

# تمرین سری چهارم مباحثی در ریاضیات

## نگار یگانه

40113437

بخش اول ) در این بخش هدف استفاده از الگوریتم isomap برای تجسم فواصل بین شهر های ایران بوده است که به دوصورت دستی و استفاده از تابع آماده پیاده سازی شد. همچنین برای هر کدام از پیاده سازی ها با استفاده از دوران وانعکاس سعی کردیم به نقشه واقعی خروجی را نزدیک کنیم.

برای پیاده سازی دستی این الگوریتم به این صورت عمل کردیم که ابتدا گراف مجاورت را با پیدا کردن کا نزدیک ترین همسایه های هر شهر ایجاد کردیم سپس با استفاده از الگوریتم دایکسترا فاصله ژئودزیک بین همه زوج شهرها محاسبه می شود و ماتریس فواصل ژئودزیک را بدست می سازیم. سپس الگوریتم MDS بر روی ماتریس فواصل ژئودزیک به کار گرفته می شود. بدین صورت که ماتریس  $H$  را که کاربردش حذف اثر میانگین نقاط و آوردن داده ها به حول مبدأ می سازیم که این ماتریس برابر است با تفاضل ماتریس همانی با ماتریسی که تمام درایه هایش  $1/n$  است. سپس ماتریس  $B = -\frac{1}{2}H(\bar{A})H$  که  $\bar{A} = (a_{ij}^2)$  ساخته می شود. مقادیر و بردار ویژه های این ماتریس را حساب کرده و ضرب حاصل از جذر مقادیر ویژه و بردار ویژا را بدست می

آوریم. فقط مقادیر ویژه های بزرگتر و بردار ویژه های متناظر با آن ها مورد استفاده قرار می گیرند.

از ماتریس های چرخش برای قابل فهم تر شدن خروجی های هر دو isomap استفاده کردیم البته زوایه هایی که استفاده شد متفاوت است.

با توجه به خروجی های تصویری می توان دید که خروجی ها بهم نزدیک بوده اند و نزدیک نقشه واقعی ایران هستند. از سه معیار هم برای نشان دادن نزدیک بودن خروجی ها استفاده شد : معیار پروکروستس که یک روش آماری برای مقایسه ی دو مجموعه از نقاط است که ممکن است تحت تبدیل های هندسی مانند انتقال، چرخش یا مقیاس دهی قرار گرفته باشند و این تحلیل بهترین هم تراز ی بین دو مجموعه نقاط را پیدا می کند و میزان تفاوت آن ها را به صورت عددی نشان می دهد با توجه به مقدار پایین بدست آمده برای این معیار می توان متوجه شد اختلاف کمی در ساختار دو خروجی وجود دارد. همچنین مقادیر بالا بدست آمده برای هم بستگی مولفه ها نیز شباهت دو خروجی را نشان می دهد. معیار دیگه ای که استفاده کردیم حفظ همسایه های نزدیک بود که جهت بررسی این است که آیا نقاطی که در فضای اصلی به هم نزدیک بوده اند، در فضای جدید نیز نزدیک باقی مانده اند یا نه. برای هر نقطه،  $k$  همسایه ی نزدیک آن در هر دو پیاده سازی استخراج می شود. سپس بررسی می کنیم چه تعداد از این همسایه ها در هر دو تعبیه مشترک هستند. نهایتاً، درصد همسایگی حفظ شده به صورت میانگین محاسبه می شود که 89 درصد بود.

بخش دوم ) در این بخش، دو روش متداول کاهش بُعد یعنی PCA و t-SNE بر روی مجموعه داده ی MNIST اعمال شده اند و مورد ارزیابی قرار گرفتند.

مجموعه داده‌ی MNIST تک بعدی شامل ۷۸۴ ویژگی برای تصاویر ارقام 0 تا 9 است. برای کاهش زمان اجرا، تنها از ۴۰,۰۰۰ نمونه تصادفی استفاده شده است که به صورت متوازن از بین کلاس‌ها انتخاب شده‌اند.

خروجی‌های بصری و زمانی هر دو روش چاپ شده‌اند که واضحاً روش PCA زمانی بسیار کمتری داشت. اما در t-SNE، هر رقم تقریباً به صورت خوشه‌ای مشخص تفکیک شده است و در نتیجه خوشه بندی واضح تری دارد.

مقدار trustworthiness برای هر دو روش بسیار بالاست اما برای t-SNE به مقدار 0.04 بیشتر است. البته در نمونه گیری‌هایی با مقدار پایین تر این اختلاف بیشتر است.

### توضیح trustworthiness :

این معیار نشان می‌دهد که ساختار محلی داده‌ها تا چه حد حفظ شده است در واقع اندازه گیری می‌کند که آیا همسایگی‌های اصلی در فضای پایین بعد نیز حفظ شده‌اند یا خیر و اگر نقاطی که در فضای اصلی دور از هم بوده‌اند، بعد از کاهش بُعد به هم نزدیک شوند، معیار Trustworthiness کاهش پیدا می‌کند.

این معیار با فرمول زیر محاسبه می‌شود :

$$T(k) = 1 - \frac{2}{nk(2n - 3k - 1)} \sum_{i=1}^n \sum_{j \in \mathcal{N}_i^k} \max(0, (r(i, j) - k))$$

که در آن برای هر نمونه  $i$ ،  $N$  نشان‌دهنده  $k$  نزدیک‌ترین همسایه‌های نمونه  $i$  در فضای خروجی است، و هر نمونه  $j$ ،  $\text{rank}(i, j)$  نشان‌دهنده رتبه همسایگی نمونه  $j$  نسبت به

نمونه  $i$  در فضای ورودی است. به عبارت دیگر، هر همسایه غیرمنتظره (نزدیک‌ترین همسایه‌ای که در فضای خروجی ظاهر شده ولی در فضای ورودی در رتبه‌های پایین‌تر قرار دارد) متناسب با رتبه‌اش در فضای ورودی جریمه می‌شود. پس اگر Trustworthiness پایین باشد، به این معنی است که روش کاهش بعد نقاطی را که در واقع دور از هم بوده‌اند، به صورت اشتباه به هم نزدیک کرده است.

ادامه توضیحات پیاده سازی :

در ادامه t-SNE برای perplexity های مختلف اجرا شد و خروجی ها چاپ شد. با افزایش perplexity مقدار انحراف KL کاهش می یابد به طوری مقدار بهینه perplexity با توجه به انحراف KL ، 100 بدست اومد البته این بدین معنا نیست که با افزایش perplexity، به کاهش بعد بهتری از لحاظ بصری می رسیم چون می توان به صورت بصری دید که خوشه در حال ادغام شدن و بدتر شدن هستند.

Perplexity معیاری است برای ایجاد تعادل میان مقیاس های local و global در نمایش داده های کاهش یافته. مقدار این پارامتر مشخص می کند که t-SNE چه تعداد همسایه های نزدیک را در نظر بگیرد. هر چه مقدار آن کم باشد تمرکز بیشتر بر همسایه های خیلی نزدیک خواهد بود که باعث خوشه بندی محلی تر می شود و خوشه ها نزدیک تر بهم شکل میگیرند و ممکن است درک روابط را کمی دشوار کند و هرچه مقدارش بالا تر باشد الگوریتم تعداد بیشتری از نقاط را در نظر می گیرد مه باعث نمایش کلی تر از ساختار می شود ولی ممکن است جزئیات محلی از بین برود.

KL divergence معیاری است که اندازه گیری می کند چقدر توزیع احتمال داده های پربعد با توزیع داده های کمبعد تفاوت دارد. در واقع، الگوریتم سعی می کند -KL Divergence را کمینه کند تا بتواند ساختار داده ها را بهتر حفظ کند. اما لزوماً کاهش این مقدار به معنای نمایش بصری بهتر نیست.

منابع :

[trustworthiness — scikit-learn 1.6.1 documentation](#)

جزوه

[t-SNE: The effect of various perplexity values on the shape — scikit-learn 1.6.1 documentation](#)