

SUBJEC – DATA MANAGEMENT, WAREHOUSING, ANALYTICS
ASSIGNMENT – 1

SUBMITTED BY,
Dhruv Doshi
B00883311
dh722257@dal.ca

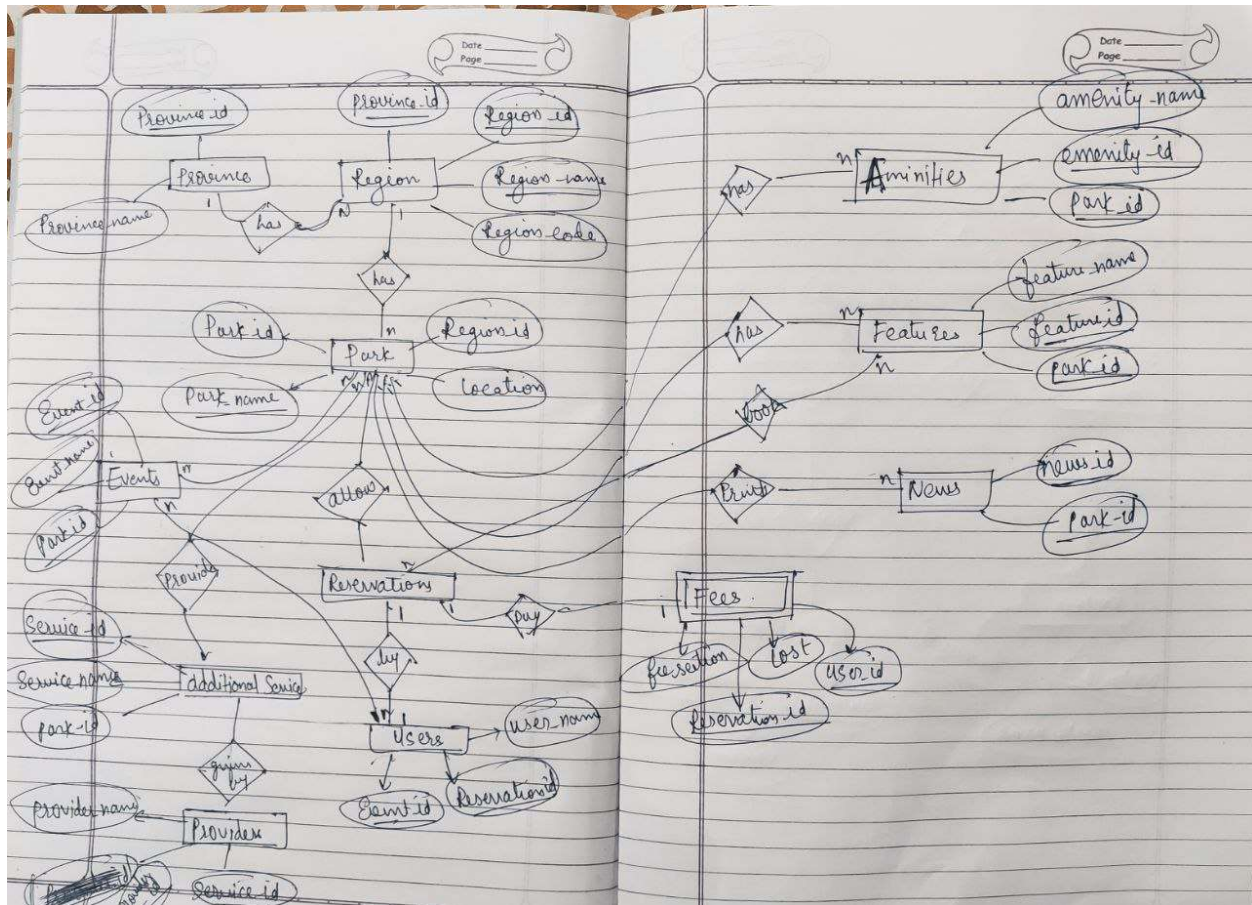
SUBMITETD TO: Prof. SAURABH DEY

Building a Data Model for Nova Scotia on its Provincial Parks.

After going through <https://parks.novascotia.ca/> I have founded these 12 entities and listed their reason for the selection in the given table.

ENTITY	REASON FOR SELECTION
Park	Parks are the center point for the whole system
Region	Region is directly connected with park and have distinct value
Province	Province could help further rectify the rules to be considered
Events	All events happening in the park.
Additional Service	Third party services
Providers	Third party services providers
Reservations	User could make reservations from here.
Users	All the data regarding guests
Features	All the additional paid features given by park
Amenities	All the free amenities given by park
Fees	Fees to be paid during registration
News	News published by the park
Staff	All the maintain staff for the park

After completion of finalizing the entities, I move forward to design the initial ER diagram for the parks. The pen and paper ERD is attached as a figure below.

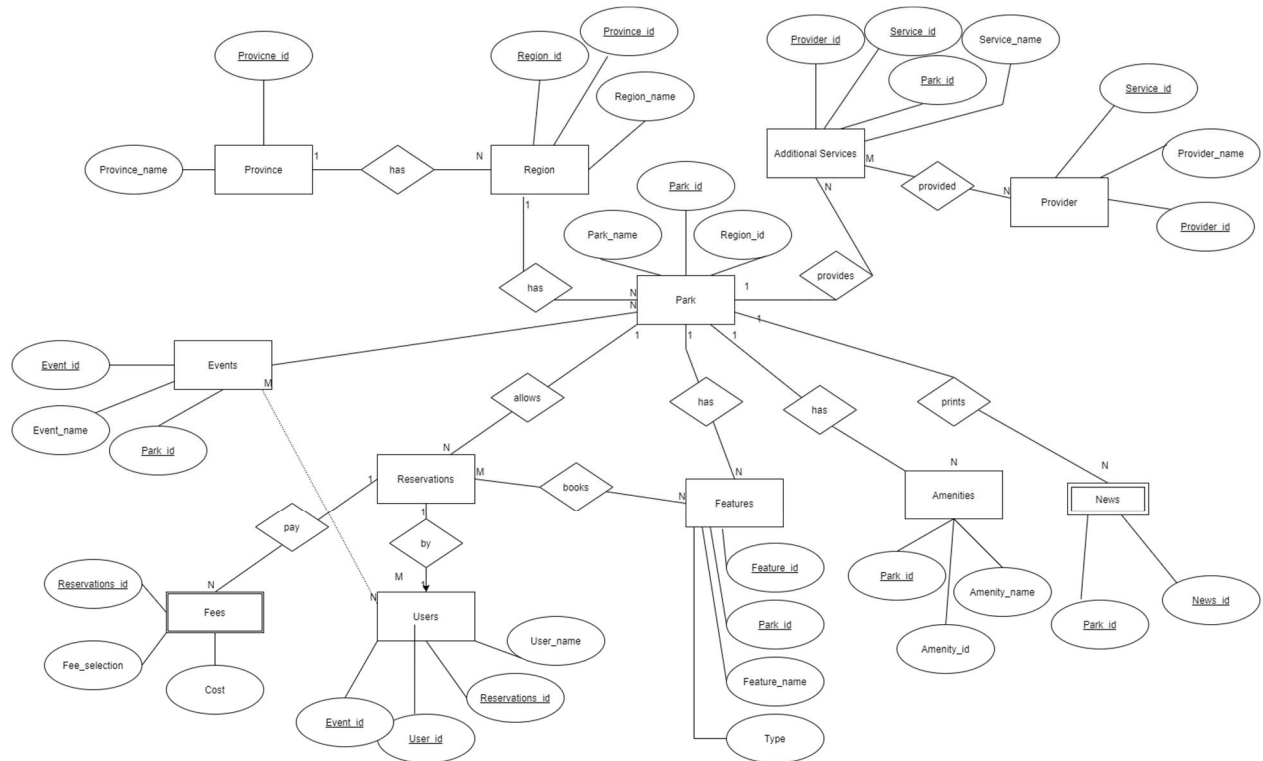


Finding the design issues, or attributes that were not considered, or entities which were found after going through the design. I have tried to cover all the aspects in the initial diagram itself still if any of the part is left then it will be corrected in this segment.

Points performed:

- Region_code is not required as it does not hold any substantial information for this use case.
- In Park we need to have Park_id as a primary key.
- Additional Service attribute should consist of the provider_id as an attribute because that will cast the connection between two tables.
- In Additional Services Park_id needed to be a primary key.
- Users entities need to have user_id attribute and that too need to be part of composite primary key.
- News will be the weak entity as it cannot be identified independently without accessing the parks along with the only way to have a primary key in it is by having the park_id and news_id being a composite key.
- There is multiple 1 to many relationships in the ER Diagram, but these all entity are individual and do not collide hence there will not be any case of fan trap here.
- There are not any many to one relationship hence no occurrence of chasm trap is identified in this ERD.

After performing the steps this will be the final ER Diagram.



Format Ocean Tracking Data and Report

The data given includes the information regarding some research work on the animals as there is vital information about the animals like their sex, age and length. Alongside that we also have information about the lab. The information about the datacenter, offload tag, detection, project and receivers is also given through the dataset. After the initial skimming I could say that the database are in bad shape and they need to be refactored which I will be doing in the later half of the assignment. In this Assignment we initially have what the changes I have performed on the datasets and why I needed to perform those steps. Continuing that it will have initial ER diagram and after that we will discuss the normalization steps and also have an ER diagram with normalization.

Datasets used for this task:

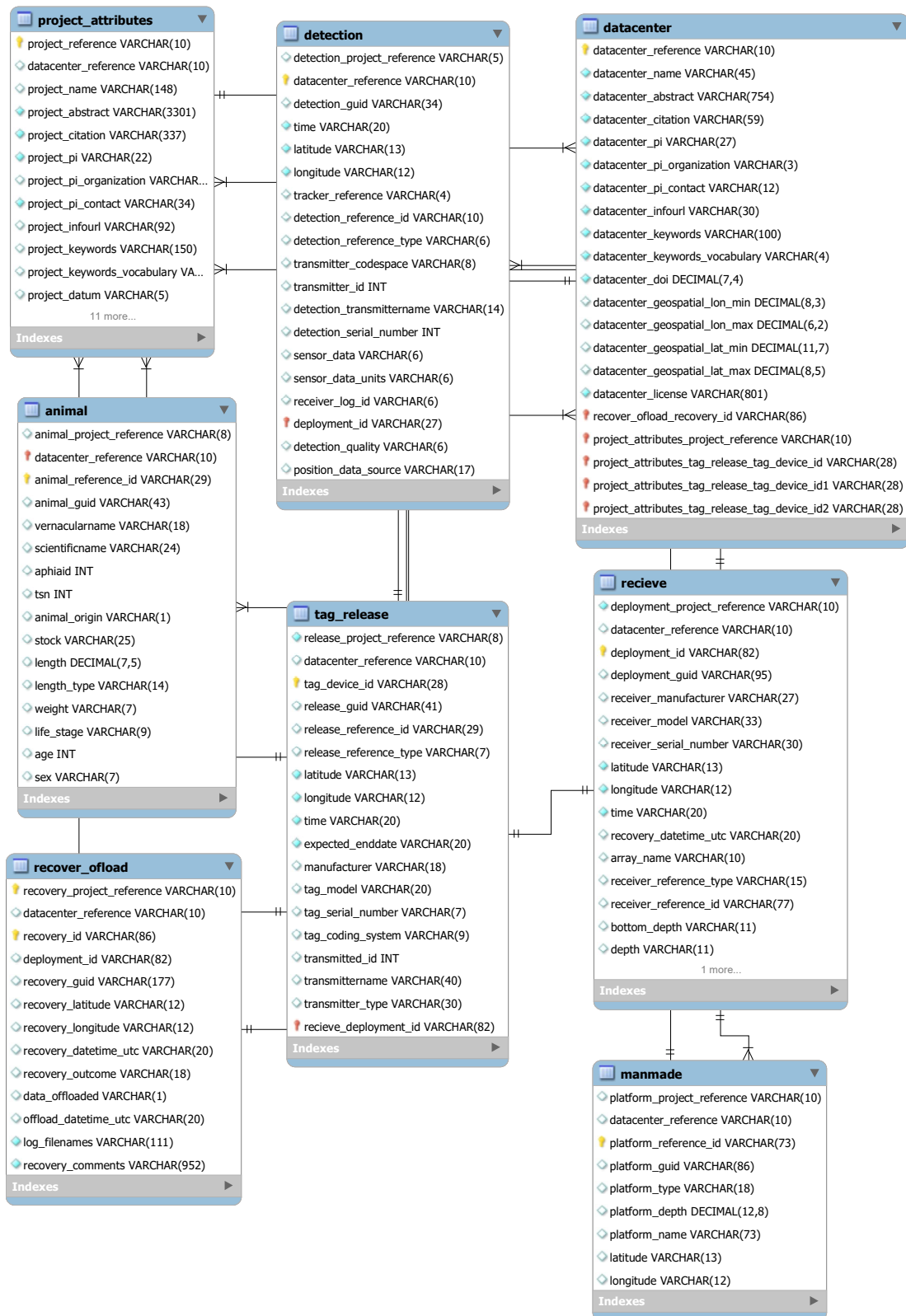
- Otnunit_att_animals
- Otnunit_aat_datacenter
- otnunit_aat_detection
- otnunit_aat_manmade
- otnunit_aat_project
- Otnunit_aat_receivers
- Otunit_aat_recover_offload_details
- Otnunit_aat_tag_releases

After cleaning and filtering the data

- Otnunit_att_animals
 - o Taxonrank – Empty column – Removed completely
 - o Animal_Origin – For every null value and aphiaid = 217628 and scientific name = Notorynchus cepedianus we will change the blank with W as per the previous records.
 - o Stock – There are three data inputs UNK, Unknown and Blanks which will be replaced by No Data.
 - o Length – The records where there is no data for stock are majorly NAN in length hence, we need to make sure there isn't any issue with any average hence we will change them by the average length of similar aphiaid. Where the information couldn't been found it had put the null values.
 - o Length_types – Blank rows with No Data.
 - o Life_stage – blank rows with No Data.
 - o Age – its int data type hence removing nan with NULL.
 - o Sex – Sex cannot be predicted hence we put UNKNOWN at blank positions.
- Otnunit_aat_datacenter
 - o Datacenter distribution statement – Empty hence removed
 - o Datacenter date modified – Empty hence removed
 - o time_coverage_start – Empty hence removed
 - o time_coverage_end – Empty hence removed
 - o datacenter_geospatial_lon_min, datacenter_geospatial_lon_max, datacenter_geospatial_lat_min, datacenter_geospatial_lat_max – wherever there is na values changes with NULL
- otnunit_aat_detection
 - o Sensor_data, Sensor_data_units – All blank is replaced with unknown.
 - o receiver_log_id, detection_quantity – All blank is replaced with unknown.
 - o Depth – the whole column does not hold any substantial values, hence removing the whole column.
 - o Uncertainty_in_latitude, Uncertainty_in_longitude – the whole column is nan hence removing the column.
 - o Depth_data_source, uncertainty_in_depth, other_position_data, dataset quantity – blank column, hence removing them.
- otnunit_aat_manmade
 - o platform_depth, latitude_degrees_north, longitude_degrees_north – wherever the data is nan change it with NULL.
- otnunit_aat_project
 - o time_coverage_end, time_coverage_start – whole column is empty hence remove the column.
 - o geospatial_vertical_positive, project_linestring – whole column is empty hence remove the column.
 - o Project_date_modified, project_distribution statement, project_doi, project_references – removing the whole column.
 - o Project_infourl – Blank and na cells are <NULL>
 - o Project_pi_contact, project_pi, project_citation, project_abstract – Blanks are converted to No Data.

- Otnunit_aat_receivers
 - Receiver_manufacturer, receiver_seriell_number, receiver_reference_id, deployment_comments – blank row converted to No Data
 - Frequencies_monitored, receiver_coding_scheme, deployed_by – blank column hence removed.
 - Recovery_datetime_utc – some blank cells are there replaced by NULL
 - Bottom_depth, depth – nan to NULL
 - Expected_reciever_life – removed because it doesn't hold any value
- Otunit_aat_recover_offload_details
 - Log_files, recovery_comments – blank replaced with No Data
 - Recovery_datetime_utc, offload_daytime_utc – nan or blank to NULL
 - Clock_synchronized, recovered_by – blank so removed
- Otnunit_aat_tag_releases
 - Tag_frequency, transmitter_type, tag_programming_id – blank hence removed

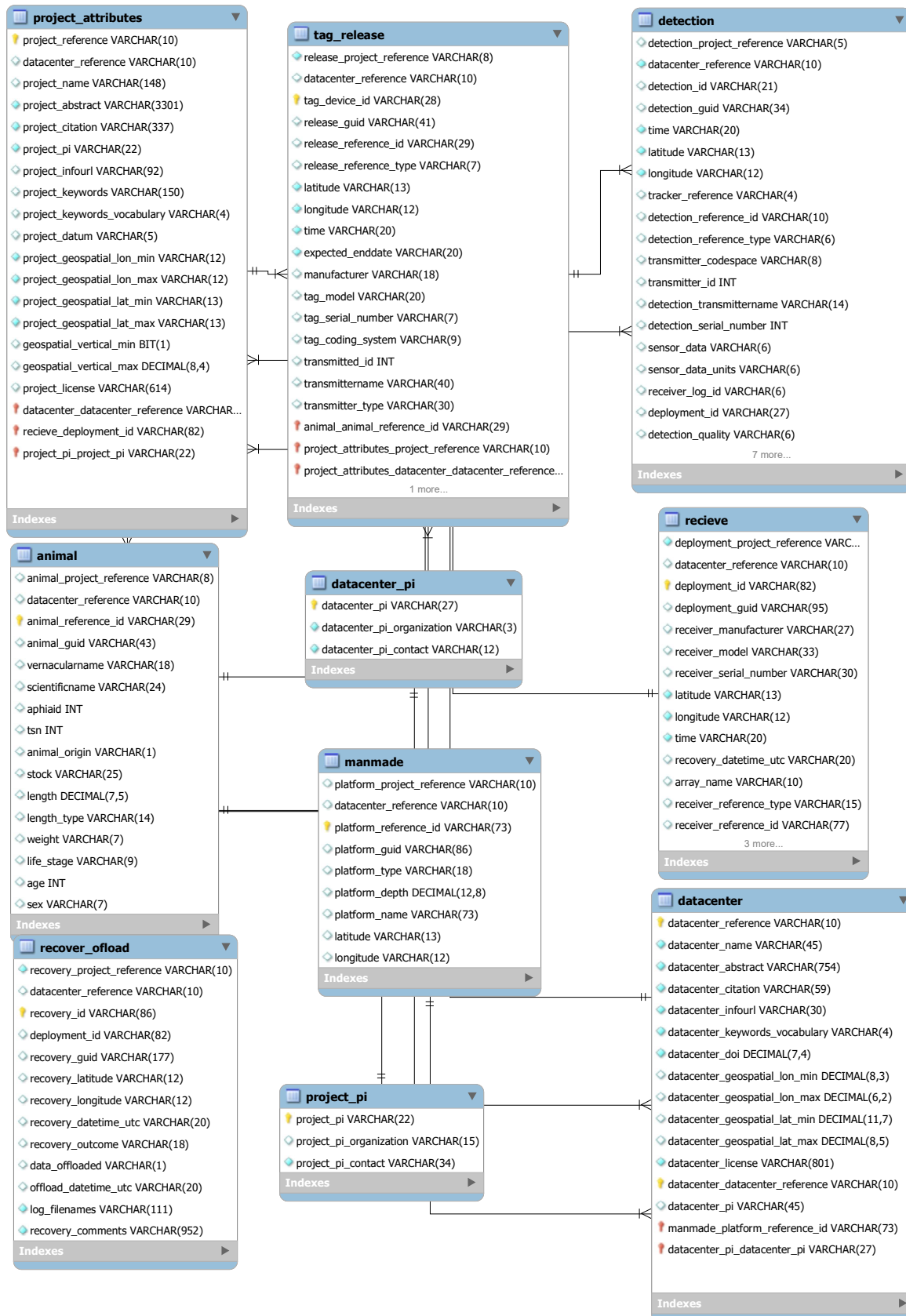
On the following page we could have a pre normalization reverse engineered ERD diagram. That consist of all the relationships and cardinalities.



Normalization:

- In datacenter table, it includes multivalued attributes keyword datacenter column, so this violated 1NF norms. To solve the issue, create a new table named keywords with two fields named references_datacenter and keywords_datacenter. After doing this step delete the multivalued attribute keyword datacenter from the table.
- Datacenter table has non – prime attributes (datacenter_pi, datacenter_pi_contact, datacenter_pi_organisation) which could uniquely identified from datacenter_pi by putting it as the primary key of the newer generated table.
 - o Table 1: Datacenter_pi_details(datacenter_pi (PK), datacenter_pi_contact, datacenter_pi_organisation)
 - o In the original table the Datacenter_pi will become primary key and datacenter_pi_contact, datacenter_pi_organisation would be removed from that table.
- Project_Attribute table is not also in 2NF, hence we will continue with the same method we proceed in the previous step.
 - o Here the issue lies with project_pi, project_pi_contact, project_pi_organisation where the tables could be designed in the form like where one table is main holding all the information with project_pi as foreign key and then the second table with project_pi as the primary key and project_pi_contact, project_pi_organisation as an attribute.

On the following page you could find the post normalization ERD diagram



Opportunities in Halifax

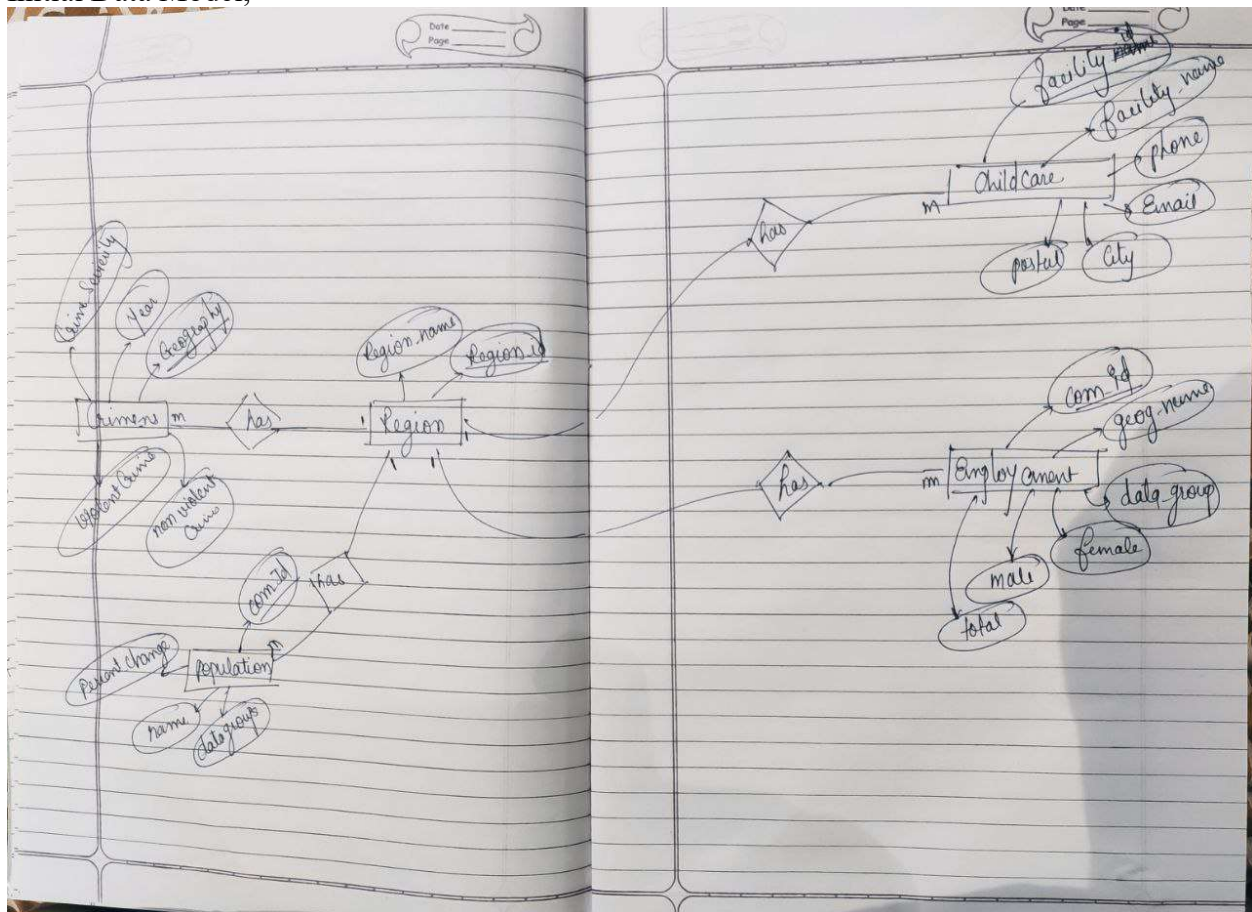
Here in this part of the assignment we need to gather 5 entities for the real-estate client who is looking to buy some property in nova scotia region.

For that work these are the five entities which had been finalized

- Region,
- Crime_ns,
- Population,
- Childcare and
- Employment

To gather the data sets <https://data.novascotia.ca/> is used as per stated in the assignment document.

Initial Data Model,



After finding the data sets I have used several Excel techniques to narrow down the data by cleaning the given data. Broadly wherever the data field in blank, it is filled by NULL values for numerical and No Data for the Varchar data types.

Here in this image you can see that there is <NULL> for each blank position in email_1

FACILITY	FACILITY	FACILITY	COUNTY	PHONE 1	EMAIL 1	EMAIL 2	TOTAL	LI	ANNUAL	ANNUAL	PROBATIC	AGE	RAN	PROG	FUI	PROG	PAI	PROG	SCI	FAMILY	H	AGE
3000898 Tadpoles	Day Care	f	Halifax Co	902-444-0421	<NULL>		16	8-Mar-19	#####		3 months	Yes	No	Yes	No	Yes	No	Yes	No	Yes		
3012474 Starting Bl	Day Care	f	Pictou Co	902-485-1957	<NULL>		33	4-Mar-19	#####		18 months	Yes	No	Yes	No	Yes	No	Yes	No	No		
2907648 Love, Laug	Day Care	f	Colcheste	902-662-3210	<NULL>		40	10-Oct-18	#####		18 months	Yes	No	Yes	No	Yes	No	Yes	No	No		
3648758 JR After S	Day Care	f	Cape Bret	902-564-8169, 902-371-8164	<NULL>		22	#####	#####		4 years - 1	No	No	Yes	No	Yes	No	No	No	No		
4089127 South End	Day Care	f	Halifax Co	902-420-1618	<NULL>		81	12-Jul-19	#####		3 months	Yes	No	No	No	Yes	No	No	Yes			
2992107 Gore Distr	Day Care	f	Hants County		<NULL>		32	8-Nov-18	9-Apr-19		18 months	Yes	No	Yes	No	Yes	No	No	No	Yes		
3626951 Montesso	Day Care	f	Halifax Co	902-865-5151	<NULL>		35	10-Jul-19	#####		8 months	Yes	No	Yes	No	Yes	No	Yes	No	Yes		
3205121 YMCA of F	Day Care	f	Pictou Co	902-752-0202 ext: 227	<NULL>		30	2-Nov-18	12-Jun-19		Pre-prima	No	No	Yes	No	Yes	No	No	No	No		
3714361 YMCA of F	Day Care	f	Pictou Co	902-752-0202 ext: 227	<NULL>		57	#####	10-Jul-18		3 months	Yes	No	No	No	No	No	No	Yes			
2901634 Garderie	Day Care	f	Halifax Co	902-865-6263	<NULL>		62	6-Nov-18	#####		3 years - 1	Yes	No	Yes	No	Yes	No	No	No	No		
3885120 Giant Step	Day Care	f	Halifax Co	902-829-3236	<NULL>		30	#####			4 years - 1	No	No	No	Yes	No	No	No	No	No		
3725593 Barrington	Day Care	f	Halifax Co	902-405-2722	<NULL>		68	#####			6 months	Yes	No	No	No	No	No	Yes				
3012472 YMCA of F	Day Care	f	Pictou Co	902-752-0202 ext: 239	<NULL>		30	#####	6-May-19		Pre-prima	No	No	No	Yes	No	No	No	No	Yes		
3028061 Garderie	Day Care	f	Inverness	902-224-1998	<NULL>		26	2-Apr-19	29-Oct-18		18 months	Yes	No	Yes	No	Yes	No	No	No	No		
2968428 Wee Bairr	Day Care	f	Halifax Co	902-477-2117	<NULL>		15	3-Oct-18	4-May-18		30 months	No	Yes	Yes	Yes	No	No	No	No	No		
3004738 Crestview	Day Care	f	Halifax Co	902-835-4295	<NULL>		14	5-Apr-19	#####		3 years - 1	No	Yes	Yes	Yes	No	No	No	No	No		
3315826 Little Lam	Day Care	f	Cumberla	902-660-3019	<NULL>		24	3-Jan-19	18-Jul-18		3 years - 1	Yes	No	Yes	No	Yes	No	No	No	No		
3002664 Red Apple	Day Care	f	Antigonis	902-863-4357, 902-863-0430	<NULL>		30	8-Apr-19	17-Oct-18		18 months	Yes	No	Yes	No	Yes	No	No	No	No		
2983227 Project EL	Day Care	f	Halifax County		<NULL>		15	8-Nov-18	#####		3 years - 5	No	Yes	No	No	No	No	No	No	No		
3039249 Apple Tre	Day Care	f	Kings Cou	902-582-3086	<NULL>		90	10-Jun-19	29-Jan-19		3 months	Yes	No	Yes	No	Yes	No	No	Yes			
3651971 Growing	Day Care	f	Halifax Co	902-434-3886	<NULL>		40	#####	#####		3 years - 1	No	Yes	Yes	Yes	No	No	No	No	No		
3644314 Cumberla	Day Care	f	Cumberla	902-667-4724	<NULL>		45	11-Oct-18	25-Jun-19		5 years - 1	No	No	No	Yes	Yes	No	No	No	No		

Whereas here Email_2 had been deleted as it does not had any single value in it.

FileHomeInsertPage LayoutFormulasDataReviewViewHelp

ClipboardFontAlignmentNumberConditional FormattingStylesCell StylesCellsEditingAnalysisSensitivity

Child_Care_Directory

Search

Dhruv Doshi

ShareComments

H1EMAIL 2

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R						
1	FACILITY	FACILITY	FACILITY	COUNTY	PHONE 1	PHONE 2	EMAIL 1	EMAIL 2	TOTAL	LI	ANNUAL	ANNUAL	PROBATIC	AGE	RAN	PROG	FUI	PROG	PAI	PROG	SCI	FAMILY	H	AGE
2	3052426 Swings Da	Day Care	f	Hants County			swingsdaycare@hotmail.com		29	26-Jun-19	11-Jan-19		18 months	Yes	No	Yes	No	Yes	No	No	No			
3	2969979 East Prest	Day Care	f	Halifax Co	902-462-0054		eastprestondacare@bellaliant.com		85	9-Oct-18	9-Jul-19		3 months	Yes	No	Yes	No	Yes	No	Yes	No	Yes		
4	2966347 Windsor	Day Care	f	Hants Cou	902-798-2001		info@windsordaycare.ca		80	#####	21-Jun-18		6 months	Yes	No	Yes	No	Yes	No	Yes	No	Yes		
5	4089503 St. Joseph	Day Care	f	Halifax Co	902-422-8441 ext: 19		damascus@stjcc.ca		104	16-Jul-19			3 months	Yes	No	No	No	No	Yes					
6	3000898 Tadpoles	Day Care	f	Halifax Co	902-444-0421				16	8-Mar-19	#####		3 months	Yes	No	Yes	No	Yes	No	No	Yes			
7	2976771 Parrsboro	Day Care	f	Cumberla	902-254-4201		woodamie@hotmail.com		19	6-Nov-18	#####		3 years - 5	No	Yes	No	No	No	No	No	No			
8	4012371 Carousel	Day Care	f	Colchester County			ctucker@ns.sympatico.ca		20	#####	1-Nov-11		3 years - 1	No	Yes	Yes	No	No	No	No	No			
9	3958837 Kids and	Day Care	f	Halifax Co	902-832-7852		larryuteck@kidsandcompany.com		76	#####	#####		3 months	Yes	No	No	No	Yes	No	No	Yes			
10	3224788 FHCC Afte	Day Care	f	Pictou Co	902-928-2211		t.e@eastlink.ca		20	#####	#####		30 months	No	Yes	Yes	Yes	No	No	No	No			
11	3117389 Connect	Day Care	f	Halifax Co	902-802-6292, 902-81		kidopolis@ic@gmail.com		24	#####	22-Jun-18		30 months	No	Yes	Yes	Yes	No	No	No	No			
12	2432411 Ponderos	Day Care	f	Antigonis	902-863-3994		effievincent@hotmail.com		30	19-Jul-18	#####		Pre-prima	No	No	Yes	No	No	No	No	No			
13	2974102 Lower Sol	Day Care	f	Antigonis	902-863-0514		lsrps@eastlink.ca		45	2-Oct-18	#####		6 months	Yes	No	Yes	No	Yes	No	Yes	No	Yes		
14	4049469 A Tiny	Lab Day	Day Care	f	Halifax Co	902-429-2539, 902-22	jillian@atinylab.ca		42	10-Jul-19			18 months	Yes	No	No	No	No	No	No	No			
15	3012474 Starting	Bl Day	Day Care	f	Pictou Co	902-485-1957			33	4-Mar-19	#####		18 months	Yes	No	Yes	No	Yes	No	No	No			
16	2984049 Christoph	Day Care	f	Halifax Co	902-468-5208		crlc@eastlink.ca		72	4-Dec-18	#####		3 months	Yes	No	No	No	No	No	No	Yes			
17	3053461 Queens	Day Care	f	Queens Co	902-354-5088		queensdaycare@ns.alliantzinc.ca		42	5-Jul-19	29-Jan-19		18 months	Yes	No	Yes	No	Yes	No	No	No			
18	3407954 Whyccor	Day Care	f	Inverness	902-756-2244		wcdc105@gmail.com		45	6-Sep-18	#####		18 months	Yes	No	Yes	No	Yes	No	No	No			
19	3652831 Little Pum	Day Care	f	Kings Cou	902-365-5137		littlepumpkinskentville@gmail.com		50	2-Apr-19	23-Oct-18		18 months	Yes	No	Yes	No	Yes	No	No	No			
20	2894476 Bright Beg	Day Care	f	Cumberla	902-667-7857		bbccc@ns.sympatico.ca		80	11-Jul-19	#####		6 months	Yes	No	Yes	No	Yes	No	No	No			
21	2992096 Lower On	Day Care	f	Colcheste	902-662-2495		regmichelin@eastlink.ca		20	10-Jan-19	5-Jun-19		30 months	No	Yes	No	Yes	No	No	No	No			
22	3216595 Winding	Day Care	f	Colcheste	902-632-2181, 902-63		pjocn@hotmail.ca		20	#####	#####		3 years - 1	No	Yes	Yes	Yes	No	No	No	No			
23	3380087 Willowbr	Day Care	f	Halifax Co	902-830-3658		lakramar@willowbraechildcare.com		132	10-Jun-19	#####		4 months	Yes	No	No	No	Yes	No	No	Yes			

Child_Care_Directory

Ready

The following page consists the ERD obtained by MySql workbench.



References

- [1]"Flowchart Maker & Online Diagram Software", *App.diagrams.net*, 2021. [Online]. Available: <https://app.diagrams.net/>. [Accessed: 01- Jun- 2021]
- [2]*Data.novascotia.ca*, 2021. [Online]. Available: <https://data.novascotia.ca/>. [Accessed: 01- Jun- 2021]
- [3]"How to import a CSV file into a MySQL database?", *Medium*, 2021. [Online]. Available: <https://medium.com/@AviGoom/how-to-import-a-csv-file-into-a-mysql-database-ef8860878a68>. [Accessed: 01- Jun- 2021]