

گزارش تکلیف ۴ درس یادگیری ماشین

کسرا سینایی

شماره دانشجویی ۸۱۰۶۹۶۲۵۴

۹ دی ۱۴۰۰

۱ سوال یک

الف

در روش KNN هرچه مقدار K را کوچکتر انتخاب کنیم پیچیدگی مدل و مرز تصمیم افزایش می‌یابد و به همین نسبت جنرالیزیشن مدل کمتر و حساسیت به نویز نیز بیشتر می‌شود. در واقع هر چه K کوچکتر باشد واریانس بیشتر خواهد بود. با افزایش K به حد بهینه‌ای برای مدل می‌رسیم و افزایش بیشتر آن منجر به خطای بایاس می‌شود. در تخمین پنجره پارزن هاپرپارامتر مدل طول پنجره در تابع کرنل می‌باشد. این پارامتر (h_n) مانند یک smoothing factor عمل می‌کند. هر چه طول پنجره بزرگتر باشد شکل تابع به دست آمده برای توزیع متغیرها نرم‌تر می‌شود. برای h های بسیار کوچک تخمین نان پارامتری پارزن به تابع دلتا دیراک شبیه خواهد شد. به زبان دیگر انتخاب طول بزرگ برای پنجره منجر به خطای بایاس و انتخاب اندازه کوچک برای طول پنجره منجر به خطای واریانس می‌شود.

ب

در روش‌های پارامتریک فرض می‌کنیم داده‌ها از یک توزیع شناخته شده و فرمول‌بندی شده مانند گوسی، پواسون و ... آمده‌اند و برای توصیف آن‌ها به تخمین پارامترهای مدل فرض شده می‌پردازیم. در روش‌های نان پارامتری هیچ فرض اضافه‌ای برای مدل حاکم بر داده‌های یادگیری قرار نمی‌دهیم و با قوانینی مانند نزدیک‌ترین همسایه، استفاده از تابع کرنل و ... توزیعی دلخواه به داده‌های یادگیری فیت می‌کنیم.

ج

چون این روش‌ها نان پارامتریک هستند لازم است برای پیش‌بینی نقاط دیده نشده تمام مجموعه یادگیری را همراه با مدل نگه داریم. (در واقع مدل مستق و به کمک یک تابع و چند پارامتر تعریف نشده که بتوان با استفاده از آن پیش‌بینی کرد، بلکه مدل با استفاده از داده‌های مجموعه یادگیری در هر نقطه تعریف می‌شود) با افزایش ابعاد فیچرها برای دست‌یابی به تخمینی دقیق تعداد نقاط مورد نیاز در دیتاست به صورت نمایی افزایش می‌یابد و در مواردی که نتوانیم حجم داده‌ها را افزایش دهیم ممکن است دچار نفرین ابعادی شویم.

د

۲ سوال دو

در این سوال می‌خواهیم نشان دهیم که سلول‌های voronoi تشکیل شده با الگوریتم نزدیک‌ترین همسایه محدب هستند. به ازای هر دو نقطه درون یک سلول، خط وصل کننده آن‌ها نیز داخل سلول قرار خواهد گرفت، فرض کنید x^* در یک سلول با لیبل مشخص واقع شده باشد و y نیز نقطه‌ای با یک لیبل متفاوت از x^* باشد. در این حالت ابر صفحه‌ای وجود دارد که فضا را طوری تقسیم می‌کند که یک طرف آن نقاط به x^* نزدیک تر هستند و طرف دیگر آن نقاط به y نزدیک ترند. اگر نقاط $x(0)$ و $x(1)$ در سلول متعلق به x^* قرار داشته باشند خطی که این دو نقطه را به هم وصل می‌کند نیز به صورت: $x(\lambda) = (1 - \lambda)x(0) + \lambda x(1); 0 \leq \lambda \leq 1$ خواهد بود. این خط نیز به طور کامل درون سلول متعلق به نقطه x^* قرار خواهد گرفت. چون فضای نصف شده توسط ابر صفحه محدب است، تمامی نقاط موجود بر روی خط تعریف شده $(x(\lambda))$ نیز به x^* نزدیک تر هستند تا y .

این موضوع به ازای تمام نقاط موجود در Voronoi Cell متعلق به نقطه x^* صادق است. همچنین برای سایر سمپل‌ها نیز صادق است (شرطی برای y نداریم) بنابراین Voronoi Cell تشکیل شده با قانون نزدیک‌ترین همسایه همواره محدب می‌باشند.

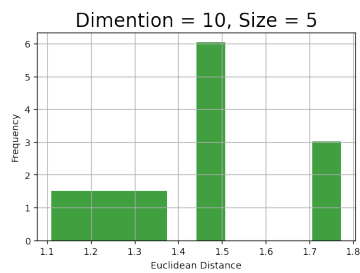
۳ سوال سه

۴ سوال چہار

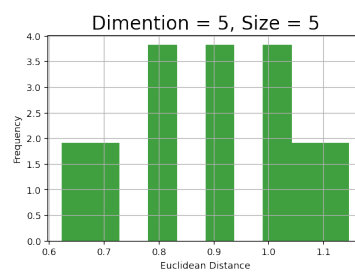
۵ سوال پنج

۶ سوال شش (شبیه سازی)

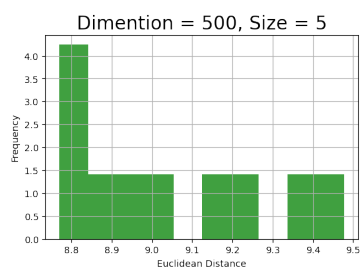
الف



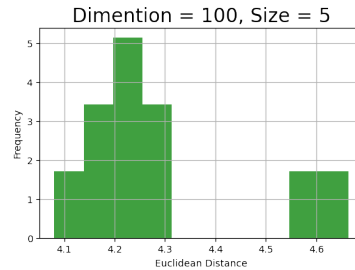
(ب) ۱۰ بعدی



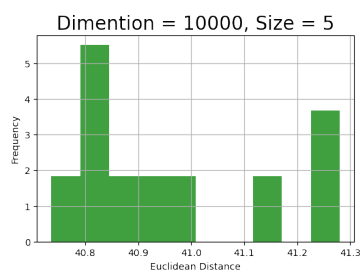
(آ) ۵ بعدی



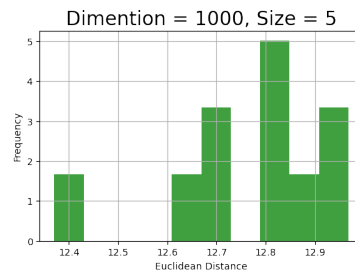
(د) ۵۰۰ بعدی



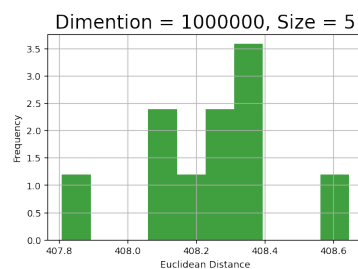
(ج) ۱۰۰ بعدی



(و) ۱۰۰۰۰ بعدی



(ه) ۱۰۰۰ بعدی

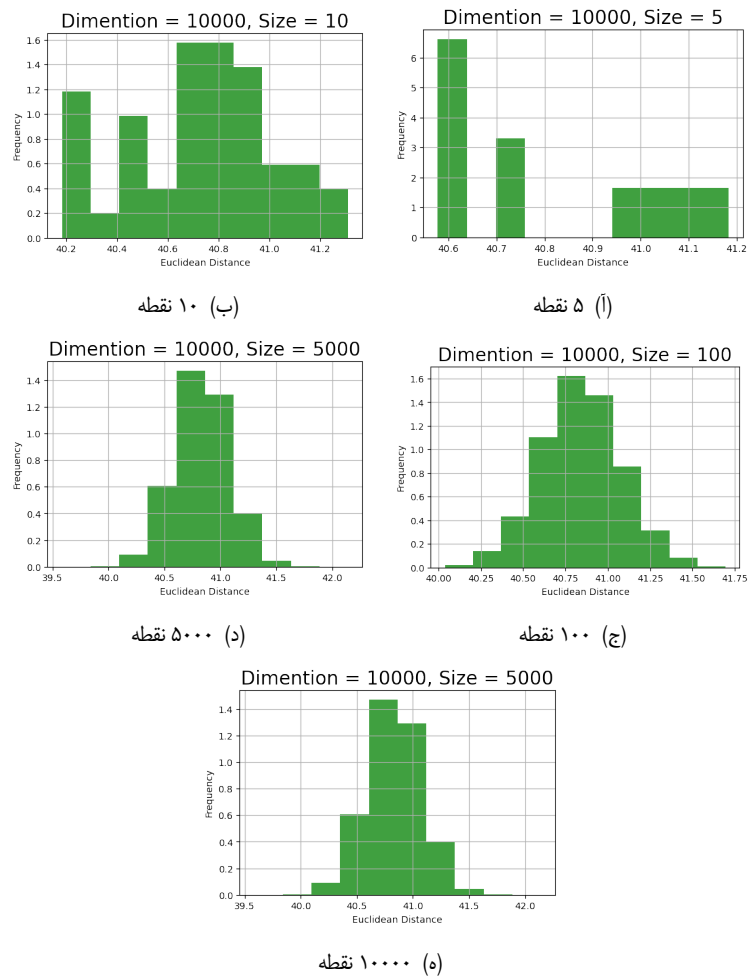


(ز) ۱۰۰۰۰۰ بعدی

شکل ۱: خواسته‌های قسمت الف

در این سوال با استفاده از تابع `rndom` کتابخانه `numpy` و تابع `hist` از کتابخانه `matplotlib` هیستوگرام فاصله نقاط را به دست می‌آوریم. برا این کار تابع `do task` پیاده سازی شده است که تعداد و بعد نقاط مورد نظر را در ورودی گرفته و نمودارها را رسم می‌کند.

ب



شکل ۲: خواسته‌های قسمت ب

ج

۷ سوال هفت (شبیه سازی)

در این سوال با استفاده از الگوریتم KNN تصاویر اعداد دست نویسی فارسی را طبقه بندی می‌کنیم.

الف

برای لود کردن دیتاست از روی فایل‌های cdb از اسکریپت پایتون معرفی شده در لینک قرار داده شده در صورت سوال استفاده می‌کنیم. یکی از توابع موجود در این اسکریپت read_hoda_dataset می‌باشد. آرگومان‌های ورودی که آدرس فایل دیتاست و سائز تصاویر می‌باشد را مطابق توضیحات پیچ گیت‌هاب به تابع می‌دهیم. آرگومان reshape را نیز True قرار می‌دهیم و آرگومان one_hot را False، در این صورت فیچرها که پیکسل عکس‌های 32×32 هستند فلت شده و لیبل‌ها نیز یک عدد بین ۰ و ۹ خواهند بود. در شکل ۳ سه نمونه از تصاویر متعلق به هر کلاس نمایش داده شده‌اند. (رنگ‌ها غیر واقعی می‌باشند، تمام تصاویر یک کاناله و gray scale هستند).

ب

در این قسمت الگوریتم KNN را به صورت دستی پیاده سازی می‌کنیم. به این منظور یک کلاس به نام KNN با سه متد زیر پیاده سازی شده‌است:

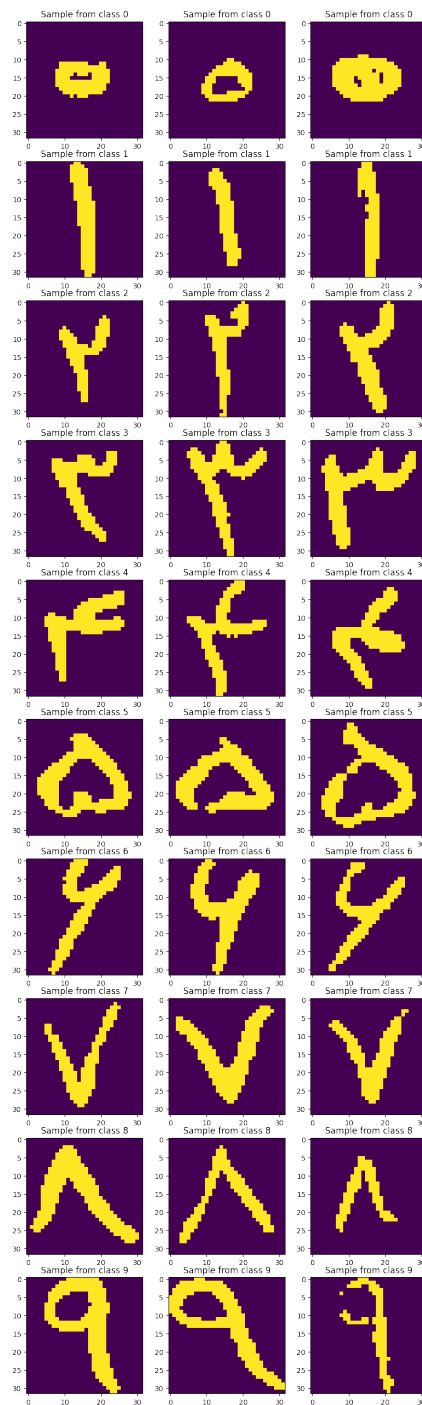
- constructor: این تابع تنها هاپیر پارامتر مدل که همان k یا تعداد همسایه‌ها می‌باشد را به همراه دیتاست و فیچرها می‌گیرد و در متغیرهای کلاس ذخیره می‌کند. تعداد سمپل‌ها و ابعاد فیچر نیز از دیگر متغیرهای این کلاس هستند که از روی ورودی‌های این تابع به دست می‌آیند. لازم به ذکر است که ورودی‌های مربوط به سمپل‌ها و لیبل‌ها باید آرایه numpy باشند.
- dist: این تابع به صورت برداری فاصله اقلیدسی دو بردار که آرگومان‌ها ورودی آن هستند را حساب می‌کند.
- predict: در این تابع برای هر سطر از متغیر ورودی (x) تخمین نان‌پارامتریک مدل را به دست می‌آوریم. ابتدا با استفاده از متد dist و تابع argsort کتابخانه numpy همسایه‌های نقاط را به دست می‌آوریم. در نهایت با استفاده از تابع unique و argmax پر تکرارترین لیبل در بین k همسایه نزدیک نقطه مورد بررسی را به دست آورده و آن نقطه را در همان طبقه قرار می‌دهیم.

سرعت محاسبه تابع predict در مقایسه با توابع آماده کتابخانه scikit learn به مقدار قابل توجهی کمتر است. لیبل کردن ۲۰۰۰ دیتا از مجموعه تست حدوداً ۳۰ دقیقه زمان برد. نتایج طبقه بندی در جدول ۱ و شکل ۴ آورده شده است.

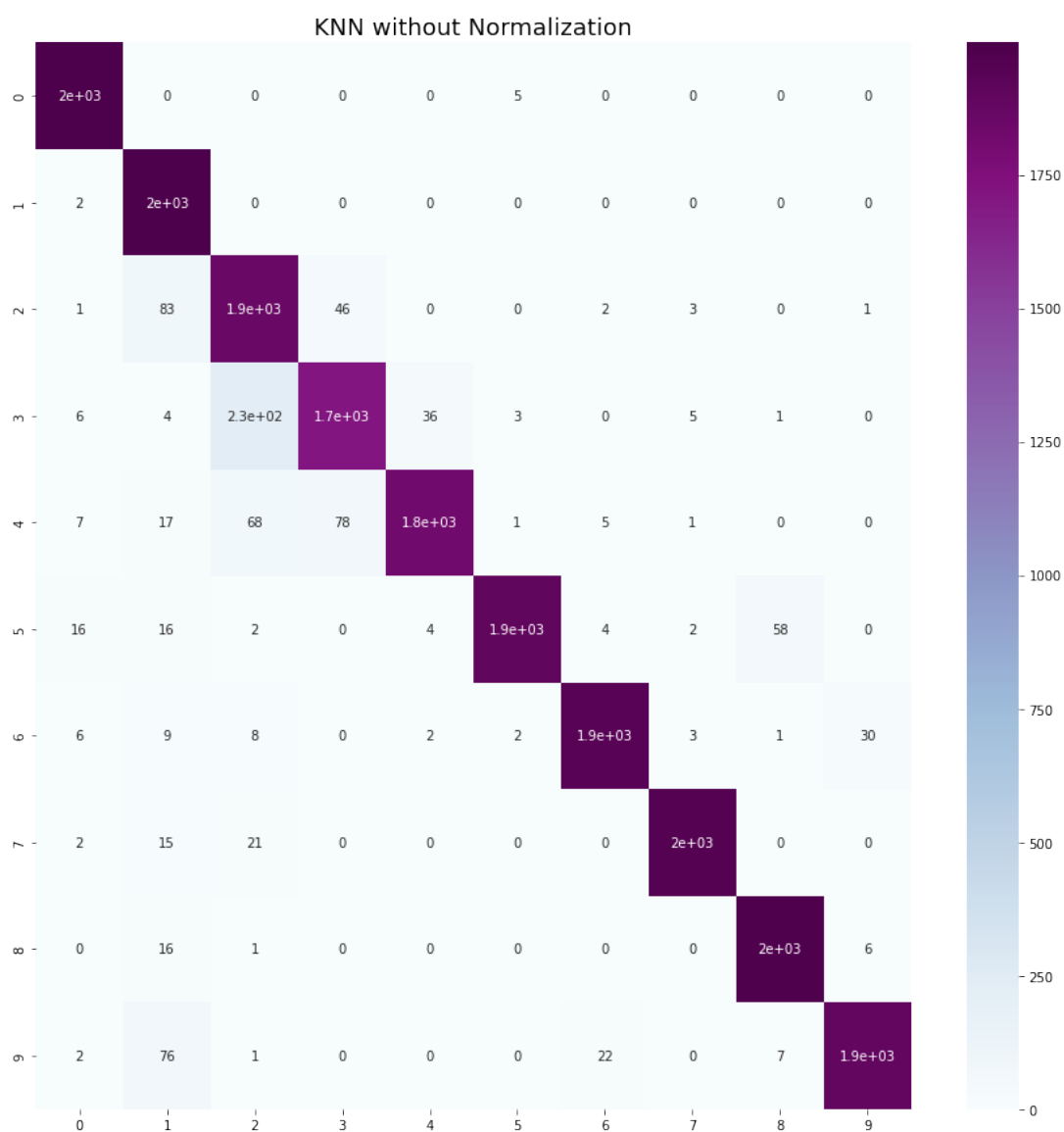
| label | ۰ | ۱ | ۲ | ۳ | ۴ | ۵ | ۶ | ۷ | ۸ | ۹ |
|----------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| f1 score | ۸۳.۹۸ | ۳۷.۹۴ | ۸۲.۸۸ | ۲۸.۸۹ | ۳۳.۹۴ | ۱۰.۹۷ | ۶۳.۹۷ | ۶۹.۹۸ | ۷۷.۹۷ | ۳۰.۹۶ |

جدول ۱: نتایج طبقه بندی KNN با داده‌های نرمالایز نشده

همانطور که از ماتریس آشفستگی و f1 score نیز مشخص است، دقت تمایز بین برخی کلاس‌ها پایین تر است. برای نمونه اعداد ۲ و ۳ فارسی شباهت زیادی به یک دیگر دارند و در مجموعه تست تعدادی از نمونه‌های این دو کلاس اشتباه طبقه بندی شده‌اند. کلاس ۹ و ۱ نیز تعداد اشتباه به نسبت بالایی دارند.



شکل ۳: چند نمونه از تصاویر موجود در دیتاست



شکل ۴: Confusion Matrix طبقه بند KNN با داده‌های نرمالایز نشده

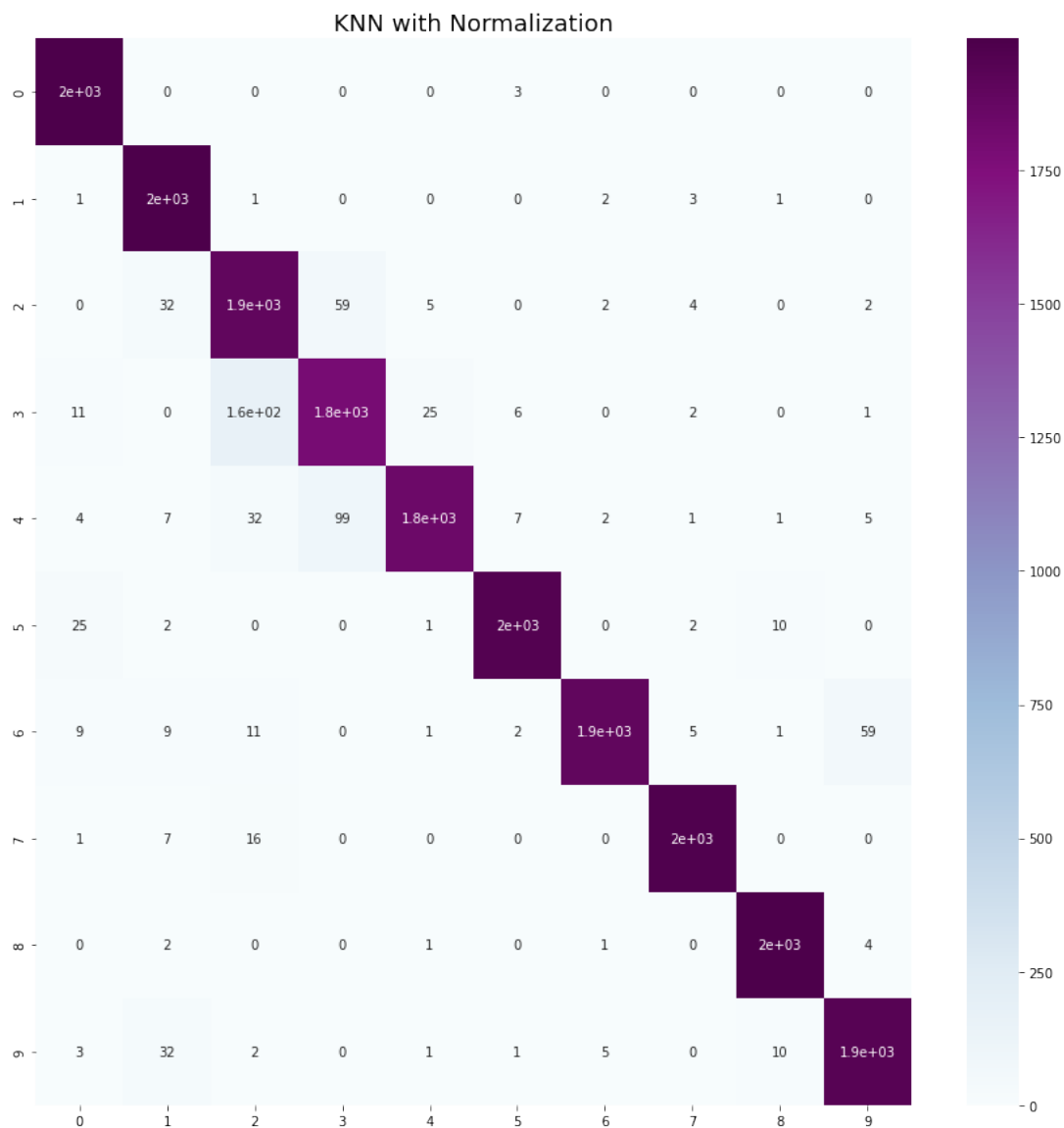
ج

در این قسمت بردار فیچر را با استفاده از رابطه زیر نرمال می‌کنیم:

$$\mathbf{D} = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\}$$

$$\vec{x}_i := \frac{\vec{x}_i}{\sum_{j=1}^d x_i^j}$$

با استفاده از رابطه فوق دو بردار X_{train} و X_{test} نرمال شده و در آرایه‌ی جدیدی ذخیره می‌کنیم و مجدداً عکس‌ها را طبقه‌بندی می‌کنیم. نتایج طبقه‌بندی با داده‌های نرمال شده بر روی مجموعه تست در شکل ۵ نشان داده شده است.



شکل ۵: Confusion Matrix طبقه‌بند KNN با داده‌های نرمال‌ایز شده

از ماتریس آشفتگی و جدول ۲ می‌توان نتیجه گرفت که در مقایسه با نتایج قسمت قبل بهبود اندکی در عملکرد طبقه‌بند به چشم می‌خورد. در تمامی قسمت‌های این تمرین از تعداد همسایه $K = 5$ در طبقه‌بند استفاده شده است.

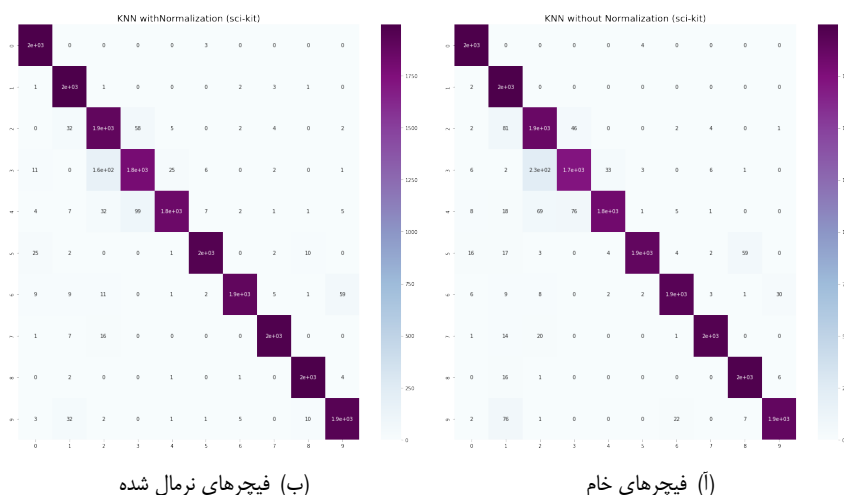
| label | ۰ | ۱ | ۲ | ۳ | ۴ | ۵ | ۶ | ۷ | ۸ | ۹ |
|-------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| F1 Score (Raw Features) | ۸۳.۹۸ | ۳۷.۹۴ | ۸۲.۸۸ | ۲۸.۸۹ | ۳۳.۹۴ | ۱۰.۹۷ | ۶۳.۹۷ | ۶۹.۹۸ | ۷۷.۹۷ | ۳۰.۹۶ |
| F1 Score (Normalized) | ۵۹.۹۸ | ۵۷.۹۷ | ۹۷.۹۱ | ۶۷.۹۰ | ۰۴.۹۵ | ۵۱.۹۸ | ۲۱.۹۷ | ۹۷.۹۸ | ۲۲.۹۹ | ۸۸.۹۶ |

جدول ۲: مقایسه عملکرد طبقه‌بند KNN با داده‌های پردازش نشده و داده‌های نرمال شده

۵

در این قسمت با استفاده از کلاس آماده KNeighborsClassifier از کتابخانه Scikit Learn طبقه‌بندی را برای داده‌های پردازش نشده و پردازش شده انجام می‌دهیم. عملکرد طبقه‌بند با متریک F1 Score و همچنین ماتریس‌های آشفتگی در شکل ۶ و جدول ۳ نشان داده شده‌اند.

برای ساخته نمونه از این کلاس تنها کافی است مقدار n neighbours را در کانستراکتور تنظیم کنیم و سپس با استفاده از تابع fit داده‌های یادگیری و لیبل آن‌ها را به آبجکت تشکیل شده بدهیم. در نهایت با استفاده از دستور predict داده‌های تست را طبقه‌بندی کنیم.



شکل ۶: طبقه‌بندی با کلاس کتابخانه scikit learn

| label | ۰ | ۱ | ۲ | ۳ | ۴ | ۵ | ۶ | ۷ | ۸ | ۹ |
|-------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| F1 Score (Raw Features) | ۸۳.۹۸ | ۴۴.۹۴ | ۸۴.۸۸ | ۵۰.۸۹ | ۳۷.۹۴ | ۰۵.۹۷ | ۶۰.۹۷ | ۶۹.۹۸ | ۷۵.۹۷ | ۳۰.۹۶ |
| F1 Score (Normalized) | ۵۹.۹۸ | ۵۷.۹۷ | ۹۹.۹۱ | ۷۰.۹۰ | ۰۴.۹۵ | ۵۱.۹۸ | ۲۱.۹۷ | ۹۷.۹۸ | ۲۲.۹۹ | ۸۸.۹۶ |

جدول ۳: نتایج طبقه‌بندی با استفاده از کتابخانه scikit learn

۸ سوال هشت (شبیه سازی)

در این تمرین به پیاده سازی روش نان پارامتری پنجره پارزن می پردازیم. داده های موجود در فایل اکسل را با دستورات کتابخانه pandas باز کرده و ستون duration را در یک آرایه numpy ذخیره می کنیم.

برای پیاده سازی طبقه بند یک کلاس به نام parzen پیاده سازی شده است. این کلاس چهار متد برای اجرای الگوریتم طبقه بندی و پیش بینی احتمال نقاط جدید دارد.

- constructor: برای ایجاد یک نمونه از این کلاس لازم است داده های یادگیری، نوع کرنل و اندازه پنجره برای فضاخوانی این تابع ارسال کرد. تعداد داده های یادگیری باید حتماً آرایه numpy باشند.

- gaussian kernel: این تابع با استفاده از رابطه $\phi(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^d h_n^d} \exp \left[-\frac{1}{2} \left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n} \right)^2 \right]$ کرنل گوسی را برای بردار ورودی حساب می کند.

- hypercube kernel: برای تست عملکرد طبقه بند قبل از یادگیری و تست بر روی داده های اصلی با استفاده از کرنل مکعبی تابع prob تست شده است. برای استفاده از این کرنل در کانستراکتور نمونه ساخته شده از این تابع باید آرگومان آخر را "cube" بگذاریم.

- prob: این تابع با استفاده از کرنل مشخص شده در کانستراکتور و رابطه $p_n(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h^d} \phi \left[\frac{\mathbf{x} - \mathbf{x}_i}{h_n} \right]$ احتمال مشاهده هر سطر از بردار ورودی را با توجه به داده های دریافت شده از تابع کانستراکتور به دست می آورد.

- plot dist: این تابع با استفاده از دستورات displot و line plot از کتاب خانه seaborn هیستوگرام داده های یادگیری را به همراه تابع توزیع احتمال تخمین زده شده را رسم می کند.

الف

با استفاده از کلاس پیاده سازی شده با کرنل گوسی و $h = 10$ توزیع تابع را برای داده های تد به دست آمده و نتیجه در شکل ۷ نشان داده شده است.

ب

مانند بخش الف این تمرین با تغییر هایپرپارامتر طبقه بند (همان طول پنجره) یادگیری و تخمین توزیع را انجام می دهیم. نتایج در شکل ۸ نشان داده شده اند.

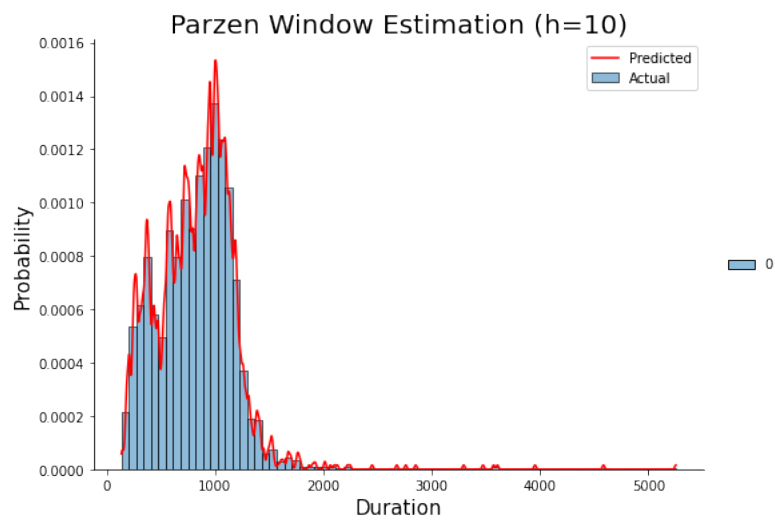
ج

در این قسمت داده ها را به طور تدریجی به تخمین زن می دهیم تا روند تغییر و همگرا شدن به توزیع اصلی را ببینیم. این نمودار در شکل ۹ آورده شده است.

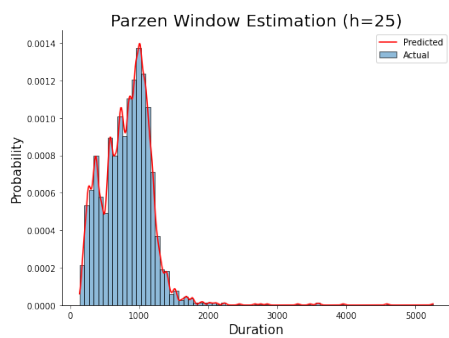
د

در این قسمت با استفاده از KernelDensity از کتابخانه scikit learn نتایج قسمت الف را صحت سنجی می کنیم.

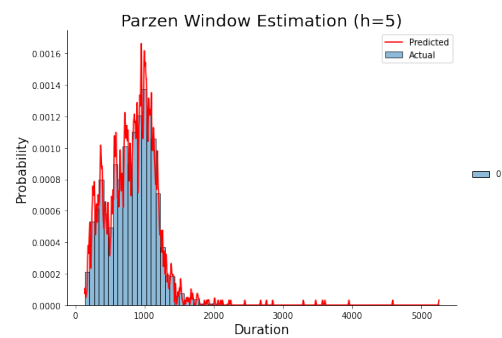
یک آبجکت KernelDensity ساخته و طول پنجره را در آرگومان ورودی (band width) تنظیم می کنیم. سپس داده های یادگیری را با فراخوانی تابع fit به مدل می دهیم. در نهایت می توان با تابع score samples احتمال داده های تست را حساب کرد. این تابع لگاریتم طبیعی احتمال را در خروجی باز می گرداند، پس لازم است آن ها را با تابع exp از کتابخانه numpy به احتمال تبدیل کرد.



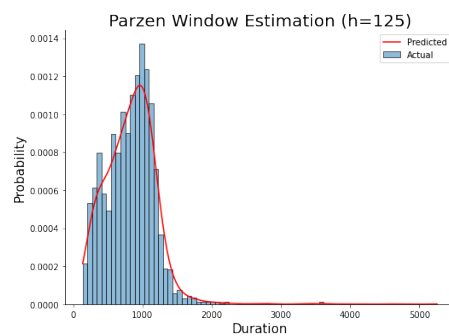
شکل ۷: توزیع تخمین زده شده با $k = 10$



(ب) $h_n = 25$

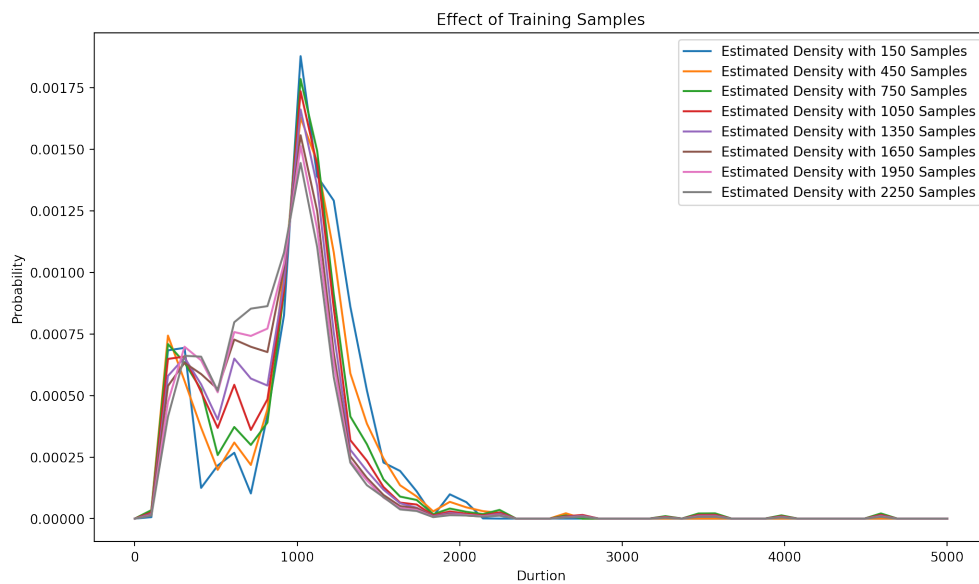


(ا) $h_n = 5$

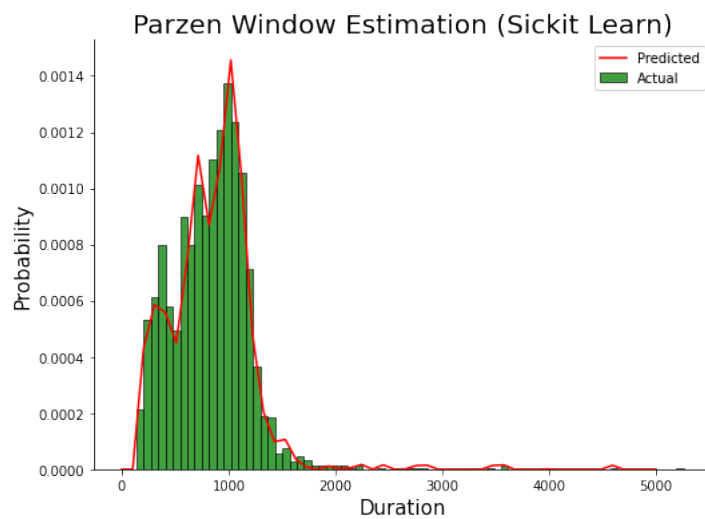


(ج) $h_n = 125$

شکل ۸: توزیع‌های تخمین زده شده با طول پنجره متفاوت



شکل ۹: روند همگرایی تخمین نان پارامتری پنجره پارزن



شکل ۱۰: تخمین نان پارامتری پارزن با استفاده از کتابخانه scikit learn