

استخراج و تحلیل احساسات نظرات کاربران

محمد سینا محلاتی
بوت کمپ هوش مصنوعی مولد

۱۴۰۴ آبان ۱۵

۱ استخراج داده‌ها از دیجی‌کالا:

در این مرحله یک خزنده (crawler) اختصاصی پیاده‌سازی و اجرا شد که از روی صفحات دسته «گوشی موبایل» در دیجی‌کالا، اطلاعات محصولات برنده سامسونگ و نظرات کاربران را جمع‌آوری می‌کند. برای هر محصول، حداقل ۵۰۰ نظر استخراج شده است. محصولات ناموجود و همچنین موارد توقف تولید نیز در استخراج لحاظ شده‌اند. کدهای مربوط به خزش در فایل digikala_crawl.py قرار دارند.

اسکریپت دو فایل CSV تولید می‌کند:

:Digikala_products.csv

ویژگی‌ها:
id, title_fa, title_en, selling_price, rrp_price, rating_avg, rating_count

:Digikala_comments.csv

ویژگی‌ها:
product_id, product_title, comment_id, created_at, rating, comment_text

خروچی نهایی دارای ابعاد زیر است:

Products: (426, 7)

Comments: (72474, 6)

در نهایت بیشتر از 70000 متن نظر استخراج شد.

۲ برچسب‌گذاری خودکار احساس با ParsBERT

برای برچسب‌گذاری از مدل از پیش‌آموزش دیده HooshvarLab/bert-fa-base-uncased-sentiment-deepsentipers-binary استفاده شد.

چرا از Rating استفاده نکردیم؟ ویژگی Rating به تنهایی معیار دقیقی برای احساس متن نیست مواردی مثل امتیاز اشتباه، تفاوت سلیقه کاربران، یا تضاد بین متن و امتیاز باعث خطا می‌شود. بنابراین، برای تولید برچسب‌های قابل اعتمادتر، از پیش‌بینی‌های مدل احساس فارسی استفاده شد.

۳ ساخت امبدینگ کامنت‌ها با ParsBERT :

در این مرحله بردارهای امبدینگ برای کامنت‌ها با استفاده از ParsBERT ساخته شد. خروجی و ابعاد داده برداری:

□ تعداد بردارها: 59464 نمونه.

□ بعد هر بردار: 768

□ شکل نهایی ماتریس امبدینگ: 768×59464 .

۴ آموزش مدل SVM (روی امبدینگ‌های ParsBERT) :

به جای جست‌وجوی شبکه‌ای، برای سرعت و سادگی، از تنظیمات زیر استفاده شد؛ این تنظیمات روی امبدینگ‌های متر acum 768 بعدی معمولاً عملکرد پایدار و خوبی می‌دهند:

rbf Kernel

2.0 C

"balanced" class_weight

False probability

True shrinking

)MB) 1000 cache_size

1e-3 tol

-1 max_iter

۵ Shekar → ParsBERT → SVM :

در این بخش یک ماژول یکپارچه پیاده‌سازی شد که متن خام کاربر را به صورت خودکار به برچسب احساس تبدیل می‌کند. ابتدا متن با Shekar نرمال می‌شود تا نویز زبانی کاهش یابد. سپس با توکن‌سازی انجام شده و با Mean Pooling روی خروجی لایه آخر و L2-Norm (یک بردار 768 بعدی پایدار برای هر کامنت ساخته می‌شود (با $\text{max_length}=256$ و پردازش دسته‌ای برای کارایی). در گام پایانی، همین بردارها وارد Pipeline ذخیره‌شده SVM می‌شوند اسکیل داخل پایپلاین انجام می‌گیرد و خروجی نهایی به صورت برچسب دودویی یا برچسب متنی ارائه می‌شود. این طراحی، مسیر را کاملاً تکرارپذیر، ماژولار و آماده استفاده می‌سازد.

