

### **7.1. Objasnite zašto se relacijske baze podataka ne smatraju dobrim rješenjem za pohranu velikih podataka (Big Data).**

- spajanje tablica je resursno zahtjevno (potreba denormalizacije čime gubimo 3. normalnu formu), nisu dizajnirane za horizontalnu skalabilnost, ocuvanje konzistentnosti postaje usko grlo (visoka razina konzistentnosti = stroga politika zaključivanja i kontrole pristupa), sekundarni indeksi usporavaju sustav kod čestih izmjena (uklanjanje istih može znatno usporiti upite)

### **7.2. Navedite prednosti stupcano-orijentiranih baza podataka**

- smanjuje se broj I/O operacija, bolja kompresija podataka, pogodno za tablice sa velikim brojem stupaca, učinkovito upravljanje rijetkim tablicama, mogu preskociti NULL vrijednosti

### **7.3. Objasnite koja je od sljedećih tvrdnji točna:**

a) podaci pohranjeni u bazu podataka HBase slijede formalno definiranu shemu,

b) baza podataka HBase dizajnirana je za laku vertikalnu skalabilnost,

c) baza podataka HBase omogućuje nisku latenciju rada s podacima i brz pristup individualnim zapisima i

d) podaci pohranjeni u bazu podataka HBase u pravilu slijede treću normalnu formu

- b) krivo (horizontalnu), a) krivo (dodaju se obitelji stupaca), d) krivo (denormaliziran oblik) c) točno

### **7.4. Objasnite način na koji se podaci fizički pohranjuju kod baze podataka HBase**

- HBase je mapa zapravo

- tablica = skup redaka, redak = vrijednosti unutar skupa obitelji stupaca, obitelj stupaca = skup stupaca, stupac = skup parova ključ-vrijednost

- Podaci se perzistiraju uz pomoć neizmjenjivih datoteka - HFile

- HFile = uređena mapa parova ključ -> vrijednost poredanih leksikografski po id-u redaka, skup parova je pohranjen u niz blokova, na kraju datoteke se pohranjuje indeks blokova pomoću kojeg HBase pronalazi traženi zapis

- kod izmjene podataka ne izmjenjuje se HFile, nego se piše prvo u HLog ili WAL (write-ahead log) a zatim se pohranjuje u memoriju (memstore), posluživanje informacija združivanjem memstora i HFile-a

- kad se memstore prepuni stvara se novi HFile uz pomoć WAL-a

### **8.1. Napišite psudokod osnovnog algoritma preporuke po suradnji korisnika**

- for every object o that user c has no preference for yet

    for every other user c' that has a preference for o

        compute a similarity s between c and c'

        add c's preference for o, weighted by s, into a running average of o

return the top objects, ranked by weighted average

**8.2. Kako osnovni algoritam preporuke temeljene na sadržaju možemo napisati u obliku matrica? Što predstavljaju elementi svake od tih matrica?**

$$U \times S = E$$

[U11 ... U1m] [S11 ... S1m] [e11 ... e1m]

[..... Uix .....] x [..... Sxj .....] = [..... eij .....]

[Un1 ... Unm] [Sm1 ... Smm] [en1 ... enm]

- uix - koliko je korisniku ci koristan objekt ox
- sxj - koliko je objekt ox sličan objektu oj
- eij - kolika je procijenjena korisnost objekta oj korisniku ci

**8.3. Objasnite kako biste ostvarili preporucivanje temeljeno na suradnji u slučaju velikog broja korisnika.**

- pri izracunu preporuke za korisnika c se uzima u obzir samo njemu slični korisnici, tako se smanjuje broj operacija na  $2nm|C'|$ , inace je  $2mn^2$

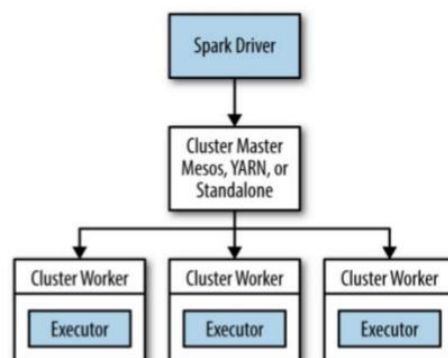
**8.4. Objasnite zašto je potrebno raditi normalizaciju matrice sličnosti objekata (S') u slučaju preporuke temeljene na sličnosti objekata (item-based recommendation)**

- jer zbog velike kolicine objekata je broj operacija jako velik, te ga je potrebno smanjiti (sta ja znam, lupam bzvz)

**9.1. Objasnite zašto kažemo da su osnovni operatori nad tokom podataka blokirajući. Kako ih možemo odblokirati?**

- da blokiraju isporuku rezultata (koja je max vrijednost beskonacnog slijeda?), mozemo ih odblokirati uvevši prozore koji racuna max vrijednost na temelju zadnjih x vrijednosti

**9.2. Skicirajte i objasnite raspodijeljeno izvođenje Sparkove aplikacije**



Driver je proces u kojem se izvodi glavna metoda aplikacije. Zadužen za podjelu posla na zadatke. Iz usmjerenog acikličkog grafa operacija stvara plan izvođenja. Koordinira raspoređivanje poslova na izvodace. Vrsi raspoređivanje poslova na izvodace temeljeno na prostornoj lokalnosti.

Izvodaci izvođe dodeljene zadatke i u memoriji pohranjuju podatke neophodne za izvođenje aplikacije.

### 9.3. Objasnite razliku između dviju vrsta operacija nad RDD-om (Resilient Distributed Dataset).

- transformacije - pretvaranja jednog toka u drugi
- izlazne operacije - operacije nad tokom podataka koje daju neki rezultat

### 9.4. Objasnite način na koji Spark obrađuje tok podataka prividno stvarnovremeno. Koji parametri su bitni za definiranje vremenskog prozora u Sparkovoj aplikaciji?

- stvarnovremenost je prividna jer se podaci iz tokova obrađuju u mikroskupinama
- windowDuration - koliko zadnjih prozora "pamti"
- slideDuration - koliko svakih prozora se pomice

### 10.1. Navedite i objasnite dva moguća načina zapisa podataka o strukturi grafa te navedite koji je pogodniji za zapis velikih mreža te objasnite zašto.

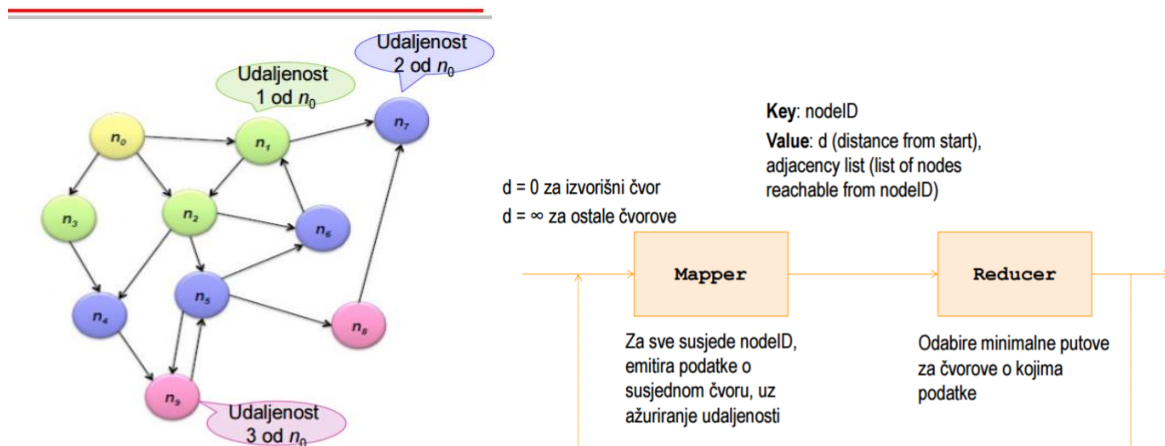
- matrica susjedstva i lista susjedstva, lista je pogodnija jer nemora svaki čvor u svojoj listi imati sve ostale čvorove nego samo one sa kojima je povezan dok je matrica full velicine  $N \times N$

### 10.2. Navedite kako možete odrediti izlazne i ulazne grane pojedinog čvora grafa ako za zapis koristite listu susjedstva

- id čvora čiju listu gledamo je čvor iz kojeg ide strelica, a čvor u njegovoj listi je čvor prema kojem ide strelica

### 10.3. Skicirajte algoritam za paralelno izvođenje Dijkstrinog algoritma u okolini MapReduce

- u jednoj iteraciji MapReduce analiziraju se čvorovi na udaljenosti 1 od izvorista, zatim na udaljenosti 2 itd



**10.4. Objasnite ideju mjeru popularnosti PageRank. Navedite kako je moguće izračunati PageRank na iterativni način u okolini MapReduce.**

- pageRank - karakterizira količinu vremena koju korisnik provede na nekoj web stranici
- mjera popularnosti web-resursa neovisna o txt sadržaju dokumenta
- ako stranica i ima hiperlink na stranicu j, tada i implicitno prenosi važnost j
- $x_i = \sum (1/N_j) * x_j \rightarrow$  rang stranice i = rang stranice j kroz broj izlaznih linkova iz j

**11.1. Što je strojno učenje?**

- strojno učenje je polje računalne znanosti koje se bavi specifičnim načinom programiranja u kojem očekujemo da računalo samostalno dođe do određenih spoznaja na osnovu dostupnih skupova i odabrane metode učenja

**11.2. Objasnite metodu linearne regresije.**

- pretpostavka je da između prediktora i predikcija postoji linearna veza
- traži se pravac kod kojeg će zbroj najmanjih kvadrata odstupanja ciljeva od pravca biti najmanji
- pogreska procjene = srednja vrijednost kvadrata odstupanja od pravca

**11.3. Ukratko navedite korake u procesu dubinske analize podataka**

- prikupljanje i integracija podataka iz različitih izvora
- definicija problema kojeg se želi riješiti provedbom dubinske analize
- priprema, čišćenje i transformacija dobivenog podatkovnog skupa
- eksploratorna analiza podataka
- odabir metoda za stvaranje prediktivnih i/ili deskriptivnih modela
- treniranje i podešavanje modela
- evaluacija i usporedba modela
- priprema izvještaja o rezultatima analize

**11.4. Navedite prednosti i nedostatke korištenja Spark + MLib arhitekture**

- prednosti: veliki podatkovni skup a postojeća analitička platforma nema dovoljno kapaciteta, snažno integracijski orijentirani, transparentno obavlja zadane zadatke na velikom skupu podataka kao da se radi o lokalnom radu, može poslužiti kao alternativni priručni stroj za obradu podataka kada se lokalni režim pokazuje nedovoljan, integracija sa sparkom otvara pristup novim funkcionalnostima kao što su Spark streaming ili stvarnovremensko izvođenje upita nad velikim skupom podataka (Spark SQL)

- nedostaci: nova tehnologija, teško procijeniti dugoročnu stabilnost API-ja, integracijska rješenja i dalje nisu na visokoj razini glede jednostavnosti, stabilnosti i transparentnosti, nepotpune funkcionalnosti, problem dugoročne stabilnosti i održivosti napisanog programskog koda, integracija sa postojećim rješenjima se može pokazati disruptivnom

**12.1. Opišite postupak dobivanja malog svijeta iz rešetke modelom Watts-Strogatz. Objasnite koja dva parametra i kako moramo pri tome promatrati da mali svijet ne bi postao slučajna mreža.**

- model Watts-Strogatz - veze u rešetki se nasumično prespajaju (vjerojatnost prespajanja  $p$ )
- precaci drastično smanjuju prosjecni najkraci put
- mali svijet - mreza s visokim koeficijentom grupiranja i malim prosjecnim najkracim putom
- moramo promatrati vjerojatnost prespajanja i koeficijent grupiranja (valjda, ne znam)

**12.2. Skicirajte proizvoljnu mrežu od 4 cvora te na primjeru jednog cvora navedite formulu objasnite što je koeficijent grupiranja cvora (local clustering coefficient). Kao iz koeficijenta grupiranja cvora možemo dobiti koeficijent grupiranja društvene mreže (network average clustering coefficient)?**

- koeficijent grupiranja cvora  $i$ ,  $C_i$ , predstavlja omjer broja međusobnih veza među susjedima cvora  $i$ ,  $E_i$ , i ukupnog broja njihovih mogućih međusobnih veza  $C_i = 2E_i / (z_i(z_i - 1))$
- koeficijent grupiranja nam kaže koliko su gusto povezani prvi susjedi nekog cvora
- koeficijent grupiranja mreže je prosjecni koeficijent grupiranja cvorova u mreži  $C = (1/N) * \sum(C_i)$

**12.3. Objasnite zašto je centralnost cvora bitna za širenje informacija u mreži. Kako definiramo centralnost po položaju (betweenness centrality)?**

- centralnost pokušava identificirati najvažnije cvorove u grafu
- centralnost po položaju, centralnost po blizini, po svojstvenom vektoru, po stupnju...

**12.4. U RDD-u prVertices se nalazi graf s pagerank-om cvorova društvene mreže (id cvora tipa Long i pagerank cvora tipa Double). Nadopunite sljedeći programski kod da ispišete razdiobu vrijednosti pagerank-a cvorova od većeg prema manjem kao niz redaka u obliku: pagerank cvora, broj cvorova s tim pagerank-om.**

```
JavaRDD<Tuple2<Long, Double>> prVertices = getVertexPagerank();
JavaPairRdd<Double, Long> prCounts = prVertices
    .mapToPair(pair -> new Tuple2<>(pair._2, 1))
    .reduceByKey((x, y) -> x+y)
    .sortByKey();
prCounts.foreach( d -> System.out.println("\n" + d._1 + ", " + d._2));
```