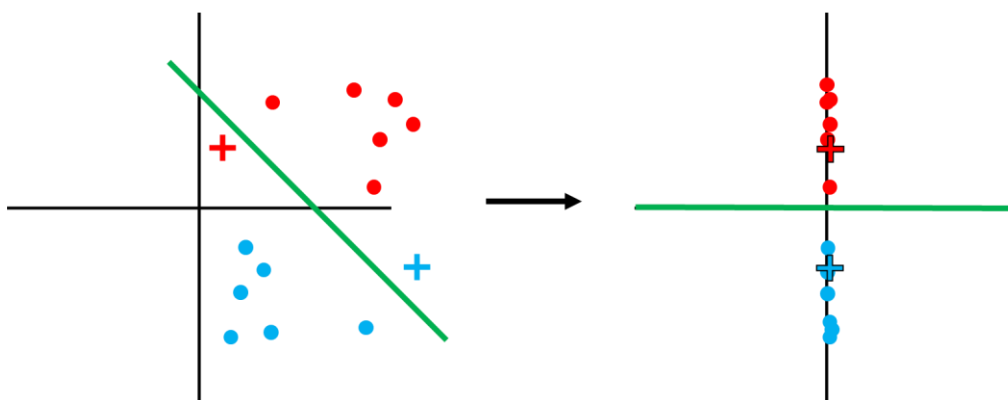




Question ONE

1.

(I)



دیتاست سمت چپ را در نظر بگیرید نقاط قرمز متعلق به کلاس 1 و نقاط آبی متعلق به کلاس 2 هستند. داده های به شکل دایره داده های آموزش و داده های + شکل داده های تست هستند.

اگر برای داده ها یک مرز تصمیم در 2 بعد پیدا کنیم انتظار داریم دسته بند ما مانند خط سبز در شکل سمت چپ شود چون به نظر دسته بند خوبی است (از این جهت که فاصله اش با هر دو کلاس تقریباً یکسان است و این فاصله را بیشینه کرده است).
با این حال این دسته بند نمیتواند داده های تست را درست دسته بندی کند.

اگر داده ها را روی محور عمودی مپ کنیم یعنی درواقع فقط همین بعد آن را در نظر بگیریم میتوان دید یک مرز تصمیم دیگر (خط افقی سبز) برای داده پیدا خواهد شد که دیگر مشکل دسته بندی غلط داده های تست را نخواهد داشت.
میتوان دید که دقت دسته بند روی داده های تست بعد از کاهش بعد بیشتر است و داده های تست درست دسته بندی شده اند.

(ب) دو مرحله ی Evaluation و subset generation.

Evaluation: هدف آن معیاری برای سنجش یک زیر مجموعه ی کاندید از ویژگی هاست.

متد های Evaluation :

Feature ranking: به هر ویژگی طبق یک معیار و استاندارد یک رنک داده میشود و اگر رنک آن به آستانه ی مشخصی نرسد، از مجموعه ویژگی ها حذف میشود. مثال هایی برای این متد : **Pearson correlation, mutual information** و **information gain**.

Probabilistic distance: محاسبه ی فاصله ی احتمالاتی بین احتمال شرطی دو کلاس یعنی $p(x | C1)$ و $p(x | C2)$. مثال هایی برای این متد : **Chernoff dissimilarity measure**.

Subset generation: هدف آن تولید یک زیر مجموعه برای سنجیده شدن توسط **Evaluation** است.

متد های **Subset generation**:

Complete search: این متد ها تضمین میکنند که به جواب بهینه (طبق معیار انتخاب شده برای **Evaluation**) خواهند رسید. برای مثال **exhaustive search** با جستجو و سنجیدن کل فضای حالت یعنی بین تمام زیر مجموعه های ممکن از ویژگی ها، پاسخ بهینه را میابد و بنابراین یک جستجوی کامل است.

Best individual N: یک رویکرد ساده این است که به هر ویژگی یک رنک (طبق معیار انتخاب شده برای **Evaluation**) داده شود و **N** بهترین آن ها به عنوان زیرمجموعه ی جواب برگردد.

$$\begin{aligned}
 Z_1 &= \omega_1^T x \\
 \text{Var}(Z_1) &= E[(\omega_1^T x - \omega_1^T \mu)^2] = E[(\omega_1^T x - \omega_1^T \mu)(\omega_1^T x - \omega_1^T \mu)^T] \\
 &= E[\omega_1^T (x - \mu)(x - \mu)^T \omega_1] = \\
 &= \omega_1^T E[(x - \mu)(x - \mu)^T] \omega_1 = \omega_1^T \Sigma \omega_1
 \end{aligned}$$

Cov matrix ←

$$\text{Var}(Z_1) = \omega_1^T \Sigma \omega_1 \quad (1)$$

← داده‌ی map شده

PCA خط و ضمیمه‌ی Var بیشینه شود :

$$\text{argmax}_{\omega_1} \omega_1^T \Sigma \omega_1 \quad \text{Subject to } \omega_1^T \omega_1 = 1$$

→ مشتق به نرم لایترال

$$\text{maximize}_{\omega_1} \omega_1^T \Sigma \omega_1 - \alpha (\omega_1^T \omega_1 - 1)$$

$$\frac{\partial}{\partial \omega_1} = 0 \rightarrow 2 \Sigma \omega_1 = 2 \alpha \omega_1 \Rightarrow \Sigma \omega_1 = \alpha \omega_1 \quad (2)$$

مشتق سری بر حسب ω_1 و α است
 معیار و قرار دادن \max شده

← بهینه‌ترین ω_1 بردار ویژه

ماتریس کواریانس (Σ) است و α نیز مقدار ویژه‌ی آن است.

$$\begin{aligned}
 \text{از (1) و (2)} \rightarrow \text{Var}(Z_1) &= \omega_1^T \Sigma \omega_1 = \alpha \omega_1^T \omega_1 = \alpha \\
 \text{Var}(Z_1) &= \alpha
 \end{aligned}$$

پس واریانس Z_1 برابر است با مقدار ویژه‌ی ماتریس Cov (Σ) و
 Z_1 هم حاصل نمایش داده‌ی x بر بردار ویژه‌ی Σ است ($Z_1 = \omega_1^T x$)
 ثابت شده بردار ویژه‌ی Σ است

$$\underset{\omega: \omega^T \omega = 1}{\operatorname{argmin}} \|x - x \omega \omega^T\|$$

استدلال: $\frac{1}{N} \operatorname{argmin}_{\omega: \omega^T \omega = 1} \|x - x \omega \omega^T\| = \operatorname{argmin}_{\omega: \omega^T \omega = 1} \sum_{i=1}^N \frac{\|x_i\|^2}{N} - \omega^T \Sigma \omega$

ماتریس کوواریانس: $\Sigma = \frac{1}{N} \sum_{i=1}^N x_i x_i^T$

$$\underset{\omega: \omega^T \omega = 1}{\operatorname{argmin}} \operatorname{var}(S) - \omega^T \Sigma \omega \Rightarrow \underset{\omega: \omega^T \omega = 1}{\operatorname{argmax}} \omega^T \Sigma \omega$$

مقدار ثابت

مسئله قبلی را به این شکل بیان می‌کنیم: $\operatorname{argmax}_{\omega: \omega^T \omega = 1} \omega^T \Sigma \omega$

$$\operatorname{argmax}_{\omega: \omega^T \omega = 1} \omega^T \Sigma \omega - \alpha (\omega^T \omega - 1)$$

$$\frac{d}{d\omega} = 0 \Rightarrow 2 \Sigma \omega = 2 \alpha \omega$$

$$\Sigma \omega = \alpha \omega$$

$$\Sigma \omega = \alpha \omega \Rightarrow \Sigma \omega = \alpha \omega$$

در notation این به شکل D از ω استفاده شده

(\bar{A}) منظور این است که ما فرض می‌کنیم یک سری فاکتورهای غیرقابل مشاهده و نهفته وجود دارند که ترکیب این فاکتورها x را حاصل می‌شود. پس در واقع هدف ما این هست که dependency ویژگی‌ها را با ترکیبی از فاکتورها که تعداد آن‌ها کمتر از ویژگی‌های اصلی است بیان کنیم. برای مثال ممکن است چندین ویژگی با همبستگی بسیار بالا داشته باشیم که این‌ها همه می‌توانند حالت‌های مختلف یک فاکتور باشند. کاربرد آن علاوه بر کاهش بعد می‌تواند knowledge extraction و بیان متغیرها با تعداد فاکتورهای کمتر باشد.

(ب)

Multidimensional Scaling (MDS): سعی میکند تمام نقاط را به یک بعد کمتر مپ کند به گونه ای که فاصله ی دو به دوی آن ها تا حد ممکن تغییر نکند. این فاصله میتواند فاصله ی اقلیدسی یا هر معیار فاصله ی دیگری باشد. در واقع MDS سعی میکند تابع زیر را کمینه کند که در آن نقاط $X_1 \dots X_n$ به نقاط $Z_1 \dots Z_n$ در ابعاد کمتری مپ شده اند و d_{ij} فاصله ی دو به دوی X_i هاست و \hat{d}_{ij} فاصله ی دو به دوی Z_i هاست و درواقع تابع دارد سعی میکند فواصل دو به دوی نقاط در فضای جدید تا حد امکان شبیه نقاط اولیه باشد و مجموع مجذورات تفاوت این فاصله ها کمینه باشد (فاصله در اینجا اقلیدسی است).

$$\min_z \sum_{i=1}^N \sum_{j=1}^N (d_{ij} - \hat{d}_{ij})^2$$

فاصله ها فواصل اقلیدسی هستند و تابع هدف MDS به فرم زیر کاهش میابد (reduce میشود) که در آن X_i یک نقطه در فضای اصلی و Z_i معادل آن نقطه در فضای جدید است. $X^T X$ نیز درواقع ماتریس Gram است.

$$\min_z \sum_{i=1}^N \sum_{j=1}^N (x_i^T x_j - z_i^T z_j)^2$$

Locally Linear Embedding: در این روش ساختار های غیرخطی در فضای سراسری بر اساس فیت های خطی در نواحی محلی به دست می آید. ایده ی اصلی LLE این است که هر نقطه جمع وزن داری از همسایه هایش است (تعریف همسایه میتواند بر اساس یک فاصله ی مشخص یا تعداد مشخصی نقاط نزدیک باشد). در مقوله ی کاهش بعد، این روش در این مورد قابل استفاده است که وزن های به دست آمده، ویژگی های هندسی ذاتی نقاط را منعکس میکنند بنابراین در ابعاد دیگر نیز معتبر خواهند بود.

در مرحله ی اول LLE ما میخواهیم وزن ها را به گونه ای پیدا کنیم که برای هر نقطه، جمع وزن دار همسایه هایش تا حد امکان به خود نقطه نزدیک باشد. برای هر نقطه x^r همسایه های آن را با نماد X_r^s نشان میدهیم. w_{rs} هم وزنی است که در همسایه ی s ام x^r ضرب میشود و باید مجذور فاصله ی این مجموع وزن دار از نقطه ی اصلی حساب شود. هدف LLE این است که مجموع این فاصله ها را کمینه کند یعنی وزن ها به گونه ای باشند که مجموع وزن دار همسایه های هر نقطه تا حد امکان به آن نقطه نزدیک باشد. و نکته ی دیگر هم اینکه برای هر نقطه باید مجموع وزن همسایگانش 1 باشد.

$$E[W|x] = \sum_{r=1}^N \left\| x^r - \sum_s w_{rs} X_r^s \right\|^2$$

$$\text{subject to } \sum_s w_{rs} = 1$$

سپس در مرحله ی دوم این وزن ها فیکس میشوند و LLE تلاش میکند با کمینه کردن تابع زیر مختصات Y (که در ابعاد جدید است) را به گونه ای بیابد که طبق همان وزن ها، این الگوهای خطی محلی (یعنی درواقع آن هندسه ی ذاتی نقاط) حفظ شوند. به این منظور تابع زیر را کمینه میکند که در آن Y نقاط در ابعاد جدید است، W_{rs} وزن های محاسبه شده در مرحله ی قبل است و Y_r^s همسایه های نقطه ی r هستند. با کمینه کردن تابع زیر میخواهیم شرط (مجموع به ازای هر نقطه ی مجذور فاصله ی مجموع وزن دار همسایه ها از نقطه ی اصلی کمینه شود) در فضای جدید نیز برقرار باشند. در ضمن شرط $Cov(Y) = I$ and $E[Y] = 0$ نیز باید برقرار باشد.

$$\mathbb{E}[Y|W] = \sum_{r=1}^N \left\| y^r - \sum_s w_{rs} y_{(r)}^s \right\|^2$$

in such a way that $Cov(Y) = I$ and $\mathbb{E}[Y] = 0$

Isomap: رویکرد isomap مانند MDS است با این تفاوت که به جای فاصله ی اقلیدسی از نوع دیگری از فاصله استفاده میکند و در واقع یک generalization غیرخطی از MDS است. ایزومپ ابتدا فواصل geodesic را محاسبه میکند و سپس MDS را روی این فواصل اجرا میکند. فاصله ی geodesic در واقع کوتاه ترین مسیر روی سطح منحنی است گویی که سطح flat باشد. برای محاسبه ی این نوع فاصله میتوان از الگوریتم هایی مانند فلویدارشال استفاده کرد.

این الگوریتم ابتدا روی فضای منحنی برای هر نقطه نقاطی را میابد که از لحاظ فاصله ی اقلیدسی فاصله ی کمتری نسبت به یک آتسانه داشته باشند ($d(i,j) < \epsilon$) سپس این نقاط را با یال به یکدیگر متصل میکند. در نتیجه یک گراف روی سطح منحنی ایجاد میشود. برای هر دو نقطه یعنی هر دو راس در این گراف کوتاه ترین مسیر محاسبه میشود و به شکل $dG(i,j)$ ذخیره میشود. حال ایزومپ عینا مانند MDS عمل میکند با این تفاوت که معیار فاصله ها را $dG(i,j)$ در نظر میگیرد. همانطور که در MDS توضیح داده شد، نقاط $X_1 \dots X_n$ به نقاط $Z_1 \dots Z_n$ در ابعاد کمتری مپ شده اند و در ایزومپ $dG(i,j)$ فاصله ی دو به دوی X_i ها در گراف است و d^G (ij) فاصله ی دو به دوی Z_i ها است و درواقع تابع دارد سعی میکند فواصل geodesic دو به دوی نقاط در فضای جدید تا حد امکان شبیه نقاط اولیه باشد و مجموع مجذورات تفاوت این فاصله ها کمینه باشد (فاصله در اینجا geodesic است که با مدل سازی به گراف و محاسبه ی کوتاه ترین مسیر بین هر دو راس به دست آمد).

$$\min_z \sum_{i=1}^N \sum_{j=1}^N \left(dG(i,j) - \hat{d}G(i,j) \right)^2$$

■

Question TWO

1.

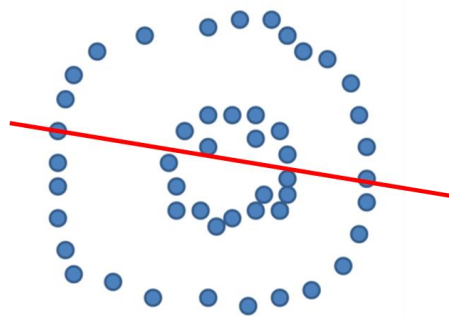
در kmeans ما می‌خواهیم برای k کلاستر، $(x - \mu_k)^2$ را کمینه کنیم که x داده‌های ما هستند و m مرکز هر کلاستر است.

در GM با k کلاستر و ماتریس کواریانس به شکل گفته شده، هدف کمینه شدن مقدار $\frac{(x - \mu_k)^2}{\sigma^2}$ است. زیرا نوع ماتریس کواریانس spherical است و این یعنی تمام خوشه‌ها دقیقاً دایره هستند (صورت کسر) که این اتفاق دقیقاً در فاصله‌ی اقلیدسی هم می‌افتد (یک دایره نقاطی هستند که فاصله‌ی اقلیدسیشان نسبت به مرکز - در اینجا مرکز خوشه - یکسان است)

2.

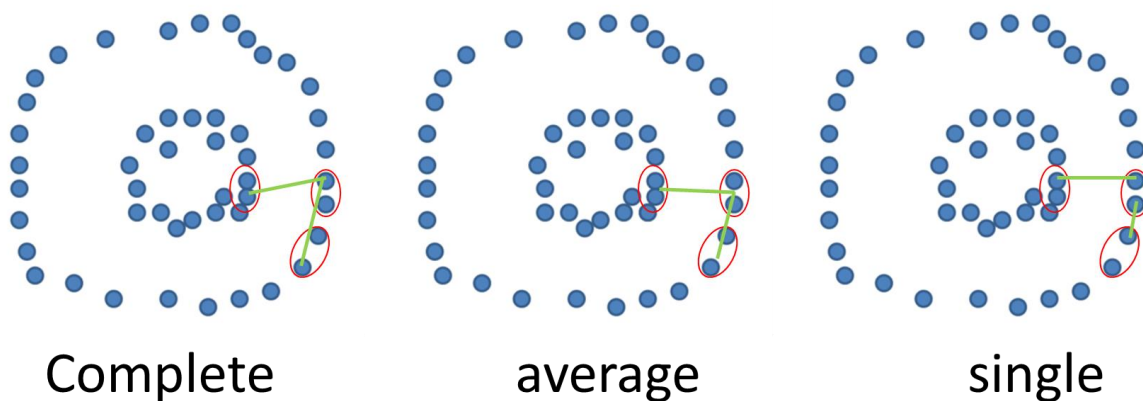
(آ) هر دو خوشه‌ها را با مرزی به این شکل (و نه با یک مرز دایره‌ای) از هم جدا می‌کنند دلیل آن این است که k-means هدفش پیدا کردن $k=2$ خوشه خواهد بود به گونه‌ای که دو مرکز خوشه پیدا کند به طوری که مجموع فاصله‌ی اجزای هر خوشه از مرکز آن خوشه کمینه شود به طور کلی یک سری مرکز خوشه و دایره‌هایی در اطراف آن در فضا رسم می‌کند (بر اساس فاصله‌ی اقلیدسی) به عنوان خوشه‌ها و بنابراین نمیتواند دو خوشه که مانند شکل زیر درون هم قرار دارند را دسته‌بندی کند زیرا مرکز این دو خیلی به هم نزدیک هستند ولی اینگونه نیست که هر کدام در دو دایره‌ی جدا نسبت به مرکز خود قرار بگیرند بلکه تفاوت آنها در فاصله‌شان از مرکز حدوداً مشترکشان است و این نوع خوشه‌ها را kmeans به شکل زیر دسته‌بندی میکند.

مدل EM هم با 2 توزیع گوسی، دقیقاً همین رفتار را خواهد داشت و مانند توضیحات بالا می‌باشد و کلا Kmeans خود حالت خاصی از EM با توزیع گوسی است.

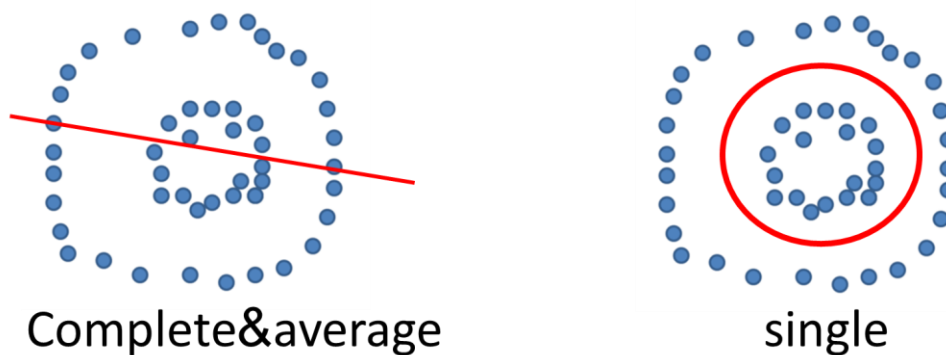


(ب) در ابتدا که هر نقطه یک کلاستر جداگانه است، در هر حالت نقاط نزدیک تر به هم یک کلاستر میشوند حالا برای مراحل بالاتر، در مورد single-link میتوان مطمئن بود که خوشه‌هایی که در ادامه با هم merge میشوند خوشه‌های تشکیل دهنده‌ی هر حلقه هستند پس دو خوشه‌ی به شکل دو حلقه از هم جدا خواهند شد. اما در دو حالت دیگر همانطور که در مثال زیر مشخص است مثلاً در مورد average ممکن است مرکز دو خوشه در یک حلقه از هم دورتر باشد تا مرکز دو خوشه در دو حلقه‌ی متفاوت هم چنین در

حالت complete نیز ممکن است ماکسیمم فاصله ی دو خوشه در یک حلقه از هم بیشتر باشد تا بیشترین فاصله ی دو خوشه در دو حلقه ی متفاوت بنابراین دو خوشه به شکل دو حلقه از هم جدا نخواهند شد.

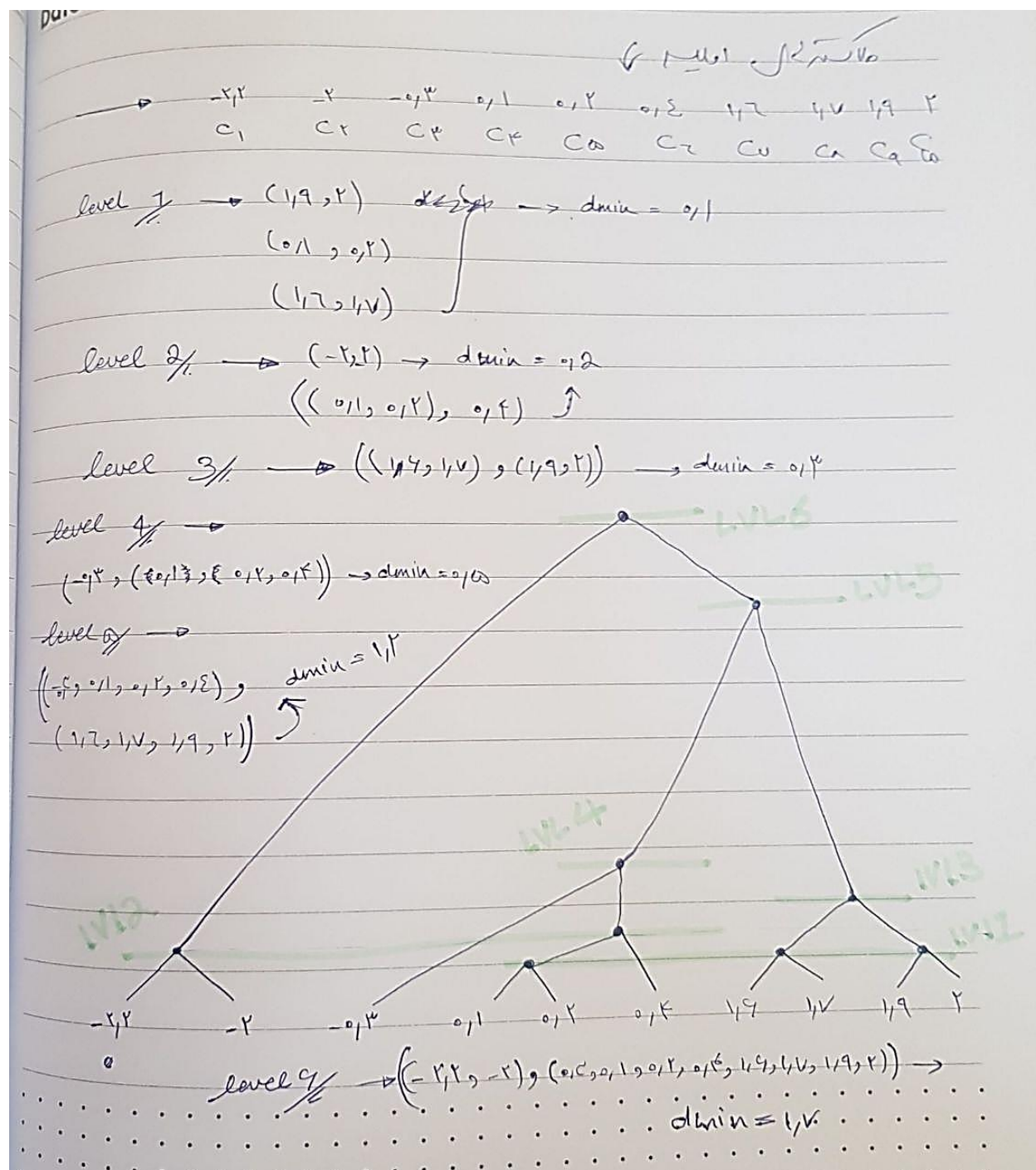


در نتیجه نتیجه ی نهایی کلاسترینگ حدودا به شکل زیر خواهد بود:



وقتی می‌خواهیم کلاسترهایی که طبق معیار گفته شده بیشترین شباهت را به هم دارند انتخاب کنیم در واقع می‌خواهیم کلاسترهایی را انتخاب کنیم که کمترین d_{min} را باهم دارند بنابراین با این روش، ابتدا تمام نقاط را یک کلاستر در نظر می‌گیریم سپس آنهایی که

کمترین d_{min} را با هم دارند با هم merge میکنیم و در level بعدی نیز همین کار را تکرار میکنیم تا به درخت نهایی برسیم. کلاستر هایی که در هر level انتخاب میشوند و d_{min} آنها نسبت به هم و هم چنین نتیجه ی نهایی در ادامه آمده است.



Question THREE.....

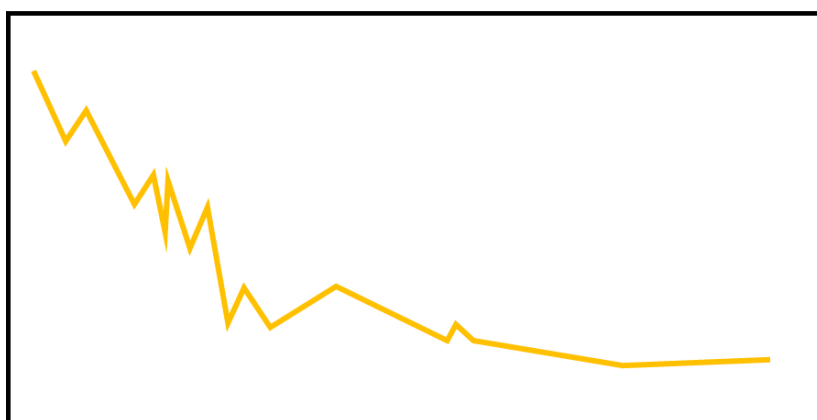
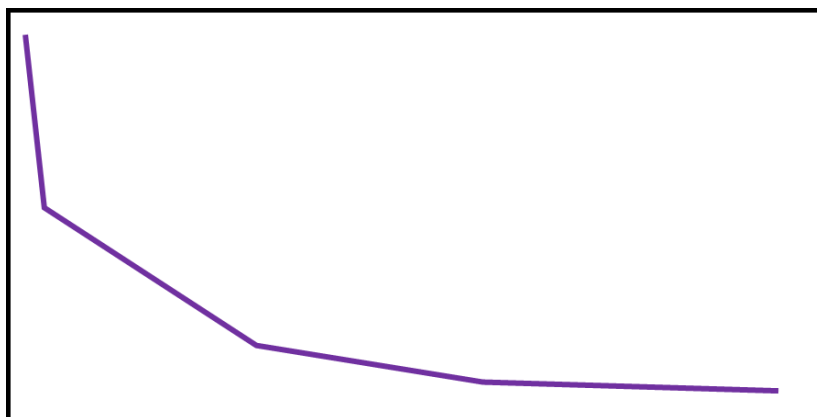
1.

(آ) پرسپترون یک الگوریتم iterative است که از یک نقطه ی شروع به سمت یک اوپتیمال محلی برای تابع هزینه حرکت میکند بنابراین نقطه ی شروع در آن تاثیر زیادی دارد و داده هایی که در ابتدا دیده میشوند تاثیرگذار تر هستند از داده هایی که در ایپاک های آخر دیده میشوند اگر چند تک داده ی outlier باشند، چون پرسپترون در هر ایپاک فقط یک داده را میبیند به سمت داده های اولیه ی outlier باایاس میشود. به همین علت بهتر است در کاربرد هایی که به نقطه ی شروع حساس هستد چندین بار الگوریتم را از نقطه های شروع متفاوت آغاز کرد.

(ب) در حالت Minibatch چون پرسپترون در هر ایپاک چندین داده میبیند و برای تصمیم جهت حرکت در فضای حالت، به گرادیان تمام آن بچ توجه میکند این باعث میشود حرکت نرم تری داشته باشد (چون یک داده ی outlier در یک بچ تاثیر کمتری دارد و اثر سایر داده های سالم آن را تا حد خوبی فیلتر میکند) و با یک شیب نسبتا نرم خطا در هر ایپاک کاهش میابد.

در حالت تک داده پرسپترون در هر ایپاک با توجه به گرادیان تنها یک داده جهت حرکت در فضای حالت را مشخص میکند این باعث میشود اگر داده outlier باشد یا نویز زیادی داشته باشد و حرکت تنها با توجه به آن صورت بگیرد، خطا روی سایر داده ها افزایش یابد پس در این حالت حرکت خیلی بیشتر نویز دارد و مانند حالت قبل یکنواخت نیست.

نمونه ای از مقدار تابع خطا در هر ایپاک برای دو حالت فوق:



2.

نکات مهم:

در این پاسخ $z_1 = \sigma(h_1)$ & $z_2 = \sigma(h_2)$ & $z_{1_final} = \sigma(o_1)$ & $z_{2_final} = \sigma(o_2)$ که σ تابع فعال ساز سیگموید است.

هم چنین در مشتق گیری ها تمام محاسبات تا 3 رقم اعشار انجام شده است.

علاوه بر لایه ی hidden، برای لایه ی خروجی نیز تابع سیگموید در نظر گرفته شده است.

$$h_1 = w_1 x_1 + w_2 x_2 + w_3 x_3 + b_1$$

$$h_1 = 0.1 \times 1 + 0.1 \times 2 + 0.1 \times 3 + 0.1 = 0.7$$

$$z_1 = \sigma(h_1) = \frac{1}{1 + e^{-h_1}} = 0.687$$

$$h_2 = w_4 x_1 + w_5 x_2 + w_6 x_3 + b_2$$

$$= 0.1 + 1.7 + 3 + 0.1 = 5.9$$

$$z_2 = \sigma(h_2) = \frac{1}{1 + e^{-h_2}} = 0.993$$

$$o_1 = w_7 z_1 + w_8 z_2 + b_3$$

$$0.1 \times 0.687 + 0.9 \times 0.993 + 0.1 = 1.089$$

$$\sigma(o_1) = \frac{1}{1 + e^{-o_1}} = 0.889 \rightarrow z_{1-final}$$

$$o_2 = w_9 z_1 + w_{10} z_2 + b_4$$

$$0.1 \times 0.687 + 0.1 \times 0.993 + 0.1 = 0.289$$

$$\sigma(o_2) = \frac{1}{1 + e^{-o_2}} = 0.58 \rightarrow z_{2-final}$$

$$SSE = (t_1 - z_{1-final})^2 + (t_2 - z_{2-final})^2$$

$$= (0.1 - 0.889)^2 + (0.05 - 0.58)^2 = 1.18$$

$$\frac{\partial E}{\partial z_{1-final}} = 2 \times (-1) (t_1 - z_{1-final}) = -2 (t_1 - z_{1-final})$$

$$= -2 (-0.789) = 1.578$$

$$\frac{\partial E}{\partial z_{2-final}} = 2 \times (-1) (t_2 - z_{2-final}) = -2 (t_2 - z_{2-final})$$

$$= -2 (-0.53) = 1.06$$

$$\frac{\partial E}{\partial \omega_v} = \frac{\partial E}{\partial z_{1-\text{final}}} \times \frac{\partial z_{1-\text{final}}}{\partial \omega_1} \times \frac{\partial \omega_1}{\partial \omega_v}$$

$$\frac{\partial E}{\partial z_{1-\text{final}}} = 1,05 \text{ V}$$

$$\frac{\partial z_{1-\text{final}}}{\partial \omega_1} = \sigma'(\omega_1) = \frac{e^{-\omega_1}}{(1+e^{-\omega_1})^2} = 0,109 \text{ V}$$

$$\frac{\partial \omega_1}{\partial \omega_v} = z_1 = 0,987$$

$$\left. \begin{aligned} dE/d\omega_v &= \\ 1,05 \text{ V} \times 0,109 \text{ V} &= \\ \times 0,987 &= \\ \underline{\underline{0,105}} \end{aligned} \right\}$$

$$\frac{\partial E}{\partial \omega_g} = \frac{\partial E}{\partial z_{1-\text{final}}} \times \frac{\partial z_{1-\text{final}}}{\partial \omega_1} \times \frac{\partial \omega_1}{\partial \omega_g}$$

$$\frac{\partial z_{1-\text{final}}}{\partial \omega_1} \frac{\partial \omega_1}{\partial \omega_g} = z_r = 0,990$$

$$\left. \begin{aligned} dE/d\omega_g &= \\ 1,05 \text{ V} \times 0,109 \text{ V} &= \\ 0,990 &= \underline{\underline{0,104}} \end{aligned} \right\}$$

$$\frac{\partial E}{\partial \omega_h} = \frac{\partial E}{\partial z_{r-\text{final}}} \times \frac{\partial z_{r-\text{final}}}{\partial \omega_r} \times \frac{\partial \omega_r}{\partial \omega_h}$$

$$\frac{\partial z_{r-\text{final}}}{\partial \omega_r} = \sigma'(\omega_r) = \sigma'(1,141 \text{ V}) = 0,109$$

$$\frac{\partial \omega_r}{\partial \omega_h} = z_1 = 0,987$$

$$\left. \begin{aligned} dE/d\omega_h &= \\ 1,05 \text{ V} \times 0,109 \text{ V} &= \\ 0,987 &= \underline{\underline{0,104}} \end{aligned} \right\}$$

$$\frac{\partial E}{\partial \omega_{10}} = \frac{\partial E}{\partial z_{r-\text{final}}} \times \frac{\partial z_{r-\text{final}}}{\partial \omega_r} \times \frac{\partial \omega_r}{\partial \omega_{10}}$$

$$\frac{\partial \omega_r}{\partial \omega_{10}} = z_r = 0,990$$

$$\left. \begin{aligned} dE/d\omega_{10} &= \\ 1,05 \text{ V} \times 0,109 \text{ V} &= \\ 0,990 &= \underline{\underline{0,104}} \end{aligned} \right\}$$

$$\frac{\partial E}{\partial \omega_1} = \frac{\partial E}{\partial z_1} \times \frac{\partial z_1}{\partial h_1} \times \frac{\partial h_1}{\partial \omega_1} = 0.191 \times 0.013 \times 1 = \underline{\underline{0.0025}}$$

$$\frac{\partial E}{\partial z_1} = \frac{\partial E_1}{\partial z_1} + \frac{\partial E_r}{\partial z_1} = (0.101 + 0.190) = 0.291$$

$$\frac{\partial E_1}{\partial z_1} = \frac{\partial E_1}{\partial o_1} \times \frac{\partial o_1}{\partial z_1} = \frac{\partial E_1}{\partial z_{1\text{final}}} \times \frac{\partial z_{1\text{final}}}{\partial o_1} \times \frac{\partial o_1}{\partial z_1}$$

$\sigma'(o_1)$

$1.0 \times 0.091 \times 0.1 = \underline{\underline{0.101}}$

$$\frac{\partial E_r}{\partial z_1} = \frac{\partial E_r}{\partial o_r} \times \frac{\partial o_r}{\partial z_1} =$$

$$\frac{\partial E_r}{\partial z_{r\text{final}}} \times \frac{\partial z_{r\text{final}}}{\partial o_r} \times \frac{\partial o_r}{\partial z_1}$$

$\sigma'(o_r)$

$1.0 \times 0.109 \times 0.1 = 0.190$

$$\frac{\partial z_1}{\partial h_1} = \sigma'(h_1) = \sigma'(1.3) = 0.013$$

→

$$\frac{\partial E}{\partial \omega_r} = \frac{\partial h_1}{\partial \omega_r} \times \frac{\partial E}{\partial z_1} \times \frac{\partial z_1}{\partial h_1} = 0.1 \times 0.291 \times 0.013 = \underline{\underline{0.0038}}$$

$$\frac{\partial E}{\partial \omega_o} = \frac{\partial h_1}{\partial \omega_o} \times \frac{\partial E}{\partial z_1} \times \frac{\partial z_1}{\partial h_1} = 0.1 \times 0.291 \times 0.013 = \underline{\underline{0.0038}}$$

$$\frac{\partial E}{\partial \omega_r} = \frac{\partial E}{\partial z_r} \times \frac{\partial z_r}{\partial h_r} \times \frac{\partial h_r}{\partial \omega_r} = 0.192 \times 0.004 \times 1 = 0.000768$$

$$\frac{\partial z_r}{\partial h_r} = \sigma'(h_r) = \sigma'(2.13) = 0.004$$

$$\frac{\partial E}{\partial z_r} = \frac{\partial E_1}{\partial z_r} + \frac{\partial E_r}{\partial z_r} = 0.139 + 0.024 = 0.163$$

$$\frac{\partial E_1}{\partial z_r} = \frac{\partial E_1}{\partial o_1} \times \frac{\partial o_1}{\partial z_r} = \frac{\partial E_1}{\partial z_{1, \text{final}}} \times \frac{\partial z_{1, \text{final}}}{\partial o_1} \times \frac{\partial o_1}{\partial z_r}$$

$1.05 \times 0.09 \times 0.19 = 0.139$

$$\frac{\partial E_r}{\partial z_r} = \frac{\partial E_r}{\partial o_r} \times \frac{\partial o_r}{\partial z_r} = \frac{\partial E_r}{\partial z_{r, \text{final}}} \times \frac{\partial z_{r, \text{final}}}{\partial o_r} \times \frac{\partial o_r}{\partial z_r}$$

$1.0 \times 0.109 \times 0.1 = 0.109$

→

$$\frac{\partial E}{\partial \omega_r} = \frac{\partial h_r}{\partial \omega_r} \times \frac{\partial E}{\partial z_r} \times \frac{\partial z_r}{\partial h_r} = 1 \times 0.163 \times 0.004 = 0.0006512$$

$$\frac{\partial E}{\partial \omega_r} = \frac{\partial h_r}{\partial \omega_r} \times \frac{\partial E}{\partial z_r} \times \frac{\partial z_r}{\partial h_r} = 1 \times 0.163 \times 0.004 = 0.0006512$$

$$\frac{dE}{db_r} = \frac{dE_1}{db_r} + \frac{dE_r}{db_r} =$$

$$\frac{dz_{1\text{final}}}{dO_1} \times \frac{dO_1}{db_r} + \frac{dz_{r\text{final}}}{dO_r} \times \frac{dO_r}{db_r}$$

$$= \sigma'(O_1) + \sigma'(O_r) = 0.099 + 0.109 = \underline{\underline{0.208}}$$

$$\frac{dE}{db_1} = \frac{dE}{dz_1} \times \frac{dz_1}{dh_1} \times \frac{dh_1}{db_1} + \frac{dE}{dz_r} \times \frac{dz_r}{dh_r} \times \frac{dh_r}{db_1}$$

$$0.099 \times 0.013 \times 1 + 0.172 \times 0.005 \times 1 = \underline{\underline{0.001}}$$

→ new values for Bias.

$$b_1 = b_1 - 0.01 \times \frac{dE}{db_1} = 0.0 - 0.01 \times 0.001 = 0.099$$

$$b_r = b_r - 0.01 \times \frac{dE}{db_r} = 0.0 - 0.01 \times 0.208 = 0.099$$

$$w_i = w_i - 0.01 \times \frac{dE}{dw_i} \Rightarrow$$

$$w_1 = 0.1 - 0.01 \times 0.001 = 0.099$$

$$w_2 = 0.2 - 0.01 \times 0 = 0.2$$

$$w_3 = 0.3 - 0.01 \times 0.010 = 0.299$$

$$w_4 = 0.5 - 0.01 \times 0.002 = 0.499$$

$$w_5 = 0.6 - 0.01 \times 0.019 = 0.591$$

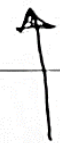
$$w_6 = 0.7 - 0.01 \times 0.001 = 0.699$$

$$w_7 = 0.8 - 0.01 \times 0.012 = 0.791$$

$$w_8 = 0.9 - 0.01 \times 0.014 = 0.891$$

$$w_9 = 0.9 - 0.01 \times 0.012 = 0.891$$

$$w_{10} = 0.1 - 0.01 \times 0.014 = 0.091$$



new values for weights.