

```

import java.io.IOException;
import java.util.ArrayList;
import java.util.Iterator;
import java.util.List;

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.io.Writable;
import org.apache.hadoop.io.WritableComparable;
import org.apache.hadoop.mapred.FileInputFormat;
import org.apache.hadoop.mapred.FileOutputFormat;
import org.apache.hadoop.mapred.JobConf;
import org.apache.hadoop.mapred.KeyValueTextInputFormat;
import org.apache.hadoop.mapred.Mapper;
import org.apache.hadoop.mapred.MapReduceBase;
import org.apache.hadoop.mapred.OutputCollector;
import org.apache.hadoop.mapred.RecordReader;
import org.apache.hadoop.mapred.Reducer;
import org.apache.hadoop.mapred.Reporter;
import org.apache.hadoop.mapred.SequenceFileInputFormat;
import org.apache.hadoop.mapred.SequenceFileOutputFormat;
import org.apache.hadoop.mapred.TextInputFormat;
import org.apache.hadoop.mapred.JobControl;
import org.apache.hadoop.mapred.lib.IdentityMapper;

public class MRExample {
    public static class LoadPages extends MapReduceBase
        implements Mapper<LongWritable, Text, Text, Text> {

        public void map(LongWritable k, Text val,
            OutputCollector<Text, Text> oc,
            Reporter reporter) throws IOException {
            // Pull the key out
            String line = val.toString();
            int firstComma = line.indexOf(',');
            String key = line.substring(0, firstComma);
            String value = line.substring(firstComma + 1);
            Text outKey = new Text(key);
            // Prepend an index to the value so we know which file
            // it came from.
            Text outVal = new Text("1" + value);
            oc.collect(outKey, outVal);
        }

        public static class LoadAndFilterUsers extends MapReduceBase
            implements Mapper<LongWritable, Text, Text, Text> {

        public void map(LongWritable k, Text val,
            OutputCollector<Text, Text> oc,
            Reporter reporter) throws IOException {
            // Pull the key out
            String line = val.toString();
            int firstComma = line.indexOf(',');
            String value = line.substring(firstComma + 1);
            int age = Integer.parseInt(value);
            if (age < 18 || age > 25) return;
            String key = line.substring(0, firstComma);
            Text outKey = new Text(key);
            // Prepend an index to the value so we know which file
            // it came from.
            Text outVal = new Text("2" + value);
            oc.collect(outKey, outVal);
        }

        public static class Join extends MapReduceBase
            implements Reducer<Text, Text, Text, Text> {

        public void reduce(Text key,
            Iterator<Text> iter,
            OutputCollector<Text, Text> oc,
            Reporter reporter) throws IOException {
            // For each value, figure out which file it's from and
            // accordingly.
            List<String> first = new ArrayList<String>();
            List<String> second = new ArrayList<String>();

            while (iter.hasNext()) {
                Text t = iter.next();
                String value = t.toString();
                if (value.charAt(0) == '1')
                    first.add(value.substring(1));
                else second.add(value.substring(1));
            }

            public static void main(String[] args) throws IOException {
                JobConf lp = new JobConf(MRExample.class);
                lp.setJobName("Load Pages");
                lp.setInputFormat(TextInputFormat.class);

                reporter.setStatus("OK");

                // Do the cross product and collect the values
                for (String s1 : first) {
                    for (String s2 : second) {
                        String outVal = key + "," + s1 + "," + s2;
                        oc.collect(null, new Text(outVal));
                        reporter.setStatus("OK");
                    }
                }
            }

            public static class LoadJoined extends MapReduceBase
                implements Mapper<Text, Text, Text, LongWritable> {

            public void map(
                Text k,
                Text val,
                OutputCollector<Text, LongWritable> oc,
                Reporter reporter) throws IOException {
                // Find the url
                String line = val.toString();
                int firstComma = line.indexOf(',');
                int secondComma = line.indexOf(',', firstComma);
                String key = line.substring(firstComma, secondComma);
                // drop the rest of the record, I don't need it anymore,
                // just pass a 1 for the combiner/reducer to sum instead.
                Text outKey = new Text(key);
                oc.collect(outKey, new LongWritable(1L));
            }

            public static class ReduceUrls extends MapReduceBase
                implements Reducer<Text, LongWritable, WritableComparable,
                Writable> {

            public void reduce(
                Text key,
                Iterator<LongWritable> iter,
                OutputCollector<WritableComparable, Writable> oc,
                Reporter reporter) throws IOException {
                // Add up all the values we see
                long sum = 0;
                while (iter.hasNext()) {
                    sum += iter.next().get();
                    reporter.setStatus("OK");
                }
                oc.collect(key, new LongWritable(sum));
            }

            public static class LoadClicks extends MapReduceBase
                implements Mapper<WritableComparable, Writable, LongWritable,
                Text> {

            public void map(
                WritableComparable key,
                Writable val,
                OutputCollector<LongWritable, Text> oc,
                Reporter reporter) throws IOException {
                    oc.collect((LongWritable)val, (Text)key);
            }

            public static class LimitClicks extends MapReduceBase
                implements Reducer<LongWritable, Text, LongWritable, Text> {

            int count = 0;
            public void reduce(
                LongWritable key,
                Iterator<Text> iter,
                OutputCollector<LongWritable, Text> oc,
                Reporter reporter) throws IOException {
                // Only output the first 100 records
                while (count < 100 && iter.hasNext()) {
                    oc.collect(key, iter.next());
                    count++;
                }
            }

            public static void main(String[] args) throws IOException {
                JobConf lp = new JobConf(MRExample.class);
                lp.setJobName("Load Pages");
                lp.setInputFormat(TextInputFormat.class);

                JobConf lfu = new JobConf(MRExample.class);
                lfu.setJobName("Load and Filter Users");
                lfu.setInputFormat(TextInputFormat.class);
                lfu.setOutputKeyClass(Text.class);
                lfu.setOutputValueClass(Text.class);
                lfu.setMapperClass(LoadAndFilterUsers.class);
                FileInputFormat.addInputPath(lfu, new
                    Path("/user/gates/users"));
                FileOutputFormat.setOutputPath(lfu,
                    new Path("/user/gates/tmp/filtered_users"));
                lfu.setNumReduceTasks(0);
                Job loadUsers = new Job(lfu);

                JobConf join = new JobConf(MRExample.class);
                join.setJobName("Join Users and Pages");
                join.setInputFormat(KeyValueTextInputFormat.class);
                join.setOutputKeyClass(Text.class);
                join.setOutputValueClass(Text.class);
                join.setMapperClass(IdentityMapper.class);
                join.setReducerClass(Join.class);
                FileInputFormat.addInputPath(join, new
                    Path("/user/gates/tmp/indexed_pages"));
                FileInputFormat.addInputPath(join, new
                    Path("/user/gates/tmp/filtered_users"));
                FileOutputFormat.setOutputPath(join, new
                    Path("/user/gates/tmp/joined"));
                join.setNumReduceTasks(50);
                Job joinJob = new Job(join);
                joinJob.addDependingJob(loadPages);
                joinJob.addDependingJob(loadUsers);

                JobConf group = new JobConf(MRExample.class);
                group.setJobName("Group URLs");
                group.setInputFormat(KeyValueTextInputFormat.class);
                group.setOutputKeyClass(Text.class);
                group.setOutputValueClass(LongWritable.class);
                group.setOutputFormat(SequenceFileOutputFormat.class);
                group.setMapperClass(LoadJoined.class);
                group.setCombinerClass(ReducerUrls.class);
                group.setReducerClass(ReducerUrls.class);
                FileInputFormat.addInputPath(group, new
                    Path("/user/gates/tmp/joined"));
                FileOutputFormat.setOutputPath(group, new
                    Path("/user/gates/tmp/grouped"));
                group.setNumReduceTasks(50);
                Job groupJob = new Job(group);
                groupJob.addDependingJob(joinJob);

                JobConf top100 = new JobConf(MRExample.class);
                top100.setJobName("Top 100 sites");
                top100.setInputFormat(SequenceFileInputFormat.class);
                top100.setOutputKeyClass(LongWritable.class);
                top100.setOutputValueClass(Text.class);
                top100.setOutputFormat(SequenceFileOutputFormat.class);
                top100.setMapperClass(LoadClicks.class);
                top100.setCombinerClass(LimitClicks.class);
                top100.setReducerClass(LimitClicks.class);
                FileInputFormat.addInputPath(top100, new
                    Path("/user/gates/tmp/grouped"));
                FileOutputFormat.setOutputPath(top100, new
                    Path("/user/gates/top100sitesforusers18to25"));
                top100.setNumReduceTasks(1);
                Job limit = new Job(top100);
                limit.addDependingJob(groupJob);

                JobControl jc = new JobControl("Find top 100 sites for users
                    18 to 25");
                jc.addJob(loadPages);
                jc.addJob(loadUsers);
                jc.addJob(joinJob);
                jc.addJob(groupJob);
                jc.addJob(limit);
                jc.run();
            }
        }
    }
}

```