

At Automattic we run a diverse array of systems and as with many companies Kafka is the glue that ties them together; letting us to shuffle data back and forth. Our experience with Kafka have thus far been fantastic, it's stable, provides excellent throughput, and the simple API makes it trivial to hook any of our systems up to it. In fact it's been so popular that we've been steadily piping more and more data through it over the past year. Now we're starting to run out of disk space necessitating an expansion of the cluster.

The expansion plan was pretty straight forward. Spin up some more brokers, have them join the same cluster, then when they are all up simply use Kafka's provided `kafka-reassign-partitions.sh` script to rebalance all topics across all brokers.

We then thought, with more available capacity why not also do a rolling upgrade of our cluster to get on the latest distro across all our servers? And seeing as we run most topics with a replication factor of 2 so it seemed better to use the reassign script to move data off each server as we upgrade as to not have partitions run with only a single copy of the data at any time during the upgrade process.

The Problem

The expansion and rolling restart method seemed like a good plan however after the upgrade we discovered that while the reassign script balances partitions between servers, those servers do not balance each topic between the configured `log.dirs`. This lead to major imbalances on each broker causing some disks to be virtually empty while other disks to be almost entirely full. A quick search showed others have also run into [similar issues](#) with no good way of fixing it.

Why Did It Happen?

The problem stems from the fact that we run our Kafka cluster in JBOD mode, meaning we don't RAID or combine the disks on the hardware or OS level. Instead we mount all disks then configure Kafka's `log.dirs` to use all drives on the server. When run this way Kafka distributes partitions assigned to each broker evenly across all log dirs available to it.

Each new partition that is created will be placed in the directory which currently has the fewest partitions.

—[Kafka Docs](#)

Unfortunately, we have some quite unbalanced topics. Some topics are huge on the order of TBs, with long retention periods. Then there are other topics that perhaps are not yet fully productionized with only a couple MB or KB in each partition. By simply balancing the *number of partitions* across `log.dirs` we end up with some disks with mostly large partitions and some with mostly small partitions.

Hack Partition Reassignment

There is currently some talk about trying to make JBOD support in Kafka better but we don't want to run with such unbalanced disk use nor did we want to bring brokers down to try and hack the metadata in order to move files around in the interim. Luckily how Kafka brokers assign partitions and moves data during partition reassignment operations is well defined which means we can force each broker to balance its data through the use of partition reassignment if we control it very carefully. The basic concept is pretty simple:

1. **Ensure the topic we want to distribute has the right number of partitions.** In order to achieve an even distribution of a topic we need to make sure we have enough partitions to evenly distribute in the first place. This means we'll need to make sure that the $\text{number of partitions} \times \text{replication factor}$ is a multiple of the $\text{number of brokers} \times \text{number of disks on each broker}$. If this is not done we will not end up with exactly the same number of partitions per disk across all brokers and disks.
2. **Balance partition counts across disks on all brokers.** We want every disks to have the same number of partitions on each broker before we move anything. Doing so ensures that when new partitions are created they will be round robined across every available disk on each broker evenly. When a broker is assigned a new partition Kafka assigns that partition to one of the directories configured in `log.dirs` that has the fewest partitions. Without first balancing partition counts across disks these assignment operations will not assign the topic being moved across disks evenly.
3. **Move one topic at a time but move all partitions of that topic together.** When Kafka reassigns partitions the new partition is created on the broker it's being moved to and then synced with the leader. By moving every replica of every partition of a single topic together Kafka will be forced to create a new copy of each partition thus distributing all those partitions evenly.

Tips and Hints

The easiest way to make sure partition counts are balanced across all disks of all brokers is to simply create a new topic and manually assign the right number of partitions to each broker. To do this go to each of your brokers and count up how many partitions you have on each disk then figure out how many more partitions you need on each broker to bring the partition count of the broker up to be the $\text{max partitions on any one disk} \times \text{number of disks}$. For example if you have 3 disks and with 10 partitions on disk A and 6 partitions on disk B & C you'll need 8 more partitions to bring the total partition count of that broker up to 30.

Once you have all the counts simply create a new topic with the number of new partitions needed for each broker. Continuing with the above example if you've determined that broker "1" needs 8 more partitions and broker "2" needs 5 more partitions you would issue the following command to create a new topic with 13 partitions, each with a single replica.

```
1 $ ./kafka-topics.sh \  
2   --zookeeper zk1.xyu.io:2181,zk2.xyu.io:2181/kafka \  
3   --create --topic xyu_gap_filler \  
4   --replica-assignment 1,1,1,1,1,1,1,1,2,2,2,2,2
```

Once this filler topic has been created the number of partitions per disk on each broker will be evened out.

Forcing a reassignment of every partition and replica for a topic is a bit harder to accomplish. First you must make sure that the current number of replicas and number of future replicas is less then the total number of brokers. (This is necessary because all brokers that you are assigning a partition to must not already be assigned that partition.) If this is not done then only some portion of partitions for a topic will actually be moved and it’s likely one of the partitions being moved will end up getting moving to the same disk as a partition not being moved leading to a unbalanced distribution once again.

With dozens of partitions and thousands of possible permutations trying to assign all partitions so that they are both balanced across all brokers and are not assigned to a broker that already contains a replica of it was not something I wanted to do by hand so I wrote a Python script to generate the necessary assignments for me.

1	#!/usr/bin/env python
2	
3	import argparse
4	import json
5	import random
6	import sys
7	
8	def get_topic(args, data):
9	topic = data['partitions'][0]['topic']
10	for partition in data['partitions']:
11	if topic != partition['topic']:
12	raise ValueError('Rebalance only works on one topic at a time')
13	return topic
14	
15	def get_brokers(args, data):
16	brokers = set()
17	if args.brokers:
18	for broker in args.brokers.split(','):
19	brokers.add(int(broker))
20	else:
21	for partition in data['partitions']:
22	for broker in partition['replicas']:
23	brokers.add(broker)
24	return brokers
25	
26	def get_target_replicas(args, data):
27	if args.replicas:
28	replicas = args.replicas
29	else:
30	replicas = get_current_replicas(args, data)
31	return replicas
32	
33	def get_current_replicas(args, data):

```

34     num_replicas = set()
35     for partition in data['partitions']:
36         num_replicas.add( len( partition['replicas'] ) )
37     return max( num_replicas )
38
39 def get_broker_list( data, replicas, brokers ):
40     broker_list = []
41     per_broker = len( data['partitions'] ) * replicas / len( brokers )
42     for broker in brokers:
43         broker_list += [ broker ] * per_broker
44     return broker_list
45
46 def rebalance( data, replicas, broker_list, verbose ):
47     random.shuffle( broker_list )
48     new_partitions = []
49     for partition in data['partitions']:
50         i = 0
51         new_partition = partition.copy()
52         new_partition['replicas'] = []
53         while True:
54             # Ran out of brokers that can be assigned
55             if len( broker_list ) == i:
56                 return None
57
58             # Don't use if matches any existing broker
59             if broker_list[i] in partition['replicas']:
60                 i += 1
61                 continue
62
63             # Don't use if matches already selected new broker
64             if broker_list[i] in new_partition['replicas']:
65                 i += 1
66                 continue
67
68             new_partition['replicas'] += [ broker_list[i] ]
69             broker_list = broker_list[0:i] + broker_list[1+i:]
70             i=0
71
72             # Found all brokers for this partition
73             if len( new_partition['replicas'] ) == replicas:
74                 break
75
76     if verbose:
77         print "Partition " + str( partition['partition'] ) + "\t" + str( partition['replicas'] ) + " -> " + str( new_partition['replicas'] )
78

```

```
79         new_partitions += [ new_partition ]
80
81     return new_partitions
82
83 # CLI args
84 parser = argparse.ArgumentParser()
85 parser.add_argument(
86     "assignment", type=str,
87     help="current partition replica assignment (JSON)"
88 )
89 parser.add_argument(
90     "--brokers", type=str,
91     help="override brokers to use"
92 )
93 parser.add_argument(
94     "--replicas", type=int,
95     help="override number of replicas after rebalance"
96 )
97 parser.add_argument(
98     "-v", "--verbose",
99     help="increase output verbosity", action="store_true"
100 )
101 args = parser.parse_args()
102 data = json.loads( args.assignment )
103
104 # Check we have enough brokers to handle rebalance
105 if get_target_replicas( args, data ) + get_current_replicas( args, data ) > len( get_brokers( args, data ) ):
106     raise ValueError( 'Cannot rebalance, not enough brokers' )
107
108 topic = get_topic( args, data )
109 replicas = get_target_replicas( args, data )
110 brokers = get_brokers( args, data )
111 broker_list = get_broker_list( data, replicas, brokers )
112
113 # Try to assign 100 times then giveup
114 print "Attempting to rebalance " + topic
115 for i in range( 0, 100 ):
116     if args.verbose:
117         print "Try " + str( i )
118     else:
119         sys.stdout.write('.')
120         if 74 == i%75:
121             print ""
122     new_partitions = rebalance( data, replicas, broker_list, args.verbose )
123     if None != new_partitions:
```

124	break
125	
126	if None == new_partitions:
127	raise RuntimeError('Could not rebalance partitions')
128	
129	print "\nDone!\n"
130	
131	print json.dumps(
132	{ "version" : 1, "partitions" : new_partitions },
133	separators=(',', ':')
134)
kafka-rebalance.py hosted with ❤ by GitHub	
view raw	

To use this script you must first get the current partition replica assignment from Kafka for the topic that you are trying to move using the `kafka-reassign-partitions.sh` tool:

```
1 $ ./kafka-reassign-partitions.sh \  
2   --zookeeper zk1.xyu.io:2181,zk2.xyu.io:2181/kafka \  
3   --broker-list '1,2,3' \  
4   --topics-to-move-json-file topic.json \  
5   --generate
```

The param specified by `broker-list` does not actually matter as we don't care about the plan generated by the tool. We're just using this to print out the current assignment in JSON format.

The `topic.json` file should look something like the following. (Remember, only move one topic at a time otherwise balance is not guaranteed!)

```
1 {  
2   "version": 1,  
3   "topics": [  
4     { "topic": "my_sample_topic" }  
5   ]  
6 }
```

The `kafka-reassign-partitions.sh` script will generate two JSON strings, one showing the current partition assignment and one showing a proposed partition assignment. Ignore the proposed assignment and instead pass the current assignment JSON string to `kafka-rebalance.py` to get a new proposed assignment, one that guarantees all partitions will relocate.

```
1 $ ./kafka-rebalance.py \  
2   '{"version":1,"partitions":[{"...},{...}]}'
```

This is a pretty convoluted way to rebalance topics across all disks on each Kafka broker however it's pretty safe as this method does not involve any metadata hacking nor does it necessitate bring any brokers down. Hopefully Kafka will become better at balancing topics across disks when run in JBOD mode by itself in the future.