arnoldHicran TermProject

April 26, 2023

MET CS 689

Hicran Arnold

Term Project: Parts Unlimited - EV Data Warehouse

1 Part 1: Are you working on your own or with a partner? If with a partner provide their name. If on your own, just state that this is the case.

I am working alone in this project

2 Part 2: Determine the project scope

- In a short paragraph, describe the topic you wish to explore –an update if any
- Update the five business questions that your data warehouse will answer.

Project Description The topic I wish to explore for my data warehouse project is the management of Parts Unlimited's EV parts business. This topic was inspired by "The Unicorn Project," which describes the challenges and opportunities of digital transformation in a large organization. Specifically, I plan to focus on storing and analyzing data related to charging stations, EV products, customer purchases, and geographic location. Parts Unlimited already sells EV parts, and my goal is to improve the organization's data management, reporting, and analysis capabilities related to this business.

Business Questions:

Business Question 1

What is the distribution of EV charging stations and EV cars across different geographic regions, and how does this relate to the popularity of EVs in those regions?

The Marketing and Sales department of your business can benefit from this query. By understanding the distribution of EV charging stations and EV cars across different geographic regions and how it relates to the popularity of EVs in those region. Additionally, the Sales department can use this information to identify regions with higher demand for EVs and focus their sales efforts in those areas to increase sales and revenue.

Business Question 2

How has the popularity of different EV models and electric vehicle types varied by geographic region over time, and how can Parts Unlimited use this information to target their marketing and sales efforts more effectively

The first question focused on the distribution of EV charging stations and EV cars across different geographic regions and its relation to the popularity of EVs. It was more general and aimed at understanding the overall trends and patterns in the adoption of EVs.

In contrast, the second question focuses specifically on the popularity of different EV models and electric vehicle types in different regions over time. It is more detailed and aimed at understanding how consumer preferences for specific EV models and types have evolved in different regions.

While both queries can help inform marketing and sales efforts, the second query can provide more granular insights into the specific EV models and types that are most popular in different regions, which can help Parts Unlimited make more targeted marketing and sales decisions.

Business Question Q3

How does the weight of EV parts relate to their price, and how does this relationship differ across different manufacturers?

The Operations and Sales departments of Parts Unlimited can benefit from this query. By understanding how the weight of EV parts relates to their price, and how this relationship differs across different manufacturers, Parts Unlimited can make more informed decisions regarding inventory management, pricing, and supplier selection.

Business Question Q4

What is the average price, weight, and quantity of products sold by Parts Unlimited over time, and how does this vary across different geographic regions in which the products have been manufactured?

The Operations and Sales departments of Parts Unlimited can benefit from this query. By understanding the average price, weight, and quantity of products sold by Parts Unlimited over time and how this varies across different geographic regions in which the products have been manufactured, Parts Unlimited can make more informed decisions regarding inventory management, pricing, and supplier selection.

Business Question Q5

How does the demand for charging stations and the number of vendors operating in different regions relate to the size of the EV market, and what trends can be identified over time?

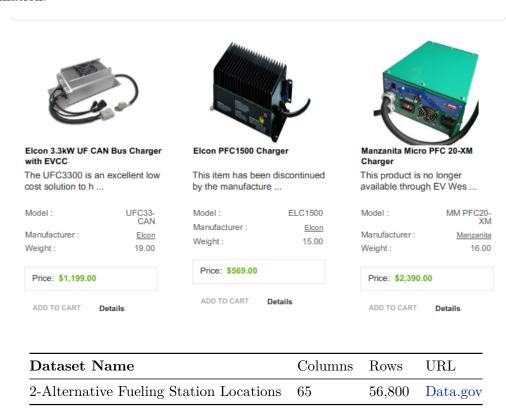
The Marketing and Sales departments of Parts Unlimited can benefit from this query. By understanding how the demand for charging stations and the number of vendors operating in different regions relates to the size of the EV market and what trends can be identified over time, Parts Unlimited can make more informed decisions regarding their marketing and sales efforts. Overall, understanding the relationship between the demand for charging stations, the number of vendors, and the size of the EV market can help Parts Unlimited make more strategic decisions regarding their marketing, sales, and product development effort

3 Part 3: Data Sources

Provide two data sources you will be using, for each data source list the number or columns and rows that are in each data source. Provide a header and first 5 rows from each source. - What is the URL or location of the data? - What information does this data provide that will help answer one or more of the above questions? - Do you see any issues in the data that will require transformation.

Dataset Name	Columns	Rows	URL
1-Product Info - EV West	7	187	EV West

The data set provides information about EV parts, manufacturer, weight and price information. The data is not available in the PDF version and requires a Python script to extract and format the information.



Provides information about the current list of charging stations and their locations. There are a lot of missing values that need to be cleaned and lat and long location info needs to be identical with the EV Car Population dataset. Price information is available in a text format and needs to be formatted so it is a number.

Inde	EV station Dataset Column exNames	Inde	EV Car Dataset	Inde	EV station Dataset Column exNames	Jnde	EV station Dataset Column exNames
1	Fuel Type Code	21	EV Other Info	41	Access Days	61	CNG PSI
2	Station Name	22	EV Network	42	Time (French) BD Blends	62	CNG Vehicle
3	Street Address	23	EV Network Web	43	(French) Groups With Access Code (French)	63	Class LNG Vehicle Class
4	Intersection Directions	24	Geocode Status	44	Hydrogen Is Retail	64	EV On-Site Renewable Source
5	City	25	Latitude	45	Access Code	65	Restricted Access
6	State	26	Longitude	46	Access Detail Code		
7	ZIP	27	Date Last Confirmed	47	Federal Agency Code		
8	Plus4	28	ID	48	Facility Type		
9	Station Phone	29	Updated At	49	CNG Dispenser Num		
10	Status Code	30	Owner Type Code	50	CNG On-Site Renewable Source		
11	Expected Date	31	Federal Agency ID	51	CNG Total Compression Capacity		
12	Groups With Access Code	32	Federal Agency Name	52	CNG Storage Capacity		
13	Access Days Time	33	Open Date	53	LNG On-Site Renewable Source		
14	Cards Accepted	34	Hydrogen Status Link	54	E85 Other Ethanol Blends		
15	BD Blends	35	NG Vehicle Class	55	EV Pricing		
16	NG Fill Type Code	36	LPG Primary	56	EV Pricing (French)		
17	NG PSI	37	E85 Blender Pump	57	LPG Nozzle Types		
18	EV Level1 EVSE Num	38	EV Connector Types	58	Hydrogen Pressures		
19	EV Level2 EVSE Num	39	Country	59	Hydrogen Standards		

	EV station				EV station	EV station
	Dataset Column		EV Car Dataset		Dataset Column	Dataset Column
$\operatorname{Ind}\epsilon$	exNames	$\operatorname{Ind}\epsilon$	exColumn Names	$\operatorname{Ind}\epsilon$	exNames	IndexNames
20	EV DC Fast	40	Intersection	60	CNG Fill Type	
	Count		Directions		Code	
			(French)			

```
[1]: import warnings
     warnings.simplefilter(action='ignore', category=Warning)
[2]: import pandas as pd
     pd.set_option("display.max_columns",10 )
[3]: import pandas as pd
     ev_station = pd.read_csv('data/alt_fuel_stations.csv')
     ev_station.head(5)
                                                             Station Name \
[3]:
       Fuel Type Code
                                     Spire - Montgomery Operations Center
                  CNG
                  CNG
                                                      PS Energy - Atlanta
     1
     2
                  CNG
                            Metropolitan Atlanta Rapid Transit Authority
                  CNG
     3
                                                    United Parcel Service
     4
                  CNG Clean Energy - Texas Department of Transportation
              Street Address
                                                         Intersection Directions \
     0
            2951 Chestnut St
     1
            340 Whitehall St From I-7585 N, exit 91 to Central Ave, left on...
         2424 Piedmont Rd NE
     2
                                                                              NaN
     3 270 Marvin Miller Dr
                                                                              NaN
         7721A Washington St I-10, Washington Ave exit, 1.5 blocks to the s...
                         CNG PSI CNG Vehicle Class
                                                     LNG Vehicle Class
              City ...
        Montgomery ...
                            3600
                                                                    NaN
                                                 MD
     1
           Atlanta ...
                            3600
                                                 MD
                                                                    NaN
     2
           Atlanta ...
                            3000
                                                 LD
                                                                    NaN
           Atlanta ...
     3
                            3600
                                                 HD
                                                                    NaN
           Houston ... 3000 3600
                                                 MD
                                                                    NaN
       EV On-Site Renewable Source Restricted Access
```

0	NaN	NaN
1	NaN	NaN
2	NaN	NaN
3	NaN	NaN
4	NaN	NaN

[5 rows x 65 columns]

Dataset Name	Columns	Rows	URL
3- EV Population Data	17	121,978	Data.gov

This dataset shows EV cars that are currently registered through Washington State Department of Licensing (DOL). Provides information about EV cars, car types, and registered locations. Location fields have to match with the format of the EV Charging dataset. There are some missing fields. Currently, I do not see any other issues related to this dataset.

Index	$Electric_Vehicle_Population$
1	VIN (1-10)
2	County
3	City
4	State
5	Postal Code
6	Model Year
7	Make
8	Model
9	Electric Vehicle Type
10	Clean Alternative Fuel Vehicle (CAFV) Eligibility
11	Electric Range
12	Base MSRP
13	Legislative District
14	DOL Vehicle ID
15	Vehicle Location
16	Electric Utility
17	2020 Census Tract

```
[4]: ev_pop = pd.read_csv('data/Electric_Vehicle_Population_Data.csv')
     ev_pop.head(5)
[4]:
        VIN (1-10)
                       County
                                    City State
                                                Postal Code
        5YJ3E1EB2J
                      Suffolk
                                Suffolk
                                            V٨
                                                     23435.0
     1
        5YJ3E1ECXL
                       Yakima
                                 Yakima
                                            WA
                                                     98908.0
     2
       WA1LAAGE7M
                       Yakima
                                  Yakima
                                                     98908.0
                                            WA
        5YJ3E1EA1K
     3
                     Danville
                               Danville
                                            VA
                                                     24541.0
        1FADP5CU9E
                      Norfolk
                                Norfolk
                                            VA
                                                     23518.0
        Legislative District DOL Vehicle ID
                                                           Vehicle Location
     0
                          NaN
                                    476647986
                                                  POINT (-76.42443 36.8752)
     1
                         14.0
                                    103490145
                                               POINT (-120.56916 46.58514)
     2
                         14.0
                                    144941534
                                               POINT (-120.56916 46.58514)
     3
                                                  POINT (-79.4172 36.58598)
                          \mathtt{NaN}
                                    168513922
     4
                          NaN
                                    150749378
                                                 POINT (-76.21549 36.92478)
```

Electric Utility 2020 Census Tract

0	NaN	5.180008e+10
1	PACIFICORP	5.307700e+10
2	PACIFICORP	5.307700e+10
3	NaN	5.159000e+10
4	NaN	5.171001e+10

[5 rows x 17 columns]

Dataset Name	Columns	Rows	URL
4- Manufacture Location	6	25	extracted by online search

Provides information about EV cars. This dataset is clean, currently I do not see any issues.

Index	Manufacture Location
1	company_name
2	latitude
3	longitude
4	city
5	state
6	country

4 Part 4: Dimensions - Review the data and the business questions from part 2.

- What fields (attributes) are in the data that will be used for the dimensions.
- Determine the dimension tables. There should be at least two non-date dimensions and one date dimension for each fact table.
- At least one (non-date) dimension in your design should have a hierarchy.
- What are the attributes that will be tracked via slowly changing dimensions?
- What attributes within the dimensions will need transformation before they are loaded into the dimension, for example it could be to build consistency or any other issues? This is where for example you might build case statements in your code to handle various scenarios. Two to three examples showing some sample data and what you think the transformation will be during your ETL would be helpful here.

In the below tables I am showing each data sets by their corresponding tables and their scd types and their data transformation information summary

EV poulation Data set			
Field Name Description	Transformat	io T able	SCD Type
	Info	Name	Info
$ev_car_pop_idthis$ is a pk	not in the	ev_car_popul	at RK
fieldssign to	original		

nfo ЬK the unique data/ has to each records be populated dol_vehicle_Id unique id for no transforev_car_population type3 each vehicle mation needed VIN no transfor $ev_car_populat$ is od type3 unique id for each make mation model needed make ev brand no transforev_car_population type3 mation needed vechicle model brand model no transforev car populat**sod** type3 mation needed current_reg_stashusws if the not in the ev_car_population type3(scd registration original field) data/ has to in current (scd type 2) be populated reg_status_prevprevious not in the ev_car_population type3(scd original registiration field) status data/ has to be populated reg_sta_updated hant the not in the ev_car_population type3(scd registiration original field) has been data/ has to updated be populated location id not in the unique location dim scd type3 location id original for each car data/ has to be populated city city name no transforlocation dim scd type3 that car mation registered needed region the region not in the location_dim scd type3 that city falls original under to data/ has to be populated country the country not in the location dim scd type3 that car original data/ has to registred be populated

EV poulation Data set				
latitude	car registered address lat info	Point field, I have split into lat and long	location_dim	scd type3
longitude	car registered address lat info	Point field, I have split into lat and long	location_dim	scd type3
city_lat	city lat info	this was not in the original I pulled via api call but I realised that Tablaue already has city lat and long info can be automaticallypopulated by city name	location_dim	scd type3
city_long	city long info	this was not in the original I pulled via api call but I realised that Tablaue already has city lat and long info can be automaticallypopulated by city name	location_dim	scd type3
latitude_prev	scd field keeping change in the lat field	not in the original data/ has to be populated	location_dim	scd type3(scd field)
longitude_pre		not in the original data/ has to be populated	location_dim	scd type3(scd field)

EV poulation Data set				
updated_at	scd field tracking adjustment date	not in the original data/ has to be populated	location_dim	scd type3(scd field)
date_id	unique id for registration date	not in the original data/ has to be populatedI have used make year to create this data	date_dim	scd type 1
day	registration date	not in the original data/ has to be populatedI have used make year to create this data	date_dim	scd type 1
month	registration month	not in the original data/ has to be populatedI have used make year to create this data	date_dim	scd type 1
year	registration year	not in the original data/ has to be populatedI have used make year to create this data	date_dim	scd type 1

EV poulation Data set				
date	registration year	not in the original data/ has to be populatedI have used make year to create this data as01/01/make year	date_dim	scd type 1

below is the charging station dataset with the final data fields information

		Transformation	of Table	SCD Type
Field Name	Description	Info	Name	Info
charging_statio	n <u>in</u> dine_idlfor each charging stations	not in the original data/ has to be populated	ev_charging_	_sta tk ions
charging_statio	n <u>m</u> ique id for each charging stations- source	no transfor- mation needed	ev_charging_	_st ationy pe 2
sk_ev_dim_id	surrogate key to keep track of scd 2 change	not in the original data/ has to be populated		SK type2(scd field)
station_name	charging station name	no transfor- mation needed	ev_charging_	_st ations pe 2
status_code	temporarily closed, open	transformed from the code (P,T) to show more readable format	ev_charging_	_st atidny pe 2
access_code	Private or Public	no transfor- mation needed	ev_charging_	_st stidnty pe 2
currency	charging station price currency	not in the original data/ has to be populated	ev_charging_	_st atidns pe 2

Field Name	Description	Transformati Info	o T able Name	SCD Type Info
row_effective_	dstd type 2 shows when the first record created	not in the original data/ has to be populated	ev_charging_	_st stidnty pe2(scd field)
row_expiration	shows if this record has been adjusted andnew recorded added	not in the original data/ has to be populated	ev_charging_	_st stidnty pe2(scd field)
row_status	scd type 2 shows if this record still effective	not in the original data/ has to be populated	ev_charging_	_st stidnty pe2(scd field)
location_id	unique location id for each charging station recorded	not in the original data/ has to be populated	location_dim	PK
city	city name of the charging station	no transfor- mation needed	location_dim	scd type3
region	the region that city falls under to	not in the original data/ has to be populated	location_dim	scd type3
country	country name	no transfor- mation needed	location_dim	scd type3
latitude	charging station address lat info	no transfor- mation needed	location_dim	scd type3
$\log tude$	charging station address long info	no transfor- mation needed	location_dim	scd type3

Field Name	Description	Transformati Info	o¶able Name	SCD Type Info
city_lat	city lat info	this was not in the original I pulled via api call but I realised that Tablaue already has city lat and long info can be automaticallypopulated by city name	location_dim	scd type3
city_long	city long info	this was not in the original I pulled via api call but I realised that Tablaue already has city lat and long info can be automaticallypopulated by city name	location_dim	scd type3
latitude_prev	v scd field keeping change in the lat field	not in the original data/ has to be populated	location_dim	scd type3(scd field)
$ m longitude_pr$		not in the original data/ has to be populated	location_dim	scd type3(scd field)
updated_at	scd field tracking adjustment date	not in the original data/ has to be populated	location_dim	scd type3(scd field)
${ m date_id}$	unique id for registration date	not in the original data/ has to be populated	date_dim	PK

To:	-1-1 NT	D	Transformatio		SCD Type
F10	eld Name	Description	Info	Name	Info
day	у	charging station open date	transformed from the open date field in source data	date_dim	scd type 1
mo	onth	charging station open month	transformed from the open date field in source data	date_dim	scd type 1
yea	ar	charging station open year	transformed from the open date field in source data	date_dim	scd type 1
dat	te	charging station open year	charging station open date	date_dim	scd type 1

EvWest data set

	EvWest data-set				
]	Field Name	Description	Transformation	o T able	SCD Type
			Info	Name	Info
r	manufacturer	unique id for	not in the	manufacturer	PK
i	id	each	original		
		manufacturer	data/ has to		
			be populated		
S	sk	surrogate key	not in the	manufacturer	SK SCD
1	manufacturer	to keep track	original		type2(scd
		of scd 2	data/ has to		field)
		change	be populated		
r	manufacturer	n man ufacturer	pdf	manufacturer	scd type 2
		station name	extraction		
			transformation		
1	manu_row_effe	estidet(spd)2	not in the	manufacturer	scd type2(scd
		tracking:	original		field)
		record	data/ has to		
		creatation	be populated		
		date			

EvWest data-set				
manu_row_	expiratity(sc2) tracking: record expiration date track	not in the original data/ has to be populated	manufacturer	scd type2(scd field)
manu_row_	_ind(scol)type 2 tracking: indicated if the record is inuse or there is an adjustment (expired/active	not in the original data/ has to be populated	manufacturer	scd type2(scd field)
product_id	unique id for each product	not in the original data/ has to be populated	Products	PK
product_co	de unique code for each product type	pdf extraction transformation	Products	scd type 3
product_na	ame product type name	pdf extraction transformation	Products	scd type 3
$ m product_de$	escrip ting description of the product	pdf extraction transformation	Products	scd type 3
product_ca	tego n roduct category	pdf extraction transformation	Products	scd type 3
produc_nar	tracks of the name change in the productPrevious name of the product	not in the original data/ has to be populated	Products	scd type3(scd field)
product_up	odatesbdattype 3: tracks of the name change in the pro- ducttracks of the adjustment date	not in the original data/ has to be populated	Products	scd type3(scd field)

EvWest				
data-set				
location	_id unique location id for each charging station recorded	not in the original data/ has to be populated	location_dim	PK
city	city name that car registered	no transfor- mation needed	location_dim	scd type3
region	the region that city falls under to	not in the original data/ has to be populated	location_dim	scd type3
country	the country that car registred	no transfor- mation needed	location_dim	scd type3
latitude	charging station address lat info	no transfor- mation needed	location_dim	scd type3
longitud	le charging station address long info	no transfor- mation needed	location_dim	scd type3
city_lat		this was not in the original I pulled via api call but I realised that Tablaue already has city lat and long info can be automaticallypopulated by city name	location_dim	scd type3

	EvWest data-set				
	city_long	city long info	this was not in the original I pulled via api call but I realised that Tablaue already has city lat and long info can be automaticallypopulated by city name	location_dim	scd type3
]	latitude_prev	scd field keeping change in the lat field	not in the original data/ has to be populated	location_dim	scd type3(scd field)
1	longitude_prev	scd field keeping change in the long field	not in the original data/ has to be populated	location_dim	scd type3(scd field)
,	updated_at	scd field tracking adjustment date	not in the original data/ has to be populated	location_dim	scd type3(scd field)
(date_id	unique id for manufacturer vendor since date	Not in the original data/ has to be populated	date_dim	PK
(day	manufacturer vendor since date -day	not in the original data/ has to be populated	date_dim	scd type 1
1	month	manufacturer vendor since date -month	not in the original data/ has to be populated	date_dim	scd type 1
	year	manufacturer vendor since date -year	not in the original data/ has to be populated	date_dim	scd type 1

EvWe data-				
date	manufacturer vendor since date	not in the original data/ has to be populated	date_dim	scd type 1

5 Part 5: Facts – Review the data and the business questions from step 1.

What measurements are in the data that will be used for the fact tables? What measures will you be calculating (i.e. using an aggregate function, or some other transformation – recall as an example some of the aggregation you did in assignment 1A)

manufacture-facts: This table is keeping track of manufacturers and manufacturer products that the company sells. The measurements are weight, price

Field Name	Description	Transformation Info	SCD Type Info
manu_fact_id	unique id for	not in the	PK
	each	original data/	
	manufacturer	has to be	
	fact record	populated	
location_id	foreign key to	extracted by	FK
	location info	merge	
$date_dim$	foreign key to	extracted by	FK
	date info	merge	
$manufacture_id$	foreign key to	extracted by	FK
	location	merge	
	manufacturer		
	table connection		
$\operatorname{product_id}$	foreign key to	extracted by	FK
	location product	merge	
	table connection		
price	price of the	pdf extraction	measurement
	product for each		
	manufacturer		
weight	weight of the	pdf extraction	measurement
	product for each		
	manufacturer		
quantity	quantity of the	not in the	measurement
	product for each	original data/	
	manufacturer	has to be	
		populated	

below is my second fact ev_charginging_facts snapshot table tracks of snapshot of the charging stations quantity and the ev population $\frac{1}{2}$

Field Name	Description	Transformation Info	SCD Type Info
ev_char_fact_id	unique id for each charging fact records	not in the original datahas to be populated	PK
location_id	fk to location info	extracted by merge	FK
$\mathrm{date_dim}$	fk to date info	extracted by merge	FK
$charging_station_id$	fk to location manufacturer table connection	extracted by merge	FK
ev_car_pop_id	fk to location ev population table connection	extracted by merge	FK
vehichle_pop_quantity	quantity of cars	pdf extraction	measurement
${ m ev_cs_quantity}$	quantity of charging stations	pdf extraction	measurement
ev_price	price of the charging stations	not in the original data/has to be populated	measurement

cumulative fact table: this tracks of the overall quantities , this is a bridge of operations and research tables

Field Name	Description	Transformation Info	SCD Type Info
manu_fact_cumu_id	unique id for	not in the original data/	PK
	each charging fact records	has to be populated	
$\operatorname{month_id}$	month (time dimension in	extracted by merge	FK
	month grain)		
location_id	fk to location info	extracted by merge	FK
${\rm manufacturer_id}$	fk to manufacturer	extracted by merge	FK
	info		
total_product_quntity	total product quantity for	aggregated with etl	measurement
total_ev_charging_statio	that month ontsotal charging stations	aggregated with etl	measurement

Field Name	Description	Transformation Info	SCD Type Info
total_ev_population	total ev car population	aggregated with etl	measurement

6 Part 7: Provide code and screenshots of loading your data into staging/data frame

For Question 7-9: I used pandas df to load data and before transformation and again I load the data with the help of the pandas sql load function. Please see my etl jupiternotebook for more details. I am including a three etl file one is for charging station facts and the other one for manufacturer facts. I grouped etl for each documents for each fact table

7 Part 8: Provide code and screenshots of transforming the data. Perhaps you are adjusting for consistency of data or calculating aggregates.

For Question 7-9: I used pandas df to load data and before transformation and again I load the data with the help of the pandas sql load function. Please see my etl jupiternotebook for more details. I am including a three etl file one is for charging station facts and the other one for manufacturer facts. I grouped etl for each documents for each fact table

8 Part 9: Provide code and screenshots of loading the data into the two dimensions and the fact. At this time, you do not need to worry about maintenance of slowly changing dimensions, the focus is on the initial data load. If you are loading into SCD2 or SCD3, make sure to show the SCD maintenance attributes populated.

I have inculuded a spearate notebook for this transformation and breakdown the process for each scd type and showed examples

8.1 Part 11: Outline one of the business questions you outlined in part 2 (can be same as part 10) and answer it with a Tableau or PowerBI visualization. In a single sentence explain why you selected this particular visualization to answer this question.

Please see my query notebook for my business questions. I answered each business questions with sql aggreagted queries. Please also see my presentation documents for visualizations for those questions