

Report big data for company

Elasticsearch & Kibana

Yuchao QIAN Ying DAI Zhen Xing Song SHI

Part 1 Installation and data preparation

Prerequisites:

Elasticsearch 2.4+Kibana 4.6.1

Data set: IMDB Movie Dataset

In order to configure the environment easily, we need to use the tools called docker to build our working space image.

And the new version of Elasticsearch is not necessary and it has too many changes and features we will not use, so we decided to use the old version (2.4), which is much easier to use.

The first step to use docker is to write a docker file to put all the configuration into itself. So it is necessary to learn how to write a Dockerfile.

As we know there are lots of images already on the cloud, so we can easily write a Dockerfile to have both Elasticsearch and Kibana using the following script:

```
FROM sebp/elk:es241_l240_k461
```

And we need to add the local file into the docker image, so we have to use the command ADD to add all the data and script from the local storage like this:

```
ADD ./start.sh /usr/local/bin/start.sh
```

And also we need to give the permission to run the script:

```
RUN chmod +x /usr/local/bin/start.sh
```

Then we have to expose some port for Elasticsearch and Kibana because it can be configured by restful API, so we have to use EXPOSE method like this:

```
EXPOSE 5601 9200 9300
```

Finally, we need to run the script we store in the image and it will only run once a script when the image is online the first time by adding this:

```
CMD [ "/usr/local/bin/start.sh" ]
```

After finish the Dockerfile, we need to prepare our dataset for Elasticsearch.

```
curl -XPUT localhost:9200/_bulk --data-binary @/data/movie.json
```

```
{"index":{"_index":"movies","_type":"record","_id":1}}
```

In order to prepare the dataset, we use ATOM and a script to create our JSON file from a csv file we download from the website. And here is the screenshot of the dataset of JSON format.

```
1{"index":{"_index":"movies","_type":"record","_id":1}}
2{"color":"Color","director_name":"James
Cameron","num_critic_for_reviews":723,"duration":178,"director_facebook_likes":0,"actor_3_facebook_
David Moore","actor_1_facebook_likes":1000,"gross":760505847,"genres":["Action|Adventure|Fantasy|
Sci-Fi"],"actor_1_name":"CCH
Pounder","movie_title":"Avatar","num_voted_users":886204,"cast_total_facebook_likes":4834,"actor_3
Studi","facenumber_in_poster":0,"plot_keywords":["avatar|future|marine|nave|
paraplegic"],"movie_imdb_link":"http://www.imdb.com/title/tt0499549/?
ref=fn_tt_tt_1","num_user_for_reviews":3054,"language":"English","country":"USA","content_rating":
3{"index":{"_index":"movies","_type":"record","_id":2}}
4{"color":"Color","director_name":"Gore
Verbinski","num_critic_for_reviews":302,"duration":169,"director_facebook_likes":563,"actor_3_faceb
Bloom","actor_1_facebook_likes":40000,"gross":309404152,"genres":["Action|Adventure|
Fantasy"],"actor_1_name":"Johnny Depp","movie_title":"Pirates of the Caribbean: At World's
End","num_voted_users":471220,"cast_total_facebook_likes":48350,"actor_3_name":"Jack
Davenport","facenumber_in_poster":0,"plot_keywords":["goddess|marriage ceremony|marriage proposal|
pirate|singapore"],"movie_imdb_link":"http://www.imdb.com/title/tt0449088/?
ref=fn_tt_tt_1","num_user_for_reviews":1238,"language":"English","country":"USA","content_rating":
5{"index":{"_index":"movies","_type":"record","_id":3}}
6{"color":"Color","director_name":"Sam
Mendes","num_critic_for_reviews":602,"duration":148,"director_facebook_likes":0,"actor_3_facebook_l
Kinneer","actor_1_facebook_likes":11000,"gross":200074175,"genres":["Action|Adventure|
Thriller"],"actor_1_name":"Christoph
Waltz","movie_title":"Spectre","num_voted_users":275868,"cast_total_facebook_likes":11700,"actor_3
Sigman","facenumber_in_poster":1,"plot_keywords":["bomb|espionage|sequel|spy|
terrorist"],"movie_imdb_link":"http://www.imdb.com/title/tt2379713/?
ref=fn_tt_tt_1","num_user_for_reviews":994,"language":"English","country":"UK","content_rating":
7{"index":{"_index":"movies","_type":"record","_id":4}}
8{"color":"Color","director_name":"Christopher
Nolan","num_critic_for_reviews":813,"duration":164,"director_facebook_likes":22000,"actor_3_facebook_
..._1_facebook_likes":22000,"gross":448230643,"genres":["Action|
..."]}
```

And here is part of the ATOM script to convert the CSV file to JSON file.

```

json-converter.coffee      Welcome Guide
1  YAML = require 'js-yaml'
2  converter = require('json-2-csv')
3  {CompositeDisposable} = require 'atom'
4  id = 0
5  unique = (array) ->
6    output = {}
7    output[array[key]] = array[key] for key in [0...array.length]
8    value for key, value of output
9
10 isEven = (val) ->
11   return val % 2 == 0
12
13 expandAction = (docs, opType='index') ->
14   metaIndex = atom.config.get('json-converter.elasticIndex')
15   metaType = atom.config.get('json-converter.elasticDocType')
16   metaId = atom.config.get('json-converter.elasticUidField')
17   metaParentId = atom.config.get('json-converter.elasticParentUidField')
18   excludeFields = atom.config.get('json-converter.elasticExcludeFields')
19
20   actions = []
21   for doc in docs
22     docKeys = Object.keys(doc)
23     id = id+1
24     action = {}
25     action[opType] = {}
26     action[opType]._index = metaIndex if metaIndex
27     action[opType]._type = metaType if metaType
28     action[opType]._id = id
29     action[opType]._parent = doc[metaParentId] if metaParentId in docKeys
30
31     if opType is 'delete'
32       actions.push(JSON.stringify(action))
33     else
34       for field, value of doc
35         delete doc[field] if field in excludeFields
36         doc = {"doc": doc} if opType is 'update'
37         actions.push(JSON.stringify(action))
38         actions.push(JSON.stringify(doc))
39
40   return actions
41
42

```

After all the above steps, we will have a Dockerfile, a JSON file.

So we need to add the data into the docker image by doing this:

```
ADD ./data /data
```

We need to modify the start.sh file from the example project and add the following code to run import the data.

```
curl -XPUT localhost:9200/_bulk --data-binary @/data/movie.json
```

Then we need to build our docker image and run it by the following command:

```
sudo docker build -t elakiba:1.0.1 .
```

```
sudo docker run -p 9200:9200 -p 5601:5601 elakiba:1.0.1
```

We also met the problem with the virtual machine that max_map_count is too small, so we need to set it to be 262144 as indicated.

```
sudo sysctl -w vm.max_map_count=262144
```

Part 2 Data analysis and dashboard

Then we have all the data in Elasticsearch server and we tried to use Kibana to plot the chart for visualization.

<http://localhost:5601> is the url of Kibana and we should Browse <http://localhost:5601> Uncheck "Index contains time-based events" Replace "logstash-*" by "mo*".

Firstly we chose to use a vertical bar chart to visualize the number of the movies per director, however, we met a big problem that the field "director_name" is analyzed, so Elasticsearch will split the name into words, so we can not have a full name of the director, which will casue a big problem for us.

So we tried to fix this problem. We go to the indices page in Kibana and we find out that all the strings are analyzed by default(like the following picture).

fields (34)		scripted fields (0)		
name	type	format	analyzed	indexed
plot_keywords	string		✓	✓
country	string		✓	✓
aspect_ratio	number			✓
actor_3_facebook_likes	number			✓
num_critic_for_reviews	number			✓

We found out that we have to change the mapping before import the data to set the filed to be not_analyzed in order to let Elasticsearch not to split the word.

So we create a mapping file we set several string filed to be not_analyzed

```
1 {
2   "mappings": {
3     "record": {
4       "properties": {
5         "color":{"type":"string","index":"analyzed"},
6         "director_name":{"type":"string","index":"not_analyzed"},
7         "num_critic_for_reviews":{"type":"long"},
8         "duration":{"type":"long"},
9         "director_facebook_likes":{"type":"long"},
10        "actor_3_facebook_likes":{"type":"long"},
11        "actor_2_name":{"type":"string","index":"not_analyzed"},
12        "actor_1_facebook_likes":{"type":"long"},
13        "gross":{"type":"double"},
14        "genres":{"type":"string"},
15        "actor_1_name":{"type":"string","index":"not_analyzed"},
16        "movie_title":{"type":"string"},
17        "num_voted_users":{"type":"long"},
18        "cast_total_facebook_likes":{"type":"long"},
19        "actor_3_name":{"type":"string","index":"not_analyzed"},
20        "facenumber_in_poster":{"type":"long"},
21        "plot_strings":{"type":"string"},
22        "movie_imdb_link":{"type":"string"},
23        "num_user_for_reviews":{"type":"long"},
24        "language":{"type":"string"},
25        "country":{"type":"string"},
26        "content_rating":{"type":"string","index":"not_analyzed"},
27        "budget":{"type":"long"},
28        "title_year":{"type":"long"},
29        "actor_2_facebook_likes":{"type":"long"},
30        "imdb_score":{"type":"double"},
31        "aspect_ratio":{"type":"double"},
32        "movie_facebook_likes":{"type":"long"}
33      }
34    }
35  }
36 }
37 ]
```

And then we do the same thing to put the file into Docker image

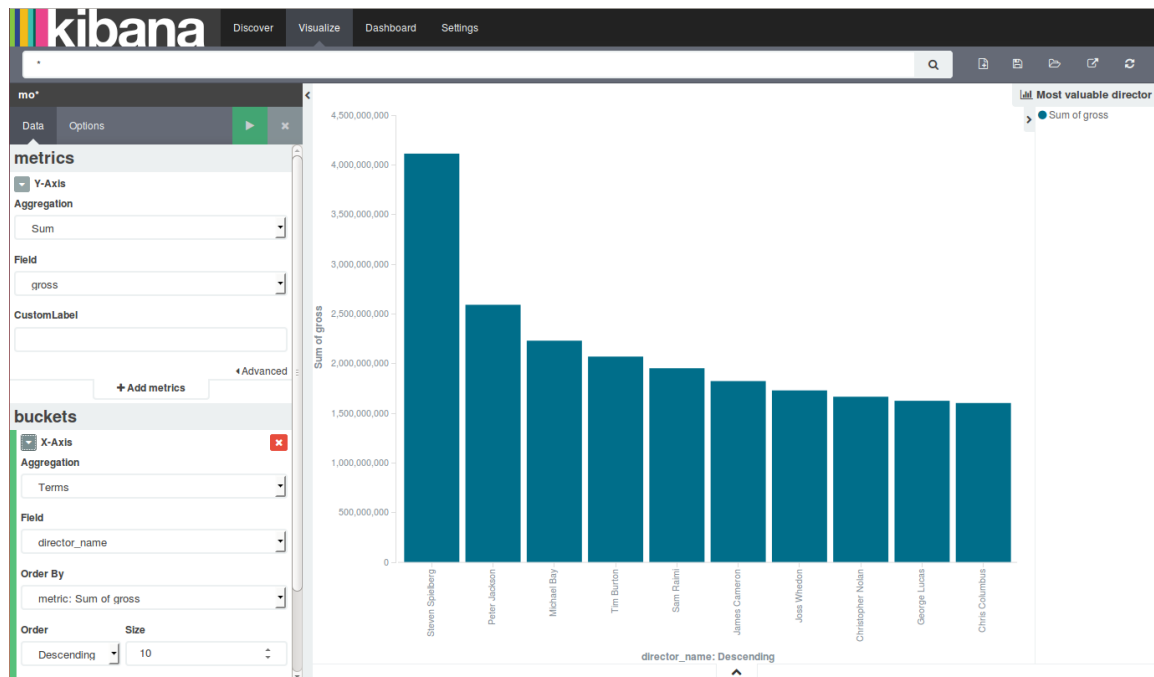
```
ADD ./mapping.json /tmp/mapping.json
```

And write following script in start.sh to put the mapping before import data.

```
curl -XPUT localhost:9200/movies -d @/tmp/mapping.json
```

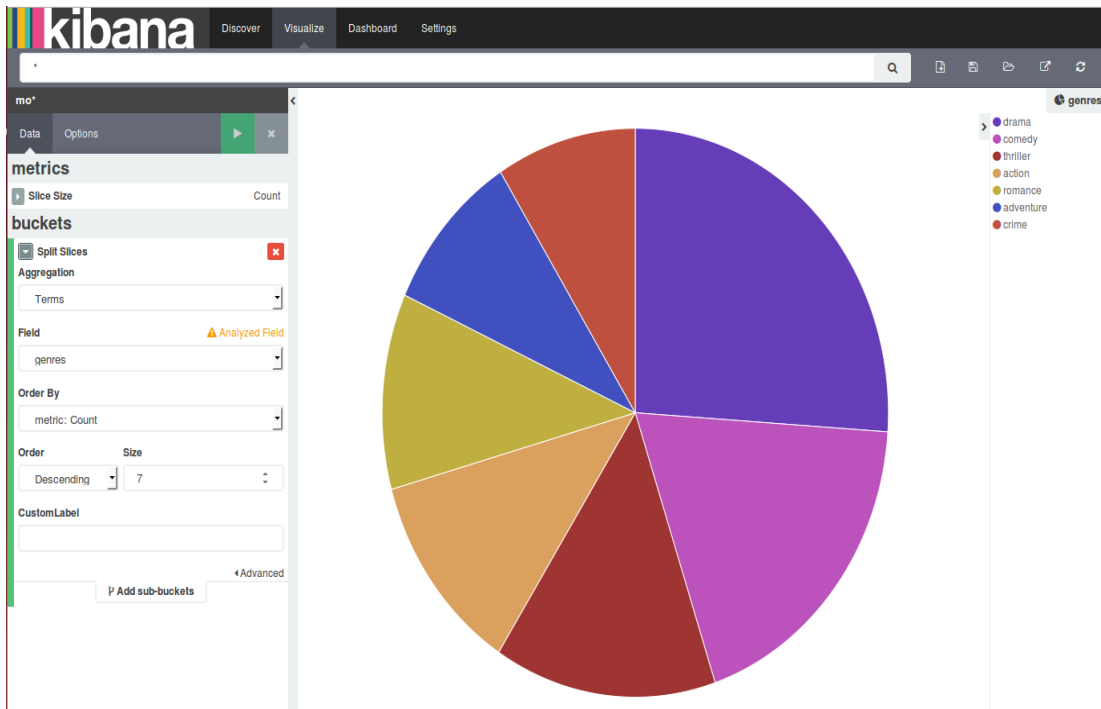
Then we can go to <http://localhost:5601> again and trying to create the vertical bar chart.

We configure the visualization with the following setting and we can have a clear idea of who is the biggest money-maker during the past years.

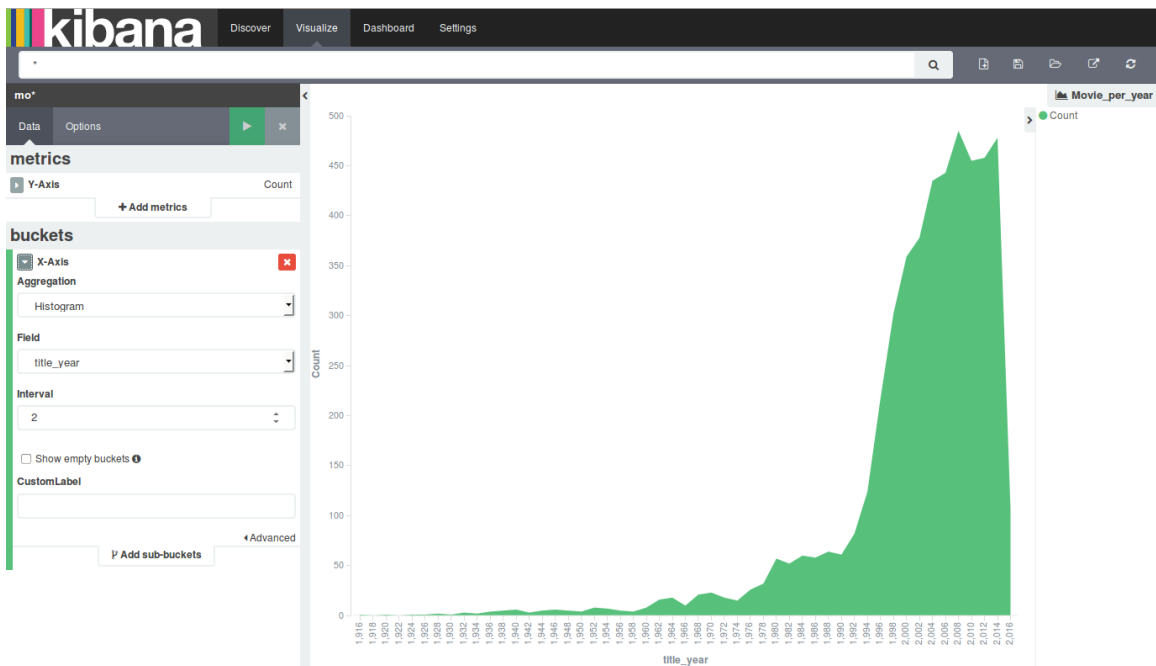


And then we create five more visualization to have a better idea of the dataset.

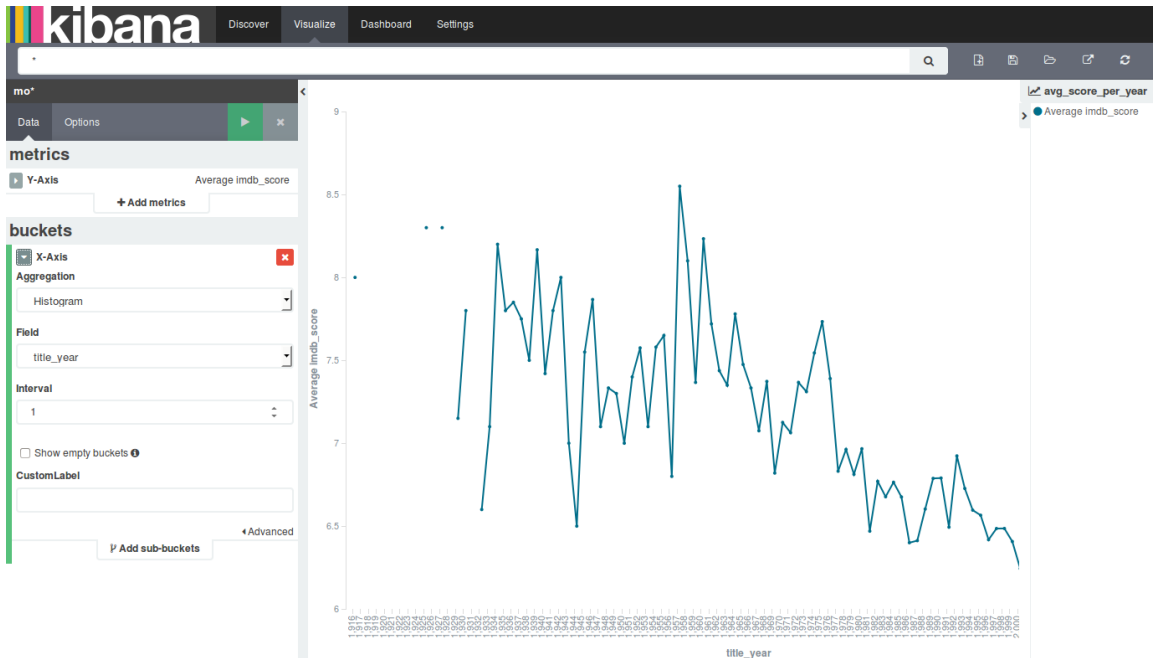
The most frequent genres of movies during the past years.



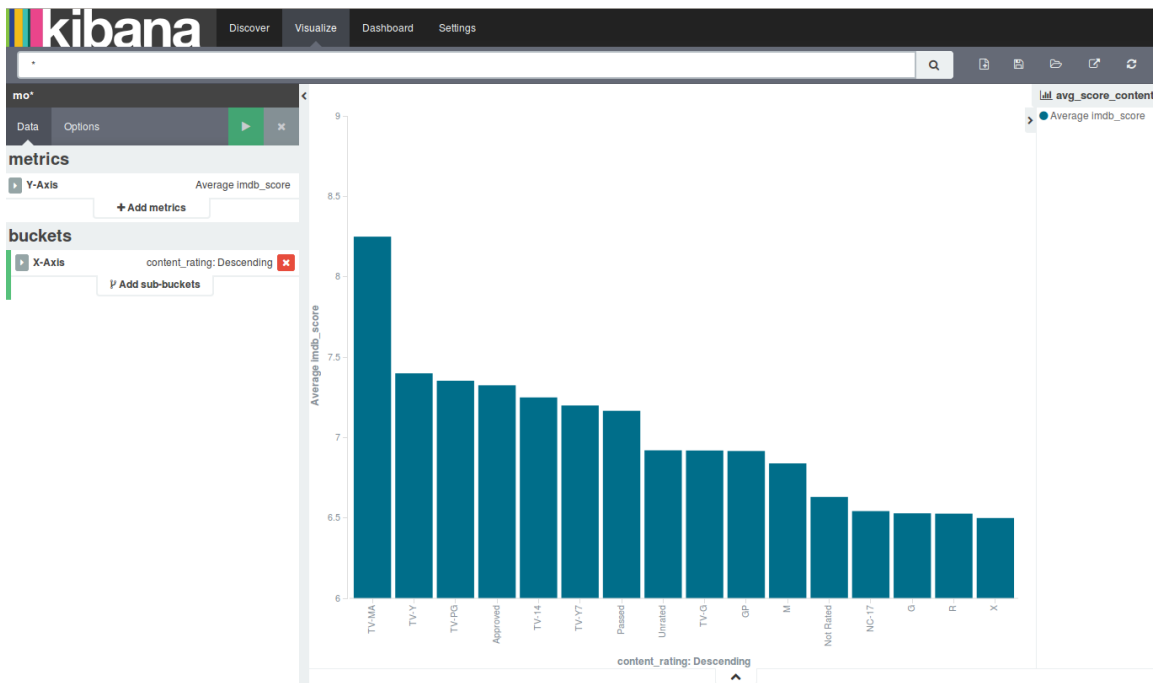
The trend of the amount of movies each year



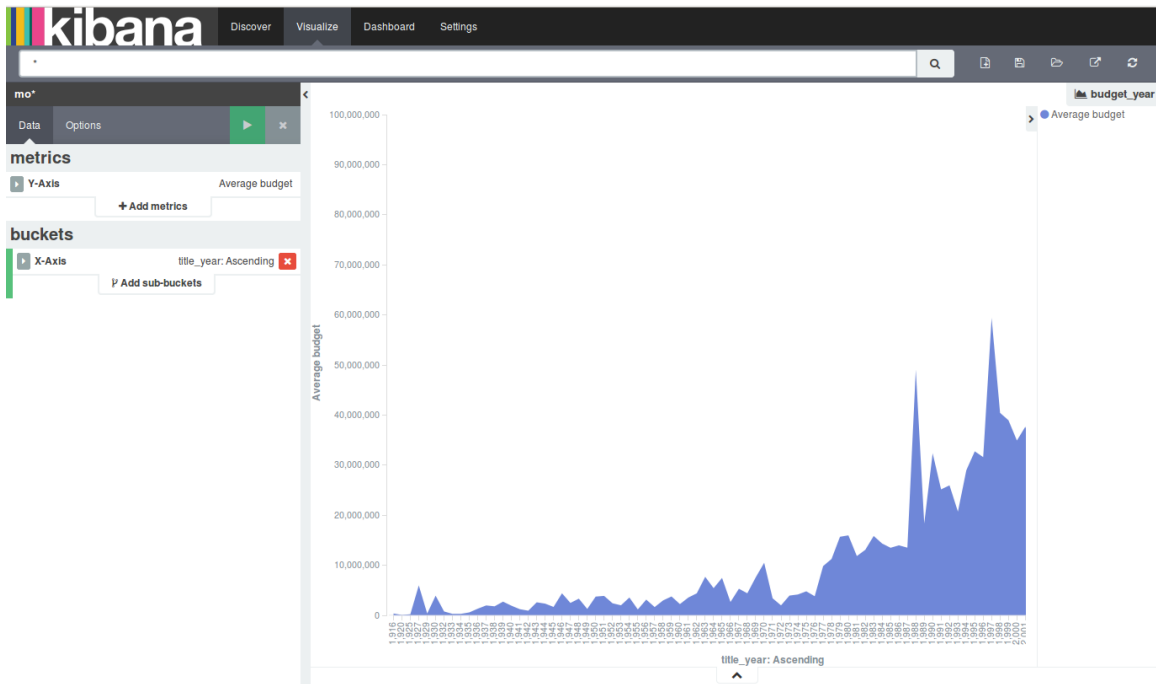
The average IMDE score by year



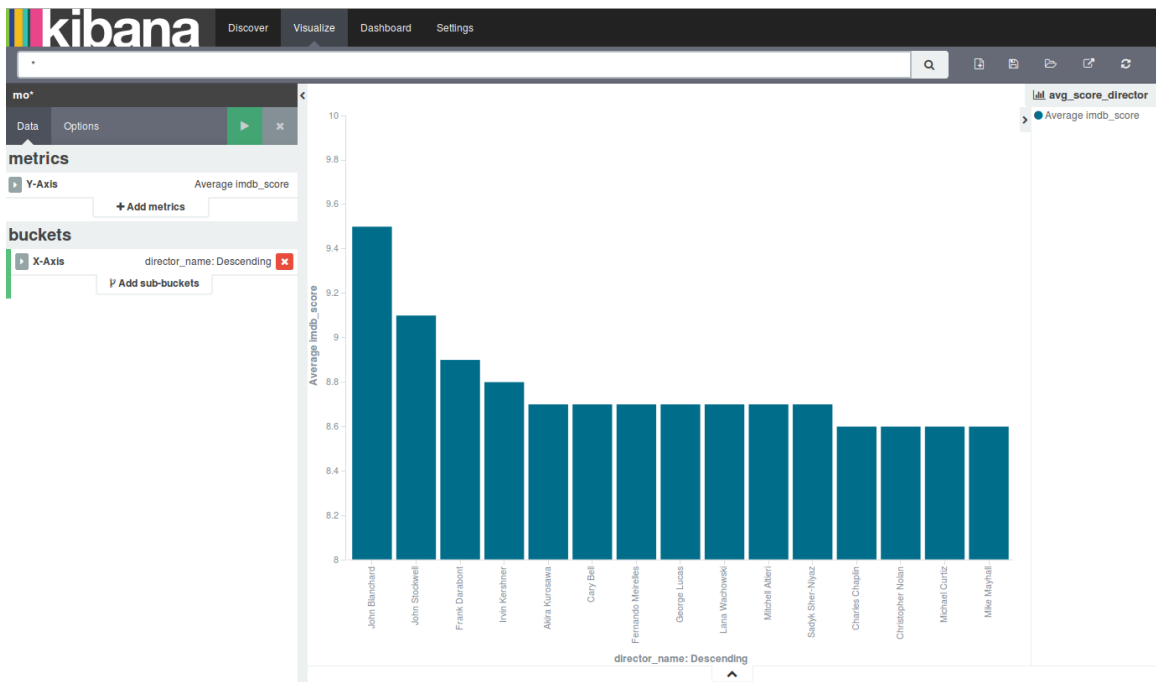
The average score by content level



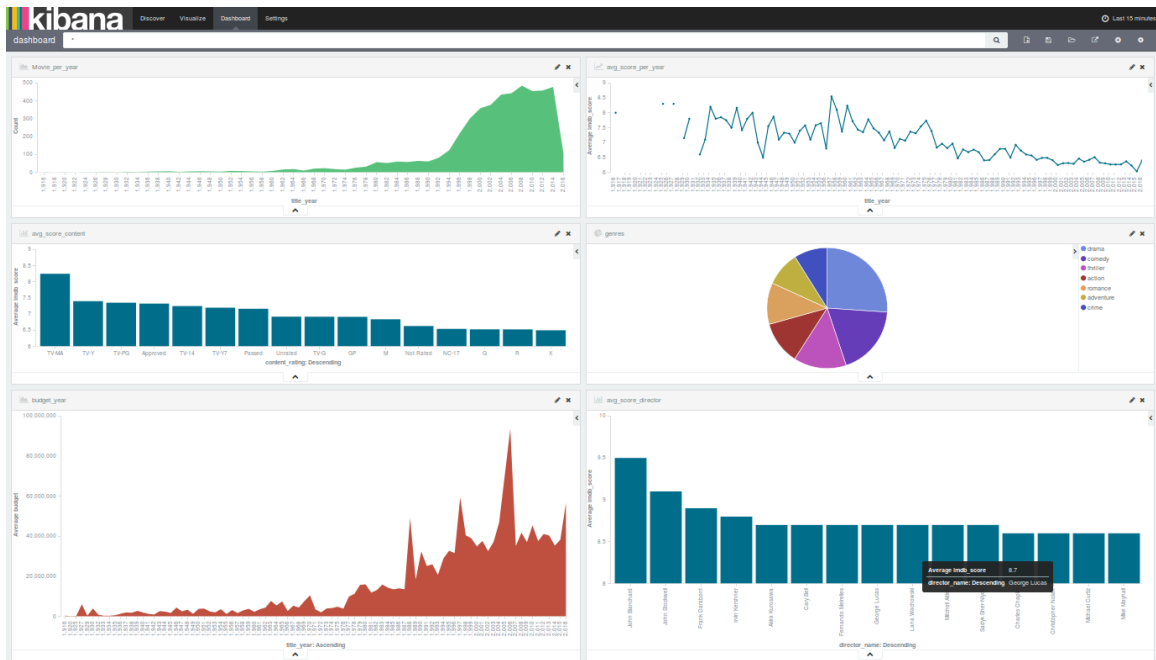
The trend of gross each year



The average IMDB score by director



Then we put all this visualization onto the dashboard.



Then the next step is to export the dashboard and visualizations also we can use the export.json to import the dashboard and visualizations.

