

کاهش ابعاد

Dimensionality Reduction

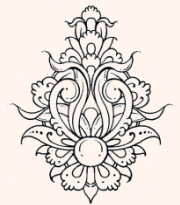
# یادگیری ماشین



دانشگاه شهید بهشتی  
پژوهشکده‌ی فضای مجازی  
پاییز ۱۴۰۱  
احمد محمودی ازناوه

# فهرست مطالب

- مزایای کاهش ابعاد
- انتخاب خصیصه
- استخراج خصیصه
- تحلیل مؤلفه‌ی اصلی
- تحلیل تفکیک خطی
- تجزیه به مقادیر تکین
- روش‌های غیر خطی
- تحلیل عاملی
- تغییر مقیاس داده‌های چند بعدی

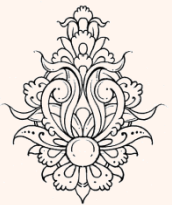
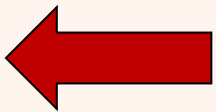


# نکبت ابعاد! (مشقت ابعاد)

- از لحاظ نظری، افزایش ابعاد منجر به بهبود عملکرد دسته‌بندی می‌شود، اما در عمل همیشه این گونه نیست.

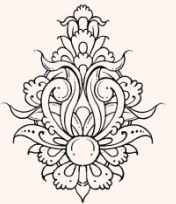
## *Curse of dimensionality*

- انتظار می‌رود در یک فرآیند ایده‌آل دسته‌بندی یا رگرسیون از خصیصه‌های بی‌اهمیت صرف‌نظر شود و فرآیند «**کاهش ابعاد**» به صورت جداگانه مورد نیاز نباشد. با این وجود کاهش ابعاد به دلایل زیر مورد توجه قرار می‌گیرد:



# مزایای کاهش ابعاد

- «کاهش حجم محاسبات»: حافظه مصرفی و حجم محاسبات به تعداد ( $N$ ) و ابعاد ( $d$ ) داده‌ها بستگی دارد.
  - زمان محاسبات
  - حافظه مورد نیاز
- «صرفه‌جویی در جمع‌آوری داده»: حذف داده‌های غیرضروری
- «مقاوم بودن» (robustness): مدل‌های ساده، هنگامی که داده‌های آموزشی کم باشد، «مقاوم‌تر» می‌باشند؛ قدرت پیش‌بینی برای تعداد مشخصی داده، با افزایش ابعاد، کاهش می‌یابد.



# مزایای کاهش ابعاد (ادامه...)

- «استفراجه دانش»: با تعداد فصیصه‌های کمتر، در مورد داده‌ها و فرآیندهای مربوط به آن درک بهتری وجود خواهد داشت. گاهی این فصیصه‌ها را می‌توان به صورت «عوامل پنهان» در نظر گرفت که متغیرهای قابل مشاهده از آنها نشأت می‌گیرند.

## *Hidden or latent factor*

- هنگامی که تعداد فصیصه‌ها (بدون از دست دادن اطلاعات) کمتر باشد، «ساختار داده‌ها» بهتر درک می‌شود. داده‌های پرت و غیرمعمول بهتر تشخیص داده می‌شود؛ قابلیت نمایش بهتری دارند.



# انتخاب - استخراج (خصیصه)

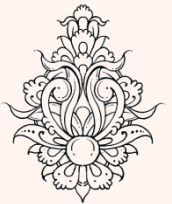
## Feature Selection vs Extraction

- انتخاب خصیصه:

- $K$  خصیصه‌ی مهم‌تر ( $k < d$ ) انتخاب می‌شود.
- الگوریتم‌های انتخاب زیرمجموعه

- استخراج خصیصه:

- $K$  خصیصه‌ی جدید، استخراج می‌شود.
- نگاشت از فضای  $n$ -بعدی به فضای  $k$ -بعدی
- روش‌های استخراج خصیصه نیز از دیدگاه‌های مختلفی قابل طبقه‌بندی هستند، روش‌های خطی در برابر روش‌های غیرخطی و یا روش‌های بی‌نظارت در برابر روش‌های بانظارت



# انتخاب زیرمجموعه

- در انتخاب زیرمجموعه، هدف انتخاب **بهترین** زیرمجموعه، زیرمجموعه‌ای با کمترین ابعاد و درست‌ترین نتیجه، می‌باشد.

- $2^d$  زیرمجموعه، در یک مجموعه  $d$ -عضوی وجود دارد، بررسی تمام حالات به جز زمانی که  $d$  کوچک باشد، امکان‌پذیر نیست.

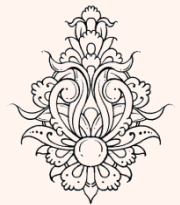
## Forward search

- جستجوی رو به جلو:

- در گام نخست، مجموعه‌ی خصیصه‌ها،  $F$  در حالت اولیه  $\emptyset$  در نظر گرفته می‌شود.

- در هر گام بهترین خصیصه به مجموعه‌ی خصیصه‌ها افزوده می‌شود. (میزان خطای  $(E(F))$  کمتر)

- برای بررسی خطا باید از داده‌های validation استفاده کرد.

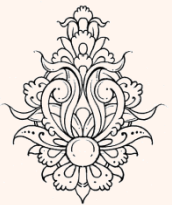


$$j = \arg \min_j E(F \cup x_j)$$
$$\text{Add } x_j \text{ to } F \text{ if } E(F \cup x_j) < E(F)$$

# انتخاب زیرمجموعه (ادامه...)

## Backward search

- جستجوی رو به عقب:
  - در گام نخست، مجموعه‌ی خصیصه‌ها،  $F$  در حالت اولیه تمامی خصیصه‌ها در نظر گرفته می‌شود.
  - در هر گام بدترین خصیصه از مجموعه‌ی خصیصه‌ها حذف می‌شود.
- هنگامی که تعداد خصیصه‌ها زیاد است، روش جستجوی رو به جلو ترجیح داده می‌شود.
- انتخاب زیرمجموعه به صورت بانظارت است.
- در کاربردهایی که یک خصیصه به تنهایی اطلاعات مفیدی ندارد، انتخاب خصیصه مفید نیست. (مانند تشخیص چهره)



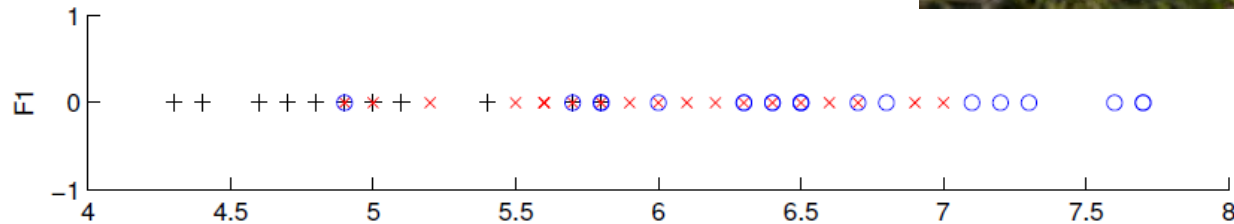


# Iris data: Single feature

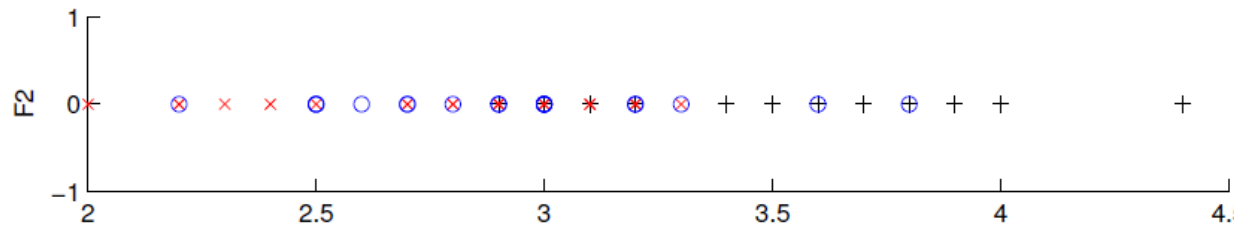
مثال



0.76



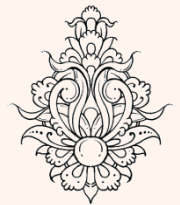
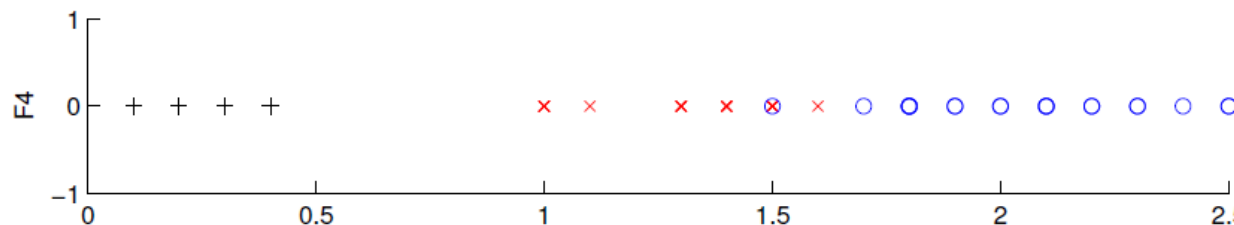
0.57



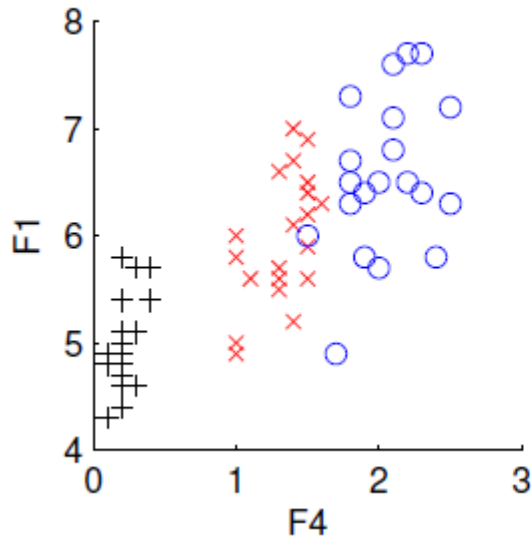
0.92



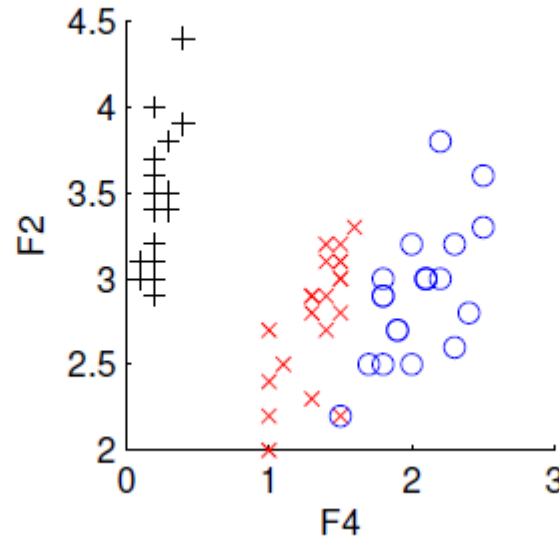
0.94



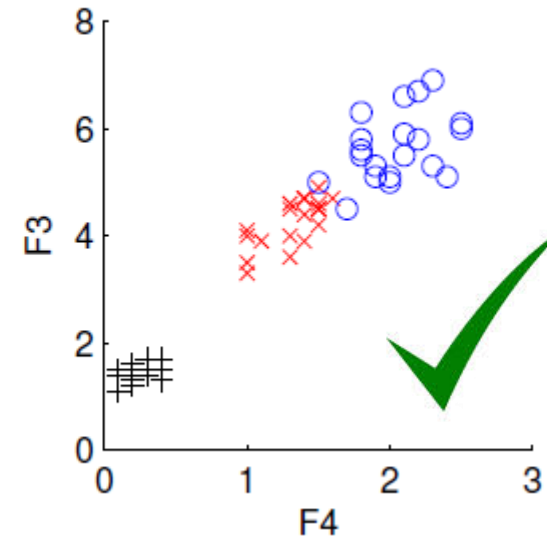
# Iris data: Add one more feature to F4



0.87



0.94

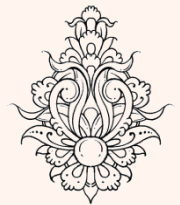


0.96

در صورت اضافه کردن فمحصه بعدی نتایج افت می‌کند!

در بسیاری موارد انتخاب فمحصه‌ها به نوع دسته‌بند بستگی دارد.

در صورت کوچک بودن دسته‌داده، فمحصه‌ی انتخاب شده، می‌تواند به نمره‌ی تقسیم پایگاه به دو دسته‌ی training و validation مربوط باشد.



## تحلیل مؤلفه‌های اصلی

- هدف نگاشت داده‌ی  $d$ -بعدی به فضای  $k$ -بعدی است ( $k < d$ )، به گونه‌ای که **کم‌ترین میزان اتلاف** رخ دهد.

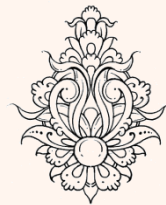
– نگاشت  $x$  در راستای  $w$ :

$$z = w^T x$$

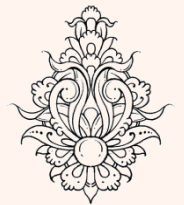
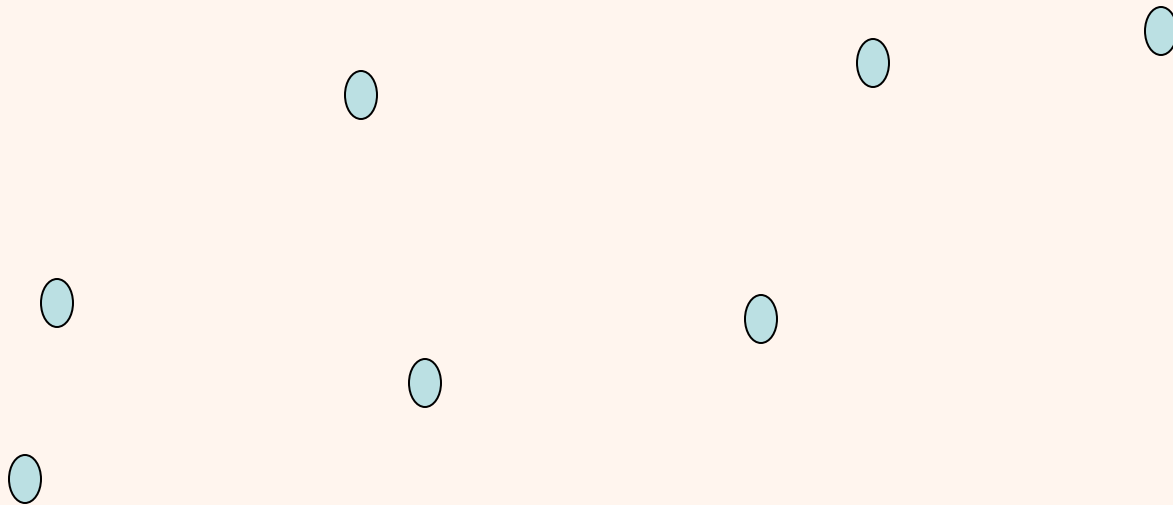
- این راستا به گونه‌ای انتخاب می‌شود که  $\text{Var}(z)$  **بیشینه** شود، راستایی که داده در امتداد آن بیشترین تغییرات را داشته باشد.

– این مسأله باعث می‌شود، تفاوت نمونه‌ها آشکارتر شود.

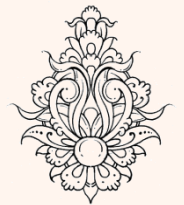
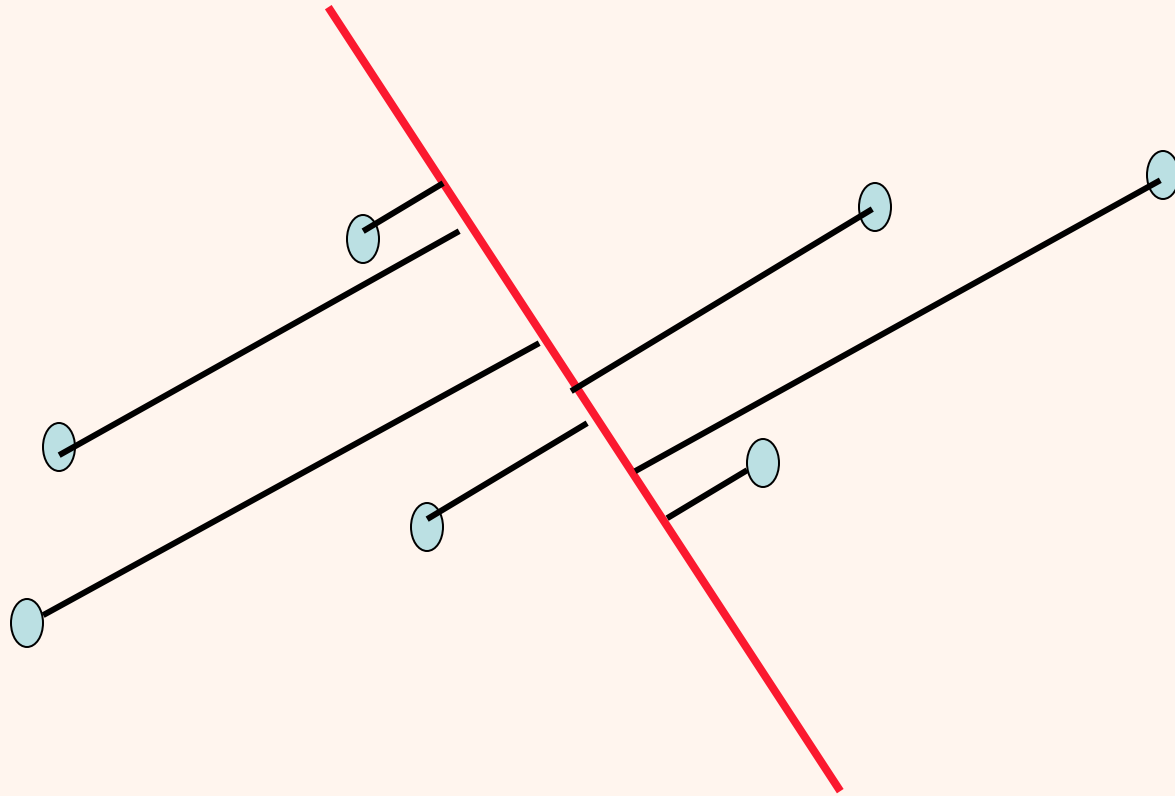
- این شیوه‌ی کاهش بعد به صورت «بی‌نظارت» است.



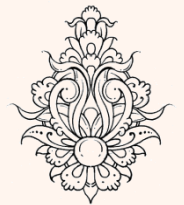
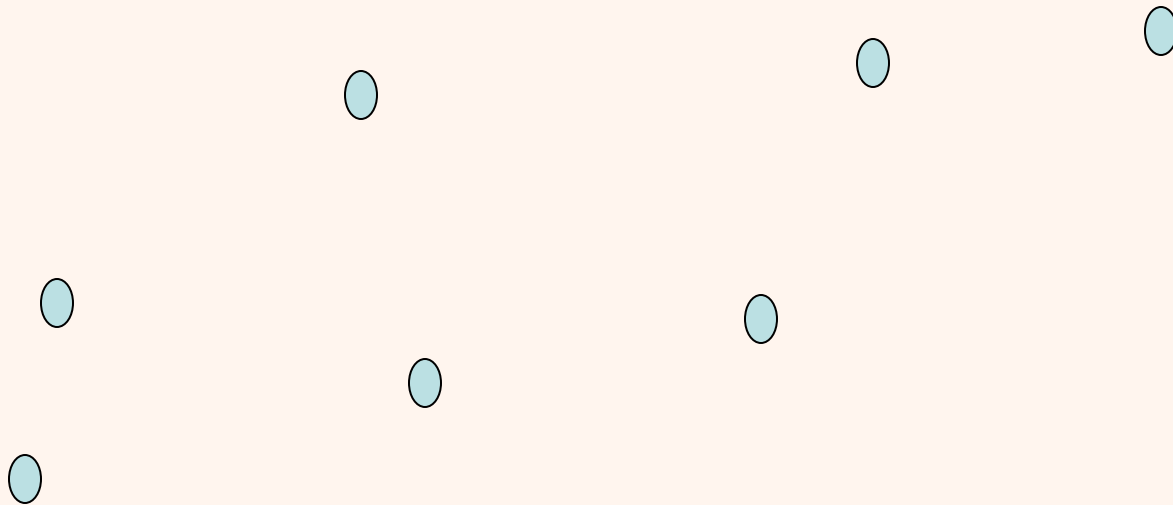
# تحلیل مؤلفه‌های اصلی (ادامه...)



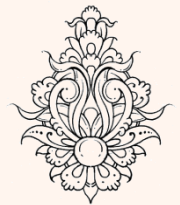
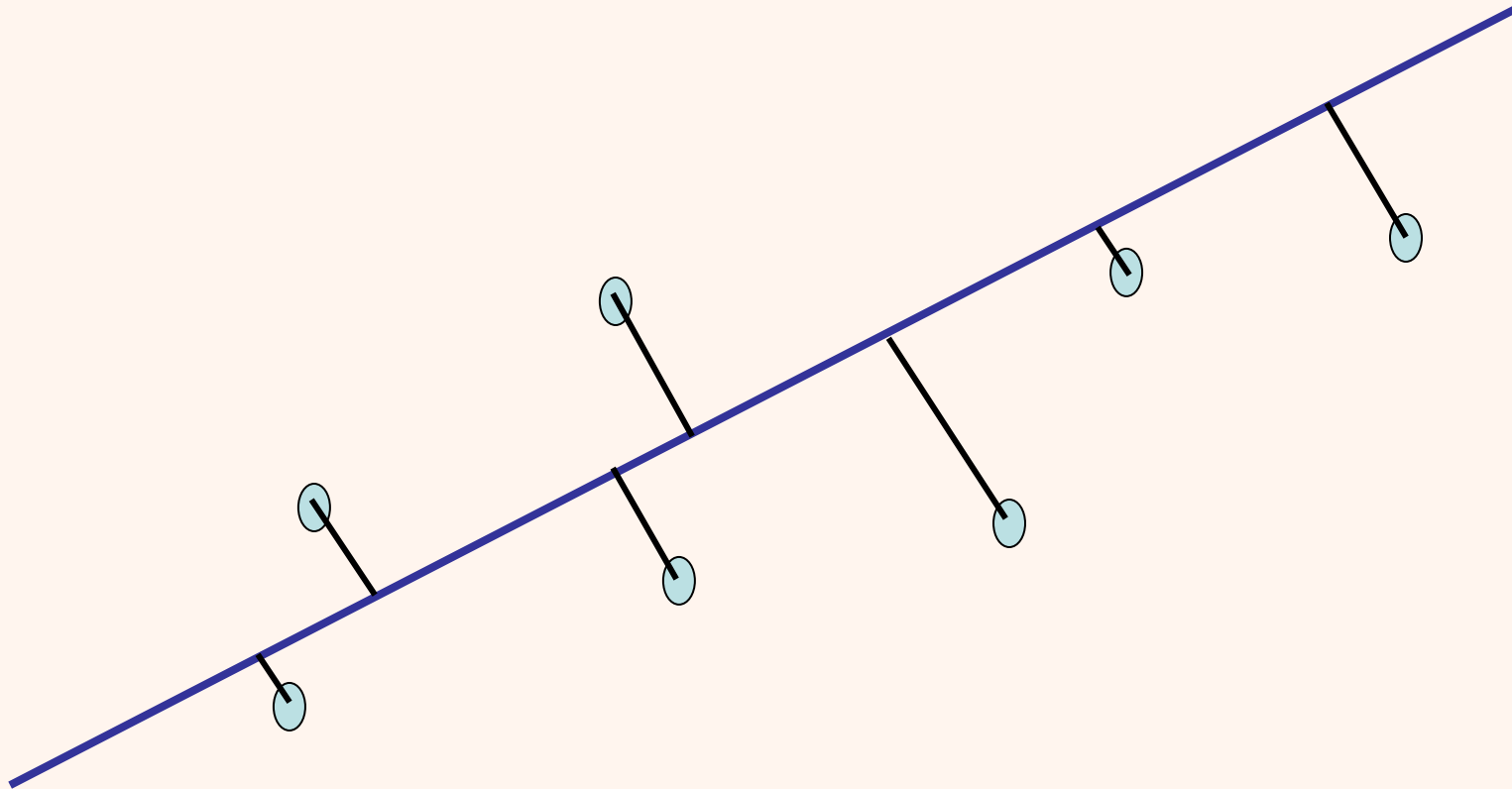
# تحلیل مؤلفه‌های اصلی (ادامه...)



# تحلیل مؤلفه‌های اصلی (ادامه...)



# تحلیل مؤلفه‌های اصلی (ادامه...)



# تحلیل مؤلفه‌های اصلی (ادامه...)

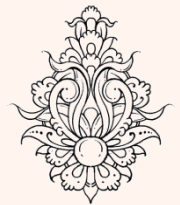
- در راستای  $w$ ، پراکندگی داده ماکزیمم می‌شود:

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

where  $\text{Cov}(x) = \Sigma$

- در این حالت تنها راستا است که اهمیت دارد، در نتیجه برای یافتن پاسخ یکتا، باید شرط زیر نیز برقرار باشد:

$$\|w\| = 1$$





# تحلیل مؤلفه‌های اصلی (ادامه...)

- در نتیجه برای اولین مؤلفه‌ی اساسی رابطه‌ی زیر به دست می‌آید:  $\max w_1^T \Sigma w_1 - \alpha (w_1^T w_1 - 1)$

- با مشتق گرفتن نسبت به  $w_1$  و برابر صفر قرار دادن آن

$$2\Sigma w_1 - 2\alpha w_1 = 0$$

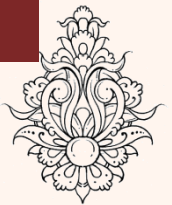
- در نتیجه

$$\Sigma w_1 = \alpha w_1$$

در نتیجه  $w_1$  یکی از بردارهای ویژه‌ی ماتریس  $\Sigma$  می‌باشد

- از طرفی  $w_1^T \Sigma w_1 = \alpha$  ، در واقع واریانس در راستای  $w_1$  برابر مقدار ویژه‌ی متناظر با آن است.

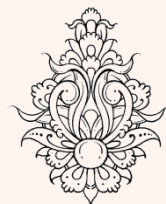
اولین مؤلفه‌ی اصلی، برابر بردار ویژه‌ی ماتریس کواریانس با بیشترین مقدار ویژه است.



# تحلیل مؤلفه‌های اصلی (ادامه...)

- برای یافتن دومین مؤلفه‌ی اصلی، علاوه بر شرایط پیش باید بر راستای اولین مؤلفه‌ی اساسی هم عمود باشد، در این حالت داده‌های نگاشت شده «**ناهمبسته**» (uncorrelated) خواهند بود.
- برای یافتن دومین مؤلفه‌ی اصلی ( $w_2$ )، باید  $\text{Var}(z_2)$  ماکزیمم شود، مشروط به متعامد بودن بر اولین مؤلفه‌ی اصلی و  $\|w_2\|=1$

$$\max_{w_2} w_2^T \Sigma w_2 - \alpha (w_2^T w_2 - 1) - \beta (w_2^T w_1 - 0)$$



# تحلیل مؤلفه‌های اصلی (ادامه...)

$$\max_{\mathbf{w}_2} \mathbf{w}_2^T \Sigma \mathbf{w}_2 - \alpha (\mathbf{w}_2^T \mathbf{w}_2 - 1) - \beta (\mathbf{w}_2^T \mathbf{w}_1 - 0)$$

• پس از مشتق گرفتن خواهیم داشت:

$$2\Sigma \mathbf{w}_2 - 2\alpha \mathbf{w}_2 - \beta \mathbf{w}_1 = 0$$

• با ضرب در  $\mathbf{w}_1^T$

$$\mathbf{w}_1^T \Sigma \mathbf{w}_2 - 2\alpha \mathbf{w}_1^T \mathbf{w}_2 - \beta \mathbf{w}_1^T \mathbf{w}_1 = 0$$

$$\mathbf{w}_1^T \mathbf{w}_2 = 0$$

$$\mathbf{w}_1^T \Sigma \mathbf{w}_2 - \beta \mathbf{w}_1^T \mathbf{w}_1 = 0$$

$$\mathbf{w}_1^T \Sigma \mathbf{w}_2 = \mathbf{w}_2^T \Sigma \mathbf{w}_1 = \lambda_1 \mathbf{w}_2^T \mathbf{w}_1 = 0$$

→  $\beta = 0$

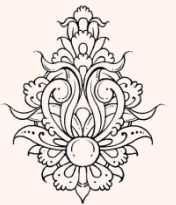
→  $\Sigma \mathbf{w}_2 = \alpha \mathbf{w}_2$

دومین مؤلفه‌ی اصلی، برابر بردار ویژه‌ی ماتریس کواریانس با بیشترین مقدار ویژه در رده‌ی دوم است، به همین ترتیب سایر مقادیر ویژه به دست می‌آیند.



# تحلیل مؤلفه‌های اصلی (ادامه...)

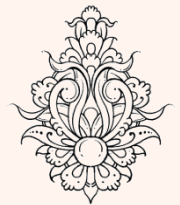
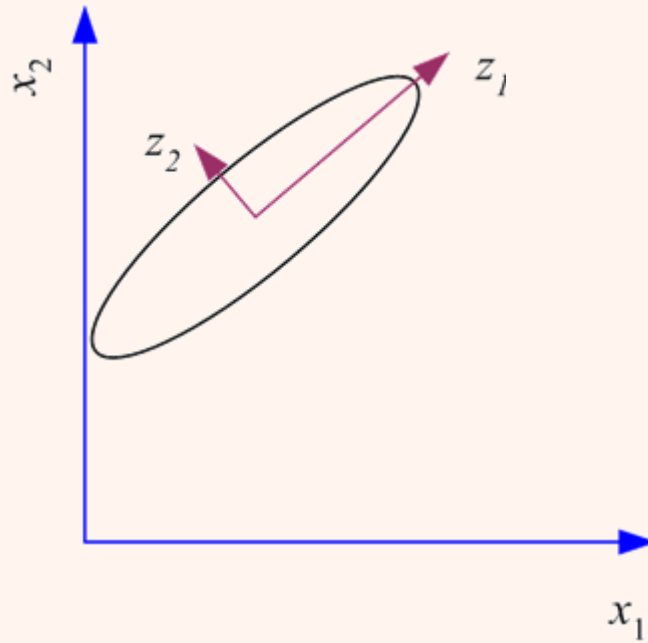
- در صورتی که ماتریس متقارن باشد، بردارهای ویژه‌ی آن متعامد هستند.
- در صورتی که ماتریس positive definite باشند، مقادیر ویژه همگی مثبت خواهند بود.
- در صورتی که ماتریس singular باشد، به اندازه‌ی rank ماتریس مقادیر ویژه غیرصفر خواهیم داشت.



# تحلیل مؤلفه‌های اصلی (ادامه...)

$$z = W^T(x - m)$$

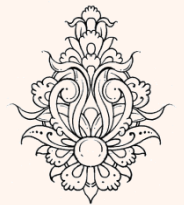
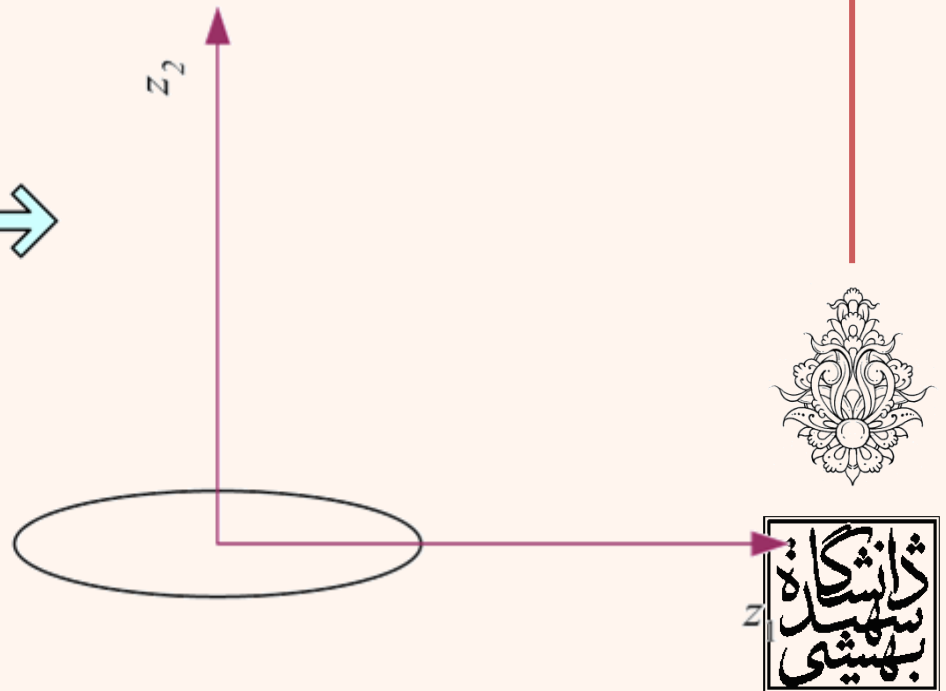
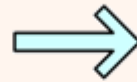
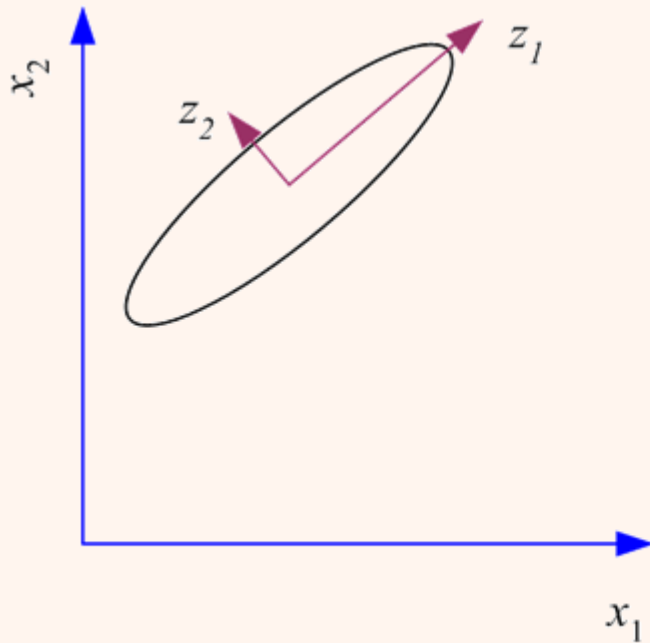
- ستون‌های  $W$ ، بردارهای ویژه‌ی ماتریس کواریانس هستند.



# تحلیل مؤلفه‌های اصلی (ادامه...)

$$z = W^T(x - m)$$

- ستون‌های  $W$ ، بردارهای ویژه‌ی ماتریس کواریانس هستند.



# کاهش بعد

• در صورتی که  $|S|$  کوچک باشد، می‌توان نتیجه گرفت برخی مقادیر ویژه، کوچک هستند. در نتیجه داده‌ها در راستای بردار ویژه‌ی متناظر با آن واریانس کمی دارد و قابل صرفنظر کردن است.

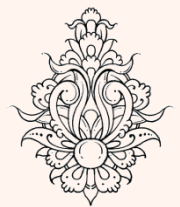
– در این حالت  $K$  مؤلفه‌ی پرارزش انتخاب می‌شوند، با فرض آن که مقادیر ویژه به صورت صعودی مرتب شده باشند.

**Proportion of Variance (PoV)**

**PoV > 0.9**

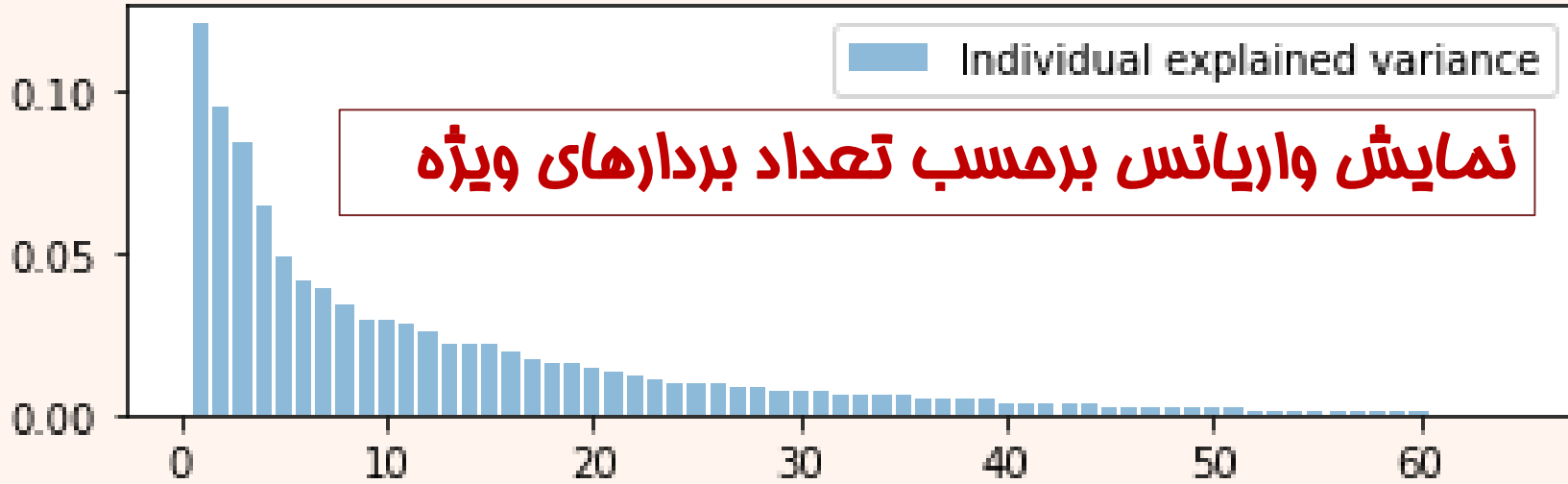
$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_k + \dots + \lambda_d}$$

– در کاربردهای نظیر پردازش تصویر یا صوت، معمولاً کاهش ابعاد قابل توجه است.

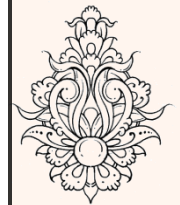
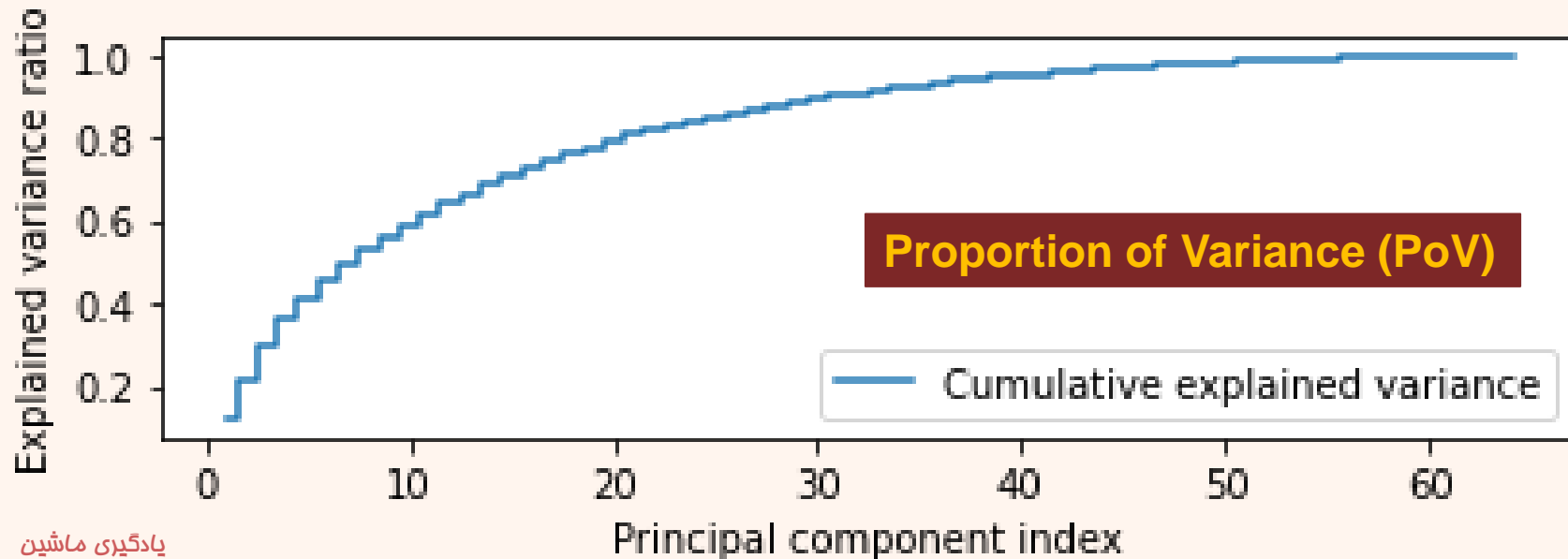


## Scree graph

### نمایش واریانس بر حسب تعداد بردارهای ویژه



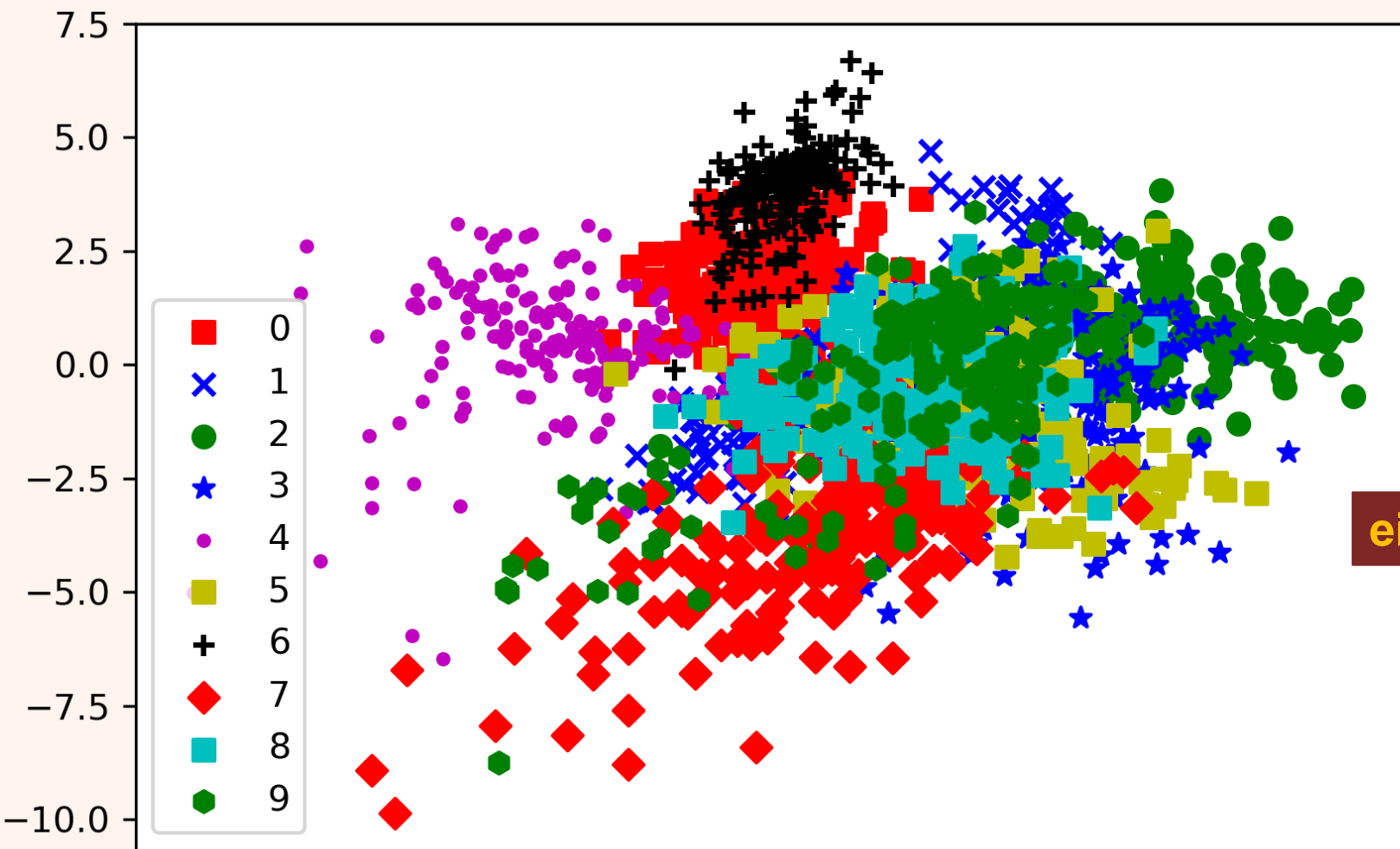
## Proportion of Variance (PoV)



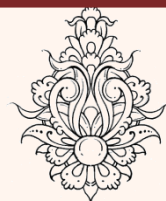
تراشگاه  
سپهر



# Visual analysis



eigendigits

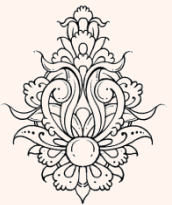


تازشگاه  
بهشتی

در صورتی که سه بعد نخست، ماوی بخش عمده‌ای از واریانس باشند، می‌توان داده‌ها از آن‌ها برای «بررسی دیداری» بهره برد.

## چند نکته

- علاوه بر در نظر گرفتن PoV، می‌توان بردارهای ویژه‌ای که مقدار ویژه‌ی متناظر آن از یک حدآستانه (به عنوان مثال میانگین واریانس) کمتر است را حذف نمود.
- در صورتی که واریانس در ابعاد مختلف تخییرات زیادی داشته باشند، بیش از مقدار همبستگی بر روی مؤلفه‌ی اصلی اثرگذار خواهد بود.
- در این شرایط می‌توان از بردارها و مقادیر ویژه‌ی «ماتریس همبستگی» (R) استفاده کرد یا این که داده‌ها را به گونه‌ای نرمال کرد که همگی واریانس یکسان داشته باشند.



# تحلیل مؤلفه‌های اصلی (ادامه...)

- از زاویه‌ی دیگری نیز می‌توان به این مسأله نگاه کرد؛ هدف یافتن ماتریس تبدیلی است که داده‌های را به گونه‌ای نگاشت کند که در فضای جدید «ناهمبسته» باشند.

$$\mathbf{z} = \mathbf{W}^T \mathbf{x} \quad \text{Cov}(\mathbf{z}) = \mathbf{D}' \quad \text{ماتریس قطری}$$

- $\mathbf{C}_{d \times d}$  ماتریسی است که ستون‌هایش بردارهای ویژه‌ی ماتریس کواریانس است:

$$\mathbf{C}^T \mathbf{C} = \mathbf{I}$$

$$\mathbf{S} = \mathbf{S} \mathbf{C} \mathbf{C}^T$$

ادامه



# تحلیل مؤلفه‌های اصلی (ادامه...)

$$\begin{aligned} \mathbf{S} &= \mathbf{S}\mathbf{C}\mathbf{C}^T \\ &= \mathbf{S}[\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_d]\mathbf{C}^T \\ &= [\mathbf{S}\mathbf{c}_1, \mathbf{S}\mathbf{c}_2, \dots, \mathbf{S}\mathbf{c}_d]\mathbf{C}^T \\ &= [\lambda_1\mathbf{c}_1, \lambda_2\mathbf{c}_2, \dots, \lambda_d\mathbf{c}_d]\mathbf{C}^T \\ &= \lambda_1\mathbf{c}_1\mathbf{c}_1^T + \lambda_2\mathbf{c}_2\mathbf{c}_2^T + \dots + \lambda_d\mathbf{c}_d\mathbf{c}_d^T \\ &= \mathbf{C}\mathbf{D}\mathbf{C}^T \end{aligned}$$

ماتریس قطری که عناصر روی قطر اصلی مقادیر ویژه‌ی ماتریس کواریانس هستند.

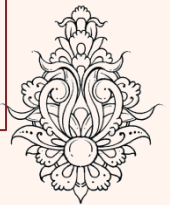
## Spectral decomposition

$$\mathbf{C}^T\mathbf{S}\mathbf{C} = \mathbf{D}$$

$$\mathbf{z} = \mathbf{W}^T \mathbf{x}, \quad \text{Cov}(\mathbf{z}) = \mathbf{W}^T \mathbf{S} \mathbf{W}$$

$$\mathbf{W} = \mathbf{C}$$

$$\text{Cov}(\mathbf{z}) = \mathbf{D}$$



# فاصله‌ی Mahalanobis

$$(x - \mu)^T \Sigma^{-1} (x - \mu) = c^2$$

$$\Sigma = \mathbf{Q} \mathbf{D} \mathbf{Q}^T$$

ماتریس شامل بردارهای ویژه  $\Sigma$  (در ستون‌ها)

$$\Sigma^{-1} = \mathbf{Q} \mathbf{D}^{-1} \mathbf{Q}^T$$

ماتریس شامل مقادیر ویژه  $\Sigma$  (قطر اصلی)

$$\Sigma^{-1} = \mathbf{Q} \mathbf{D}^{-1/2} \mathbf{D}^{-1/2} \mathbf{Q}^T$$

$$C^2 = (x - \mu)^T \mathbf{Q} \mathbf{D}^{-1/2} \mathbf{D}^{-1/2} \mathbf{Q}^T (x - \mu)$$

$$C^2 = [(x - \mu)^T \mathbf{Q} \mathbf{D}^{-1/2}] [\mathbf{D}^{-1/2} \mathbf{Q}^T (x - \mu)]$$

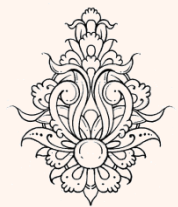
$$C^2 = [(x - \mu)^T \mathbf{Q} \mathbf{D}^{-1/2}] [(x - \mu)^T \mathbf{Q} \mathbf{D}^{-1/2}]^T$$

$$[(x - \mu)^T \mathbf{Q} \mathbf{D}^{-1/2}]$$

تبدیل:

۱- چرخش در راستای بردارهای ویژه

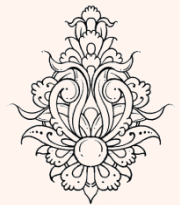
۲- تخریب مقیاس بر اساس مقادیر ویژه (واریانس در راستای بردار ویژه)



# چند نکته

- PCA، نسبت به نویز به شدت حساس است.  
– یک روش ساده حذف داده‌های پرت با استفاده از فاصله‌ی Mahalanobis پیش از محاسبه‌ی ماتریس کواریانس است.
- در میان تمام بردارهای متعامد، PCA کم‌ترین میزان خطا را دارد.  
**Reconstruction error**  $\sum_t \|\hat{\mathbf{x}}^t - \mathbf{x}^t\|$
- Hotelling transform و Karhunen-Loève expansion نام‌های دیگری هستند برای مفاهیم مشابه به کار می‌روند.
- در common principal components برای هم‌بندی کلاس‌ها مؤلفه‌های اساسی یکسانی در نظر گرفته می‌شود، با این تفاوت که کواریانس هر کلاس متفاوت در نظر گرفته می‌شود.

$$S_i = \mathbf{C} \mathbf{D}_i \mathbf{C}^T$$



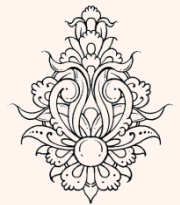
# کاربرد PCA در شناسایی چهره



پایگاه داده‌ی ORL

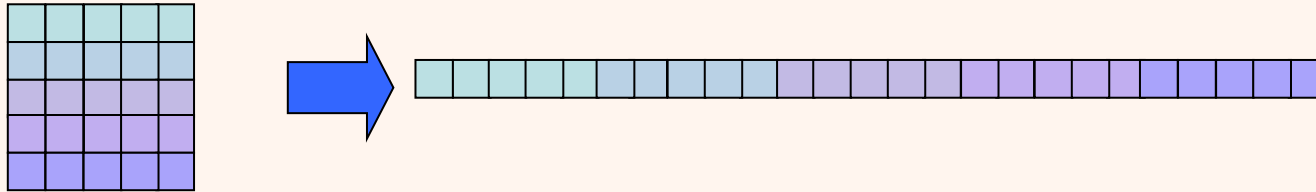


میانگین چهره‌ها

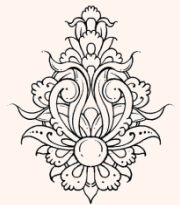
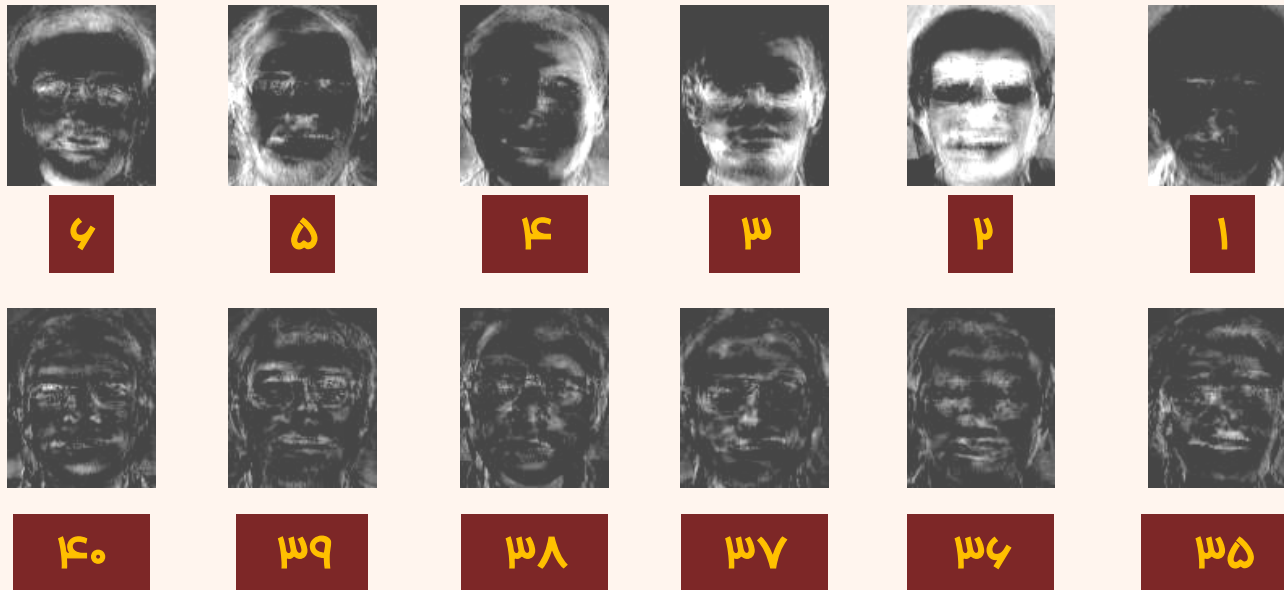


M. Turk, A. Pentland, "Eigenfaces for Recognition", *Journal of Cognitive Neuroscience*, vol. 3, no. 1, pp. 71-86, 1991.

# کاربرد PCA در شناسایی چهره (ادامه...)

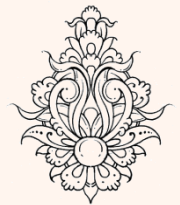
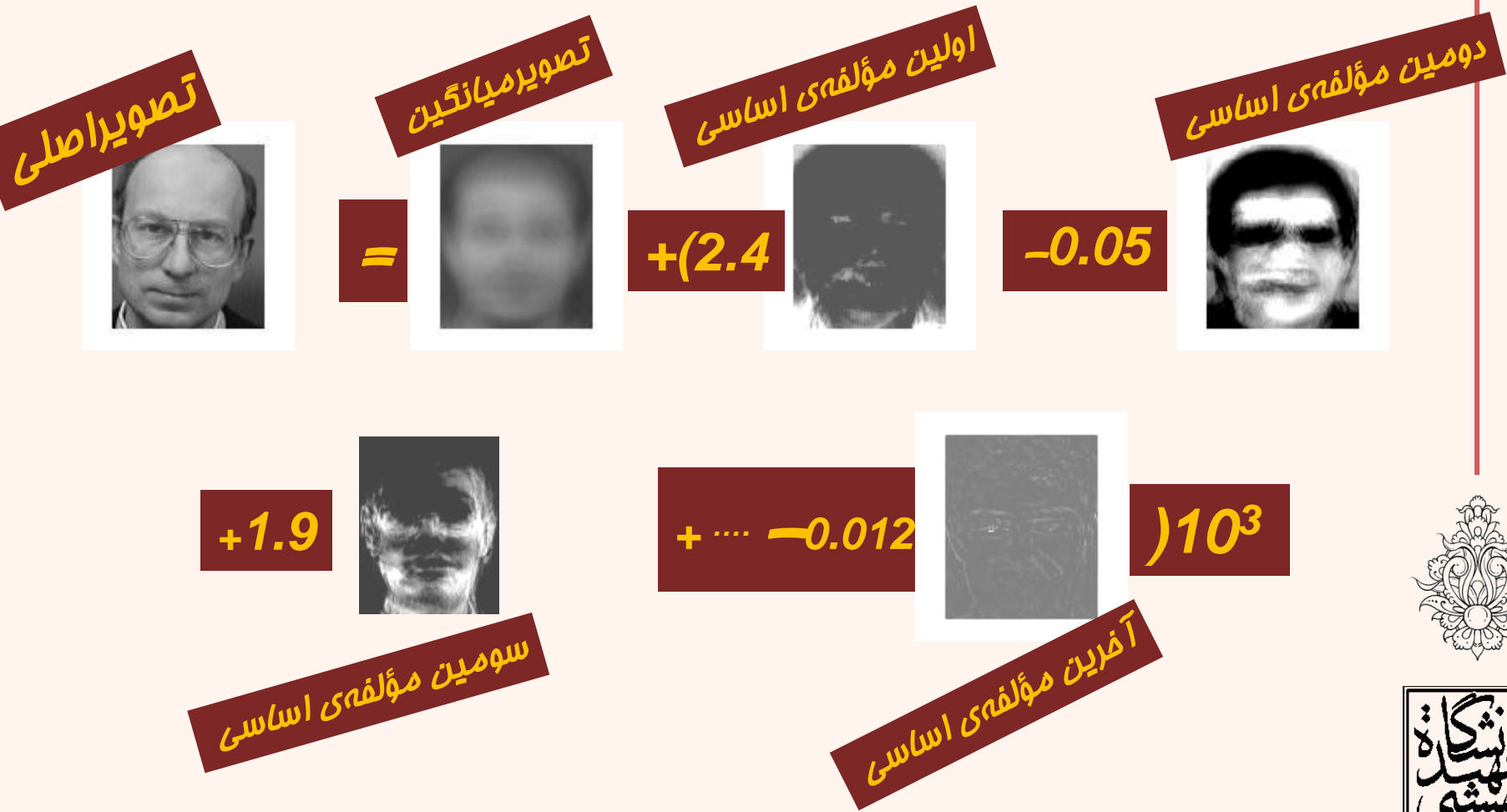


## Eigenfaces

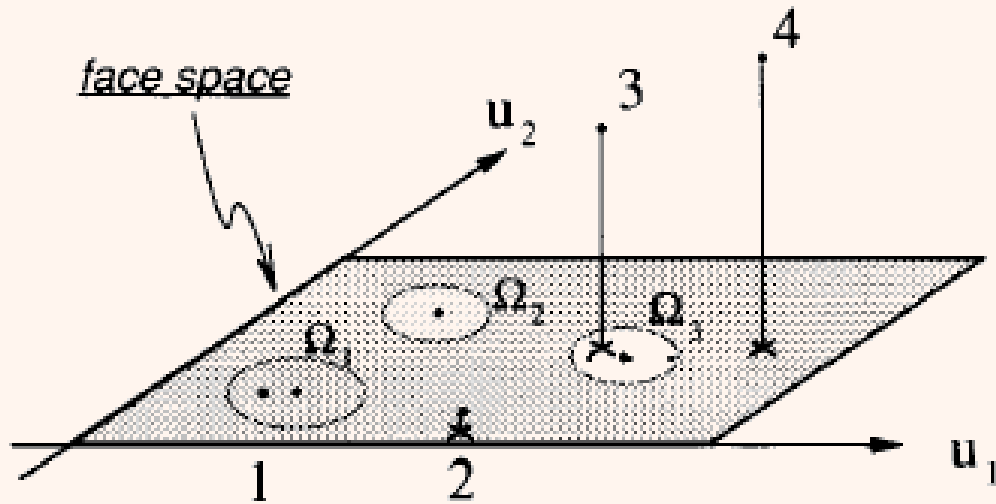




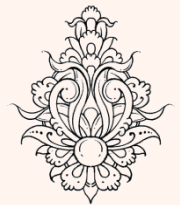
# کاربرد PCA در شناسایی چهره (ادامه...)



# تشخیص چهره

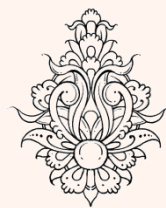


$$\|\hat{\mathbf{x}} - \mathbf{x}\|$$



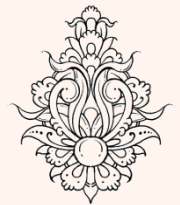
# Feature Embedding

- ماتریس داده‌ها به صورت  $\mathbf{X}_{N \times d}$  است.
- در صورت این که میانگین داده‌ها صفر باشد، ماتریس کواریانس خصیصه‌ها  $\mathbf{X}^T \mathbf{X}_{d \times d}$  می‌باشد، در نتیجه
- با ضرب طرفین در  $\mathbf{X}$   
$$(\mathbf{X}^T \mathbf{X}) \mathbf{w}_i = \lambda_i \mathbf{w}_i$$
- در نتیجه  $\mathbf{X} \mathbf{w}_i$  بردار ویژه  $\mathbf{X} \mathbf{X}^T_{N \times N}$  با مقدار ویژه  $\lambda_i$  است.
- در این حالت بردار ویژه، مختصات نمونه‌ها در راستای  $\mathbf{w}_i$  را نشان می‌دهد.  
$$(\mathbf{X} \mathbf{X}^T) \mathbf{X} \mathbf{w}_i = \lambda_i \mathbf{X} \mathbf{w}_i$$



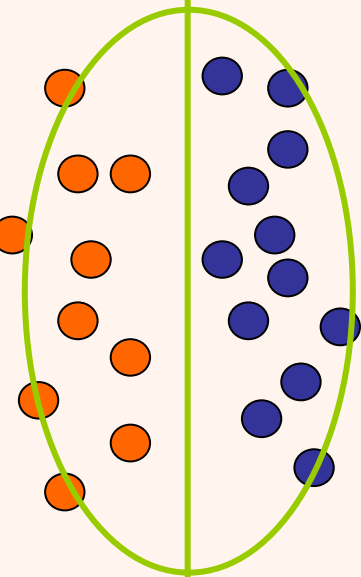
# Feature Embedding

- ثابت می‌شود، رتبه‌ی ماتریس حداکثر  $\min(d, N)$  می‌باشد.
- برای یک پایگاه داده حاوی چهل تصویر  $۲۵۶ \times ۲۵۶$ 
  - ماتریس کواریانس فصیصه‌ها  $۶۵۵۳۶ \times ۶۵۵۳۶$  خواهد بود.
  - در حالی که ماتریس شباهت نمونه‌ها  $۴۰ \times ۴۰$  می‌باشد.
  - این ماتریس، شباهت دوبه‌دو نمونه‌ها را نشان می‌دهد، از این نظر می‌توان گفت این شیوه داده‌های  $d$  بعدی را در فضای  $k$  بعدی به گونه‌ای قرار می‌دهد که فاصله‌ی بین آن‌ها حفظ می‌شود.
  - این ماتریس شامل ضرب داخلی دو به دو همه نمونه‌های آموزش است.





## دسته‌بندی دو کلاسه



• آیا PCA برای دسته‌بندی مناسب است؟

- راستای نگاشت بر اساس واریانس، انتخاب می‌شود.

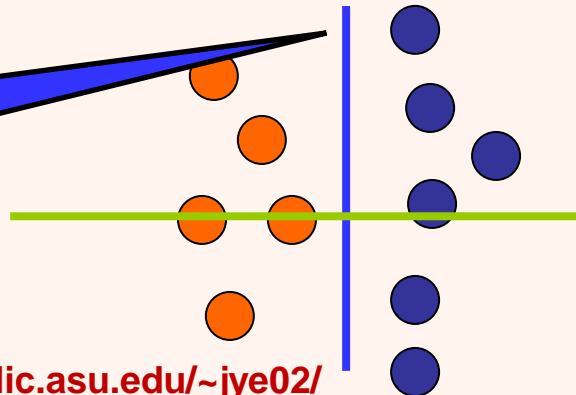
- در این میان ممکن است اطلاعات دسته‌ها از بین بروند.

• تحلیل تفکیک خطی، «بانظارت» است و برای دسته‌بندی به کار می‌رود.

- هدف آن کاهش بعد همراه با حفظ اطلاعاتی است که بین دسته‌ها تمایز قائل می‌شود.



در این راستا دو کلاس همپوشانی دارند

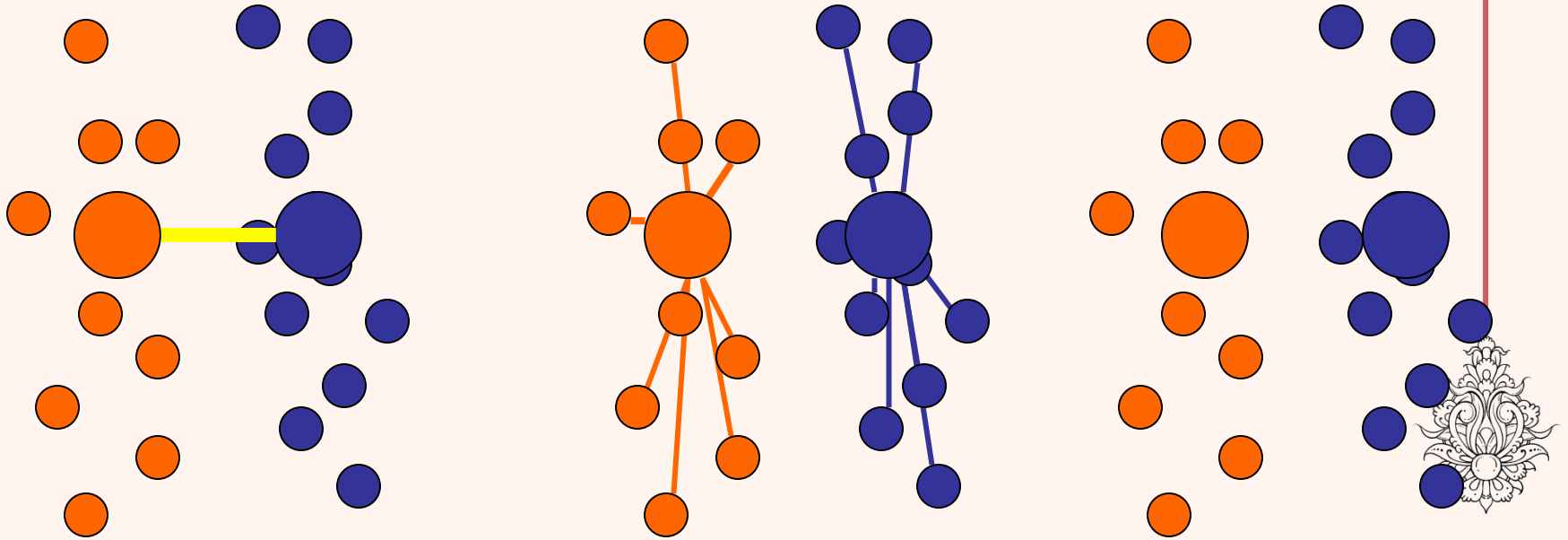


در این راستا دو کلاس بدون خطا دسته‌بندی می‌شوند

# کاهش ابعاد برای دسته‌بندی

دسته‌بندی دو کلاس

برای انتخاب راستای مناسب برای نگاشت، باید اطلاعات دسته‌ها نیز در نظر گرفته شود.



Between-class distance

Within-class distance



# تحلیل تفکیک خطی (ادامه...)

در LDA، نداشتن به گونه‌ای انجام می‌شود که فاصله‌ی بین دو کلاس حداکثر شده و فاصله‌ی نمونه‌های متعلق به یک کلاس مینیمم گردد.

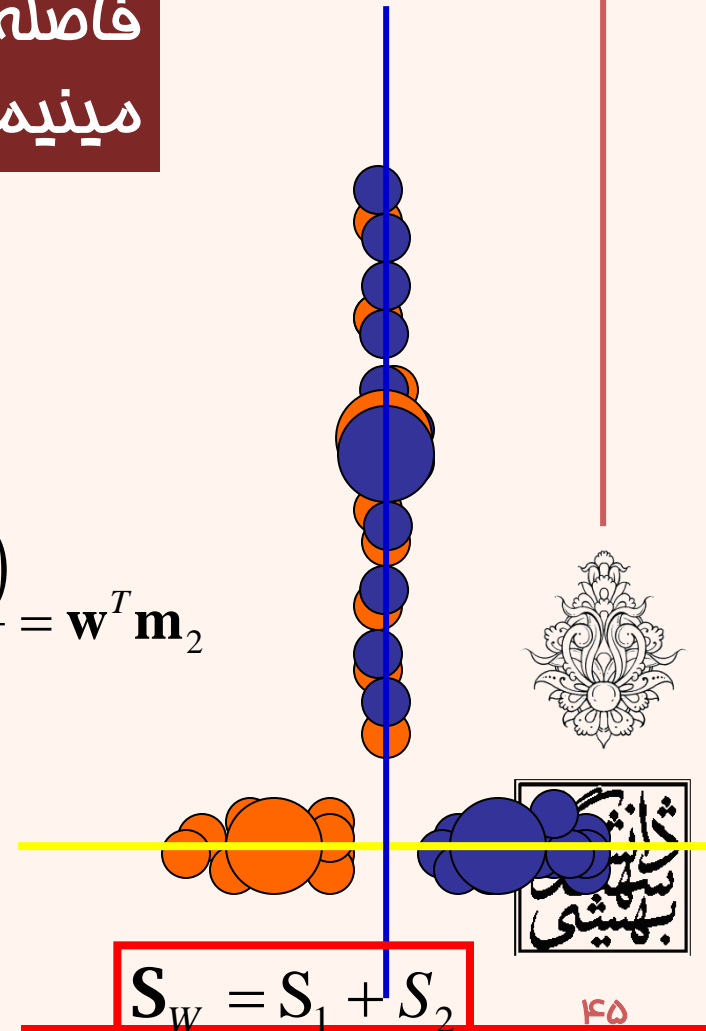
دسته‌بندی دو کلاس

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

$$m_1 = \frac{\sum_t \mathbf{w}^T \mathbf{x}^t r^t}{\sum_t r^t} = \mathbf{w}^T \mathbf{m}_1, \quad m_2 = \frac{\sum_t \mathbf{w}^T \mathbf{x}^t (1 - r^t)}{\sum_t (1 - r^t)} = \mathbf{w}^T \mathbf{m}_2$$

$$s_1^2 = \sum_t (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t$$

$$s_2^2 = \sum_t (\mathbf{w}^T \mathbf{x}^t - m_2)^2 (1 - r^t)$$



$$S_W = S_1 + S_2$$

۴۵

Total within Class scatter

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

# تحلیل تفکیک خطی (ادامه...)

دسته‌بندی دو کلاس

$$\begin{aligned} (m_1 - m_2)^2 &= (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 \\ &= \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} \\ &= \mathbf{w}^T \mathbf{S}_B \mathbf{w} \end{aligned}$$

Between class scatter matrix

$$s_1^2 = \sum_t (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t$$

$$\mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$$

$$= \sum_t \mathbf{w}^T (\mathbf{x}^t - \mathbf{m}_1)(\mathbf{x}^t - \mathbf{m}_1)^T \mathbf{w} r^t$$

Class scatter matrix for C1

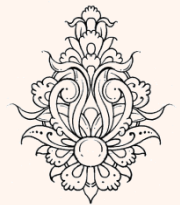
$$= \sum_t \mathbf{w}^T \mathbf{S}_1 \mathbf{w} r^t$$

$$\mathbf{S}_1 = \sum_t r^t (\mathbf{x}^t - \mathbf{m}_1)(\mathbf{x}^t - \mathbf{m}_1)^T$$

$$s_1^2 + s_2^2 = \mathbf{w}^T \mathbf{S}_W \mathbf{w}$$

$$\mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$$

Total within Class scatter





# تحلیل تفکیک خطی (ادامه...)

دسته‌بندی دو کلاس

- هدف ماکزیمیم کردن رابطه‌ی زیر است:

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} = \frac{|\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)|^2}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

- با مشتق گرفتن:

$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \left( 2(\mathbf{m}_1 - \mathbf{m}_2) - \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \mathbf{S}_W \mathbf{w} \right) = 0$$

Scalar

$$\mathbf{w} = c \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

یادآوری: جداساز خطی

$$\mathbf{w} = \Sigma^{-1} (\mu_1 - \mu_2)$$

when  $p(\mathbf{x} | C_i) \sim \mathcal{N}(\mu_i, \Sigma)$

بدین ترتیب، برای دسته‌های نرمال، LDA جداساز بهینه است.



# دسته‌بندی برای بیش از دو کلاس

- زمانی که تعداد کلاس‌ها بیشتر از دو باشد: برای کاهش ابعاد، ماتریس  $W_{d \times k-1}$  برای نگاشت مورد استفاده قرار می‌گیرد:

$$\mathbf{z} = \mathbf{W}^T \mathbf{x}$$

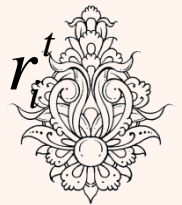
$$\mathbf{S}_W = \sum_{i=1}^K \mathbf{S}_i \quad \mathbf{S}_i = \sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T$$

Within-class scatter

Between-class scatter:

$$\mathbf{S}_B = \sum_{i=1}^K N_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^T \quad \mathbf{m} = \frac{1}{K} \sum_{i=1}^K \mathbf{m}_i \quad N_i = \sum_t r_i^t$$

- پس از نگاشت،  $W^T \mathbf{S}_W W$  و  $W^T \mathbf{S}_B W$  ماتریس‌های پراکندگی داده «بین‌دسته‌ها» و «درون‌دسته‌ها» خواهند بود.

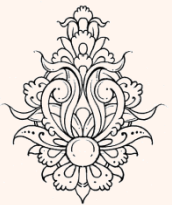


# دسته‌بندی برای بیش از دو کلاس

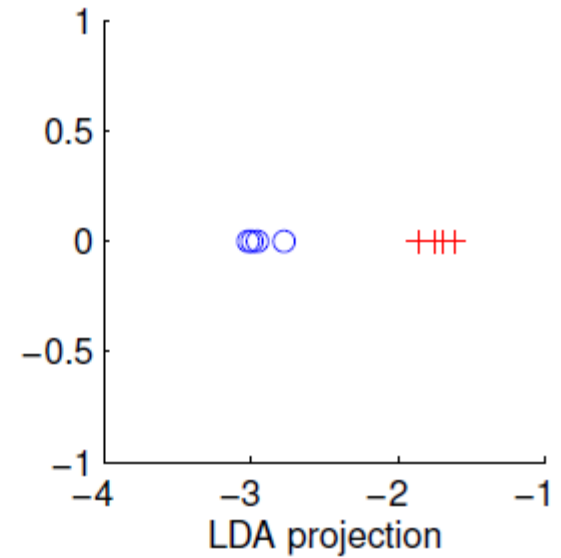
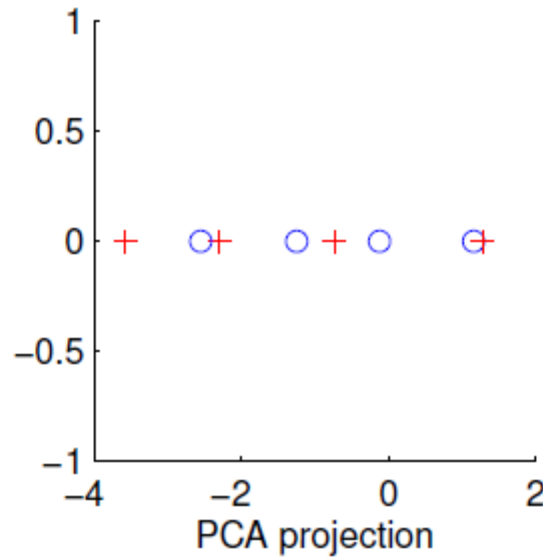
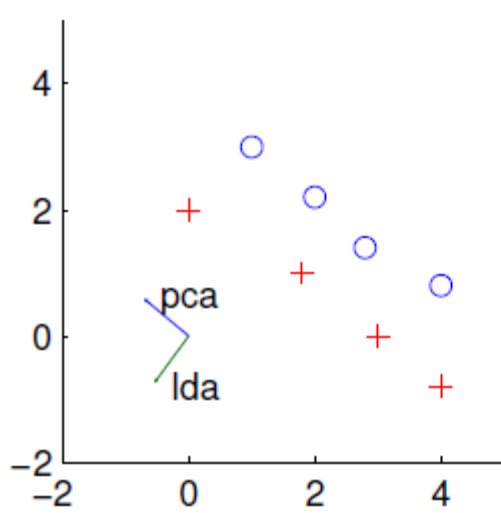
- در نتیجه در صورت پیشینه شدن عبارت زیر، دسته‌بندی به بهترین شکل انجام می‌شود.
- برای ماتریس کواریانس، دترمینان معیاری است که پراکندگی داده را نشان می‌دهد.

$$J(W) = \frac{|W^T S_B W|}{|W^T S_W W|}$$

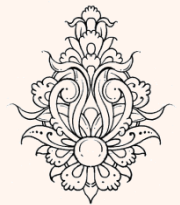
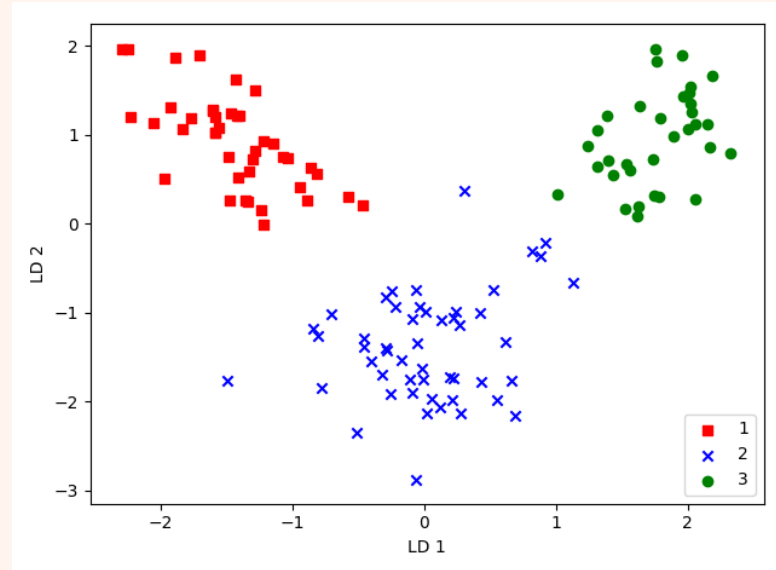
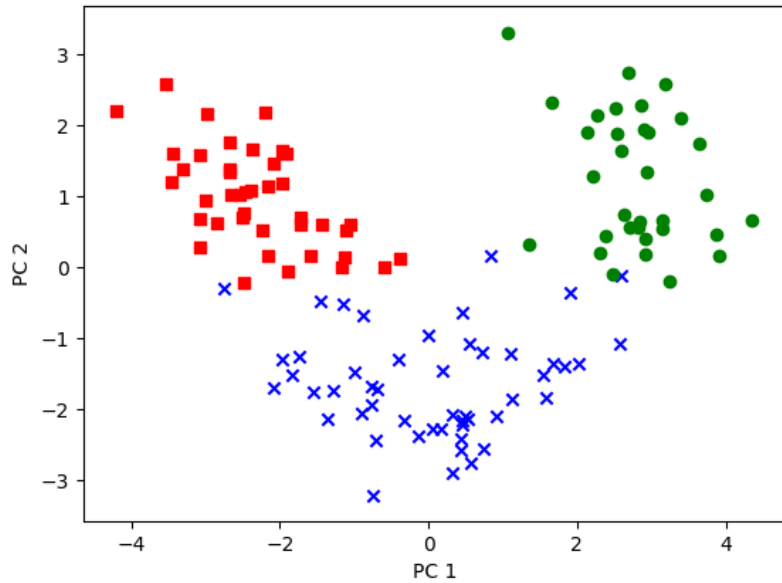
- در این حالت پاسخ، بردارهای ویژه متناظر با بزرگ‌ترین مقادیر ویژه‌ی ماتریس  $S_W^{-1} S_B$  خواهد بود.

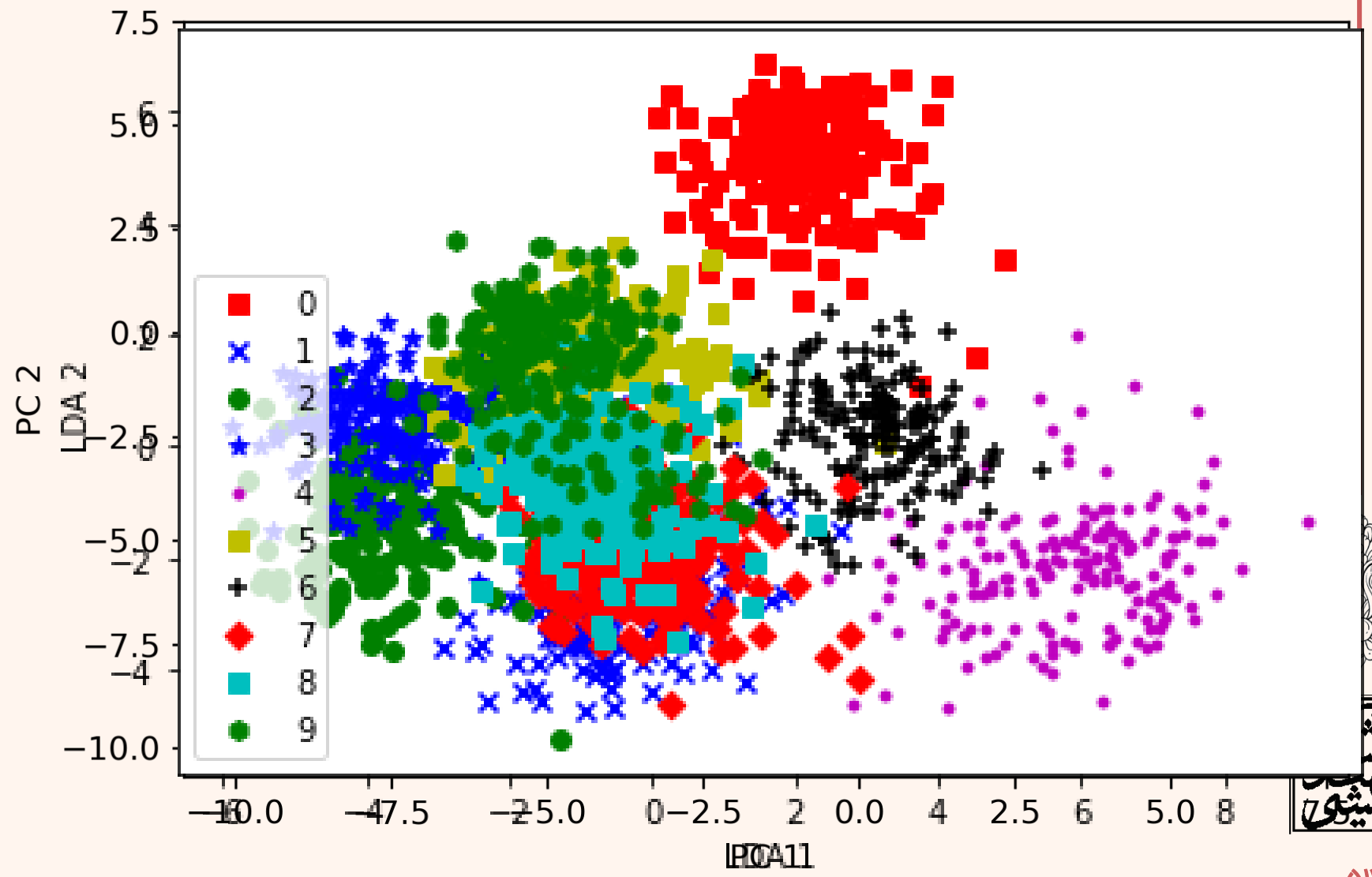


# PCA vs LDA



# PCA vs LDA





## تجزیه‌ی مقادیرهای تکین

- با استفاده از SVD، یک ماتریس به سه ماتریس تجزیه می‌شود:

$$\mathbf{X}_{N \times d} = \mathbf{V}_{N \times N} \mathbf{A}_{N \times d} \mathbf{W}_{d \times d}^T$$

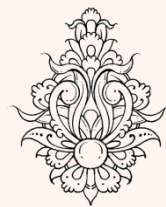
- $\mathbf{V}$  شامل بردارهای ویژه‌ی  $\mathbf{X}\mathbf{X}^T$  می‌باشد،  $\mathbf{W}$  شامل بردارهای ویژه‌ی  $\mathbf{X}^T\mathbf{X}$  است و  $\mathbf{A}$  مقادیر ویژه را در  $k$  عنصر قطری خود دارد.

$$\mathbf{X}\mathbf{X}^T = (\mathbf{V}\mathbf{A}\mathbf{W}^T)(\mathbf{V}\mathbf{A}\mathbf{W}^T)^T = \mathbf{V}\mathbf{A}\mathbf{W}^T\mathbf{W}\mathbf{A}^T\mathbf{V}^T = \mathbf{V}\mathbf{E}\mathbf{V}^T$$

$$\mathbf{X}^T\mathbf{X} = (\mathbf{V}\mathbf{A}\mathbf{W}^T)^T(\mathbf{V}\mathbf{A}\mathbf{W}^T) = \mathbf{W}\mathbf{A}^T\mathbf{V}^T\mathbf{V}\mathbf{A}\mathbf{W}^T = \mathbf{W}\mathbf{D}\mathbf{W}^T$$

$$\mathbf{E} = \mathbf{A}\mathbf{A}^T$$

$$\mathbf{D} = \mathbf{A}^T\mathbf{A}$$

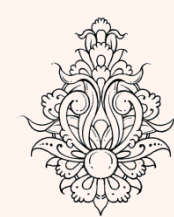
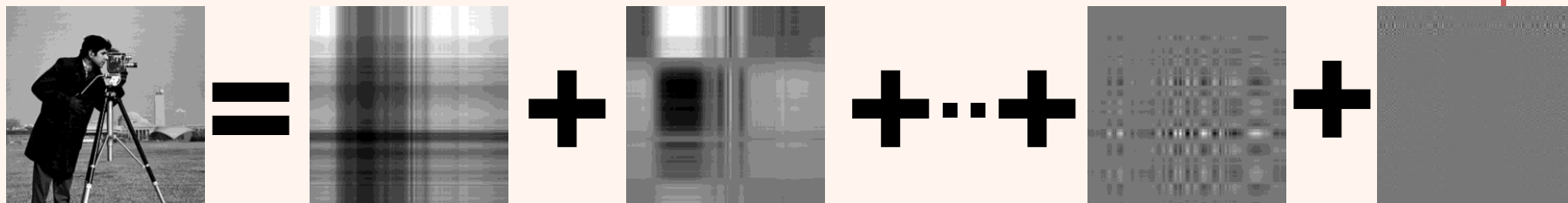


$$\mathbf{X}_{N \times d} = \mathbf{V}_{N \times N} \mathbf{A}_{N \times d} \mathbf{W}_{d \times d}^T$$

# تجزیه‌ی مقدارهای تکین

$$\mathbf{X} = a_1 \mathbf{v}_1 \mathbf{w}_1^T + \dots + a_k \mathbf{v}_k \mathbf{w}_k^T$$

$$\mathbf{X} = \sum_{i=1}^k a_i \mathbf{v}_i \mathbf{w}_i^T$$





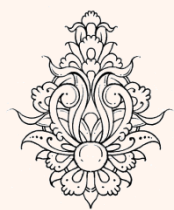


|          | 0             | 1             | 2             | 3             | 4             | 5             | 6             |
|----------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| NLP      | -3.471975e-01 | 0.000000e+00  | -1.942887e-01 | 0.000000e+00  | 0.000000e+00  | -4.183768e-01 | 8.164966e-01  |
| flying   | -1.396622e-01 | 0.000000e+00  | 6.739877e-01  | -9.992007e-16 | -4.440892e-16 | 5.996402e-01  | 4.082483e-01  |
| enjoy    | 1.249001e-16  | -3.687707e-01 | -2.220446e-16 | -8.041284e-01 | 4.662464e-01  | -1.110223e-16 | -5.551115e-17 |
| I        | -8.340573e-01 | -5.277047e-18 | 2.854104e-01  | -7.948171e-16 | -1.162003e-15 | -2.371134e-01 | -4.082483e-01 |
| deep     | -4.053355e-01 | 0.000000e+00  | -6.530953e-01 | 0.000000e+00  | 0.000000e+00  | 6.396638e-01  | -1.110223e-16 |
| like     | 1.665335e-16  | -9.167567e-01 | -2.775558e-17 | 2.318040e-01  | -3.253062e-01 | 5.551115e-17  | 1.110223e-16  |
| learning | -1.318390e-15 | -1.535102e-01 | 1.221245e-15  | 5.473979e-01  | 8.226726e-01  | 3.885781e-16  | -5.551115e-17 |

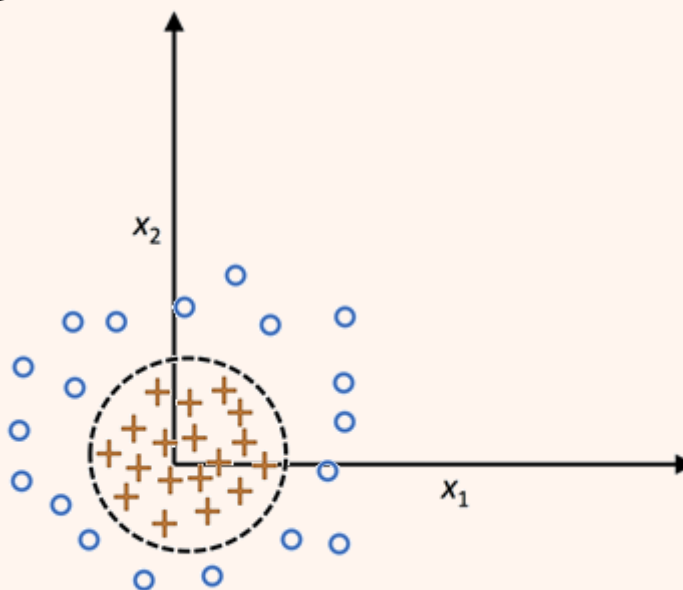
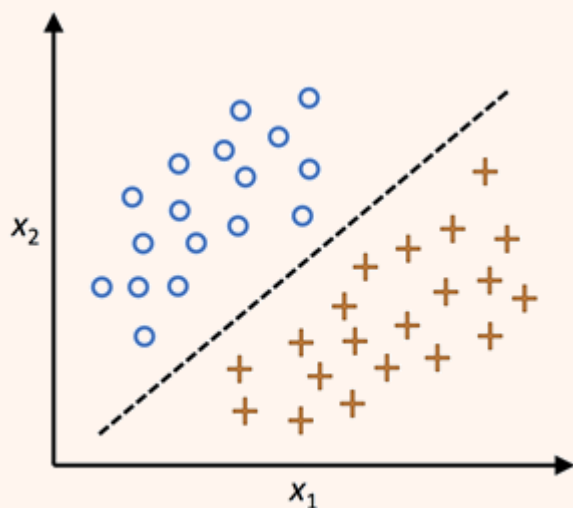
Matrix U (before selecting k singular features)

|          | 0             | 1             | 2             |
|----------|---------------|---------------|---------------|
| NLP      | -3.471975e-01 | 0.000000e+00  | -1.942887e-01 |
| flying   | -1.396622e-01 | 0.000000e+00  | 6.739877e-01  |
| enjoy    | 1.249001e-16  | -3.687707e-01 | -2.220446e-16 |
| I        | -8.340573e-01 | -5.277047e-18 | 2.854104e-01  |
| deep     | -4.053355e-01 | 0.000000e+00  | -6.530953e-01 |
| like     | 1.665335e-16  | -9.167567e-01 | -2.775558e-17 |
| learning | -1.318390e-15 | -1.535102e-01 | 1.221245e-15  |

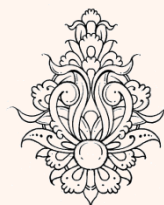
Matrix U (after selecting k singular features)



- زمانی که داده‌ها جدایی‌پذیر خطی نیستند، استفاده از روش‌ها خطی برای کاهش بعد موجب از دست رفتن اطلاعات کلاس‌ها می‌شود.



در این حالت نیز می‌توان داده‌ها را به فضایی با ابعاد بالاتر نگاشت کرد به گونه‌ای که در آن فضا به صورت خطی قابل جداسازی باشند و پس از آن از روش‌های سنتی کاهش بعد بهره برد.



## Kernel principal component analysis

• در این حالت نیز می‌توان از **kernel trick** استفاده کرد.

– با توجه به این که در Feature embedding از ماتریس تشابه (ضرب داخلی) نمونه‌های آموزشی استفاده می‌شود با استفاده از توابع هسته‌ی مناسب به راحتی و بدون مواجه شدن با نگاشت به فضایی دیگر، کاهش ابعاد به شیوه‌ی غیرخطی قابل انجام خواهد بود.

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^k \quad (k \gg d)$$

توابع کرنل ضرب داخلی دو بردار ویژگی در فضایی با ابعاد بالا را بدون نگاشت با آن فضا به دست می‌آورد:

$$k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)})$$



## توابع کرنل متعارف:

### Radial Basis function(RBF)

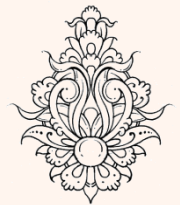
$$\kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp\left(-\gamma \|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2\right)$$

$$\kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \left(\mathbf{x}^{(i)T} \mathbf{x}^{(j)} + \theta\right)^p$$

### polynomial kernel

### *Hyperbolic Tangent Kernel*

$$\kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \tanh\left(\eta \mathbf{x}^{(i)T} \mathbf{x}^{(j)} + \theta\right)$$



# حذف میانگین

$$k_{ij} = \phi_i^T \phi_j$$

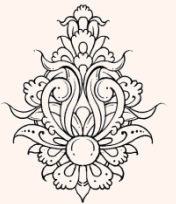
- باید توجه داشت که لازم است نمونه‌ها در فضای جدید به گونه‌ای نرمال شود که میانگین آنها

صفر باشد:

$$k'_{ij} = \left( \phi_i - \frac{1}{N} \sum_m \phi_m \right)^T \left( \phi_j - \frac{1}{N} \sum_n \phi_n \right)$$

$$k'_{ij} = \phi_i^T \phi_j - \phi_i^T \left[ \frac{1}{N} \sum_n \phi_n \right] - \left[ \frac{1}{N} \sum_m \phi_m^T \right] \phi_j + \left[ \frac{1}{N} \sum_n \phi_n \right] \left[ \frac{1}{N} \sum_m \phi_m^T \right]$$

$$k'_{ij} = k_{ij} - \left[ \frac{1}{N} \sum_n k_{in} \right] - \left[ \frac{1}{N} \sum_m k_{mj} \right] + \frac{1}{N^2} \sum_m \sum_n k_{mn}$$



$$\mathbf{1}_{N[i,j]} = [1/N]_{i,j}$$

$$\mathbf{K}' = \mathbf{K} - \mathbf{1}_N \mathbf{K} - \mathbf{K} \mathbf{1}_N + \mathbf{1}_N \mathbf{K} \mathbf{1}_N$$



# مثال گام به گام

```
def rbf_kernel_pca(X, gamma, n_components):
```

```
    sq_dists = pdist(X, 'sqeuclidean')
    mat_sq_dists = squareform(sq_dists)
    K = exp(-gamma * mat_sq_dists)
```

$$\kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma \|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2)$$

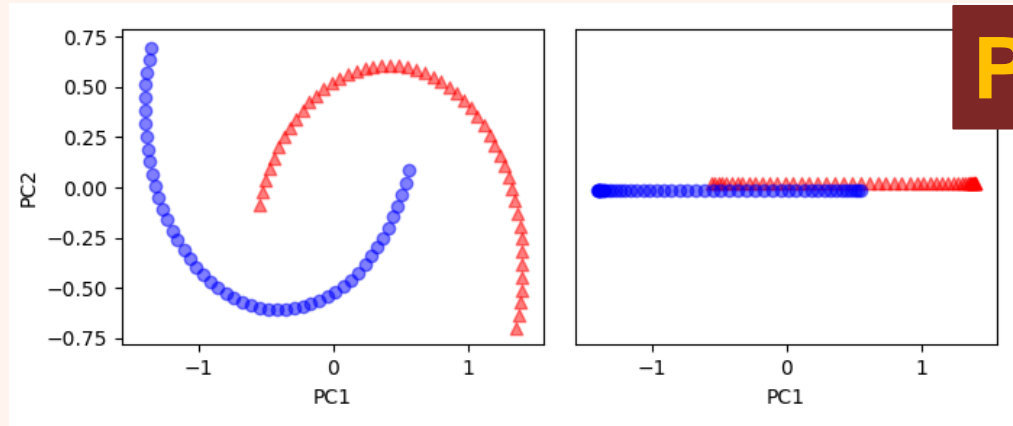
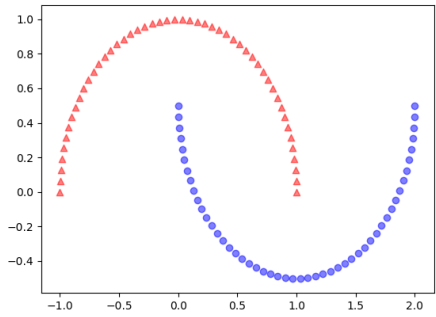
```
    N = K.shape[0]
    one_n = np.ones((N, N)) / N
    K = K - one_n.dot(K) - K.dot(one_n) + one_n.dot(K).dot(one_n)
```

$$\mathbf{K}' = \mathbf{K} - \mathbf{1}_n \mathbf{K} - \mathbf{K} \mathbf{1}_n + \mathbf{1}_n \mathbf{K} \mathbf{1}_n$$

```
    eigvals, eigvecs = eigh(K)
    eigvals, eigvecs = eigvals[::-1], eigvecs[:, ::-1]
```

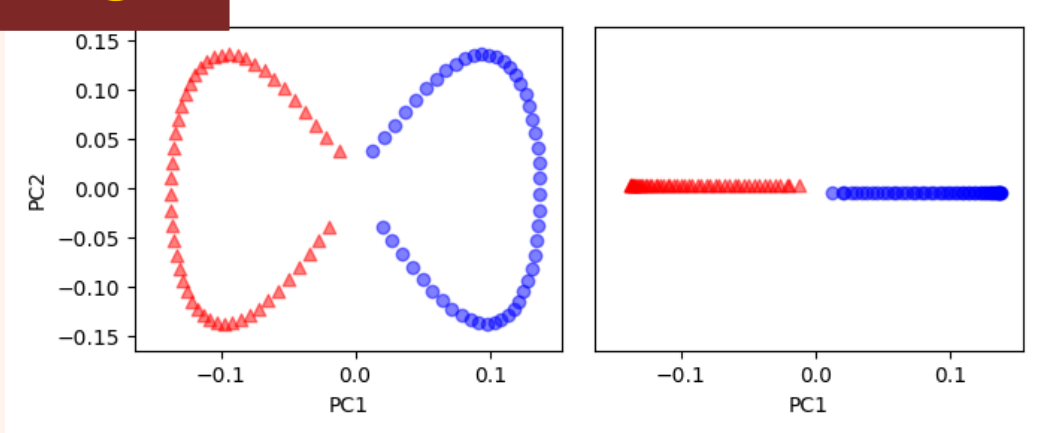
```
    X_pc = np.column_stack((eigvecs[:, i]
                             for i in range(n_components)))
```

```
    return X_pc
```



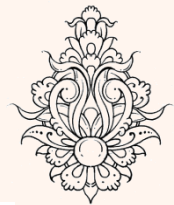
**PCA**

**KPCA**



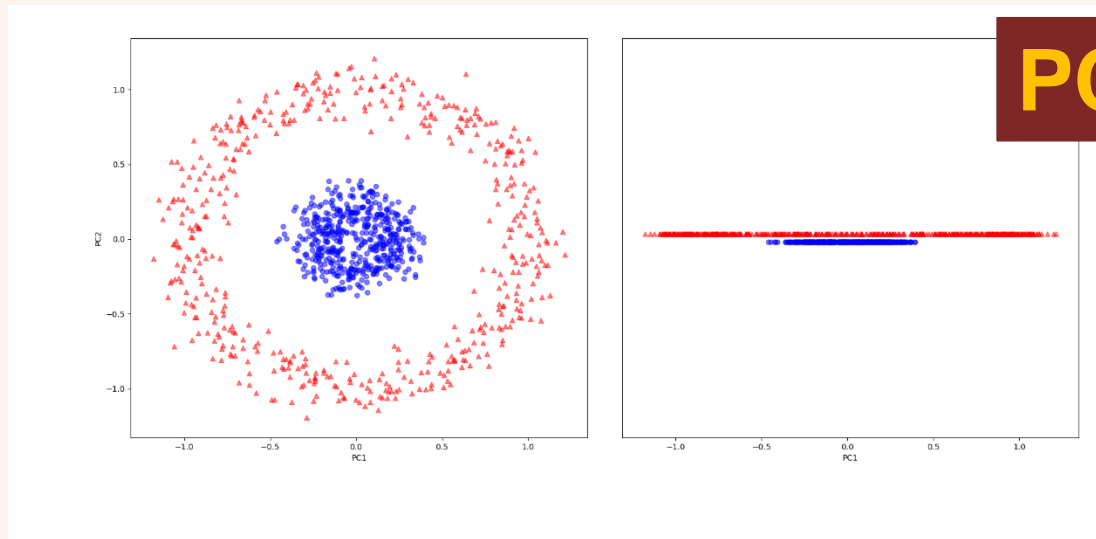
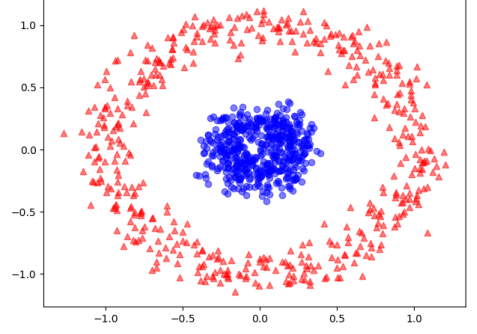
```
from sklearn.decomposition import KernelPCA
```

```
X, y = make_moons(n_samples=100, random_state=123)
scikit_kpca = KernelPCA(n_components=2, kernel='rbf', gamma=15)
X_skernpca = scikit_kpca.fit_transform(X)
```



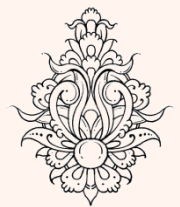
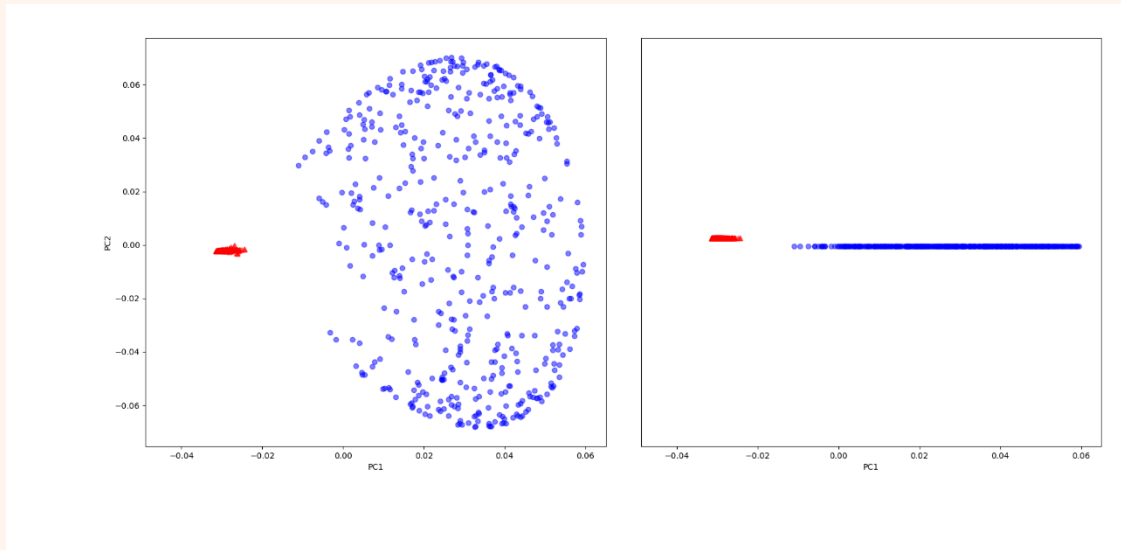
تراشگاه  
سپهر  
بهشتی





**PCA**

**KPCA**



تراشگاه  
سپهر  
بهشتی