



دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

پروژه پردازش زبان‌های طبیعی

## هوش مصنوعی

بهار ۱۴۰۳

استاد: محمدحسین رهبان

گردآورندگان:

مهلت ارسال: ۱ خرداد

موضوع

## ۱ شرح مختصر پروژه

در این پروژه قرار است تا مقدمات پردازش زبان‌های طبیعی (NLP) آشنا شوید و به کمک یکی از مدل‌های زبانی (LM) مبتنی بر ترنسفورمر یک مسئله‌ی واقعی دسته‌بندی متن را حل کنید. مجموعه داده‌ای در این پروژه استفاده می‌شود یک مجموعه داده‌ی فیلم‌های فارسی است که با نام «persianmovies.csv» در مستندات پروژه قابل دسترس است. این مجموعه داده شامل توضیح مختصری از فیلم‌های فارسی و ویژگی‌های آن‌ها نظیر ژانر و سال تولیدشان به هر دو زبان فارسی و انگلیسی است. در طی این پروژه با موارد زیر دست و پنجه نرم می‌کنید:

۱. کار با داده‌ها

۲. آموزش مدل دسته‌بندی متن

## ۲ کار با داده‌ها

یکی از مهم‌ترین کارها قبل از آموزش مدل و پیش‌بینی، جمع‌آوری داده و تمیز کردن داده است! قبل از هر چیزی نیاز است که داده‌هایمان را بشناسیم. در نوت‌بوک داده شده از شما خواسته شده که EDA و گام‌های مختلف تحلیل داده را برای دو زبان فارسی و انگلیسی به صورت جداگانه انجام دهید. توصیه می‌شود که علاوه بر کارهای گفته شده، تحلیل‌ها و تنظیم‌های بیشتری را نیز روی مجموعه داده‌ها انجام دهید. گام‌های خواسته شده از شما به طور کلی به صورت زیر هستند:

۱. پیش پردازش داده

۲. توصیف آماری

۳. تصویرسازی داده

۴. تحلیل ژانر

۵. تحلیل مبتنی بر زمان

۶. تحلیل رتبه بندی

۷. همبستگی و تحلیل چندمتغیره

۸. شناسایی مقادیر پرت

۹. متعادل سازی داده‌ها

دقت کنید که حتماً اگر برای Balance کردن دیتاست، تصمیم به حذف دادگان یا اضافه کردن آنها می‌گیرید (منظور از حذف دادگان داده‌های null نیست و البته دقت کنید که داده‌های null نیز نباید در دیتاست باشند)، داده‌تست را از قبل جدا کنید که بعداً در اعتبارسنجی دچار خطای overfit نشوید! سپس با توجه به تمام تحلیل‌ها و نمودارها و تقسیم داده‌ها به دو بخش train و test مدل را آموزش دهید. (random\_seed=42)

**از هرگونه دستکاری عمدی و بی‌دلیل دادگان test برای بالاتر بردن دقت جواب نهایی بپرهیزید!**

### ۳ آموزش مدل دسته‌بندی متن

در این بخش قرار است به کمک مدل زبانی BERT یک مدل دسته‌بندی برای ژانر فیلم‌ها آموزش دهید. به صورت خلاصه، ورودی مدل شما، خلاصه‌ای از فیلم است و خروجی آن یک دسته‌بندی بین ژانرهای موجود است. روند کار به این صورت است که مدل زبانی، ورودی داده شده را به یک بازنمایی (Representation) می‌برد و سپس به کمک یک لایه Fully Connected دسته‌بندی مورد نظر انجام می‌شود. روش‌های متنوعی برای انجام این کار وجود دارد که ما از Full Fine-tuning مدل به صورت end-to-end استفاده می‌کنیم. به این صورت که تمامی وزن‌های مدل زبانی BERT به همراه دسته‌بند موجود روی CLS head آموزش و به‌روز می‌شوند. در این بخش ۴ مرتبه مدل BERT را روی داده‌هایی که در قسمت قبل آماده کرده‌اید، آموزش می‌دهید و نتایج را مشاهده می‌کنید:

#### ۱. آموزش مدل ParsBERT روی داده‌های فارسی

(آ) داده‌های فارسی پیش‌پردازش شده

(ب) داده‌های فارسی پیش‌پردازش نشده

#### ۲. آموزش مدل BERT روی داده‌های انگلیسی

(آ) داده‌های انگلیسی پیش‌پردازش شده

(ب) داده‌های انگلیسی پیش‌پردازش نشده

برای مدل BERT از نسخه bert-base-uncased و برای مدل ParsBERT از نسخه HooshvareLab/bert-fa-zwnj-base استفاده کنید؛ همچنین برای آموزش مدل می‌توانید از Trainer موجود در Huggingface استفاده کنید.

**توجه: توصیه می‌شود برای آموزش مدل‌های حتماً از Google Colab استفاده کنید.**