



دانشگاه صنعتی امیرکبیر  
دانشکده مهندسی کامپیوتر

مبانی هوش محاسباتی  
تمرین پیاده‌سازی فازی  
(Fuzzy C-Means Clustering)

استاد درس: دکتر عبادزاده  
بهار ۱۴۰۰

توی درس با منطق فازی آشنا شدید و مباحث تئوری اش رو خوندید. ایده‌های کلی قضیه این بود که بیایم مجموعه‌هایی تعریف کنیم که عناصر مختلف بتونن به درجات مختلف عضو اون مجموعه‌ها باشن، گزاره‌هایی تعریف کنیم که بتونن به طور نسبی درست یا غلط باشن، بتونیم پلی ایجاد کنیم بین قوانین منطقی و محاسبات ریاضی به طوری که این قوانین بتونن در درجات مختلفی صادق باشن و مواردی از این قبیل.

بعد از جلو رفتن مباحث تئوری این حوزه، خیلی‌ها شروع کردن این ایده‌ها رو بهرن توی سایر بخش‌های علوم کامپیوتر و به الگوریتم‌ها و مسائل مختلف، از این منظر نگاه کنن. حاصلش همیشه حل مسائل Classification به کمک الگوریتم‌های مبتنی بر منطق فازی، تعریف شدن شبکه‌های عصبی-فازی که بر طبق قوانین منطقی فازی کار می‌کنن و...

توی این تمرین پیاده‌سازی قراره که سراغ الگوریتم خوشه‌بندی K-Means<sup>1</sup> بریم و نسخه‌ی فازی اون یعنی Fuzzy C-Means رو پیاده‌سازی کنیم. یکی از بهبودهایی که این الگوریتم نسبت به حالت غیرفازی اش داره اینه که داده‌ها فورس نمی‌شن که فقط به یک خوشه‌ی خاص تعلق پیدا کنن بلکه می‌تونن به درجات مختلف به خوشه‌ها تعلق پیدا کنن. این مورد برای نقاطی که در حالت مرزی بین خوشه‌ها قرار دارن، اتفاق می‌افته.

---

<sup>1</sup> الگوریتم K-Means رو می‌تونید از [این لینک](#) مرور کنید.

## الگوریتم Fuzzy C-Means

در قدم اول، توی این الگوریتم، مثل K-Means تعداد خوشه‌هایی که می‌خوایم رو مشخص می‌کنیم. بعد به اون تعداد، مرکز خوشه (Centroid) اولیه به صورت رندوم تولید می‌کنیم. در ادامه، توی یک حلقه باید دوتا کار رو انجام بدیم:

۱- پیدا کردن اینکه هر داده به کدام خوشه (یا خوشه‌ها) تعلق داره.

۲- آپدیت کردن مرکز خوشه‌ها براساس داده‌های متعلق بهشون.

برای کار اول، توی C-Means به هر خوشه به چشم یک مجموعه فازی نگاه میشه. در نتیجه، هر داده به تمام خوشه‌ها تعلق داره ولی به اندازه‌های مختلف. میزان تعلق داده‌ی  $k$  ام به خوشه‌ی  $i$  ام از رابطه‌ی زیر محاسبه میشه:

$$u_{ik} = \frac{1}{\sum_{j=1}^c \left( \frac{\|X_k - V_i\|}{\|X_k - V_j\|} \right)^{\frac{2}{m-1}}}$$

$X_k$  داده‌ی  $k$  ام،  $c$  تعداد خوشه‌ها،  $V_i$  مرکز خوشه‌ی  $i$  ام و  $m$  یک پارامتره که باید برای الگوریتم مشخص کنیم. کاری که رابطه‌ی بالا انجام میده اینه که میزان تعلق داده به یک خوشه رو با توجه به نزدیکی به مرکز اون خوشه و مقایسه‌اش با نزدیکی به سایر مراکز خوشه‌ها حساب می‌کنه.

برای کار دوم، باید میانگین نقاطی که به خوشه تعلق دارن رو حساب کنیم و اون رو به عنوان مرکز خوشه‌ی جدید در نظر بگیریم. از اونجا که همه‌ی نقاط عملاً عضو تمامی خوشه‌ها هستن، پس باید میانگین وزن‌دار بگیریم.

$$V_i = \frac{\sum_{k=1}^N u_{ik}^m X_k}{\sum_{k=1}^N u_{ik}^m}$$

اینطوری، اون داده‌هایی که تعلق بیشتری دارن به یک خوشه، بیشتر نقش دارن توی تعیین مرکز اون خوشه؛ که منطقی هم هست.

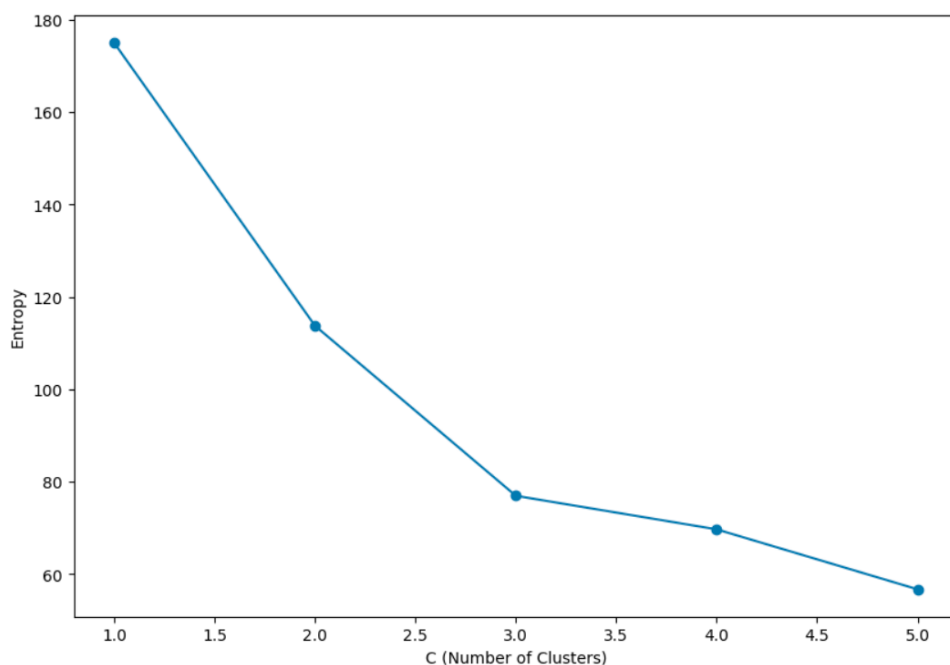
پس بدین شکل این دوتا کار رو توی یک حلقه باید انجام بدیم. این حلقه رو مثلاً به ازای ۱۰۰ بار اجرا می‌کنیم که مطمئن بشیم خوشه‌ها به ثبات رسیدن و دیگه تغییر زیادی نمی‌کنن.

## نحوه انتخاب C بهینه

سوالی که مطرح هست اینه که تعداد خوشه‌ها رو برای مجموعه‌ی داده‌ای که در اختیار داریم، چی قرار بدیم؟ خیلی راه ایده‌آلی برای انتخاب C بهینه وجود نداره، اما چندتا روش هست که معمولا از اون‌ها استفاده میشه. یکی از اون‌ها، محاسبه‌ی Entropy داده‌ها و بعد استفاده از روش Elbow هستش. آنتروپی از فرمول زیر محاسبه میشه:

$$entropy = - \sum_i \sum_k u_{ik} \ln(u_{ik})$$

همانطور که می‌دونید، مقدار آنتروپی در معنای کلی‌اش نشون‌دهنده میزان بی‌نظمی و پراکندگی داده‌هاست. توی بحث ما، آنتروپی نشون‌دهنده این هست که به طور میانگین چقدر داده‌ها از مراکز خوشه‌هایی که پیدا شده فاصله دارن. اگر به ازای C های مختلف این مقدار رو پلات کنیم، یک چنین نموداری خواهیم داشت:



منطقی هست که با افزایش  $C$ ، میزان آنتروپی حتما کاهش پیدا کنه (چونکه مراکز زیاد میشن و فاصله نقاط از مراکز هم کمتر میشه). حالا روش Elbow می‌گه که توی این سیر نزولی، اون  $C$  ای رو انتخاب کنیم که از اونجا به بعد، دیگه خیلی آنتروپی مون کاهش چشمگیری پیدا نکنه<sup>2</sup>. در نتیجه با توجه به شکل بالا،  $C=3$  می‌تونه گزینه مناسبی باشه.

---

<sup>2</sup> توضیحات بیشتر: [لینک](#)

## موارد تحویلی

الگوریتم C-Means رو پیاده‌سازی کنید و بر روی ۴ دیتاستی که در اختیار دارید اجرا کنید. نمودار آنتروپی برحسب C رو پلات کرده و از طریق روش Elbow، تعداد خوشه‌های بهینه رو تعیین کنید. همچنین، برای دیتاست اول، چند مقدار مختلف برای m را امتحان کرده و در مورد تاثیر آن بر روی خوشه‌های ایجاد شده بحث کنید.

از اونجایی که داده‌ها به همگی خوشه‌ها تعلق دارن، در نتیجه نمی‌تونیم خیلی ساده بیایم برای هر خوشه یه رنگ خاص در نظر بگیریم و داده‌های متعلق به اون خوشه رو با اون رنگ نشون بدیم و پلات کنیم (در واقع برای نمایش دقیق خروجی خوشه‌بندی، نیازه که یه Color Gradient طوری رو در نظر بگیریم). اما صرفاً برای اینکه یک شهودی داشته باشیم، می‌تونیم خوشه‌ها و داده‌ها رو به دنیای Crisp برگردونیم. کافیه که ببینیم هر داده به کدوم خوشه بیشتر از بقیه خوشه‌ها تعلق داره، و از اون طریق رنگ اون خوشه رو بهش بدیم. برای داده‌ها دوبعدی دیتاست، این کار رو انجام بدید و پلات کنید (مراکز خوشه‌ها را نیز پلات کنید).