

Sparse PCA

The first principal component is $z_1 = X v_1$

$$z_1 = \begin{array}{|c|c|c|c|} \hline | & | & | & | \\ \hline f_1 & f_2 & \dots & f_{p-1} & f_p \\ \hline | & | & & | & | \\ \hline \end{array} \begin{array}{c} v_{11} \\ v_{21} \\ \vdots \\ v_{p1} \end{array} = v_{11} f_1 + v_{21} f_2 + \dots + v_{p-1} f_{p-1} + v_{p1} f_p$$

X v_1

What if we want to select a subset of $\{f_1, \dots, f_p\}$ important for summarizing the information in X ?

SCoTLASS by Joliffe et. al. (2003):

$$\text{maximize } v^T S v \quad \text{subject to } \begin{cases} \|v\|^2 = 1 \\ \|v\|_1 \leq c \end{cases}$$

If $c \uparrow \infty \Rightarrow$ PCA. If $c \downarrow 0 \Rightarrow v = 0$.

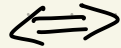
$$1 \leq c \leq \sqrt{p} \quad \text{as } \|v\|_2 \leq \|v\|_1 \leq \sqrt{p} \|v\|_2$$

If $v = (1, 0, \dots, 0) \Rightarrow \|v\|_1 = 1$. If $v = (\frac{1}{\sqrt{p}}, \dots, \frac{1}{\sqrt{p}}) \Rightarrow \|v\|_1 = \sqrt{p}$

Power iteration method

Step 1 $\tilde{v} = \frac{X^T X}{n-1} v$

Step 2 $v = \frac{\tilde{v}}{\|\tilde{v}\|_2}$



At iteration $t+1$:

$$v^{(t+1)} = \frac{X^T X v^{(t)}}{\|X^T X v^{(t)}\|_2}$$

Penalized matrix decomposition by Witten et al (2009):

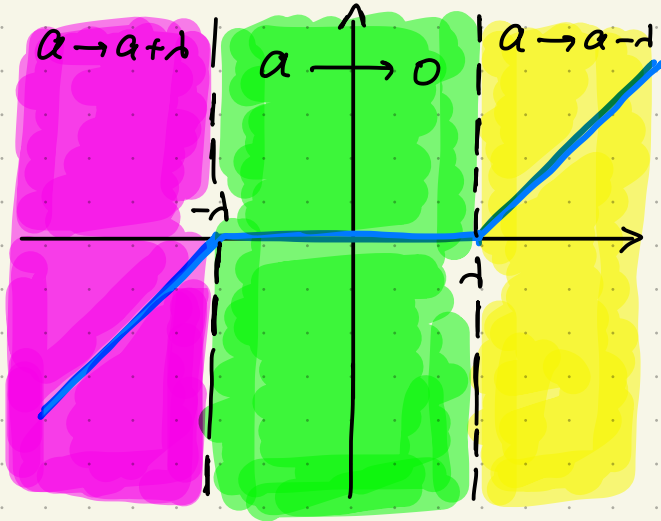
At iteration $t+1$: $v^{(t+1)} = \frac{S_\lambda(X^T X v)}{\|S_\lambda(X^T X v)\|_2}$

where $S_\lambda(a) = \text{sign}(a)(|a| - \lambda)_+$ is

soft-thresholding operator applied coordinatewise.

and λ is such that $\|v^{(t+1)}\|_1 = C$.

$$S_d(a) = \text{sign}(a) (|a| - d)_+$$



Note that

- $S_d(a) = a$ for $d = 0$
- $S_d(a/d) = S_d(a)$

$$a \geq d \quad \text{sign}(a) = 1, \quad |a| = a > d \Rightarrow S_d(a) = 1 \cdot (a - d)$$

$$-d \leq a \leq d \quad |a| < d \Rightarrow |a| - d < 0 \Rightarrow S_d(a) = 0$$

$$a \leq -d \quad \text{sign}(a) = -1, \quad |a| = -a > d \Rightarrow S_d(a) = -1 \cdot (-a - d) = a + d$$

Sparse SVD

Penalized matrix decomposition by Witten et al (2009):

Given $X \in \mathbb{R}^{n \times p}$ find $d \in \mathbb{R}$, $u \in \mathbb{R}^n$, $v \in \mathbb{R}^p$

minimize $\|X - d u v^T\|_F$ subject to $\begin{cases} \|u\|_2 = 1 \\ \|u\|_1 \leq C_1 \\ \|v\|_2 = 1 \\ \|v\|_1 \leq C_2 \end{cases}$
 u, d, v

- If C_1 and C_2 are very large then u and v are singular vectors (see HW2)
- Why not just $\|u\|_1 = 1$ and $\|v\|_1 = 1$?

Let's denote by (i, j) the index with the largest $|x_{ij}|$. Then $u = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}_i$ $v = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}_j$ $d = x_{ij}$

$X = \begin{pmatrix} - & x_{ij} & - \\ | & | & | \end{pmatrix} \approx \begin{pmatrix} 0 & \dots & 0 & \dots & 0 \\ 0 & \dots & x_{ij} & \dots & 0 \\ 0 & \dots & 0 & \dots & 0 \end{pmatrix} = x_{ij} \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}_i (0 \dots 0 \underset{j}{1} 0 \dots 0)$

- For fixed u you need to

$$\underset{v}{\text{maximize}} \quad u^T X v \quad \text{subject to} \quad \begin{cases} \|v\|_2 = 1 \\ \|v\|_1 \leq c_2 \end{cases}$$

- Lagrangian is

$$\mathcal{L}(v, \lambda, \mu) = u^T X v - \lambda (\|v\|_1 - c_2) - \frac{\mu}{2} (v^T v - 1)$$

- One can show that optimal u and v are:

$$v = \frac{S_\lambda(X^T u)}{\mu} \quad \text{where } S_\lambda(a) \text{ is soft-thresholding}$$

- To enforce $\|v\|_2 = 1$ we need

$$\mu = \|S_\lambda(X^T u)\|_2 \quad \text{so} \quad v = \frac{S_\lambda(X^T u)}{\|S_\lambda(X^T u)\|_2} = v(\lambda)$$

- We need to find λ such that $\|v(\lambda)\|_1 \leq c_2$

- If we don't have additional constraint $\|u\|_1 \leq c$, we will get **sparse PCA**.

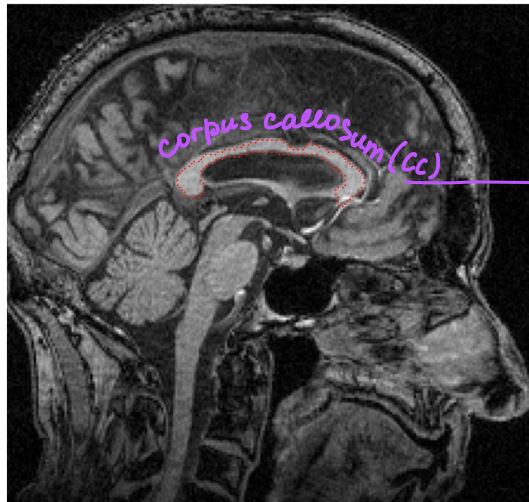
At iteration $t+1$: Step 1 $u^{(t+1)} = \frac{XV^{(t)}}{\|XV^{(t)}\|_2}$

Step 2 $v^{(t+1)} = \frac{S_\lambda(X^T u^{(t+1)})}{\|S_\lambda(X^T u^{(t+1)})\|_2} = \frac{S_{\lambda'}(X^T X v^{(t)})}{\|S_{\lambda'}(X^T X v^{(t)})\|_2}$

$$| S_\lambda(X^T X v / \|Xv\|_2) = S_{\lambda'}(X^T X v) \text{ for } \lambda' = \lambda \cdot \|Xv\|_2$$

- Note that if $\lambda = 0$ then $v^{(t+1)} = \frac{X^T X v^{(t)}}{\|X^T X v^{(t)}\|_2}$ that is, power iteration method.
- After finding u_1, d_1, v_1 you can apply method to $\tilde{X} = X - d_1 u_1 v_1^T$ and find $u_2, d_2, v_2 \dots$

Sparse PCA on shapes (from ESL11)

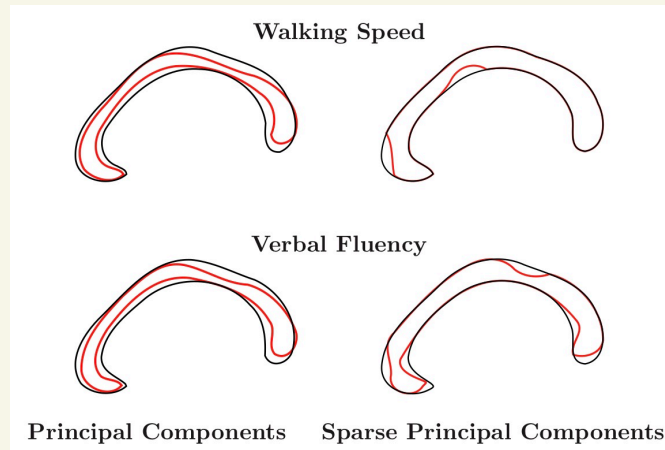


CC annotated with landmarks

$$\begin{pmatrix} x_1, y_1 \\ \dots \\ x_N, y_N \end{pmatrix} \rightarrow (x_1, y_1, \dots, x_N, y_N)$$
$$X = \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{pmatrix} \leftarrow n=569 \text{ elderly people}$$

black: the mean CC shape

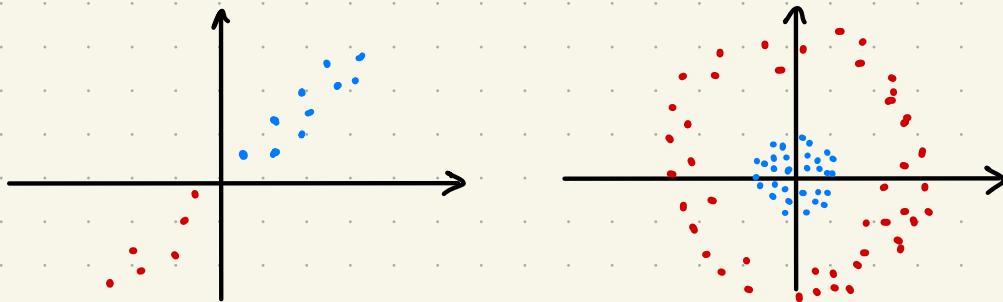
red: PC loading vectors



Kernel PCA

The main limitation of PCA is **linearity**

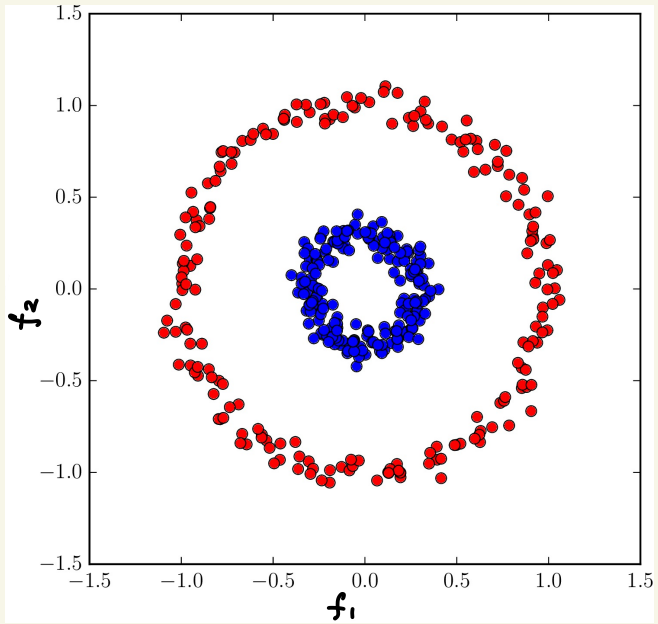
$$| Z_i = X V_i = \begin{pmatrix} | & & | \\ f_1 & \dots & f_p \\ | & & | \end{pmatrix} \begin{pmatrix} v_{1i} \\ \vdots \\ v_{pi} \end{pmatrix} = f_1 v_{1i} + \dots + f_p v_{pi} \quad (\text{linear function of } f_1, \dots, f_p)$$



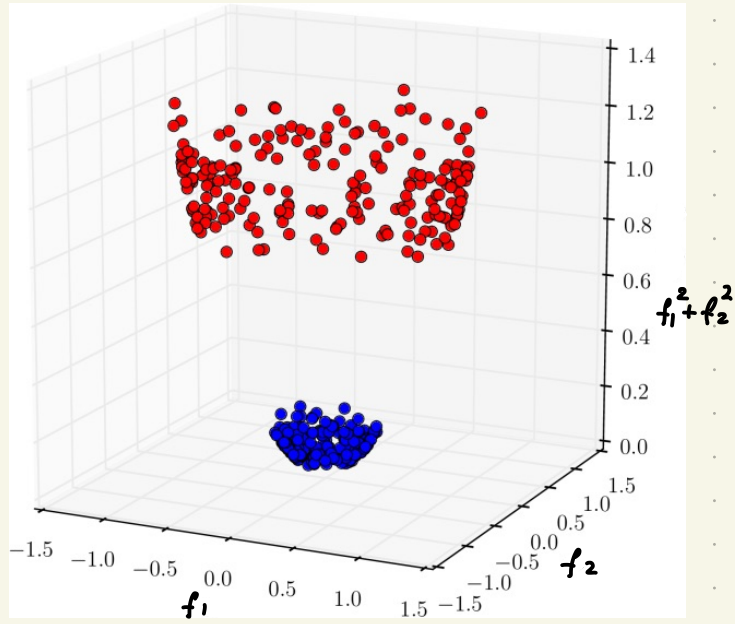
Idea: transform the feature space, each observation $x_i \in \mathbb{R}^p$ becomes $\phi(x_i) \in \mathbb{R}^q$ (typically, $q > p$)

Then PCA for $\phi(x_1), \dots, \phi(x_n)$ is **non-linear** for x_1, \dots, x_n .

Example: $\Phi(f_1, f_2) = (f_1, f_2, f_1^2 + f_2^2)$



Φ \rightarrow



Denote by $\Phi = \begin{pmatrix} -\Phi(x_1)^T \\ \dots \\ -\Phi(x_n)^T \end{pmatrix} \in \mathbb{R}^{n \times q}$ the transformed data

and $K = \Phi\Phi^T \in \mathbb{R}^{n \times n}$ the inner product matrix,
i.e. $K_{ij} = \langle \Phi(x_i), \Phi(x_j) \rangle$.

① Assume that Φ is centered.

We can compute PCs using K only (not Φ)!

• Eigen vectors of $S_\Phi = \frac{1}{N-1} \Phi^T \Phi$ live in the row-space of Φ , i.e. $v = \Phi^T d$ for some $d \in \mathbb{R}^n$.

$$\left| S_\Phi v = \frac{1}{N-1} \Phi^T (\Phi v) = \lambda v \right.$$

• d is an eigenvector of K .

$$\left| \frac{1}{N-1} \Phi\Phi^T \Phi\Phi^T d = \lambda \Phi\Phi^T d \Rightarrow K^2 d = \underbrace{(N-1)}_{\lambda'} \lambda K d \right.$$

Solve $Kd = \lambda' d$

• d should be normalized by the $\sqrt{\lambda}$ of the e value

$$| \|v\|^2 = d^T \Phi \Phi^T d = d^T K d = \lambda' \|d\|^2 = 1 \Rightarrow \|d\| = \frac{1}{\sqrt{\lambda'}}$$

• The PC scores for Φ are just Kd

$$| z = \Phi v = \Phi \Phi^T d = Kd.$$

• To project $\Phi(x)$ onto the PC direction v we need to know d and $\langle \Phi(x_i), \Phi(x) \rangle$ only.

$$| v^T \Phi(x) = d^T \Phi \cdot \Phi(x) = d^T \begin{pmatrix} \langle \Phi(x_i), \Phi(x) \rangle \\ \langle \Phi(x_n), \Phi(x) \rangle \end{pmatrix}$$

② If Φ is not centered then replace K by $\tilde{K} = C K C$ where $C = I - \frac{11^T}{n}$.

$$| \hat{\Phi} = \Phi C \quad \text{so} \quad \tilde{K} = \hat{\Phi} \hat{\Phi}^T = C \Phi \Phi^T C = C K C$$

Kernel PCA

- compute K
- compute the top eigenvalue of K (d')
- find the top eigenvector of K (d) and scale $\|d\| = \frac{1}{\sqrt{d'}}$
- find scores $Z = Kd$.

KPCA relies only on $K(x, y) = \langle \Phi(x), \Phi(y) \rangle$, that is called **kernel function**.

Examples:

Quadratic kernel: $K(x, y) = (1 + \langle x, y \rangle)^2$

$| K(x, y) = 1 + x_1^2 y_1^2 + \dots + x_p^2 y_p^2 + 2x_1 y_1 + \dots + 2x_p y_p = \langle \Phi(x), \Phi(y) \rangle$ for $\Phi(x) = \begin{pmatrix} 1 \\ x_1^2 \\ \vdots \\ x_p^2 \\ \sqrt{2} x_1 \\ \vdots \\ \sqrt{2} x_p \end{pmatrix}$

Polynomial kernel: $K(x, y) = (1 + \langle x, y \rangle)^d$

Radial kernel: $K(x, y) = e^{-\gamma \|x - y\|^2}$

Local dimension reduction methods

tSNE (t-distributed stochastic neighbor embedding)

Given points $x_1, \dots, x_n \in \mathbb{R}^p$

- compute distances $\|x_i - x_j\|^2$
- compute probabilities p_{ij} of selecting neighbors (i, j)
(use Gaussian distribution)

Given embedding $z_1, \dots, z_n \in \mathbb{R}^q$ ($q < p$)

- compute distances $\|z_i - z_j\|^2$
- compute probabilities q_{ij} of selecting neighbors (i, j)
(use t-distribution)

Find z_1, \dots, z_n such that p_{ij} and q_{ij} are "similar"
(use KL divergence)

Main parameter :

- **perplexity**, balances local and global attention

tSNE vs PCA :

+ Non-linear, good for visualization, captures local neighbours

- Slow, struggles with noise, less interpretable :

- no meaning of the tSNE coordinates and distances
- Small distances are informative
- cluster sizes are not informative
- distances between clusters are not informative

UMAP (Uniform Manifold Approximation and Projection)

Given points $x_1, \dots, x_n \in \mathbb{R}^p$

- constructs a weighted k -neighbour graph
(use "fuzzy simplicial complex")

Given embedding $z_1, \dots, z_n \in \mathbb{R}^q$ ($q < p$)

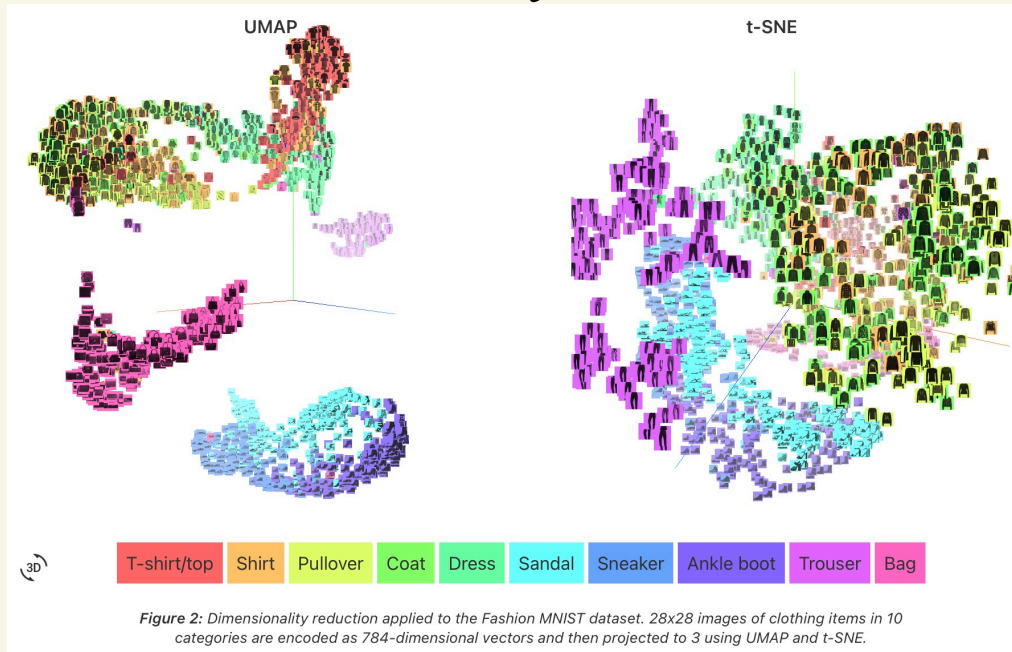
- constructs a weighted k -neighbour graph
(use "fuzzy simplicial complex")

Find z_1, \dots, z_n such that the graphs are "similar"
(use cross-entropy)

Two main parameters:

- **n_neighbors**, the number of nearest neighbours
- **min_dist**, how tightly UMAP packs neighbours

Comparing to t-SNE, UMAP is faster and better at preserving more global structure



Practical aspects:

Sometimes it's better to combine PCA & tSNE/UMAP.

- Filter the data
- Do PCA, reduce dimensionality and noise
- Plot with UMAP/tSNE, try various hyperparameters.

Perturbation theory for PCA

Given $S \in \mathbb{R}^{p \times p}$, consider $\hat{S} = S + E \in \mathbb{R}^{p \times p}$ where $E \in \mathbb{R}^{p \times p}$ is a symmetric noise matrix.

Denote the eigendecompositions by $S = U \Lambda U^T$, $\hat{S} = \hat{U} \hat{\Lambda} \hat{U}^T$.

Recall the definition of the spectral norm:

$$\|A\|_2 = \sqrt{\lambda_1(A^T A)} = d_1(A).$$

Thm (Weyl's) $\max_{i=1..p} |d_i - \hat{d}_i| = \|\Lambda - \hat{\Lambda}\|_2 \leq \|E\|_2$, i.e.

eigenvalues are stable under perturbation.

$$\begin{aligned} \hat{d}_1 &= \max_{\|v\|=1} v^T \hat{S} v = \max_{\|v\|=1} (v^T S v + v^T E v) \leq \\ &\leq d_1 + \max_{\|v\|=1} |v^T E v| = d_1 + \|E\|_2 \end{aligned}$$

$$|\hat{\lambda}_1 - \lambda| \leq \|E\|_2$$

Eigenvectors are not stable!

Example: $S = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ and $E = \begin{pmatrix} 0 & \epsilon \\ \epsilon & 0 \end{pmatrix}$, $\hat{S} = \begin{pmatrix} 1 & \epsilon \\ \epsilon & 1 \end{pmatrix}$

$u_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ and $\lambda_1 = 1$; $u_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ and $\lambda_2 = 1$

$\hat{u}_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ and $\lambda_1 = 1 + \epsilon$; $\hat{u}_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ and $\lambda_2 = 1 - \epsilon$

Measuring agreement between \mathcal{U} and $\hat{\mathcal{U}}$ is tricky.

Recall $U_{(r)} = (u'_1 \dots u'_r) \in \mathbb{R}^{p \times r}$ and $\hat{U}_{(r)} = (\hat{u}'_1 \dots \hat{u}'_r) \in \mathbb{R}^{p \times r}$

• Why not $\|U_{(r)} - \hat{U}_{(r)}\|_F$?

| If $\lambda_1 = \dots = \lambda_p$ then U is any orthogonal matrix.

• We need to measure the agreement between

$\mathcal{U} = \text{span}(u_1, \dots, u_r)$ and $\hat{\mathcal{U}} = \text{span}(\hat{u}_1, \dots, \hat{u}_r)$

Principal angles between subspaces

Consider $A, B \in \mathbb{R}^{n \times r}$ such that $A^T A = B^T B = I$.

Denote $\mathcal{A} = \text{span}(a_1, \dots, a_r)$ $\mathcal{B} = \text{span}(b_1, \dots, b_r)$

Then **principal angle** between \mathcal{A} and \mathcal{B} is

$$\theta_1 = \angle(\mathcal{A}, \mathcal{B}) = \arccos(d_1(A^T B))$$

$$\left| \angle(\mathcal{A}, \mathcal{B}) = \min_{\substack{a \in \mathcal{A} \quad b \in \mathcal{B} \\ \|a\| = \|b\| = 1}} \arccos(a^T b) = \min_{\substack{x, y \\ \|x\| = \|y\| = 1}} \arccos(x^T A^T B y) \right.$$



The general statement is

$$A^T B = U \cos \Theta V^T \quad \text{where } \cos \Theta = \begin{pmatrix} \cos \theta_1 & & \\ & \ddots & \\ & & \cos \theta_r \end{pmatrix}$$

and $\theta_1, \dots, \theta_r$ are called **principal angles**.

Distance Between Subspaces

Define the distance between \mathcal{A} and \mathcal{B} as

$$d(\mathcal{A}, \mathcal{B}) = \|\sin \theta\|_F$$

Let $P_A = AA^T$, $P_B = BB^T$ denote the projection operators and A_\perp and B_\perp are orthogonal complements.

$$\text{Then } d(\mathcal{A}, \mathcal{B}) = \frac{1}{\sqrt{2}} \|P_A - P_B\|_F = \|A^T B_\perp\|_F.$$

$$I = A^T A = A^T (BB^T + B_\perp B_\perp^T) A = U \cos^2 \theta U^T + A^T B_\perp B_\perp^T A$$

$$A^T B_\perp B_\perp^T A = I - U \cos^2 \theta U^T = U (I - \cos^2 \theta) U^T = U \sin^2 \theta U^T$$

$$\text{tr}(U \sin^2 \theta U^T) = \text{tr}(\sin^2 \theta) = \|\sin \theta\|_F^2 = \text{tr}(A^T B_\perp B_\perp^T A) = \|A^T B_\perp\|_F^2$$

$$\|P_A - P_B\|_F^2 = \|AA^T - BB^T\|_F^2 = \text{tr}(AA^T AA^T) - 2 \text{tr}(AA^T BB^T) + \text{tr}(BB^T BB^T)$$

$$= r - 2 \text{tr}(A^T B B^T A) + r = 2r - 2 \text{tr}(I - A^T B_\perp B_\perp^T A) = 2 \text{tr}(A^T B_\perp B_\perp^T A)$$

Davis-Kahan theory

Denote by δ the eigengap, i.e.

$$\delta = \min_{\substack{1 \leq i \leq r, r+1 \leq j \leq p}} |d_i(S) - d_j(\tilde{S})| > 0$$

Then $d(U_{(r)}, \hat{U}_{(r)}) \leq \frac{\|E\|_F}{\delta}$.