

From the post (<http://datascience.la/benchmarking-random-forest-implementations/>) we can see that Tree methods, especially the Random Forest, in Spark consumes more time and gives lower accuracy compare with those counter-parts in H2O and XGBoost. Since XGBoost is written in C++ with better tuned server performance, we'll spend most of our effort comparing Spark with H2O this time. We'll try to figure out that to what extend (speed and accuracy) does H2O better than Spark, why those advantages exist and how to fix Spark.

Random forest is a combination of decision trees. Decision tree, as a basic learner, suffers from relatively lower accuracy due to its high variance. Gradient boosted tree (GBT) reduces the variance, in the meantime it reduces the bias of a single tree. While random forest (RF) only reduces the variance through a method called "combination" or "bagging". Usually a GBT could reach the best performance among GBT, RF, and decision tree. However, the RF is faster and easy to parallelize in comparison with GBT.

In order to answer the question we mentioned above, we'll first to reproduce the experiment result of that post. Note that interestingly, the blog author uses different hyper-parameters in different methods, which leads to uneven results. So after recreating the previous result, we will try to construct a series of experiments bottom-up, with the same hyper-parameters as long as possible.

Dataset

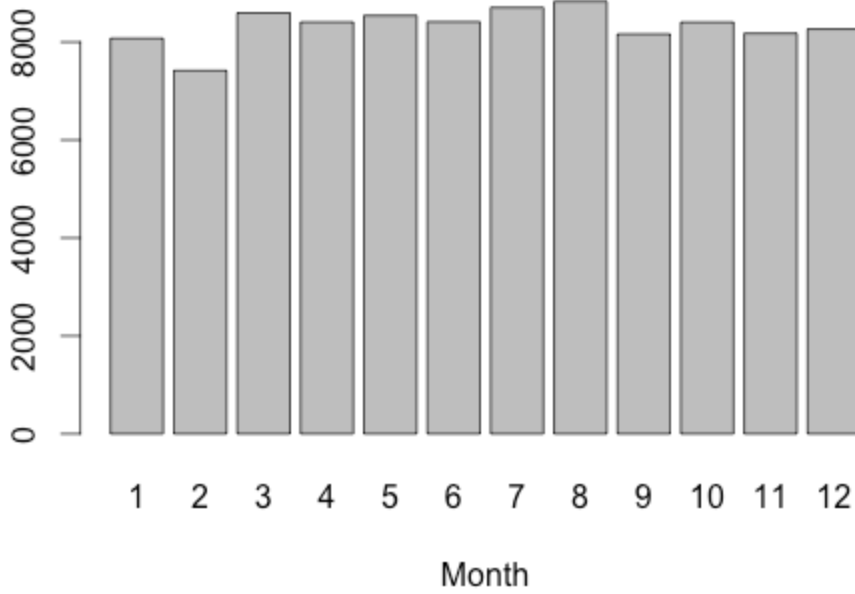
We use the airline on-time performance dataset, which is the same with the dataset used in that blog. Training data is sampled from airline dataset from year 2005 to 2006, while test data is sampled from 2007 airline dataset.

Columns of the data:

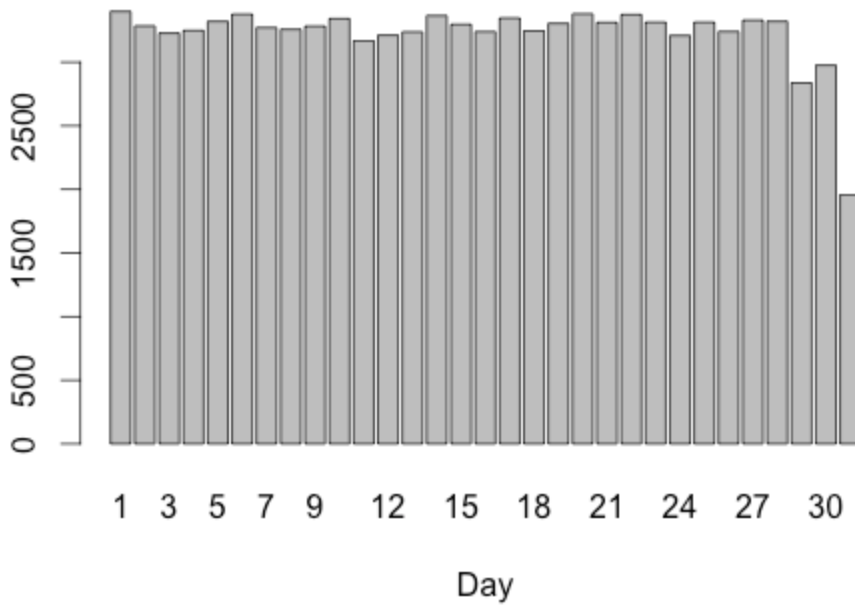
Label	Feature 1	Feature 2	Feature 3	Feature 4	Feature 5	Feature 6	Feature 7	Feature 8
dep_delayed_15min	Month	DayofMonth	DayOfWeek	DepTime	UniqueCarrier	Origin	Dest	Distance

Here are some statistics results of train-0.1m.csv. According to the training set statistics, it is a highly biased dataset with 80956 "N" and 19044 "Y" for the target column "dep_delayed_15min". Our test set has similar distribution with 78383 "N" and 21617 "Y".

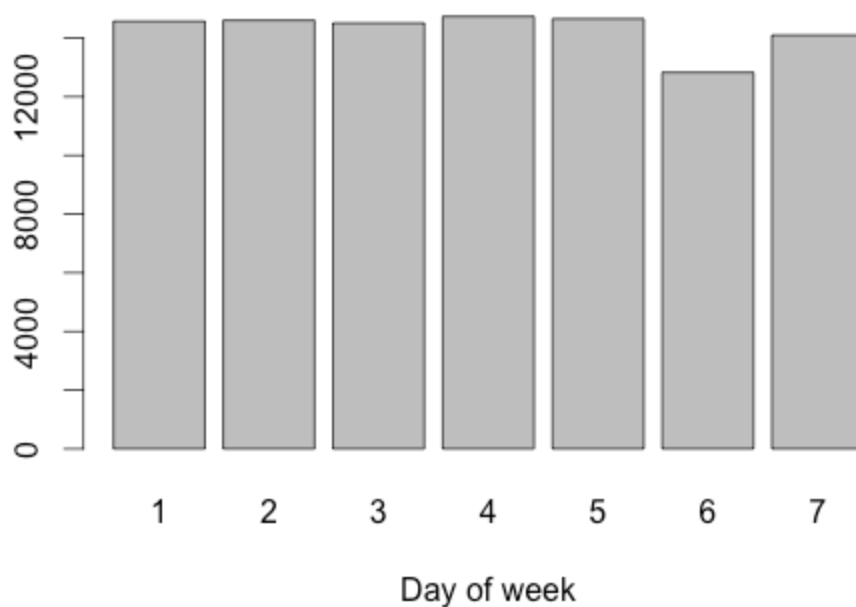
Distribution of Months



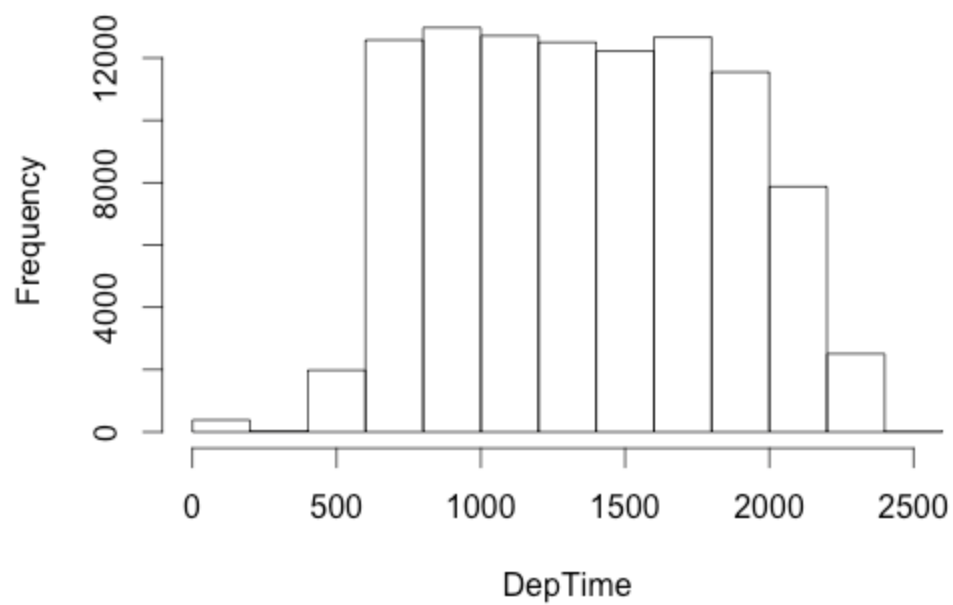
Distribution of Days



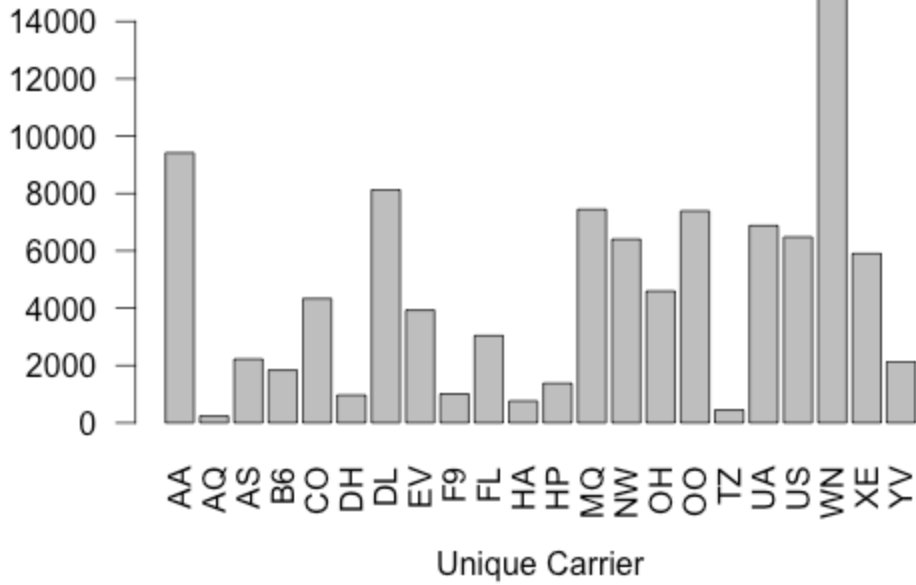
Distribution of Day of week



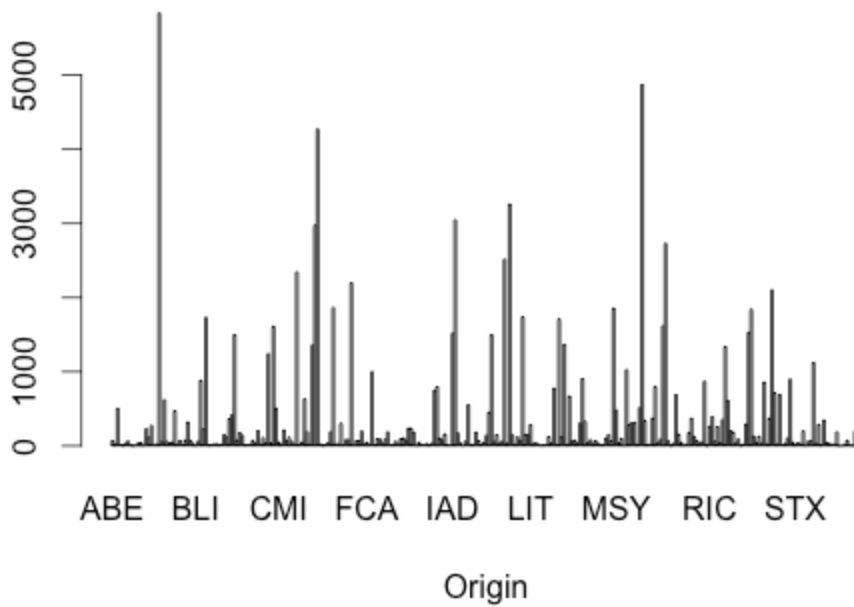
Histogram of DepTime



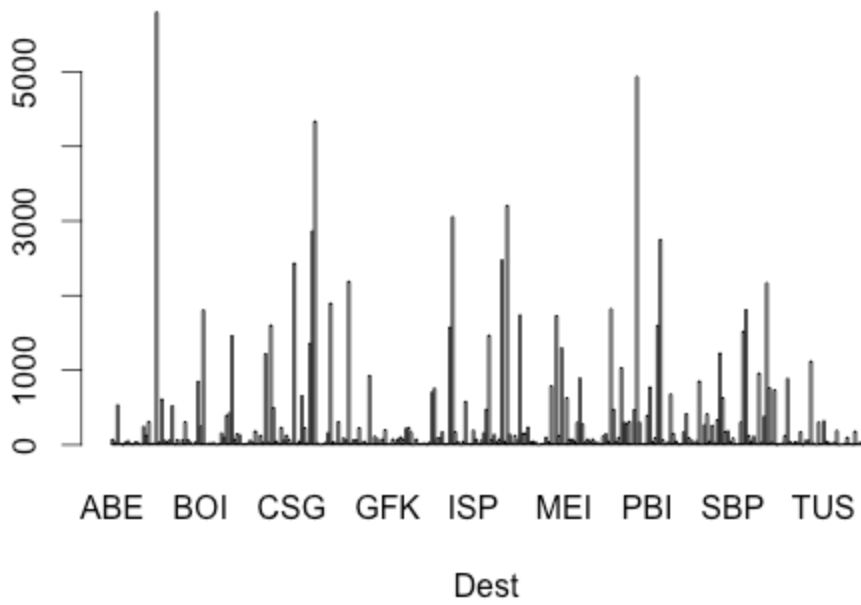
Distribution of Unique Carrier



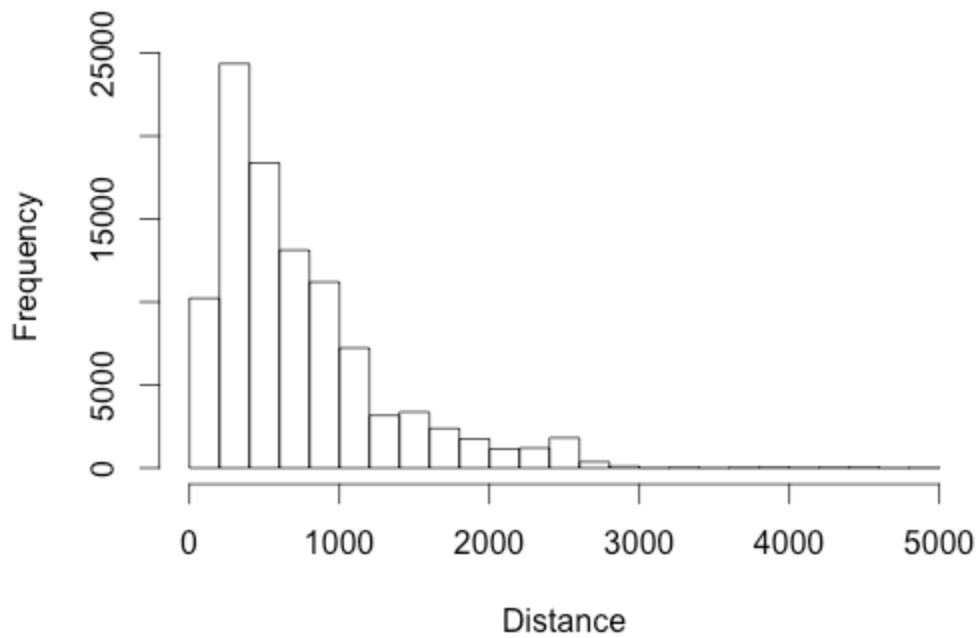
Distribution of Original Places



Distribution of Destinations



Histogram of Distance



How to handle discrete values

In this dataset, all columns except DepTime and Distance are all categorical features. The easiest way to transform the data into LabeledPoint style is RFormula. However, RFormula is not suitable here because it produces different shapes of dataset in comparison with the original one. RFormula uses One-hot encoder so it expands the original dataset into thousands of columns. It brings two drawbacks:

1. The volume of dataset is expanded, which may hurt the performance.
2. One-hot encoder splits one column into cardinality size of new columns, while Random Forest cannot take groups of features into consideration, so it may hurt the accuracy.

The RFormula also recognizes DepTime and Distance into categorical features, so it brings more unnecessary new columns and reduces the accuracy a step further. Because DepTime and Distance are the two most important features for this task.

On the contrary, H2O uses the original dataset, without further preprocessing. As a result, I re-processed the data with the StringIndexer for each categorical columns.

Some results

Parameters for Spark

impurity	List(gini)
maxDepth	List(5)
maxBins	List(400)
minInstancesPerNode	List(1)
numTrees	ArraySeq(1)
featureSubsetStrategy	List(auto)
subsamplingRate	List(1.0)
maxMemoryInMB	32
cacheNodeIds	false
checkpointDir	None
checkpointInterval	10

Spark

dataset	base	maxDepth=20	ntrees=10	ntrees=10	ntrees=10	100, 20	500	500, 20
---------	------	-------------	-----------	-----------	-----------	---------	-----	---------

				maxDepth =20				
0.01	0.651820 952	0.556367 746	0.663224 047	0.633448 682	0.678784 341	0.669790 915	0.684695 681	NA
0.1	0.693036 09	0.586464 084	0.702090 021	0.674857 269	0.706579 952	0.722670 43	0.707556 013	NA
1	0.701479 433	0.624999 02	0.705408 545	0.744321 364	0.710382 107	NA	NA	NA
10	0.698504 332	0.688066 209	0.705764 649	NA	NA	NA	NA	NA

H2O + R

dataset	base	maxDepth =20	ntrees=10	ntrees=10 , maxDepth =20	ntrees=10 0	100, 20	500	500, 20
0.01	0.618594 7	0.544936	0.625409 2	0.615516 8	0.660197 1	0.659832 3	0.663743 3	0.661779
0.1	0.688368 3	0.569516 9	0.694882 7	0.664632 5	0.708888 2	0.707817 7	0.709294 4	0.714758 3
1	0.642678 6	0.602974 1	0.710957 7	0.695598	0.713014 4	0.735657	0.712805	0.741455 4
10	0.677531 8	0.645212 7	0.701887	0.707783 2	0.713378 2	0.747154 1	0.713035	0.752933 5

Conclusions for now

OOM caused by new array

Repeatedly generate new arrays Endless OOM for heap

Wrong computation of working tasks in the same time

Parameters not good, like maxBins, 只能对所有的列统一设定.

调整一下base learner, 看看不同数据集的情况下base learner能达到的最好水平, 然后看看H2O和Spark哪个可以通过ensemble提升的更多。现在看来是H2O更多。所以Spark ensemble的方面可能存在问题。

1. Parameters differences (different params between H2O and Spark)

2. Data ingestion differences
(categorical/continuous/OHE/input->treePoint->baggedInput->nodeIdCache)
3. Base learner differences (sampling bins, order/unorder categorical features, tree pruning)
4. Ensembling differences (weighted point, subsampling, feature subset, voting)
5. Design and implementation differences (histogram collecting/OOM, reduce problem, ser/de problem)