Setup Yahoo SAMOA on S4

On Virtual Machine with Ubuntu Linux

Li Huang

Computer Engineering and Computer Science

J.B. Speed School of Engineering

University of Louisville

l0huan08@gmail.com

2014/1/29 Edit: 2014/2/5

Statement:

This article only reflects the author's personal opinion, and may contain some views not right. If you find something wrong in this article or have advices please contact the author. This article is for non-commercial purpose, which helps the people who are interested to SAMOA. Some figures are referenced from some other sources.

Contents

1.	Setup Environment	2
	Preparation	
	Setup SAMOA	
	Test SAMOA	
	3.1 Simple Test Cluster Evaluation	4
	3.2 Forest Cover Type Classification	
4	Recommend Readings	

SAMOA(Scalable Advanced Massive Online Analysis) is a new distributed streaming data mining platform which is now being developing by Yahoo, and S4 is one of the support platform for running SAMOA(another is Twitter Storm). Since currently there is very few material can be find to introducing the Setup process of SAMOA and S4, in this article I will introduce my setup process of Yahoo! SAMOA based on S4. This article assumes you have already setup the Ubuntu system and Yahoo S4. For the detail of setup Ubuntu and S4, please see my another article "Yahoo SAMOA Whitepaper 2-Setup Yahoo S4".

Comparing with setup S4, setup SAMOA is much easier! The whole process takes about 1 hour.

1. Setup Environment

I use my daily-use laptop to setup SAMOA and S4. Its configuration is below:

Hardware:

Computer: Lenovo Laptop U310

CPU: Intel Core i3-2367M, with 2 cores at 1.40Ghz

Memory(RAM): 2.00 GB System Type: 64 bit

Software:

Base operating system: Windows 7 Professional Service Pack 1 - 64 bit

Virtual Machine to host S4: VMware Player 5.0.2

Operating system on Virtual Machine: Ubuntu Linux Desktop 12.04 - 64 bit

Stream Platform: Yahoo S4 0.6.0

1. Preparation

Before setup SAMOA, you need setup a Linux and a Stream Process Engine (SPE) in your computer. For operating system, I setup Ubuntu 12.04 desktop 64bit in a virtual machine. Currently SAMOA supported Twitter Storm and Yahoo S4. In this article, I use Yahoo S4 as my SPE, since I failed to setup Storm on my virtual machine (it is hard to configuration).

Note: Another reason I gave up Storm is that: To test Storm in Cluster environment, we need at least 3 nodes--a nimbus(center server) and 2 computing nodes, but my computer has not so much memory for running 3 VMs at same time.

2. Setup SAMOA

This process is simple and as same as the official document. The only things you need to note are:

Firstly you need install a tool called "maven" to build SAMOA.

SAMOA support three modes: local, S4, Storm. Local mode run on a single machine without any SPE; S4 mode run SAMOA on S4; Storm mode run SAMOA on Storm.

The process of setup SAMOA:

1. Setup maven. (I use maven 3.1.1. It is a tool to build SAMOA)

apt-get install maven

Or: find "maven" in "ubuntu software center"

Or: download it from the official website

2. Download SAMOA from its website

Go to the dir "cd /opt ".

Clone SAMOA 0.0.1 from GitHub.

git clone git://github.com/yahoo/samoa.git

And set its write permission to all user by "sudo chmod -R 777 /opt/samoa"

3. In the "samoa" folder,

Run "mvn package",

to build the software package of "local mode" will be build. The result file is

"target/SAMOA-Local-0.0.1-SNAPSHOT.jar"

Run "mvn -Ps4 package"

to build the softwarepackage of "S4 mode" will be build. The result file is

"target/SAMOA-S4-0.0.1-SNAPSHOT.jar"

Note: SAMOA local mode do not need any SPE, but can only run a single machine rather than a cluster.

```
2
                           hl@hlnode1: /etc
                                                                  + _ □ X
 File
      Edit View
                Search Terminal
                               Help
[INFO] Copying jmockit-1.5.jar to /opt/samoa/samoa-s4/target/lib/jmockit-1.5.jar
[INFO] Copying samoa-api-0.0.1-SNAPSHOT.jar to /opt/samoa/samoa-s4/target/lib/sa
moa-api-0.0.1-SNAPSHOT.jar
[INFO]
[INFO]
      --- maven-assembly-plugin:2.4:single (make-assembly) @ samoa-s4 ---
[INFO] Reading assembly descriptor: src/main/assembly/samoa-s4.xml
     Building jar: /opt/samoa/target/SAMOA-S4-0.0.1-SNAPSHOT.jar
INFO]
INF0]
INFO] Reactor Summary:
[INFO]
INFO] SAMOA ...... SUCCESS [7.319s]
INFO] samoa-api ...... SUCCESS
                                                             [6.649s]
INFO] samoa-s4 ...... SUCCESS [30.980s]
INFO]
INFO] BUILD SUCCESS
INF01
[INFO] Total time: 49.889s
[INFO] Finished at: Wed Jan 29 01:30:33 EST 2014
[INFO] Final Memory: 22M/54M
INF0]
hl@hlnode1:/opt/samoa$
```

4. Edit configure file, "<samoa folder>/bin/samoa-s4.properties".

Same as the official documents, but the document is not clear with this part:

- # Deployment strategy
 samoa.deploy.mode=local
- Note: When I saw this setting at first time, I did not know what was it mean. Now I know, it means that: If you are run SAMOA in a single (virtual) machine, you should set deploy mode as "local". When you want to run SAMOA in several (virtual)machines, you need to change "local" into "cluster"
- 5. Add SAMOA to environment file by add these line to the end of "/etc/profile":

```
export SAMOA_HOME=/opt/samoa
```

And run "source /etc/profile" to update your environment variables.

3. Test SAMOA

After setup SAMOA, you should test SAMOA. We have two simple test app for SAMOA, "Cluster Evaluation" and "Forest Cover Type". So far (by 2014/1/29).

Note

SAMOA Bug: However, SAMOA 0.0.1 has a bug that you can only run these two examples in "local" mode, but cannot run them on S4 mode. Currently the developers are trying to solve this problem.

3.1 Simple Test -- Cluster Evaluation

This task will do a cluster and evaluate the performance. This example can be found in the last part of :

https://github.com/yahoo/samoa/wiki/Executing-SAMOA-with-Apache-S4

Syntax: bin/samoa <platform> <jar-location> <task & options>

(1) Run in Local Mode

You can run in local mode:

bin/samoa local target/SAMOA-Local-0.0.1-SNAPSHOT.jar "ClusteringEvaluation"

You will see output result below:

```
hl@hlnode1: /opt/samoa
                                                                               ↑ _ □ X
      Edit View Search Terminal
                                     Help
hl@hlnode1:/opt/samoa$ bin/samoa local target/SAMOA-Local-0.0.1-SNAPSHOT.jar
usteringEvaluation'
bin/samoa
Deploying to LOCAL
Command line string = ClusteringEvaluation
[main] INFO com.yahoo.labs.samoa.LocalDoTask - Sucessfully instantiating com.yah
oo.labs.samoa.tasks.ClusteringEvaluation
[main] INFO com.yahoo.labs.samoa.learners.clusterers.LocalClustererProcessor -
rained model using 1000 events with classifier id 1
[main] INFO com.yahoo.labs.samoa.learners.clusterers.LocalClustererProcessor -
rained model using 1000 events with classifier id 0
[main] INFO com.yahoo.labs.samoa.learners.clusterers.LocalClustererProcessor -
rained model using 2000 events with classifier id 1
[main] INFO com.yahoo.labs.samoa.learners.clusterers.LocalClustererProcessor -
rained model using 2000 events with classifier id 0
[main] INFO com.yahoo.labs.samoa.learners.clusterers.LocalClustererProcessor - rained model using 3000 events with classifier id 1
[main] INFO com.yahoo.labs.samoa.learners.clusterers.LocalClustererProcessor -
rained model using 3000 events with classifier id 0
[main] INFO com.yahoo.labs.samoa.learners.clusterers.LocalClustererProcessor - rained model using 4000 events with classifier id 1
[main] INFO com.yahoo.labs.samoa.learners.clusterers.LocalClustererProcessor -
```

```
hl@hlnode1: /opt/samoa
        Edit
                      Search
                              Terminal
  File
              View
[main] INFO com.yahoo.labs.samoa.learners.clusterers.LocalClustererProcessor -
 ained model using 50000 events with classifier id 0
[main] INFO com.yahoo.labs.samoa.evaluation.EvaluatorProcessor - last event is r
eceived!
[main] INFO com.yahoo.labs.samoa.evaluation.EvaluatorProcessor - total count: 9
[main] INFO com.yahoo.labs.samoa.evaluation.EvaluatorProcessor - com.yahoo.labs.
samoa.evaluation.EvaluatorProcessorid = 0
evaluation instances,SSQ,van Dongen,Rand statistic
0.0,20.129163178128866,0.7962790697674419,0.8867091074444589
1000.0,20.03877460326178,0./616504854368932,0.8586902167246927
2000.0,15.895277164346854,0.5974304068522484,0.7548343020291628
3000.0,19.866134960895106,0.7133726647000983,0.7746014602157031
4000.0,28.631325987623143,0.640625,0.7221671693235674
5000.0,7.576155342546159,0.7707292707292708,0.846001998001998
5000.0,13.624210160959136,0.7117003367003367,0.7781919163553536
7000.0,21.683682246891404,0.7345505617977528,0.7629708412750229
8000.0,14.718398888658736,0.6845238095238095,0.8125018885788822
9000.0,15.758046878655744,0.7526595744<u>68085</u>1,0.8267854782986275
[main] INFO com.yahoo.labs.samoa.evaluation.EvaluatorProcessor - total evaluatio
n time: 11 seconds for <u>9</u> instances
hl@hlnode1:/opt/samoa$
```

(2) Run in S4 Mode

You can also run it in \$4 mode:

bin/samoa S4 target/SAMOA-S4-0.0.1-SNAPSHOT.jar "ClusteringEvaluation"

Important Note:

However, running ClusteringEvaluation in **S4 mode** has a bug (by January.29.2014). Please see https://github.com/yahoo/samoa/issues/33 to find the detail.

In my computer, it shows error message like:

```
hl@hlnode1: /etc
                                                                                                                                                                                                     ↑ _ □ X
    File
                  Edit
                                                 Search
                                View
                                                                      Terminal
                                                                                              Help
149) ~[samoa-api-0.0.1-SNAPSHOT.jar:na]
                     at com.yahoo.labs.samoa.evaluation.ClusteringEvaluatorProcessor.process(
ClusteringEvaluatorProcessor.java:101) ~[samoa-api-0.0.1-SNAPSHOT.jar:na]
                     at com.yahoo.labs.samoa.evaluation.ClusteringEvaluatorProcessor.process(
ClusteringEvaluatorProcessor.java:91) ~[samoa-api-0.0.1-SNAPSHOT.jar:na]
at com.yahoo.labs.samoa.topology.impl.S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4ProcessingItem.onEvent(S4Proces
                     at OverloadDispatcher668.dispatchEvent(Unknown Source) ~[na:na]
                     at org.apache.s4.core.ProcessingElement.handleInputEvent(ProcessingEleme
nt.java:461) ~[s4-core-0.6.0-incubating.jar:0.6.0-incubating]
                     at org.apache.s4.core.Stream$StreamEventProcessingTask.run(Stream.java:3
33) ~[s4-core-0.6.0-incubating.jar:0.6.0-incubating]
at org.apache.s4.comm.staging.BlockingThreadPoolExecutorService$Runnable
WithPermitRelease.run(BlockingThreadPoolExecutorService.java:178) ~[s4-comm-0.6.
0-incubating.jar:0.6.0-incubating]
                     at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.
java:1145) ~[na:1.7.0_51]
                     at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor
  java:615) ~[na:1.7.0_51]
                     at java.lang.Thread.run(Thread.java:744) ~[na:1.7.0_51]
```

3.2 Forest Cover Type Classification

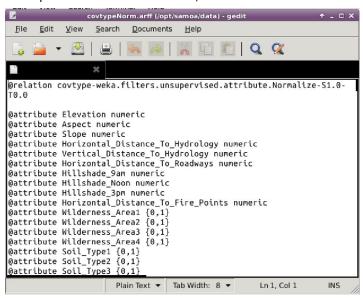
The second example does a classification job on a datasource contains many points of a forest. This example can be found in:

https://github.com/yahoo/samoa/wiki/Getting%20Started

1. Download dataset and unzip it to "/opt/samoa/data":

http://downloads.sourceforge.net/project/moa-datastream/Datasets/Classification/covtypeNorm.arff.zip

If we open this file with a text editor, we can see the file like this:



- 2. classifying the CoverType dataset with the bagging algorithm.
- Local Mode

```
bin/samoa local target/SAMOA-Local-0.0.1-SNAPSHOT.jar
"PrequentialEvaluation -l classifiers.ensemble.Bagging -s
(ArffFileStream -f /opt/samoa/data/covtypeNorm.arff) -f 100000"
```

The output will be a list of the evaluation results, plotted each 100,000 instances.

The output will be a list of the evaluation of classification performance results, plotted each 100,000 instances. As the figure below shows:

You can also add "-d `pwd`/result.csv" in the command to output the result into a .csv file.

S4 mode

You can also try it in **S4 mode**, by typing:

bin/samoa S4 target/SAMOA-S4-0.0.1-SNAPSHOT.jar

"PrequentialEvaluation -l classifiers.ensemble.Bagging -s

(ArffFileStream -f /opt/samoa/data/covtypeNorm.arff) -f 100000"

Note:

* The command parameters can be seen at:

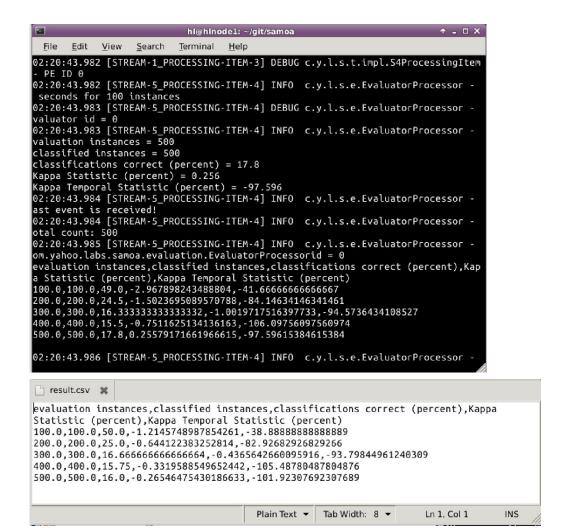
https://github.com/yahoo/samoa/wiki/Prequential%20Evaluation%20Task

- **-f** 100000 means the result will be plotted every 100,000 instances
- * You can add -i 200000 to limit the maximum number of instances to test/train on (as 200000)

I recommend use this command below, which only evaluate 500 instances to save your time, and output the result file to "result.csv":

```
bin/samoa S4 target/SAMOA-S4-0.0.1-SNAPSHOT.jar
"PrequentialEvaluation -1 classifiers.ensemble.Bagging -s
(ArffFileStream -f `pwd`/covtypeNorm.arff) -f 100 -i 500 -d
`pwd`/result.csv"
```

You will see the result like below, with the result.csv in your SAMOA folder:



This result.csv files shows the **performance** (correct rate, Kappa Statistic) of classification of the first 500 instances.

4. Recommend Readings

- Official website of Yahoo SAMOA, including a document tell how to setup SAMOA on S4: http://yahoo.github.io/samoa/
 https://github.com/yahoo/samoa/wiki/Executing-SAMOA-with-Apache-S4
- Official website of Install Maven

http://maven.apache.org/download.cgi#Installation