

به نام خدا



فاز اول پروژه درس پردازش زبان و گفتار

قسمت دیتا

دکتر بهروز مینایی

بهار ۱۴۰۱

تیم حل تمرین:

محمد یارمقدم

امیرحسین امینی مهر

ثمین حیدریان

مهسا انوریان

توحید عابدینی

رضا قهرمانی

تحويل مرحله اول پروژه در قالب يك ريبازيتوري گيتهاپ ميباشد. ساختار ريبازيتوري بصورت زير ميباشد. چنانچه ريبازيتوري خصوصي است به اساتيد حل تمرين دسترسي بدهيد.

- توضيح كلي پروژه در ReadMe.md
- يك فايل گزارش در ريشه ريبازيتوري تحت عنوان P1_Report.pdf
- پوشه src شامل كد هاي پروژه
- پوشه data شامل داده هاي پروژه

مجموعه داده

لازم است مجموعه داده به شكل خام و پيش پردازش شده موجود باشد.

- داده خام بايد بشكلي ذخيره شده باشد كه با اجراي مجدد دستور / كد / اسكريپت جمع آوري داده، داده خام بصورت خودكار بروز رساني شود (مثلا در پوشه‌اي به نام raw)
- داده پيش پردازش شده لازم است بصورت مرحله-مرحله در پوشه/فايل هاي جداگانه ذخيره شود (مثلا در پوشه هايي با نام هاي (word_broken، sentence_broken، cleaned
- داده ها بايد بشكلي ذخيره شده باشند كه برچسب هاي براحتي در هر مرحله قابل تفكيك باشند.
- چنانچه منابع مختلفي براي جمع آوري داده استفاده شده، ساختار/ نام فايل / پوشه به گونه‌اي باشد كه منبع جمع آوري داده در آن مشخص باشد.

كد جمع آوري و پردازش داده

كدهاي لازم براي پروژه به سه منظور نوشته شده‌اند.

- جمع آوري / استخراج داده (كرال)
 - پيش پردازش داده (شامل تميز كردن داده، شكستن جملات، شكستن كلمات)
 - استخراج آمار
- لازم است كد جمع آوري داده و پردازش آن بصورت ماجولار نوشته شده باشد. بشكلي كه از خط فرمان بتوان مراحل مختلف پيش پردازش و دريافت داده را اجرا كرد. همچنين لازم است يك اسكريپت/نرم افزاري براي اجراي تمام مراحل جمع آوري داده، پيشپردازش و استخراج آمار داده داشته باشيد بطوريكه محقق/ دانشجوي ديگر بتواند با اجراي اين اسكريپت داد هاي مشابه پوشه data بدست آورد.

گزارش

در گزارش موارد زیر را قید کنید.

- منبع دقیق داده بطوریکه بازیابی آن با روش مشابه برای یک محقق دیگر قابل انجام و راس تأزمایی باشد.
- روش جمع‌آوری، مراحل و ابزارهای استفاده شده برای جمع‌آوری داده.
- فرمت داده‌ها (فایل و ساختار پوشه). ساختار هر فایل به چه صورت است و برچسب‌های مختلف چگونه از هم متمایز هستند.
- پیش پردازش‌های انجام شده
 - روش / ابزار تفکیک جملات
 - روش / ابزار تفکیک توکن‌ها / کلمات
 - روش / معیارهای تمیز کردن داده
 - اندازه داده قبل / بعد تمیز کردن داده
- واحد برچسب‌گذاری (جمله، توییت، صفحه وب، ...) و روش برچسب‌گذاری
- آمار داده به تفکیک برچسب در قالب جدول "و" نمودار
 - تعداد "واحد" داده
 - تعداد جملات
 - تعداد کلمات
 - تعداد کلمات منحصر به فرد
 - هیستوگرام تعداد تکرار هر کلمه منحصر به فرد به ترتیب از فرکانس بالا به پایین

نکته مهم:

در انتها حتما حداقل داکيومنت و کد جمع‌آوری و مرتب‌سازی داده‌های خود را در یک فایل زیپ قرار داده و در کوئرا آپلود نمایید.

موفق باشید 😊