

شیوه‌نامه‌ی پروژه تشخیص زبان فارسی از غیرفارسی

هدف از این پروژه، ارائه یک مدل کلاسه‌بندی برای تشخیص زبان فارسی از غیرفارسی در بستر شبکه‌های اجتماعی است.

فایل‌های برنامه:

language_classification.ipynb: این فایل شامل کتابخانه‌ها، کلاس و متدهای مورد نیاز است که مدل اصلی تشخیص زبان در آن پیاده‌سازی شده است.

run.py: این فایل داده‌ها را می‌خواند، آن‌ها را پیش‌پردازش می‌کند، بردارهای جاسازی کلمات را به دست می‌آورد و سپس زبان متن را با استفاده از یک مدل از پیش آموزش‌دیده KNN پیش‌بینی می‌کند.

new_final_dataset.xlsx: مجموعه داده مورد استفاده در این پروژه دارای ۲۲۵۹۳ نمونه از زبان‌های فارسی، پشتو، دری، عربی، اردو، انگلیسی، روسی و اوکراینی است. برچسب‌گذاری این مجموعه داده به این صورت است که برچسب ۱ نشان‌دهنده متن فارسی و برچسب ۰ نشان‌دهنده متن غیرفارسی است.

insta_wchr.vec: این فایل حاوی بردارهای جاسازی کلمات است که برای تبدیل متن به بردار استفاده می‌شود.

KNN_model.pkl: این فایل مدل K-Nearest Neighbors از پیش آموزش‌دیده، حاصل از اجرای فایل **language_classification.ipynb** برای کلاسه‌بندی است.

requirements.txt: این فایل دارای پکیج‌های مورد نیاز پروژه، برای نصب است.

README.md: این فایل دارای نسخه پایتون مورد استفاده و راهنمایی برای اجرای فایل **run.py** است.

معرفی کلاس ها و توابع farsivsnfarsiimplementation.py:

۱. تابع **preprocess**: این تابع برای پیش پردازش متنی که در یک دیتافریم ورودی قرار دارد، استفاده می شود.
۲. کلاس **Word2VecVectorizer**: یک کلاس تعریف شده است که برای تبدیل متن به بردارهای واژه با استفاده از فایل `insta_wchr.vec` می باشد. این کلاس سه متد اصلی دارد:
 - تابع **__init__**: در این تابع، مدل `Word2Vec` و اندازه جاسازی (`embedding size`) به عنوان ورودی گرفته می شود. این متد وقتی اجرا می شود، ابتدا یک پیام چاپ می کند که نشان می دهد داده ها در حال بارگذاری هستند، سپس مدل `Word2Vec` و اندازه جاسازی را ذخیره می کند و در نهایت یک پیام دیگر چاپ می شود که نشان می دهد بارگذاری داده ها تمام شده است.
 - تابع **fit**: این تابع نیازی به انجام عملیات خاصی ندارد.
 - تابع **transform**: این تابع برای تبدیل داده های ورودی به بردارهای عددی کلمات استفاده می شود.
 - تابع **fit_transform**: این تابع ابتدا `fit` را صدا می زند تا داده ها آماده شوند و سپس `transform` را صدا می زند تا داده ها تبدیل شوند. در نتیجه، داده های ورودی از دو مرحله `fit` و `transform` عبور می کنند و بردارهای معادل کلمات تولید می شود.
۳. تابع **language_detect**: این تابع برای تشخیص زبان متن ها در یک مجموعه داده استفاده می شود. این تابع شامل مراحل پیش پردازش داده، آموزش و ارزیابی مدل برای تشخیص زبان متن ها است. این تابع یک فایل اکسل به نام `final_dataset.xlsx` را می خواند و دو ستون به نام `textcontent`` و `target`` را از آن استخراج می کند. سپس از تابع `preprocess()` برای پیش پردازش متن ها استفاده می کند و موارد تکراری و مقادیر تهی را هم حذف می کند. داده ها به دو بخش آموزش و آزمون تقسیم می شوند. متن ها با استفاده از مدل به بردارهای عددی تبدیل می شوند. مدل پیشنهادی آموزش و ارزیابی می شود. این مدل به صورت فایل `Pickle` ذخیره می شود. همچنین برای استفاده از این تابع، باید فایل اکسل مربوطه به همراه فایل `Word2Vec` موجود باشد و این تابع یک سری وابستگی به کتابخانه های مختلف دارد که باید قبل از اجرا نصب شده باشند.