

دانشگاه تهران – دانشکده مهندسی برق و کامپیوتر

درس علم داده: تمرین نهایی جامع (نسخه حرفه‌ای توسعه‌یافته)

تیم درس علم داده – بهار ۱۴۰۴

ویرایش ۰.

عنوان پژوهش: تحلیل مهاجرت جهانی استعدادهای فنی با رویکرد داده محور
داده اصلی: GlobalTechTalent_50k.csv (۵۰،۰۰۰ رکورد)
مسئله اصلی: پیش‌بینی Migration_Status (۰ = مهاجرت، ۱ = عدم مهاجرت)

(۱) هدف آموزشی و فلسفه ارزیابی

این تمرین نهایی برای سنجش همزمان درک ریاضی، پیاده‌سازی مهندسی، تحلیل انتقادی، و گزارش‌نویسی حرفه‌ای طراحی شده است. انتظار می‌رود دانشجو در این پژوهش نشان دهد که می‌تواند:

• مسئله واقعی را به یک مسئله داده محور قابل ارزیابی تبدیل کند؛

• از داده خام تا استقرار پیشنهادی، یک خط لوله کامل و بازنویسیدن بسازد؛

• از مدل‌های ساده تا مدل‌های غیرخطی/پیشرفته را با استدلال علمی مقایسه کند؛

• با روش‌های XAI و تحلیل عدالت، خروجی مدل را قابل حسابرسی نماید.

(۲) اقلام تحولی الزامی

۱. نوتبوک اصلی پژوهش با ساختار Q1 تا Q17 و بخش Capstone.

۲. گزارش نهایی (حداکثر ۲۵ صفحه بدون پیوست) شامل شکل‌ها، جدول‌ها، تفسیرها و محدودیت‌ها.

۳. بسته کد شامل اسکریپت‌های قابل اجرا، فایل وابستگی‌ها، و راهنمای اجرای کامل.

۴. پاسخنامه تشریحی برای همه سوالات (فرمول، استدلال، خروجی).

۵. خلاصه مدیریتی ۱ تا ۲ صفحه برای مخاطب غیرتخصصی.

حداقل استاندارد فنی:

• ثبت random seed در تمام بخش‌های تصادفی.

• تقسیم شفاف train/validation/test.

• کنترل صریح نشت داده و نشت زمانی.

• ثبت نسخه کتابخانه‌ها و محیط اجرا.

• توان بازنویسیدن نتایج صرفاً با اجرای دستورات اعلام شده.

(۳) ساختار نمره‌دهی (۰ ۲۳۰ نمره)

| نمره | حوزه مهارتی | بلوک |
|------|--|---|
| ۲۰ | مبانی: چرخه عمر علم داده، عملیات پایتونی، EDA | A |
| ۲۰ | استباط آماری، طراحی بصری و روایت داده | B |
| ۲۵ | SQL پیشرفت و مهندسی داده در مقیاس | C |
| ۴۵ | مدل‌سازی نظارت شده و بهینه‌سازی | D |
| ۲۰ | یادگیری بدون ناظارت و کاهش بُعد | E |
| ۳۰ | یادگیری عمیق، NLP، مدل‌های زبانی/عامل‌ها | F |
| ۱۵ | عدالت، سوگیری، حاکمیت و استقرار مسئولانه | G |
| ۲۵ | کپسون یکپارچه (پیدامه‌سازی + تبیین + گزارش) | H |
| ۳۰ | توسعه حرفه‌ای پایش تولید (کالیبراسیون، درفت، ریکورس) | I |
| ۲۰+ | افزونه‌های پژوهشی/تولیدی پیشرفت | (Bonus) J |
| ۲۳۰ | | جمع بلوک‌های اصلی امتیاز تشویقی |
| ۲۰+ | | |

(۴) صورت سوال توسعه‌یافته

بلوک A – مبانی (۰ ۲۰ نمره)

(۱) چرخه عمر علم داده و صورت‌بندی مسئله (۰ ۱۰ نمره)

یک سند فنی کوتاه ارائه دهید که شامل موارد زیر باشد:

۱. تعریف مسئله، ذی‌نفعان، و تصمیم‌هایی که مدل قرار است پشتیبانی کند.
 ۲. تعریف KPI‌های فنی (مثل C_{AUC}، Recall@K، Calibration) و عملیاتی.
 ۳. فهرست ریسک‌ها (نشت داده، drift، تغییر سیاست مهاجرت، خطای برچسب).
 ۴. برنامه پایش پس از استقرار (آستانه هشدار، فرکانس بازآموزی، مسیر ارجاع انسانی).
- خروجی مورد انتظار: نمودار چرخه عمر + یک جدول ریسک.

(۲) عملیات پایتونی، کنترل کیفیت و EDA (۰ ۱۰ نمره)

۱. پروفایل‌گیری داده: نوع ستون‌ها، مقادیر گمشده، تکراری، بازه‌های نامعتبر.
۲. شناسایی پرت‌ها با حداقل دو روش (مثالاً IQR و z-score) و تحلیل پیامد.
۳. حداقل ۸ نمودار اکتشافی با تفسیر کاربردی.
۴. پیدامه‌سازی یکتابع پیش‌پردازش مازولار و قابل آزمون.

بلوک B – استباط آماری و مصورسازی (۰ ۲۰ نمره)

(۳) طراحی مطالعه و استباط (۰ ۱۰ نمره)

- ۰. تفکیک مطالعه مشاهده‌ای از ادعای علی برای این مسئله.
- ۰. ارائه حداقل یک بازه اطمینان معنی‌داری و تفسیر دقیق آن.
- ۰. تعریف یک آزمون فرض با H_0/H_1 ، سطح معنی‌داری، کنترل خطای نوع اول.
- ۰. بررسی اعتبار پیش‌فرض‌ها (نرمال‌بودن تقریبی/حجم نمونه/استقلال).

(Q۴) طراحی بصری و روایت برای تصمیمگیر (۱۰ نمره)

- طراحی یک بخش داشبوردی برای مخاطب غیرتخصصی.
- توجیه انتخاب رنگ، مقیاس، ترتیب، و annotation‌ها.
- ارائه حداقل یک مثال نموذار گمراهنده + نسخه اصلاح شده.

بلوک C – SQL و مهندسی داده (۲۵ نمره)

(SQL Q۵) SQL پیشرفته (۱۵ نمره)

سه کوئری مستقل ارائه دهید:

۱. میانگین متحرک ۳ ساله Research_Citations با window frame .
۲. رتبه‌بندی/دهکبندی کشورها یا کاربران با RANK/DENSE_RANK/NTILE .
۳. یک تحلیل cohort (ماندگاری یا انتقال وضعیت در زمان) با CTE .

(Q۶) نشت داده و معماری داده در مقیاس (۱۰ نمره)

- ویژگی‌های نشتدهنده را شناسایی و با منطق زمانی حذف کنید.
- معماری Bronze/Silver/Gold یا معادل آن را برای داده آفلاین/آنلاین پیشنهاد دهید.
- راهکار feature store با train-serving consistency توضیح دهید.

بلوک D – مدل‌سازی نظارت شده و بهینه‌سازی (۴۵ نمره)

(Q۷) مدل‌های خطی/لجستیک + Elastic Net (۱۵ نمره)

۱. مدل پایه خطی/لجستیک را اجرا کنید.
۲. تابع هزینه Elastic Net را بنویسید و گرادیان/زیرگرادیان را استخراج کنید.
۳. تفسیر ضریب p-value، بازه اطمینان، و پایداری ضرایب را ارائه دهید.

(Q۸) تحلیل بهینه‌سازی در Ravine (۱۰ نمره)

- رفتار SGD، Momentum، و Adam را روی یک تابع دره‌ای مقایسه کنید.
- مسیر همگرایی و نوسان را رسم و تحلیل کنید.
- برای ناهمگنی مقیاس ویژگی‌ها پیشنهاد عملی بدهید.

(Q۹) مقایسه خانواده مدل‌های غیرخطی (۲۰ نمره)

• KNN و SVM

• Random Forest و Decision Tree

• یک مدل Boosting (ترجیحاً XGBoost)

الزامات:

۱. پروتکل cross-validation و تنظیم ابرپارامتر.
۲. جدول متریک‌ها (حداقل: Accuracy, ROC-AUC, F1, Precision, Recall).
۳. تحلیل خطا (الگوهای اشتباه، ماتریس در هم ریختگی، حساسیت تصمیم).

بلوک E – بدون نظارت (۲۰ نمره)

(Q1۰) کاهش بُعد (۱۰ نمره)

- PCA: محاسبه و تفسیر Explained Variance Ratio
- یک روش مکمل (UMAP/t-SNE/Random Projection).
- تفسیر هندسی ابعاد نهفته.

(Q1۱) خوشبندی (۱۰ نمره)

- Silhouette + Elbow + KMeans
- DBSCAN یا روش چگالی محور معادل.
- تحلیل پایداری خوشها و معنای کاربردی آنها.

بلوک F – یادگیری عمیق، NLP و مدل‌های زبانی (۳۰ نمره)

(Q1۲) مدل عصبی جدولی و مدل توالی/متن (۱۵ نمره)

- یک MLP روی داده جدولی آموزش دهید.
- یک مدل توالی/متنی (RNN/LSTM/GRU/CNN) اجرا کنید.
- عملکرد را با baseline کلاسیک مقایسه کنید.

(Q1۳) مدل‌های زبانی و LLM Agent (۱۵ نمره)

جريان عامل محور طراحی کنید:

plan -> retrieve -> reason -> verify

و برای آن:

- معیار ارزیابی صحت/وفاداری/ایمنی تعریف کنید.
- سیاست مهار hallucination و محدودیت دسترسی ابزارها مشخص کنید.

بلوک G – اخلاق، عدالت و حاکمیت (۱۵ نمره)

(Q1۴) عدالت، سوگیری، و استقرار مسئولانه (۱۵ نمره)

- تحلیل زیرگروهی بر اساس کشور/تحصیلات/سابقه.
- بررسی proxy discrimination و سوگیری تاریخی سیاست‌ها.
- ارائه سیاست human-in-the-loop و مسیر اعتراض/بازبینی.

بلوک H – کپستون یکپارچه (۲۵ نمره)

خروجی کپستون باید شامل موارد زیر باشد:

۱. خط لوله کامل داده تا پیش‌بینی با کنترل نشت.
۲. یک گزارش model card شامل مفروضات، محدودیت‌ها، متريک‌ها.
۳. تبيين محلی و سراسری با SHAP.
۴. جدول عدالت زیرگروهی + توصيه استقرار + مانيتورينگ.

بلوک I – توسعه حرفه‌ای پایش تولید (۳۰ نمره)

(Q15) کالیبراسیون احتمال و سیاست آستانه (۱۰ نمره)

- منحنی کالیبراسیون/قابلیت اطمینان را برای بهترین مدل رسم کنید.
- حداقل یک معیار کالیبراسیون (مثل Brier یا ECE) گزارش کنید.
- دو سیاست آستانه ارائه دهید:

۱. آستانه بهینه بر اساس F1

۲. آستانه بهینه بر اساس هزینه نامتقارن خطأ (Mثلاً FN) پر هزینه‌تر از FP

(Q16) تشخیص درفت و طراحی مانیتورینگ (۱۰ نمره)

- داده را به دو پنجره مرجع/جاری تقسیم کنید (ترجمانی).
- برای ویژگی‌های عددی، PSI محاسبه و رتبه‌بندی کنید.
- حداقل یک شاخص درفت دسته‌ای (مثل JS divergence) گزارش کنید.
- سیاست هشدار/حرانی و محرک بازآموزی را تعریف کنید.

(Q17) تحلیل Counterfactual Recourse (۱۰ نمره)

- برای پیش‌بینی‌های منفی نزدیک آستانه، کمینه تغییر لازم برای تغییر تصمیم را برآورد کنید.
- حداقل دو ویژگی عملیاتی قابل مداخله را بررسی کنید.
- نرخ موفقیت ریکورس و میانه مقدار مداخله برای هر ویژگی را گزارش کنید.
- درباره امکان‌پذیری/اخلاقی بودن این مداخلات بحث کنید.

بلوک J (Bonus ۲۰+) – افزونه‌های پیچیده

- **DAG Causal و شناسایی (۵):** گراف، مجموعه تعديل، محدودیت‌های علی را صریح بیان کنید.
- **عدم قطعیت / Conformal (۵):** بازه یا امتیاز اطمینان با پوشش تجربی گزارش شده.
- اعتبارسنجی زمانی (۵): اسپلیت زمانی در برابر تصادفی، تحلیل افت عملکرد و درفت.
- سروینگ آنلاین/استریمینگ (۵): طرح ویژگی‌های تازه، SLA نگهبان درفت/OOD و مسیر rollback.

۵) قالب‌های تحویل و استاندارد مستندسازی

- همه شکل‌ها باید عنوان، واحد، و caption دقیق داشته باشد.
- تمام ادعاهای باید با خروجی کد/جدول پشتیبانی شوند.
- کد باید مازو لار، دارای نامگذاری استاندارد، و قابل اجرا با دستور واحد باشد.
- گزارش باید شامل بخش «محدودیت‌ها» و «کارهای آینده» باشد.

۶) معیار کسر نمره

- وجود نشت داده بدون گزارش: تا ۵۰٪ کسر در بخش مرتبط.
- نبود بازتولیدپذیری: تا ۳۰٪ کسر در پروژه.
- تفسیر نادرست آماری (مثلاً برداشت اشتباه از p-value): کسر ۲۰٪ در سوال.
- فقدان تحلیل عدالت یا اخلاق: کسر کامل نمره بلوک G.

۷) امتیاز تشویقی (تا +۱۰ نمره)

- توسعه تحلیل علیٰ با DAG و بحث شناسایی.
- افزودن drift و تحلیل temporal validation بین سال‌ها.
- اضافه کردن conformal prediction یا برآورد عدم قطعیت پیشرفت.

۸) اصول اخلاق علمی

- استفاده از منبع بیرونی بدون ارجاع، تخلف آموزشی محسوب می‌شود.
- نتایج ضعیف/منفی نیز باید صادقانه گزارش شوند.
- استفاده از مدل‌های زبانی باید شفاف و مستند باشد.

جمع‌بندی: این تمرین صرفاً «ساخت مدل» نیست؛ هدف، ارائه یک سامانه تصمیم‌گیر قابل دفاع است که از نظر فنی، آماری، اخلاقی و مهندسی در سطح حرفه‌ای دانشگاهی قابل ارزیابی باشد.