Q1) a) $P_{ij} = \dfrac{\left(1 + \|x_i - x_j\|^2\right)^{-1}}{\sum_{k=1}^{N} \sum_{l \neq k} \left(1 + \|x_k - x_l\|^2\right)^{-1}}$

Consider $Q_{ij}$, where $y = x_i$ : $\dfrac{\left(1 + \|y_i - y_j\|^2\right)^{-1}}{\sum_{k=1}^{N} \sum_{l \neq k} \left(1 + \|y_k - y_l\|^2\right)^{-1}}$

if $y_i = x_i$ based on the assumption, $P_{ij} = Q_{ij}$

Thus KL Divergence $= \sum_{i \neq j} P_{ij} \log\left(\dfrac{P_{ij}}{Q_{ij}}\right) = \sum_{i \neq j} P_{ij} \underbrace{(0)}_{\boxed{= 0}}$

$\underset{1}{\underbrace{\qquad}}$

b) Doing the substitution: $q_{ij} = P_{ij}$

$KL(P \| Q) = \sum_{k=1}^{N} \sum_{i \neq j} P_{ij} \log\left(\dfrac{P_{ij}}{q_{ij}}\right)$    , $z_i = y_i - c$

$q_{ij}(y) = \dfrac{\left(1 + \|y_i - y_j\|^2\right)^{-1}}{\sum_{k=1}^{N} \sum_{l \neq k} \left(1 + \|y_k - y_l\|^2\right)^{-1}}$    $q_{ij}(z) = \dfrac{\left(1 + \|z_i - z_j\|^2\right)^{-1}}{\sum_{k=1}^{N} \sum_{l \neq k} \left(1 + \|y_k - y_l\|^2\right)^{-1}}$

$z_i - z_j = (y_i - c)(y_j - c) = y_i - y_j = z_i - z_j$

$z_k - z_l = (y_k - c) - (y_l - c) = y_k - y_l = z_k - z_l$.

Thus, both distribution are equal and hence

KL Divergence $KL(P \| Q(y)) = KL(P \| Q(z))$

This illustrates that shifting a fixed set with lower dimension space Q, doesn't alter the overall relationship and mapping between higher & and lower dimesional space. Thus, relationship depends on the relative positions and distances between points, unaffected by such translations.

c)   Using formula.  $\|x_1 - x_2\|^2 = \alpha$   $\|x_1 - x_3\|^2 = \beta = 100\alpha$

$$P_{3|1} = \frac{e^{-\beta/2\sigma_1^2}}{e^{-\alpha/2\sigma_1^2} + e^{-\beta/2\sigma_1^2}} = \frac{e^{-50\alpha/\sigma_1^2}}{e^{-\alpha/2\sigma_1^2} + e^{-50\alpha/\sigma_1^2}}$$

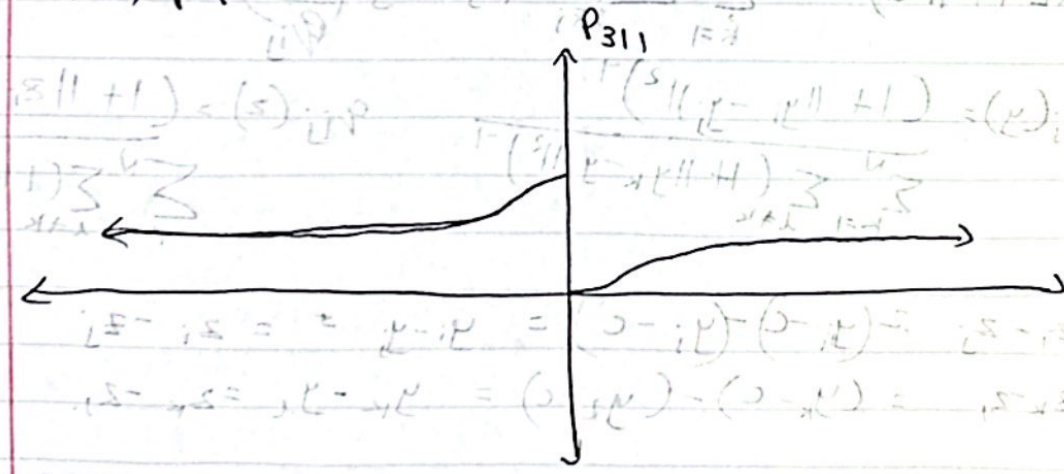Limits   $\sigma_1^2 \to \infty$   and   $\sigma_1^2 \to 0$.

Both terms in denom approach 1 as $e^0 = 1$.

$P_{3|1} = \frac{1}{2}$

The term $e^{-50\alpha/\sigma_1^2}$ approaches 0 much faster than $e^{-\alpha/2\sigma_1^2}$
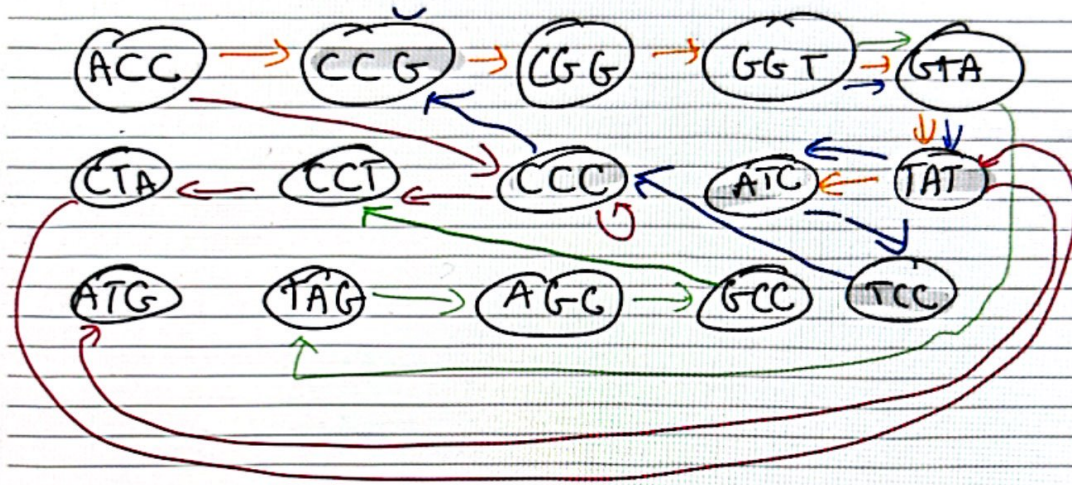
therefore, $P_{3|1} \to 0$

Thus, graph looks like



Implications

When $\sigma_1^2$ is large, model considers both nearby and distant points equally, leading to less distinction between $x_2$ and $x_3$ relative to $x_1$. When $\sigma_1^2$ is small, distinction becomes sharp, effectively ignoring distant points. $x_3$ relative to $x_1$ and $x_2$ preserving more local structure. Changing $\sigma_1^2$ impacts relative influence of distant points in t-SNE. greater $\sigma_1^2$ leads to a more global perspective, while smaller $\sigma_1^2$ values emphasizes local relationship

Q2)

E₁:  ACC    CCG    CGG    GGT  GTA  TAT  ATC

E₂:  ACC    CCC    CCC    CCT  CTA  TAT  ATG

E₃:  GGT  GTA    TAG    AGC  GCC  CCT

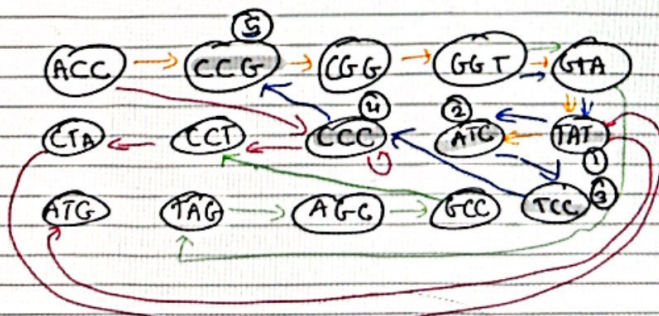t₄:  GGT  GTA    TAT    ATC  TCC  CCC  CCG



b) Pseudoalign with Skipping

TATCCCG → TA

TAT → ATC → TCC → CCC → CCG

forward complement.

Find TAT. Find ATC
se then find connection
between two kmer nodes
we do this and note
down all



Completely aligns.
with true

t₄

c) Pseudoalign with skipping

CGGGATC

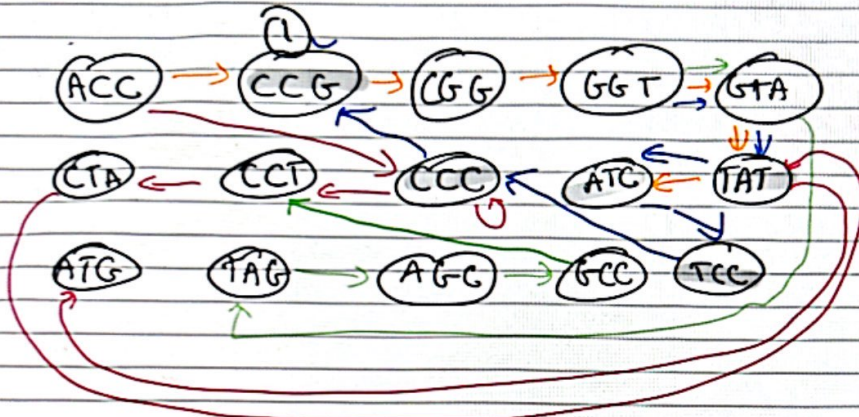CGG → GGG → GGA → GAT → ATC

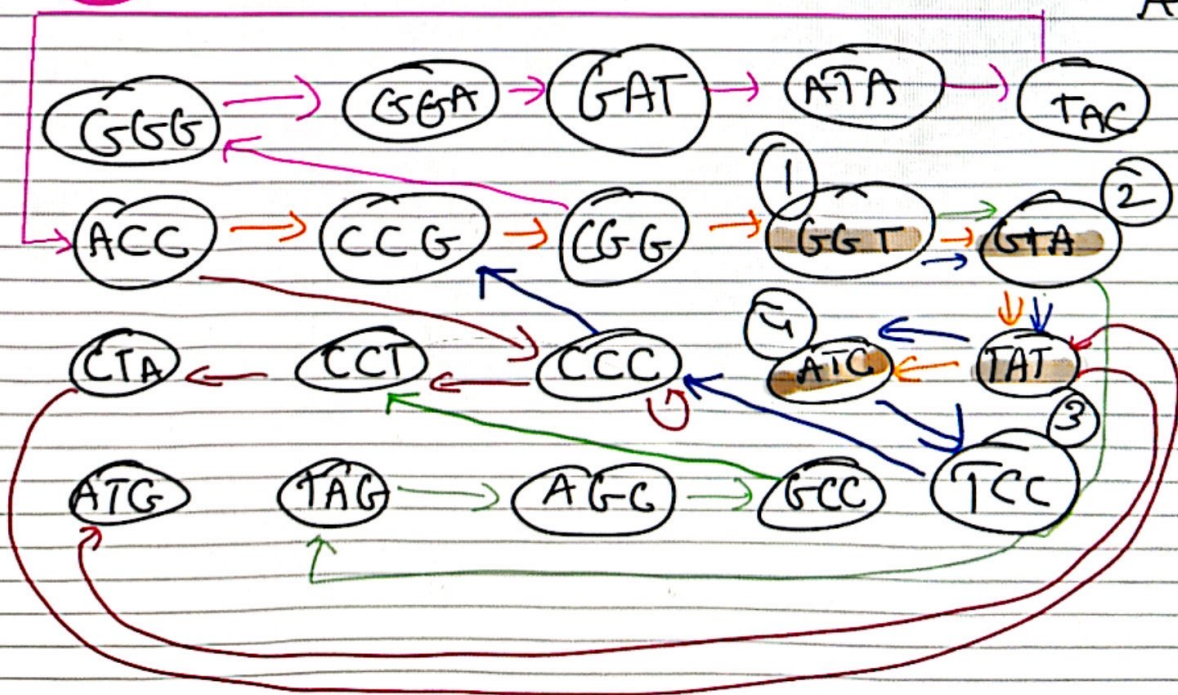c) Pseudoalign with skipping

    C GGG ATC

CGG → GGG → GGA → GAT → ATC

GGG Not found. GGA Not found. and GAT not found. Thus, we can't pseudoalign the sequence because no match is found from the first kmer to the last kmer (CGG → ATC).



first we find CCG. We find GTA which doesn't exist in the graph. After finding nullset, we terminate different alignment.

Reverse complement → GATCCCG → GAT Not found, thus terminate.

Union of forward and reverse is thus the null set
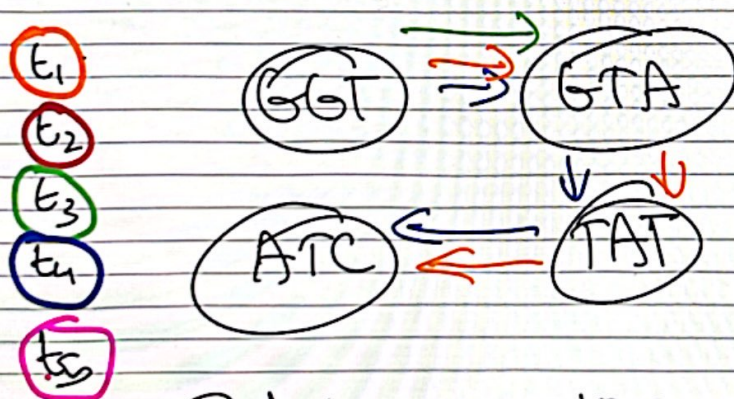
d) (t_5)    CGG   GGG   GGA  GAT ATA TAC
                                      ACC

e) Pseudoalign with skipping:

$$GGT \rightarrow GTA \rightarrow TAT \rightarrow ATC$$

From the graph above, we first visit GGT. , then, we visit GTA which aligns with $t_3, t_1$ and $t_4$. We then visit TAT, which gives a match with $t_1$ and $t_4$. Same with ATC.

Nodes visited

① → ② → ③ → ④

$t_1$
$t_2$
$t_3$
$t_4$
$t_5$

GGT → GTA
ATC ← TAT

Taking union

We get read = { }

for the entire length of the read

Q3a) ⟹ First , we choose the cell state. $(G_1, S, G_2, m)$
⟹ After this step. the machine labels the cell with the specific barcode corresponding to its shape.
→ labelled cells are then integrated in the barcode via droplet creation. No perturbations means that the cell's natural state is preserved throughout.

→ At the end of the process, cells in $m$ state die with $\delta$ probability. Observed gene expression of that cell droplet goes to zero.

→ For each cell, the only the expression of 1 gene is sampled. complete dataset consists of a vector of length C, where each entry represents the expression level of a single gene in 1 cell.

3b) Random Variables to consider.

→ Cell state → $S_c \sim$ Categorical $(G_1, S, G_2, M)$

→ Cell label → $B_c \sim$ Categorical $(P',)$

distribution of Barcode labels.

→ ~~As cell~~ Cell death →

for all cells death rate is same

cell only dies in state $M$. → ~~than~~

$P(\text{dead cell} | M) = \delta$

(~~cells don't~~ die $P(\text{dead cell} | (G_1, S, G_2) = 0$

$(i = 0 \to c)$

→ ~~Observation~~. $P(D_i) = \begin{cases} 1 - \delta^{E}_X(A) & P_c = \sim D(\text{Alive}) \text{ where } A \text{ is} \\ \delta(A) & D_c = \delta(\text{dead}) \end{cases}$ # cells in state $M$

Observation:

→ Gene expression → $D_i = \begin{cases} 0 & \text{if cell } i \text{ is alive} \\ \Theta_1 & \text{if cell } i \text{ is dead} \end{cases}$

$O_{\text{GE}} \sim$ GE

~~For~~ Gene expression

For $D_i = 1 \to$ Gene expression is zero

Ilse, $D_i = 0 \to$ ~~gene expr~~

$O_c = $ Distribution $\&(\overset{\text{Expected}}{\text{Expression}} | S_c) \cdot \delta(F)$

# cells are that alive

$O_1 = 0$ if cell is dead.

$$= \frac{0.8(0.01)}{0.0575} \cong 0.1391$$

proportioned
$$\propto (G_1, S, G_2, M)$$

$\delta$ (Probability of death in state M)

cell state $S_c$ → $D_i$ → $O_c$

for $i = 0 \to c$