

PCA

Suppose we have n data samples
 $x_1, \dots, x_n \in \mathbb{R}^p$. We denote

1) the sample mean $\mu^* = \frac{1}{n} \sum_{k=1}^n x_k, \mu^* \in \mathbb{R}^p$

$\underbrace{\phantom{\sum_{k=1}^n x_k}}$
mean done element-wise,
i.e., $\mu_i^* = \frac{1}{n} \sum_{k=1}^n x_{ki}$

2) the sample covariance

$$C^* = \frac{1}{n-1} \sum_{k=1}^n (x_k - \mu^*)(x_k - \mu^*)^T$$

$\underbrace{\phantom{\sum_{k=1}^n (x_k - \mu^*)(x_k - \mu^*)^T}}$

this is an outer product!

↳ for $a, b \in \mathbb{R}^p$, $(a \cdot b^T)_{ij} = a_i b_j$

$$a \cdot b^T = \begin{pmatrix} a_1 b_1 & a_1 b_2 & \dots & a_1 b_p \\ \vdots & \ddots & & \vdots \\ a_p b_1 & \dots & \ddots & a_p b_p \end{pmatrix}$$

First approach

1) We want to find a lower-dim. representation of our data samples.

Idea: Project the samples into a d-dim.

Example: $20 \rightarrow 10$

Subspace with basis vectors $v_i \in \mathbb{R}^p$,
 $i \in [1, d]$ with $d < p$. For convenience, we choose an orthonormal basis $\rightarrow v_i^T v_j = \delta_{ij}$,

If V has the vectors as columns,
 $V = (v_1, \dots, v_d)$ then $V^T V = \mathbf{I}$.

(note that $V V^T \neq \mathbf{I}$!!)

2) A good representation has to be close to the original

$\mu_k + V \beta_k \stackrel{!}{\approx} x_k$

offset

basis vectors

vector components

reconstruct the original sample

β_k is our low-dim. representation of x_k !

We phrase this as an optimisation problem:

$$\min_{\mu, V, \beta_k, \quad \nabla V = 1} \sum_{k=1}^n \|x_k - (\mu + V \beta_k)\|_Q^2$$

Euclidean norm

Find μ : W.l.o.g. we assume $\sum_{k=1}^n \beta_k = 0$.

To find μ , we set the gradient to 0:

$$\nabla_\mu \sum_{k=1}^n \|x_k - (\mu + V \beta_k)\|^2 \stackrel{!}{=} 0$$

$$\Leftrightarrow \sum_{k=1}^n x_k - (\mu + V \beta_k) = 0$$

$$\Leftrightarrow \mu = \frac{1}{n} \sum_{k=1}^n x_k \quad \text{since} \quad \sum_{k=1}^n \beta_k = 0$$

$$= \mu^* \quad (\text{for simplicity, we write } \mu)$$

Find β_k : Same approach! $\nabla_{\beta_k} \sum_{k=1}^n \|x_k - (\mu + V \beta_k)\|^2$

note: $\|\cdot\|^2$ terms \Rightarrow decoupler in k !

Problem: $\nabla_{\beta_k} \|\underbrace{V \beta_k}_{\text{matrix-vector product}}\|^2 = ?$

For simplicity, let's write this with indices:

$$\begin{aligned}\frac{\partial}{\partial \beta_{k,m}} \sum_i \left(\sum_j V_{ij} \beta_{kj} \right)^2 &= 2 \sum_i \left(\sum_j V_{ij} \beta_{kj} \right) V_{im} \\ &= 2 \sum_{ij} V_{mi}^T V_{ij} \beta_{kj} \\ \Rightarrow \nabla_{\beta_k} \|\underbrace{V \beta_k}_{\text{matrix-vector product}}\|^2 &= V^T V \beta_k\end{aligned}$$

Lets apply this: $\nabla_{\beta_k} \|x_k - (\mu + V \beta_k)\|^2 = 0$

$$\Leftrightarrow V^T (x_k - (\mu + V \beta_k)) = 0$$

$$\Leftrightarrow V^T (x_k - \mu) = \underbrace{V^T V}_{=1} \beta_k = \beta_k$$

Quickly check our assumpt.: $\sum_{k=1}^n \beta_k = \sum_{k=1}^n V^T (x_k - \mu) = 0$

Thus, β_k is obtained by projecting x_k to the low.-dim. subspace!

The last step is to find the basis vectors!

↪ insert β_k : $\min_{\substack{V \\ V^T V = \mathbb{1}}} \|x_k - \underbrace{\tilde{V} V^T}_{\neq \mathbb{1}}(x_k - \mu)\|^2$

Write this out: $(x_k - \mu - \underbrace{\tilde{V} V^T}_{=1}(x_k - \mu))^2$

$$= (x_k - \mu)^2 + (x_k - \mu)^T \underbrace{\tilde{V} V^T}_{\neq \mathbb{1}} \tilde{V} V^T (x_k - \mu)$$
$$- 2(x_k - \mu)^T \tilde{V} V^T (x_k - \mu)$$
$$= (x_k - \mu)^2 - (x_k - \mu)^T \tilde{V} V^T (x_k - \mu)$$

Note: $(x_k - \mu)^2$ does not depend on V

$$\Rightarrow \min \text{ turns to } \max_{\substack{V \\ V^T V = \mathbb{1}}} \sum_{k=1}^n (x_k - \mu)^T \tilde{V} V^T (x_k - \mu)$$

\hookrightarrow
due to the minus sign,
 \min becomes \max

We can rewrite this to arrive at a nice expression.

Substitute $a_k = (x_k - \mu)$ for simplicity.

$$\begin{aligned} \rightarrow \sum_k a_k^T V V^T a_k &= \sum_{k,n} \sum_i \sum_j a_{kn} V_{ni} V_{ij}^T a_{kj} \\ &= \sum_{k,j \in n} V_{in}^T a_{kn} a_{kj} V_{ji}^T \\ &= \sum_{i,j \in n} V_{in}^T \underbrace{\sum_k a_{kn} a_{kj}}_{= [a_k a_k^T]_{n \times 1}} V_{ji}^T \\ &= \sum_i [V^T (\sum_k a_k a_k^T) V]_{ii} \\ &= \text{Tr}(V^T (\sum_k a_k a_k^T) V) \end{aligned}$$

Substitute back a_k :

$$\begin{aligned} \text{Tr}(V^T [\sum_k (x_k - \mu)(x_k - \mu)^T] V) \\ = \underbrace{(n-1)}_{\substack{\uparrow \\ \text{Covariance matrix!}}} C^* \end{aligned}$$

Thus, our problem becomes

$$\max_{V,} \operatorname{Tr}(V^T C^* V)$$

s.t. $V^T V = 1$

We can do some more rewriting:

$$\begin{aligned} \operatorname{Tr}(V^T C^* V) &= \sum_{i, m, j} V_{im}^T C_{mj}^* V_{ji} = \sum_{i, m, j} V_{mi} C_{mj}^* V_{ji} \\ &= \sum_i V_i^T C^* V_i \end{aligned}$$

↑ ↑
basis vectors

product sums over column entries!

C^* is a symmetric and real-valued matrix.

Thus, it has an eigenvalue decomp.

$$\text{In the exercises, we showed that } \max_{V,} \operatorname{Tr}(V^T C^* V) \quad (\lambda=1)$$

s.t. $\|V\|^2=1$

is solved by the eigenvector of C^* with the largest eigenvalue. For $\lambda > 1$, V is given by the eigenvectors of C^* with the λ largest eigenvalues. \square

Second approach

Idea: Lets find a projection of the data with maximum variance.

The projection is given by $V \in \mathbb{R}^{n \times d}$, $V^T V = I$, and the projected points are $y_k = V^T x_k$

Formulate problem:

Sample mean: $\bar{\mu} = \frac{1}{n} \sum_{k=1}^n y_k = \frac{1}{n} \sum_{k=1}^n V^T x_k = V^T \mu^*$

Sample variance: $\frac{1}{n-1} \sum_{k=1}^n \|y_k - \bar{\mu}\|^2$
 $= \frac{1}{n-1} \sum_{k=1}^n \|V^T x_k - V^T \mu^*\|^2$
 $= \frac{1}{n-1} \sum_{k=1}^n \|V^T(x_k - \mu^*)\|^2$

Maximize the sample variance:

$$\max_{V^T V = I} \sum_{k=1}^n \|V^T(x_k - \mu)\|^2$$

→ same objective as in

same solution! ← the prev. approach!



Notes on Marčenko-Pastur

$x_k = g + \sqrt{\beta} g_0 u \rightarrow x_k$ is a sum of scaled Gaussian, and hence also a Gaussian RV.

The moments are:

$$1) E(x_k) = \underbrace{E(g)}_{=0} + \sqrt{\beta} u \underbrace{E(g_0)}_{=0} = 0 \Rightarrow \mu = 0$$

$$\begin{aligned} 2) \Sigma &= E((x_k - \mu)(x_k - \mu)^T) \\ &\stackrel{*}{=} E(x_k x_k^T) = \underbrace{E(g g^T)}_{=1} + \underbrace{E(g g_0)}_{=0} \sqrt{\beta} u u^T \\ &\quad + \underbrace{E(g_0 g^T)}_{=0} \sqrt{\beta} u + E(g_0^2) \underbrace{\beta u u^T}_{=1} \\ &= 1 + \beta u u^T \end{aligned}$$

previously

Johnson-Lindenstrauss lemma

Given: $x_1, \dots, x_m \in \mathbb{R}^n$ are our data samples

$y_1, \dots, y_m \in \mathbb{R}^d$. $y_i = \frac{1}{\sqrt{d}} G x_i$ is the proj. data
 $G \in \mathbb{R}^{d \times n}$, with $G_{ij} \sim N(0, 1)$

y_i is a rescaled unit Gaussian random variable

$\hookrightarrow y_{ik} = \frac{1}{\sqrt{d}} \sum_j G_{kj} x_{ij} = \text{sum of Gaussian RVs } (G_{kj}) \text{ scaled by } \frac{1}{\sqrt{d}} x_{ij}!$

$\Rightarrow y_{ik}$ is a Gaussian RV!

Lets calculate the moments:

$$E(y_i) = \frac{1}{\sqrt{d}} E(G) x_i = 0 \text{ since } E(G)=0$$

$$\begin{aligned} \text{Var}(y_{ik}) &= \text{Var}\left(\frac{1}{\sqrt{d}} \sum_j G_{kj} x_{ij}\right) \\ &= \frac{1}{d} \sum_j x_{ij}^2 \underbrace{\text{Var}(G_{kj})}_{=1} = \frac{1}{d} \|x_i\|^2 \end{aligned}$$

Thus, we can represent y_i as a rescaled unit Gaussian $g_i \in \mathbb{R}^d$, $E(g_i)=0$ and $\text{Var}(g_i)=1$:

$$y_i = \frac{\|x_i\|}{\sqrt{d}} g_i$$

↑
1-vector

Proof of one side

We will only do the proof for one case, as the technique is the same for proving the whole theorem.

Moreover, we will prove

$$P(\|y_i\|^2 \geq (1+\epsilon)^2 \|x_i\|^2) \leq \frac{1}{m^3},$$

which is equivalent to what we are looking for:

$$P(\|y_i\| \geq (1+\epsilon) \|x_i\|) \leq \frac{1}{m^3}$$

? the prob. that our projection deviates too much!

We first write out $\|y_i\|^2 = \frac{\|x_i\|^2}{d} \sum_j g_{ij}^2$.

Thus, we want to bound

$$P(\|y_i\|^2 \geq (1+\epsilon)^2 \|x_i\|^2) = P\left(\sum_j g_{ij}^2 \geq (1+\epsilon)^2 d\right),$$

which can be done via the Chernoff bound!

$$P\left(\sum_j g_{ij}^2 \geq (1+\epsilon)^2 d\right) \leq e^{-t(1+\epsilon)^2 d} E(e^{t \sum_j g_{ij}^2})$$

↑ moment-gen.
function!

The probability
we would like
to bound

Chernoff

We need to calculate the moment-gen. fn.

First, we have $E(e^{t \sum g_{ij}^2}) = \prod_{i,j} E(e^{t g_{ij}^2})$
↑
independence

Then, we expand the expectation value: $\mathcal{N}(0,1)$

$$\begin{aligned} E(e^{t g_{ij}^2}) &= \frac{1}{\sqrt{2\pi}} \int e^{tz^2} \cdot e^{-\frac{z^2}{2}} dz \\ &\stackrel{\text{rewrite as } z \text{ in integral}}{=} \frac{1}{\sqrt{2\pi}} \int e^{-\frac{z^2}{2} + t z^2} dz \\ &= \frac{1}{\sqrt{2\pi}} \int e^{-\frac{z^2}{2} \cdot (1-2t)} dz \\ &= \frac{1}{\sqrt{1-2t}} \cdot \underbrace{\frac{1}{\sqrt{2\pi}} \int e^{-\frac{z^2}{2}} dz'}_{=1} \\ &= \frac{1}{\sqrt{1-2t}} \end{aligned}$$

Substitute with
 $z = z \sqrt{1-2t}$
 $dz' = dz \sqrt{1-2t}$

Thus, we have: $E(e^{-t \sum g_{ij}^2}) = (1-2t)^{-\frac{d}{2}}$

and thus $P(|Y_i| \geq (1+\epsilon)|X_i|) \leq e^{-t(1+\epsilon)^2} \cdot (1-2t)^{-\frac{d}{2}}$

We are not yet finished, as we have to choose t ! t is arbitrary, so after some meditation (joking), you have to work backwards from the result to find this: $t = \frac{1}{2} \left(\frac{(1+\epsilon)^2 - 1}{(1+\epsilon)^2} \right)$

$$\begin{aligned}
 &\Rightarrow P(\|y_i\| \geq (1+\epsilon)\|x_i\|) \\
 &\leq e^{-\frac{\alpha}{2}[(1+\epsilon)^2 - 1]} \cdot \left(\frac{1}{(1+\epsilon)^2} \right)^{-\frac{\alpha}{2}} \\
 &= e^{-\frac{\alpha}{2}[(1+\epsilon)^2 - 1] - \frac{\alpha}{2} \ln\left(\frac{1}{(1+\epsilon)^2}\right)} \\
 &\stackrel{\text{use } \ln(1+x) \leq x - \frac{x^2}{4}}{=} e^{-\alpha \left[\epsilon + \frac{\epsilon^2}{2} - \ln(1+\epsilon) \right]} \\
 &\leq e^{-\alpha \left[\epsilon + \frac{\epsilon^2}{2} - \epsilon + \frac{\epsilon^2}{4} \right]} \\
 &\text{for } x \in [0, 1] \quad = e^{-\alpha \frac{3\epsilon^2}{4}} \\
 &\text{(just Taylor expansion)}
 \end{aligned}$$

Now set $\alpha = 4 \ln(m)/\epsilon^2$

$$\Rightarrow P(\|y_i\| \geq (1+\epsilon)\|x_i\|) \leq \frac{1}{m^3}$$

Proof for the left tail ($1-\varepsilon$) $\|x_i\|$ is similar!

Note that our result is valid for any vector

$v \in \mathbb{R}^n$. Hence, it also holds for, e.g.,

$$v = x_i - x_j$$

$$\hookrightarrow P(L(v) \geq (1+\varepsilon)\|v\|)$$

$$= P(\|y_i - y_j\| \geq (1+\varepsilon)\|x_i - x_j\|) \leq \frac{1}{m^3}$$

since $L(v) = \frac{1}{f_d} G(x_i - x_j) = \frac{1}{f_d} 6x_i - \frac{1}{f_d} 6x_j = y_i - y_j$

The last step

The lemma states that lengths and distances do not deviate too much for all data points in both directions (being too large or too small).

We have

$$\begin{aligned} P_i &= P(\|y_i\| \geq (1+\varepsilon)\|x_i\| \text{ or } \|y_i\| \leq (1-\varepsilon)\|x_i\|) \\ &\leq \frac{2}{m^3} \end{aligned}$$

$$\begin{aligned} \text{and } P_{ij} &= P(\|y_i - y_j\| \geq (1+\varepsilon)\|x_i - x_j\| \text{ or } \|y_i - y_j\| \leq (1-\varepsilon)\|x_i - x_j\|) \\ &\leq \frac{2}{m^3} \end{aligned}$$

What is the probability that at least one of these inequalities is true?

↳ remember Kolmogorov! For event A, B,

$$\begin{aligned} P(A \cup B) &= P(A) + P(B) - P(A \cap B) \\ &\leq P(A) + P(B) \end{aligned}$$

$$\Rightarrow P(\bigvee_i A_i) \leq \sum_i P(A_i)$$

\Rightarrow P(at least one deviation)

$$\leq \sum_i P_i + \sum_{i \neq j} P_{ij} \leq \frac{2}{m}$$

$\leq \frac{2}{m^3}$ $\leq \frac{2}{m^3}$ \uparrow
in total
we have
 m^2 vectors

\Rightarrow P(no deviation) = $1 - P(\text{at least one deviation})$

$$\geq 1 - \frac{2}{m}$$

II

Spectral clustering

Laplacian

We define it as $L = D - W$, $L \in \mathbb{R}^{n \times n}$

with W_{ij} = "weight" between node i and j

and $D_{ii} = \sum_j W_{ij}$, $D_{ij} = 0$ if $i \neq j$
with $W_{ij} \geq 0$

Properties of L :

1) For every $f \in \mathbb{R}^n$, we have $f^T L f = \frac{1}{2} \sum_i^n W_{ii} (f_i - f_j)^2$

Proof:
$$\begin{aligned} f^T L f &= f^T D f - f^T W f \\ &= \sum_i D_{ii} f_i^2 - \sum_{ij} f_i W_{ij} f_j \\ &= \frac{1}{2} \left(\sum_i D_{ii} f_i^2 - 2 \sum_{ij} f_i f_j W_{ij} + \sum_j D_{jj} f_j^2 \right) \end{aligned}$$

$$= \frac{1}{2} \sum_{ij} W_{ij} (f_i^2 - 2f_i f_j + f_j^2) = \frac{1}{2} \sum_{ij} W_{ij} (f_i - f_j)^2$$

2) L is symmetric and pos. def.

\hookrightarrow We have $D^T = D$ and $W^T = W \Rightarrow L^T = L$

\hookrightarrow Pos. def. follows from 1)

- 3) Smallest eigenvalue is 0 with the one-vector
 $\mathbf{1}$ as eigenvector.
 $\text{Lo } (1, 1, 1, \dots, 1)$

Lets look at the product $(L\mathbf{1})_i$:

$$(L\mathbf{1})_i = D_{ii} - \sum_j w_{ij} = \sum w_{ij} - \sum w_{ij} = 0 \quad \text{II}$$

- 4) L has non-negative, real-valued eigenvalues
 Lo follows from symmetry (real-valued)
 and pos def. (pos. eigenvalues)

Connected components of G = multiplicity of 0-EV of L

- 1) Lets assume a graph with one connected component. Let f be the eigenvector with eigenvalue 0.

We know that $f^T L f = \frac{1}{2} \sum_{ij} w_{ij} (f_i - f_j)^2 = 0$
 since $Lf = 0$

Since $w_{ij} \geq 0$, for all $w_{ij} \neq 0$ we have to demand $f_i = f_j$.

$$\begin{array}{c} v_{ij} > 0 \\ \text{---} \\ \text{---} \end{array} \Rightarrow f_i = f_j = c \quad | \quad \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \Rightarrow f_i = f_j = f_k = c$$

Basically: for all nodes i, k connected by a path,
we have $f_i = f_k$.

Since G is connected, we find $f = \underline{1} \cdot c$

2) Now assume we have k connected components.

A_k contains all node indices of the k^{th} comp.

$\hookrightarrow V = \bigcup_k A_k$ and $A_i \cap A_j = \emptyset$ if $i \neq j$

We can then decompose the sum:

$$f^T L f = \frac{1}{2} \sum_k \sum_{i,j \in A_k} w_{ij} (f_i - f_j)^2 = 0$$

For $k=0$, this is solved by setting $f_i = \begin{cases} 1 & \text{if } i \in A_0 \\ 0 & \text{else} \end{cases}$

For $k=1, \dots, n$ $f_i = \begin{cases} 1 & \text{if } i \in A_k \\ 0 & \text{else} \end{cases}$

Let's call the found EV for k : f^k

From the def. follows: $(f^L)^T f^m = 0$ if $L \neq m$

\Rightarrow We have k orthogonal eigenvectors
with eigenvalue $0!$

□

Graph cut

Lets assume we have again a graph
 $G = (V, E)$ with edge weights $W_{ij} \leftarrow$ measure of similarity

For clustering, we want to separate
the nodes into k sets A_k such that

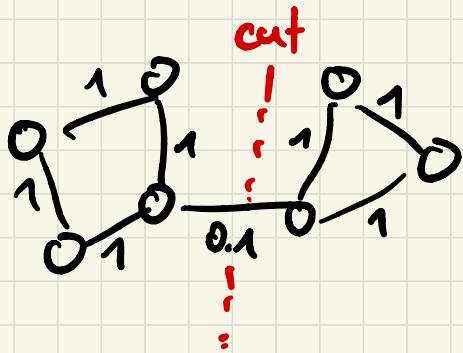
- 1) Edges between nodes from different sets
have low weight

Lets denote by $W(A, B) = \sum_{\substack{i \in A \\ j \in B}} W_{ij}$

We can formalize this as a mincut problem:

The problem of finding k cuts is to
minimize $\text{cut}(A_1, \dots, A_k) := \frac{1}{2} \sum_{i=1}^k W(A_i, \bar{A}_i)$

Example for $k=2$:



This formulation has one problem:

we can solve it by separating individual nodes
from the graph... those are bad clusters
(it's just outliers!!)

Thus, adjust to also include the set sizes!

We take the so-called Ratiocut.

$$\text{RatioCut}(A_1, \dots, A_k) := \sum_{i=1}^k \frac{\text{cut}(A_i, \bar{A}_i)}{|A_i|}$$

Note: look at $f = \frac{1}{\alpha} + \frac{1}{\alpha-1} \rightsquigarrow \frac{d}{d\alpha} f = -\frac{1}{\alpha^2} - \frac{1}{(\alpha-1)^2} = 0$

$$\begin{aligned}\alpha^2 - 2\alpha + 1 &= \alpha^2 \\ \Rightarrow \alpha &= \frac{1}{2}\end{aligned}$$

Connecting RatioCut to graph Laplacian

k=2: Our goal is to solve the following

$$\min_{A \in V} \text{RatioCut}(A, C(A))$$

First, we introduce the notation $\bar{A} = C(A)$

Then, we define the vector $f \in \mathbb{R}^n$

$$f_i = \begin{cases} \sqrt{|A|/|\bar{A}|} & \text{if } i \in A \\ -\sqrt{|A|/|\bar{A}|} & \text{else} \end{cases}$$

Basically, this is an indicator fct.

\hookrightarrow if $f_i > 0$, then $i \in A$

$f_i < 0$, then $i \in \bar{A}$

Note that:

$$\begin{aligned} 1) \sum_i f_i &= \sum_{i \in A} \sqrt{\frac{|A|}{|\bar{A}|}} - \sum_{i \in \bar{A}} \sqrt{\frac{|A|}{|\bar{A}|}} \\ &= \sqrt{|A||\bar{A}|} - \sqrt{|A||\bar{A}|} = 0 \end{aligned}$$

$\Rightarrow f$ is orth. to the 1-vector!

$$2) \sum_i f_i^2 = |\bar{A}| + |A| = \text{number vertices}$$

Let's evaluate $f^T L f$:

$$f^T L f = \frac{1}{2} \sum_{i,j} w_{ij} (f_i - f_j)^2$$

$$= \frac{1}{2} \underbrace{\sum_{i \in A, j \in \bar{A}} w_{ij} \left(\sqrt{\frac{|A|}{|\bar{A}|}} + \sqrt{\frac{|\bar{A}|}{|A|}} \right)^2}_{= \text{cut}(A, \bar{A})} + \frac{1}{2} \underbrace{\sum_{i \in \bar{A}, j \in A} w_{ij} \left(-\sqrt{\frac{|A|}{|\bar{A}|}} - \sqrt{\frac{|\bar{A}|}{|A|}} \right)^2}_{= \text{cut}(\bar{A}, A)} \text{ since } w_{ij} = w_{ji}$$

$$= \text{cut}(A, \bar{A}) \cdot \left(\frac{|A|}{|\bar{A}|} + \frac{|\bar{A}|}{|A|} + 2 \right)$$

$$= \text{cut}(A, \bar{A}) \cdot \left(\frac{|A|}{|\bar{A}|} + \frac{|\bar{A}|}{|A|} + \frac{|A|}{|\bar{A}|} + \frac{|\bar{A}|}{|A|} \right)$$

$$= \text{cut}(A, \bar{A}) \cdot \left(\frac{|A|+|\bar{A}|}{|\bar{A}|} + \frac{|\bar{A}|+|A|}{|A|} \right) \quad \begin{matrix} \# \text{ vertices} \\ V = |A| + |\bar{A}| \end{matrix}$$

$$= V \cdot \text{cut}(A, \bar{A}) \cdot \left(\frac{1}{|A|} + \frac{1}{|\bar{A}|} \right) \rightarrow \text{that's RatioCut!}$$

Thus: Optimising RatioCut is equivalent to

$\min(f^T L f)$ with $f \perp \mathbf{1}$ and $\|f\|_2^2 = V$
with f_i defined as before.

Solving this is NP-hard!

BUT we can relax it by allowing general vectors f !

$$\hookrightarrow \min_f (f^T L f) \text{ such that } f \perp \mathbf{1} \\ \|f\|_2^2 = V$$

This is similar to PCA, but instead of max, we have min!

\hookrightarrow Instead of the largest eigenvalue, we look at the lowest!

\hookrightarrow But $f \perp \mathbf{1}$, and $\mathbf{1}$ is the eigenvector of the lowest eigenvalue 0

\hookrightarrow Thus, this is solved by the rescaled eigenvector corresponding to the second-smallest eigenvalue!

Note: through the relaxation, we lose the simple rule $f_i \geq 0$ to assign clusters. Thus, we have to use k-means!

general k: For $k > 2$, we follow a similar approach. Let's assume a partitioning into k sets $A_i, i \in \{1, k\}$. Then define k indicator vectors

$$F_{ij} = \begin{cases} \frac{1}{\sqrt{|A_j|}} & \text{if } i \in A_j \\ 0 & \text{otherwise} \end{cases}$$

$$F \in \mathbb{R}^{v \times k}; v = \# \text{nodes}$$

F is orthonormal! $F^T F = I$

$$\hookrightarrow (F^T F)_{il} = \sum_j F_{ij}^T F_{jl} = \sum_j F_{ji} F_{jl}$$

$$\text{case 1 } (i=l): \sum_j F_{ji} F_{jl} = \sum_{j \in A_i} \frac{1}{\sqrt{|A_i|}} = 1$$

$$\text{case 2 } (i \neq l): \sum_i F_{ji} F_{jl} = 0$$

as both cannot be
>> 0 at the same time

In the following, let's define h_i as the i th column of F .

Let's look again at $h_k^T L h_k$

$$h_k^T L h_k = \frac{1}{2} \sum_{i,j} w_{ij} (\underbrace{h_{ki} - h_{kj}}_{\rightarrow f_{ik}})^2 = \overline{f_{jk}}$$

$$= \frac{1}{2} \sum_{\substack{i \in A_k \\ j \in \bar{A}_k}} w_{ij} \left(\frac{1}{|A_k|} - 0 \right)^2$$

$$+ \frac{1}{2} \sum_{\substack{i \in \bar{A}_k \\ j \in A_k}} w_{ij} \left(\frac{1}{|\bar{A}_k|} - \frac{1}{|A_k|} \right)^2 \underbrace{= 0}_{=0}$$

$$+ \frac{1}{2} \sum_{\substack{i \in \bar{A}_k \\ j \in \bar{A}_k}} w_{ij} (\underbrace{0 - 0}_{=0})^2$$

$$+ \frac{1}{2} \sum_{\substack{i \in \bar{A}_k \\ j \in A_k}} w_{ij} \left(0 - \frac{1}{|A_k|} \right)^2$$

$$= \text{Cent}(A_k, \bar{A}_k) \frac{1}{|A_k|} = \text{RatioCent}(A_k, \bar{A}_k)$$

In addition, we have:

$$(F^T L F)_{ii} = \sum_n F_{in}^T \sum_m L_{nm} F_{mi} = \sum_{nm} L_{nm} F_{ni} F_{mi}$$

$$= h_i^T L h_i$$

Thus: $\sum_i h_i^T L h_i = \sum_i (F^T L F)_{ii} = \text{Tr}(F^T L F)$

and we arrive at the following statement:

Minimizing $\text{RatioCut}(A_1, \dots, A_k)$

can be rewritten as

$$\min_{A_1, \dots, A_k} \text{Tr}(F^T L F) \text{ such that } F^T F = I$$

F defined as previously

which we can again relax to

$$\min_F \text{Tr}(F^T L F) \text{ such that } F^T F = I$$

Again, this is like PCA but with min instead of max. So the columns of F are given by the k eigenvectors of L with the smallest eigenvalues (including the 0 eigenvalue)!

Diffusion maps

The first step we have to do is construct a random walk on our data. For this, we have to turn it again into graph form!

Let $x_1, \dots, x_n \in \mathbb{R}^p$ be our n data samples. As for spectral clustering, each data point is represented by a node in the graph, and edges are weighted by the similarity of data points, e.g., $w_{ij} = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$.

$n=3$:



Next, we define the transition matrix:

$$M_{ij}^{ii} = \frac{w_{ij}}{D_{ii}}, \text{ with } D_{ii} = \sum_j w_{ij}$$

and $D_{ij}^{ii} = 0$ if $i \neq j$

We can write this in a more compact form:

$$M = D^{-1}W$$

This transition matrix represents the probability of moving from node i to node j :

$$P(\text{walk to node } j \text{ currently at node } i) = M_{ij}$$

Note: we construct the random walk like this!

Here, our random walker prefers walking between data points that are more similar!

There are a few properties:

1) $P(\text{walk to any node} | \text{currently at } i) = \sum_j M_{ij} = \frac{\sum_j w_{ij}}{\sum_j D_{ii}} = 1$

2) The probability to end up at node j when starting at node i and doing t steps is:
↳ expect!

$$P(\text{at } j \text{ after } t \text{ steps} | \text{start at } i) = (M^t)_{ij}$$

Why? Let's look at $t=2$.

$$(M^2)_{i,j} = (M \cdot M)_{i,j} = \sum_k M_{ik} M_{kj}$$

jump from i to k jump from k to j
sum over all k
= all possible paths!

Is this still normed?

$$\hookrightarrow \sum_j \sum_k M_{ik} M_{kj} = \sum_k \left(M_{ik} \underbrace{\sum_j M_{kj}}_{=1} \right) = \underbrace{\sum_k M_{ik}}_{=1} = 1$$

Basically, in "every" step we have to end up somewhere \Rightarrow probabilities are conserved.

Starting at node i , we can therefore extract the probabilities of where we end up after t steps:

$$P_{i,t} = e_i^T M^t = ([M^t]_{i,1}, [M^t]_{i,2}, \dots, [M^t]_{i,n})$$

↑
unit vector

Thus, we assign each data point a vector $P_{i,t}$ which tells us where a random walk starting at node i typically ends up after t steps.

Insight: If two data points i, j are "close", i.e., strongly connected in the graph, then $P_{i,t}$ and $P_{j,t}$ will be similar!

BUT $P_{i,t}$ and $P_{j,t}$ are not low-dimensional!
To find such a representation, we will look at the spectrum of M_t .

Spectral decomposition

Note that $M = D^{-1}W$ is, unfortunately, not symmetric. However, the matrix

$$S = D^{1/2} M D^{-1/2}$$

$$\text{is symmetric, since } S = D^{1/2} M D^{-1/2} = D^{1/2} D^{-1} W D^{1/2} \\ = D^{-1/2} W D^{-1/2}$$

S is symmetric and real-valued, so its spectral decomposition is given by

$$S = V \Lambda V^T \quad \text{with} \quad V^T V = \mathbb{1} \\ \Lambda = \text{diag}(\underbrace{\lambda_1, \dots, \lambda_n}_{\text{Eigenvalues}})$$

From this, we can recover M :

$$M = D^{-1/2} S D^{1/2} = \underbrace{D^{-1/2} V}_{\Phi} \Lambda \underbrace{V^T D^{1/2}}_{\Psi^T}$$

$$\Phi = (\Phi_1, \dots, \Phi_n) \\ \begin{matrix} \uparrow & \text{columns} \\ \vdots & \text{= vectors} \end{matrix}$$

$$= \Phi \Lambda \Psi^T \\ \begin{matrix} \uparrow & \text{rescaled} \\ \downarrow & \text{Eigenvectors of } S! \end{matrix}$$

$$\Phi \text{ and } \Psi \text{ are bi-orthogonal} \rightsquigarrow \Phi^T \Psi = \mathbb{1} \\ \text{or } \Phi_i^T \Psi_j = \delta_{ij}$$

Moreover, we have

$$\begin{aligned} M \Phi_k &= \underbrace{\phi \Delta \psi^T}_{\text{e}_k} \underbrace{\phi_k}_{\text{e}_k^T} = \phi \Delta e_k = \phi_k \lambda_k \\ &= e_k \underbrace{\gamma}_{(0, \dots, 1, \dots, 0)} \\ &\quad \underbrace{\text{e}_k^T \text{ entry}}_{\text{e}_k^T \gamma} \end{aligned}$$

and similarly: $\psi_k^T M = \underbrace{\psi_k^T \phi}_{\text{e}_k^T} \Delta \psi^T = e_k^T \Delta \psi^T = \lambda_k \psi_k^T$

$\Rightarrow \psi_k^T = \left. \text{left-eigenvector} \right\}$ of M with
 $\phi_k = \left. \text{right-eigenvector} \right\}$ eigenvalue λ_k

$\Rightarrow M$ can be represented as

$$M = \sum_{k=1}^n \lambda_k \underbrace{\phi_k \psi_k^T}_{\text{after product!}}$$

and $M^t = \sum_{k=1}^n \lambda_k^t \phi_k \psi_k^T$

Low-dimensional embeddings

We originally aimed for $e_i^T M^t$. We can represent this now using our new eigenbasis:

$$\begin{aligned} e_i^T M^t &= \sum_{k=1}^n \lambda_k^t e_i^T \phi_k \psi_k^T \\ &= \sum_{k=1}^n \underbrace{\lambda_k^t \phi_{k,i}}_{\text{vector components}} \underbrace{\psi_k^T}_{\text{eigenbasis}} \end{aligned}$$

With this, we define the diffusion map at step t as:

$$Y_{i,t} = \begin{bmatrix} \lambda_1^t \phi_{1,i} \\ \vdots \\ \lambda_n^t \phi_{n,i} \end{bmatrix} \quad \text{for } i=1, \dots, n$$

We can show that: 1) $\lambda_1 = 1$ and $\phi_1 = \underbrace{1}_{1 \times d}$
 2) $|\lambda_k| \leq 1$

\Rightarrow We

1) drop ϕ_1

and for large t , we 2) can truncate the vector as $\lambda_k \rightarrow 0$!

\Rightarrow Final embedding

$$y_{c,it}^{(d)} = \begin{bmatrix} \lambda_2^t \phi_{2,i} \\ \vdots \\ \lambda_{d+1}^t \phi_{d+1,i} \end{bmatrix}$$

These are our low-dim. embeddings of the data!

We can ground our initial intuition between similarity of random walks and being close in the embedded space some more:

Theorem: For any pair of nodes i, j , we have

$$\|y_{i,it}^{(n)} - y_{j,it}^{(n)}\|^2 = \sum_{l=1}^n \frac{1}{D_u} [P(\text{at } l \text{ after } t \text{ steps) start at } i) - P(\text{at } l \text{ after } t \text{ steps) start at } j)]^2$$

distance in
Eigenspace

difference of prob. to end up
at the same node after t steps.

Proof: write $P(\text{at } L \text{ after } t \text{ steps} | \text{start at } i)$

$$= P(Z(t) = L | Z(0) = i)$$

First, some rewriting:

$$\sum_{l=1}^n \frac{1}{D_{ll}} \left[P(X(t)=l | Z(0)=i) - P(Z(t)=l | X(0)=i) \right]^2$$

$$= \sum_{l=1}^n \frac{1}{D_{ll}} \left[e_i^T M^t e_l - e_j^T M^t e_l \right]^2$$

$$= \sum_{L=1}^n \frac{1}{D_{LL}} \left[\sum_{k=1}^n \lambda_k^t [\phi_{k,i} - \phi_{k,j}] \psi_{k,L}^T \right]^2$$

$$= \sum_{l=1}^n \left[\sum_{k=1}^n \lambda_k^t [\phi_{k,l}^i - \phi_{k,l}] \frac{\psi_{k,l}^T}{\sigma_{k,l}^2} \right]^2$$

$$\mathbf{V}^T \underbrace{\mathbf{D}}_{1/2} \mathbf{V}$$

$$= \sum_{l=1}^n \left[\sum_{k=1}^n \lambda_k^t [\phi_{k,i} - \phi_{k,j}] \underbrace{\sqrt{v_{k,l}}}_{\text{pref.}} \right]^2$$

sum of vectors
pref.

and $v_{k,l} = \sqrt{v_{k,k} D_{ll}}$

$$= \left\| \sum_{k \in \Lambda} \lambda_k^t [\phi_{k,i} - \phi_{k,j}] v_k^\top \right\|^2$$

$$\begin{aligned}
 &= \sum_{k=1}^n (\lambda_k^t [\phi_{k,i} - \phi_{k,j}])^2 \\
 \text{✓ } k \text{ is orthogonal!} \quad &\stackrel{\phi_1 = 1}{=} \sum_{k=2}^n \left(\underbrace{\lambda_k^t \phi_{k,i}}_{y_{irt,k}^{(n)}} - \underbrace{\lambda_k^t \phi_{k,j}}_{y_{jrt,k}^{(n)}} \right)^2 \\
 &= \| y_{irt}^{(n)} - y_{jrt}^{(n)} \|^2 \quad \square
 \end{aligned}$$

Note on spectral clustering:

$$\text{Let } L_{rw} = I - D^{-1}W$$

$$\text{and } L_N = D^{-1/2} L D^{-1/2} = I - D^{-1/2} W D^{-1/2}$$

If v is an eigenvector of L_{rw} with eigenvalue λ , then $u = D^{1/2}v$ is an eigenvector of L_N with eigenvalue λ .

$$\begin{aligned}
 \underline{\text{Proof:}} \quad L_N u &= \lambda u \Leftrightarrow D^{-1/2} L_N u = \lambda D^{-1/2} u = \lambda v \\
 \text{and } D^{-1/2} L_N u &= D^{-1/2} L_N D^{1/2} v = L_{rw} v \\
 &\Leftrightarrow L_{rw} v = \lambda v \quad \square
 \end{aligned}$$