



پردازش زبان طبیعی



آزمونک اول - پاسخ نامه

زمان: ۴۰ دقیقه

۱- در جدول زیر در سمت چپ عبارات منظم نوشته شده است و در سمت راست رشته هایی که ممکن است توسط این عبارات منظم پذیرفته شوند. مقابل هر عبارت منظم شماره رشته ای که می پذیرد را بنویسید. (یا وصل کنید) (جمعا ۵ نمره - هر یک ۰,۵ نمره)

شماره	رشته	عبارت منظم
1	ABODE	AB?B?C?A
2	BBC6	AB?C+E
3	ACE	B*C[2-6]
4	ABBA	A..C
5	ABDC	^A[BC]OD[EF]\$
6	DAB	ABC[0-9]
7	ABCD	(A B)(B C)[^D](C D)\$
8	CAB	[A-C][^AC][A-C]
9	ABC1	[^A-C][^AC][A-C]
10	ABCDE	A+B+C+D+E+

۲- در مورد perplexity به سوالات زیر پاسخ دهید:

a. معیار perplexity یا سرگشتی یک مدل زبانی برابر با ۷۸۲ است، این عدد یعنی چه؟ مدل زبانی که perplexity بیشتری دارد

مدل بهتری است یا مدلی که perplexity آن کمتر است؟ مقدار ایده آل برای perplexity چند است؟ (۵ نمره)

Perplexity در واقع همان **weighted branching factor** است و بیان می کند بعد از یک کلمه به طور میانگین چند کلمه مناسب می تواند بیاید. بنابراین هر چه میزان **perplexity** کمتر باشد به این معنی است که مدل توانسته است انتخاب های بعدی را محدود کند و در نتیجه مدل با **perplexity** کمتر مدل بهتری است. همچنین مقدار ایده آل برای **perplexity** برابر یک است که هیچگاه قابل دسترس نیست.

b. به نظر شما **perplexity** زبان فارسی بیشتر است یا انگلیسی؟ چرا؟ (۵ نمره)

از یک جهت انگلیسی میتواند بیشتر باشد، چون تعداد کلمات زبان انگلیسی خیلی بیشتر از فارسی است و در نتیجه تنوع ظهور کلمات

بیشتر و **perplexity** بیشتر میشود. از طرف دیگر، فارسی چون **free word order** و همچنین **highly inflectional**

است، به این دو دلیل میتواند تنوع کلمات بیشتر دیده شود. در مجموع بنظر فارسی **Perplexity** بیشتری دارد.



پردازش زبان طبیعی

آزمونک اول - پاسخ نامه

زمان: ۴۰ دقیقه

۳- یکی از مشکلات اصلی سامانه های خطایاب املایی، عدم وجود اسامی همه اشخاص در دیکشنری آن است. در این صورت در خطایابی یک نامه (یا خبر) اتفاقی که می افتد، آن است که اسامی اشخاصی که در آن نامه (یا خبر) به عنوان خطا تشخیص می دهد. چه روشی پیشنهاد می دهید؟ (۵ نمره)

استفاده از UNK handling و همچنین استفاده از Adaptive LM که در آن یکبار از روی کل سند تست نیز یک unigram آموزش می یابد و به دیکشنری اصلی با یک وزنی اضافه میشود