

$$1. (1) p(x; \mu) = \sum_z p(x|z; \mu) p(z; \mu)$$

$$p(z=k; \mu) = \pi_k$$

$$p(x_i | \mu_k) = \prod_{j=1}^D \mu_{kj}^{x_{ij}} (1 - \mu_{kj})^{(1-x_{ij})}$$

$$\log p(z_i, x_i | \mu_k) = \sum_k \mathbb{I}_{\{z_i=k\}} \left[\log \pi_k + \sum_{j=1}^D \log (\mu_{kj}^{x_{ij}} (1 - \mu_{kj})^{(1-x_{ij})}) \right]$$

$$= \sum_k \mathbb{I}_{\{z_i=k\}} \left[\log \pi_k + \sum_{j=1}^D (x_{ij} \log \mu_{kj} + (1-x_{ij}) \log (1 - \mu_{kj})) \right]$$

$$q_i(z_i=k) = p(z_i=k | x_i; \mu^{old}) = r_{ik}$$

$$\mu_{new} = \arg \max_{\mu} \sum_i E \left[\sum_k \mathbb{I}_{\{z_i=k\}} (\log \pi_k + \sum_{j=1}^D \log \mu_{kj}^{x_{ij}} (1 - \mu_{kj})^{(1-x_{ij})}) \right]$$

$$= \arg \max_{\mu} \sum_i \sum_k r_{ik} (\log \pi_k + \sum_{j=1}^D (x_{ij} \log \mu_{kj} + (1-x_{ij}) \log (1 - \mu_{kj})))$$

Take derivatives and set it to zero:

$$\frac{\partial}{\partial \mu_{kj}} \sum_i \sum_k r_{ik} (\log \pi_k + \sum_{j=1}^D (x_{ij} \log \mu_{kj} + (1-x_{ij}) \log (1 - \mu_{kj}))) = 0$$

$$\sum_i r_{ik} \left(\frac{x_{ij}}{\mu_{kj}} - \frac{1-x_{ij}}{1-\mu_{kj}} \right) = 0$$

$$(1 - \mu_{kj}) \sum_i r_{ik} x_{ij} = \mu_{kj} \sum_i r_{ik} (1 - x_{ij})$$

$$\mu_{kj} = \frac{\sum_i r_{ik} x_{ij}}{\sum_i r_{ik}}$$

$$(2) \beta(\alpha, \beta, \mu_{kj}) = A \mu_{kj}^{\alpha-1} (1 - \mu_{kj})^{\beta-1} \text{ where } A = \frac{(\alpha+\beta-1)!}{(\alpha-1)!(\beta-1)!}$$

$$\mu_{kj} = \arg \max_{\mu_{kj}} [\mathcal{L}(q; \mu) + \sum_j \sum_k \log \beta(\alpha, \beta, \mu_{kj})]$$

Take derivatives and set it to zero

$$\frac{\partial}{\partial \mu_{kj}} \mathcal{L}(q; \mu) + \frac{\partial}{\partial \mu_{kj}} \sum_j \sum_k \log \beta(\alpha, \beta, \mu_{kj}) = 0$$

$$\sum_i r_{ik} \left(\frac{x_{ij}}{\mu_{kj}} - \frac{1-x_{ij}}{1-\mu_{kj}} \right) + \frac{\partial}{\partial \mu_{kj}} \sum_j \sum_k [\log A + (\alpha-1) \log \mu_{kj} + (\beta-1) \log (1 - \mu_{kj})] = 0$$

$$\sum_i r_{ik} \left(\frac{x_{ij}}{\mu_{kj}} - \frac{1-x_{ij}}{1-\mu_{kj}} \right) + \frac{\alpha-1}{\mu_{kj}} - \frac{\beta-1}{1-\mu_{kj}} = 0$$

$$(1 - \mu_{kj})(\alpha - 1 + \sum_i r_{ik} x_{ij}) = \mu_{kj} (\beta - 1 + \sum_i r_{ik} (1 - x_{ij}))$$

$$\mu_{kj} = \frac{\alpha - 1 + \sum_i r_{ik} x_{ij}}{\alpha + \beta - 2 + \sum_i r_{ik}}$$

2. (1) T_1 :

$\begin{array}{c} \text{Speed distance} \\ \text{Centroid} \end{array}$	A_1	A_2	A_3	A_4	A_5	A_6	A_7	A_8
(2,10)	0 ✓	25	72	13	50	52	65	5
(5,8)	13	18	25 ✓	0 ✓	13 ✓	17 ✓	52	2 ✓
(1,2)	65	10 ✓	53	52	45	29	0 ✓	58

"✓" means optimal of the column.

Cluster1	(2,10)	Center
Cluster2	(8,4) (5,8) (7,5) (6,4) (4,9)	(2,10)
Cluster3	(2,5) (1,2)	(6,6)
		(1.5, 3.5)

T_2 :

$\begin{array}{c} \text{Speed distance} \\ \text{Centroid} \end{array}$	A_1	A_2	A_3	A_4	A_5	A_6	A_7	A_8
(2,10)	0 ✓	25	72	13	50	52	65	5 ✓
(6,6)	32	17	8 ✓	5 ✓	2 ✓	4 ✓	41	13
(1.5, 3.5)	42.5	2.5 ✓	42.5	32.5	32.5	20.5	2.5 ✓	36.5

Cluster1	(2,10) (4,9)	Center
Cluster2	(8,4) (5,8) (7,5) (6,4) (6.5, 5.25)	(3, 9.5)
Cluster3	(2,5) (1,2)	(1.5, 3.5)

T_3 :

$\begin{array}{c} \text{Speed distance} \\ \text{Centroid} \end{array}$	A_1	A_2	A_3	A_4	A_5	A_6	A_7	A_8
(3, 9.5)	1.25 ✓	21.25	55.25	6.25 ✓	36.25	39.25	60.25	1.25 ✓
(6.5, 5.25)	42.81	20.31	3.81 ✓	9.81	0.31 ✓	1.81 ✓	40.81	20.31
(1.5, 3.5)	42.5	2.5 ✓	42.5	32.5	32.5	20.5	2.5 ✓	36.5

Cluster1	(2,10) (5,8) (4,9)	center
Cluster2	(8,4) (7,5) (6,4)	($\frac{11}{3}, 9$)
Cluster3	(2,5) (1,2)	($7, \frac{13}{3}$)
		(1.5, 3.5)

T_4 :

$\begin{array}{c} \text{Speed distance} \\ \text{Centroid} \end{array}$	A_1	A_2	A_3	A_4	A_5	A_6	A_7	A_8
($\frac{11}{3}, 9$)	$\frac{34}{9}$ ✓	$\frac{168}{9}$	$\frac{394}{9}$	$\frac{25}{9}$ ✓	$\frac{224}{9}$	$\frac{274}{9}$	$\frac{505}{9}$	$\frac{1}{9}$ ✓
($7, \frac{13}{3}$)	$\frac{514}{9}$	$\frac{229}{9}$	$\frac{10}{9}$ ✓	$\frac{157}{9}$	$\frac{4}{9}$ ✓	$\frac{10}{9}$ ✓	$\frac{373}{9}$	$\frac{277}{9}$
(1.5, 3.5)	$\frac{85}{2}$	$\frac{5}{2}$ ✓	$\frac{85}{2}$	$\frac{65}{2}$	$\frac{65}{2}$	$\frac{41}{2}$	$\frac{5}{2}$ ✓	$\frac{73}{2}$

Assignments don't change
converge

Cluster1	$A_1 A_4 A_8$
Cluster2	$A_3 A_5 A_6$
Cluster3	$A_2 A_7$

(2) As the figure, A_1 must link A_8 , A_2 must link A_6 , A_4 can't link A_5 .

$\begin{array}{c} \text{Speed distance} \\ \text{Centroid} \end{array}$	A_1	A_2	A_3	A_4	A_5	A_6	A_7	A_8
(2,10)	0 ✓	25	72	13	50	52	65	5 ($\rightarrow A_8$)
(5,8)	13	18	25 ✓	0 ✓	13 ($\rightarrow A_4$)	17	52	2
(1,2)	65	10 ✓	53	52	45 ✓	29 ($\rightarrow A_5$)	0 ✓	58

	Center			
Cluster 1	(2, 10)	(4, 9)	(3, 9.5)	
Cluster 2	(8, 4)	(5, 8)	(6.5, 6)	
Cluster 3	(2, 5)	(7, 5)	(6, 4)	(1, 2)

Control \ data	A ₁	A ₂	A ₃	A ₄	A ₅	A ₆	A ₇	A ₈
C ₃ , (9.5)	1.25 ✓	21.25	55.25	6.25	36.25	39.25	60.25	1.25 (→ A ₁)
(6.5, 6)	36.25	21.25	6.25 ✓	6.25 ✓	125 (→ A ₂)	4.25	46.25	15.25
(4, 4)	40	5 ✓	16	17	10 ✓	4 (→ A ₂)	13 ✓	25

converge

Cluster1	A ₁	A ₃		
Cluster2	A ₃	A ₄		
Cluster3	A ₂	A ₅	A ₆	A ₇

3. Given a dataset $D = \{x^{(1)}, \dots, x^{(N)}\} \subset \mathbb{R}^D$

Let the mean be $\mu := \frac{1}{N} \sum_{i=1}^N x^{(i)}$

K-dimensional subspace S is spanned by an orthonormal basis $\{u_k\}_{k=1}^K$ where $u_k \in \mathbb{R}^D$, $u_i^T u_j = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}$

Approximate data x : $\tilde{x} = \mu + \text{Proj}_S(x - \mu) = \mu + \sum_{k=1}^K z_k u_k$ where $z_k = u_k^T (x - \mu)$

Let $U \in \mathbb{R}^{D \times K}$ be a matrix with columns $\{u_k\}_{k=1}^K$

reconstruction of x : $\tilde{x} = \mu + U z$

representation of x : $z = U^T (x - \mu)$

the mean of the reconstructions $\tilde{\mu} = \frac{1}{N} \sum_{i=1}^N \tilde{x}^{(i)}$

we need to $\max_{U \in \mathbb{R}^{D \times K}} \frac{1}{N} \sum_{i=1}^N \|\tilde{x}^{(i)} - \tilde{\mu}\|^2$

$$\begin{aligned} \tilde{\mu} &= \frac{1}{N} \sum_{i=1}^N \tilde{x}^{(i)} \\ &= \frac{1}{N} \sum_{i=1}^N (\mu + U z^{(i)}) \\ &= \frac{1}{N} \sum_{i=1}^N (\mu + U U^T (x^{(i)} - \mu)) \\ &= \frac{1}{N} \sum_{i=1}^N (\mu + x^{(i)} - \mu) \\ &= \frac{1}{N} \sum_{i=1}^N x^{(i)} \\ &= \mu \end{aligned}$$

$$\therefore \tilde{x}^{(i)} - \tilde{\mu} = (\mu + U z^{(i)}) - \mu = U z^{(i)}$$

$$\|U z^{(i)}\|^2 = (U z^{(i)})^T (U z^{(i)}) = z^{(i)T} U^T U z^{(i)} = z^{(i)T} z^{(i)} = \|z^{(i)}\|^2$$

$$\max_{U \in \mathbb{R}^{D \times K}} \frac{1}{N} \sum_{i=1}^N \|\tilde{x}^{(i)} - \tilde{\mu}\|^2 = \max_{U \in \mathbb{R}^{D \times K}} \frac{1}{N} \sum_{i=1}^N \|U z^{(i)}\|^2 = \max_{U \in \mathbb{R}^{D \times K}} \frac{1}{N} \sum_{i=1}^N \|z^{(i)}\|^2$$

$$\therefore z^{(i)} = U^T (x^{(i)} - \mu)$$

$$\begin{aligned}\max_{\substack{U \\ U^T U = I}} \frac{1}{N} \sum_{i=1}^N \|z^{(i)}\|^2 &= \max_{\substack{U \\ U^T U = I}} \frac{1}{N} \sum_{i=1}^N \|U^T (x^{(i)} - \mu)\|^2 \\ &= \max_{\substack{U \\ U^T U = I}} \frac{1}{N} \sum_{i=1}^N \text{Trace}(U^T (x^{(i)} - \mu)(x^{(i)} - \mu)^T U)\end{aligned}$$

the empirical covariance matrix $\Sigma = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \mu)(x^{(i)} - \mu)^T$

$$\max_{\substack{U \\ U^T U = I}} \frac{1}{N} \sum_{i=1}^N \text{Trace}(U^T (x^{(i)} - \mu)(x^{(i)} - \mu)^T U) = \max_{\substack{U \\ U^T U = I}} \text{Trace}(U^T \Sigma U) = \max_{\substack{U \\ U^T U = I}} \sum_{k=1}^K u_k^T \Sigma u_k$$

The Lagrangian function $\mathcal{L}(U, \Lambda_K) = \text{Trace}(U^T \Sigma U) + \text{Trace}(\Lambda_K^T (I - U^T U))$ where $\Lambda_K = \text{diag}([\hat{\lambda}_1, \dots, \hat{\lambda}_K])$

$$\text{Let } \frac{\partial \mathcal{L}(U, \Lambda_K)}{\partial U} = 2\Sigma U - 2U\Lambda_K = 0$$

$$\Sigma u_k = \hat{\lambda}_k u_k, \quad k=1, \dots, K$$

u_k is eigenvector, $\hat{\lambda}_k$ is eigenvalue

utilizing SVD decomposition, $\Sigma = Q\Lambda_D Q^T = \sum_{i=1}^D \lambda_i q_i q_i^T$ where $Q = [q_1, \dots, q_D] \in \mathbb{R}^{D \times D}$, q_i is the eigenvector corresponding to i th largest eigenvalue λ_i , $\Lambda_D = \text{diag}([\lambda_1, \dots, \lambda_D])$ with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D$

$$\sum_{k=1}^K u_k^T \Sigma u_k = \sum_{k=1}^K \sum_{i=1}^D \lambda_i (u_k^T q_i)(q_i^T u_k) = \sum_{i \in T} \lambda_i$$

$$T \subset \{1, 2, \dots, D\}, |T| = K$$

we should pick top- K eigenvalues, correspondingly, the first K columns of Q are the optimal solution to U

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i = \begin{bmatrix} 6.9 \\ 3.5 \\ 5.1 \end{bmatrix}$$

$$\Sigma = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T = \begin{bmatrix} 2.09 & 1.45 & -0.39 \\ 1.45 & 2.25 & -1.15 \\ -0.39 & -1.15 & 7.09 \end{bmatrix}$$

$$\text{SVD decomposition } \Sigma = Q\Lambda_D Q^T = \begin{bmatrix} -0.1376 & 0.6990 & 0.7017 \\ -0.2505 & 0.6609 & -0.7075 \\ 0.9583 & 0.2731 & -0.0842 \end{bmatrix} \begin{bmatrix} 7.4465 & 0 & 0 \\ 0 & 3.3085 & 0 \\ 0 & 0 & 0.6749 \end{bmatrix} \begin{bmatrix} -0.1376 & -0.2505 & 0.9583 \\ 0.6990 & 0.6609 & 0.2731 \\ 0.7017 & -0.7075 & -0.0842 \end{bmatrix}$$

$$\text{pick top 2 eigenvalues, } U = [q_1, q_2] = \begin{bmatrix} -0.1376 & 0.6990 \\ -0.2505 & 0.6609 \\ 0.9583 & 0.2731 \end{bmatrix}$$

The new representation $z = U^T (x_i - \mu)$

$$z_1 = \begin{bmatrix} -2.1514 \\ -0.1731 \end{bmatrix} \quad z_2 = \begin{bmatrix} 3.8042 \\ -2.8875 \end{bmatrix} \quad z_3 = \begin{bmatrix} 0.1522 \\ -0.9869 \end{bmatrix} \quad z_4 = \begin{bmatrix} -4.7065 \\ 1.3015 \end{bmatrix} \quad z_5 = \begin{bmatrix} 1.2938 \\ 2.2791 \end{bmatrix}$$

$$z_6 = \begin{bmatrix} 4.0993 \\ -0.1436 \end{bmatrix} \quad z_7 = \begin{bmatrix} -1.6258 \\ -2.2321 \end{bmatrix} \quad z_8 = \begin{bmatrix} 2.1145 \\ 3.2512 \end{bmatrix} \quad z_9 = \begin{bmatrix} -0.2348 \\ 0.3730 \end{bmatrix} \quad z_{10} = \begin{bmatrix} -2.7464 \\ -1.0689 \end{bmatrix}$$

Programming Problems

1. PCA

- (1) partition the dataset

```
X = dataset[:, :-1] is the features
```

```
y = dataset[:, -1] is the labels
```

- (2) calculate the mean, covariance, SVD decomposition

```
mu = np.mean(X, axis=0) is the mean of dataset.
```

```
sigma = np.cov(X.T) is the covariance matrix
```

```
u, s, v = np.linalg.svd(sigma) this is the process of SVD decomposition
```

```
U = u[:, :2] is the top 2 eigenvectors
```

- (3) calculate the new representation

```
z = U.T @ (X - mu).T is the new representation of dataset X
```

2. K-means

- (1) randomly pick 3 points as centers

```
initialCenter = np.random.choice(dataset.shape[0], size=3, replace=False)
```

```
x = z.T
```

```
center = (x[initialCenter[0]], x[initialCenter[1]], x[initialCenter[2]])
```

- (2) repeat “assignment” and “refitting” until convergence

```
while True:
```

```
    center = refitting(x, labelList)
```

```
    if all(assignment(x, center) == labelList):
```

```
        break # converge
```

```
    labelList = assignment(x, center) every element corresponds to a label
```

3. Silhouette Coefficient

- (1) For each data point i , let a_i be the mean distance between i and all other points in the same cluster, and let b_i be the mean distance between i and all other points in the next nearest cluster. The Silhouette Coefficient for data point i is defined as $s_i = (b_i - a_i) / \max(a_i, b_i)$. The overall silhouette coefficient is the average of s_i overall data points.

(2) silhouette coefficient: 0.4802142699427175

4. Rand Index

- (1) Given two clusterings of these data points, let “a” be the number of pairs of elements that are in the same cluster in both clusterings, and let “b” be the number of pairs of elements that are in different clusters in both clusterings. The Rand Index is defined as $RI = 2(a+b)/(n(n-1))$ where n is the number of data.

(2) rand index: 0.8743677375256322