

L'indice di Davies-Bouldin.

Sono dati N punti nel piano Euclideo e si vuole definire una loro partizione in un dato numero K di clusters, in modo da ottimizzare l'indice di Davies Bouldin, che è definito come segue.

Per ogni cluster C_k , sia T_k il numero di punti assegnati al cluster. Usando la norma Euclidea, si definisce un indicatore di compattezza del cluster C_k come

$$S_k = \sqrt{\frac{1}{T_k} \sum_{j \in C_k} ((x_j - x_k)^2 + (y_j - y_k)^2)}$$

dove si sono indicate con indice j le coordinate dei punti del cluster e con indice k le coordinate del centroide del cluster.

Si definisce poi per ogni coppia di clusters C_h e C_k un indicatore di separazione, M_{hk} , che è dato dalla distanza Euclidea tra i due centroidi.

La misura di quanto bene due clusters C_h e C_k sono distinti è data da $R_{hk} = \frac{S_h + S_k}{M_{hk}}$, che si vorrebbe basso; in tal modo, infatti, la partizione risulta buona quando i due clusters sono compatti (valori di S bassi) e distanti tra loro (valore di M alto).

L'indice di Davies Bouldin è definito come $DB = \frac{1}{N} \sum_{k=1}^K \max_{h \neq k} \{R_{kh}\}$.

Formulare il problema, classificarlo e risolvere l'esempio riportato con un numero a scelta di clusters (ad es. $K = 5$). Discutere ottimalità e unicità della soluzione ottenuta.

Punto	x	y
1	24	9
2	16	33
3	8	32
4	42	31
5	40	45
6	41	89
7	13	71
8	37	64
9	34	66
10	50	58
11	91	43
12	68	27
13	63	29
14	61	45
15	54	50
16	62	79
17	65	75
18	80	81
19	85	67
20	51	56

Tabella 1: Posizioni dei punti.

Soluzione. I dati del problema sono il numero N di punti in posizione $(x_i, y_i) \forall i = 1, \dots, N$ ed il numero K di clusters.

Una soluzione è definita da una partizione dei punti in clusters e dalla posizione dei centroidi dei clusters. La partizione può essere rappresentata da variabili binarie di assegnamento w_{ik} per ogni punto i ed ogni cluster k . La posizione dei centroidi viene rappresentata da variabili continue e libere (\bar{x}_k, \bar{y}_k) per ogni cluster k .

Per comodità si possono definire le variabili ausiliarie implicate nella definizione dell'indice.

$$S_k = \sqrt{\frac{1}{\sum_{j=1}^N w_{jk}} \sum_{i=1}^N w_{ik} ((\bar{x}_k - x_i)^2 + (\bar{y}_k - y_i)^2)} \quad \forall k = 1, \dots, K.$$

$$M_{k'k''} = \sqrt{(\bar{x}_{k'} - \bar{x}_{k''})^2 + (\bar{y}_{k'} - \bar{y}_{k''})^2} \quad \forall k' \neq k'' = 1, \dots, K.$$

$$R_{k'k''} = \frac{S_{k'} + S_{k''}}{M_{k'k''}} \quad \forall k' \neq k'' = 1, \dots, K.$$

Infine si può introdurre una variabile continua D_k per ogni cluster ed imporre il vincolo seguente:

$$D_{k'} \geq R_{k'k''} \quad \forall k' \neq k'' = 1, \dots, K$$

per linearizzare la definizione che implica la funzione max.

Gli altri vincoli del problema sono quelli di assegnamento

$$\sum_{k=1}^K w_{ik} = 1 \quad \forall i = 1, \dots, N.$$

Per evitare che il denominatore nella definizione di S_k possa valere 0 provocando problemi numerici al solutore, può essere utile aggiungere vincoli sulla cardinalità minima dei clusters:

$$\sum_{i=1}^N w_{ik} \geq 1 \quad \forall k = 1, \dots, K.$$

L'obiettivo è

$$\text{minimize } z = \frac{1}{N} \sum_{k=1}^K D_k.$$

Il modello risultante è di programmazione non-lineare intera.

Fissate le variabili binarie, la minimizzazione di S_k in ogni cluster dà luogo ad un sotto-problema convesso, ma la presenza dei termini M inficia questa proprietà. Per studiare la convessità sarebbe necessario studiare le derivate dell'indice di Davies-Bouldin rispetto alla coordinata (non importa se \bar{x} o \bar{y} , data la simmetria) del generico centroide. L'obiettivo z è una funzione convessa delle variabili D . Ogni variabile D è una funzione convessa delle variabili R . Si tratta di determinare se ogni $R_{k'k''}$ è una funzione convessa di \bar{x}_i (o \bar{y}_i) quando $w_{ik'} = 1$ (o $w_{ik''} = 1$). Ai fini della soluzione dell'esercizio non è richiesto di addentrarsi in calcoli ulteriori.