

MCMD L16 : MapReduce | triangle count

MapReduce

M = Massive Data

Mapper(M) \rightarrow {(key,value)}

Shuffle({(key,value)}) \rightarrow group by "key"

Reducer ({"key,value_i"}) \rightarrow ("key, f(value_i))

Can repeat, constant # of rounds

Given graph $G=(V,E)$

Assume $|V|=n$ and $|E| = m = n^{\{1+c\}}$
typical large graphs have c in $[0.08, 0.5]$

EXAMPLE NUMBERS (via SNAP)

	n	m	c
FaceBook subset:	4,039	88,234	0.37
Twitter subset:	81,306	1,768,149	0.27
GPlus subset:	107,614	13,673,453	0.42
Live Journal:	4,847,571	68,993,773	0.17
Friendster:	65,608,336	1,806,067,135	0.185

$N(v)$ = neighbors of v

cluster coefficient $cc(V)$
= fraction $N(v)$, neighbors themselves
How dense a subgraph is

** need to find all triangles for each v in V **

[based on paper by Sure + Vassilvitskii 11]

```
(sequential)
for each v in V
  for each (u,w) in N(v)
    if (u,w) in E  $\rightarrow$  Triangle[v]++
```

$T = \sum_{\{v \text{ in } V\}} |N(v)|^2$
 $O(n^2)$ if some v $N(v) = O(n)$

(parallel)

Map 1: $G=(V,E) \rightarrow (v,u),(u,v)$ for (v,u) in E

Reduce 1: $(v, N(v)) \rightarrow ((u,w),v)$ s.t. u,w in $N(v)$

Map 2: $\rightarrow ((u,w),v)$ (output of R1)
 $\rightarrow ((u,w),\$)$ for (u,w) in E

Reduce 2: $((u,w),\{v_1,v_2,v_3,\dots,v_t,\$\})$
iff $\$,$ then $\rightarrow (v_i,1/3)$

Map 3: identity
Red 3: aggregate

:(running time still $\max_{\{v \text{ in } V\}} |N(v)|^2$

LiveJournal

80% reducers done in 5 min
99% reducers done in 35 min
some 60 minutes

Idea 1: count each triangle once, with lowest degree

(sequential)
for each v in V
 for each (u,w) in $N(v)$
 if $\text{deg}(u) > \text{deg}(v) \ \&\& \ \text{deg}(w) > \text{deg}(v)$
 if (u,w) in $E \rightarrow \{\text{Tri}[v]++, \text{Tri}[u]++, \text{Tri}[w]++\}$

In Reduce 1, add if condition.
In Reduce 2, $\rightarrow (v_i,1)$
 $\rightarrow (u,1) , (w,1)$

Works better!
Last reducer at around 2 minutes

two types of nodes:
 $L = \{v \mid N(v) \leq \sqrt{m}\}$
 $H = \{v \mid N(v) > \sqrt{m}\}$

$|L| \leq n \rightarrow$ produce $O(m)$ pairs
 $n * \sqrt{m} = m \rightarrow n < \sqrt{m}$ $n=m^a \quad m = m^a * m^b$ where $a+b < 1$
 $b < 1/2$
 $n * (\sqrt{m})^2 = O(m^{3/2})$ work = $m^a + m^{2b}$
 $|H| \leq 2\sqrt{m} \rightarrow$ produce $O(m)$ pairs
 $2 \sqrt{m} (2 \sqrt{m})^2 = O(m^{3/2})$ since only pairs with H

if $m = O(n^2)$ (very dense)
 $n \sim \sqrt{m}$
→ $O(m^{\{3/2\}})$ work (optimal!)

Idea 2 : Graph Split

partition V into p equal-size sets $\{V_1, V_2, \dots, V_p\}$

For triples (V_i, V_j, V_k) → subgraph $G_{\{ijk\}} = G[V_i + V_j + V_k]$

 compute triangles on $G_{\{ijk\}}$

 triangles counted $\{1, p-2, \text{ or } (p-1 \text{ choose } 2) \sim p^2\}$ times

 figure out and adjust

subgraph has $O(m/p^2)$ edges in expectation

 → $V_{\{i,j,k\}}$ has $3n/p$ vertices, random edge in set w.p. $9/p^2$

work: $p^3 * O((m/p^2)^{\{3/2\}}) = O(m^{\{3/2\}})$

total space used: $O(pm)$ expected

p about 20 worked best on LiveJournal graph