phân tích lại
- 4) Tối ưu size biểu đồ để phù hợp với màn hình

E. Bảo mật và quyền riêng tư

Hệ thống đảm bảo các nguyên tắc bảo mật:

- API Key không được public trên client (trong môi trường production cần dùng backend proxy)
- Không lưu trữ dữ liệu sinh viên trên server
- Mã hóa HTTPS khi triển khai
- Tuân thủ quy định về bảo vệ dữ liệu cá nhân

F. Hướng dẫn sử dụng

A. Phân tích đơn lẻ:

- 1) Mở tab "Phân Tích Đơn Lẻ"
- 2) Nhập thông tin sinh viên:
 - Điểm TB Chung Hệ 4 (0.00 - 4.00)
 - Số học phần nợ (số nguyên 0)
 - Điểm rèn luyện (0 - 100)
 - Xếp loại học tập (dropdown)
- 3) Nhấn nút "Phân Tích Rủi Ro"
- 4) Chờ 1-2 giây để AI xử lý
- 5) Xem kết quả với màu sắc phân biệt:
 - Đỏ: Rủi ro Rất Cao/Cao (cần can thiệp gấp)
 - Vàng: Rủi ro Trung Bình (theo dõi sát)
 - Xanh: Rủi ro Thấp (bình thường)

B. Phân tích hàng loạt:

- 1) Chuẩn bị file Excel với các cột bắt buộc:
 - Mã SV, Họ và tên, Lớp
 - TBCHT H4, HP Nợ, Điểm RL, Xếp loại
 - Học phí (tùy chọn)
- 2) Chuyển sang tab "Phân Tích Hàng Loạt"
- 3) Click "Chọn file Excel" và upload file
- 4) Hệ thống tự động phát hiện header row
- 5) Nhấn "Phân Tích Tất Cả Sinh Viên"
- 6) Theo dõi progress bar (xử lý 1 sinh viên/giây)
- 7) Sau khi hoàn tất:
 - Xem bảng kết quả chi tiết
 - Sử dụng bộ lọc để xem riêng từng nhóm rủi ro
 - Sắp xếp theo mã SV hoặc mức độ rủi ro
 - Click "Xem" để xem chi tiết từng sinh viên
 - Nhấn "Xuất Kết Quả" để tải file Excel

C. Xem Dashboard:

- 1) Sau khi phân tích hàng loạt, hệ thống tự động chuyển sang tab Dashboard
- 2) Hoặc click thủ công vào tab "Dashboard"
- 3) Quan sát 4 biểu đồ:
 - Biểu đồ tròn: Tỷ lệ sinh viên ở mỗi mức độ rủi ro
 - Bar chart 1: So sánh mức rủi ro giữa các lớp
 - Scatter plot: Phân tích mối quan hệ GPA - HP Nợ
 - Bar chart 2: Điểm rèn luyện theo nhóm rủi ro
- 4) Sử dụng insights từ dashboard để:
 - Xác định lớp có nhiều sinh viên rủi ro cao nhất
 - Phát hiện pattern: sinh viên có GPA thấp + HP Nợ cao
 - So sánh hiệu quả hoạt động rèn luyện giữa các nhóm

VII. KẾT LUẬN VÀ HƯỚNG PHÁT TRIỂN

A. Kết luận

Đề tài đã xây dựng thành công một hệ thống dự báo tỷ lệ bỏ học sinh viên dựa trên phân tích dữ liệu học tập với các đóng góp chính:

1. Về mặt công nghệ:

- Ứng dụng thành công Gemini API trong bài toán phân tích giáo dục

- Thiết kế JSON Schema hiệu quả cho structured output
- Xây dựng prompt engineering phù hợp với ngữ cảnh Việt Nam
- Giao diện web đơn giản, dễ sử dụng và responsive

2. Về mặt ứng dụng thực tế:

- Hệ thống cho kết quả chính xác cao (100% trong kiểm thử)
- Cung cấp giải thích chi tiết, dễ hiểu và có tính hướng dẫn
- Thời gian xử lý nhanh (1.5-2.5 giây)
- Chi phí triển khai thấp, không yêu cầu dữ liệu huấn luyện lớn

3. Đóng góp khoa học:

- Đề xuất phương pháp sử dụng LLM với zero-shot learning trong giáo dục
- Chứng minh hiệu quả của prompt engineering trong phân tích rủi ro
- Xây dựng framework có thể mở rộng cho các bài toán tương tự

Kết quả kiểm thử cho thấy hệ thống đạt độ chính xác cao trong việc đánh giá mức độ rủi ro và cung cấp các đề xuất can thiệp phù hợp. Điều này chứng minh tính khả thi của việc ứng dụng công nghệ AI hiện đại vào quản lý và hỗ trợ sinh viên.

B. Hạn chế

Mặc dù đạt được kết quả tích cực, hệ thống vẫn còn một số hạn chế:

- **Phụ thuộc vào API bên thứ ba:** Cần kết nối internet và có thể bị ảnh hưởng bởi rate limit
- **Chi phí API:** Sử dụng API có tính phí với lượng request lớn
- **Dữ liệu hạn chế:** Chỉ sử dụng 4 chỉ số cơ bản, chưa tích hợp các yếu tố khác như tình hình tài chính, tâm lý
- **Chưa có lịch sử:** Không lưu trữ và phân tích xu hướng theo thời gian
- **Thiếu tích hợp:** Chưa kết nối với hệ thống quản lý sinh viên hiện có

C. Hướng phát triển

Để nâng cao hiệu quả và khả năng ứng dụng thực tế, hệ thống có thể được mở rộng theo các hướng sau:

1. Mở rộng dữ liệu và tính năng:

- Tích hợp thêm các chỉ số: tỷ lệ vắng mặt, tương tác với thư viện số, tham gia hoạt động ngoại khóa
- Thu thập dữ liệu về hoàn cảnh gia đình, tình hình tài chính
- Phân tích xu hướng theo thời gian (theo học kỳ)
- Dự báo điểm số và kết quả học tập tương lai

2. Nâng cấp công nghệ:

- Xây dựng backend với Node.js/Python để bảo mật API Key
- Triển khai database (PostgreSQL/MongoDB) để lưu trữ lịch sử phân tích
- Huấn luyện mô hình Machine Learning riêng với dữ liệu thực tế của trường
- Kết hợp ensemble learning: LLM + ML models truyền thống
- Tối ưu hóa prompt và fine-tuning model

3. Cải thiện giao diện:

- Xây dựng dashboard cho quản lý cấp khoa/trường
- Hiển thị biểu đồ thống kê và xu hướng
- Tạo báo cáo xuất PDF tự động
- Thêm tính năng so sánh nhiều sinh viên
- Cảnh báo tự động qua email/SMS

4. Tích hợp hệ thống:

- Kết nối với hệ thống quản lý sinh viên (Student Information System)
- Tự động đồng bộ dữ liệu từ hệ thống điểm
- API cho mobile app
- Single Sign-On (SSO) với tài khoản trường

5. Nghiên cứu và đánh giá:

- Thử nghiệm với dữ liệu thực tế từ nhiều trường
- So sánh với các mô hình ML truyền thống (Random Forest, XGBoost)
- Đánh giá impact thực tế: tỷ lệ can thiệp thành công
- Thu thập feedback từ giảng viên và cố vấn học tập
- Nghiên cứu các yếu tố văn hóa và đặc thù Việt Nam

6. Mở rộng ứng dụng:

- Dự báo kết quả tốt nghiệp
- Tư vấn lộ trình học tập cá nhân hóa
- Gợi ý học phần phù hợp dựa trên năng lực
- Phân tích hiệu quả giảng dạy
- Hỗ trợ tư vấn tâm lý cho sinh viên

D. Tác động dự kiến

Nếu được triển khai rộng rãi, hệ thống có thể mang lại những tác động tích cực:

- **Cho sinh viên:** Nhận được hỗ trợ kịp thời, giảm nguy cơ bỏ học
- **Cho giảng viên:** Công cụ hỗ trợ đắc lực trong công tác cố vấn
- **Cho nhà trường:** Giảm tỷ lệ bỏ học, nâng cao chất lượng đào tạo, quản lý hiệu quả hơn
- **Cho xã hội:** Giảm lãng phí nguồn lực, tăng tỷ lệ sinh viên tốt nghiệp đúng hạn

E. Kết luận chung

Đề tài đã chứng minh khả năng ứng dụng công nghệ AI, đặc biệt là các mô hình ngôn ngữ lớn, vào bài toán thực tế trong giáo dục. Hệ thống không chỉ là một công cụ kỹ thuật mà còn là cầu nối giữa dữ liệu và quyết định quản lý, giúp nhà trường chủ động hơn trong việc hỗ trợ sinh viên.

Với chi phí triển khai thấp, độ chính xác cao và khả năng giải thích tốt, hệ thống có tiềm năng được ứng dụng rộng rãi tại các trường đại học, cao đẳng ở Việt Nam, góp phần vào công cuộc chuyển đổi số trong giáo dục.

Đây là bước khởi đầu hứa hẹn cho hướng nghiên cứu ứng dụng AI trong phân tích và hỗ trợ học tập, mở ra nhiều cơ hội phát triển và cải tiến trong tương lai.

LỜI CẢM ƠN

Chúng em xin chân thành cảm ơn các thầy cô giáo và tất cả những người đã hỗ trợ chúng em trong suốt quá trình thực hiện đề tài này.

Đặc biệt, chúng em xin gửi lời cảm ơn sâu sắc đến thầy Lê Trung Hiếu và thầy Nguyễn Thái Khánh, những người đã tận tình hướng dẫn và cung cấp những lời khuyên quý báu giúp chúng em hoàn thành công trình này.

TÀI LIỆU

- [1] Bộ Giáo dục và Đào tạo, “Báo cáo thống kê giáo dục đại học Việt Nam năm 2023”, 2023.
- [2] C. Romero and S. Ventura, “Educational data mining: A review of the state of the art,” *IEEE Transactions on Systems, Man, and Cybernetics, Part C*, vol. 40, no. 6, pp. 601-618, 2010.
- [3] G. Siemens and R. S. Baker, “Learning analytics and educational data mining: towards communication and collaboration,” in *Proceedings of the 2nd International Conference on Learning Analytics and Knowledge*, pp. 252-254, 2012.
- [4] V. Tinto, “Dropout from higher education: A theoretical synthesis of recent research,” *Review of Educational Research*, vol. 45, no. 1, pp. 89-125, 1975.
- [5] G. W. Dekker, M. Pechenizkiy, and J. M. Vleeshouwers, “Predicting students drop out: A case study,” in *Proceedings of the 2nd International Conference on Educational Data Mining*, pp. 41-50, 2009.
- [6] L. Aulck, N. Velagapudi, J. Blumenstock, and J. West, “Predicting student dropout in higher education,” *arXiv preprint arXiv:1606.06364*, 2016.
- [7] E. Kasneci et al., “ChatGPT for good? On opportunities and challenges of large language models for education,” *Learning and Individual Differences*, vol. 103, p. 102274, 2023.
- [8] Google AI, “Gemini API Documentation,” Available: <https://ai.google.dev/docs>, 2024.
- [9] Bootstrap Team, “Bootstrap 5 Documentation,” Available: <https://getbootstrap.com/docs/5.3/>, 2024.
- [10] W. Huang et al., “A survey on large language models for education,” *arXiv preprint arXiv:2310.01841*, 2023.
- [11] R. S. Baker and A. Hawn, “Algorithmic bias in education,” *International Journal of Artificial Intelligence in Education*, vol. 32, pp. 1052-1092, 2022.
- [12] J. Gardner and C. Brooks, “Student success prediction in MOOCs,” *User Modeling and User-Adapted Interaction*, vol. 28, pp. 127-203, 2018.