

# پروژه نهایی درس داده‌کاوی

عنوان پروژه: «علموصیار» – داشبورد تحلیلی و سیستم پیشنهاددهنده اساتید

## ۱. ایده‌ی کلی پروژه «علموصیار»

داستان از کجا شروع می‌شود؟  
همه‌مون موقع انتخاب واحد می‌ریم سراغ کanal «اساتید علموص» تو تلگرام تا ببینیم بچه‌ها درباره‌ی هر استاد چی گفتن 😊 خود همین کanal در واقع یه دیتاست خفن و واقعی‌ه، فقط خیلی شلخته و پراکنده‌ست.

کاری که تو این پروژه باید بکنید اینه که:

- پیام‌های خام کanal (فایل JSON) رو تبدیل کنید به یه داشبورد تحلیلی + یه سیستم پیشنهاددهنده‌ی استاد/درس 

هدف نهایی:

- به دانشجو کمک کنه راحت‌تر استاد/درس مناسب خودش رو پیدا کنه 
- به اساتید و مدیرگروه‌ها یه دید کلی از بازخورد دانشجوها بده 

تمرکز پروژه روی این سه تا محوره:

1. **مهندسي داده + NLP**  
... Parsing, Cleaning, Normalization, Sentiment, Keywords, Clustering
2. **مدلهای داده‌کاوی**  
نظرارتی / غیرنظرارتی، خوشبندی، insight، روابطه بین متغیرها
3. **محصول نهایی**  
یه داشبورد تعاملی که یه دانشجوی عادی هم بفهمه چی به چیه + سیستم پیشنهاددهنده

## ۲. خروجی نهایی قراره چی باشه؟

در نهایت به یک وب‌اپلیکیشن تحلیلی (داشبورد) + سیستم پیشنهاددهنده استاد/درس برسیم که روی داده‌های واقعی کانال کار کنه.

کاربری که میاد تو سیستم (دانشجو یا مدیر آموزشی) باید بتونه:

### جستجو (Search)

- بر اساس نام استاد، فامیل یا اسم درس سرچ کنه.

### فیلتر (Filter)

فیلتر کنه بر اساس:

0 دانشکده

0 درس

0 ترم

0 میانگین امتیازها

0 وضعیت نمره‌دهی

0 ...

### پروفایل استاد (Professor Profile)

برای هر استاد یه نمای کلی (۳۶۰ درجه) ببینه، شامل:

- میانگین ۶ امتیاز عددی (تدريس، مدیریت کلاس، تسلط، تعامل، پاسخ‌گویی و...)
- تحلیل وضعیت نمره‌دهی (مثلاً: ۷۰٪ منصفانه، ۲۰٪ سختگیر، ۱۰٪ آسان)
- تحلیل متن کامنت‌ها:

0 امتیاز احساسات (Sentiment)

0 ابر کلمات (Word Cloud)

- روند امتیاز و رضایت در ترم‌های مختلف (اگه اطلاعات زمانی داشته باشین)

### مقایسه‌ی اساتید (Compare)

• چندتا استاد/درس رو انتخاب کنه و کنار هم مقایسه‌شون کنه:

0 میانگین امتیازها

0 وضعیت نمره‌دهی

... و sentiment 0

## (Recommender) سیستم پیشنهادهندۀ

بر اساس ترجیحات دانشجو، مثلًا:

- استاد/درس‌هایی که دوست داره
- اولویت‌ها:
  - نمره‌دهی آسان
  - امتحان سبک
  - پژوهش‌محور بودن
  - تدریس قوی
  - ... و ...

یه لیست اساتید/درس‌های پیشنهادی برگردانید.

مدل پیشنهادهندۀ می‌تونه مثلًا:

- باشه Content-Based
- ساده Collaborative Filtering
- ساده Rule-based

مهم اینه که:

- منطقی باشه
- بشه توضیح داد چرا این استادها پیشنهاد شدن 
- تو گزارش مستند شده باشه 

---

## ۳. داده از کجا می‌آید؟

### ۳.۱. منبع اصلی

بهتون یه فایل JSON می‌دیم که خروجی کامل پیام‌های کانال علموص اساتیده.  
این می‌شه دیتای خام پروژه‌تون. ([دانلود فایل JSON](#))

### ۳.۲. ساختار کلی

هر رکورد JSON = یک پیام کانال. تو ش اینا رو داریم:

- فراداده (Metadata)

- 0 شناسه پیام، تاریخ/ترم (اگه هست)، اطلاعات کانال
- متن پیام (**message\_text**) :
  - 0 اسم استاد
  - 0 دانشکده
  - 0 اسم درس
  - 0 ۶ تا امتیاز عددی (از ۱۰)
  - 0 توضیحات درباره نمره‌دهی، حضور و غیاب و... (متن آزاد)
- کامنت/توضیح دانشجو (متن بدون ساختار)

### ✳️ ۳. امتیاز اضافه

اگه خواستین بهترش کنین:

- به جای استفاده از JSON آماده، خودتون با **API تلگرام** (Telethon / Pyrogram) داده رو بکشید بیرون. این جوری داده‌های داشبورد همیشه به روز می‌مونن.
- اگه این کار رو می‌کنید، تو گزارش باید توضیح بدین:
  - 0 چطوری وصل شدین
  - 0 چطوری داده جمع کردین
  - 0 چطور ذخیره و مدیریت کردین

## ۴. فازهای پروژه؛ قدم به قدم

این یه **roadmap** پیشنهادی‌ه؛ لازم نیست مو به مو همین باشه، ولی در این حد کار باید انجام بشه.

 **فاز ۱: بارگذاری، اکتشاف و Parsing (Load, EDA & Parsing)**

 **تسک ۱ - بارگذاری و EDA**

- خوندن JSON و تبدیلش به **DataFrame** تو Pandas

• فهمیدن ساختار داده:

- 0 چند تا پیام داریم؟
- 0 حدوداً چندتا استاد مختلف داریم؟
- 0 توزیع امتیازها چطوریه؟
- 0 چقدر داده‌ی گمشده داریم؟

• کشیدن چند نمودار اکتشافی:

- هیستوگرام امتیازها
- تعداد نظرات به تفکیک دانشکده و ...
- مستندسازی EDA تو یه نوتبوک جدا (با توضیحات متنی)

### تسک ۲ – نوشتن Parser

یه تابع/کلاس پایتون که:

- متن خام پیام رو بگیره
- در عوض یه دیکشنری مرتب بد، شامل حداقل:

  - (بعداً ساخته میشه) professor\_id ◦
  - professor\_name\_raw ◦
  - department ◦
  - course\_name ◦
  - rating\_6 تا rating\_1 ◦
  - grading\_status\_raw ◦
  - attendance\_status\_raw ◦
  - comment\_text ◦
  - اگه قابل استخراج باشه date یا term ◦

نکات مهم:

- احتمالا لازمه از **Regex** و الگوهای متنی استفاده کنید.
- پارسر باید نسبت به تغییرات ریز (ایموجی، جایه‌جایی خطها، فاصله‌ها و...) مقاوم باشه.
- اگه جایی نتوونست درست `parse` کنه، اون پیام رو تو یه ستون مثل `parse_error` لاغ کنه.
- برای پارسر چندتا **Unit Test** بنویسید (روی نمونه‌ها و چندین حالت).



## فار ۲: پاکسازی و استانداردسازی (Data Cleaning & Standardization)

### تسک ۳ – یکی‌کردن اسامی اساتید

چالش:

- «دکتر الف. رضایی»
- «علیرضا رضایی»
- «رضایی»

اینا شاید همه یک نفر باشن 😊

راه حل پیشنهادی:

- نرمال سازی اولیه از جمله:
    - حذف فاصله های اضافی
    - یکسان سازی حروف عربی / فارسی
    - lowercase برای اسمی لاتین
  - ترکیب دانشکده + الگوریتم های شباهت رشتہ ای (Levenshtein, FuzzyWuzzy) و ...
  - در نهایت گروه بندی اسمی مشابه و ساختن یک professor\_id یکتا برای هر استاد.
- تسک ۴ - یکی کردن اسم درس و برچسب های متنی 
- یکی کردن اسم درس ها:
    - «مبانی طراحی ۲» و «مبانی طراحی ۲» ← یه درس واحد
  - استانداردسازی وضعیت نمره دهی و حضور و غیاب:
    - یه سری برچسب تعریف کنید مثل:
      - نمره دهی: {آسان، منصفانه، سخت‌گیر، نامشخص}
      - حضور و غیاب: {سخت‌گیر، متوسط، آزاد، نامشخص}
  - متن های آزاد مثل «نمره نمیده»، «نمره اش خوبه» رو مپ کنید به این برچسبها.
- 

### فاز ۳: Feature Engineering, NLP و مدل ها

اینجا پروژه از حالت «گزارش ساده» می‌رود تو مود تحلیل داده کاوانه جدی.

#### ۱.۴.۳.۱. تحلیل احساسات (Sentiment Analysis)

ورودی: ستون comment\_text

حداقل کار:

1. پیش‌پردازش متن فارسی ([Shekar](#) یا موارد مشابه):

◦ نرمال سازی

◦ توکن سازی

◦ حذف stopword ها

2. استفاده از یک روش ساده برای Sentiment

0 روش‌های Ruled-Based (به شرطی که عملکرد و دقت قابل تحمیل روی داده واقعی داشته باشند)

0 یا مدل آماده‌ی سبک (مثلاً یه مدل HuggingFace که روی CPU جواب بد).  
خروجی:  
یه ستون جدید مثل sentiment\_score (مثلاً [-1,1]) یا برچسب {منفی، خنثی، مثبت}.  
امتیاز اضافه: استفاده از مدل‌های ParsBERT و گزارش نتایج. ✨

---

#### ۴.۳.۲. کلمات کلیدی و Topic ها

برای هر استاد:

- متن‌ها رو نرمال‌سازی کنید + حذف کنید.
- با TF-IDF یا Frequency کلمات کلیدی رو دریارید.
- ازش استفاده کنید برای:
  - 0 ساخت ابر کلمات
  - 0 شناخت کلی حال و هوای کامنت‌ها

امتیاز اضافه: استفاده از Topic Modeling مثل LDA یا BERTopic برای پیدا کردن موضوعات اصلی نظرات. ✨

---

#### ۴.۳.۳. یک مدل نظارتی (Supervised) – اجباری

باید حداقل یه مسئله‌ی Supervised تعریف و حل کنید. مثلاً:

- پیش‌بینی overall\_score از روی:
  - 0 میانگین ۶ امتیاز
  - 0 وضعیت نمره‌دهی (کدگذاری شده)
  - 0 میانگین sentiment
- طبقه‌بندی نظرات به {مثبت/خنثی/منفی} از روی ویژگی‌های عددی و متند.
  - 0 حداقل‌ها:

- (یا یه Cross-Validation ساده) Train/Test
  - حداقل یک مدل ساده مثل:
    - Logistic Regression
    - Random Forest
    - SVM ساده
  - گزارش معیارهای ارزیابی مناسب (Accuracy، F1، MAE، RMSE و...)
  - یه تحلیل کوچیک از Feature Importance یا حداقل تفسیر منطقی نتایج.
- 

#### ۴.۳.۴. یک مدل غیرنظرارتی (*Unsupervised*) – اجباری

مثال خوب: خوشبندی استادید بر اساس پروفایلشون:

- ورودی:
  - میانگین ۶ امتیاز + وضعیت نمره‌دهی + میانگین sentiment feature ساده دیگه
- روش‌ها:
  - KMeans
  - Agglomerative Clustering
- خروجی و نمایش:
  - بعده رو کم کنید (مثلًاً PCA)
  - هر استاد رو به صورت یه نقطه روی نمودار دو بعدی نشون بدین
  - برای هر خوشه یه تفسیر انسانی بدین مثل:
    - «سختگیر ولی محبوب»
    - «آسانگیر با رضایت متوسط»
    - ...

#### ۴.۳.۵. سیستم پیشنهادهندی استاد/درس – اجباری

باید یه سیستم پیشنهادهندی کارا و منطقی بسازید. می‌تونه ساده باشه.

چند سناریو:

- **:Content-Based Filtering**
  - برای هر استاد یه بردار ویژگی بسازید (امتیازها، sentiment، کلمات کلیدی)

۰ با شباهت (cosine similarity) اساتید مشابه یه استاد/درس مورد علاقه رو

پیشنهاد بدین

**:Rule-Based •**

۰ بر اساس فیلترهای انتخابی کاربر (امتیاز بالای ۸، نمره‌دهی منصفانه، پروژه‌محور بودن

و...) اساتید رو رتبه‌بندی و پیشنهاد بدین

**: Hybrid ساده •**

۰ ترکیبی از امتیازات عددی، sentiment و شباهت کلمات

مهم ایناست که:

• سیستم منطقی باشه

• تو داشبورد شفاف باشه «چرا» این استادها پیشنهاد شدن

• تو گزارش طراحی و جایگزین‌هاش توضیح داده بشه 

---

#### Web App / فاز ۴: ساخت داشبورد

ابزار پیشنهادی

• پیشنهاد اصلی: **Streamlit**

چون پایتون محور، سریع، راحت و خوب برای DataFrame.

صفحات حداقلی داشبورد

1. صفحه‌ی اصلی (Overview):

۰ چند تا کارت شاخص:

▪ تعداد کل نظرات

▪ تعداد اساتید پوشش داده شده

▪ میانگین امتیاز کلی

۰ یکی دو نمودار:

▪ هیستوگرام توزیع امتیازها

▪ توزیع وضعیت نمره‌دهی

2. صفحه‌ی جستجو و فیلتر:

۰ جستجوی استاد/درس

۰ فیلتر بر اساس دانشکده، ترم، بازه‌ی امتیاز، وضعیت نمره‌دهی

۰ لیست اساتید مطابق جستجو، با یه خلاصه‌ی کوچیک از هرکدام

### ۳. صفحه‌ی پروفایل استاد:

۰ نمودار رادری برای ۶ امتیاز عددی

۰ نمودار دایره‌ای/دونات برای نمره‌دهی و حضور و غیاب

۰ Word Cloud مثبت و منفی

۰ خلاصه‌ی sentiment در ترم‌های مختلف (اگه داده باشه)

### ۴. صفحه‌ی مقایسه (Compare):

۰ انتخاب چند استاد

۰ نمودار میله‌ای برای مقایسه امتیازها

۰ یه جدول خلاصه مقایسه

### ۵. صفحه‌ی پیشنهادهنده (Recommender):

۰ کاربر ترجیح‌هاش رو وارد کنه (مثلاً تدریس مهم‌تره یا نمره‌دهی آسان)

۰ سیستم یه لیست مرتب از اساتید/درس‌های پیشنهادی بده

۰ کنار هر پیشنهاد یه توضیح کوتاه مثل:

- «میانگین تدریس ۹.۲، نمره‌دهی منصفانه، sentiment بالا»

## ۵. حداقل‌ها در مقابل کارهای خفن‌تر 😎

### ۱. حداقل کارهایی که باید انجام بشه

برای نمره‌ی خوب (نه لزوماً فول‌مارک)، حداقل:

- تبدیل DataFrame به JSON تمیز با ستون‌های اصلی
- یکپارچه‌سازی نسبی اسم اساتید، درس‌ها، وضعیت نمره‌دهی
- یه تحلیل Sentiment حداقلی روی نظرات
- حداقل یک مدل Supervised + یک مدل Unsupervised
- یک سیستم پیشنهادهندۀ ساده و منطقی
- داشبورد با صفحات حداقلی که گفتیم
- README واضح و قابل اجرا
- کار با Git و commit معنادار

## ۵.۲. چیزهایی که نمره‌ی طلایی می‌آره ✨

امتیاز اضافه برای تیم‌هایی که:

- از مدل‌های NLP پیشرفت‌هه مثل ParsBERT استفاده کنن و مستند کنن
- سیستم پیشنهاده‌نده خلاقانه‌تر/Hybrid واقعی پیاده کنن
- تحلیل‌های تحقیقی جذاب ارائه بدن مثل:
  - رابطه‌ی سختگیری نمره‌دهی و رضایت دانشجو
  - پیدا کردن «اساتید underrated» یا «اساتید پررسیک»
  - خوشبندی اساتید و تفسیر باکیفیت خوشه‌ها
- پروژه رو با Docker یا اسکریپت‌های تمیز بسته‌بندی کنن

## ۶. بخش «۳ تصمیم مهم طراحی»

تو گزارش نهایی یه بخشی باید باشه به اسم:

«۳ تصمیم مهم طراحی (Design Decisions) که در پروژه گرفتیم»

برای هر تصمیم باید بنویسید:

- چیکار کردین و کجای سیستم بوده (پارسر، مدل، داشبورد، Recommender و...)
- چه گزینه‌های جایگزینی داشتین (مثلًا fuzzy matching، Streamlit vs Dash، BERT vs TF-IDF و...)
- چرا این گزینه رو انتخاب کردین و اون یکی رو رد کردین

استدلالهاتون هم باید:

- به نتایج تجربی، محدودیت‌ها (زمان، منابع، پیچیدگی) یا تجربه‌ی خودتون تکیه کنه
- توضیح خیلی ساده مثل «این راحت‌تر بود» کافی نیست 😊

## ۷. ملاحظات اخلاقی و Bias

چون داده‌ها مربوط به آدم‌های واقعی (اساتید) هست، حتماً تو گزارش نهایی حداقل ۱-۲ پاراگراف درباره‌ی اینا بنویسید:

- چه هایی ممکنه تو داده باشه؟ Bias

- مثلاً ناراضی‌ها بیشتر نظر می‌دن؟
  - بعضی اساتید کمتر دیده می‌شن تو کانال؟
  - خطرات استفاده‌ی نادرست از این سیستم چیه؟
    - برچسب‌زن ناعادلانه
    - تقویت شایعات
  - اگه این محصول واقعاً مستقر بشه، چی‌کار می‌شه کرد که:
    - استفاده ازش منصفانه‌تر و مسئولانه‌تر باشه؟
- 

## ۸. تحویل‌ها و ارزیابی

### ۱. چی باید تحویل بدین؟

#### • پروپوزال اولیه (۲-۱ صفحه):

- ایده‌ی تیم
- ابزارهایی که قراره استفاده کنید و چرا بیش
- طرح کلی Data Pipeline و ساختار اولیه داشبورد
- نوتبوک‌های پایتون:

Supervised/Unsupervised، مدل‌های EDA، Parsing، Cleaning، NLP ◦

◦ با توضیحات متنی (Markdown) کافی

#### • کد نهایی اپ (App):

◦ سورس داشبورد + سیستم پیشنهاددهنده (.../Streamlit/Dash)

◦ ساختار تمیز (پوشه‌های /data/, src/, app و ...)

#### • :README.md

◦ قدم به قدم راه‌اندازی محیط (virtualenv/conda، نصب پکیج‌ها)

◦ چطور نوتبوک‌ها و اپ رو اجرا کنیم

◦ پیش‌نیازها (نسخه پایتون و ...)

#### • گزارش نهایی (Report):

◦ خلاصه EDA و چالش‌های داده

◦ توضیح پارسر و Cleaning

◦ مدل‌های داده‌کاوی (Supervised / Unsupervised + Recommender)

◦ بخش «۳ تصمیم مهم طراحی»

◦ بخش Bias و ملاحظات اخلاقی

## • ریپو Git (ترجیحاً GitHub):

- لینک ریپو + README + تاریخچهی منطقی commitها

## ۸.۲. نمره‌دهی تقریبی

### • ۳۵٪ - محصول نهایی:

- آیا داشبورد و Recommender واقعاً کار می‌کنن؟

- UI UX قابل قبوله؟

- برای یه دانشجوی معمولی مفیده؟

### • ۳۰٪ - داده، NLP و مدل‌ها:

- کیفیت پارسر و Cleaning

- هوشمندی در یکسان‌سازی اسامی و دسته‌بندی‌ها

- کیفیت تحلیل NLP

- طراحی و ارزیابی مدل‌های Recommender و Supervised/Unsupervised

### • ۱۵٪ - کیفیت کد و ساختار پروژه:

- تمیزی و مازولار بودن

- وجود تست‌ها

- قابلیت اجرا و تکرارپذیری

### • ۱۰٪ - Insight و EDA:

- درک اولیه از داده

- نمودارهای اکتشافی

- چند Insight جالب و معنادار

### • ۱۰٪ - بخش Design Decisions و ملاحظات اخلاقی

## ۹. چند قانون مهم آخر

### • توصیه‌ی خیلی جدی:

🚫 وارد کردن دستی داده (کپی/پیست تو اکسل و...) ممنوعه

كل فرآيند باید خودکار باشه: از خوندن JSON (يا API تلگرام) تا ساخت داشبورد و پيشنهادها.

### • توصیه‌ی جدی:

از همون اول روی Git کار کنيد.

## • روش کار پیشنهادی (Iterative) :

1. اول مطمئن شید پارسربخش خوبی از داده درست کار می‌کنه.
2. بعد یه داشبورد ساده بسازید که فقط امتیازهای عددی رو نشون بده.
3. بعدهش کم‌کم NLP، مدل‌های داده‌کاوی و آخر سر سیستم پیشنهاددهنده رو اضافه کنید.

**هدف نهایی** این نیست که پیچیده‌ترین مدل دنیا رو بسازین؛  
هدف اینه که:

- **یه سیستم کاربردی، پایدار، قابل توضیح و منصفانه** داشته باشین
  - **و بتونین همه‌ی تصمیم‌های مهم مسیر رو با زبان داده و شواهد توضیح بدین**
- 

## ۱۰. زمانبندی

برای این‌که کار پروژه رو منظم‌تر پیش ببرید، این جدول زمانبندی رو رعایت کنید:

فعالیت	تاریخ	توضیحات
مشخص کردن اعضای تیم، آماده‌سازی تیم، تحويل پروپوزال اولیه برای انجام پروژه و ساخت مخزن گیت‌هاب	جمعه ۱۴ آذر	مشخص کردن اعضای تیم، ایده‌ی اولیه شامل ابزارهای پیشنهادی و طرح کلی داشبورد و مدل‌ها، ساخت مخزن iust-supervisor و دعوت گیت‌هاب به مخزن پروژه و پر کردن <a href="#">این فرم</a>
انجام پروژه و تحويل کدها، نتایج و گیت‌هاب پروژه	جمعه ۱۹ دی	تحویل نوت‌بوک‌ها، کد اپ/داشبورد، گزارش نهایی و لینک مخزن گیت‌هاب پروژه
برگزاری ارائه‌ها (در صورت نیاز)	زمان متعاقباً اعلام می‌شود	در صورت نیاز به ارائه‌ی حضوری/آنلاین، تاریخ و جزئیات از طریق کanal اطلاع‌رسانی می‌شود