



## بازیابی پیشرفته اطلاعات

نیم سال اول ۱۴۰۲-۰۳

مدرس: دکتر حمید بیگی

۴ آذر

زمان: ۲۵ دقیقه

کوییز سوم (۱۰۰ نمره)

## سوال ۱: تخمین احتمال تولید کوثری (۴۰)

در نظر بگیرید که سه داکيومنت زیر ( $d_1, d_2, d_3$ ) را در Collection داریم. به کمک روش هموارسازی Jelinek-Mercer smoothing احتمال تولید کوثری را برای هر داکيومنت برای کوثری داده شده محاسبه کنید و در انتها سندها را بر اساس این احتمالات رتبه بندی کنید.  $\lambda$  را برابر  $1/4$  در نظر بگیرید.

$d_1$ : Shakespeare was a renowned playwright and wrote poets  
 $d_2$ : the famous playwright William wrote the tragedy of Hamlet  
 $d_3$ : Works of Shakespeare are celebrated worldwide  
 query: Shakespeare wrote Hamlet

تعداد کل کلمات در کلکسیون:  $N = 23$ 

تعداد تکرار Shakespeare: ۲

تعداد تکرار wrote: ۲

تعداد تکرار Hamlet: ۱

 $d_1$ :

تعداد کل کلمات: ۸

$$P(\text{Shakespeare} | M_{d_1}) = \frac{1}{8}$$

$$P(\text{wrote} | M_{d_1}) = \frac{1}{8}$$

$$P(\text{Hamlet} | M_{d_1}) = 0$$

$$P(\text{Shakespeare} | M_c) = \frac{2}{23}$$

$$P(\text{wrote} | M_c) = \frac{2}{23}$$

$$P(\text{Hamlet} | M_c) = \frac{1}{23}$$

$$P(\text{query} | d_1) = (\frac{1}{4} \frac{1}{8} + \frac{3}{4} \frac{2}{23})(\frac{1}{4} \frac{1}{8} + \frac{3}{4} \frac{2}{23})(\frac{1}{4} 0 + \frac{3}{4} \frac{1}{23})$$

 $d_2$ :

تعداد کل کلمات: ۹

$$P(\text{Shakespeare} | M_{d_2}) = 0$$

$$P(\text{wrote} | M_{d_2}) = \frac{1}{9}$$

$$P(\text{Hamlet} | M_{d_2}) = \frac{1}{9}$$

$$P(\text{query} | d_2) = (\frac{1}{4} 0 + \frac{3}{4} \frac{2}{23})(\frac{1}{4} \frac{1}{9} + \frac{3}{4} \frac{2}{23})(\frac{1}{4} \frac{1}{9} + \frac{3}{4} \frac{1}{23})$$

 $d_3$ : تعداد کل کلمات: ۶

$$P(\text{Shakespeare} | M_{d_3}) = \frac{1}{6}$$

$$P(\text{wrote} | M_{d_3}) = 0$$

$$P(\text{Hamlet} | M_{d_3}) = 0$$

$$P(\text{query} | d_3) = (\frac{1}{4} \frac{1}{6} + \frac{3}{4} \frac{2}{23})(\frac{1}{4} 0 + \frac{3}{4} \frac{2}{23})(\frac{1}{4} 0 + \frac{3}{4} \frac{1}{23})$$

بنابراین با توجه به مقادیر احتمالات محاسبه شده، رتبه بندی سندها به شرح زیر خواهد بود:

$$d_2 > d_1 > d_3$$

محاسبات بیشتر که برای جواب نیاز نیست و همان ۳ عبارت فوق برای جواب درست کافی است.

$$P(\text{query} | d_1) \approx 0.00030$$

$$P(\text{query} | d_2) \approx 0.00036$$

$$P(\text{query} | d_3) \approx 0.00026$$

## سوال ۲: naive bayes (۴۰)

با توجه به جدول زیر به سوالات الف-د پاسخ دهید.

الف) یک طبقه بندی چندجمله ای با استفاده از روش naive bayes را روی داده های جدول تخمین بزنید.

ب) طبقه‌بندی به دست آمده را روی داده تست اعمال کنید.  
 ج) یک طبقه‌بندی برنولی با استفاده از روش naive bayes را روی داده‌های جدول تخمین بزنید.  
 د) طبقه‌بندی به دست آمده را روی داده تست اعمال کنید.

نکته: در جدول زیر داک ۱، داک ۲، داک ۳ و داک ۴ مربوط به داده آموزش است و داک ۵ مربوط به داده تست است.

آیدی داکيومنت	کلمات موجود در داکيومنت	آیا مهم است؟
داک ۱	باران تند	بله
داک ۲	برف تند شد	بله
داک ۳	حيوان آمد	خير
داک ۴	آمد پلنگ تند	خير
داک ۵	تند تند آمد	؟

### سوال ۳: kNN vs SVM (۱۰)

روش k-Nearest Neighbors (kNN) با روش Support Vector Machines (SVMs) از نظر مرزهای تصمیم‌گیری (decision boundaries) مقایسه کنید و تفاوت‌ها را ذکر نمایید.

kNN و SVM دو روش کاملاً متفاوت برای تعیین مرزهای تصمیم در طبقه‌بندی داده‌ها هستند. kNN یک روش غیرپارامتری است که به طور صریح مرزهای تصمیم را یاد نمی‌گیرد. بلکه، طبقه‌بندی‌ها را بر اساس مقیاس شباهت بین موارد آزمایشی جدید و داده‌های آموزش دیده انجام می‌دهد. مرزهای تصمیم به طور ضمنی بر اساس همسایگی‌های محلی اطراف هر نقطه داده پدید می‌آیند. بنابراین، مرزهای تصمیم kNN وابسته به داده و بسیار پیچیده هستند و می‌توانند با شکل داده‌های آموزش تطبیق پیدا کنند. برعکس، SVM مرزهای تصمیم صریح را به شکل هایپرپلین‌های حاشیه حداکثر یاد می‌گیرند. SVMها با هدف به دست آوردن یک جداسازی سراسری و واحد بین کلاس‌ها که حاشیه یا فاصله بین هایپرپلین و نزدیک‌ترین نقاط آموزشی (بردارهای پشتیبان) را به حداکثر می‌رساند، بهینه می‌شوند. بنابراین، مرزهای تصمیم SVM تابع‌های نرم خطی یا غیرخطی صاف و مستقل از نقاط داده فردی هستند. پس می‌توان گفت که، مرزهای تصمیم kNN ضمنی و وابسته به داده هستند، در حالی که مرزهای تصمیم SVM تابع‌های صاف صریح و مستقل از توزیع داده هستند. SVMها به دنبال مرزهای تصمیم ساده و مقاوم هستند، در حالی که kNN می‌تواند مرزهای بسیار پیچیده را یاد بگیرد.

با توجه به جدول زیر به سوالات الف تا د پاسخ دهید.

الف) یک طبقه‌بندی چندجمله‌ای با استفاده از روش naive bayes را روی داده‌های جدول تخمین بزنید.  
 ب) طبقه‌بندی به دست آمده را روی داده تست اعمال کنید.  
 ج) یک طبقه‌بندی برنولی با استفاده از روش naive bayes را روی داده‌های جدول تخمین بزنید.  
 د) طبقه‌بندی به دست آمده را روی داده تست اعمال کنید.

نکته: در جدول زیر داک ۱، داک ۲، داک ۳ و داک ۴ مربوط به داده آموزش است و داک ۵ مربوط به داده تست است.

ایندی داکيومنت	کلمات موجود در داکيومنت	آیا مهم است؟
داک ۱	باران تند	بله
داک ۲	برق تند شد	بله
داک ۳	خیوان آمد	خیر
داک ۴	آمد پلنگ تند	خیر
داک ۵	تند تند آمد	؟

$V =$  تعداد ویژگی‌ها

$10 = 2m$  عدد داده‌ها

$$P(C) = \frac{1}{7}$$

$$P(\bar{C}) = \frac{1}{7}$$

$$P(\text{باران} | C) = \frac{1+1}{5+7} = \frac{2}{12}$$

$$P(\text{تند} | C) = \frac{2+1}{5+7} = \frac{3}{12}$$

$$P(\text{شد} | C) = \frac{1+1}{5+7} = \frac{2}{12}$$

$$P(\text{خیوان} | C) = \frac{1+1}{5+7} = \frac{2}{12}$$

$$P(\text{آمد} | C) = P(\text{پلنگ} | C) = \frac{0+1}{5+7} = \frac{1}{12} = P(\text{آمد} | C)$$

$$P(\text{باران} | \bar{C}) = \frac{0+1}{5+7} = \frac{1}{12}$$

$$P(\text{تند} | \bar{C}) = \frac{1+1}{5+7} = \frac{2}{12}$$

$$P(\text{شد} | \bar{C}) = \frac{0+1}{5+7} = \frac{1}{12}$$

$$P(\text{خیوان} | \bar{C}) = \frac{0+1}{5+7} = \frac{1}{12}$$

$$P(\text{باران} | \bar{C}) = \frac{1+1}{5+7} = \frac{2}{12}$$

$$P(\text{تند} | \bar{C}) = \frac{1+1}{5+7} = \frac{2}{12}$$

$$P(\text{آمد} | \bar{C}) = \frac{2+1}{5+7} = \frac{3}{12}$$