

```
1 Indicacions sobre les solucions d'alguns problemes del segon parcial i el f
  inal
2 Gener 2015
3
4 Exercici 1 del parcial, 4 del final
5
6 Amb una sola fase de mapreduce, suposant que es dona una línia
7 a cada instància de mapreduce)
8
9 map(string linia)
10     l = l.split(" ")
11     autor = l[0];
12     l = l[1..final];
13     eliminar repetits de l;
14     per cada x de l, output (autor,x)
15
16 combiner(autor,l)
17     eliminar repetits de l;
18     output (autor,l);
19
20 reduce(autor,l)
21     // L és una llista de llistes
22     eliminar repetits de l;
23     output (autor,longitud(l));
24
25
26 Eliminar repetits al map i al combiner
27 és opcional, però molt convenient per eficiència.
28 Notem que un combiner rep elements de molts maps,
29 i per tant és *més* eficient fer-ho a tots dos
30 llocs que només al combiner.
31
32 Una alternativa, probablement més eficient, és
33 no emetre una tupla per cada comentari, sinó una
34 per cada post:
35
36 map(string linia)
37     l = l.split(" ")
38     autor = l[0];
39     l = l[1..final];
40     eliminar repetits de l;
41     output (autor,l);
42
43 combiner(autor,L)
44     // L és una llista de llistes
45     fusionar totes les llistes de L
46     i eliminar repetits; sigui l el resultat
47     output (autor,l);
48
49 reduce(autor,L)
50     // L és una llista de llistes
51     fusionar totes les llistes de L
52     i eliminar repetits; sigui l el resultat;
53     output (autor,l);
54
55 Exercici 4 del parcial
56
57 Manera 1: executar els dos algorismes separatament,
58 i combinar els ratings (p.ex. fent la mitjana, triant el màxim)
59 en un nou ranking, o bé triant els top k/2 de cadascun
```

```
60 i per donar k recomanacions.
61
62 Manera 2: executar primer CB per seleccionar un conjunt
63 d'items prometedors, i llavors fer CF només sobre aquest subset.
64 O en l'ordre invers (es consideren el mateix mètode).
65
66 En aquest problema hi ha hagut la tendència de donar solucions
67 força vagues.
68
69 Exercici 5 del parcial, 6 del final
70
71 No s'hi valien declaracions generals que no eren
72 escenaris sinó problemes, en especial si només
73 es repetia el que diuen les descripcions de spark
74 en les primeres línies
75 - "per a problemes de machine learning": hadoop va molt
76 bé per a molts problemes de machine learning. Problemes
77 de machine learning en temps real és una altra cosa.
78 - "per a problemes molt grans": idem, hadoop va bé
79 per a problemes molt grans també.
80 - "per a problemes de molta velocitat, spark diu que
81 pot ser fins a 100 vegades més ràpid que hadoop"...
82 en alguns problemes, es tractava de saber quins i dir-ho
83 (o si no se sap, buscar altres raons).
84
85 Escenaris concrets podria ser:
86
87 - anàlisi de seguretat en temps real
88 - anàlisi de sèries financeres en temps real
89 - recomanacions on productes i preferències varien en el temps
90 - servei transaccional alhora que anàlisi de tendències
91 en temps real d'un lloc com Twitter
92 - ...
93
94 Exercici 1 del final:
95 1. fals (els tf no cal recalcular-los)
96 2. cert
97 3. cert
98 4. fals
99 5. (no comptava, però fals)
100 6. cert
101 7. fals (de fet, el que es diu no té sentit)
102 8. fals (no té solució)
103 9. fals: afegir damping factor NO arregla gens aquest problema.
104 se seguirà perdent pagerank i només la solució amb
105 tots els pageranks 0 és factible
106 10. cert
107
108
109
110
111
```