

At. 1.1. f94

سُ بَدائِش وَاَعْرَاجِ دُرُكِي سَا

۱- استخراج ویژگی ها از تصاویر: استفاده مستقیم از ویژگی‌ها (مانند رنگ، موقعیت، اندازه، ...)

۵ ۱) کامن پیکینی دادها: در معاری خام، مثال، مقادیر زیادی از پیکل هستند که بسیاری از آن ها چکنات برای تحلیل  
ملی سطح بالا غیر ضروری یا غیر مفید باشد، با استخراج و حذف پیکینی داده ها کامن و یا به طور معمول جداول  
از ویژگی ملی مختص و سترز نیست. (آید)

۲) تالیف بر سرِ رُگی (بازی) هم یعنی: بدلِ می (بعضی) مانند ۷۶۶۱۶ با اختلاص از لایِ می (لایِ لوس) و آگوه و سافتر می

هم بعضی را تسلی می کنند این و آن می ما اغلب مواردی هستند بعد از آن، باعث آسایش و سایر اطلاعاتی در باره ایشان

10. سنہ: ایس۔ مزائیدہ الکومرٹسلی Machina Learning اجازت سے تا جلی ریلل سلی خام برقی سلی برجستہ تر و حندلرتہ

(۳) استفاده از ویژگی های استخراج شده برای اسراف مختلف ، با استخراج ویژگی های توانمند از آن ، با در نظر گرفتن مختلف  
مانند : طبعه بنویس ، کلاس بعد از ظهر ، استفاده کنید

15 (4) با عبارتی نسبت به تغییرات و متنرنگی ها: در رنگی سلی استوار جمع شده اغلب در برابر تغییرات جزئی در تقاضای مقولست  
بیشتر در این سبیل ها کمک می کند تا از در رنگی سلی استوار برای تصمیم گیری استفاده کنند.

۵) ارتباط با دولت قبلی: استخراج و یکی با بر اساس محل بلی. حالت ۶۶۶۶ که با تائید غیبی از طرفین آمدنی  
نعمت شطرنج، ب. محابین، کلان ران، دعد از دولت که این محل از قبل دارند استفاده کنند. این یک راه کار آمد بلی استغله  
از و یکی بلی استخراج دولت.

در این استخراج برتری ما از تصویب شما الحاح دارد تا با داره ملی سرگرم، قتل دولت و منصب و برتری اگر کنید

۲- استخراج رگي از تقاویر بي زائید اليسی در زمان تقویر بيانی ماسک الکتریکی آن اطلاعات معناداری

و قابل استفاده از تصاویر خام استخراج می شود. این فرآیند به کمک الگوریتم های یکپارچه میسر می شود و دستاورد اصلی کارها  
داده های خام یکپارچه با جبره ای لازم و ویژگی های مفید تر و غیر تکرار کننده و ترکیب های زیاده نداشتن در نتیجه اکتفا به بافت ها  
لبه ها و رنگ ها و سایر اطلاعات بهتری هم باشند.

۲۵. ۳. تکیه ایج یون استراجه و رنگی ما انزعا و یون

① استخراج ویژگی‌های لانه‌گوش، Convolutional Feature Extraction، این روش برای به دست آوردن ویژگی‌های لانه‌گوش (CNN)







«پایانه سازی» خوشه یابی»

۴- Clustering یک روش یادگیری بدون نظارت است که هدف آن گروه بندی اشیاء داده ها بر اساس ویژگی های مشابه آن ها است. در روش محبوب در Clustering، K-Means و DBSCAN است. در ادامه توضیحی را می توانیم در مورد آن ها بگوییم.

«K-Means»

K-Means یکی از روش های پرکاربرد خوشه بندی است که در آن هدف این است که داده ها به K خوشه تقسیم شوند. الگوریتم به طور مکرر به صورت زیر عمل می کند:

- ۱- انتخاب اولیه مراکز خوشه ها: الگوریتم با انتخاب K نقطه به عنوان مراکز خوشه شروع می شود.
- ۲- اختصاص داده ها: هر نقطه از داده ها به نزدیک ترین مرکز خوشه اختصاص داده می شود.
- ۳- تعیین مراکز خوشه ها: مراکز خوشه ها بر اساس میانگین نقاط تخصیص یافته به هر خوشه به روز می شوند.
- ۴- تکرار تا زمانی که خوشه ها تثبیت شوند: این فرآیند تکرار می شود تا زمانی که تغییرات مراکز خوشه ها صاف باشد.

- \* مزایا: ① سادگی: K-Means یک الگوریتم نسبتاً ساده و سریع است که به راحتی قابل اجرا است.
- ② کارایی: به ویژه برای داده های بزرگ، K-Means به دلیل سادگی و کارایی آن محبوب است.
- ③ قابلیت تفسیر: خوشه های ایجاد شده توسط K-Means به راحتی قابل تفسیر هستند.
- \* معایب: ① نیاز به تعیین تعداد خوشه ها (K): نیاز به مشخص سازی تعداد خوشه ها دارد که این امر گاهی اوقات از روشی ساخته شده نباشد.

15

- ② حساسیت به مقدار اولیه: نتایج K-Means می تواند به انتخاب اولیه مراکز اولیه خوشه ها حساس باشد.
- ③ فرضیه ساختار کروی: K-Means فرض می کند که خوشه ها به صورت کروی هستند، که ممکن است در همه موارد صادق نباشد.
- ④ حساسیت به داده های Noisy: داده های Noisy می تواند به طور قابل توجهی روی نتایج K-Means تأثیر بگذارد.

«DBSCAN»

- DBSCAN یک الگوریتم خوشه بندی مبتنی بر چگالی است. در این روش، Clusters ما بر اساس متعلق یا چگالی بالا شکل می گیرند. الگوریتم به طور مکرر به صورت زیر عمل می کند:
- ۱- شناسایی نقاط مرکزی: نقاطی که دارای تعداد معینی از همسایگان در یک ناحیه مشخص (E) هستند به عنوان نقاط مرکزی شناخته می شوند.
  - ۲- گسترش خوشه ها: خوشه ها با افزودن نقاطی که به همسایگان آن ها متعلق می دارند گسترش می یابند.
  - ۳- نقاط Noisy: داده های که به هیچ خوشه ای متعلق ندارند به عنوان داده های Noisy شناسایی می شوند.

25

- مثالی: ① عدم نیاز به تعداد خوشه‌ها از پیش: DBSCAN نیاز به تعیین تعداد خوشه‌ها ندارد.
- ② مدیریت داده‌های Noise: DBSCAN توانایی حذف داده‌های Noise را از خوشه‌ها دارد.
- ③ قابلیت خوشه‌بندی: DBSCAN قادر به شناسایی خوشه‌های نامنظم و غیر استاندارد است.
- ④ حساسیت: DBSCAN به تنظیم پارامتر  $\epsilon$  (شماره) و تعداد نقاط محلی حساس است. نیاز دارد که مشخص است ابتدا مشخص باشد.
- ⑤ حساسیت به نویز: اگر نویز داده‌ها به طور قابل توجهی متفاوت باشد، ممکن است DBSCAN نتواند خوشه‌های مناسب را شناسایی کند.
- ⑥ کارایی: برای داده‌های بسیار بزرگ، برای داده‌های بسیار بزرگ، PBSCAN ممکن است کمتر از K-Means باشد.
- انتخاب بین DBSCAN و K-Means:
- انتخاب بین این دو روش بستگی به طبیعت داده‌ها و نیازهای پروژه دارد. K-Means برای داده‌های با ساختار گردی و خوشه‌های با اندازه‌های بسیار مناسب است. در مقابل DBSCAN برای داده‌های با ساختار نامنظم و خوشه‌های با شکل‌های نامنتظم انتخاب بهتری دارد.

برای پیدا کردن تعداد مناسب خوشه‌ها (K) در K-Means روش‌های مختلفی وجود دارد: یکی از رایج‌ترین روش‌ها Elbow است. این روش شامل اجرای K-Means با تعدادهای مختلف از خوشه‌ها است و سپس بررسی نمودار تغییرات و پیدا کردن نقطه‌ای که تغییرات ناگهانی در مقدار هزینه (Cost) رخ دهد. در روش Elbow، شما مقادیر مختلف K را در الگوریتم K-Means قرار می‌دهید و سپس مقدار SSE (Sum of Squared Errors) را برای هر مقدار K محاسبه می‌کنید. وقتی که کاهش SSE به طور قابل توجهی کم می‌شود، این نقطه ممکن است مقدار مناسب K باشد.

★ من در اکثر استفاده از الگوریتم Silhouette، من و سیلوئت را هم اضافه کردم که بهتر بتوانم K مناسب را انتخاب کنم. حل این مسئله سیلوئت چیست؟

سیلوئت یک معیار برای ارزیابی کیفیت خوشه‌بندی است که نشان می‌دهد چگونه داده‌های موجود در یک خوشه به هم نزدیک هستند و چقدر از خوشه‌های دیگر فاصله دارند. برای هر یک از این روش‌ها یک عدد بین ۰ تا ۱ محاسبه می‌شود. اگر عدد به ۱ نزدیک باشد، خوشه‌ها از یکدیگر جدا شده‌اند.

نقشه سیلوئت:  $s = \frac{b-a}{\max(a,b)}$ ، جاییکه  $a$  فاصله هر یک از نقاط خوشه‌ها از یکدیگر است و  $b$  فاصله هر یک از نقاط خوشه‌ها از نزدیک‌ترین خوشه دیگر.

این روش بین ۰ و ۱ است. سیلوئت بالا (نزدیک ۱) نشان می‌دهد که خوشه‌ها به طور کامل از خوشه‌های دیگر جدا شده‌اند.



از یکدیگر جدا شده اند. نو ترکیب به ۵ شکل دسته این است که نمونه در هر خوشه یک رنگ است (نمونه این است که نمونه به خوشه دیگری نزدیک تر است. کاربرد نو سیلونت به ارزیابی کیفیت خوشه بندی این نو شکل و ده که چقدر خوشه ها از یکدیگر متجانس هستند و چقدر داده ها در خوشه خود پاینده هستند. انتخاب تعداد خوشه های مناسب و انتخاب این نو سیلونت برای پیدا کردن تعداد بهینه خوشه ها استفاده کرد. مقدار  $k$  که بیشترین نو سیلونت را دارد ممکن است تعداد مناسب خوشه ها باشد.

### homogeneity

معیار یکپارچگی یا Homogeneity یکی از معیارهای است که برای ارزیابی کیفیت خوشه بندی مورد استفاده قرار می گیرد. این معیار به ما می گوید که چقدر خوشه ها یکپارچگی هستند و تا چه حد هر خوشه شامل عناصر یکدست است (یا دسته بندی) مناسب است. مفهوم homogeneity: معیار یکپارچگی نشان می دهد که هر خوشه چقدر شامل داده هایی است که فقط از یک دسته بندی یکدست می آیند. به عبارت دیگر، این معیار می گوید که آیا همه عناصر در یک خوشه متعلق به یک کلاس هستند یا خیر. نحوه محاسبه یکپارچگی: معیار یکپارچگی با استفاده از اطلاعات متقابل (Mutual Information) بین خوشه ها و کلاس های واقعی محاسبه می شود و پس با اطلاعات کلاس های واقعی (نتیجه یک عدد بین ۰ و ۱ است). \* یک مقدار نزدیک به ۱ نشان می دهد که هر خوشه شامل داده هایی از یک کلاس است (خوشه ها یکپارچگی هستند). \* یک مقدار نزدیک به ۰ نشان می دهد که خوشه ها حاوی داده هایی از کلاس های مختلف هستند (خوشه ها متنوع هستند). کاربرد homogeneity برای ارزیابی خوشه بندی: اگر داده های شما دارای کلاس های خاصی هستند (حتی اگر این کلاس ها برای مدل به عنوان برچسب های آموزشی نبوده اند) معیار یکپارچگی می تواند به شما کمک کند تا کیفیت خوشه بندی را ارزیابی کنید. \* برای استفاده از این معیار، به داده های با برچسب نیاز داریم تا بتوانیم یکپارچگی را اندازه گیری کنیم. این برچسب ها کلاس های واقعی و اقربا داده ها را نشان می دهد. محدودیت ها: معیار یکپارچگی متوازن نیست معنی است که داده های شما دارای برچسب های واقعی باشند. این معیار برای داده های بدون برچسب کاربرد ندارد.

### DBSCAN clustering

برای پیدا کردن ۴ در این الگوریتم ۲ روش وجود دارد که اولین روشی که ما به دست می آوریم  $k$  که از نو سیلونت استفاده کنیم و در این حالت بیشترین نو سیلونت در یک  $k$  را می بینیم و در آن روش  $k$  را انتخاب می کنیم.

دوم: رسم نمودار فاصله بین نقاط نزدیک ترین همایه: با رسم نمودار فاصله بین نقاط نزدیک ترین همایه و تعیین میزان تیرات در فواصل بین نقاط را سه کنید. به طور مثال، یک به یک یا تیر ناگهانی در فاصله مانع از همایه بر روی یک است.

۶- نتیجه فرقی K-Means با K-Bray ۱-۸ به صورت آن که تعدادی خوب در یک دسته مناسب قرار گرفته اند و همین cluster یکی در میان هستند و به صورت خودی به نقلی آید که آن تولد این مارا cluster بنویسند البته به نقلی رسد که انتخاب K در این حالت تأثیر مهمی داشته است. در نمودار DBSCAN او تولد فقط چند cluster کم را تعیین داده و بر حسب visualization و به عمل گردان. دلایل هم می باشد به E دارد.

«کلاس بعد»

تبدیل مولفه های اصلی (PCA) یک تکنیک برای کاهش ابعاد داده ها می باشد که برای کلاسیک بعد داده ها می باشد و تبدیل ویژگی ها و صورت های استفاده می شود. PCA با شناسایی مولفه های اصلی که بیشترین واریانس را در داده ها را توضیح می دهند کار می کند. این تکنیک به ویژه زمانی مفید است که داده ها دارای مقیاس های مختلف و ویژگی ها باشند و می توانیم آن ها را به تعداد کمی از ویژگی های اصلی کاهش دهیم.

نقشه عملکرد PCA: PCA به این صورت عمل می کند که مولفه های اصلی (شناختی) می کند به طور خطی ترکیبی از ویژگی های اصلی هستند و بیشترین واریانس را دارند در اینجا چند هم کلاسیک که PCA انجام می دهد آورده شده است. ۱- مرکزیت داده ها: داده ها معمولاً مرکزیت می یابند یعنی میانگین هر ویژگی از داده ها می شود تا داده ها حول محور صفر توزیع شوند.

۲- محاسبه ماتریس کواریانس: ماتریس کواریانس نشان می دهد که چگونه ویژگی ها نسبت به یکدیگر تغییر می کنند. PCA با محاسبه این ماتریس شروع می کند.

۳- محاسبه بردارها و مقادیر ویژه: PCA بردارهای ویژه و مقادیر ویژه ماتریس کواریانس را محاسبه می کند این بردارها جهت های مولفه های اصلی هستند و مقادیر ویژه میزان واریانس توضیح داده شده توسط هر مولفه اصلی را نشان می دهد.

۴- انتخاب مولفه های اصلی: مولفه های اصلی به ترتیب کمترین واریانس مرتب می شوند. با انتخاب میز مولفه اول، می توانیم بیشترین واریانس داده ها را توضیح دهیم.

۵- تبدیل داده ها به مولفه های اصلی: با استفاده از مولفه های اصلی انتخاب شده، داده ها به فضای جدیدی با ابعاد کم تبدیل می شوند.



چرا از PCA استفاده می‌کنیم؟ PCA یکی کلاس ابعاد داده‌ها بعد از دست دادن اطلاعات زیاد است این کار برای  
مختلف دارد ① کلاسی بی‌بیمبگی، با کلاسی تعاد و ترکی ما، حسابات سریع تر و آسان تر می‌شوند  
② معود سازی، PCA و ترکی معود سازی داده‌ها در فضای دو بعدی یا سه بعدی ممکن است  
③ کلاسی نویز، حذف مولفه‌هایی که دارای نویز کمتری دارند ممکن است به کلاسی نویز در داده‌ها کمک کند.

محدودیت‌های PCA

فقدان تغییر مرکزی، مولفه‌های اصلی ممکن است تغییر مرکزی بعضی اطلاعات باشد زیرا ترکیب فضای از ویژگی‌های اصلی هستند  
فرض بر فضای بدون PCA، و اساسی فرض فضای بدون روابط بین ویژگی‌ها عمل کند اگر در اینجا فرض باشد این روش  
ممکن است سوئیچ باشد  
حسیت به مقیاس، اگر ویژگی‌ها در مقیاس‌های مختلفی باشند نتایج PCA ممکن است تحت تأثیر قرار گیرد.

در داده‌ها بعد از انجام PCA می‌بینیم که بسیار بهتر شده‌اند cluster می‌انجام شده و این به ما کمک کرده است.

مزایای و معایب

- 1- در صورتی که قبل از این دو معیار صحبت کرده‌ایم.
- 4- می‌تواند برای درمان نیز به‌یونیت انجام داده‌ام
- ط- برای بهبود عملکرد مدل ما در سبک‌های مختلف و می‌تواند از روش‌ها و استراتژی‌های مختلفی استفاده کرد در اینجا به چند  
روش هم برای ارزیابی عملکرد مدل می‌توانیم استفاده کنیم و در زیر
- ① روشی برای ارزیابی داده‌ها، می‌تواند برای داده‌ها و سبک‌های مختلف و می‌تواند برای داده‌ها و سبک‌های مختلف
- ⑤ روشی برای ارزیابی داده‌ها، می‌تواند برای داده‌ها و سبک‌های مختلف و می‌تواند برای داده‌ها و سبک‌های مختلف
- ⑥ انتخاب و تنظیم مدل، تنظیم یادگیری، انتخاب مدل مناسب، استفاده از ترکیب مدل‌ها
- ④ ارزیابی و اعتبارسنجی
- ⑤ ارزیابی داده‌ها
- ⑥ استفاده از مدل‌های پیچیده‌تر