

# Кластеровање вискодимензионих података

Ненад Митић

Математички факултет  
[nenad@matf.bg.ac.rs](mailto:nenad@matf.bg.ac.rs)





























































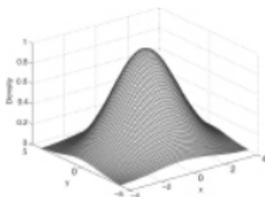




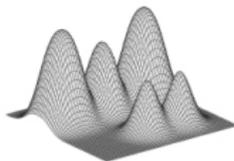


# DENCLUE - илустрација

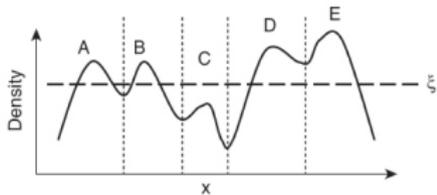
$$K(y) = e^{-\text{distance}(x,y)^2/2\sigma^2}$$



Set of 12 points.



Overall density—surface plot.





# DENCLUE - кораци

- У првом кораку препроцесирања формира се мрежа за податке дељењем минималног граничног хиперправоугаоника у  $d$ -димензионалне хиперправоугаонике дужине ивице  $2\sigma$
- Одређују се мрежне ћелије које садрже тачке. Ћелије су нумерисане на једној ивици граничног хиперправоугаоника
- Вредности нумерације се чувају у дрвету претраживања ради ради касније обраде.
- За сваку сачувану мрежну ћелију чува се број тачака, сума тачака у ћелији и везе са суседним ћелијама
- У другом кораку се врши кластеровање узимајући у обзир само ћелије са великим бројем тачака и ћелије које су повезане са њима
- За сваку тачку  $k$ , функција локалне густине израчунава се узимајући у обзир само оне тачке које су из мрежних ћелија које су у њеној близини
- DENCLUE спаја густе ћелије које се могу спојити путем тачака које имају густину већу од задатог параметра  $\xi$
- Крајњи резултат су само кластри чија је густина већа или једнака  $\xi$

# SNN

## SNN (Shared Near Neighbour Graph)

- Иницијално, алгоритам одређује  $K$  најближих суседа сваке тачке у улазном скупу
- Сличност парова тачака се рачуна као број заједничких најближих суседа обе тачке у пару
- Сличност пара је 0 ако друга тачка није у скупу најближих суседа првој тачки
- Густина сваке тачке је једнака броју тачака унутар растојања  $\epsilon$
- Тачке су класификоване као тачке језгра ако је густина тачака  $\geq$  већа од задатог прага
- Кластери се формирају од тачака у језгру, при чему истом кластеру припадају тачке које су на растојању  $\leq \epsilon$
- Тачка је шум ако је на растојању већем од  $\epsilon$  од било које тачке језгра
- Преостале тачке се доделе кластеру који садржи најсличније тачке језгра

# Хибридни алгоритми

BRIDGE - комбинација К-средина и приступа заснованог на густини

- Изврши се алг. К-средина и свакој тачки додели идентификација кластера
- За сваки кластер К-средина изврши се DBSCAN са параметрима наслеђеним од К-средина
- Помоћу DBSCAN се одреде тачке које јесу и нису су шум
- Изврши се К-средина са првобитним центроидима уз искључивање тачака које су шум

CUBN - комбинација К-средина и приступа заснованог на растојању

- У првом кораку се одређују тачке на граници
- Методом најближих суседа се кластерују тачке на граници. Добија се скуп кластера са граничним тачкама
- Методом најближих суседа се кластерују унутрашње тачке и тачке из скупа кластера са граничним тачкама

# GRIDCLUS

- Алгоритам дели простор података помоћу координатне мреже у облику решетке сачињене од хиперкоцки
- Тачке се смештају у неку од хиперкоцки у зависности од вредности њихових атрибута
- Кластеровање се врши методом тражења суседа (хиперкоцки које се додирују)

# BANG

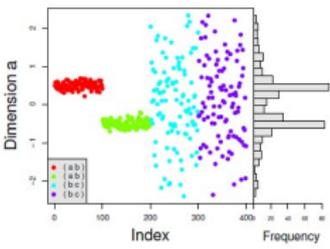
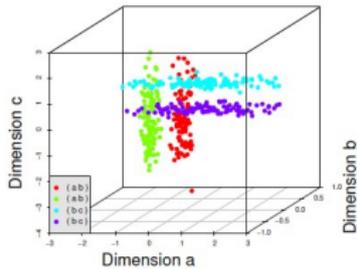
- Проширење GRIDCLUS алгоритма - простор се дели у хијерархијски скуп блокова
- Блокови су сортирани у опадајућем редоследу према густини
- Блокови са највећом густином постају центри кластера
- Остали блокови се итеративно кластерују према њиховој густини тако што формирају нови кластер или се додељују постојећим

# WaveCluster

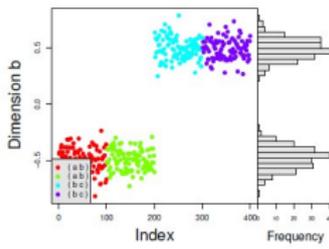
- Може да кластерује, поред осталих, просторне и мултимедијалне податке
- Користи се трансформација улазних података таласићима и начужење кластера у тако трансформисаним подацима
- Алгоритам је отпоран на шум, може да открије кластере произвољног облика са различитим нивоом детаља
- Може да ради са великом количином улазних података



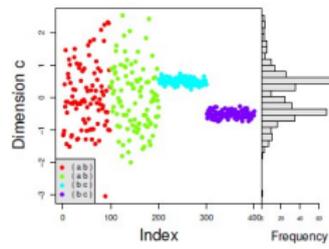
# Кластеровање по подпросторима



(a) Dimension a



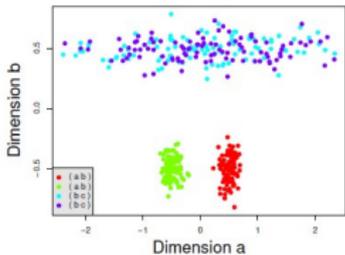
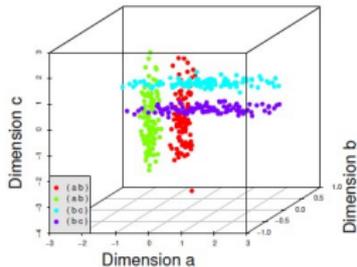
(b) Dimension b



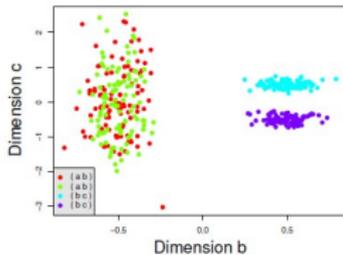
(c) Dimension c



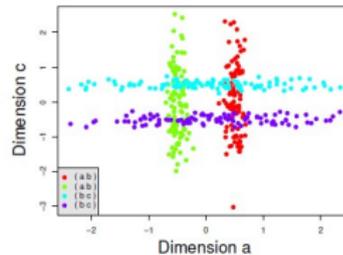
# Кластеровање по подпросторима



(a) Dims  $a$  &  $b$



(b) Dims  $b$  &  $c$



(c) Dims  $a$  &  $c$

# CLIQUE

## CLIQUE (CLustering In QUEst)

- Један од првих алгоритама који проналази кластере унутар потпростора података
- Хибридни алгоритам заснован на густини и мрежи који врши кластеровање по потпросторима високодимензионих нумеричких података
- Одређује кластере угнежене у потпросторима без велике интервенције корисника
- Користи APRIORI технику за кластеровање потпростора
- Пронађени густе потпростори се сортирају по покривености, где је покривеност дефинисана као део скупа података који покривају густе јединице у потпростору
- Алгоритам затим проналази густе суседне ћелије у сваком од изабраних потпростора користећи претрагу у дубину. Кластери настају комбиновањем таквих мрежних ћелија користећи похлепну стратегију
- CLIQUE је у стању да пронађе различите врсте и облике кластера, као и било који број кластера у било којем броју димензија



# ENCLUS

## ENCLUS (ENTropy-based CLUStering)

- Модификација CLIQUE алгоритма која не мери директно густину или покривеност, већ ентропију
- Потпростор са кластерима обично има нижу ентропију од потпростора без кластера
- Ентропија се смањује како се повећава густина ћелија
- Користи АПРИОРИ принцип и приступ одоздо према горе као и CLIQUE за одређивање значајних потпростора
- Може да пронађе кластере произвољног облика који су угнеждени у потпросторима оригиналног скупа података

# MAFIA

## MAFIA (Merging of Adaptive Finite IntervAls)

- Модификација CLIQUE алгоритма која не користи фиксну величину ћелија решетке са једнаким бројем делова у свакој од димензија, већ конструише адаптивну мрежу ћелија за сваку димензију
- Циљ модификација је побољшање кластеровања у потпросторима и омогућавање паралелизације процеса кластеровања ради обраде великих скупова података
- Алгоритам формира хистограм за одређивање минималног броја ћелија неке димензије, а затим комбинује суседне ћелије сличне густине да би се формирале веће
- Када су ћелије дефинисане, MAFIA наставља слично као CLIQUE, користећи АПРИОРИ приступ формира кластере у простору веће димензије користећи густе ћелије
- MAFIA проналази произвољан број кластера произвољног облика у потпросторима различите величине

# OPTIGRID

## OptiGrid (OPTimal GRID-Clustering)

- Користи процену густине ради одређивања центара кластера као да је кластеровање рађено DENCLUE алгоритмом
- Кластер је регион са концентрисаном густином центриран око јаког атрактора густине или локалног максимума функције густине, при чему је густина већа од задатог прага
- Рекурзивном поделом простора атрибута у вишедимензионе решетке формира се оптимална мрежа партиција конструисањем хиперравни
- Конструисане хиперравни деле простор на области мале густине и чувају просторе високе густине (кластере), и посебно центре кластера
- Хиперравни се одређују коришћењем скупа скупих договорених линеарних пројекција простора атрибута

# OPTIGRID варијанте

## O-Cluster (Orthogonal partitioning CLUSTERing)

- Уклања ограничење OPTIGRID-а које се односи на скалабилност у односу на величину меморије
- Користи технику случајног узорковања и бафер релативно мале величине
- Користи стратегију деобе на делове паралелно са координатним осама ради налажења густих региона
- Поред хиперравни користи и статистички тест ради провере квалитета хиперравни које секу простор

## CBF (Cell-Based Filtering)

- Фокусира се на скалабилност мрежне структуре
- Може да ради са великим датотекама у меморији
- Омогућава ефикасно уношење и дохватање кластера из мрежне структуре
- Дефинише јасне критеријуме за конструкцију хиперравни које деле простор



















