

گزارش تمرین پیاده سازی دوم درس هوش محاسباتی

بخش اجباری (FCM.py)

شرح مختصری راجع به نحوه ی پیاده سازی الگوریتم و برخی موارد جزئی تر مانند تعیین شرط خاتمه و یا مقداردهی های اولیه، ارائه دهید.

پس از آنکه فایل سمپل تست مورد نظر از کاربر دریافت شد، برنامه ابتدا نقاط را پلات می کنیم (اگر داده ها دوبعدی باشند) و سپس به ازای $c=2$ تا $c=10$ ، الگوریتم را اجرا می کنیم.

روند اجرای الگوریتم بدین صورت بوده است که ابتدا مراکز خوشه ها (V_i ها) را به صورت رندوم جنریت می کنیم.

سپس در یک حلقه، توابع تعلق (u_{ik} ها) را حساب کرده و مجدد مراکز خوشه ها را از روی توابع تعلق حساب می کنیم.

در نهایت و پس از خاتمه حلقه (حلقه 10 دفعه اجرا می شود، این عدد تجربی است و برای دسته بندی داده ها مناسب بوده است)، آرایه توابع تعلق مربوط به آن C را به یک آرایه به نام $c_memberships$ (آرایه تمام آرایه های توابع تعلق) می افزاییم.

در ادامه، سراغ تعیین بهترین C می رویم. از فرمول زیر برای محاسبه آنتروپی هر C به خصوص استفاده می کنیم:

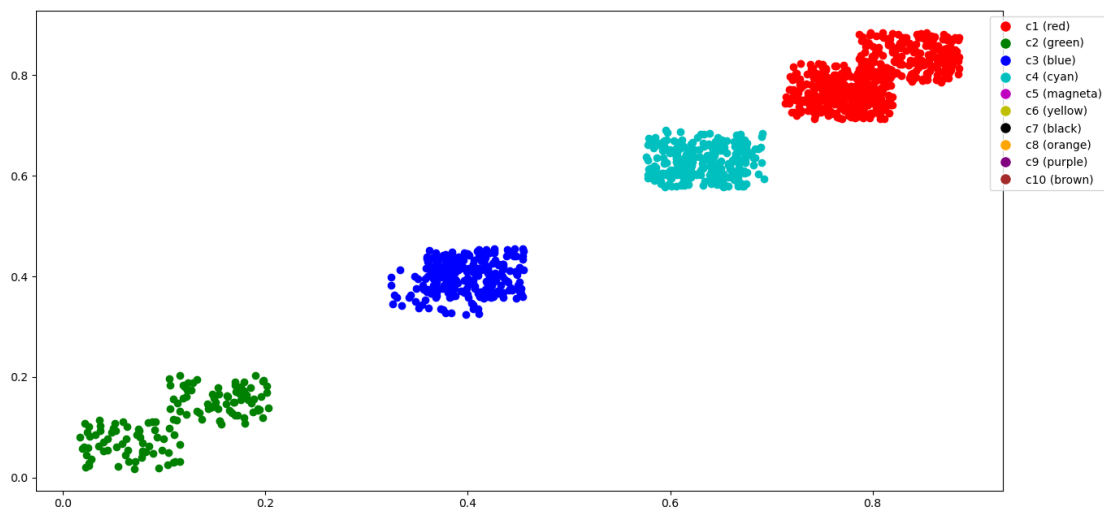
$$Entropy_c = - \sum_i \sum_k \frac{u_{ik} \ln(u_{ik})}{\ln(c)}$$

پس از اینکه بهترین C تعیین شد، نقاط را دسته بندی می کنیم. دسته بندی بدین صورت است که هر داده به خوشه ای نگاشت می شود که تابع تعلق آن ماکزیمم باشد. نهایتاً، در کیس داده های دوبعدی، این داده ها رنگ خوشه مورد نظر خود را می گیرند و پلات می شوند.

پارامتر m را چه مقداری در نظر گرفتید ؟ مقدار این پارامتر چه تأثیری بر نتایج دارد؟

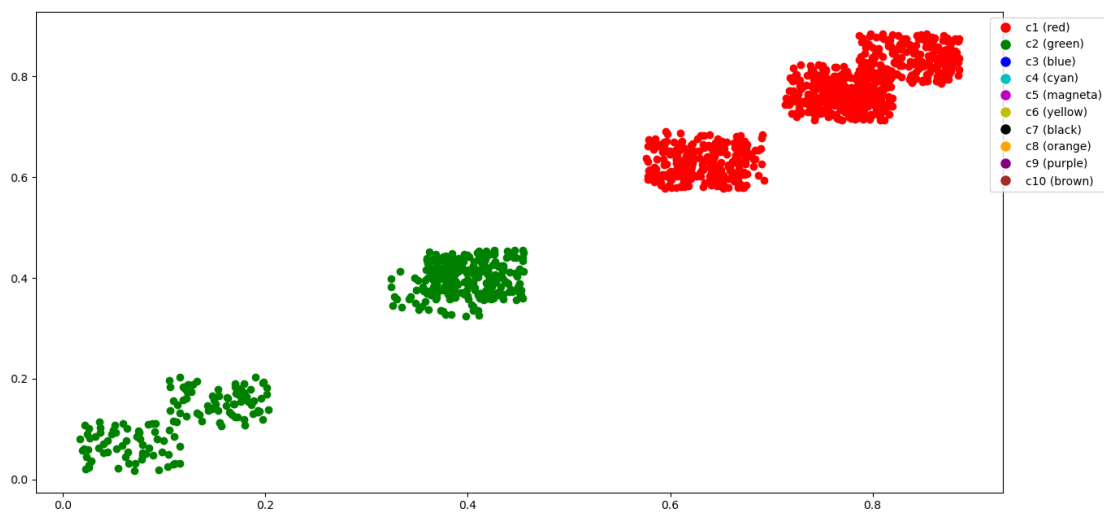
مقدار 3 در نظر گرفته شده است. این مقدار به صورت تجربی به دست آمده است. با افزایش m ، شاهد کم شدن تعداد خوشه ها می شویم. به عنوان مثال:

به ازای $m = 3$ داریم:



sample2. $m = 3$. $c = 4$

به ازای $m = 4$ داریم:



sample2. $m = 4$. $c = 2$

رابطه تابع هزینه در الگوریتم FCM را بنویسید و ارتباط آن را با پارامتر c شرح دهید

تابع هزینه در الگوریتم FCM به صورت زیر می باشد:

$$\sum_i \sum_k u_{ik}^m ||x_k - c_i||^2$$

به طور کلی، این تابع هزینه فاصله داده ها از مراکز خوشه ها را در u_{ik}^m ضرب کرده و از حاصل جمع می گیرد.

ارتباط این پارامتر با c بدین شکل است که با افزایش مکرر c ، چون احتمال نزدیک بودن یک داده به یکی از مراکز خوشه بیشتر می شود، در نتیجه این سیگما نیز مقدار کمتری پیدا می کند. در حالتی که تعداد خوشه ها بسیار زیاد باشد (مثلا برابر تعداد داده ها باشد) آنگاه هر داده مرکز خوشه تکی خود می شود؛ در نتیجه این هزینه برابر صفر می شود.

شایان ذکر است که برای ما افزایش تعداد خوشه ها لزوما امری مثبت نیست؛ و بنابراین از معیارهای دیگر مانند آنتروپی داده ها برای تعیین بهترین c استفاده می کنیم.

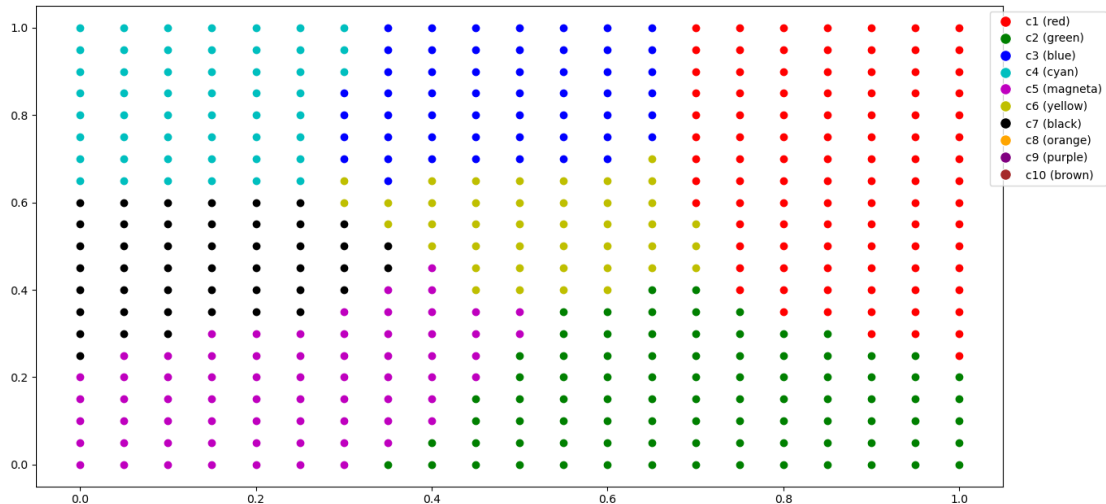
بخش امتیازی (FkNN.py)

در این بخش نیز، پس از آنکه از بخش قبل بهترین C را پیدا کرده ایم، می‌آییم و تعدادی داده تست (test data) جنریت می‌کنیم. (این داده ها به صورت evenly spaced و منظم، کل فضا $1*1$ مورد نظر را می‌پوشانند)

سپس می‌آییم و برای هر کدام از داده‌های تست، نزدیک ترین k نقطه داده‌های اصلی مان به این داده را پیدا می‌کنیم. (این کار با k بار پیمایش روی داده‌های اصلی صورت می‌گیرد)

پس از آنکه برای هر داده تست، نزدیک ترین k داده به آن را پیدا کردیم، می‌آییم و تابع تعلق داده تست را نسبت به خوشه‌هایمان (که در بخش قبل یافتیم) محاسبه می‌کنیم.

در آخر، بسته به اینکه هر داده به کدام خوشه از همه بیشتر تعلق دارد، رنگ آن را به رنگ خوشه متناظرش در می‌آوریم و به ازای همه داده‌های تست، پلات می‌کنیم. نمونه ای از این مرزبندی آورده شده است:



sample3. c = 7