

## تصویر کلی پروژه: دقیقاً قراره چه چیزی بسازی؟

تو این پروژه، شما یک دیتاست واقعی‌نما از ۵۰ هزار متخصص فنی دارید و باید پیش‌بینی کنید چه کسی احتمالاً مهاجرت کاری می‌کند و چه کسی نه.

اما نکته مهم اینجاست:

نمود فقط برای ساختن یک مدل با عدد خوب نیست.

نمود اصلی برای اینه که نشان بدی:

1. مسئله را درست فهمیدی و درست تعریف کردي

2. داده را علمی و تمیز آماده کردي

3. مدل‌ها را درست مقایسه کردي (نه سلیقه‌ای)

4. نتیجه را قابل توضیح و قابل دفاع کردي

5. به عدالت، خطأ، ریسک و استفاده واقعی فکر کردي

یعنی در نهایت باید یک «سیستم تصمیم‌پار قابل اتکا» تحويل بدی، نه صرفاً یک فایل کد.

---

## خروجی‌هایی که باید تحويل بدھی (و چرا مهم‌اند)

### (Q20 + Capstone Q1) نوت‌بوک اصلی (1)

این نوت‌بوک قلب پروژه است. هر سوال باید یک بخش مجزا داشته باشد.

کسی که نوت‌بوک را باز می‌کند باید بتواند مرحله‌به‌مرحله بفهمد:

سؤال چی بوده •

شما چه روشی انتخاب کردید •

نتیجه چه شده •

چرا این نتیجه منطقی است •

## (2) گزارش نهایی (حداکثر ۲۵ صفحه)

این برای «دادستان علمی پروژه» است.

یعنی حتی اگر کسی کد را نبیند، از روی گزارش بفهمد کارتان چقدر درست و حرفه‌ای بوده.

## (3) پسته کد

کد باید مازولار و قابل اجرا باشد.

یعنی با یک دستور مشخص (یا چند دستور واضح) پروژه اجرا شود.

اگر کسی مجبور شود با دست فایل‌ها را تغییر دهد، کیفیت مهندسی پایین حساب می‌شود.

## (4) پاسخ تشریحی سوالات

برای سوال‌های مفهومی و ریاضی (مثل Elastic Net یا استنباط) باید توضیح تشریحی روشن بدھی، نه فقط خروجی عددی.

## (5) خلاصه مدیریتی (۱-۲ صفحه)

برای مخاطب غیر فنی:

«مدل چه می‌گوید؟ چقدر قابل اعتماد است؟ کجا باید احتیاط کنیم؟»

---

**حداکل استاندارد فنی (چیزهایی که اگر رعایت نکنی، پروژه از نظر حرفه‌ای ناقص است)**

## ثابت Seed

اگر امروز و فردا نتایج فرق کند، باز تولید پذیر نیست.  
پس همه‌جا seed مشخص داشته باش.

## تفکیک train/validation/test

این بخش خیلی مهمه.  
اگر مرز این‌ها قاطع شود، نتیجه‌ها بیش برآورده می‌شوند و در واقعیت خراب می‌شوند.

## کنترل نشت داده (Data Leakage)

خیلی از پروژه‌ها همینجا نابود می‌شوند.  
هر فیچری که بعد از زمان تصمیم‌گیری به وجود آمده باشد نباید وارد مدل شود.

## ثبت محیط اجرا

نسخه پایتون و پکیج‌ها باید مشخص باشد تا خروجی شما در سیستم دیگر هم تکرار شود.

---

# توضیح ساده و دقیق هر بلوک

---

## بلوک A – مبانی

---

### (Q1) صورت‌بندی مسئله و چرخه عمر داده‌محور

اینجا باید نشان بدھی می‌فهمی «مدل» فقط یک قطعه از کل سیستم است.

باید روشن کنی:

- دقیقاً چه تصمیمی می‌خواهیم با مدل پشتیبانی کنیم؟  
(مثلًاً شناسایی افراد با ریسک مهاجرت بالا برای برنامه نگهداشت)
- موفقیت را با چه معیارهایی می‌سنجیم؟  
(فقط Accuracy کافی نیست؛ AUC، F1، Calibration مهم‌اند)
- چه ریسکهایی داریم؟  
(نشت داده، تغییر رفتار کاربران، تغییر سیاست مهاجرت، خطای لیبل)
- بعد از استقرار چه می‌شود؟  
(چه چیزی را پایش می‌کنیم؟ هر چند وقت؟ چه زمانی بازآموزی؟)

**خروجی خوب:** یک دیاگرام چرخه عمر + جدول ریسک‌ها و راه حل.

---

## EDA) کیفیت داده و Q2

این بخش یعنی قبل از مدل‌سازی باید «چشم باز» روی داده داشته باشی.

چه کارهایی؟

- نوع ستون‌ها، مقادیر گمشده، تکراری‌ها، مقادیر غیرمنطقی
- پرت‌ها را با حداقل دو روش پیدا کن (IQR و z-score مثلاً)
- حداقل ۸ نمودار معنادار بکش (نه صرفاً برای پر کردن)
- یکتابع preprocessing قابل استفاده مجدد بساز (و تست کن)

چرا مهم است؟

چون اگر داده را نفهمی، هر مدلی هم بسازی در بهترین حالت عدد ظاهری خوب می‌دهد، نه نتیجه قابل اعتماد.

---

## بلوک B – استنباط آماری و روایت بصری

---

### (Q3) استنباط

اینجا باید نشان بدھی که فرق همبستگی، معناداری آماری و علیت را می‌فهمی.

- آیا مطالعه مشاهده‌ای است یا آزمایشی؟ •
- یک CI معتبر بده (و درست تفسیر کن) •
- یک آزمون فرض کامل انجام بده ( $H_0/H_1$ ، سطح معنی‌داری، نتیجه) •
- پیش‌فرضهای آزمون را بررسی کن •

#### نکته مهم:

p-value کوچک ≠ اثر بزرگ یا علت قطعی.  
فقط می‌گوید داده با فرض صفر سازگار نیست.

---

### (Q4) تصمیم‌گیری برای تصمیم‌گیری

این بخش درباره «انتقال درست معنا» است.

- یک داشبورد روایت‌دار بساز؛ KPI‌ها واضح باشند •
- رنگ، مقیاس، ترتیب نمودارها دلیل داشته باشد •
- یک نمودار گمراه‌کننده نشان بده و نسخه درستش را ارائه کن •

**مثال خطای رایج:**

محور `u` بریده شده و اختلاف را مصنوعی بزرگ نشان می‌دهد.

---

## بلوک SQL – C و مهندسی داده

---

### Q5) پیشرفت‌های SQL

باید نشان بدهی query تحلیلی بلد هستی، نه فقط `select` ساده.

- window function سه‌ساله با moving average
- رتبه‌بندی/دهک‌بندی با rank/ntile
- CTE با cohort analysis

اینجا کیفیت تحلیل زمانی و ساختار query خیلی مهم است.

---

### Q6) نشت داده + معماری مقیاس‌پذیر

دو چیز هم‌زمان می‌خواهند:

1. منطق علمی: حذف فیچرهای نشتشی با دید زمانی
2. منطق مهندسی: طراحی معماری داده (Bronze/Silver/Gold + Feature Store)

باید توضیح بدهی:

- داده خام کجاست؟
- داده تمیز کجاست؟

- فیچر نهایی برای مدل کجاست؟
  - چطور مطمئن می‌شوی همان فیچری که در train بوده در serving هم همان است؟
- 

## بلوک D – مدل‌سازی نظارت شده و بهینه‌سازی

---

### Elastic Net + لجستیک (Q7)

هدف اینه که از پایه علمی دور نشی.

- یک baseline معقول بساز
- تابع هزینه Elastic Net را بنویس و توضیح بده
- ضرایب را تفسیر کن (علامت، اندازه، p-value، CI)
- درباره پایداری ضرایب صحبت کن

### چرا مهم؟

چون baseline قوی و قابل تفسیر بهترین نقطه مرجع برای مدل‌های پیچیده‌تر است.

---

### (Q8) مقایسه SGD / Momentum / Adam / Ravine روی

اینجا می‌خواهند بفهمند optimization را واقعاً می‌فهمی.

- مسیر همگرایی را نشان بده

نوسان در راستای شب تند را توضیح بده •

بگو برای ناهمگنی مقیاس فیچرها کدام optimizer مناسب‌تر است و چرا •

---

### (Q9) مقایسه خانواده مدل‌ها

اینجا باید «منصفانه» مقایسه کنی:

SVM/KNN •

Tree/RF •

(XGBoost (ترجیحاً Boosting •

و الزاماً:

درست CV •

درست tuning •

• تحلیل خطأ (فقط جدول عدد کافی نیست)

نکته حرفه‌ای:

مدلی بهتر است که علاوه بر عملکرد، پایداری و تفسیر و هزینه محاسباتی مناسبی هم داشته باشد.

---

## بلوک E – بدون نظارت

---

## (Q10) کاهش بُعد

- PCA را درست انجام بده و EVR را توضیح بده
  - یک روش دیگر (t-SNE/UMAP) اضافه کن
  - در مورد معنی فضای نهفته با احتیاط حرف بزن (خصوصاً t-SNE)
- 

## (Q11) خوشه‌بندی

- KMeans + Elbow + Silhouette
- هم DBSCAN اجرا کن
- پایداری خوشه‌ها را بررسی کن

### هدف:

فقط عدد خوشه ندهی؛ باید بگویی خوشه‌ها در دنیای واقعی چه معنایی دارند.

---

## F – Deep Learning, NLP, LLM

---

## (Q12) یک مدل عصبی جدولی + یک مدل توالی/متن

می‌خواهند ببینند مدل عمیق را کورکورانه استفاده نمی‌کنی.

- MLP روی جدول

یک RNN/LSTM/GRU/CNN روی داده توالی/متن •

مقایسه واقعی با baseline کلاسیک •

### نتیجه‌گیری خوب:

اگر مدل کلاسیک بهتر بود، باید صادقانه بگویی و دلیلش را تحلیل کنی.

---

## LLM Agent (Q13) طراحی

فقط «استفاده از LLM» نیست، طراحی فرایند است:

**plan** → **retrieve** → **reason** → **verify**

باید تعریف کنی:

چطور خطای hallucination را کم می‌کنی؟ •

معیار ارزیابی کیفیت چیست؟ •

چه محدودیت امنیتی/دسترسی برای ابزارها می‌گذاری؟ •

---

## بلوک G – عدالت و حاکمیت

### عدالت و سوگیری (Q14)

اینجا خیلی مهم است چون خروجی مدل روی آدم‌ها اثر دارد.

تحلیل زیرگروهی انجام بده (کشور، تحصیلات، سابقه) •

• تبعیض غیرمستقیم (proxy discrimination) را بررسی کن

• روند human-in-the-loop تعریف کن

• مسیر اعتراض/بازبینی تصمیم داشته باش

### پیام کلیدی:

مدل دقیق ولی ناعادلانه، برای تصمیم‌گیری واقعی قابل قبول نیست.

---

## بلوک H – Capstone (یکپارچه‌سازی)

اینجا همه چیز باید به هم وصل شود:

.1 leakage-safe pipeline

.2 model card

.3 SHAP محلی + سراسری

.4 جدول عدالت + توصیه استقرار + برنامه پایش

این بخش نشان می‌دهد شما فقط «مدل‌ساز» نیستی، بلکه «طراح سیستم» هستی.

---

## بلوک I – پایش تولید و تحلیل پیشرفته (Q15–Q20)

---

### Q15) Calibration + Threshold Policy

منحنی کالیبراسیون •

ECE یا Brier •

دو آستانه: •

یکی برای بهترین F1 ○

یکی بر اساس هزینه کسب و کاری (FN vs FP) ○

**هدف:** تصمیم‌گیری بهتر، نه صرفاً خروجی احتمال خام.

---

### Q16) Drift Detection

پنجره مرجع و جاری بساز •

PSI برای عددی‌ها •

JS divergence برای دسته‌ای‌ها •

آستانه‌های هشدار/بحرانی و trigger بازآموزی تعریف کن •

---

### Q17) Counterfactual Recourse

برای موارد نزدیک مرز تصمیم، کمترین تغییر لازم را پیدا کن •

فقط تغییرهای واقع‌بینانه را پیشنهاد بده •

نرخ موفقیت recourse و هزینه تغییر را گزارش کن •

بعد اخلاقی اش را هم تحلیل کن •

---

### Q18) Temporal Validation

ارزیابی زمانی غلطان انجام بده •

افت عملکرد در زمان را بسنج •

ارتباط با drift را بررسی کن •

اگر ستون زمان نداری: fallback منطقی و مستند بده.

---

### Q19) Uncertainty Quantification

روش conformal یا مشابه •

پوشش تجربی در چند سطح اطمینان •

تحلیل پهنا و under-coverage •

هدف: بدانی مدل کجا «مطمئن نیست» تا آن موارد به انسان ارجاع شوند.

---

### Q20) Fairness Mitigation

• fairness بگیر اندازه را پایه مدل

• کن اجرا مداخله روشن یک (مثل reweighing)

• کن مقایسه را بعد قبل

• بده بده نهایی تصمیم سیاستی قید با

---

## کجاها معمولاً نمره کم می‌شود؟

1. توضیح بی توپی نشده داده نشست

2. CV نادرست split یا

3. اشتباه p-value یا CI تفسیر

4. عدالت تحلیل تحلیل نبود

5. اجرا دستور seed، پکیج نسخه بازتولید مستندات نداشتند

---

## اصول اخلاقی که باید رعایت شود

• باشد داشته ارجاع باید بیرونی منبع هر

• نکن حذف منفی یا بد نتیجه

- اگر از مدل زبانی کمک گرفتی، شفاف بگو کجا و چطور
- 

## جمع‌بندی خیلی ساده

این پروژه یعنی ساختن یک سامانه تصمیم‌گیر که:

- فنی درست باشد
- آماری درست باشد
- اخلاقی قابل دفاع باشد
- مهندسی اش قابل اجرا در عمل باشد