

CLUSTER CON VINCOLI

Potrebbe essere necessario tracciare dei vincoli per forzare l'attribuzione dei punti a determinati cluster, perché magari ci sono vincoli fisici o vincoli richiesti dall'utente (Una banca che vuole piazzare i suoi ATM in città deve tener conto degli ostacoli come fiumi e colline)

CONSTRAINTS ON INSTANCES

Questa tecnica consiste nello specificare quando una coppia o un set di oggetti dovrebbero essere messi insieme

1) MUST-LINK(x, y)

x e y devono stare nello stesso cluster.

2) CANNOT-LINK(x, y) se $\text{dist}(x, y) > \text{Threshold}$

CONSTRAINTS ON CLUSTERS

Specificare dei requisiti per i cluster stessi, come per esempio il numero minimo di oggetti nel cluster, il diametro massimo, la forma, numero di cluster, ...

1) d -constraints

Per ogni coppia di cluster C_i e C_j e $\forall i, j$

Per ogni coppia di istanze $x \in C_i$ e $y \in C_j$ $\forall x, y$

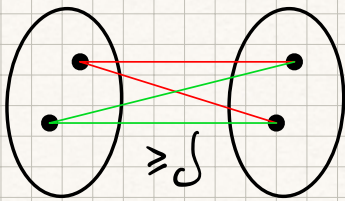
$$D(x, y) \geq d$$

2) ϵ -Constraints

Per ogni cluster $C_i : |C_i| > 1$

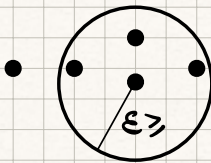
\forall punto $x \in C_i, \exists y \in C_i : \epsilon \geq D(x, y), x \neq y$.

Il d -Constraint può essere convertito in un vincolo di tipo **must-link**:



Per ogni punto x applica **must-link**(x, y) & $\text{dist}(x, y) < d$.

Mentre l' ϵ -constraint può essere visto come un **must-link** del tipo:



Per ogni punto x **must-link** almeno un punto $y : \text{dist}(x, y) \leq \epsilon$.

CONSTRAINTS ON SIMILARITY MEASURES

Specificare dei requisiti che il calcolo della similarità fra gli oggetti deve rispettare.

HARD vs SOFT

1) **Hard Constraints**: Un cluster che viola questo vincolo è inaccettabile

2) **Soft constraints**: Un cluster che viola questo vincolo è accettabile quando non ci sono soluzioni migliori.

Come capire se un vincolo è utile o meno?

1) INFORMATIVENESS

È la quantità di informazione che ha il vincolo e che va oltre il metodo di cluster usato. Si può vedere come, preso un modello di cluster A e un insieme di vincoli C, il numero di vincoli che non sono soddisfatti dal metodo A.

2) COHERENCE OF A SET OF CONSTRAINTS

Quanto i vincoli non si ostacolano fra loro e quindi siano coerenti.

COP K-MEANS : Hard Constraints

Praticamente crea delle SUPER-ISTANZE per i vari must-link:

1) **Transitive closure** dei must-link, dove se esistono due must-link : $\text{must-link}(A, B)$ e $\text{must-link}(B, C)$ di conseguenza esiste il $\text{must-link}(A, C)$ implicito.

Sostituire gli oggetti della transitive closure con un oggetto dato dalla media degli oggetti nella transitive closure, questa è detta **super-istanza**.

Questa super-istanza è pesata col numero di oggetti che rappresenta.

Usare ora un **k-means modificato** per la gestione dei cannot-link. Quando bisogna assegnare il punto a un centroide bisogna vedere se c'è un cannot-link che lo impedisce.

```
1 def COPkmeansHC(Su = unlabeled_data, Sl = labeled_data, K = number_of_clusters, q = number_of_constraints):
2     ML = []
3     CL = []
4     # Generazione dei q constraints
5     for i in range(0, q):
6         # Scegli due punti a caso dal set labeled.
7         # Se appartengono alla stessa classe, aggiungi la coppia nei MustLink. Altrimenti, aggiungi la coppia nei CannotLink.
8         x = randomly select point from Sl
9         y = randomly select point from Sl
10        if (label(x) == label(y)):
11            ML = ML U (x, y)
12        else:
13            CL = CL U (x, y)
14
15    # CC sono i punti uniti dalle varie transitive closure
16    CC = transitive_closure(ML)
17
18    # calcola il punto medio di ogni CC
19    r = CC.size()
20    for i in range(1, r):
21        mean = calculate_mean(CC[i])
22
23    # genera k centroidi iniziali a caso
24    centroids = []
25    for i in range(0, K):
26        centroids = centroids U randomly select point from Su
27
28    while convergence is not reached:
29        # assegna ogni punto al centroide più vicino
30        for i in range(0, Su.size()):
31            assign_point_to_nearest_centroid(Su[i], centroids)
32        # aggiorna i centroidi
33        for i in range(0, centroids.size()):
34            centroids[i] = calculate_mean(centroids[i])
35
36    return centroids
```


Handling Soft-Constraints

Quando un cluster non rispetta un soft-constraints allora viene penalizzato.

L'obiettivo è massimizzare la qualità dei cluster e minimizzare il numero di vincoli violati.

Objective Function

Somma delle distanze usata da K-means per calcolare la media dei punti di un cluster, per poi ricavarne i nuovi centroidi:

$$\frac{1}{|C_i|} \sum_{p \in C_i} p + \text{penalità}$$

Così da spostare il centroide.

1) Penalizzare un must-link violato

Viene aggiunta la distanza $\text{dist}(C_1, C_2)$ alla funzione obiettivo, ovvero la distanza dei due cluster C_1 e C_2 che hanno oggetti che violano il must-link.

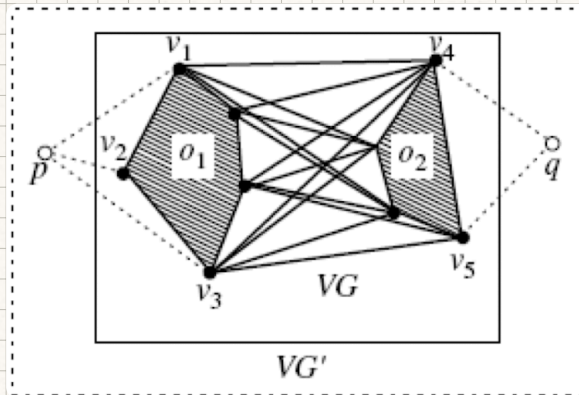
$$\text{must-link}(x, y) \text{ ma } x \in C_1 \text{ e } y \in C_2.$$

2) Penalty of a cannot-link violation

Quando due oggetti x, y con vincolo cannot-link sono associati allo stesso centroide C , allora viene aggiunta alla funzione obiettivo la distanza fra C e il suo più vicino altro centroide C' .

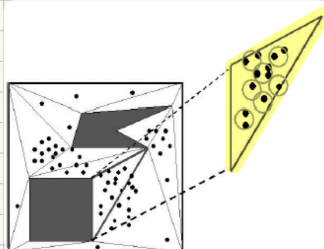
VELOCIZZARE LE PERFORMANCE

Si utilizza un grafo detto **visibility graph** che rappresenta gli ostacoli. Ogni vertice è un "angolo" di un ostacolo, e i nodi vengono uniti da archi se e solo se si possono "vedere" fra loro. Per esempio se vogliamo rappresentare una casa, allora collegheremo solo i nodi corrispondenti agli angoli della casa che si vedono:



Creiamo il grafo VG' partendo da VG , aggiungendo i due nodi p e q . Colleghiamo i nodi mutuamente visibili fra loro.

Possiamo ora raggruppare i punti del dataset in **MICRO-CLUSTER**, ovvero gruppi di punti ottenuti dividendo la regione in triangoli fra i vertici del grafo, e poi utilizzare un qualsiasi algoritmo di cluster per cercare i micro-cluster



Usando i micro-cluster la complessità computazionale viene ridotta.