

Analyst Workshop

In this workshop, we will be working with some dataset that contains insurance rates data from across the United States, providing insights into the premiums charged by insurers, the underlying factors that affect those rates, and claims history analysis. The data is designed to help researchers understand the inner workings of the insurance industry, and how rates are calculated. It includes information on premiums, underlying factors, current premium prices, indicated premium prices, selected premium prices, fixed expenses, and more

There are 3 tables with the schemas and column definitions below.

- Cgr_definitions_table - Dimension table with the premium rates for a Combined grade rating (CGR)
- Cgr_premiums_table - Historic dataset on customer's premium, expenses and CGR factor
- Territory_definitions_table - Location data such as county, country_code, zipcode related to the customer's territory

Cgr_definitions_table

Column name	Description
cgr	Combined grade rating. (Numeric)
aa	Average annual premium. (Numeric)
bb	Base premium. (Numeric)
cc	Cost of capital. (Numeric)
va	Value of assets. (Numeric)
dd	Direct written premium. (Numeric)
hh	Homeownership. (Categorical)

Cgr_premiums_table

Column name	Description
territory	The territory in which the person lives. (String)
gender	The person's gender. (String)
birthdate	The person's birthdate. (Date)
ypc	The person's years of prior coverage. (Integer)
current_premium	The person's current premium. (Float)
indicated_premium	The person's indicated premium. (Float)
selected_premium	The person's selected premium. (Float)
underlying_premium	The person's underlying premium. (Float)

fixed_expenses	The person's fixed expenses. (Float)
underlying_total_premium	The person's underlying total premium. (Float)
cgr_factor	The person's CGR factor. (Float)

Territory_definitions_table

Column name	Description
territory	The territory in which the person lives. (String)
county	The county in which the person lives. (String)
county_code	The county code for the county in which the person lives. (String)
zipcode	The zip code for the county in which the person lives. (String)
town	The town in which the person lives. (String)

Data Upload / Ingestion via UI

1. Create a schema with your name / unique ID
2. Upload each csv via the UI to create a new table within your schema
3. Once the tables are created, navigate to the table and there would be an AI generated description for the table. Accept or edit as necessary
4. This makes it easy for the search function to locate the right data object with the right metadata
5. Search for CGR with the full definition

Data upload via Volumes

1. Besides creating a direct delta table, you can also upload files into a volume as csv. This can be done either via the Databricks UI or a direct upload into ADLS.
2. First create a volume pointing towards an external location that you have set up on Databricks.
3. Create subdirectories within for the different files, as this will future proof it if there are additional new files that should be appended to the table

Create a new volume

×

Volume name

vehicle_data

Volume type [Learn more](#)

☐ Managed volume
☒ External volume

External location

beatrice_liew_ext

This external location will act as a source for your volume

Path

s3://one-env-uc-external-location/beatrice-liew

The prefix where the volume will be created

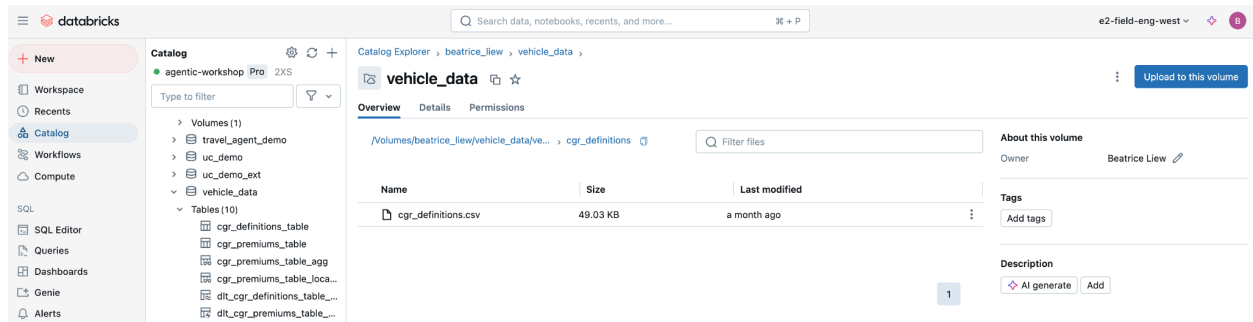
Comment (optional)

Cancel

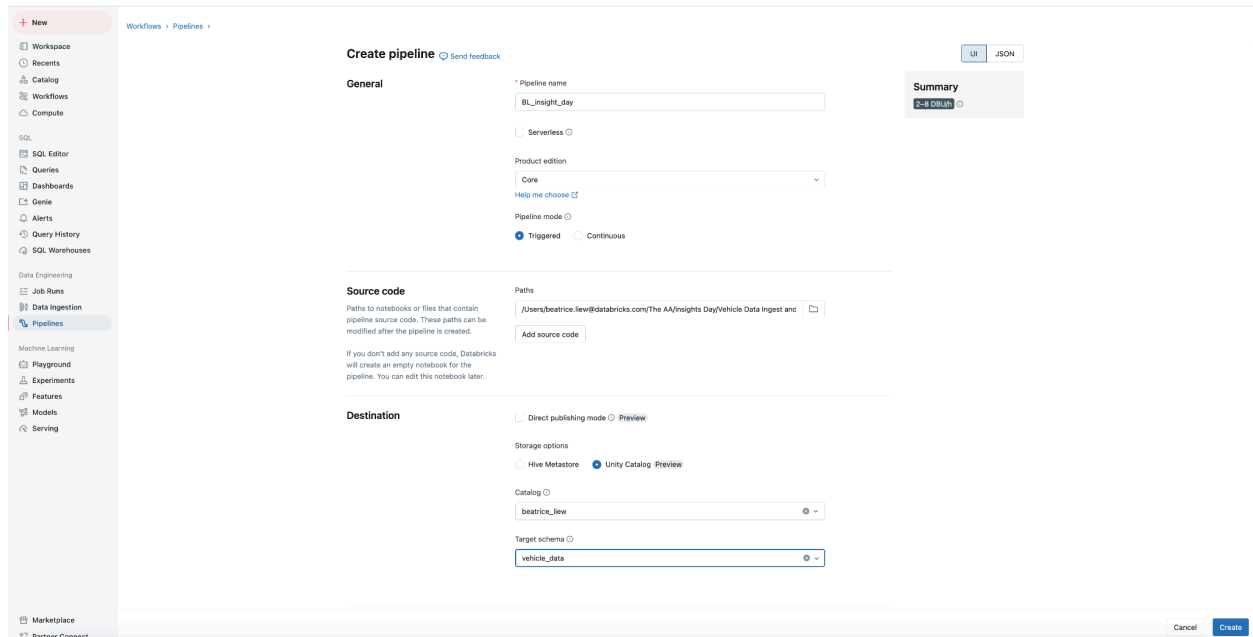
Create

Data Engineering

1. Three csv files are already uploaded in your workspace
2. Download any one (Python or SQL) of the [notebooks](#) from Github
3. Change the volume paths in the notebook to point towards the path in your workspace (ie: /Volumes/beatrice_liew/vehicle_data/vehicle_data/cgr_definitions/)



4. Create a new pipeline under Delta Live Tables and fill out the details below:
 - a. Pipeline name - use a unique name
 - b. Product edition - Core
 - c. Pipeline mode - Triggered
 - d. Paths - The path to the notebook under workspaces
 - e. Destination - Choose Unity Catalog and write to your schema
 - f. Click create



5. Once the pipeline is created, you would need to click on Run Now to process the data. Note that there are 2 options - 1) Full refresh all - all the files will be reprocessed 2) Start - this will process incremental new data
6. To schedule this pipeline, click on Jobs on the left pane, click on create Job
7. In create task, fill out:
 - a. Task name - give the task a name (eg: Vehicle data DLT pipeline)

- b. Type - Pipeline (Or DLT)
- c. Pipeline - Select the DLT pipeline you've created

Runs

Tasks

DLT_Pipeline

Pipeline: vehicle_data_sql

+ Add task

Task name*

DLT_Pipeline

Type*

Pipeline

Pipeline*

vehicle_data_sql

Trigger a full refresh on the pipeline

Notifications

+ Add

Retries

+ Add

Metric thresholds

+ Add

Cancel

Save task

8. On the right pane, under Schedules & Triggers, add a trigger:
- a. Choose file arrival - this means that the pipeline is triggered within 1 min if a file is added to a volume
 - b. Storage location - add the volume path for cgr_premiums directory (eg: /Volumes/beatrice_liew/vehicle_data/vehicle_data/cgr_premiums/)

Schedules & Triggers

Trigger Status

Active

Paced

Trigger type

File arrival

This job does not have any failure notifications. Consider adding email or webhook notifications to alert if trigger evaluation fails.

File arrival triggers monitor cloud storage paths of up to 10,000 files for new files. These paths are either Unity Catalog volumes or external locations managed through Unity Catalog.

Storage location

/Volumes/beatrice_liew/vehicle_data/vehicle_data/cgr_premiums/

Advanced

Test connection

Cancel

Save

Job notifications

No notifications

Edit notifications

Duration and streaming backlog thresholds

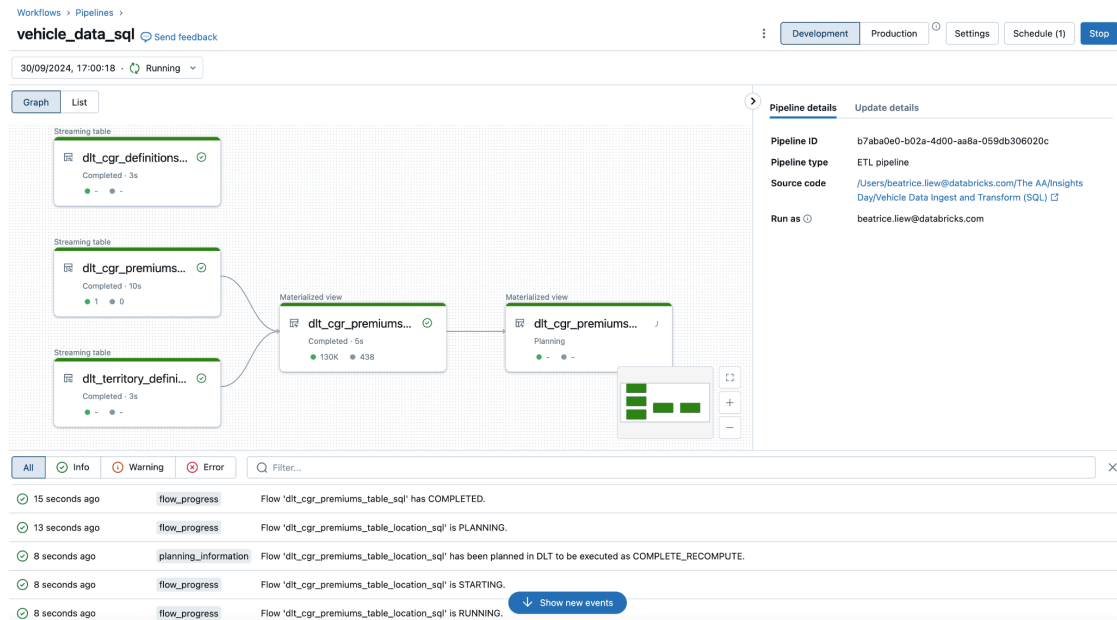
No thresholds defined

Add metric thresholds

- Click on run now - Click on the job and under runs, you will see that it's found no new files

Trigger status: Monday, September 30, 2024 4:54 PM - no run triggered
Found no new files and 1 file in total.

- Add an additional file named ([cgr_premiums_additional_data.csv](#)) to the path
- In around 1 minute, the pipeline would have picked up the new file and you can view the DLT pipeline to see that there is 1 row processed



Data modelling

- Switch over to the SQL editor
- Create the primary keys by running the following command. As the columns cannot be nulls for primary keys, we need to set the constraint to non nullable.

```
ALTER TABLE cgr_definitions_table ALTER COLUMN cgr SET NOT NULL;  
ALTER TABLE cgr_definitions_table ADD CONSTRAINT cgr_pk PRIMARY KEY(cgr);  
ALTER TABLE territory_definitions_table ALTER COLUMN territory SET NOT NULL;  
ALTER TABLE territory_definitions_table ADD CONSTRAINT territory_pk PRIMARY  
KEY(territory);
```

- Create the foreign keys by running the following command. This links the foreign keys to the primary keys in the dimension tables

```
ALTER TABLE Cgr_premiums_table ADD CONSTRAINT cgr_fk FOREIGN KEY(cgr) REFERENCES  
cgr_definitions_table;  
ALTER TABLE Cgr_premiums_table ADD CONSTRAINT territory_fk FOREIGN KEY(territory)  
REFERENCES territory_definitions_table;
```

- Flip back to the catalog view of the tables and there would be a PK / FK icon next to the columns

5. Click on view relationships to see the model / ERD diagram (If available)

Data analysis

1. Return to SQL editor
2. Let's add an age column to the Cgr_premiums_table using the query below (Or try to use Databricks assistant for this if not familiar with the syntax: Return a query to calculate the age in years based on birthdate)

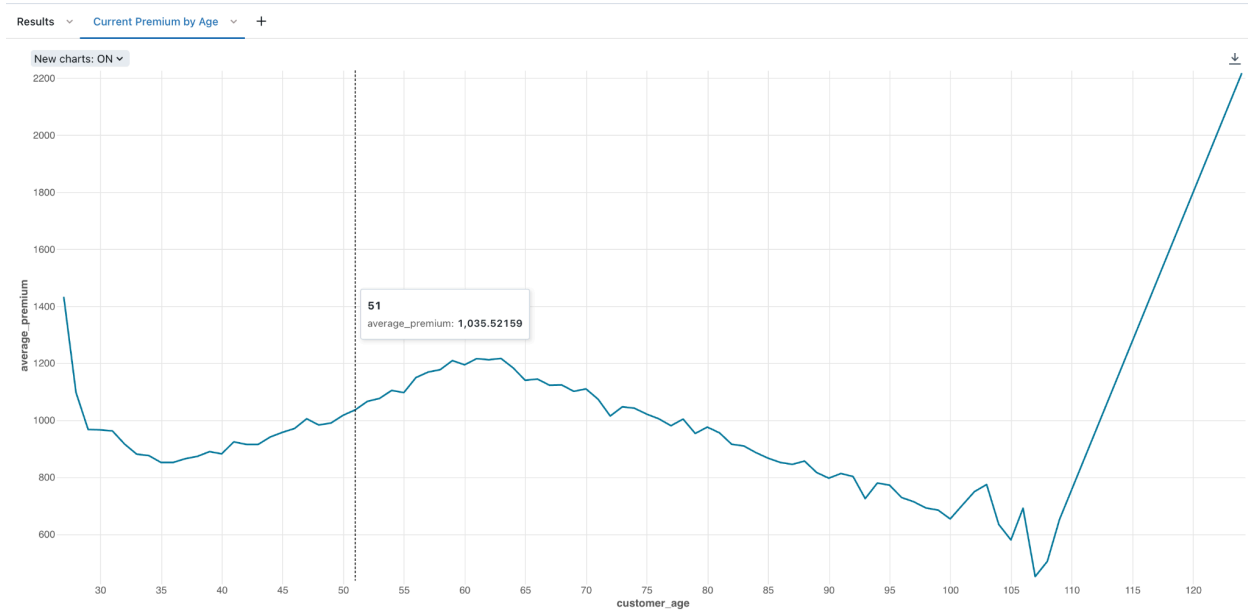
```
SELECT *, datediff(year, birthdate, CURRENT_DATE()) as customer_age
FROM cgr_premiums_table
```

3. Find the trend between average current_premium and the age of the customer (Or ask assistant: what is the relationship between average of current premium and the customer's age)

```
select customer_age, avg(current_premium) as average_premium from (
SELECT *, datediff(year, birthdate, CURRENT_DATE()) as customer_age
FROM cgr_premiums_table)
group by customer_age
order by customer_age asc
```

Results ▾ +			
#	customer_age	average_premium	
1	27	1429.78	
2	28	10 1429.78	
3	29	966.67	
4	30	965.13	
5	31	961.63	
6	32	915.50	
7	33	880.06	
8	34	875.09	
9	35	851.06	
10	36	851.05	
11	37	864.15	
12	38	872.46	
13	39	888.96	
14	40	881.24	
15	41	923.36	
16	42	914.42	

4. Create a visualisation to visualise this using a line chart with customer_age on the x axis and average_premium on the y axis.



5. Notice that data seems odd after >100 years. Let's use a where clause to filter it out.

```
select customer_age, avg(current_premium) as average_premium from (  
SELECT *, datediff(year, birthdate, CURRENT_DATE()) as customer_age  
FROM cgr_premiums_table  
)  
where customer_age < 100  
group by customer_age  
order by customer_age asc
```

6. Enrich the historic customer premium information with location data by joining to the territory_definitions_table. Create a view with this.

```
SELECT a.*, datediff(year, a.birthdate, CURRENT_DATE()) as customer_age, b.*  
FROM cgr_premiums_table a  
LEFT JOIN territory_definitions_table b  
on a.territory = b.territory
```

7. Find the average current_premium by town, which town has the highest premium? How many customers are there? (Assistant: What is the average current premium by town, using the table territory_definitions_table. How many customers are in each town, alongside their average premium?)

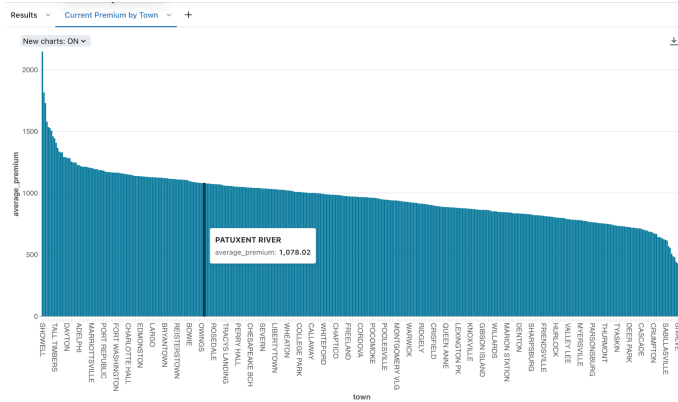
```

select town, avg(current_premium) as average_premium, count(*) as number_of_customers
from
(
SELECT a.*, datediff(year, a.birthdate, CURRENT_DATE()) as customer_age, b.*
FROM cgr_premiums_table a
LEFT JOIN territory_definitions_table b
on a.territory = b.territory
)
group by town
order by average_premium desc

```

Results ▾ Current Premium by Age ▾ +				
#	town	average_premium	number_of_customers	
1	SHOWELL	2145.03	1	
2	HARMANS	1813.14	1	
3	KITZMILLER	1728.24	4	
4	BARNESVILLE	1576.99	3	
5	DOWELL	1533.72	3	
6	STEVENSON	1523.91	7	
7	CHELTENHAM	1503.08	64	
8	BUCKEYSTOWN	1457.93	2	
9	TALL TIMBERS	1441.36	6	
10	UPPER FALLS	1406.38	8	
11	BEL ALTON	1363.64	16	
12	CLARKSVILLE	1334.80	133	
13	WELCOME	1328.34	23	
14	ISSUE	1327.59	7	

8. Create a chart to visualise the relationship between premium and town



9. To reuse the query, save it on SQL editor into your own workspace. You can also share the query with others

Scheduling queries

1. To schedule this query to run at a set time, save the query first.
2. Use the query below to create a view and save the query

```
CREATE OR REPLACE VIEW cgr_premiums_table_location as
SELECT a.*, datediff(year, a.birthdate, CURRENT_DATE()) as customer_age, b.county,
b.county_code, b.zipcode, b.town, b.area
FROM cgr_premiums_table a
LEFT JOIN territory_definitions_table b
on a.territory = b.territory
```

3. Create an aggregated view based on this view to get the average premium by town. Save the query

```
CREATE OR REPLACE VIEW cgr_premiums_table_agg as
select town, avg(current_premium) as average_premium, count(*) as number_of_customers
from cgr_premiums_table_location
group by town
order by average_premium desc
```

4. If the SQL queries are not dependent on each other / anything else, click on schedule and set the frequency required.
5. To build dependencies between scheduled queries, use workflows
6. Click create job
7. Click on create task and add the SQL query in for creating cgr_premiums_table_location
8. Add another task and add the SQL query for creating cgr_premiums_table_agg
9. Add in a schedule or trigger for the update
10. There's also an option to add in notification via email to notify in the event of failure / success

Premium View Job ☆

Runs

Tasks

Create_premium_location

Premium location view

Serverless Shared Endpoint

Create_premium_agg_view

Premium agg view

Serverless Shared Endpoint

+ Add task

Task name*

Create_premium_location

Type*

SQL

SQL task*

Query

SQL query*

Premium location view

SQL warehouse*

Serverless Shared Endpoint (M)

Parameters

The selected SQL query does not have parameters.

Depends on

Select task dependencies...

Notifications

+ Add

Retries

+ Add

Metric thresholds

+ Add

Job details

Job ID

563646870708398

Creator

Beatrice Liew

Run as

Beatrice Liew

Tags

Add tag

Description

Add description

Lineage

No lineage information for this job.

Schedules & Triggers

None

Add trigger

Compute

Serverless Shared Endpoint

Medium, 1-4 clusters, 1 active cluster

View SQL warehouse

Swap

Job parameters

No job parameters are defined for this job

Edit parameters

Job run settings

Queue

Maximum concurrent runs

1

Edit concurrent runs

Duration and streaming backlog thresholds

No thresholds defined

Alerts

1. There's an option to create an alert to listen in on the data
2. Let's say we would like to be alerted if there's any town's with an average premium > 2000
3. Set up the alert as shown below and click on create alert

New alert

Alert name

Average premium by town

Query

Average premium by town

Trigger condition

Value column

average_premium

First row

Operator

≥

Threshold value

2000

When query result has no rows, set state to

UNKNOWN

[Preview alert](#)

Notification: When alert is triggered

Send notification

Just once

When alert returns back to normal

☒ Send notification

Template

Use default template

Create alert

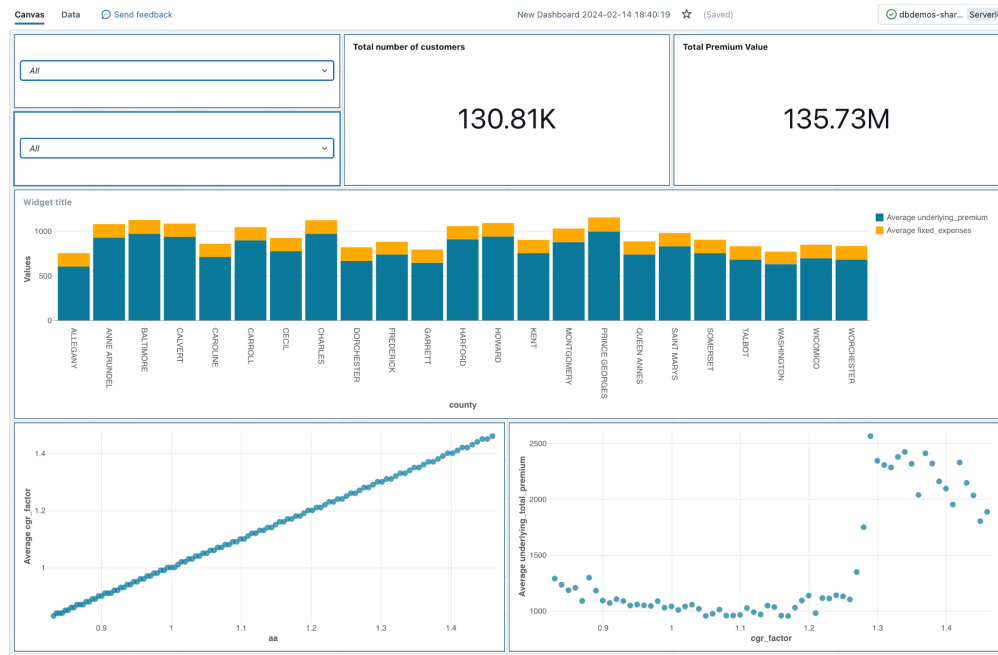
4. Set up a schedule for when you want the check to be done, which would either trigger an alert if the trigger condition is met, or not if not met
5. Add your email as the notification destination
6. Click on run once to test that it is working. You should receive an email when it's triggered

Lakeview dashboards

1. Navigate to Dashboards and click on lakeview dashboard
2. First, add data that we want to visualise in the data tab. Click on select a table and choose the view with data enriched with location data cgr_premiums_table_location
3. There's also an option to define transformations using create from SQL. Enrich the cgr_premiums_table using the table cgr_definitions_table with the query below

```
select *  
from beatrice_liew.vehicle_data.cgr_premiums_table a  
left join beatrice_liew.vehicle_data.cgr_definitions_table b  
on a.cgr = b.cgr
```

4. On the canvas, use the cgr_premiums_table_location dataset
5. To create a card with the total number of customers in the dataset, add a viz from the bottom blue bar. Under visualisation, select counter. Add count(*) as the value. Add in widget title to name it total number of customers
6. There's also an option to use text to vis, which allows you to use natural language to define the visual. Let's use what is the total value of current premium for another counter. Add in widget title to name it total premium value
7. To visualise the average underlying premium by county, use a bar visualisation. Select county on the x axis and underlying premium for the y axis (use AVG here). Add in avg for fixed expenses in the y axis to show the ratio of cost/premium
8. To visualise the relationship between the average cgr factor and the underlying premium, use a scatter plot. Select cgr_factor on the x axis and underlying_total_premium for the y axis (use AVG here).
9. Add a filter to allow us to filter the data with county. Notice that the other dataset will not be filtered as the column does not exist there.
10. Add another filter to allow us to filter the data on gender. As the column exists in both datasets, we will need to add both columns
11. This dashboard can be scheduled and shared



Connect to PBI

1. There are a few options to connect to PBI
 - a. Click on open in PBI desktop. This will download a connection file that you could open directly on PBI desktop
 - b. Click on publish to PBI service which can either publish the entire schema as a semantic model with relationships

Catalogs > beatrice_liew >
 beatrice_liew.vehicle_data
 Owner: beatrice.liew@databricks.com
 Tags: Add tags
 Comment: Add comment

Publish to Power BI workspace Create

Open in Tableau
 Open in Power BI Desktop
 Publish to Power BI workspace

Tables Volumes Models Functions Details Permissions

Filter tables 5 tables

Name	Created at	Owner	Popularity
cgr_definitions_table	2024-02-14 13:17:29	beatrice.liew@databricks.com	all
cgr_premiums_table	2024-02-14 13:21:40	beatrice.liew@databricks.com	all
cgr_premiums_table_agg	2024-02-14 18:00:34	beatrice.liew@databricks.com	----
cgr_premiums_table_location	2024-02-14 18:00:31	beatrice.liew@databricks.com	all
territory_definitions_table	2024-02-14 13:24:17	beatrice.liew@databricks.com	all

- c. Open PBI, click on Get data -> sign in under > Azure Databricks
 - i. Fill out the information below, you can find the connection details on your DB SQL warehouse, under connection details
 1. Server hostname
 2. HTTP path
 3. Catalog (optional)
 4. Schema(optional)
 - ii. Select the table to query

SQL Warehouses

Shared Endpoint

Overview **Connection details** Monitoring

Server hostname
adb-984752964297111.11.azuredatabricks.net

Port
443

Protocol
https

HTTP path
/sql/1.0/warehouses/2e78d4a8647e3893

JDBC URL
2.6.25 or later

Databricks supports drivers released within the last two years. [Download drivers here](#)

jdbc:atabricks://adb-984752964297111.11.azuredatabricks.net:443/default;transportMode=http;ssl=1;AuthMech=3;httpPath=/sql/1.0/warehouses/2e78d4a8647e3893;

Use these details to connect to this warehouse

Tableau Power BI dbt Python Java Node.js Go More tools

Azure Databricks

Server Hostname
adb-984752964297111.11.azuredatabricks.net

HTTP Path
/sql/1.0/warehouses/2e78d4a8647e3893

Advanced Options (optional)

Default catalog (optional)
beatrice_liew

Database (optional)
retail_data

Automatic Proxy Discovery (optional)

Native query (Requires: Default catalog) (optional)
Example: select * from db.schemaname.tablename

OK Cancel

Genie demo

1. Using the Cgr_premiums_table table, click on create -> Genie Data Room
2. Within Data room, select all three tables used in this exercise and create room
3. Click on explain dataset to show how Genie explains the dataset from column names. Note that the explanation includes relationships mapped out during data modelling. But without modelling, data rooms can typically infer relationships if the col names are similar
4. Ask some example questions or use surprise me to get some insights from the data
5. Pick an example question and run it
 - a. What is the age distribution of customers?
 - b. Find the trend between average current_premium and the age of the customer

- c. Find the average current premium by town, which town has the highest premium?
 - d. What is the average current premium for each gender across all territories?
 - e. What is the average current premium for each gender split by county?
 - f. Which quarters are customers born in? - Create and use get_quarter function
6. Ask some questions such as which age group has the highest current premium
7. Create a table function, log it to unity catalog and ask similar questions for Genie to trigger this function, which will be tagged as a trusted asset as end users can see that it's a curated answer that they can rely on

To create get_quarter function

```
create or replace function get_quarter(date DATE) RETURNS INT
RETURN
CASE WHEN MONTH(date) BETWEEN 2 AND 4 THEN 1
WHEN MONTH(date) BETWEEN 5 AND 7 THEN 2
WHEN MONTH(date) BETWEEN 8 AND 10 THEN 3
ELSE 4
END
```

To create a table function as a trusted asset:

```
CREATE OR REPLACE FUNCTION average_current_premium()
returns table (town string, average_current_premium double)
comment "What is the average current premium by town?"
return
SELECT
    t.town,
    AVG(p.current_premium) AS average_current_premium
FROM
    beatrice_liew.vehicle_data.cgr_premiums_table p
    JOIN beatrice_liew.vehicle_data.territory_definitions_table t ON p.territory =
t.territory
GROUP BY
    t.town
ORDER BY
    average_current_premium DESC;
```