

DIPLOMAMUNKA

**METRIKUS TEREK BEÁGYAZÁSA,
ALKALMAZÁS KLASZTEREZÉSRE**

Zsbán Ambrus

témavezető: **Lukács András**

MTA Számítástechnikai

és Automatizálási Kutatóintézet

belső témavezető: **Rónyai Lajos**

BME Matematikai Intézet, Algebra Tanszék

2008

BME

DIPLOMAMUNKA KIVONAT

Metrikus terek beágyazása, alkalmazás klaszterezésre

Zsbán Ambrus

Témavezető: Lukács András, MTA Számítástechnikai és Automatizálási Kutatóintézet.

Belső témavezető: Rónyai Lajos, BME Matematikai Intézet, Algebra Tanszék.

Dolgozatomban először rövid áttekintést nyújtok a véges metrikus beágyazások elméletéről. Ez a módszer szélsőérték-feladatokra ad közelítő megoldást olyan módon, hogy egy metrikus teret egy másik, egyszerűbb szerkezetű metrikus térbe ágyaz be. Bemutatok néhány tételt, amelyek ilyen beágyazást konstruálnak; és egy lemmát, amelynek segítségével az ilyen beágyazásokat fel lehet használni.

Részletesen ismertetem a metrikus beágyazások elméletének egy alkalmazását: egy közelítő megoldást adó algoritmust egy konkrét klaszterezési feladatra, a min-sum k -clustering problémára. Ez az algoritmus egy általános metrikus teret először egy speciális formájú véletlen fa által meghatározott térbe képezi, majd ezen a metrikán rekurzív módon számítja ki egy jó közelítést.

A klaszterezést – vagyis adatok előre nem meghatározott, közeli pontokból álló csoportokba sorolását – az adatbányászatban többek között arra használják, hogy egy nagy méretű magas dimenziós adathalmaz viselkedését az egyes klaszterek néhány tulajdonságának vizsgálatával könnyebben megérthessük.

My thesis first gives a short review of the theory of finite metric embeddings. This method gives approximating solutions to optimization problems by embedding a metric space to another, simple metric space. I describe a few theorems that construct such embeddings, and a lemma allowing us to use these embeddings.

I comprehensively present an application of this theory, namely an algorithm that approximates a concrete clustering problem, min-sum k -clustering. This algorithm works by first transforming a general metric space to a space induced by a random tree of a certain special form, then gets a good approximation on this new metric with a recursive calculation.

Clustering, which means dividing points to groups of points close to each other without groups known in advance, is used in data mining for understanding the behaviour of large high-dimension data sets by examining few properties of each cluster.

Budapest, 2008. május 21.

Tartalomjegyzék

1. Bevezetés a metrikus terek beágyazásának elméletébe	3
2. Metrikus terek közelítése véletlen fákkal	10
3. A min-sum k -clustering közelítő megoldása	22
Hivatkozások	47

1. Bevezetés a metrikus terek beágyazásának elméletébe

A véges metrikus terek beágyazásának elméletével kapcsolatos vizsgálatokat Bourgain 1985-ös eredménye indította el; azóta annyi eredmény született, hogy ezeknek csak kis részét tudom ismertetni. Ezekről az eredményekről jó áttekintést ad Indyk, Matoušek [6] fejezete egy kézikönyvből, és valamivel bővebben Matoušek könyvének [7] fejezete. Ebben a diplomamunkában ismeretlem ennek az elméletnek egy konkrét alkalmazását a min-sum k -clustering probléma közelítő megoldására.

1.1. definíció. Legyen a továbbiakban V pontoknak egy adott véges halmaza, és n a pontok száma. Egy $d : V \times V \rightarrow \mathbb{R}$ függvényt **távolságfüggvénynek** nevezzük a ponthalmazon, ha szimmetrikus, tehát ha minden $u, v \in V$ pontpárra $d(u, v) = d(v, u)$, és $d(u, u) = 0$. Ha a távolságfüggvény pozitív, vagyis minden $u, v \in V$ pontpárra ha $u \neq v$ akkor $0 < d(u, v)$; és ha a háromszögeyenlőtlenség is teljesül, azaz minden $u, v, w \in V$ pontokra $d(u, w) \leq d(u, v) + d(v, w)$; akkor a távolságfüggvény **metrika**. Ha a háromszögeyenlőtlenség teljesül, de a pozitivitás helyett csak nemnegativitást ($0 \leq d(u, v)$) kötünk ki, akkor a távolságfüggvényt **szemimetrikának** hívjuk.

1.2. definíció. Tekintsünk egy optimalizálási problémát, amely a távolságfüggvénnyel van paraméterezve. Pontosabban legyen T a **megengedett megoldások** halmaza; g pedig a **célfüggvény**, amely a távolságfüggvényből és megengedett megoldásból álló párokat képezi a valós számokba. Az alapfeladatunk, hogy adott d távolságfüggvényre keressünk meg egy $t_0 \in T$ optimális megoldást, amelyre $g(d, t_0)$ minimális. Bizonyos esetben megelégedhetünk közelítő megoldással is: adott d távolságfüggvény mellett egy t megoldásról akkor mondjuk, hogy β közelítéssel optimális, ha $g(d, t) \leq \beta \cdot g(d, t_0)$ (itt t_0 az optimális megoldás, és $1 \leq \beta$ valós paraméter; továbbá ezt a definíciót abban az esetben fogjuk használni, ha a célfüggvény értéke nemnegatív).

Az optimalizálási problémákat általában csak szemimetrikákra vagy metrikákra megszorítva fogjuk vizsgálni, de a következő definícióhoz technikailag

egyszerűbb bármely távolságfüggvényre kiterjeszteni őket, mivel az összes távolságfüggvény együtt valós vektorteret alkot.

1.3. definíció. A fenti optimalizálási problémát (vagy a célfüggvényét) akkor hívjuk a távolságfüggvényben **pozitív lineárisnak**, ha bármely rögzített $t \in T$ megoldásra a g függvény előáll $g(d, t) = \sum_{u,v \in V} c_{uv} \cdot d(u, v)$ alakban valamilyen nemnegatív c_{uv} konstansokra.

Nyilvánvaló, hogy ha a célfüggvény pozitív lineáris, akkor a távolságfüggvénytől (amelyet a vektortér elemeinek tekintünk) lineárisan függ. Fontos továbbá észrevenni, hogy ha a célfüggvény pozitív lineáris, akkor szemimetrika távolságokon az értéke mindig nemnegatív. Abból azonban, hogy a célfüggvény a távolságfüggvényben lineáris és minden szemimetrikán nemnegatív értéket vesz fel, még nem következik, hogy pozitív lineáris is: példa erre a $g(d, t) = d(u, v) + d(v, w) - d(v, w)$ függvény.

Pozitív lineáris problémára példa az utazó ügynök problémája, a legkisebb súlyú feszítőfa keresése, vagy metrikus térben az egymáshoz legközelebbi két pont megkeresése.

1.4. definíció. Legyen most d és d^* két szemimetrika ugyanazon a csúcshalmazon. Azt mondjuk, hogy a d^* szemimetrika α **torzítással** közelíti d -t, ha minden u, v pontpárra $d(u, v) \leq d^*(u, v) \leq \alpha \cdot d(u, v)$ (ahol $1 < \alpha$ valós).

A torzítás fogalmának a következő egyszerű lemma adja meg a jelentőségét. Ez a lemma lehetővé teszi, hogy egy pozitív lineáris optimalizálási problémára közelítő megoldást adjunk, ha tudunk találni a bemenő metrikához egy másik metrikát, amely az eredetit kis torzítással közelíti, és amely elég speciális formájú ahhoz, hogy az optimalizálási problémát könnyebb megoldani rajta.

1.5. lemma. Tegyük fel, hogy a d^* szemimetrika α torzítással közelíti a d szemimetrikát, és a g célfüggvény pozitív lineáris. Ha t_1 optimális megoldása a d^* szemimetrika mellett problémának, azaz minimalizálja a $g(d^*, t_1)$ függvényt, akkor α közelítéssel optimalizálja a problémát az eredeti d szemimetrikával, azaz $g(d, t_1) \leq \alpha \cdot g(d, t_0)$ az optimális t_0 megoldásra. Hasonlóan,

ha a t megoldás β közelítéssel optimális a d^* szemimetrikával, akkor $\alpha\beta$ közelítéssel optimális a d szemimetrikával.

Bizonyítás. Vegyük észre, hogy a lemma feltételei mellett bármely $t \in T$ megengedett megoldásra $g(d, t) \leq g(d^*, t) \leq \alpha \cdot g(d, t)$. Valóban, ha a t megoldást rögzítve a g függvényt felírjuk az 1.3. definíciónak megfelelő alakban, akkor, mivel $0 \leq c_{uv}$,

$$g(d, t) = \sum_{u,v} c_{uv} \cdot d(u, v) \leq \sum_{u,v} c_{uv} \cdot d^*(u, v) = g(d^*, t),$$

$$g(d^*, t) = \sum_{u,v} c_{uv} \cdot d^*(u, v) \leq \sum_{u,v} c_{uv} \cdot \alpha \cdot d(u, v) = \alpha \cdot g(d, t).$$

Most ha t_1 az optimális megoldás a d^* szemimetrika mellett, akkor bármely $t \in T$ megoldásra $g(d^*, t_1) \leq g(d^*, t)$. Speciálisan ha t_0 az optimális megoldás a d szemimetrikával,

$$g(d, t_1) \leq g(d^*, t_1) \leq g(d^*, t_0) \leq \alpha \cdot g(d, t_0),$$

azaz t_1 valóban α -közelítés az eredeti problémára.

Ha most t a d^* szemimetrikával β faktor erejéig optimális, akkor

$$g(d, t) \leq g(d^*, t) \leq \beta \cdot g(d^*, t_1) \leq \beta \cdot g(d^*, t_0) \leq \alpha\beta \cdot g(d, t_0),$$

tehát t valóban $\alpha\beta$ -közelítése az optimumnak.

Metrikus beágyazásnak azt az általános módszert hívjuk, amikor egy metrikus teret egy másik hasonló térrel helyettesítünk. A második metrikus tér lehet egyszerűbb szerkezetű az elsőnél, rövidebben leadható, vagy akár csak a megoldandó feladat szempontjából valamilyen okból kényelmesebb. Látuk, hogy egy pozitív lineáris feladat közelítő megoldásában segíthet, ha a metrikus teret alacsony torzítással tudjuk beágyazni.

Nézzünk most néhány példát olyan ismertebb tételekre, amelyekkel egy valamilyen értelemben jól közelítő metrikus teret megkonstruálhatunk.

1.6. definíció. Legyen k pozitív egész, $1 \leq p$ valós. Az k hosszúságú valós vektorokból álló teret a

$$d((u_0, \dots, u_{k-1}), (v_0, \dots, v_{k-1})) = ((y_0 - x_0)^p + \dots + (y_{k-1} - x_{k-1})^p)^{1/p}$$

metrikával l_k^p jelöli; ugyanezt a vektorteret a

$$d((u_0, \dots, u_{k-1}), (v_0, \dots, v_{k-1})) = \max(|y_0 - x_0|, \dots, |y_{k-1} - x_{k-1}|)$$

metrikával l_k^∞ jelöli.

1.7. állítás. Bármely n elemű metrikus teret izometrikusan be lehet ágyazni az l_{n-1}^∞ térbe, azaz bármely n elemű vektortérhez van az l_{n-1}^∞ térnek olyan részhalmaza, amelyre megszorított metrikus tér 1 torzítással közelíti az eredeti metrikus teret.

Ezt az állítást nagyon egyszerű belátni, de inkább csak elméleti jelentősége van, mivel a kapott metrikus teret általában nem könnyebb kezelni az eredeténél.

Jegyezzük meg rögtön, hogy nem lehet bármely véges metrikus teret egy l^p térbe izometrikusan beágyazni. Például azt a metrikus teret, amely az A, B, C, D pontokból áll, és

$$\begin{aligned} 1 &= d(A, B) = d(B, C) = d(C, D) = d(A, D) \\ 2 &= d(A, C) = d(B, D), \end{aligned}$$

nem lehet l^2 térbe beágyazni, bármi is legyen a dimenzió. Hasonló példákat lehet készíteni bármilyen p -re.

Jíří Matousek 1996-os eredménye a következő tétel.

1.8. tétel. Legyen $D = 2b - 1$, ahol $1 < b$ egész. Tetszőleges n elemű metrikus tér beágyazható D torzítással az l_k^∞ térbe, ahol a dimenzió

$$k = O(bn^{1/b} \log n)$$

Bourgain 1985-ös beágyazási tétele a következő.

1.9. tétel. Bármely n pontból álló metrikus tér $O(\log n)$ torzítással beágyazható egy l^2 térbe.

Adatbányászati feladatokban előfordul, hogy a metrikus tér egy magas dimenziós l^p térrel van megadva. Ilyen magas dimenziós nagy adathalmazokon különféle feladatokat akarhatunk elvégezni, mint például klaszterezés, klasszifikáció, outlierok keresése, statisztikai modellezés stb. Az algoritmusok szempontjából sok esetben segíthet, ha alacsonyabb dimenziós térrel dolgozunk. Ebben segíthet a következő Johnson–Lindenstrauss dimenziócsökkentési lemma.

1.10. tétel. Egy l^2 tér bármely n pontból álló részét $1 + \varepsilon$ torzítással be lehet ágyazni a $k = O(\varepsilon^{-2} \log n)$ dimenziós l_2^k térbe bármely adott $0 < \varepsilon \leq 1$ esetén.

A három előző tétel bizonyítását az olvasó megtalálhatja [7]-ben. Ezek a bizonyítások tanulságos példát adnak a valószínűségi módszerek alkalmazására, ami egyébként a metrikus beágyazások konstruálására tipikus módszer.

A fentiekén kívül nagyon sok beágyazási tétel létezik, lehet például foglalkozni különféle tulajdonságú gráfok által indukált metrikákkal, vagy nullákat és egyeseket tartalmazó vektorok terén Hamming- vagy Jaccard-távolsággal. Vannak olyan tételek is, amelyek a torzítás helyett más mérőszámmal jellemzik egy beágyazás jóságát.

Most rátérünk tulajdonképpen célunkra, a min-sum k -clustering probléma vizsgálatára.

1.11. definíció. Legyen az előzőekhez hasonlóan rögzítve az n elemű V csúcshalmaz, és bemenetnek adott a d metrika ezen a halmazon. Legyen továbbá rögzítve egy k pozitív egész paraméter. A **min-sum k -clustering probléma** feladata megkeresni V egy olyan partícióját (legfeljebb) k részre, ami az egy partíción belüli pontpárok távolságának összegét minimalizálja. Megengedett megoldás tehát a csúcshalmaz tetszőleges $\{P_0, \dots, P_{N-1}\}$ diszjunkt partíciója (bármely N -re; $V = P_0 \dot{\cup} \dots \dot{\cup} P_{N-1}$). A célfüggvény értéke

$$g(d, \{P_i\}) = \frac{1}{2} \sum_{0 \leq i < N} \sum_{u \in P_i} \sum_{v \in P_i} d(u, v).$$

Bartal, Charikar, Raz [3] cikkében egy új közelítő algoritmust ad a problémára. Ez a metrikus beágyazás módszerét alkalmazza, vagyis a metrikát egy speciális formájú (egyszerűbben kezelhető) véletlen metrikával helyettesíti, majd ezen ad meg egy közelítő megoldást. Ennek a bizonyítást fogom jelen diplomamunkában részletesen ismertetni.

A min-sum k -clustering problémát korábban is vizsgálták. A fenti eredmény azért újszerű, mert a legtöbb másik algoritmussal ellentétben k -tól függetlenül polinomiális idejű. Csak rögzített k mellett polinomiális konstans faktossal közelítő algoritmust mutat be [8].

Dolgozatunk fő tétele tehát a következő.

1.12. tétel. Adható olyan polinomiális időkorlátú algoritmus, amelynek bemenete tetszőleges – páronkénti távolságokkal megadott – véges metrikus tér, kimenete pedig a min-sum k -clustering problémának ezen a téren az optimálisnál várhatóan legfeljebb $O(\log^2 n)$ -szer rosszabb megoldása.

Először is szögezzük le a következőt.

1.13. állítás. A min-sum k -clustering probléma pozitív lineáris.

Bizonyítás. A célfüggvény fenti felírásából nyilvánvaló, ha a szummákat felcseréljük:

$$c_{uv} = \frac{1}{2}d(u, v) \quad \text{ha } u, v \in P_i \text{ valamely } i\text{-re,}$$

$$c_{uv} = 0 \quad \text{különben.}$$

Nem az 1.5. lemmát fogjuk használni azonban, hanem egy valamivel bonyolultabb változatot.

1.14. definíció. Legyen d egy szemimetrika, és d^* egy szemimetrika értékű valószínűségi változó (mindkettő a V ponthalmazon). Azt mondjuk, hogy d^* egy $1 < \alpha$ paraméterrel valószínűségben közelíti a d szemimetrikát, ha minden $u, v \in V$ -re biztosan $d(u, v) \leq d^*(u, v)$ (azaz d_N dominálja d_M -et), de minden $u, v \in V$ -re $d^*(u, v)$ várható értéke legfeljebb $\alpha \cdot d(u, v)$.

1.15. lemma. Tegyük fel, hogy a d^* véletlen szemimetrika α paraméterrel valószínűségben közelíti a d szemimetrikát, és a g célfüggvény pozitív lineáris. Ha t_1 optimális megoldása a d^* -ra vonatkozó problémának, azaz bármely d^* kimenet esetén minimalizálja a $g(d^*, t_1)$ függvényt, akkor $\mathbf{E}(g(d, t_1)) \leq \alpha \cdot g(d, t_0)$ az optimális t_0 megoldásra. Hasonlóan, ha a t megoldás mindig β közelítéssel optimális megoldása a d^* szemimetrikához tartozó feladatnak, akkor $\mathbf{E}(g(d, t_1)) \leq \alpha\beta \cdot g(d, t_0)$.

Megjegyzem, hogy a Markov-egyenlőtlenség miatt ebből az is következik, hogy ha $\alpha\beta < \gamma$, akkor annak, hogy a kapott t megoldás $g(d, t)$ célfüggvénye nagyobb az optimális $g(d, t_0)$ célfüggvény γ -szorosánál, legfeljebb $(\alpha\beta - 1)/(\gamma - 1)$ lehet a valószínűsége.

Bizonyítás. Hasonlóan járunk el, mint az 1.5. lemmánál.

Először állítom, hogy bármely $t \in T$ megengedett megoldásra $g(d, t) \leq g(d^*, t)$ biztosan teljesül, és $\mathbf{E}(g(d^*, t)) \leq \alpha \cdot g(d, t)$. Valóban, hiszen

$$g(d, t) = \sum_{u,v} c_{uv} \cdot d(u, v) \leq \sum_{u,v} c_{uv} \cdot d^*(u, v) = g(d^*, t);$$

$$\mathbf{E}(g(d^*, t)) = \sum_{u,v} c_{uv} \cdot \mathbf{E}(d^*(u, v)) \leq \sum_{u,v} c_{uv} \cdot \alpha \cdot d(u, v) = \alpha \cdot g(d, t).$$

Ebből pedig egyrészt ha t_1 optimális a d^* szemimetrikával, akkor

$$g(d, t_1) \leq g(d^*, t_1) \leq g(d^*, t_0)$$

$$\mathbf{E}(g(d^*, t_0)) \leq \alpha \cdot g(d, t_0);$$

így tehát t_1 valóban átlagosan legfeljebb α -szor rosszabb az optimális megoldásnál:

$$\mathbf{E}(g(d, t_1)) \leq \alpha \cdot g(d, t_0).$$

Ha pedig csak annyit tudunk, hogy t a módosított $g(d^*, t)$ problémának β -közelítése, akkor

$$g(d, t) \leq g(d^*, t) \leq \beta \cdot g(d^*, t_1) \leq \beta \cdot g(d^*, t_0)$$

$$\mathbf{E}(g(d^*, t_0)) \leq \alpha \cdot g(d, t_0);$$

így tehát t az eredeti problémára átlagosan legfeljebb $\alpha\beta$ -szor rosszabb az optimálisnál:

$$\mathbf{E}(g(d, t)) \leq \alpha\beta \cdot g(d, t_0).$$

A második fejezetben definiálunk egy metrika-osztályt, a μ paraméterű hierarchikusan elválasztott fák által indukált metrikák osztályát, és bebizonyítjuk, hogy bármely metrikát $O(\mu(2 \ln n + 2)(1 + \log_\mu n))$ torzítással közelíteni lehet ilyen metrikával. A harmadik fejezetben megmutatjuk, hogyan lehet ezen metrikákon hatékonyan konstans közelítéssel megoldani a min-sum k -clustering problémát.

2. Metrikus terek közelítése véletlen fákkal

Vegyünk egy tetszőleges összefüggő gráfot, aminek minden éléhez egy pozitív élsúly van rendelve, röviden egy súlyozott élű gráfot. Egy ilyen súlyozott élű gráfból kaphatunk egy metrikát a gráf csúcshalmazán vagy ennek egy részén: távolságnak a két csúc közötti legrövidebb út hosszát választjuk. Ez a konstrukció még nagyon általános: minden metrikát meg lehet kapni súlyozott élű gráfból. Értelmes fogalmakat kaphatunk azonban, ha a gráfra megszorításokat teszünk.

Ha speciálisan a gráf egy fa (összefüggő körmentes gráf), akkor bármely két levél között pontosan egy út van, így a definíciót fogalmazzuk egyszerűen így.

2.1. definíció. Egy súlyozott élű fa leveleinek V halmazára felett **indukálja** azt a metrikus teret, amiben a távolság a két levél közötti út élsúlyainak összege.

Az általunk vizsgált egyszerűbb metrikus tereket az úgynevezett hierarchikusan elválasztott fákból (HST) kapjuk ilyen módon.

2.2. definíció. Egy súlyozott élű gyökeres fát μ -HST-nek hívunk az $1 < \mu$ paraméterrel és Δ átmérővel, ha minden levele azonos $\Delta/2$ távolságra van a gyökértől, és minden gyereke vagy levél, vagy szintén μ -HST pontosan Δ/μ átmérővel.

A definícióból következik, hogy egy Δ átmérőjű μ -HST-ben az i -edik szinten lévő csúcstól az összes levél leszármazottja pontosan $\mu^{-i}\Delta/2$ távolságra van.

Vegyük észre, hogy az itt használt Δ átmérő nagyobb lehet a kapott metrikus tér átmérőjénél (a maximális távolságnál) ha a gyökérnek csak egy gyereke van, de kisebb sosem lehet annál. A gyökeret nem tekintjük levélnek, még akkor sem, ha csak egy él indul ki belőle. (Megjegyzem, hogy némely cikk ennél általánosabban definiálja a HST-t.)

2.3. definíció. Egy HST-t speciális formájúnak hívunk, ha minden levél egy szinten van, azaz mindegyiket ugyanannyi él választja el a gyökértől.

A levelekbe menő élek felosztásával könnyen beláthatjuk a következő állítást.

2.4. állítás. Bármely HST-t speciális formára hozhatunk, és ez a művelet megtartja a Δ és μ értékeit, a levelek halmazát, és az ezen indukált gráf-metrikát.

Bartal [1] megmutatja, hogy bármely (véges) metrikus teret bármely adott μ -re valószínűségben közelíteni lehet egy, a tér pontjaira kifeszített μ -HST-vel, $\alpha = O(\mu \log n / \log_{\mu} n)$ paraméterrel. Ezt a módszert fogom ebben a fejezetben kifejteni.

A közelítés paraméterét $\alpha = O(\mu \log n \log \log n)$ -re javítja [2]. A módszert derandomizálni is lehet: [4] egy determinisztikus algoritmust ad, amellyel elég $O(n \log n)$ fát megvizsgálni. Rosszabb korlátú, de egyszerűbb bizonyítást ad meg [5].

E fejezet fő tétele tehát a következő.

2.5. tétel. Legyen adott egy d véges metrikus tér a V ponthalmazon. Hatékony véletlen algoritmussal konstruálhatunk egy véletlen speciális formájú μ -HST-t úgy, hogy a HST által a leveleken indukált d^* véletlen gráf-metrika valószínűségben α paraméterrel közelíti d -t, ahol

$$\alpha = \mu(2 \ln n + 2)(1 + \log_{\mu} n)$$

Itt $1 < \mu$ előre rögzített paraméter.

Először kimondunk egy lemmát, amely a bizonyítás egyik lépését technikailag segíti.

2.6. lemma. Legyenek F_0, \dots, F_{N-1} illetve E_0, \dots, E_{N-1} olyan véletlen események, hogy bármely $0 \leq k < N$ indexre F_k és E_k kizárja egymást. Tegyük továbbá fel, hogy van egy olyan $0 \leq U \leq N - 2$ index, amelyre F_U biztosan nem következik be.

Definiáljuk az m valószínűségi változót, mint a legkisebb olyan $0 \leq m < N$ indexet, amelyre F_m hamis (az előbbi feltevés miatt mindig van ilyen index). A G eseményt definiáljuk, mint E_m .

Tegyük fel, hogy minden $0 \leq k < N$ indexre van olyan ξ_k szám, hogy a következő három egyenlőtlenség teljesül:

$$1/n \leq \xi_k \quad (1)$$

$$\mathbf{P}(E_k \wedge F_0 \wedge \dots \wedge F_{k-1}) \leq \mathbf{P}(F_0 \wedge \dots \wedge F_{k-1}) \cdot \frac{n}{n-1} \cdot p \cdot \xi_k \quad (2)$$

$$\mathbf{P}(F_k \wedge F_0 \wedge \dots \wedge F_{k-1}) \leq \mathbf{P}(F_0 \wedge \dots \wedge F_{k-1}) \cdot \frac{n}{n-1} \cdot (1 - \xi_k). \quad (3)$$

Állítjuk, hogy ekkor $\mathbf{P}(G) \leq 2p$.

Bizonyítás. Teljes indukcióval látjuk be a következő egyenlőtlenséget

$$\mathbf{P}(G \mid F_0 \wedge \dots \wedge F_{k-1}) \leq \frac{2n-2-k}{n-1} \cdot p$$

egyenlőtlenséget. Pontosabban a teljes indukció k szerint visszafele halad az $0 \leq k \leq U + 1$ intervallumon, és azt állítjuk, hogy

$$\mathbf{P}(G \wedge F_0 \wedge \dots \wedge F_{k-1}) \leq \frac{2n-2-k}{n-1} \cdot p \cdot \mathbf{P}(F_0 \wedge \dots \wedge F_{k-1}). \quad (4)$$

A $k = U + 1$ alapesetben az $F_0 \wedge \dots \wedge F_i$ feltétel triviálisan sosem teljesül, így az egyenlőtlenség mindkét oldala nulla. Másrészt az $k = 0$ esetben a feltétel triviálisan mindig igaz, így az egyenlőtlenség az $\mathbf{P}(G) \leq 2p$ alakra egyszerűsödik, és ez éppen az, amit igazolni akarunk.

Mármost lássuk be az indukciós állítást k -ra, ha feltesszük $k - 1$ -re. Itt $k \leq U - 1 \leq N - 2$. Tudjuk tehát, hogy

$$\mathbf{P}(G \wedge F_0 \wedge \cdots \wedge F_{k-1} \wedge F_k) \leq \frac{2n - 3 - k}{n - 1} \cdot p \cdot \mathbf{P}(F_0 \wedge \cdots \wedge F_{k-1} \wedge F_k). \quad (5)$$

Mármost

$$\begin{aligned} \mathbf{P}(G \wedge F_0 \wedge \cdots \wedge F_{k-1}) &= \\ &= \mathbf{P}(G \wedge F_0 \wedge \cdots \wedge F_{k-1} \wedge F_k) + \mathbf{P}(G \wedge F_0 \wedge \cdots \wedge F_{k-1} \wedge \neg F_k) \end{aligned} \quad (6)$$

Csak hogy ha $F_0 \wedge \cdots \wedge F_{k-1} \wedge \neg F_k$ igaz, akkor $m = k$, és így G ekvivalens E_k -vel. Ezért aztán a második tagot átírhatjuk a következőképpen, és használhatjuk a (3) feltételt:

$$\mathbf{P}(G \wedge F_0 \wedge \cdots \wedge F_{k-1} \wedge \neg F_k) = \mathbf{P}(F_0 \wedge \cdots \wedge F_{k-1} \wedge E_k) \quad (7)$$

Helyettesítsük be a (6) második tagjába a (7) kifejezést, az első tagjába pedig az (5) indukciós feltételt.

$$\begin{aligned} \mathbf{P}(G \wedge F_0 \wedge \cdots \wedge F_{k-1}) &\leq \\ &\leq \frac{2n - 3 - k}{n - 1} \cdot p \cdot \mathbf{P}(F_0 \wedge \cdots \wedge F_{k-1} \wedge F_k) + \mathbf{P}(F_0 \wedge \cdots \wedge F_{k-1} \wedge E_k) \end{aligned}$$

Összevetve ezt az a (3) és a (2) feltételekkel, ezt kapjuk:

$$\begin{aligned} \mathbf{P}(G \wedge F_0 \wedge \cdots \wedge F_{k-1}) &\leq \\ &\leq \mathbf{P}(F_0 \wedge \cdots \wedge F_{k-1}) \cdot \left(\frac{2n - 3 - k}{n - 1} \cdot p \cdot \frac{n}{n - 1} \cdot (1 - \xi_k) + \frac{n}{n - 1} \cdot p \cdot \xi_k \right) \end{aligned}$$

Viszont ebben az együttható

$$\begin{aligned}
& \frac{2n-3-k}{n-1} \cdot p \cdot \frac{n}{n-1} \cdot (1-\xi_k) + \frac{n}{n-1} \cdot p \cdot \xi_k = \\
& = \left(\frac{2n^2-3n-nk}{(n-1)^2} + \frac{-n^2-2n+nk}{(n-1)^2} \xi_k \right) \cdot p \leq \\
& \leq \left(\frac{2n^2-3n-nk}{(n-1)^2} + \frac{n(-n-2+k)}{(n-1)^2} \cdot \frac{1}{n} \right) \cdot p = \\
& = \frac{2n^2-4n-nk-2+k}{(n-1)^2} \cdot p = \frac{2n-2-k}{n-1} \cdot p
\end{aligned}$$

(Felhasználtuk, hogy $k \leq n-2$, és az (1) feltételt.) Ez pedig az indukció a (4) állítása, amit bizonyítani akartunk.

Most a tényleges bizonyítás első lépéseként megadunk egy módszert a csúcshalmaz bizonyos tulajdonságú véletlen partícionálására. (Noha az elnevezés hasonló, ne tévessze össze ezt a partícionálást a min-sum clustering probléma megoldásával: ez utóbbi partícionálásról ebben a szakaszban nem ejtünk szót.) Pontosan a következőt állítjuk.

2.7. lemma. A V_0, \dots, V_{N-1} valószínűségi változókat megválaszthatjuk úgy, hogy a V csúcshalmaz ezen halmazok (klaszterek) $V_0 \dot{\cup} \dots \dot{\cup} V_{N-1}$ diszjunkt uniójaként áll elő, és a következő három tulajdonság teljesül. Egyrészt minden V_k klaszter átmérője legfeljebb $2r \ln n$; másrészt ha $d(u, w) = 0$, akkor u és w biztosan egy klaszterben van; harmadrészt bármely $u, w \in V$ csúcsokra annak a valószínűsége, hogy u és w külön klaszterben van, legfeljebb $2p$, ahol $p = d(u, w)/r$. (Itt $0 < r$ egy tetszőleges paraméter; d pedig adott szemimetrika.)

Bizonyítás. A bizonyításhoz megkonstruálunk egy megfelelő (V_0, \dots, V_{N-1}) valószínűségi változót. Feltehetjük, hogy $2 \leq n$, különben az állítás triviális.

2.8. definíció. Először rögzítsük a csúcsok egy tetszőleges sorrendjét: $V = \{v_0, \dots, v_{n-1}\}$. Minden $0 \leq k < n$ indexhez válasszunk egy z_k véletlenszámot, amelyre $0 \leq z_k < r \ln n$, az eloszlása

$$\mathbf{P}(z_k \leq t) = \frac{n}{n-1} (1 - e^{-t/r}),$$

ha $0 \leq t \leq r \ln n$, és a z_0, \dots, z_{n-1} számok teljesen függetlenek. Ezután legyen $N = n$, és legyen

$$V_k = \{u \in V \mid d(u, v_k) \leq z_k \text{ és } u \notin V_0 \cup \dots \cup V_{k-1}\}.$$

Másképpen a k sorszámú klaszter a v_k körüli z_k sugarú gömbben lévő csúcsokat tartalmazza, kivéve azokat, amelyek valamely korábbi klaszterben vannak.

Most belátjuk, hogy az így definiált véletlen partícióval a lemma állításai igazak.

Előrebocsájtjuk, hogy a z_k -ra megadott eloszlás folytonos, és hogy eloszlásfüggvényét $0 \leq t$ esetén felírhatjuk a következő egységes alakban:

$$\mathbf{P}(z_k \leq t) = \frac{n}{n-1} (1 - e^{-\min(t, r \ln n)/r}). \quad (8)$$

Mivel $d(v_k, v_k) \leq z_k$ mindenképp igaz minden k -ra, ezért a partíció valóban lefedi az egész ponthalmazt: $V = V_0 \cup \dots \cup V_{n-1}$. Az is látható, hogy a klaszterek diszjunktak. Bármely klaszter sugara legfeljebb $r \ln n$, hiszen ha $u \in V_k$ akkor $d(u, v_k) \leq z_k < r \ln n$, ezért az átmérője is legfeljebb kétszer ekkora: ha $u, w \in V_k$ akkor $d(u, w) \leq d(u, v_k) + d(v_k, w) \leq 2r \ln n$. A második tulajdonság is nyilvánvalóan igaz: ha $d(u, v) = 0$ akkor u és v ugyanabba a klaszterbe kerül.

A harmadik tulajdonság igazolásához legyen $u, w \in V$ két rögzített csúcs. A 2.6. lemmát fogjuk felhasználni. Jelölje F_k azt az eseményt, hogy $z_k < d(u, v_k)$ és $z_k < d(v_k, w)$ is teljesül; E_k pedig azt jelöli, hogy $z_k < d(u, v_k)$ vagy $z_k < d(v_k, w)$ közül pontosan az egyik teljesül.

Az nyilvánvaló, hogy F_k és E_k kizárja egymást. Mivel $u, w \in V$, vannak olyan $U_0, U_1 < N$ indexek, hogy $u = v_{U_0}$ illetve $w = v_{U_1}$, ezekre az indexekre pedig $d(u, v_{U_0}) = 0 \leq z_k$ illetve $d(v_{U_1}, w) = 0 \leq z_k$ miatt F_k nem teljesül. Ráadásul $U_0 \neq U_1$, így az $U = \min(U_0, U_1) \leq N - 2$ indexre is igaz mindez.

Vegyük észre, hogy mivel F_k és E_k függvénye z_k -nak, és z_0, \dots, z_{N-1} függetlenek, ezért az

$$(F_0, E_0), \dots, (F_k, E_k)$$

párok teljesen függetlenek egymástól. Ezért aztán a (2), (3) feltételeket egyszerűsíthetjük $\mathbf{P}(F_0 \wedge \dots \wedge F_{k-1})$ -gyel. Ahhoz, hogy a lemmát alkalmazhassuk, azt kell tehát igazolni, hogy minden k indexre van olyan ξ_k szám, hogy

$$\begin{aligned} 1/n &\leq \xi_k \\ \mathbf{P}(E_k) &\leq \frac{n}{n-1} \cdot p \cdot \xi_k \\ \mathbf{P}(F_k) &\leq \frac{n}{n-1} \cdot (1 - \xi_k) \end{aligned}$$

Vezessük be a következő jelöléseket:

$$\begin{aligned} \varrho_0 &= \min(d(u, v_k), d(v_k, w)); & \varrho_0^* &= \min(\varrho_0, r \ln n); \\ \varrho_1 &= \max(d(u, v_k), d(v_k, w)); & \varrho_1^* &= \min(\varrho_1, r \ln n). \end{aligned}$$

Ezzel a jelöléssel az E_k eseményt úgy is átfogalmazhatjuk, hogy $\varrho_0 \leq z_k < \varrho_1$. Használjuk fel a (8) eloszlásfüggvényt (kihasználva, hogy folytonos).

$$\begin{aligned} \mathbf{P}(E_k) &= \mathbf{P}(\varrho_0^* \leq z_k < \varrho_1^*) = \frac{n}{n-1} (e^{-\varrho_0^*/r} - e^{-\varrho_1^*/r}) = \\ &= \frac{n}{n-1} e^{-\varrho_0^*/r} (1 - e^{(\varrho_0^* - \varrho_1^*)/r}) \quad (9) \end{aligned}$$

Teljesen hasonlóan F_k akkor és csak akkor teljesül, ha $z_k < \varrho_0$. Tehát

$$\mathbf{P}(F_k) = \mathbf{P}(z_k < \varrho_0) = \frac{n}{n-1} (1 - e^{-\varrho_0^*/r}) \quad (10)$$

Vegyük most észre, hogy

$$\varrho_1^* - \varrho_0^* \leq \varrho_1 - \varrho_0 = |d(u, v_k) - d(v_k, w)| \leq d(u, w),$$

amiből, mivel $1 + x \leq e^x$ minden x valós számra igaz,

$$1 - e^{(\varrho_1^* - \varrho_0^*)/r} \leq (\varrho_1^* - \varrho_0^*)/r \leq d(u, w)/r = p$$

Másrészt $\varrho_0^* \leq r \ln n$ miatt az is igaz, hogy

$$1/n \leq e^{-\varrho_0^*/r}$$

Így aztán $\xi_k = e^{-\varrho_0^*}$ választással (9) és (10) adja a kívánt feltételeket. Valóban alkalmazhatjuk tehát a 2.6. lemmát. Vizsgáljuk meg, mit ad az állítása.

A 2.8. definícióból következik, hogy az u pont abban a V_{k_0} klaszterben van benne, amelyre $d(u, v_{k_0}) \leq z_k$, de minden $k < k_0$ indexre $z < d(u, v_k)$. Ugyanígy, w abban a V_{k_1} klaszterben van, amelyre $d(v_{k_1}, w) \leq z_k$, de minden $k < k_1$ indexre $z < d(v_k, w)$. Viszont az m valószínűségi változót úgy definiáltuk, mint azt az indexet, amelyre F_m hamis, de minden $k < m$ indexre F_m igaz. Ezért aztán F_k definíciójából minden $k < m$ indexre $z_k < d(u, v_k)$ és $z_k < d(v_k, w)$; viszont $d(u, v_m) \leq z_k$ vagy $d(v_m, w) \leq z_k$.

Együtt ez azt jelenti, hogy u és v közül legalább az egyik a V_m klaszterben van. Ráadásul látható, hogy G azaz E_m akkor és csak akkor igaz, ha nem mindkét pont van ebben a klaszterben. Azt kaptuk, hogy G akkor és csak akkor következik be, ha u és v különböző klaszterben van. A lemma azt állítja, hogy ennek a valószínűsége legfeljebb $2p$. Ezzel a 2.7. lemmát beláttuk.

Következő lépésként módosítjuk a fenti partícionálási módszert, hogy egy újabb követelményt is kielégítsen.

2.9. lemma. A V_0, \dots, V_{N-1} valószínűségi változókat megválaszthatjuk úgy, hogy a V csúcshalmaz ezen halmazok (klaszterek) $V_0 \dot{\cup} \dots \dot{\cup} V_{N-1}$ diszjunkt uniójaként áll elő, és a következő három tulajdonság teljesül. Egyrészt minden V_k klaszter átmérője legfeljebb $r(2 \ln n + 2)$; másrészt ha $d(u, w) \leq r/n$, akkor u és v biztosan egy klaszterben van; harmadrészt bármely $u, w \in V$ csúcsokra annak a valószínűsége, hogy u és w külön klaszterben van, legfeljebb $2p$, ahol $p = d(u, w)/r$ ($0 < r$ egy tetszőleges paraméter).

(Vegyük észre, hogy az átmérőre adott korlát kissé gyengébb az előzőnél, míg a harmadik tulajdonság változatlan.)

Bizonyítás. Először képezzünk a metrikus térből egy új szemimetrikát úgy, hogy összehúzzuk az r/n -nél közelebbi csúcsokat egy csúcsba. Pontosabban tekintsük azt az ekvivalenciarelációt, amely szerint az u és v csúcs ekvivalens egymással, ha $d(u, v) \leq r/n$, vagy ha más V -beli csúcsokon keresztül u -ból el lehet jutni v -be legfeljebb r/n hosszú ugrások sorozatával. Definiáljunk egy új d_C szemimetrikus teret a V halmazon, amelyben a $d_C(u, v)$ távolság egyenlő az u ekvivalenciaosztálya és v ekvivalenciaosztálya közti legkisebb távolsággal.

Alkalmazzuk a 2.7. lemmát az új d_C szemimetrikára (de változatlan r paraméterrel), így V egy (V_0, \dots, V_{N-1}) partícióját kapjuk. Állítom, hogy erre a partícióra teljesülnek a lemma követelményei.

Mivel összesen n csúcs van, ezért bármely pontból bármely vele ekvivalens pontba legfeljebb $n - 1$ darab, egyenként r/n hosszú ugrással el lehet jutni, így két ekvivalens pont távolsága kisebb r -nél. Ebből az következik, hogy bármely két csúcsra $d_C(u, v) \leq d(u, v) + 2r$, így az egyes klaszterek átmérője a d metrika szerint legfeljebb $2r$ -rel nagyobb, mint a d_C szerint, azaz legfeljebb $r(2 \ln n + 2)$. Ha az u, v pontokra $d(u, v) \leq r/n$, akkor ez a két pont ekvivalens, ezért $d_C = 0$, így ez a két pont biztosan egy klaszterbe kerül. Végül mivel bármely két u, v csúcsra $d_C(u, v) \leq d(u, v)$, ezért a harmadik tulajdonság is teljesül: u és v legfeljebb $d(u, v)/r$ valószínűséggel kerül külön klaszterbe.

Valóban a fenti tulajdonságú véletlen partícionálást kaptunk tehát.

Most már készen állunk arra, hogy megadjuk a μ -HST-t előállító algoritmust.

2.10. definíció. Legyen bemenetként adott a V (nem üres) csúcshalmaz, a d metrika ezen a halmazon, az $1 < \mu$ paraméter, és a Δ átmérő, amelyről feltesszük, hogy nem kisebb a V átmérőjénél a d metrika szerint. A bemenethez tartozó **közelítő fát**, amely egy súlyozott élű gyökeres véletlen fa V levélhalmazzal, rekurzívan definiáljuk a következőképpen.

Az alapeset, ha $1 = n$ (n jelölje V elemszámát): ekkor a közelítő fa álljon egy levélből és egy gyökérből köztük egy $\Delta/2$ súlyú éllel.

Különben partíciónáljuk a V csúcshalmazt a 2.9. lemma szerint a d metrika alapján az $r = \Delta/(\mu(2 \ln n + 2))$ paraméterrel. Dobjuk el az üres klasztereket. Ezután minden V_k klaszterhez rekurzív módon készítsük el a közelítő fát, metrikának a d metrika megszorítását használva, a μ paraméterrel és a Δ/μ átmérővel. (Az egyes véletlen közelítő fákat a partíciónálástól függetlenül készítsük el.) Végül vegyünk fel egy új gyökércsúcsot, amelynek gyerekei az így kapott kisebb közelítő fákat, melyek bármelyikének a gyökerét $(1 - 1/\mu)\Delta/2$ súlyú él köti össze az új gyökérrel.

2.11. lemma. Ha a V, d, μ, Δ bemenetek teljesítik a feltételeket, akkor a közelítő fa rekurzív definíciója értelmes, kimenete egy Δ átmérőjű μ paraméterű HST, melynek levélhalmaza V , és ez a fa olyan d^* metrikát indukál, amely az 1.14. definíció értelmében

$$\alpha = \mu(2 \ln n + 2)(1 + \log_\mu n)$$

paraméterrel közelíti d -t.

Bizonyítás. Ahhoz, hogy a definíció értelmes, csak két dolgot kell belátunk: azt, hogy a rekurzív hívás bemenetei teljesítik a kikötéseket, és azt, hogy az eljárás terminál, tehát nem hívja magát közvetetten végtelen sokszor.

A rekurzív hívás argumentumairól rögtön adódik, hogy V_k nem üres halmaz, d megszorítása ugyanerre pedig metrika rajta, és hogy $1 < \mu$. Azt kell csak belátni, hogy a Δ/μ átmérő nem kisebb a V_k átmérőjénél, ez is igaz azonban, mert a 2.9. lemma garantálja, hogy a kapott klaszterek átmérője legfeljebb $2(r \ln n + 2) = \Delta/\mu$.

Az l mélységű rekurzív hívásnál az átmérő bemenet egyenlő Δ/μ^l -nel így a csúcshalmaz átmérője is legfeljebb ennyi. A használt metrika mindig ugyanannak a d -nek a megszorítása, így elég nagy l -re ez a korlát kisebb a legkisebb d -beli távolságnál is, ekkor pedig a csúcshalmaz csak egy csúcsból állhat, ilyenkor pedig az alapeset érvényesül. Ebből aztán látható, hogy a rekurzív terminál.

A konstrukció és a 2.2. definíció összevetéséből nyilvánvaló, hogy a kimenet valóban egy HST a μ paraméterrel, Δ átmérővel, és V levélhalmazzal.

Vizsgáljuk most a valószínűségben közelítés érvényét. Legyen $u, v \in V$ két csúcs, jelölje a fában mért távolságukat $d^* = d^*(u, v)$, az eredeti metrikában mért távolságukat $d = d(u, v)$. Annyit kell csak belátnunk, hogy $d \leq d^*$ biztos, de $\mathbf{E}(d^*) \leq \alpha d$. Nyilván feltehetjük, hogy $u \neq v$.

Az első egyenlőtlenséget indukcióval láthatjuk be a rekurzió mélysége szerint. Két esetet kell vizsgálnunk: ha u és v a partíciónálásnál különböző klaszterbe került, akkor a kapott fában gyökérnek különböző gyerekeitől származnak le, $d^* = \Delta$, ez pedig a feltételből nagyobb vagy egyenlő V átmérőjénél, így d -nél is. Ha viszont azonos klaszterbe kerültek, akkor ugyanannak a részfának a levelei, ezt a részfát pedig egy rekurzív hívás állította elő, amelyre alkalmazhatjuk az indukciós feltételt, így ebben a részfában a távolságuk legalább d .

A második állításhoz jelölje i azt a nemnegatív egészt, amelyre $\Delta/\mu^{i+1} < d \leq \Delta/\mu^i$ (nyilván $0 \leq i$). Bizonyítsuk i -re vonatkozó teljes indukcióval azt az erősebb állítást, hogy $\mathbf{E}(d^*) < d\mu(2 \ln n + 2)(1 + \min(i, \log_\mu n))$. Vegyük észre, hogy ha $\log_\mu n \leq i$, akkor $d \leq \Delta/n$, így a partíciónálásról kikötött harmadik tulajdonság miatt u és v biztosan ugyanabba a klaszterbe kerül, így az indukciós feltételt alkalmazhatjuk erre a klaszterre, ezért valóban $\mathbf{E}(d^*) \leq d\mu(2 \ln n + 2)(1 + \log_\mu n)$. Másrészt ha $i = 0$, akkor $\Delta/\mu < d$, ezért u és v nem kerülhet azonos klaszterbe a 2.9. lemma első tulajdonsága miatt, így $d^* = \Delta < d\mu < d\mu(2 \log n + 2)$. Végül ha $0 < i < \log_\mu n$, alkalmazzuk a második tulajdonságot: ez azt mondja ki, hogy annak a valószínűsége, hogy a két csúcs külön klaszterbe kerül, legfeljebb $d/r = d\mu(2 \ln n + 2)/\Delta$. Ha a két csúcs külön klaszterbe kerül, akkor a $d^* = \Delta$. Ha ugyanabba a klaszterbe, akkor az indukciót alkalmazhatjuk a megfelelő részfára, ennek az átmérője Δ/μ , míg d változatlan, így erre a rekurzív hívásra i eggyel kisebb mint a teljes fára, tehát az állításból d^* feltételes várható értéke legfeljebb $d\mu(2 \log n + 2)i$. A két lehetőséget összerakva éppen azt kapjuk, hogy

$$\mathbf{E}(d^*) \leq \frac{d\mu(2 \ln n + 2)}{\Delta} \cdot \Delta + d\mu(2 \log n + 2)i = d\mu(2 \log n + 2)(i + 1).$$

Meg kell még jegyezni, hogy a fenti számításban n -et állandónak tekintettük, holott a részfákra n értéke kisebb az eredeti fa elemszámánál. El-

lenőrizhető azonban, hogy ez nem rontja el a becsléseket. Ezzel beláttuk a 2.11. lemmát.

A fentiek már közvetlenül igazolják a fejezet fő tételét.

Bizonyítás (2.5. tétel). Vegyük a V, d metrikus tér 2.10. definíció szerinti közelítő fáját, ahol a paraméter μ , és a Δ átmérő egyenlő a metrikus tér átmérőjével. A 2.11. lemma szerint ez valóban egy μ -HST, amely által indukált d^* metrika α paraméterrel valószínűségben közelíti d -t. Ezt a HST-t a 2.4. állítás szerint speciális alakra hozhatjuk, és ez nem rontja el a többi tulajdonságot.

Nem fogjuk teljes részletességgel bebizonyítani, hogy ez a konstrukció polinom idejű véletlenített algoritmussal elkészíthető, de bizonyos részleteket kiemelünk.

Egyrészt a bemenet méretéhez meg kell adni pontosan, milyen formájú bemenetet kap az algoritmus. Azt feltételezem, hogy μ értékét a bemenet tartalmazza, de egy 1-nél nagyobb alsó korlátot rögzítettünk rá. hogy a d metrikában minden távolság egész szám, és a metrikát egyszerűen az összes $d(u, v)$ távolság felsorolásával adjuk meg.

Azt nem nehéz látni, hogy a 2.9. lemmában, a 2.10. definícióban, és a 2.4. állításban szereplő konstrukciókat polinom időben el lehet végezni. Csak annyit jegyezni meg, hogy a közelítő fa rekurziója összesen legfeljebb $l \cdot n$ -szer hívja meg magát, ahol l a rekurzió legnagyobb mélysége (egyben a kapott fa mélysége), és $l = O(\log \Delta)$, ahol Δ a d -beli legnagyobb távolság (a legkisebb távolság legalább 1).

Gondot csak az okoz, hogy a 2.8. konstrukcióját is végre lehet-e hajtani polinomiális idő alatt – sőt, egyáltalán algoritmikusan. Ehhez a definícióhoz ugyanis elő kellene állítani egy véletlen valós számot egy megadott eloszlással. Vegyük azonban észre, hogy ha a z_k véletlen szám helyett egy másikat választunk, amely pontosan ugyanazoknál a $d(u, v)$ távolságoknál kisebb, mint az igazi, akkor az eredmény ugyanaz lesz, ilyen $d(u, v)$ távolság pedig csak n^2 féle van. Azt nyilván megtehetjük P -ben polinomiális időben (és polinomiális sok véletlen bittel), hogy z_k helyett ennek egy közelítését számítjuk ki, amelynek az eloszlása z_k -étől mindenhol legfeljebb 2^{-P} -nel tér el. Ez a közelítés

legfeljebb $O(2^{-P}) \cdot n^2$ valószínűséggel kerül másik intervallumba azok közül, amelyeket a $d(u, v)$ távolságok határolnak. Mivel legfeljebb $n^{O(1)} \cdot \log^{O(1)} \Delta$ ilyen véletlenszámot használunk, így legfeljebb $O(2^{-P}) \cdot N^{O(1)}$ valószínűséggel kapunk rossz eredményt, ha N a bemenet mérete. Azonban még ha rossz eredményt is kapunk, a végső eredmény akkor is egy HST, és az ebből kapott d^* indukált gráf-metrika dominálja d -t, mivel pedig $d^*(u, v)$ sosem lehet nagyobb, mint Δ , így $\mathbf{E}(d^*(u, v)) \leq \alpha \cdot d(u, v) + \Delta \cdot O(2^{-P}) \cdot N^{O(1)}$. Ezért aztán az ilyen közelítéssel kapott d^* mérték $\alpha + \varepsilon$ paraméterrel valószínűségben közelíti d -t, ahol $\varepsilon = N^{O(1)} \cdot O(2^{-P'})$, és a futásidő P' -nek és N -nek polinomjával korlátos.

3. A min-sum k -clustering közelítő megoldása

Ebben a fejezetben megadjuk a HST-kre vonatkozó közelítő algoritmust [3] alapján. Az algoritmus bemenetként már egy HST-t kap, kimenete a min-sum k -clustering problémának (lásd 1.11. definíció) olyan megoldása, amely a célfüggvény szerint az optimálistól csak legfeljebb konstans faktoriall tér el, determinisztikus, és polinom időben lefut. Kiemeljük, hogy a k számot bemenetnek tekintjük, tehát a közelítésre és az időre adott korlát minden k -ra egyszerre érvényes. A közelítéshez azonban szükséges, hogy a HST μ paraméterét n -től függően válasszuk meg. A választást később fogjuk megadni.

HST alatt ebben a fejezetben végig a 2.3. definíció szerinti speciális formájú HST-t értünk – ezt megtehetjük, mert a 2.5. tétel ilyet ad. A fejezet fő tétele tehát a következő.

3.1. tétel. Bármely bemenetként kapott speciális formájú μ -HST-re és $0 \leq k$ egészre polinomiális idő alatt konstans faktoriall optimális közelítést tudunk adni a HST indukált metrikáján a min-sum k -clustering problémára. Ehhez azonban feltesszük, hogy μ később meghatározandó, csak n -től függő korlátok között van, ahol az alsó korlát nagyobb egy n -től független 1-nél nagyobb konstansnál, és a korlátok bármely n -re kielégíthetőek.

Most az első előkészítő lépésben visszavezetjük a feladatot arra az esetre, amikor a HST-nek legfeljebb $O(\log_\mu n)$ szintje van.

3.2. lemma. Legyen μ mint az előző tételben. Tegyük fel, hogy van egy polinomiális idejű algoritmusunk, amelynek a bemenete egy $\alpha_L \log_\mu n$ szintű n levelű μ -HST, és egy $0 \leq k$ egész, kimenete pedig egy β konstans közelítéssel optimális megoldás a min-sum k -clustering problémára az indukált metrikán. Ekkor szintén polinom időben 2β konstans közelítéssel optimális megoldást tudunk adni a min-sum k -clustering problémára tetszőlegesen magas μ -HST-re is. Itt α_L egy később meghatározott β -től függő pozitív konstans.

Bizonyítás. Legyen a bemenetünk egy Δ átmérőjű μ -HST az n elemű V levélhalmaz felett, és a k egész szám. Tegyük fel, hogy az optimális megoldás célfüggvénye g_0 . Legyen l az a legkisebb egész, amelyre $\beta \cdot g_0 < \Delta/\mu^l$. Ekkor egy β faktorig optimális megoldásban a klasztereken belül biztosan nem szerepelhet $\Delta/(2\mu^l)$ -nél hosszabb távolság, így a fa első l szintjének nincs sok jelentősége: ha $2 \leq l$, ezeket akár össze is vonhatjuk egy szintté, azaz felveszünk egy új gyökeret, és az l -edik szinten lévő összes csúcsot közvetlenül ehhez kötjük egy-egy $(1 - 1/\mu)\Delta/(2\mu^l)$ hosszú éllel.

Másrészt legyen l' a legnagyobb egész, amelyre $\Delta/\mu^{l'} \leq g_0$. Ekkor nincs jelentősége a $\Delta/(n^2 \cdot \mu^{l'+1})$ -nél rövidebb távolságoknak sem, mert ha ezeket megnöveljük egységesen $\Delta/(n^2 \cdot \mu^{l'+1})$ -re, akkor ez a célfüggvényt legfeljebb $\Delta/\mu^{l'+1}$ -nel növeli meg, így legfeljebb kétszer rosszabb közelítést kaphatunk. Vonjuk tehát össze a fában az (első összevonás előtt) l' -edik és alatta lévő szinteket, azaz minden levelet az l' -edik szinten lévő csúcsokhoz közvetlenül egy-egy $\Delta/(2\mu^{l'})$ hosszú éllel.

Mivel nyilván $l' - l = O(\log_\mu n)$, ezért az így kapott fának legfeljebb ennyi szintje van, és ha a feltételezett algoritmust futtatjuk rajta, akkor a kapott közelítő megoldás az eredeti fán is 2β paraméterű közelítés.

Csak egy probléma van: g_0 értékét nem tudjuk előre meghatározni. Sebj, egyszerűen kipróbálhatunk minden szóba jövő l, l' értéket, és mindegyikre meghívhatjuk az algoritmust, majd a kapott megoldások helyességét könnyen ellenőrizhetjük, és kiválaszthatjuk a legjobbat. Ez legfeljebb annyi hívást igényel, ahány szintes a bemenetként adott fa, ez pedig még mindig csak polinomiális időt jelent.

Második előkészítő lépésként a min-sum k -clustering problémáról átté-

rünk egy hasonló problémára.

3.3. definíció. Legyen adott egy (V, d) metrikus tér és egy k nemnegatív egész. A **balanced k -median** probléma egy megengedett megoldása a középpontok K halmazából, és a klaszterek P_v halmazaiból áll: ahol $K \subseteq V$ egy pontosan k elemű halmaz, és minden $v \in K$ középpontra P_v a V egy részhalmaza és $v \in P_v$, a klaszterek diszjunktak, uniójuk a teljes V . A minimalizálandó célfüggvény a klaszterek elemeinek a középponttól való távolságának a klaszterek elemszámával súlyozott összege:

$$g(d, (K, (P_v))) = \sum_{v \in K} |P_v| \sum_{w \in P_v} d(v, w).$$

Megjegyzem, hogy a balanced k -median problémát másképpen is szokták definiálni, mégpedig a $v \in P_v$ megszorítás nélkül.

Nyilvánvaló a következő állítás.

3.4. állítás. A balanced k -median probléma pozitív lineáris.

3.5. lemma. Ha adott metrikus téren és adott k mellett a $(K, (P_v))$ megengedett megoldás β -közelíti a balanced k -median problémát, akkor a $\{P_v\}$ partíció 2β -közelíti a min-sum k -clustering problémát.

(A jelöléssel kapcsolatban hadd jegyezzem meg, hogy míg a balanced k -median probléma megoldásában a (P_v) csúshalmazokból álló rendszer K elemeivel van indexelve, a min-sum k -clustering problémában az ugyanígy jelölt $\{P_v\}$ nincsenek indexelve.)

Bizonyítás. Legyen $\{P'_0, \dots, P'_{k-1}\}$ a min-sum k -clustering probléma optimális megoldása. Ennek min-sum k clustering szerinti célfüggvénye

$$g'(d, \{P'_i\}) = \frac{1}{2} \sum_i \sum_{u \in P'_i} \sum_{w \in P'_i} d(u, w).$$

Legyen most minden i -re v'_i a P'_i halmaz azon eleme, amelyre $\sum_{w \in P'_i} d(v'_i, w)$ minimális. Ez után legyen minden i -re $P'_{v'_i} = P'_i$. Az így kapott $(\{v'_i\}, (P'_{v'_i}))$

megengedett megoldása a balanced k -mediannak, és célfüggvénye

$$\begin{aligned} g(d, (\{v'_i\}, (P'_{v'_i}))) &= \sum_i |P'_i| \cdot \sum_{w \in P'_i} d(v'_i, w) = \sum_i \sum_{u \in P'_i} \sum_{w \in P'_i} d(v'_i, w) \leq \\ &\leq \sum_i \sum_{u \in P'_i} \sum_{w \in P'_i} d(u, w) = 2g'(d, \{P'_i\}). \end{aligned}$$

Másrészt a $(K, (P_v))$ a balanced k -medianra β -közelítés, ezért a fenti célfüggvénynek legfeljebb β -szorosa:

$$\sum_{v \in K} |P_v| \cdot \sum_{w \in P_v} d(v, w) = g(d, (K, (P_v))) \leq \beta \cdot g(d, (\{v'_i\}, (P'_{v'_i}))).$$

Viszont a $\{P_v\}$ k elemű partíciója V -nek, így megengedett megoldása a minimum k -clustering problémának, célfüggvénye pedig

$$\begin{aligned} g'(d, \{P_v\}) &= \frac{1}{2} \sum_{v \in K} \sum_{u \in P_v} \sum_{w \in P_v} d(u, w) \leq \frac{1}{2} \sum_{v \in K} \sum_{u \in P_v} \sum_{w \in P_v} d(u, v) + d(v, w) = \\ &= \sum_{v \in K} |P_v| \sum_{w \in P_v} d(v, w) = g(d, (K, (P_v))) \quad (11) \end{aligned}$$

A három egyenlőtlenségből együtt pedig $g'(d, \{P_v\}) \leq 2\beta \cdot g'(d, \{P'_v\})$, tehát $\{P_v\}$ valóban 2β faktorig optimális megoldás.

A 3.2. lemma és a 3.5. lemma miatt a továbbiakban csak egy adott n levelű, legfeljebb $O(\log_\mu n)$ szintű μ -HST által indukált gráf-metrikán kell polinom időben konstans közelítéssel megoldást adnunk a balanced k -median feladatra.

Erre egy dinamikus programozáson alapuló algoritmust fogunk bemutatni, amely a problémának bizonyos, egyre nagyobb részmegoldásait keresi. Először leírunk egy olyan változatot, amely optimális megoldást ad, de polinomiálnál több időt használ. Ezután megadjuk ennek egy javítását, ami már polinomiális idejű, de cserébe egy konstans faktorialább rosszabb megoldást is adhat.

Először definiálni fogjuk a probléma részleges megoldásának a fogalmát. A részmegoldást szemléletesen úgy képzeljük el, mint a balanced k -medián

megoldásának részben meghatározott változata, de ennek a képnek nem adunk precíz jelentést. Speciálisan nem feltétlenül tartozik minden rész megoldáshoz teljes megoldás.

Definiálunk egy nemnegatív értékű célfüggvényt (kiértékelő függvényt) is a rész megoldásokon, valamint definiáljuk a rész megoldás lenyomatát. A dinamikus programozási algoritmus lényege az lesz, hogy bizonyos lenyomatokhoz előállítjuk a célfüggvényt minimalizáló részleges megoldást (és a célfüggvény itt felvett értékét is). A bizonyítás sok technikai részletből áll. Meg kell ugyanis adni a legegyszerűbb rész megoldásokat, amelyekből kiindulunk; azt a műveletet, amivel két kisebb rész megoldásból egy nagyobbat kaphatunk, és azt, hogy számíthatjuk ki a célfüggvényét és a lenyomatát egy ilyen összetételnek; azt, hogy minden rész megoldást felépíthetünk ilyen módon; és hogy minden teljes megoldást meg lehet kapni egy bizonyos lenyomatú rész megoldásból.

3.6. definíció. Legyen tehát a továbbiakban adva a Δ átmérőjű speciális formájú μ -HST a V leveleken. Legyen adott a k nemnegatív egész paraméter; és jelölje n a V elemszámát. Jelölje L a fa szintjeinek számát leszámítva a levelek szintjét: a gyökér a nulladik, a gyökér gyerekei az első, stb, a levelek az L -edik **szinten** vannak. Ha két levél legfiatalabb közös őse az i -edik szinten van, akkor azt mondjuk, hogy a két levél az i szinten **kapcsolódik**. Természetesen bármely két különböző levélhez van olyan i ($0 \leq i < L$), hogy a két levél az i -edik szinten kapcsolódik. Idézzük fel a HST-nek azt a tulajdonságát, hogy ilyenkor a két levél távolsága $\mu^{-i} \Delta$

3.7. definíció. Részleges megoldásnak bizonyos $(U, C, (p_{v,i}), (A_v))$ négyeseket hívunk, ahol $U \subseteq V$ és $C \subseteq U$ csúcshalmazok, minden $v \in C$ csúcsra $A_v \subset U$, és minden $v \in C$ csúcsra és $0 \leq i < L$ egészre $p_{v,i}$ nemnegatív egész. A négyesről annyit kötünk még ki, hogy az A_v halmazok páronként diszjunktak legyenek; minden $v \in C$ -re $v \in A_v$; és hogy ha $a_{v,i}$ jelöli az A_v olyan csúcsainak számát, amelyek v -vel az i -edik szinten kapcsolódnak, akkor minden $v \in C$ -re és minden $0 \leq i < L$ -re $a_{v,i} \leq p_{v,i}$.

Szemléletesen a részleges megoldás a balanced k -median megoldásának egy részleges meghatározását jelenti. Az U , amelyet egy részfa leveleinek

kell elképzelni, azt adja meg, hogy a rész megoldás mekkora részét rögzíti a megoldásnak: ha $U = V$, akkor a teljes megoldást egyértelműen rögzíti. A másik három változó jelentése: C az U -ban lévő középpontok halmaza, és A_v v klaszteréből az U -beli pontok halmaza, és $p_{v,i}$ a klaszter olyan pontjainak száma, amelyek v -vel az i -edik szinten kapcsolódnak. Speciálisan $1 + \sum_{0 \leq i < L} p_{v,i}$ a v klaszterének elemszáma.

3.8. definíció. Egy $(U, C, (p_{v,i}), (A_v))$ részleges megoldás **célfüggvényének** értéke legyen

$$\sum_{v \in C} \left(1 + \sum_{0 \leq i < L} p_{v,i} \right) \sum_{0 \leq i < L} p_{v,i} \mu^{-i} \Delta.$$

Ez szemléletesen a balanced k -medián célfüggvényéből az U -beli középpontú klaszterekhez tartozó részösszeg. (Megjegyzem, hogy ez nem az 1.2. definíció értelmében vett célfüggvény.)

3.9. definíció. Egy $(U, C, (p_{v,i}), (A_v))$ részleges megoldáshoz definiáljuk ennek (egyértelmű) **lenyomatát**: ez legyen az $(U, c, s, (q_i))$ négyes; ahol $c = |C|$; $s = |U| - \sum_{v \in C} |A_v|$, és minden $0 \leq i < L$ indexre $q_i = \sum_{v \in C} (p_{v,i} - a_{v,i})$, ahol itt $a_{v,i}$ ismét az A_v olyan csúcsainak száma, amelyek v -vel az i -edik szinten kapcsolódnak.

A lenyomat paramétereinek szemléletes jelentése a következő: c az U -beli középpontok száma; s az olyan U -beli csúcsok száma, amelyek U -n kívüli középponthez tartoznak; és q_i az olyan U -n kívüli csúcsok száma, amelyek U -beli középponthez kapcsolódnak az i -edik szinten át.

A definícióból könnyen következik néhány korlát.

3.10. állítás. Minden részleges megoldás lenyomatában $0 \leq s \leq n$, és minden i -re $0 \leq q_i$.

Ismertetek néhány megszorítást, ami az összes, az algoritmus által vizsgált lenyomatra teljesülni fog. Mivel a rekurziót még nem adtuk meg, ezeket csak később látjuk be.

Az első megszorítás U -ra vonatkozik.

3.11. definíció. Az U levélhalmazt **szabályosnak** nevezzük, ha nem üres, és V -nek van egy olyan w belső csúcsa, és a belső csúcsnak néhány gyereke, hogy U pontosan ezen gyermekekből leszarmazott leveleket tartalmazza.

(A gyerek csúcsok lehetnek belső csúcsok vagy levelek w -tól függően. Megjegyzem, hogy ez az w nem feltétlenül egyértelmű: ha U a w összes leszarmazottjának halmaza, akkor w helyett vehetnénk ennek szülőjét is – kivéve ha w a gyökér csúcs.) Ebből következik, hogy ha w az l -edik szinten van, akkor U bármely elemével együtt az összes olyan levelet is tartalmazza, amelynek vele az $l + 1$ -edik szinten közös őse van.

3.12. állítás. Az összes, a rekurzió által vizsgált lenyomatban U szabályos levélhalmaz. Ez az állítás a rekurzió megadásából világosan következni fog.

A második és a harmadik megszorítás q_i -re vonatkozik.

3.13. állítás. Minden, a rekurzió során vizsgált lenyomatban minden $0 \leq i < L$ szintre $q_i \leq 2n(n - |U|)$

A harmadik megszorítás a bizonyításhoz ugyan nem szükséges, ezért nem fogjuk igazolni, mégis a bizonyítás megértését segítheti. Ez a következő mondja ki: a 3.11. definícióban definiált l -re $0 = q_{l+1} = q_{l+2} = \dots = q_{L-1}$. Szemléletesen ennek az oka, hogy az első tulajdonság miatt U -n kívüli csúcs az U -n belüli klaszter-középponthoz csak l -edik vagy felsőbb szintű csúcson keresztül kapcsolódhat.

A rekurzió célja az lesz, hogy a $(V, k, 0, (0, \dots, 0))$ lenyomathoz találjuk meg az optimális (azaz minimális célfüggvényű) részmegoldást. Ez elegendő a következő lemma miatt.

3.14. lemma. A balanced k -median probléma egy tetszőleges megoldásához elkészíthetünk egy olyan részmegoldást, amelynek lenyomata éppen

$$(V, k, 0, (0, \dots, 0))$$

és amelynek a célfüggvénye megegyezik a részmegoldás célfüggvényével. Megfordítva, egy olyan részmegoldáshoz, amelynek a lenyomata a fentivel egyezik meg, elkészíthetünk egy megoldást ugyanolyan célfüggvénnyel.

Bizonyítás. Legyen adott a $(K, (P_v))$ megoldás. Legyen $p_{v,i}$ a P_v klaszter olyan pontjainak száma, amelyek v -hez az i -edik szinten kapcsolódnak. Vizsgáljuk a $(V, K, (p_{v,i}), (P_v))$ részmegoldást. Ez valóban részmegoldás, mivel minden $v \in K$ -ra és $0 \leq i < L$ szintre a P_v halmazok páronként diszjunktak, $v \in P_v$, és $a_{v,i} = p_{v,i}$. A részmegoldás célfüggvénye

$$\sum_{v \in K} \left(1 + \sum_{0 \leq i < L} p_{v,i} \right) \sum_{0 \leq i < L} p_{v,i} \mu^{-i} \Delta,$$

csak hogy $p_{v,i}$ megválasztása miatt $1 + \sum_i p_{v,i} = |P_v|$, és mivel ha két csúc az i -edik szinten kapcsolódik, akkor a távolságuk éppen $\mu^{-i} \Delta$, ezért

$$\sum_i p_{v,i} \mu^{-i} \Delta = \sum_{u \in P_v} d(u, v);$$

így ez a célfüggvény éppen megegyezik az eredeti megoldás min-sum k -clustering célfüggvényével.

Visszafele legyen adott a $(V, C, (p_{v,i}), (A_v))$ megoldás. Vizsgáljuk ekkor a $(C, (A_v))$ megoldást. Tudjuk, hogy az A_v halmazok diszjunktak, így a klaszterek is; és $c = k$ darab klaszterközpont van. A lenyomatban $s = 0$, ezért a részmegoldás olyan, hogy $|U| = \sum_{v \in C} |A_v|$, tehát minden csúc bekerült valamelyik klaszterbe. Ez tehát megengedett megoldása a problémának. Másrészt mivel a lenyomatban $q_i = 0$, ezért a lenyomat definíciójából $0 = \sum_{v \in C} (p_{v,i} - a_{v,i})$, csak hogy a részmegoldás definíciójából ennek az összegnek minden tagja nemnegatív, így $p_{v,i} = a_{v,i}$ minden v -re. Az előzőhöz hasonló érveléssel látható az is, hogy a részmegoldás célfüggvényének értéke egyenlő a balanced k -medián célfüggvényével a megadott megoldáson.

Most pedig lássuk a rekurzió lépéseit. Ez, mint már említettük, bizonyos $(U, c, s, (q_i))$ lenyomatokra megkeresi a célfüggvény minimumát az ilyen lenyomatú részmegoldásokon, és egy ilyen optimális részmegoldást.

3.15. definíció. A $(U, c, s, (q_i))$ lenyomatú részmegoldások minimumának értékét $D(U, c, s, (q_i))$ jelöli. (Előfordul, hogy valamilyen lenyomathoz nincs részmegoldás, ilyenkor a minimumot $+\infty$ -nek tekintjük.)

Tegyük tehát fel, hogy adott leveleknek egy U halmaza, amely teljesíti az első megszorítást, és adottak a $0 \leq c, s, q_0, \dots, q_{L-1} \leq n$ egészek. A számítási módszernek két esete van: az alapeset az lesz, ha U egyetlen levélből áll, a többi esetben az optimum számítását visszavezetjük kisebb U -ra.

Legyen tehát először $U = \{u\}$ egyelemű. Az optimum ilyenkor könnyen számítható.

3.16. lemma. Egyelemű U halmazzal a lenyomatokhoz tartozó optimumokat a következő képletek adják meg, és a megfelelő részmegoldást is könnyen meg lehet konstruálni.

$$D(\{u\}, 0, 1, (0, \dots, 0)) = 0$$

$$D(\{u\}, 1, 0, (q_0, \dots, q_{L-1})) = \sum_i \left(1 + \sum_i q_i \right) q_i \mu^{-i} \Delta$$

$$D(\{u\}, c, s, (q_0, \dots, q_{L-1})) = +\infty \quad \text{minden más esetben}$$

Bizonyítás. Ha $c = 0$ és $s = 1$ és $q_0 = \dots = q_{L-1} = 0$, akkor az optimum 0, és az ehhez tartozó részmegoldás triviális: $C = \emptyset$, így p és A üresek. Ha $c = 1$ és $s = 0$, akkor bármilyen (q_0, \dots, q_{L-1}) paraméterekre van egy egyértelmű részmegoldás: $C = \{u\}$, $A_u = \{u\}$, és $p_{u,i} = q_i$ minden i -re. Erre a célfüggvény értéke valóban a fent megadott. Ha $c = 0$ de s vagy q nem a fenti alakú, akkor könnyen látható, hogy nincs megfelelő lenyomatú részmegoldás. Ugyanez a helyzet, ha $1 < c$, vagy ha $1 = c$ mellett $0 < s$.

Legyen most U legalább kételemű.

3.17. állítás. Ha az U levélhalmaz szabályos a 3.11. definíció szerint, és legalább kételemű, akkor felírható az $U^{(0)}$ és $U^{(1)}$ szabályos levélhalmazok uniójaként úgy, hogy van a gráfnak egy w belső csúcsa, amelyre mind U , mind $U^{(0)}$ és $U^{(1)}$ előáll a w néhány gyerekének összes leszármazott leveleként. Jelöljük w szintjét l -l. Bármely $U^{(0)}$ -beli csúcs bármely $U^{(1)}$ -beli csúcshoz az l -edik szinten kapcsolódik.

Bizonyítás. Mivel az U reguláris, van egy olyan w csúcs és w gyerekeinek egy W halmaza, amelyre U pontosan a W pontjaiból leszármazott levelek

halmaza. Feltehetjük, hogy W legalább két elemből áll, mert ha W egyelemű lenne, akkor w helyett válasszuk W egyetlen elemét, és egyidejűleg W helyett az új w összes gyerekének halmazát (ezt a lépést legfeljebb véges sokszor kell megismételni, és mivel U nem egyelemű, végül nem kaphatunk levelet). Osszuk a W halmazt két diszjunkt nem üres részre: $W = W^{(0)} \dot{\cup} W^{(1)}$. A rekurziós lépéshez $W^{(0)}$ -t tetszőlegesen választhatjuk W nem üres valódi részalmazai közül, de a bizonyítás későbbi részében további megszorítást fogunk adni rá. Legyen $U^{(0)}$ a $W^{(0)}$ -tól leszármazott összes levél halmaza, és képezzük hasonlóan $W^{(1)}$ -ből $U^{(1)}$ -et. Ezekre a halmazokra az állítás triviális.

Az U -hoz tartozó optimumokat a fenti felosztás után az $U^{(0)}$ -hoz illetve $U^{(1)}$ -hez tartozó optimumokból fogjuk rekurzív módon megkapni. Mivel ezek a levélhalmazok U -nál kisebbek, a rekurzió véges lesz. Pontosan azt állítom, hogy a következő képlet érvényes.

3.18. lemma. Legyen $(U, c, s, (q_i))$ egy lenyomat, melyben az U legalább kételemű szabályos levélhalmaz a 3.17. állítás szerinti módon fel van osztva két részre. Ekkor az ilyen lenyomatú részmegoldások minimális célfüggvényét a következő képlettel számolhatjuk, és könnyen konstruálhatunk egy megfelelő részmegoldást is, ha a képlet jobb oldalán szereplő optimumokhoz ismerünk egy megfelelő részmegoldást.

$$D(U, c, s, (q_0, \dots, q_{L-1})) = \min(D(U^{(0)}, c^{(0)}, s^{(0)}, (q_0^{(0)}, \dots, q_{L-1}^{(0)})) + \quad (*) \\ + D(U^{(1)}, c^{(1)}, s^{(1)}, (q_0^{(1)}, \dots, q_{L-1}^{(1)}))$$

ahol a minimum a

$$c^{(0)}, c^{(1)}, s^{(0)}, s^{(1)}, q_0^{(0)}, \dots, q_{L-1}^{(0)}, q_0^{(1)}, \dots, q_{L-1}^{(1)}, \sigma^{(0)}, \sigma^{(1)}$$

egész értékű változók minden olyan értékén fut, ahol teljesülnek a következő

egyenlőségek

$$\begin{aligned} c &= c^{(0)} + c^{(1)} \\ s &= s^{(0)} + s^{(1)} - \sigma^{(0)} - \sigma^{(1)} \\ q_l &= q_l^{(0)} + q_l^{(1)} - \sigma^{(0)} - \sigma^{(1)} \\ q_i &= q_i^{(0)} + q_i^{(1)} \quad \text{ha } i \neq l \end{aligned}$$

és a következő kiegészítő egyenlőtlenségek

$$\begin{aligned} 0 &\leq \sigma^{(0)} \leq \min(s^{(1)}, q_l^{(0)}) \\ 0 &\leq \sigma^{(1)} \leq \min(s^{(0)}, q_l^{(1)}) \\ 0 &\leq c^{(0)}, c^{(1)}, s^{(0)}, s^{(1)} \leq n \\ 0 &\leq q_i^{(0)}, q_i^{(1)}. \end{aligned}$$

Mielőtt elkezdenénk a bizonyítást, ellenőrizzük a rekurzióra korábban megadott két megszorítást. A 3.12. állítás nyilvánvaló, csak a 3.13. állítás szorul indoklásra.

Bizonyítás (3.13. állítás). A 3.14. lemmában megadott kiindulásra az állítás nyilvánvaló, mivel ott $q_i = 0$ és $|U| = n$. Tegyük most fel, hogy az $(U, c, s, (q_i))$ lenyomatára igaz a korlát. Be kell látnunk ebből, hogy a (*) egyenlőtlenség jobb oldalán minden, az összegben szereplő lenyomatra is igaz. Mivel $|U^{(g)}| < |U|$, ezért a korlát

$$q_i + 2n \leq 2n(n - |U|) + 2n \leq 2n(n - |U^{(g)}|).$$

Viszont mivel $0 \leq \sigma^{(0)} + \sigma^{(1)} \leq s^{(1)} + s^{(0)} \leq 2n$, ezért

$$q_i^{(0)} + q_i^{(1)} \leq q_i + \sigma^{(0)} + \sigma^{(1)} \leq q_i + 2n,$$

és $0 < q_i^{(0)}, q_i^{(1)}$, ezért a korlát valóban teljesül.

Bizonyítás. Először belátjuk, hogy a (*) bal oldala nem kisebb a jobb oldalnál. Ehhez vegyük az eredeti $(U, c, s, (q_i))$ lenyomathoz tartozó optimális

$(U, C, (p_{v,i}), (A_v))$ részmegoldást: ennek a célfüggvénye $D(U, c, s, (q_i))$. Szedjük ezt szét két kisebb részmegoldásra, amelyet $g = 0, 1$ indexekkel jelölünk. Az $(U^{(g)}, C^{(g)}, (p_{v,i}^{(g)}), (A_v^{(g)}))$ részmegoldásban $U^{(g)}$ a korábban rögzített halmaz; $C^{(g)} = C \cap U^{(g)}$; továbbá minden $v \in C^{(g)}$ középpontra $p_{v,i}^{(g)} = p_{v,i}$ és $A_v^{(g)} = A_v \cap U^{(g)}$. (Ez a képzési szabály összhangban van a részmegoldás szemléletes jelentésével: a két indexelt részmegoldás ugyanazt a megoldást határozza meg kevésbé szigorúan, mint az eredeti.)

Legyen $(U^{(g)}, c^{(g)}, s^{(g)}, (q_i^{(g)}))$ a megfelelő részmegoldás lenyomata. Legyen továbbá

$$\sigma^{(g)} = \sum_{v \in C^{(g)}} |A_v \setminus U^{(g)}|, \quad (12)$$

azaz szemléletesen az olyan csúcsok száma, amelyek $U^{(g)}$ -beli középponthez tartoznak, de maguk $U^{(1-g)}$ -beliek. Belátjuk, hogy ezekre a paraméterekre a kikötött egyenletek és egyenlőtlenségek teljesülnek, így a két részmegoldás célfüggvénye szerepel a $(*)$ -beli minimumban.

Nyilvánvaló, hogy a paraméterek egészek, hogy a lenyomat paramétereire adott korlátok teljesülnek, és hogy $\sigma^{(g)}$ nemnegatív. Az első két egyenlőség is könnyen adódik: $c = |C| = |C^{(0)} \dot{\cup} C^{(1)}| = c^{(0)} + c^{(1)}$. Hasonlóan a második egyenlőséghez

$$\begin{aligned} s &= \left| U \setminus \dot{\bigcup}_{v \in C} A_v \right| = \\ &= \left| U^{(0)} \setminus \dot{\bigcup}_{v \in C^{(0)}} (A_v \cap U^{(0)}) \setminus \dot{\bigcup}_{v \in C^{(1)}} (A_v \cap U^{(0)}) \right| + \\ &+ \left| U^{(1)} \setminus \dot{\bigcup}_{v \in C^{(1)}} (A_v \cap U^{(1)}) \setminus \dot{\bigcup}_{v \in C^{(0)}} (A_v \cap U^{(1)}) \right| = \\ &= s^{(0)} - \sigma^{(1)} + s^{(1)} - \sigma^{(0)}. \end{aligned}$$

A szereplő részösszegekből kapjuk azt is, hogy $\sigma^{(1)} \leq s^{(0)}$ illetve $\sigma^{(0)} \leq s^{(1)}$.

A harmadik és negyedik egyenlőséghez először használjuk a definíciókat.

$$q_i = \sum_{v \in C} p_{v,i} - \sum_{v \in C} a_{v,i}$$

$$q_i^{(g)} = \sum_{v \in C^{(g)}} p_{v,i}^{(g)} - \sum_{v \in C^{(g)}} a_{v,i}^{(g)}$$

Az első tagra $C = C^{(0)} \cup C^{(1)}$ és $p_{v,i}^{(g)} = p_{v,i}$ miatt nyilván

$$\sum_{v \in C} p_{v,i} = \sum_{v \in C^{(0)}} p_{v,i}^{(0)} + \sum_{v \in C^{(1)}} p_{v,i}^{(1)}.$$

A második tagot hasonlóan kettébontjuk:

$$\sum_{v \in C} a_{v,i} = \sum_{v \in C^{(0)}} a_{v,i} + \sum_{v \in C^{(1)}} a_{v,i}.$$

Idézzük fel, hogy $a_{v,i}$ az A_v olyan pontjainak a száma, amelyek v -hez az i szinten kapcsolódnak; analóg módon $a_{v,i}^{(g)}$ az $A_v^{(g)} = A_v \cap U_v$ halmaz ilyen elemeinek a száma. Ha tehát $v \in C^{(g)}$, akkor $a_{v,i}^{(g)}$ -be csak olyan pontokat számolunk, amelyeket $a_{v,i}$ -ben is. Az utóbbiban még olyan pontok szerepelnek, amelyek $U^{(1-g)}$ -be esnek: és $\sigma^{(g)}$ épp az ilyen különleges pontok számát adja meg az összes $v \in C^{(g)}$ -re együtt. Mivel egy ilyen pont v -hez csak az l -edik szinten kapcsolódhat, ezért az eltérés teljes egészében az $i = l$ indexű esetbe számít bele. Így tehát

$$\sum_{v \in C^{(g)}} a_{v,l} = \sum_{v \in C^{(g)}} a_{v,l}^{(g)} + \sigma^{(g)}$$

$$a_{v,i} = a_{v,i}^{(g)} \quad \text{ha } i \neq l \text{ és } v \in C^{(g)},$$

ami összegezve az előzővel éppen a kívánt egyenlőségeket adja. Végül mivel a fentiek miatt

$$0 \leq \sum_{v \in C^{(g)}} (p_{v,l} - a_{v,l}) = \sum_{v \in C^{(g)}} (p_{v,l}^{(g)} - a_{v,l}^{(g)}) - \sigma^{(g)} = q_l^{(g)} - \sigma^{(g)}$$

ezért a $\sigma^{(g)} \leq q_l^{(g)}$ korlát is teljesül. Az utolsó két egyenlőtlenség a 3.10. állítás

miatt igaz: ha ez a korlát nem teljesül, akkor a megfelelő lenyomathoz az optimum $+\infty$, tehát az ilyen eseteket nem kell figyelembe venni.

Be kell még látnunk, hogy a részmegoldás célfüggvénye a két kisebb részmegoldás célfüggvényének összege. Ez viszont triviális, hiszen $p_{v,i} = p_{v,i}^{(g)}$ és a célfüggvényt definiáló összeget szét lehet választani a $v \in C^{(0)}$ és $v \in C^{(1)}$ esetekhez tartozó összegekre.

A (*) egyenlet másik irányú becsléséhez azt kell belátni, hogy akárhogy adottak a jobb oldali összeg paraméterei, és két ilyen lenyomatú részmegoldás, a kettőből össze lehet rakni egy részmegoldást a teljes U -n a megfelelő lenyomattal, és ennek a célfüggvénye az előző kettő célfüggvényének az összege.

A bizonyítást jelentősen egyszerűsíti, hogy felhasználhatjuk az előző bizonyítás állításait, ha belátjuk, hogy a kapott U -ra vonatkozó részmegoldást a fenti módszerrel szétválasztva az $U^{(0)}$ -ra és $U^{(1)}$ -re eredetileg adott részmegoldást kapjuk. Ez azonnal adja az összeg paramétereire adott egyenleteket és egyenlőségeket, és a célfüggvények additivitását is; ezek az egyenletek pedig biztosítják, hogy a részmegoldás lenyomata épp a kívánt lenyomat lesz. Be kell azonban látni a (12) egyenlőséget, mivel ezt felhasználtuk a bizonyításban.

Legyenek tehát adottak az

$$(U^{(0)}, C^{(0)}, (p_{v,i}^{(0)}), (A_v^{(0)})), (U^{(1)}, C^{(1)}, (p_{v,i}^{(1)}), (A_v^{(1)}))$$

részmegoldások és a $\sigma^{(g)}$ paraméterek úgy, hogy ezek a részmegoldások lenyomatával és az adott $(U, c, s, (q_i))$ lenyomattal együtt teljesítsék a minimum paramétereire adott összes feltételt.

Ekkor az $(U, C, (p_{v,i}), (A_v))$ részmegoldást a következőképpen állítjuk elő. Legyen $C = C^{(0)} \cup C^{(1)}$ és legyen $p_{v,i} = p_{v,i}^{(g)}$ ahol $v \in C^{(g)}$. Kezdetben legyen $A_v = A_v^{(g)}$ ha $v \in C$. Mivel $U^{(1)}$ -ben $s^{(1)}$ olyan pont van, ami egyik $A_v^{(1)}$ -ben sincs benne, ezek közül kiválaszthatunk $\sigma^{(0)}$ darabot. Ezeket belevesszük tetszőleges $v \in C^{(0)}$ -hoz tartozó A_v halmazhoz, csak arra kell vigyázni, hogy A_v -ben ne legyen $p_{v,l}$ -nél több olyan pont, ami v -hez az l szinten kapcsolódik. Ezt mindig megtehetjük, mert a $\sigma^{(0)} \leq q_i^{(0)}$ feltétel miatt a $p_{v,l} - a_{v,l}$ többletek

összege az összes $v \in C^{(0)}$ -ra éppen $q_i^{(0)}$. Hasonlóan vegyünk hozzá $v \in C^{(1)}$ -hez tartozó A_v -khez pontokat $U^{(0)}$ -ból.

Látható, hogy ez a konstrukció biztosítja az a (12) egyenlet teljesülését. Az is nyilvánvaló, hogy a részmegoldás szétválasztásával visszacapjuk a $C^{(g)}$, $p_{v,i}^{(g)}$, $A_v^{(g)}$ értékeket.

Ezzel tehát a 3.18. lemma bizonyítását befejeztük. A bizonyítás konstruktívan megadta azt is, hogyan állítjuk elő az U -hoz tartozó optimális részmegoldást olyan $U^{(g)}$ -khez tartozó részmegoldásokból, amelyekre a rekurzió alkalmazható.

A módszer a balanced k -median megoldására ezután egyszerűen a következő lenne: a 3.18. és a 3.16. lemmák ismételt használatával kiszámoljuk a $(V, k, 0, (0, \dots, 0))$ lenyomathoz tartozó egy optimális részmegoldást, majd a 3.14. lemmával ebből megkonstruálunk egy optimális megoldást. Ez az eljárás nyilván terminál, mert mindig egyre kisebb U halmazokra vezetjük vissza a nagyobbakat, és a lenyomatban a $0 \leq c, s, q_0, \dots, q_{L-1} \leq 2n^2$ egészek csak véges sok értéket vehetnek fel. E közben nyilván érdemes feljegyezni minden lenyomathoz a kiszámított optimumot, hogy egynél többször ne kelljen ugyanazt kiszámolni.

Sajnos azonban ez az eljárás nem polinom idejű, ugyanis nem tudjuk, hogy e közben nem kell-e polinomiálnál több féle lenyomathoz optimumot számolni. A fejezet hátralevő részében a módszert úgy módosítjuk, hogy ez a futásidő k -tól függetlenül n -ben polinomiális legyen, és be kell látnunk, hogy ez a módosítás nem ad az optimálisnál lényegesen rosszabb közelítést.

3.19. lemma. Legyen adott az U szabályos levélhalmaz és a

$$0 \leq c, s, q_0, \dots, q_{L-1}, q'_0, \dots, q'_{L-1}$$

egészek. Tegyük fel, hogy minden i -re $q_i \leq q'_i$. Ekkor

$$D(U, c, s, (q_i)) \leq D(U, c, s, (q'_i)).$$

Ha még azt is feltesszük, hogy a $0 \leq \delta$ valós szám olyan, hogy minden i -re

$q'_i \leq (1 + \delta) \cdot q_i$, akkor

$$D(U, c, s, (q'_i)) \leq (1 + \delta)^3 D(U, c, s, (q_i)).$$

Bizonyítás. Az első egyenlőtlenséghez egyszerűen vegyünk az $(U, c, s, (q'_i))$ lenyomathoz tartozó optimális részmegoldást, majd ebben csökkentjük a $p_{v,i}$ értékeket úgy, hogy a lenyomat éppen $(U, c, s, (q_i))$ legyen (csak arra kell vigyázni, hogy $0 \leq p_{v,i}$ igaz maradjon). Ettől a célfüggvény nyilvánvalóan nem nőhet.

A másik iránynak elég különös bizonyítása van. Legyen az $(U, c, s, (q_i))$ lenyomatú részmegoldások közül az optimális $(U, C, (p_{v,i}), (A_v))$. Erre tehát a célfüggvény (amit röviden jelöljünk G -vel)

$$G = D(U, c, s, (q_i)) = \sum_{v \in C} \left(1 + \sum_i p_{v,i} \right) \sum_i p_{v,i} \mu^{-i} \Delta.$$

Vegyünk most fel minden $0 \leq i < L$ szintre és minden $0 \leq j < q'_i - q_i$ indexre egy $\varepsilon_{i,j}$ valószínűségi változót, amelynek értéke egy C -beli pont, eloszlása

$$\mathbf{P}(v = \varepsilon_{i,j}) = \frac{p_{v,i}}{\sum_{v \in C} p_{v,i}},$$

és az összes $\varepsilon_{i,j}$ teljesen független egymástól. Képezzük ezekből minden $0 \leq i < L$ szinthez és $v \in C$ középponthoz a $b_{v,i}$ változót, amelynek értéke azon $0 \leq j < q'_i - q_i$ indexek száma, amelyre $v = \varepsilon_{i,j}$. Erre nyilván igaz $\sum_{v \in C} b_{v,i} = q'_i - q_i$, ezért az $(U, C, (p_{v,i} + b_{v,i}), (A_v))$ egy véletlen részmegoldás, és ennek a lenyomata biztosan $(U, c, s, (q_i))$. Be fogjuk látni, hogy ezen részmegoldások átlagos célfüggvénye legfeljebb a megadott felső korlát, ebből pedig következik, hogy legalább egy részmegoldás célfüggvénye is legfeljebb ennyi, tehát az egyenlőtlenség igaz.

Vegyünk észre, hogy mivel $q_i \leq \sum_{v \in C} p_{v,i}$, és $q'_i - q_i \leq \delta q_i$ ezért

$$\mathbf{E}(b_i) = (q'_i - q_i) \cdot \frac{p_{v,i}}{\sum_{v \in C} p_{v,i}} \leq \delta p_{v,i}.$$

Számítsuk ki tehát a rész megoldás célfüggvényét, jelöljük ezt G' -vel.

$$\begin{aligned} G' &= \sum_{v \in C} \left(1 + \sum_i (p_{v,i} + b_{v,i}) \right) \sum_i (p_{v,i} + b_{v,i}) \mu^{-i} \Delta = \\ &= \sum_{v \in C} \left(1 + \sum_i p_{v,i} \right) \sum_i p_{v,i} \mu^{-i} \Delta + \sum_{v \in C} \left(1 + \sum_i p_{v,i} \right) \sum_i b_{v,i} \mu^{-i} \Delta + \\ &\quad + \sum_{v \in C} \left(\sum_i b_{v,i} \right) \sum_i p_{v,i} \mu^{-i} \Delta + \sum_{v \in C} \left(\sum_i b_{v,i} \right) \sum_i b_{v,i} \mu^{-i} \Delta \end{aligned}$$

Számítsuk ki, illetve becsljük meg ennek a várható értékét. Az első három taggal egyszerűen boldogulhatunk. Az első tag konstans G ; a második tag várható értéke

$$\sum_{v \in C} \left(1 + \sum_{i'} p_{v,i} \right) \sum_i \mathbf{E}(b_{v,i}) \mu^{-i} \Delta \leq \sum_{v \in C} \left(1 + \sum_{i'} p_{v,i} \right) \sum_i \delta p_{v,i} \mu^{-i} \Delta = \delta G;$$

hasonlóan a harmadik tagé

$$\sum_{v \in C} \left(\sum_{i'} \mathbf{E}(b_{v,i}) \right) \sum_i p_{v,i} \mu^{-i} \Delta \leq \sum_{v \in C} \left(\sum_{i'} \delta p_{v,i} \right) \sum_i p_{v,i} \mu^{-i} \Delta \leq \delta G.$$

A negyedik tag kicsit nehezebb. Az $\varepsilon_{i,j}$ -k teljes függetlensége miatt ha $i \neq i'$, akkor $b_{v,i}$ és $b_{v,i'}$ függetlenek. Így a negyedik tag várható értéke

$$\begin{aligned} &\sum_{v \in C} \sum_{i'} \sum_i \mathbf{E}(b_{v,i'} b_{v,i}) \mu^{-i} \Delta = \\ &= \sum_{v \in C} \sum_{i'} \sum_i \mathbf{E}(b_{v,i'}) \mathbf{E}(b_{v,i}) \mu^{-i} \Delta + \sum_{v \in C} \sum_i (\mathbf{E}(b_{v,i}^2) - \mathbf{E}(b_{v,i})^2) \mu^{-i} \Delta \end{aligned}$$

Ennek az első tagja az előzőkhöz hasonlóan nem nagyobb $\delta^2 G$ -nél. A második tagban viszont $b_{v,i}$ binomiális eloszlású, ebből a szórás

$$\mathbf{E}(b_{v,i}^2) - \mathbf{E}(b_{v,i})^2 = (q'_i - q_i) \cdot \frac{p_{v,i}}{\sum_{v \in C} p_{v,i}} \cdot \left(1 - \frac{p_{v,i}}{\sum_{v \in C} p_{v,i}} \right) \leq \delta p_{v,i};$$

így hát ez az utolsó tag

$$\begin{aligned} \sum_{v \in C} \sum_i (\mathbf{E}(b_{v,i}^2) - \mathbf{E}(b_{v,i})^2) \mu^{-i} \Delta &\leq \sum_{v \in C} \sum_i \delta p_{v,i} \mu^{-i} \Delta \leq \\ &\leq \sum_{v \in C} \left(1 + \sum_{i'} p_{v,i'} \right) \sum_i \delta p_{v,i} \mu^{-i} \Delta = \delta G. \end{aligned}$$

Összegezve az előzőeket, mivel $0 \leq \delta$, valóban azt kaptuk, hogy

$$\mathbf{E}(G') \leq G + \delta G + \delta G + \delta^2 G + \delta G \leq (1 + \delta)^3 G,$$

amiből, ha egy, az átlagnál nem nagyobb lehetséges G' értéket választunk, a lemma állítása adódik.

Most definiálni fogunk egy függvényt bizonyos lenyomatokon, amit a 3.18. és 3.16. lemmákhoz hasonló rekurzió definiál, közel van D -hez, de gyorsabban ki lehet számítani.

3.20. definíció. Legyen adott a később meghatározandó $1 < \gamma$ szám. Nevezzük megengedett egészeknek a 0-t, és az összes $\lceil \gamma^j \rceil$ számot, ahol $0 \leq j$. Jelölje bármely q nemnegatív egészre $f(q)$ a q -nál nem kisebb számok közül a legkisebb megengedett egészt.

3.21. állítás. Ha $0 \leq q$ egész, akkor $q \leq f(q) < \gamma \cdot q$.

Bizonyítás. Feltehetjük, hogy $2 \leq q$. Nyilván $q \leq f(q)$. Ha viszont $\gamma \cdot q \leq f(q) = \lceil \gamma^j \rceil$ lenne, akkor $f(q) - 1 \leq \gamma^j$, így

$$q - 1 < \frac{f(q)}{\gamma} - \frac{1}{\gamma} \leq \frac{f(q) - 1}{\gamma} \leq \gamma^{j-1}$$

tehát $q \leq \lceil \gamma^{j-1} \rceil$

3.22. definíció. Egy (q_i) vektort megengedett vektornak hívunk, ha a q_i -k mind megengedett egészek.

A módszerünk az lesz, hogy a lenyomatokat olyan lenyomatokkal fogjuk közelíteni, amelyekben (q_i) megengedett vektor. Ilyen lenyomatból már csak polinomiálisan sokat kell megvizsgálni. Ehhez meg kell adnunk a γ és μ paramétert megfelelően.

3.23. definíció. Rögzítsünk egy $0 < \varepsilon < 1$ és egy $1 < \Gamma$ valós konstanst. Rögzítsük a két korábban definiált paramétert úgy, hogy csak n -től függjenek, és teljesüljön a következő két aszimptotika.

$$\begin{aligned}\mu &= \Theta((\ln n)^\varepsilon) \\ \gamma &= \Gamma^{1/O(\ln^2 n)}\end{aligned}$$

Ez összhangban van a korábbi kikötéssel, hogy μ -nek n -től független 1-nél nagyobb alsó korlátja legyen, és hogy $1 < \gamma$.

3.24. definíció. Tegyük fel, hogy az U legalább kételemű szabályos levélhalmaz, amely pontosan a W -beli csúcsok összes leszármazottját tartalmazza, ahol W egy (l -edik szinten lévő) w belső csúcs legalább két gyerekéből áll. Ekkor kettéoszthatjuk a W halmazt $W = W^{(0)} \dot{\cup} W^{(1)}$ módon úgy, hogy $||W^{(0)}| - |W^{(1)}|| \leq 1$ legyen. Legyen $U^{(0)}$ és $U^{(1)}$ a $W^{(0)}$ illetve $W^{(1)}$ leszármazottainak halmaza. Ez egy a 3.17. állításnak megfelelő felbontás. Rögzítsünk minden U halmazhoz egyetlen ilyen $(U^{(0)}, U^{(1)})$ felosztást.

3.25. állítás. Ha bármely U -ból kiindulva újra és újra felosztjuk két részre, akkor legfeljebb $\log_2 n$ mélységben oszthatjuk fel a nélkül, hogy l változna, ezért legfeljebb $L \cdot \log_2 n = O(\log^2 n)$ lépésben egyelemű halmazhoz jutunk.

3.26. definíció. A B függvény olyan $(U, c, s, (q_i))$ lenyomatokon van értelmezve, ahol U szabályos levélhalmaz, és (q_i) megengedett vektor. Ha U egyelemű, akkor

$$B(U, c, s, (q_i)) = D(U, c, s, (q_i)).$$

Ha U legalább kételemű, akkor legyen $U^{(0)}, U^{(1)}$ és l megadva a 3.24. definícióban rögzített módon.

$$\begin{aligned}B(U, c, s, (q_0, \dots, q_{L-1})) &= \min(B(U^{(0)}, c^{(0)}, s^{(0)}, (q_0^{(0)}, \dots, q_{L-1}^{(0)})) + \\ &\quad + B(U^{(1)}, c^{(1)}, s^{(1)}, (q_0^{(1)}, \dots, q_{L-1}^{(1)}))\end{aligned}$$

ahol a minimum a

$$c^{(0)}, c^{(1)}, s^{(0)}, s^{(1)}, q_0^{(0)}, \dots, q_{L-1}^{(0)}, q_0^{(1)}, \dots, q_{L-1}^{(1)}, \sigma^{(0)}, \sigma^{(1)}$$

egész értékű változók minden olyan értékén fut, ahol $q_i^{(g)}$ megengedett egészek, és teljesülnek a következő egyenlőségek

$$\begin{aligned} c &= c^{(0)} + c^{(1)} \\ s &= s^{(0)} + s^{(1)} - \sigma^{(0)} - \sigma^{(1)} \\ q_l + \sigma^{(0)} + \sigma^{(1)} &\leq q_l^{(0)} + q_l^{(1)} < \gamma \cdot (q_l + \sigma^{(0)} + \sigma^{(1)}) \\ q_i &\leq q_i^{(0)} + q_i^{(1)} < \gamma q_i \quad \text{ha } i \neq l \end{aligned}$$

és a következő kiegészítő egyenlőtlenségek

$$\begin{aligned} 0 &\leq \sigma^{(0)} \leq \min(s^{(1)}, q_l^{(0)}) \\ 0 &\leq \sigma^{(1)} \leq \min(s^{(0)}, q_l^{(1)}) \\ 0 &\leq c^{(0)}, c^{(1)}, s^{(0)}, s^{(1)} \leq n. \end{aligned}$$

Látható, hogy ez a rekurzió nagyon hasonlít a 3.18. lemmában D -re megadott rekurzióra. Az egyetlen eltérés az, hogy a q_i számot megengedett egészre kerekítjük. Vizsgáljuk meg a rekurzió futásidejét. Ehhez be kell látnunk a 3.13. állítás analógiáját.

3.27. állítás. Minden $(U, c, s, (q_0, \dots, q_{L-1}))$ lenyomatban, amit a fenti rekurzió megvizsgál $B(V, k, 0, (0, \dots, 0))$ kiszámításához, minden i -re

$$q_i \leq \alpha_q \cdot n \ln^2 n$$

valamely $0 < \alpha_q$ konstansra.

Bizonyítás. A bizonyítás is hasonlít a 3.13. állításéhoz.

Azt igazoljuk H -ra vonatkozó indukcióval, hogy ha a B -t definiáló rekurzió a $(U, c, s, (q_0, \dots, q_{L-1}))$ lenyomatot a $(B, k, 0, (0, \dots, 0))$ lenyomattól

számított H -adik rekurzív hívási szinten vizsgálja, akkor minden i -re

$$q_i \leq 2nH\gamma^H.$$

Ez nyilván igaz a $0 = H$ alapesetre. Tegyük most fel, hogy az $(U, c, s, (q_i))$ lenyomathoz a H -adik szinten jutottunk el, és erre igaz az indukciós feltétel. Be kell látnunk, hogy ekkor az összes $(U^{(g)}, c^{(g)}, s^{(g)}, (q_i^{(g)}))$ lenyomatra, amely a minimumban szerepel, szintén igaz. De $\sigma^{(g')} \leq s^{(1-g')} \leq n$ miatt valóban

$$q_i^{(g)} < \gamma \cdot (q_i + \sigma^{(0)} + \sigma^{(1)}) \leq \gamma \cdot (2n(H-1)\gamma^{H-1} + 2n) \leq 2nH\gamma^H.$$

Mármost a 3.25. állítás miatt H sosem lehet több $O(\log^2 n)$, ezért valóban

$$\begin{aligned} q_i &= O(nH\gamma^H) = O(n \log^2 n \cdot \gamma^{\log^2 n}) = \\ &= O(n \log^2 n \Gamma^{\log^2 n / O(\log^2 n)}) = O(n \log^2 n). \end{aligned}$$

3.28. állítás. $B(V, k, 0, (0, \dots, 0))$ értékét n -ben polinomiális időben ki lehet számolni.

Bizonyítás. Adjunk becslést arra, hány különböző $(U, c, s, (q_i))$ lenyomatot vizsgál meg a rekurzió. Nyilvánvaló, hogy U pontosan $2n-1$ különféle értéket vesz fel. Azt tudjuk, hogy $0 \leq c \leq n$, $0 \leq s \leq n$, és $0 \leq q_i$, és a 3.27. állításban beláttuk, hogy $q_i \leq \alpha_q n \ln^2 n$, ennél kisebb megengedett szám pedig csak $O(\log^2 n)$ darab van. Mivel a szintek száma $L = O(\log_\mu n) = O(\log n / \log^\epsilon n)$, ezért a lehetséges lenyomatok száma valóban

$$O(n^3 \cdot (\log^2 n)^{\log^{1-\epsilon} n}) = n^{O(1)}.$$

Most a rekurzió elvégzéséhez egyszerűen minden egyes lenyomathoz kiszámíthatunk B értékét úgy, hogy minden két másik lenyomatról és minden $0 \leq \sigma^{(0)}, \sigma^{(1)} \leq n$ értékekről ellenőrizzük, teljesítik-e ezek a kikötött egyenleteket és egyenlőtlenségeket, és kiszámítjuk az összegek minimumát. Ha gondoskodunk róla, hogy az egyes lenyomatokra a B függvényt egynél többször ne számoljuk ki, akkor ez az eljárás valóban polinom időben elvégezhető.

Most pedig lássuk be, hogy a rekurzió eredménye is közel van D -hez.

3.29. lemma. Ha az $(U, c, s, (q_i))$ lenyomatban U szabályos levélhalmaz, és (q_i) megengedett vektor, akkor

$$D(U, c, s, (q_i)) \leq B(U, c, s, (q_i)) \leq \alpha_B \cdot D(U, c, s, (q_i)),$$

valamely $0 < \alpha_B$ konstansra.

3.30. definíció. Jelölje a továbbiakban $h(U)$ azt a mélységet, ahány lépés után U -ból legkésőbb elérhetünk egy egyelemű halmazhoz, ha választhatunk, hogy U -ról $U^{(0)}$ -ra vagy $U^{(1)}$ -re lépünk. Másképpen ha U egyelemű, akkor $h(U) = 0$; különben $h(U) = \max(h(U^{(0)}), h(U^{(1)}))$.

Bizonyítás. Az alsó korlátot teljes indukcióval bizonyítjuk az U -ra vonatkozó rekurzió mélysége, azaz $h(U)$ szerint.

Ha $0 = h(U)$, vagyis U egyelemű, akkor a lemma nyilvánvaló, ezért most tegyük fel, hogy legalább kételemű.

Legyenek $U^{(g)}, c^{(g)}, s^{(g)}, q_i^{(g)}, \sigma^{(g)}$ azok a paraméterek, amelyekre a B -t definiáló rekurzió a minimumot felveszi, tehát

$$B(U, c, s, (q_i)) = B(U^{(0)}, c^{(0)}, s^{(0)}, (q_i^{(0)})) + B(U^{(1)}, c^{(1)}, s^{(1)}, (q_i^{(1)})).$$

Az indukciós feltétel szerint

$$D(U^{(g)}, c^{(g)}, s^{(g)}, (q_i^{(g)})) \leq B(U^{(g)}, c^{(g)}, s^{(g)}, (q_i^{(g)})),$$

ezért

$$D(U^{(0)}, c^{(0)}, s^{(0)}, (q_i^{(0)})) + D(U^{(1)}, c^{(1)}, s^{(1)}, (q_i^{(1)})) \leq B(U, c, s, (q_i))$$

Legyen

$$\begin{aligned} q'_l &= q_l^{(0)} + q_l^{(1)} - \sigma^{(0)} - \sigma^{(1)} \\ q'_i &= q_i^{(0)} + q_i^{(1)} \quad \text{ha } i \neq l. \end{aligned}$$

Ekkor

$$D(U, c, s, (q'_i)) \leq D(U^{(0)}, c^{(0)}, s^{(0)}, (q_i^{(0)})) + D(U^{(1)}, c^{(1)}, s^{(1)}, (q_i^{(1)})),$$

mivel könnyen látható, hogy a 3.18. lemmában felsorolt összes egyenlőtlenség teljesül. Viszont minden i -re $q_i \leq q'_i$, így alkalmazhatjuk a 3.19. lemmát, e szerint

$$D(U, c, s, (q_i)) \leq D(U, c, s, (q'_i)),$$

így valóban

$$D(U, c, s, (q_i)) \leq B(U, c, s, (q_i)),$$

A felső korláthoz azt az erősebb állítást látjuk be $h(U)$ -ra vonatkozó indukcióval, hogy

$$B(U, c, s, (q_i)) \leq \gamma^{3h(U)} \cdot D(U, c, s, (q_i)).$$

Írjuk fel $D(U, c, s, (q_i))$ -re a (*) szerinti rekurziót: legyenek ebben a rekurzióban a minimumnál felvett paraméterek $U^{(g)}, c^{(g)}, s^{(g)}, q_i^{(g)}, \sigma^{(g)}$. Tehát

$$D(U, c, s, (q_i)) = D(U^{(0)}, c^{(0)}, s^{(0)}, (q_i^{(0)})) + D(U^{(1)}, c^{(1)}, s^{(1)}, (q_i^{(1)})).$$

Legyen $Q_i^{(g)} = f(q_i^{(g)})$. Ekkor ha $i \neq l$, akkor

$$q_i = q_i^{(0)} + q_i^{(1)} \leq Q_i^{(0)} + Q_i^{(1)} < \gamma q_i^{(0)} + \gamma q_i^{(1)} = \gamma q_i,$$

és

$$\begin{aligned} q_l + \sigma^{(0)} + \sigma^{(1)} &= q_l^{(0)} + q_l^{(1)} \leq Q_l^{(0)} + Q_l^{(1)} \\ Q_l^{(0)} + Q_l^{(1)} &< \gamma q_l^{(0)} + \gamma q_l^{(1)} = \gamma(q_l + \sigma^{(0)} + \sigma^{(1)}), \end{aligned}$$

így az összes feltétel teljesül ahhoz, hogy

$$B(U, c, s, (q_i)) \leq B(U^{(0)}, c^{(0)}, s^{(0)}, (Q_i^{(0)})) + B(U^{(1)}, c^{(1)}, s^{(1)}, (Q_i^{(1)})).$$

Az indukciós feltételből

$$B(U^{(g)}, c^{(g)}, s^{(g)}, (Q_i^{(g)})) \leq \gamma^{3h(U)-3} D(U^{(g)}, c^{(g)}, s^{(g)}, (Q_i^{(g)})).$$

Viszont a 3.19. lemma miatt, mivel $Q_i^{(g)} \leq \gamma q_i^{(g)}$,

$$D(U^{(g)}, c^{(g)}, s^{(g)}, (Q_i^{(g)})) \leq \gamma^3 \cdot D(U^{(g)}, c^{(g)}, s^{(g)}, (q_i^{(g)}));$$

ebből pedig valóban

$$B(U, c, s, (q_i)) \leq \gamma^{3h(U)} \cdot D(U, c, s, (q_i)).$$

Végül az indukciós feltételből következik a lemma állítása, mert a 3.25. állítás szerint bármely U szabályos levélhalmazra $h(U) = O(\log^2 n)$, így

$$\gamma^{3h(U)} = \Gamma^{\log^2 n / O(\log^2 n)} = O(1).$$

Annyit kell még belátnunk, hogy a B függvényt a rekurzió alapján nem csak ki tudjuk számolni, hanem a megfelelő részmegoldást is meg tudjuk találni.

3.31. állítás. Bármely $(U, c, s, (q_i))$ lenyomathoz, amelyhez eljuthatunk a $B(V, k, 0, (0, \dots, 0))$ rekurzív kiszámításánál, polinomiális időben meg tudunk adni egy olyan részmegoldást, amelynek lenyomata megegyezik ezzel a lenyomattal, és célfüggvénye nem nagyobb, mint $B(U, c, s, (q_i))$.

Bizonyítás. Az állítást természetesen indukcióval látjuk be U mérete szerint. A 3.16. lemma szerint egyelemű U -ra az állítás igaz. Tegyük fel tehát, hogy a megadott lenyomatban U legalább kételemű szabályos halmaz, és hogy minden kisebb U -ra igaz az állítás.

Írjuk fel a 3.26. definíció szerinti rekurziót a megadott lenyomatra, és keressük meg a minimum paramétereit. A 3.28. állításban láttuk, hogy ezt polinomiális idő alatt meg lehet tenni. A kapott $c^{(g)}, s^{(g)}, q_i^{(g)}, \sigma^{(g)}$ paraméterekkel tehát

$$B(U, c, s, (q_i)) = B(U^{(0)}, c^{(0)}, s^{(0)}, (q_i^{(0)})) + B(U^{(1)}, c^{(1)}, s^{(1)}, (q_i^{(1)})).$$

Számítsuk ki az indukciós feltétel alapján az $(U^{(g)}, c^{(g)}, s^{(g)}, (q_i^{(g)}))$ lenyomatokhoz a megfelelő részmegoldásokat. Ezeknek a célfüggvénye külön-külön legfeljebb $B(U^{(g)}, c^{(g)}, s^{(g)}, (q_i^{(g)}))$, így összesen legfeljebb $B(U, c, s, (q_i))$.

A két részmegoldásból most a 3.18. lemma bizonyításával teljesen megegyező módon összerakhatunk egy nagyobb részmegoldást. A kapott részmegoldás lenyomata $(U, c, s, (q'_i))$, ahol ha $i \neq l$, akkor

$$q_i = q_i^{(0)} + q_i^{(1)} = q'_i,$$

míg

$$q_l \leq q_l^{(0)} + q_l^{(1)} - \sigma^{(0)} - \sigma^{(1)} = q'_l,$$

tehát minden i -re $q_i \leq q'_i$. Így aztán ezt a részmegoldást a 3.19. lemma alsó korlátjához hasonlóan átalakíthatjuk egy $(U, c, s, (q_i))$ lenyomatú részmegoldássá, ezzel nem növelve a célfüggvényt. Ennek a részmegoldásnak tehát $B(U, c, s, (q_i))$ -nél nem nagyobb a célfüggvénye.

Az, hogy mindezt valóban polinomiális időben végrehajthatjuk, a 3.28. állításhoz teljesen hasonló számítással ellenőrizhető.

Most már csak össze kell fűznünk a lemmákat a bizonyításhoz.

3.32. állítás. Olyan n csúcsú μ -HST fákon, amelyeknek $O(\log_\mu n)$ szintje van, polinomiális idő alatt konstans közelítést tudunk adni az indukált metrikán felírt min-sum k -clustering problémára; feltéve, ha mi határozhatjuk meg a μ paramétert előre n függvényében (de 1-nél nagyobb alsó korláttal).

Bizonyítás. Alkalmazzuk a 3.31. állítást a $(V, k, 0, (0, \dots, 0))$ lenyomatra. Az így kapott részmegoldás célfüggvénye legfeljebb $B(V, k, 0, (0, \dots, 0))$, de a 3.29. lemma szerint ez nem nagyobb $\alpha_B \cdot D(V, k, 0, (0, \dots, 0))$ -nál, tehát ez a részmegoldás a legjobb ugyanilyen lenyomatú részmegoldásnál is csak legfeljebb konstans faktorialrosszabb.

A 3.14. lemma szerint ebből kapunk egy konstans faktor erejéig optimális megoldást a balanced k -medián problémára. Ez a 3.5. lemma miatt a min-sum k -clustering problémának is konstans faktorig optimális megoldása.

Bizonyítás (3.1. tétel). A 3.2. lemmából és a 3.32. állításból következik.

Végül fő tételünk bizonyítása.

Bizonyítás (1.12. tétel). Alkalmazzuk a 2.5. tételt a bemenetként kapott metrikus térre a 3.23. definícióban választott μ paraméterrel. Ez polinomiális idő alatt megad egy μ paraméterű HST-t, amelyen az indukált metrika az eredeti teret $\alpha = \mu(2 \ln n + 2)(1 + \log_\mu n)$ paraméterrel valószínűségben közelíti.

Futtassuk le a 3.1. tétel algoritmusát ezen a HST-n a bemenetként megadott k -val. Ez megadja egy β konstans faktor erejéig közelítő megoldását a min-sum k clustering problémának az indukált metrikus téren.

Az 1.15. lemma szerint így a min-sum k -clustering problémára az eredeti metrikus téren átlagosan legfeljebb $\alpha\beta$ faktor erejéig közelítő megoldást kapunk. Viszont μ értékét behelyettesítve látható, hogy $\alpha\beta = O(\log^2 n)$.

Hivatkozások

- [1] Yair Bartal, Probabilistic approximation of metric spaces and its algorithmic application. *Annual Symposium on Foundations of Computer Science*, **37** (1996), 184–193.
- [2] Yair Bartal, On Approximating Arbitrary Metrics by Tree Metrics. *Annual ACM symposium on Theory of Computing*, **30** (1998), 161–168.
- [3] Yair Bartal, Moses Charikar, Danny Raz, Approximating min-sum k -Clustering in Metric Spaces. *Annual ACM Symposium on Theory of Computing*, **33** (2001), 11–20.
- [4] M. Charikar, C. Chekuri, A. Goel, S. Guha, and S. Plotkin, Approximating a Finite Metric by a Small Number of Tree Metrics. *Annual Symposium on Foundations of Computer Science*, **39** (1998), 379–388.
- [5] Piotr Indyk, Algorithmic applications of low-distortion geometric embeddings. *Annual Symposium on Foundations of Computer Science*, **42** (2001), 10–35.

- [6] Piotr Indyk, Jiří Matoušek, Low-distortion embeddings of finite metric spaces. Chapter 8 in *Handbook of Discrete and Computational Geometry*, Jacob. E. Goodman and Joseph. O'Rourke (editors). CRC Press, 2004.
- [7] Jiří Matoušek, *Lectures on discrete geometry*. Springer, 2002, Series: Graduate Texts in Mathematics. Chapter 15: Embedding Finite Metric Spaces into Normed Spaces.
- [8] W. Fernandez de la Vega, Marek Karpinski, Claire Kenyon, Yuval Rabani, Polynomial Time Approximation Schemes for Metric Min-Sum Clustering. *Electronic Colloquium on Computational Complexity (ECCC)*, Rep. 025 (2002)