



روش‌های یادگیری ماشین در پردازش زبان طبیعی (۸۳۰۴۳۶۸)  
نیم‌سال اول ۱۴۰۳-۱۴۰۲

تاریخ تحویل:

۱۴۰۲/۰۸/۱۳

تمرین شماره ۲

۱. (۵٪) [پژوهش - مدل‌های زبانی بزرگ و یادگیری ماشین] هدف از این بخش استفاده از مدل‌های زبانی بزرگ برای انجام وظایف مرتبط با مدل‌های مولد و جداساز<sup>۱</sup> است. LLMها در بسیاری از وظایف رایج در پردازش زبان طبیعی نسبت به مدل‌های پیشین بهتر عمل کرده‌اند.

قصد داریم که بدانیم که مدل‌های زبانی بزرگ چگونه وظایف مدل‌های مولد و جداساز را به صورت دقیق حل می‌کنند. با بررسی مقاله‌های اخیر (از سال ۲۰۲۰ تا کنون) و انتخاب دو وظیفه (یک وظیفه توسط مدل‌های مولد و یک وظیفه توسط مدل‌های جداساز) به خوبی انجام می‌شود) بیان کنید که چگونه مدل‌های زبانی بزرگ این وظایف را استدلال و پاسخ می‌دهند. توضیحات لازم را به صورت مطالب قابل ارائه برای کلاس درس آماده کنید.

۲. (۱۵٪) [نظری - بیز ساده و توزیع احتمالاتی] هدف از این بخش، مطالعه مبحث‌های بیز ساده و توزیع احتمالاتی است.  
۱. دسته‌بندی برای تفکیک ایمیل‌های غیر تبلیغاتی از تبلیغاتی طراحی کرده‌ایم، این دسته‌بند در ۹۰ درصد مواقع پاسخ درست می‌دهد. در صورتی که فرض کنیم تنها یک درصد ایمیل‌ها تبلیغاتی هستند، اگر ایمیل تبلیغاتی تشخیص داده شود، با چه احتمالی تبلیغاتی است؟

۲. یک پیشامد موفقیت یا شکست با احتمال موفقیت  $P$  مرتباً تکرار می‌شود. پیشامدها در صورت مشاهده‌ی دو موفقیت یا دو شکست پیاپی به پایان می‌رسند. احتمال دو موفقیت پیاپی به شرط آن که دو شکست پیاپی اتفاق نیفتد چیست؟

۳. در یک آزمایش پزشکی سه پارامتر  $x$ ،  $y$  و  $z$  اندازه‌گیری می‌شود. اگر شخصی سالم باشد، توزیع احتمالی این سه پارامتر به صورت گوسی با واریانس‌های به ترتیب ۱، ۴ و ۹ و میانگین‌های به ترتیب ۱، ۴ و ۶ است. برای شخص بیمار نیز توزیع احتمالی این سه پارامتر به صورت گوسی با واریانس‌های به ترتیب ۱، ۴ و ۹ و میانگین‌هایی به ترتیب برابر با ۸، ۶ و ۱۰ است. با فرض اینکه قصد داریم از روش بیز ساده (روشی که در آن فرض می‌شود سه پارامتر  $x$ ،  $y$  و  $z$  مستقل از یکدیگر هستند) برای تفسیر نتایج این آزمایش استفاده کنیم، نشان دهید این روش معادل این است که حاصل عبارت  $\alpha x + \beta y + \gamma z + \lambda$  را بدست آورده و اگر این حاصل مثبت بود، فرد بیمار و در غیر اینصورت فرد سالم در نظر گرفته می‌شود. با فرض اینکه احتمال بیمار بودن فردی که به طور تصادفی از جامعه انتخاب می‌گردد برابر  $\frac{1}{2}$  است، مقادیر  $\alpha$ ،  $\beta$ ،  $\gamma$  و  $\lambda$  را محاسبه کنید.

۳. (۲۰٪) [نظری و پیاده‌سازی - تحلیل احساسات] هدف از این بخش، شناخت نقاط قوت و ضعف و تاثیر پارامترهای مدل بر دقت بیز ساده و استفاده از آن برای دسته‌بندی نظرهای ثبت شده در اسنپ‌فود<sup>۲</sup> است. برچسب هر نظر می‌تواند مثبت (خوشحال یا راضی) و منفی (ناراحت یا ناراضی) است. بایستی احتمال خوشحال بودن نظرهای موجود در دادگان آزمایش را به کمک بیز ساده پیش‌بینی کنید. به این منظور قصد داریم تا با پیاده‌سازی مدل بیز ساده نظرهای مجموعه‌دادگان آزمون را دسته‌بندی کنیم (ما ساختاری را برای شروع پیاده‌سازی در اختیار شما قرار داده‌ایم `naive_bayes.py` لازم است آن را طبق بخش‌های مختلف سوال تکمیل کرده و موارد توضیحی خواسته شده را در گزارش خود ارائه کنید، دقت کنید که برای آموزش مدل خود به صورت تصادفی ۸۰ درصد از دادگان آموزش را جدا کرده و بر آن‌ها مدل را آموزش دهید).

<sup>1</sup> Generative and Discriminative Models

<sup>2</sup> [Snappfood](https://www.snappfood.com/)



روش‌های یادگیری ماشین در پردازش زبان طبیعی (۸۳۰۴۳۶۸)  
نیم‌سال اول ۱۴۰۲-۱۴۰۳

دانشکده علوم و فنون نوین

تاریخ تحویل:  
۱۴۰۲/۰۸/۱۳

تمرین شماره ۲

۱. واحدسازی<sup>۱</sup>، سبب کلمات<sup>۲</sup> و شمارش<sup>۳</sup>: پیش از ساخت مدل باید متن را به نمایشی تبدیل کنیم که برای مدل قابل فهم است. در اینجا از نمایش BOW استفاده می‌کنیم. مدل Naïve Bayes به ترتیب کلمات در سند اهمیت نمی‌دهد و فقط تعداد دفعه‌هایی که کلمه‌ها تکرار می‌شوند مهم است.

۱. تابع `tokenize_doc(doc)` را پیاده‌سازی کنید. این تابع یک سند (یک رشته) را به عنوان ورودی دریافت می‌کند و آن را به کلمات مختلف جدا می‌کند (از فاصله برای جداسازی استفاده کنید). سپس هر کلمه را به تعداد دفعاتی که در آن سند تکرار شده، توسط یک دیکشنری تصویر کرده و این دیکشنری را به عنوان خروجی باز می‌گرداند. خطوط مرتبط را در فایل برای ارزیابی کد خود از کامنت در بیاورید و اجرا کنید.

۲. تابع `update_model` را پیاده‌سازی کنید. پیش از شروع، حتما کامنت‌های مرتبط با این تابع را مطالعه کنید، تا بدانید چه چیزی را باید به‌روزرسانی کنید. همچنین به آرگومان‌هایی که کلاس NaiveBayes دریافت می‌کند دقت کنید و در پیاده‌سازی‌های خود در نظر بگیرید. تابع `train_model` را اجرا کنید. تعداد کلمات موجود در دادگان آموزش چقدر است؟ بیاید به بررسی شمار آمارگان ذخیره شده توسط `update_model` بپردازیم. به این منظور تابع `top_n` را پیاده‌سازی کنید.

۳. این تابع یک دسته (مقادیر برچسب‌های احساس، برای مثال؛ دسته/کلاس مثبت یا منفی) و یک عدد  $n$  را به عنوان ورودی دریافت می‌کند و  $n$  کلمه پرتکرار که در سند نسبت به دسته مورد نظر تکرار شده‌اند را به عنوان خروجی باز می‌گرداند. به کمک این تابع، ۱۰ رایج‌ترین کلمه در دسته احساس مثبت و منفی را بدست آورده و در گزارش خود ارائه کنید.

۴. آیا ۱۰ رایج‌ترین کلمه در دسته مثبت/منفی کمکی به تفکیک این دو دسته انجام می‌دهد؟ آیا این فرایند برای دیگر مجموعه‌دادگان در زبان فارسی به نتایج مشابهی منجر می‌شود؟

۲. احتمال کلمات و شبه‌شماره‌ها<sup>۴</sup>: مدل بیز ساده فرض می‌کند که همه ویژگی‌ها با توجه به برچسب کلاس مستقل هستند. در این مساله، این بدان معنی است که احتمال دیدن یک کلمه خاص در یک سند با برچسب کلاس  $y$  مستقل از بقیه کلمات در آن سند است.

۱. تابع `p_word_given_label` را پیاده‌سازی کنید. این تابع باید مقدار  $p(w|y)$  را حساب کند (به عبارت دیگر؛ احتمال مشاهده کلمه  $w$  با فرض اینکه برچسب سند  $y$  است).

۲. به کمک این تابع احتمال کلمه «عالی» و «افتضاح» را برای هر دو برچسب مثبت و منفی حساب کنید. کدام کلمه احتمال بیشتری در دسته مثبت دارد؟ کدام کلمه احتمال بیشتری در دسته منفی دارد؟ آیا این نتایج مطابق انتظار است؟

۳. توضیح دهید که اگر کلمه‌ای در دسته مثبت داده‌های آموزش است ولی در دسته منفی نیست (و حتی بالعکس)، احتمال آن را چگونه محاسبه می‌کنید؟ شرح دهید که چه چیزی مشکل ایجاد می‌کند؟

۴. این مشکل را می‌توان به کمک شبه‌شماره‌ها حل کرد. شبه‌شماره یک مقدار ثابت است که به تعداد تکرار هر کلمه در مدل اضافه می‌شود. از این مقدار برای هموارسازی<sup>۵</sup> محاسبه کلماتی استفاده می‌شود که اطلاعات کمی از آن‌ها در دادگان موجود است. به این منظور تابع `p_word_given_label_and_pseudocount` را پیاده‌سازی کنید.

<sup>1</sup> Tokenization

<sup>2</sup> Bag-of-Words (BOW)

<sup>3</sup> Counting

<sup>4</sup> Pseudocounts

<sup>5</sup> Smoothing



روش‌های یادگیری ماشین در پردازش زبان طبیعی (۸۳۰۴۳۶۸)  
نیم‌سال اول ۱۴۰۳-۱۴۰۲

تاریخ تحویل:  
۱۴۰۲/۰۸/۱۳

تمرین شماره ۲

۳. احتمال پیشین<sup>۱</sup> و درست‌نمایی<sup>۲</sup>: همانطور که پیش‌تر گفته شد، مدل بیز ساده فرض می‌کند که همه کلمات در یک سند با توجه به برچسب اسناد مستقل از یکدیگر هستند. به همین دلیل می‌توانیم درست‌نمایی یک سند را به صورت زیر تعریف کنیم:

$$p(w_{d_1}, \dots, w_{d_n} | y_d) = \prod_{i=1}^n P(w_{d_i} | y_i)$$

در اینجا  $w_{d_i}$ ، امین کلمه در سند  $d$  و  $y_d$ ، برچسب سند  $d$  است. با این حال اگر سندی دارای کلمه‌های زیادی باشد، احتمال آن بسیار کوچک می‌شود و با کمبود حسابی<sup>۳</sup> مواجه خواهیم شد. هنگام کار با مدل‌های احتمالی، کمبود حسابی یک مشکل رایج است. اگر با آن آشنا نیستید، می‌توانید یک نمای کلی از [ویکی‌پدیا](#) دریافت کنید. یکی از راه‌حل‌های رایج در برخورد با این مشکل، تبدیل فضا به فضای لگاریتمی<sup>۴</sup> است.

۱. لگاریتم درست‌نمایی (همانطور که در عبارت بالا آمده است) را حساب کنید.

۲. تابع `log_likelihood` را پیاده‌سازی کنید (دقت کنید که باید از فراخوانی تابع `p_word_given_label_and_pseudocount` استفاده کنید).

۳. تابع `log-prior` را پیاده‌سازی کنید. این تابع یک برچسب را به عنوان ورودی می‌گیرد و لگاریتم قسمتی از دادگان آموزش که این برچسب را دارند، باز می‌گرداند.

۴. نرمال‌سازی<sup>۵</sup> و قانون تصمیم‌گیری<sup>۶</sup>: بیز ساده مدلی است که به ما می‌گوید چگونه احتمال پسین یک سند با کلاسی (به عبارت دیگر؛  $p(y_d | w_d)$ ) را محاسبه کنیم. به کمک قانون بیز داریم که:

$$p(y_d | w_d) = \frac{p(y_d)p(w_d | y_d)}{p(w_d)}$$

در بخش قبل تابع‌های لازم برای محاسبه  $\log \text{prior}(\log[P(y_d)])$  و  $\log \text{likelihood}(\log[P(w_d | y_d)])$  را پیاده‌سازی کرده‌اید. حال لازم است که نرمال‌ساز را محاسبه کنید ( $P(w_d)$ ).

۱. نرمال‌ساز<sup>۷</sup>  $p(w_d)$  را بدست آورید.

۲. لگاریتم احتمال پسین را با لگاریتم گرفتن از تساوی بالا، بدست آورید.<sup>۸</sup>

۳. یک راه برای دسته‌بندی یک سند این است که لگاریتم احتمال پسین نرمال‌نشده برای هر دو برچسب محاسبه کنیم و سپس `argmax` بگیریم (به عبارت دیگر؛ برچسبی که منجر به لگاریتم احتمال پیشین نرمال‌نشده بیشتری می‌شود). لگاریتم احتمال پسین نرمال‌نشده از جمع لگاریتم احتمال پیشین و لگاریتم درست‌نمایی یک سند بدست می‌آید. چرا به محاسبه لگاریتم نرمال‌ساز نیازی نداریم؟

۴. تابع `unnormalized_log_posterior` را پیاده‌سازی کنید.

۵. تابع `classify` را پیاده‌سازی کنید (این تابع باید از لگاریتم نرمال‌نشده احتمال پسین استفاده کند، اما نباید نرمال‌ساز را محاسبه کند).

<sup>1</sup> Prior

<sup>2</sup> Likelihood

<sup>3</sup> Arithmetic Underflow

<sup>4</sup> Log-Space

<sup>5</sup> Normalize

<sup>6</sup> Decision Rule

<sup>7</sup> Normalizer is denoted as  $P(w_d)$ , which is the probability of observing the document regardless of its label.

<sup>8</sup> The log of the posterior probability is the logarithm of the expression representing the probability of a certain label given the observed data. The log transformation is often applied to simplify calculations and avoid numerical underflow issues.

روش‌های یادگیری ماشین در پردازش زبان طبیعی (۸۳۰۴۳۶۸)  
نیم‌سال اول ۱۴۰۳-۱۴۰۲



تاریخ تحویل:  
۱۴۰۲/۰۸/۱۳

تمرین شماره ۲

۵. ارزیابی: پس از آموزش مدل و پیاده‌سازی عملکرد دسته‌بندی، قصد داریم که دقت مدل را ارزیابی کنیم.
۱. تابع `evaluate_classifier_accuracy` را پیاده‌سازی کنید. این تابع باید تمام نمونه‌های موجود در مجموعه‌داده‌گان آزمایش را برچسب‌گذاری کند. سپس قسمتی از داده‌گان آزمایش را که درست برچسب‌گذاری شده‌اند را گزارش کنید (شبه‌شماره را برابر با ۱ در نظر بگیرید).
  ۲. در این بخش با مقادیر مختلف شبه‌شماره (مقادیر؛ [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1]) دقت مدل را بررسی کنید. نتایج را به صورت نمودار نمایش دهید (از تابع `plot_pseudocounts_vs_accuracy` استفاده کنید).
  ۳. یک نمونه بیابید که دسته‌بند آن را اشتباهی تشخیص داده باشد. فکر می‌کنید که چرا مدل اینگونه عمل کرده است؟ چه راهکارهایی پیشنهاد می‌دهید که مدل این نمونه را درست تشخیص دهد.
  ۴. تحلیل کاوشگرانه<sup>۱</sup>: حال از آنجا که مدل ما آموزش دیده است، می‌توان به تحلیل کاوشگرانه داده‌ها پرداخت. دیدیم که ۱۰ رایج‌ترین کلمه<sup>۲</sup> برای هر کلاس چندان تبعیض آمیز نیستند. اغلب اوقات، یک آمار توصیفی‌تر، نرخ درست‌نمایی کلمات<sup>۳</sup> است. نرخ درست‌نمایی کلمات را به صورت تعریف می‌کنیم:

$$LR(w) = \frac{P(w | y = positive)}{P(w | y = negative)}$$

۱. حال به کمک مدل آماده شده، به با توجه به تعریف، اگر LR برای کلمه‌ای برابر با ۵ بدست آید، به این معنی است که احتمال تعلق داشتن آن کلمه به دسته مثبت ۵ برابر محتمل‌تر از تعلق داشتن به دسته منفی است. همچنین اگر LR برای کلمه‌ای برابر با ۰.۳۳ بدست آید، به این معنی است که احتمال تعلق داشتن آن کلمه به دسته مثبت یک سوم احتمال تعلق آن به دسته منفی است.
۲. خروجی تابع LR در چه بازه‌ای است؟
۳. تابع `likelihood_ratio` را پیاده‌سازی کنید (این تابع یک کلمه را به عنوان ورودی می‌گیرد و نرخ درست‌نمایی کلمات را مطابق فرمول بالا محاسبه می‌کند).
۴. LR را برای دو کلمه «عالی» و «افتضاح» حساب کنید. مقادیر بدست آمده را با نرخ ۱۰ رایج‌ترین کلمه که در بخش‌های قبل بدست آوردید، مقایسه کنید.
۵. توضیح دهید که چگونه مقادیر LR کلمات به مدل بیز ساده مرتبط هستند؟ برای مثال اگر LR کلمه‌ای ۱ باشد، آیا به این معنی است این کلمه برای دسته‌بند بیز ساده با اهمیت است یا خیر؟ اگر LR برای کلمه‌ای بسیار با مقدار ۱ اختلاف داشته باشد (برای مثال؛ ۰.۰۱ و ۱۰۰)، آیا به این معناست که این کلمه برای دسته‌بند بیز ساده با اهمیت/بی اهمیت است؟ همچنین بیان کنید که کلمه‌ای که دارای  $LR=0.01$  است، در مقایسه با کلمه‌ای که  $LR=100$  دارد، در دسته‌بندی چه تاثیری دارد؟ (موارد خواسته شده را به صورت دقیق توضیح دهید)
۷. تعمیم به چند کلاس: در دنیای واقعی مساله‌های دسته‌بندی چند کلاسه<sup>۵</sup> نسبت به دسته‌بندی دو کلاسه به مراتب بیشتر هستند.

<sup>1</sup> Evaluation

<sup>2</sup> Exploratory Analysis

<sup>3</sup> Top 10 Most Common Words

<sup>4</sup> Words Likelihood Ratio (LR)

<sup>5</sup> Multi-Class Classification



روش‌های یادگیری ماشین در پردازش زبان طبیعی (۸۳۰۴۳۶۸)  
نیم‌سال اول ۱۴۰۳-۱۴۰۲

تاریخ تحویل:  
۱۴۰۲/۰۸/۱۳

تمرین شماره ۲

۱. اگر قصد داشته باشیم مدل حال حاضر را به منظور دسته‌بندی چند کلاسه تغییر دهیم، آمار شمارشی<sup>۱</sup> که ذخیره می‌کنیم چگونه تغییر می‌کند؟ تمام نمونه‌های موجود در مجموعه‌دادگان آزمایش را دسته‌بندی کنید و بخشی از نمونه‌هایی را که به درستی دسته‌بندی شده‌اند را در گزارش خود ارائه نمایید. دقت دسته‌بندی خود را با شبه‌شماره برابر با مقدار ۱، گزارش کنید.
۲. نرمال‌ساز چگونه تغییر می‌کند؟
۳. قانون تصمیم‌گیری جدید چه خواهد بود (به عبارت دیگر، تابع دسته‌بندی چگونه تغییر می‌کند)؟
۸. ارزیابی مجموعه‌دادگان آزمایش: دادگان آزمایش test.csv در اختیار به شما قرار گرفته است، لازم است که شما به کمک مدل بیز ساده که تکمیل کردید به دادگان تست برچسب بزنید. قصد داریم به میزان خوب بودن مدل شما امتیازی بدهیم، امتیاز نهایی مدل شما تابعی از سطح زیر نمودار منحنی مشخصه عملکرد سیستم است. AUC-ROC مشخصه‌ای از میزان تطابق بین برچسب‌های داده‌شده و برچسب‌های پیش‌بینی شده است. امتیاز نهایی مدل به کمک رابطه زیر محاسبه می‌شود.

$$\text{score} = 4 \times (\text{AUC-ROC} \times 100 - 50)$$

۴. (۴۰٪) [پایاده‌سازی - سخنان تنفرپراکنانه] هدف از این بخش تمرین سه روش Naïve Bayes، Logistic Regression و KNN است. دادگان سخنان تنفرآمیز<sup>۲</sup> در اختیار شما قرار گرفته است. قصد داریم به کمک الگوریتم‌های یادگیری ماشین، از دادگان آموزش برای یادگیری مدل و دادگان آزمایش برای بررسی وضعیت یادگیری مدل استفاده کنیم. اگر سخنی دارای سه ویژگی زیر باشد، به عنوان سخن تنفرآمیز در نظر گرفته می‌شود:
  - حمله عمدی<sup>۳</sup> باشد.
  - به سمت یک گروه خاصی از افراد باشد.
  - توسط ویژگی‌های آن گروه منشأ بگیرد.

برچسب هر عبارت با مقدار ۱ (به معنای تنفرآمیز بودن) و ۰ (به معنای تنفرآمیز نبودن) نشان داده می‌شود. ساختار پیشنهادی برای طراحی در قالب فایل‌های utils.py، classifier.py و main.py در اختیار شما قرار گرفته است. برای پایاده‌سازی دسته‌بندی‌های Naïve Bayes، Logistic Regression و KNN شما نیاز به پایاده‌سازی الگوریتم‌های خواسته شده در utils.py و classifier.py دارید. توضیح‌های مورد نیاز برای هر پایاده‌سازی در ادامه برای شما ارائه می‌شود.

۱. استخراج ویژگی: در utils.py تابع tokenize عمل جداسازی کلمه‌ها با الگوی داده شده انجام می‌شود. برای مثال جمله "I love nlp" را به لیستی از کلمه‌های ["I", "love", "nlp"] تبدیل می‌کند. به کمک این تابع می‌خواهیم استخراج ویژگی‌های یک‌تایی‌ها<sup>۴</sup> و دو تایی‌ها<sup>۵</sup> را پایاده‌سازی کنیم. برای هر «استخراج‌کننده ویژگی» چهار تابع fit، init، transform و transform\_list تعریف شده است که لازم است آن‌ها را تکمیل کنید. در ادامه، تابع‌های اشاره شده برای روش Unigram توضیح داده می‌شود.

`__init__(self)`

در این تابع یک دیکشنری به نام unigram را مقداردهی اولیه می‌کنید.

`fit(self, text_set: list)`

در این قسمت لازم است تا دیکشنری unigram را به کمک `text_set` که لیستی از جمله‌های جداسازی شده است، مقداردهی کنید. برای مثال اگر:

<sup>1</sup> Count Statistics

<sup>2</sup> Hatespeech

<sup>3</sup> Deliberate attack

<sup>4</sup> UnigramFeature

<sup>5</sup> BigramFeature



```
Text_list = [{"I", "love", "nlp"}, {"I", "like", "python"}]
```

در این صورت پس از اجرای تابع بر روی text\_list و self.unigram، نتیجه به صورت زیر خواهد بود:

```
{"I": 0, "love": 1, "nlp": 2, "like": 3, "python": 4}
```

توجه داشته باشید که اندازه حروف<sup>۱</sup> اهمیتی ندارد، برای مثال دو کلمه "Python" و "pYthON" یک کلمه در نظر گرفته می‌شوند.

```
transform(self, text: list)
```

در این تابع به کمک دیکشنری بدست آمده از بخش قبل، یک بردار ویژگی برای جمله text بدست آورید. برای مثال اگر:

```
Text = ["I", "love", "nlp"]
```

در این صورت پس از اجرای تابع بر روی text، بردار ویژگی زیر را به عنوان خروجی بازمی‌گردانیم:

```
[1, 1, 1, 0, 0]
```

```
transform_list(self, text_set: list)
```

برای لیستی از جمله‌های جداسازی شده، لیستی از بردارهای ویژگی را به عنوان خروجی بازگردانید. برای مثال:

```
Text_list = [{"I", "love", "nlp"}, {"I", "like", "python"}]
```

در این صورت خروجی این تابع باید یک آرایه NumPy به شکل زیر باشد:

```
array([[1, 1, 1, 0, 0], [1, 0, 0, 1, 1]])
```

توضیحات داده شده برای هر چهار تابع را در کلاس UnigramFeature تکمیل کنید و به کمک این توابع یک استخراج‌کننده bigram در کلاس BigramFeature نیز بنویسید.

۲. پیاده‌سازی دسته‌بندهای Naïve Bayes، Logistic Regression و KNN: در Classifier.py دو کلاس NaiveBayesClassifier و

LogisticRegressionClassifier و KNNClassifier داده شده است، برای هر الگوریتم لازم است که سه تابع fit، predict و predict را تکمیل کنید. در ادامه جزئیات تابع‌های گفته شده برای دسته‌بند Naïve Bayes ارائه می‌شود:

```
__init__(self)
```

دسته‌بند Naïve Bayes از احتمال هر دسته و احتمال ویژگی‌ها برای هر دسته برای دسته‌بندی استفاده می‌کند، پس دو دیکشنری برای ذخیره این احتمالات تعریف کنید.

```
fit(self, X, Y)
```

در این قسمت باید احتمال‌های تعریف شده را برای دسته‌ها و ویژگی‌ها برای داده X و برچسب y حساب کنید (از laplace smoothing استفاده کنید).

```
predict(self, X)
```

برای ارزیابی مدل بدست آمده، تابعی بنویسید که دادگان آزمایش را توسط مدلی که آموزش داده‌اید، برچسب‌گذاری کند (خروجی این تابع باید یک آرایه NumPy است).

توابع گفته شده را در کلاس NaiveBayesClassifier کامل کنید. سپس همین کار را برای مدل LogisticRegression در کلاس LogisticRegressionClassifier انجام دهید.

در رگرسیون لاجیستیک بردار وزنی به اندازه تعداد ویژگی‌ها تعریف کنید. سپس این وزن‌ها را در تابع fit به کمک فرمول‌های زیر، بروزرسانی کنید.

محاسبه تابع هزینه:

<sup>۱</sup> Lowercase or Uppercase



$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^i * \log(h\theta(x^i)) + (1 - y^i) * \log(1 - h\theta(x^i))]$$

محاسبه گرادیان:

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h\theta(x^i) - y^i) * x_j^i$$

بروزرسانی وزن‌ها به کمک گرادیان:

$$\theta_j = \theta_j - \alpha * \frac{\partial J(\theta)}{\partial \theta_j}$$

به صورت پیش‌فرض از نرخ یادگیری 0.01 و تعداد اپیاک 100 استفاده کنید و به کمک وزن‌های بدست آمده، دادگان آزمایش را پیش‌بینی کنید.

$$Predict = \sigma(XW)$$

۳. اجرای مدل‌ها و ارزیابی نتایج: در این قسمت دقت مدل‌های ساخته شده را بررسی می‌کنیم. به این منظور، خط فرمان زیر را

در محیط Colab اجرا کنید (اگر از cmd/shell/bash استفاده می‌کنید، به جای python از python استفاده کنید).

`!python main.py -model AlwaysPredictZero`

در صورت موفق بودن اجرای این فرمان، خروجی به صورت زیر خواهد بود.

```
Namespace(model='AlwaysPredictZero', feature='unigram', path='./data/')
==== Train Accuracy ====
Accuracy : 956 / 1913 = 0.4997
```

جهت ارزیابی مدل NaiveBayesClassifier با استخراج‌کننده ویژگی یک‌تایی، دستور زیر را اجرا کنید:

`!python main.py -model NaiveBayes -feature unigram`

دقت بدست آمده برای مدل‌های Naive Bayes، Logistic Regression و KNN به کمک دو روش استخراج ویژگی یک‌تایی و دو تایی در جدول زیر گزارش کنید.

جدول ۱: نتایج مدل‌های بدست آمده

| استخراج‌کننده<br>ویژگی دو تایی | استخراج‌کننده<br>ویژگی یک تایی | استخراج ویژگی<br>مدل |
|--------------------------------|--------------------------------|----------------------|
|                                |                                | بیز ساده             |
|                                |                                | رگرسیون لاجستیک      |
|                                |                                | k-نزدیک‌ترین همسایه  |

۵. (۲۰٪) [پیاده‌سازی - نزدیکی اعداد] هدف از این تمرین پیاده‌سازی الگوریتم نزدیک‌ترین همسایه است. در این بخش شما می‌بایست به طراحی

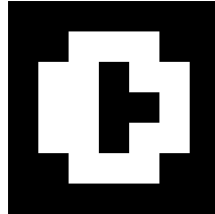
و پیاده‌سازی طبقه‌بند k-نزدیک‌ترین همسایه<sup>۱</sup> بپردازید، همچنین دادگان ارقام دست‌نویس دودویی<sup>۲</sup> در اختیار شما قرار گرفته است. مجموعه داده ارائه شده شامل اعداد ۰ تا ۹ بوده که در قالب مجموعه‌ای از صفرها و یک‌ها در ابعاد ۲۸ در ۲۸ نمایش داده می‌شوند.

مجموعه دادگان: هر نمونه در مجموعه دادگان دارای ارتفاعی با ۲۸ پیکسل و عرضی با ۲۸ پیکسل، که در مجموع ۷۸۴ پیکسل است. هر پیکسل می‌تواند دو مقدار صفر و یک را بگیرد، شکلی که پیکسل‌های دارای مقدار «۱» نشان می‌دهند، همان عددی است که توسط کاربر نوشته شده است

<sup>۱</sup> KNN Classifier

<sup>۲</sup> Digit-recognizer

و بدین صورت مجموعه دادگان آموزشی، دارای ۷۸۵ ستون است، که ستون برچسب<sup>۱</sup> رقمی است که توسط کاربر ترسیم شده است.



```
00000000
00111000
01101100
01100100
01101100
00111000
00000000
```

شکل ۱: نمایش عدد صفر در مجموعه دادگان

۴. پیاده‌سازی: کلاس KNN که پیاده‌سازی می‌کنید، باید حداقل دارای متدهای زیر باشد.

```
__init__(self, k=3)
```

طبقه‌بند KNN را با مقدار پیش‌فرض  $K$  (تعداد همسایگان) به صورت ۳ مقداردهی کنید.

```
fit(self, X_train, y_train)
```

این متد مدل KNN را بر ماتریس مشخصه  $X_{train}$  و بردار هدف  $y_{train}$  آموزش می‌دهد.

```
predict(self, X_test)
```

این متد داده‌های آزمون  $X_{test}$  را می‌پذیرد و برچسب‌های پیش‌بینی شده را بر می‌گرداند.

هر گونه متد یا کلاس اضافی را که برای راه‌حل خود ضروری می‌دانید، اضافه کنید (هدف و استفاده از هر جزء اضافی را مستند کنید).

۵. ارزیابی: لازم است که گزارش کاملی از رویکرد خود، تصمیمات طراحی، چالش‌های پیش‌روی و راه‌حل اجرا شده، تاثیر مقادیر مختلف  $k$  بر

عملکرد مدل را مورد بحث قرار دهید. عملکرد شما بر اساس معیارهای زیر ارزیابی خواهد شد: «دقت، صحت، فراخوانی و ماتریس درهم‌ریختگی

دسته‌بند KNN» و «نشان‌دادن درک تاثیر پارامترهای کلیدی، مانند  $K$  در مدل خود».

<sup>۱</sup> Label