

ETL AND TEXTUAL ANALYTICS

Data Analytics for Professional Accountants (ACCTG 522)

Class 3 | MPAcc class of 2026

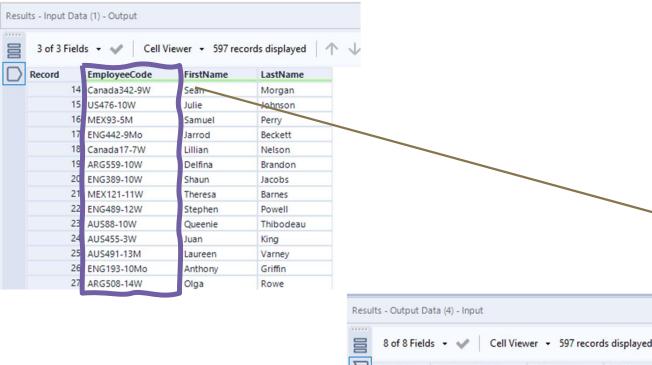


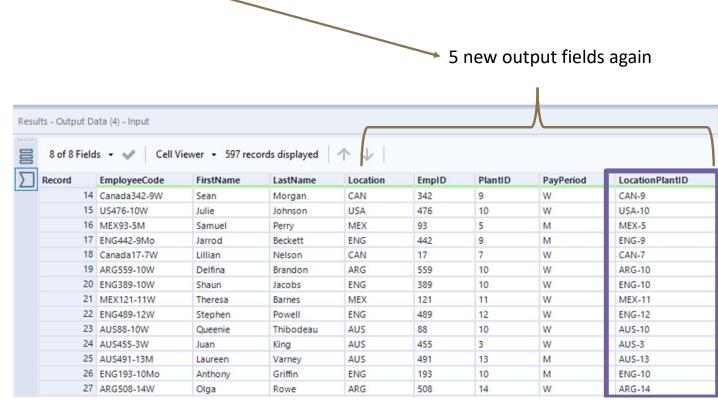
Welcome!

Class Agenda

- Review and ETL Overview
- ETL Cases in the FSB Remote Labs
- Conclusion and look ahead

Review: ETL Overview





True/False: Our regular expression (\u*)(\d*)-(\d*)(\u*) only worked beacuse of an Alteryx-specific feature?



True: Alteryx treats \u, \l, and \w the same

False: \d is any digit

True: Alteryx has an option to ignore case

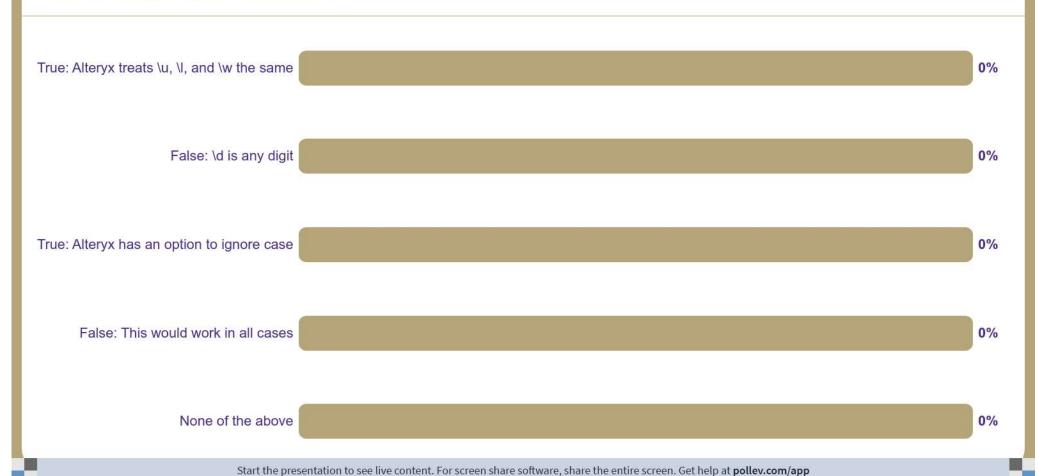
False: This would work in all cases

None of the above

Start the presentation to see live content. For screen share software, share the entire screen. Get help at pollev.com/app

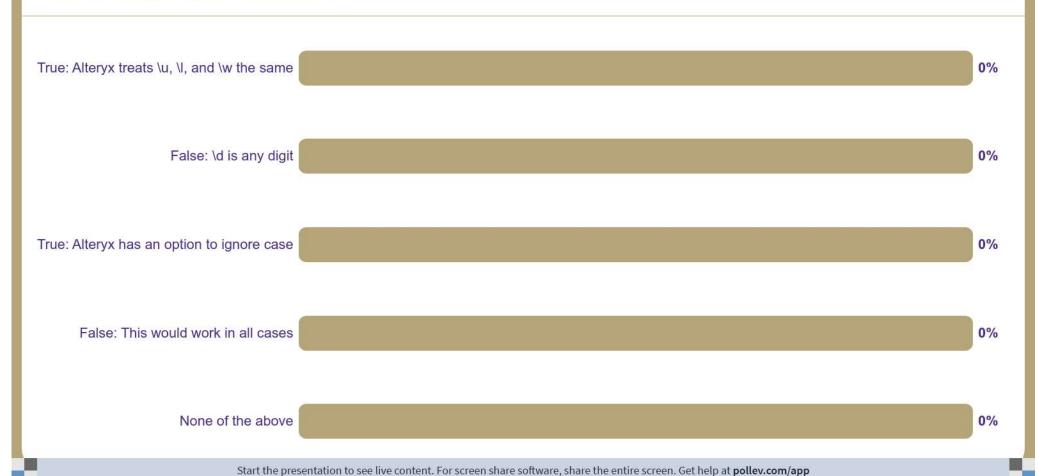
True/False: Our regular expression (\u*)(\d*)-(\d*)(\u*) only worked beacuse of an Alteryx-specific feature?





True/False: Our regular expression (\u*)(\d*)-(\d*)(\u*) only worked beacuse of an Alteryx-specific feature?





Extension

Using RegEx101 tool attempt to Transform these similar EmployeeCodes:

Canada1234-1Mo Mexico0001234-121Monthly Australia45789-13Annual AUS123-13M FRA334455-145weekly GRE67 1a

Take aways

```
Regular expressions:

[A-Za-z] lists of acceptable terms (alphabetic)

[A-Za-z0-9] as above, alphanumeric

*,+ zero, one or more

^,$ beginning, end of the line

| or

\d any digit, \d{N} a specific "N" number of digits

. Any character, \. A full stop (escaping)
```

Extension Exercise

```
^[A-Za-z]+\d+(-|\s)\d+[MAWmaw]
^[A-Za-z]+\d+[-\s]\d+[MAWmaw]
```

After Class Solution:

```
^([A-Za-z]+)(\d+)[-\s](\d+)([MAWmaw])[A-Za-z]*$
Captured (line 1, 4 separate outputs):
Canada,1234,1,M
Discarded (line 1):
-,0
```

```
REGULAR EXPRESSION

! / ^([A-Za-z]+)(\d+)[-\s](\d+)([MAWmaw])[A-Za-z]*$

TEST STRING

Canada1234-1Mo !

Mexico0001234-121Monthly!

Australia45789-13Annual!

AUS123-13M!

FRA334455-145weekly!

GRE67 1a
```

Extra RegEx examples

Dates:

- Invoice Date: 03/15/2025
- Payment Due: 04/14/2025

^\D*\s*(\S*)

\d*\/\d*\/\d*

Invoices:

- INV-100234 ^(?!Reference :)INV-\d+\$
- INV-567890
- Reference: INV-456

 $^([A-Z]+)-(\d^*)$

INV-(\d{6})

Extra RegEx examples

Firm Identifiers:

- CIK: 0000320193
- Ticker: TSLA
- Ticker: AAPL

[CIK|Ticker]+: (\d+|[A-Z]+)

Matching dollar amounts and accounts:

- Total: \$12,450.00 ^Total:\s([\$]\d+,*\d+.*\d*)\$
- Late Fee: \$25
- Refund: (\$1,200)
- .+: (.+)

Extra RegEx examples

Phone numbers, email, credit card numbers:

- 206-123-4567
- ^\d{3}-\d{3}-\d{4}\$
- real_person@example.com
- Potential fraud detected: 4532-9876-1234-5678

```
(\d+-\d+-\d+)($|-(\d+))|([A-Za-z]+)([@_.]) ([A-Za-z]+)([@_.]) ([A-Za-z]+)([@_.])([A-Za-z]+)
```

ETL Cases in the FSB Remote Labs

ETL Cases, Continued

Today we will examine "Joining Data":

- Joining data is another **fundamental step** in the ETL process, as in most cases, developing insight will require data from different databases.
- Joining data can sometimes be referred to as merging data, blending data, or combining data.

Why is joining data important?

Joining data is part of the transforming step in ETL.

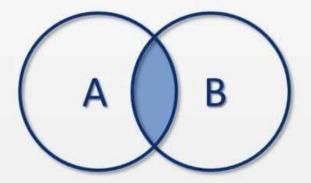
- Incorrect joins are a data quality issue, incorrect joins can:
 - crash (or really slow down) our programs!
 - have us lose important data.
 - Have us generate incorrect conclusions/insight from our analysis.

Joining using unique keys

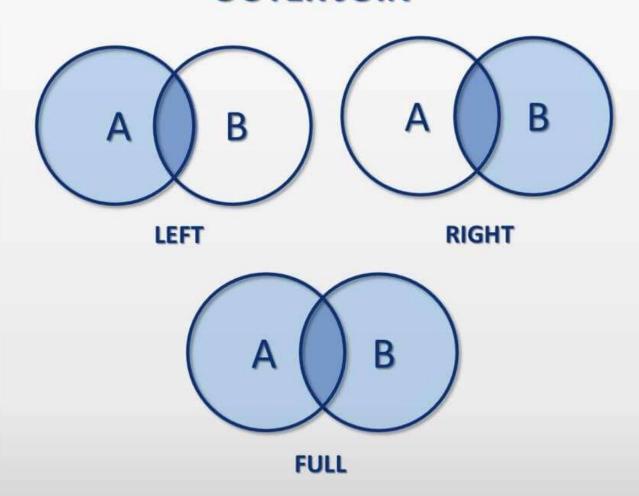
In all cases we will use unique keys to join the data.

- All the unique keys are described in the case
- One of them was tricky, can you remember which one?

INNER JOIN



OUTER JOIN



A B

SELECT <select_list> FROM TableA A LEFT JOIN TableB B

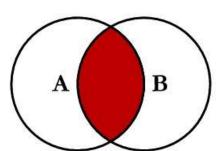
ON A.Key = B.Key

A B

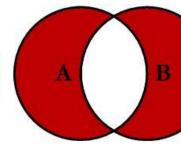
SELECT <select_list>
FROM TableA A
LEFT JOIN TableB B
ON A.Key = B.Key
WHERE B.Key IS NULL

SELECT <select_list>
FROM TableA A
FULL OUTER JOIN TableB B
ON A.Key = B.Key

SQL JOINS

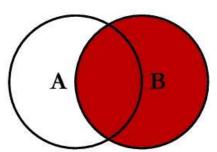


SELECT <select_list>
FROM TableA A
INNER JOIN TableB B
ON A.Key = B.Key

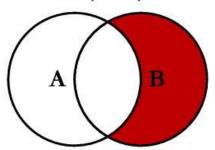


© C.L. Moffatt, 2008

B



SELECT <select_list>
FROM TableA A
RIGHT JOIN TableB B
ON A.Key = B.Key



SELECT <select_list>
FROM TableA A
RIGHT JOIN TableB B
ON A.Key = B.Key
WHERE A.Key IS NULL

SELECT <select_list>
FROM TableA A
FULL OUTER JOIN TableB B
ON A.Key = B.Key
WHERE A.Key IS NULL
OR B.Key IS NULL

Joining using unique keys

In all joins we will use "one to many" matches.

 Joining the large information source ("JELineItems") to the smaller set of information in the other data sheets.

Why join across multiple software?

Because joining data is so fundamental, we want to work towards becoming software agnostic.

Working with cutting-edge software allows you to:

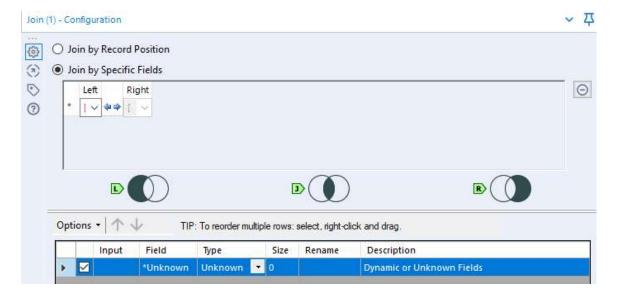
- Become adaptable and resilient to new software.
- Impress recruiters with practical discussions of cutting-edge software.

Joining across different tools:

Some terminolgy:

- Alteryx: drag and drop "Join" and "Union" tools
- PowerBI: data relationships tab
- Python: code-level with package "pandas"
- SQL: base code-level
- Tableau: Both relationships and joins.

Alteryx Joins:





SQL/Python:

```
SELECT
JEL.*,
BU.*

FROM JELineitems AS JEL
LEFT JOIN BusinessUnits AS BU
ON JEL.BusinessUnitID = BU.BusinessUnitID;

SELECT
JEL.*,
BU.*

FROM JELineitems AS JEL
RIGHT JOIN BusinessUnits AS BU
ON JEL.BusinessUnitID = BU.BusinessUnitID;
```

```
import pandas as pd
# Example: read from Excel file (each sheet as
DataFrame)
file path = "your file.xlsx"
# Load the two sheets
JELineitems = pd.read excel(file path,
sheet name="JELineitems")
BusinessUnits = pd.read excel(file path,
sheet_name="BusinessUnits")
# Perform SQL-style join on BusinessUnitID
# Equivalent to: SELECT * FROM JELineitems
#
         INNER JOIN BusinessUnits
         ON JELineitems.BusinessUnitID =
BusinessUnits.BusinessUnitID;
merged df = JELineitems.merge(
  BusinessUnits,
  on="BusinessUnitID",
  how="inner" # can also use 'left', 'right', or 'outer'
```

Lab Agenda

- 1. Open files in excel / open case and identify the keys
- 2. Compare Tableau and Power BI

Tableau Relationships

☐ JELineItems+ (Analytics_mindset_case_studies_ET...

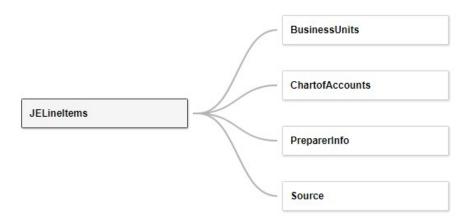
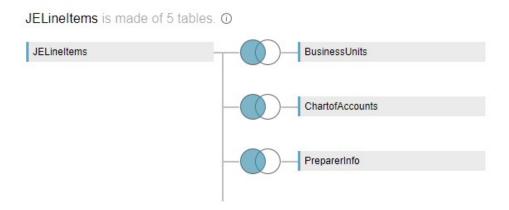
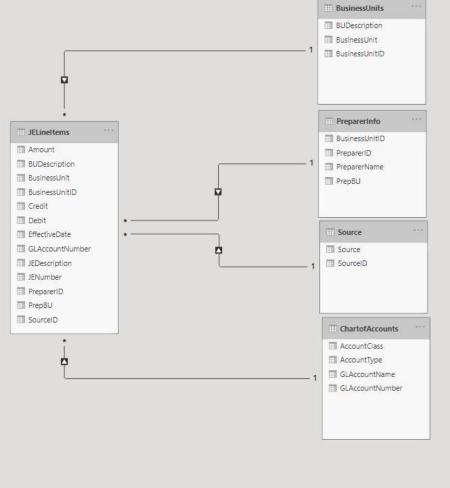


Tableau Joins



Power BI Data Relationships



To FSB remote labs

Again I recommend **saving screenshots** of your work in Alteryx (I've added blank slides if you want to save them in this ppt document, ctrl+m adds more slides)

Two Potential Issues Explained:

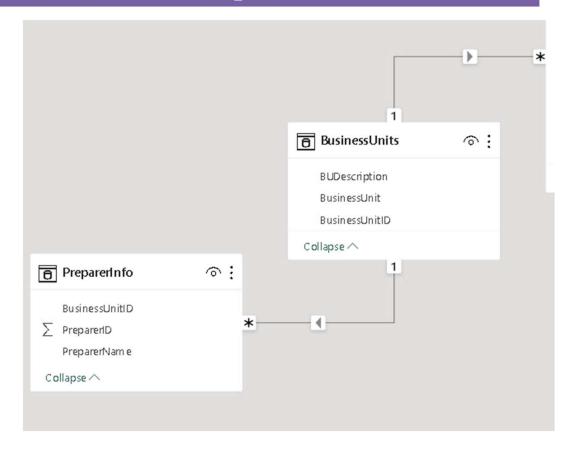
Issue 1 In the underlying dataset PreparerID is **not unique:**

For example ID 101 appears 6 times, repeating once for each BusinessUnit

PreparerName	PreparerID BusinessUnitID	
Carolyn Slater Food Operations Manager	101	1
Alexander Knox Assistant Manager	101	2
Trevor Roberts Housekeeping Supervisor	101	4
Gavin McGrath Dining Commons Manager	101	3
Rodney Wilkins Maintenance Supervisor	101	6

Two Potential Issues Explained:

PowerBI: maps the multiple prepareIDs to a unique BusinessUnitID, allowing it to join both datasets to JELineItems via BusinessUnitID.



Two Potential Issues Explained:

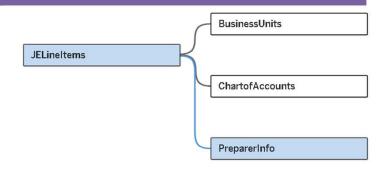
Tableau: we can resolve the alert by filling out the keys which are

BusinessUnitID= BusinessUnitID

AND

PreparerID = PreparerID.

Which means it has to match both keys.



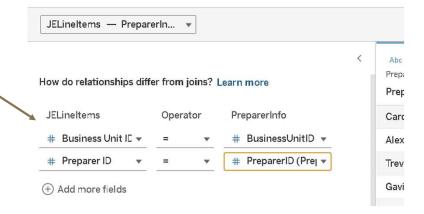


Tableau Completed

Add the Source dataset and the **data** relationship is complete.

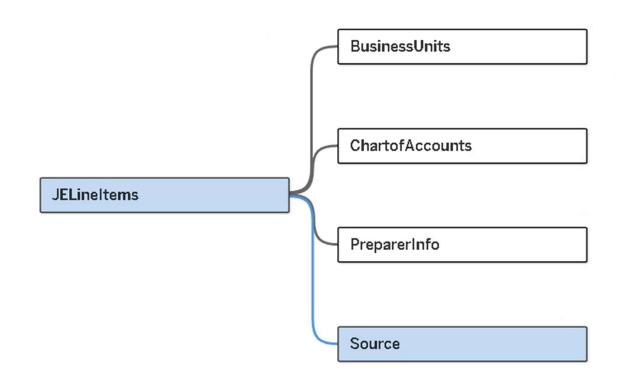


Tableau Join Version

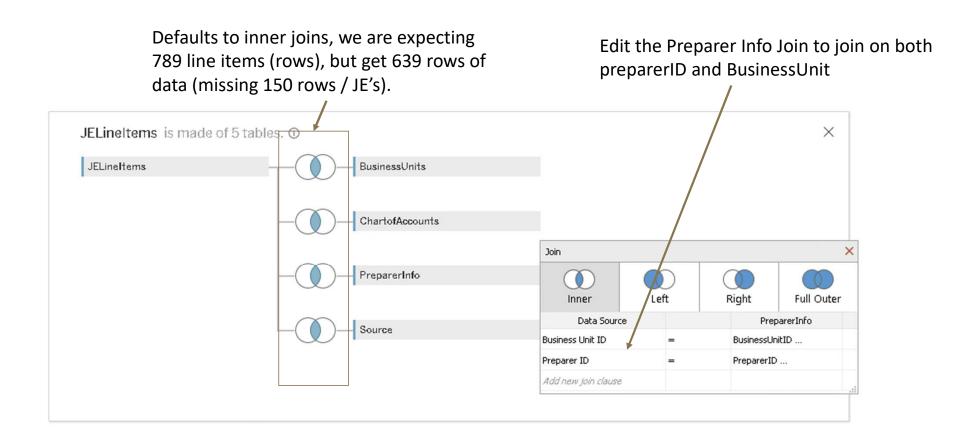


Tableau Join Version

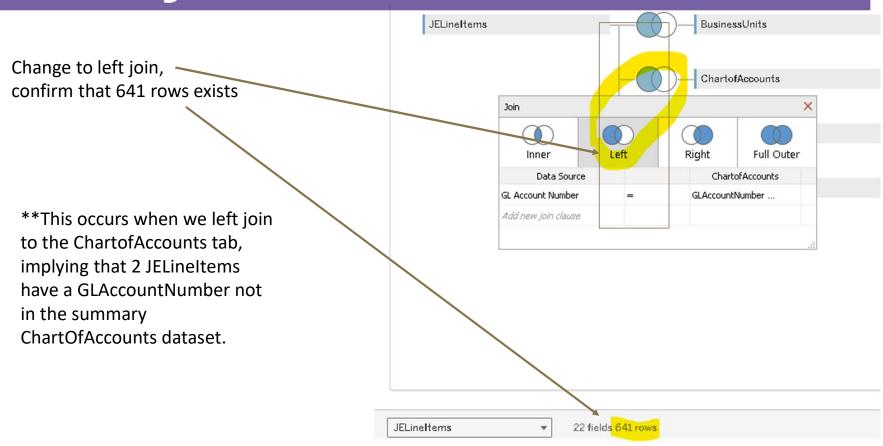
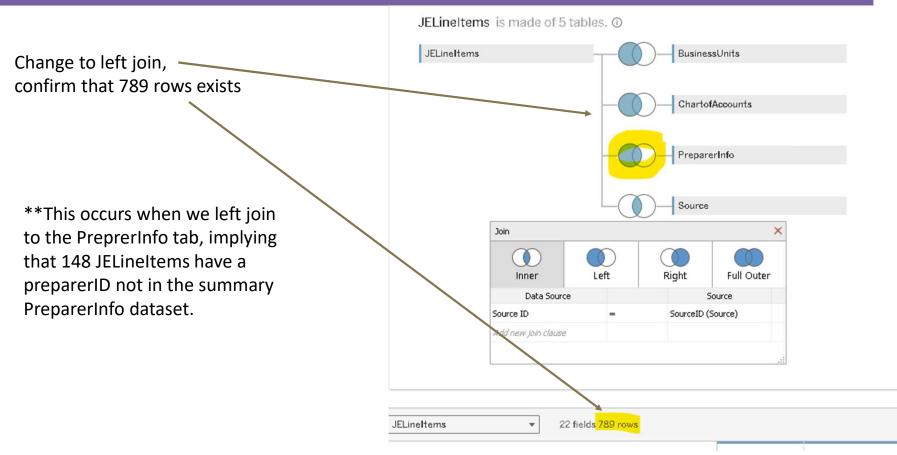
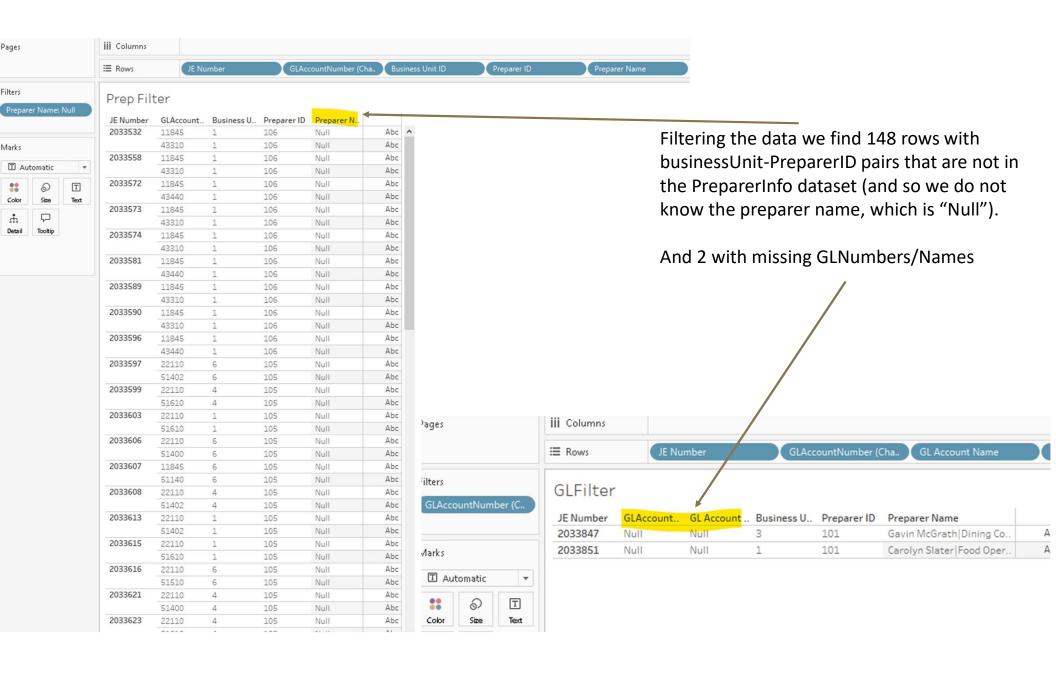


Tableau Join Version





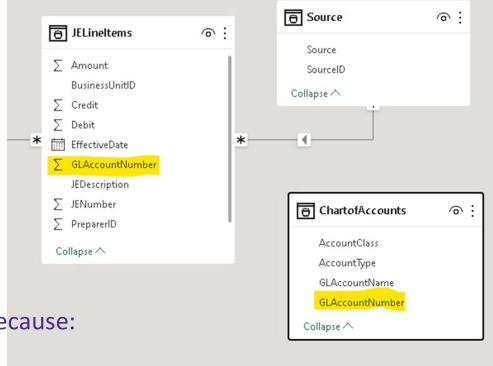
Two Potential Issues Explained:

Issue 2 GLAccounts not mapped in PowerBI:

PowerBI read GLAccountNumbers as numbers instead of string format (warned us with 2 read errors)

PowerBI Workflow

Notice that GLACccountNumber is summed (numeric) in JELineItems but not in Chart of accounts (means it is string).

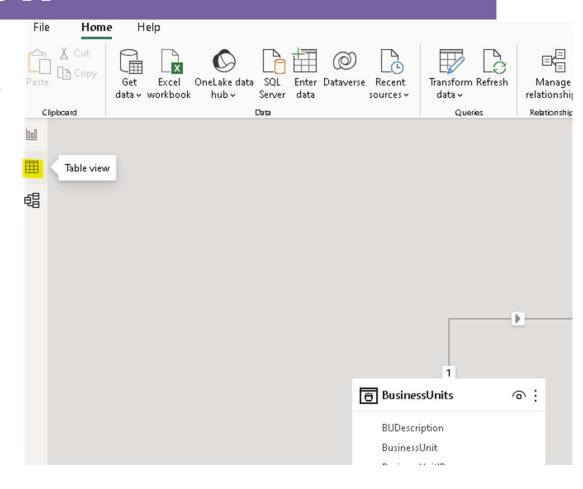


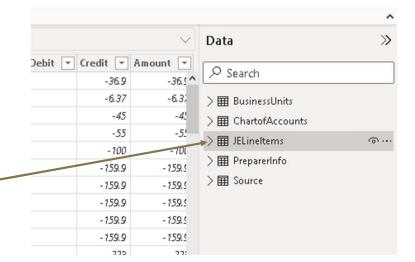
*Joins can't be between numbers and string because:

"1"≠1

PowerBI Workflow

Go to Table view on left index

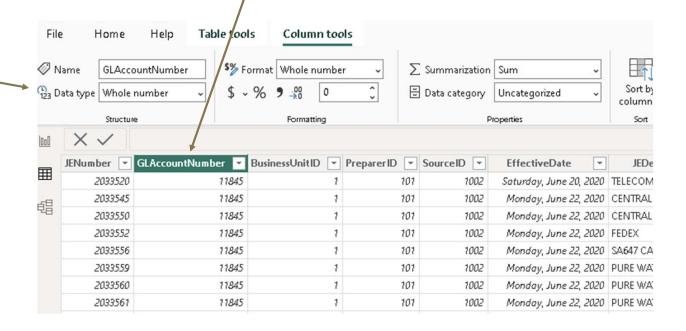


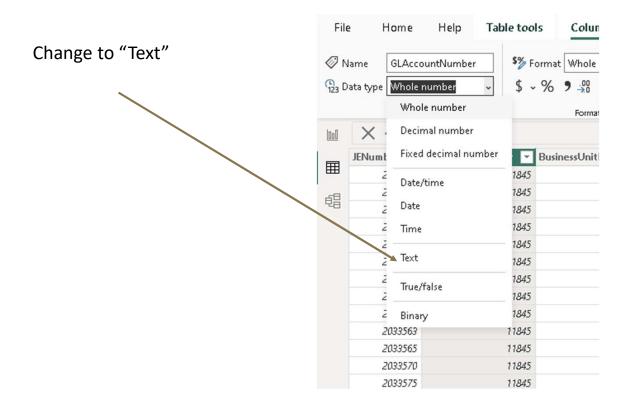


On the right hand menu select the JELineitems dataset

Select the GLAccountNumber Field / Column

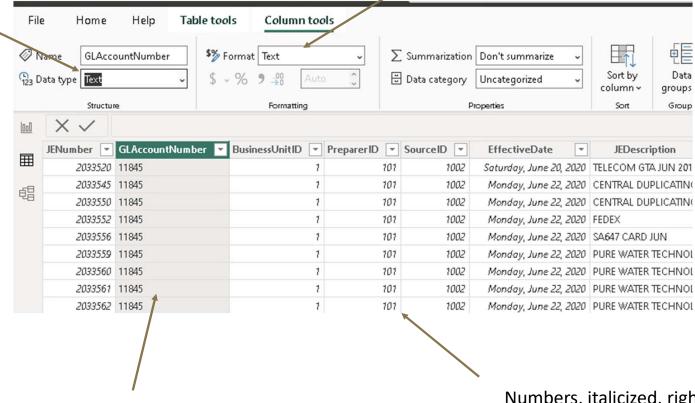
Note Dat type is whole number





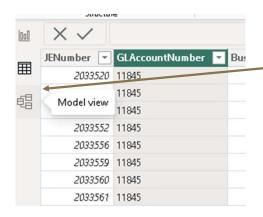
Confirm the change and how it affects format

Note format is just how the data is displayed, and it updates automatically to text (Data Type is how it is recognized by PowerBI).



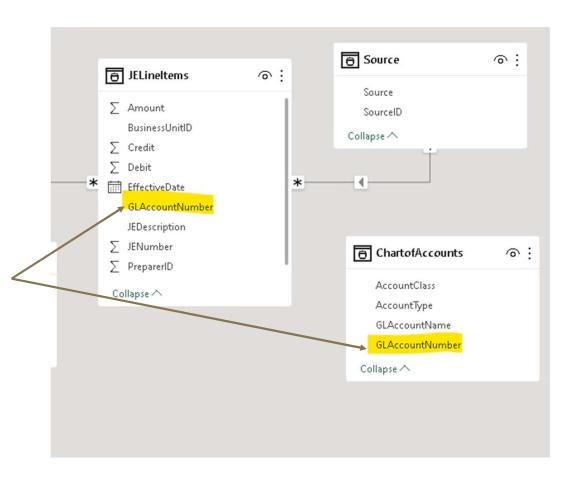
Text, left aligned

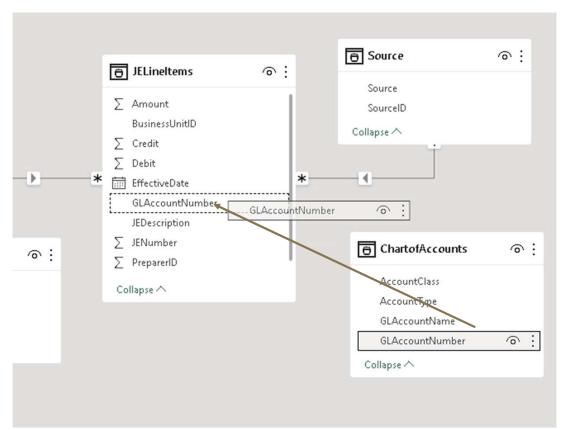
Numbers, italicized, right aligned



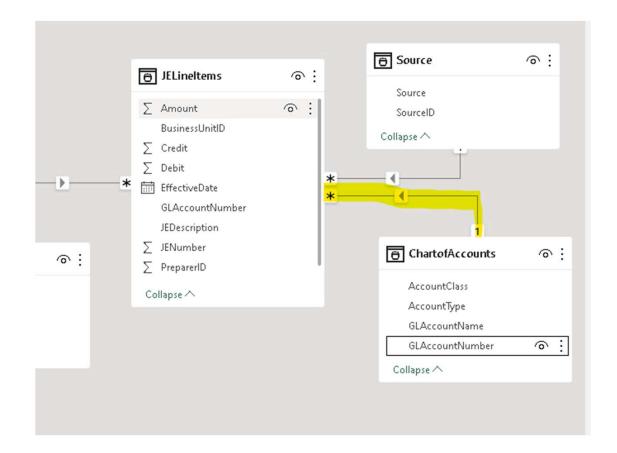
1 return to model view

Confirm
GLAccountNumber is text
in both (neither include
the summation sign)

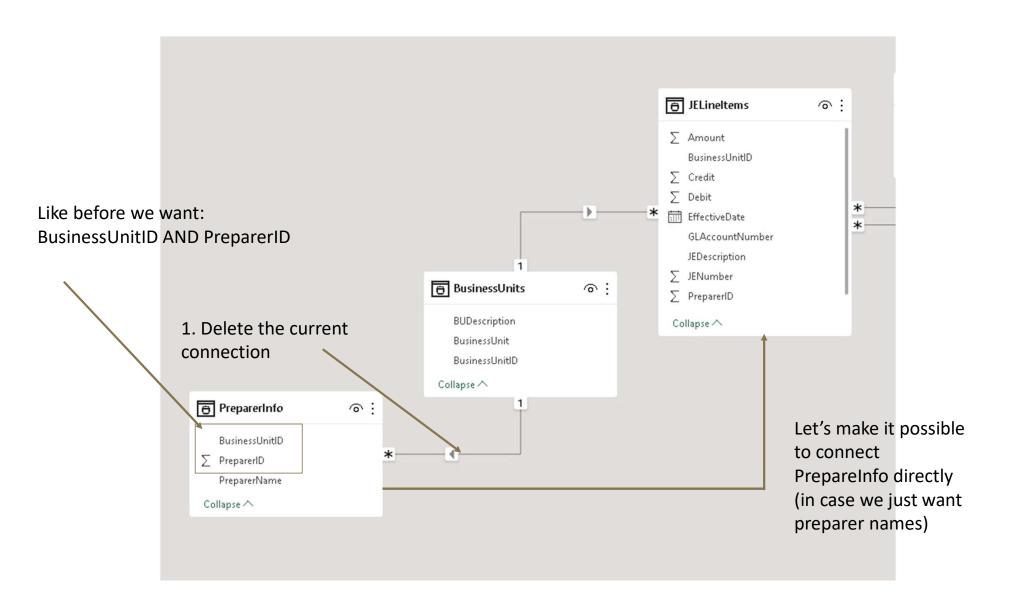


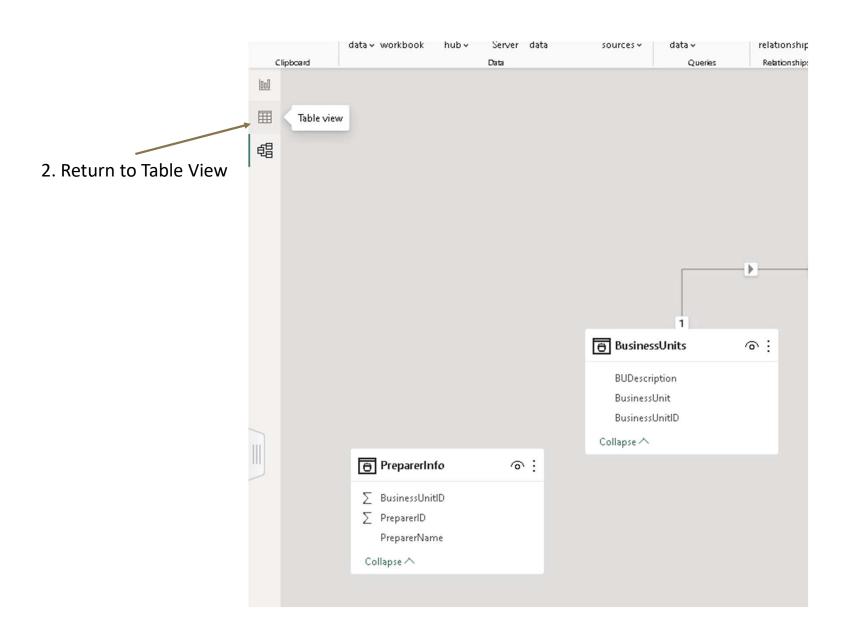


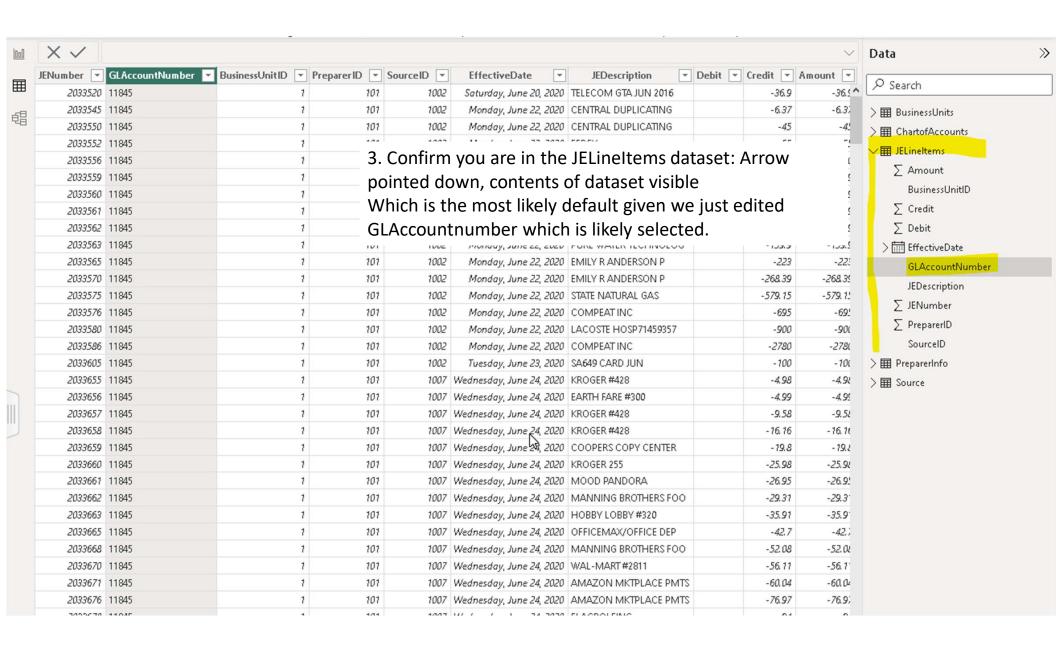
Drag GLAccountNumber from ChartOfAccounts to JELineItems and match it with GLAccountNumber

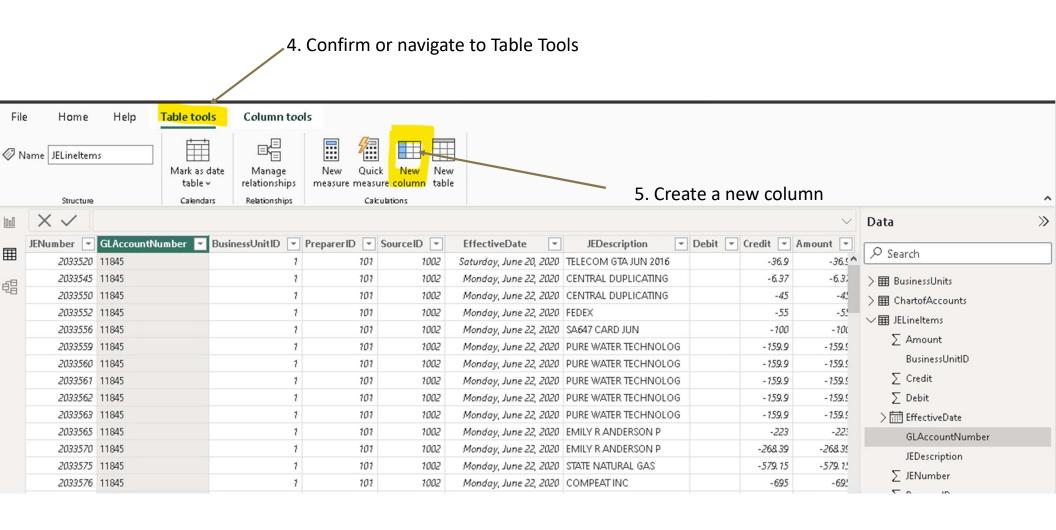


New Relationship Arrow appears in the data model

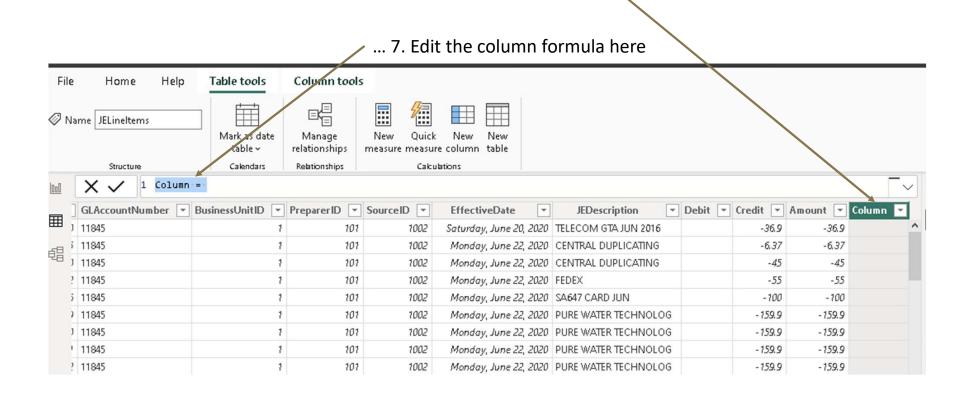








6. Confirm new column and ...



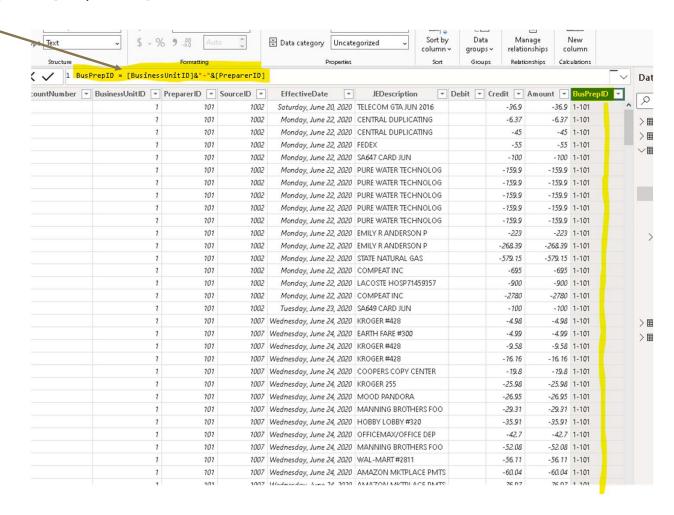
Our formula goal is to concatenate BusinessUnitID and PreparerID to create a unique key i.e., 1-101, 2-101 etc.

Our formula goal is to concatenate BusinessUnitID and PreparerID to create a unique key i.e., 1-101, 2-101 etc.

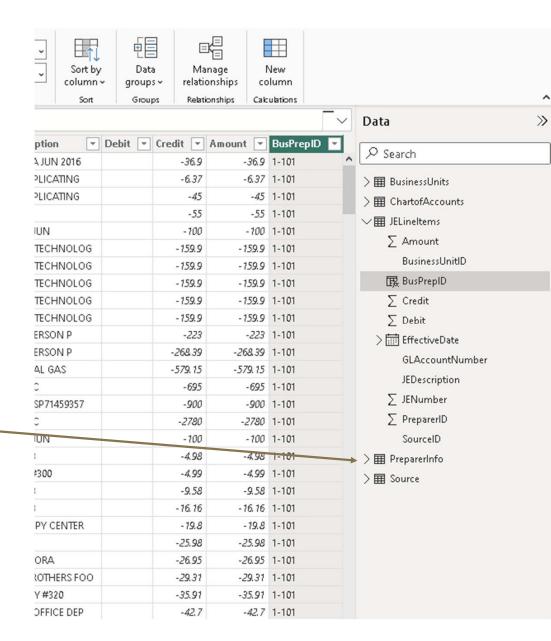
8. Replace Column = with:

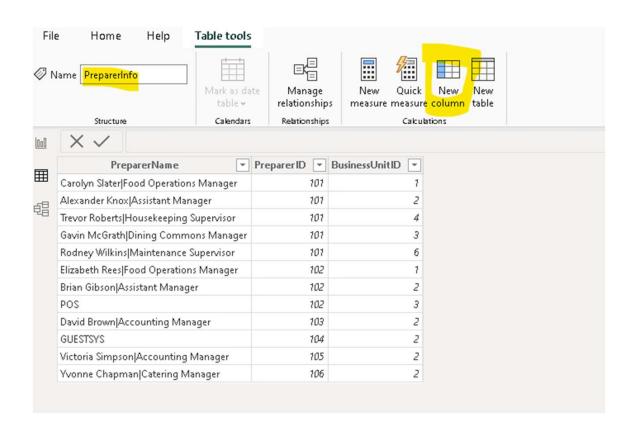
BusPrepID = [BusinessUnitID]&"-"&[PreparerID]

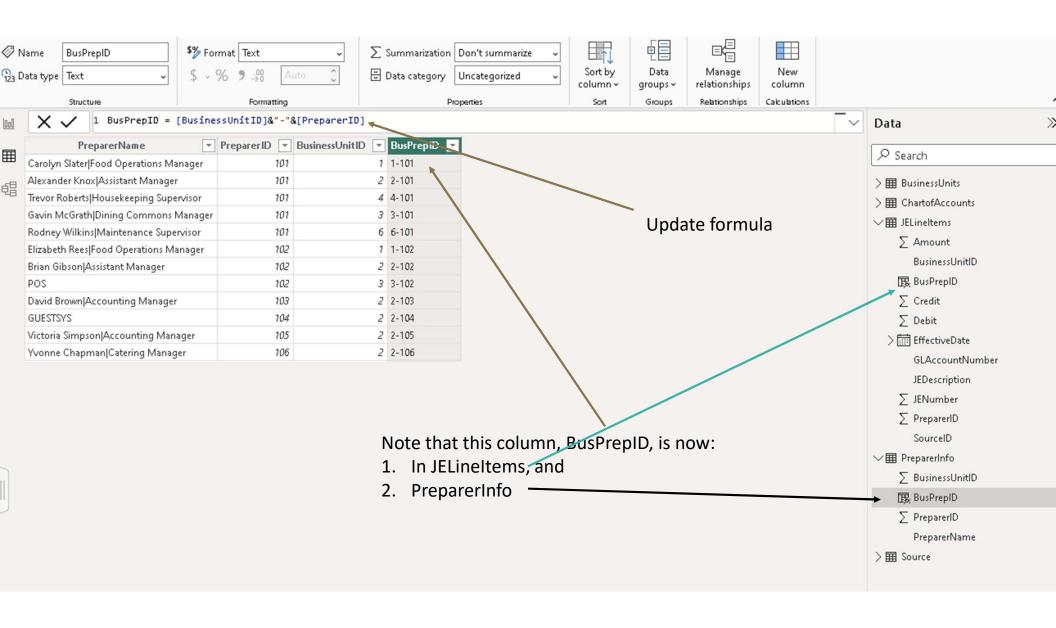
Then press Enter

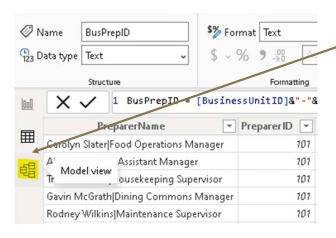


10. Still in Tables, select the PreparerInfo Dataset and create the same new column "BusPrepID"



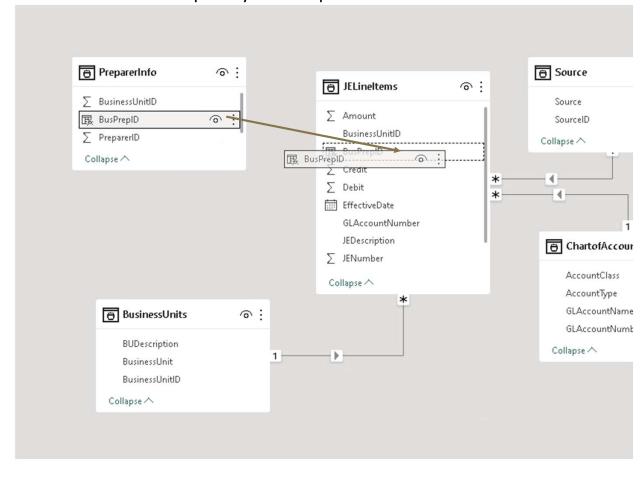






Return to model view

Connect PreparerInfo and JELineItems by dragging the new unique key "BusPrepID"



PowerBI Completed

one (1)—

Confirm the new BusPrepID connection and the data relationship **model** is complete.

6 Source @: Source SourceID Collapse ^ AccountClass (a) JELine tems @: AccountType 民 BusPrepID GLAccountName ∑ Credit GLAccountNumber ∑ Debit Collapse ^ Effective Date GLAccountNumber **JEDescription** ∑ JENumber BusinessUnits ◎: ∑ PreparerID BUDescription SourceID BusinessUnit Collanse A BusinessUnitID Collapse ^ Note it should be PreparerInfo @: to many (*) BusinessUnitID -(1)-民 BusPrepID PreparerID Collapse ^

Conclusion and look-ahead

ETL is the first step in the analytical process – we cannot undertake data analytics without data!

ETL is especially important when the data needs to be cleaned or transformed for use in analytics.

Conclusion and look-ahead

Wednesday:

Cases: Enron Emails Case – RegEx Heavy

