

گزارش نهایی کامل پروژه در قالب IEEE

تحلیل داده‌محور مهاجرت جهانی استعداد های فنی

تیم دستیاران آموزشی درس علم داده
دانشکده مهندسی برق و کامپیوتر، دانشگاه تهران
بسته کیستون کارشناسی ارشد – بهار ۱۴۰۴

چکیده

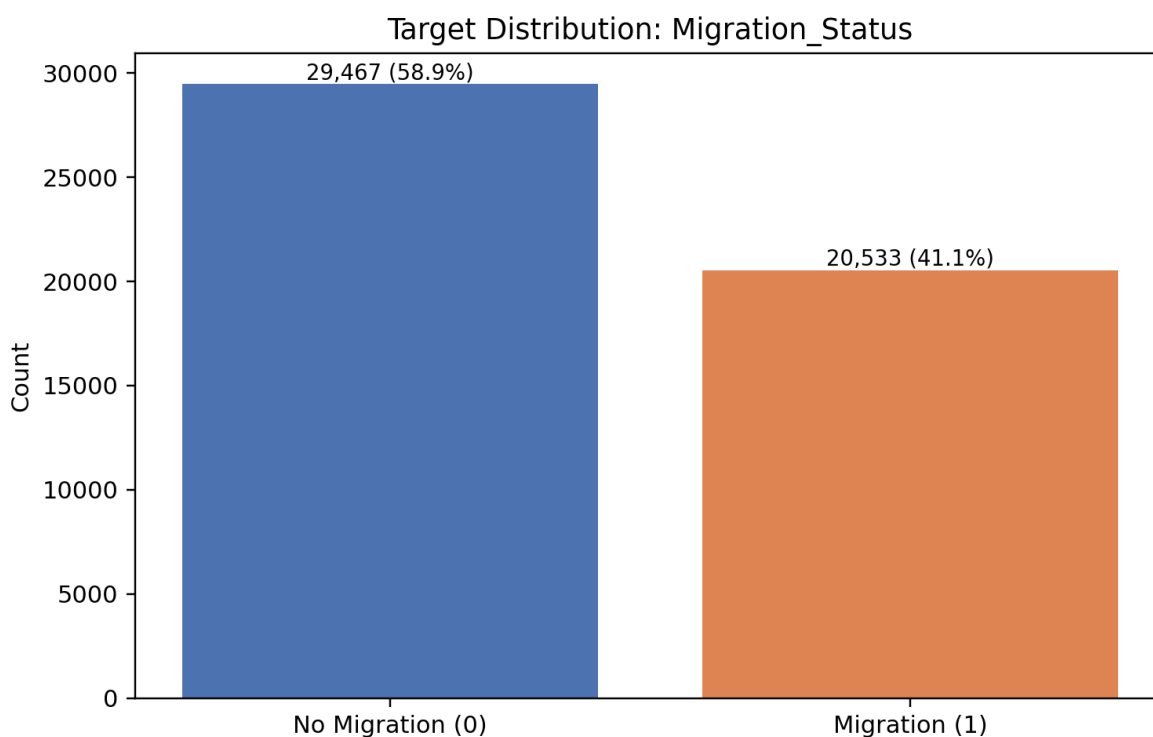
این گزارش نسخه کامل و قابل داوری پروژه نهایی درس علم داده است. پیاده‌سازی حاضر کل مسیر مهندسی داده تا تحلیل‌های تولیدی را پوشش می‌دهد: کنترل‌نشت، مدل‌سازی خطی و غیرخطی، بهینه‌سازی، یادگیری بدون نظارت، تبیین‌پذیری با SHAP، و سه افزونه پیشرفته Q18-Q20 برای پایداری زمانی، عدم قطعیت و عدالت الگوریتمی.

تعریف مسئله و دامنه پروژه

هدف، پیش‌بینی Migration_Status برای ۵۰ هزار متخصص فناوری است. خروجی نهایی صرفاً یک مدل نیست؛ بلکه یک بسته کامل دانشگاهی-مهندسی است که شامل کد، آزمون، نوت‌بوک، گزارش دوزبانه و خروجی‌های استاندارد CSV/JSON/PNG است.

تشخیص اولیه داده

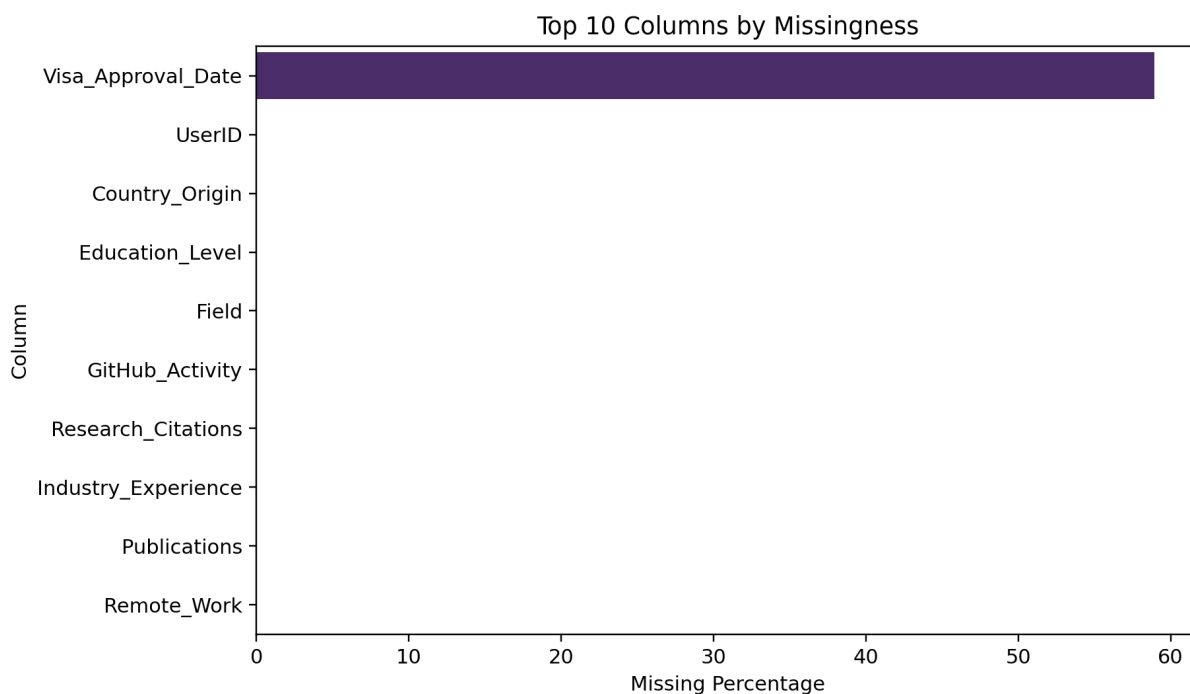
توازن کلاس هدف



شکل ۱: توزیع کلاس‌های متغیر هدف Migration_Status.

تفسیر: توزیع کلاس‌ها نامتوازن خفیف است و ارزیابی باید فراتر از دقت خام باشد.
اثر تصمیمی: معیارهای AUC/F1 و تحلیل آستانه، معیار اصلی انتخاب مدل قرار می‌گیرند.
محدودیت/تهدید: نسبت کلاس‌ها ممکن است در زمان استقرار تغییر کند و ثابت فرض کردن آن پرریسک است.

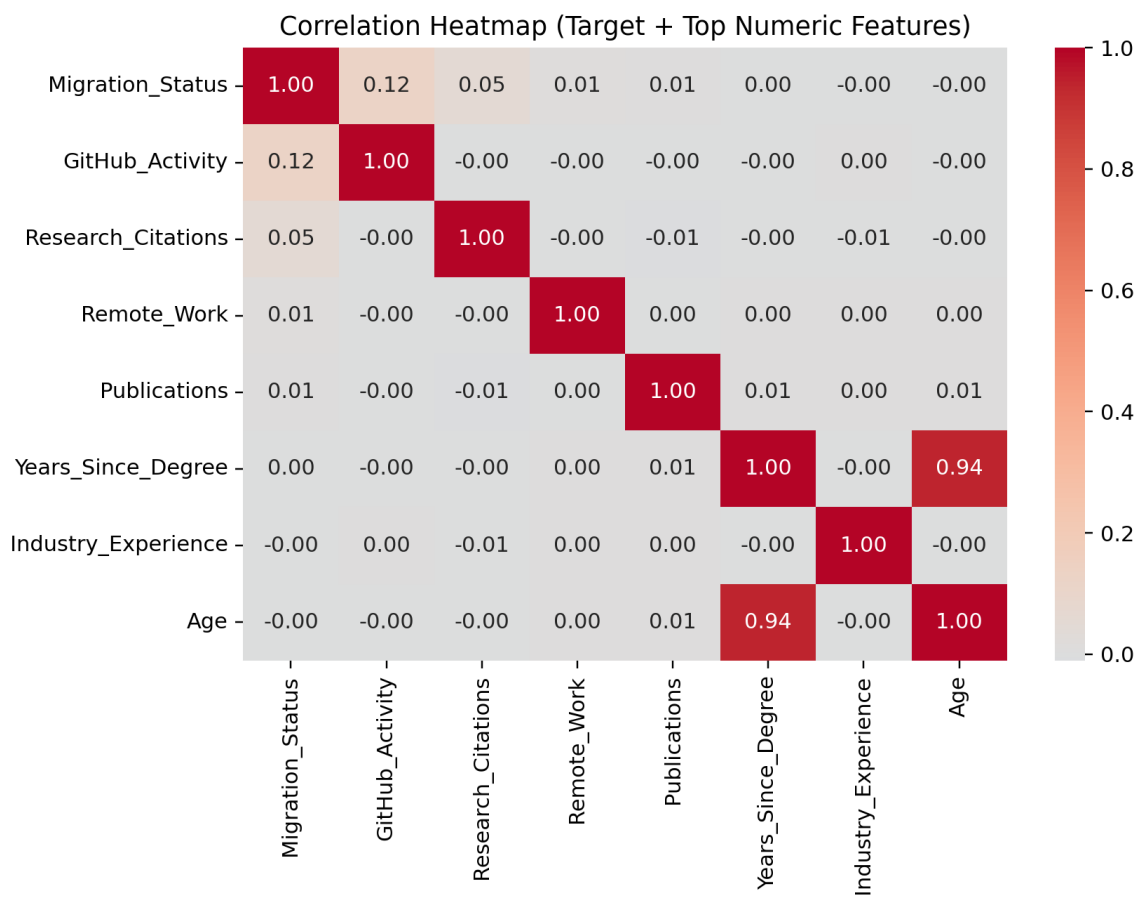
الگوی داده‌های گمشده



شکل ۲: ده ستون با بیشترین نرخ داده گمشده.

تفسیر: گمشدگی در ستون‌های فرایندی متمرکز است و بخشی از آن با نشت داده همپوشانی دارد.
اثر تصمیمی: ویژگی‌های پس‌رخداد، قبل از آموزش حذف می‌شوند یا به‌صورت کنترل‌شده وارد تحلیل می‌شوند.
محدودیت/تهدید: اگر گمشدگی تصادفی نباشد، خود الگوی گمشدگی می‌تواند منشأ بایاس شود.

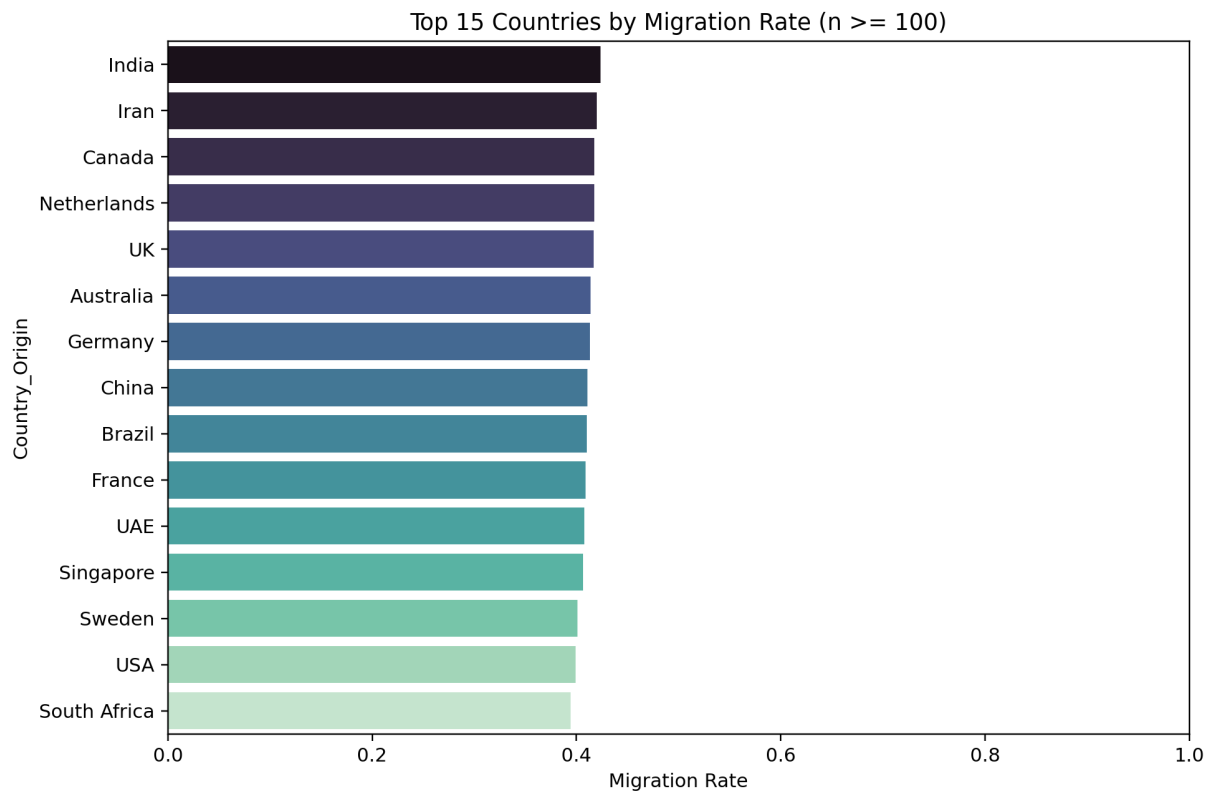
ساختار همبستگی



شکل ۳: نقشه همبستگی ویژگی‌های عددی کلیدی با هدف.

تفسیر: چند ویژگی مهم سیگنال دارند اما هیچ متغیری به‌تنهایی پاسخ کامل نمی‌دهد.
 اثر تصمیمی: استفاده از مدل‌های چندمتغیره و غیرخطی به‌صورت هم‌زمان توجیه می‌شود.
 محدودیت/تهدید: همبستگی دلالت علی ندارد و می‌تواند بازتاب عوامل پنهان سیاستی باشد.

نرخ مهاجرت در سطح کشور



شکل ۴: مقایسه نرخ مهاجرت بین کشورها پس از اعمال حداقل حجم نمونه.

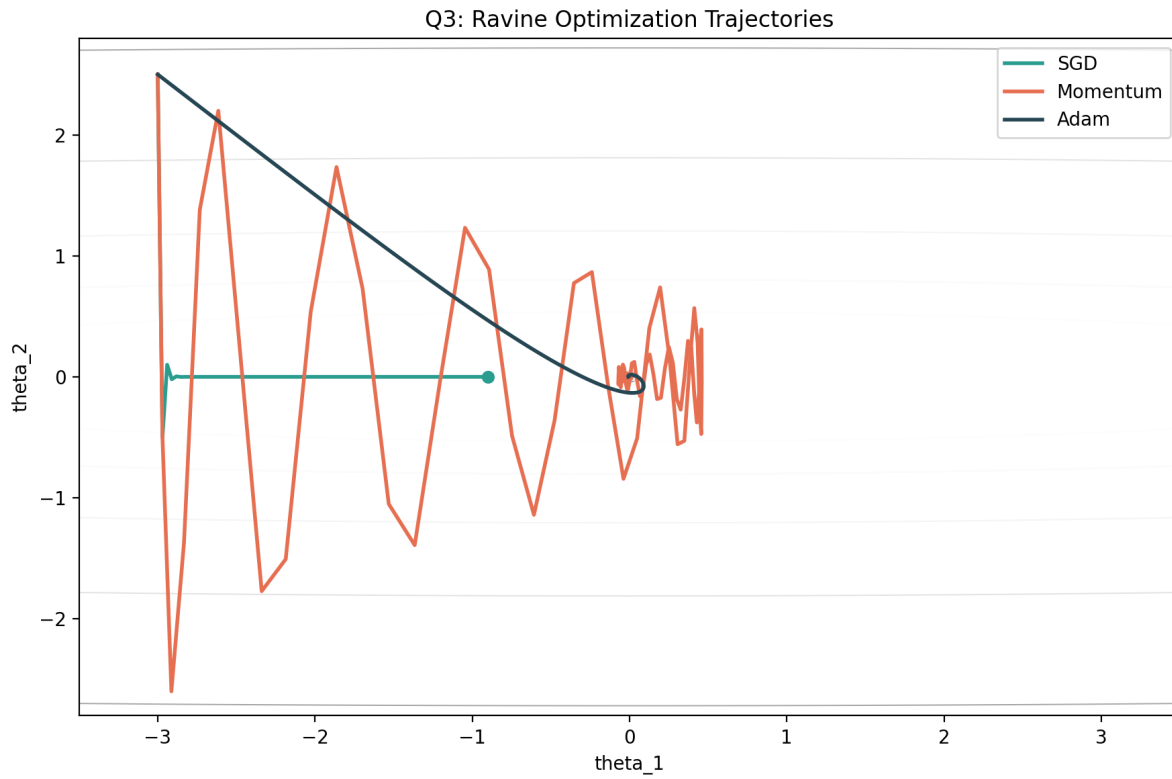
تفسیر: تفاوت گروهی بین کشورها معنادار است و بی‌توجهی به آن ارزیابی را ناقص می‌کند.
اثر تصمیمی: تحلیل عدالت قبل و بعد از مداخله الزامی است.
محدودیت/تهدید: تفاوت‌ها می‌توانند ناشی از محدودیت‌های سیاستی باشند نه توان واقعی افراد.

نتایج هسته اصلی (Q۱) تا (Q۶)

Q۱: مهندسی داده و نشئت

خروجی SQL در `code/solutions/q1_moving_average.sql` ذخیره شده است. ویژگی `Visa_Approval_Date` به‌عنوان نشئت مستقیم حذف می‌شود.

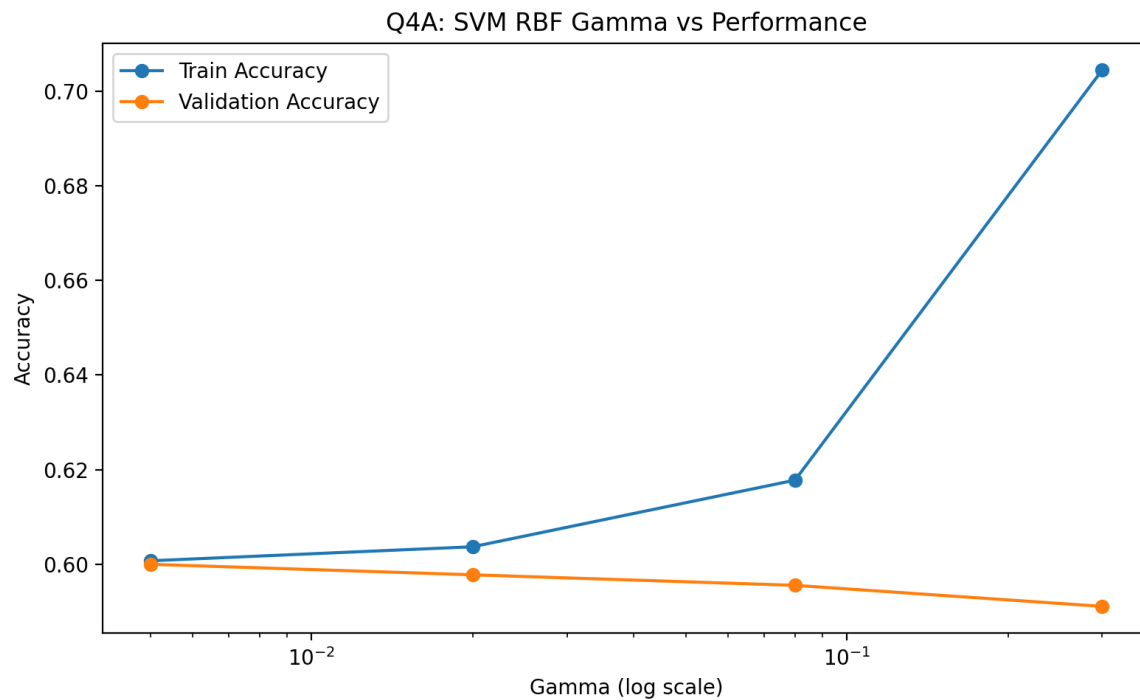
Q۳: تحلیل بهینه‌سازها



شکل ۵: مسیر همگرایی SGD، Momentum و Adam روی سطح دره‌ای.

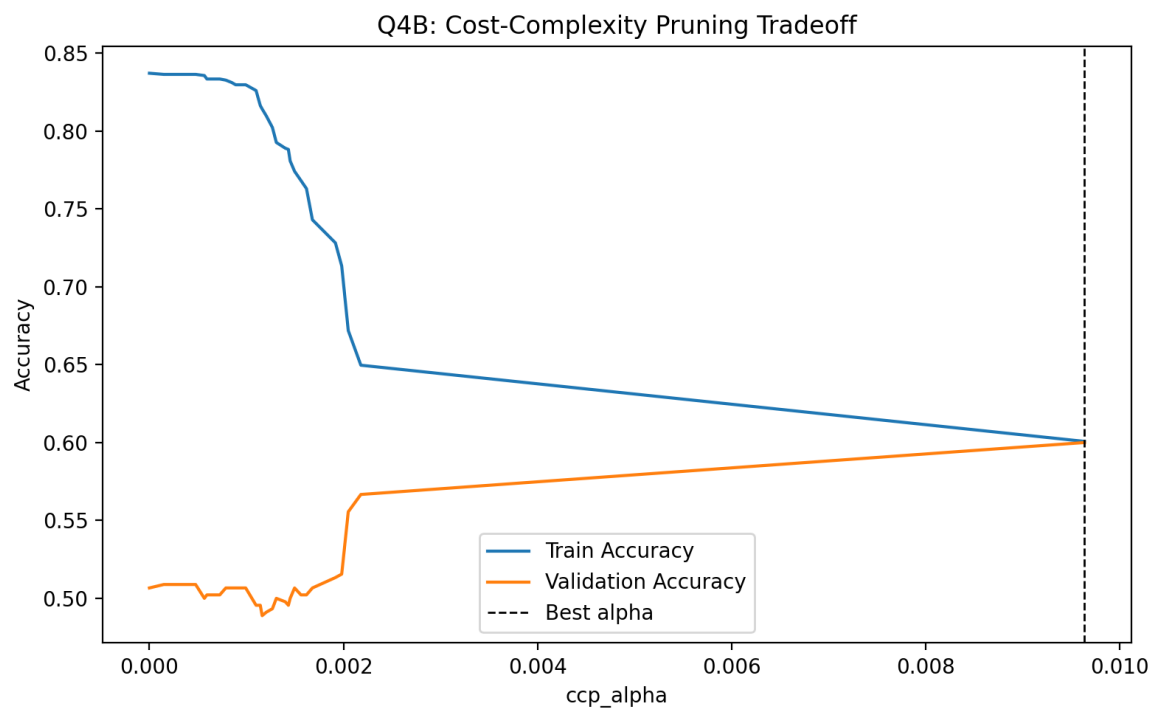
تفسیر: Momentum/Adam نوسان عرضی را کاهش می‌دهند و همگرایی سریع‌تری دارند.
اثر تصمیمی: برای توابع بدشروط، بهینه‌ساز شتابدار انتخاب پیش‌فرض مناسب‌تری است.
محدودیت/تهدید: رفتار روی تابع آزمایشی، کل پیچیدگی اهداف غیرمحدب را پوشش نمی‌دهد.

Q۴: مدل غیرخطی و کنترل پیچیدگی



شکل ۶: حساسیت دقت اعتبارسنجی نسبت به γ در SVM-RBF.

تفسیر: افزایش γ می‌تواند مرز تصمیم را بیش‌ازحد محلی و پرنوسان کند. **اثر تصمیمی:** در حالت بیش‌برازش، γ کاهش می‌یابد و روی اعتبارسنجی تنظیم می‌شود. **محدودیت/تهدید:** نتیجه نهایی به مقیاس‌بندی و همپوشانی کلاس‌ها وابسته است.

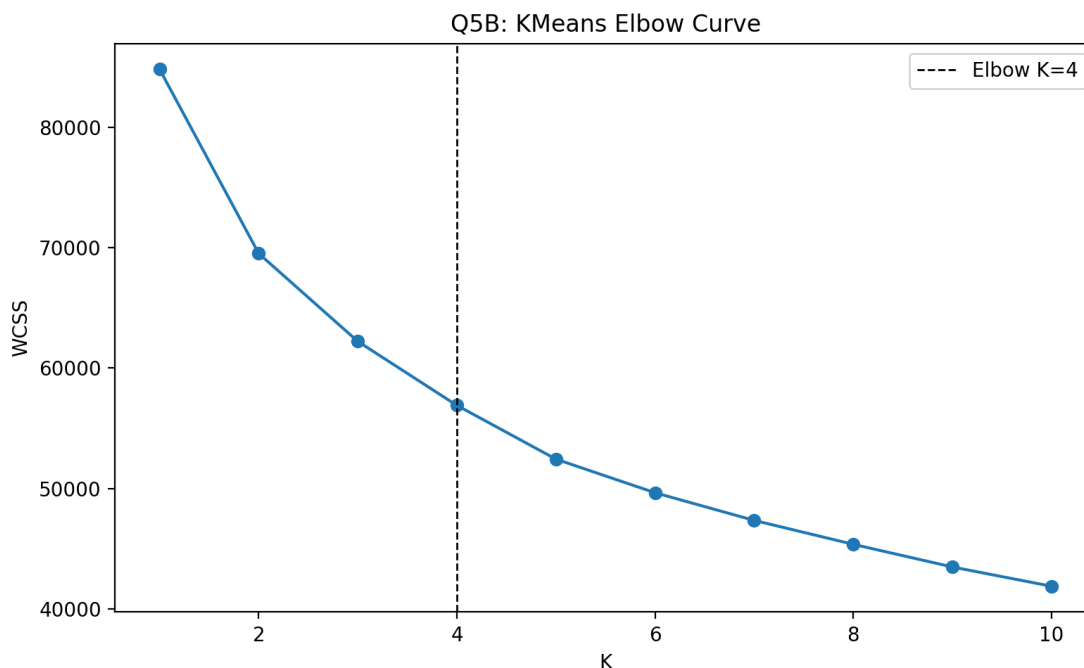


شکل ۷: مصالحه خطا-پیچیدگی در هرس Decision Tree.

تفسیر: α بزرگ‌تر اندازه درخت را کم می‌کند و واریانس را پایین می‌آورد. **اثر تصمیمی:** انتخاب α بر اساس عملکرد اعتبارسنجی انجام می‌شود، نه خطای آموزش.

محدودیت/تهدید: منحنی هرس نسبت به نحوه تقسیم داده حساس است.

Q۵: ساختار بدون نظارت



شکل ۸: منحنی آرنج WCSS برای انتخاب تعداد خوشه‌ها.

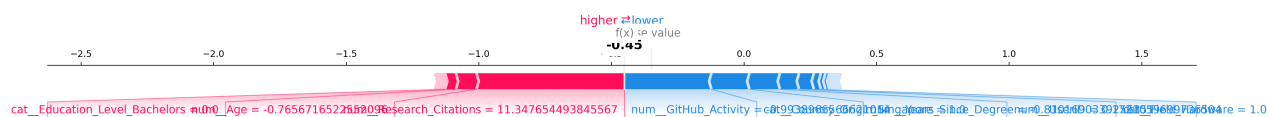
تفسیر: پس از یک K میانی، کاهش WCSS بازده نزولی پیدا می‌کند. اثر تصمیمی: K به‌صورت heuristic انتخاب و سپس با تفسیرپذیری خوشه‌ها اعتبارسنجی می‌شود. محدودیت/تهدید: در منحنی‌های کم‌انحنا، محل آرنج می‌تواند مبهم باشد.

Q۶: مدل کیستون و SHAP

پرو فایل اجرا: **balanced**.

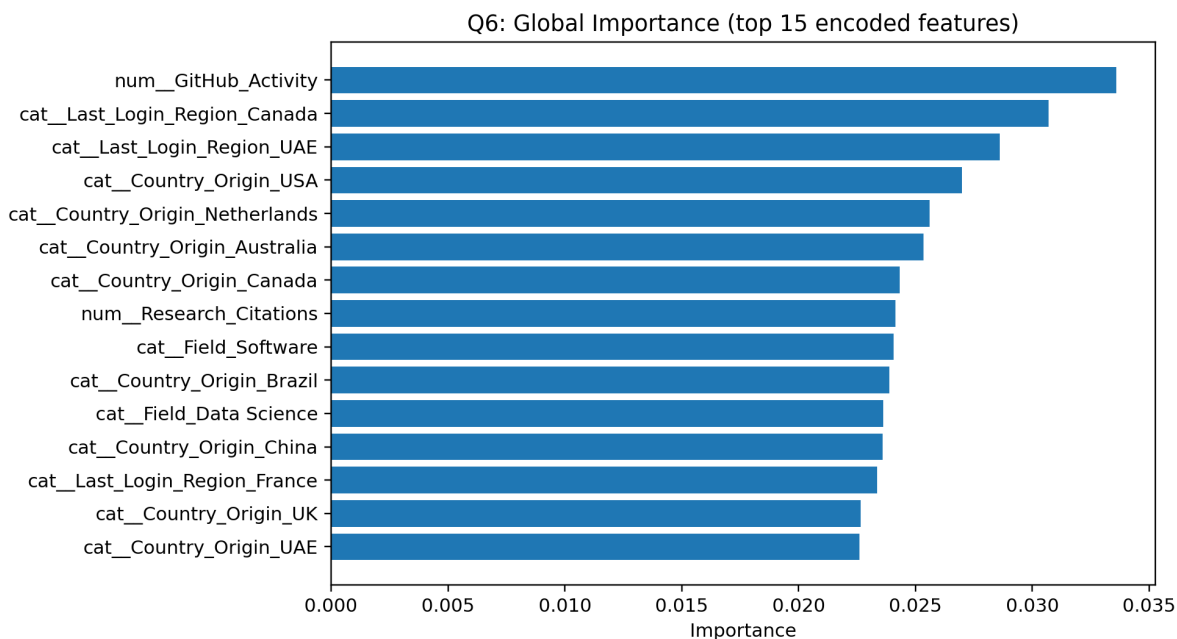
مدل کیستون: **XGBoost**.

AUC: ۰.۵۴۹۵, Accuracy: ۰.۵۸۳۵, F1: ۰.۲۴۷۵.



شکل ۹: توضیح محلی SHAP برای نمونه منتخب.

تفسیر: اختلاف base value و خروجی نهایی از جمع سهم ویژگی‌ها ساخته می‌شود. اثر تصمیمی: بازبینی پرونده روی عوامل غالب منفی/مثبت متمرکز می‌شود. محدودیت/تهدید: SHAP رفتار مدل را توضیح می‌دهد، نه رابطه علی دنیای واقعی را.



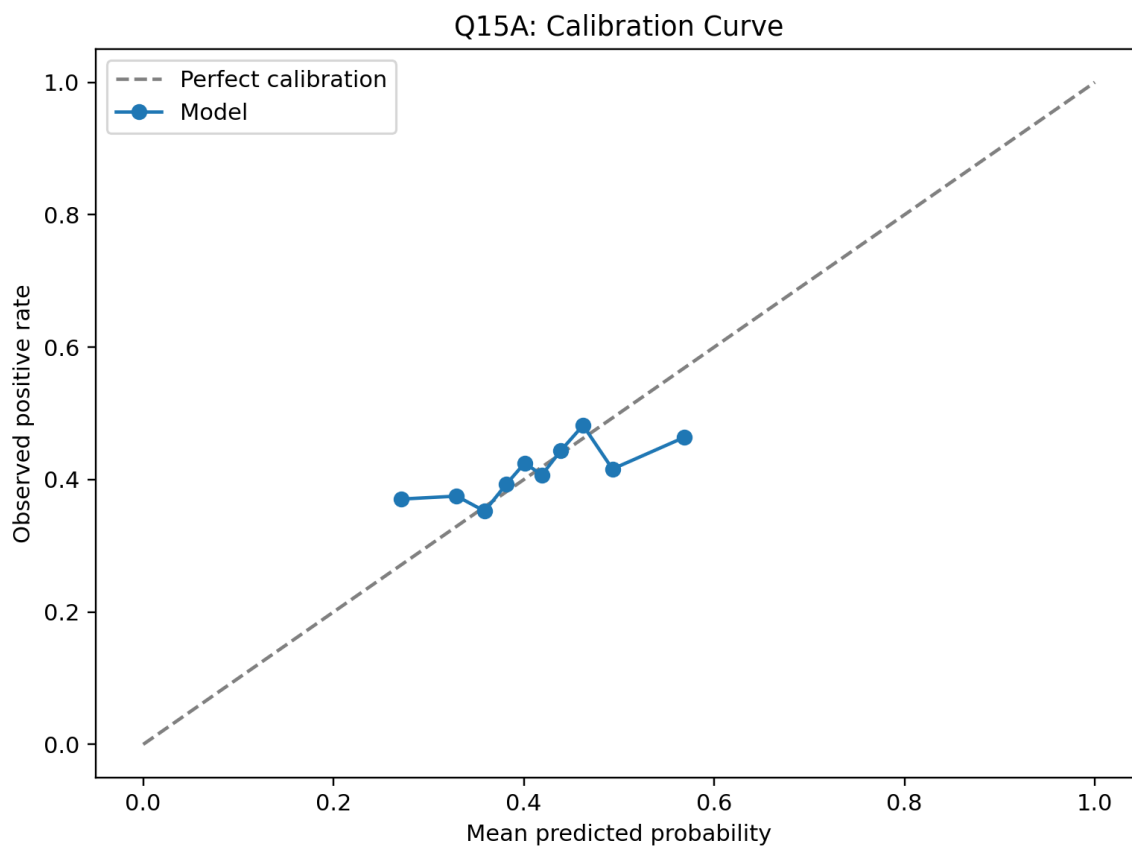
شکل ۱۰: نمای سراسری اهمیت ویژگی‌ها در مدل کپستون.

تفسیر: ویژگی‌های مرتبط با فعالیت فنی و استاد، وزن بیشتری در پیش‌بینی دارند.
اثر تصمیمی: کنترل کیفیت داده برای ویژگی‌های پراثر، اولویت حاکمیتی می‌گیرد.
محدودیت/تهدید: اهمیت سراسری، ناهمگنی اثر در همه زیرگروه‌ها را نشان نمی‌دهد.

بلوک پیشرفته تولیدی (Q۱۵) تا (Q۲۰)

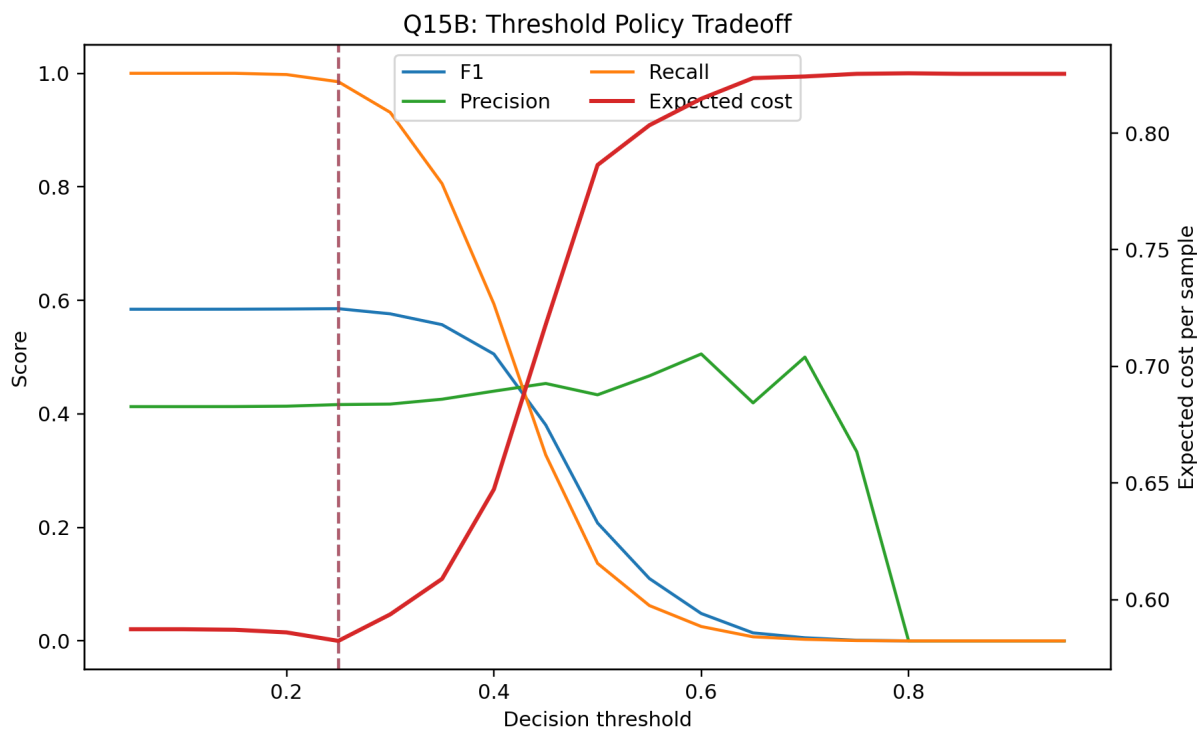
Q۱۵: کالیبراسیون و سیاست آستانه

Brier: ۰.۲۴۳۶، ECE: ۰.۳۲۷، آستانه بهینه F1: ۰.۲۵۰۰.



شکل ۱۱: منحنی قابلیت اطمینان کالیبراسیون احتمال.

تفسیر: کالیبراسیون نشان می‌دهد احتمال پیش‌بینی‌شده تا چه حد با فراوانی واقعی هم‌راستا است.
اثر تصمیمی: تصمیم‌های حساس با احتمال کالیبره‌شده گرفته می‌شوند، نه امتیاز خام.
محدودیت/تهدید: کالیبراسیون در گذر زمان ممکن است فرسوده شود.

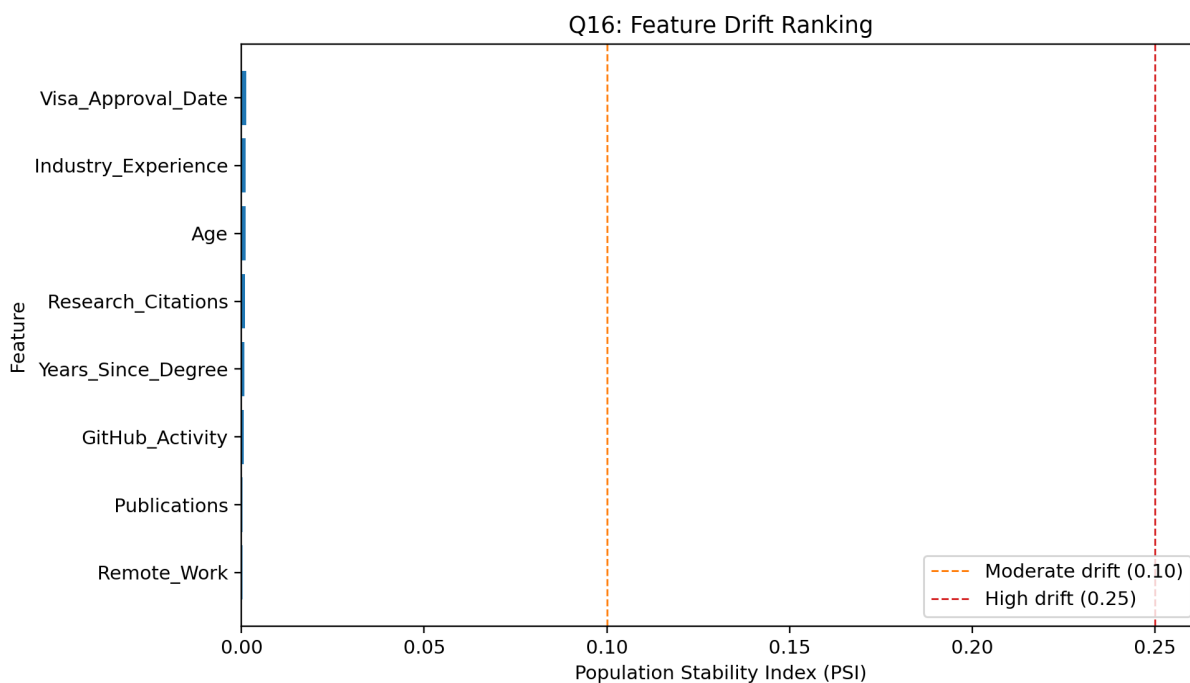


شکل ۱۲: مصالحه آستانه بین دقت، بازخوانی، F1 و هزینه مورد انتظار.

تفسیر: هر آستانه، توازن متفاوتی بین خطاهای نوع اول/دوم ایجاد می‌کند.
اثر تصمیمی: آستانه نهایی با ماتریس هزینه سیاستی انتخاب می‌شود.
محدودیت/تهدید: مفروضات هزینه ممکن است بین کشور/سازمان متفاوت باشد.

Q۱۶: پایش درفت

ویژگی با بیشترین درفت: **Visa_Approval_Date** با $PSI = 0.13$.

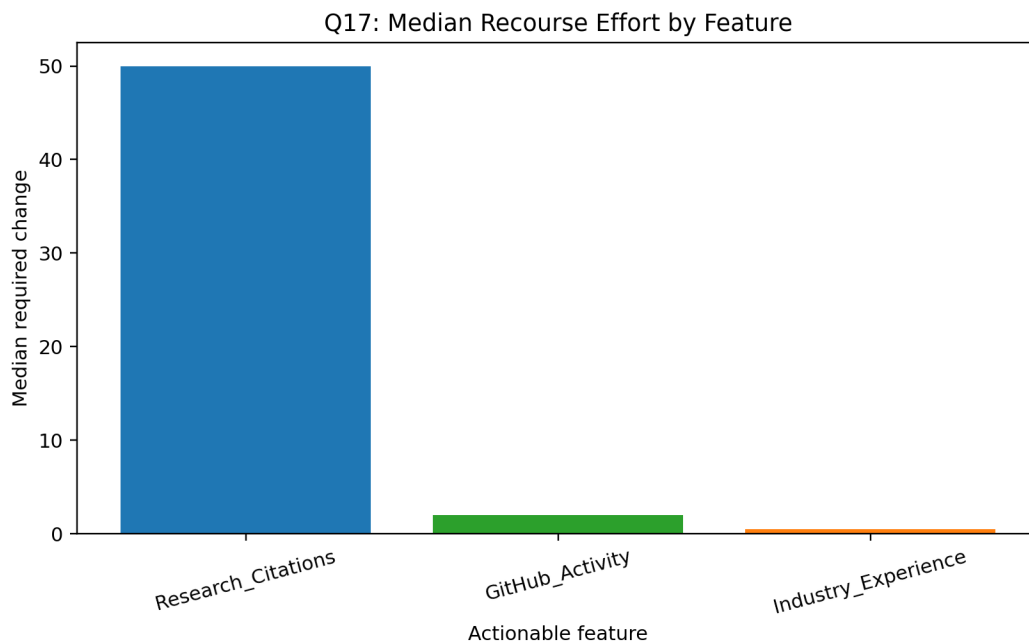


شکل ۱۳: رتبه‌بندی درفت ویژگی‌ها بر اساس PSI.

تفسیر: بخشی از ویژگی‌ها ناپایداری توزیعی قابل‌توجه دارند.
اثر تصمیمی: آستانه‌های هشدار PSI به پایش دوره‌ای و بازآموزی متصل می‌شود.
محدودیت/تهدید: PSI تغییر رابطه ویژگی-هدف را مستقیماً اندازه نمی‌گیرد.

Q۱۷: ریکورس مقابله‌ای

نرخ موفقیت ریکورس: ۰.۰۰۰۰.۱

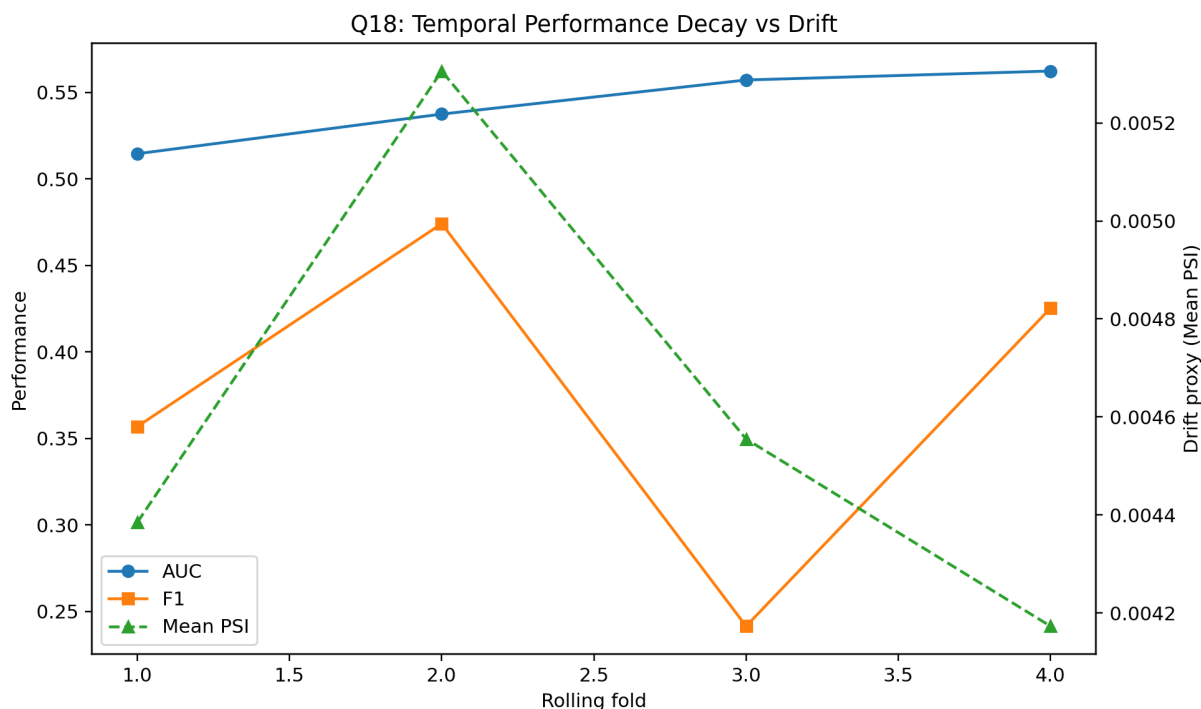


شکل ۱۴: میانه تغییر لازم برای عبور از مرز تصمیم در هر ویژگی قابل اقدام.

تفسیر: هزینه تغییر در ویژگی‌های عملیاتی یکسان نیست و قابل کمی‌سازی است.
اثر تصمیمی: پیشنهاد اقدام به متقاضی بر اساس کم‌هزینه‌ترین مسیر ممکن داده می‌شود.
محدودیت/تهدید: عملی بودن ریکورس به محدودیت‌های واقعی خارج از داده وابسته است.

Q۱۸: اعتبارسنجی زمانی و افت عملکرد

میانگین AUC زمانی: ۰.۵۴۲۸، افت AUC: ۰.۴۷۸.

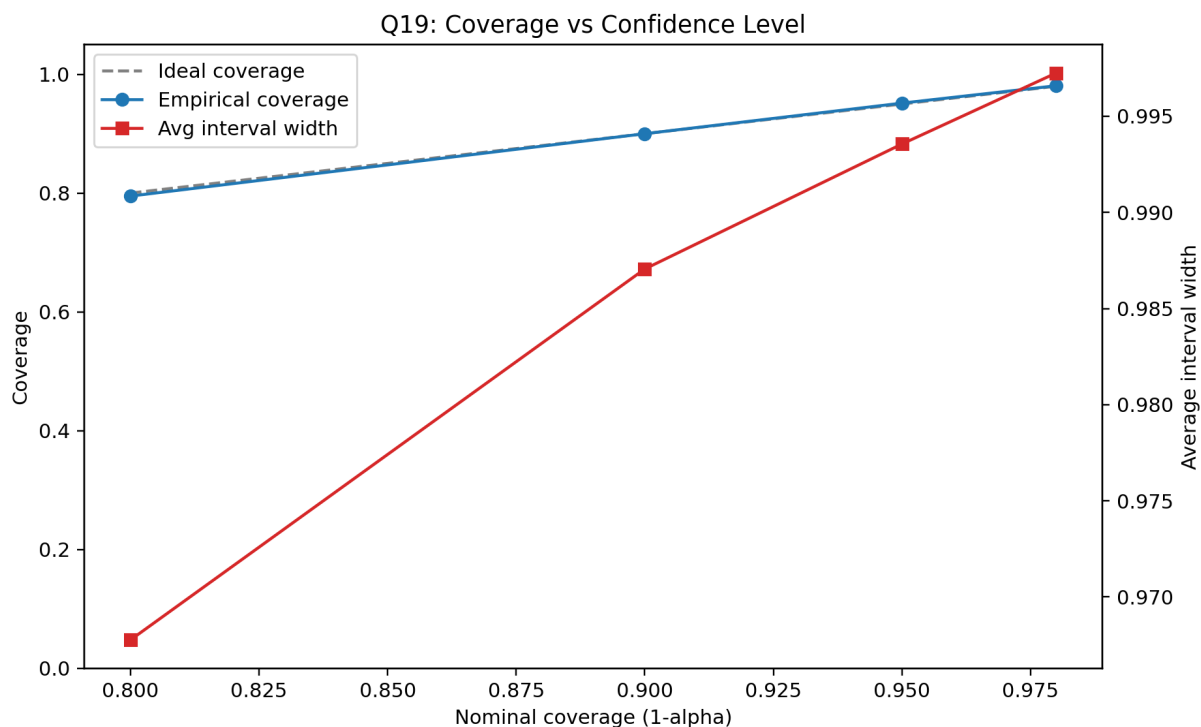


شکل ۱۵: افت عملکرد در پنجره‌های زمانی غلطان در کنار شاخص درخت.

تفسیر: عملکرد بین foldهای متوالی تغییر می‌کند و با شدت درخت قابل مقایسه است.
 اثر تصمیمی: قبل از استقرار، سنجش زمانی اجباری است و صرف random split کافی نیست.
 محدودیت/تهدید: در نبود ستون زمانی معتبر، fallback باید صریح گزارش شود.

Q۱۹: کمی‌سازی عدم قطعیت

Coverage@90: ۹۰.۰۰.۰، بیشینه کپوششی: ۰.۰۵۰.۰.

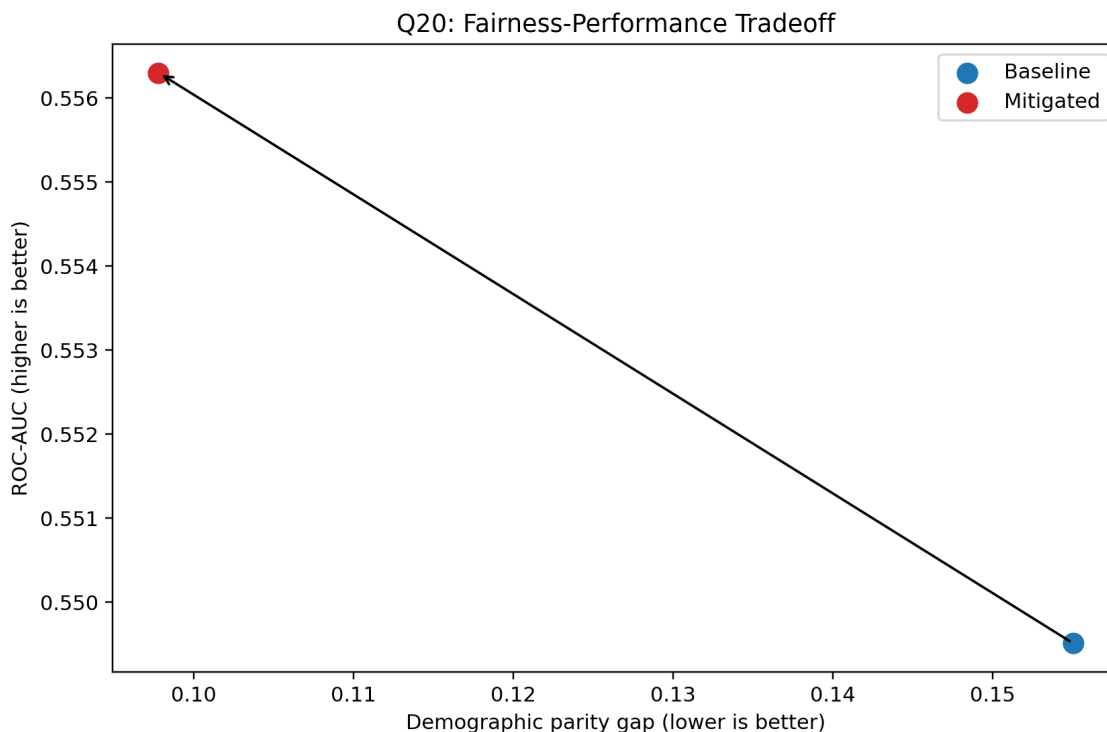


شکل ۱۶: پوشش اسمی در برابر پوشش تجربی و پهنای بازه عدم قطعیت.

تفسیر: بازه‌های conformal توازن بین اطمینان پیش‌بینی و عرض بازه را روشن می‌کنند.
اثر تصمیمی: موارد کم‌اطمینان به مسیر بررسی انسانی ارجاع می‌شوند.
محدودیت/تهدید: تضمین‌های پوشش در صورت جابه‌جایی توزیع تضعیف می‌شوند.

Q20: مداخله عدالت الگوریتمی

شکاف DP پایه: ۰.۱۵۵۰، شکاف DP پس از مداخله: ۰.۰۹۷۸، نتیجه قید سیاستی: true.



شکل ۱۷: مصالحه عدالت-عملکرد از مدل پایه تا مدل مداخله‌شده.

تفسیر: بازوزنده‌ی، نقطه عملکردی مدل را روی صفحه عدالت-کارایی جابه‌جا می‌کند.
اثر تصمیمی: استقرار منوط به عبور هم‌زمان از قید افت عملکرد و بهبود عدالت است.
محدودیت/تهدید: بهبود یک معیار عدالت ممکن است برای زیرگروه‌های ریزتر کافی نباشد.

بازتولیدپذیری و خروجی‌ها

اجرای کامل با دستور زیر انجام می‌شود:

```
python code/scripts/full_solution_pipeline.py --profile {fast,balanced,heavy}
```

خروجی‌ها در code/solutions و code/figures ذخیره می‌شوند، شامل run_summary.json نسخه ۲، خروجی‌های Q18-Q20 و فایل‌های latex_metrics برای گزارش خودکار.

جمع‌بندی

پروژه در این نسخه به یک بسته حرفه‌ای کیستون تبدیل شده است: پیچیدگی علمی، استاندارد مهندسی، و قابلیت داوری آموزشی را به‌صورت یکپارچه و قابل بازتولید پوشش می‌دهد.