



## یادگیری ماشین

پاییز ۱۴۰۳

استاد: علی شریفی زارچی

مسئول تمرین: نیکی سپاسیان

مهلت ارسال نهایی: ۱۸ آبان

تمرین دوم

مهلت ارسال امتیازی: ۱۱ آبان

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روزهای مشخص شده است.
- در طول ترم، برای هر تمرین می‌توانید تا ۵ روز تأخیر مجاز داشته باشید و در مجموع حداکثر ۱۵ روز تأخیر مجاز خواهید داشت. توجه داشته باشید که تأخیر در تمرین‌های عملی و تئوری به صورت جداگانه محاسبه می‌شود و مجموع تأخیر هر دو نباید بیشتر از ۱۵ روز شود. پس از اتمام زمان مجاز، دو روز اضافی برای آپلود غیرمجاز در نظر گرفته شده است که در این بازه به ازای هر ساعت تأخیر، ۲ درصد از نمره تمرین کسر خواهد شد.
- اگر بخش عملی یا تئوری تمرین را قبل از مهلت ارسال امتیازی آپلود کنید، ۲۰ درصد نمره اضافی به آن بخش تعلق خواهد گرفت و پس از آن، ویدئویی تحت عنوان راهنمایی برای حل تمرین منتشر خواهد شد.
- حتماً تمرین‌ها را بر اساس موارد ذکر شده در صورت سوالات حل کنید. در صورت وجود هرگونه ابهام، آن را در صفحه تمرین در سایت کوئرا مطرح کنید و به پاسخ‌هایی که از سوی دستیار آموزشی مربوطه ارائه می‌شود، توجه کنید.
- در صورت هم‌فکری و یا استفاده از هر منابع خارج درسی، نام هم‌فکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.
- فایل پاسخ‌های سوالات نظری را در قالب یک فایل pdf به فرمت  $HW2\_T\_ [STD\_ID].pdf$  آماده کنید و برای سوالات عملی، هریک را در یک فایل zip جداگانه قرار دهید و فایل zip اول را به فرمت  $HW2\_P1\_ [STD\_ID].zip$  و فایل zip دوم را به فرمت  $HW2\_P2\_ [STD\_ID].zip$  نامگذاری کرده و هرکدام را به صورت جداگانه آپلود کنید.
- گردآورندگان تمرین: محمدپارسا بشری، عرفان جعفری، مهدی طباطبایی، فاطمه السادات موسوی، محمد مولوی

## سوالات نظری (۱۰۰ نمره)

۱. (۲۰ نمره) در یک مسئله طبقه‌بندی دو کلاسه با دو ویژگی، از هر کلاس دو داده داریم:

$$\omega_1 : \begin{pmatrix} 4 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \omega_2 : \begin{pmatrix} 2 \\ 5 \end{pmatrix}, \begin{pmatrix} 3 \\ 4 \end{pmatrix}$$

- الف) با محاسبه بردار میانگین و ماتریس کوواریانس، داده‌ها را با استفاده از PCA به فضای یک بعدی برده و تمایزپذیری کلاس‌ها را بررسی کنید. برای هر دو راستا مسئله را حل کنید.
- ب) در قسمت الف یک بردار پیشنهاد دهید که اگر به همه داده‌ها اضافه شود مولفه‌ی اساسی اول تغییر نکند.
- حل.  
(الف)

برای محاسبه بردار میانگین باید میانگین تمام داده‌ها را مستقل از کلاس آن‌ها حساب کرد:

$$\bar{x} = \frac{1}{4} \left( \begin{pmatrix} 4 \\ 1 \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} 2 \\ 5 \end{pmatrix} + \begin{pmatrix} 3 \\ 4 \end{pmatrix} \right) = \begin{pmatrix} 2.5 \\ 2.5 \end{pmatrix}$$

سپس باید ماتریس کوواریانس را حساب کنیم:

$$\tilde{X} = \begin{pmatrix} 4 - 2/5 & 1 - 2/5 \\ 1 - 2/5 & 0 - 2/5 \\ 2 - 2/5 & 5 - 2/5 \\ 3 - 2/5 & 4 - 2/5 \end{pmatrix} = \begin{pmatrix} 1/5 & -1/5 \\ -1/5 & -2/5 \\ -1/5 & 2/5 \\ 0/5 & 1/5 \end{pmatrix}$$

$$S = \frac{1}{N} \tilde{X}^T \tilde{X} = \begin{pmatrix} 1/66 & 0/33 \\ 0/33 & 5/66 \end{pmatrix}$$

حال باید بردارهای ویژه و مقادیر ویژه متناظر با ماتریس کواریانس را پیدا کنیم:

$$Sv = \lambda v$$

$$\lambda_1 = 5/7, \lambda_2 = 1/63$$

$$v_1 = \begin{pmatrix} -0/08 \\ -0/99 \end{pmatrix}, v_2 = \begin{pmatrix} 0/99 \\ -0/08 \end{pmatrix}$$

حال داده‌ها را به راستای اول project می‌کنیم:

$$X' = XA = Xv_1$$

$$X' = \begin{pmatrix} 4 & 1 \\ 1 & 0 \\ 2 & 5 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} -0/08 \\ -0/99 \end{pmatrix} = \begin{pmatrix} -1/31 \\ -0/08 \\ -5/11 \\ -4/21 \end{pmatrix}$$

همانگونه که مشخص است، دو کلاس تمایزپذیر هستند. حال همین کار را برای راستای دوم انجام می‌دهیم:

$$X' = XA = Xv_2$$

$$X' = \begin{pmatrix} 4 & 1 \\ 1 & 0 \\ 2 & 5 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} 0/99 \\ -0/08 \end{pmatrix} = \begin{pmatrix} 3/88 \\ 0/99 \\ 1/58 \\ 2/65 \end{pmatrix}$$

همانگونه که می‌بینید، دو کلاس تمایزپذیر نیستند.

ب) نشان می‌دهیم که اضافه شدن هر بردار دلخواه  $a$  از فضای ویژگی‌ها به تمام داده‌ها مقدار مولفه اساسی اول را ثابت نگه می‌دارد.

$$\text{scalar} \rightarrow a = [10 \ -6]^T$$

$$\bar{X}_{new} = \frac{1}{n} \sum (x_i + a)$$

$$= \frac{1}{n} (\sum x_i + na) = \bar{X} + a$$

$$\bar{X}_{new} = \bar{X} + a$$

بنابراین داریم:

$$Cov(X_{new}) = \frac{1}{n-1} \sum (x_i + a - \bar{X}_{new})(x_i + a - \bar{X}_{new})^T$$

$$\begin{aligned}
&= \frac{1}{n-1} \sum (x_i + a - \bar{X} - a)(x_i + a - \bar{X} - a)^T \\
&= \frac{1}{n-1} \sum (x_i - \bar{X})(x_i - \bar{X})^T \xrightarrow{\text{ثابت}} \\
&= Cov(x)
\end{aligned}$$

بنابراین بردارهای ویژه تغییر نخواهند کرد و مقدار مولفه اساسی اول تغییر نخواهد کرد.

۲. (۲۰ نمره) قصد داریم مدلی بسازیم که با گرفتن ورودی  $X$ ، متغیر خروجی  $Y$  را تخمین بزند. فرض کنید  $M$  مدل ضعیف به نام‌های  $f_1(X), \dots, f_M(X)$  داریم که روی نمونه‌های bootstrap شده‌ای از دیتاست اصلی آموزش داده شده‌اند و دارای بایاس و واریانس یکسان هستند. مدل نهایی به شکل زیر تعریف می‌شود:

$$f_{\text{ensemble}}(X) = \frac{1}{M} \sum_{i=1}^M f_i(X)$$

الف) بایاس و واریانس مدل  $f_{\text{ensemble}}(X)$  را بر حسب بایاس و واریانس مدل‌های ضعیف محاسبه کنید. در این قسمت فرض کنید مدل‌های ضعیف از یکدیگر مستقل هستند. تحلیل کنید که افزایش تعداد مدل‌های ضعیف ( $M$ ) چه تاثیری روی بایاس و واریانس مدل نهایی دارد.

ب) حال فرض کنید مدل‌های ضعیف مستقل نیستند و correlation بین هر دو مدل  $f_i(X)$  و  $f_j(X)$  ( $i \neq j$ ) برابر  $\rho$  است. حال مانند قسمت الف، بایاس و واریانس مدل نهایی را محاسبه کرده و تاثیر تعداد مدل‌ها ( $M$ ) و وابستگی بین آن‌ها ( $\rho$ ) را روی این مقادیر بررسی کنید.

پ) به سوالات زیر به طور خلاصه پاسخ دهید.

- آیا یادگیرنده‌های ضعیف در adaboost نیاز به مشتق‌پذیر بودن دارند؟ چرا؟
- بین bagging و boosting، کدام یک از نظر محاسباتی گران‌تر می‌باشد؟ چرا؟

حل.  
الف) بایاس مدل جدید به شکل زیر محاسبه می‌شود:

$$\text{Bias}(f_{\text{ensemble}}(X)) = \frac{1}{M} \sum_{i=1}^M \text{Bias}(f_i(X)) = \text{Bias}(f_i(X))$$

واریانس مدل جدید نیز با توجه به مستقل بودن مدل‌های ضعیف به شکل زیر محاسبه می‌شود:

$$\text{Var}(f_{\text{ensemble}}(X)) = \frac{1}{M^2} \sum_{i=1}^M \text{Var}(f_i(X)) = \frac{1}{M} \text{Var}(f_i(X))$$

بنابراین می‌توان دید که با افزایش تعداد مدل‌ها ( $M$ ) واریانس مدل نهایی کاهش می‌یابد در حالیکه بایاس ثابت می‌ماند. برای مدل‌های ضعیف که معمولاً واریانس بالا و بایاس پایین دارند، این کار منجر به کاهش چشمگیری در خطای مدل نهایی می‌شود.

ب) بایاس مدل نهایی در حالت مستقل نبودن مدل‌های ضعیف فرقی با قسمت قبل ندارد. واریانس مدل جدید

به شکل زیر محاسبه خواهد شد:

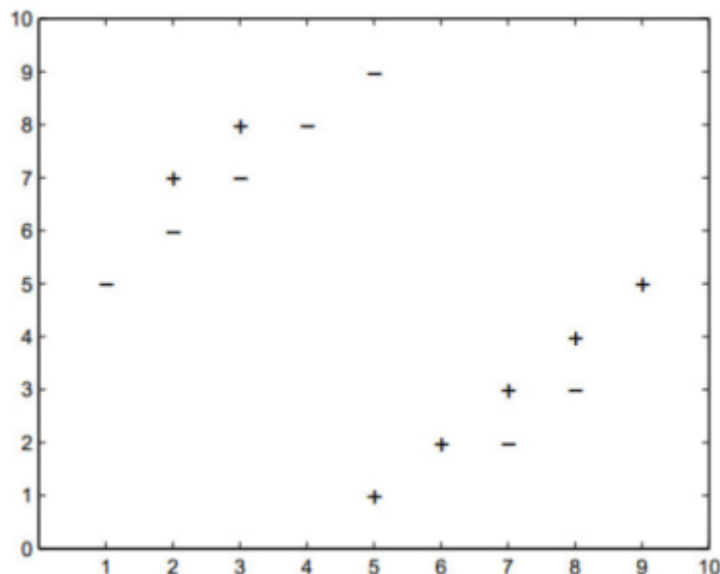
$$\begin{aligned}
 \text{Var}(f_{\text{ensemble}}(X)) &= \text{Var}\left(\frac{1}{M} \sum_{i=1}^M f_i(X)\right) \\
 &= \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \text{Cov}(f_i(X), f_j(X)) \\
 &= \frac{1}{M^2} \sum_{i=j} \underbrace{\text{Cov}(f_i(X), f_i(X))}_{\text{Var}(f_i(X))} + \frac{1}{M^2} \sum_{i \neq j} \text{Cov}(f_i(X), f_j(X)) \\
 &= \frac{1}{M^2} \sum_{i=j} \text{Var}(f_i(X)) + \frac{1}{M^2} \sum_{i \neq j} \rho \text{Var}(f_i(X)) \\
 &= \frac{1}{M} \text{Var}(f_i(X)) + \frac{M-1}{M} \rho \text{Var}(f_i(X))
 \end{aligned}$$

بنابراین می‌توان دید که با افزایش  $\rho$  تاثیر تعداد مدل‌های ضعیف ( $M$ ) کمتر می‌شود. به عبارتی اضافه کردن تعداد بیشتری مدل ضعیف، واریانس مدل را به قدر کافی کاهش نمی‌دهد؛ زیرا به دلیل وابستگی زیاد بین این مدل‌ها، در واقع مدل‌های اضافه شده دارند پیش‌بینی‌های مشابهی با مدل‌های موجود ارائه می‌دهند.

(پ)

- خیر، adaboost نیازی به مشتق‌پذیر بودن یادگیرنده‌های ضعیف ندارد زیرا بر وزن‌دهی مجدد و رای‌گیری متکی است و نه بر روش‌های بهینه‌سازی مبتنی بر گرادیان.
- Boosting، زیرا ماهیت دنباله‌دار boosting و وابسته بودن هر مدل به نتایج مدل قبلی اجازه موازی‌سازی را نمی‌دهد. درحالی که در bagging یادگیرنده‌های پایه به دلیل استقلال می‌توانند به صورت موازی آموزش داده شوند.

۳. (۲۰ نمره) در این سوال یک دسته‌بند KNN با متریک فاصله  $L_2$  در نظر بگیرید. کلاس‌ها را تماماً دو حالت (+/-) در نظر خواهیم گرفت. به سوالات زیر با توجه به مجموعه داده مشخص شده در تصویر پاسخ دهید.



- الف) به ازای چه مقدار  $k$  خطای این دسته‌بند کمینه می‌شود؟ مقدار این خطا چقدر است؟  
 ب) چرا استفاده از مقادیر بسیار زیاد یا بسیار کم برای  $k$  می‌تواند منجر به خطا شود؟

پ) فرض کنید از روش Leave One Out Cross Validation استفاده کنیم. به ازای چه مقدار  $k$  خطای دسته‌بندی کمینه می‌شود؟ مقدار این خطا چقدر است؟ (استفاده از ابزارهای sklearn برای به دست آوردن  $k$  بهینه مجاز است.)

ت) مرز تصمیم برای دسته‌بند 1-NN را برای این مجموعه داده در تصویر نشان دهید.

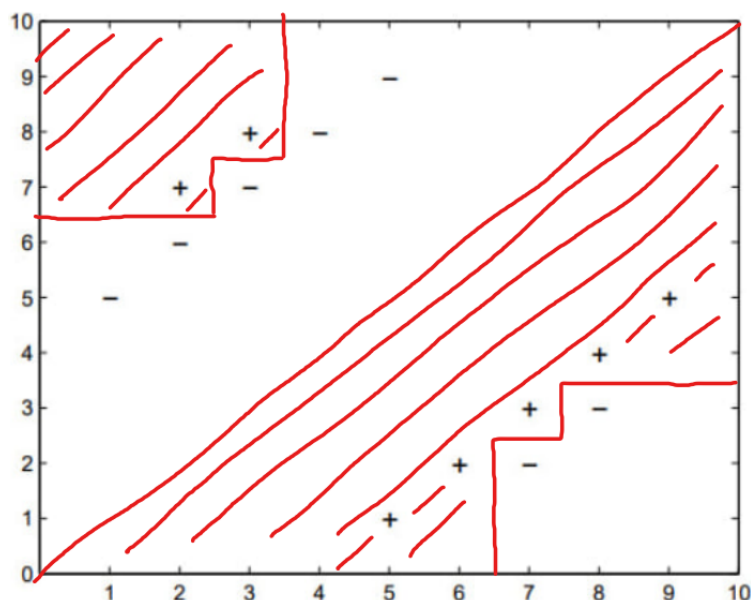
الف) چون هر نقطه‌ای همسایه خودش (نزدیکترین همسایه خودش) به شمار می‌رود، 1-NN کمترین خطا را بدست می‌دهد که همان صفر است.

ب) **مقادیر زیاد  $k$** : خط پایین سمت راست تصویر که دو داده منفی روی آن قرار دارند را در نظر بگیرید؛ این خط بوسیله خطی متشکل از داده‌های مثبت، از سایر داده‌های منفی جدا می‌شود. در صورت افزایش  $k$ ، دادگان مثبت نیز در دسته‌بندی در نظر گرفته می‌شوند که موجب افزایش خطا خواهند شد.

**مقادیر اندک  $k$** : در دو سمت مجموعه داده که دادگان dense هستند، می‌توان مشاهده کرد که اکثر داده‌هایی که فاصله اقلیدسی کمی از یکدیگر دارند، عضو کلاس‌های متفاوت هستند، که در صورت کوچک گرفتن اندازه همسایگی این امر موجب افزایش خطا خواهد شد.

پ)  $k = 5$ , error = 0.28571 برای بدست آوردن  $k$  بهینه می‌توانید از ابزارهای sklearn استفاده کنید (کد ضمیمه پاسخنامه شده است.)

ت)



۴. (۲۰ نمره)

تابع هزینه الگوریتم خوشه‌بندی k-means به شکل زیر تعریف می‌شود:

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

که در آن  $x_1, x_2, \dots, x_n$  نمونه‌ها و  $\mu_1, \mu_2, \dots, \mu_k$  مراکز خوشه‌ها هستند. منظور از  $S_j$  نیز مجموعه‌ای از نمونه‌هاست که به مرکز  $\mu_j$  نزدیک‌تر از هر خوشه‌ی دیگر هستند.

الف) یک مرحله از الگوریتم را در نظر بگیرید که برچسب داده‌ها  $y_j$  ثابت است و مراکز خوشه‌ها  $\mu_i$  به‌روزرسانی می‌شوند. نشان دهید که برای کمینه کردن تابع هزینه در این مرحله، کافی است میانگین هر خوشه به‌عنوان مرکز آن خوشه قرار گیرد.

ب) آیا الگوریتم k-means نسبت به مقداردهی اولیه مراکز خوشه‌ها حساس است؟ آیا این الگوریتم به‌طور تضمینی همگرا می‌شود؟

پ) در مرحله‌ای از الگوریتم k-means که در آن میانگین خوشه‌ها  $\mu_i$  ثابت‌اند و برچسب‌های نقاط داده  $y_j$  به‌روزرسانی می‌شوند، گاهی ممکن است یک نقطه  $X_j$  به چندین مرکز خوشه با فاصله مساوی نزدیک باشد. اگر یکی از گزینه‌ها این باشد که  $X_j$  در همان خوشه‌ای که در تکرار قبلی بوده باقی بماند، توضیح دهید چرا بهتر است این گزینه انتخاب شود؟ اگر این اصل رعایت نشود، چه مشکلی ممکن است پیش بیاید؟

ح. الف) اگر  $n_j$  تعداد نقاط نمونه اختصاص داده شده به خوشه  $j$  باشد، و از تابع هزینه  $L$  استفاده کنیم، مشتق تابع هزینه نسبت به  $\mu_i$  به صورت زیر خواهد بود:

$$\frac{\partial L}{\partial \mu_i} = \sum_{y_j=i} (\mu_i - X_j)$$

با برابر صفر قرار دادن این مشتق، خواهیم داشت:

$$\sum_{y_j=i} \mu_i = \sum_{y_j=i} X_j$$

که معادل است با:

$$n_j \mu_i = \sum_{y_j=i} X_j$$

بنابراین:

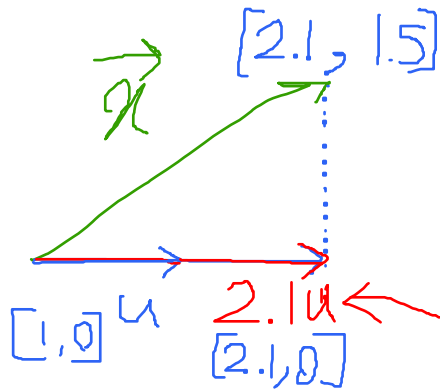
$$\mu_i = \frac{1}{n_j} \sum_{y_j=i} X_j$$

که همان میانگین نقاط نمونه اختصاص داده شده به خوشه  $i$  است.

ب) الگوریتم k-means به صورت تضمینی به حداقل محل (نه لزوماً سراسری) همگرا می‌شود و به انتخاب مقدار  $k$  و محل مراکز اولیه‌ی خوشه‌ها حساس است زیرا این دو مقدار نقش تعیین کننده‌ای در همگرایی به یک نقطه بهینه سراسری دارند.

پ) اگر یک نقطه نمونه بین دو یا چند میانگین خوشه که فاصله‌ای مساوی دارند، دائماً جابه‌جا شود، الگوریتم k-means ممکن است به یک چرخه بی‌پایان از تغییرات برسد. در این حالت، هر بار که برچسب آن نقطه به‌روزرسانی می‌شود، بدون اینکه هیچ بهبودی در کاهش تابع هزینه حاصل شود، نقطه به خوشه دیگری منتقل می‌شود. این باعث می‌شود الگوریتم نتواند به حالت پایدار برسد. در نتیجه، ممکن است الگوریتم به پایان نرسد، حتی با وجود اینکه تابع هزینه دیگر کاهش پیدا نمی‌کند و پیشرفتی در خوشه‌بندی حاصل نمی‌شود.

۵. (۲۰ نمره) فرض کنید مجموعه‌ای از نقاط  $x^{(1)}, \dots, x^{(m)}$  داده شده‌اند. فرض کنید داده‌ها نرمال شده‌اند و دارای میانگین صفر و واریانس یک در هر بعد هستند. همچنین فرض کنید  $f_u(x)$  تصویر نقطه  $x$  در جهت بردار یکه  $u$  باشد. به عبارتی اگر داشته باشیم:



$$V = \{au : a \in \mathbb{R}\}$$

آنگاه:

$$f_u(x) = \arg \min_{v \in V} \|x - v\|^2$$

نشان دهید بردار  $u$  که خطای MSE بین نقاط و تصویر آن‌ها را کمینه می‌کند، همان مولفه اصلی اول ( $PC_1$ ) است. به عبارتی نشان دهید که:

$$\arg \min_{u: u^T u = 1} \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - f_u(x^{(i)})\|^2$$

برابر اولین مولفه اصلی است.

حل.  
ابتدا داریم:

$$f_u(x) = \arg \min_{v=au} \|x - au\|^2$$

$$\begin{aligned} \frac{\partial}{\partial a} \|x - au\|^2 &= \frac{\partial}{\partial a} (x - au)^T (x - au) \\ &= \frac{\partial}{\partial a} (x^T x - 2au^T x + a^2 u^T u) \\ &= -2u^T x + 2au^T u = 0 \quad \Rightarrow \quad a = \frac{u^T x}{u^T u} \end{aligned}$$

$$f_u(x) = au = \frac{u^T x}{u^T u} u$$

حال عبارت را بازنویسی می‌کنیم:

$$\begin{aligned}
\arg \min_{\|u\|=1} \sum_{i=1}^m \|x^{(i)} - f_u(x^{(i)})\|^2 &= \arg \min_{\|u\|=1} \sum_{i=1}^m \|x^{(i)} - uu^T x^{(i)}\|^2 \\
&= \arg \min_{\|u\|=1} \sum_{i=1}^m (x^{(i)} - uu^T x^{(i)})^T (x^{(i)} - uu^T x^{(i)}) \\
&= \arg \min_{\|u\|=1} \sum_{i=1}^m (x^{(i)T} x^{(i)} - u^T x^{(i)} x^{(i)T} u) \\
&= \arg \min_{\|u\|=1} \sum_{i=1}^m (x^{(i)T} x^{(i)} - \mathbf{1}(u^T x^{(i)})^2 + (u^T x^{(i)})^2) \\
&= \arg \min_{\|u\|=1} \sum_{i=1}^m -(u^T x^{(i)})^2 \\
&= \arg \max_{\|u\|=1} \sum_{i=1}^m (u^T x^{(i)})^2 \\
&= \arg \max_{\|u\|=1} u^T \left( \sum_{i=1}^m x^{(i)} x^{(i)T} \right) u
\end{aligned}$$

حل معادله زیر را با استفاده از روش لاگرانژ انجام می‌دهیم:

$$\arg \max_{\|u\|=1} u^T \Sigma u$$

conditional optimization

حال برای حل این مسئله از ضرایب لاگرانژ استفاده می‌کنیم:

$$\mathcal{L}(u, \lambda) = u^T \Sigma u - \lambda(u^T u - 1)$$

حال برای حل این معادله مشتق لاگرانژین را محاسبه می‌کنیم:

$$\frac{\partial}{\partial u} \mathcal{L}(u, \lambda) = 2\Sigma u - 2\lambda u = 0 \quad \Rightarrow \quad \Sigma u = \lambda u$$

پس می‌توان نتیجه گرفت که جواب این مسئله همان بردارهای ویژه ماتریس  $\Sigma$  هستند.

$$u^T \Sigma u = u^T \lambda u = \lambda u^T u = \lambda$$

در نتیجه این مقدار جواب بهینه‌ای است که معادل با بردارهای ویژه ماتریس  $\Sigma$  می‌باشد که همان اولین مولفه PCA است.