

## پروژه ی سوم درس NLP ترجمه ماشینی

### (۱) جمع آوری داده:

برای ترجمه رمان « صد سال تنهایی » را به دو زبان انگلیسی و اسپانیایی انتخاب کردم. پی دی اف هر کدام از زبان های مورد نظر را از اینترنت دانلود کردم و تبدیل آن به text به خاطر تایپی بودن متن کتاب کار سختی نبود. فایل کتاب ها را با نام های \_one hundred years of solitude\_ English.txt و Gabriel García Márquez.txt مشاهده می کنید.

### (۲) نصب برنامه ی mooses:

معمولا برای نصب mooses طبق گفته ی سایت بهتر است از سیستم عامل لینوکس استفاده شود. من هم این کار را کردم ولی در نصب boost به مشکل برخورددم. بنابراین بر روی ویندوز نصب کردم. که از دو لینک زیر برای نصب آن استفاده کردم.

<https://jon.dehdari.org/teaching/uds/moses/Oxk\el6https://www.youtube.com/watch?v=hh-V>

### (۳) کار با mooses و اجرای آن:

<http://www.statmt.org/moses/?n=Moses.Baseline> مطابق با این لینک مراحل ترجمه ی ماشینی را برای داده ی نمونه ای که در این سایت معرفی شده بود را جلو رفتیم. پس از نصب giza++ برای align کردن corpus ها و همچنین دانلود داده نمونه ، متاسفانه در مرحله ی tokenization با مشکل زیر مواجه شدم که process با اینکه مدت زیادی را در حال ران گذاشتم، به پایان نمیرسید و نمیتوانستم ادامه ی مراحل را انجام بدم. بعد از search و پرسیدن از بچه های دیگر هم نتوانستم مشکل را برطرف کنم.

