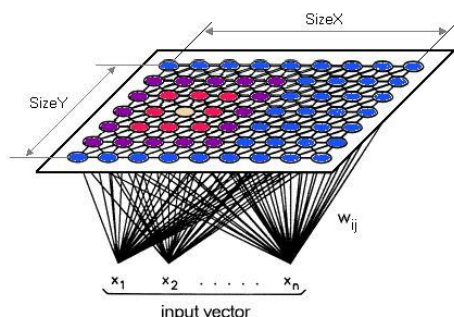


تمرین دوم درس شبکه های عصبی

شبکه خود سازمانده

شبکه خودسازمانده یک شبکه دو لایه است که لایه اول دارای تعداد نورون متناسب با ابعاد داده ورودی است و لایه خروجی یک لایه است که میتوان آنرا به صورت خطی، دوبعدی و یا چند بعدی نمایش داد که به منظور یادگیری بدون نظارت و برای کارهایی مثل دسته بندی، کاهش ابعاد داده، بصری سازی اطلاعات و ... میتوان به کار برد.



شکل ۱: شبکه خودسازمانده کوهونن

در این شبکه طی یک فرآیند یادگیری رقابتی، نورون برنده مشخص می شود و وزن نورون بنده و نورون های همسایه آن، بروز رسانی می شوند. در ابتدای فرآیند سعی بر این است که شعاع همسایگی زیاد در نظر گرفته شود و در دفعات بعدی این شعاع با یک تابع خطی و یا گوسی و ... کاهش یابد. همچنین برای نرخ یادگیری نیز از همین قانون استفاده می کنیم و در ابتدا نرخ بالایی را در نظر می گیریم و سپس آن را در هر دفعه کاهش بر اساس یک رابطه خطی و یا یک رابطه گوسی کاهش می دهیم. رابطه زیر برای به روز رسانی وزن ها در نظر گرفته می شود.

$$W'_j = W_j + \beta NS [X - W_j] \quad (1)$$

که در رابطه فوق NS تابع همسایگی و β نرخ یادگیری است که در زیر نمونه های خطی این دو نشان داده شده است.

$$\beta_t = \beta_0 \left(1 - \frac{t}{T}\right)$$

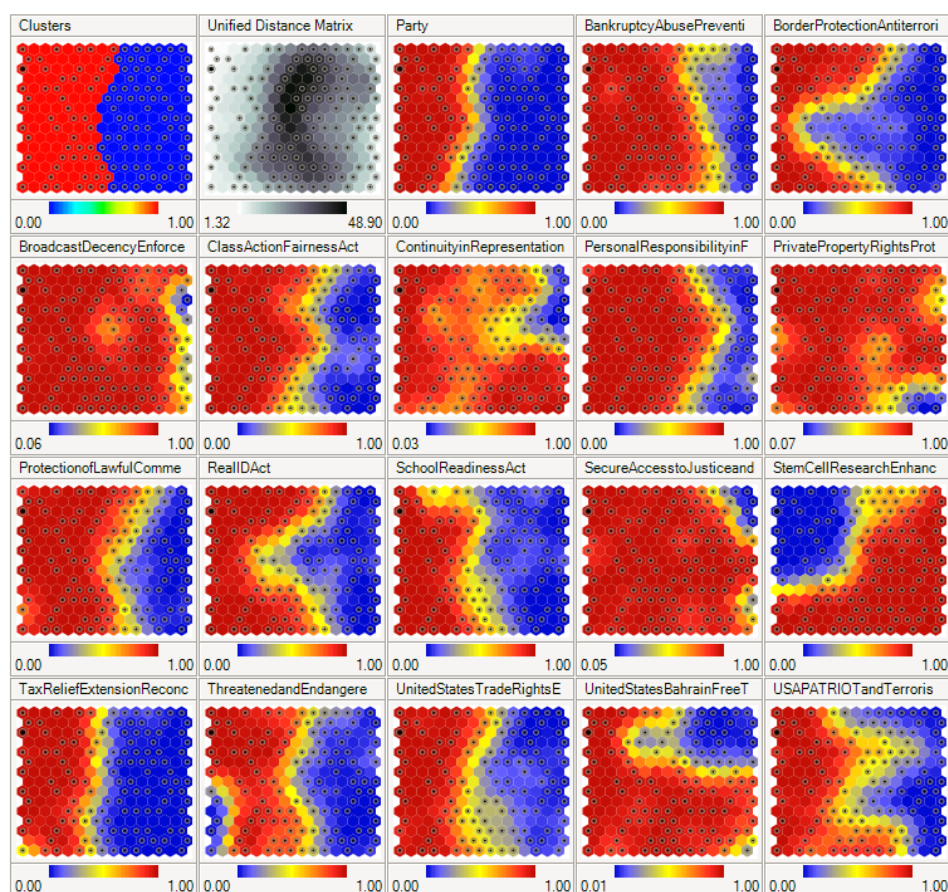
$$NS = e^{-\frac{d_{ij}^2}{2\sigma^2}}$$

$$\sigma_t = \sigma_0 \left(1 - \frac{t}{T}\right)$$

در رابطه فوق T تعداد دفعات تکرار است.

یکی از ویژگی های شبکه SOM این است که داده های با ابعاد بالا را می توان به یک محور یک بعدی و یا یک صفحه دوبعدی و یا سه بعدی به وسیله یک نگاشت غیرخطی تصویر نمود. از این رو یکی از کاربردهای این شبکه کاهش ابعاد داده است که از این طریق می توان تحلیل خوبی از داده های با حجم بالا به وسیله بصری سازی پیدا کرد.

SOM یک الگوریتم خوشه بندی بدون نظارت است. به این صورت که نمونه های مشابه را در نقشه ویژگی نزدیک تر نشان می دهد (یعنی نمونه های مشابه را به گره هایی که به هم نزدیک تر هستند نگاشت می کند). بنابراین آنچه که SOM در واقع تولید می کند، نگاشتی از فضای ورودی X به فضای کاهش یافته Y است (متداول ترین آنها یک شبکه دو بعدی است که Y را یک فضای دو بعدی می کند). حال از این نقشه نقشه های ویژگی به دست آمده می توانیم برای تحلیل و بصری سازی داده ها استفاده کنیم.



شکل ۲: نقشه های ویژگی ارزیابی آرا نمایندگان کنگره آمریکا

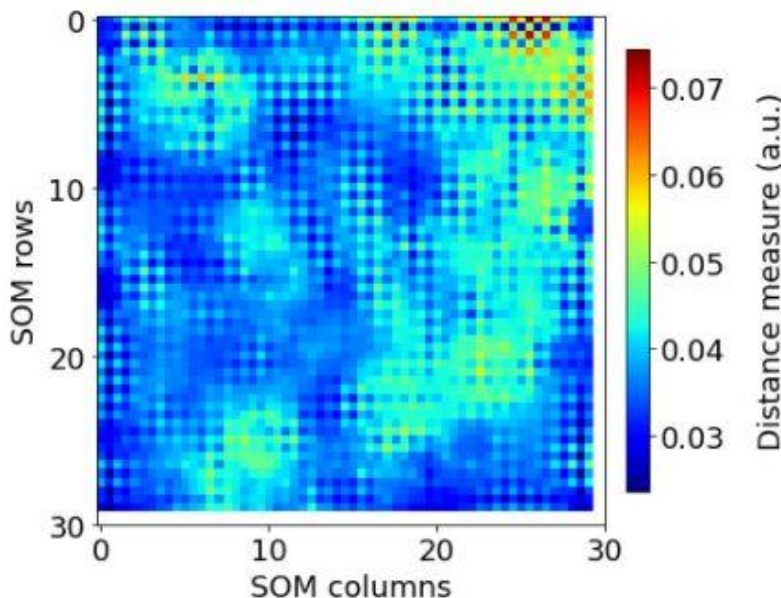
در شکل بالا شبکه خود سازمانده روی مجموعه داده رای های مثبت و منفی نمایندگان کنگره به لایحه های مختلف نمایش داده شده است. نمودار اول خوشه بندی را هنگامی که داده ها به دو خوشه تقسیم می شوند نشان می دهد. نمودار دوم میانگین فاصله را با همسایگان نشان می دهد: فواصل بزرگتر تاریک تر هستند. طرح سوم عضویت حزب جمهوری خواه (قرمز) یا دموکرات (آبی) را پیش بینی می کند. نمودارهای دیگر هر کدام نقشه حاصل را با مقادیر پیش بینی شده در یک بعد ورودی پوشش می دهند: قرمز به معنای رای «بله» پیش بینی شده در آن لایحه، آبی به معنای رای «نه» است. در واقع می توان بررسی کرد که هر طرح به چه میزانی در جدا کردن نمایندگان دو حزب از هم تاثیر دارد.

آماده سازی مجموعه داده‌ها

مجموعه داده‌ها شامل ۸۰۰۰۰ نمونه با ابعاد ۴۷۲۳۶ می‌باشد. هدف ما از این تمرین استفاده از شبکه خودسازمانده کوهونن برای کاهش ابعاد داده‌ها و بصری سازی آن است. فرم داده‌ها sparse matrix است که تنها درایه‌های غیر صفر را ذخیره می‌کند و بخش زیادی از مجموعه داده‌ها صفر است. برای اعمال som روی مجموعه داده‌ها نیاز است که آن‌ها را به numpy array تبدیل کنیم که به این ترتیب علاوه بر مقادیر غیر صفر باید حجم زیادی صفر نیز اضافه شود. به دلیل اینکه حجم RAM در colab جوابگوی پردازش این حجم از دیتا نیست ما به کمک resample از میانه نمونه‌های داده شده بخش از آن‌ها را انتخاب می‌کنیم.

SUSI

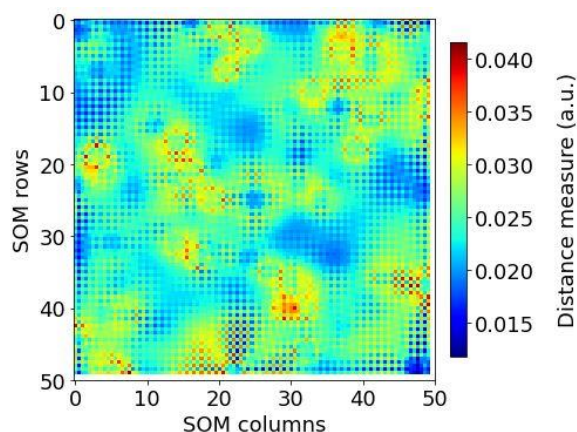
برای اعمال som از کتابخانه susi استفاده می‌کنیم. این کتابخانه به نسبت سایر کتابخانه‌ها همچون minisom عملکرد مناسب‌تری دارد و برای اجرا به RAM و زمان کمتری نیاز دارد. پس از نصب این کتابخانه، شبکه som را با ابعاد 30*30*47236 ایجاد می‌کنیم. تعداد تکرارهای را برابر ۱۰۰۰ و نرم یادگیری اولیه را ۰,۷ در نظر می‌گیریم. برای اعمال همسایگی نیز از تابع گوسی استفاده می‌کنیم. در شکل زیر نمودار u-matrix را مشاهده می‌کنید.



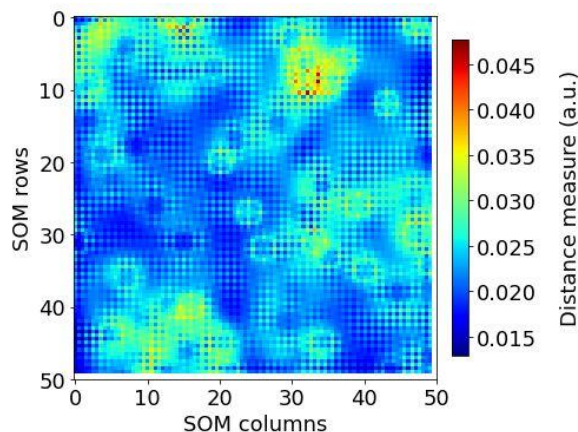
شکل ۳: نقشه u-matrix برای 30*30 SOM

در شکل بالا نورون‌های آبی نزدیک به هم هر کدام یک منطقه را تشکیل می‌دهند. در واقع ناحیه‌های زرد رنگ را می‌توان مرزهای جداکننده قسمت‌های مختلف در نظر گرفت. نورون‌های آبی رنگ نشان‌دهنده فاصل کم بین بردار وزن‌های این نورون‌ها با نورون‌های همسایه آن است از این موضوع می‌توان نتیجه گرفت تراکم داده‌ها در قسمت آبی رنگ بیشتر است. به بیان دیگر شبکه SOM یک تصویر از داده‌ها می‌دهد که داده‌های نزدیک به هم را به نورون‌های مجاور آبی رنگ تصویر می‌کند. و داده‌های با فاصله زیاد را به نورون‌های متمایل به قرمز تصویر می‌کند. دقت شود حتی اگر داده‌ها به نورون‌های مجاور تصویر شوند اما به فرض یک نورون آبی و دیگری قرمز باشد. نشان‌دهنده فاصله زیاد داده‌ها از هم می‌باشد. در واقع با این نحوه رنگ آمیزی، شبکه SOM چگالی توزیع داده‌ها را به صورت گسسته به نمایش می‌گذارد.

در بررسی یک نقشه دیگر، شبکه som را با ابعاد $50 \times 50 \times 47236$ ایجاد می‌کنیم. تعداد تکرارهای را برابر ۱۰۰۰ و نرم یادگیری اولیه را ۰,۷ در نظر می‌گیریم. و این بار ۷۰۰۰ نمونه را به صورت تصادفی انتخاب می‌کنیم تا RAM داده شده توسط colab جواب گوی حجم محاسبات باشد. برای اعمال همسایگی نیز از تابع گوسی استفاده می‌کنیم. همچنین، این نقشه را بار دیگر با نرخ یادگیری اولیه ۰,۵ به فرض ثابت بودن سایر پارامترها اجرا کردیم. در شکل زیر نمودارهای u-matrix را مشاهده می‌کنید.



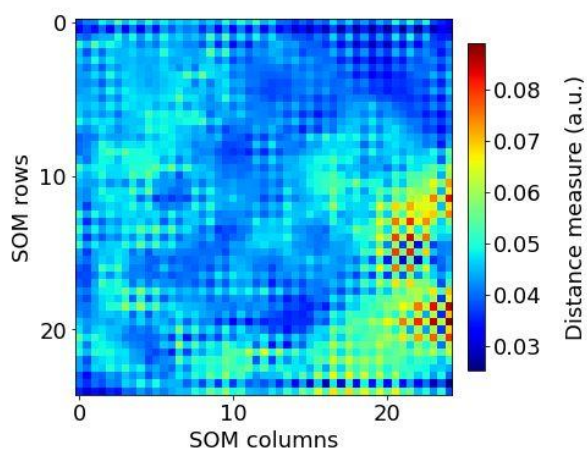
شکل ۵: نقشه u-matrix برای SOM 50×50 با نرخ یادگیری اولیه ۰.۷



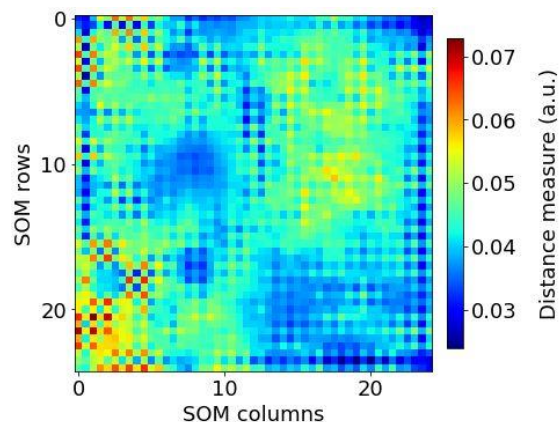
شکل ۴: نقشه u-matrix برای SOM 50×50 با نرخ یادگیری اولیه ۰.۵

همان طور که مشاهده می‌کنید نرخ یادگیری اولیه تاثیر زیادی در عملکرد الگوریتم دارد. و اجرای الگوریتم با نرخ ۰,۵ باعث متراکم شدن نقشه شده است. در شکل ۳ به خوبی می‌توان داده‌های تصویر شده به نودهای مجاور آبی رنگ دسته‌بندی کرد و از نودهای زرد رنگ نیز به عنوان مرز استفاده کرد.

همچنین در زیر، شبکه som را با ابعاد $25 \times 25 \times 47236$ ایجاد می‌کنیم. تعداد تکرارهای را برابر ۱۰۰۰ و نرخ یادگیری اولیه را ۰,۷ و ۰,۵ در نظر می‌گیریم. و این بار ۱۲۰۰۰ نمونه را به صورت تصادفی انتخاب می‌کنیم تا RAM داده شده توسط colab جواب گوی حجم محاسبات باشد. برای اعمال همسایگی نیز از تابع گوسی استفاده می‌کنیم. در شکل زیر نمودار u-matrix را مشاهده می‌کنید.



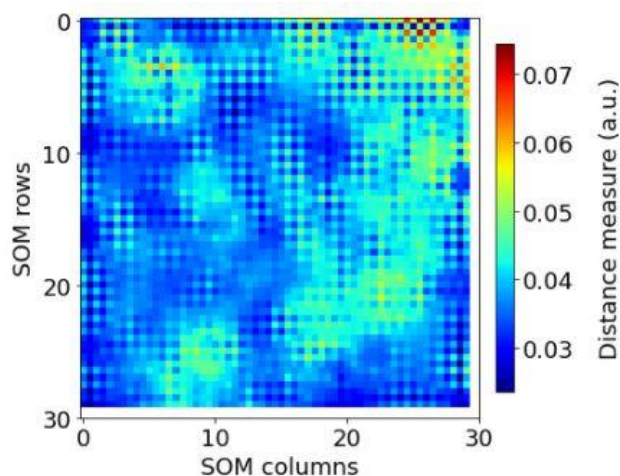
شکل ۷: نقشه u-matrix برای SOM 25×25 با نرخ یادگیری اولیه ۰.۷



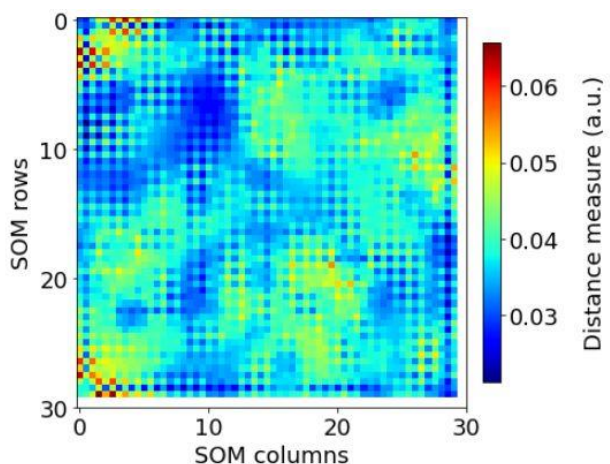
شکل ۶: نقشه u-matrix برای SOM 25×25 با نرخ یادگیری اولیه ۰.۵

برای نقشه 25*25 بر خلاف نقشه 50*50 نرخ یادگیری 0.7 عملکرد بهتری نسبت به مقدار 0.5 دارد.

همچنین در زیر، شبکه som را با ابعاد 30*30*47236 ایجاد می‌کنیم. تعداد تکرارهای را برابر ۱۰۰۰ و نرخ یادگیری اولیه را ۰.۷ و ۰.۵ در نظر می‌گیریم. و این بار ۱۰۰۰۰ نمونه را به صورت تصادفی انتخاب می‌کنیم. برای اعمال همسایگی نیز از تابع گوسی استفاده می‌کنیم. در شکل زیر نمودار u-matrix را مشاهده می‌کنید.



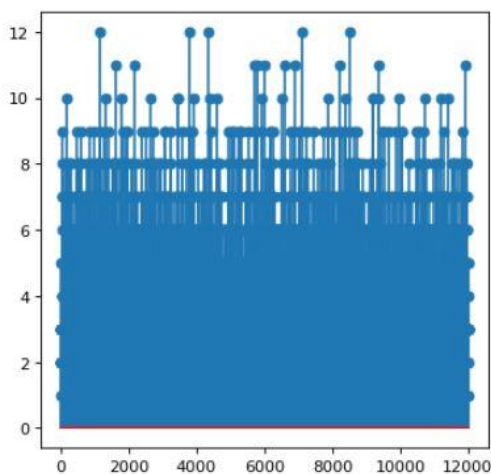
شکل ۹: نقشه u-matrix برای SOM 30*30 با نرخ یادگیری اولیه ۰.۷



شکل ۸: نقشه u-matrix برای SOM 30*30 با نرخ یادگیری اولیه 0.5

خوشه بندی

مجموعه داده‌های داده شده multi-label هستند و هر داده به چند کلاس از ۱۰۳ کلاس موجود متعلق است. به عنوان مثال ممکن است یک خبر یا یک داستان هم در کلاس موضوعات فرهنگی باشد و هم در کلاس موضوعات بین‌المللی. در زیر نموداری از تعداد کلاس‌هایی که هر داده به آن متعلق است آمده است.

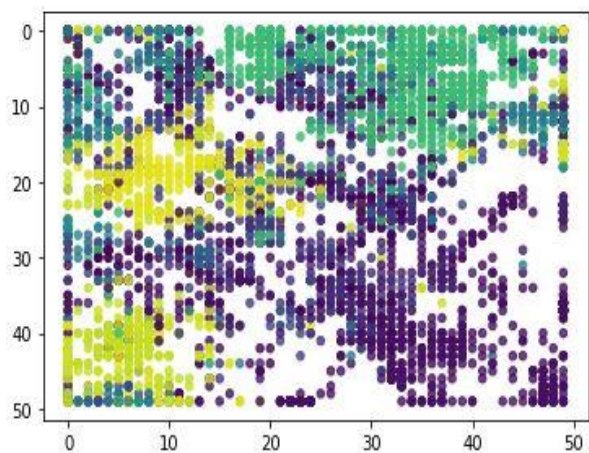


شکل ۱۰: نمایش تعداد کلاس‌های هر نمونه

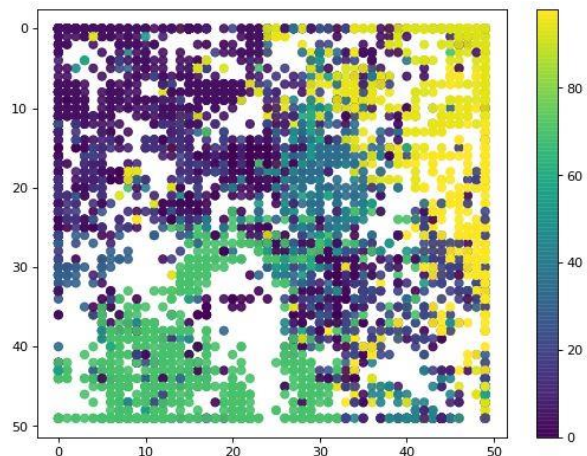
برای خوشه بندی از کتابخانه `sus` استفاده کردیم. تابع `get_cluster` در این کتابخانه به این صورت عمل می کند که هر داده را با توجه به تنسور های به وجود آمده در نقشه `som` هر یک از ورودی ها را به یک نورون از نقشه نگاشت می کند بر مبنای ترتیب داده های داده شده نورون های مورد نظر را به عنوان خروجی می دهد. در واقع هر نورون را یک خوشه در نظر می گیریم. به کمک `som` به دست آمده در قسمت قبل داده ها را به نورون، خوشه ها، نگاشت می کنیم. برای نمایش خوشه بندی و رنگ آمیزی خوشه ها به کمک `label` ها باید در نظر داشته باشیم که داده ها `multi-label` هستند و نمایش تمام `label` های داده ممکن نیست. ما برای نمایش خوشه ها با `label` ها آمدیم و برای هر داده اولین `label` داده را در نظر گرفتیم و به کمک `scatter` مجموع آن ها را نمایش دادیم. در زیر نمایش های خوشه بندی `som` های با اولین `label` داده ها در بخش قبل آورده شده است.

تحلیل توزیع داده ها

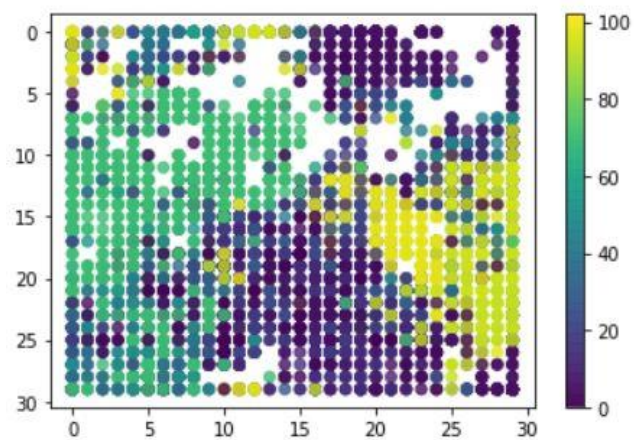
همان طور که مشاهده می کنید، عمل خوشه بندی به خوبی انجام شده و داده ها با ابعاد بسیار بالا به خوبی توسط `som` تفکیک شده است. حال به تحلیل نمودار نقشه `som 30*30` با نرخ یادگیری `۰,۷` که نمودار `u-matrix` آن نسبت به بقیه حالت ها داده ها را متراکم تر کرده است می پردازیم. همان طور که از همه شکل ها از جمله شکل `۱۴` که مورد نظر ماست آشکار است می توان با `4` خوشه داده ها را خوشه بندی کرد که با `4` رنگ بنفش، آبی، سبز و زرد نمایش داده شده است. در رابطه با جدایی پذیری داده ها در شکل `۱۴`، تا حد قابل قبولی جدایی پذیر هستند به جز داده خوشه آبی رنگ که با خوشه سبز رنگ در بخش هایی هم پوشانی دارند و جداپذیر نیست. برای پیچیدگی داده ها می توان گفت با توجه به ابعاد بسیار بالای داده ها، نقشه شکل نسبتاً ساده ایی از داده ها داده است و می توان گفت داده ها پیچیدگی زیادی همچون اشکال مارپیچی و `XOR` ندارند و با اعمال یک `MLP` به خروجی `som` آن ها را می توان دسته بندی کرد. داده های نویزی به داده هایی گفته می شود که در هیچ خوشه ایی جای نگیرند. با توجه به شکل `۱۴` می توان گفت داده ها دارای نویز کمی هستند.



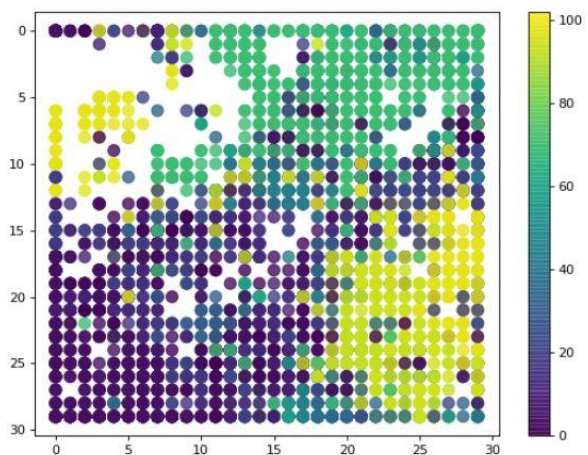
شکل ۱۲: نقشه 50*50 SOM با نرخ یادگیری اولیه 0.7



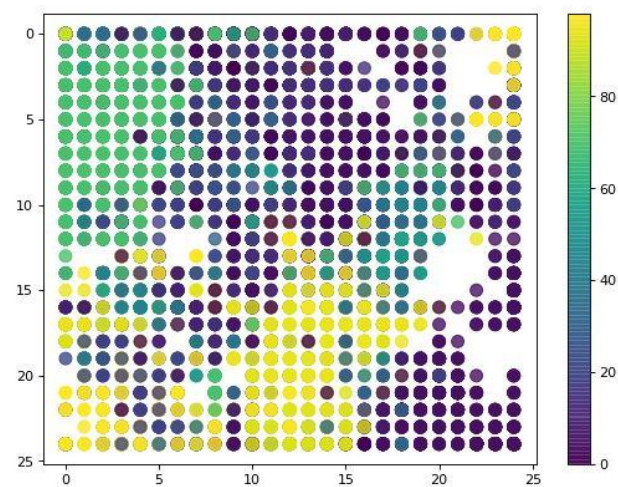
شکل ۱۱: نقشه 50*50 SOM با نرخ یادگیری اولیه 0.5



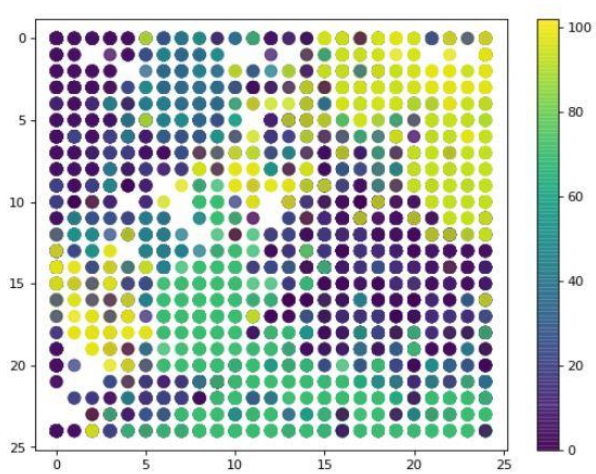
شکل ۱۴: نقشه 30*30 SOM با نرخ یادگیری اولیه 0.7



شکل ۱۳: نقشه 30*30 SOM با نرخ یادگیری اولیه 0.5



شکل ۱۵: نقشه 25*25 SOM با نرخ یادگیری اولیه 0.7



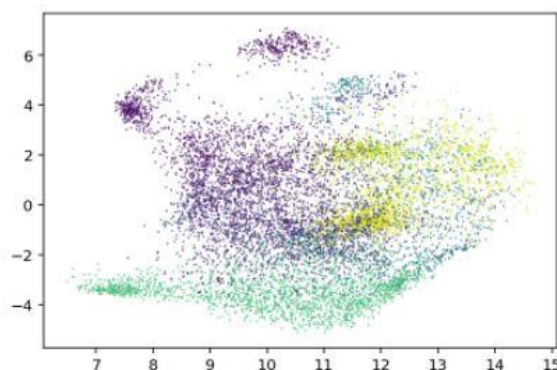
شکل ۱۵: نقشه 25*25 SOM با نرخ یادگیری اولیه 0.5

Purity

به طور کلی **purity** یک معیار ارزیابی مناسب برای داده های تک **label** می باشد و معیارهای ارزیابی دیگری برای داده های **multi-label** مورد استفاده می گردد. **Purity** را به این صورت محاسبه می کنیم: ابتدا یک تنسور با ابعاد نقشه و با عمق کلاس ها یا برجسب ها در نظر می گیریم. برای هر نود، داده هایی که به آن نگاشت شده را پیدا می کنیم. حال با بررسی کلاس های داده های نگاشت شده به آن خانه بردار تنسور هر خانه را به ازای هر برجسب آن داده یکی اضافه می کنیم. در نهایت نیز کلاسی که دارای بیشترین مقدار است را مرجع در نظر می گیریم. با جمع کلاس های مرجع هر نود و تقسیم آن ها به کل داده ها مقدار **purity** را به دست می آوریم. مقدار **purity** با این روش برای 25×25 som و نرخ یادگیری ۰,۷ برابر ۸۸,۸۱٪ شد. در روش دیگر برای محاسبه **purity** همچون کشیدن نقشه های بخش قبل از اولین کلاس هر داده نگاشت شده استفاده می شود که در این روش **purity** برابر ۶۲,۷۵٪ می شود. همچنین برای 50×50 som نیز **purity** برای دو روش فوق به ترتیب برابر ۹۱,۹۵٪ و ۷۱,۷۱٪ می باشد. دقت شود مقایسه **purity** برای som ها با ابعاد متفاوت کار درستی نیست.

لازم به ذکر است پیاده سازی som نیز پیگیری شد که با توجه به حجم بالای داده ها و ابعاد بسیار زیاد آن استفاده از این روش در محیط **colab** بسیار زمان بر بوده و نیاز به **RAM** بالایی برای تعداد داده های محدود دارد که عملاً تحلیل با تعداد نمونه های ۱۰۰ یا ۲۰۰ تا دقیق نخواهد بود.

همچنین روش دیگری برای صحت سنجی نتایج تمرین به این صورت انجام شده که به کمک روش **trunked SVD** ابعاد داده ها را ابتدا به ۱۵۰ و سپس به کمک **ivis** به دو بعد کاهش دهیم که در شکل زیر به نمایش داده شده است. همان طور که مشاهده می کنید تحلیل های بیان شده در بخش های بالا در شکل زیر نیز صادق است.



شکل ۱۵: ادغام دو روش **SVD** و **IVIS** برای کاهش بعد و نمایش نقشه ویژگی حاصل شده