SQL Elapsed Time Collection Method Comparision Analysis

Author: Craig Shallahamer (craig@orapub.com), Version 1a 22-July-2011

Background and Purpose

The purpose of this notepad is to investigate three methods of gathering SQL elapsed times; SQL trace, source code instrumentation, and OraPub's Elapsed Time Sampler (beta v4).

Experimental Data

Below is all the experimental data. The experiment was run on a Dell single four-core CPU, Oracle 11.2G. According to "cat /proc/version": Linux version 2.6.18-164.el5PAE (mockbuild@ca-build10.us.oracle.com) (gcc version 4.1.2 20080704 (Red Hat 4.1.2-46)) #1 SMP Thu Sep 3 02:28:20 EDT 2009. There was a tremendous CBC latch contention load, the OS was CPU bottlenecked at 100% utilization. The sample set interval was 5 minutes.

The order of sample data is elapsed time (seconds).

The three sets are:

In[1]:=

ssToolV3 data was gathered using the OraPub SQL Elapsed Sampler version 3a. **ssIntrumentedV1** data was gathered using time gathered just before and after the sql statement was executed each time. **ssTraceV4** data was gathered using a shell script to parse through a trace file, which only had data for the specific sqlid under observation.

```
ssToolV3 = {1.767666, 1.518642, 1.518561, 1.784412, 1.808626, 2.009722, 1.518732, 1.520677, 1.51853,
   1.764539, 1.771699, 1.525678, 1.810673, 1.777627, 6.89338, 2.288473, 2.72366, 2.964467,
   2.489718, 1.333654, 1.518478, 2.245655, 1.827613, 2.008537, 2.002821, 4.410587, 2.731605,
   2.542492, 2.241653, 2.01862, 2.009639, 1.762681, 1.555664, 1.759525, 1.569629, 1.283757,
   2.735525, 2.000659, 1.518548, 2.002676, 2.723594, 1.759653, 2.491581, 2.05056, 2.735544,
   2.000499, 2.001705, 2.725593, 3.513618, 2.516472, 2.022676, 1.759562, 1.518653, 2.26362,
   4.42956, 2.241625, 1.759642, 1.521695, 1.759607, 1.763891, 2.000644, 2.73557, 1.762558,
   1.840116, 2.000518, 2.241591, 2.482646, 1.767594, 1.525654, 1.283704, 2.241546, 2.753659,
   1.759537, 2.010594, 1.760708, 1.758692, 2.009638, 2.290598, 2.492627, 2.24156, 1.767647};
ssInstrumentedV1 = {1.743327, 1.664634, 1.620656, 1.621973, 2.28016, 1.809512, 1.911768, 1.635249,
   1.626635, 1.80816, 1.656637, 1.789185, 1.624409, 1.832635, 1.803162, 7.255142, 2.330535,
   2.830604, 3.104761, 2.507755, 1.632316, 1.691193, 2.12539, 1.835945, 1.943198, 1.90315,
   4.694471, 2.863778, 2.841089, 2.205007, 1.881308, 1.859121, 1.753866, 1.687689, 1.627337,
   1.715177, 1.586555, 2.636796, 2.024178, 2.281286, 2.023567, 2.630976, 1.875818, 2.555945,
   2.224021, 2.8679, 1.929357, 2.196423, 2.724408, 3.699935, 2.974219, 1.964847, 1.86812, 2.092287,
   2.467177, 4.894718, 2.202867, 1.911851, 1.608843, 1.725006, 1.771269, 2.288353, 2.800475,
   1.588898, 1.891323, 2.437246, 2.552358, 2.457707, 1.683852, 1.632384, 1.715668, 2.061057,
   2.701076, 1.858633, 1.800178, 1.782631, 1.631529, 2.14331, 2.284766, 2.628907, 2.59879};
ssTraceV4 = {1.660351, 1.618885, 1.620529, 2.079340, 1.808099, 1.908893, 1.633803, 1.625207,
   1.805905, 1.654301, 1.786923, 1.622185, 1.831073, 1.801623, 7.077461, 2.329004, 2.829103,
   3.102503, 2.505485, 1.630009, 1.688804, 2.123024, 1.833691, 1.940932, 1.901513, 4.493080,
   2.860241, 2.838904, 2.202768, 1.879802, 1.857663, 1.752419, 1.686225, 1.625099, 1.673690,
   1.584356, 2.634501, 2.022832, 2.128306, 2.022080, 2.629479, 1.874318, 2.554511, 2.222008,
   2.865609, 1.927211, 2.095340, 2.722211, 3.598867, 2.972704, 1.962327, 1.865951, 1.866360,
   2.266289, 4.397910, 2.201021, 1.909574, 1.606611, 1.722349, 1.769086, 2.286116, 2.699535,
   1.586294, 1.889070, 2.435653, 2.550875, 2.356181, 1.682314, 1.630129, 1.713436, 2.059333,
   2.699629, 1.857081, 1.797896, 1.780331, 1.629258, 2.140887, 2.282376, 2.553857, 2.596437};
```

Basic Statistics

In this section I calculate the basic statistics, such as the mean and median. My objective is to ensure the data has been collected and entered correctly and also to compare the two datasets to see if they appear to be different.

```
In[4]:=
       ssTool = ssToolV3;
       ssInstrumented = ssInstrumentedV1;
       ssTrace = ssTraceV4:
      myData =
         {"Trace", Mean[ssTrace], Median[ssTrace], StandardDeviation[ssTrace], Length[ssTrace]
         },
         {"Instr", Mean[ssInstrumented], Median[ssInstrumented],
          StandardDeviation[ssInstrumented], Length[ssInstrumented]
         },
         {"Tool", Mean[ssTool], Median[ssTool], StandardDeviation[ssTool], Length[ssTool]
         }
        3
       toGrid = Prepend[myData, {"Method", "Mean", "Median", "Std Dev", "Samples"}];
       Grid[toGrid, Frame \rightarrow All]
Out[7]=
       {{Trace, 2.19921, 1.91839, 0.795444, 80},
        {Instr, 2.22215, 1.92936, 0.834282, 81}, {Tool, 2.12768, 2.00064, 0.786436, 81}}
```

	Method	Mean	Median	Std Dev	Samples
9]=	Trace	2.19921	1.91839	0.795444	80
	Instr	2.22215	1.92936	0.834282	81
	Tool	2.12768	2.00064	0.786436	81

Sample Comparison Tests (when normality does NOT exist)

If our sample sets are **not normally distributed**, we can not perform a simple t-test. We can perform what are called location tests. I did some research on significance testing when non-normal distributions exists. I found a very nice reference:

http://www.statsoft.com/textbook/nonparametric-statistics

The paragraph below (which is from the reference above) is a key reference to what we're doing here:

...the need is evident for statistical procedures that enable us to process data of "low quality," from small samples, on variables about which nothing is known (concerning their distribution). Specifically, nonparametric methods were developed to be used in cases when the researcher knows nothing about the parameters of the variable of interest in the population (hence the name nonparametric). In more technical terms, nonparametric methods do not rely on the estimation of parameters (such as the mean or the standard deviation) describing the distribution of the variable of interest in the population. Therefore, these methods are also sometimes (and more appropriately) called parameter-free methods or distribution-free methods.

Being that I'm not a statistician but still need to determine if these sample sets are significant different, I let *Mathematica* determine the appropriate test. Notice that one of the above mentioned tests will probably be the test *Mathematica* chooses.

Note: If we run our normally distributed data through this analysis (speically, the "LocationEquivalenceTest"), *Mathematica* should detect this and use a more appropriate significant test, like a t-test.

Here we go with the hypothesis testing (assuming our sample sets are not normally distributed):

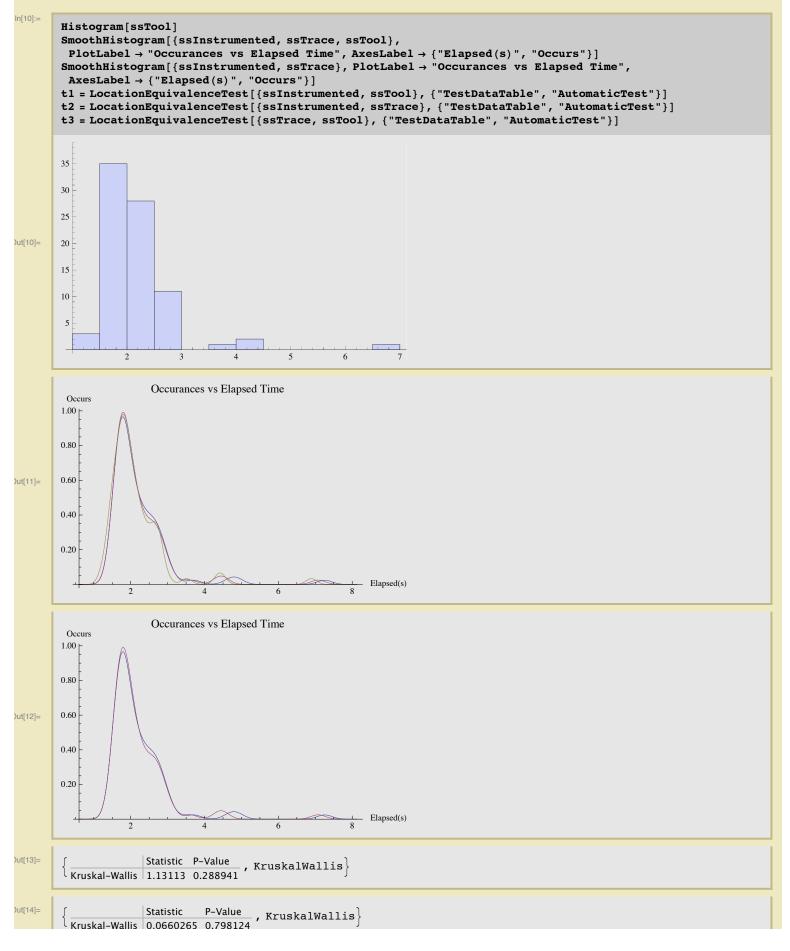
1. Our P value threshold is 0.05, which is our alpha.

2. The null hypotheses is the two populations have the same mean. (Remember we have to sample sets, which is not the population.)

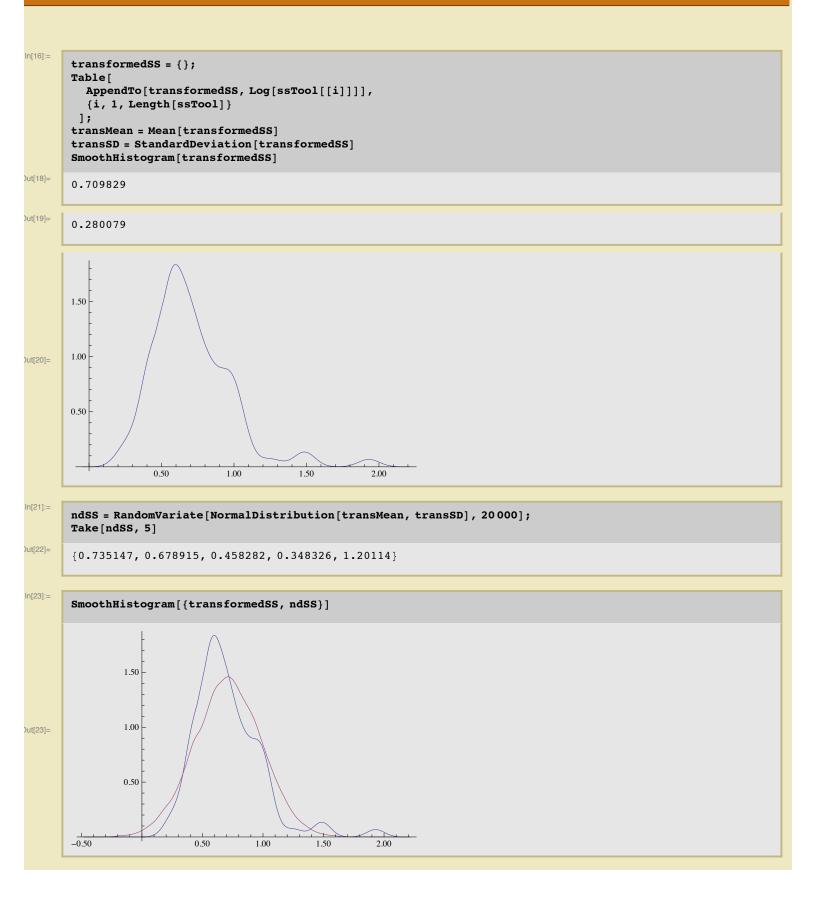
Do the statistical test to compute the P value.

Out

Elapsed hypomes is and say the difference between our samples is significant. (Which is what I'm hoping to see.) However, if the P³ value is greater than the threshold, we cannot reject the null hypothesis and any difference between our samples are not statistically significant; randomness, picked the "wrong" samples, etc.



Log Normal Distribution Fit Test



Dut[24]=