

→ ITC-Exp → Sampling results  $\Rightarrow$   
Candidate

C-L Pool size  $\Rightarrow \underline{\underline{100}} \Rightarrow$

New layer pts = 3469

Val, test = 433, 435 resp

• Feature size  $\Rightarrow$  Graph Embedding size  $\Rightarrow \cancel{20} [1 \times 20]$   
 $\hookrightarrow$  Rest G(N-Layer's output  
 size = 20

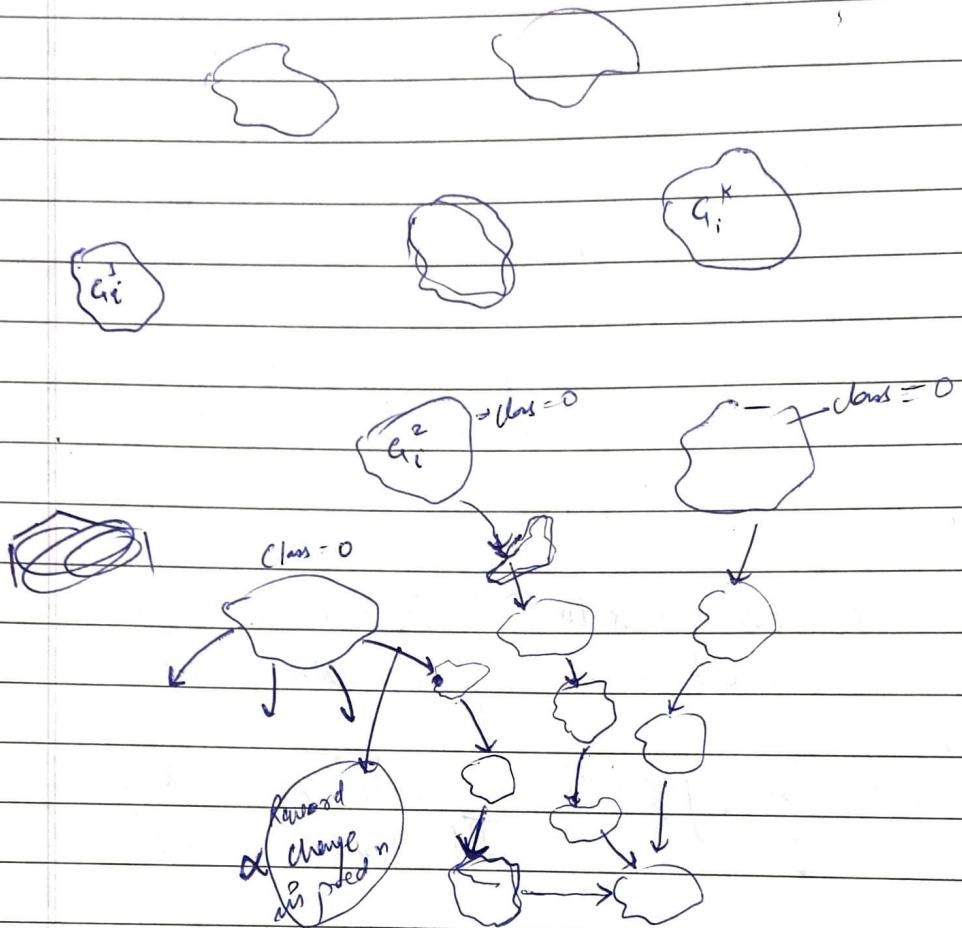
• Restricting the Random Walk to Delta Domain Dis  $\Rightarrow$  Not  
 going out of data manifold or expanding the delta manifold too much

$\hookrightarrow$  In GCFE  $\rightarrow$  they only go to a distance  $\theta$  from any input graph  
 $\hookrightarrow$  fixed beforehand (hyperplane)  
 ie.  $\rightarrow$  similar w.r.t classification

$\Rightarrow$  simply perturbing the set of similar graphs (similar as per trained model, not necessarily based on similarity of just graph structure & features) till we reach a common cf  $\Rightarrow$  why (we check cf by the pred of trained model on proposed cf)  $\Rightarrow$  does not seem

to be very appropriate in terms of explaining the model, as it may not fully reflect the representations learnt by GNN that we may want to explain.

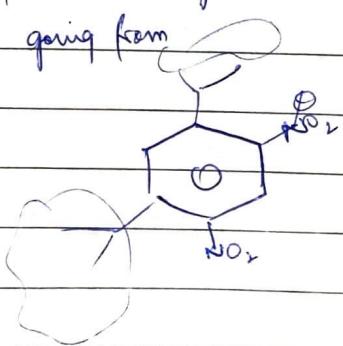
$\Rightarrow$  It would generate VCFs  $\rightarrow$  but these may still not be CF explanations (in the true sense)



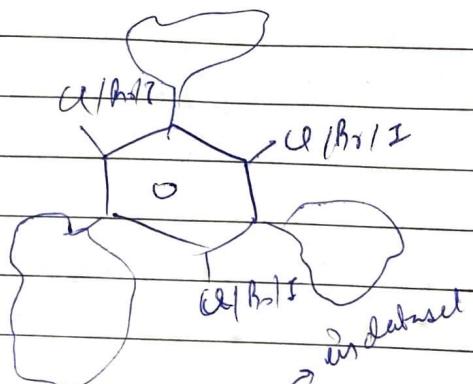
\* How to make sure that the Random walk does not go to far from the data set / manifold?

$\epsilon_n \rightarrow$  No point in changing

going from



to



If most graphs don't have a Cl/Br/I directly attached to the benzene

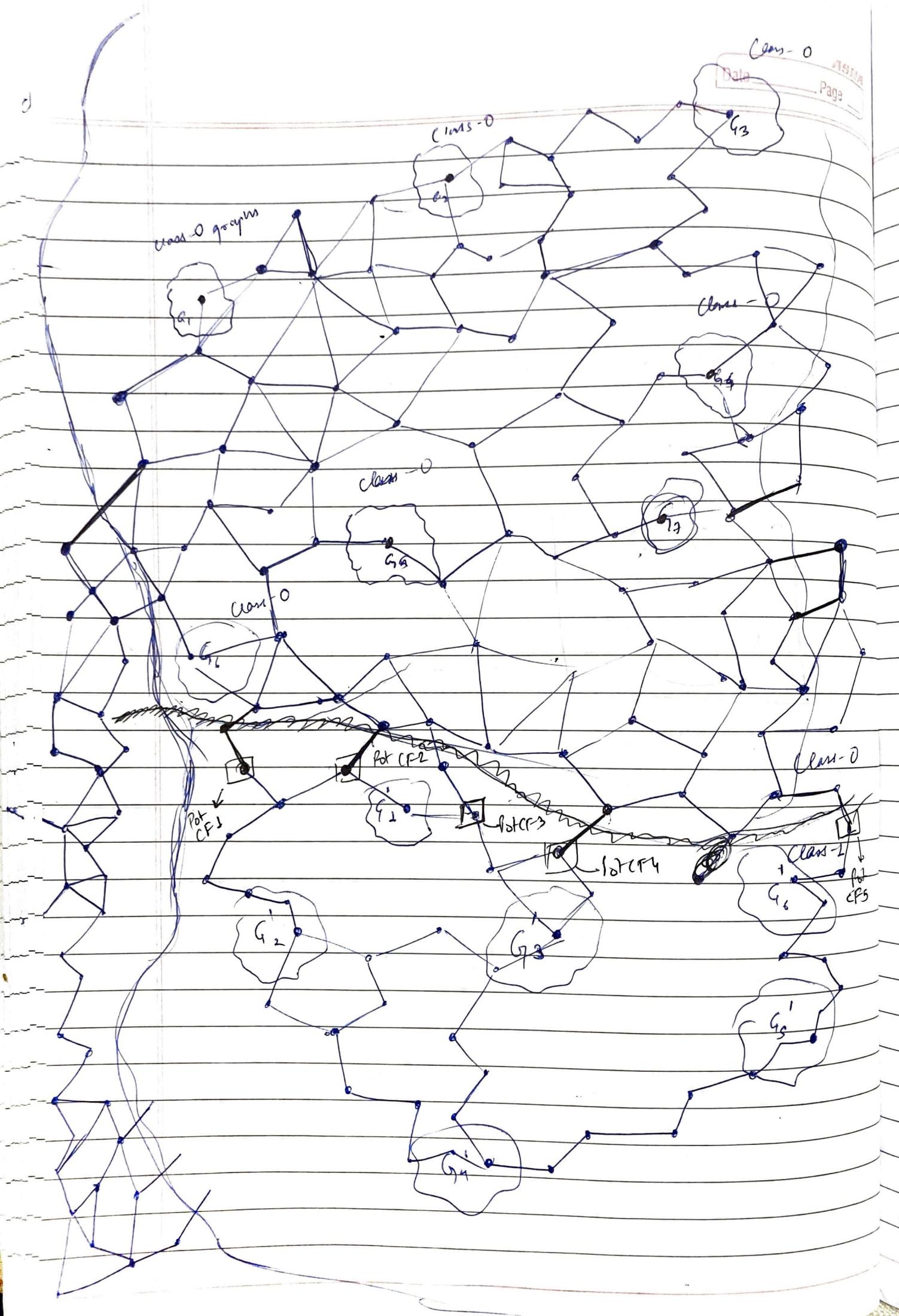
→ An check key computing GNN embedding of the new or added graph → if the GNN embedding is outlier → we may not want to explore it further → But has expl<sup>n</sup> vs exploitation tradeoff

An intermediate graph may be too out of data domain but the fine cf may be in ~~domain~~ domain. → cf needs May keep this learnable by using RL (introducing a Reward setting) to automatically figure out what is good, what's not

↓  
① Approach 1: Re-compute clusters of repr<sup>n</sup> of input graphs,  
& check if the new graph belongs to  
any one cluster within a dist threshold  
→ ~~can be~~ Feasible to compute  
(Takes same time as  
checking if a new the  
new graph is CF / not CF)  
+ Distance from all clusters  $O(k)$

② Approach 2: After each new graph is formed →  
embeddings of  
recluster all the graphs till now → & sum  
do local outlier factor to check if it is outlier or  
not? → ~~Infeasible Approach~~ (Not possible to recluster  
at each time step)





→ Consider the convex hull of all embeddings (both class 0 & class 1) →

Then, a new graph can be checked to be an outlier  
if it is more than  $\theta$  dist from convex hull

↓  
Can be  
pre-computed

(ideally, it's embedding should lie within the convex hull)

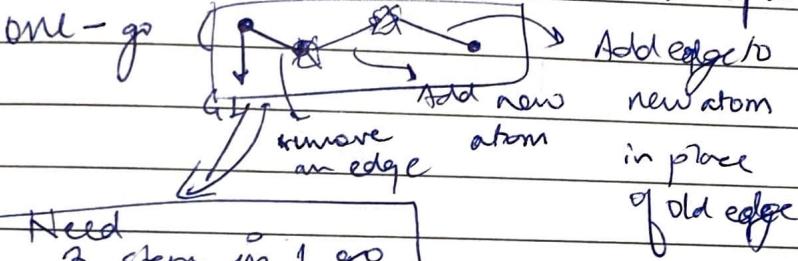
→ Complexity of this computation?

→ Penalise → Such potential (P) which are at high distance from this convex hull ⇒

(Also → read about some k-nearest nbrs like approaches for finding whether a new point is within data-manifold.)

\* Also → sometimes can't just remove an edge (for ex-  
molecular graphs → if removing an edge → will need  
to add some other edge to satisfy the valence)

↳ how to handle that ⇒ Need multiple steps in



→ Can select

among these pot CF & arrange them in ↑ order of  
their distance to avg distance to all input graphs to  
be explained.

Can tackle this

using immutable features ⇒ If a transition  
leads to a change in an immutable feature →  
keep making sequential changes till that immutable feature  
change is restored.

↳ if we simply don't allow such transitions  
we won't be able to generate new graphs ⇒  
As we need to remove an edge first at any  
step for an edge removal action.

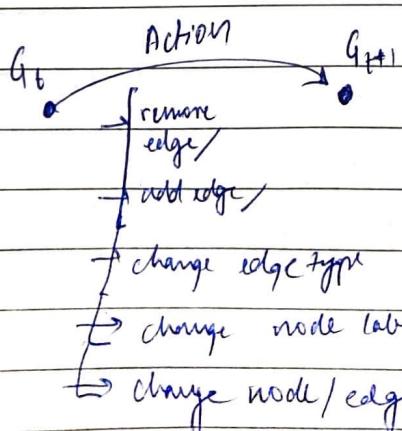
i Reward for transition to a new graph  $\Rightarrow$  ~~Prediction~~  
 $\propto \uparrow$  in  $p_{\text{pred}}^n$  for altered label /  
Decrease in  $p_{\text{pred}}^n$  for current  
(undesired label)

$\alpha$   $\downarrow$   $\frac{1}{\text{dist}}$  embedding  
Distance from convex hull  
of embeddings (data manifold)

(maybe)  $\propto \frac{1}{\text{dist}}$   
Any dist from  
all graphs to be  
explained  
↓ how to quantify this  
compute  
efficiently

\* Still all this is not using the Representations learnt  
by the GNN or its learnt decision boundaries

~~Reward~~



$$R_{t+1} = [P(\hat{y} = \text{y desired}) - P(\hat{y} = \text{y current/non-dm})]$$

$$d_1 \quad d_2 \quad d_3$$

$$\cancel{e} \cdot \cancel{e} \cdot \cancel{e} \quad (d_1 + \epsilon) \quad (d_2 + \epsilon) \quad (d_3 + \epsilon)$$

$d_i = \text{distance from}$   
boundary of convex hull of  
all embeddings

$d_2 \rightarrow \text{avg distance of } G_{t+1}$   
from all graphs to be  
explained

We can have an estimate of  $d_2$  using the dist of GNN Embedding of  $G_{t+1}$  to all graphs to be explained.

Date \_\_\_\_\_ Page \_\_\_\_\_

ASIAN

Also  $\propto$  Data-density at the new pt  $G_{t+1}$

(How to estimate that  $\rightarrow$  One proxy could be

1

$$\frac{\text{sum of dist of } k\text{-nearest}}{\text{nbr}}$$

$d_3 = \frac{\text{sum of}}{\text{avg}} \text{dist of } k\text{-nearest nbrs (ie Actual graphs)} \quad k\text{-nearest}$   
 in the dataset  $\rightarrow$  not generated graphs

But this will take  $O(N \log N + N)$  complexity  
 where  $N = \text{no of past train samples} \rightarrow$  for each transition  $\rightarrow$  Can't go this  $\rightarrow$  Choose this approach

$$R_{t+1} = \frac{[p(\hat{y} = y_{\text{desired}}) - p(\hat{y} = y_{\text{current}} / \text{non desired})]}{e^{d_1 (d_2 + \epsilon) (d_3 + \epsilon)}}$$

$d_1 \Rightarrow$  distance from boundary of convex hull  $\rightarrow d_1 \geq 0$  if pt outside C.H

$d_2 \Rightarrow$  avg dist of embedding of  $G_{t+1}$  from all graphs  $\quad 0 \leq d_2 \leq w$   
 to be explained

$d_3 \Rightarrow$  Proxy for data density at  $G_{t+1}$   $\rightarrow$  sum of dist of  $k$  nearest nbrs