

4.10 model card + evaluation report

کارت مدل و گزارش ارزیابی - ProDecks

نسخه: 1.0

تاریخ: 1403/11/17

ارزیابی شده - MVP وضعیت: مدل‌های

فصل ۱: مرور کلی مدل‌ها

هدف این سند ۱.۱.

شامل نحوه آموزش، ارزیابی، ProDecks ارائه اطلاعات جامع درباره مدل‌های یادگیری ماشین استفاده شده در محدودیت‌ها و شرایط عملکرد بهینه.

۱.۲. MVP مدل‌های موجود در

۱. (Card Recommender) مدل پیشنهاد کارت CR-MODEL-v1

۲. (Completion Time Predictor) مدل پیش‌بینی زمان تکمیل CTP-MODEL-v1

۳. blockage (Blocker Detector) مدل تشخیص BD-MODEL-v1

۱.۳. اصول اخلاقی و مسئولیت

شفافیت: افشاری کامل قابلیت‌ها و محدودیت‌ها

عدالت: بررسی و کاهش bias

مسئولیت‌پذیری: تعیین مسئول برای هر مدل

امنیت: محافظت از مدل‌ها در برابر سوءاستفاده

۲.۱. هدف مدل فصل ۲: مدل پیشنهاد کارت (CR-MODEL-v1)

۲.۲. معماری مدل

پیشنهاد خودکار کارت‌های مشابه هنگام ایجاد کارت جدید، بر اساس الگوهای تاریخی کاربر و تیم

۲.۳. نویسنده

• Collaborative Filtering با Neural Network Enhancement

• لایه‌ها

- Input Layer: ۵۰۰ ویژگی

- Embedding Layer: ۱۰۰ بعد

تفصیر اهمیت ویژگی‌ها ۳.۵.

- تجربه مسئول: ۳۵٪ اهمیت ۱.
- طول متن: ۲۰٪ اهمیت ۲.
- اهمیت ۱۵٪: subcards تعداد ۳.
- اولویت: ۱۰٪ اهمیت ۴.
- عوامل دیگر: ۲۰٪ اهمیت ۵.

محدودیت‌ها ۳.۶.

- (> 5 subcards) دقت پایین برای کارت‌های بسیار پیچیده •
- عدم در نظر گرفتن عوامل خارجی (مرخصی، تعطیلات) •
- نیاز به بهروزرسانی مداوم با داده‌های جدید •

تست bias ۳.۷.

- $\Delta < 2\%$ ناچیز bias: بر اساس جنسیت bias بررسی •
- به نفع کاربران با تجربه bias بر اساس تجربه: کمی bias بررسی •
- بر اساس زبان: عملکرد مشابه فارسی و انگلیسی bias بررسی •

blockage (BD-MODEL-v1) فصل ۴:

هدف مدل ۴.۱.

دارند قبل از وقوع blockage تشخیص خودکار کارت‌هایی که احتمال

معماری مدل ۴.۲.

- نوع Binary Classification با Neural Network
- + برای دنباله‌های زمانی LSTM: لایه‌ها
- پارامترها: ... ۸۵ پارامتر •

داده‌های آموزشی ۴.۳.

- کارت blocked (blocked): ۲,۰۰۰ مثبت •
- کارت عادی (not blocked): ۱۸,۰۰۰ منفی •
- کلاس مثبت oversampling: توازن •

ارزیابی عملکرد ۴.۴.

- دقت: ۸۷.۳٪
- Precision: ۷۶.۵٪
- Recall: ۷۲.۱٪

۵.۲. ارزیابی robustness

- تغییرات کوچک در داده: عملکرد پایدار
- retraining تغییرات بزرگ در توزیع: نیاز به
- مناسب ورودی‌های 极端: handling

۵.۳. ارزیابی bias

- بررسی بر اساس زیرگروه‌ها
- متوسط bias: تجربه کاربران -
- کم bias: اندازه تیم -
- قابل توجه bias: نوع صنعت -

- اقدامات کاهش bias:
- داده‌های متعادل -
- regularization
- post-processing

فصل ۶: تفسیرپذیری و توضیحپذیری

۶.۱. روش‌های تفسیر

۱. SHAP Values: برای مدل‌های tree-based
۲. LIME: برای مدل‌های neural
۳. Feature Importance: برای تمام مدل‌ها
۴. Attention Weights: برای مدل‌های sequence

۶.۲. نمونه تفسیر

CTP: برای مدل subcard، این کارت به دلیل طولانی بودن توضیحات (۴۰۰ کلمه) و داشتن "۳" و اختصاص به کاربری با تجربه متوسط، احتمالاً ۴۰ ساعت زمان خواهد برد

۶.۳. محدودیت‌های تفسیر

- تفسیرهای تقریبی، نه قطعی
- عدم توانایی در تفسیر تعاملات پیچیده
- نیاز به تخصص برای درک برخی تفسیرها

فصل ۷: امنیت و حریم خصوصی

۷.۱. حفاظت از مدل‌ها

- Watermarking: برای تشخیص سرقت

- در حالت استراحت
- دسترسی محدود

حريم خصوصی داده‌ها ۷.۲.

- Differential Privacy: $\epsilon = 1.0$
- برای آینده برنامه‌ریزی شده
- قبل از آموزش

حملات احتمالی و دفاع ۷.۳.

- حملات Evasion:
- روش: تغییر ورودی برای فریب مدل
- دفاع: adversarial training

- حملات Poisoning:
- روش: تزریق داده بد برای تخریب مدل
- دفاع: outlier detection

- حملات Extraction:
- روش query: استخراج مدل با
- دفاع: rate limiting, query auditing

فصل ۸: کنترل‌های ایمنی

کنترل‌های قبل از استقرار ۸.۱.

- توسط کمیته اخلاق
- تست‌های جامع
- stakeholders تأیید

کنترل‌های حین اجرا ۸.۲.

- برای تصمیمات مهم Human-in-the-loop
- توسط کاربر override امکان
- های متغیر بر اساس ریسک threshold

کنترل‌های نظارتی ۸.۳.

- لاگ تمام تصمیمات
- امکان audit
- گزارش‌های منظم

فصل ۹: برنامه بهبود

بهبودهای کوتاه‌مدت (۳ ماه)

- به CR ۸۵ افزایش دقت مدل
- به CTP ۲۰ کاهش خطای مدل
- به recall BD ۸۰ بهبود مدل

بهبودهای میان‌مدت (۶ ماه)

- افزودن مدل‌های جدید
- بهبود تفسیری‌ذیری
- bias کاهش

بهبودهای بلند‌مدت (۱۲ ماه)

- multimodal مدل‌های
- reinforcement یادگیری
- personalization پیشرفته

فصل ۱۰: نتیجه‌گیری

با استانداردهای بالای اخلاقی، امنیتی و عملکردی توسعه یافته‌اند ProDecks مدل‌های در حالی که عملکرد خوبی دارند، محدودیت‌های آنها به‌طور شفاف بیان شده و برنامه‌ای برای بهبود مستمر وجود دارد. کنترل‌های ایمنی و تفسیری‌ذیری تضمین می‌کنند که مدل‌ها به صورت مسئولانه‌ای مورد استفاده قرار گیرند.

ضمیمه‌ها

ضمیمه ۱: نتایج کامل ارزیابی

ضمیمه ۲: کد آموزش مدل‌ها

ضمیمه ۳: گزارش‌های تفسیری‌ذیری

ضمیمه ۴: نتایج تست‌های امنیتی

تهیه‌کنندگان

ProDecks تیم یادگیری ماشین

ML سرپرست) حامد کوهی

AI مشاوران اخلاق

