# CS 327E Lecture 13
## Shirley Cohen

November 21, 2016

# Plan for Today

- Reading Quiz

- MySQL + JSON

- Final Project Assignment

# Readings for Today

- JSON Data Type from the MySQL Reference Manual (section 12.6)

- JSON Functions from the MySQL Reference Manual (section 13.16)

# Question 1

Which of the following statements is **false** about MySQL support for JSON:

A.    MySQL has a native JSON datatype

B.    MySQL converts JSON documents to a binary format

C.    MySQL indexes JSON documents directly

D.    MySQL supplies a number of functions for operating on JSON data

E.    MySQL automatically validates JSON documents that are stored in JSON columns

# Question 2

In the context of JSON, `normalization` means the process of removing duplicate keys (or names) from a document.

A.  True

B.  False

# Question 3

```
 1 ▾ {
 2     "text": "RT @ESW_UTEXAS: Meeting tonight at 7pm in JGB 2.216! We'll
           have Chick-fil-A and a rep from @airliquidegroup will present!
           #utexas",
 3     "retweeted_status":
 4 ▾   {
 5       "entities":
 6 ▾     {
 7         "user_mentions":
 8 ▾       [
 9 ▾         {
10             "id": 104481993,
11             "indices":
12 ▾           [
13               75,
14               91
15             ]
16           }
17         ]
18       }
19     }
20   }
```

What does the path expression below evaluate to?

`$.retweeted_status.entities.`
`user_mentions[0].indices[0]`

A.   {}

B.   75

C.   NULL

D.   91

E.   "id" : 104481993

# Question 4

```
 1 ▾ {
 2      "text": "RT @ESW_UTEXAS: Meeting tonight at 7pm in JGB 2.216! We'll
          have Chick-fil-A and a rep from @airliquidegroup will present!
          #utexas",
 3      "retweeted_status":
 4 ▾    {
 5        "entities":
 6 ▾      {
 7          "user_mentions":
 8 ▾        [
 9 ▾          {
10             "id": 104481993,
11             "indices":
12 ▾           [
13               75,
14               91
15             ]
16           }
17         ]
18       }
19     }
20  }
```

Let `j` be to the JSON document shown. What does the `select` statement below evaluate to?

```
select JSON_EXTRACT(j,
'$.retweeted_status.entities.
user_mentions[0].id');
```

A.    "indices" : [75, 91]

B.    75

C.    91

D.    104481993

E.    None of the above

# Question 5

```
 1 ▾ {
 2      "text": "RT @hiangieee: The students of @UTAustin took it to the streets today.
         #utprotest #Morningafter https://t.co/JvTLajTqT8",
 3      "retweet_count": 6236,
 4      "id_str": "796541058192736256",
 5      "favorited": false,
 6 ▾    "user": {
 7        "profile_background_image_url_https": "https://pbs.twimg.com
           /profile_background_images/569242807726915584/WDTmWx2F.jpeg",
 8 ▾      "entities": {
 9 ▾        "description": {
10            "urls": []
11          }
12        },
13        "followers_count": 73,
14        "utc_offset": -21600
15      },
16 ▾    "metadata": {
17        "iso_language_code": "en",
18        "result_type": "recent"
19      }
20  }
```

Let `j` be the JSON document shown. What does the `select` statement below evaluate to?

```
select JSON_SEARCH(j, 'one', 'RT');
```

A. "$.text"

B. "$.*"

C. {}

D. NULL

E. None of the above

# Demo 1

# UT Class Enrollment & Twitter



Logical ERD - UT Class Enrollment - CS 327E Fall 2016

# New DDL

```
 3   drop table if exists Major;
 4   create table Major (
 5      code int auto_increment primary key,
 6      name varchar(32) not null,
 7      college varchar(32) not null
 8   );
```

```
26   alter table Student add column major_code int;
27   alter table Student add constraint fk_major_code
28      foreign key (major_code) references Major(code);
```

```
42   drop table if exists Tweet;
43   create table Tweet (
44      tweet_id varchar(32) generated always
45         as (json_unquote(json_extract(tweet_doc, '$.id_str'))) stored primary key,
46      screen_name varchar(32) generated always
47         as (json_unquote(json_extract(tweet_doc, '$.user.screen_name'))) stored,
48      created_at datetime generated always
49         as (str_to_date(json_unquote(json_extract(tweet_doc, '$.created_at')),
50            '%a %b %d %H:%i:%s +0000 %Y')) stored,
51      tweet_doc json,
52      major_code int,
53      foreign key (major_code) references Major(code)
54   );
```

# Twitter Client

```python
16  def do_data_pull(api_inst):
17
18      sql_query = "select code, name from Major order by name"
19
20      try:
21          conn = create_connection()
22          db_cursor = conn.cursor()
23          query_status = run_stmt(db_cursor, sql_query)
24          resultset = db_cursor.fetchall()
25
26          for record in resultset:
27              major_code = record[0]
28              major_name = record[1]
29
30              utexas_query = "(#UTexas OR @UTAustin OR url:utexas.edu) AND "
31              twitter_query = utexas_query + "'" + major_name + "'"
32              print "twitter_query: " + twitter_query
33              twitter_cursor = tweepy.Cursor(api_inst.search, q=twitter_query, lang="en")
34
35              for page in twitter_cursor.pages():
36                  for item in page:
37                      json_str = json.dumps(item._json)
38                      print "found a " + major_name + " tweet"
39                      insert_stmt = "insert into Tweet(tweet_doc, major_code) values(%s, %s)"
40                      run_prepared_stmt(db_cursor, insert_stmt, (json_str, major_code))
41                      do_commit(conn)
42
43      except pymysql.Error as error:
44          is_success = False
45          print "do_data_pull: " + e.strerror
```

# Demo 2

# Concept Question 1

We want to extend the Twitter Client to check for duplicate tweets before doing the insert into MySQL. Assume that in Python we extract the id of the tweet and store the value in the variable `$id`. How can we formulate a SQL query that checks for duplicate tweets given `$id`?

```
43  create table Tweet (
44     tweet_id varchar(32) generated always
45        as (json_unquote(json_extract(tweet_doc, '$.id_str'))) stored primary key,
46     screen_name varchar(32) generated always
47        as (json_unquote(json_extract(tweet_doc, '$.user.screen_name'))) stored,
48     created_at datetime generated always
49        as (str_to_date(json_unquote(json_extract(tweet_doc, '$.created_at')),
50           '%a %b %d %H:%i:%s +0000 %Y')) stored,
51     tweet_doc json,
52     major_code int,
53     foreign key (major_code) references Major(code)
54  );
```

A. select count(*) from Tweet where tweet_id = `$id`

B. select * from Tweet order by tweet_id

C. select count(distinct tweet_id) from Tweet where tweet_id = `$id`

D. select * from Tweet where tweet_id = `$id`

E. select count(tweet_id) from Tweet where tweet_id = `$id`

# Concept Question 2

We want to implement a more accurate count of tweets per UT major. More specifically, we want to filter out all retweets and only add up the origin tweets. Assume that for all tweets, we extract the origin tweet id from the tweet and we store this value in a new field called `Tweet.origin_tweet_id`. How can we modify the query below to only count unique tweets?

```
select m.name, m.code, count(t.tweet_id) as tweet_count
from Major m left outer join Tweet t
on m.code = t.major_code
group by m.name, m.code
order by tweet_count desc;
```

A. Replace: `count(t.tweet_id)` with: `count(t.origin_tweet_id)`
B. Replace: `count(t.tweet_id)` with: `count(distinct t.origin_tweet_id)`
C. Change the outer join to an inner join
D. Change the left outer join to a right outer join
E. None of the above

# Final Project

http://www.cs.utexas.edu/~scohen/project/final_project.pdf