



Pattern Recognition

HW8 : clustering

Spring 99



۱. L^p Norm

نورم p ام یک بردار در حالت کلی به شکل زیر تعریف می شود که در آن p یک عدد طبیعی است.

$$\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} \quad p \in \mathbb{N}.$$

با استفاده از تعریف بالا موارد زیر را اثبات کنید.

$$\|x\|_\infty = \max_j (|x_j|) \quad (۱)$$

$$\|x\|_\infty \leq \|x\|_2 \leq \sqrt{n} \|x\|_\infty \quad (ب)$$

$$\|x\|_\infty \leq \|x\|_2 \leq \|x\|_1 \quad (ج)$$

* در سوالات ۲ تا ۵ از مجموعه داده DS1 استفاده کنید.

۲. Agglomerative Hierarchical Clustering

در این سوال الگوریتم خوشه بندی سلسله مراتبی تجمعی را بر روی مجموعه داده DS1 که همراه تمرین قرار گرفته است، اعمال می کنیم. در این سوال می توانید از کتابخانه های آماده استفاده کنید. در مجموعه داده علاوه بر نمونه ها، شماره کلاستر هر نمونه نیز جدا گانه آورده شده است. برای مقایسه عملکرد الگوریتم از این شماره ها استفاده می کنیم.

(۱) با استفاده از الگوریتم خوشه بندی تجمعی داده ها را به ۳ خوشه دسته بندی کنید.

(ب) برای هر یک از خوشه ها متوسط فاصله از میانگین دسته را محاسبه کنید. با تحلیل اعداد بدست آمده انسجام درونی خوشه ها را بایک دیگر مقایسه کنید.

$$Mean\ Dist\ Mean(S_i) = \frac{1}{Q_i} \sum_{x \in S_i} d(x, \mu_i)$$

در رابطه بالا منظور از $d(x, y)$ فاصله اقلیدسی دو بردار است، همچنین S_i نشانه خوشه i ام می باشد.

(ج) برای این خوشه بندی، درخت Dendrogram را رسم نمایید. سپس با استفاده از مفهوم life time تعیین کنید که چه تعداد خوشه برای این داده ها مناسب است.

(د) با استفاده از شماره خوشه هایی که در اختیار شما قرار گرفته است، دقت خوشه بندی را بدست آورید. توجه کنید که الگوریتم خوشه بندی تناظر یک به یک میان خوشه ها و شماره هایی که در اختیار شما قرار داده شده است برقرار نمی کند و این نگاهت را خود شما باید ایجاد کنید. در جفت کردن خوشه ها به اعداد حالتی را در نظر بگیرید که دقت خوشه بندی بالاتر می شود.

۳. KMeans

الگوریتم KMeans از جمله الگوریتم هایی است که بر روی داده هایی با خوشه های compact بسیار موثر واقع می شود. در این الگوریتم هر یک از نمونه ها به یکی از خوشه های $S = \{S_1, S_2, \dots, S_k\}$ تخصیص داده می شوند که در آن $k \leq Q$ است. تابع هدفی که در این الگوریتم بهینه می شود within-cluster sum of squares (WCSS) نام دارد که رابطه آن در زیر آورده شده است.

$$\min_s \left(\sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|^2 \right)$$

مراحل الگوریتم:

۱. مقداردهی اولیه Q, n, k و μ_i . برای مقداردهی اولیه نشان‌دسته‌ها می‌توانید k نمونه از مجموعه داده را انتخاب کنید.

۲. تا زمانی که مقادیر نشان‌دسته‌ها ثابت نشده اند مراحل زیر را تکرار می‌کنیم:

۱.۲ هر یک از نمونه‌ها را به نزدیک‌ترین نشان‌دسته اختصاص می‌دهیم.

۲.۲ مقدار نشان‌دسته را به میانگین نمونه‌های اختصاص یافته به آن، بروزرسانی می‌کنیم.

۳. نشان دسته‌ها را باز می‌گردانیم.

برای اطلاعات بیشتر به [اینجا](#) مراجعه کنید.

(آ) با استفاده از الگوریتم بالا داده‌ها را به ۳ خوشه تقسیم کنید. (در این قسمت مجاز به استفاده از کتابخانه‌های آماده نیستید.)

(ب) با استفاده از کتابخانه **sklearn** مرحله قبل را تکرار کنید.

(ج) برای هر یک از خوشه‌بند‌های بالا مقدار **Mean Dist Mean** را برای هر یک از خوشه‌ها محاسبه کنید. با تحلیل این اعداد عملکرد پیاده‌سازی خود و کتابخانه را مقایسه کنید.

(د) دقت هر یک از خوشه‌بند‌های بالا را محاسبه کنید. توجه کنید که الگوریتم خوشه‌بندی تناظر یک‌به‌یک میان خوشه‌ها و شماره‌هایی که در اختیار شما قرار داده شده است برقرار نمی‌کند و این نگاشت را خود شما باید ایجاد کنید. در جفت کردن خوشه‌ها به اعداد حالتی را در نظر بگیرید که دقت خوشه‌بندی بالاتر می‌شود.

۴. Separation Index

شاخص **SI** یک متریک برای **cluster validity** می‌باشد. این شاخص متناسب با نسبت فاصله‌برون خوشه‌ای بر فاصله درون خوشه‌ای می‌باشد، بنابراین هرچه خوشه‌ها از یک دیگر تمیزپذیرتر باشند و انسجام درونی بیشتری داشته باشند، مقدار **SI** بزرگتری دارند. رابطه این شاخص در زیر آورده شده است.

$$SI = \min_j \left\{ \min_{i(i \neq j)} \left\{ \frac{d(S_i, S_j)}{\max_l d(S_l, S_l)} \right\} \right\}$$

$$d(S_i, S_j) = \min_{x,y} \{d(x, y | x \in S_i, y \in S_j)\}$$

$$d(S_l, S_l) = \min_{x,y} \{d(x, y | x, y \in S_l)\}$$

شاخص **SI** را برای هر یک از خوشه‌بند‌های سوال ۲ و ۳ محاسبه کنید. کدام یک از خوشه‌بند‌ها اعتبار بیشتری دارد؟

۵. Fisher's Discriminant Index

یکی دیگر از شاخص‌های جدایی‌پذیری **FDI** می‌باشد که با آن در مبحث **multi-class LDA** آشنا شدید. رابطه این شاخص در زیر آمده است.

$$FDI = \text{trace}(S_W^{-1} S_B)$$

$$S_W = \sum_{i=1}^k S_i \quad S_B = \sum_{i=1}^k Q_i (\hat{\mu}_i - \hat{\mu})(\hat{\mu}_i - \hat{\mu})^T \quad S_i = \sum_{q=1}^{Q_i} (x^q - \hat{\mu}_i)(x^q - \hat{\mu}_i)^T$$

$$\hat{\mu}_i = \sum_{x^q \in S_i} x^q \quad \mu = \sum_{q=1}^Q x^q$$

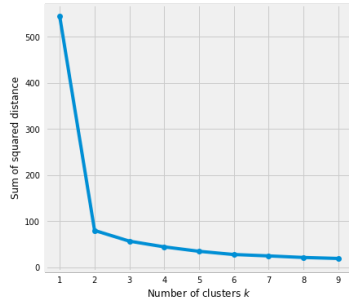
شاخص **FDI** را برای هر یک از خوشه‌بند‌های سوال ۲ و ۳ محاسبه کنید. کدام یک از خوشه‌بند‌ها اعتبار بیشتری دارد؟

۶. Image Compression using KMeans

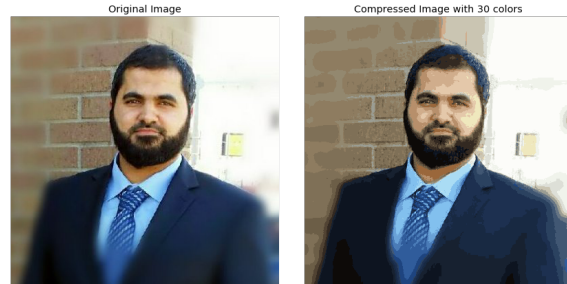
از الگوریتم **KMeans** میتوان برای فشرده سازی تصاویر استفاده کرد. به این صورت که k رنگ از رنگ های تصویر انتخاب می شود و هر یک از رنگ های تصویر به نزدیک ترین رنگ از k رنگ تغییر پیدا می کند. در حالت معمولی برای ذخیره سازی تصویر به $8 \times 3 \times \text{rows} \times \text{columns}$ بیت نیاز است (هر پیکسل ۳ کانال رنگی دارد که هر رنگ با ۸ بیت نمایش داده می شود). ، در حالی که با استفاده از **KMeans** می توان تصویر را با اندازه $8 \times 3 \times k + \log_2 k \times \text{rows} \times \text{columns}$ بیت نمایش داد.

با استفاده از الگوریتم **KMeans** تصویر همراه تمرین را فشرده کنید و تصویر حاصل را در گزارش خود بیاورید.

برای انتخاب k مناسب از **Elbow Method** استفاده می کنیم. در این روش خطای **MSE** میان تصویر فشرده شده و تصویر اصلی را برای k های مختلف رسم می کنیم. نقطه شکستگی نمودار را به عنوان k مناسب انتخاب می کنیم.



(ب) روش elbow برای انتخاب k مناسب



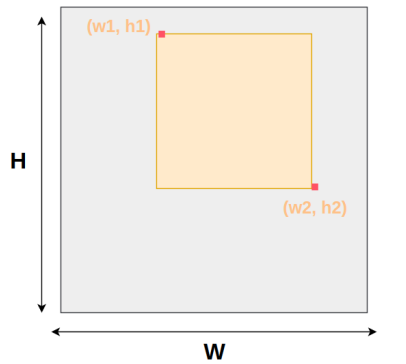
(آ) نمونه از فشرده سازی تصویر با استفاده از خوشه بندی

۷. Bounding Box Clustering

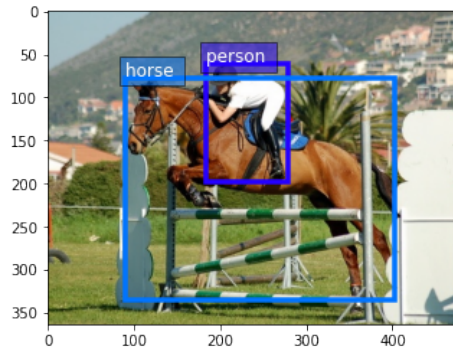
استفاده اصلی **bounding box** ها در سیستم های **object detection** می باشد. این باکس ها حدود یک شی را در تصویر مشخص می کنند. اهمیت دسته بندی این باکس ها در آنجا است که سیستم هایی همانند **YOLO** مشخصاتی از چندین باکس به عنوان هاپر پارامتر دارند. در این تمرین بر مبنای **KMeans Clustering** برای تعیین این هاپر پارامتر ها معرفی می شود.

در این سوال از مجموعه داده **PASCAL Visual Object Classes (VOC)** استفاده می شود که برای سادگی بیشتر نشان گذاری های این مجموعه داده در فایل **bboxes.csv** در اختیار شما قرار گرفته است.

یک باکس در حالت کلی با ۴ عدد (w_1, h_1, w_2, h_2) نشان داده می شود که در آن مختصات گوشه بالا-چپ و گوشه پایین-راست می باشد. از آنجایی که موقعیت و اندازه نسبی باکس در خوشه بندی اهمیت ندارد، هر باکس را به صورت نرمالیزه شده $(\frac{w_2-w_1}{W}, \frac{h_2-h_1}{H})$ نشان می دهیم که در آن W و H طول و عرض تصویر می باشد.



(ب) مشخصات یک باکس در تصویر



(آ) نمونه تصویر از مجموعه داده

(آ) با استفاده از الگوریتم **KMeans** و فاصله اقلیدسی، داده ها را یک بار به ۵ خوشه و بار دیگر به ۹ خوشه تقسیم کنید. با استفاده از **scatter plot** چینش خوشه ها را نمایش دهید.

متریکی که در قسمت قبل برای فاصله سنجی استفاده شد، متریک مناسبی نیست چرا که به باکس هایی با اندازه بزرگ تر ارزش بیشتری می دهد. (اگر از پارامتر α در **satter plot** استفاده کنید، چگالی نقاط بهتر دیده می شوند و این ارزش دهی نادرست فاصله اقلیدسی مشهود تر است.) به همین دلیل از متریک **Intersection over Union** برای شباهت سنجی دو باکس استفاده خواهیم کرد. برای اطلاعات بیشتر به [اینجا](#) مراجعه کنید.

$$IoU(b_1, b_2) = \frac{\min(x_1, x_2) \times \min(y_1, y_2)}{x_1 y_1 + x_2 y_2 - \min(x_1, x_2) \times \min(y_1, y_2)}$$

$$b_1 = (x_1, y_1) \quad b_2 = (x_2, y_2) \quad d(b_1, b_2) = 1 - IoU(b_1, b_2)$$

(ب) این بار با استفاده از متریک **IoU** داده ها را یک بار به ۵ خوشه و بار دیگر به ۹ خوشه تقسیم کنید. همچنین با استفاده از **scatter plot** تقسیم بندی داده ها را نشان دهید.



شکل ۳: متریک **IoU** برای چند باکس

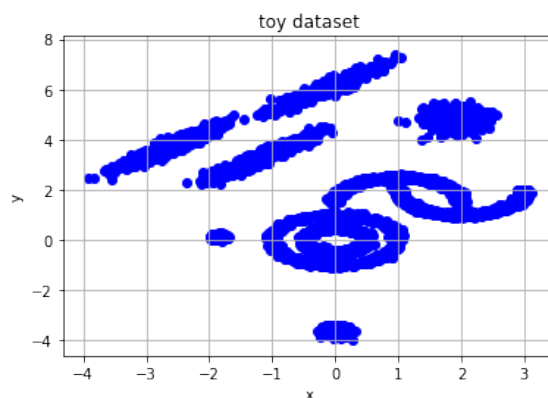
۸. DBSCAN vs KMeans

دسته دیگری از الگوریتم های خوشه بندی بر اساس چگالی داده ها خوشه ها را انتخاب می کنند، از این دسته از الگوریتم ها می توان به روش خوشه بندی **Density-Based Spatial Clustering of Applications with Noise** اشاره کرد. برای آشنایی بیشتر به [اینجا](#) و [اینجا](#) مراجعه کنید.

(آ) به صورت مختصر عملکرد این الگوریتم را توضیح دهید

(ب) در این قسمت می خواهیم عملکرد **DBSCAN** و **KMeans** را بر روی یک مجموعه داده ساختگی (**ToyDataSet.csv**) با یک دیگر مقایسه کنیم. این مجموعه داده دارای ۱۰ خوشه مجزا است. با استفاده از هر یک از روش ها داده ها را خوشه بندی کنید و با استفاده از **Scatter plot** خوشه ها را نمایش دهید.

(ج) عملکرد **DBSCAN** بر چه نوع خوشه هایی بهتر از عملکرد **KMeans** می باشد؟



شکل ۴: مجموعه داده ساختگی

نکات پایانی:

۱. شما باید پاسخ های خود را با الگو PATREC-HW8-SID.zip در محل تعیین شده آپلود کنید
 ۲. گزارش شما معیار اصلی ارزیابی خواهد بود، در نتیجه دقت کنید کیفیت عکس ها مناسب باشند.
 ۳. کدهای خود را به تفکیک سوال ارسال کنید و استفاده از دیگر زبان های برنامه نویسی ممانعتی ندارد.
 ۴. هدف از انجام تمرین یادگیری مباحث درس می باشد، بنابر این تمرین را خودتان انجام دهید. در صورت کشف مشابهت بلا توجیه، با توجه به قوانین درس عمل خواهد شد.
 ۵. شما میتوانید سوالات خود را از طریق ایمیل sj.pakdaman@ut.ac.ir بپرسید
- موفق باشید