

به نام خدا

پروژه ی سوم هوش مصنوعی:

برای انجام این پروژه باید ابتدا فایل های ترین داده شده را بخوانیم و در صورت نیاز روی آن ها تغییراتی اعمال کنیم.

ابتدا فایل های مورد نظر را باز میکنیم. برای اینکه دشمنی درستی از کلمات داشته باشیم نیاز است که علائم نگارشی را از فایل حذف کنیم و این کار را انجام میدهیم و برای اینکه هر خط را مشخص کنیم از علامت `<s>/s>` استفاده میکنیم و هر مصرع را معادل يك خط در نظر میگیریم.

در مرحله ی بعد باید يك لیست از همه ی کلمات در نظر بگیریم که آن را با `parse` کردن در هر خط انجام میدهیم و کلمات آن را به لیست اضافه کنیم.

برای یافتن تعداد تکرار هر کلمه باید ابتدا يك دیکشنری از کل کلمات درست کنیم و مقدار `value` برای هر کلمه تعداد تکرار آن جمله را ست کنیم. (این کار برای روش یونیگرام لازم است.)

برای روش بایگرام چون باید هر کلمه با توجه به کلمه ی آن یعنی به صورت جفت کلمه سنجیده شود نیاز به دیکشنری دیگری است که کلید آن جفت کلمات کنار هم باشد و مقدار آن تعداد تکرار آن ها باشد که این کار مشابه حالت قبلی ابتدا کلید ها را یافته در لیستی قرار میدهیم سپس دیکشنری را با مقدار های تکرار هر کدام تشکیل میدهیم.

سپس برای مدل یونیگرام و بایگرام باید احتمال هر کدام را حساب کنیم و آن مقدار را برای کلید های به دست آمده در دیکشنری جدید قرار میدهیم.

با توجه به اینکه ممکن است احتمال تکرار بعضی جفت ها و کلمات صفر باشد بنابراین باید از مدل `backoffmodel` استفاده کنیم. که برای اینکار با توجه به ضرایبی که تعیین میکنیم از نتایج یونیگرام و بایگرام استفاده میکنیم.

برای نسبت دادن مقادیر به ضرایب باید در نظر بگیریم که مقدار اپسیلون عدد بسیار کمی باشد زیرا که آن عدد برای وقت هایی است که آن زوج در داده های آموزشی وجود نداشته باشد. از طرفی وجود یک لغت به صورت بایگرام نیز اهمیت زیادی دارد پس لاندا 2 مهم تر از لاندا 3 می باشد.

برای استفاده از فایل تست باید نتیجه ی مدل را به ازای هر سه شاعر به دست بیاوریم و شاعری که احتمال بالاتری داشته باشد را به عنوان جواب باز میگردانیم.