

Домашна задача – Hadoop и Map-Reduce

Група 37

Опис на проблемот

Податоците од log датотеките можат да послужат за стекнување на многу информации кои се од големо значење. Во склоп на оваа задача ние сакаме да извлечеме корисни информации обработувајќи log датотека превземена од материјалите на udacity.org од курсот за Hadoop. Датотеката содржи 44477844 линии во следниот формат:

Ip_address	--	timestamp	"request type	File "	bytes	cs
На пример:						
10.223.157.186	--	[15/Jul/2009:15:50:36]	"GET	RagingPhoenix_2DSleeve.jpeg "	200	168

Цел: да ги излистаме сите IP адреси од кои бил отворен секој фајл во текот на еден месец.

Ова значи дека сакаме да добиеме излез во следниот формат:

month	file	[ip_address_1	ip_address2	ip_address_3	...]
-------	------	---------------	-------------	--------------	------

Секој клик на веб страна креира околу 100 бајти податоци во типичен лог. Како последица на ова може да се генерираат стотици гигабајти податоци за еден ден. Издвојувањето на информациите кои ни требаат од оваа огромна маса податоци побарува повеќе од обичен секвенцијален алгоритам.

Изворен код

Бидејќи целта ни е клуч во мапата да се месецот и пристапениот фајл, креиравме класа FileAccess која ги содржи истите. За секој ред во влезната датотека, секој mapper креира објект од FileAccess како клуч и му доделува вредност еднаква на IP адресата од која бил направен пристапот. Потоа reducer-от ги зема сите IP адреси кои имаат ист клуч и ги става во еден TreeSet.

```
public class Logs {
    public static class FileAccess implements Comparable<FileAccess> {
        private Text file, month;
        public FileAccess() {
            this.file = new Text();
            this.month = new Text();
        }
        public void set(Text file, Text month) {
            this.file = file;
            this.month = month;
        }
        public int compareTo(FileAccess o) {
            if (file.compareTo(o.file) == 0) return (month.compareTo(o.month));
            else return (file.compareTo(o.file));
        }
        public String toString() {
            return month.toString() + "\t" + file.toString();
        }
    }
    public static class Map extends Mapper<LongWritable, Text, FileAccess, Text> {
        private final FileAccess fileAccess = new FileAccess();
        private Text file = new Text();
        private Text ipAddress = new Text();
        private Text month = new Text();
        public void map(LongWritable key, Text value, Context context)
            throws IOException, InterruptedException {
            final String[] data = value.toString().trim().split(" ");
            if (data.length > 6) {
                ipAddress.set(data[0]);
                String[] datetime = data[3].trim().split(":");
                if (datetime.length > 0 && datetime[0].length() > 1) {
                    String date = datetime[0].substring(1);
                    file.set(data[6]);
                    String[] nums = date.split("/");
                }
            }
        }
    }
}
```

