

Metrický prostor

$M = (D, d)$

D – datová doména

d – funkce vzdálenosti/metrická funkce/metrika: $D \times D \rightarrow \mathbb{R}$

-funkce může být:

-diskrétní (vrací hodnoty pouze z malé, předem definované množiny)

-spojitá (kardinalita množiny je velmi velká až nekonečná)

Metrický prostor musí splňovat:

(a) nezáporná funkce vzdálenosti:

$$\forall x, y \in D, d(x, y) \geq 0$$

(b) symetrie:

$$\forall x, y \in D, d(x, y) = d(y, x)$$

(c) identita:

$$\forall x, y \in D, x = y \Leftrightarrow d(x, y) = 0$$

(d) trojúhelníková nerovnost:

$$\forall x, y, z \in D, d(x, z) \leq d(x, y) + d(y, z)$$

Alternativní definice:

(p1) nezáporná funkce vzdálenosti:

$$\forall x, y \in D, d(x, y) \geq 0$$

(p2) symetrie:

$$\forall x, y \in D, d(x, y) = d(y, x)$$

(p3) reflexivita:

$$\forall x \in D, d(x, x) = 0$$

(p4) kladnost (positiveness):

$$\forall x, y \in D, x \neq y \Rightarrow d(x, y) > 0$$

(p5) trojúhelníková nerovnost:

$$\forall x, y, z \in D, d(x, z) \leq d(x, y) + d(y, z)$$

Pseudo-metrika

když vlastnost (p4) neplatí?

Quasi-metrika

kdyz vlastnost (p2) neplati?

Super-metrika

kdyz je na vlastnost (p5) kladen silnejsi pozadavek:

$$\forall x, y, z \in \mathcal{D} : d(x, z) \leq \max \{d(x, y), d(y, z)\}$$

Minkowskeho vzorec/vzdalenosti/metriky

-tez nazyvane jako L_p metriky

-definovane na n dimenzionalnich vektorech:

$$L_p[(x_1, \dots, x_n), (y_1, \dots, y_n)] = \sqrt[p]{\sum_{i=1}^n |x_i - y_i|^p}$$

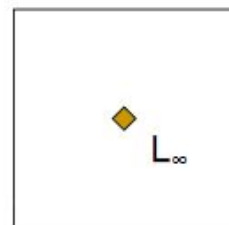
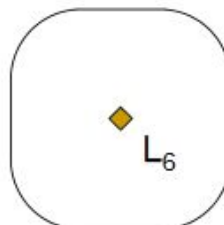
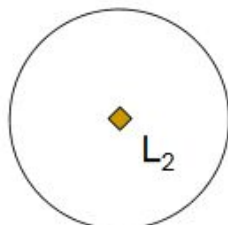
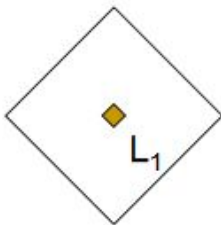
specialni pripady:

L_1 – Manhattanska metrika ($p=1$)

L_2 – Euklidovska metrika ($p=2$)

L_∞ – Maximalni (nekonecna) metrika:

$$L_\infty = \max_{i=1}^n |x_i - y_i|$$



na obrazku je znazornena mnozina bodu ktere maji v dane metrice stejnou vzdalenost od pivota

Jaccarduv koeficient

-index podobnosti

-mereni "vzdalenosti" mezi dvemi mnoziny:

$$d(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

"Tanimoto similarity" definovana na vektorech:

$$d_{TS}(\vec{x}, \vec{y}) = 1 - \frac{\vec{x} \cdot \vec{y}}{\|\vec{x}\|^2 + \|\vec{y}\|^2 - \vec{x} \cdot \vec{y}}$$

$\vec{x} \cdot \vec{y}$ je skalarni soucin

$\|\vec{x}\|$ je velikost vektoru (eukleidovska norma)

Vyjmenovat 3 typy dotazu

-range query

-nearest neighbor query

-k-nearest neighbor query

-reversed k-nearest neighbor query

-similarity join

-kombinovane dotazy (combined query)

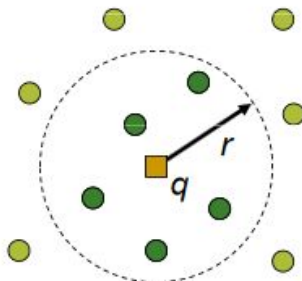
Range query

"vsechny muzea az do vzdalenosti 2km od meho hotelu"

$$R(q, r) = \{x \in X \mid d(q, x) < r\}$$

q – pivot

r – radius



Nearest neighbor (NN) query

"nejblizsi muzeum od meho hotelu"

$$NN(q) = x$$

$$x \in X, \forall y \in X, d(q, x) \leq d(q, y)$$

x ma mensi vzdalenost k pivotovi nez kdokoli jiny

k-NN query

"pet nejblizsich muzei od meho hotelu"

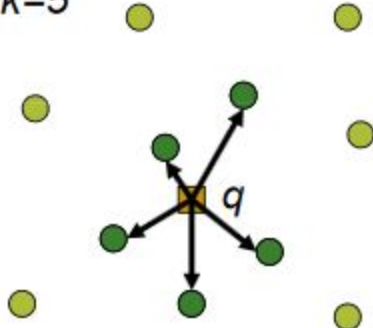
$$k\text{-}NN(q, k) = A$$

$$A \subseteq X, |A| = k$$

$$\forall x \in A, y \in X - A, d(q, x) \leq d(q, y)$$

vysledek je **A**, coz je podmnozina domeny **X** o velikosti **k**, pricemz kazdy prvek z **A** ma mensi vzdalenost k pivotovi **q** nez kazdy prvek z doplnku

k=5



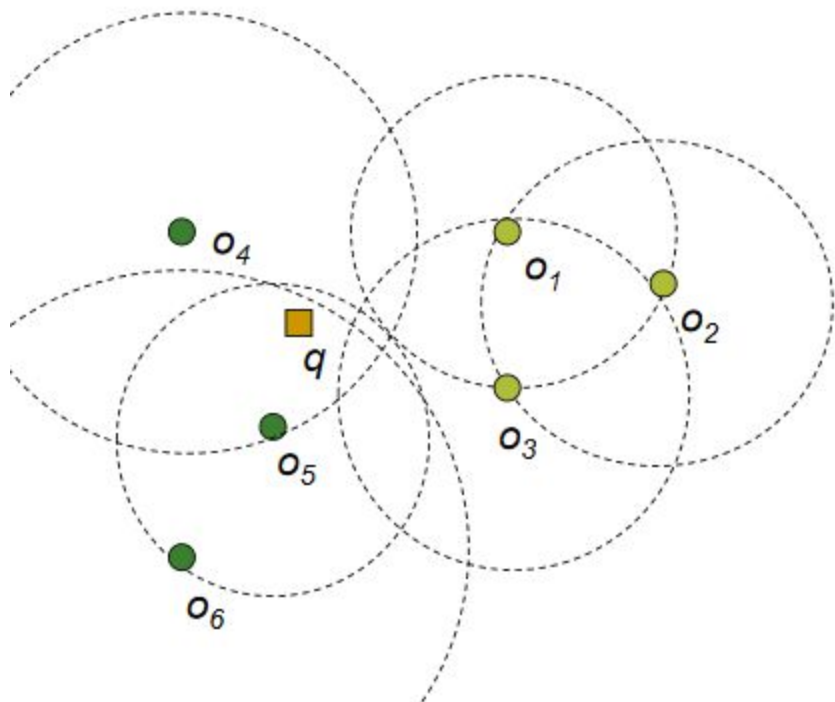
Reversed k-NN query

"vsechny hotely, které mají dané muzeum mezi svými K nejbližšími muzei"

$$kRNN(q) = \{R \subseteq X, \forall x \in R : q \in kNN(x) \wedge \\ \forall x \in X - R : q \notin kNN(x)\}$$

Kazdy prvek **x** který vybereme, musi mit pivota **q** mezi svými *k*-nejblizsimi sousedy a do vysledku zaradime kazde takove **x**. Tj. pokud jej nevybereme tak jen proto ze pivot **q** mezi *k*-nejblizsimi sousedy prvku **x** není. => **VYSLEDEK NEMUSI BYT VELIKOSTI k!!!**

příklad pro 2-RNN



objekty: o4, o5 a o6, mají pivota q mezi svými dvěma nejbližšími sousedy => výsledek nemá velikost k.

Similarity join

"pár hotelů a muzeí, které jsou od sebe 5 minut pěšky"

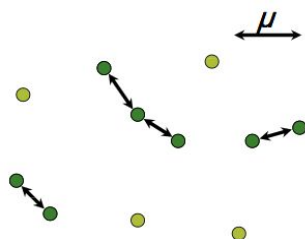
$$X \subseteq \mathcal{D}, Y \subseteq \mathcal{D}, \mu \geq 0$$

$$J(X, Y, \mu) = \{(x, y) \in X \times Y : d(x, y) \leq \mu\}$$

obecně definované na dvou různých množinách (X, Y, muzea, hotely..)

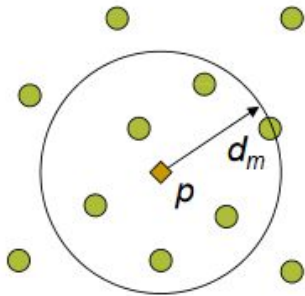
pokud X=Y pak se jedná o Similarity Self-join

μ je vzdálenost, kterou od sebe prvky páru mohou maximálně mít



Ball partitioning

-rozdeluje datovou mnozinu na dve casti podle pivota (p) a radiusu (d_m)



Vnitřní množina (inner set)

$\{x \in X \mid d(p, x) \leq d_m\}$ všechny prvky uvnitř a na hranici koule

Vnější množina (outer set)

$\{x \in X \mid d(p, x) > d_m\}$ všechny prvky mimo kouli

Multi-way ball partitioning

-rozšíření ball partitioningu který domenu rozdeluje pomocí dvou radiusu (d_{m1} , d_{m2}) na 3 množiny:

Vnitřní množina (inner set)

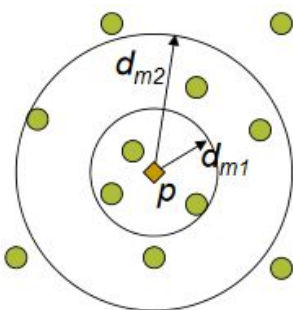
$\{x \in X \mid d(p, x) \leq d_{m1}\}$ všechny prvky ležící uvnitř

Vnější množina (outer set)

$\{x \in X \mid d(p, x) > d_{m2}\}$ všechny prvky ležící vně

Prostřední množina (middle set)

$\{x \in X \mid d(p, x) > d_{m1} \text{ \& } d(p, x) \leq d_{m2}\}$ všechny prvky ve "středním pasmu"



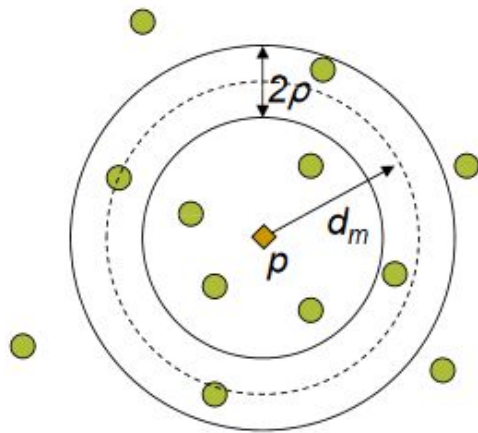
prvky na hranici se berou ze patří dovnitř koule, proto je rovnost u vnitřní a střední množiny

Excluded middle partitioning

-rozsireni bezneho ball partitioningu

motivace:

kvuli objektum ktere lezi blizko hranice muzeme casto byt nuceni pristupovat k obema mnozinam (inner i outer set). Proto datovou mnozinu rozdeline na 3 casti, pricemz minimalne do jedne z nich nebude nutne pristupovat – objekty na hranici zahrneme do teto treti mnoziny.



Vnitřní množina (inner set)

$$\{x \in X \mid d(p, x) \leq d_m - p\}$$

Vnější množina (outer set)

$$\{x \in X \mid d(p, x) > d_m + p\}$$

Vynechaná množina (excluded set)

vsechny ostatni prvky

Generalized hyperplane partitioning

-rozdělování datové množiny pomocí nadrovin

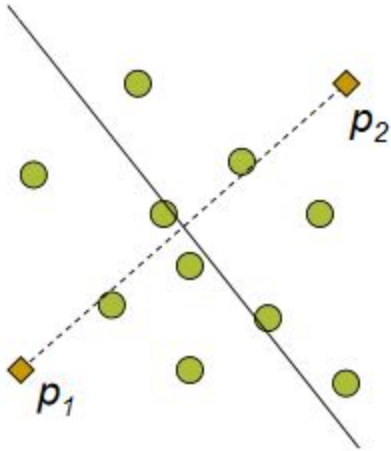
-zvolíme dva pivoty a rozdělíme datovou množinu na dvě, podle toho ke kterému pivotu má který prvek blíže.

Jedna množina:

$$\{x \in X \mid d(p_1, x) \leq d(p_2, x)\} \text{ (rovnost - musíme si zvolit kam padnou objekty na hranici)}$$

Druhá množina:

$$\{x \in X \mid d(p_1, x) > d(p_2, x)\}$$



Pivot filtering

- strategie pouzivana k redukcii poctu kontrolovaných prvku → rychlost ← to je dostatečná motivace ne?
- používá trojúhelníkovou nerovnost pro “odstřel” kontrolovaných prvku
- predpokládáme že máme strukturu ve které jsou spočítány vzdálenosti všech prvku od pivota p

1)příklad filtrování s jedním pivotem

p – pivot

q,r – range query

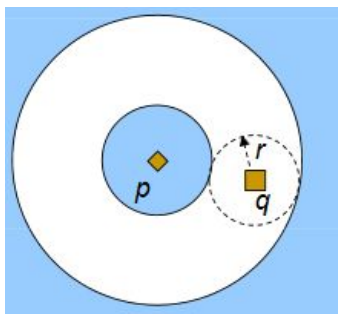
o – kontrolovaný prvek

Algoritmus

odstřel prvek o pokud něco z tohoto platí:

$d(p,o) < d(p,q) - r$ (vnitřní modrý kruh)

$d(p,o) > d(p,q) + r$ (vnější modrá oblast)



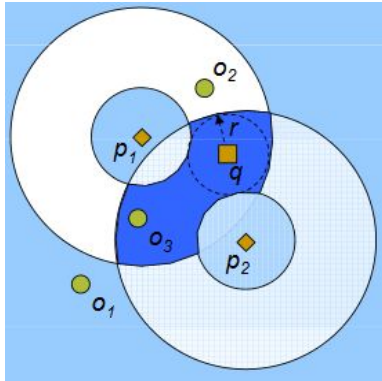
prvky v modré oblasti se nemohou kvalifikovat pro dotaz q,r a to jen díky znalosti jejich vzdálenosti od p

2)příklad filtrovani se dvema pivoty

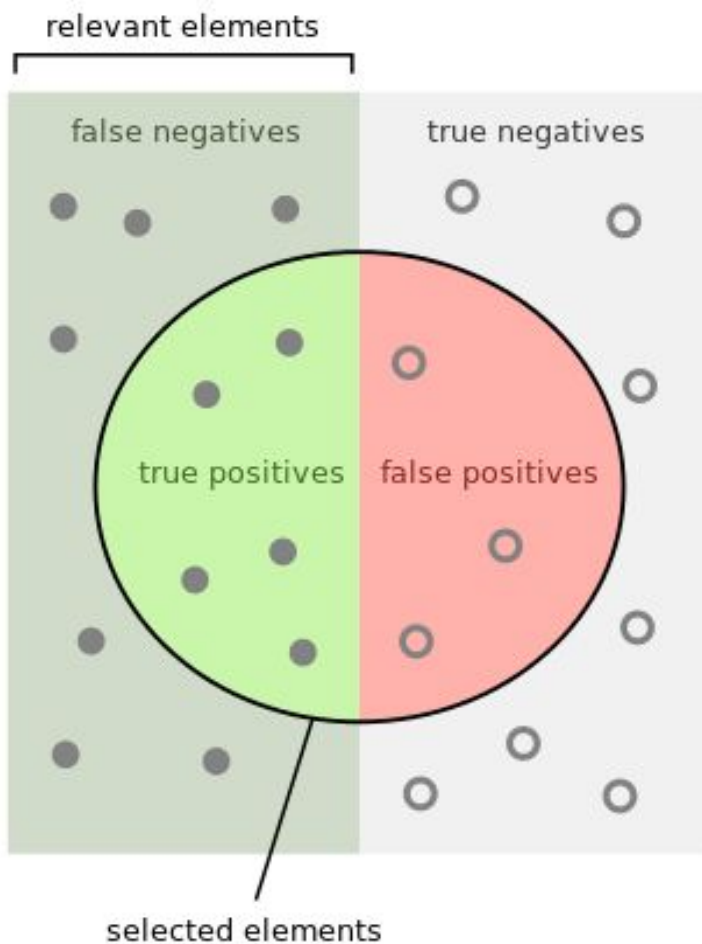
-jen tmave modrou oblast ma cenu uvazovat.

-efektivita se zvysuje s pouzitim vice pivotu

(oblast kterou je potreba kontrolovat se zmensuje – delame prunik)



Precision, Recall



Precision (presnost)

pomer mezi poctem vybranych relevantnich prvku k poctu vseh vybranych.

Jinak receno

jaka cast vybranych prvku je relevantni? (jsou spravnou odpovedi na dotaz)

jeste jinak receno:



rika nam to chybu jakou jsme udelali.

1 → vsechny prvky, které jsme vybrali jsme vybrali opravnene

0.9 → nektère prvky jsme vybrali spatne

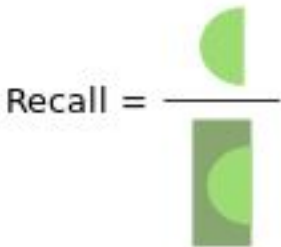
Recall (senzitivita)

pomer mezi poctem vybranych relevantnich prvku ku vsem RELEVANTNIM prvkum

jinak receno:

jakou cast z relevantnich prvku jsme vybrali

jeste jinak receno:



rika nam to jak velkou cast spravne odpovedi na dotaz jsme vybrali.

1 → vsechny prvky které jsou spravne jsme vybrali

0.7 → nektère prvky jsme nevybrali presteze jsou spravne

DUSLEDEK

Muzeme mit velkou **precision** (klidne 1) ale pokud vracime v odpovedi na dotaz jen malo prvku, i kdyz jich tam je relevantnich hodne (mame nizky **recall**) tak vysledek nestoji za nic.

Stejně tak když máme kvalitní **recall** (klidně 1), což znamená že žádný z relevantních prvků ve výsledku nevynecháme, ale zase máme špatnou **precision** (vrátíme skoro všechny prvky – tzn. i spoustu špatných), tak opět výsledek nestojí za nic.

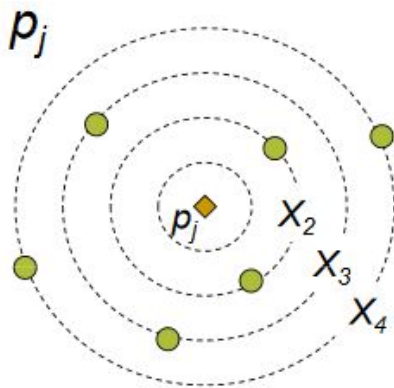
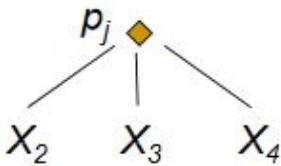
Burkhard-Keller Tree, popsat, nakreslit příklad

2/5

- použitelný pouze při použití diskrétní distance funkce
- rekurzivně dělí datovou množinu
- začneme s výběrem libovolného prvku p_j z datové množiny a sestrojíme množiny:

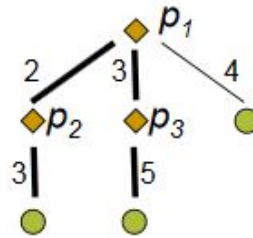
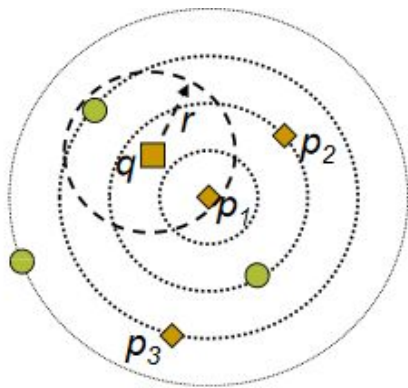
$$X_i = \{o \in X \mid d(o, p_j) = i\}$$

- tzn. že například množina X_2 obsahuje všechny objekty které mají od p_j vzdálenost 2
- z prvku p_j udeláme kořen stromu, z každého X_i se stane podstrom. Prázdné množiny X_i vynecháme (v případě že žádný objekt nemá od p_j vzdálenost i).



Range query $R(q, r)$

- začneme od kořene
- v každém vnitřním uzlu p_j deláme:
 - zarad p_j do výsledku když: $\text{if } d(q, p_j) \leq r$
 - vejdi do potomka p_i když: $\text{if } \max\{d(q, p_i) - r, 0\} \leq i \leq d(q, p_i) + r$



zlute prvky jsou klasifikovane, tluste cary znazornuji uzly které jsme navstivili

Vantage point tree

- Používá ball partitioning
- rekurzivně dělí datovou množinu
- vyvázený binární strom

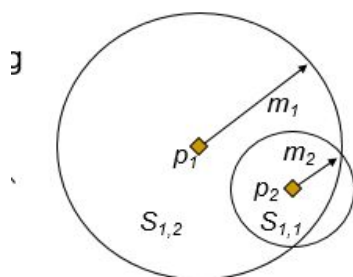
Algoritmus

1. Je potřeba vybrat vztyčný bod (vantage point) a poté určit median(**median ze vzdaleností vzhledem k vantage pointu**)
2. Rozdělíme datovou množinu na dvě podmnožiny podle mediánu

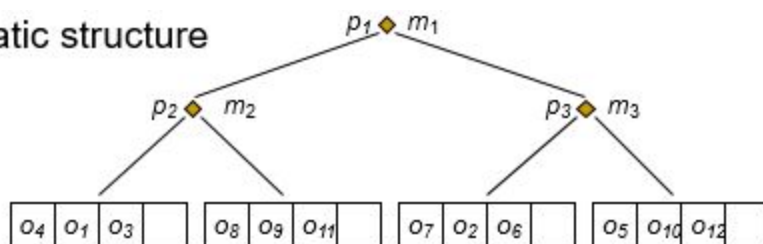
$$S_1 = \{x \in X - \{p\} \mid d(x, p) \leq m\}$$

$$S_2 = \{x \in X - \{p\} \mid d(x, p) \geq m\}$$

- Rovnitko je použito v obou množinách a proto je strom vyvážený
- vantage point (p) je uložen vždy v nelistovém uzlu
- Jeden nebo více objektů jsou umístěny v listech



Static structure



Range query $R(q,r)$

-(uvazujeme pripad ze v kazdem listu je jen jeden objekt)

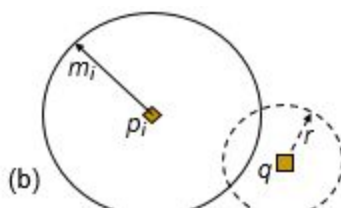
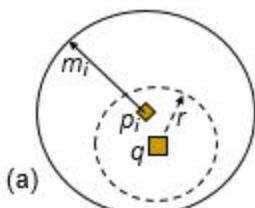
-prochazime strom od korene

-V každém vnitřním uzlu je potřeba udělat:

if $d(q, p_i) \leq r$ zaradime p_i do vysledku

if $d(q, p_i) - r \leq m_i$ prohledame levý podstrom – inner set (a, b)

if $d(q, p_i) + r \geq m_i$ prohledame pravý podstrom – outer set (b)



Jak lze vidět v případě (b), může se stát že bude nutné prohledat levý i pravý podstrom

Nearest neighbor NN(q)

- inicializujeme: $d_{NN} = d_{max}$ $NN = nil$
- (d_{nn} =dosavadní nejmenší vzdálenost od q , NN =prvek s touto vzdáleností)
- procházíme strom od kořene
- v každém interním uzlu (p_i, m_i) , děláme:

if $d(q, p_i) \leq d_{NN}$	nastav $d_{NN} = d(q, p_i)$, $NN = p_i$
if $d(q, p_i) - d_{NN} \leq m_i$	prohlédáme levý podstrom
if $d(q, p_i) + d_{NN} \geq m_i$	prohlédáme pravý podstrom

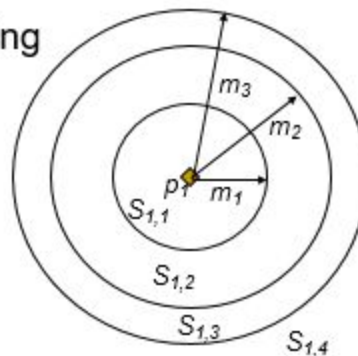
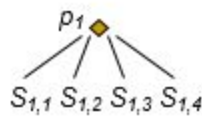
k-NN(q)

- jedina změna je že potřebujeme pole: $d_{NN}[k]$ a $NN[k]$
- dana pole udržujeme setříděná vzhledem k vzdálenosti od pivotu q

Multi Way VPT

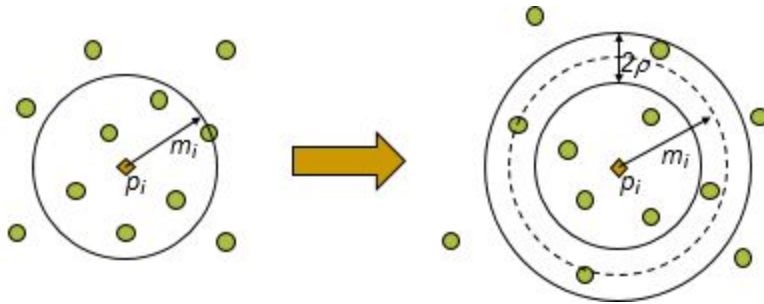
- další typ VPT
- Veškeré principy dědí z klasického VPT
- jedná se o m-ární strom oproti klasickému VPT (který je binární).
- místo bezneho ball partitioningu používáme Multi-way ball partitioning (nemáme jen jeden radius, ale máme jich více => rozdělujeme na více množin, nejen vnitřek a vnějšek)

- applies multi-way ball partitioning

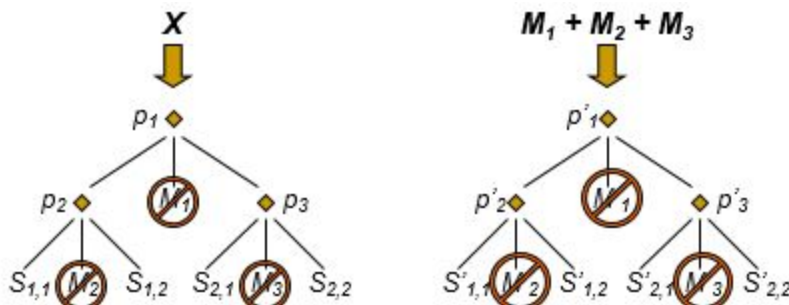


Vantage point forest (VPF)

- „les binárních stromů“
- oproti beznemu VPT (který používá ball partitioning) používá VPF excluded middle partitioning



- excluded middle partitioning rozdeluje datovou množinu na 3 podmnožiny, pričom ale my tu vytvárame binárny strom...Jak? Tak že prostě vytváříme strom jen z dvou podmnožin: vnitru (inner set) a vnějšku (outer set).
- třetí podmnožinu (excluded set - prvky na hranici koule) si dáme bokem a ze všech těchto vynechaných částí pak zkonstruujeme další (opět binární) VPT (rekurzí).



Range query $R(q,r)$

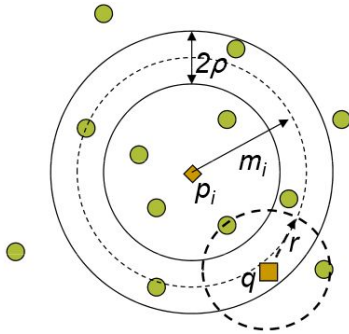
- procházíme strom od kořene
- $\rho(p)$ udává polovici sírky hranice excluded middle setu
- v každém vnitřním uzlu (p_i, m_i) , děláme:

if $d(q, p_i) \leq r$
 if $d(q, p_i) - r \leq m_i - \rho$
 if $d(q, p_i) + r \geq m_i - \rho$
 if $d(q, p_i) + r \geq m_i + \rho$
 if $d(q, p_i) - r \leq m_i + \rho$
 If $d(q, p_i) - r \geq m_i - \rho$
 &&
 $d(q, p_i) + r \leq m_i + \rho$

zaradíme p_i do výsledku
 prohledáme levý podstrom
 prohledáme následující strom
 prohledáme pravý podstrom
 prohledáme následující strom

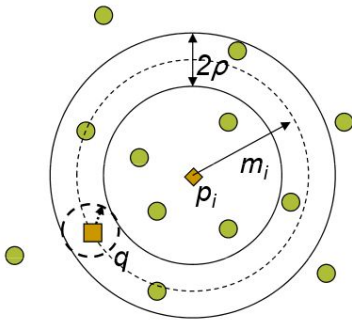
prohledáme POUZE následující strom

PRIKLAD1



Dotaz $R(q, r)$ zasahnuv vsechny 3 množiny \Rightarrow prohledame oba podstromy (levy i pravy) a jeste k tomu jdeme dal hledat v nasledujicim stromu (protoze ten byl vyvoren z excluded middle)

PRIKLAD2



Dotaz $R(q, r)$ zasahl jen prostredni (excluded middle) množinu \Rightarrow levy ani pravy podstrom neprohledavame, jdeme rovnou hledat do nasledujiciho strom (protoze ten byl vyvoren z excluded middle)

Bisector Tree

- zaklad pro GHT (generalized hyperplane tree)
- vyuzivame deleni pomoci nadrovin (viz generalized hyperplane partitioning)

Algoritmus

- zvolime dva pivoty $p_1, p_2 \in X$
- ty nam rozdeli datovou množinu na dve podmnožiny, podle toho, kteremu pivotu ma dany objekt bliz. Delici uzel (nelistovy) ma tedy tvar: $\langle p_1, p_2 \rangle$

$$S_1 = \{o \in X, d(o, p_1) \leq d(o, p_2)\}$$

$$S_2 = \{o \in X, d(o, p_1) > d(o, p_2)\}$$

-zvolíme hodnoty r_1^c a r_2^c tak, aby:

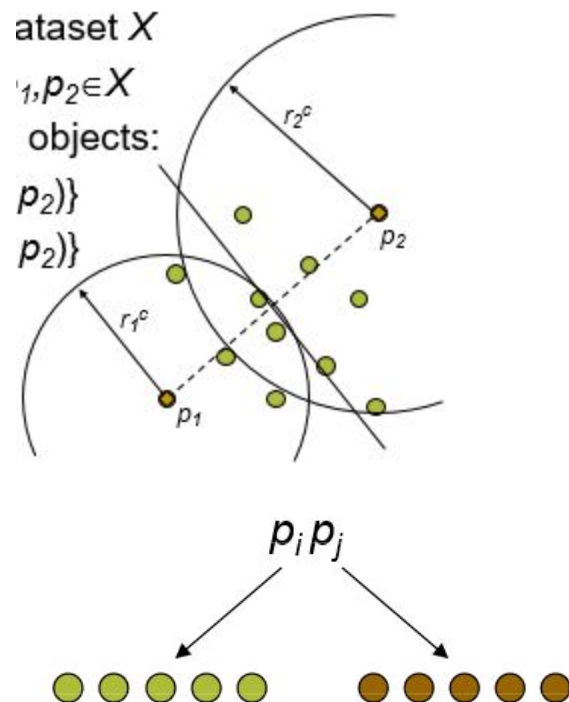
Všechny prvky z množiny S1 byly uvnitř koule se středem v p_1 a poloměrem r_1^c
a

Všechny prvky z množiny S2 byly uvnitř koule se středem v p_2 a poloměrem r_2^c

-zvoleným hodnotám r_1^c a r_2^c se říká 'covering radii'

-koule se mohou překrývat (no problemo)

-rekurzivně opakujeme na množině S1 a S2



Range query $R(q, r)$

-procházíme strom od kořene

-v každém interním uzlu $\langle p_i, p_j \rangle$ děláme tohle (jak pro p_i , tak pro p_j):

Zařadíme p_x do výsledku if $d(q, p_x) \leq r$ (p_x je přímo v r -okolí q)

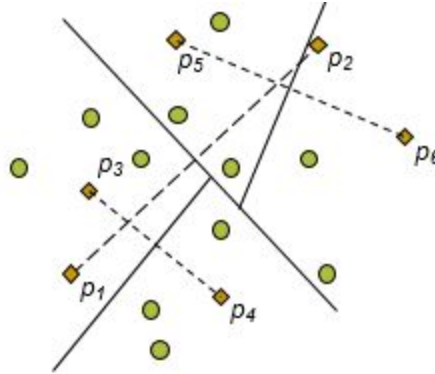
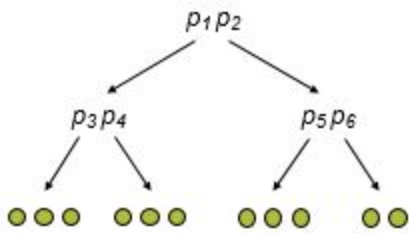
Vlezeme do podstromu p_x if $d(q, p_x) - r \leq r_x^c$ (range query koule zasahuje do covering radiusu $p_x \Rightarrow$ musíme prohledat)

GHT (Generalized hyperplane tree)

-podobný jako bisector tree

-dělení se děje pomocí nadrovin

-je vybrána dvojice objektu a následně je rozdělena datová množina na dvě podmnožiny podle toho ke kterému objektu má kdo blíž.



Range query $R(q, r)$

-procházíme strom od kořene

-V každém vnitřním uzlu $\langle p_i, p_j \rangle$ děláme:

-zaradíme p_x do výsledku

if $d(q, p_x) \leq r$

(jsme v dosahu radiusu)

-jdeme do levého podstromu

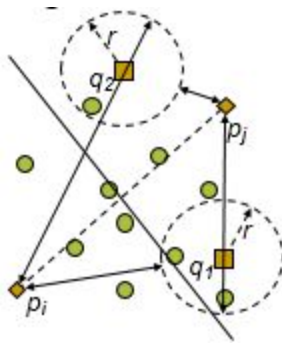
if $d(q, p_i) - r \leq d(q, p_j) + r$

(jsme blíže p_i)

-jdeme do pravého podstromu

if $d(q, p_i) + r \geq d(q, p_j) - r$

(jsme blíže p_j)



M-tree vlastnosti, idea, naco je dobry, jake tam jsou typy uzlu, co je v nich ulozene

2/65

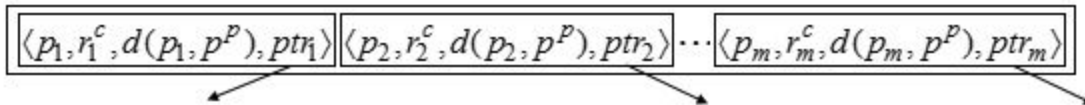
3/6

M-tree

- převážně dynamická struktura
- Disk-oriented (fixní velikost uzlů)
- Vytvořeno pomocí bottom-up architektury
 - Inspirováno R-trees a B-trees
- Veškerá data jsou uložena v listech
- Interní uzly: odkazy na podstromy a další informace
- Podobné jako GNAT, ale objekty jsou uloženy v listech.

M-tree: Vnitřní uzly

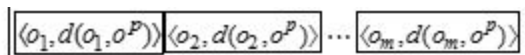
- Vnitřní uzel obsahuje přístup do každého podstromu
- Obsah vnitřního uzlu:
 - Pivot: p
 - Pokrývací poloměr podstromu: r^c
 - Vzdálenost od p do parent pivotu p^p : $d(p, p^p)$
 - Ukazatel na podstrom: ptr



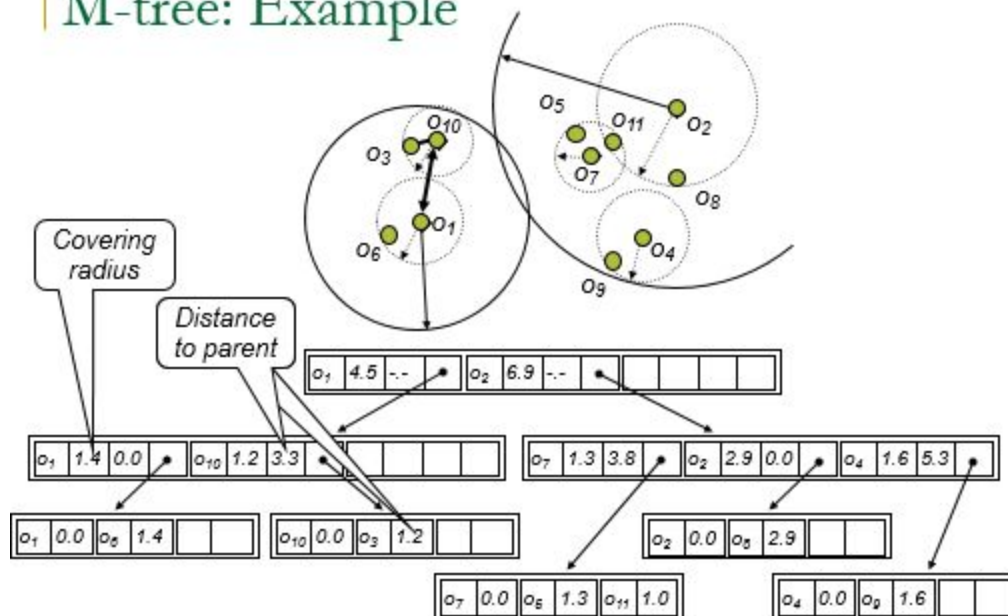
- Všechny objekty v podstromu ptr jsou v mezích vzdálenosti r^c z p .

M-tree: Listové uzly

- Obsahují data
- Každá každý záznam je složený z dvojice:
 - object (jeho identifikátor): o
 - vzdálenost mezi o a jeho nadřazeným pivotem: $d(o, o^p)$



M-tree: Example



M-tree: Vložení prvku

- vložíme nový objekt o_N :
- rekurzivně sestupně projdeme strom a určíme nejlepší list pro uložení o_N
- v každém kroku vstoupíme do podstromu s pivotem p pro který:
 - není potřeba zvětšit poloměr r^c i.e., $d(o_N, p) \leq r^c$
 - v případě, že existuje více listů, vybereme ten s nejbližším p v o_N

- minimalizujeme zvětšení r^c
- když dosáhneme listu N pak je třeba rozlišit tyto situace:
 - if N není plný then ulož o_N do N
 - else **Split**(N, o_N).

M-tree: Split

Split(N, o_N):

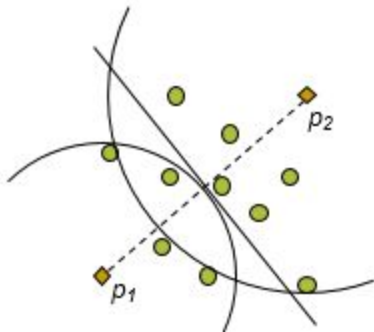
- Necht' S je množina obsahující všechny prvky z N a o_N
- Vybereme dva pivoty p_1 a p_2 z S
- Partition S do S_1 a S_2 v závislosti na p_1 a p_2
- Ulož S_1 do N a S_2 do nově vytvořeného uzlu N'
- If N is root
 - Allocate a new root and store entries for p_1, p_2 there
- else (let N^p and p^p be the parent node and parent pivot of N)
 - Replace entry p^p with p_1
 - If N^p is full, then Split(N^p, p_2)
 - else store p_2 in node N^p

M-tree: Split Policy

- Partition S to S_1 and S_2 according to p_1 and p_2

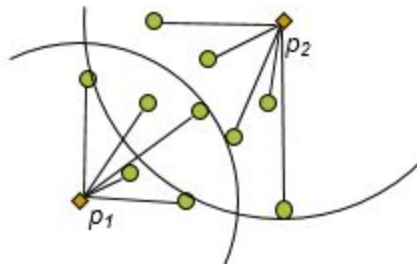
■ Unbalanced

- Generalized hyperplane



■ Balanced

- Larger covering radii
- Worse than unbalanced one



Leaf node splitting a metody jakými se to dělá, rozdíl mezi splitováním v Slim a M-tree

3/36

Distance-index (D-index)

3/72

-Hybridní struktura pro velké databáze.

-Kombinace dvou věcí: pivot filtering a partitioning