

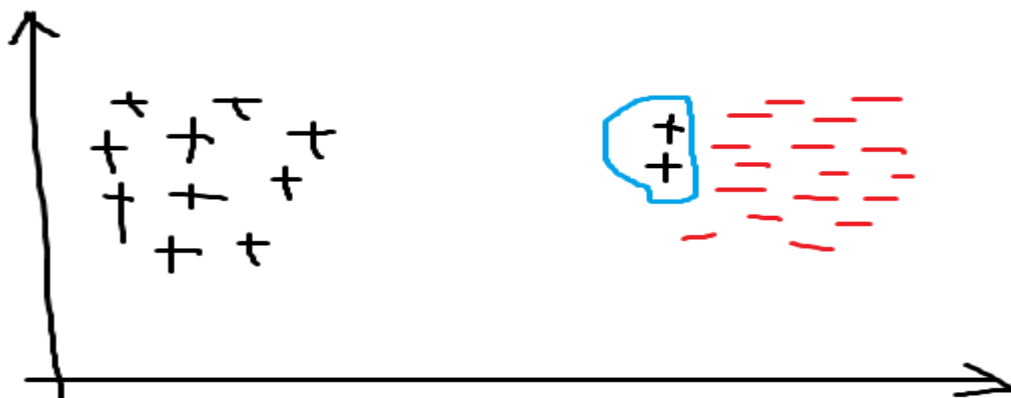


به نام پروردگار
پردازش زبان طبیعی تمرین اول
موعد تحویل :



۱- به سوالات زیر پاسخ مناسب دهید

۱. می‌دانیم Naïve Bayes و Logistic regression هر دو جزء طبقه‌بندهای خطی هستند، هنگامی که دادگان آموزش ما کم هستند به نظر شما کدام طبقه‌بند بهتر عمل می‌کند؟ هنگامی که دادگان آموزش زیاد هستند چگونه؟ چرا؟
۲. مدل‌سازی Naïve Bayes به این صورت است که سند را به صورت BOW مدل‌سازی می‌کند، مزایا و معایب این روش مدل‌سازی را عنوان کنید. اگر در یک تسک شرط Conditional Independence وجود داشته باشد (یعنی ویژگی‌ها به شرط برچسب مستقل از هم باشند) آیا Naïve Bayes می‌تواند یک طبقه‌بند خوب باشد؟
۳. می‌توانید یک مثال بزنید که Naïve Bayes باعث تخمین بیش از اندازه (Overestimate) می‌شود؟ به عبارت دیگر: $P(f1, f2|class) < P(f1|class) \times P(f2|class)$
۴. گفته شد که Naïve Bayes یک طبقه‌بند خطی است. به نظر شما در شکل زیر این طبقه‌بند دو ایتمی را که دور آن خط آبی کشیده شده است جزء کلاس مثبت در نظر می‌گیرد یا کلاس منفی؟ کمی توضیح دهید.



۵. راجع به مشکلات مربوط به نامتوازن بودن دیتاست (imbalanced dataset) و روش‌های مقابله با آن

تحقیق کنید و چند مورد را ذکر کنید؟

۶. تفاوت development test set (devtest) را با test set بیان کنید؟

۲- فرض کنید که اسناد زیر داده شده‌اند:

برچسب	سند
Comedy	fun, couple, love, love
Action	fast, furious, shoot
Comedy	couple, fly, fast, fun, fun
Action	furious, shoot, shoot, fun
Action	fly, fast, shoot, love

در صورت استفاده از طبقه‌بند Naïve Bayes به همراه هموارسازی add-1 به نظر شما سند زیر چه برچسبی خواهد گرفت؟

fast, couple, shoot, fly

۳- جدول زیر ماتریس آشفته‌گی را برای سه کلاس نشان می‌دهد دقت و فراخوانی را برای هر کدام از کلاس‌ها حساب کنید:

		برچسب واقعی		
		A	B	C
خروجی سیستم	A	۱۰	۱۲	۲
	B	۵	۷۰	۸۰
	C	۸	۹	۲۰۰

لطفا به قواعد حل تمرین که در CECM قرار داده شده است توجه کنید.
