

YELP DOCUMENTATION

Firstly I created the default schema and relationship by the documentation provided by the YELP. Here, I have not used any kind of constraints.

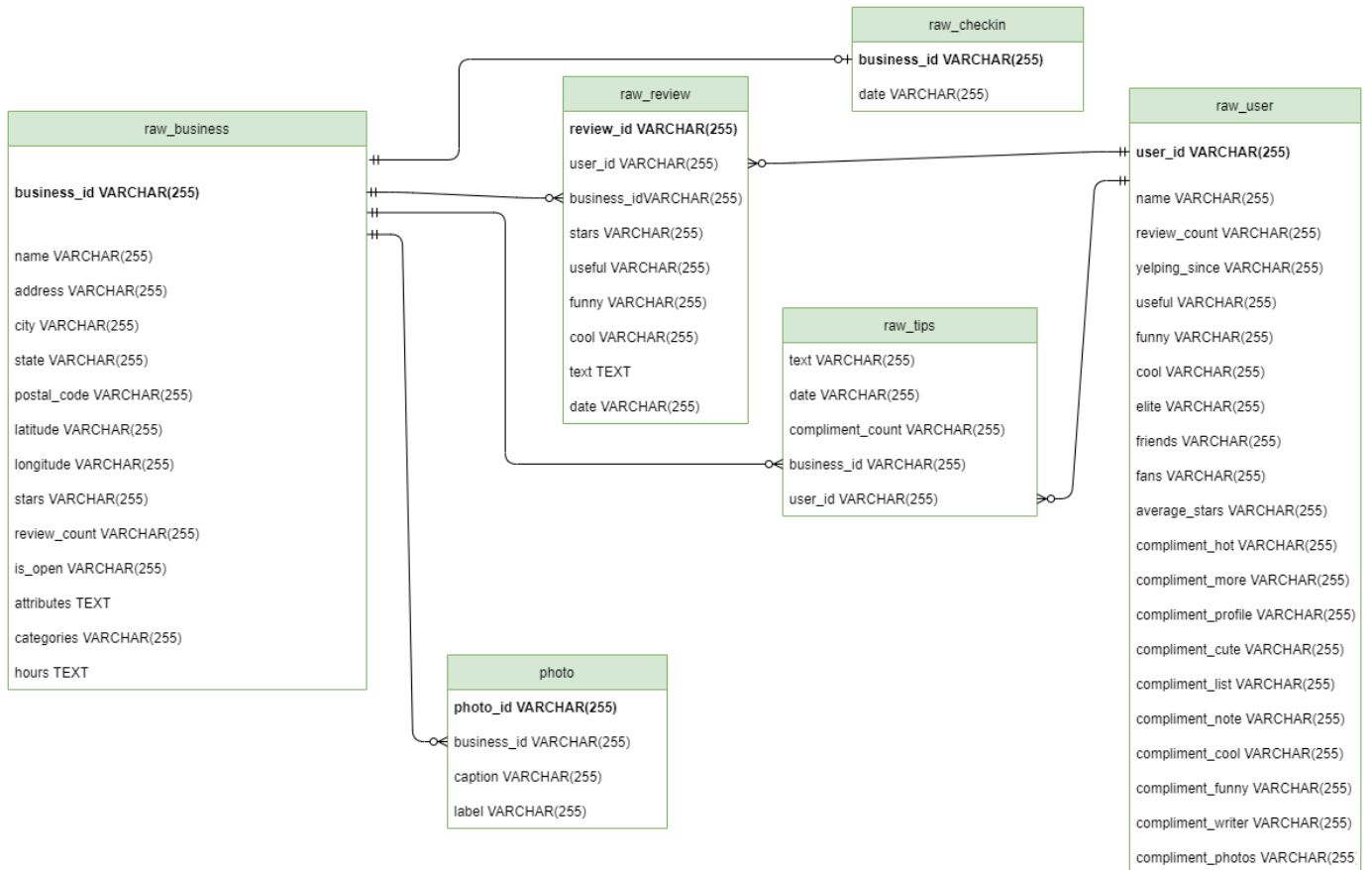


Fig:ER diagram of Raw tables

Logical Modelling:

After the data was pushed into the raw table, I had a lot of insight about the data and I started further exploring its entities and attributes and the data types for it.

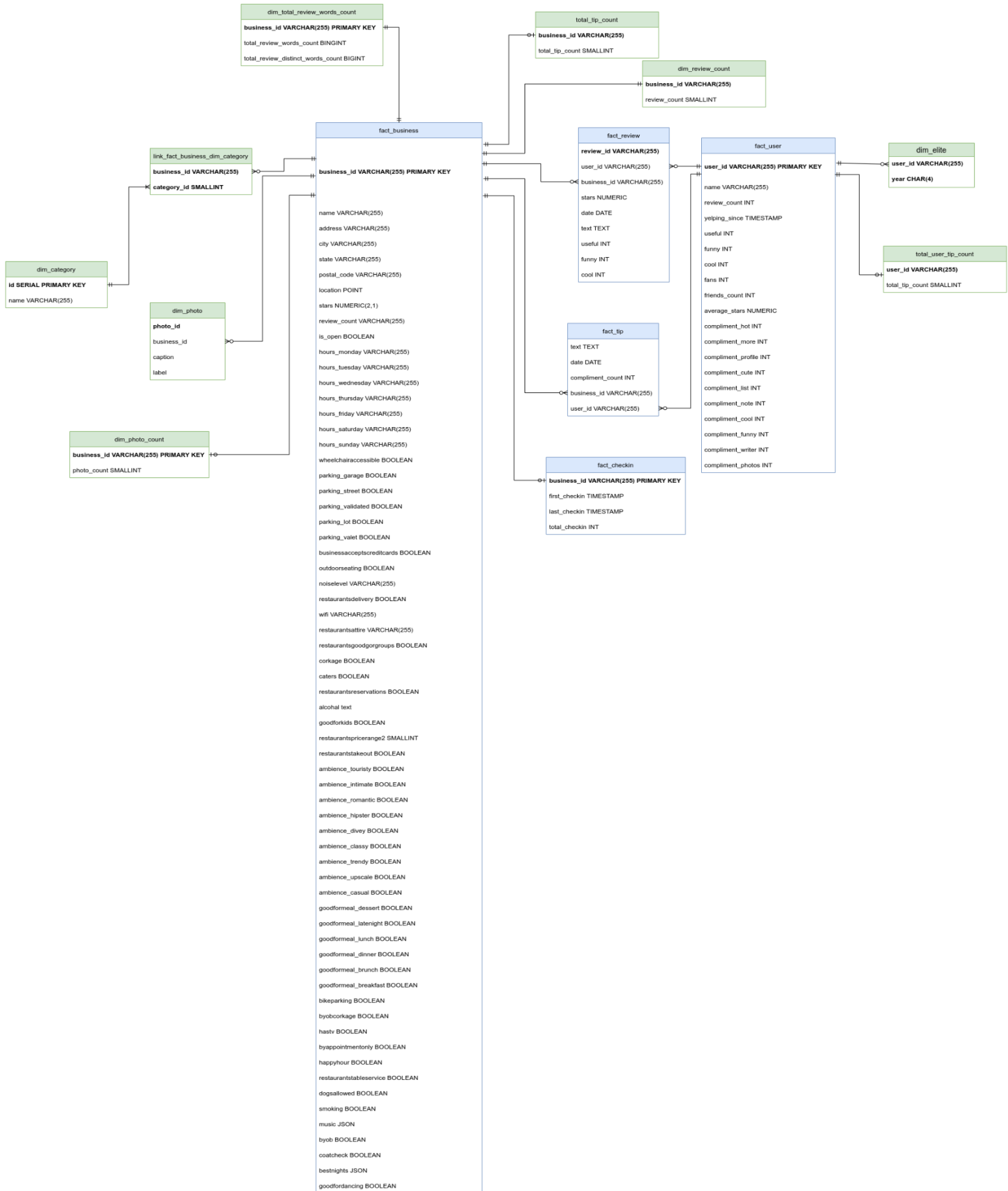
Below is the list of entities , description and domain of the model.

Entity	Description	Domain
dim_category	The categories of the business	
Attributes: id	Identifier for the entity,SK,FK	Auto generated,Serial

Attributes: review_id user_id business_id stars date text useful funny cool	Identifier of the entity,PK Identifier of the entity user,FK Identifier of the entity business,FK The stars given by the user to the business The date at which the review was given The review text The number of useful reaction by users The number of funny reaction by users The number of cool reaction by users	Text Text Text Numeric Date Text INT INT INT
fact_tip	The tip given by the user to the business.	
Attributes: text date compliment_count business_id user_id	The text about the tip given to the business by the user The date at which the tip was given The count of the compliment The business to which the tip is given The user who gives the tips .	Text Date INT Valid ID from table fact_business Valid ID from table fact_user
fact_checkin	The checkin done on the business.	
Attributes: business_id first_checkin last_checkin total_checkin	The identifier which references the business,FK The TIMESTAMP at which the first checkin was done The TIMESTAMP at which the last checkin was done The total number of checkin done	Valid ID from the table fact_business TIMESTAMP TIMESTAMP INT
fact_user	The info about the user	
Attributes: user_id name review_count useful funny cool fans friends_count average_stars compliment_hot compliment_more compliment_cute compliment_list compliment_note compliment_cool compliment_funny Compliment_writer compliment_photos	Identifier of the entity fact_user,PK The name of the user The total review count of the user The total number of useful reaction received to their reviews The total number of funny reaction received to their reviews The total number of cool reaction received to their reviews The total number of fans of the user The total number of friends of the user The number of compliments_hot The number of compliments_hot The number of compliment_more The number of compliment_cute The number of compliment_list The number of compliment_note The number of compliment_cool The number of compliment_funny The number of Compliment_writer The number of compliment_photos	PK Text INT INT INT INT NUMERIC INT INT INT INT INT INT INT INT INT INT INT
dim_elite	The year in which the user was elite.	

Attributes: user_id year	Identifier which references the user.PK The valid year at which the user was elite,PK	Valid ID referencing the fact_user table CHAR
total_user_tip_count		
Attributes: user_id total_tip_count	Identifier which references the user.PK The total number of tips count by user	Valid ID from fact_user table. INT

Proposed ER model
Below is the proposed ER diagram of the warehouse.



Validation:

Fan-Trap

It looks like there is a fan trap everywhere throughout the ER model, but the above fan trap is not going to affect our model as per my design.

Chasm-Trap

Since, the pathway exists between all of the entities, so there's no occurrence of the chasm-trap.

Physical Implementation:

1. Making raw tables
2. Pushing the data from the pipeline into the raw tables.
3. Data cleaning of the table.
 - Making the 'None' value as NULL.
 - Changing the data types of the attributes as per the model proposed.
 - Changing the values like 'no', 'uno', 'yes_corkage', 'yes_free', 'yes_free' to 0 and 1 and casting it to Boolean.
 - While splitting the words by comma on the category field of the business entity, spaced at the beginning was trimmed off to make consistency among the same category name.
 - Making the elite years '20,20' as the 2020.
4. Creating the schemas as proposed above in the ER model.
5. Pushing the data into the fact and dimension table by further cleaning the data if necessary.
6. Validation of different aspects of the data, such as
 - Checking the total photo_id and the total distinct photo_id is 0.
 - Checking all the unique photos are associated with the unique business_id
 - Checking if all the friends are not the user
 - Checking if the yelping_since is not in the future
 - Checking if the average_stars is not in between 0 and 5
 - The review count of a business is not equal to the provided reviews_count at fact_businesss.
7. Exporting the data into flat files for the visualization of data on power BI as I have DBMS on the linux system on my computer, and PowerBI is only supported on the windows.
8. Finally uploading the data on PowerBI for visualization.

CODE DESCRIPTION LINK:

https://github.com/callingsandesh/yelp/blob/final_project/docs/code_explanation.md

Visualization:

LINK:

https://github.com/callingsandesh/yelp/blob/final_project/docs/Dashboard%20and%20Report.pdf