

Final Project Documentation

Yelp-Database

The Yelp dataset is a subset of our businesses, reviews, and user data for use in personal, educational, and academic purposes. Available as JSON files.["<https://www.yelp.com/dataset>"]. This is a big database with more than 6M rows of data so we need to make sure our scripts can handle such requirements.

Scope of the project:

- a) The scope of this project is to design and create a clean relational database for the publicly available yelp database which is given in JSON format.
- b) Querying this dataset to analyze service based businesses and calculate key metrics to further understand the market.

Example of Yelp Database JSON:

Each file is composed of a single object type, one JSON-object per-line.

Take a look at some examples to get you started: <https://github.com/Yelp/dataset-examples>.

Note: the follow examples contain inline comments, which are technically not valid JSON. This is done here to simplify the documentation and explaining the structure, the JSON files you download will not contain any comments and will be fully valid JSON.

business.json

Contains business data including location data, attributes, and categories.

```
{  
  // string, 22 character unique string business id  
  "business_id": "tnhfDv5lI8EaGSXZGiuQGg",  
  
  // string, the business's name  
  "name": "Garaje",  
  
  // string, the full address of the business  
  "address": "475 3rd St",  
  
  // string, the city  
  "city": "San Francisco",  
  
  // string, 2 character state code, if applicable  
  "state": "CA",  
  
  // string, the postal code  
  "postal code": "94107",  
  
  // float, latitude  
  "latitude": 37.7817529521,  
  
  // float, longitude  
  "longitude": -122.39612197,
```

```

// float, star rating, rounded to half-stars
"stars": 4.5,

// integer, number of reviews
"review_count": 1198,

// integer, 0 or 1 for closed or open, respectively
"is_open": 1,

// object, business attributes to values. note: some attribute values might be objects
"attributes": {
  "RestaurantsTakeOut": true,
  "BusinessParking": {
    "garage": false,
    "street": true,
    "validated": false,
    "lot": false,
    "valet": false
  },
},

// an array of strings of business categories
"categories": [
  "Mexican",
  "Burgers",
  "Gastropubs"
],

// an object of key day to value hours, hours are using a 24hr clock
"hours": {
  "Monday": "10:00-21:00",
  "Tuesday": "10:00-21:00",
  "Friday": "10:00-21:00",
  "Wednesday": "10:00-21:00",
  "Thursday": "10:00-21:00",
  "Sunday": "11:00-18:00",
  "Saturday": "10:00-21:00"
}
}

```

review.json

Contains full review text data including the user_id that wrote the review and the business_id the review is written for.

```

{
  // string, 22 character unique review id
  "review_id": "zdSx_SD6obEhz9VrW9uAWA",

  // string, 22 character unique user id, maps to the user in user.json
  "user_id": "Ha3iJu77CxlrFm-vQRs_8g",

  // string, 22 character business id, maps to business in business.json
  "business_id": "tnhfDv5lI8EaGSXZGiuQQg",
}

```

```

// integer, star rating
"stars": 4,

// string, date formatted YYYY-MM-DD
"date": "2016-03-09",

// string, the review itself
"text": "Great place to hang out after work: the prices are decent, and the ambience is fun.
It's a bit loud, but very lively. The staff is friendly, and the food is good. They have a good
selection of drinks.",

// integer, number of useful votes received
"useful": 0,

// integer, number of funny votes received
"funny": 0,

// integer, number of cool votes received
"cool": 0
}

```

user.json

User data including the user's friend mapping and all the metadata associated with the user.

```

{
  // string, 22 character unique user id, maps to the user in user.json
  "user_id": "Ha3iJu77CxlrFm-vQRs_8g",

  // string, the user's first name
  "name": "Sebastien",

  // integer, the number of reviews they've written
  "review_count": 56,

  // string, when the user joined Yelp, formatted like YYYY-MM-DD
  "yelping_since": "2011-01-01",

  // array of strings, an array of the user's friend as user_ids
  "friends": [
    "wqoXYLWmpkEH0YvTmHBsJQ",
    "KUXLLiJGrtSsapmxmpvTA",
    "6e9rJKQC3n0RSKyHLViL-Q"
  ],

  // integer, number of useful votes sent by the user
  "useful": 21,

  // integer, number of funny votes sent by the user
  "funny": 88,

  // integer, number of cool votes sent by the user
  "cool": 15,

  // integer, number of fans the user has
  "fans": 1032,
}

```

```

// array of integers, the years the user was elite
"elite": [
    2012,
    2013
],

// float, average rating of all reviews
"average_stars": 4.31,

// integer, number of hot compliments received by the user
"compliment_hot": 339,

// integer, number of more compliments received by the user
"compliment_more": 668,

// integer, number of profile compliments received by the user
"compliment_profile": 42,

// integer, number of cute compliments received by the user
"compliment_cute": 62,

// integer, number of list compliments received by the user
"compliment_list": 37,

// integer, number of note compliments received by the user
"compliment_note": 356,

// integer, number of plain compliments received by the user
"compliment_plain": 68,

// integer, number of cool compliments received by the user
"compliment_cool": 91,

// integer, number of funny compliments received by the user
"compliment_funny": 99,

// integer, number of writer compliments received by the user
"compliment_writer": 95,

// integer, number of photo compliments received by the user
"compliment_photos": 50
}

```

checkin.json

Checkins on a business.

```

{
    // string, 22 character business id, maps to business in business.json
    "business_id": "tnhfDv5lI8EaGSXZGiuQGg"

    // string which is a comma-separated list of timestamps for each checkin, each with format
    // YYYY-MM-DD HH:MM:SS
    "date": "2016-04-26 19:49:16, 2016-08-30 18:36:57, 2016-10-15 02:45:18, 2016-11-18
    01:54:50, 2017-04-20 18:39:06, 2017-05-03 17:58:02"
}

```

```
}
```

tip.json

Tips written by a user on a business. Tips are shorter than reviews and tend to convey quick suggestions.

```
{
  // string, text of the tip
  "text": "Secret menu - fried chicken sando is da bombbbbbbb Their zapatos are good too.",

  // string, when the tip was written, formatted like YYYY-MM-DD
  "date": "2013-09-20",

  // integer, how many compliments it has
  "compliment_count": 172,

  // string, 22 character business id, maps to business in business.json
  "business_id": "tnhfDv5lI8EaGSXZGiuQGg",

  // string, 22 character unique user id, maps to the user in user.json
  "user_id": "49JhAJh8vSQ-vM4Aourl0g"
}
```

photo.json

Contains photo data including the caption and classification (one of "food", "drink", "menu", "inside" or "outside").

```
{
  // string, 22 character unique photo id
  "photo_id": "_nN_DhLXkfwEkwPNxne9hw",
  // string, 22 character business id, maps to business in business.json
  "business_id": "tnhfDv5lI8EaGSXZGiuQGg",
  // string, the photo caption, if any
  "caption": "carne asada fries",
  // string, the category the photo belongs to, if any
  "label": "food"
}
```

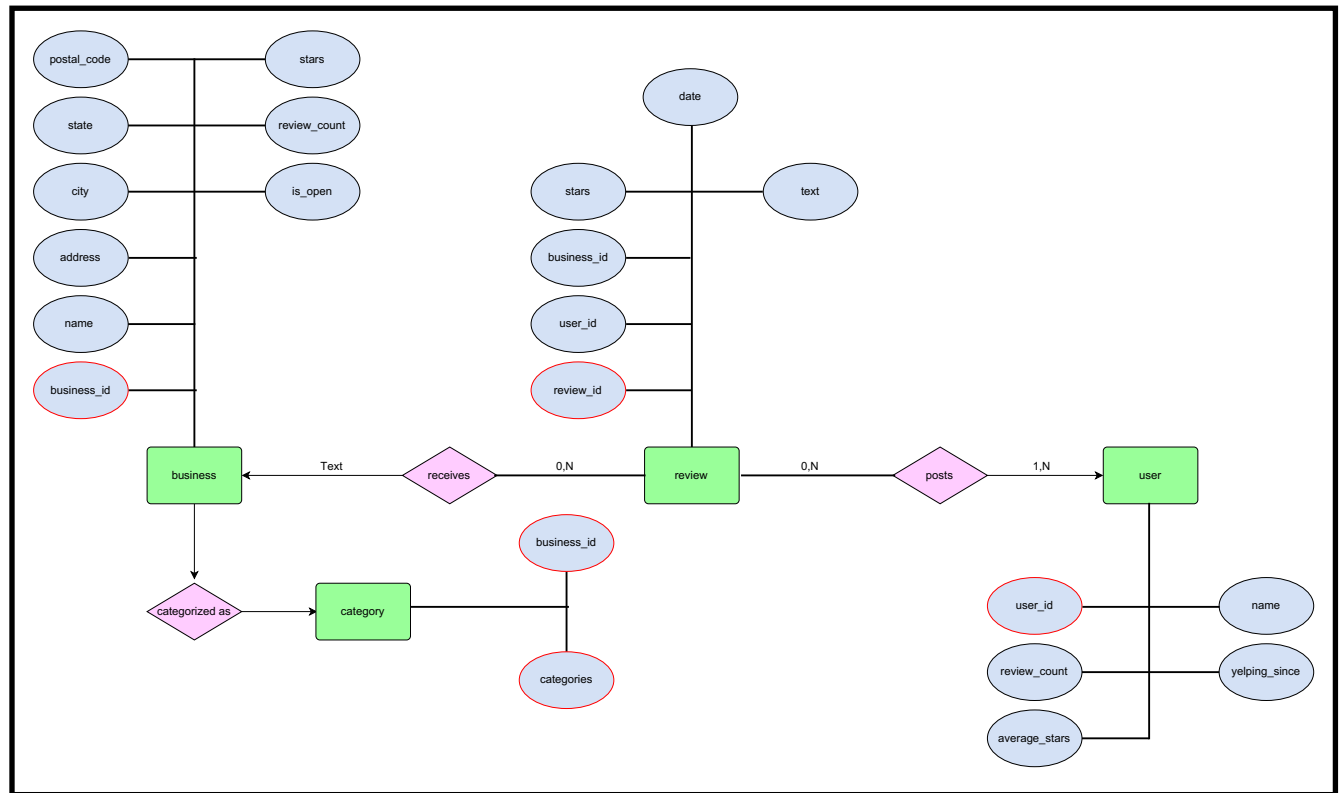
Questions we aim to answer:

1. Name of the user that has posted the most reviews.
2. Number of businesses in each category
3. Average stars of business in the state of FLORIDA in each postal code
4. Category with the most reviews
5. What percentage of existing Home service businesses on yelp that have happy customers ?
6. What is the number of reviews per year in the home services category ?
7. What is the average rate of reviews written per year ?
8. What is the number of happy customers in each category ?
9. How many businesses have shutdown per category?
10. What is the number of businesses per category in MA ?

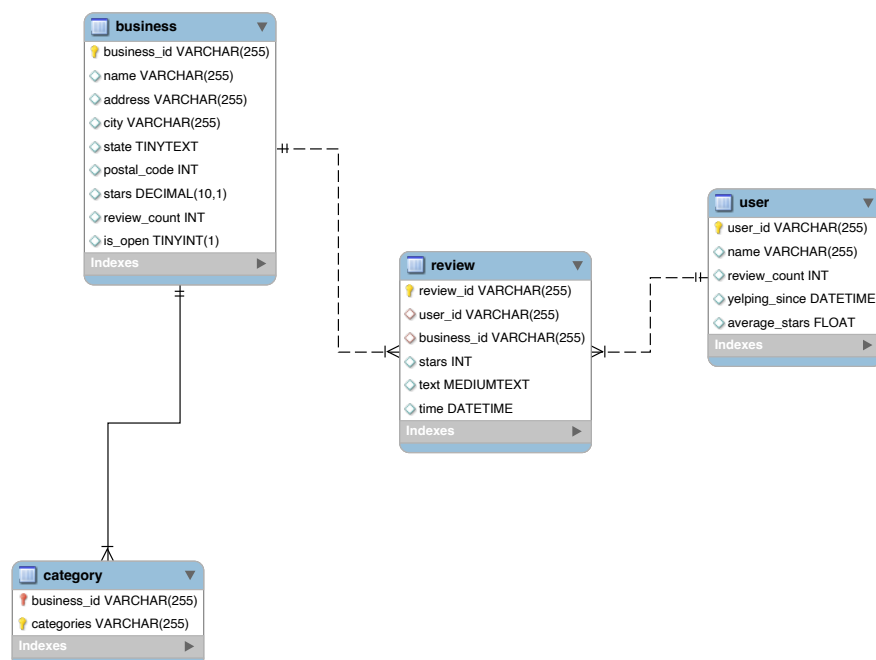
Designing the SQL Database

Keeping in mind the scope of the project and the analysis we aim to perform, we structured the data fulfill these objectives.

Entity-Relationship Diagram



UML Diagram



Note:

As is evident from the ERD and UML the unstructured JSON data available to us is converted into a clean structured form by dropping redundant and irrelevant data. All tables made are in 3rd Normal Form that means:

- a) All data is atomic.
- b) There are no partial dependencies.
- c) There are no transitive dependencies.

SQL Table Creation and Data Insertion

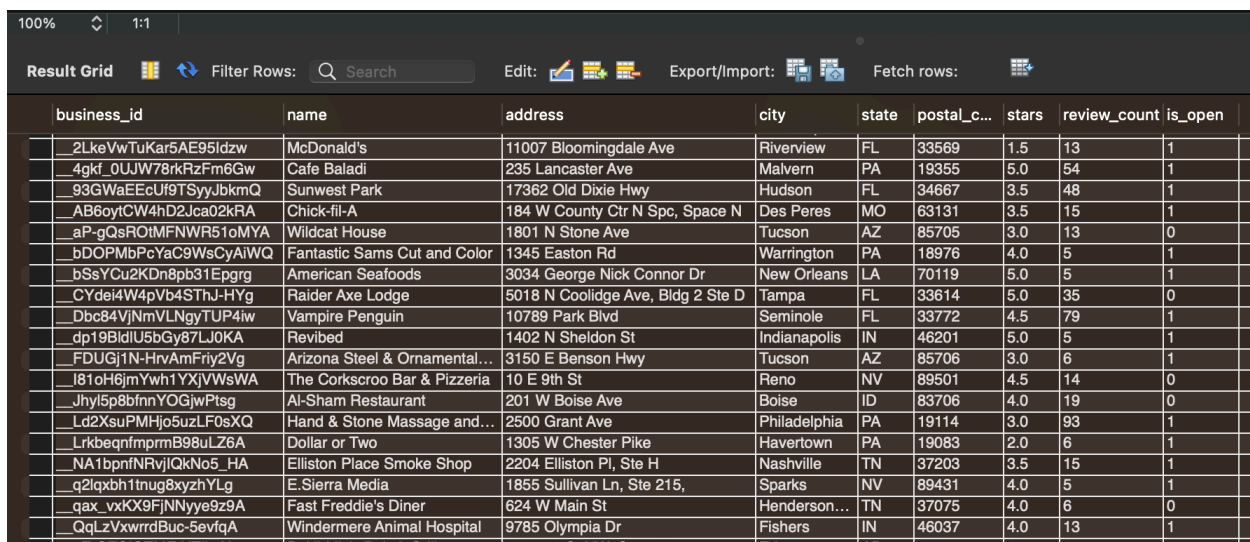
[data_insertion.ipynb](#): script that converts the json dataset into structured data and inserts the data into the designed mysql database.

This way of inserting data takes a lot of time and isn't efficient to handle big datasets such as ours so we dig deeper and optimized the process and reduced the run time significantly.

[optimized_data_insertion.ipynb](#): script that is highly optimized to convert yelp JSON dataset which contains more 6M rows into structured and inserting that data into the database much faster than data_insertion.ipynb.

Output:

business Table



business_id	name	address	city	state	postal_c...	stars	review_count	is_open
2LkeVwTuKar5AE95ldzw	McDonald's	11007 Bloomingdale Ave	Riverview	FL	33569	1.5	13	1
4gkf_0UJW78rkRzFm6Gw	Cafe Baladi	235 Lancaster Ave	Malvern	PA	19355	5.0	54	1
93GWaEEcUf9TSyyJbkMQ	Sunwest Park	17362 Old Dixie Hwy	Hudson	FL	34667	3.5	48	1
AB6oytCW4hD2Jca02kRA	Chick-fil-A	184 W County Ctr N Spc, Space N	Des Peres	MO	63131	3.5	15	1
aP-gQsROtMFNWR51oMYA	Wildcat House	1801 N Stone Ave	Tucson	AZ	85705	3.0	13	0
bDOPMbPcYaC9WsCyAiWQ	Fantastic Sams Cut and Color	1345 Easton Rd	Warrington	PA	18976	4.0	5	1
bSsYCu2KDn8pb31Epgrg	American Seafoods	3034 George Nick Connor Dr	New Orleans	LA	70119	5.0	5	1
CYdei4W4pVb4SThJ-HYg	Raider Axe Lodge	5018 N Coolidge Ave, Bldg 2 Ste D	Tampa	FL	33614	5.0	35	0
Dbc84VjNmVLNgYTUP4iw	Vampire Penguin	10789 Park Blvd	Seminole	FL	33772	4.5	79	1
dp19BldIU5bGy87LJ0KA	Revived	1402 N Sheldon St	Indianapolis	IN	46201	5.0	5	1
FDUGj1N-HrvAmFriy2Vg	Arizona Steel & Ornamental...	3150 E Benson Hwy	Tucson	AZ	85706	3.0	6	1
J81oH6jmYwh1YXjVWsWA	The Corkscroo Bar & Pizzeria	10 E 9th St	Reno	NV	89501	4.5	14	0
Jhyl5p8bfnnYOGjwPtsg	Al-Sham Restaurant	201 W Boise Ave	Boise	ID	83706	4.0	19	0
Ld2XsuPMHjo5uzLF0sXQ	Hand & Stone Massage and...	2500 Grant Ave	Philadelphia	PA	19114	3.0	93	1
LrkbeqnfmpmB98uLZ6A	Dollar or Two	1305 W Chester Pike	Havertown	PA	19083	2.0	6	1
NA1bpnfNRvjIQkNo5_HA	Ellistn Place Smoke Shop	2204 Ellistn Pl, Ste H	Nashville	TN	37203	3.5	15	1
q2lqxbh1tnug8xyzhYLg	E.Sierra Media	1855 Sullivan Ln, Ste 215,	Sparks	NV	89431	4.0	5	1
qax_vxKX9fJNNyye9z9A	Fast Freddie's Diner	624 W Main St	Henderson...	TN	37075	4.0	6	0
QqLzVxwrrdBuc-5evlqA	Windermere Animal Hospital	9785 Olympia Dr	Fishers	IN	46037	4.0	13	1

100% 1:1

Result Grid Filter Rows: Search Edit: Export/Import:

	business_id	categories	
▶	___UdvaxCnwsQ7nA1eKZAQ	Hotels & Travel	
	___UdvaxCnwsQ7nA1eKZAQ	Transportation	
	___UdvaxCnwsQ7nA1eKZAQ	Buses	
	__2LkeVwTuKar5AE95ldzw	Burgers	
	__2LkeVwTuKar5AE95ldzw	Coffee & Tea	
	__2LkeVwTuKar5AE95ldzw	Food	
	__2LkeVwTuKar5AE95ldzw	Restaurants	
	__2LkeVwTuKar5AE95ldzw	Fast Food	
	__4gkf_0UJW78rkRzFm6Gw	Halal	
	__4gkf_0UJW78rkRzFm6Gw	Lebanese	
	__4gkf_0UJW78rkRzFm6Gw	Middle Eastern	
	__4gkf_0UJW78rkRzFm6Gw	Restaurants	
	__4gkf_0UJW78rkRzFm6Gw	Mediterranean	
	__93GWaEEcUf9TSyyJbkmQ	Active Life	
	__93GWaEEcUf9TSyyJbkmQ	Beaches	
	__93GWaEEcUf9TSyyJbkmQ	Parks	
	__93GWaEEcUf9TSyyJbkmQ	Rafting/Kayak...	

	user_id	name	review_count	yelping_since	average_stars
▶	___6aix-XvFcQz3GauAPpw	John	7	2011-09-26 18:10:05	2.75
	___9Jl-8aF7z58JdTWlyjw	Erin	10	2011-04-23 23:47:24	3.9
	___aO4EhZULBsJPvJ4rMhA	David	14	2010-12-04 22:21:55	1.69
	___DjEwDb9e7Ny-NOiezZw	Patama	10	2013-10-28 19:29:04	4.8
	___hQj63mwwgFDhehNwv3ZQ	Taleena	11	2014-10-25 21:26:56	4.09
	___l9ZYdYGkZ6dMYxwJEIQ	Jim	239	2011-08-17 17:35:52	3.99
	___Pul-lcE0y9hjAer8GrQ	DeAnna	1	2015-02-15 15:49:36	1
	___tX0MgAQYPaWssEjSxKw	Mary	2	2016-02-22 00:12:38	3
	___-6h87PZrkaT0SAuQc-w	Jiban	2	2016-10-03 03:57:50	5
	___-2EyiraNzLPrq6o1QDIA	Dave	2	2018-07-06 17:09:08	4.5
	___-6iHRGqL_K9oM9KNs6pg	Morgan	22	2014-12-21 19:58:02	2.39
	___-8WEg7xD0CwCDu04MzBA	Maria	39	2010-05-03 16:50:07	4
	___-Kt26YrtJxGdWs8FqKCg	Shane	11	2012-08-06 04:44:19	4
	___-LqRflgxlTHOxsGMBt_Q	Craig	9	2012-11-20 22:57:48	3.67
	___-uwr6nLywqDa8d_QD4-A	Rose...	18	2010-03-28 14:46:49	3.5
	___-w2mUbDTIC6u1lklEWwg	Zach	13	2012-12-24 04:54:21	4.46
	___-WWzo7OEEed01pOA-rhPw	Ron	8	2011-11-21 13:53:41	1
	___-YOsZp7ilfYVwD8Wdszg	David	2	2014-10-03 14:11:46	3
	___-zs2o-TGiRzK4WqmgI2A	John	2	2017-04-12 04:12:03	1
	___05rytNjsye9MBhqB0DMA	Marck	1247	2006-02-12 20:20:32	3.65
	___0cufkRcZf12oTT04xGLg	Mat	36	2014-10-31 20:18:36	3.65
	___0D94KGQl7dBCcA2MmH...	Chris	48	2012-07-13 11:32:20	3.5
	___0TajwzDW-qJ2301ooWUg	B	11	2015-02-25 21:07:52	3.17

100%

1:1

Result Grid

Filter Rows:

Search

Edit:

Export/Import:

Fetch rows:

review_id	user_id	business_id	stars	text	time
7mnNcWBXBn76xXXsYyw	geqxRkOpMg7S3XmPx1LySA	I2vZuyG1KjtBtwINfHQA	5	Such great food and service ..not to mention AM...	2019-08-23 19:27:18
881evFAzgZxmEvHy-v	NbOt9ikm5Y14d_ZiGyIf7w	lqW5vulkrfF9a2hiLgrRw2s	5	Excellent. Juicy, well-seasoned salmon. Even th...	2018-04-08 17:40:51
B0ROhOE7Y_krtlls_bQ	1gm_7imVolEXgyY4Q_1mY6uA	9n-1QLX3ntBIjBlMWgsSpig	5	We are originally from northern New Jersey acr...	2012-02-08 16:16:56
R4r3-3EICmiWHBQ5Www	-kef2nuId6Cnhf5drLR07Sw	8wD5iNDnrJqI-TBtAO1LSQ	5	When I first entered King and I, I wondered, "wh...	2011-06-30 21:25:36
eElPf3S0sm6fmhTrpD1MA	YQnLRLCvc0-dTTgaH2qw1Q	krcJ4KkAsRymBJJT3CKauQ	4	Have eaten here twice for lunch. Food is tasty t...	2015-01-26 17:23:53
fGeCuKXRtMu6hmJkJ3u6A	qiOmSEFFEEsVWVtsIIoaaFw	pwiFIUAwzDXbXMGOOnrB-g	4	Delicious. Harvest quile possibly has the best f...	2006-11-14 00:40:00
GhU2U_3cmcbYy5ylfwg	3uwDTFPwohqh8Q3yYVs_FA	qupFYcm_e2lO5MPP9ZOAA	5	Amazing experience here. Will definitely buy mo...	2015-10-11 03:45:52
gu3aPxpLzQ-J8bi680Q	Z90K8RRx5jgAtCHvwzckJFA	UkyODD3LU4CTe7yNDhpnmOXg	4	I love brunch at Day By Day :) We go all the tim...	2010-10-19 04:00:42
hjNQv6lBHndhwG4Z4Ug	pd4vedOOUoh-iKAxSTdCXA	hxXiJG1Qlae0Dasfb-tZA	3	Drunches from the friends required me to make...	2017-08-31 18:19:30
ijx4svckVGMcPb73dTQ	EGrX9pbnd1s8z1uj1L_vdQ	wzkrZJ6Auiz1_WbOXvpTmA	1	I should have taken Matthew's comment more s...	2014-07-01 22:56:16
ld1ela-jAtpmCCCNxwg	qtldV6LwnXSbOTa7IC4P6w	ZYUxmZLH8le21j3e42Ufw	1	Cara was at an event the ended at a certain tim...	2016-09-19 17:15:12
qb8Z8oxTrXDder9Q5A	citrE57NKQMGPY3AHAKPdOQ	iuCCHFWmjIFbGYvgB9fw	5	This post is a year over due so in the summer of...	2020-08-20 23:25:16
qda8FYIGIXUSWcproA	J-Ta786tsYzz39KpeysgtQ	u.LhNGXB15Yln-ue70AG	5	After buying our home in South Philly, we broug...	2018-04-14 10:28:44
LQl4ph-5Eq506ypocqg	uD4ssusSuH1PgnaadpZsSQ	3a10MaqYLLUH8B1uEm02PW	5	I'm a fan for life after what happened with an or...	2015-12-28 19:34:37
RvluxpfVXgbg7zMEHKRA	3IkCuUlrgnC19jiXkBzSLA	feafesMR_MbxaOOS55HYUGA	2	I eat at Wasabi frequently for dinner... though it'	2012-06-15 19:40:38
S17UEEbotmf9sbbl1Q	JEA472SCepG4h1p6dNLwYpA	BLIQNJzmItG779flSHMd4gfA	5	The BEST Mexican food around, and prices. Ev...	2021-07-15 14:45:10
WG7GdQPccos8NYO0Ye	YkL5tsQVNI_wHkdrdNWGMp	i4GBHXV2E8ScODmx9xA	4	Great food and drinks and atmosphere. Had br...	2021-02-26 19:11:05
wXTTh9ak3AFYCr1_sbw	SYllSUtLiGbXYQRJhALw	jrcSiBNlw2zdVbnRTelSGF	5	Well Jack I'm in love with you once again! I'm n...	2018-11-25 20:29:06

Use Cases

use-cases.sql : SQL file that contains 5 use cases and its SQL statements

Example output:

```
20
21 -- use case 3: number of businesses in each state by category
22 SELECT
23     categories,
24     state,
25     COUNT(business.business_id) AS num_business
26 FROM
27     business
28 JOIN category ON business.business_id = category.business_id
29 GROUP BY categories, state
30 ORDER BY num_business desc;
```

100% 28:30

Result Grid

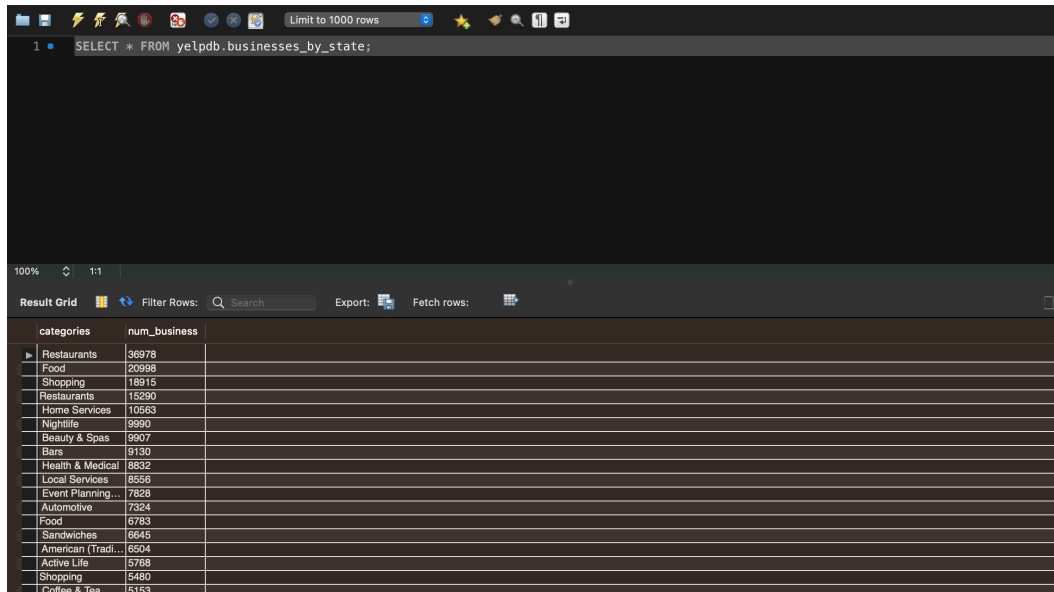
Filter Rows: Search Export: Fetch rows:

categories	num_business
Restaurants	36978
Food	20998
Shopping	18915
Restaurants	15290
Home Services	10563
Nightlife	9990
Beauty & Spas	9907
Bars	9130
Health & Medical	8832
Local Services	8556
Event Planning...	7828
Automotive	7324
Food	6783
Sandwiches	6645
American (Trad...	6504
Active Life	5768
Shopping	5480
Coffee & Tea	5159

Views

model_views.sql: This file creates views for all the use cases for our database.

Example Output:



The screenshot shows a database query tool interface. At the top, a SQL query is entered: `SELECT * FROM yelpdb.businesses_by_state;`. Below the query editor, the results are displayed in a grid format. The grid has two columns: `categories` and `num_business`. The results list various business categories and the number of businesses in each.

categories	num_business
Restaurants	36978
Food	20998
Shopping	18916
Restaurants	18290
Home Services	10583
Nightlife	9990
Beauty & Spas	9907
Bars	9130
Health & Medical	8832
Local Services	8556
Event Planning...	7828
Automotive	7324
Food	6783
Sandwiches	6645
American (Tradi...	6504
Active Life	5768
Shopping	5480
Coffee & Tea	5153

Caption