

مستندات

علی علی محمدی

۹۶۱۳۰۲۷

برای استفاده از برنامه لازم است ابتدا یک `instance` از کلاس `LanguageModel` بسازیم. به این منظور باید مقادیر زیر را به عنوان پارامترهای مدل، تعیین نماییم:

- آدرس فولدری را که دو فایل `train.json` و `valid.json` در آن قرار دارند، به عنوان ورودی اول (`corpus_dir`)
- عدد `N` به عنوان مشخص کننده مدل (`Unigram/Bigram`)
- یک متغیر از جنس `Boolean` به عنوان `smoothing`، که اگر `True` باشد، `Laplace Smoothing` اعمال می شود و اگر `False` باشد، بدون `Smoothing` خواهد بود. مقدار پیشفرض این پارامتر، `False` است.

پس از آن تابع `constructor` این کلاس فراخوانی می شود که در آن اعمال زیر به ترتیب انجام می شوند:

- ✓ هر دو فایل ورودی مسئله، در مدل بارگزاری می شود.
- ✓ کاراکترهای اضافی حذف می شوند. به این منظور، حروف الفبای فارسی، اعداد، علامت سؤال و نقطه را درون یک `dict` قرار می دهیم تا بتوان در $O(1)$ به آن ها دسترسی داشت و مشخص کرد که آیا کاراکتری به این مجموعه تعلق دارد یا خیر. به قبل و بعد نقطه ها و علامت های سؤال، یک `space` اضافه می کنیم تا در ادامه که متن را به توکن ها تجزیه می کنیم، در توکن ها لحاظ شوند و به کلمه ی قبل خود نچسبند.
- ✓ متن را به توکن ها تجزیه می شود.
- ✓ توکن های عددی با حرف `N` جایگزین می شود.
- ✓ پیشنهاد مدل های `n-gram`، حتماً شامل `Unigram` می شود. به این منظور، ابتدا `Unigram` را آموزش می دهیم. به این ترتیب، کلمات پر تکرار قابل شناسایی می شوند. کلمات را بر اساس تعداد تکرار مرتب می کنیم.
- ✓ به جای مابقی توکن های خارج از ده هزار توکن اول، `UNK` قرار می دهیم.
- ✓ کلمات پرتکرار در فایل `most_frequent.txt` ذخیره می شوند.
- ✓ مقادیر مطلوب مسئله گزارش می شود.

- ✓ جملات در فایل جداگانه به نام `sentences.txt` ذخیره می‌شوند و میانگین تعداد کلمات درون جملات محاسبه و گزارش می‌شود تا در ادامه به عنوان مقدار حد طول جمله مورد استفاده قرار گیرد.
- ✓ نمودار توزیع قانون توانی رسم می‌شود. خط رگرسیون را برای درک بهتر نمودار به آن اضافه می‌شود.

به طور مشابه، تابع آموزش **Bigram** نیز پیاده‌سازی می‌شود.

در تابع **prob**، که ورودی آن یک جمله است؛ با توجه به نوع مدل، احتمال وقوع آن جمله محاسبه می‌شود. دقت شود که **smoothing** برای **Unigram** معنای خاصی ندارد و در صورت نیاز، برای **Bigram** از **Laplace Smoothing** استفاده می‌شود.

در تابع **generate**، با توجه به نوع مدل، با استفاده از تعداد رخدادها هر توکن که از قبل محاسبه شده بود، کلمات جمله انتخاب می‌شوند و یک جمله تولید می‌شود. ورودی این تابع، کلمه‌ی اول جمله و حد طول جمله است. در حالت پیشفرض، این حد برابر میانگین طول جملات در نظر گرفته شده است که برابر ۳۶ کلمه است.

در تابع **evaluate**، با استفاده از داده‌های درون `valid.json`، اعتبار جملات تولید شده توسط مدل سنجیده می‌شود. هر جمله‌ی این فایل، با جمله‌ی متناظر تولید شده توسط مدل آموزش‌دیده که کلمه‌ی ابتدایی آن‌ها یکسان است، مقایسه می‌شود. با استفاده از متریک فاصله‌ی **Levenshtein**، نرخ خطای کلمه (**WER**) محاسبه می‌شود. دقت شود که این معیار بر اساس مقدار اولیه‌ی **N** که مشخص‌کننده‌ی نوع مدل است گزارش داده می‌شود؛ یعنی اگر مدل **Bigram** باشد، این نرخ مربوط به **Bigram** خواهد بود. در ادامه میانگین این نرخ، با توجه به تعداد کل جملات محاسبه می‌شود و به عنوان خروجی تابع بازگردانده می‌شود.

❖ دقت شود که این تابع تنها دارای یک ورودی است که مشخص می‌کند که جملات تولید شده برای مقایسه، چه طولی خواهند داشت. به طور پیشفرض، این مقدار برابر میانگین تعداد کلمات جملات (۳۶) است.

- ✓ پیشنهاد می‌شود عدد کوچک‌تری برای این پارامتر در نظر گرفته شود؛ مثلاً ۱۰.

❖ توجه کنید که تعداد جملات درون `valid.json` خیلی زیاد است و این تابع از نظر محاسباتی، هزینه‌ی زیادی دارد.

- ✓ به این منظور پیشنهاد می‌شود که در حلقه‌ی **for** این تابع، مقدار **i** محدودتر شود؛ مثلاً:

`for i in range(100):`

به این ترتیب، زمان اجرای برنامه کاهش می‌یابد.

در پایان، **Driver Code**ی نوشته شده است که اجرای توابع را ممکن می‌سازد.