

# NATURAL LANGUAGE PROCESSING

## گزارش پروژه - پردازش زبان طبیعی



محمدرضا اخگری زیری

۹۶۳۱۰۰۱

زمستان ۱۳۹۸

## دسته بندی متون



در این پروژه قصد داریم تا به وسیله‌ی پردازش زبان طبیعی، متون فارسی را دسته‌بندی کنیم.

الگوریتم استفاده شده در این پروژه naive bayse میباشد و مدل‌های زبانی استفاده شده، شامل: unigram و bigram میباشد.

در ابتدای کار، اقدام به خواندن محتوای فایل HAM-Train.txt کردیم. در هر خط از این فایل جمله‌ای مربوط به دسته بندی‌ای نوشته شده بود. با پردازش این خطوط توانستیم تعداد کلمات و تعداد زوج کلمات را در هر دسته حساب کنیم (فایل‌های آنها به همراه پروژه قرار داده شده است).

از نکات قابل توجه میتوان به این اشاره کرد که برای بالا بردن دقت علامت <S> را به ابتدای هر جمله اضافه کردم تا در زوج کلمات، کلمات اول جمله نیز حساب شود.

سپس اقدام به خواندن محتوای فایل HAM-Test.txt میکنیم. این فایل، برای آزمایش الگوریتم ماست. برای آزمایش هر خط را از فایل خوانده و در ابتدای آن نماد شروع <S> را میگذاریم. حال باید برای دسته بندی‌های مختلف‌ای که از قسمت قبل یافتیم، احتمال وجود تک کلمات و زوج کلمات را حساب کنیم.

برای هموار سازی از ترکیب خطی این احتمالات استفاده شده و با فرمول زیر حساب میشود:

$$p(i|i-1) = \lambda p_i + (1 - \lambda)p(i|i-1)$$

مقدار  $\lambda$  سه مقدار متفاوت ۰,۲, ۰,۵ و ۰,۹ داده شد. و نتیجه آنها به همراه پروژه قرار داده شده است ولی بهترین حالت ۰,۵ شد.

برای پیدا کردن مقادیر احتمال از الگوریتم naive bayse استفاده شده است، که در این الگوریتم به دلیل کوچک بودن الگوریتم‌ها از ضرب آنها خودداری کردیم و مجموع log های آنها را حساب میکنیم. کلماتی که در training set قرار ندارند را با احتمال -۱۰۰۰ در نظر میگیریم.

برای ارزیابی نهایی پروژه از سه متغیر گفته شده در سوال استفاده میکنیم که به شرح زیر است:

$$Precision = \frac{TP}{TP + FP}, recall = \frac{TP}{TP + FN}, F = 2 * r * \frac{p}{r + p}$$

در صورت نزدیک بودن به یک یعنی پیاده سازی بهتری داشتیم.

