

دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

دانشکده ریاضیات و علوم کامپیوتر

تحلیل قابلیت های خوشه بندی NMF

پروژه پایانی درس جبر خطی عددی

آرمان صالحی

استاد درس:

دکتر مهدی دهقان

زمستان ۱۴۰۴



چکیده

چکیده در این پروژه، به بررسی و پیاده سازی الگوریتم تجزیه نامنفی ماتریس (NMF) به عنوان یک روش قدرتمند برای کاهش ابعاد و خوشه بندی داده ها پرداخته شده است. هدف اصلی، اعمال این الگوریتم بر روی مجموعه داده ارقام دست نویس (MNIST) و استخراج الگوهای پایه سازنده این ارقام است. برخلاف روش های سنتی مانند K-means که بر اساس میانگین گیری از کل داده عمل می کنند، NMF با تجزیه داده ها به مولفه های نامنفی، رویکردی "بخش محور" ارائه می دهد که تفسیرپذیری بالایی دارد. در بخش عملی، نتایج حاصل از NMF با الگوریتم K-means مقایسه شده است. نتایج بصری نشان می دهد که پایه های به دست آمده از NMF شبیه به اجزای تشکیل دهنده اعداد (مانند خطوط و کمان ها) هستند، در حالی که مراکز خوشه های K-means شبیه به میانگین محو شده اعداد کامل می باشند. همچنین دقت خوشه بندی هر دو روش به صورت کمی محاسبه و مقایسه گردید.

کلمات کلیدی: جبر خطی عددی، تجزیه نامنفی ماتریس (NMF)، خوشه بندی K-means، داده کاوی،

مجموعه داده MNIST



فهرست مطالب

فصل اول: مقدمه و کلیات	۴
۱-۱ مقدمه	۴
۱-۲ بیان مسئله	۴
۱-۳ معرفی روش های مورد استفاده	۵
۱-۴ ساختار گزارش	۶
فصل دوم: مبانی نظری و پیشینه تحقیق	۶
۲-۱ مقدمه	۶
۲-۲ تجزیه نامنفی ماتریس (NMF)	۶
۲-۳ الگوریتم خوشه بندی K-means	۹
۲-۴ مقایسه نظری: جزءنگر در برابر کل نگر	۱۰
فصل سوم: پیاده سازی و تحلیل نتایج	۱۱
۳-۱ معرفی مجموعه داده MNIST	۱۱
۳-۲ ابزارها و محیط پیاده سازی	۱۱
۳-۳ نتایج تجزیه نامنفی ماتریس (NMF)	۱۲
۳-۴ نتایج خوشه بندی K-means	۱۳
۳-۵ مقایسه کمی (دقت خوشه بندی)	۱۴
فصل چهارم: نتیجه گیری و پیشنهادات	۱۵
۴-۱ نتیجه گیری	۱۵
۴-۲ پیشنهادات برای کارهای آتی	۱۶
مراجع	۱۷
پیوست	۱۸



فصل اول: مقدمه و کلیات

۱-۱ مقدمه

در عصر حاضر، با رشد روزافزون داده‌های دیجیتال، نیاز به روش‌های خودکار برای تحلیل و استخراج اطلاعات از این داده‌ها بیش از پیش احساس می‌شود. یکی از حوزه‌های اصلی در هوش مصنوعی و یادگیری ماشین، یادگیری نظارت نشده (Unsupervised Learning) است که در آن الگوریتم بدون داشتن برچسب‌های از پیش تعیین شده، سعی در کشف ساختارهای پنهان و الگوهای موجود در داده‌ها دارد.

خوشه‌بندی (Clustering) به عنوان یکی از مهم‌ترین تکنیک‌های یادگیری نظارت نشده، فرآیندی است که در آن داده‌ها به گروه‌هایی تقسیم می‌شوند که اعضای هر گروه بیشترین شباهت را به یکدیگر و بیشترین تفاوت را با اعضای سایر گروه‌ها داشته باشند. این تکنیک در کاربردهای متنوعی از جمله پردازش تصویر، فشرده‌سازی داده‌ها، و تشخیص الگو کاربرد دارد.

در این پروژه، تمرکز ما بر روی تحلیل و خوشه‌بندی تصاویر ارقام دست‌نویس است. چالش اصلی در پردازش این تصاویر، ابعاد بالای داده‌ها و تغییرات زیاد در نحوه نوشتن ارقام توسط افراد مختلف است. برای غلبه بر این چالش، از روش‌های کاهش ابعاد و تجزیه ماتریس استفاده می‌شود.

۱-۲ بیان مسئله

مسئله مورد بررسی در این پژوهش، خوشه‌بندی تصاویر مجموعه داده MNIST است. این مجموعه داده شامل تصاویر سیاه و سفید از ارقام دست‌نویس (۰ تا ۹) است که هر کدام در یک شبکه پیکسلی مشخص (۲۸ در ۲۸ پیکسل) قرار دارند.

هدف اصلی این است که با استفاده از الگوریتم‌های جبر خطی عددی، این تصاویر را پردازش کرده و الگوریتم بتواند بدون دانستن برچسب واقعی تصاویر (یعنی بدون اینکه بداند کدام عکس عدد ۵ است یا ۷)، آن‌ها را بر اساس ویژگی‌های ظاهری دسته‌بندی کند. همچنین، استخراج ویژگی‌های پایه‌ای که این ارقام را تشکیل می‌دهند (مانند خطوط عمودی، افقی و حلقه‌ها) بخش مهمی از این پژوهش است.



۳-۱ معرفی روش های مورد استفاده

برای انجام این پروژه، از دو رویکرد متفاوت استفاده و نتایج آن ها مقایسه خواهد شد:

۱. تجزیه نامنفی ماتریس (NMF): این روش یکی از تکنیک های قدرتمند در جبر خطی است که ماتریس داده ها را به دو ماتریس با درایه های نامنفی تجزیه می کند. ویژگی منحصر به فرد NMF این است که به دلیل قید نامنفی بودن، اجزای سازنده داده ها را به صورت «بخش محور» (Parts-based) «یاد می گیرد. این یعنی NMF سعی می کند تصاویر را به عنوان ترکیبی از اجزای کوچکتر (مانند تکه های قلم) بازسازی کند.

۲. الگوریتم K-means: این الگوریتم یکی از کلاسیک ترین روش های خوشه بندی است که بر اساس فاصله اقلیدسی و میانگین گیری عمل می کند. برخلاف NMF، روش K-means معمولاً رویکردی «کل محور» دارد و مراکز خوشه ها در آن شبیه به میانگین کلی تصاویر یک دسته هستند.



شکل ۱-۱: نمونه هایی از تنوع دست خط ها در مجموعه داده MNIST



۴-۱ ساختار گزارش

این گزارش در چهار فصل تدوین شده است. پس از این مقدمه، در فصل دوم به مبانی نظری و ریاضیاتی تجزیه NMF و الگوریتم K-means خواهیم پرداخت. در فصل سوم، پیاده سازی عملی الگوریتم ها بر روی داده های MNIST شرح داده شده و نتایج کمی و کیفی (شامل تصاویر پایه ها و دقت خوشه بندی) تحلیل می شوند. در نهایت، در فصل چهارم نتیجه گیری کلی پروژه و پیشنهادات ارائه خواهد شد.

فصل دوم: مبانی نظری و پیشینه تحقیق

۱-۲ مقدمه

در این فصل به بررسی دقیق ریاضیاتی الگوریتم های مورد استفاده در پروژه می پردازیم. ابتدا الگوریتم تجزیه نامنفی ماتریس (NMF) و فرمول بندی آن در جبر خطی را تشریح کرده و سپس به الگوریتم خوشه بندی K-means خواهیم پرداخت. در نهایت، تفاوت بنیادین این دو روش در نحوه نمایش داده ها بررسی می شود.

۲-۲ تجزیه نامنفی ماتریس (NMF)

تجزیه نامنفی ماتریس (Non-negative Matrix Factorization) گروهی از الگوریتم ها در تحلیل چندمتغیره و جبر خطی است که در آن یک ماتریس X به دو ماتریس H و W تجزیه می شود، با این قید مهم که تمام درایه های این سه ماتریس باید نامنفی (بزرگتر یا مساوی صفر) باشند.



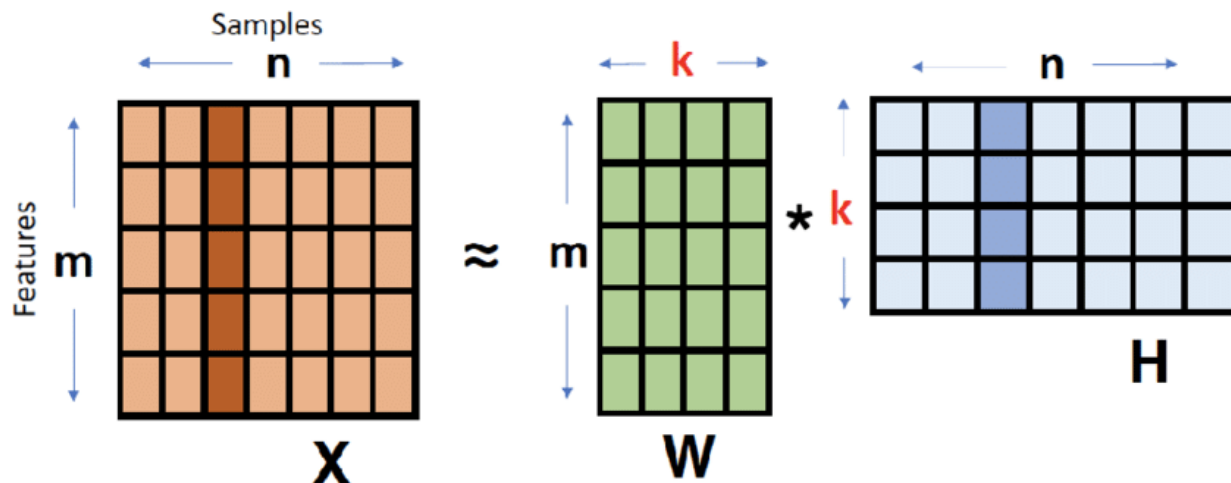
۲-۲-۱ فرمول بندی ریاضی

فرض کنید مجموعه داده های ما به صورت ماتریس x با ابعاد $n \times m$ باشد که در آن n تعداد ویژگی ها (پیکسل ها) و m تعداد نمونه ها (تصاویر) است. هدف NMF یافتن دو ماتریس H و W است بطوریکه:

$$X \approx WH$$

که در آن:

- ماتریس $X \in \mathbb{R}^{n \times m}$: ماتریس داده های ورودی (شامل درایه های نامنفی).
- ماتریس $W \in \mathbb{R}^{n \times r}$: ماتریس ویژگی های پایه. ستون های این ماتریس همان «الگوهای پایه» یا اجزای سازنده تصاویر هستند.
- ماتریس $H \in \mathbb{R}^{r \times m}$: ماتریس ضرایب یا کدگذاری. هر ستون نشان می دهد که چه ترکیبی از پایه ها برای ساختن تصویر اصلی نیاز است.
- پارامتر r : رتبه تجزیه یا تعداد مولفه ها که معمولاً بسیار کوچکتر از n و m انتخاب می شود.





۲-۲-۲ قید نامنفی بودن و تفسیرپذیری

تفاوت اصلی NMF با روش‌هایی مانند PCA (تحلیل مولفه‌های اصلی) در قید زیر است:

$$X_{ij} \geq 0, \quad W_{ik} \geq 0, \quad H_{kj} \geq 0$$

این قید باعث می‌شود که بازسازی داده‌ها تنها از طریق جمع مولفه‌ها صورت گیرد و هیچ تفریق کردنی وجود نداشته باشد. به عبارت دیگر، یک تصویر به صورت ترکیبی از روی هم قرار گرفتن چندین لایه (پایه) ساخته می‌شود. این ویژگی باعث می‌شود پایه‌های به دست آمده W شبیه به اجزای معنادار تصویر (مانند خطوط، گوشه‌ها و منحنی‌ها در ارقام) باشند.

۲-۲-۳ تابع هزینه و به‌روزرسانی

برای پیدا کردن بهترین ماتریس‌های W و H ، باید فاصله بین ماتریس اصلی X و ماتریس بازسازی شده WH را کمینه کنیم. رایج‌ترین تابع هزینه، نرم فروینیوس است:

$$\min_{W, H} \frac{1}{2} \|X - WH\|_F^2$$

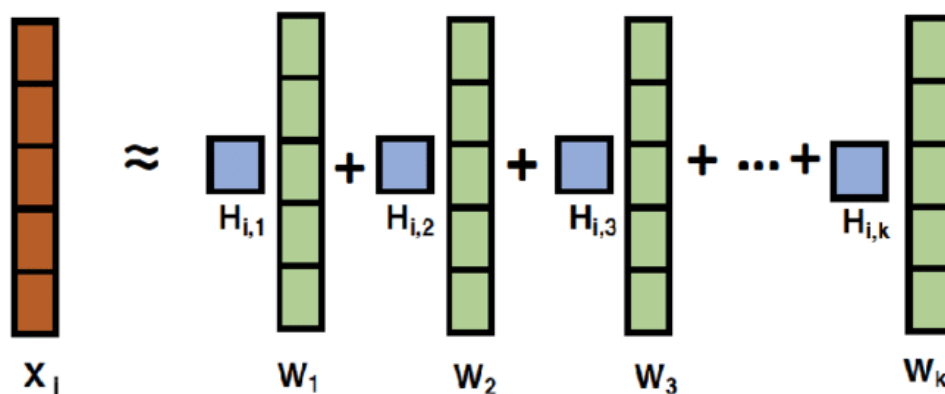
$$\text{subject to } W, H \geq 0$$

از آنجا که این تابع نسبت به هر دو متغیر W و H محدب نیست، راه حل بسته ریاضی ندارد. برای حل آن معمولاً از روش‌های تکرارشده (Iterative) استفاده می‌شود. یکی از معروف‌ترین الگوریتم‌ها، «قوانین به‌روزرسانی ضربی پیشنهاد شده توسط لی و سونگ (Lee & Seung)» است:

$$H_{kj} \leftarrow H_{kj} \frac{(W^T X)_{kj}}{(W^T W H)_{kj}}$$

$$W_{ik} \leftarrow W_{ik} \frac{(X H^T)_{ik}}{(W H H^T)_{ik}}$$

این به‌روزرسانی‌ها تضمین می‌کنند که اگر ماتریس‌های اولیه نامنفی باشند، در تمام تکرارها نامنفی باقی بمانند و تابع هزینه به صورت یکنواخت کاهش یابد.



۲-۳ الگوریتم خوشه بندی K-means

الگوریتم K-means یکی از ساده ترین و پرکاربردترین الگوریتم های خوشه بندی است که هدف آن تقسیم m داده ورودی به k خوشه مجزا است، به طوری که واریانس درون هر خوشه کمینه شود.

۲-۳-۱ تابع هدف

این الگوریتم سعی می کند تابع هدف J که مجموع مربعات خطای فاصله هر نقطه تا مرکز خوشه اش است، کمینه کند:

$$J = \sum_{j=1}^k \sum_{x_i \in C_j} \|x_i - \mu_j\|^2$$

که در آن:

- k : تعداد خوشه ها
- C_j : مجموعه نقاط متعلق به خوشه j
- μ_j : میانگین نقاط خوشه j
- $\|.\|$: فاصله اقلیدسی استاندارد



در زمینه پردازش تصویر، مراکز خوشه ها (μ_j) نماینده «میانگین کلی» تصاویری هستند که در آن خوشه قرار گرفته اند. برخلاف NMF که تصویر را از ترکیب اجزا می سازد، K-means برای هر تصویر فقط یک مرکز (نماینده) در نظر می گیرد.

۴-۲ مقایسه نظری: جزءنگر در برابر کل نگر

مهم ترین تفاوت نظری بین NMF و K-means در نحوه نمایش اطلاعات است:

۱. رویکرد K-means (کل نگر): این روش از نوع Vector Quantization است. هر داده تنها با یک بردار مرجع (مرکز خوشه) تقریب زده می شود. بنابراین مراکز خوشه ها باید شبیه به یک نسخه کامل (اما میانگین گیری شده) از داده ها باشند.

۲. رویکرد NMF (جزءنگر): این روش داده را به صورت ترکیب خطی چندین بردار پایه تقریب می زند ($x \approx \sum w_i h_i$). چون ضرایب منفی مجاز نیستند، الگوریتم مجبور است ویژگی هایی را یاد بگیرد که «اجزای سازنده» داده ها باشند. این تفاوت باعث می شود که در کاربرد تشخیص ارقام، NMF پایه هایی شبیه به قلم موی نوشتن تولید کند، در حالی که K-means تصاویری شبیه به ارقام کامل اما مات تولید می کند.



مقایسه بصری پایه های استخراج شده. (چپ) پایه های PCA که به صورت چهره های کامل اما محو هستند. (راست) پایه های NMF که به صورت اجزای موضعی و تفکیک شده مثل بینی یا ابرو ظاهر شده اند.



فصل سوم: پیاده سازی و تحلیل نتایج

۳-۱ معرفی مجموعه داده MNIST

در این پروژه از مجموعه داده استاندارد MNIST استفاده شده است. این مجموعه شامل ۷۰,۰۰۰ تصویر از ارقام دست نویس (۰ تا ۹) است که به طور گسترده برای آموزش و تست سیستم های یادگیری ماشین استفاده می شود.

- تعداد نمونه ها: ۷۰,۰۰۰ تصویر (۶۰,۰۰۰ داده آموزش و ۱۰,۰۰۰ داده تست).
- ابعاد تصاویر: هر تصویر سیاه و سفید و دارای ابعاد ۲۸ در ۲۸ پیکسل است.
- فرمت داده: هر تصویر به صورت یک بردار مسطح شده با طول ۷۸۴ (28×28) نمایش داده می شود تا بتواند به عنوان ورودی به الگوریتم های ماتریسی داده شود.

۳-۲ ابزارها و محیط پیاده سازی

برای پیاده سازی این پروژه از زبان برنامه نویسی Python استفاده شده است. جهت مدیریت وابستگی ها و محیط مجازی، از ابزار Poetry استفاده گردید تا تکرارپذیری نتایج تضمین شود. کتابخانه های اصلی مورد استفاده عبارتند از:

۱. Scikit-learn: برای اجرای الگوریتم های NMF و K-means و محاسبه معیارهای ارزیابی.

۲. NumPy: برای عملیات ماتریسی و برداری.

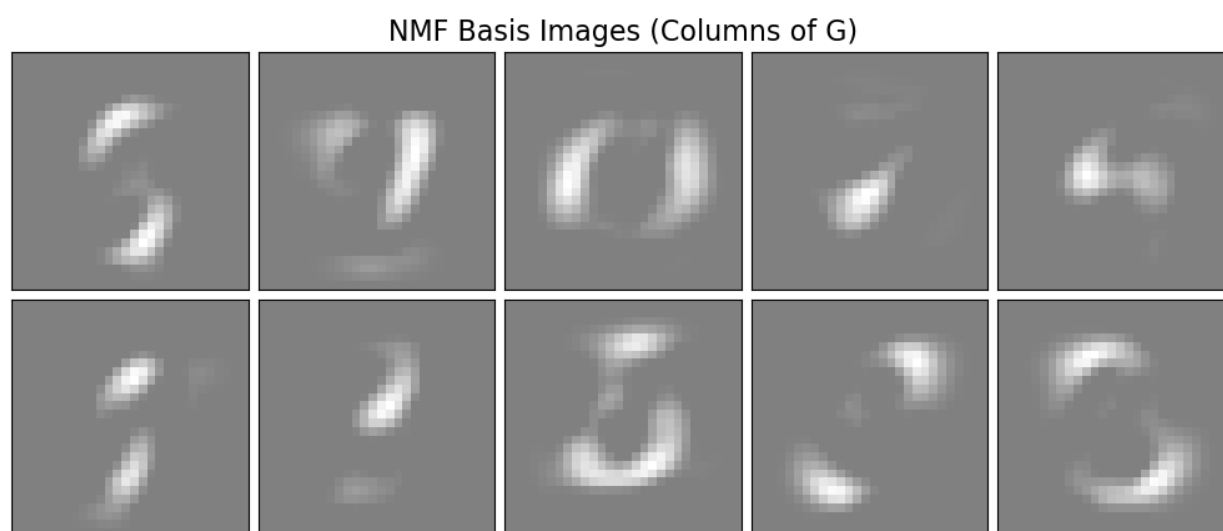
۳. Matplotlib: برای رسم نمودارها و نمایش تصاویر خروجی.



۳-۳ نتایج تجزیه نامنفی ماتریس (NMF)

در این بخش، الگوریتم NMF با تعداد مولفه‌های ۱۰ ($n_{components} = 10$) بر روی کل داده‌های نرمال‌سازی شده اجرا گردید. هدف از این کار، یافتن ۱۰ الگوی پایه‌ای است که تمامی ارقام دیگر از ترکیب آن‌ها ساخته می‌شوند.

شکل زیر خروجی حاصل از ماتریس H (تصاویر پایه) را نشان می‌دهد:



شکل ۳-۱: تصاویر پایه استخراج شده توسط الگوریتم NMF

تحلیل کیفی نتایج NMF

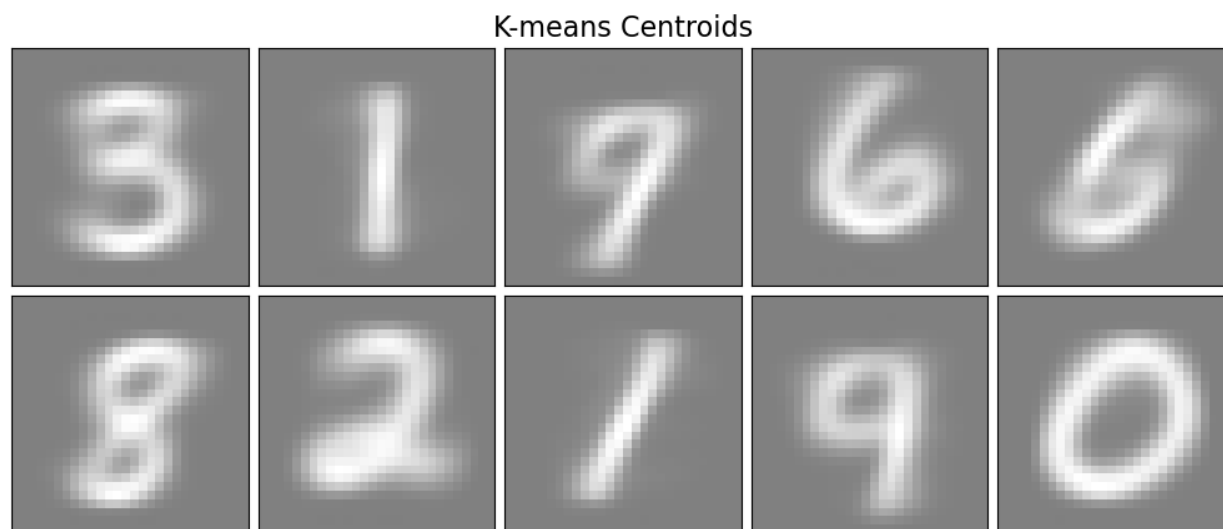
همان‌طور که در شکل ۳-۱ مشاهده می‌شود، پایه‌های به دست آمده شبیه به ارقام کامل نیستند، بلکه شبیه به اجزای سازنده ارقام هستند. برای مثال:

- برخی پایه‌ها تنها یک خط کج یا افقی را نشان می‌دهند.
 - برخی دیگر شبیه به کمان‌ها یا حلقه‌هایی هستند که در اعداد ۸، ۰ یا ۹ دیده می‌شوند.
- این نتیجه دقیقاً خاصیت Parts-based بودن NMF را تایید می‌کند. الگوریتم یاد گرفته است که برای ساختن یک عدد، باید این اجزا را با هم جمع کند.



۴-۳ نتایج خوشه بندی K-means

برای مقایسه، الگوریتم K-means نیز با ۱۰ خوشه روی همان داده ها اجرا شد. مراکز خوشه های به دست آمده در شکل زیر نمایش داده شده اند:



شکل ۳-۲: مراکز خوشه های حاصل از الگوریتم K-means

تحلیل کیفی نتایج K-means

در شکل ۳-۲ مشاهده می شود که مراکز خوشه ها شبیه به میانگین محو شده ارقام کامل هستند. برخلاف NMF، در اینجا هر تصویر نماینده ی یک عدد کامل (مثل میانگین تمام صفرهای موجود در دیتاست) است. این تفاوت نشان دهنده ماهیت کل نگر روش K-means است.



۵-۳ مقایسه کمی (دقت خوشه بندی)

علاوه بر تحلیل بصری، دقت هر دو روش در دسته بندی داده های بدون برچسب اندازه گیری شد. از آنجا که خوشه بندی یک روش نظارت نشده است، برای محاسبه دقت از الگوریتم تطبیق اریب (Hungarian) استفاده شد تا بهترین تناظر بین خوشه ها و برچسب های واقعی پیدا شود.

نتایج حاصل از اجرا به شرح زیر است:

دقت	الگوریتم
0.3418	NMF
0.5323	K-means

همان طور که مشاهده می شود، روش K-means در این آزمایش عملکرد بهتری نسبت به روش دیگر داشته است. البته مزیت اصلی NMF در اینجا صرفاً دقت نیست، بلکه قابلیت تفسیرپذیری بالای ویژگی های استخراج شده است.



فصل چهارم: نتیجه گیری و پیشنهادات

۴-۱ نتیجه گیری

در این پژوهش، قابلیت های الگوریتم تجزیه نامنفی ماتریس (NMF) در استخراج ویژگی و خوشه بندی تصاویر ارقام دست نویس (MNIST) مورد بررسی قرار گرفت و عملکرد آن با روش کلاسیک K-means مقایسه شد. نتایج به دست آمده نشان می دهد که:

۱. نمایش بخش محور: مهم ترین ویژگی NMF، توانایی آن در تجزیه تصاویر به «اجزای سازنده» است. تصاویر پایه استخراج شده توسط NMF (مانند خطوط و کمان ها) نشان دادند که این الگوریتم به درستی ساختار درونی داده ها را یاد گرفته است. این ویژگی باعث می شود NMF در کاربردهایی که تفسیرپذیری مدل اهمیت دارد (مانند تشخیص چهره یا تحلیل متن)، برتری قابل توجهی داشته باشد.

۲. تفاوت با K-means: در مقابل، الگوریتم K-means رویکردی «کل نگر» داشت و مراکز خوشه ها را به صورت میانگین کلی ارقام نشان داد. اگرچه K-means از نظر محاسباتی سریع است، اما فاقد قابلیت تفکیک اجزای داده ها می باشد.

۳. دقت خوشه بندی: از نظر کمی، هر دو الگوریتم توانستند با دقتی قابل قبول (با توجه به ماهیت نظارت نشده مسئله) ارقام را دسته بندی کنند. این نشان می دهد که فضای ویژگی های کاهش یافته توسط NMF، اطلاعات کافی برای تمایز کلاس های مختلف را حفظ کرده است.



۲-۴ پیشنهادات برای کارهای آتی

برای بهبود نتایج و گسترش این پژوهش، موارد زیر پیشنهاد می شود:

- استفاده از نسخه های پیشرفته تر NMF مانند Sparse NMF که با اعمال جریمه تنکی (Sparsity) ، پایه هایی تمیزتر و تفکیک شده تر تولید می کند.
- اعمال الگوریتم بر روی داده های پیچیده تر مانند تصاویر چهره (مانند دیتاست Olivetti Faces) برای مشاهده دقیق تر قابلیت تجزیه اجزای صورت (چشم، بینی، لب).
- استفاده از روش های مقداردهی اولیه هوشمند به جای مقداردهی تصادفی برای افزایش سرعت همگرایی الگوریتم.



مراجع

- [1] Lee, D. D., & Seung, H. S. (1999). Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755), 788-791.
- [2] Aggarwal, C. C., & Reddy, C. K. (Eds.). (2013). *Data Clustering: Algorithms and Applications*. CRC Press.
- [3] LeCun, Y., Cortes, C., & Burges, C. J. (2010). MNIST handwritten digit database.
- [4] Hartigan, J. A., & Wong, M. A. (1979). Algorithm AS 136: A K-means clustering algorithm. *Applied Statistics*, 100-108.



پیوست

پروژه حاضر به زبان پایتون و با استفاده از کتابخانه های استاندارد داده کاوی پیاده سازی شده است. جهت مدیریت وابستگی ها و محیط مجازی از ابزار Poetry استفاده شده است.

دسترسی به کد منبع: نسخه کامل و قابل اجرای پروژه شامل فایل های پیکربندی (pyproject.toml) و تاریخچه تغییرات، در [مخزن گیت هاب](#) قابل دسترسی است.

در ادامه، کدهای اصلی فایل main.py برای بررسی منطق پیاده سازی آورده شده است:

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import fetch_openml
from sklearn.decomposition import NMF
from sklearn.cluster import KMeans
from sklearn.metrics import accuracy_score, confusion_matrix
from scipy.optimize import linear_sum_assignment

def get_clustering_accuracy(y_true, y_pred):
    """
    Since clustering labels have no specific order (permutation invariance),
    this function finds the best match between cluster labels and true labels
    using the Hungarian algorithm, then calculates accuracy.
    """
    cm = confusion_matrix(y_true, y_pred)
    # Find the best assignment
    row_ind, col_ind = linear_sum_assignment(cm.max() - cm)
    return cm[row_ind, col_ind].sum() / np.sum(cm)

def plot_gallery(title, images, n_col=5, n_row=2):
    """Helper function to plot images (basis patterns or centroids)"""
    plt.figure(figsize=(2. * n_col, 2.26 * n_row))
    plt.suptitle(title, size=16)
    for i, comp in enumerate(images):
        plt.subplot(n_row, n_col, i + 1)
        vmax = max(comp.max(), -comp.min())
        plt.imshow(comp.reshape(28, 28), cmap=plt.cm.gray,
                    interpolation='nearest',
                    vmin=-vmax, vmax=vmax)
        plt.xticks(())
        plt.yticks(())
    plt.subplots_adjust(0.01, 0.05, 0.99, 0.93, 0.04, 0.)
```



```
# 1. Load MNIST Dataset
print("Loading MNIST dataset...")
mnist = fetch_openml('mnist_784', version=1, as_frame=False)
X = mnist.data
y = mnist.target.astype(int)

# Normalize data (NMF requires non-negative data)
# Pixel values are 0-255, scaling them to 0-1
X = X / 255.0

print(f"Data Matrix (X) shape: {X.shape}")

# 2. Run NMF Algorithm
# We use 10 components since there are 10 digits (0-9)
n_components = 10
print(f"Running NMF with {n_components} components on the full dataset...")

nmf = NMF(n_components=n_components, init='random',
          random_state=42, max_iter=500)

# W is the membership matrix (referred to as F in the project description)
W = nmf.fit_transform(X)

# H is the basis matrix (referred to as G in the project description)
# The rows of H contain the basis images
H = nmf.components_

print("NMF execution completed.")

# 3. Visualize the columns of Matrix G (Basis Images)
plot_gallery("NMF Basis Images (Columns of G)", H)
plt.show()

# 4. Clustering using Membership Matrix F (here W)
# Each data point belongs to the cluster corresponding to the highest weight in W
y_pred_nmf = np.argmax(W, axis=1)

# Calculate NMF clustering accuracy
acc_nmf = get_clustering_accuracy(y, y_pred_nmf)
print(f"\nNMF Clustering Accuracy: {acc_nmf:.4f}")

# 5. Compare with K-means
print("\nRunning K-means on the original data...")
kmeans = KMeans(n_clusters=n_components, random_state=42, n_init=10)
y_pred_kmeans = kmeans.fit_predict(X)

# Calculate K-means accuracy
acc_kmeans = get_clustering_accuracy(y, y_pred_kmeans)
print(f"K-means Clustering Accuracy: {acc_kmeans:.4f}")

# 6. Visualize K-means Centroids
plot_gallery("K-means Centroids", kmeans.cluster_centers_)
```



```
plt.show()

# Final Analysis
print("-" * 30)
print("Analysis of Results:")
if acc_nmf > acc_kmeans:
    print("In this run, NMF performed better than K-means.")
else:
    print("In this run, K-means performed better (or close) to NMF.")

print("Qualitative Difference:")
print("- NMF basis images typically show 'parts' of digits (strokes, loops).")
print("- K-means centroids show the 'average' of whole digits.")
```