Scan the QR to record your attendance
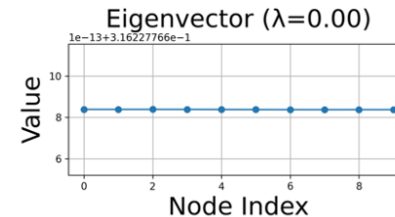
# Quick review – An illustration of EVD on L
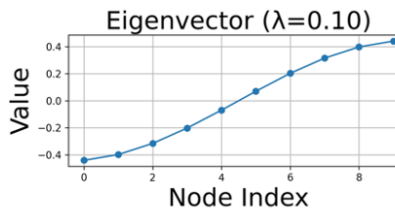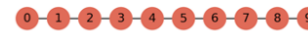


A line graph with 10 nodes

A ➔ L=D-A ➔ EVD

➔ $\lambda \in R^{10}$ & $10 * u \in R^{10}$

# Quick review -- Introduction to Graph Science

- Unnormalized / Normalized / Random Walk Graph Laplacian
- EVD on L
  - $u_k$ : Fourier functions, i.e. vibration modes of the graph.
  - $\lambda_k$ : frequencies of the Fourier functions, i.e. how fast $u_k$ vibrate.


- Distance
  - Euclidean/L2, Cosine/dot product, KL divergence, Wasserstein distance…
- Pre-processing
  - Center, Z-scoring, Project on $L_2$ sphere, Min-Max normalization

$$X =$$

# Quick review -- Graph Clustering

- k-means: partition the dataset into k clusters, minimize the least-squares
  - $L = \frac{1}{kn} \sum_{q=1}^{k} \sum_{i \in S_q} \left\| x_i - m_q \right\|^2$

- EM for k-means: cluster update (expectation) $\leftarrow\rightarrow$ mean update (maximization)
  - Monotonic loss, Guaranteed convergence, Speed complexity O(n·d·k· $n_i$)
  - NP-hard, Good initialization

- Non-liner k-means: $x \rightarrow \phi(x)$, higher dim, more separable, $L = \frac{1}{kn} \sum_{q=1}^{k} \sum_{i \in S_q} \theta_i \left\| \phi(x_i) - m_q \right\|^2$

- Indicator function: $F \in \{0,1\}^{n \times k} \rightarrow L = \frac{1}{kn} tr(F^\top \Theta D), \Theta = diag(\theta_1, \dots, \theta_n) \in R^{n \times n}, D \in R^{n \times k}$ distance
  - $D = diag(K)1_K^\top - 2K\Theta FZ + 1_n diag(ZF^\top \Theta K \Theta FZ) \in R^{n \times k}$
  - $K = \phi(x)\phi(x)^\top \in R^{n \times n}$ (mapping, similarity) $\rightarrow$ direct similarity $\rightarrow K(x_i, x_j)$ linear, Gaussian, Polynomial...
  - Now, EM becomes
    - $F^0 \rightarrow D^0 \rightarrow F^1 \rightarrow D^1 \rightarrow \cdots$
    - No $m_q$ in $D$ $\rightarrow$ No maximization !!!

# Quick review -- Graph Clustering

$$D = diag(K)1_K^\top - \mathbf{2K\Theta FZ} + \mathbf{1_n diag(ZF^\top \Theta K \Theta FZ})$$

$$\frac{\phi(x_i)^\top \sum_{j \in S_q} \theta_j \phi(x_j)}{\sum_{j \in S_q} \theta_j}$$

$$Z^{-1} = diag(1_n^T \Theta F)$$

$$\left( \frac{\sum_{j \in S_q} \theta_j \phi(x_j)}{\sum_{j \in S_q} \theta_j} \right) \left( \frac{\sum_{j \in S_q} \theta_j \phi(x_j)}{\sum_{j \in S_q} \theta_j} \right)^\top$$

$$D_{iq} = \left\| \phi(x_i) - m_q \right\|^2 = K_{ii} - \mathbf{2A_{iq}} + \mathbf{B_{qq}}$$

$$\phi(x_i)^\top m_q \qquad m_q m_q^\top$$

# Quick review -- Graph Clustering

- Non-liner k-means: $L = \frac{1}{kn} \sum_{q=1}^{k} \sum_{i \in S_q} \theta_i \left|\left| \phi(x_i) - m_q \right|\right|^2 \rightarrow m_q \in R^{k \times d}$  EM means

- Spectral mean $m_i \in R^{n \times d}$ : cluster-centric $\rightarrow$ sample-centric

  - $L = \frac{1}{kn} \sum_{q=1}^{k} \sum_{i \in S_q} \theta_i \left|\left| \phi(x_i) - m_q \right|\right|^2 = \frac{1}{kn} \sum_{q=1}^{k} \sum_{i \in S_q} \left|\left| \theta_i^{1/2} \phi(x_i) - \theta_i^{1/2} m_q \right|\right|^2 = \frac{1}{kn} \left|\left| \theta^{\frac{1}{2}} \phi - \theta^{\frac{1}{2}} M \right|\right|^2, M = FZF^{\top} \Theta \phi$

  - Indicator $Y = \Theta^{1/2} F Z^{1/2} \rightarrow L = \frac{1}{kn} \left|\left| \theta^{\frac{1}{2}} \phi - YY^{\top} \theta^{\frac{1}{2}} \phi \right|\right|^2 \rightarrow \max_{Y \in R^{n \times k}} tr(Y^{\top} \Theta^{1/2} K \Theta^{1/2} Y) \rightarrow \max_{Y \in R^{n \times k}} tr(Y^{\top} A Y)$

  - A is symmetric & positive semi-definite $\rightarrow$ Y are the k largest eigenvectors, $\max_{Y \in R^{n \times k}} tr(Y^{\top} A Y) = \sum_{q=1}^{k} \lambda_q$

- Spectral clustering

  - A=$\Theta^{1/2} K \Theta^{1/2}$ $\rightarrow$ EVD $\rightarrow$ pick the largest k eigenvectors as $Y^*$ $\rightarrow$ binarize $Y^*$ (standard k-means on $Y^*$ )

  - Independent of initialization, O(n^2k)

# Core API

- Plot scatter
  - plt.scatter(x, y, c, s), c is label / value …
- Compute pairwise distance
  - sklearn.metrics.pairwise.pairwise_distances(X, metric='euclidean')
- Numpy
  - sort(a, axis=-1), argsort(a, axis=-1), ascending order