

به نام خدا

الگوریتم k نزدیکترین همسایه

فاطمه دلدار- ملیحه رضایی

الگوریتم k نزدیکترین همسایه

الگوریتم k نزدیکترین همسایه یک الگوریتم تعلیم با سرپرستی است. در حالت کلی از این الگوریتم به دو منظور استفاده می‌شود: برای تخمین تابع چگالی توزیع داده‌های تعلیم و برای طبقه‌بندی داده‌های تست بر اساس الگوهای تعلیم.

تخمین چگالی توزیع داده‌ها با استفاده از الگوریتم K_n نزدیک‌ترین همسایه

برای تخمین $p(x)$ از روی n نمونه‌ی تعلیم توسط الگوریتم k نزدیک‌ترین همسایه می‌توانیم یک سلول به مرکزیت x ایجاد کرده و اجازه دهیم این شعاع این سلول تا حدی گسترش پیدا کند که k_n نمونه‌ی تعلیم را در بر گیرد. این نمونه‌ها k_n نزدیک‌ترین همسایه‌های x هستند.

• در حالت کلی k را به صورت k_n در نظر می‌گیریم که k_n تابعی تعریف شده از n است.

اگر چگالی نقاط تعلیم اطراف x زیاد باشد سلول کوچک می‌شود و بنابراین نتیجه‌ی به دست آمده نتیجه‌ی بهتری است و در صورتی که چگالی نقاط تعلیم اطراف x کم باشد سلول بزرگ می‌شود. در حالت کلی چگالی توزیع به ازای هر نقطه‌ی x توسط رابطه‌ی زیر محاسبه می‌شود:

$$p_n(x) = \frac{k_n / n}{V_n}$$

اگر با رشد n ، k_n نیز افزایش پیدا کند به طوری که با رفتن n به سمت بی-نهایت k_n نیز به بی‌نهایت میل کند، آنگاه می‌توان مطمئن بود که k_n/n یک تخمین خوب از این احتمال است که یک نقطه در یک سلول به حجم V_n قرار بگیرد. بنابراین دو شرط زیر شروط لازم و کافی برای این است که $p_n(x)$ در نهایت به $p(x)$ همگرا شود:

$$\lim_{n \rightarrow \infty} k_n = \infty, \quad \lim_{n \rightarrow \infty} k_n / n = 0$$

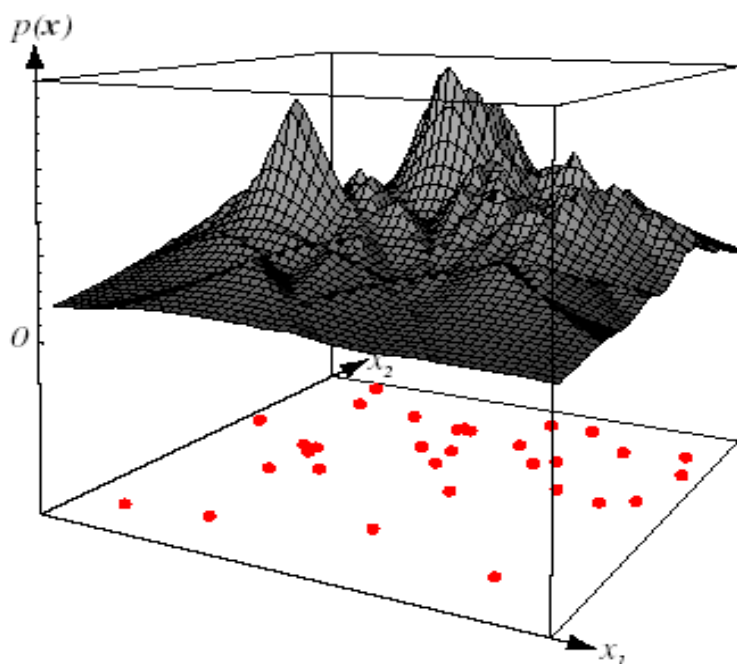
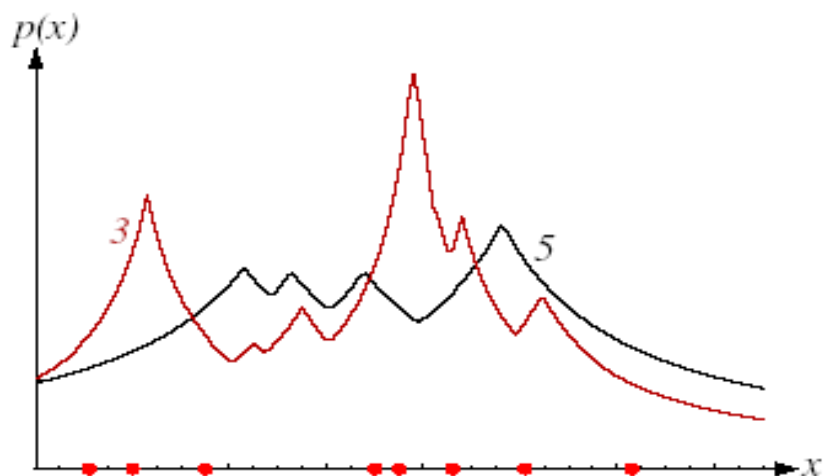
به عنوان مثال اگر k_n را \sqrt{n} در نظر بگیریم داریم:

$$k_n = \sqrt{n}$$

$$V_n \approx \frac{1}{\sqrt{n} p(x)}, \quad V_n \approx \frac{V_1}{\sqrt{n}} \rightarrow 0 \text{ as } n \rightarrow \infty$$

و بنابراین شروط فوق برای k_n داده شده برقرار بوده و با افزایش n به سمت بی‌نهایت $p_n(x)$ به $p(x)$ میل می‌کند. شکل‌های زیر تخمینی از تابع چگالی

توزیع نمونه‌های داده شده را به ترتیب برای نمونه‌های یک‌بعدی و دوبعدی نشان می‌دهند.



تخمین احتمال‌های posteriori

تکنیک بیان شده در قسمت قبل می‌تواند برای تخمین احتمال‌های posteriori، $P(w_i|x)$ از یک مجموعه از n نمونه‌ی برچسب‌گذاری شده استفاده شود. فرض می‌کنیم که یک سلول به حجم v حول x رسم کرده‌ایم که k نمونه‌ی

ورودي را در بر گرفته است و k_i تا از این نمونه‌ها در کلاس w_i قرار دارند. بنابراین یک تخمین برای احتمال $p(x, w_i)$ به صورت

$$p_n(\mathbf{x}, \omega_i) = \frac{k_i / n}{V}$$

است و بنابراین داریم:

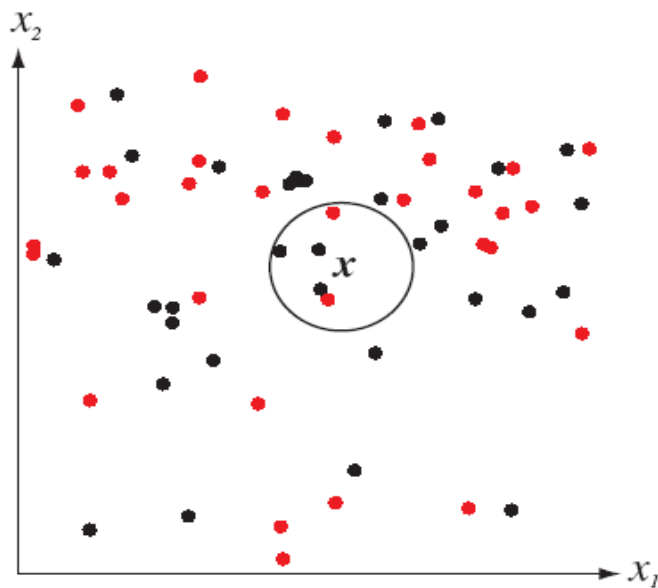
$$p_n(\omega_i | \mathbf{x}) = \frac{p_n(\mathbf{x}, \omega_i)}{\sum_{j=1}^c p_n(\mathbf{x}, \omega_j)} = \frac{\frac{k_i / n}{V}}{\sum_{j=1}^c \frac{k_j / n}{V}} = \frac{k_i}{k}$$

یعنی برای اینکه نرخ خطا حداقل باشد، نمونه‌ی تست را در دسته‌ای قرار می‌دهیم که بیشترین تکرار را در سلول داشته باشد. اگر تعداد نمونه‌ها به اندازه‌ی کافی زیاد باشد و سلول به اندازه‌ی کافی کوچک باشد، این انتخاب از لحاظ کارایی بهترین انتخاب است.

قانون نزدیک‌ترین همسایه

در تعریف قانون نزدیک‌ترین همسایه k را مساوی یک در نظر گرفته شده است. فرض می‌کنیم $D^n = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ مجموعه‌ای از n الگوی ورودی باشد و $\mathbf{x}' \in D^n$ نزدیک‌ترین الگوی ورودی به نقطه‌ی تست \mathbf{x} باشد. قانون نزدیک‌ترین همسایه برای طبقه‌بندی \mathbf{x} آن را در کلاسی مشابه با کلاس \mathbf{x}' قرار می‌دهد. قانون نزدیک‌ترین همسایه یک روال sub-optimal است یعنی نرخ خطای آن معمولاً بیشتر از حداقل نرخ خطای ممکن یعنی نرخ خطای الگوریتم بیز است. ولی ثابت می‌شود که در صورت استفاده از تعداد نامحدودی از الگوهای ورودی نرخ خطا در بدترین حالت بیشتر از دو برابر نرخ خطای الگوریتم بیز نخواهد شد.

قانون k نزدیک‌ترین همسایه گسترشی از قانون نزدیک‌ترین همسایه است و همان‌طور که واضح است این قانون \mathbf{x} را در دسته‌ای طبقه‌بندی می‌کند که بیشترین تکرار را در بین k نزدیک‌ترین همسایه‌ی \mathbf{x} دارد. به عنوان مثال در شکل زیر \mathbf{x} توسط این قانون در کلاس سیاه قرار می‌گیرد.



پیچیدگی محاسباتی در الگوریتم نزدیک‌ترین همسایه

برای محاسبه‌ی پیچیدگی الگوریتم k نزدیک‌ترین همسایه فرض می‌کنیم که n نمونه‌ی تعلیم d بعدی داریم و می‌خواهیم نزدیک‌ترین آنها به نقطه‌ی تست x ($k=1$) را به دست آوریم. در ساده‌ترین روش برای هر نقطه‌ی تعلیم فاصله‌ی آن تا نقطه‌ی تست را محاسبه کرده و نزدیک‌ترین آنها به نقطه تست را برمی‌گردانیم. محاسبه‌ی فاصله‌ی اقلیدس برای هر نقطه در $O(d)$ انجام می‌شود و چون n نقطه‌ی تعلیم ذخیره شده داریم این محاسبه $O(nd)$ زمان می‌برد و در نهایت انتخاب کوتاه‌ترین فاصله نیز در زمان $O(n)$ انجام می‌شود. بنابراین پیچیدگی زمانی الگوریتم در حالت کلی $O(nd)$ می‌باشد. سه تکنیک الگوریتمی برای کاهش پیچیدگی محاسباتی الگوریتم نزدیک‌ترین همسایه وجود دارد:

- محاسبه‌ی فاصله‌های جزئی
- Prestructuring
- ویرایش الگوهای ذخیره شده

محاسبه‌ی فاصله‌های جزئی: در این روش فاصله‌ی بین دو نقطه را ابتدا به صورت کامل محاسبه نمی‌کنیم بلکه آن را به صورت تدریجی و با افزودن ابعاد محاسبه می‌کنیم. اگر فضا را d بعدی در نظر بگیریم، فاصله‌ی جزئی بر مبنای r بعد انتخاب شده از این d بعد به صورت زیر محاسبه می‌شود:

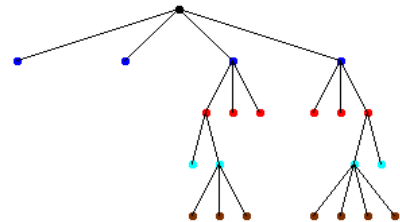
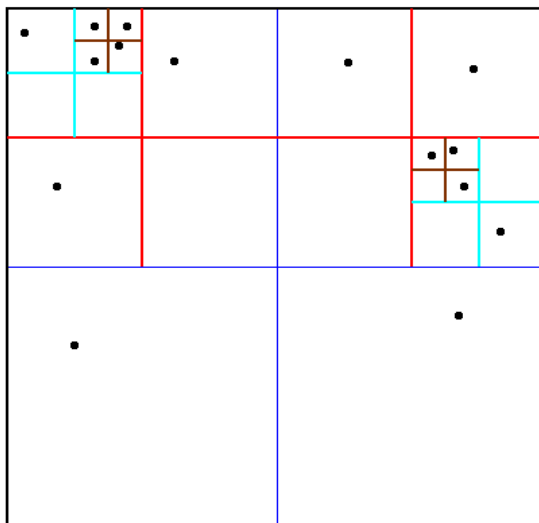
$$D_r(\mathbf{a}, \mathbf{b}) = \left(\sum_{k=1}^r (a_k - b_k)^2 \right)^{1/2}, \quad r < d$$

می‌توانیم محاسبه‌ی فاصله‌ی جزئی را تا زمانی ادامه دهیم که بزرگتر از فاصله‌ی فعلی نقطه‌ی تست از نزدیک‌ترین الگوی فعلی شود. در این حالت

مطمئن می‌شویم که این الگو نزدیکترین الگو به نقطه‌ی تست مورد نظر نیست و محاسبه‌ی فاصله را متوقف می‌کنیم.

Prestructuring: در این روش بر روی الگوهای تعلیم مورد جستجو شکلی از درخت جستجو را ایجاد می‌کنیم. در حین فرآیند طبقه‌بندی کردن فقط فاصله‌ی نقطه‌ی تست را تا یک یا تعداد کمی از الگوهای ذخیره شده در درخت محاسبه می‌کنیم. به این معنی که یک الگو که نزدیک‌تر از بقیه‌ی الگوها به نقطه‌ی تست به نظر می‌رسد را انتخاب کرده و به صورت بازگشتی فقط الگوهای متصل به آن را در درخت جستجو می‌کنیم. اگر درخت به درستی ساخته شده باشد تعداد کل الگوهای موردنیاز برای جستجو کاهش پیدا می‌کند.

یکی از این ساختارهای درختی QuadTree است. در این درخت هر نود میانی چهار فرزند دارد. این درخت اغلب به این منظور استفاده می‌شود که فضاهای دوبعدی را به صورت بازگشتی به چهار ناحیه تقسیم کند. این نواحی می‌توانند مربع، مستطیل یا هر شکل دلخواه دیگری داشته باشد. شکل زیر ساخته شدن یک QuadTree را نشان می‌دهد.



این روش به هیچ عنوان تضمین نمی‌کند که نزدیکترین الگو را پیدا کند. در جستجوی الگوهای تعلیم از طریق این درخت در هر مرحله $\frac{3}{4}$ فضای جستجو حذف می‌شود و فقط $\frac{1}{4}$ ای باقی می‌ماند که فاصله‌ی نقطه‌ی تست تا مرکز آن ناحیه کمتر از سه ناحیه‌ی دیگر است. بنابراین ممکن است نزدیکترین الگو به نقطه‌ی تست در یکی از نواحی دیگر قرار داشته باشد (این مشکل بیشتر در مواقعی پیش می‌آید که نقطه‌ی تست نزدیک به مرز ناحیه قرار دارد) و بنابراین با حذف آن ناحیه نزدیکترین الگو به نقطه‌ی تست را نیز از دست

بدهیم. با وجود اینکه این روش روش تضمین‌شده‌ای نیست ولی برای کاهش فضای جستجو روش مناسبی است.

ویرایش الگوهای ذخیره‌شده: روش سوم برای کاهش پیچیدگی جستجوی نزدیکترین همسایه حذف الگوهای غیرمفید در حین تعلیم است. یک روش برای کاهش پیچیدگی فضایی $O(n)$ این است که الگوهایی که همگی نقاط تعلیم اطراف آنها در کلاسی مشابه با خود الگو قرار دارند را حذف کنیم. با این کار مرزهای تعلیم تغییر نمی‌کند در حالی که فضای جستجو کاهش پیدا می‌کند. یک الگوریتم ساده برای این کار به صورت زیر است:

Algorithm. (Nearest-Neighbor Editing)

1. **begin initialize** $j \leftarrow 0$, $D \leftarrow$ data set, $n \leftarrow$ # prototypes
2. construct the full Voronoi diagram of D
3. **do** $j \leftarrow j + 1$; for each prototype \mathbf{x}'_j
4. find the Voronoi neighbors of \mathbf{x}'_j
5. if any neighbor is not from the same class as , **then** mark \mathbf{x}'_j
6. **until** $j = n$
7. discard all points that are not marked
8. construct the Voronoi diagram of the remaining (marked) prototypes
9. **end**

از آنجایی که الگوهایی که در مرز تصمیم قرار دارند حداقل یکی از همسایه‌هایشان از دسته‌ای متفاوت هستند، بنابراین با اجرای الگوریتم بالا حذف نمی‌شوند و مرز تصمیم تغییری پیدا نمی‌کند. این الگوریتم تضمین نمی‌دهد که حداقل مجموعه الگوهای تعلیم مفید باقی بمانند. یکی از مشکلات چنین سیستم نزدیک‌ترین همسایه‌ای این است که نمی‌توان بعداً الگوی تعلیمی به این الگوها اضافه کرد چون برای اجرای این الگوریتم نیاز است که همگی الگوهای تعلیم وجود داشته باشند.

می‌توان این سه روش کاهش پیچیدگی را با هم ادغام کرد، یعنی ابتدا الگوها را ویرایش کرد، بعد بر روی الگوهای باقیمانده یک درخت جستجو ایجاد کرد و در نهایت در هنگام انجام عمل طبقه‌بندی فاصله‌های جزئی را محاسبه کرد.

پیاده‌سازی الگوریتم k نزدیکترین همسایه

در پیاده‌سازی الگوریتم k نزدیک‌ترین همسایه در نرم‌افزار مطلب کلاسهای زیر به کار رفته است:

کلاس CalDistance.m

این کلاس فاصله‌ی بین همه‌ی نقاط تعلیم از همه‌ی نقاط تست را بر اساس یکی از معیارهای فاصله محاسبه کرده و ماتریس فاصله‌ها را برمی‌گرداند. ورودی این کلاس train (نقاط تعلیم)، test (نقاط تست) و distance (معیار فاصله-ی موردنظر) است که این معیار فاصله می‌تواند یکی از موارد زیر باشد:

- euclidean: فاصله‌ی اقلیدس
- cityblock: مجموع قدرمطلق فاصله‌ها
- cosine: یک منهای کسینوس زاویه‌ی محصور بین دو نقطه (نقاط را به عنوان بردار در نظر می‌گیریم)
- correlation: یک منهای کوررولیشن بین دو نقطه
- hamming: این فاصله فقط برای داده‌های باینری مناسب است و درصدی از بیت‌ها را که با هم متفاوت هستند نشان می‌دهد.

کلاس KNNClassification.m

این کلاس کلاس اصلی برنامه می‌باشد و k نزدیکترین همسایه به هر یک از نقاط تست را محاسبه کرده و بر اساس آنها مشخص می‌کند که هر یک از نقاط تست در چه دسته یا کلاسی قرار می‌گیرند.

ورودی‌های این کلاس train (نقاط تعلیم)، test (نقاط تست)، group (ماتریسی ستونی که تعداد سطرهای آن مساوی با تعداد سطرهای train است و مشخص می‌کند که هر یک از الگوهای تعلیم در چه کلاسی قرار می‌گیرند)، k که تعداد نزدیکترین همسایه مورد نظر برای الگوریتم KNN را مشخص می‌کند، dist که معیار فاصله‌ی مورد نظر برای محاسبه‌ی فاصله‌ی بین نقاط تعلیم و نقاط تست را بیان می‌کند و rule که می‌تواند یکی از موارد زیر باشد:

- consensus: وقتی این گزینه به عنوان rule انتخاب شود به این معنی است فقط در صورتی که همه‌ی k نزدیکترین همسایه‌ی یک نقطه‌ی تست به یک کلاس مشابه تعلق داشته باشند آن نقطه‌ی تست نیز در همان کلاس طبقه‌بندی می‌شود ولی اگر حتی یکی از نقاط k نزدیکترین همسایه‌ی یک نقطه‌ی تست در کلاسی متفاوت با بقیه قرار داشته باشد، آن نقطه‌ی تست قابل طبقه‌بندی با این الگوریتم نیست. در سایر موارد به غیر از مواقعی که این گزینه به عنوان rule قرار می‌گیرد از قانون اکثریت استفاده می‌شود یعنی یک نقطه‌ی تست در کلاسی قرار می‌گیرد که تعداد نقاطی از k نزدیکترین همسایه‌ی آن که در آن کلاس قرار دارند نسبت به سایر کلاس‌ها بیشتر باشد. در مواقعی که گره ایجاد می‌شود یعنی بیشتر نقاط در یک کلاس قرار نمی‌گیرند بلکه نقاط قرار گرفته در بیش از یک کلاس با هم برابر هستند از یکی از سه قانون زیر برای شکستن این گره استفاده می‌شود:

- nearest: یعنی از بین گروه‌های به هم گره خورده کلاس نزدیکترین همسایه به عنوان کلاس خروجی انتخاب می‌شود.
- farthest: یعنی از بین گروه‌های به هم گره خورده کلاس دورترین همسایه به عنوان کلاس خروجی انتخاب می‌شود.
- random: یعنی برای شکستن گره به صورت تصادفی یکی از همسایه‌هایی که به هم گره خورده‌اند انتخاب شده و کلاس آن همسایه به عنوان کلاس خروجی برگردانده می‌شود.

در کلاس KNNClassification، ابتدا توسط تابع CalDistance فاصله‌ی بین نقاط تعلیم و تست محاسبه شده و در ماتریس d قرار می‌گیرد. سپس این ماتریس را مرتب می‌کنیم یعنی برای هر نقطه‌ی تست نقاط تعلیم را به صورت صعودی بر حسب فاصله از آن نقطه‌ی تست در سطر مربوط به آن قرار داده و سپس k نقطه‌ی تعلیم اول هر یک از این سطرها را به عنوان k نزدیکترین همسایه‌ی هر یک از نقاط تست جدا می‌کنیم. پس از آن در یک حلقه به ازای هر یک از نقاط تست محاسبه می‌کنیم که هر یک از k نزدیکترین همسایه‌ی آن در چه کلاسی قرار دارند و کلاسی که ماکزیمم تعداد همسایه‌ها در آن قرار دارند را به عنوان کلاس خروجی انتخاب می‌کنیم. در صورت به وجود آمدن گره نیز با توجه به نوع rule گره را باز می‌کنیم.

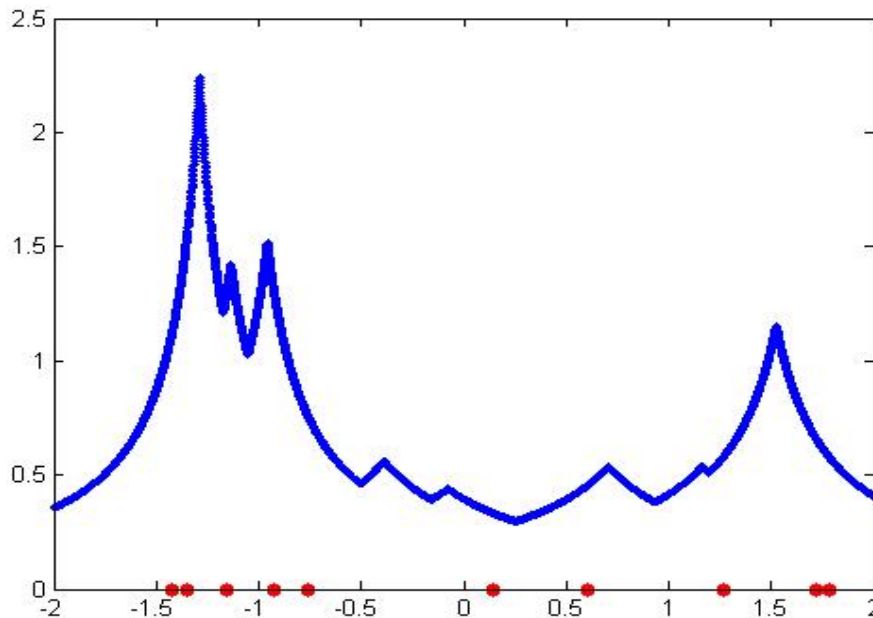
کلاس classification.m

این کلاس کلاس اجرایی برنامه است که در آن الگوهای تعلیم و گروه‌های آنها، الگوهای تست، مقدار k ، معیار فاصله و rule را انتخاب کرده و با توجه به آنها کلاس KNNClassification را فراخوانی می‌کنیم.

کلاس KNN1.m

این کلاس برای تخمین تابع چگالی توزیع نقاط تعلیم با توجه به الگوریتم k نزدیکترین همسایه استفاده می‌شود. به این منظور نقاط تست را از ابتدا تا انتهای بازه‌ای که نقاط تعلیم در آن قرار دارند با فاصله کم از هم در نظر گرفته و برای هر کدام از آنها تابع CalDistance را فراخوانی می‌کنیم تا فاصله‌ی آن از تمام نقاط تعلیم را به دست آوریم. سپس این فاصله‌ها را مرتب می‌کنیم. از آنجایی که در استفاده از الگوریتم k نزدیکترین همسایه برای تخمین چگالی توزیع الگوهای ورودی در نقطه‌ی x به مرکزیت آن نقطه سلولی ایجاد می‌کنیم و اجازه می‌دهیم شعاع این سلول تا حدی گسترش پیدا کند که k نزدیکترین همسایه‌ی نقطه‌ی x را در بر گیرد، بنابراین در لیست مرتب شده k امین عدد را به عنوان شعاع در نظر می‌گیریم و در نهایت با استفاده از رابطه‌ی $p=k/n.v$ مقدار تابع چگالی در نقطه‌ی مورد نظر را به دست می‌آوریم. در کلاس KNN1 داده‌ها به صورت یک بعدی در نظر گرفته شده‌اند و بنابراین به جای v طول خط یا همان k را در نظر می‌گیریم. در شکل زیر نمونه‌ای از

خروجي اين کلاس نشان داده شده است. (الگوهاي ورودی با نقاط قرمز مشخص شده‌اند)



کلاس KNN2.m

این کلاس کاملاً مشابه کلاس KNN1 است با این تفاوت که در این کلاس داده‌ها به صورت دو بعدی در نظر گرفته شده‌اند و به همین دلیل به جای v مساحت دایره‌ای به مرکزیت x و شعاع k را محاسبه می‌کنیم. در شکل زیر نمونه‌ای از خروجی این کلاس آمده است.

