



Tutorial 5: KDDCup'99 Analysis using HIVE

CN7022 - Big Data Analytics

Dr Amin Karami (a.karami@uel.ac.uk)

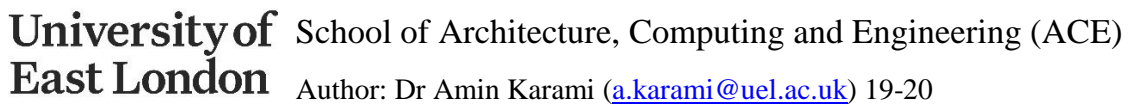
LEARNING OUTCOMES: After completing this tutorial, you should:

- Have gotten a hands-on experience in using HIVE for data analysis
- Understand the real case studies applied in the Big Data platform for analysis
- Visualize the outcomes of queries into the graphical representations

Phase 1: Getting Ready

- This is the data set used for The Third International Knowledge Discovery and Data Mining Tools Competition, which was held in conjunction with KDD-99 The Fifth International Conference on Knowledge Discovery and Data Mining. The competition task was to build a network intrusion detector, a predictive model capable of distinguishing between “bad” connections, called intrusions or attacks, and “good” normal connections. This database contains a standard set of data to be audited, which includes a wide variety of intrusions simulated in a military network environment.
- The first task for big data analytics is to understand the data carefully. The featured are described [here](#).
- This dataset has 22 types of bad connections: `back`, `buffer_overflow`, `ftp_write`, `guess_passwd`, `imap`, `ipsweep`, `land`, `loadmodule`, `multihop`, `neptune`, `nmap`, `perl`, `phf`, `pod`, `portsweep`, `rootkit`, `satan`, `smurf`, `spy`, `teardrop`, `warezclient`, `warezmaster`.
- The whole overview of the features are as follows:

<i>feature name</i>	<i>description</i>	<i>type</i>
<code>duration</code>	length (number of seconds) of the connection	continuous
<code>protocol_type</code>	type of the protocol, e.g. tcp, udp, etc.	discrete
<code>service</code>	network service on the destination, e.g., http, telnet, etc.	discrete
<code>src_bytes</code>	number of data bytes from source to destination	continuous
<code>dst_bytes</code>	number of data bytes from destination to source	continuous
<code>flag</code>	normal or error status of the connection	discrete
<code>land</code>	1 if connection is from/to the same host/port; 0 otherwise	discrete
<code>wrong_fragment</code>	number of “wrong” fragments	continuous
<code>urgent</code>	number of urgent packets	continuous
<code>host</code>	number of “host” indicators	continuous
<code>num_failed_logins</code>	number of failed login attempts	continuous



- A sample of data is:

- Download the data from [here](#). This is the KDDCup'99 data with 708MB. Then, copy this file into the Cloudera machine in **Desktop**.

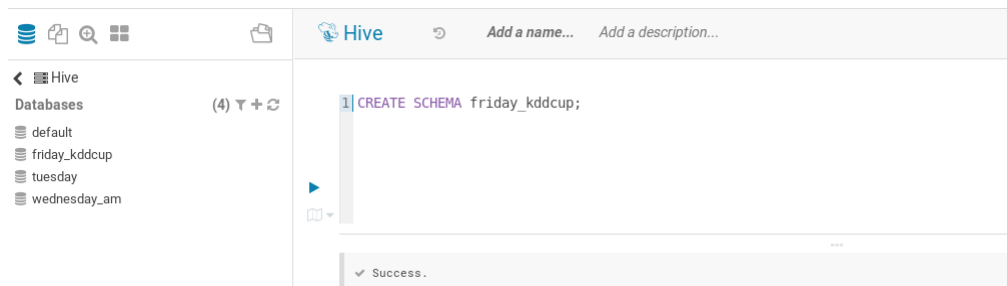


Phase 2: Create a Database and Table in HIVE

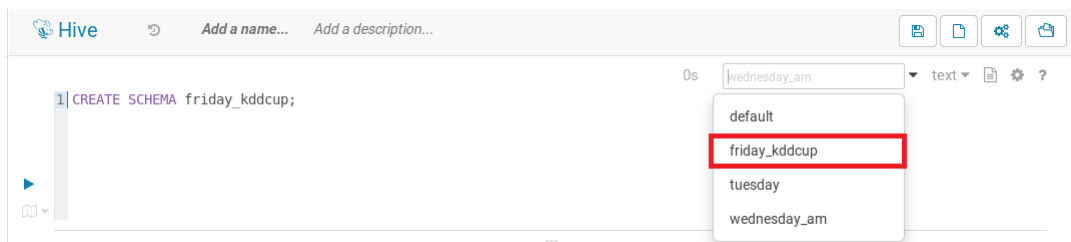
- a) Create a database (**Select your own database using the tutorial day**):

Note: every group should create their own database: **friday**, **tuesday**, **wednesday_am**, **wednesday_pm**

```
CREATE SCHEMA friday_kddcup;
```



- b) Refresh the page, then select the **friday_kddcup** database to work with it (**Select your own database using the tutorial day**):



- c) Create a table:

```
CREATE TABLE IF NOT EXISTS kddcup (duration int, protocol_type string, service string, flag string, src_bytes int, dst_bytes int, land int, wrong_fragment int, urgent int, host int, num_failed_logins int, logged_in int, num_compromised int, root_shell int, su_attempted int, num_root int, num_file_creations int, num_shells int, num_access_files int, num_outbound_cmds int, is_host_login int, is_guest_login int, count int, srv_count int, serror_rate float, srv_serror_rate float, rerror_rate float, srv_rerror_rate float, same_srv_rate float, diff_srv_rate float, srv_diff_host_rate float, dst_host_count int, dst_host_srv_count int, dst_host_same_srv_rate float, dst_host_diff_srv_rate float, dst_host_same_src_port_rate float, dst_host_srv_diff_host_rate float, dst_host_serror_rate float, dst_host_srv_serror_rate float, dst_host_rerror_rate float, dst_host_srv_rerror_rate float, connection_status string) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' LINES TERMINATED BY '\n' STORED AS TEXTFILE;
```



d) Load the data into table:

```
LOAD DATA LOCAL INPATH '/home/cloudera/Desktop/kddcup.data'  
OVERWRITE INTO TABLE kddcup;
```

e) Let's show the first 10 records of data.

```
select * from kddcup limit 20;
```

Query History Saved Queries Results (20) Q, ✓

	kddcup.duration	kddcup.protocol_type	kddcup.service	kddcup.flag	kddcup.src_bytes	kddcup.dst_bytes	kddcup.land	kddcup.wrong_fragment	kddcup.urgent	kddcup.host	kddcup.num_failed_logins	kddcup.logged_in	kddcup.num_compromised	kddcup
1	0	tcp	http	SF	215	45076	0	0	0	0	1	0	0	
2	0	tcp	http	SF	162	4528	0	0	0	0	1	0	0	
3	0	tcp	http	SF	236	1228	0	0	0	0	1	0	0	
4	0	tcp	http	SF	233	2032	0	0	0	0	1	0	0	
5	0	tcp	http	SF	239	486	0	0	0	0	1	0	0	
6	0	tcp	http	SF	238	1282	0	0	0	0	1	0	0	
7	0	tcp	http	SF	235	1337	0	0	0	0	1	0	0	
8	0	tcp	http	SF	234	1364	0	0	0	0	1	0	0	
9	0	tcp	http	SF	239	1295	0	0	0	0	1	0	0	
10	0	tcp	http	SF	181	5490	0	0	0	0	1	0	0	
11	0	tcp	http	SF	184	124	0	0	0	0	1	0	0	
12	0	tcp	http	SF	185	9020	0	0	0	0	1	0	0	
13	0	tcp	http	SF	239	1295	0	0	0	0	1	0	0	
14	0	tcp	http	SF	181	5450	0	0	0	0	1	0	0	
15	0	tcp	http	SF	236	1228	0	0	0	0	1	0	0	
16	0	tcp	http	SF	233	2032	0	0	0	0	1	0	0	
17	0	tcp	http	SF	238	1282	0	0	0	0	1	0	0	
18	0	tcp	http	SF	235	1337	0	0	0	0	1	0	0	
19	0	tcp	http	SF	234	1364	0	0	0	0	1	0	0	
20	0	tcp	http	SF	239	486	0	0	0	0	1	0	0	

Phase 3: KDDCup Data Querying using HIVE

a) List the number of connections for different connections' statuses

```
SELECT connection_status, protocol_type, count(*) as Counts FROM  
kddcup GROUP BY connection_status, protocol_type;
```

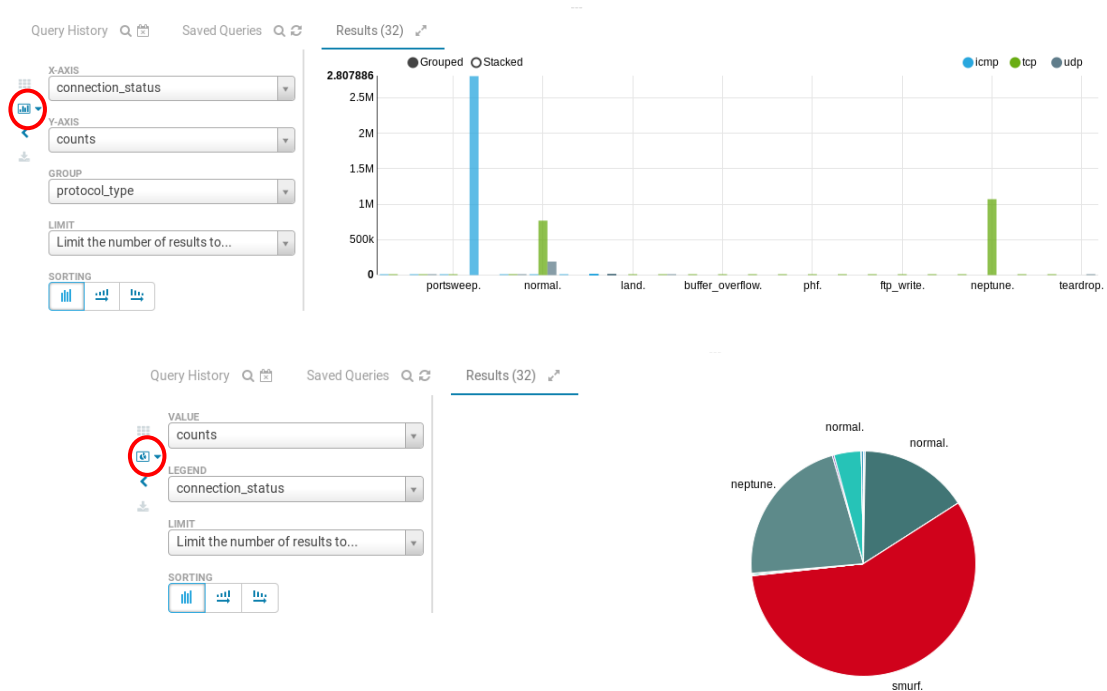
GROUP BY takes a while to be completed.

Query History Saved Queries Results (32) Q, ✓

	connection_status	protocol_type	counts
1	ipsweep.	icmp	11557
2	land.	tcp	21
3	nmap.	icmp	1032
4	nmap.	udp	250
5	normal.	tcp	768670
6	portsweep.	icmp	6
7	rootkit.	tcp	7
8	smurf.	icmp	2807886
9	teardrop.	udp	979
10	back.	tcp	2203



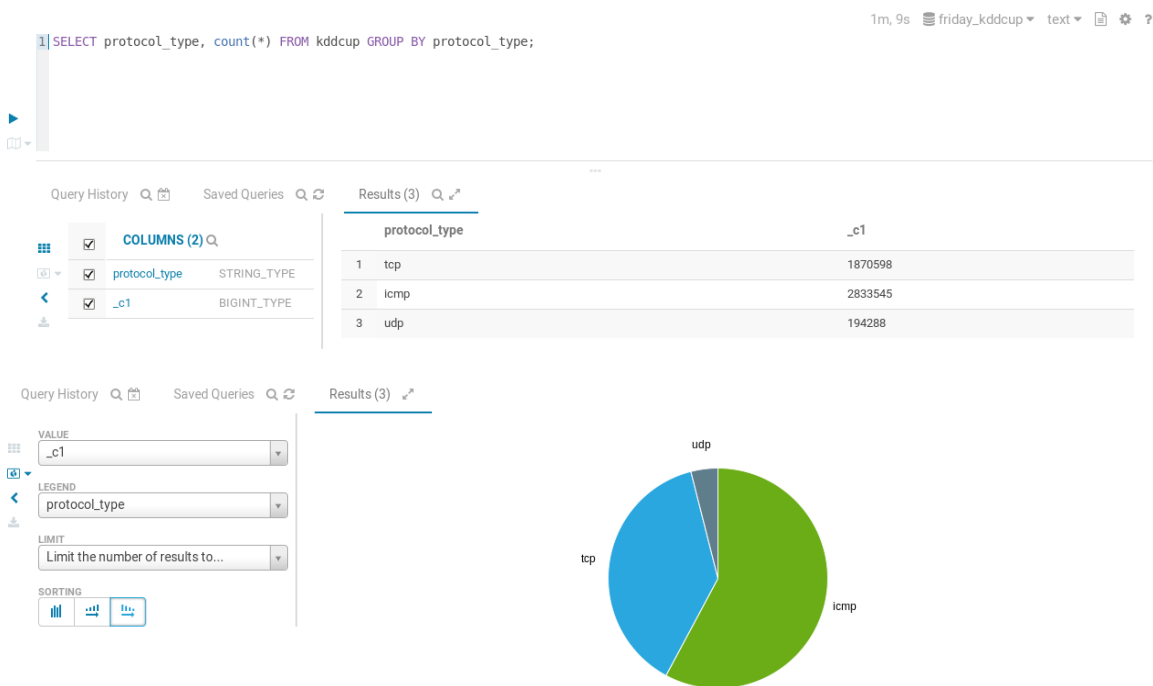
Now, it is a time to present our work visually:

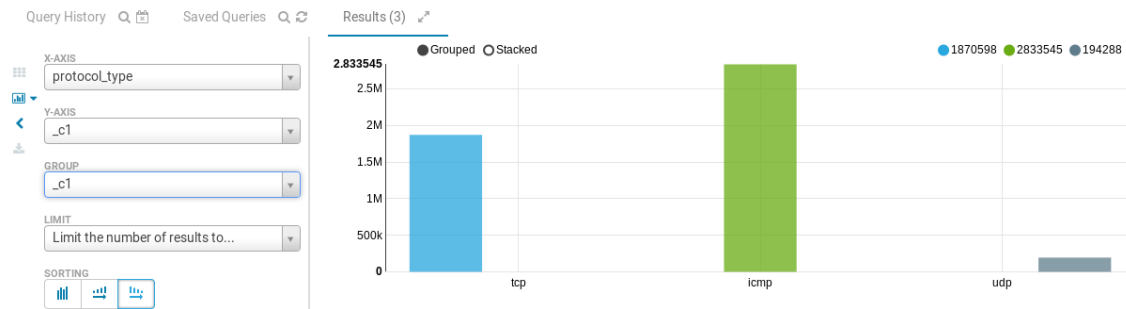


b) List the number of connections.

```
SELECT protocol_type, count(*) FROM kddcup GROUP BY protocol_type;
```

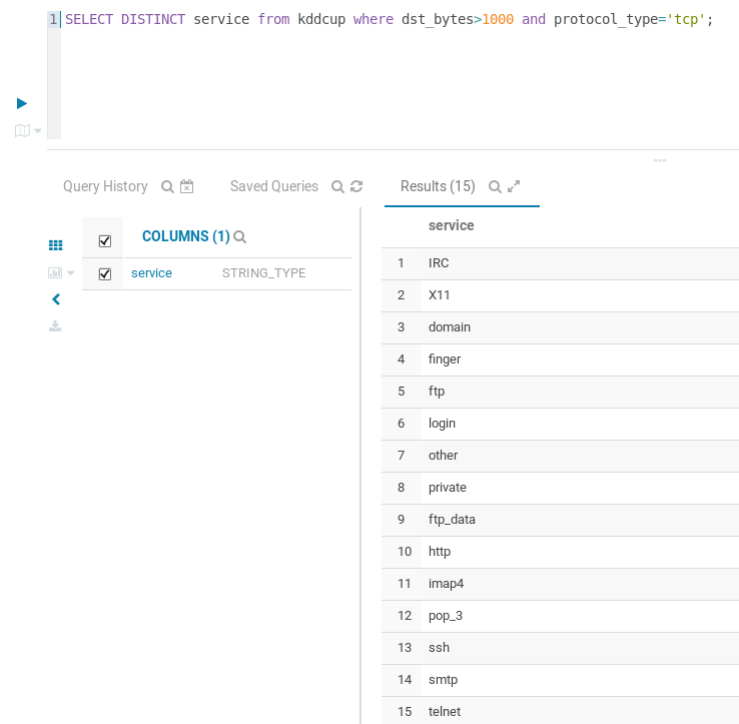
GROUP BY takes a while to be completed.





- c) List the services based on the received destination packet's bytes and the tcp protocol.

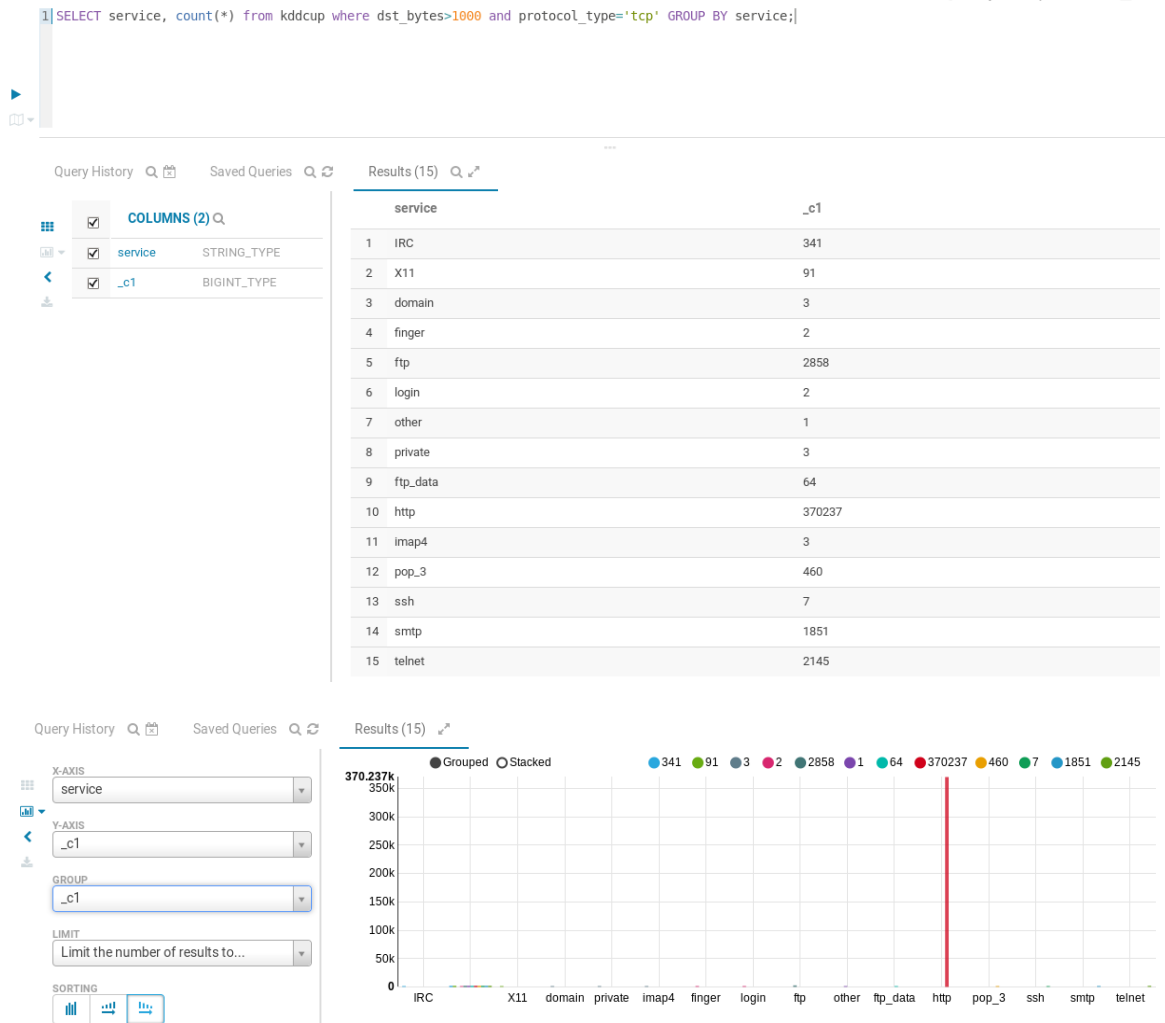
```
SELECT DISTINCT service from kddcup where dst_bytes>1000 and protocol_type='tcp';
```



Since there is only one value/column, we cannot use 2D plots.

- d) List the number of services based on the received destination packet's bytes and the tcp protocol.

```
SELECT service, count(*) from kddcup where dst_bytes>1000 and protocol_type='tcp' GROUP BY service;
```





Phase 4: Your turn to make more queries

It is your turn to practice HIVE queries using Analytic and Mathematical and many more functions. You should inquire the KDDCup dataset to create more information from this big sized dataset. For instance, provide useful information about the number of attacks (and sub-attacks) and normal connections. You can summarize the data based on the several features, such as the status of the connections, protocol types, services, server error rates, etc.

