

From Unsupervised Clustering Results to MLP Classification

- استفاده از نتایج Clustering برای آموزش MLP

مقاله مرجع (Clustering)

**Performance Evaluation of K-Means and Hierarchical Clustering

in Terms of Accuracy and Running Time**

در این پژوهه، هدف بررسی یک مسیر یادگیری ترکیبی است که از **Unsupervised Learning** آغاز می‌شود و در ادامه، خروجی آن برای آموزش یک **Multilayer Perceptron (MLP)** مورد استفاده قرار می‌گیرد.

در تمام مراحل پژوهه، تحلیل‌ها و تصمیم‌گیری‌ها باید صرفاً بر اساس مطالب، نتایج عددی و نتیجه‌گیری‌های مقاله مرجع (cluster.pdf) Clustering محاسبه فرمول‌ها نیست.

در مقاله مرجع، عملکرد دو الگوریتم **Hierarchical Clustering** و **K-Means** روی دیتاست‌هایی مانند **Diabetes** و **Iris** مورد بررسی قرار گرفته است. اگرچه این دیتاست‌ها دارای Label واقعی هستند، اما در فرآیند Clustering، این Label‌ها به صورت کامل نادیده گرفته شده‌اند و تنها برای ارزیابی نهایی **Accuracy** مورد استفاده قرار گرفته‌اند. بنابراین، انجام شده در مقاله ماهیت **Unsupervised Clustering** دارد.

در ابتدای پژوهه، دانشجو باید یکی از دیتاست‌های معرفی‌شده در مقاله Iris (یا Diabetes) را انتخاب کرده و همان دیتاست را در تمام مراحل پژوهه دنبال نماید. سپس، باید هر دو الگوریتم **K-Means** و

Hierarchical Clustering را روی دیتاست انتخاب شده اجرا کند. هدف از اجرای هر دو الگوریتم، مقایسه عملکرد آن‌ها و تحلیل تفاوت نتایج است.

در این پژوهه، دانشجو ملزم به پیاده‌سازی دقیق، گام‌به‌گام یا ابزار محور مقاله نیست. با این حال، اجرای Clustering با ایده، منطق کلی و چارچوب آزمایشی مقاله انجام شود؛ به این معنا که:

- نوع دیتاست‌ها و مسئله مطابق مقاله انتخاب شود؛
- تعداد خوش‌ها با ساختار مسئله هم‌خوان بوده و مشابه تنظیمات مقاله در نظر گرفته شود؛
- Label‌های واقعی کلاس در مرحله Clustering استفاده نشوند؛
- Accuracy صرفاً برای تحلیل و مقایسه با نتایج مقاله محاسبه شود.

هدف از این رویکرد آن است که نتایج به دست آمده توسط دانشجو از نظر رفتار کلی الگوریتم‌ها و الگوی عملکرد، قابل مقایسه با نتایج گزارش شده در **Table 2** مقاله باشد. تأکید می‌شود که یکسان بودن عددی Accuracy با مقاله مدنظر نیست؛ بلکه تحلیل این موضوع اهمیت دارد که کدام الگوریتم، مطابق نتیجه‌گیری مقاله (**Conclusion – Page 3**)، برای دیتاست انتخاب شده مناسب‌تر است.

پس از تحلیل نتایج Clustering و با تکیه بر نتیجه‌گیری مقاله، دانشجو باید الگوریتم مناسب‌تر را برای دیتاست انتخاب شده مشخص کند. خروجی این الگوریتم شامل شماره خوش‌های نمونه داده است. این شماره خوش‌ها در ادامه پژوهه به عنوان **Pseudo-Label** جدید داده‌ها در نظر گرفته می‌شوند که به آن‌ها گفته می‌شود.

در مرحله بعد، یک شبکه **Feedforward Multilayer Perceptron (MLP)** با استفاده از همین Pseudo-Label‌ها آموزش داده می‌شود. در این مرحله:

- ورودی شبکه، Feature‌های اصلی دیتاست است؛
- خروجی شبکه، Pseudo-Label‌های تولیدشده توسط Clustering است.

در این پژوهش، شبکه MLP وظیفه یادگیری Label‌های واقعی دیتاست را ندارد، بلکه هدف آن یادگیری و مدل‌سازی ساختار خوشبندی کشف شده در مرحله Unsupervised است. بنابراین، آموزش MLP به عنوان ادامه مفهومی فرآیند Clustering مقاله تلقی می‌شود.

در نهایت، کل این فرآیند نمونه‌ای از Clustering-based Label Propagation محسوب می‌شود (نه Semi-Supervised به معنای دقیق کلمه، چون در هیچ مرحله‌ای از Label واقعی در آموزش استفاده نمی‌شود)؛ زیرا:

- کشف ساختار داده‌ها به صورت Unsupervised انجام شده است؛
- آموزش MLP به صورت Supervised صورت می‌گیرد، اما با Label‌هایی که واقعی نیستند و به صورت Pseudo تولید شده‌اند.

هدف نهایی پژوهش این است که دانشجو درک کند چگونه می‌توان، با تکیه بر تحلیل نتایج یک مقاله Unsupervised Clustering برای آموزش یک مدل MLP استفاده کرد و این فرآیند را به صورت مفهومی تحلیل نمود.

مراحل پروژه

مرحله ۱ — انتخاب دیتاست و اجرای Clustering

۱-الف) یکی از دو دیتاست **Diabetes** یا **Iris** را که در مقاله مرجع بررسی شده‌اند انتخاب کنید. دیتاست انتخاب شده را در تمام مراحل پروژه ثابت نگه دارید.

۱-ب) هر دو الگوریتم **Hierarchical Clustering** و **K-Means** را روی دیتاست انتخاب شده اجرا کنید.
در اجرای هر دو الگوریتم:

- تعداد خوشها باید با ساختار مسئله (تعداد کلاس‌های واقعی دیتاست) همخوانی داشته باشد
- واقعی داده‌ها در مرحله Clustering استفاده نشود Label
- Accuracy نهایی Clustering را برای هر دو الگوریتم محاسبه و گزارش دهید

۱-ج) نتایج Accuracy به دست آمده را با **Table 2** مقاله مقایسه کنید. آیا الگوی رفتاری الگوریتم‌ها (کدام بهتر است) با نتایج مقاله همخوانی دارد؟ توضیح دهید.

۱-د) بر اساس **Conclusion** (صفحه ۳) مقاله، کدام الگوریتم برای دیتاست انتخاب شده مناسب‌تر است؟ دلیل بیاورید.

مرحله ۲ — تولید Pseudo Labels

۲-الف) الگوریتم برتری که در مرحله قبل مشخص کردید را انتخاب کنید. خروجی این الگوریتم (شماره خوشه هر نمونه) را به عنوان **Pseudo Label** در نظر بگیرید.

۲-ب) توضیح دهید که چرا به این Label‌ها «Pseudo» می‌گویند و چه تفاوتی با Label واقعی دیتاست دارند؟

مرحله ۳ — آموزش MLP با Pseudo Labels

در این مرحله یک شبکه Feedforward MLP (Multi-Layer Perceptron) با استفاده از Feature‌های تولیدشده در مرحله قبل آموزش می‌دهید. ورودی شبکه Feature‌های اصلی دیتاست و خروجی آن Pseudo Label‌ها هستند.

۳-الف) آماده‌سازی داده‌ها (مقیاس‌بندی)

پیش از آموزش مدل MLP، لازم است که داده‌های ورودی (Feature‌ها) را مقیاس‌بندی کنید. مدل‌های MLP به داده‌های مقیاس‌بندی شده حساس هستند و این کار به همگرایی بهتر و سریع‌تر مدل کمک می‌کند.

- پیشنهاد: از StandardScaler در کتابخانه scikit-learn برای استانداردسازی داده‌ها استفاده کنید. (یعنی میانگین Feature‌ها صفر و واریانس آن‌ها یک شود).

۳-ب) طراحی و آموزش مدل MLP

یک مدل MLP از نوع **Feedforward** با استفاده از MLPClassifier در کتابخانه scikit-learn طراحی و آموزش دهید. در زمان تعریف مدل، به نکات زیر توجه کنید:

۱. اندازه و تعداد لایه‌های پنهان (`hidden_layer_sizes`):
 - این پارامتر، تعداد لایه‌های میانی (پنهان) و تعداد نورون‌های هر لایه را مشخص می‌کند. این لایه‌ها وظیفه‌ی اصلی یادگیری الگوهای پیچیده در داده‌ها را بر عهده دارند.
 - پیشنهاد: برای شروع، از ساختار `hidden_layer_sizes=(10, 10)` استفاده کنید. این یعنی دو لایه پنهان که هر کدام ۱۰ نورون دارند. این یک نقطه شروع مناسب برای این پروژه است.

۲. تعداد Epoch‌ها :(max_iter)

- پارامتر max_iter (که همان تعداد Epoch‌هاست) نشان می‌دهد که مدل چند بار کل مجموعه داده‌های آموزشی را ببیند و پارامترهای داخلی خود را بهروزرسانی کند.
- پیشنهاد: مقدار max_iter=500 را تنظیم کنید. این تعداد معمولاً برای مدل‌های ساده MLP کافی است تا به یک همگرایی مناسب برسند و مدل بتواند Pseudo Label‌ها را یاد بگیرد.

۳. تابع فعال‌ساز (Activation Function)

- لایه‌های داخلی (مخفي) شبکه: مدل MLPClassifier به صورت پیش‌فرض از تابع فعال‌ساز relu برای لایه‌های داخلی استفاده می‌کند که یک انتخاب مناسب برای یادگیری الگوهاست. نیازی به تغییر این تنظیم نیست.
- لایه خروجی شبکه: در مورد لایه خروجی (که جواب نهایی را تولید می‌کند)، مدل MLPClassifier به صورت هوشمند و خودکار تابع فعال‌ساز مناسب (مثل Softmax برای مسائل چندکلاسه مانند Iris، یا Sigmoid برای مسائل دودویی مانند Diabetes) را متناسب با تعداد دسته‌های دیتای شما تنظیم می‌کند.
- بنابراین، نیازی نیست شما هیچ تابعی برای لایه خروجی انتخاب کنید و تمرکزتان صرفاً بر استفاده از مدل باشد.

مدل را با Pseudo Label‌های تولیدشده آموزش داده و عملکرد آن را روی داده‌های آزمایشی گزارش کنید.

۴-ج) تحلیل مفهومی خطای آموزش با افزایش Epoch

به صورت مفهومی توضیح دهید که به طور کلی با افزایش تعداد Epoch در آموزش یک شبکه MLP ، انتظار دارید چه اتفاقی برای خطای مدل بیفتد؟

۴) ارزیابی مدل

پس از آموزش، مدل را روی داده‌های Test ارزیابی کنید:

Accuracy

در صد نمونه‌هایی که MLP خوش آن‌ها را درست پیش‌بینی کرده است را محاسبه و گزارش دهید. این عدد نشان می‌دهد مدل تا چه حد توانسته ساختار خوشه‌بندی را یاد بگیرد.

Confusion Matrix

دانشجو باید **Confusion Matrix** مدل روی داده‌های تست را محاسبه و نمایش دهد

سوالات تحلیلی

سوال ۱) عدد Accuracy به دست‌آمده از MLP را با Clustering Accuracy مرحله قبل مقایسه کنید. آیا این دو عدد باید به هم نزدیک باشند؟ چرا؟

سوال ۲) از روی Confusion Matrix تحلیل کنید که MLP در تشخیص کدام خوشه بیشترین خطأ را داشته. این خطأ با کیفیت Pseudo Label های آن خوشه ارتباط دارد؟