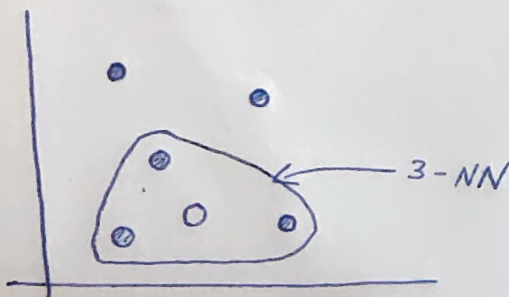## Q1) Graph learning with missing data

In real world scenario, missing node attributes is a common problem that we encounter while learning graph. From the problem statement, two things can be inferred,

1) For a given node, the data is partially missing & not entirely

2) Multiple attributes are missing from a single node.

This is a problem of multivariate missing value imputation. There are several techniques in Machine learning and Data mining literature for this problem. Here we will be making use of famous k-NN imputation to find missing attributes.

The key idea of k-NN is that we make use of similar nodes which are 'nearby' and impute the missing values. Here, Euclidean distance is calculated between nodes to find the nearest neighbour by ignoring the missing values.



3-NN

There are some Graph Neural Network Architectures which discuss about missing value imputation along with graph learning. {" Wasserstein diffusion on graphs with missing attributes" by Chen, Zhixian and Ma, Tengi & Song, 2021}

But due to time limit, I have decided to go with k-NN imputation.

→ Now coming to graph learning, there are 3 main kind of graph learning

1) Linear combination model

2) Smooth signal model

3) Probabilistic Graphical model

Here, we will be learning graph under the assumption of smoothess. { Paper : "How to learn a graph from smooth Signals" Vassilis Kalofolias"}

→ A graph signal is smooth if the signal values associated with two end vertices of edge with large weights in the graph tend to be similar

→ Define, combenatorial graph laplacian $L = D - W$

    $W \to$ Weight matrix

    $D \to$ Diagonal matrix with each element equal to sum of weights connected to that node

    $D = \sum_j W_{ij}$

→ For the smoothness, we need need to minimize,

    $\min \ \text{trace}(X^T L X)$

→ $\text{tr}(X^T L X) \ = \ \text{tr}(X^T D X) - \text{tr}(X^T W X)$

$$= \ \sum_{i=1}^{m} \sum_{j=1}^{m} x_i^T W_{ij} x_i - \ \sum_{i=1}^{m} \sum_{j=1}^{m} x_i^T W_{ij} x_j$$

$$= \tfrac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} (x_i - x_j)^T W_{ij} (x_i - x_j)$$

$$= \tfrac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} W_{ij} \| x_i - x_j \|_2^2$$

$$= \tfrac{1}{2} \| W \circ Z \|_{1,1} \ = \ \tfrac{1}{2} \ \text{tr}(W, Z)$$

where $Z = \| x_i - x_j \|_2^2$   an
    $ij$

and   $\circ \to$ Hadmard product

So, the objective func. is

$$\min_{W} \quad \frac{1}{2} \sum_{f} W_{ij} \| x_i - x_j \|^2 + f(W)$$

Role of $f(W)$ : $f(W)$ has 2 ~~l~~ roles to play

1) ~~Using l₁ norm~~

1) To prevent $W$ from going into trivial solution $W = 0$

2) To impose a prior belief in the graph signal

→ Using $l_1$ norm as $f(W)$ will be usefull in creating a sparse $W$, but it will translate into adding a constant to square distance in $Z$, so we make use of $l_2$ norm.

$$tr(x^T L x) + \eta \| W \|_{1,1} = \frac{1}{2} \| W \circ (2\eta + Z) \|_{1,1}$$

→ The final objective function :

$$\min_{W} \quad \| W \circ Z \|_{1,1} - \alpha \, 1^T \log (W \cdot 1) + \beta \| W \|_F^2$$

where $\quad 1 = [1, 1 \ldots 1]^T$

$W 1$ → Diagonal elements of $L$

Taking the log prevents ~~from~~ the weights of a node going to zero, or in other words, it prevents formation of isolated nodes.

→ The addition of the Frobnius norm will penalize larger weights but do not penalize smaller ones.

→ ~~The well~~ $\alpha$ and $\beta$ are the ~~type~~ parameters which can be tuned.

⇒ Furthemore, if we do a change of variable from
$$\tilde{w} = w/\gamma, \quad \text{we can show that}$$

$$\beta \quad F(z, \alpha, \beta) = \gamma F(z, \alpha/\gamma, \beta\gamma)$$

$$= \alpha F(z, 1, \alpha\beta)$$

→

This implies, if we ~~want~~ to ~~learn~~ obtain a $w$ with fixed scale, then we can solve by tuning only one parameter $\beta$ (Take $\alpha = 1$).

→ For a the optimization, point of view, we use vector form representation of $w$

⇒ $\min\limits_{w \in W} \quad f_1(w) + f_2(Kw) + f_3(w)$

where $W1 \geq Kw$ to impose a fixed scale on $w$

$$f_1(w) = 1$$

where,  $f_1(w) = 1\{w \geqslant 0\} + 2 w^T z$

$$1\{w \geqslant 0\} = \begin{cases} 0 & \text{if } w \geqslant 0 \\ \infty & \text{for otherwise} \end{cases}$$

$f_2(Kw) = f_2(d) = -\alpha \, 1^T \log(d)$

where  $Kw = W1$  to impose a fixed scale on $W$

$f_3(w) = \beta \|w\|^2$