

به نام خدا
درس مبانی یادگیری عمیق
گزارش پروژه پایانی

استاد درس : دکتر مرضیه داوودآبادی
دستیاران : مرتضی حاجی آبادی، سحر سرکار، فائزه
صادقی، مهسا موفق بهروزی، الناز رضایی، پریسا ظفری،
حسن حماد، سید محمد موسوی، کمیل فتحی، شایان
موسوی نیا، امیررضا ویشه

دانشگاه علم و صنعت ایران، دانشکده مهندسی کامپیوتر
نیمسال اول تحصیلی ۱۴۰۲ - ۱۴۰۳



موضوع:

تحلیل احساسات در متن فارسی

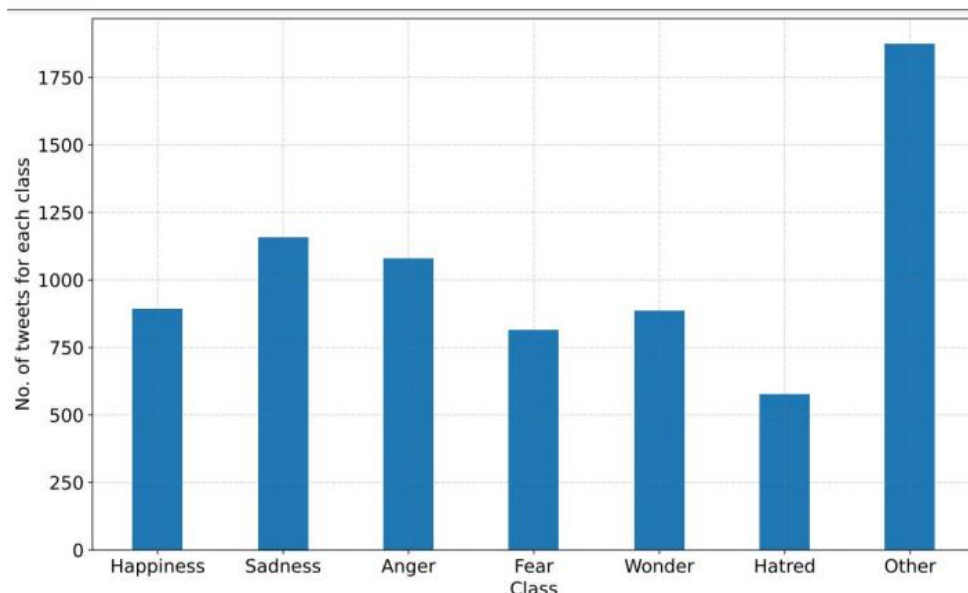
ردیف	نام و نام خانوادگی	شماره دانشجویی
۱	ملیکا محمدی فخار	99522086
۲	ستاره باباجانی	99521109

جدول ۱: مشخصات اعضای گروه

۱ شرح موضوع و مجموعه دادگان

تجزیه و تحلیل احساسات، همچنین به عنوان opinion mining شناخته می شود، یک تکنیک پردازش زبان طبیعی (NLP) است که شامل تعیین احساسات بیان شده در یک قطعه متن است. هدف این است که احساسات و نگرش هایی را که متن منتقل می کند، چه مثبت، چه منفی یا خنثی، درک کنیم. تجزیه و تحلیل احساسات به طور گسترده در برنامه های مختلف، از جمله نظارت بر رسانه های اجتماعی، تجزیه و تحلیل بازخورد مشتری، تحقیقات بازار و تعدیل خودکار محتوا استفاده می شود.

در این پروژه از مجموعه داده ArmanEmo استفاده شده است که شامل بیش از ۷۰۰۰ جمله فارسی با هفت دسته احساسات است. این مجموعه داده از منابع مختلف از جمله توییتر، اینستاگرام و نظرات دیجی کالا جمع آوری شده است. دسته بندی ها بر اساس شش احساس اصلی (خشم، ترس، شادی، نفرت، غم، تعجب) و یک دسته دیگر انجام شده است.



۲ پیش پردازش داده ها

پیش پردازش متن یک گام مهم در وظایف پردازش زبان طبیعی (NLP) برای تمیز کردن و تبدیل داده های متن خام به قالبی مناسب برای تجزیه و تحلیل و یادگیری ماشین است. روش های اصلی استفاده از پیش پردازش متن در NLP عبارتند از:

کاهش noise: حذف اطلاعات نامربوط مانند کاراکترهای خاص، علائم نگارشی، و ارقام عددی را که ممکن است به تجزیه و تحلیل کمک نکنند.

Normalization: تبدیل متن به حروف کوچک و جلوگیری از تکرار کلمات.

توکن سازی: تقسیم متن به کلمات جداگانه.

Stemming: حذف پیشوندها یا پسوندها برای به دست آوردن ریشه کلمات.

Spell Checking: اصلاح اشتباهات املائی و تایپی.

حذف تگ های HTML و کاراکترهای خاص: حذف تگ های HTML و URL ها و هشتگ ها از متن.

Vectorization: تبدیل کلمات یا عبارات به بردارهای عددی.

در این پروژه از روش های کاهش نویز، نرمالسازی، توکن سازی و حذف تگ های html برای پیش پردازش داده استفاده شده است.

۳ انتخاب مدل

مدل های موجود برای تحلیل احساسات به شرح زیر هستند:

BERT (Bidirectional Encoder Representations from Transformers): این مدل توسط گوگل

توسعه داده شده است، یک مدل مبتنی بر transformer ها است که به نتایج پیشرفته‌ای در کارهای مختلف NLP، از جمله تجزیه و تحلیل احساسات دست یافته است.

GPT (Generative Pre-trained Transformer): در حالی که مدل های GPT، مانند GPT-3، در درجه

اول برای تولید متن شناخته شده‌اند، این مدل برای تحلیل احساسات نیز قابل تنظیم است.

XLNet: یکی دیگر از مدل های مبتنی بر transformer است که نقاط قوت مدل های auto-regressive و

auto-encouoding را ترکیب می کند. برای درک context و وابستگی های داخل متن، استفاده میشود.

RoBERTa (Robustly optimized BERT approach): نسخه بهبود یافته BERT است که برخی از

محدودیت های آن را برطرف می کند و به نتایج بهتری می رسد.

DistilBERT: یک نسخه مقطر BERT است که بیشتر عملکرد خود را حفظ می کند اما تعداد پارامترها به میزان قابل توجهی کاهش یافته است. این برای استنتاج سریعتر و نیازهای کمتر منابع و در عین حال حفظ دقت طراحی شده است.

در این پروژه، با توجه به منابع محدود و نیاز به دقت بالا از مدل پیش آموخته DistilBert استفاده شده است و خروجی آن به یک pre-classifier داده شد و پس از آن از تابع فعالسازی tanh استفاده شد. سپس با استفاده از روش dropout منظم سازی صورت گرفت و در انتها، یک لایه classifier خطی قرار داده شد تا ابعاد خروجی به 7 کلاس احساسات برسد.

۴ اقدامات انجام شده

مقاله ArmanEmo بر روی تشخیص احساسات از متن فارسی تمرکز دارد. این مطالعه مجموعه دادگان «ArmanEmo» را معرفی می‌کند که یک مجموعه داده جدید شامل بیش از ۷۰۰۰ جمله فارسی برچسب‌گذاری شده توسط انسان که به هفت احساس بر اساس مدل Ekman و یک دسته «دیگر» طبقه‌بندی شده‌اند. این مقاله به جزئیات جمع‌آوری داده‌ها و فرآیند annotation می‌پردازد، که از منابعی مانند نظرات توییتر، اینستاگرام و دیجی کالا استفاده میکند و ترکیبی از روش‌های annotation دستی و خودکار را به کار می‌برد. این مطالعه همچنین چندین مدل پایه را برای طبقه‌بندی احساسات با تمرکز بر مدل‌های زبان مبتنی بر transformer بررسی می‌کند. مدل با بهترین عملکرد به میانگین امتیاز F1، 75.39%، در کل مجموعه داده آزمایشی دست یافت.

در نهایت، مقاله چالش‌ها و پیچیدگی‌های تشخیص احساسات مبتنی بر متن، از جمله دشواری طبقه‌بندی دقیق احساسات مختلط و محدودیت‌های مدل‌های فعلی در گرفتن حالات هیجانی ظریف را برجسته می‌کند. برای آشنایی بیشتر با کتابخانه Parsivar از این [لینک](#) استفاده شده است که نحوه استفاده توابع آن را با مثال آموزش داده است.

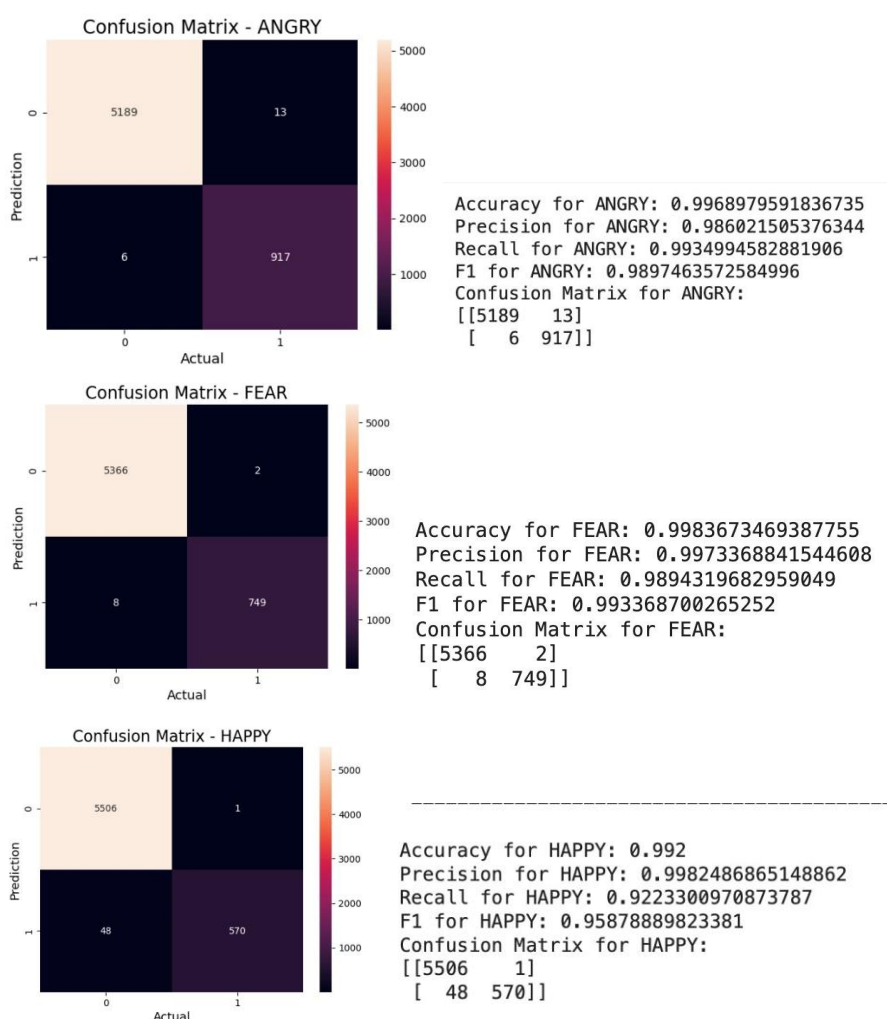
اقدامات انجام شده به شرح زیر هستند:

۱. پیش پردازش با روش‌های ذکر شده در بخش پیش پردازش داده‌ها
۲. خواندن فایل‌های آموزش و تست و اجرای پیش پردازش بر روی آنها
۳. تبدیل هر یک از برچسب‌های احساسات به یک one-hot vector
۴. توکن‌سازی داده آموزش
۵. مقداردهی اولیه ابرپارامترها
۶. استخراج id, mask, token_type_id از دیتا آموزش و تست
۷. طراحی مدل طبق توضیحات ذکر شده در بخش انتخاب مدل
۸. انتخاب تابع ضرر و تابع بهینه‌ساز
۹. آموزش مدل بر روی داده آموزش و تست آن روی داده آزمون
۱۰. تغییر فرمت خروجی اصلی و خروجی به دست آمده برای مقایسه
۱۱. محاسبه دقت و confusion matrix
۱۲. پیش‌بینی خروجی به ازای 5 داده تصادفی تست
۱۳. تعریف تابع پیش‌بینی خروجی به ازای یک متن

۵ ارزیابی مدل

از توابع recall, precision, f1-score در کنار confusion matrix که مربوط به کتابخانه sklearn.metrics هستند، استفاده شد. confusion matrix امکان نمایش بصری performance مدل با مقایسه برچسب ها و مقادیر پیش بینی شده را میدهد.

تصاویر زیر خروجی معیارها بر روی برخی کلاس ها و دقت کل مدل بر روی داده های آزمون میباشد:



Total accuracy is 0.9958483965014577

سپس از میان داده های آزمون، ۵ داده زیر را بصورت تصادفی انتخاب کرده و با استفاده از مدل، احساس موجود در آنن این چنین پیش بینی شده است:

ببینید * اینقدر بد صدا و سیما و مسئولان ما عمل کردند که خود من در ابتدا جبهه گرفته بودم و پس از پایان متوجه شدم اشتباه کردم * دیدنش خالی از لطف نیست:
 The text is :
 The True Label is : OTHER
 The Predicted label is OTHER
 The text is : شاید همین باشد: « رویای آمریکایی »
 The True Label is : HATE
 The Predicted label is HATE
 The text is : الان عادل توی تود مفتحه بعد یه کلیپ حماسی از کیروش پخش میکنه و این موفقیترو ربط میده به کیروش:
 The True Label is : OTHER
 The Predicted label is OTHER
 The text is : ز گذشته انتصایاتی را در شهرداری دیدم که طرف حتی انقبای کار را نمیدونه !!! بدون هیچ سابقه کاری و دارای قرارداد کارگری (خنده حضار)
 The True Label is : SURPRISE
 The Predicted label is SURPRISE
 The text is : فیلم ترسناک شد که:
 The True Label is : FEAR
 The Predicted label is FEAR

خروجی مدل به ازای متن ورودی زیر، این چنین میباشد:

the text is: امروز برای من خیلی روز خوبی بود
 the label is: HAPPY

۶ بخش امتیازی

۷ مراجع

۱. [ARMANEMO: A PERSIAN DATASET FOR TEXT-BASED EMOTION DETECTION](#) مقاله

۲. [data-hub](#)

۳. [Confusion matrix\(geeks for geeks\)](#)

پیاده‌سازی این پروژه با همکاری گروه آیسامیاهی‌نیا و ریحانه شاهرخیان صورت گرفته است.