

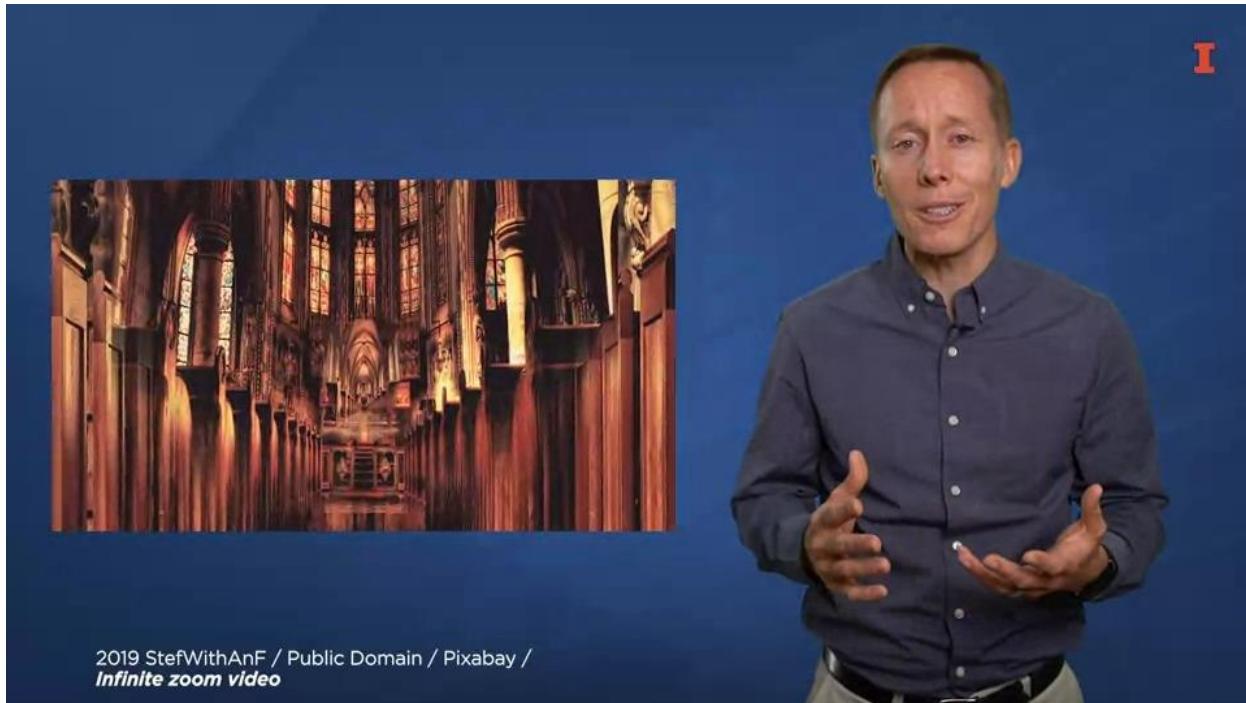
Module 2: ETL and EDA Using PowerBI

Table of Contents

Module 2: ETL and EDA Using PowerBI.....	1
Module 2 Overview.....	2
Module 2 Introduction.....	2
Lesson 2-1: Business Problem.....	11
Lesson 2-1.1 Data.....	11
Lesson 2-1.2 Business Problem.....	19
Lesson 2-2: Introduction to Power BI.....	28
Lesson 2-2.1 Introduction to Power BI	28
Lesson 2-2.2 Installing Power BI	40
Lesson 2-2.3 ETL 1: Examine Data with Power Query Editor.....	53
Lesson 2-2.4 ET2: Dates and Calculated Columns.....	67
Lesson 2-2.5 ETL 3: Checking for and Eliminating Outliers	75
Lesson 2-2.6 ETL 4: Data Models and Joins.....	89
Lesson 2-3: EDA with Power BI	108
Lesson 2-3.1 EDA 1: Univariate Plots for Numeric Data: Histograms and Boxplots	108
Lesson 2-3.2 EDA 2 - Univariate, Bivariate, and Multivariate Plots with Categorical	122
Lesson 2-3.3 EDA 3: Filters, Slicers, and Drill Through.....	136
Lesson 2-3.4 EDA 4: Multivariate Plots: Scatter Plots and Line Plots.....	150
Lesson 2-3.5 EDA 5: Publishing a Power BI Report	162
Module 2 Review	167
Module 2 Conclusion	167

Module 2 Overview

Module 2 Introduction



2019 StefWithAnF / Public Domain / Pixabay /
Infinite zoom video

Stephen Few, a data visualization guru, shares how organizations like the CIA teach new spy recruits to make observations. They're trained to first get an overview of what's going on around them and then when they spot something abnormal, they shift from a broad observational alertness to more focused perspective, in which they analyze the details. In this module, you'll learn how to use power BI to do with data, what spies do when observing their surroundings, get an overview of the data and then narrow in on certain aspects of the data that seem abnormal and analyze them.

I

“Having an overview is very important. It reduces search, allows detection of overall patterns, and aids the user in choosing the next move. A general heuristic of visualization design, therefore, is to start with an overview. But it is also necessary for the user to access details rapidly. One solution is overview + detail: to provide multiple views, an overview for orientation, and a detailed view for further work.”

- Few, Stephen. Now You See It, p 84

Stephen Few goes on to further quote some data visualization colleagues about why this approach is so useful. Having an overview is very important, it reduces search, allows detection of overall patterns, and aids the user in choosing the next move. A general heuristic of visualization design, therefore, is to start with an overview. But it is also necessary for the user to access details rapidly. One solution is overview plus detail; to provide multiple views, an overview for orientation and a detailed view for further work.

Overview first, zoom and filter,
then details-on-demand



In fact, there's something called Schneiderman's Mantra, which is overview first, zoom and filter, then details-on-demand.



Now Powered BI is a great tool for implementing Schneiderman's Mantra because it is optimized for visualizing data, and so making charts is very quick and intuitive. Power BI is so intuitive it empowers many people, including those who do not have coding skills to make a wide variety of charts and then to use those charts to zoom in and filter the data to get a more focused perspective.

Business Analytic Workflow

Frame a question

Identify data sources and how to acquire the data

Data management and extraction, transformation, and loading (ETL)



With respect to the Business Analytic Workflow, Power BI excels Exploratory Data Analysis. In some sense, Power BI is like a grown up version of Excel's pivot table and pivot chart functionality, like excel, Power BI is pretty intuitive to use and you don't need to know much about programming to use it.

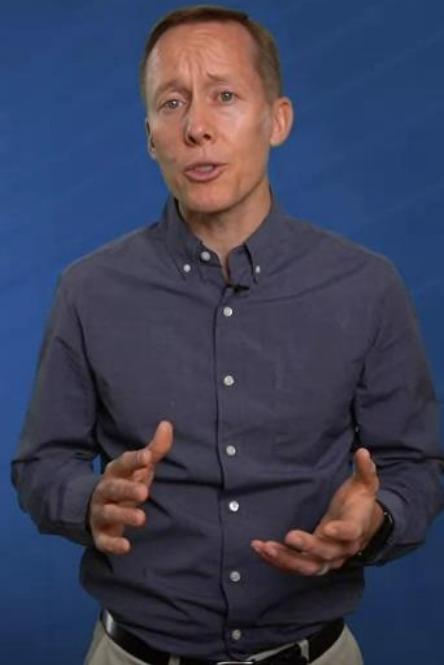
Business Analytic Workflow

Exploratory data analysis (EDA)

Data modeling

Results analysis and business insight

Visualizing and communicating findings and solutions



Also, like Excel, Power BI has a really powerful tool for interacting with the data. The Power Query Editor, you can manipulate the shape of the data from wide to long. You can create new fields that are based on calculations of other fields. You can split or combine text fields. You can also perform Joynes and Lookups.

FACT Framework:

Frame the question

Assemble the data

Calculate the results

Tell others the results



Thus, many of the key assembly processes that can be done in excel, can also be done in the Power Query Editor. The method for doing so is different, because you can't point and click on cells to make changes. However, the benefit is that you can handle much larger amounts of data. While Power BI is a very powerful tool, it is only a tool unless you have first framed a question and then have access to relevant data, which are the first two parts of the fact framework, then you're not as likely to convert data into action even when you do have a well framed question and good data, obtaining actionable insights is not a straightforward process. Quoting his data visualization colleagues, Stephen Few writes,

I

“Users often try to make a “good” choice by deciding first what they do not want, i.e., they first try to reduce the data set to a smaller, more manageable size. After some iterations, it is easier to make the final selections from the reduced data set. This iterative refinement of progressive querying of data sets is sometimes known as hierarchical decision-making.”

- Few, Stephen. Now You See It, p 84

"users often try to make a good choice by deciding first what they do not want, that is, they first try to reduce the data set to a smaller, more manageable size. After some iterations, it is easier to make the final selections from the reduced data set. This iterative refinement of progressive querying of data sets is sometimes known as hierarchical decision making."



Now, in my experience, I found that this quote is accurate. Power BI helps with this hierarchical decision making process by using visualizations to identify and communicate useful relationships in the data. This lessons are not meant to be a comprehensive tutorial on Power BI. There are many excellent tutorials for becoming familiar with all the bells and whistles of Power BI.

Objectives

1. Identify the strengths and weaknesses of Power BI for assembling data and exploratory data analysis through hands on practice.
2. Use the Power Query Editor to assemble and preprocess data.



The objectives of this module are for you to; one, identify the strengths and weaknesses of Power BI, for assembling data and Exploratory Data Analysis through hands on practice. Two use the Power Query Editor to assemble and preprocess data.

Objectives

3. Create basic charts in Power BI for exploring and communicating business data.
4. Combine multiple charts into a report that can be published for use by others.



Three, create basic charts in Power BI for exploring and communicating business data, and four, combine multiple charts into a report that can be published for use by others.

Lesson 2-1: Business Problem

Lesson 2-1.1 Data

The slide features a dark blue background. On the left side, there is a graphic element consisting of four colored squares (red, green, yellow, blue) forming a stylized 'T' shape, followed by the word 'TECA' in a large, bold, dark blue sans-serif font. To the right of this graphic, a man with short brown hair, wearing a dark blue button-down shirt, stands with his hands slightly open as if gesturing while speaking. In the top right corner of the slide, there is a small red letter 'I'.

TECA is company that owns over 150 convenience stores and gas stations.

Sells gas, candy, soda, chips, lottery tickets, etc.

In these lessons, we will be using point-of-sale data from Teca, which is the name of a company that owns over 150 convenience stores and gas stations throughout the middle of the United States. These stores sell typical convenience store items; gas, candy, soda, chips, lottery tickets and so on.



This data is known as point-of-sale data because it is generated at the point of the sale, like the cash register or the gas pump. You can probably imagine some of the information that is in this dataset.

Point-of-Sale Data

5 main categories of transaction:

- When item sold
- Customer sold to
- What product sold
- Store location
- Revenue and profit data



The columns in this dataset contain information about five main categories of the transaction. One, information that identifies when the transaction takes place and the identifier of the transaction. Two, information about the customer. Three, information about the product or products that are sold. Four, information about the store where the transaction took place, and five, revenue and profit data.

TECA Data Can Be Aggregated By:

Product/Year

Customer/Parent

Customer/Day

Customer/Parent/Year

Parent/Year



This type of granular data is awesome because we can aggregate it in many different ways, such as by day, customer, product, product category, and even by hour or minute of the day. We will be using a random sample of the data during the years of 2017, 2018, and 2019. Specifically, we will use only 3,000 observations. It would definitely be more informative to look at the full dataset,



but that full data set is so large that we would need a special machine with hundreds of gigabytes of RAM to analyze all of it at once. Now most machines can handle a whole lot more than 3,000 transactions. But since there is bound to be a wide variety of machine computing power among you learners, we wanted to keep it small enough that hopefully, even the weakest machine will be able to process this data.

Benefits of a Small Dataset:

Reduce wait times

Use as experimental sample that you can apply to larger dataset



Another benefit of using such a small data set is that it allows you to drastically reduce wait times as you experiment with different capabilities of the software. Now, even if you had lots of experience with the data analytics software, it would still be good practice to conduct data assembly and exploratory data analytic tasks using a sample of the data. Again, it will reduce the amount of time that you spend waiting for computations to be performed. Once you zero in on the specific data preparation tasks that need to be performed, then you can very easily set those processes to run on the full dataset.



Description: This is a sample of 3,000 rows and 23 variables of TEGA's point-of-sale data during the years 2017-2019. TEGA is a company that owns over 350 convenience stores and gas stations throughout the middle of the United States. These stores sell typical convenience store items: gas, candy, soda, chips, lottery tickets, and so on. Each row of this dataset represents a unique product of a single transaction. To protect the actual convenience stores, the coordinates have been changed to match those of a nearby post office.

Column Name	Example	Description
unique_id	2612027	A number that uniquely identifies each row.
transaction_id	20181219 562 3 2 4909048	A unique identifier for each transaction.
unformatted_date	3/14/19	Date in mdy format
customer_id	6977.63	A unique identifier for each loyalty customer. Blank for customers that are not loyalty customers.
product_id	2179	A unique code for each product
product_name	HOT DOG FCS	A descriptive name for each product
category_id	280	A unique code for each product category
category_name	Roller Grill Food	A descriptive name for each product category
parent_id	279	A unique code for each category parent
parent_name	Roller Grill	A descriptive name for each category parent
product_count	9	
site_id	297	A unique code for each store
site_name	562 Columbia	A descriptive name for each store
address	1015 Providence Rd	The street number and name of the store's location
city	Columbia	The city of the store's location
zip	65203	The five-digit zip code of the store's location
latitude	38.9509	Latitude coordinates describing how far north the store sits above the equator
longitude	-92.3362	Longitude coordinates describing how far west the store sits from the prime meridian.
site_status	ACTIVE	An indicator of whether the store is active or not.
revenue	2.09	The revenue generated from that line item. Equal to the price times the units.
gross_profit	0.732	The revenue minus the costs.
costs	1.358	The value of the inventory that was sold. Equal to the cost per unit times the units.
units	2	The number of units sold of that product.

Now there are two other files that go along with this dataset. One is a data dictionary that summarizes the contents of the dataset. Feel free to refer to it to get more specific information of what's in the dataset, especially as you first start working with it. Once you get some experience assembling and exploring the data, you'll learn a lot more about it than what is described in the data dictionary.

postal_code	state_province	state_province_code	country_name
35630	Alabama	AL	USA
72712-2500	Arkansas	AR	USA
80012-9997	Colorado	CO	USA
80640-9998	Colorado	CO	USA
36206	Alabama	AL	USA
51436-8702	Iowa	IA	USA
64401-9998	Missouri	MO	USA
64745-9998	Missouri	MO	USA
50469-1073	Iowa	IA	USA
73001-9998	Oklahoma	OK	USA
50665-7732	Iowa	IA	USA
36606	Alabama	AL	USA
81128-9990	Colorado	CO	USA
68660-4511	Nebraska	NE	USA
35124	Alabama	AL	USA



The other file contains information about states and zip codes. We will use this file to join the state data to the point-of-sale data using the data analytics software. There's not a data dictionary for this file, but it's pretty self-explanatory, so you should be able to interpret it without much help.

Lesson 2-1.2 Business Problem



In business analytics, sometimes you start with the business problem and then gather data that will help solve the specific problem. In this set of lessons, we are taking the opposite approach and we will explore the data to see what kind of business problems for which it can provide actionable insight.



In this video, we will elaborate on business problems for which we can gain insight by using the TECA data as a reminder that TECA data is point of sale data that comes from over 150 convenience stores in the middle of the United States.



Each row represents a line item of a transaction. The features of each row are related to five areas of the transaction,

Point-of-Sale Data

5 main categories of transaction:

- When item sold
- Customer sold to
- What product sold
- Store location
- Revenue and profit data



one identifying information about the transaction, including the hour and minute of the day when the transaction took place, two identifying information about the customer if they are loyal to a customer or not, three information about the product or products that were purchased by the customer, four information that identifies the store where the transaction took place, and five revenue and profit data.

About Data Gathering

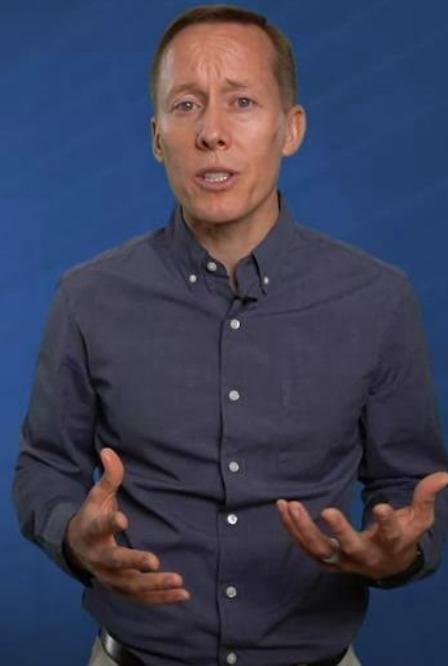
Need actionable business insight
that will help achieve a goal



With such granular data, there is a multitude of potential business insights that we could obtain. Remember that we want a specific kind of business insight, actionable insight that will help us achieve some goal. Broadly speaking, those goals could be related to tactical business issues such as increased customer or employee satisfaction, less frequent stock-outs, less frequent stocking of unpopular goods that spoil or become obsolete, reduction of labor costs, and improved employee scheduling.

Questions About Business Insights

1. When should we order more of a specific candy bar so that we don't run out of stock?
2. What times of the year and day should we schedule extra employees to be working to keep up with customers?



Specifically, some of these questions are, one, when should we order more of a specific candy bar so that we don't run out of stock? Two what times of the year and day should we schedule extra employees to be working to keep up with customers?

Questions About Business Insights

3. What are the high margin products that we can place near the cash register?
4. What are complementary products that customers often purchase together so that we can place them next to each other?



Three what are the high margin products that we can place near the cash register? Four what are complementary products that customers often purchase together so that we can place them next to each other?

Questions About Business Insights

5. How can we price our products to obtain the highest profit?
6. Which customers do we expect to see on a regular basis so that we can offer them a reward?



Five how can we price our products to obtain the highest profit? Six which customers do we expect to see on a regular basis so that we can offer them a reward?

Strategic Questions

1. Should we operate and manage a bakery, or lease out the space to an external bakery or fast food restaurant?
2. If we have the financial means to invest in a new store, where should that store be located?



This data can also be used for strategic decisions, such as introducing new products by entering into a long-term relationship with a supplier, investing in another store, and employee training. Some of these specific questions might be the following. One should we operate and manage a bakery or lease out the space to an external bakery or fast food restaurant? Two if we have the financial means to invest in a new store, where should that store be located?

Strategic Questions

3. Which employees would be best-qualified for providing training to other employees?



Three, which employees would be best qualified for providing training to other employees?



This is just a sample of the questions that we could seek to gain insight about from the TECA data. I hope this gets you thinking about additional questions that you could ask and how you would use the data to get insight to those questions. Coming up with questions that you can answer with the data is probably more important than the data itself,

Remember:

Data itself does not provide insight

Mind is most important tool



I want to emphasize that the data itself does not provide insight. While it's true that we are going to focus on using business analytic tools to summarize the data and look for patterns in the data, remember that your mind is the most important business analytic tool. Try not to lose sight of the overall business problem as we get into the details of cleaning and exploring the data.

Lesson 2-2: Introduction to Power BI

Lesson 2-2.1 Introduction to Power BI

The image consists of two parts. On the left is a screenshot of the Gartner Magic Quadrant for Analytics and Business Intelligence Platforms. The chart is a 2x2 matrix with 'ABILITY TO EXECUTE' on the vertical axis and 'COMPLETENESS OF VISION' on the horizontal axis. The quadrants are: CHALLENGERS (top-left), LEADERS (top-right), VISIONARIES (bottom-right), and NICHE PLAYERS (bottom-left). Several companies are plotted: Microsoft (LEADERS), Tableau (LEADERS), Oracle (VISIONARIES), ThoughtSpot (VISIONARIES), TIBCO Software (VISIONARIES), SAP (VISIONARIES), SAS (VISIONARIES), Yellowfin (VISIONARIES), IBM (VISIONARIES), Domo (CHALLENGERS), Google (Looker) (CHALLENGERS), MicroStrategy (CHALLENGERS), Amazon Web Services (NICHE PLAYERS), Alibaba Cloud (NICHE PLAYERS), Pyramid Analytics (NICHE PLAYERS), Board (NICHE PLAYERS), Inform (Information Builders) (NICHE PLAYERS). A red 'I' logo is in the top right corner of the slide area. On the right is a video frame of a man with short brown hair, wearing a dark blue button-down shirt, standing against a blue background. He is gesturing with his hands clasped together at waist level. The video frame has a thin black border.

In this lesson, we want to introduce you to a business analytics software tool Power BI, which has a robust set of point and click functions for data exploration and visualization. You may have heard of Power BI because it was developed by the tech giant Microsoft. In recent years, Power BI has become a leader in the business analytics space and there are many reasons why that's the case.



Power BI has a wide variety of stunning data visualizations that can be created very quickly. This makes it a great tool for exploratory data analysis

Why Power BI?

- Wide variety of charts
- Charts are easy to create
- Easy to filter and drill-down
- Can add in custom visualizations
- Intuitive

because you can quickly create a variety of multidimensional charts and then filter and drill through the data when you spot a distinctive pattern so that you can zoom in on what is causing that pattern. If you're looking for a chart that isn't included in Power BI, you can easily add in custom visualizations from external sources or create some with R and python.



Power BI is one of the most intuitive business analytic tools that I've used, especially if you're already familiar with other Microsoft products like Word, PowerPoint and especially Excel. The environment of Power BI will be very familiar to those accustomed to other Microsoft products.



For instance, there's a ribbon along the top with different tabs and familiar icons. New tabs appear when you interact with visualizations and data. You can right click on charts and other elements on the screen to bring up additional options. You can perform common tasks in more than one way, which makes it easy for beginners, but also more efficient for advanced users. Saving and opening files, copying and pasting, formatting, creating shapes, importing data and many other such tasks are all completed in a way that is very similar to how you would do in Excel. Not only is the environment familiar, but many of the common data assembly tasks have also been made available with the click of a button

using the Power Query editor, which is built into Power BI. The data transformation tools available in the Power Query editor are just as important as the data visualization tools. However, it's harder to show data transformation tasks than it is to show a chart. So the visualization tools are probably more memorable to most people than the Power Query editor.

Power Query Functions

Easy to get a quick overview
of the data

Rename, delete, relocate columns

Create calculated columns
from numeric and text data

Evaluate the quality of a
column of data



Some of the tasks that you can perform with the Power Query editor include getting a quick overview of the rows and columns, as well as visually exploring the data as you would in Excel, renaming columns, as well as deleting and relocating them, or creating calculated columns based on mathematical combinations of values and numeric columns, or from sub strings of text columns. Quickly evaluating the quality of a column of data to see how many rows have missing values or errors, and then replacing missing values and errors with other values.

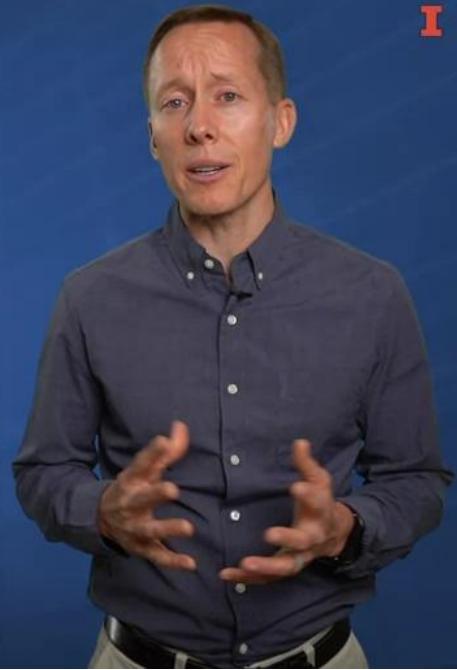
Power Query Functions

Sort rows

Pivot the shape

Group and summarize

Join datasets



Sorting the rows of data based on values in a column. Pivoting the shape of the data from wide to long or from long to wide. Grouping the data and then calculating summary tables. Joining two data sets together based on a primary key.



There are many other tasks that you can easily perform. We won't have time to go over them, but with a little effort, you should easily be able to find out how to do almost anything you want.



The screenshot shows the 'Applied Steps' pane of the Power Query editor. The pane lists the following steps:

- Source
- Promoted Headers
- Changed Type
- Renamed Columns
- Changed Type1
- Replaced Errors
- Replaced Errors1
- Replaced Errors2
- Inserted Day Name
- Renamed Columns1
- Added Custom
- Sorted Rows
- Removed Top Rows
- Sorted Rows1
- Filtered Rows
- Merged Queries
- Expanded states
- Renamed Columns2

The step 'Changed Type2' is highlighted with a yellow background.

There's one thing that I really like about the Power Query editor that has been missing from Excel, and that is a record of your steps are saved so that you can remember the process for the future, or easily revert back to an earlier version of the data. In my opinion, this has made Excel difficult to use for data analytics because it's hard to replicate a data manipulation process.

Data Analysis Expressions (DAX)

Used for custom data
transformations



While there are a lot of point and click data visualization and transformation tasks that you can perform, there's also an option to use code if there isn't already a button for it. Data analysis expressions or D-A-X can be used for custom data transformations. As I mentioned before, you can use R or Python to create custom visualizations. However, just as a relatively small proportion of users take advantage of VBA in Excel, I suspect that a relatively small proportion of users take advantage of DAX or D-A-X because there's so much you can do without learning to type out code.



Power BI is also a great way to tell others about data analytic results, but you have to pay for it. If your organization has a subscription, you can publish your reports to a space where others in your organization can access and interact with it. This allows multiple people to analyze a single source of data at the same time, rather than having different versions of the data floating around. Power BI has many other powerful built-in features, such as the ability to send alerts and to apply a limited number of advanced analytic algorithms to make predictions or identify key influencers. There's a lot to like about Power BI.

Why NOT Power BI?

You have to pay for full functionality

Only works on Windows operating system

It's not meant for creating customized models

Harder to deploy automated results



To be fair, I should also mention the weaknesses relative to other commonly used data analytic tools. One weakness relative to open source tools, is that you have to pay for the full functionality of Power BI. While you can do a lot with the free version you have to pay to unlock the most powerful features. Another weakness is that Power BI only works on Windows operating systems. There is a way to use Power BI for those who have Mac or Linux operating systems, but you have to go through a few extra steps, such as partitioning your hard drive or installing some of their software like Virtual Box, which allow you to run the Windows operating system. There's also a cloud based version of Power BI, but it doesn't contain all of the functionality. Even if it did, it's not ideal because you would have to load your data to a cloud based server.



Now, how does Power BI compared to a data analytic language like R or Python? Well, as is often the case with point and click software, there's a trade off in that it's easier to use, but limited in what it can do. For instance, Power BI is not the best tool if you want to create customized models using specific analytic algorithms that have been rigorously tuned and trained. Another limitation relative to a data analytic language is that, I suspect Power BI is limited in its ability to deploy automated results. For example, I suspect that it would be difficult or impossible to set up Power BI to acquire data, send it to Google's image detection model for classifying what is in the picture, and then send the results to some other output platform. Thus, it's probably not the tool of choice for creating customized apps that depend on a lot of sophisticated data analysis. All in all, I would say Power BI is a wonderful tool for allowing a broad range of people in an organization to conduct exploratory data analysis because it's very intuitive to use. But it's not the tool of choice for those who want to get deep into data analytics or who want to deploy automated data analytics solutions.

Lesson 2-2.2 Installing Power BI

Microsoft Power BI Desktop

Software Application

Visualize simple to complex sets of data, organize them and generate professional-looking reports that can be shared via the Power BI platform



Size: 293 MB

Category: Businessapplication

Operating System: Windows

Download Microsoft Power BI Desktop 2.85.985.0
<https://www.softpedia.com/get/Office-tools/Other-Office-Tools/Microsoft-Power-BI-Desktop.shtml>

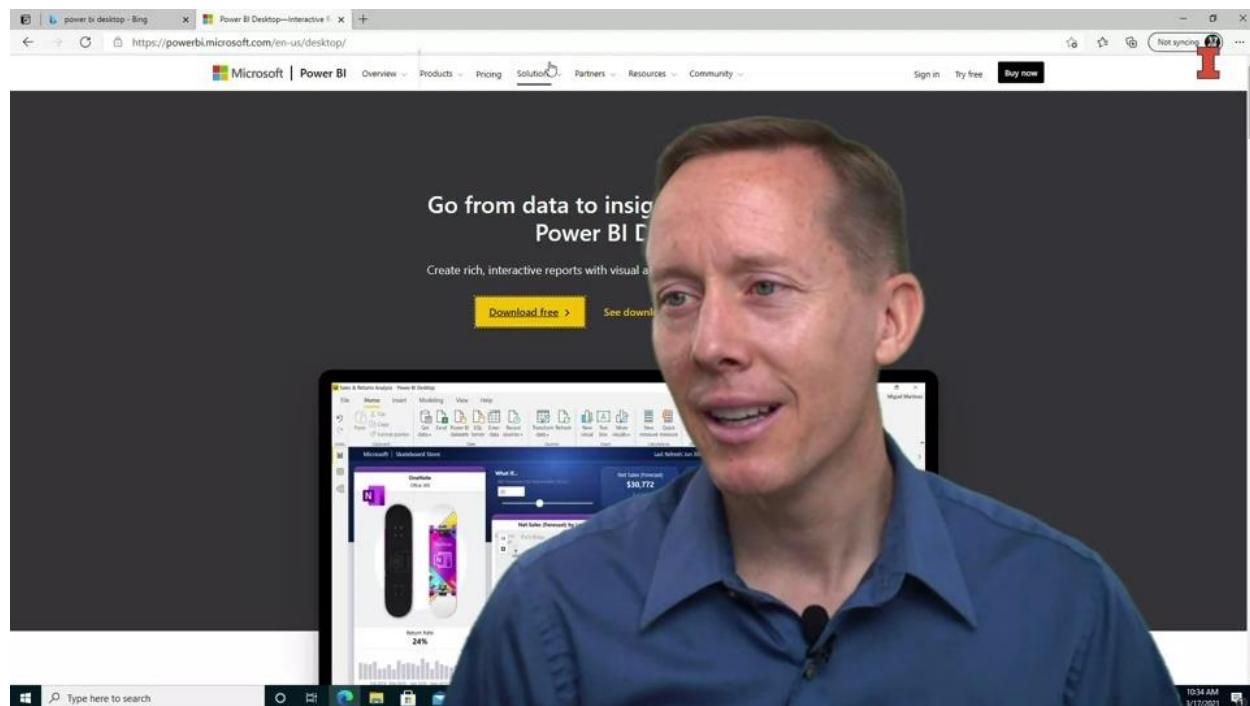
S 3.2/5
softpedia.com

Suggest an edit

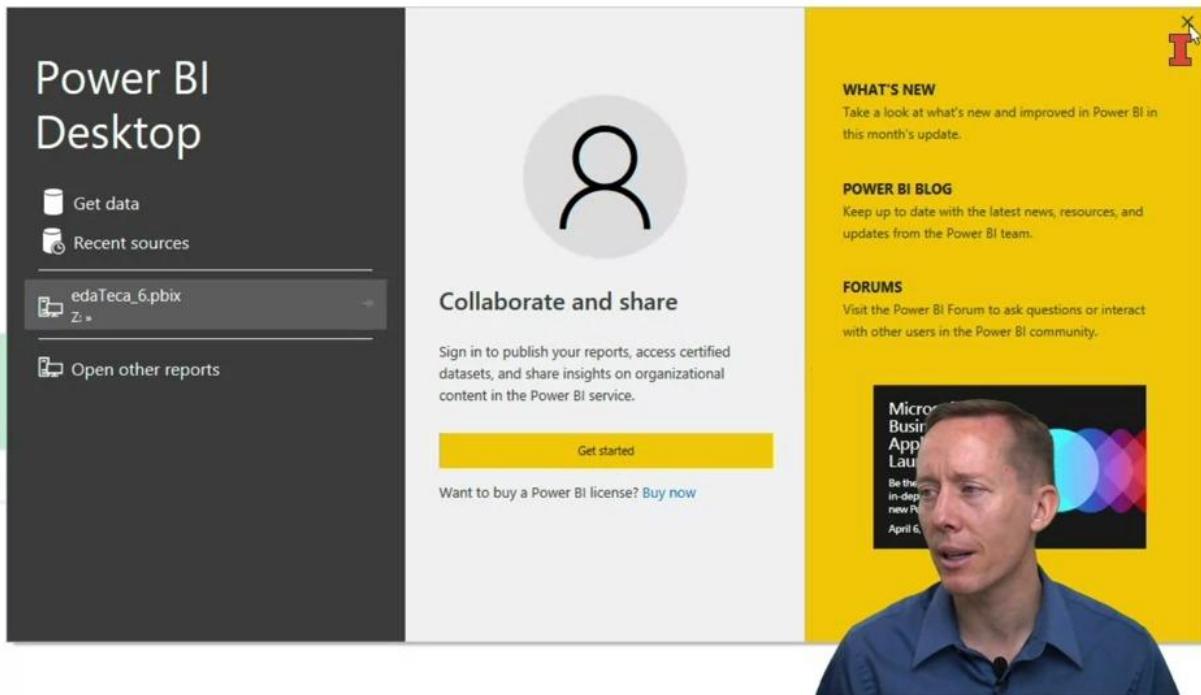


The image shows a video player interface for a Microsoft Power BI Desktop software page. On the left, there's a white box containing product details: title, category, operating system requirements, download link, and a user rating of 3.2/5 from softpedia.com. On the right, there's a video frame showing a man with short brown hair, wearing a blue button-down shirt, speaking directly to the camera. He appears to be in a studio setting with a plain white background.

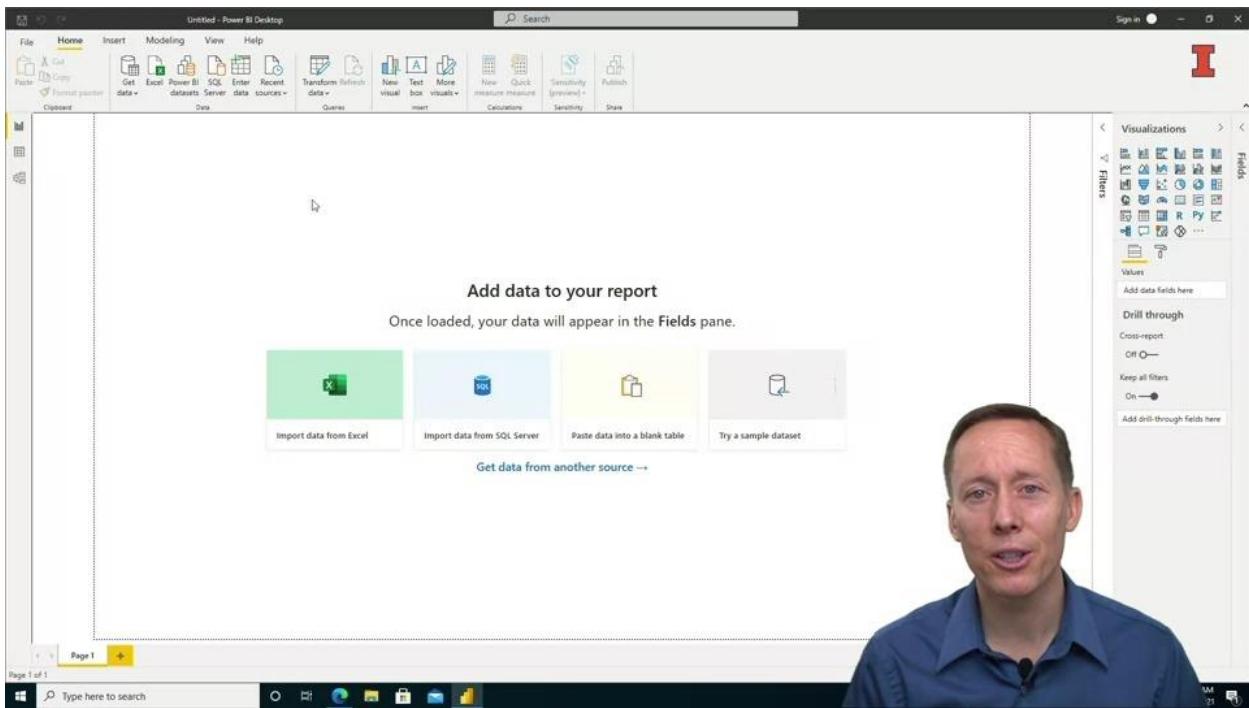
In this video, we want to help you get started using Power BI on your own machine and then give you a high level overview of the Power BI environment. The first thing to recognize is that Power BI only works on Windows operating systems. So if you're using a machine that has a Mac operating system or a Linux operating system, then you'll have to install some other software that allows you to emulate a Windows operating system. And that's what I'm doing.



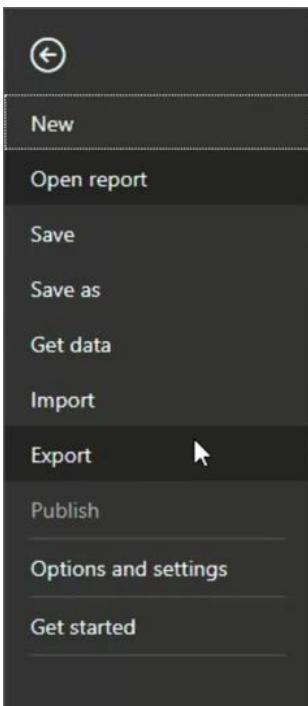
This demonstration is based on a Mac and so I have virtual box running, which allows me to emulate a Windows operating system. All right, once you get that going, go ahead and open a browser and just do a search for Power BI desktop, you don't want the online version it doesn't have as many features. One of the top links should be Power BI desktop and if you click on that it'll take you to a website where you can download Power BI I desktop for free. And if you click on that it will do everything for you. At least it did for me and you can install it and then launch it and there you go. So using power bi desktop is very easy to install on your own machine. There may be some differences for some of you, but I can't anticipate what all of those would be. Now is Power bi desktop really free? Yes, it is. But if you want the full functionality, especially to use it with your organization, you will have to pay for additional features.



Once you've installed Power bi desktop, you can go ahead and launch it. All right, once Power BI I starts up, you'll get this welcome window that intends to help you get started with the most recent project that you're working on and gives you other information about updates. I'm going to go ahead and close that and what I want to do here is just explain to you the overall environment so that you get a feel for how it works.

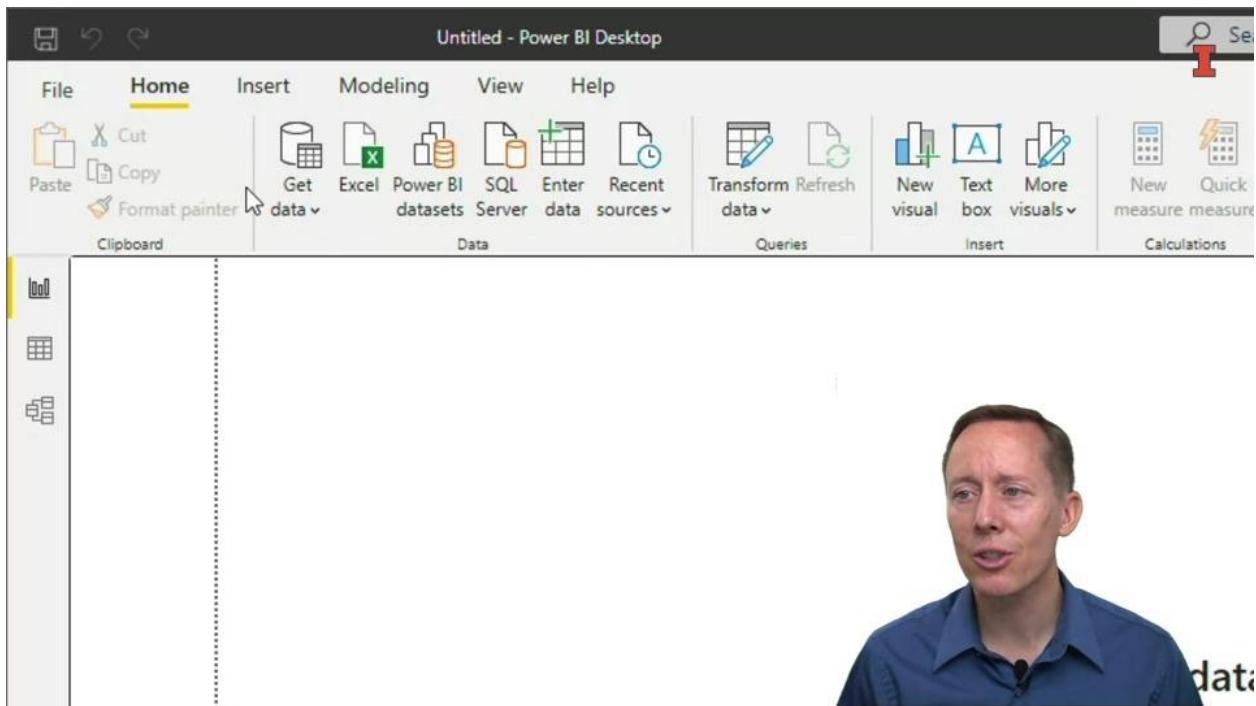


So the first thing to recognize is that this environment is very similar to an environment such as Excel. This main area here is the canvas where you will build visualizations and interact with the data and we'll get to that in a little bit.

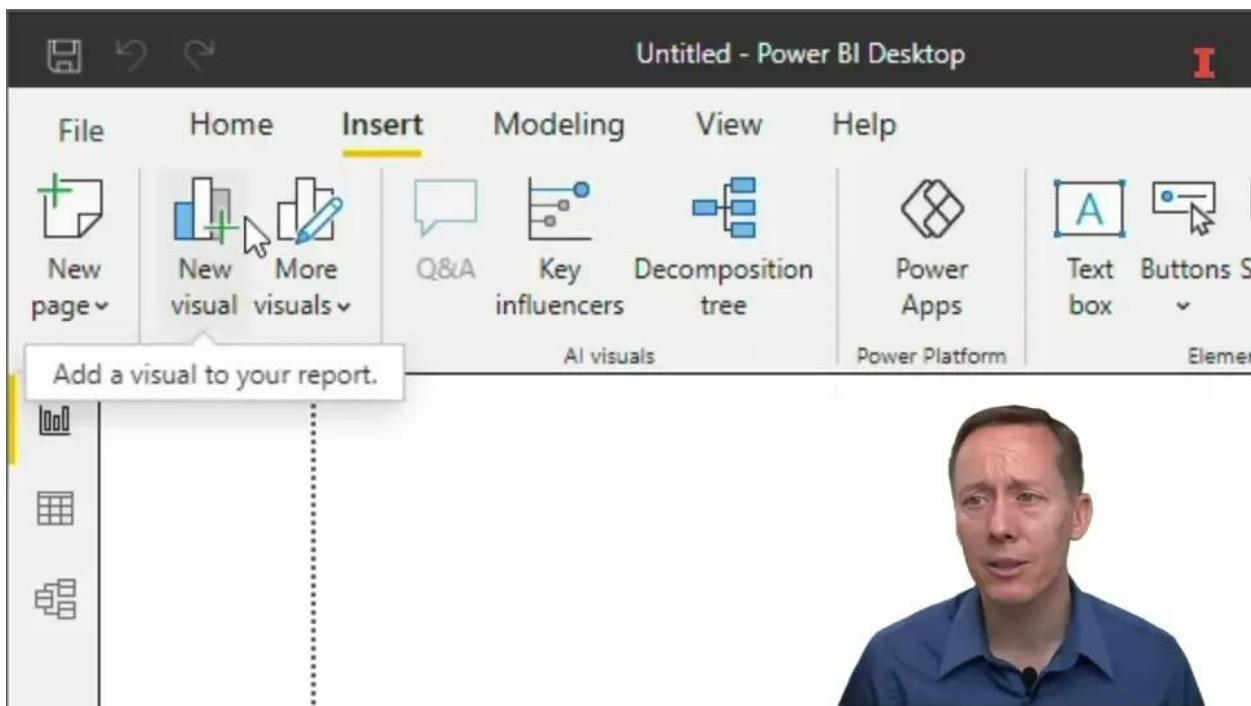


The screenshot shows a Microsoft Power BI application window. On the left, a vertical file menu is open, displaying options such as New, Open report, Save, Save as, Get data, Import, Export (which is highlighted with a cursor), Publish, Options and settings, and Get started. To the right of the menu, the main workspace displays a recent report titled "edaTeca_6.pbix" from the Z: drive, which was opened 2 days ago. A "Browse reports" button is located at the top right of the workspace area. In the bottom right corner of the slide, there is a video thumbnail of a man with short brown hair, wearing a blue shirt, speaking.

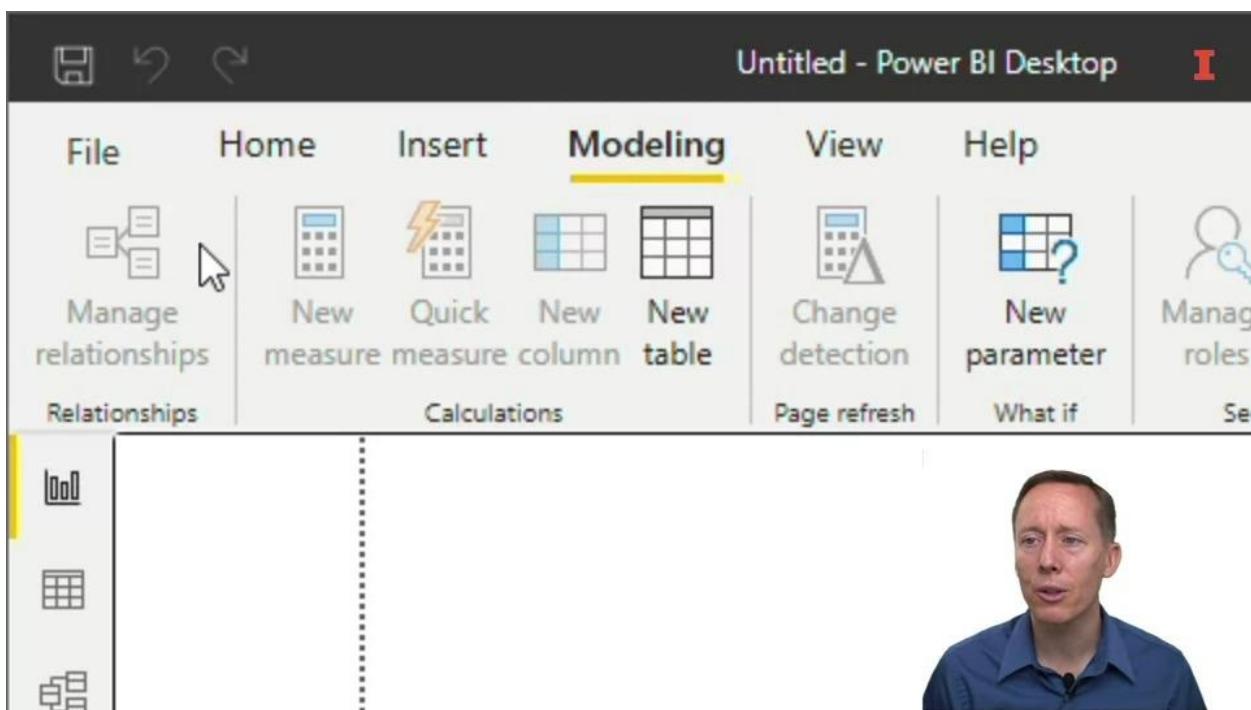
Let's talk about what everything on the side and top are all about. So let's start with this top menu up here. You've got the file menu which if you're a user of other Microsoft products, this will be very familiar where you can open a new report, open an existing report, save your progress, save as, get data, import, export and so on. So you can see here is a report that I've recently been working on.



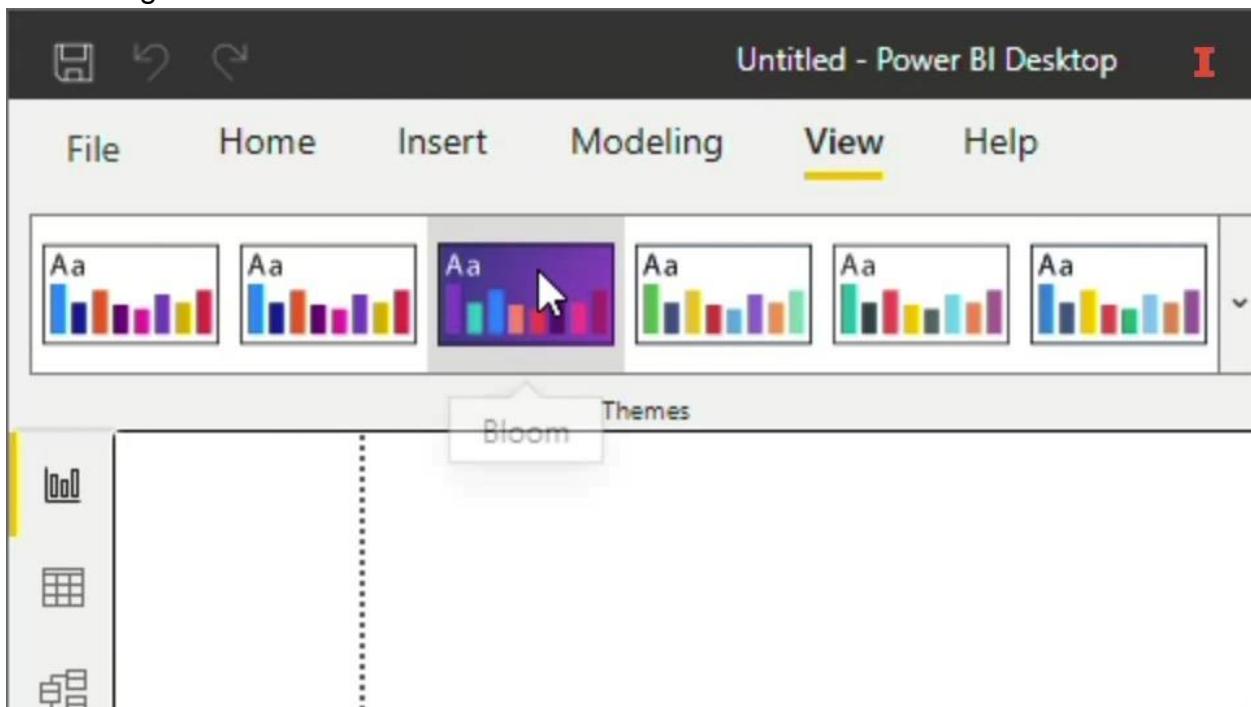
Then you've got the home menu which is kind of the default menu. This has different sections in here for commonly used features in Power BI. So this first section is the clipboard and this again if you're familiar with Microsoft products, this is very familiar. You've got paste, cut, copy and a format painter here. And then you've got a data section. And this section is very useful for accessing data. So you can get data from a CSV file or an Excel file from Power BI data set or a SQL server, so this allows you to connect to a remote SQL server. You could enter data manually and then just recent sources. Then you've got this query section and we will be using this a lot more in other lessons. This allows you to manipulate the data. Then you've got an insert section for inserting some new visualizations and text boxes, some of this you'll be able to use in more detail on this insert tab here. You can also create some quick calculations, some new measures or columns in the data and then a couple of other things, you've got the share here if you want to publish a report online for others to access. And this is one of the features that you have to pay for if you want to allow your whole team to access it.



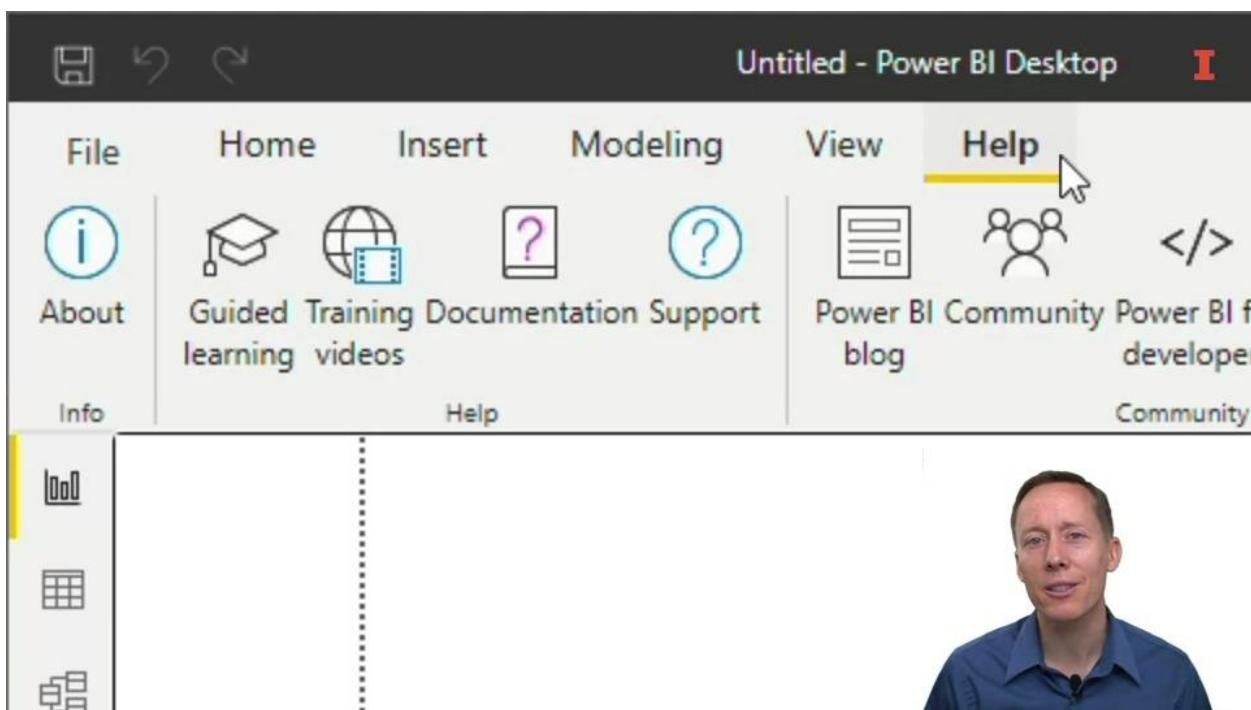
All right then we've got the insert tab at the top and this allows you to insert a new page of your report down here at the bottom or insert new visualizations and a lot of these things you can do in other areas as well. So you can go to the sidebar on the right here and insert visualizations that way. So some of this is redundant, just depends on how you want to use it. And then you have this AI visual section and this is actually really cool, we're not going to go into this. But key influencers allows you to evaluate what is influencing categorical data, whereas the decomposition tree I find is more useful for continuous or numeric data. Then you've got a power app section where you can actually create applications for your business that are based on power BI. And finally you've got this Element section where you can add in a text box or buttons, shapes, images. Again, this should be very similar to what you may have seen in Excel.



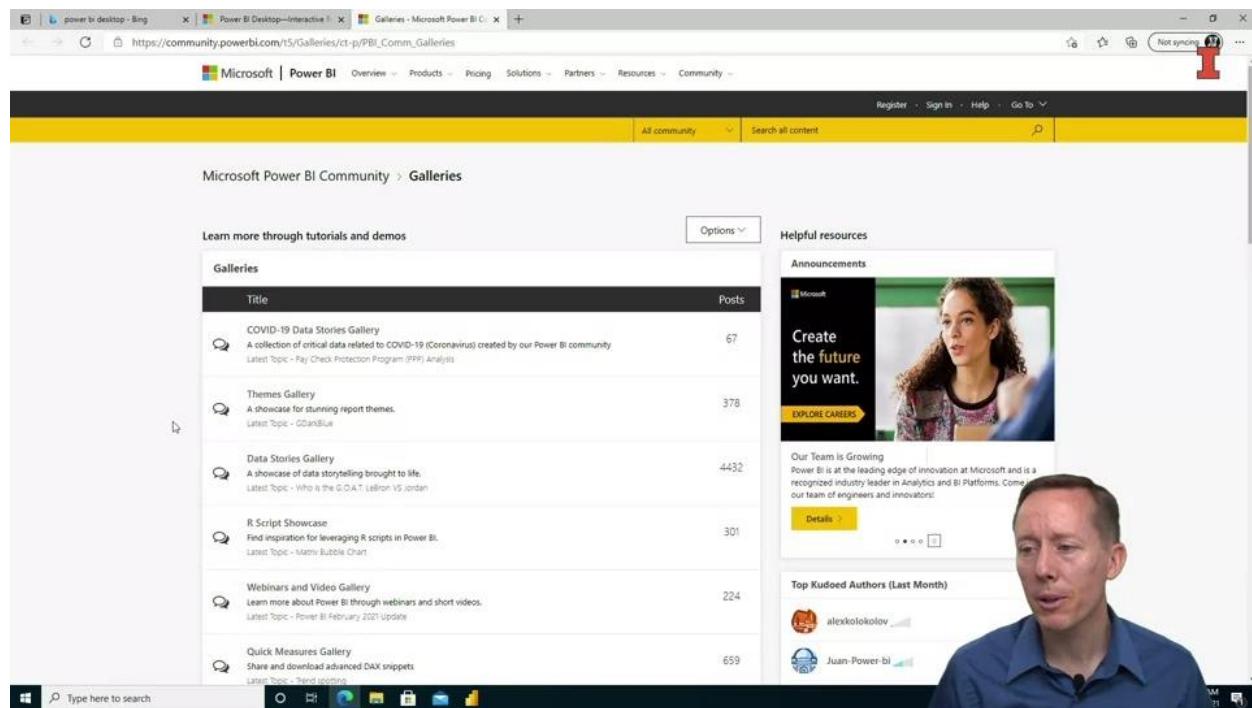
All right, then you've got the Modeling section at the top, which we're not going to use a whole lot, but you can manage relationships here, which you can also access in another area. You can again create calculations or new columns in your data and a variety of other things.



Then we've got the View tab here which allows you to change the aesthetics of what your reports looks like, so the Color scheme, the Layout.



And then one thing that I think is very helpful in any software tool is learning how to use the Help features. So I would take some time exploring what is in the help section here. If you know how to use this, you'll be able to extend your application much beyond what we'll talk about in these lessons. So you can see that Microsoft has done a great job of providing Help. You got Guided learning, Training videos, Documentation, Support, where you can chat in. There's a Power BI blog, there is a large Community of people that use Power BI. And then for developers you can Submit ideas, External tools that you can import and then what I think is super useful are these Examples. So let me click on this and just show you a couple of things. You can use a Sample data set if you want, that is included with Power BI. You can also look at Sample reports and then the thing that I really want to focus on in this lesson here is this community galleries and partner showcase.



Title	Posts
COVID-19 Data Stories Gallery A collection of critical data related to COVID-19 (Coronavirus) created by our Power BI community	67
Themes Gallery A showcase for stunning report themes.	378
Data Stories Gallery A showcase of data storytelling brought to life.	4432
R Script Showcase Find inspiration for leveraging R scripts in Power BI.	301
Webinars and Video Gallery Learn more about Power BI through webinars and short videos.	224
Quick Measures Gallery Share and download advanced DAX snippets.	659

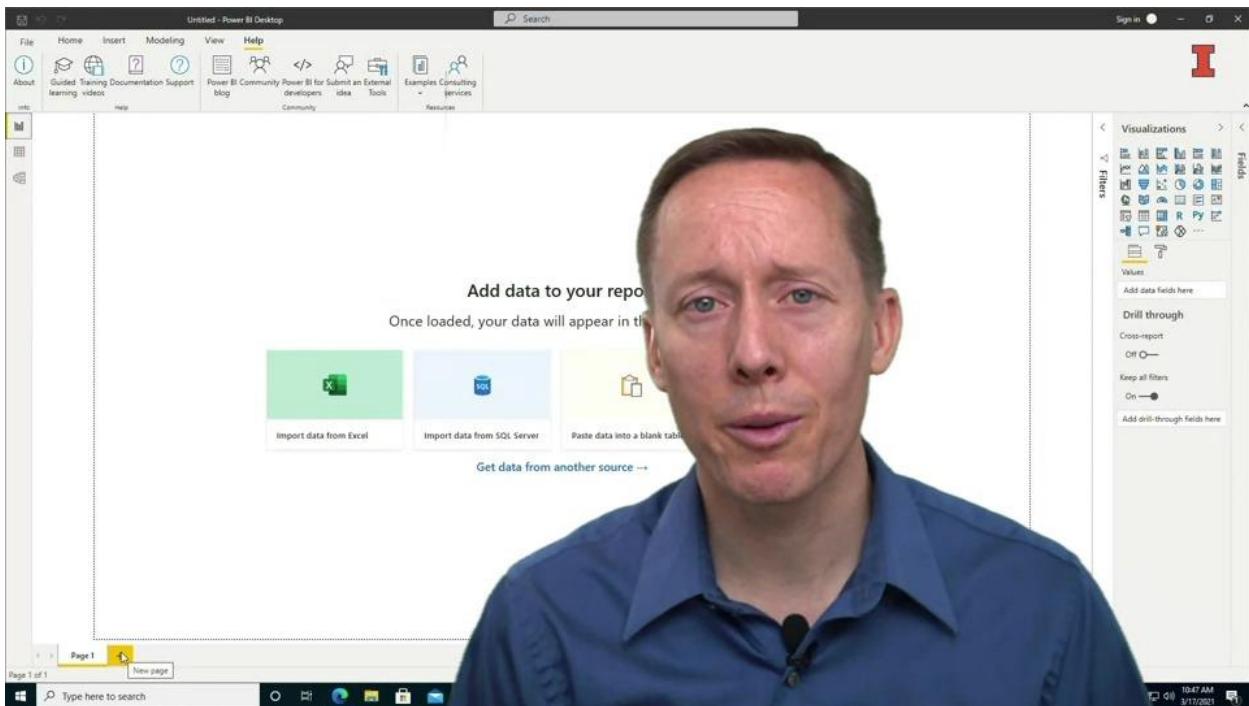
So if I click on community galleries, this will bring up a website with a variety of reports that people have made for a specific theme. So let me click on this top one here, COVID-19 Data Stories Gallery, and this will be a collection of reports, that people have made using the COVID-19 data. I'll go ahead and just click on this first one here and you can see that this report comes up that I can now interact with and explore all online. I don't have to have it on my desktop per se. But I think it's really helpful for getting an idea of what power BI is capable of doing. And can do all these hover overs that interact with the charts and then it gives you an explanation of how to use this dashboard down below.

The screenshot shows a Microsoft Edge browser window displaying the Power BI Partner Showcase at <https://powerbi.microsoft.com/en-us/partner-showcase/?term=Bcountry=UnitedStates&industry=Bdepartment>. The page has a yellow header with the text "See what our partners are doing with Power BI". Below the header, there's a search bar labeled "Search Keywords" and three dropdown menus: "United States", "All Industries", and "All Departments". A main content area features a video of a man speaking, with the text "World-class BI solutions, customized for your business" above it. Below the video, there are six thumbnail cards for partner solutions: Hitachi Consulting - Legal Services Analytics, Agile Analytics - Agile CRM Analytics 365, Agile Analytics - HR Analytics, MANAGILITY - Agriculture and Farming Analytics & Dynamics 365, FreshBI - Cash to Cash Dash, and M.R.E. Business Intelligence. The bottom of the screen shows a Windows taskbar with a search bar and several pinned icons.

The other section in this example that I want to show you is the Partner showcase. So if you click on Partner showcase it will take you to a website where you can explore reports made from businesses that use power BI. So let's go ahead and look at the Explorer Solutions and then from here I can explore Power BI reports from a variety of specific industries. So if you click on this all industries here, you can see that there are many different industries I could choose from. I'll go ahead and choose education and then I can even go into different departments and look at how they're using Power BI. Let me just go ahead and select accounting, and here are some companies that are using Power BI. And I'll click on this Hitachi Consulting and this gives me an overview of the report that they use. I can even click on View Report and interact with their report and read more about it.

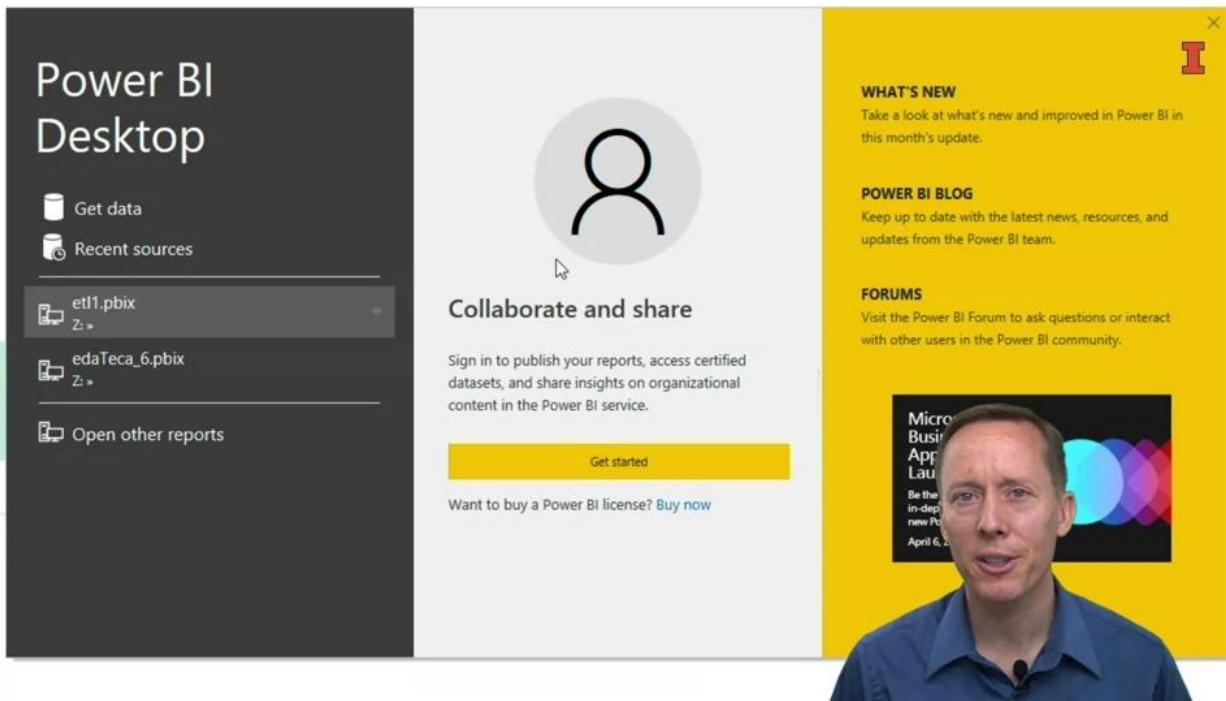
The screenshot shows the Microsoft Power BI desktop application. On the left, there is a video feed of a man in a blue shirt. The main workspace is divided into several sections: 'Visualizations' (containing icons for various chart types like bar, line, pie, and funnel charts), 'Filters' (with a dropdown menu), 'Fields' (with a dropdown menu), and 'Values' (with a section for 'Add data fields here'). A dotted vertical line separates the video feed from the workspace.

All right, let's go ahead and go back to Power BI and there are just a few other elements here that I want to illustrate to you. So on the right side bar we've got the main controls where we'll spend a lot of time for creating visualizations. So there are all sorts of different charts that you can create, Bar charts, Funnel charts, Line charts and many others. Then on the left sidebar we have menu items that allow us to quickly toggle between a Report view, a Data view, or the Data Relationship view. And then at the bottom you can see that you can add different pages either to the Data, Modeling, Relationships, or to the Reports.



So there's a quick overview of the Power BI environment of course, as we get into working with data, all of this will become much more meaningful to you.

Lesson 2-2.3 ETL 1: Examine Data with Power Query Editor

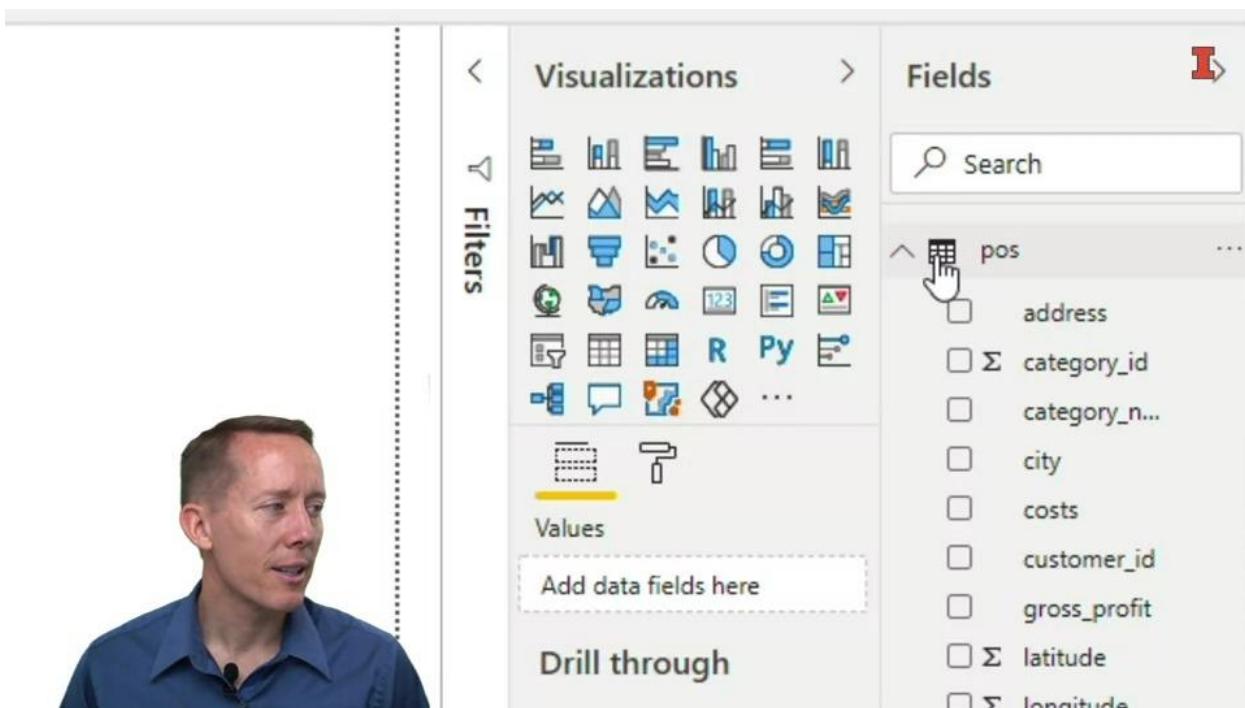


In this video, we want to show you how you can connect Power BI to a data source and then open that day to explore it and start making some transformations with it. So the first thing I need to do is start up Power BI so that I don't have to continue going to this menu and searching for it this way I'm going to pin it to my task bar and then just click on it from there. All right, and once Power BI has started up, I can access some things very quickly from here. But I want to show you how to do this without this. So I won't use that.

The screenshot shows the 'Get Data' wizard in Power BI. At the top, there's a preview of a CSV file with columns: Date, Item ID, and Description. The data includes entries like '7/9/2018 7670.214 456 KICKSTART ORIG UL 16OZ CAN' and '1/6/2017 NA 18953 PrepayFuel'. Below the preview are three buttons: 'Load' (highlighted in yellow), 'Transform Data', and 'Cancel'. A video overlay of a man in a blue shirt is visible on the right side of the screen.

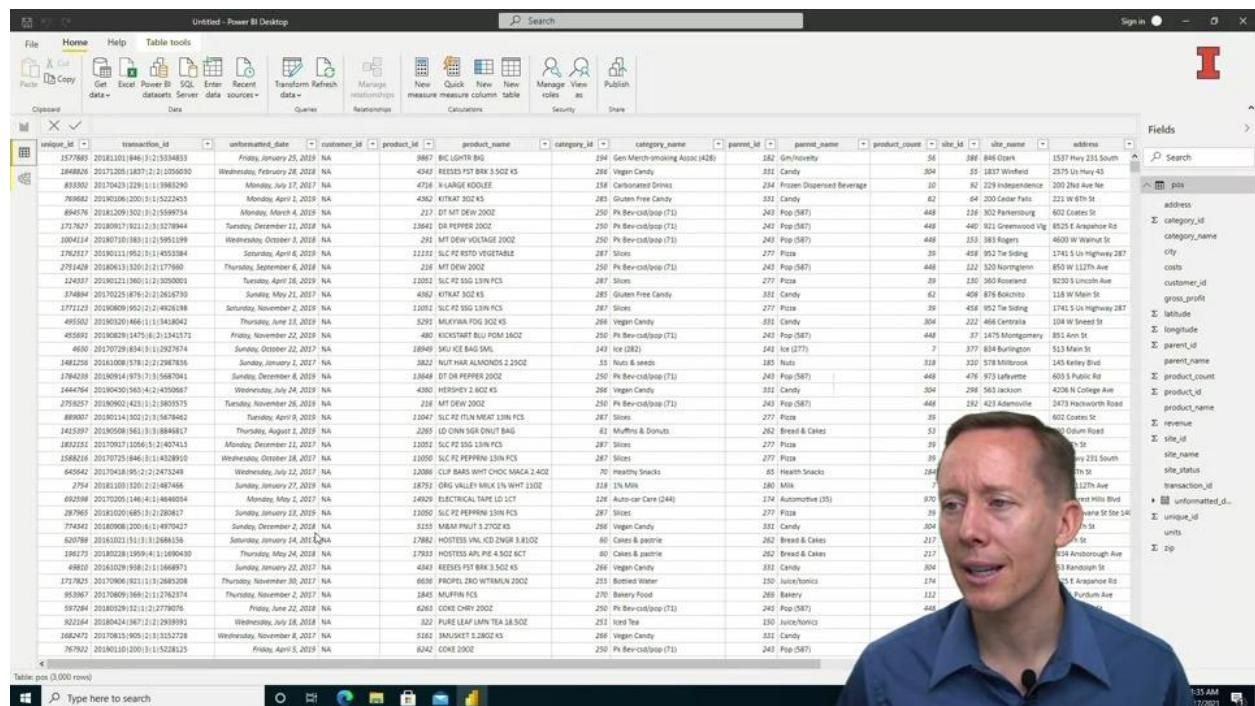
Date	Item ID	Description
7/9/2018	7670.214	456 KICKSTART ORIG UL 16OZ CAN
28/2018	NA	18963 DIESELBIO
14/2017	NA	993 IA QUICK \$50 321
31/2017	NA	5367 KOOLEE
18/2019	NA	1177 IA BNS CRSSWRD 953
23/2018	NA	2738 OLD GLD GLD BX
31/2018	7185.684	10812 SIGNATURE CHKN SAND
13/2019	NA	10059 CO \$50 SERIES II CASH CANNON BNKRL BLAST 856
1/6/2017	NA	18953 PrepayFuel

And so what I want to do is get some data. So I'm going to go to the home tab and click get data. Now, if I knew I wanted to connect to an Excel file, I could click on that Excel file. What type of file are we going to use? We're going to use a file, I'll show you where it is on this machine. I go to this tecaData folder, we're going to connect to this tecaPosSmall.csv file. So if you're following along, you want to download that data to your machine and then you'll point Power BI to the folder where it's located. And so anyway, that's the file that I'm looking for. And since it's a csv file, I want to click on this get data icon. So this wizard pops up and I can see the different options for data sources that I can connect to. And there are a whole bunch of different data sources which is pretty cool. We're doing a very basic one which is text/CSV. So I'll click on that and then connect. And then I will navigate to where this file is located and there it is. TecaPosSmall.csv, click open. And I get this preview of the data and I can see the type of data that it's assuming this is and so we're assuming it's Western European data that's being eliminated with commas. And the data type detection is based on 200 rows, I could change that if I wanted to. We could also go directly to transforming the data from here. But I want to click load so that I can give you a quick overview of where the data is in the general Power BI environment.



The screenshot shows a user interface for data analysis. On the left, there is a video feed of a man speaking. To his right is a panel titled "Visualizations" containing various chart and data representation icons. Below this is a section titled "Values" with a button "Add data fields here". To the right of the "Visualizations" panel is a "Fields" panel. At the top of the "Fields" panel is a search bar labeled "Search". Below the search bar is a list of fields, each with a checkbox and an icon: pos (hand icon), address (square icon), category_id (sum icon), category_n... (square icon), city (square icon), costs (square icon), customer_id (square icon), gross_profit (square icon), latitude (sum icon), and longitude (sum icon). A red edit icon is located at the top right of the "Fields" panel.

And then the first thing you'll notice is that on the far right hand side, you'll see that data set and a list of each column in that data set. So you can see that I've got the address, category_id and so forth. And you'll see that address doesn't have an icon next to it. Category_id has this summation icon next to it, which means it's a numeric field. So I've got several numeric fields. So, if it doesn't have an icon it means it's a text field or a character string. And I've got this field down here that has a calendar icon next to it, meaning it's a date field. Now from here, I could actually click on these three dots and do a lot of other things with this data. But the main thing I want to do from here is rename this by double clicking on it from tecaposSmall to something simpler, just like pos, for point of sale data and click enter. And it does some updates. Now I could connect to another data source and if I did that, I would see that other data source listed here as well.



The screenshot shows the Power BI Desktop interface. On the left, there's a grid view of data rows with columns like transaction_id, product_name, category_name, parent_name, and address. A sidebar on the right lists various fields with their data types and some filtering options. The address field is currently selected.

transaction_id	product_name	category_name	parent_name	address
1577883	BIC LIGHTER BOX 2.5OZ KS	Gm-Merch-smoking Assoc	182 Gm/merch/smoking Assoc	1537 Hwy 231 South
1846803	4342 KIRKLAND SIGNATURE 100CT 1.5OZ KS	Vegan Candy	331 Candy	304 1837 Winfield
8308473	Monstar, July 17, 2018 NA	Carbonated Drinks	234 Frozen Dispensed Beverage	302 229 Independence
7688642	4784 KIRKLAND SIGNATURE 100CT 1.5OZ KS	Gluten Free Candy	303 Candy	42 202 2nd Ave NE
4554579	2018012093021(1):12559734	Monstar, April 10, 2018 NA	237 KIRKLAND SIGNATURE	307 102 2nd Ave NE
1727627	201809071921(1):12987944	Monday, December 11, 2018 NA	231 DA PEPPER 200Z	308 92 229 Independence
10084114	2018071901801(1):15951199	Tuesday, December 11, 2018 NA	230 DA PEPPER 200Z	309 102 2nd Ave NE
1762517	2018051111952(1):14553884	Wednesday, October 3, 2018 NA	230 DA PEPPER 200Z	310 92 229 Independence
2781428	201805111301(1):1177600	Saturday, September 4, 2018 NA	231 MT DEW 200Z	311 102 2nd Ave NE
1234537	2019012211960(1):13000003	Tuesday, April 16, 2019 NA	231 MT DEW 200Z	312 92 229 Independence
3748484	201702251876(1):121265710	Sunday, May 21, 2017 NA	231 SLC PZ SIG 1.5OZ PCS	313 102 2nd Ave NE
1771221	201908091952(1):14021638	Saturday, November 2, 2018 NA	231 SLC PZ SIG 1.5OZ PCS	314 102 2nd Ave NE
495502	201903201466(1):13418042	Thursday, June 13, 2019 NA	232 MULYANA FOOL DOZ KS	315 Candy
455852	201903201475(1):12114315	Friday, November 22, 2019 NA	230 KICKSTART BLU POM 16OZ	316 Frozen Dispensed Beverage
4620	201707291834(1):12927974	Sunday, October 22, 2017 NA	230 KICKSTART BLU POM 16OZ	317 92 2nd Ave NE
3481258	201610081578(1):12987833	Sunday, January 1, 2017 NA	232 NUT HAR ALMONDS 2.5OZ	318 92 2nd Ave NE
1784203	201909141973(1):15687041	Sunday, December 8, 2019 NA	230 NUT HAR ALMONDS 2.5OZ	319 92 2nd Ave NE
1444764	201904501563(1):14370867	Wednesday, July 24, 2019 NA	230 OAT DR PEPPER 200Z	320 92 2nd Ave NE
2759257	201909021423(1):13803575	Tuesday, November 26, 2019 NA	230 OAT DR PEPPER 200Z	321 92 2nd Ave NE
889007	201911141302(1):15479446	Tuesday, April 9, 2019 NA	230 OAT DR PEPPER 200Z	322 92 2nd Ave NE
1425397	201905081561(1):18848817	Thursday, August 1, 2019 NA	230 OAT DR PEPPER 200Z	323 92 2nd Ave NE
1832151	201709071056(1):12104741	Monday, December 11, 2017 NA	230 OAT DR PEPPER 200Z	324 92 2nd Ave NE
1588216	201707251846(1):14328950	Wednesday, October 18, 2017 NA	230 OAT DR PEPPER 200Z	325 92 2nd Ave NE
645642	201704181950(1):12475248	Wednesday, July 12, 2017 NA	230 OAT DR PEPPER 200Z	326 92 2nd Ave NE
2754	201811011320(1):12475248	Sunday, January 27, 2019 NA	230 OAT DR PEPPER 200Z	327 92 2nd Ave NE
652598	201702051148(1):14640304	Monday, May 1, 2017 NA	1429 ELECTRICAL TAPE LD 1CT	328 92 2nd Ave NE
287985	201812011685(1):12808167	Sunday, January 13, 2019 NA	11050 SLC PZ PLN MEAT 13IN PCS	329 92 2nd Ave NE
774541	20180101200(1):14970427	Sunday, December 2, 2018 NA	5235 MBM PRUIT 2.7OZ KS	330 Candy
620788	20161021151(1):12686158	Saturday, January 14, 2017 NA	17882 HOSTESS VNL ICED ZNGR 3.81OZ	331 Candy
183372	201805281399(1):11690400	Thursday, May 24, 2018 NA	17933 HOSTESS API PIE 3.4OZ 6CT	332 Candy
498102	201610291958(1):11668973	Sunday, January 22, 2017 NA	4043 REESES PST BKR 3.5OZ KS	333 Candy
1717823	201709061821(1):12685202	Thursday, November 30, 2017 NA	6630 PROPEL ZERO WTRMLN 200Z	334 Candy
953867	201708091369(1):12762374	Thursday, November 2, 2017 NA	3845 MUFFIN PCS	335 Candy
597294	20180329131(1):12770076	Friday, June 22, 2018 NA	6263 COKE CHRY 200Z	336 Candy
922144	201804241367(1):12199893	Wednesday, July 18, 2018 NA	322 PURE LIFE LMRN TEA 18.5OZ	337 Candy
1662477	20170815105(1):13155728	Wednesday, November 8, 2017 NA	5183 SMV/SKET 3.28OZ KS	338 Candy
76792	201601101200(1):15228125	Friday, April 5, 2019 NA	6242 COKE CHRY	339 Candy

Now for a quick exploration of the data, I can go over to the left sidebar and click on this grid icon which is the data tab. And from here I can explore the data as if it were a worksheet in Excel. So it's got that very familiar feel to it. I can even click on this down arrow next to a column name and see the number of different unique values in that field. And I can do some filtering and sorting from here.

The screenshot of the Power BI Desktop interface shows a data table with the following columns and some sample data:

id	transaction_id	unformatted_d	customer_id	product_id	product_name	category_id	cat	product_count	site_id	site_name	address
1577883	20181011011846(1:1)15334833	Friday, January 25, 2019 NA	9867	842	BIG LIGHT BURGER 2.5CT KS	294	Gen.Mkt.	58	846 Ozark	1537 Hwy 231 South	
1846603	2017120518717(1:1)2050607	Wednesday, February 6, 2019 NA	4342	842	BIG LIGHT BURGER 2.5CT KS	266	Vegan	334	1337 Winfield	2375 W Hwy 43	
831084	2017010710011(1:1)19863294	Monday, July 17, 2017 NA	4781	842	BIG LIGHT BURGER 2.5CT KS	253	Carbs	30	82 229 Independence	200 2nd Ave NE	
760884	2018090610011(1:1)19863294	Monday, September 6, 2018 NA	4262	842	BIG LIGHT BURGER 2.5CT KS	245	Carbs	40	82 229 Independence	200 2nd Ave NE	
494579	201812091302(1:1)15399734	Monday, April 4, 2019 NA	217	8441	DIA PEPPER 200Z	250	Pw	448	115 302 Greenberg	200 2nd Ave NE	
1727607	201809171921(1:1)12728044	Tuesday, December 11, 2018 NA	28441	8441	DIA PEPPER 200Z	250	Pw	448	440 121 Greenberg	2125 W 4th Street Rd	
10004114	201807101380(1:1)15951199	Wednesday, October 3, 2018 NA	293	571	MIT DEW VOLCANO 200Z	250	Pw	448	235 383 Rogers	4600 W 4th Street Rd	
1782517	201802111952(1:1)15353884	Saturday, April 6, 2019 NA	21151	8441	SUC PZ RSTD VEGETABLE	287	Pw	20	456 192 7m Seling	1741 S Dixie Highway 287	
2751429	20180531130(1:1)177660	Thursday, September 6, 2018 NA	216	571	MIT DEW VOLCANO	260	Pw	448	332 320 Northgate	850 W 112th Ave	
1245317	201803211901(1:1)10000001	Tuesday, April 16, 2019 NA	13001	845	SUC PZ SGD 12IN PCS	287	Snack	29	330 365 Keweenaw	9231 S Lincoln Ave	
3748484	201702251876(1:1)20616760	Sunday, May 21, 2017 NA	4062	845	KITKAT 50Z	285	Grocery	62	456 676 Bekonto	118 W Main St	
1771212	201908091932(1:1)21402138	Saturday, November 2, 2019 NA	13001	845	SUC PZ SGD 12IN PCS	287	Slices	29	458 192 Tie Seling	1741 S Dixie Highway 287	
495502	20190301046(1:1)13418042	Thursday, June 13, 2019 NA	5293	845	MULVANIC FOZ DOZ K3	264	Vegan	304	222 466 Centralia	154 W 6th Street	
455867	201902091475(1:1)21145371	Friday, November 22, 2019 NA	480	845	KICKSTART BLU POM 16OZ	250	Pw	448	37 1475 Montgomery	895 Ann St	
4630	201707291834(1:1)20297974	Sunday, October 22, 2017 NA	18803	845	SKU ICE BAG 5MIL	241	Ice (287)	7	377 834 Burlington	513 Main St	
3481254	201610081578(1:1)21298783	Sunday, January 1, 2017 NA	3822	845	NUT HAR ALMONDS 2.5OZ	33	Achts & Snacks	318	322 378 Millbrook	145 Kelley Blvd	
1784230	201909141973(1:1)15687041	Sunday, December 8, 2019 NA	33648	845	DT DR PEPPER 200Z	250	Pw	448	476 973 Lafayette	6005 S Public Rd	
1444764	201904301563(1:1)14350687	Wednesday, July 24, 2019 NA	4580	845	HERSHEY 2 BOZ K5	266	Vegan Candy	304	298 563 Jackson	4206 N College Ave	
2759257	201909021433(1:1)13803575	Tuesday, November 26, 2019 NA	218	571	MIT DEW 200Z	250	Pw	448	192 423 Ademarino	2473 Haciworth Road	
889007	201901141302(1:1)15679446	Tuesday, April 9, 2019 NA	11047	845	SUC PZ PLN MEAT 13IN PCS	287	Slices	29	216 302 Pemberburg	602 Coates St	
1425397	201905081561(1:1)18845817	Thursday, August 1, 2019 NA	2265	571	LD OVN SRN DNLNT BAG	63	Muffins & Donuts	53	296 561 Gardendale	890 Custer Road	
1832251	201709171056(1:1)21047413	Monday, December 11, 2017 NA	11051	845	SUC PZ SGD 12IN PCS	287	Slices	29	2 2056 Rockwell City	520 4th St	
1588216	201707251846(1:1)14328910	Wednesday, October 18, 2017 NA	11050	845	SUC PZ PEPPERN 12IN PCS	287	Slices	39	386 846 Ozark	1537 Hwy 231 South	
645642	201707181950(1:1)212475248	Wednesday, July 12, 2017 NA	12098	571	CLIP BARS WHIT CHOC MACA 2.4OZ	70	Healthy Snacks	184	457 95 Sonoma	211 E 18th St	
2754	201811011320(1:1)1497466	Sunday, January 27, 2019 NA	18751	845	ORD VALLEY MILK 1% WHIT 12OZ	338	1% Milk	7	322 320 Northgate	850 W 112th Ave	
652598	201702051148(1:1)14648034	Monday, May 1, 2017 NA	14929	121	ELECTRICAL TAPE LD 1CT	121	Auto/Car	324	324 148 Belia Vtts	1750 Forest Hills Blvd	
287985	201810201685(1:1)208987	Sunday, January 13, 2019 NA	11050	845	SUC PZ PEPPERN 12IN PCS	287	Slices	325	165 Centennial	7255 S Havana St Ste 14	
774341	201801201200(1:1)14970427	Sunday, December 2, 2018 NA	5155	845	M&M PRNT 3.27OZ K5	287		326	221 47th St	608 47th St	
620798	20181021151(1:1)12688156	Saturday, January 14, 2017 NA	17882	17882	HOTTEST APIK VIN ICO ZINGER 3.81OZ	2703		327	2314 Anthonough Ave	2314 Anthonough Ave	
193372	201802281999(1:1)11490030	Thursday, May 24, 2018 NA	17933	17882	HOTTEST APIK VIN ICO ZINGER 3.81OZ	2703		328	148 Belia Vtts	1750 Forest Hills Blvd	
498102	201610291598(1:1)11668973	Sunday, January 22, 2017 NA	4043	845	REESES PST BXR 1.8OZ	6636	PROPEL ZX2	329	148 Belia Vtts	1750 Forest Hills Blvd	
177782	201709061821(1:1)2085320	Thursday, November 30, 2017 NA	6636	845	REESES PST BXR 1.8OZ	6642	COKE C	330	148 Belia Vtts	1750 Forest Hills Blvd	
953867	201708091369(1:1)12762374	Thursday, November 2, 2017 NA	3845	845	MUPPIN FC	287		331	148 Belia Vtts	1750 Forest Hills Blvd	
597284	20180229132(1:1)21779276	Friday, June 22, 2018 NA	6265	845	COKE C	287		332	148 Belia Vtts	1750 Forest Hills Blvd	
202184	201804241367(1:1)2099993	Wednesday, July 18, 2018 NA	322	845	COKE C	287		333	148 Belia Vtts	1750 Forest Hills Blvd	
1662471	201708151905(1:1)13152728	Wednesday, November 8, 2017 NA	5382	845	IMB C	287		334	148 Belia Vtts	1750 Forest Hills Blvd	
767927	20190110200(1:1)15228625	Fridays, April 5, 2019 NA	6242	845	IMB C	287		335	148 Belia Vtts	1750 Forest Hills Blvd	

Table: pss (3,000 rows) Columns: product_name (1,167 distinct values)

So in my mind this is more of a quick exploration or quick access to some basic tools for manipulating the data. What we want to do is really get into a setting where we can do a lot of transformations with the data.

- Power Query Editor

unique_id	transaction_id	unformatted_date	customer_id
1	2612027	20181219 562 3 2 4909048	3/14/2019 6977.63
2	1281537	20170721 497 1 4 1771950	10/14/2017 1240.056
3	2339438	20170628 302 2 2 4410239	9/21/2017 NA
4	126643	20181015 478 2 1 3513035	1/8/2019 NA
5	74185	20190817 953 2 2 2344053	11/10/2019 NA
6	498571	20171210 446 1 1 1938749	3/5/2018 NA
7	71994	20161023 219 1 1 1500134	1/16/2017 NA
8	2868683	20190517 391 2 1 7433220	8/10/2019 NA
9	2523329	20161027 368 1 100 2845898	1/19/2017 NA
10	2557308	20170304 208 1 1 1959504	5/21/2018 NA
11	111699	20180519 473 1 1 3285930	NA
12	1699367	20180415 914 2 1 3497322	NA

So if we go to the home tab, click on transform data icon, this will bring up the Power Query Editor. Now the Power Query Editor is Microsoft's powerful tool for transforming data. And it brings up this window on top of Power BI and it's got a familiar feel to it, just like all Microsoft products. We've got this hierarchical relationship between the functions in here where we start with this menu icon along the top and then for each menu icon in the ribbon here we've got a grouping of different functions. And then on the left sidebar we've got a list of the different data sets. I've only read in one data set but you can see it's pos, point of sale. And then in the main area here we've got the data that we can explore. And then on the right side we've got the name of the data set. And what I think is super cool is this applied steps area? This will keep track of the steps that we make for transforming the data so that we have a log of what has been done. And we can undo any transformations that we have made to the data. All right, and then along the bottom we've got an overview of the structure of the data set and we also have that in the Power BI environment as well. We can see that I've got 3000 rows and anyway, it tells us some of the same information in the Power Query Editor. Where I've got 23 columns and over 999 rows of data. And then the profiling that we see in here is just based on the top 1000 rows and I'll come back to that in a little bit.

led - Power Query Editor

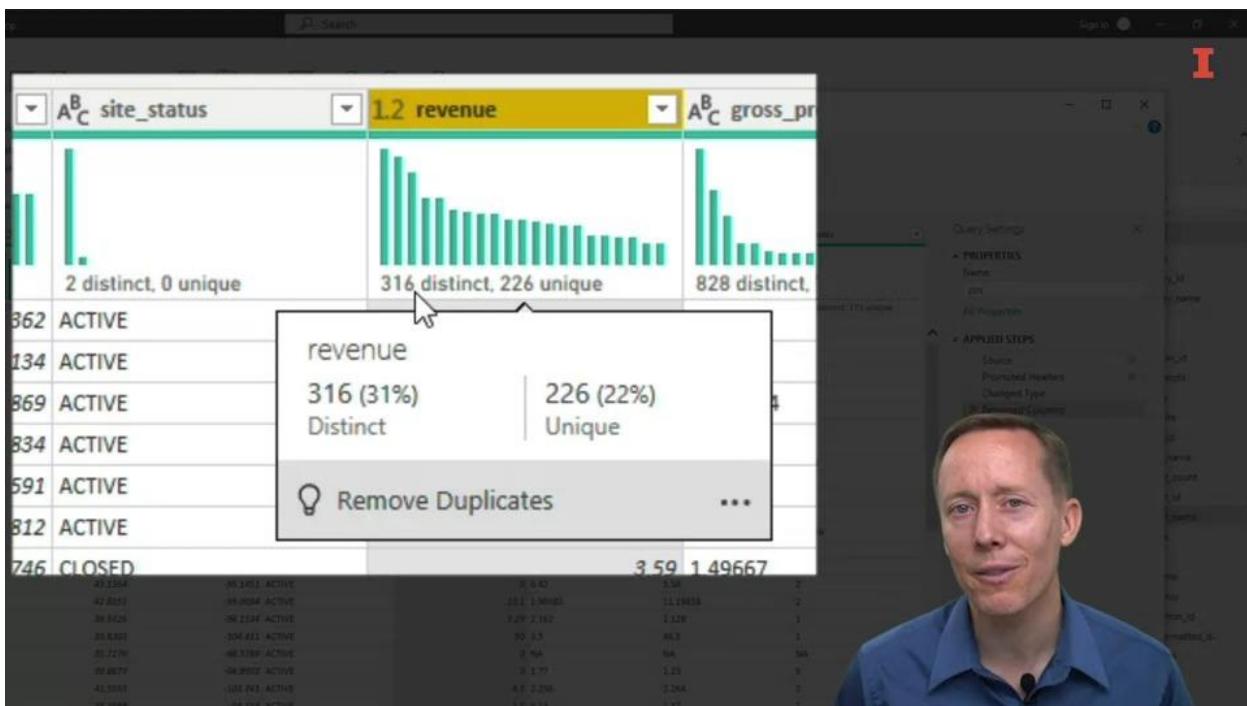
The screenshot shows the Microsoft Power Query Editor interface. At the top, there's a ribbon with tabs: Transform (selected), Add Column, View, Tools, and Help. Below the ribbon is a toolbar with various icons for data manipulation like Transpose, Reverse Rows, Count Rows, Detect Data Type, Rename, Unpivot Columns, Pivot Column, Move, Split Column, Format, Merge Columns, Extract, Parse, and Statistics. The main area displays a table with four columns: unique_id, transaction_id, date, and customer_id. The unique_id column has a numeric value and a small icon indicating it's a whole number. The transaction_id column has a text value starting with '2612027'. The date column has a date value '3/14/2019'. The customer_id column has a numeric value '6977.63'. A portrait of a man is overlaid on the left side of the table.

unique_id	transaction_id	date	customer_id
1	2612027	20181219 562 3 2 4909048	3/14/2019 6977.63
2	1281537	20170721 497 1 4 1771950	10/14/2017 1240.056
3	2339438	20170628 302 2 2 4410239	9/21/2017 NA
4	126643	20181015 473 2 1 3513035	1/8/2019 NA
5	74185	20190817 953 2 2 2344053	11/10/2019 NA
6	498571	20171210 446 1 1 1938749	3/5/2018 NA
7	71994	20161023 219 1 1 1500134	1/16/2017 NA
8	2868683	20190517 391 2 1 7433220	8/10/2019 NA
	523329	20161027 368 1 100 2845898	1/19/2017 NA
	308	20170304 208 1 1 1959504	5/28/2017 NA
	200	20180610 470 1 1 20000000	6/10/2018 544.77

All right, so let's start looking at the data in here. We can see the different columns here I've got unique_id, I've got this icon next to it. And this means it is a numeric field and I can get more information about that if I click on the transform tab and I can see that the data type is actually a whole number. And if I wanted to change this to a different data type, I could simply click on that menu button and select one of these data types that make sense. Then I've got transaction_id. This data type, you can see it's ABC. That means it's a text or character string data type, unformatted date has a calendar icon. So that means it's a date data type and so on and so forth. Now let's say that I don't like the name of a particular column of data. I can simply double click on that and type in a new name and it will update.

1.2 revenue	A _C gross_profit	A _C costs	A _C units
revenue			
1000 (100%)	0 (0%)	0 (0%)	
Valid	Error	Empty	
...			
	3.59 1.49667	2.09333	1
	1.49 0.610297	0.879703	1
	29.9 4.194	25.706	13.98
	8 8	0	1
	0.99 0.36818	0.62182	
	1.99 0.84625	1.14375	
	84.01 8.1327	75.8773	
	2 0.14	1.86	
	1 1	0	
	6 0.42	5.58	
	13.1 1.90182	11.19818	
	3.29 2.162	1.128	

I'll now scroll over to the numeric columns that I am really interested in which are revenue, gross_profit, costs and units. So below each column name, you'll see this green bar here and that tells us some information about the quality of that column of data and I'll come back to that in a minute.



Let me just click on the revenue column here and let's go ahead and go to the view tab and let's click on column distribution. When we do that, you can see we get this really quick overview of the number of distinct values in a column and how frequently they occur. So it doesn't give us a bar for every unique value, but it does give us a bar for the top 20 or so unique values and it tells us that there are 316 distinct values and 226 are unique. What that means is that 226 values show up only one time.

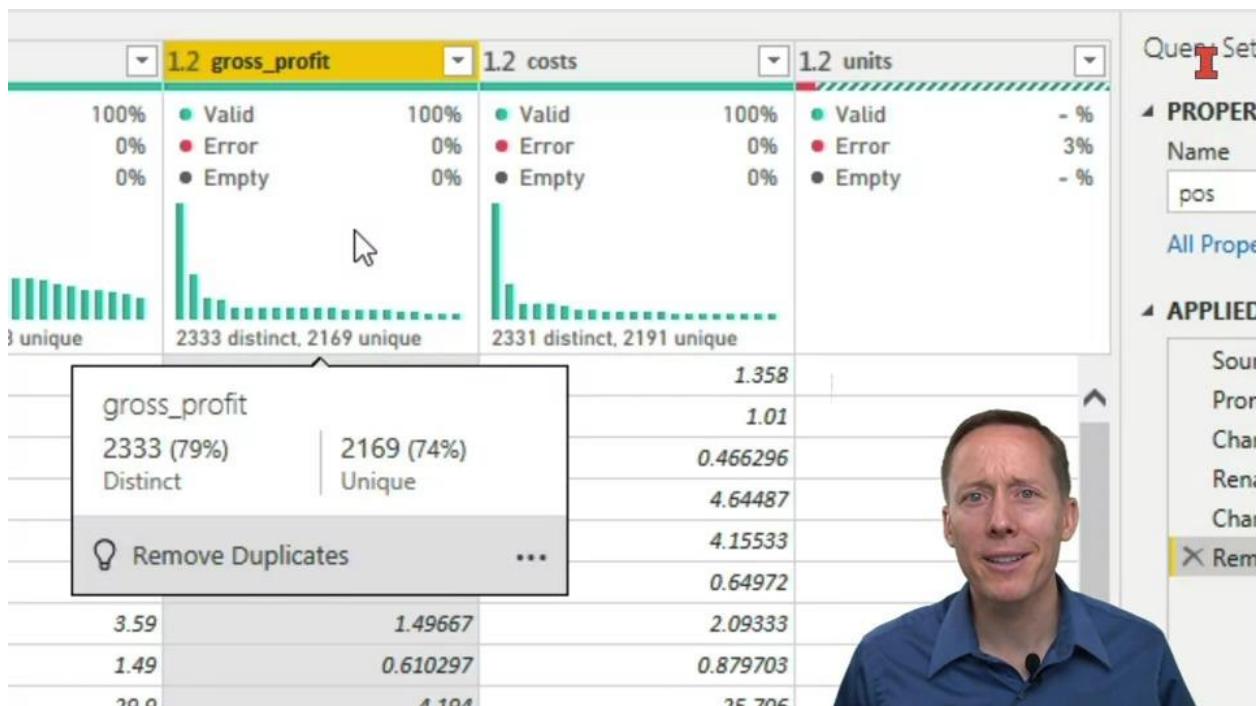


Column profiling based on top 1000 rows

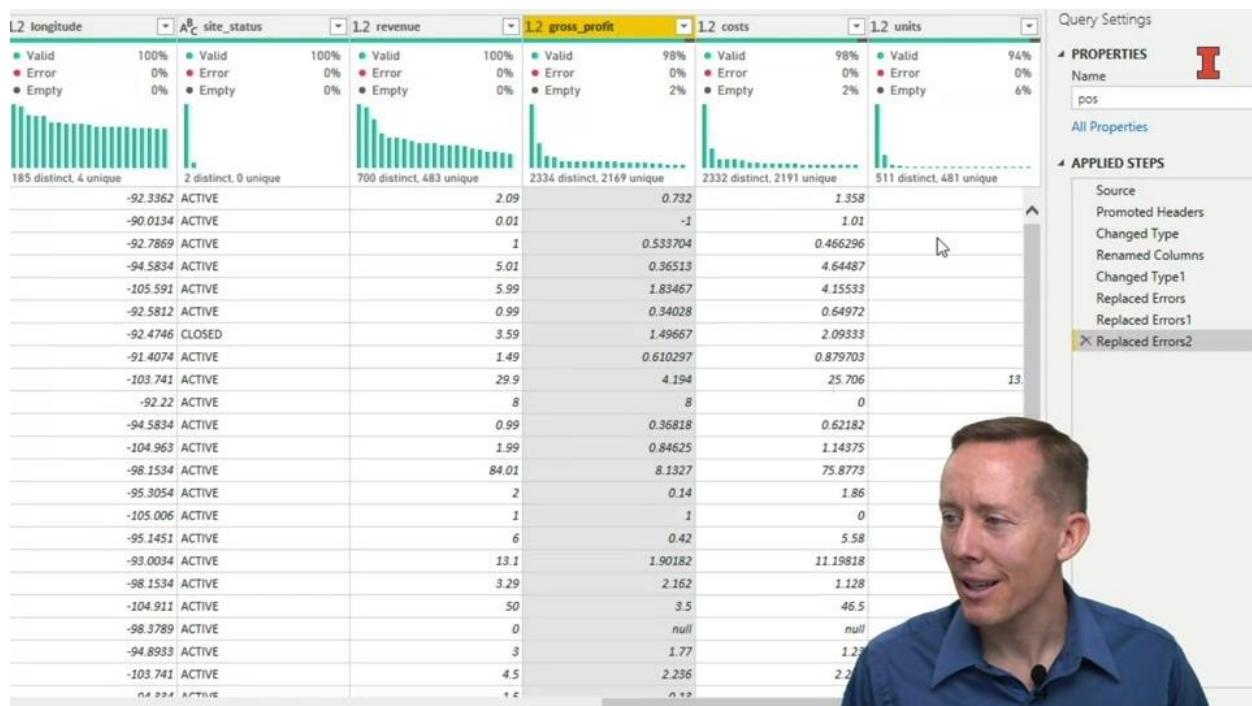
Now let's get a little bit more information and click on the column profile box here. And now down here in this bottom section, I can see more specifically what these values are. So I can see that a dollar and 99 shows up five percent approximately for these 1000 rows of data and how often these other values show up. This is convenient store data. So you would expect a lot of these numbers, these smaller numbers to show up associated with snacks and drinks and so forth. I can also see the column statistics. So I can see the number of rows of data that this is based on. So this is just based on 1000 rows, the number of error entries, empty values, distinct and unique again. NaN, so not a number of values, such as if you divide something by zero. Zero values, so zero shows up 32 times. The min is negative 20 the max is 1123.15, the average and standard deviation. Now this may be sufficient, but if you really want to get an idea for these numeric columns, what the min and max is. And we don't want to just base it off of a sample of 1000 rows.

The screenshot shows the Power BI Desktop interface. In the center is a 'Query Editor' window displaying a table with 3,000 rows and 23 columns. The columns include latitude, longitude, revenue, gross_profit, costs, units, and various date/time and ID fields. The 'Column Tools' ribbon is visible at the top. To the right of the editor is a 'Properties' pane showing the 'Name' is 'pos'. Below it is a 'Applied Steps' pane with a step labeled 'Renamed Columns'. A video overlay of a man in a blue shirt is visible on the right side of the screen.

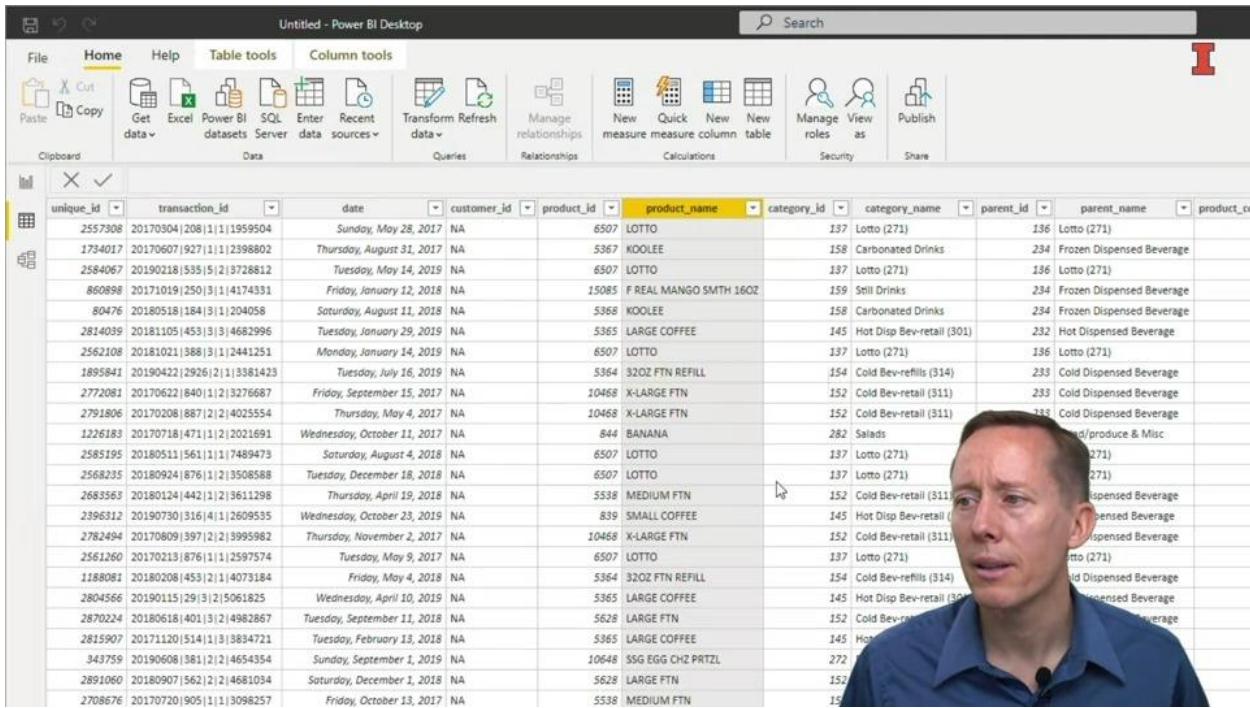
We can go ahead and change this column profiling to be based on the entire data set. And the trade off is that this will now give us more information, but it takes a little bit longer. In this instance here, I only have 3000 versus 1000 rows, it's not going to make much of a difference. If you're dealing with a million rows, then that would be a significant trade off in terms of time that you'd have to deal with. Anyway now I can see that I've actually got 700 distinct values, 483 are unique and the min is negative 100, the max is 1123.15. Alright, from here, let's go ahead and explore the quality of the state of a little more.



And we'll go up to the column quality box and check that and now along the top for each column, you can see the number and with the percentage of valid entries, error entries and empty values in here. Now as I look at revenue, it looks great. What about gross_profit? Gross_profit looks great. However, it is a text data type. So I want to change that. So if I click on gross_profit and go to transform, I can select a data type here. And I want this to be a decimal number. And when I do that it now says that there are 71 errors in here, or two percent are errors. And I could actually change the data type if I just click on that icon there and go to decimal number for costs, units, I will also change this to a decimal number. And now you can see that for all of these. I've got some values which are valid and others which are errors in here. And this bar along the top tries to indicate, send a quick signal about the quality of the colon. Now, what if I want to change? Maybe I just want to remove the errors, I can click on this remove errors button. And when I do that, it will remove all the rows that have errors in this gross_profit column. What if I decide that was a rush decision, I didn't actually want to remove those rows. I want to replace them with something else.



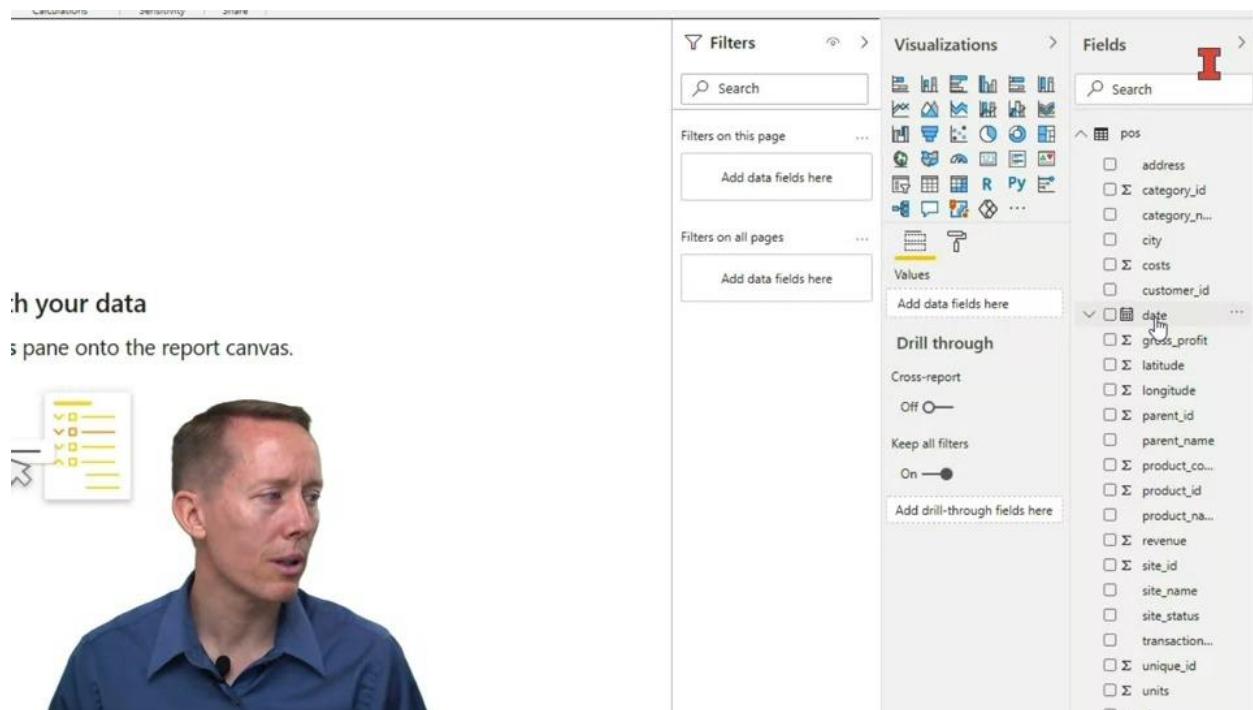
Well I could go over here to these applied steps and remember it's keeping track of everything I do. So I can see that I've renamed columns, I've changed the type and I've removed errors. So I'm going to just unclear or click that undo button and now I've got these errors that are back in the gross_profit column. And if I decided hey, I want to do something else with these missing values, these error values, I want to replace the errors with something else. I can do that, click on replace errors and I want to replace them with null. Now null means it's missing. And you'll see down here that it just replaced that error with a null. So I could go ahead and clean up these other columns as well. I'll go ahead and do that for these remaining columns.



The screenshot shows the Power BI Desktop interface with the 'Home' tab selected. A data grid is displayed, showing a table with columns: unique_id, transaction_id, date, customer_id, product_id, product_name, category_id, category_name, parent_id, parent_name, and product_cd. The data consists of various transactions and products, such as 'LOTTO', 'KOOLEE', and 'MEDIUM FTM'. The 'date' column contains dates like 'Sunday, May 28, 2017' and 'Tuesday, April 10, 2018'. The 'product_name' column includes items like 'F REAL MANGO SMTH 16OZ' and 'LARGE COFFEE'. The 'category_name' column lists categories like 'Carbonated Drinks' and 'Large Coffee'. The 'parent_name' column shows parent categories like 'Lotto' and 'Cold Dispensed Beverage'. The 'product_cd' column contains codes like '137' and '152'. The top ribbon menu shows 'File', 'Home', 'Help', 'Table tools', 'Column tools', and various icons for data import, relationships, calculations, security, and publishing.

All right, now that I've cleaned up those numeric columns, I feel pretty good about where things are at and so what if I want to close my project and apply these changes? I can't simply click on this X box here. What I need to do is go to the home tab here and click close and apply and it will apply those changes to the data and then close out of that and then it will update any visualizations that I may have that are based on that data. And at this point if I am done with my session of Power BI I can save this. As a specific file and I'll save it as etl1 and click save. And I will now have access to this when I want to use Power BI. Next time I can quickly open up and it will preserve these changes. And to close out Power BI I'll click on this X icon at the top. That's hopefully a quick introduction to how you can start using Power BI to making transformations to your dataset.

Lesson 2-2.4 ET2: Dates and Calculated Columns



The screenshot shows the Power BI desktop interface. On the left, there is a small thumbnail of a video player. The main workspace is mostly empty. On the right, there are three side panes: 'Filters' (with sections for 'Filters on this page' and 'Filters on all pages'), 'Visualizations' (with various chart icons), and 'Fields'. The 'Fields' pane is expanded, showing a hierarchical list of fields under the 'pos' category. The 'date' field is currently selected, indicated by a red box around its icon. Other visible fields include address, category_id, category_name, city, costs, customer_id, gross_profit, latitude, longitude, parent_id, parent_name, product_code, product_id, product_name, revenue, site_id, site_name, site_status, transaction_id, unique_id, and units.

In your data

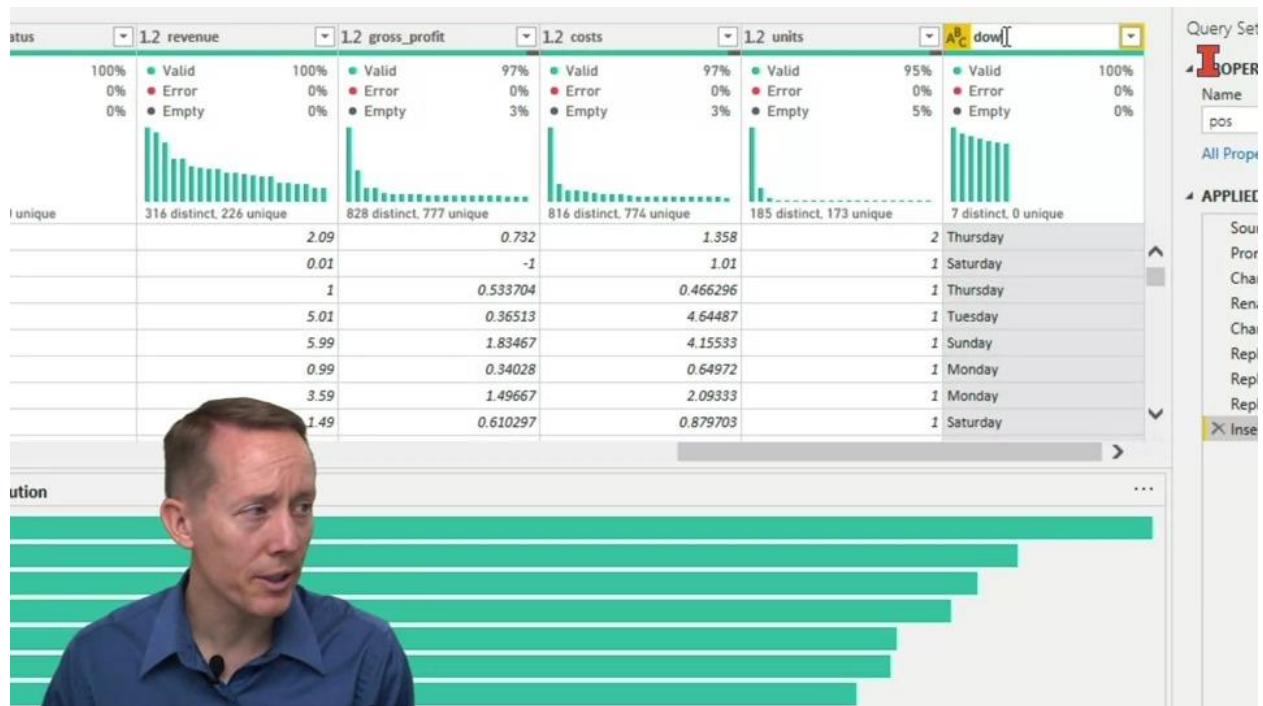
pane onto the report canvas.



In this video, we will focus on dealing with dates and calculated columns in Power BI. Let's go ahead and start up power BI. At this point, let's go ahead and just go with the suggestion to open this recent Power BI file, PBIX file, ETL1. Double click on that, and it will load the data and any other reporting that I have done with this file. Now that the data is loaded, let's just go to the right sidebar and explore this date column here.

The screenshot shows a video player interface on the left with a man in a blue shirt speaking. On the right is the Power BI 'Fields' pane. At the top, there are tabs for 'Filters', 'Visualizations', 'Fields', and a search bar. Below the tabs, there's a search bar and a section for 'Filters on this page'. A 'Date Hier...' node is expanded, showing a hierarchy with 'Year', 'Quarter', 'Month', and 'Day' under it. Other visible fields include 'customer_id', 'gross_profit', and 'latitude'. A cursor is hovering over the 'Day' node.

You'll notice that it has this drop-down arrow next to it. Now, Power BI is really awesome at dealing with dates. It does a lot of things automatically. If I click on that down-arrow, I see this date hierarchy column, and click on that down-arrow and you can see that I've got a year, quarter, month and day that I could use for creating visualizations with my data. Now this is really cool because it's not a new column that has been created, and that's important because it doesn't make our data any bigger by adding new columns. It just is a quick measure in here that it can use as if it were a new column, and then Power BI will make those transformations. If I wanted to create a chart that shows maybe revenue by year, I can do that, or by quarter or by month or by day. Now, this day is not day of the week, it's day of the month. Numbers one through 31.



What if I did want to create another column such as day of the week? I can definitely do that. The way I'll do it is I'll illustrate that by going into the power query editor, by clicking on this transform data icon. Now I can click on the date column here, and if I go to the add column menu item at the top, I can go to the column from example's button. I just want to do it for this column from selection, and now if I double click on this new column here, you can see it's pointing to the date column here, it gives me a list of possible new columns that I might want to create from this date column here. I've got the date in this format 3/14/2019. The month, the year format, the number of days that have passed from a specific date, the day of the month, day of the week, Thursday, and a bunch of other things, maybe the first day of the month, first day of the quarter, just the months name, a lot of different options that I can choose from here. Now I want day of the week, so I'm going to click Thursday, double click that and then click "Okay". You can see I've got this new column day and name, and I've got the day of the week for every observation. I'm going to simplify this and change the state name to dow, make it a little bit simpler.

The screenshot shows the Microsoft Power Query Editor interface. At the top, there's a ribbon with tabs like 'File', 'Transform', 'Add Column', 'View', 'Tools', and 'Help'. Below the ribbon is a toolbar with various icons for data manipulation. The main area displays a preview of a dataset with several columns: 'longitude', 'site_status', 'revenue', 'gross_profit', 'costs', 'units', and 'date'. Each column has a summary bar at the top showing counts for 'Valid', 'Error', and 'Empty' values. Below the preview is a 'Column statistics' table and a 'Value distribution' chart for the 'date' column, which shows the count of distinct days from Monday to Sunday. On the right side, there's a 'Query Settings' pane with a 'PROPERTIES' section where 'Name' is set to 'pos', and an 'APPLIED STEPS' pane listing various data transformation steps. A man in a blue shirt is visible on the right, likely the professor speaking.

You notice that when I click on this here or when I just look at these quick summary information here at the top, I can see that there are seven distinct values, zeros are unique, meaning that all of them show up more than once. Down here at the bottom, I can see that Thursday shows up the most than Saturday than Sunday, and so forth.

This screenshot is nearly identical to the one above, showing the Power Query Editor with the 'pos' dataset. The data preview, column statistics, and value distribution chart for the 'date' column are all present. The man in the blue shirt is again visible on the right, continuing his explanation.

That's pretty awesome. Power query is very efficient at dealing with dates. Now let's illustrate how to create a calculated column based on two or more columns in the data set.



If I look at the different numeric values, a couple of that are relevant here would be revenue and gross profit. As I explore the data, I can see that gross profit is a number that ranges from negative 15-132. This tells me it's not a percentage number. What if I want to create a gross margin or at gross profit as a percent of revenue? I can easily do that.

Custom Column

Add a column that is computed from the other columns.

New column name: gp_margin

Custom column formula: = [gross_profit]/[revenue]

Available columns:

- longitude
- site_status
- revenue
- gross_profit**
- costs
- units
- dow

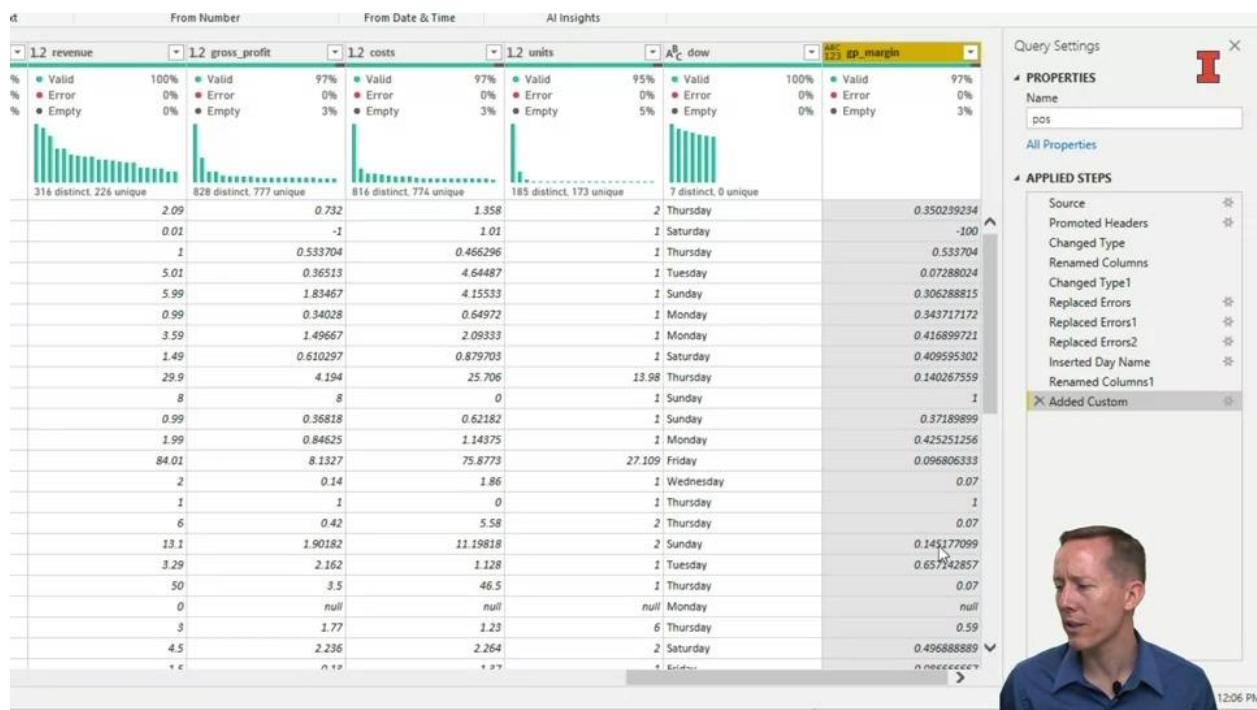
<< Insert

Learn about Power Query formulas

OK Cancel

No syntax errors have been detected.

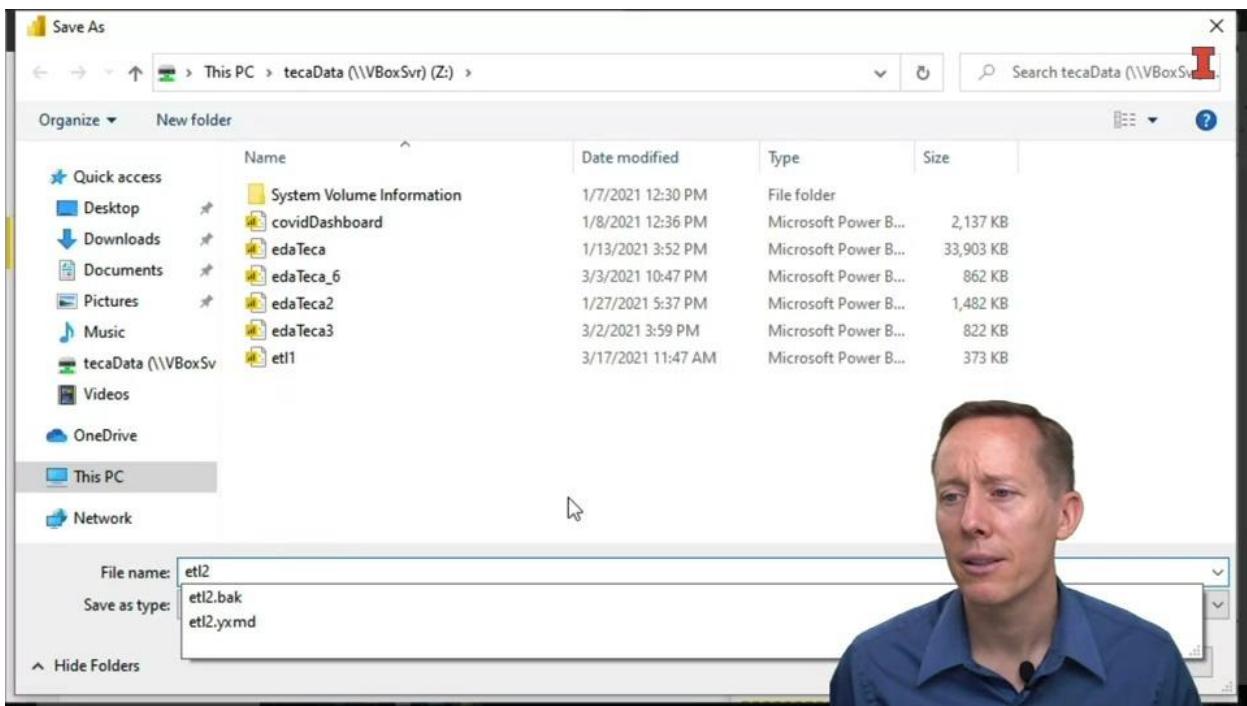
The way I would do that, is I would go in this add column tab, or go to this custom column here and click on that, and it will bring up a wizard that will allow me to type in a formula for creating a new column. Now, first thing I'll do is I'll give this a new column name and I'll label it gp_margin. Then from here, I can go down into this custom column formula box and start typing things. If I just type a letter d, it'll bring up a suggestion of possible functions that I might want to use on a column of data. Now, I don't want to use this function here, but I think that's important to recognize that it has this auto complete feature built into it. What I want to do, is I want to calculate gross profit as a percent of revenue. I'm actually going to go down to the available columns in here, and you can see they're ordered in terms of how they show up in the data set, not alphabetical. I'll double click gross profit, and the way it indicates a column in this formula is it simply put square brackets around it. If I want to take gross profit and divide it by revenue, I can click on the forward slash arrow. At this point you see this red squiggly line and this red line over here, which is indicating that there's an error with this formula. If I were to play this now, it actually won't let me, if I were to do that, it had just result in an error. I could navigate to revenue here. But let me just show you, if I know the column name, I can simply click the open square bracket and I'll lists all the columns in here. If I start typing revenue, it will narrow it down to the column that I'd want. Anyway, there's revenue, gross profit divided by revenue, there's no longer a syntax error. I can go ahead and click "Okay".



There I've got this new column, gp_margin, and I can see that there are 97 percent valid entries, no errors and three percent that are empty. If I want to see the empty values, I could scroll down and look at them and I can see that I've got no here. That's one of the empty values. Why is it null? I might want to understand why. Well, it's because the gross profit value is null here, and so null divide by zero is going to just give us null.



There you go. That's how you can very easily deal with date columns and create calculated columns in power query editor.



I'll go ahead and save this. I'll close and apply this. Once those changes have been applied, I can go over and verify that they're here. If I look at the different columns, I've got now the dow, for day of week and I also have gp_margin here that I just created, which I can now use to create visualizations. Now that I've done that, I will save this file. I could just click save, but I want to save it as a new file, and so I will save it as etl2, to separate what we did from a prior lesson. Click "Save" and then exit out of Power BI.

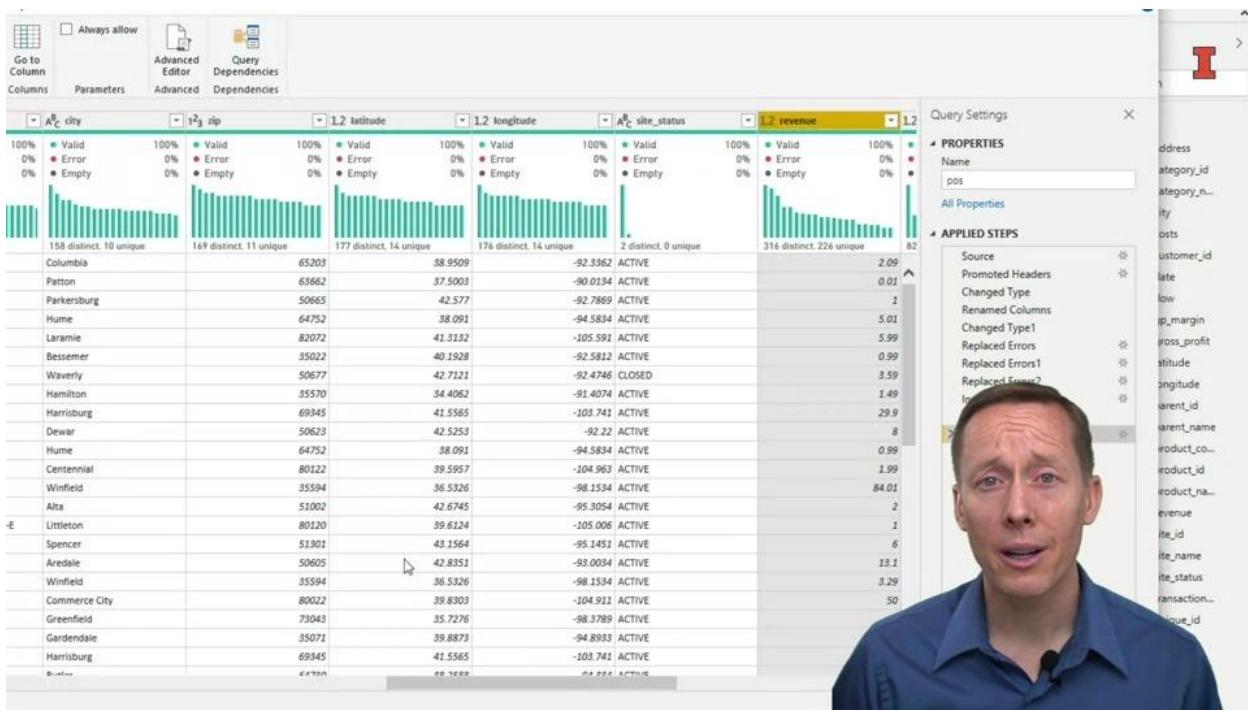


There's an example of how you can deal with date columns and Power BI, as well as create calculated columns.

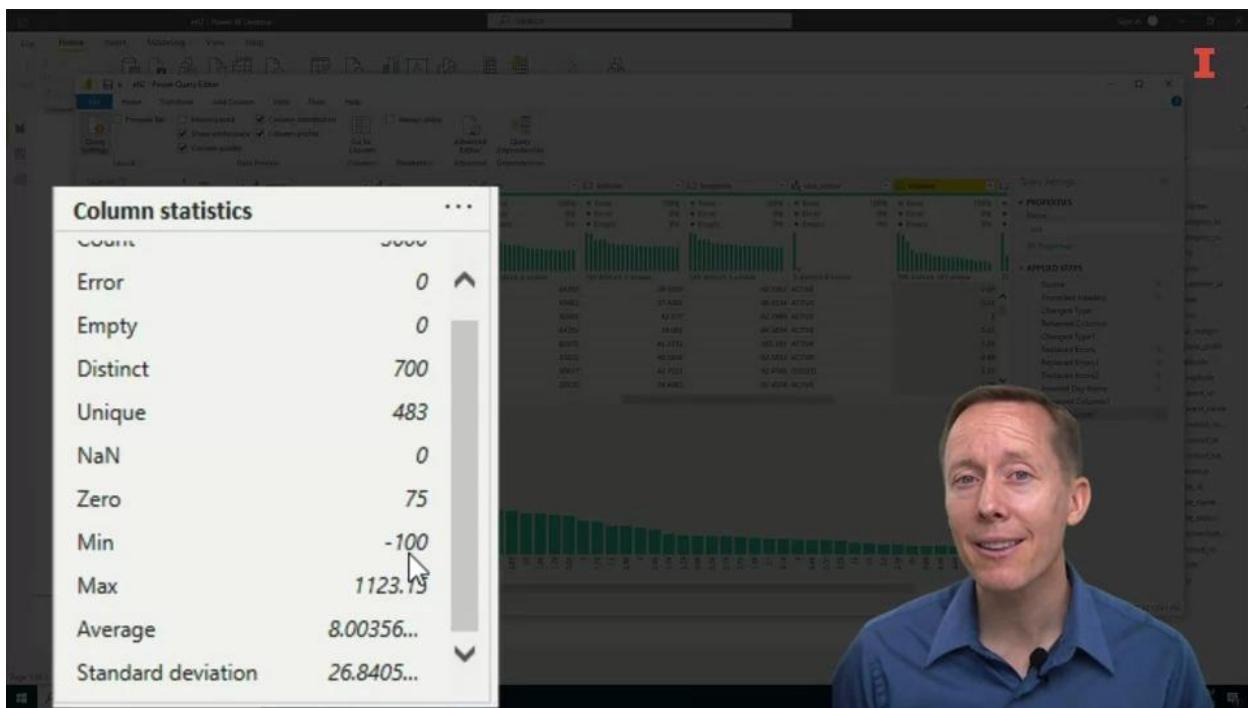
Lesson 2-2.5 ETL 3: Checking for and Eliminating Outliers

The screenshot shows the Power BI Desktop application with the 'etl2 - Power BI Desktop' report open. The main area displays a table with several columns: transaction_id, date, customer_id, product_id, product_name, and categories_id. The 'Applied Steps' pane on the right lists various transformations applied to the data, such as Promoted Headers, Changed Type, Renamed Columns, and Replaced Errors. A man's face is overlaid on the right side of the interface.

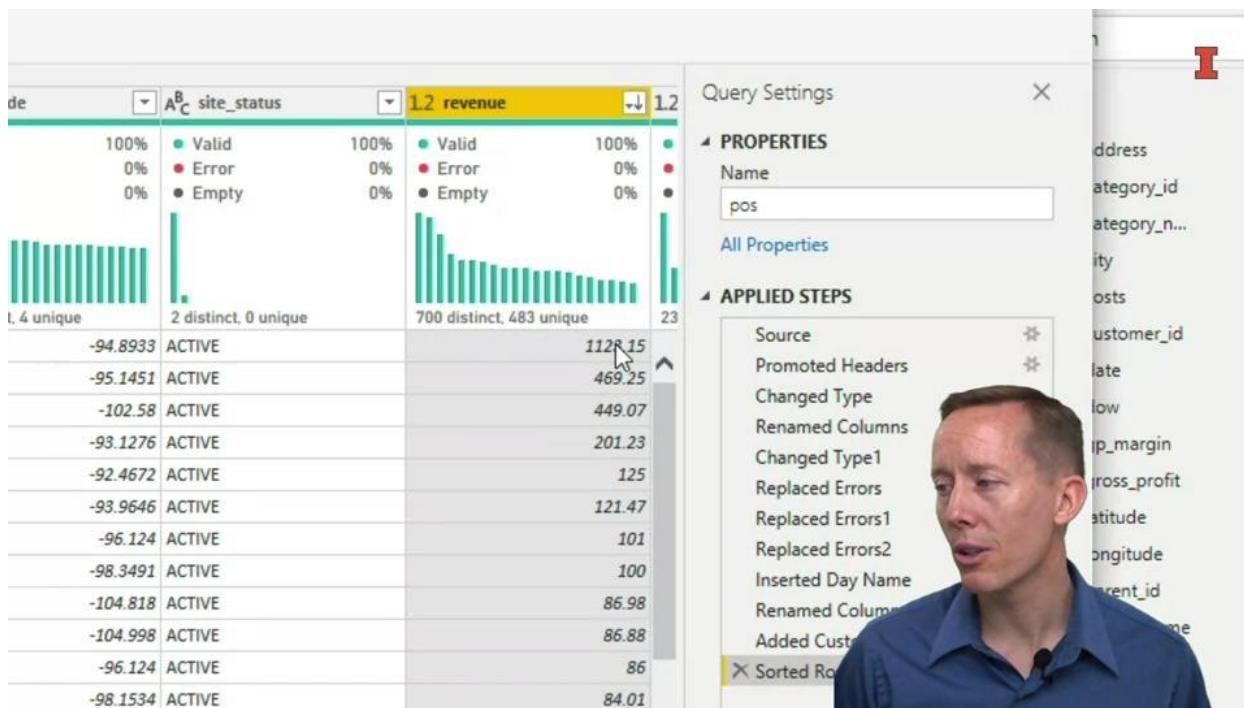
In this video, we want to show you how to check for outliers in both numeric and categorical columns and then remove outliers or deal with them in some other way. The first thing I'll do is once I open up Power BI as I will select the report that I've been working on, etl2, and once it has loaded the data and any visualizations that I may have I will click on the transform data icon, and now we are in the power query editor. Now, I want to focus first on looking at the revenue column because that is a column that we really want to understand more.



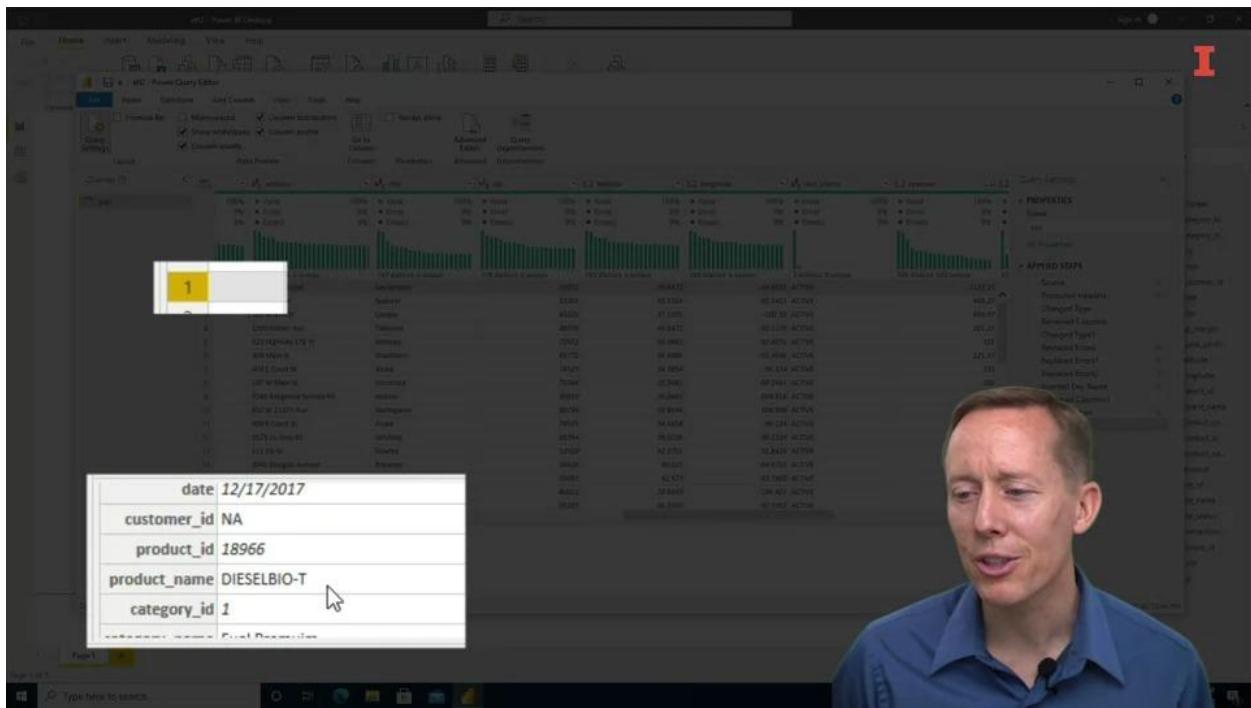
So to find that column, I know that it's on the far right, but I'll show you another thing that you can do. If you've got a wide data set and you don't know where the column is located, if you go into the View tab, you can click on go to column and then you can either scroll through the different columns or just start typing in its name and it will pop up. You can click on it, click "Okay", and it will move over, scroll to where that column is located. Now, let's explore whether or not there are outliers in the revenue column.



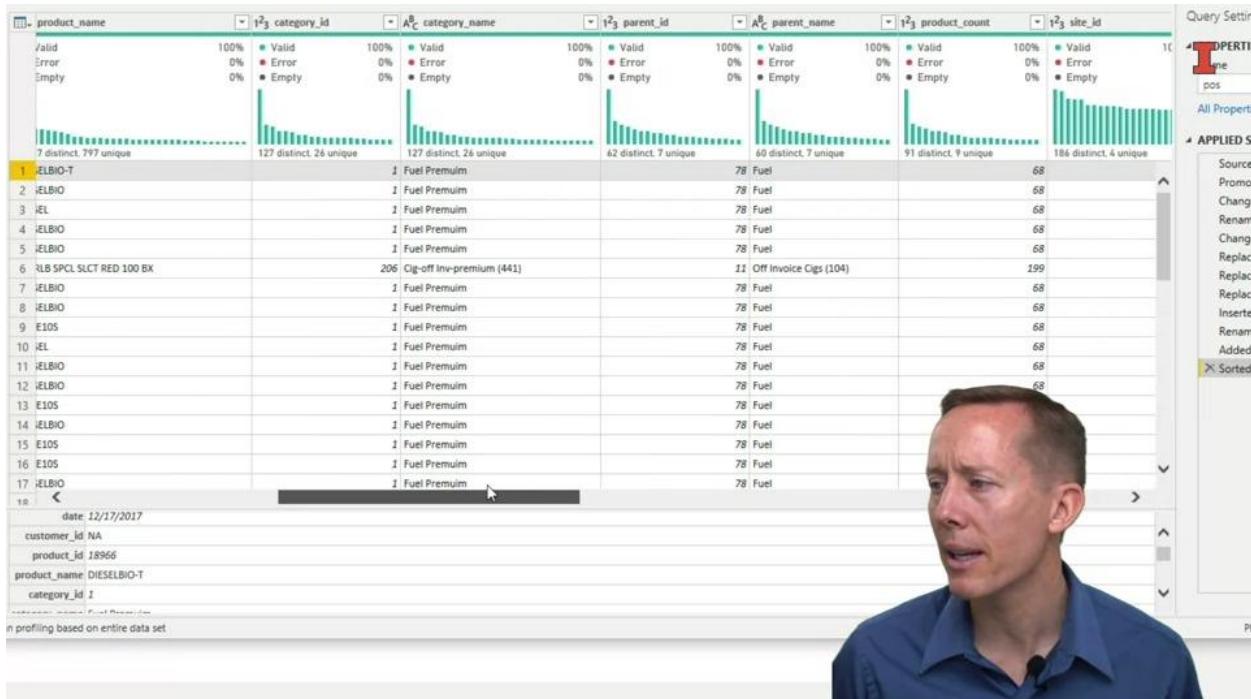
So I'll click on that column header and make sure that I have the column profile selected, and also make sure that I'm profiling this on the entire data set. Now I can look at the Min, max, average, and standard deviation to get an idea of whether or not there are outliers. So let's first focus on the average. The average is eight, the standard deviation is about 26, and if this were a normally distributed column of data then any outliers are typically either above three standard deviations above the mean, or three standard deviations below the mean. So if I take 26 and multiply it by 3, that's about 78 and then add that to the average of 8, that means anything that's above about 86 would be considered an outlier, and if I take that 78 and subtract it from 8, that means anything below a negative 70 would be considered an outlier as well. So now I can compare those two numbers, negative 70 and 86 to the mean and the max and see that indeed there are observations that extend beyond those points.



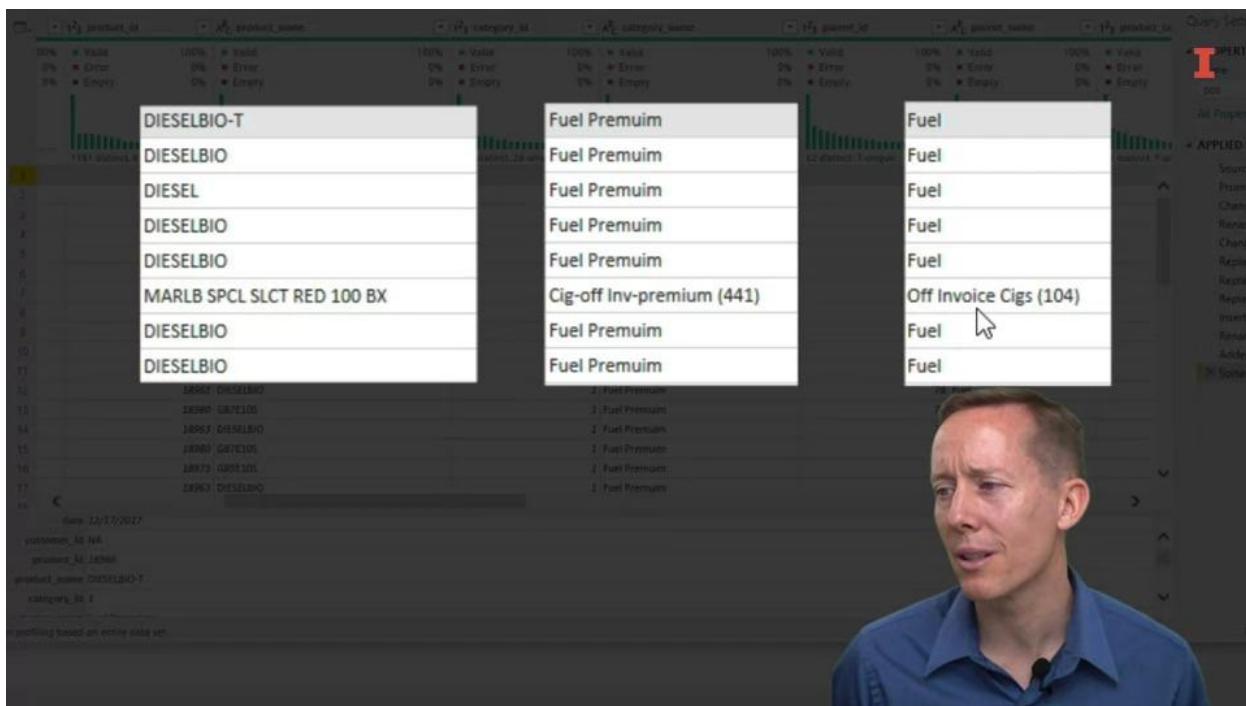
So let's first identify those observations and see if we can figure out what's going on. So I can go to the revenue column, click on the drop down arrow next to the column name and let's sort descending first. Now I can see that the data has been sorted, the rows have been sorted based on the value of revenue with the highest value at the top, and you can visually explore and verify that these are the highest values.



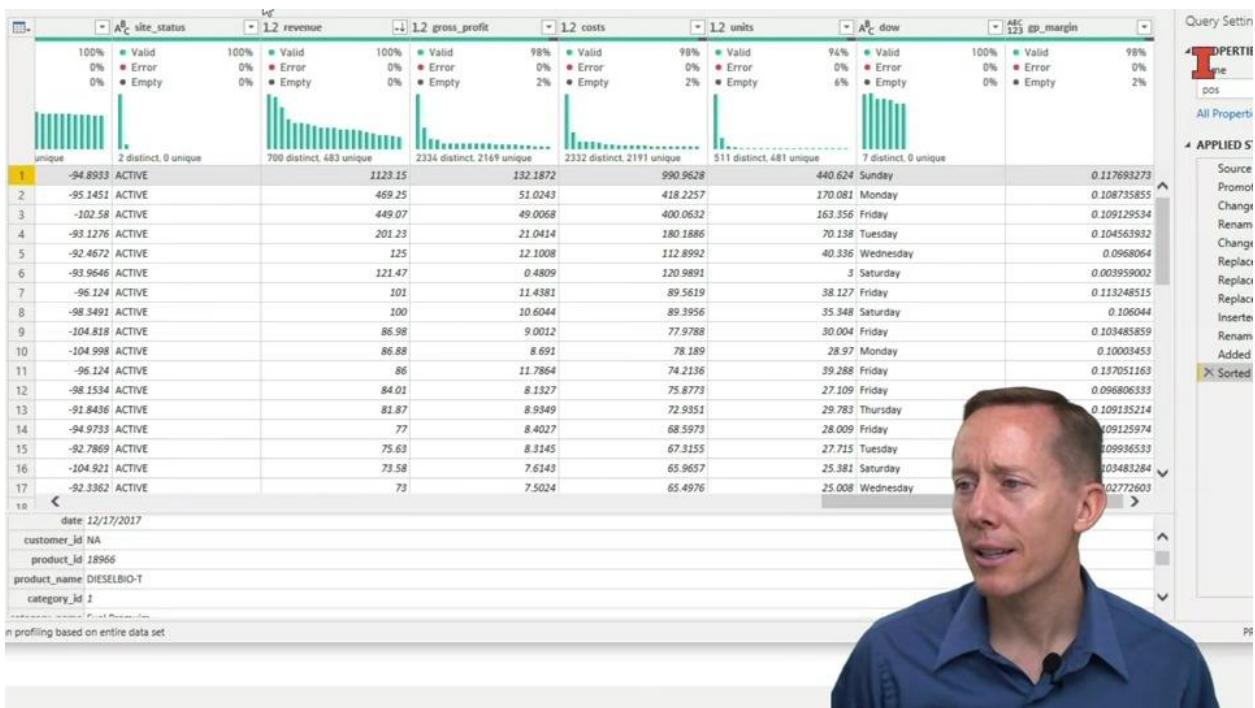
Now, if I click on a column label here, I can get information about the values in the other columns for this observation, and I can see that for this observation, it came from diesel bio, so it was a fuel observation.



Anyway, so that gives me a way to quickly see what this observation is about. If I want to look at all of these high observations, it looks like about the first eight are above 86, we'll say 87. Then I may want to figure out why, what's going on with these? I could just scroll over until I get to the parent name and category name



and product name and basically real quickly see that these are all related to fuel except for one, which is a large purchase of cigarettes, apparently.



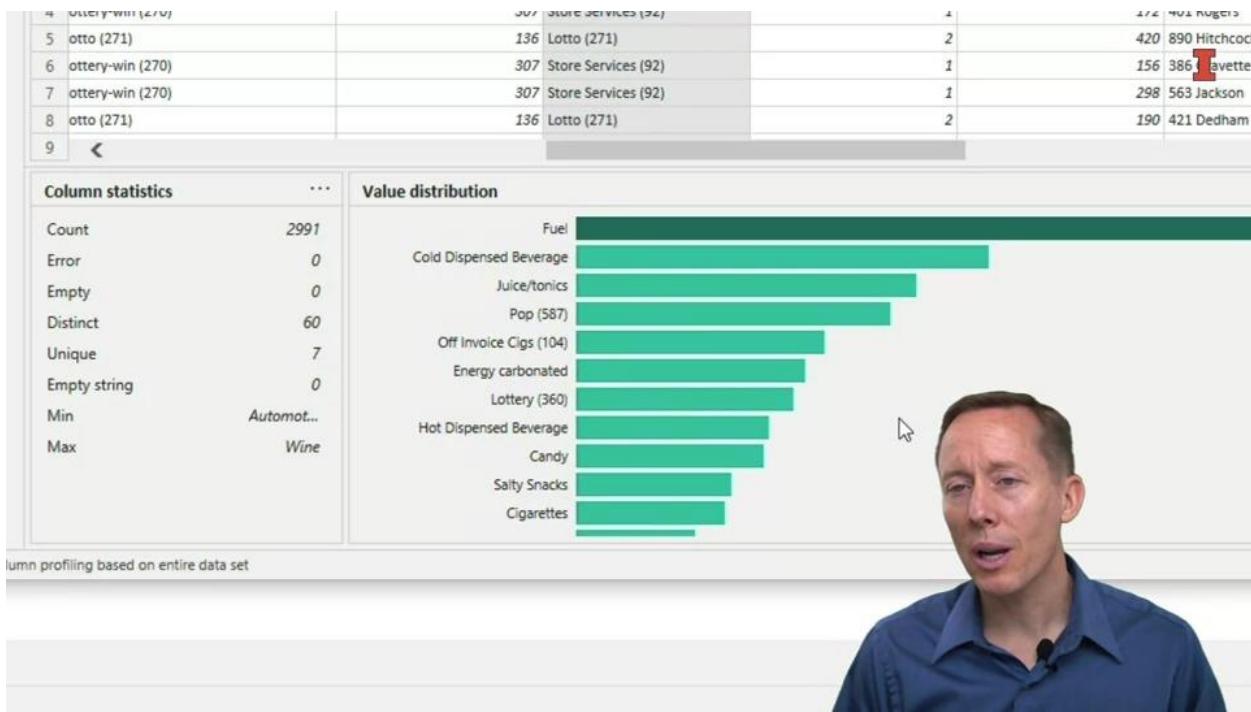
So at this point, I may think to myself, okay, I want to get rid of these. They're not entered in correctly. If they were, then I'd want to go deal with that, but if they're not, then I may not want them though, because they're not representative of most of the other transactions. So one way that I can get rid of these observations is I can go back to the revenue column, and since we've already got it sorted in descending order of revenue, I can go down to this and see the unique values and I can just manually say, "Hey, I don't want these top eight values in here, anything that's above 86 and click Okay," and that will remove those rows from this data set from any additional analysis that I do with this data set. Now, there's another way that I can deal with these and to show you that I'm going to undo that way of filtering. I'll go to the Home tab. In here, I can go to the Reduce Rows section of the ribbon and click on this "Remove Rows" button. I want to remove the top X number of rows. All I have to do is say, hey, since I know these top eight rows are the ones that I don't want, I can click eight there and click "OK". Now this only works because I first sorted the data. If I want to remove the bottom rows of data, I could scroll down to the bottom or I could resort it again, but sort in ascending order. Notice, again, that in this Applied Steps area, it's keeping track of all these different transformations I do with the data. I sorted the data. Now I can see I've got this negative 100 value here. There's really just one observation that's below negative 70.

	de	A _B C site_status	1.2 revenue	-7	1.2 gross_profit	-7	1.2 costs	-7	1.2 units	-7	A _B C dow	-7
		100% 0% 0%	100% 0% 0%	100% 0% 0%	98% 0% 2%	98% 0% 2%	98% 0% 2%	98% 0% 2%	94% 0% 6%	94% 0% 6%	ABC 123 E	ABC 123 E
		Valid Error Empty	Valid Error Empty	Valid Error Empty	Valid Error Empty	Valid Error Empty	Valid Error Empty	Valid Error Empty	Valid Error Empty	Valid Error Empty	Valid Error Empty	Valid Error Empty
		1, 4 unique	2 distinct, 0 unique	692 distinct, 475 unique	2326 distinct, 2161 unique	2324 distinct, 2183 unique	504 distinct, 474 unique	7 distinct, 0 unique	7 distinct, 0 unique	7 distinct, 0 unique	7 distinct, 0 unique	7 distinct, 0 unique
1	-94.0176	ACTIVE		-100	0	0			100	Wednesday		
2	-91.3771	ACTIVE		-50	0	0			50	Thursday		
3	-95.0022	CLOSED		-24	0	0			24	Friday		
4	-94.9895	ACTIVE		-20	0	0			20	Sunday		
5	-94.1246	ACTIVE		-20	0	0			20	Sunday		
6	-98.3491	ACTIVE		-20	0	0			20	Monday		
7	-94.4614	ACTIVE		-18	0	0			18	Tuesday		
8	-95.3812	ACTIVE		-13	0	0			13	Saturday		
9	-94.823	CLOSED		-12	0	0			12	Tuesday		
10	-94.0928	ACTIVE		-12	0	0			12	Monday		
11	-92.3313	ACTIVE		-10	0	0			10	Monday		
12	-92.22	ACTIVE		-10	0	0			10	Monday		
13	-104.891	ACTIVE		-10	0	0			10	Friday		
14	-94.4614	ACTIVE		-10	0	0			10	Thursday		
15	-107.009	CLOSED		-10	0	0			10	Wednesday		
16	-92.3734	CLOSED		-10	0	0			10	Sunday		
17	-104.963	ACTIVE		-10	0	0			10	Thursday		
18	-94.9733	ACTIVE		-10	0	0			10	Tuesday		
19	-105.036	ACTIVE		-10	0	0			10	Sunday		
20	-104.771	ACTIVE		-8	0	0			8	Wednesday		
21	-92.2822	ACTIVE		-6	0	0			6	Sunday		
22	-91.3257	ACTIVE		-6	0	0			6	Friday		
23	-94.079	ACTIVE		-6	0	0			6	Wednesday		

I could explore this and see what it's related to. In this way it's store services. Again, let's assume that this is entered correctly, but it's not representative of what we normally do, so we want to get rid of it. At this point, I might just want to click on this and uncheck the negative 100 and click "OK".



Now I have dealt with my outliers for the numeric columns. Now, this is just one way of dealing with them. Oftentimes, in Power BI, we don't want to remove those rows here, but we may want to remove them using slicers or filters after we create observations. But anyway, you can remove them from the data in the Power Query editor. That's how you do with numeric columns of data. Let's say we want to look at categorical columns of data.



Let's go to the parent name column and let's explore this column a little bit. We can see that in the parent name column, I've got 60 distinct values and seven of those are unique. If I click on that column header and I have the column profile checked in the view menu item, I can get an idea of what those values are in the 60 columns and how often they occur. I can very quickly see that fuel shows up most often, cold dispense beverages, then juice/tonics, and then pop. Hey, three of these top four are beverages. That might stand out to me, so I keep that in mind. Anyway, I can go down through here. If I think about visualizations and how I want to explore things, having 60 different values, seven of which only occur once, that could make it problematic for color coding things. There'd be too many colors. The technical name for this is high cardinality. High cardinality means there's a lot of different values in a set. Oftentimes, it just makes it too cumbersome and hides some of what's going on. I want to reduce the cardinality somehow.

The screenshot shows the Power BI Desktop application with the 'Power Query Editor' open. The main area displays a table with 25 columns and over 999 rows. The columns include 'category_name', 'parent_id', 'parent_name', 'product_count', 'site_id', 'site_name', and 'address'. The 'Applied Steps' pane on the right lists several steps taken during the data transformation process, such as 'Remove Duplicates', 'Changed Type', and 'Replaced Errors'. A man in a blue shirt is overlaid on the right side of the interface.

One thing I could do is click on the parent name or hover over it and go to "Remove Duplicates". If I do that and ask myself, does that work? Well, actually, what happened is now I say, I've just got 60 rows. It removed any observation for which the parent name occurred a second time. That's not what I wanted to do. Fortunately, we don't have to worry about, how do I undo that? I can just go over to the Applied Steps and I click on that "Remove Duplicates" button. Now they're all back there. I've got all the data. I want to show you a cool way for dealing with high cardinality in Power BI. First, I need to exit the Power Query editor, so I'll close and apply these changes. Great.

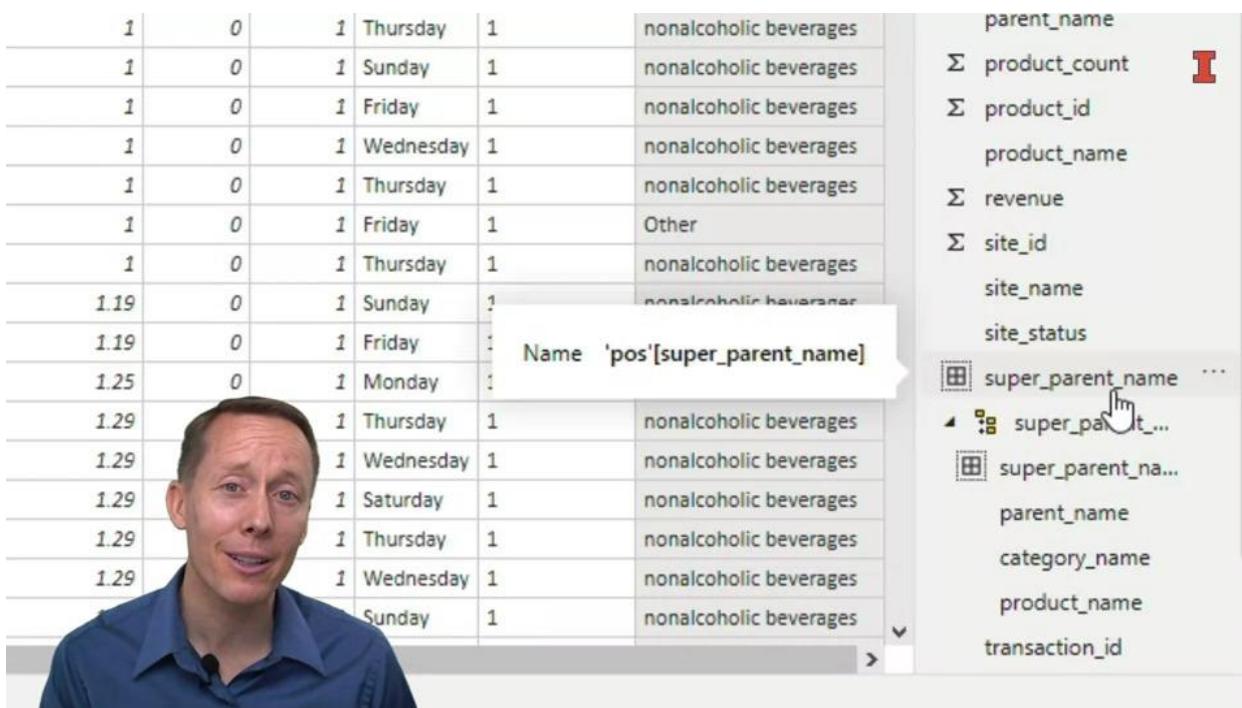
The screenshot shows a software window titled "Groups". At the top, there are fields for "Name" (set to "super_parent_name") and "Field" (set to "parent_name"). Below these are dropdown menus for "Group type" (set to "List") and "Group by" (set to "parent_name"). The main area is divided into two sections: "Ungrouped values" on the left and "Groups and members" on the right. In the "Ungrouped values" section, there is a list of items: Lottery (300), Lotto (271), Milk, Newspapers, No Tax Meds (246), Off Invoice Cigs (104), Packaged Sandwiches, Pizza, Pop Deposits (588), and Postage Stamps (233). In the "Groups and members" section, there is a list of groups: alcoholic beverages, Fuel, nonalcoholic beverages, snacks, and Other. A checkbox labeled "Include Other group" is checked. At the bottom of the window are "Group" and "Ungroup" buttons, and "OK" and "Cancel" buttons.

Now, let's go ahead and go over to the fields section and let's go to the parent name and click on the three dots next to it. Let's go to New group. What we will do is we will manually create groups for all of these different parent names. I'm actually going to create a new column. This column I'll call super parent name. Now what I can do is group together, perhaps any parent name that is a beverage. As I start exploring these parent names, now they're listed in alphabetical order, I might, real quickly see here I've got beer. I'll click on that and I can see that this group comes up. I could create a group, but I want to add more to that. Let's go down and identify any other perhaps alcoholic beverages in here. I've already done this, so I know that there's a liquor in here. I'm holding Control Down and clicking on "Liquor" and then there's wine. I've got those three parent names selected. I'll click "Group", and now I've got this new group and the members associated with it and I want to rename this to alcoholic beverages. That's one group. Now, I also may want to create just a beverages group. Let's go down and identify all the beverages that are nonalcoholic and put those into a group. Now I've got this nonalcoholic beverages group here that's made up of these items. Now let's go ahead and create another group that is made up of snacks. I created a new snacks group made up of these parent name items and I can collapse that. I can say I've got these three groups so far, alcoholic beverages, nonalcoholic beverages, snacks. I also want to create a group that's made up just of fuel. I'm doing that because I know fuel is a high-volume item. Now, all these other things may not be as important to us. They may make up a smaller amount of revenue, and so I don't necessarily want to leave those out or leave them separate. I can real quickly, click this "Include Other Groups" and everything else falls into this other category. Then I'll click "Okay".



	site_status	revenue	gross_profit	costs	units	dow	gp_margin	super_parent_name
5	ACTIVE	0.5	0.5	0	1	Wednesday	1	nonalcoholic beverages
5	ACTIVE	0.99	0.99	0	1	Wednesday	1	Other
3	ACTIVE	0.99	0.99	0	1	Wednesday	1	Other
3	ACTIVE	0.99	0.99	0	1	Wednesday	1	Other
3	ACTIVE	1	1	0	1	Saturday	1	nonalcoholic beverages
7	ACTIVE	1	1	0	1	Tuesday	1	nonalcoholic beverages
1	ACTIVE	1	1	0	1	Thursday	1	nonalcoholic beverages
1	ACTIVE	1	1	0	1	Tuesday	1	nonalcoholic beverages
5	ACTIVE	1	1	0	1	Tuesday	1	nonalcoholic beverages
4	ACTIVE		1	0	1	Friday	1	nonalcoholic beverages
5	ACTIVE		1	0	1	Wednesday	1	nonalcoholic beverages
5	ACTIVE		1	0	1	Wednesday	1	Other
5	ACTIVE		1	0	1	Tuesday	1	nonalcoholic beverages
4	ACTIVE		1	0	1	Friday	1	nonalcoholic beverages
5	ACTIVE		1	0	1	Tuesday	1	nonalcoholic beverages
3	ACTIVE		1	0	1	Saturday	1	nonalcoholic beverages
4	ACTIVE		1	0	1	Saturday	1	nonalcoholic beverages

Now I can say I've got this super parent name column over here and if I click on the data and explore it quickly, I can see that this column is added to the end, and if I click on the Drop Down arrow, I can see that I've only got these five different categories in here; alcoholic beverages, fuel, nonalcoholic beverages, other and snacks. That's one way of dealing with high cardinality columns. Now, here's another cool thing you can do with Power BI.

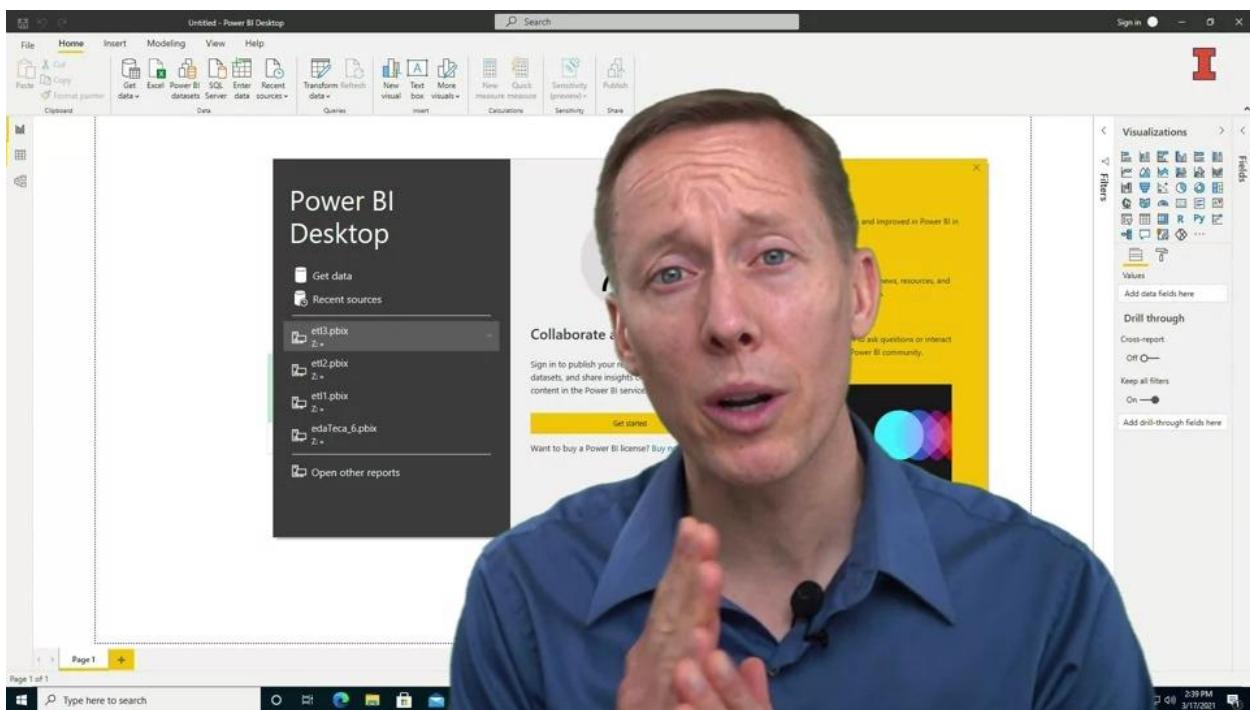


The screenshot shows a data grid and a sidebar menu. The data grid contains columns for various variables. A tooltip is visible over the 'super_parent_name' column, showing the path 'Name > 'pos'[super_parent_name]'. The sidebar menu lists several columns: parent_name, product_count, product_id, product_name, revenue, site_id, site_name, site_status, super_parent_name (with a red highlighted icon), super_pa..., parent_name, category_name, product_name, and transaction_id.

1	0	1	Thursday	1	nonalcoholic beverages	parent_name
1	0	1	Sunday	1	nonalcoholic beverages	product_count
1	0	1	Friday	1	nonalcoholic beverages	product_id
1	0	1	Wednesday	1	nonalcoholic beverages	product_name
1	0	1	Thursday	1	nonalcoholic beverages	revenue
1	0	1	Friday	1	Other	site_id
1	0	1	Thursday	1	nonalcoholic beverages	site_name
1.19	0	1	Sunday	1	nonalcoholic beverages	site_status
1.19	0	1	Friday	1	Name > 'pos'[super_parent_name]	
1.25	0	1	Monday	1	super_parent_name	
1.29		1	Thursday	1	nonalcoholic beverages	super_pa...
1.29		1	Wednesday	1	nonalcoholic beverages	super_pa...
1.29		1	Saturday	1	nonalcoholic beverages	parent_name
1.29		1	Thursday	1	nonalcoholic beverages	category_name
1.29		1	Wednesday	1	nonalcoholic beverages	product_name
			Sunday	1	nonalcoholic beverages	transaction_id

Remember how with the date column, it created automatically this date hierarchy so we could scroll down through the data and start with year then go down to quarter and then month and date, we can do that with any column of data or columns of data. This will be very helpful as you create visualizations that you may want to drill down into. Let me show you how to do that. I can go to super parent name, click on the three dots there and click on New hierarchy. Now I've got this super parent name hierarchy created and I can add different columns to that hierarchy, so I can drill down. After the super parent name, comes the parent name. I'll go up to the parent name column, click on the three dots next to it, and click "Add to Hierarchy", super parent name hierarchy. If there were more than one hierarchy, I can choose which hierarchy I'd want to add it to. Now, under the super parent name hierarchy, I've got parent name and after parent name, I would want to add category name. I go up to category name, click on the three dots, and add that to the super parent name hierarchy. There is category name, and then I can go to product name and add that to the hierarchy as well. Creating this hierarchy is a bonus, it's something that'll be useful as we create visualizations. Sometimes you want to create these hierarchies after exploring the data a little bit, as well as dealing with outliers. Once you start exploring the data visually, you'll see they're outliers. We'll come back to dealing with outliers later. But this is one way to deal with outliers in both numeric and categorical columns.

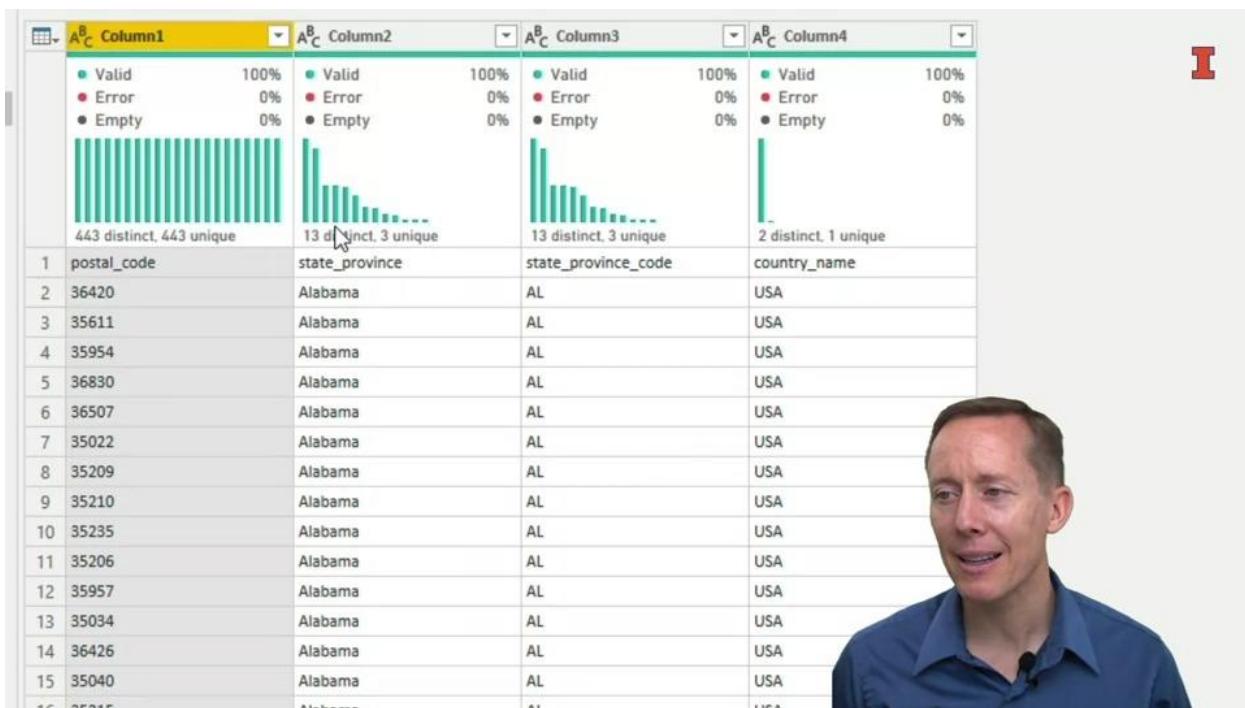
Lesson 2-2.6 ETL 4: Data Models and Joins



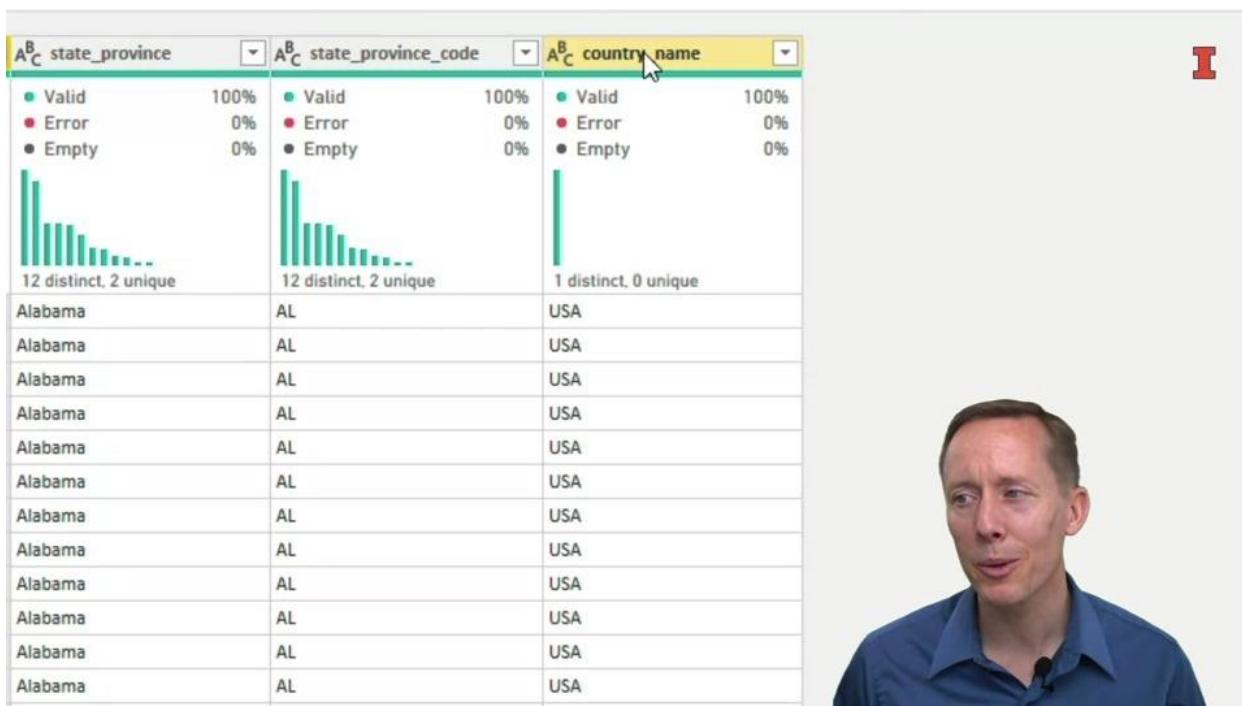
In this video, we will focus on joining data together in Power BI. This is important because oftentimes you may have one table or a data set that you're focused on, but you want to bring in information from other data sets. We'll need to do that by matching values on one or more columns. We will focus on how to do that in Power BI, and along the way, we will also illustrate how to create a simple chart and how to do some data wrangling as well.

The screenshot shows the Power BI desktop application. In the top navigation bar, there are three main sections: 'Filters', 'Visualizations', and 'Fields'. The 'Fields' section is currently active, indicated by a red icon in the top right corner. Below the navigation bar, there is a search bar and a 'Search' button. Under the 'Fields' section, there is a 'Search' bar and a list of datasets. Two datasets are visible: 'pos' and 'states'. The 'states' dataset is selected, as indicated by a hand cursor icon pointing at it. A tooltip for the 'states' dataset provides details: Name: states, Storage mode: Import, and Last refresh: 3/17/2021 2:41:57 PM. At the bottom of the Fields section, there is a 'Values' button and a 'Add data fields here' button.

The first thing that we should do is open up Power BI and I will open this project that we've been working etl3. Now we've got this project loaded and we can see in the field section that we've got one data set, pos, or the point of sale data, and within this data set, we have a column called Zip, for zip code, also known as postal code. We can see it's a numeric value which will be important a little bit later on. Now, what if I want to aggregate data and visualize data based on state or province? I don't currently have that information in here. But if I had a list of zip code and which state they belong to, I could bring in that state information and then do that type of aggregation. Fortunately, I have that data and that data is located in the same folder as this data, it doesn't have to be, but it's called states.csv. Let's first read in that data by going to the home tab and clicking on "Get Data." Since it's a.csv file, we'll click "Text/CSV" and Connect. Then indicate we want the states.csv file and here is a preview of the data. I'll go ahead and load it. Now if I look over in the field section, you can see I've got pos and states. Now, I like the name states, so I'll just leave it at that. But let's go into the Power Query Editor to transform and further analyze what this data set looks like.



Here on the left-hand side, we can see we've got pos and the states, I'll click on "States," and the first thing that I think we should do is just look at the structure of this data set. We can see that it has four columns and 443 rows. Now, if you browse this data, you can see that the first row appears to be the column header, but Power BI did not know that it should be the column header. So we want to somehow get this up to the column name up here. Fortunately, there's a really easy way to do that because this happens quite a bit.



In the home tab, we can go over here and click "Use First Row as Headers." Now we can see that there are only 442 rows instead of 443, and these column names are more appropriate. Let's start looking at the values in the columns, inspect the quality, and so forth. For the postal code, we can see that there are 442 distinct values and 442 unique values. This means that each postal code only shows up one time. If we look at the state_province column, we can see there are 12 distinct values and two are unique. So two of the states only show up once. We can click on the "state_province_code" or hover over that and see that it's very similar to a state or province, it's just that these are two letter abbreviations. Then the country_name has only one distinct value and zero unique, meaning everything is USA. Well, this column probably, really isn't that important. Let's go ahead and get rid of it.

	1 ² 3 site_id	A ^B site_name	A ^B address	A ^B city	1 ² 3 zip	1.2 latitude
10%	Valid	100%	Valid	100%	Valid	100%
0%	Error	0%	Error	0%	Error	0%
0%	Empty	0%	Empty	0%	Empty	0%
	178 distinct, 16 unique	178 distinct, 16 unique	178 distinct, 16 unique	160 distinct, 11 unique	170 distinct, 14 unique	178 distinct, 16 unique
1	27 143 Halevville	42417 Hwy 195	Halevville		35565	
2	100 248 Newell	119 S Fulton St	Newell		50568	
1	274 520 Atlantic	1630 E 7Th St	Atlantic		50022	
1	172 401 Rogers	401 N 8Th St	Rogers		72756	
2	420 890 Hitchcock	107 W Main St	Hitchcock		73744	
1	156 386 Gravette	613 Main St Sw	Gravette		72736	
1	298 563 Jackson	4206 N College Ave	Jackson		36545	
2	190 421 Dedham	415 3Rd St	Dedham		5110	
1	199 435 Golden City	519 Main St	Golden City		64748	
1	434 914 Centennial	2221 E Arapahoe Rd	Centennial		80122	
1	436 916 Pagosa Springs	250 Hot Springs Blvd	Pagosa Springs		81147	
1	252 493 Brewton	2041 Douglas Avenue	Brewton		36426	

We'll click on that column header and then in the home tab here, click on "Remove Columns," and we got rid of that column. As always, if we make a mistake or want to revert back to a previous step, we can go over to applied steps and remove that. But I'm going to go ahead and leave that here. The idea here is that we want to take these postal codes here and maybe we just browse them in. As we go down, look at these, you can see that some of these observations have a nine-digit zip code. We've got hyphenate the first five digits, so hyphen in the last four digits. What we want to do is take these state values and link it to each transaction here based on the zip code. If we scroll over here a little bit, we can see we've got the zip code in here and there are the unique values. These are only five-digit zip codes. At some point, we're going to have to get rid of those last four digits.

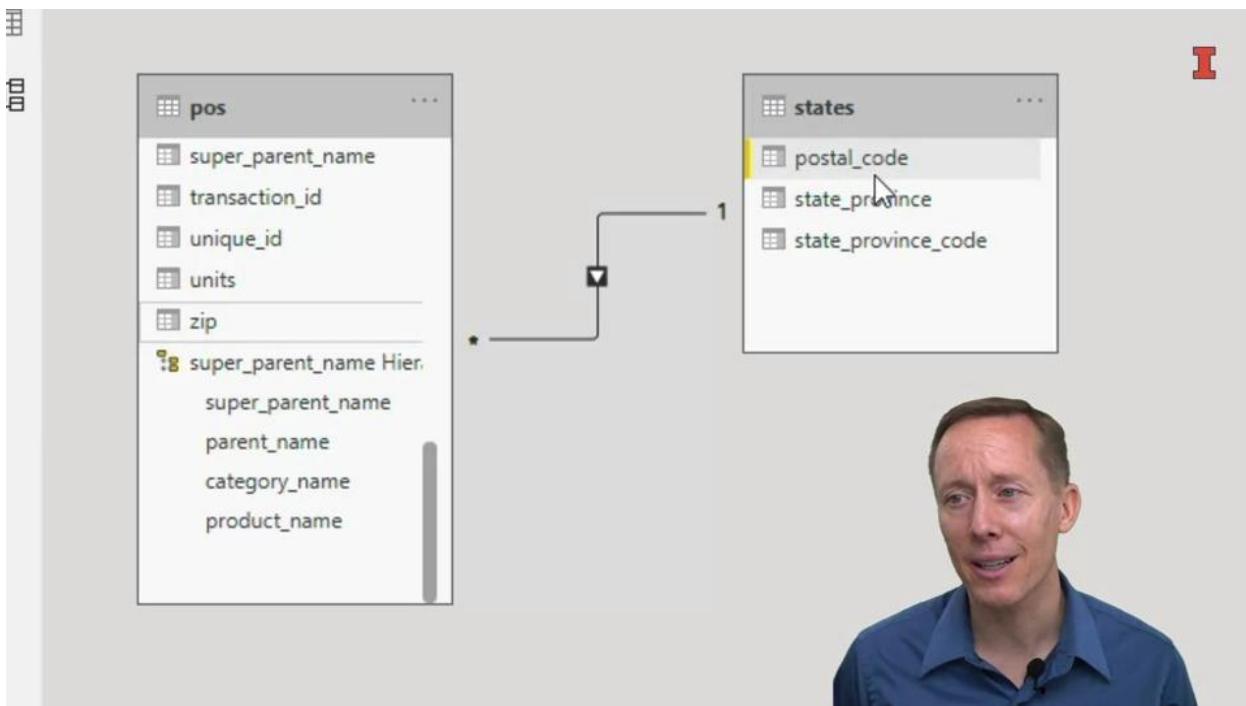
The screenshot shows a data modeling interface with two datasets displayed as rectangles. The left dataset is 'pos' and the right is 'states'. The 'pos' dataset contains columns: address, category_id, category_name, city, costs, customer_id, date, dow, gp_margin, gross_profit, and latitude. The 'states' dataset contains columns: postal_code, state_province, and state_province_code. A sidebar on the left shows icons for different data types. A message at the top says: 'Upgrade to the new model view for an improved design that makes it easy to identify data sources and manage relationships. Keep in mind, I'.

But once we've done that with the data, let's just go ahead and close and apply these changes that we made to the states' data set. Now that we've made those changes, let's talk about how to join data using the data modeling tool. I'll click on this Data Model, on the left sidebar. This gives me a rectangle for each data set that I have. Now, if there are any data sets that have the same column name and the same type, it will automatically recognize that these can be joined together.

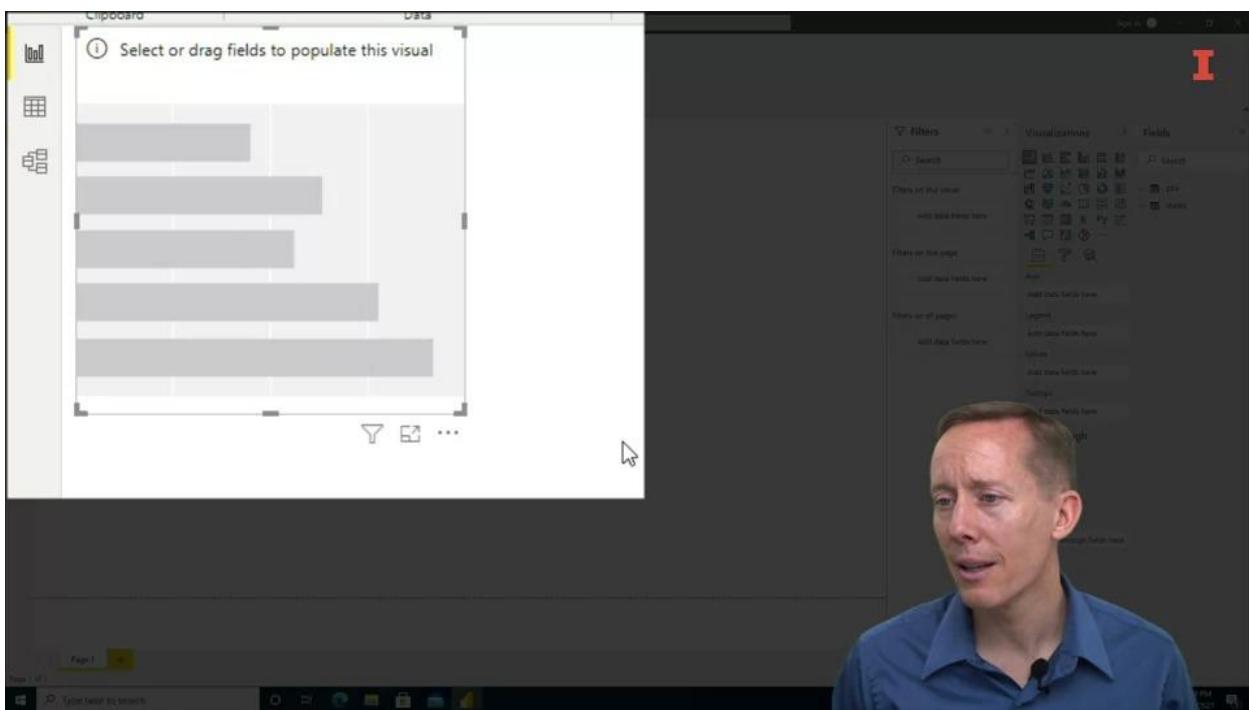
The screenshot shows a data modeling interface with a sidebar containing icons for various data types like date, dimension, fact, and more. A specific column's properties are being edited:

- Synonyms:** postal code, postal_code
- Display folder:** Enter the display folder
- Is hidden:** No (radio button selected)
- Formatting:** A section with a dropdown menu currently set to "Text".
- Data type:** Text (selected from a dropdown menu)
- Format:** Text (selected from a dropdown menu)

That's not the case here, because in the POS data set, we've got the zip column. If I click on that, I can see that the data type is a whole number. In the states column I've got the postal code, and while it's really the same thing, this is a text value here. Now, that doesn't really make a difference in this data model here.



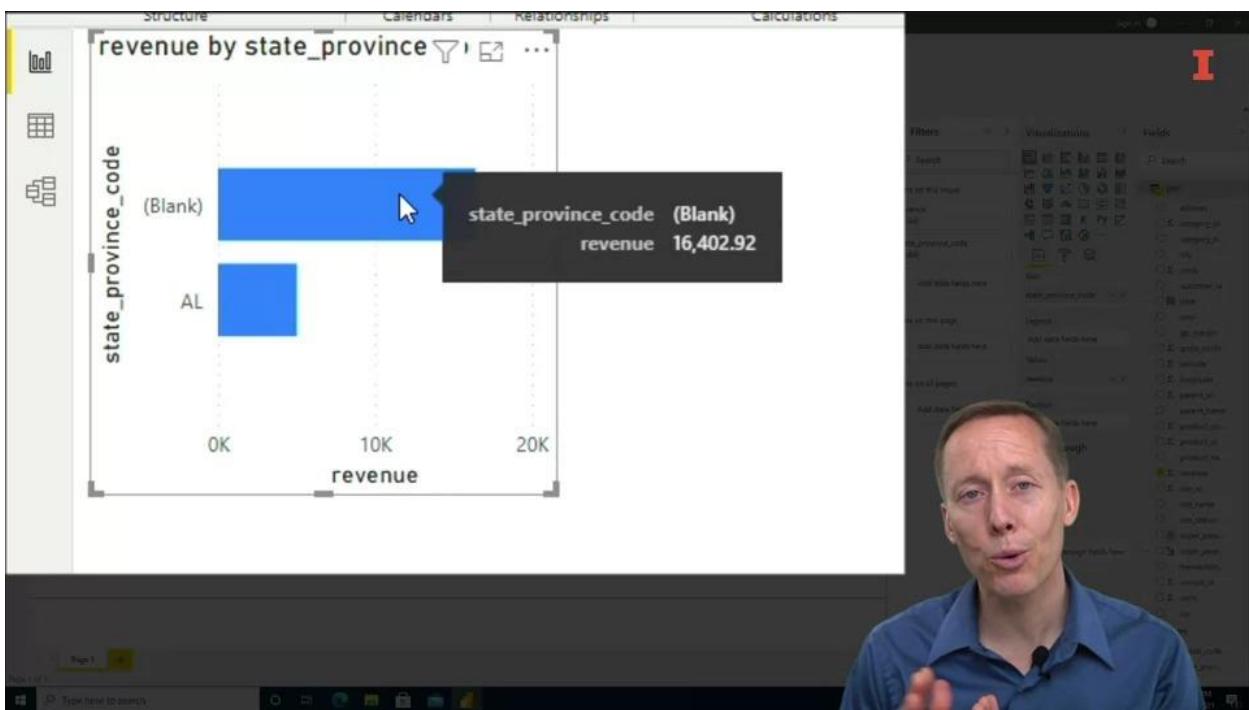
Let's just go ahead and pretend we don't know there's any difference or that there are the last four digits in the states data set. What I can do is I can explicitly define a relationship by clicking on a column name and dragging it over to a column name on another table. I can click on this arrow here and see that these tables are joined together using the zip column and the postal code column. I can also see a star on the side of this arrow that's closest to the POS table and a 1 on the side of this arrow that's closest to the states table. What that means is that there's a 1 to many relationship. In other words, each transaction in each row in POS matches up with only one observation in the states data set, but the observations in the states data set based on postal code could match up with many observations in the POS data set.



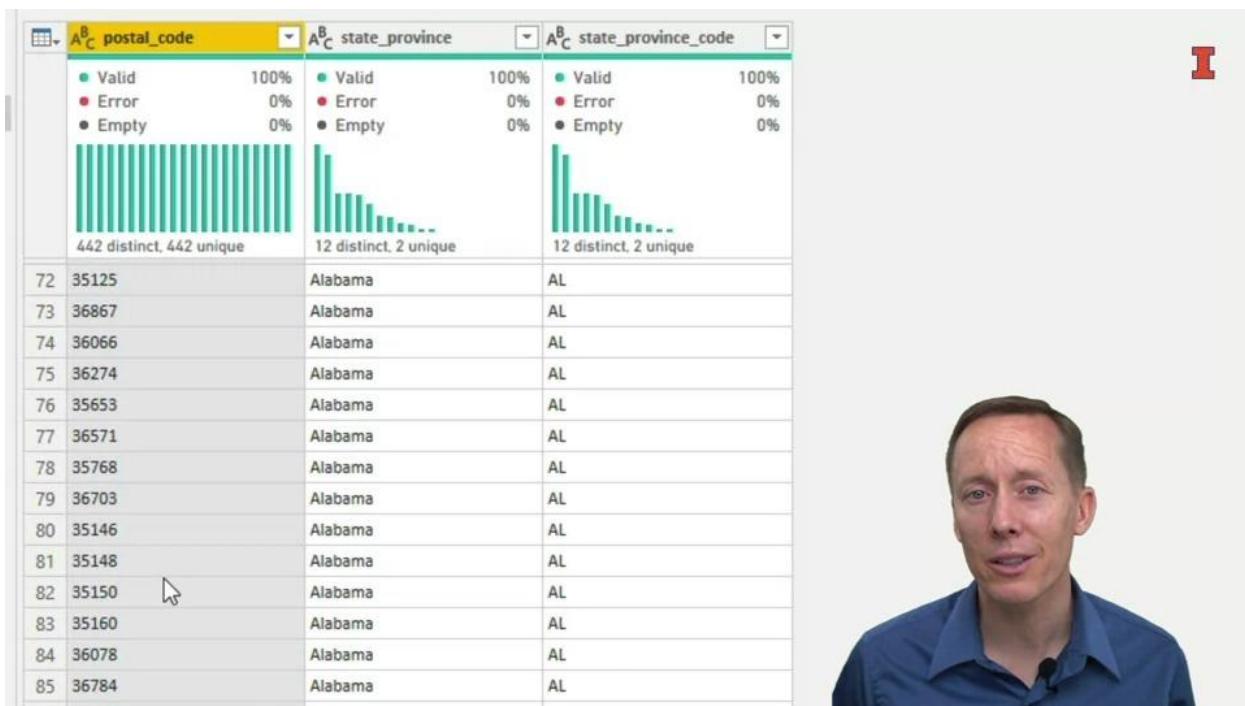
Right now, since we've done that, let's go ahead and just try creating a basic chart. Now, I realize you may not know how to create charts yet in Power BI, but whether you do or you don't it's really easy. Let's just go ahead and create a bar chart here. I will click on that stacked bar chart icon. You can see it creates this folder for the chart here,

A screenshot of the Power BI Fields pane. It shows a search bar and a list of fields under the 'Fields' section. The 'states' field is selected, indicated by a yellow checkmark. Other fields listed include 'pos', 'postal_code', and 'state_provi...'. Below the fields, there are sections for 'Axis', 'Legend', and 'Values', each with a 'Add data fields here' button. To the left of the Fields pane, a video player window shows a man in a blue shirt looking towards the camera.

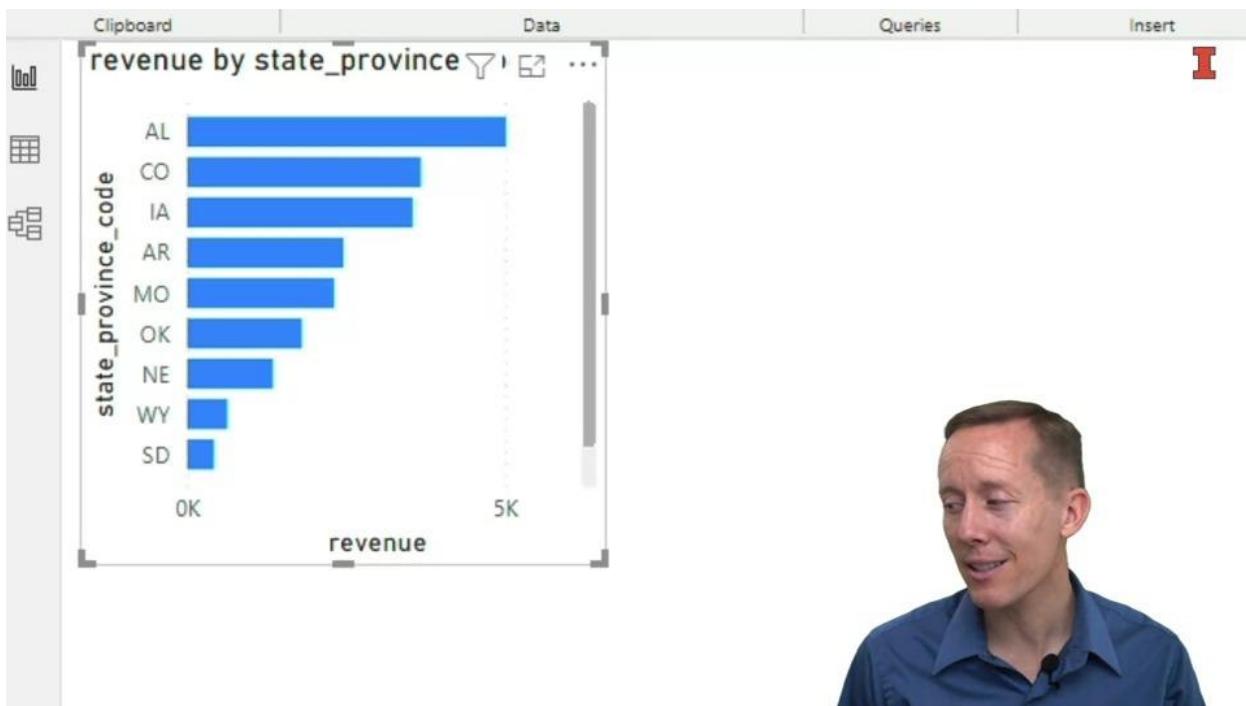
and then if I go into the states data set in here, I will drag state or province to the axis, and then in the POS data, I will drag revenue to the values.



I'm not going to focus on what that means at this point here, but it's pretty easy to see that we created a bar for each state and we added up the revenue. Now you can see that I've got a lot of blanks and one AL for Alabama. I presume there's more than just one state in the this data set. What's going on here? What are all these blanks? Remember how some of the postal codes in the states data set have the four additional digits for the postal code, so we need to go in and clean that up.



Let's go ahead and go to the home button and click on the transform data. Now that we're in the power query editor, we want to go into the states data set and recall that this data set has some observations that have the additional four digits in there. It would be awesome if we could just somehow extract only the information before that hyphen. There is a way to do that. Make sure we've highlighted that column. We'll go to transform, and then we'll go over to the extract area here for text columns, and we're going to select text before delimiter and our delimiter is the hyphen. Click "Hyphen" click "Okay" and now as I go down and look at these different postal codes, you can see that there are no longer the last four digits in that hyphen.



There are a lot of other ways that you can wrangle text data or character strings in here. You can split it, format it and so forth, anyway that's just one example. But now that we've done that let's go ahead and close and apply this, and notice how this chart will be updated here once that has finished processing. Now, we've got more than just Alabama in here and there are no blank states in there. That data model is really cool because even though states doesn't show up in one data set, we can use it from another data set.

The screenshot shows a 'Merge' dialog in the foreground and an 'AI Insights' card in the background.

Merge Dialog:

- Table:** pos
- Columns:** product_count, site_id, site_name, address, city, zip, latitude, longitude, site_status
- Join Kind:** Left Outer (all from first, matching from second)
- Options:** Use fuzzy matching to perform the merge, Fuzzy matching options
- Warning:** Select the same number of columns from both visible tables to continue.

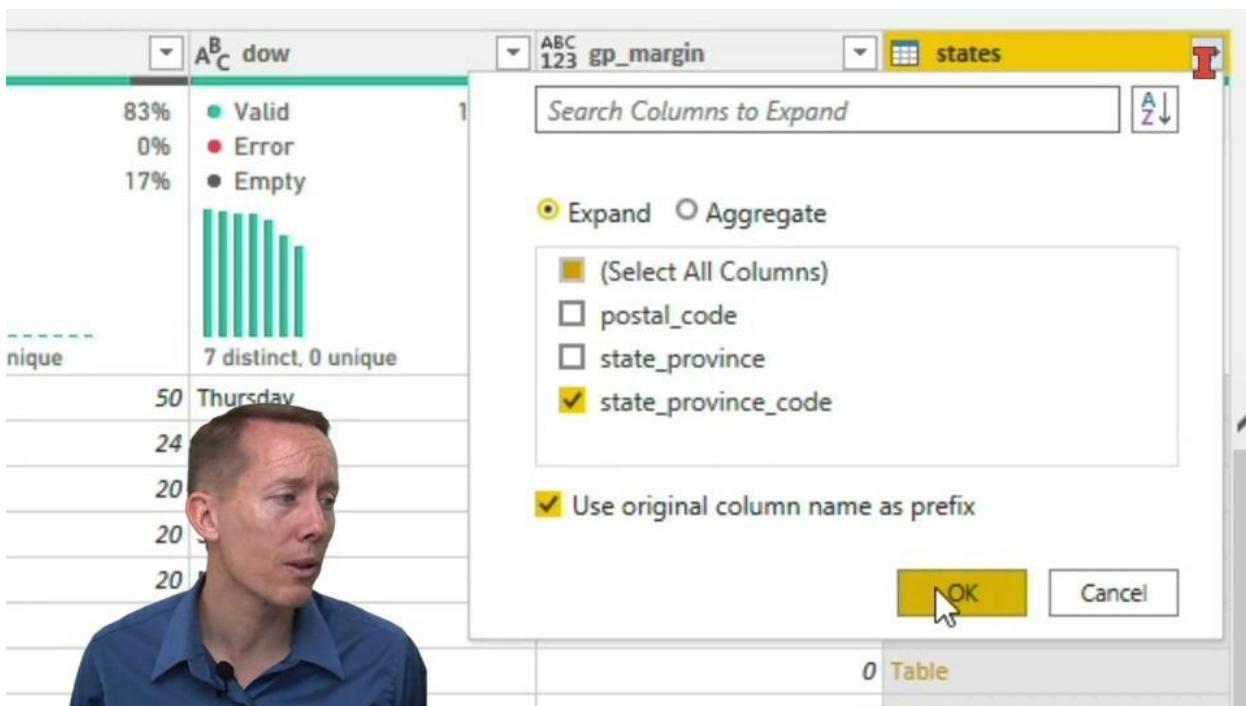
AI Insights Card:

- Category:** category_id
- Properties:**
 - Valid: 100%
 - Error: 0%
 - Empty: 0%
- Visual:** A bar chart showing the distribution of category_ids. 68 distinct, 10 unique.
- Applied Steps:**
 - Source
 - Promoted Headers
 - Changed Type
 - Renamed Columns
 - Changed Type1
 - Replaced Errors
 - Replaced Errors1
 - Replaced Errors2
 - Inserted Day Name
 - Renamed Columns1
 - Added Custom
 - Sorted Rows
 - Removed Top Rows
 - Sorted Rows1
 - Filtered Rows

Sometimes, though, we do want to create a new column for this POS data set. Let's talk about how to do that. I'll click on the "Transform data" and in the home tab here, I can go over to the combine section and just click on "Merge queries" and now what I have is a representation of the POS data set. I could change this data set if I wanted to, but I have this data set. I'm going to scroll over to the zip column, click on that. Down here, I indicate the other data set I want to connect it to and that will be states, and I want to link these together based on the postal code and the zip code. Now we get this prompt to indicate a privacy may be an issue here; we can ignore it for this data set. Now, once we do that, it's not allowing us to link these data sets together. Well, the reason is because it says select columns of the same type.

The screenshot shows a 'Merge' dialog in Microsoft Power BI Data Studio. On the left, there's a preview of two tables: 'pos' and 'states'. The 'pos' table has columns like t_name, product_count, site_id, site_name, address, city, zip, latitude, longitude, and site. The 'states' table has columns like postal_code, state_province, and state_province_code. A dropdown menu shows 'Left Outer (all from first, matching from second)'. Below it is a checkbox for 'Use fuzzy matching to perform the merge'. A note says 'The selection matches 2991 of 2991 rows from the first table.' On the right, an 'AI Insights' card displays metrics for 'category_id': Valid (100%), Error (0%), and Empty (0%). It also shows a histogram of category IDs with 68 distinct, 10 unique values. A list of applied steps is on the far right, including actions like 'Source', 'Promoted Headers', 'Changed Type', etc.

Let's go ahead and make that translation. We can see that the zip column and the point of sale data is a numeric value, whereas in states it is a text or character string value. We need to do convert one to the other. Let's just go ahead and convert postal code to a whole number and then we'll go back over to the point of sale data, click on Merge Queries, and now I can scroll over to zip, select that, select the States, and now select Postal code. This time it says it will match up, it will allow us to move forward. But before we do that, let's just look at a couple of other options here. We can do different types of joins here; left, right, full, inner, and then anti-joins as well. You don't have to necessarily start with POS, you could start with states and then link it to the POS data and then use a right join. In this case, we're going to stick with that left join and click Okay.



Now you can see in this POS data we've got this new column that states and it says we've got a table in every row here, so let's click on that arrow here. This allows us to indicate which columns we want to keep. Well, I really only want to keep the state province code. I don't want to keep the postal code again. We've already got that in the zip, and I don't want the state province; that's a longer name of the code. I'll just make sure I'm only selecting the state province code, click Okay.



	A ^B C dow	ABC 123 gp_margin	A ^B C state
83%	● Valid	100%	● Valid
0%	● Error	0%	● Error
17%	● Empty	0%	● Empty
			
unique	7 distinct, 0 unique		10 distinct, 1 unique
50	Thursday	0	AL
24	F	0	IA
20		0	OK
20	S	0	AR
20	S	0	IA
		0	AR
		0	AL
		0	AL

Now I've got that two digit abbreviation for each state in here, and it's a long name so we double-click on that and rename it just to state. All right, perfect. We have now created a new column in the POS data set. I can close and apply this.



The screenshot shows a Microsoft Power BI interface. On the left, there is a bar chart titled "Revenue by state_province" showing revenue for various states/provinces. In the center, a "Drill-through" pane is open with the following settings:

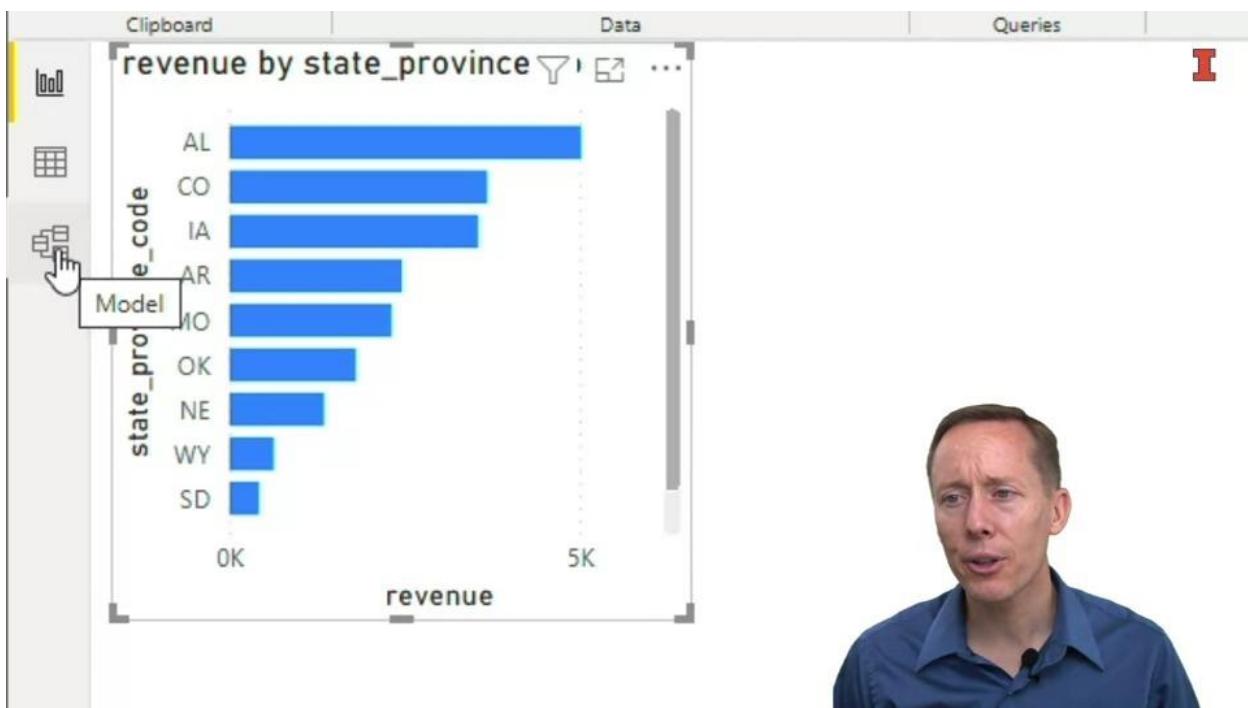
- Cross-report:** Off
- Keep all filters:** On
- Add drill-through fields here:** (empty)

On the right, a list of fields is displayed with checkboxes:

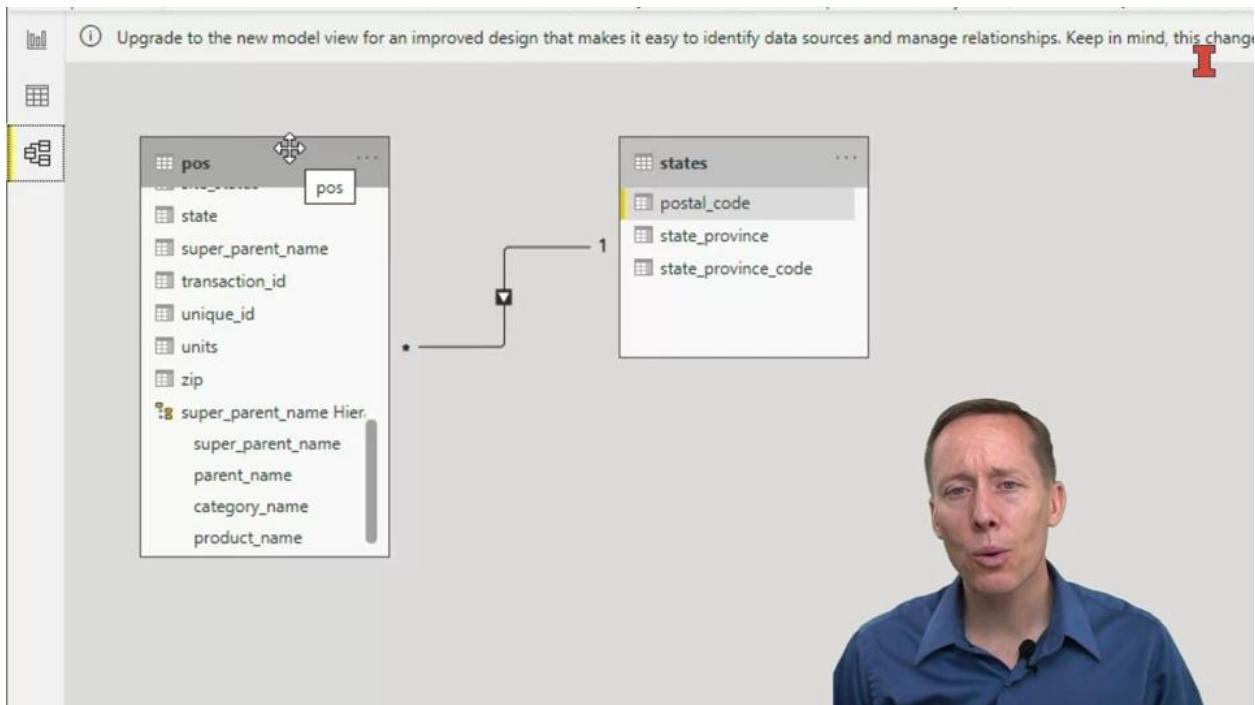
- product_na...
- Σ revenue
- Σ site_id
- site_name
- site_status
- state
- super_pare...
- transaction...
- Σ unique_id
- Σ units
- zip

Now you can go over to the POS data and see that state is in here so we don't necessarily have to join it together using a data model.

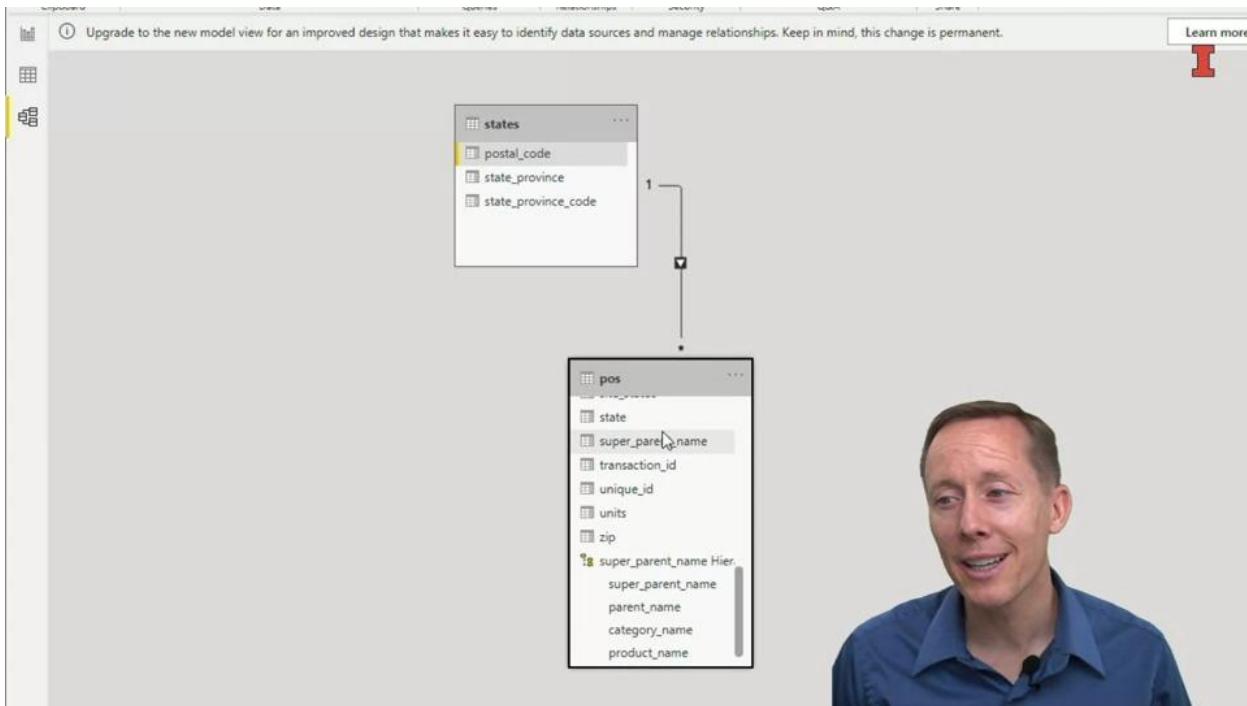
At this point you might be wondering, why are there two different ways to join data? Well, there are advantages and disadvantages to both approaches. As I've already mentioned, adding a new column to a data set can increase the storage size of that data set and all of this file that we're working on and that's not a good idea especially if we're dealing with large amounts of data.



The benefit of the data model is that we don't have to create more data to store. On the other hand,



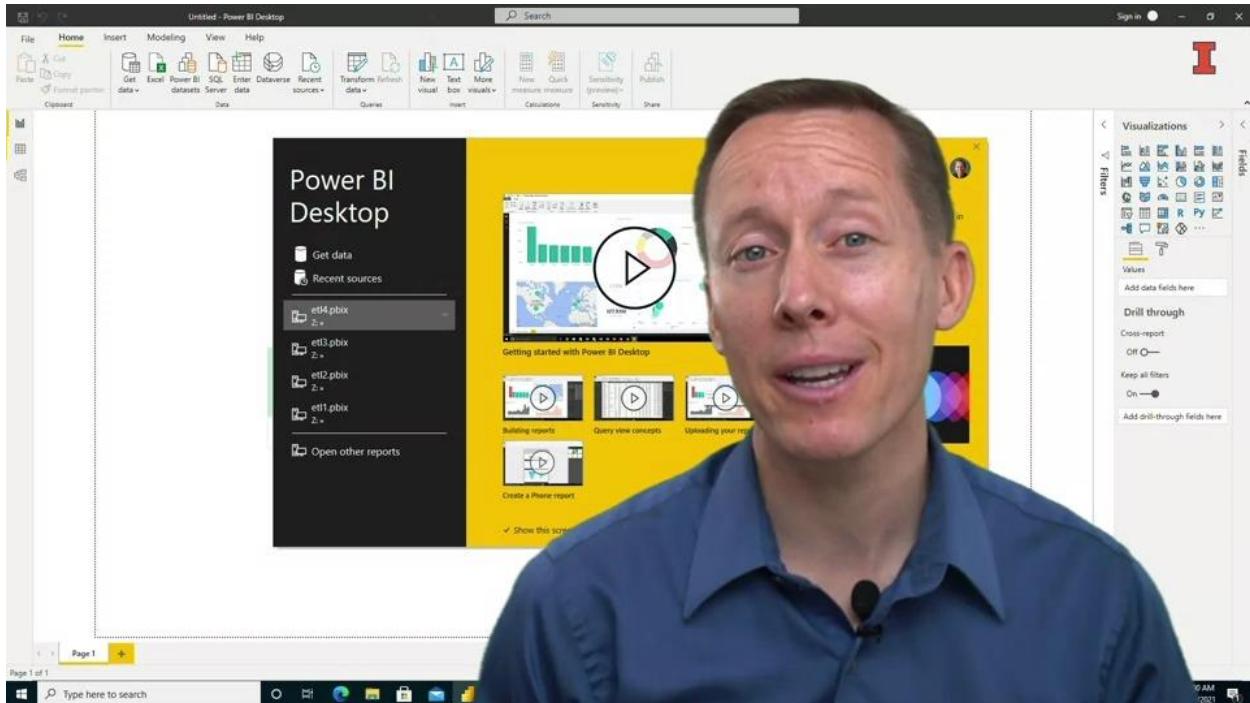
if you are using a data model to join data and you've got four or five different tables in here, you may have this long hierarchical relationship and that can take a long time for queries to run and also just be overwhelming. What's often done is a combination approach, and it's called a Star Schema in which you've got one main data set



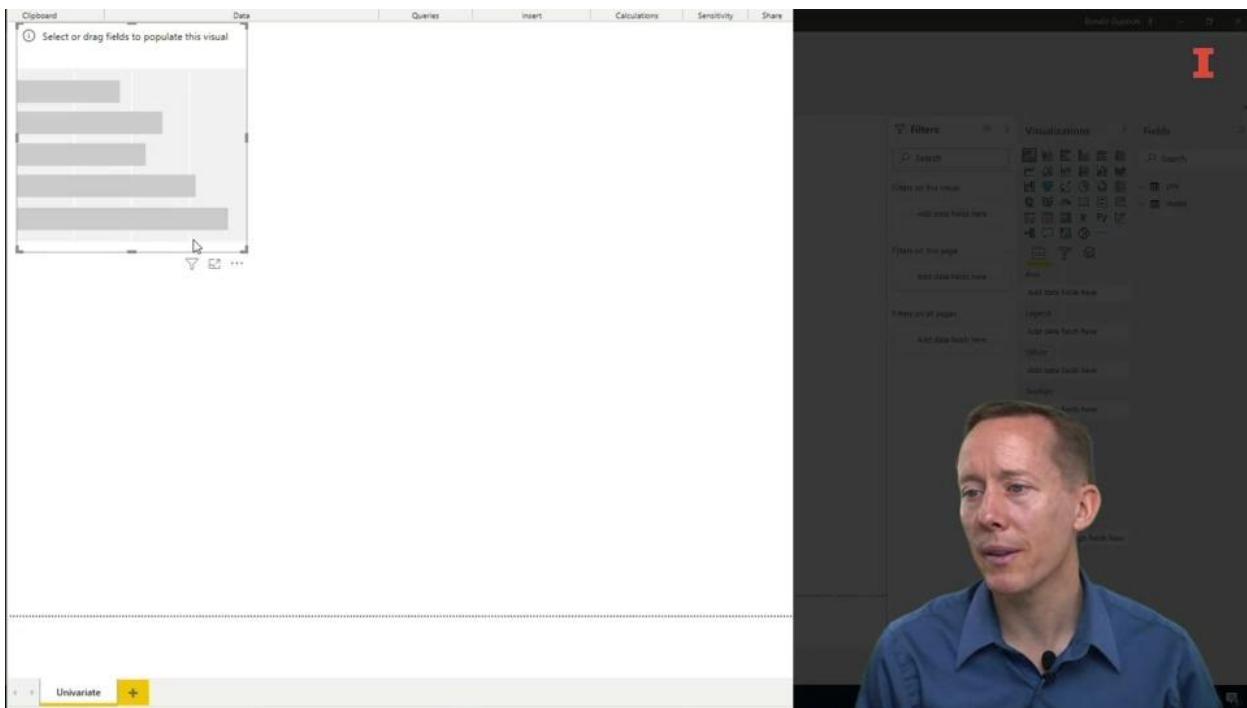
and any other data sets that you want to join together to it, you'll just make sure that they join directly to that key data set. In that way it won't have to query a lot of different hierarchical relationships in here. If you need to add in a column of data from, so you can match one table with another that is not this key table then you'll add in a new column into this key table in the middle. Anyway, that's a star schema. There is an introduction for joining data sets together in Power BI. It's really pretty easy. The benefit is that by bringing data together from different sources, it can really open up the insights that you can gain.

Lesson 2-3: EDA with Power BI

Lesson 2-3.1 EDA 1: Univariate Plots for Numeric Data: Histograms and Boxplots



In this video, we focus on creating visualizations with Power BI. Specifically, we will demonstrate how to create visualizations for exploring one column of data at a time, which is known as univariate data analysis. In this video we'll focus on creating histograms and box plots, which are used for columns that have a numeric and data type.



The first thing that we should do is open up Power BI and I'm going to open up a project that I've been working on which is called etl4. In this project, you can see I've got two data sets, I've got POS for point of sale data and states data. Now we're going to create visualizations in this main area here in the Canvas. But before we do that, let's go ahead and just name this tab, Univariate, so that we can remember as we add more tabs, what is on each tab. Creating visualizations is super simple in Power BI. We've got a variety of default visualizations that come built into Power BI. This first one here is a stacked bar chart, you can see we've got line charts, we've got pie charts, doughnut charts, tree maps and many others. Let's just go ahead and illustrate how to create a basic chart. You simply click on it and you'll see that a placeholder shows up on the Canvas, which we can resize, and that's often a good thing to do when you're creating the chart



The screenshot shows the Microsoft Power BI ribbon. The 'Visualizations' tab is active. On the left, there are sections for 'Filters', 'Visualizations' (with various chart icons), and 'Fields' (listing 'pos' and 'states'). On the right, there are sections for 'Axis', 'Legend', 'Values', 'Toolips', and 'Drill through'. A context menu is open over the 'Visualizations' tab, with the 'Get more visuals' option highlighted.

and then what we will do is we will add fields to the access legend values tool tips and any other placeholders that show up in this area below the visualization. Now we want to create a histogram, so as you look through here, you'll notice that there isn't a histogram that comes built in to Power BI. Fortunately, we can add in more charts by clicking on these three dots here that says get more visuals and then select get more visuals.

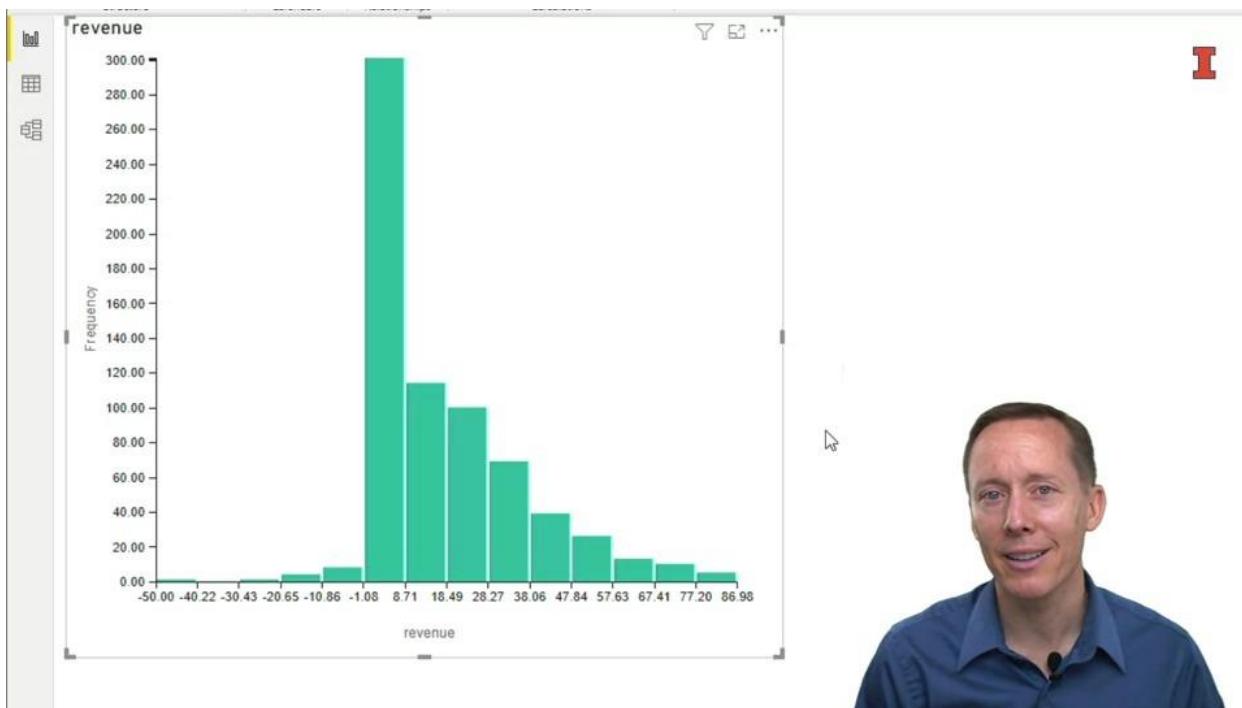
The screenshot shows a user interface for the Power BI AppSource. On the left, there's a sidebar with a search bar containing 'histogram' and a magnifying glass icon. Below the search bar are several filter categories: Category (selected), Editor's Picks, All, Advanced Analytics, Data Visualizations, Filters, Gauges, Infographics, KPIs, Maps, Power BI Certified, and Time. The main area displays four search results for 'Histogram Chart': 1. 'Histogram Chart' by MAQ Software, which visualizes data distribution over a continuous interval or time period. It has a yellow 'Add' button with a cursor hovering over it. 2. 'Histogram with points by MAQ Software...' which displays density of distribution using bars along with actual values represented by points. It also has an 'Add' button. 3. 'Histogram by PQ Systems' which is a bar chart to visualize data distribution and statistics. It includes a note about requiring an additional purchase and has an 'Add' button. 4. 'Violin Plot' which uses a violin shape to visualize data distribution. It has an 'Add' button. At the top right, there's a 'Sort by: Recommended' dropdown.

Now, in order to get more visuals using this approach, you'll have to sign in, and if your organization has a Power BI account, you can sign in with those credentials. I'm signed in with my organization's credentials. Now that I'm in here, I can go ahead and just search for a specific type of chart. I'll type in histogram, there are a variety of histogram charts that I can choose from. I'll just select this first one and click add and it tells me that this visual was successfully imported into this report so I'll click "Okay".

The image shows a user interface for data analysis, likely Power BI. On the left, there is a portrait of a man in a blue shirt. To his right is a large panel containing three main sections: 'Filters', 'Visualizations', and 'Fields'.

- Filters:** This section includes search fields for 'Search' and 'Add data fields here'. It also has sections for 'Filters on this visual', 'Filters on this page', and 'Filters on all pages', each with an 'Add data fields here' button.
- Visualizations:** This section displays a grid of icons representing different chart types. A histogram icon is highlighted with a yellow border. Below the grid are buttons for 'Values', 'Frequency', and 'Drill through', along with dropdowns for 'Cross-report' (set to 'Off') and 'Keep all filters' (set to 'On').
- Fields:** This section shows a tree view under the heading 'pos'. The 'revenue' field is selected and highlighted with a yellow border. Other visible fields include address, category_id, category_n..., city, costs, customer_id, date, dow, gp_margin, gross_profit, latitude, longitude, parent_id, parent_name, product_co..., product_id, product_na..., revenue, site_id, and site_name.

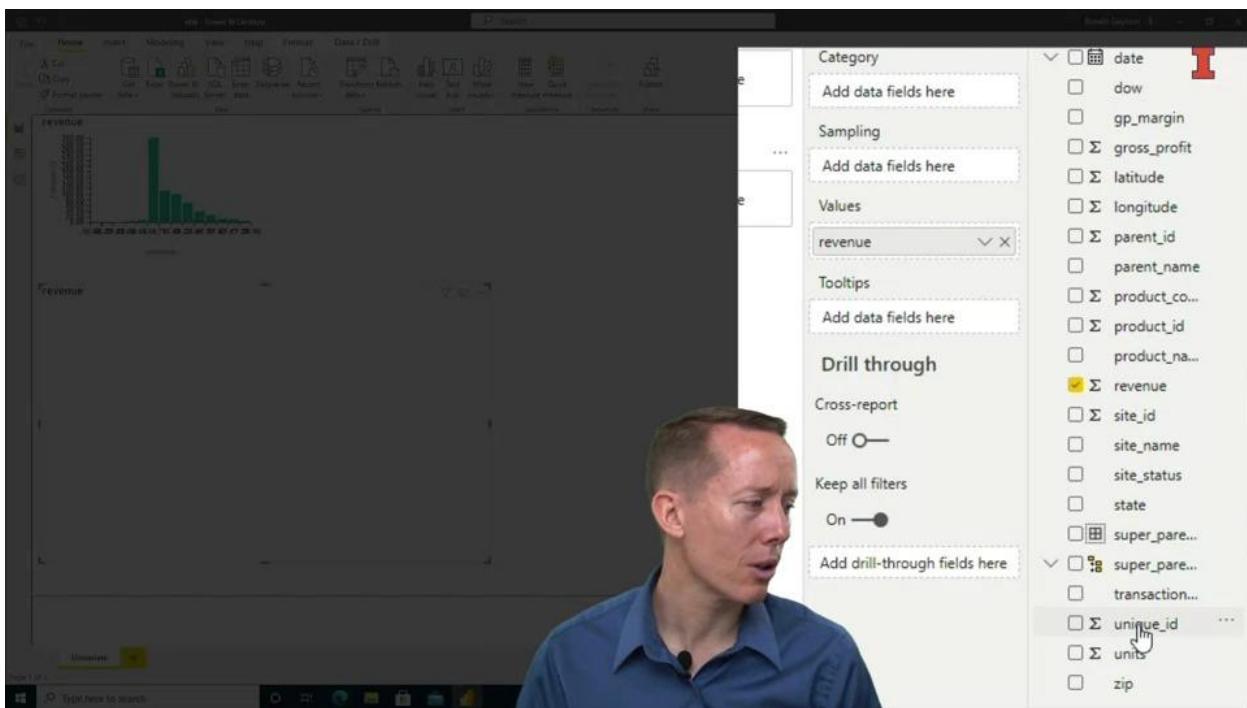
Now you can see I've got the imported visualizations down here below. Let's go ahead and just change the type of this chart by simply clicking on "Histogram", and you can see that updates to a histogram type chart or a different template here and now what I need to do is add column names to the values and perhaps frequency box as well. I'm interested in exploring revenue, so let's go ahead and open this POS data and go down to revenue and drag that into the Values box



and there we have a histogram. Let's take a minute to just identify what a histogram is doing. A histogram takes a continuous variable or a column of data that is numeric and it divides it into equally sized bins and in this case here, the default looks like it's about 12 different bins or so. Let's just hover over this top one, it stands out the most. You can see that it says frequency is 301, meaning there are 301 of approximately 3,000 observations that fall into this range of negative one dollar and eight cents to eight dollars and 71 cents. That's not for a transaction necessarily, it could be but if a transaction has more than one line item, then it will just be for that one line item. You can see that the range here is negative dollar eight cents to eight dollars and 71 cents. The next being right above that only has 114 observations and that ranges from eight dollars 71 cents to \$18.49. It gives you a real quick overview not only of where most of the observations fall, but also how that is distributed. You can see there are a lot more positive values relative to negative values. That's a histogram. Another great chart for exploring numeric data is a Box and Whisker Plot.



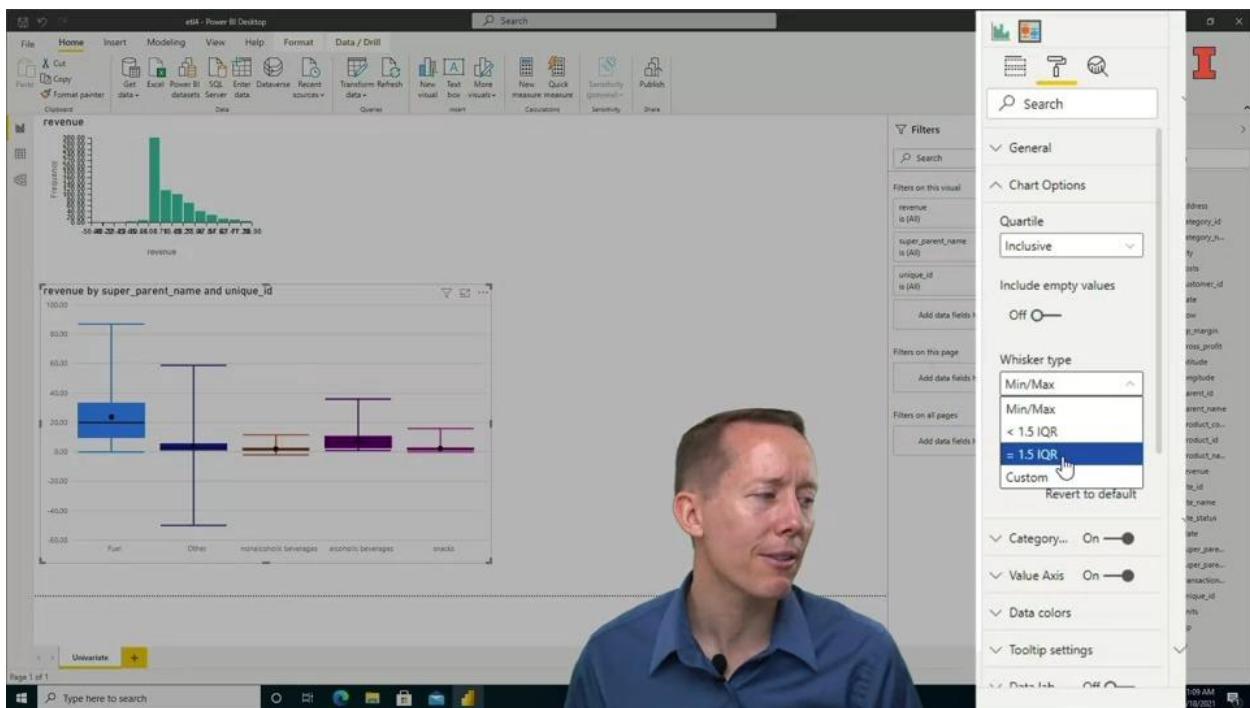
You'll notice that there's not a Box and Whisker Plot built into Power BI by default, so let's go ahead and add one in. I will search for just Box and sure enough, I get at least a couple of different Box and Whisker Plots here. I'm going to go with this one that has a checkmark by it. Click "Add" and we can see that this was successfully imported. We'll click "Okay", and I'll go ahead and make sure I'm not clicking on the histogram click off of that and actually, I might want to even resize this now that this looks how I want it to look, I'll just shrink it down. Now I will add in this Box and Whisker Chart and let's go ahead and make this bigger



and we can see that this chart has four different potential variables that I could add in here. Let's start with the values. I want to focus on revenue, so I'll drag revenue into there and nothing shows up and that's because Power BI is really good at aggregating data and if we click on this down arrow next to revenue, we can see that what it's doing is it's summing up all of the revenue observations into one data point so that will not allow us to create a Box and Whisker Plot. We need more than one observation. If we wanted to, we could average or minimum, maximum, you can do other aggregations, but sometimes there will be an option that just says don't aggregate and if it had it, then I would select that. In this case here though, let's look at the sampling. Let's try using sampling and I suspect that what we need to do is identify that there is a real label for each of these observations and we want to identify that we're using a row label as a unique observation. Let's go and take unique ID, drag that to sampling.



Still nothing shows up, and we can see that by the icon here, it looks like this is intended to create a box and whisker plot for different categories. Let's add something to the category box there. I'm going to select super parent name, drag that into there. All right, there we are. This now is a box and whisker plot for each of the five different super parent name categories that we created, fuel, other non-alcoholic beverages, alcoholic beverages and snacks. It's pretty cool, we can see the median by the line there, and then the dot is the average, then the whiskers extend out to looks like the extreme observations. Now, in all of these plots, we can customize them.



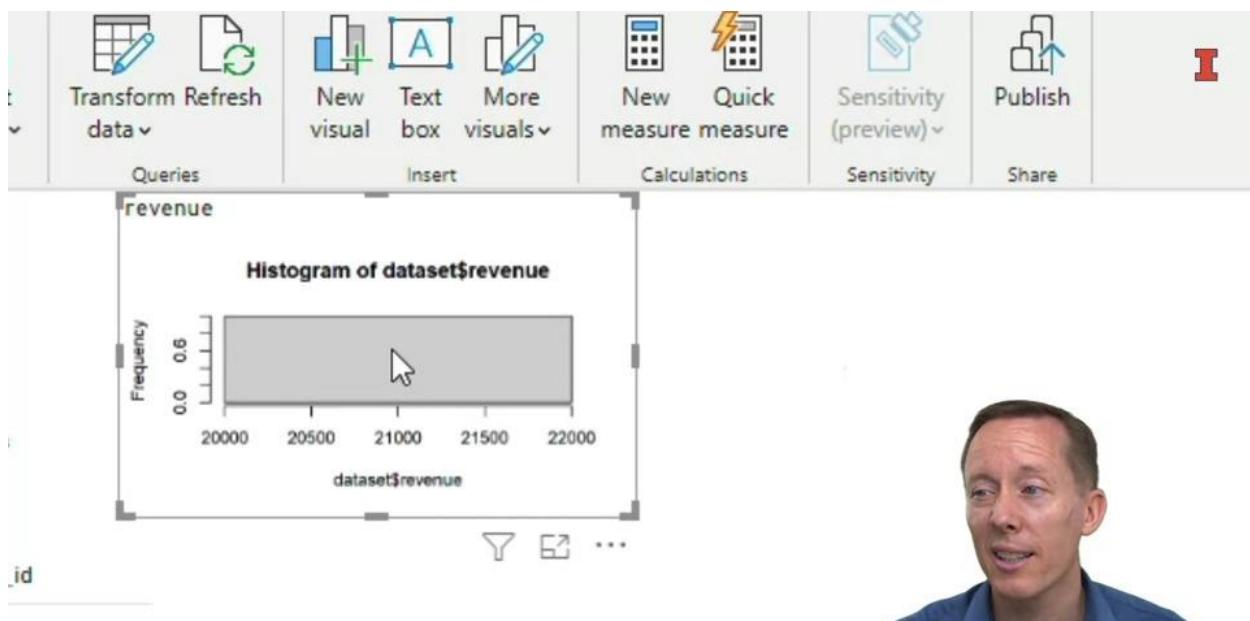
If you go over to this ruler paintbrush, this is how you format plots. With that box and whisker plot selected, I will click format. You can see all the different ways that I can format this, I can add a title, I can change the background, the aspect ratio, data labels, so on and so forth. I'm just going to go into chart options here and I want to maybe identify that these whisker types, I don't want to go all the way to the min the max, I want to do what is often done, which is equal to one and a half times the interquartile range.



I'll select that, you can see this changed a little bit. These whiskers now, the interquartile range is the distance between the top of the box and the bottom of the box, and that's the twenty third percentile and the seventy five percentile or third quartile in first quarter. We're taking one and a half times that distance and making these whiskers go out that far. This just shows, that's why it's symmetric now, it's going out there. It just shows us that there are observations in that range.

The screenshot shows the Microsoft Power BI desktop interface. A histogram titled 'Revenue' is displayed on the left, showing a distribution of values. On the right, the 'Fields' pane is open, listing various data columns such as costs, customer_id, date, dow, gp_margin, gross_profit, latitude, longitude, parent_id, parent_name, product_co..., product_id, product_na..., revenue, site_id, site_name, and site_status. The 'revenue' field is selected. Below the histogram, there is a script editor window containing Python code related to data processing.

Now, we can also go in here and say, "Hey, I want to toggle on the outliers." When I do that, you can see that now I have the stubble. You can see that there are observations that go beyond that range and that allows us to get a sense of how the data is distributed in each of these columns here. We could also customize the histogram if we wanted to by indicating. We're going to general hair, we can say, "Hey, beans instead of maybe 12 or so, I want 20 beans." I type in 20 and hit enter and you can see that it added more beans in there. Creating and editing charts, formatting charts and power BI is so simple, it's awesome. Let me show you one other type of chart that we can add in here. Let me just resize this, make it a little bit smaller. Make room for another chart, and this other type of chart is in our chart. You can see you can add in our scripts or Python scripts let me just click r. You'll get this prompt to enable it, because it could damage your computer, someone is doing something nefarious, so I'm going to go ahead and click enable. You can see it creates a placeholder here and there's just one item for this type of chart so far, and that is a value. Now, I could actually add in more than one value, but I'm just going to start with revenue. If I wanted, I could add in longitude or some other column as well.

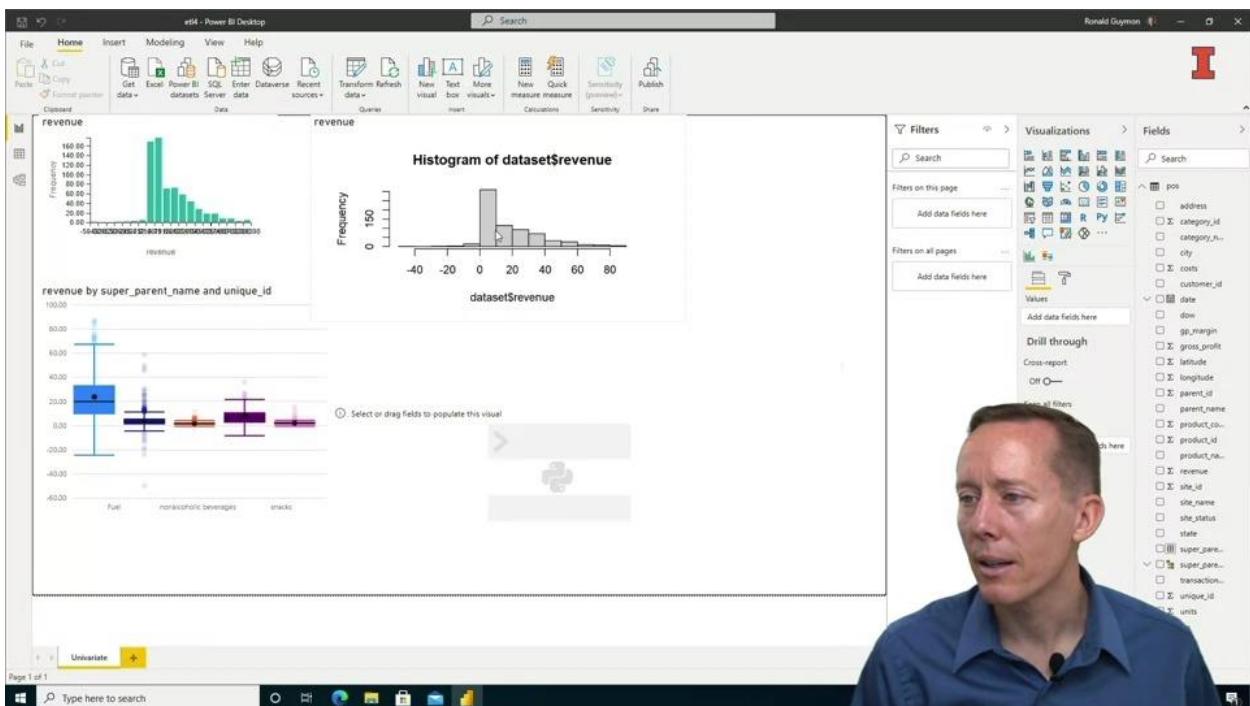


When I do that, it populates this r script editor and it's basically saying, "Hey, we're going to take this revenue, create a new data set and it will have revenue as a column in there." Now, I can click in there and I can type our code to create any type of plot I want. In this case, I'm going to type, I'll just do hist for a histogram and I will say data set and the column is revenue, and click run. You can see it doesn't look very pretty there,

The screenshot shows the Power BI r script editor. On the left, there's a sidebar with sections for 'revenue is (All)', 'Filters on this page', 'Filters on all pages', and 'Drill through'. On the right, there's a list of fields: address, category_id, category_n... (with a red 'I' icon), city, costs, customer_id, date, and dow. A context menu is open over the 'revenue' field, listing options like 'Remove field', 'Rename for this visual', 'Don't summarize' (which is highlighted with a grey background), 'Sum', 'Average', 'Minimum', 'Maximum', 'Count (Distinct)', 'Count', 'Standard deviation', 'Variance', 'Median', 'Show value as', 'New quick measure', and 'super_pare...'. The 'Don't summarize' option is currently selected.



and let's go ahead and explore why by clicking on this down arrow next to revenue. Sure enough, it is summarizing it all into one data point, one observation, so let's click on the don't summarize option.

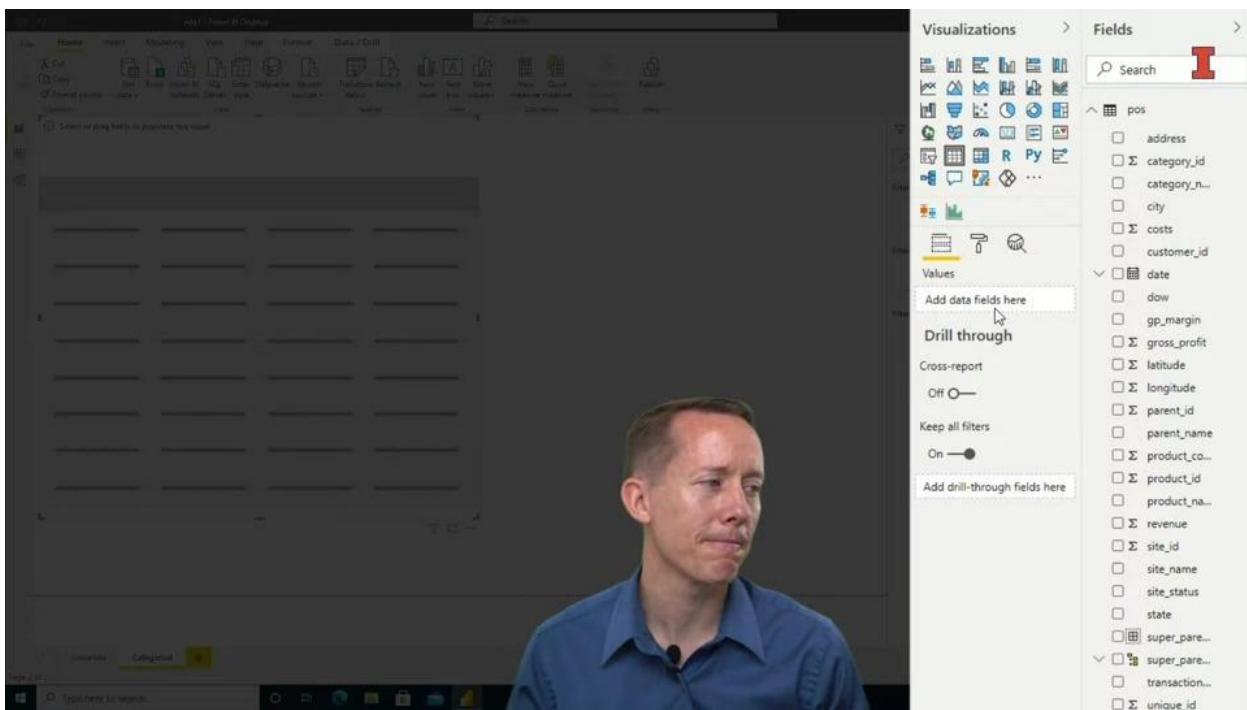


There we go, this looks like a much better histogram that's somewhat similar to this other histogram that was added. The point is, if there is a type of plot that you really like using from our Python, you can add that in. Now, one thing that you didn't see is that before adding in or being able to use our scripts in Python scripts, you need to make sure that you have these loaded on your machine. If you try to add in something for Python here, it will tell you, "Hey, first make sure it's loaded on your machine." If it is, and you probably won't get that prompt, but if it isn't, then you'll need to add it in and there will be help that will allow you to point you into how to add that software. You can see that creating visualizations is very simple. We looked at some uni-variate visualizations, you can see it's very easy to resize them. We can also move them around on the canvas here and display them in a variety of different ways. We can get rid of them by clicking on it, clicking on these three dots and clicking remove or simply hitting the backspace or delete key. It's so intuitive, creating charts in power BI. Those are introduction to creating charts as well as how to create uni-variate charts.

[Lesson 2-3.2 EDA 2 - Univariate, Bivariate, and Multivariate Plots with Categorical](#)



In this video, we'll talk about exploratory data analysis that focuses on categorical columns of data. When exploring categorical data, it's often useful to identify the cardinality or the number of unique values. We can do that in the Power Query Editor, but we can also do that with visualizations. In this lesson, we will first talk about data visualization tools for univariate data analysis with categorical columns. But we will also show you bivariate data analysis in which you analyze two columns of data together.



I'm going to open up Power BI here and open up this project I've been working on eda1. When I open it, I have to identify that I want to enable script visuals, meaning I have some visualizations from are in here. I'll click that. Here is the tab that has the univariate analysis on it. Let's go ahead and create a new tab at the bottom. Click that "Plus Button" and we'll name this categorical. Let's start by analyzing the super parent name in the point of sale data, so we've got the super parent name here. We need to decide how we want to analyze it. What I want to do is start with the table. Click on table. This adds a table visualization to my canvas here and I'll just drag super parent name to the values.

This screenshot shows a categorical field named 'super_parent_name' in Power BI Desktop. The field has five distinct values: 'alcoholic beverages', 'Fuel', 'nonalcoholic beverages', 'Other', and 'snacks'. The 'snacks' value is currently selected, as indicated by a mouse cursor pointing at its corresponding row in the list.

This is a categorical field because it's not numeric data. You can see that we've got five different super parent names, alcoholic beverages, fuel, non-alcoholic beverages, other and snacks.

That's the first thing, is just looking at the number of different values in a column of data. This is relatively low cardinality because it only has five unique values.

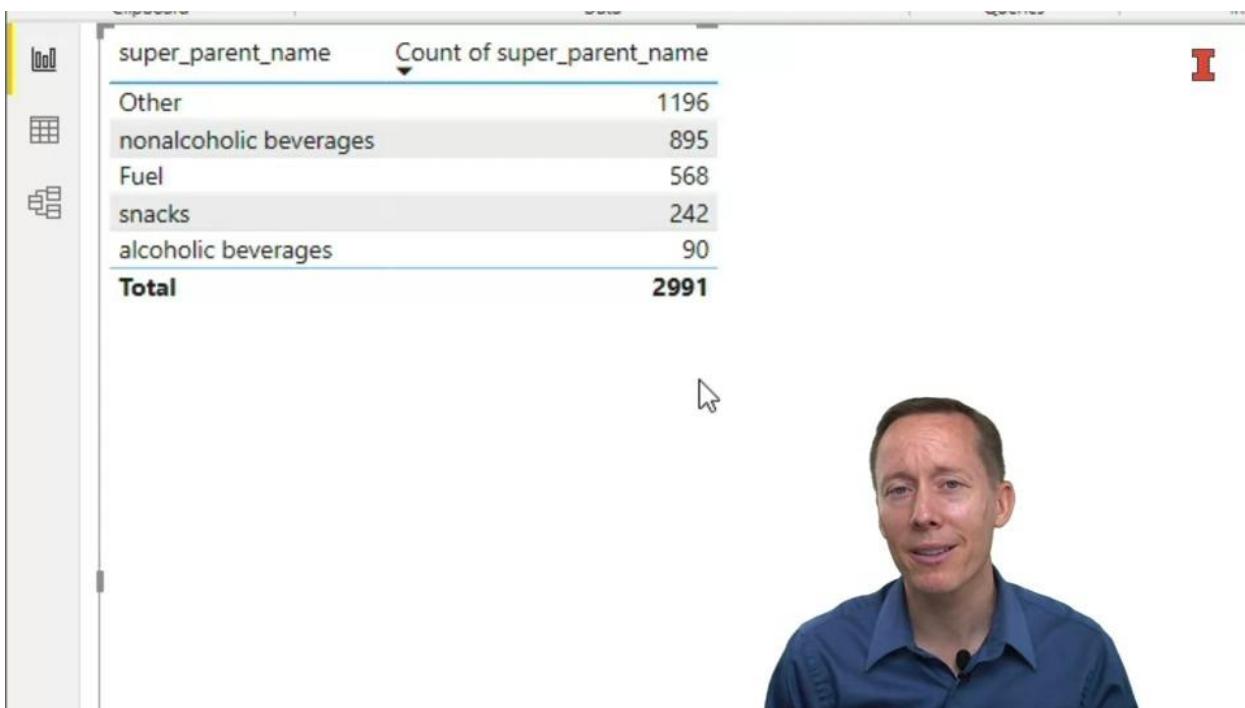
This screenshot shows a large list of product names in a Categorical field in Power BI Desktop. The list includes various items such as '1 JOHNS WEDGE', '200Z COFFEE REFILL', and '3 HOUR EX STRENGTH GRP 1.93OZ'. The right side of the screen shows the Power BI Fields pane, which lists many fields categorized under various entities like address, category_id, city, costs, customer_id, state, etc.

If we were to compare this to product name, let's take out super parent name, compare it to product name. You can see by the scroll bar here just scrolling down, we've got a lot of different values here. That can be overwhelming.

The screenshot shows a Microsoft Power BI desktop application. On the left, there is a data grid with two columns: 'super_parent_name' and 'parent_name'. The data includes rows like 'alcoholic beverages', 'Fuel', and 'snacks'. To the right of the grid is a 'Visualizations' pane with various chart icons. Below it is a 'Fields' pane containing a search bar and a tree view of fields under a 'pos' category. A context menu is open over the 'super_parent_name' field in the values list, with 'First' highlighted.

super_parent_name	parent_name
alcoholic beverages	alcoholic beverages
Fuel	Fuel
nonalcoholic beverages	nonalcoholic beverages
Other	Other
snacks	snacks

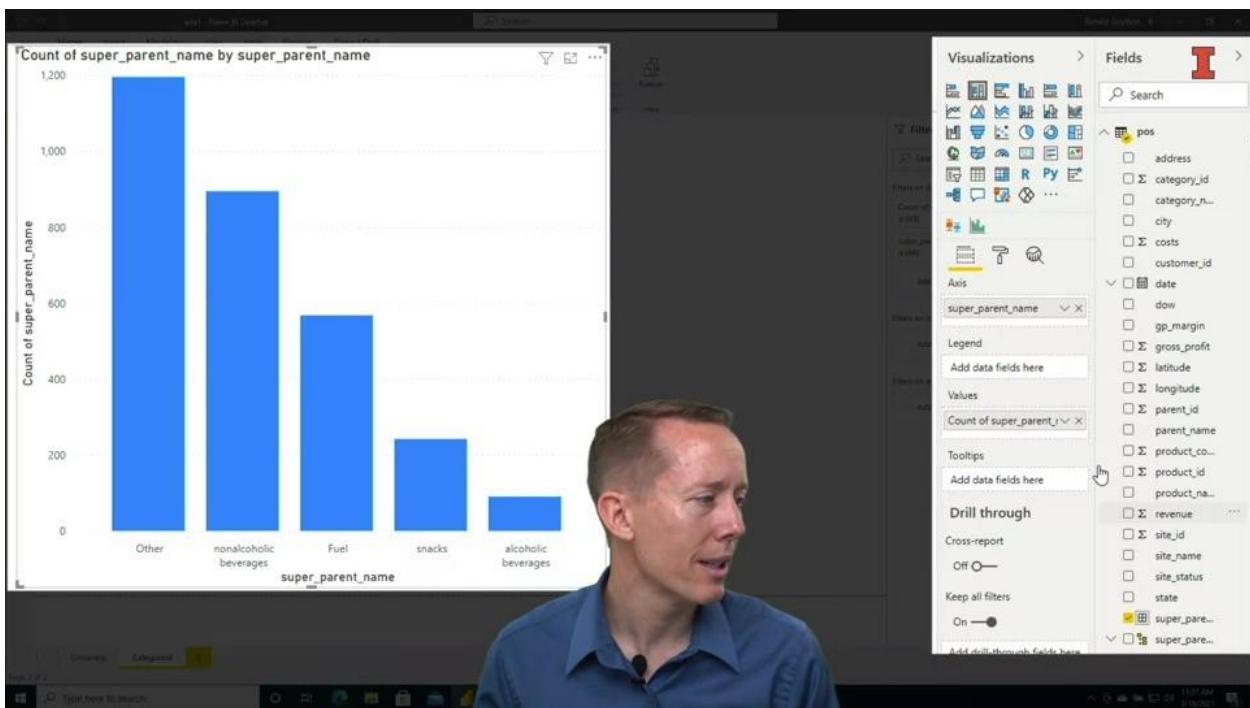
Let's start with a low cardinality category, which is super parent name. Once we have that, we can start by just looking at the count of observations for each of those values. The way we can do that is by simply dragging super parent name again to the values box, and we can see now it's just two columns with the same thing. But if you click the drop-down arrow next to the second one, we can indicate that different ways of summarizing this data and what we want is count. Click on "Count".



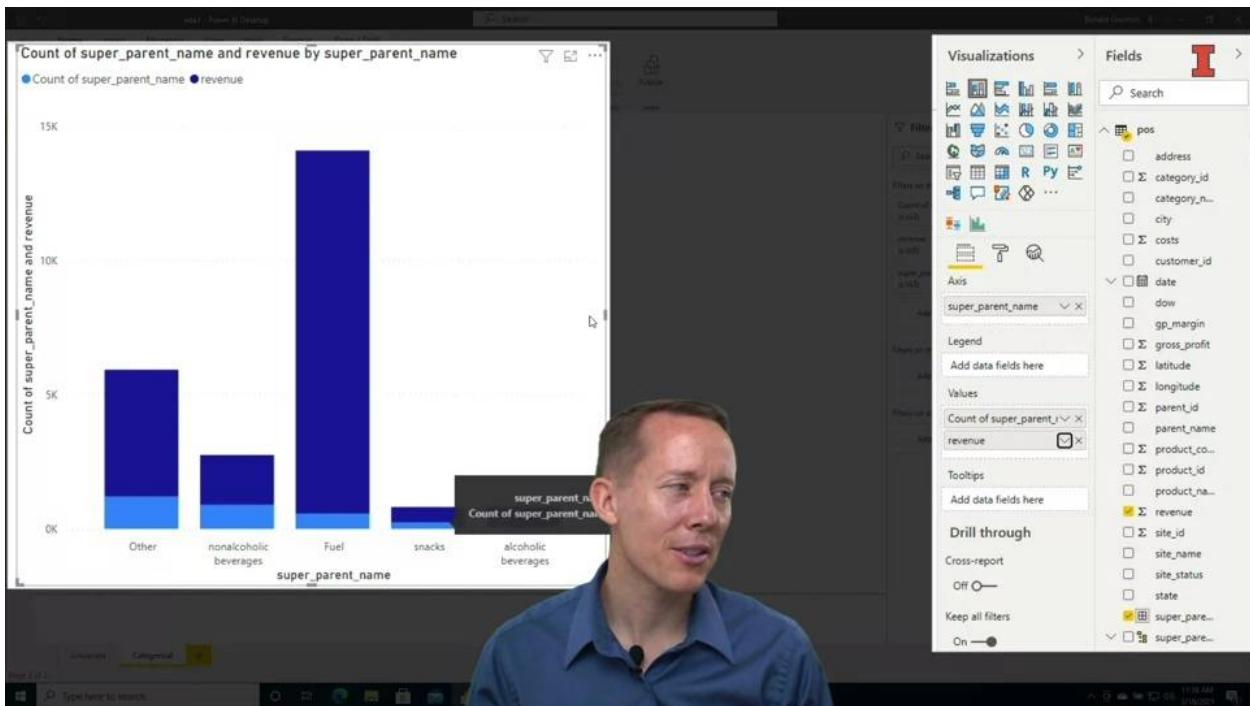
Now this table updates so that we can see that we have 90 observations that are alcoholic beverages, 568 that are fuel and so on, and we can see the total there. Basically in this, If you remember, each line of data is a line item that is purchased in a transaction. We've got 2991 items. We started with 3000 then we got rid of some outliers. You can see this is by default, ordered alphabetically. We could order it, if you it over the column name and click on the "Arrow", we can switch the ordering so that it's in a descending alphabetical order, or we can order it by the count of super parent name, which might be more useful in this case here. We can see the other makes up most of the observations. Then non-alcoholic beverages and fuel and so forth. This, I would say, is a univariate way to analyze categorical data.

The screenshot shows the Power BI Data Explorer interface. On the left, there's a sidebar with icons for different data types like tables, charts, and maps. Below that is a list of columns under the heading 'pos'. Some columns have checkmarks indicating they are selected: 'revenue' (sum), 'super_parent_name' (count), and 'super_parent_name' (average). Other columns listed include 'address', 'category_id', 'category_n...', 'city', 'costs', 'customer_id', 'date', 'dow', 'gp_margin', 'gross_profit', 'latitude', 'longitude', 'parent_id', 'parent_name', 'product_co...', 'product_id', 'product_na...', 'site_id', 'site_name', 'site_status', 'state', 'super_pare...', and 'super_pare...'. There are also sections for 'Filters on this visual', 'Values', 'Drill through', 'Cross-report', and 'Keep all filters'.

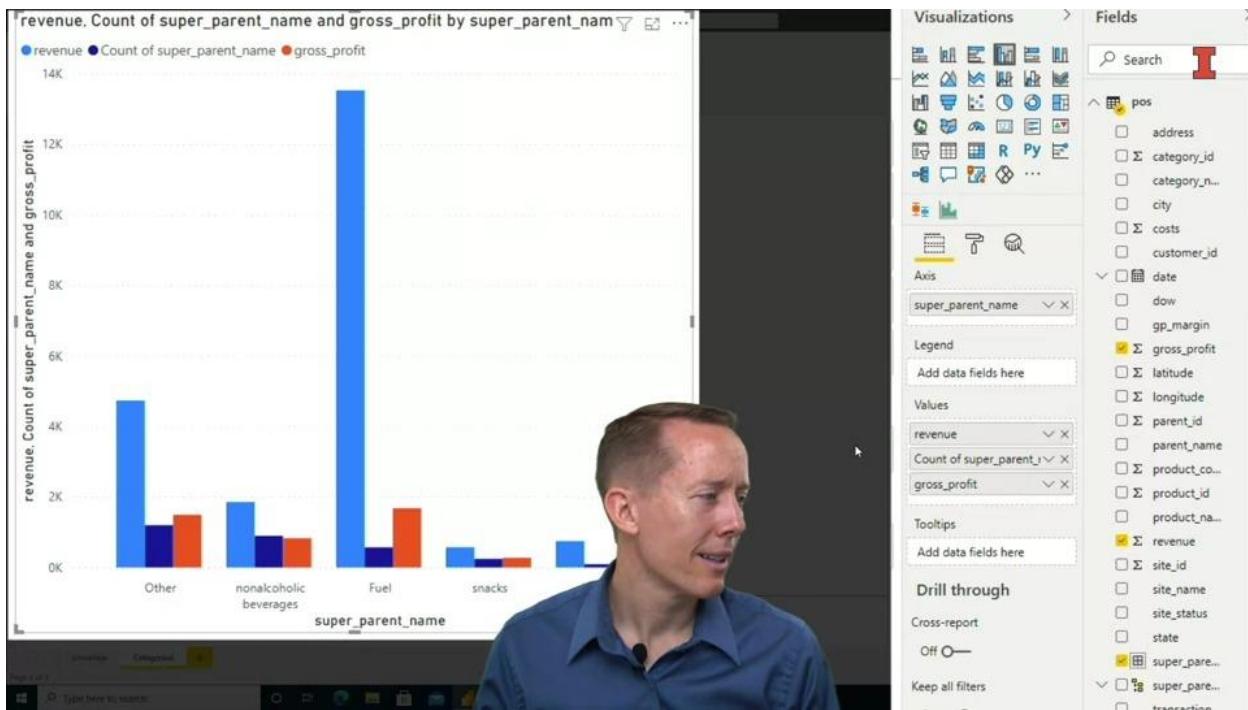
If we wanted to make this a bivariate analysis, we could drag in some other column and let's drag in revenue, since that's something that we're interested in. Now you can see that we've got a column that summarizes the revenue. You can see if you click on revenue, if you want to look at the average, we could use average instead, aggregate it to see on average what are people spending for each of these different super parent names. That's pretty cool. I'm going to go ahead and just get out of the average. Click that out of the values and start with this.



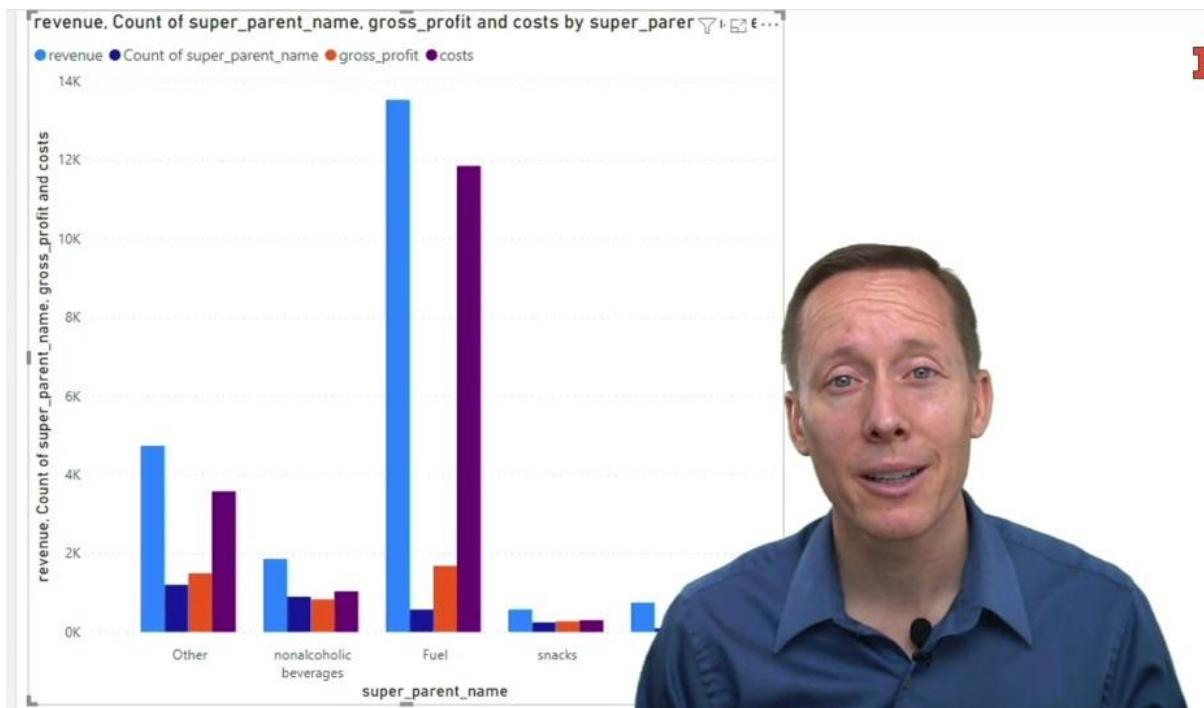
Now let's convert this table into a bar chart. The reason is because when we look at relationships between numeric values, they're much easier to see using visualizations. So I'm going to just click on this stacked bar chart. You can see it's horizontal, and now it shows us very quickly that other is about a fifth more than non-alcoholic beverages. We can convert this to a stacked or a vertical content very easily. You can see it's so simple to rotate among different types of charts here. What if I do want to look at revenue and count together? I can very easily do that by dragging revenue to this values box here.



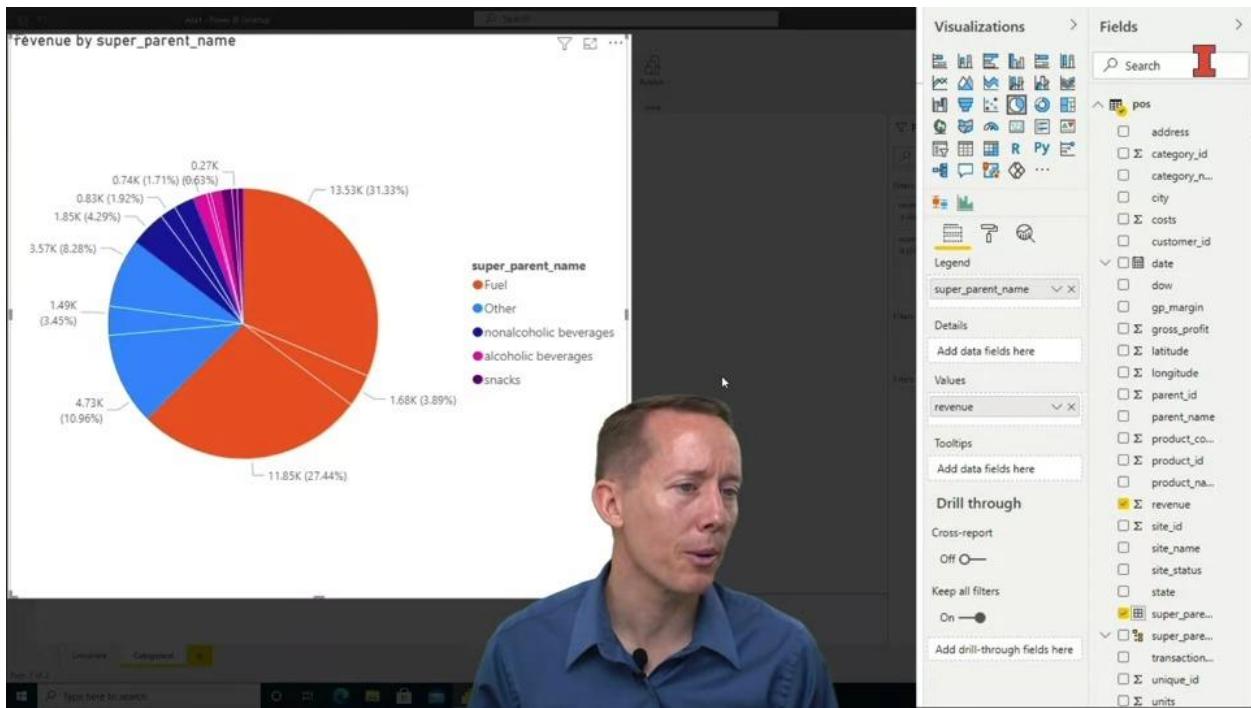
Now I've got a stacked bar chart where the light blue is the count and the dark blue is the sum of the revenue. Again, we could change these to aggregate it in different ways if we wanted to look at average amount spent on revenue. You can see the stacked bar charts can be helpful. I don't know that in this case, stacking count and revenue on top of each other is useful. Let's say, hey, we want to look at a grouped bar chart, clustered bar chart.



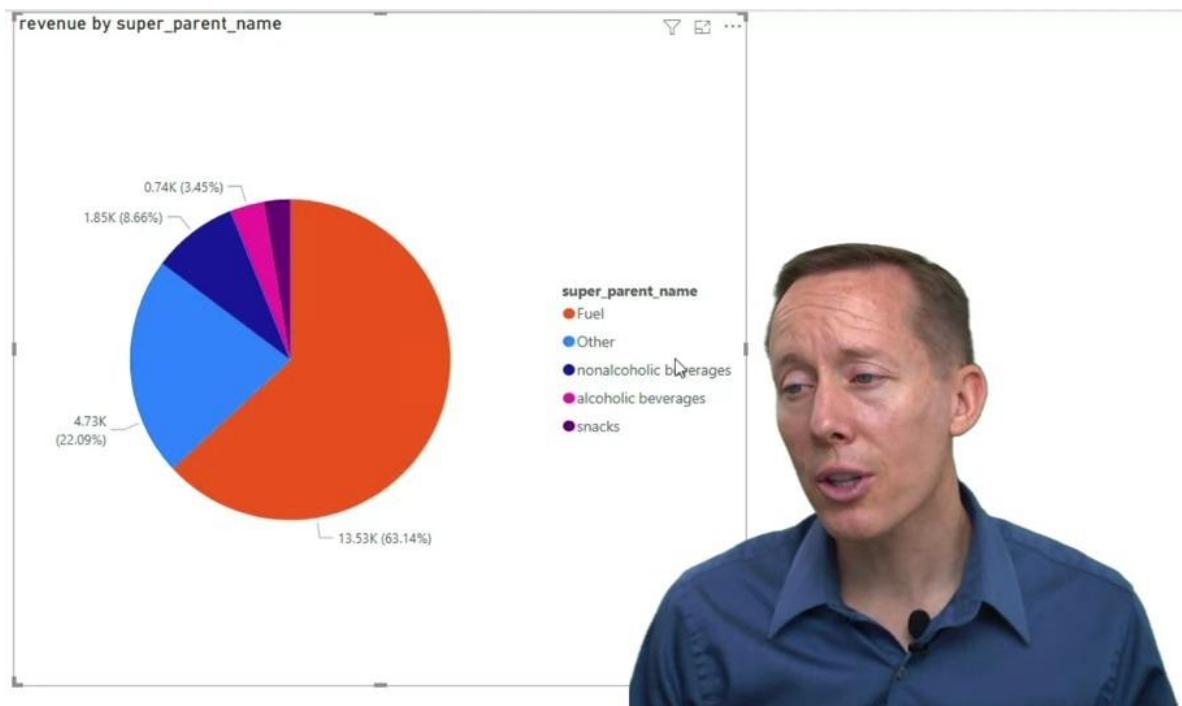
I can navigate that way to cluster bar chart and very quickly convert it to a clustered bar chart. If I want to change the order so that I have revenue before the count, I can just drag the revenue pill up above the count pill and that changes the order there. I can obviously add in a lot of other columns in here, so I can look at gross profit, I could look at the costs.



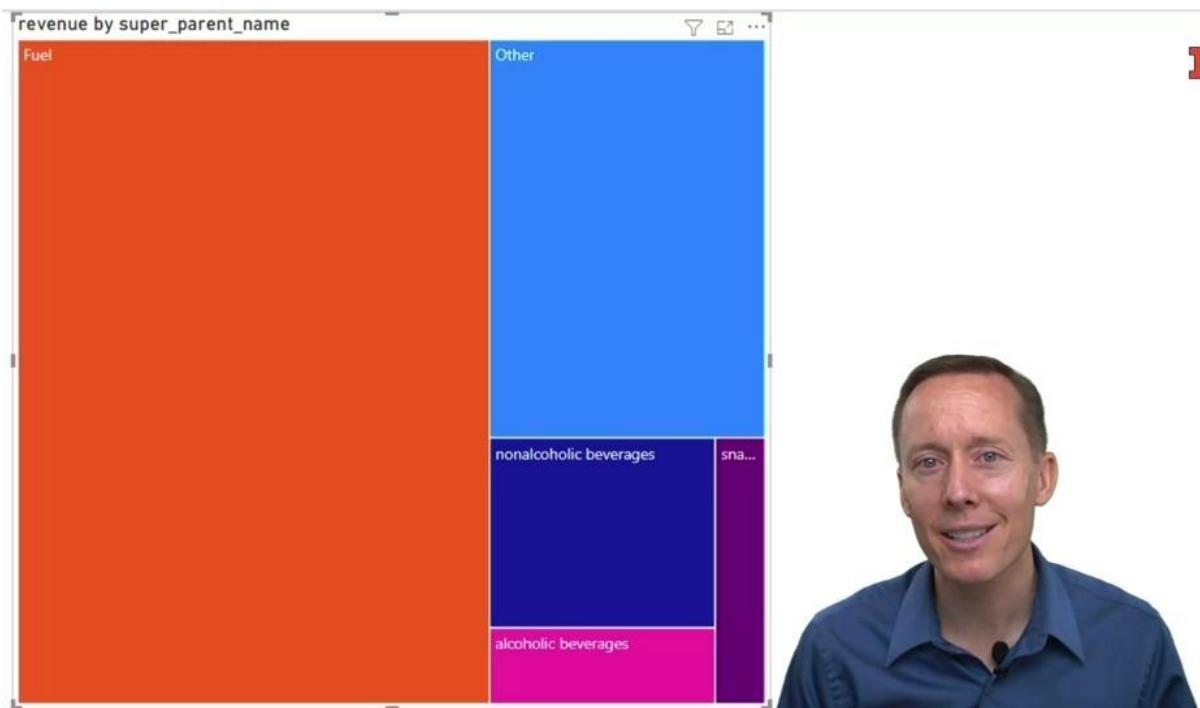
Pretty cool how we can see all of that data. I encourage you to think about the use of these visualizations. They can very easily become overwhelming and may communicate more than what you intend to communicate. The hope is that there are some intraocular traumatic impact. There's a conclusion that hits you right between the eyes, oftentimes we want to communicate that. Although if you want something for people to explore and really take time to look at, maybe you do want something that's busy like this. So I encourage you to play around with the different types of charts and the different formats for these charts and to think about how quickly each chart is communicating what you hope it's communicating.



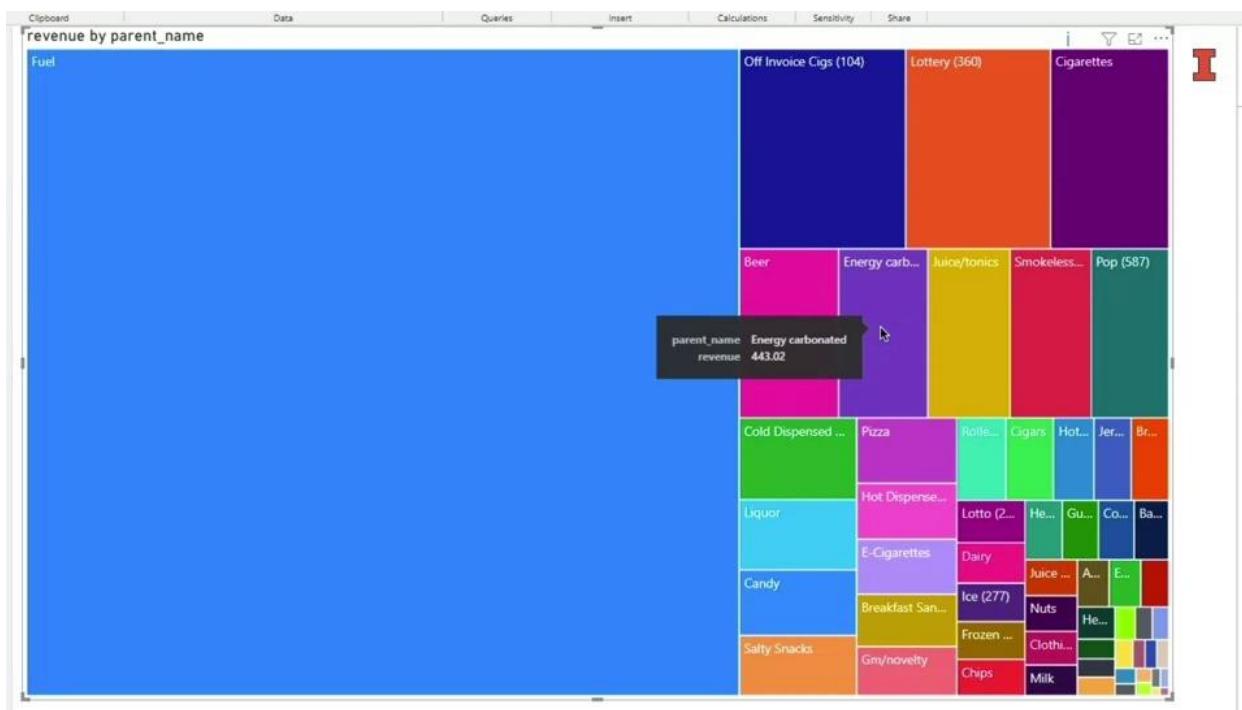
For exploratory data analysis, it may be useful to rotate between different types of charts. Bar charts are helpful to see relative amounts, a pie chart is often useful to show percentages. This pie chart is very messy. Pie charts aren't usually intended to show more than one different type of value. Let's narrow it down to just revenue. We'll get rid of these other ones here.



This is helpful. This allows us to see that fuel makes up 63.14 percent of the revenue or \$13.530. Pie charts are very helpful, my rule of thumb is six or fewer slices or categories in a pie chart can be useful, especially when you're looking at the percent of the total. The downside with pie charts is that they don't display lots of pie slices very well. Another downside is that it's typically not ideal to have the legend separate from the slice of the pie. We would like real close so your eyes don't have to keep moving back and forth.

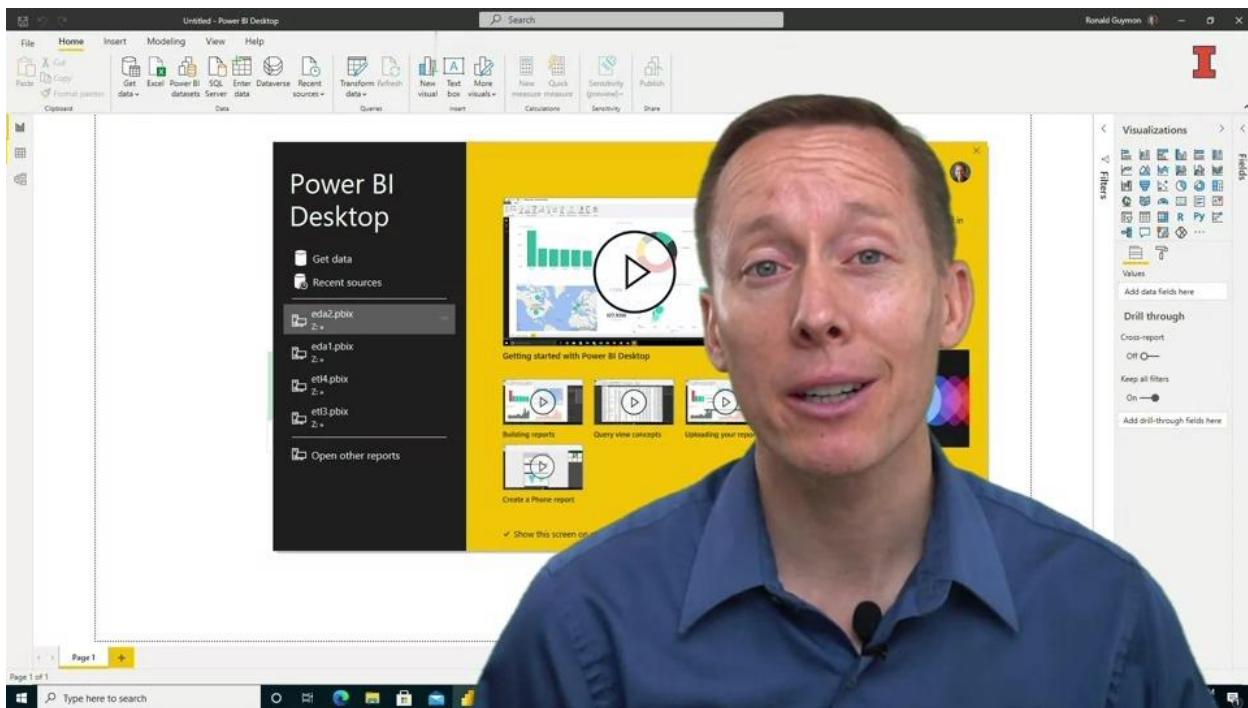


Another type of chart that you may be aware of is called a tree map. A tree map is like a pie chart, but in rectangular form. The thing I like about tree maps is that often the pieces of the tree map have the label on them. They are also, I think, much more helpful for displaying categorical data that have more than six different values in them.

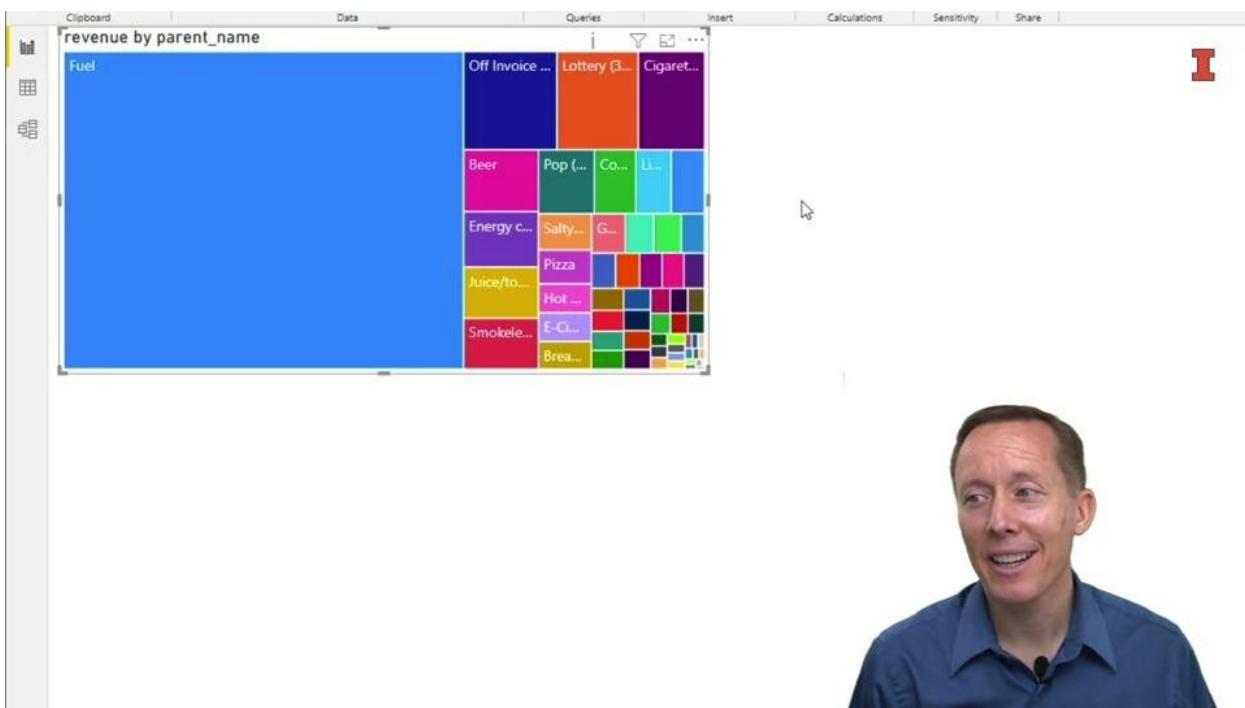


If we wanted to look at maybe instead of super parent name, we want to look at parent name, which has 60 different unique values in there. You can see how quick it is to update a chart and we can make this bigger. Now we can see that fuel makes up about the majority, but it's easy to see some of these other parent names that are populating the total revenue. That's how you can do some exploratory data analysis with categorical data in Power BI.

Lesson 2-3.3 EDA 3: Filters, Slicers, and Drill Through



In this video, we want to show you some tools for implementing Schneiderman's mantra, which is overview first, Zoom filter, then details on demand. Specifically, we'll show you how easy it is to use charts and tables as filters and also how to use the filter pane slicers and finally, how to use the drill-down features, in Power Bi for hierarchical data.



The first thing I will do to start with is load this project I've been working on EDA2, and I will go ahead and click on this and able script visuals. You can see in this report, I've got a univariate tab and a categorical tab and the univariate tab has histograms and a box and whisker plot, the categorical tab right now just has a tree map. First of all, I'm going to resize this tree map, and then I will add in a table,

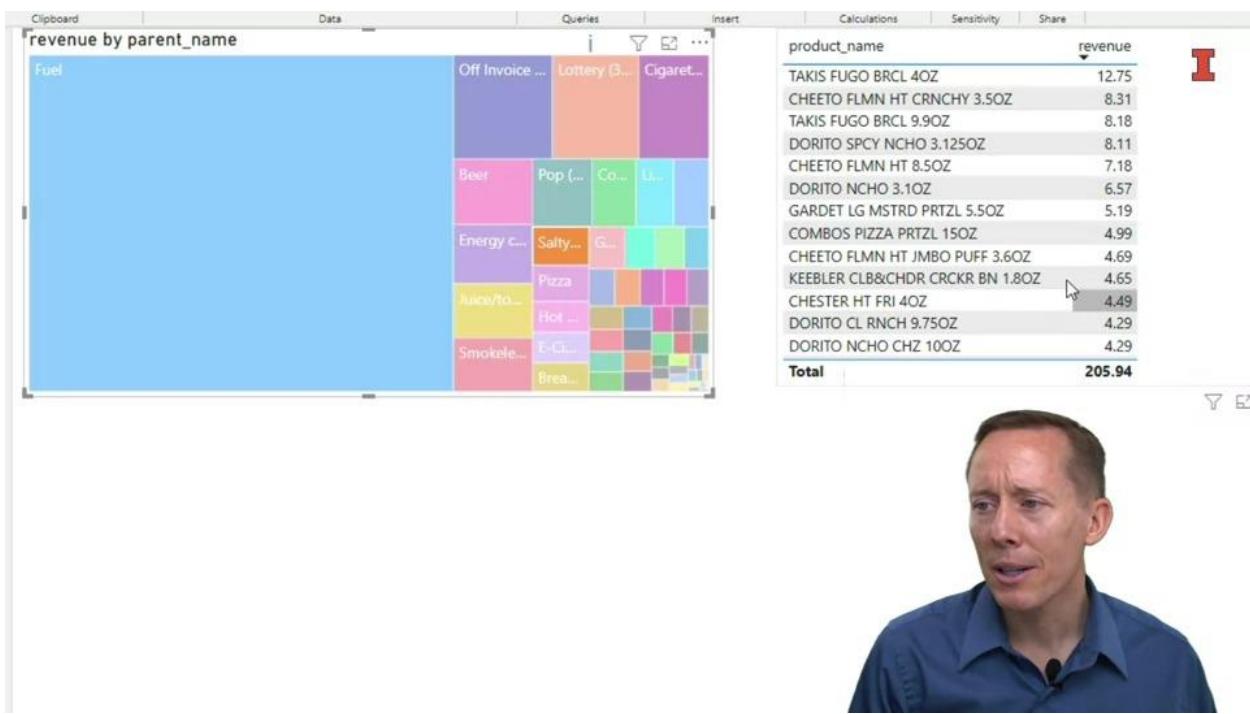
The image shows a data analysis interface with a sidebar on the right and a main content area. In the sidebar, under 'Fields', the 'product_name' field is selected, indicated by a yellow highlight. Other fields listed include address, category_id, category_name, city, costs, customer_id, date, done, gpm_margin, gross_profit, latitude, longitude, parent_id, parent_name, product_code, product_id, revenue, site_id, site_name, user_status, state, and super_page. In the main content area, there is a list of products under the heading 'product_name'. Some items in the list are highlighted with gray bars, such as '1000Z MUG 2017' and '1000Z COFFEE REFILL'. A cursor is visible over the '1000Z MUG 2017' item.

so I'm going to just click on this table icon in the visualizations paint, and I'm going to resize this and move it up next to the tree map and in this table, values filled, I'm going to add in product name. This gives me a list of all the different products that are sold or that are in the data center,

The screenshot shows a Power BI interface with a table visualization. The table has 'product_name' in the first column and 'revenue' in the second column. The data includes various products like G87E10S, No Lead Plus, G85E10S, etc., with their respective revenues. A 'Total' row at the bottom shows a revenue of 21,420.54. To the right of the table is a video player window showing a man in a blue shirt. The video player has a 'Filters or' section with two dropdown menus, both currently set to 'A'.

product_name	revenue
G87E10S	6,309.89
No Lead Plus	1,635.62
G85E10S	1,385.27
DIESELBIO	768.10
DIESEL	559.05
G91E10P	544.25
GE85	411.87
G87E00M	241.68
G91E00P	221.61
G87E10M-B34	190.75
G87E10M	182.71
G89E10M-B50	171.62
LOTTO	159.50
Total	21,420.54

and I want to see what the total revenue is for each product, so I will also drag revenue, to the value's filled. You can see maybe it doesn't appear to be summarized, but it is, if you click on the download extra revenue, you can see that it's summing up the total amount of revenue for each of these product names, and it might make more sense to sort this in terms of, the highest revenue items to the lowest revenue items. The first thing I want to show you is that inherent in power BI plots is the ability to filter.



If I click on one square in this tree map, let's say salty snacks. Notice how this table updates and includes the product names only for the salty snacks, and I can see the total is 2,005 dollars 94 cents in salty snacks. If I'm ever curious as to what is filtering this table,

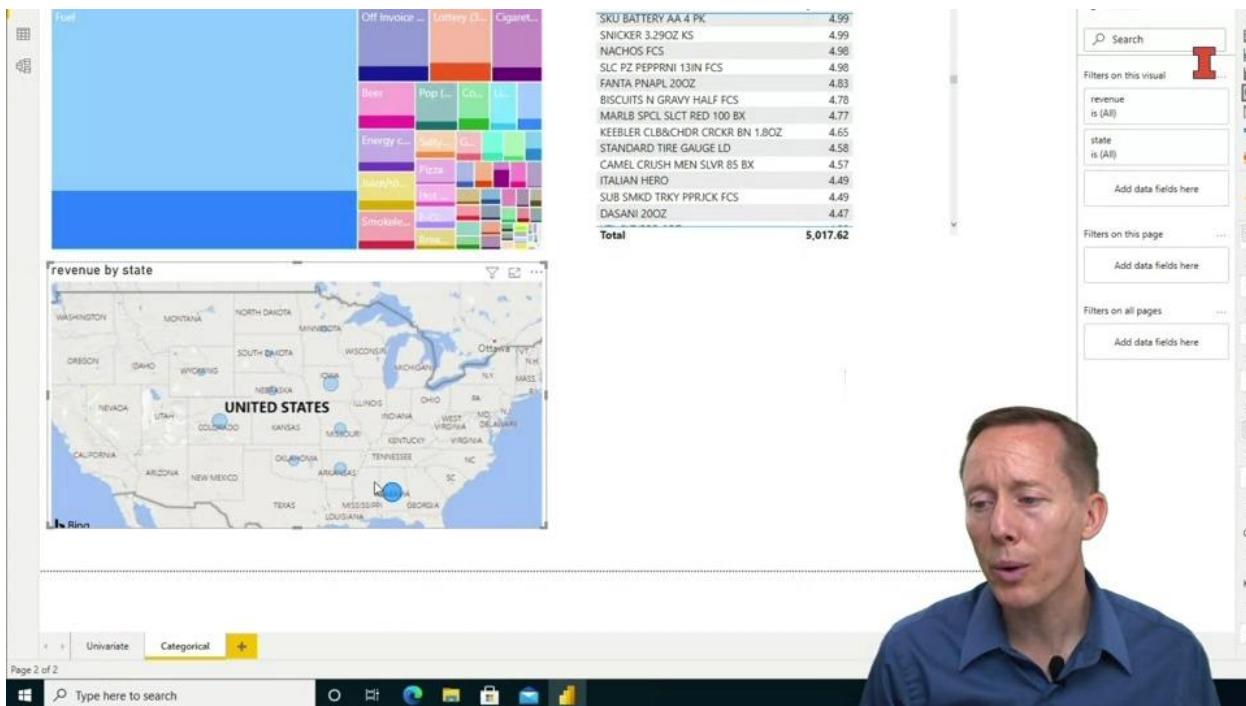
The screenshot shows the Power BI interface with the filters pane open. It displays filters for "product_name" (All) and "revenue" (All). A tooltip "Filters and slicers affecting this visual" is shown, indicating that the filter "Salty Snacks (parent_name)" is applied.

product_name	revenue
TAKIS FUGO BRCL 4OZ	12.75
CHEETO FLMN HT CRNCHY 3.5OZ	8.31
TAKIS FUGO BRCL 9.9OZ	8.18
DORITO SPCY NCHO 3.125OZ	8.11
CHEETO FLMN HT 8.5OZ	7.18
DORITO NCHO 3.1OZ	6.57
GARDET LG MSTRD PRTZL 5.5OZ	5.19
COMBOS PIZZA PRTZL 15OZ	4.99
CHEETO FLMN HT JMBO PUFF 3.6OZ	4.69
KEEBLER CLB&CHDR CRCKR BN 1.8OZ	4.65
CHESTER HT FRI 4OZ	4.49
DORITO CL RNCH 9.75OZ	4.29
DORITO NCHO CHZ 10OZ	4.29
Total	205.94

I can hover on this, click on the table and then hover over it and I see the funnel button, and it will tell me what filters are currently applied. I'll go ahead and click salty snacks again and click on the funnel and I can see that salty snacks is being used to filter this table.

product_name	revenue
'S PIN CHRG SYTRC CBLPSKU	17.99
SUB SMOD TRK PPLRCK FCS	17.95
FIL 1LT	17.47
200Z COFFEE REFILL	16.99
DUAL USB W/ MICRO CBL	16.99
FLAT CAP	16.99
HOTTESS PWD SGR MNNT BAG 10.5OZ	16.96
COKE 200Z	16.73
MARLB SPL BLND MEN BLK BX	16.69
MARLB SLVR 100 BX	16.40
WTR 1LT	16.39
WHL PZ 1 TOPPING 13IN PCS	16.00
BUSCH LT 18PK CAN	15.99
BUCKLER TINNED CAN	15.99
Total	21,420.54

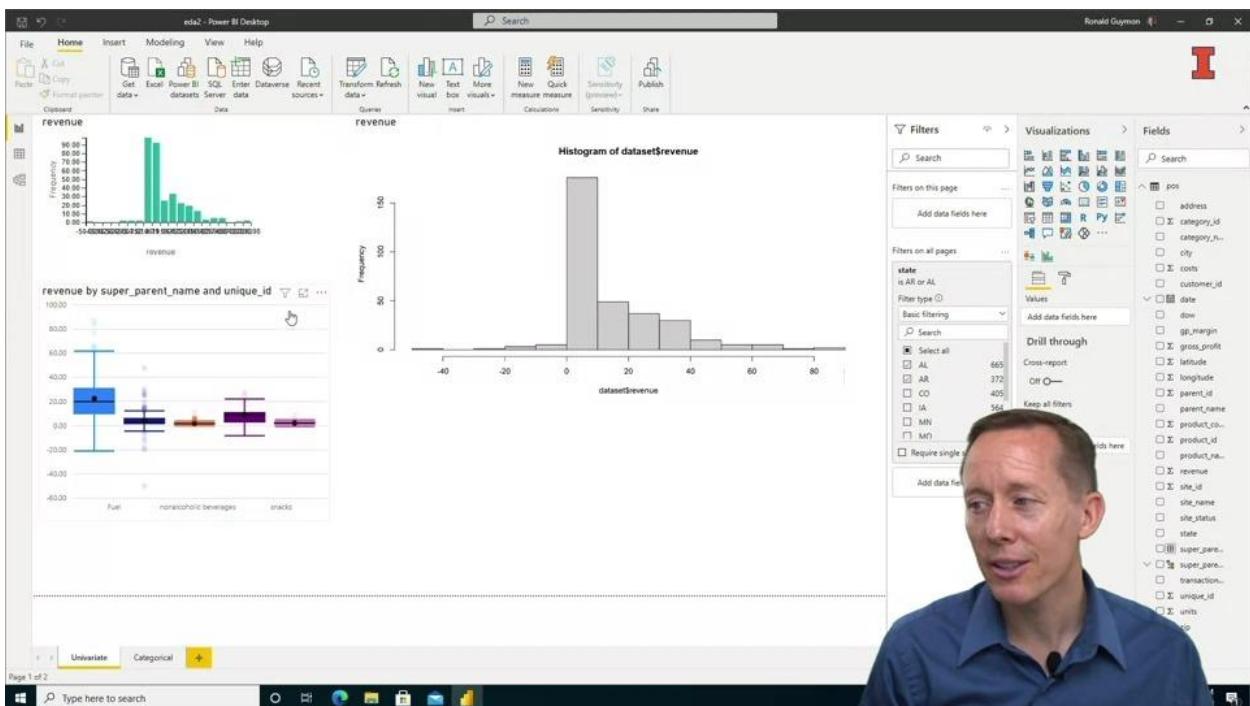
Now, I can also use the table to filter the tree map as well. I'll turn off the filter on the tree map by just clicking on the square that had been clicked on. Let's go ahead, and I will click on, coffee refill and this is a very small segment of the hot drinks item there, so it highlights it for me and I actually like it. Sometimes it's hard to find if it's a small item, but I like that it leaves the context of everything else there rather than just show me one rectangle. I really like how it leaves the context in there, so anyway, filtering is really very easy and powered BI, but there are some other ways to filter that. I'd like to illustrate for you. To do that, let's go ahead and turn off this filter and let's add in a new chart and let's add in a map.



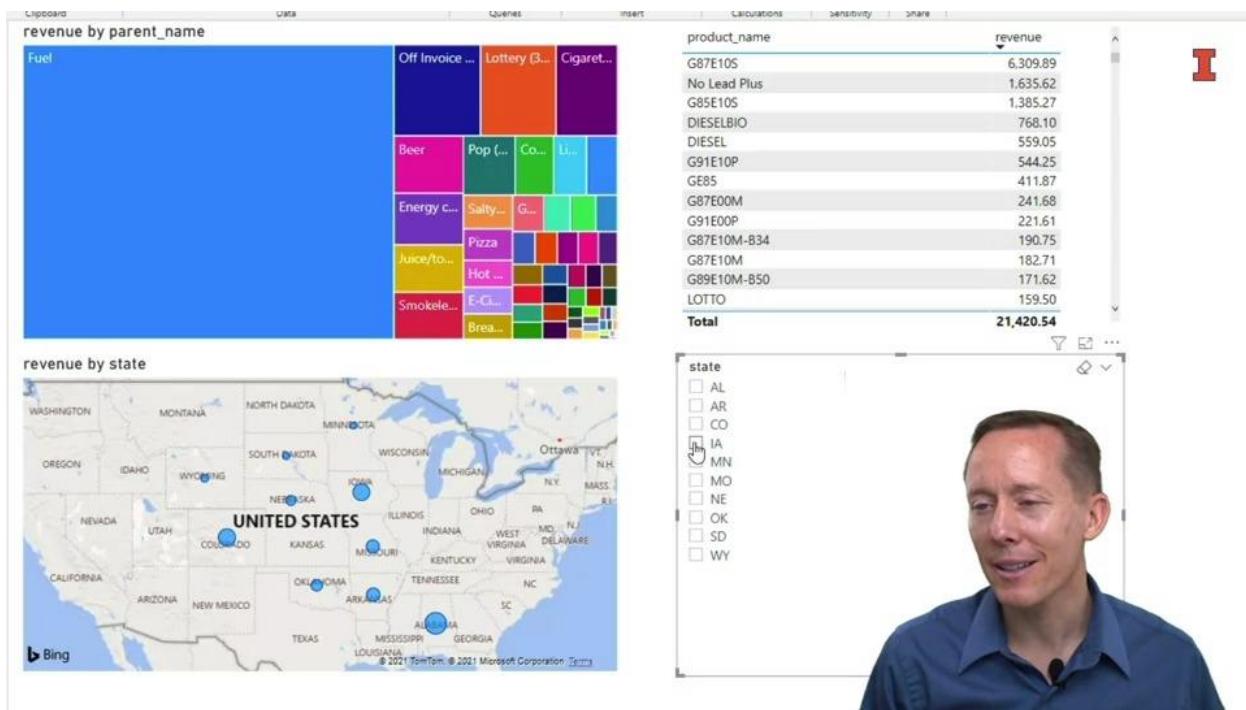
Now, power bi is really good at recognizing locations by name. You can also add in latitude and longitude coordinates, but let's go ahead and just use state, to add into the location box for this map and see what happens. Beautiful. You can see that we've got a dot on each of these states, so it recognizes the states and where it should go on a map. Now, I can resize these dots. Let's say I want to see which state sells the most. I can drag revenue to the size field. Now, I can very easily see that it looks like Alabama, [inaudible] Colorado has the most sales, and now I could use this as a filter for the other plots as well, if I want to see the total sales just for Alabama, also I click on that, the tree map and the table update to show me the sales just in Alabama. Again, I really like how the tree map includes the context of everything else. I can see how much of the total sells fuel cells Alabama makes up out of all the sales and all these other apparent names as well. Now, what if I want to filter the visualizations, not just on this tab, but also on the univariate tab as well?

The screenshot shows a Power BI desktop interface with a dashboard containing two visualizations: a treemap chart titled "Revenue by parent_name" and a map titled "Revenue by state". A video overlay of a man speaking is overlaid on the dashboard. The Power BI interface shows the "Filters" pane on the right, which includes a search bar, filters for "revenue" (is (All)), "state" (is (All)), and "Filter type" set to "Basic filtering". The "Fields" pane also shows a search bar and lists various fields like "pos", "address", "category_id", etc., with "revenue" checked.

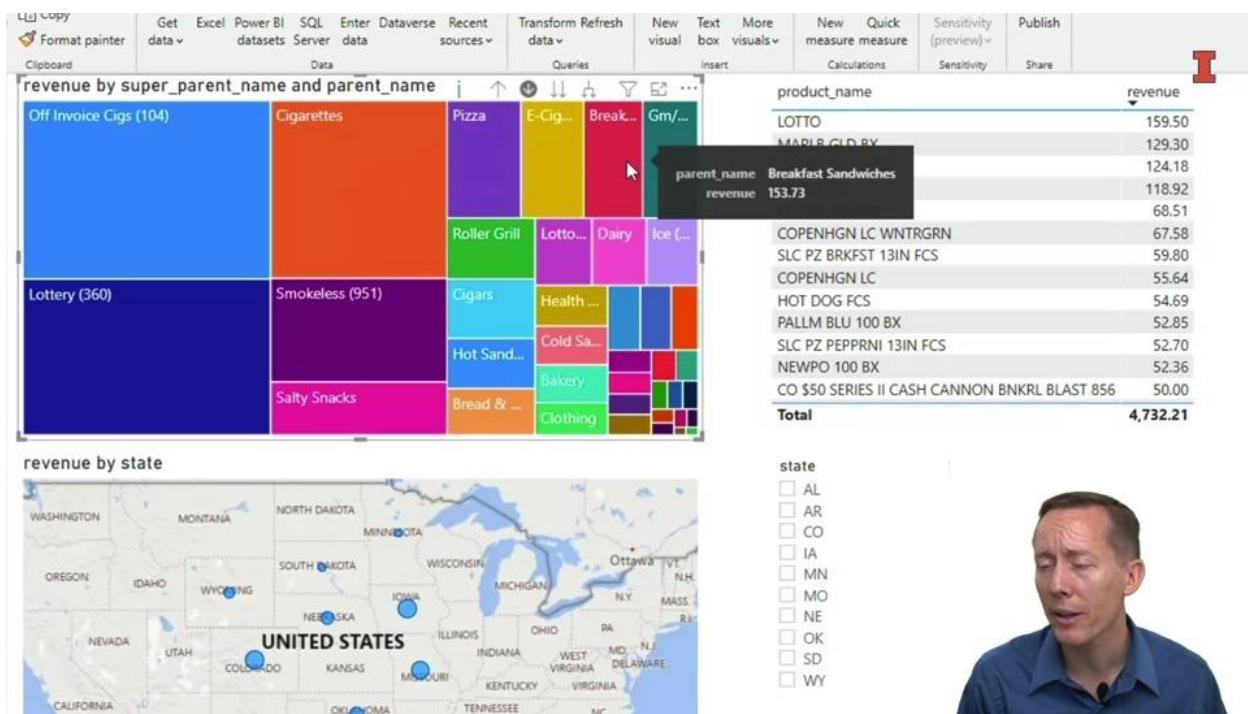
One way I can do that, is by using this filters bar over here and if I drag state into this first box here, you'll see this will allow me to filter only on a particular visual, which is the map here.



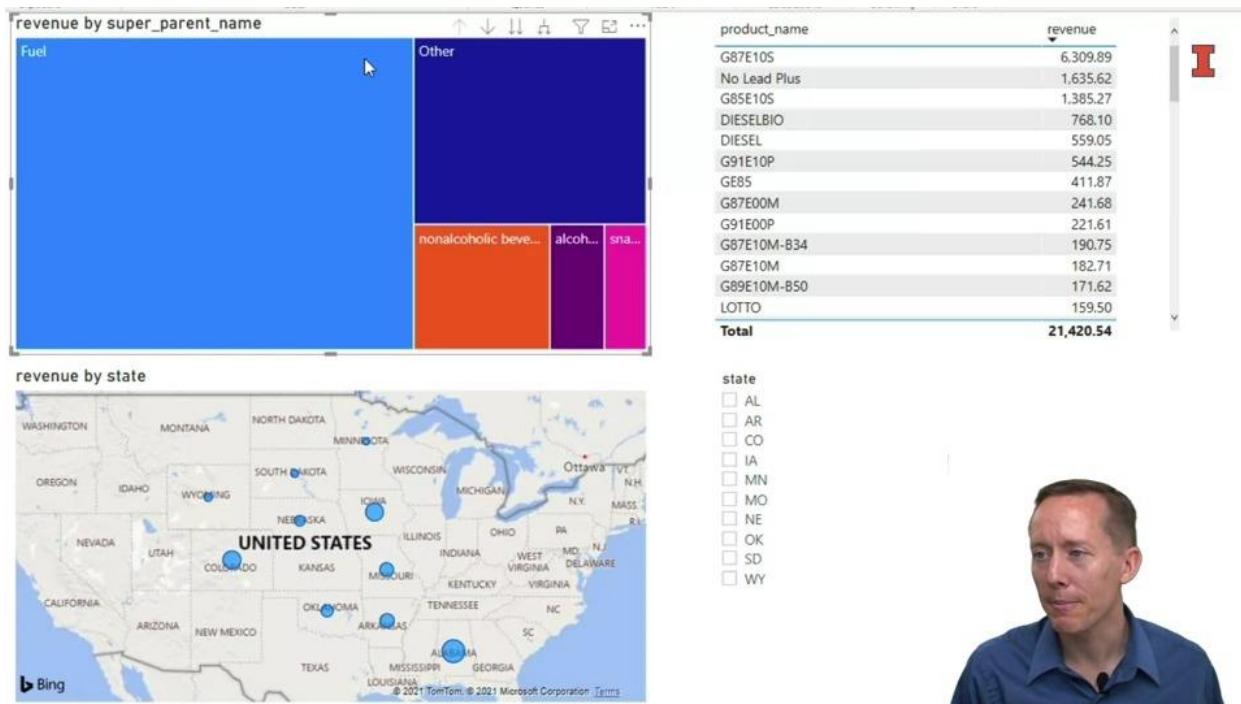
I can just click on Ar for Arkansas and you can see it zooms in on Arkansas and filters it just to that. I'll remove that filter and then drag state to filters on this page and I can just click on Arkansas again, now that not only filtered the map, but the tree map and the table as well. Finally, I could. You state to filter the observations on all the pages. If I click on Arkansas here and go over to the Univariate tab, you probably can't tell but these have been filtered. Let me click on another state here AL for Alabama. Notice how these charts also are updated. That's another way that you can explore the data using filters.



Now, another way to filter the data is by using slicers. If you've ever used Excel's pivot tables and added slicers to those, it's very much the same way. I'll go ahead and remove this filter from the filter bar and I'll go back to the Categorical tab, make sure I haven't selected any of the visualizations and I'll click on this slicer icon here. Now, I can drag state into the field and now I can essentially have that what was in the Filter bar on the right here within the chart itself. This is helpful if we want to publish a report for others to use and that filter bar for me is more like a design mode, whereas this is kind of production where this will allow other users who are not able to change the layout of this report to filter the data. Anyway, we can filter in those three different ways, using the visualizations and tables, using the filter bar, and using slicers. The last thing I want to show you in this lesson is how you can use the drill-down functionality and Power BI. This is really exciting. I think this is super cool. To illustrate this, I'm going to play with the treemap



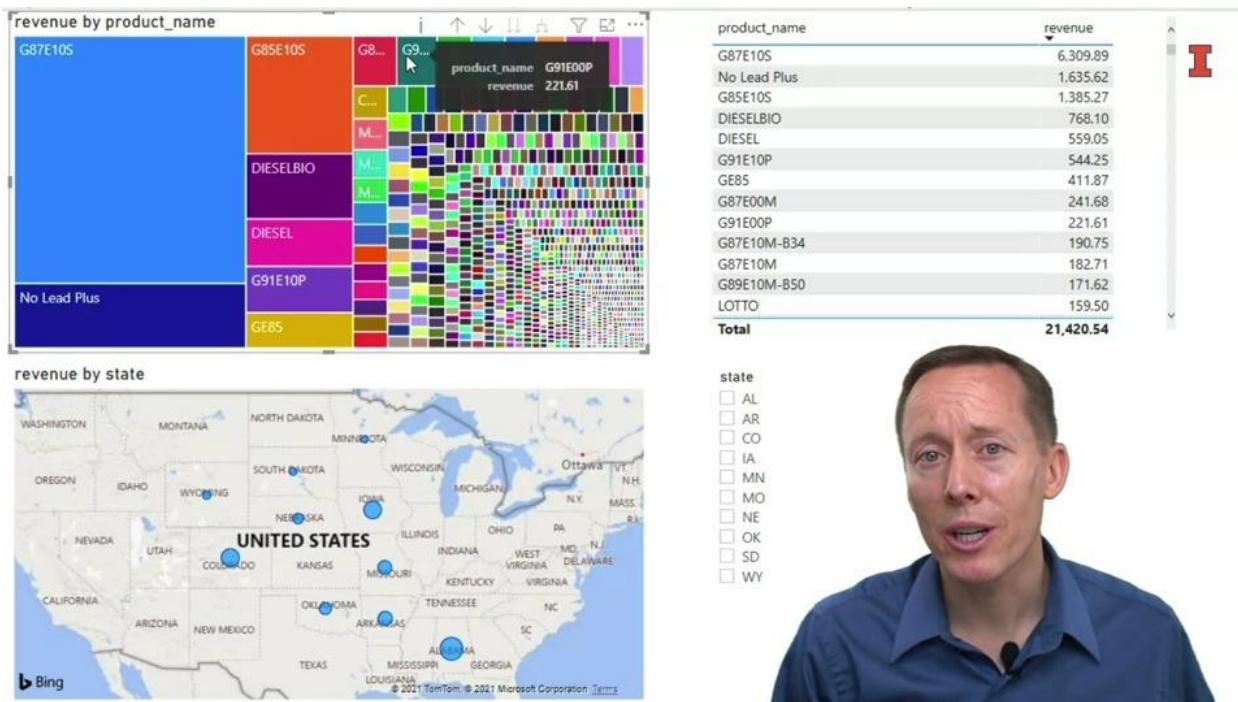
and I'm going to remove the parent name and replace it with this super parent name hierarchy. Now, if you remember the super parent name hierarchy has not just superpower name, but then below that, the parent name, then the category name, and then the product name. This is a hierarchical relationship in the way Teka organizes the categories of products that it sells. Now that I've done that, notice that in this treemap I've got some additional icons up here. Let me turn off the filters. If I hover over these, this allows me to click, to turn on, to "Drill down". Let's start with this. If I click on that now, if I click on one of these super parent names, let's try Other. I know it has a lot. Now, you can see that it's showing me what makes up the total amount in the Other category. I can see the largest amount is off-invoice cigs and then lottery, cigarette smokeless, so on and so forth. Now, I can drill down even further. Let's say I want to go to breakfast sandwiches and when I click on that also note what happens to the dots on the map.



All right, so those are changing as well and we can see the amount of breakfast food that is now purchased in each of the states, relatively speaking, as well as in this table over here. Now, let's click on breakfast food one more time to see what makes up breakfast food. We can see the different breakfast food items and they're laid out in this table very well. Now that we're in this lowest level, I may want to use this as a filter, not as a drill down. I can turn off that drill down mode and now I can click on, for instance, this pretzel bun and now everything will filter to show me just the pretzel buttons. I can see those are sold mostly in Iowa and Alabama, total of \$18.05, a little bit in South Dakota as well. Then I can undrill if I click on this up arrow button and go back up to the top layer or I guess that's called drilling up, not on drilling. I can turn off that drill mode and turn off all the filters.



Now, another way that you can drill down to other levels is by clicking on this double arrows, and watch what happens here if I click on that. Watch what happens to fuel and other. They then now are broken down into the various components still on the same treemap. So fuel, its product name or the parent name category is just fuel. But these other super parent names are made up of more than one parent name, so those got broken down into more categories. I can keep going down to the next level and you can see that now this is getting overwhelming in here. I've got all these different product names basically. Now I've got all the product names that I'm at and it's just overwhelming. Now, if you do have something like this and you really want to explore it more, you can click on this focus mode here and that will make it so that it fills up the whole screen and allows you to see a little bit more detail. This is still way too much, though, to really learn much from other than maybe these ones that make up the majority. Then I can click on "Back to report" to go back to the report.

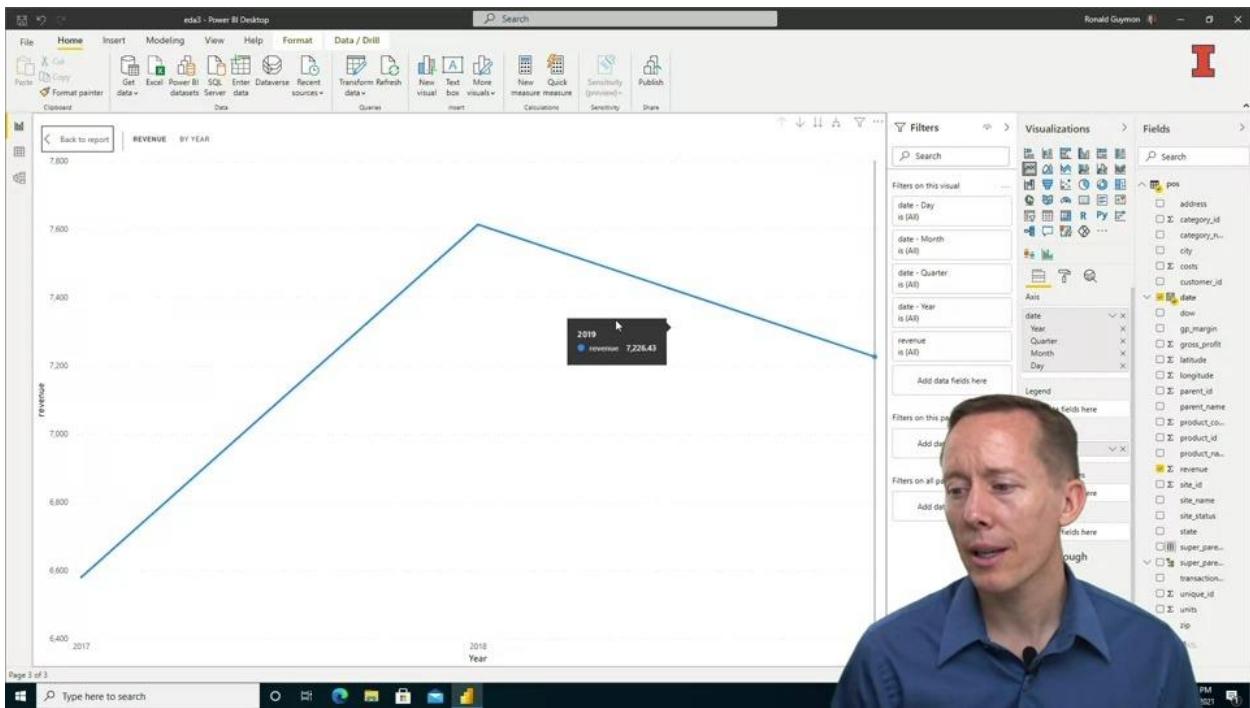


These are demonstration of how to use the drill-down feature and Power BI. I think it's super awesome. In conclusion, I'll leave you with the challenge, see if you can replicate this state's map so that you can drill down from state, down to zip code and then down to city.

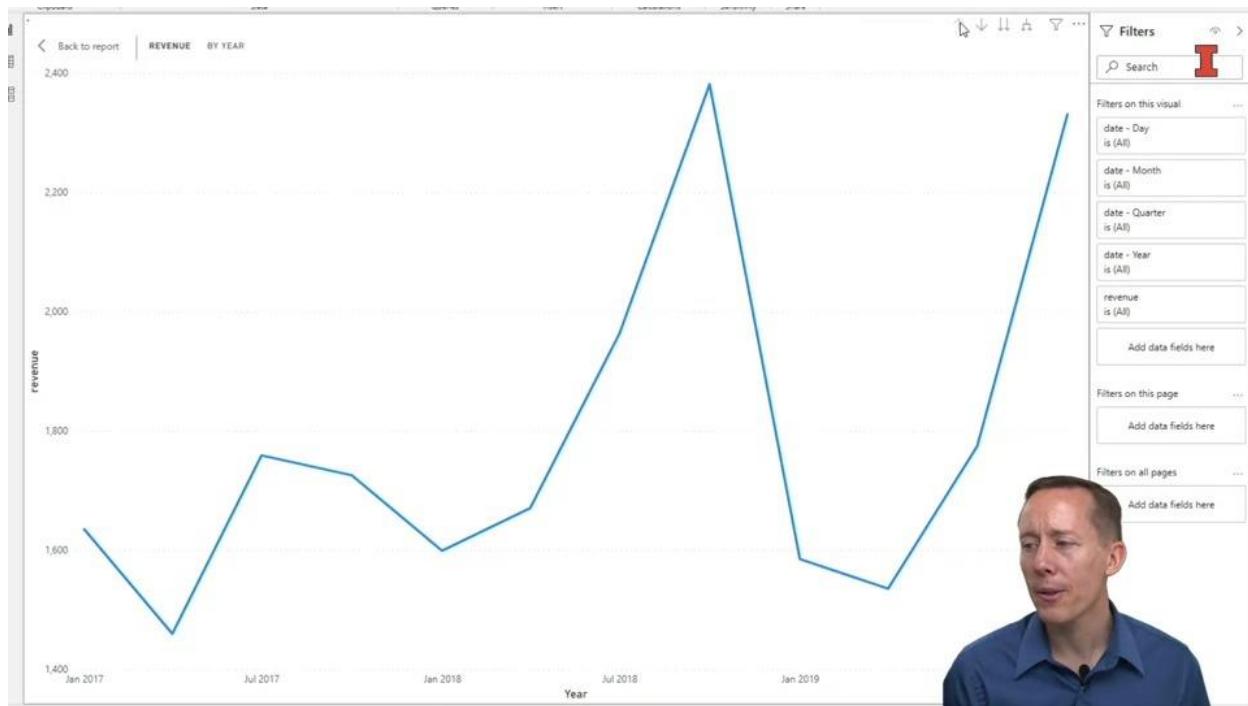
Lesson 2-3-4 EDA 4: Multivariate Plots: Scatter Plots and Line Plots



In this lesson, we are going to focus on multivariate plots, in contrast to univariate plots to look at one column of data and bivariate plots to look at two columns of data. Multivariate plots look at three or more columns of data. And these types of plots can include a mixture of categorical and numeric data. Now the focus will be on scatter plots and line plots.



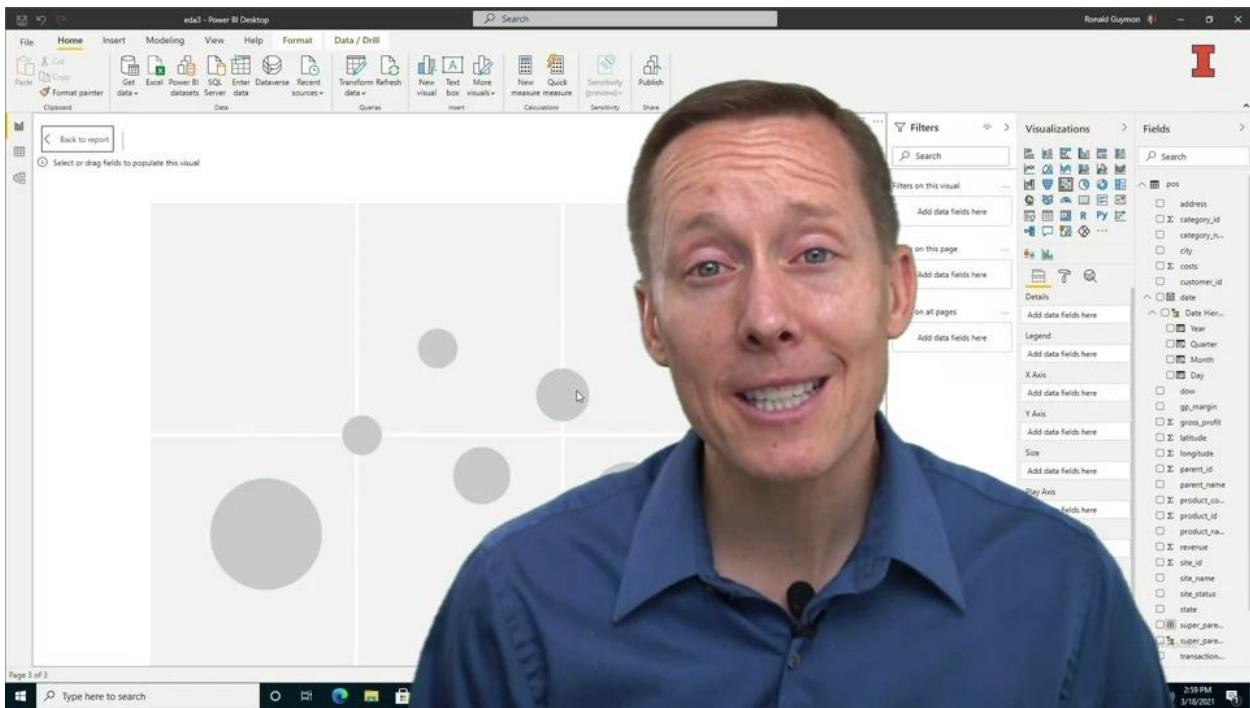
So first let's open up our Power BI and I'm going to open up this project I've been working on ETA 3. And on this report you can see that we have a tab for univariate plots and then a tab for categorical plots, I'm going to create a new tab for multivariate plots. And in this tab I am going to create a line chart. Now line charts are really helpful for connecting points that are related to each other in terms of time. So I'm going to go ahead and drag onto the axis here, the date and notice how it includes the hierarchies of the date, we'll come back to that. And then in the value section I'm going to drag revenue and I'm going to focus on this plot, I'll blow it up so we can see it better. Now, notice what it does is it's taking the date at the highest level of the hierarchy year and adding up the values so that we've got a total revenue for each year. And if I look at revenue, I can see that it's summing it up there. So that's pretty awesome, how quickly you can create a line plot in Power BI. What if I want to look at total revenue in terms of quarter?



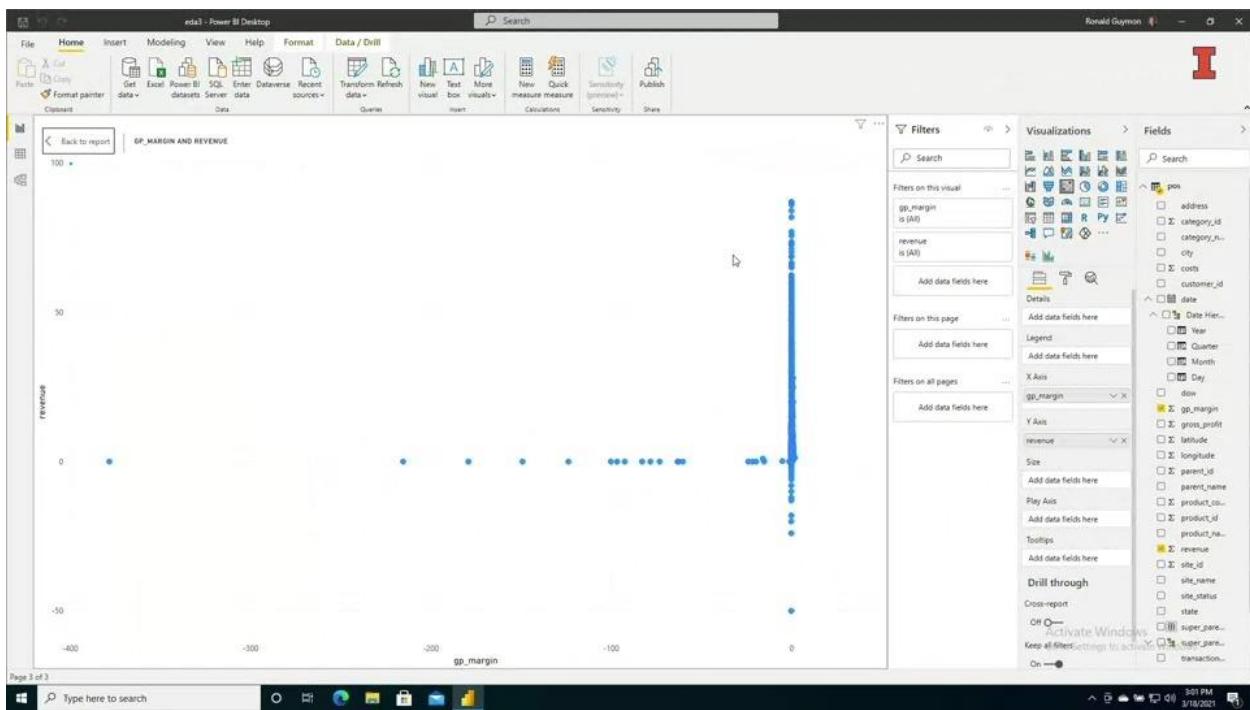
Well, there's a couple of ways I can do that, I can go up here and since there's a hierarchical relationship, I can just drill down to the next level in the hierarchy. And you can see that it's grouping it by a quarter, now, it's not by quarter and years, just by quarter. So all quarter one overall three years has being grouped together and I can continue to drill down to see month and then to see day, which is really day of the month. All right, I'll go back up to the top, I'll drill up about two year again. And if I want to, I could expand all down one level in the hierarchy and when I do this, it will break it up into quarter and year and then month by year and finally day and year. So that's using the drill down feature, but you may not want the drill down feature.



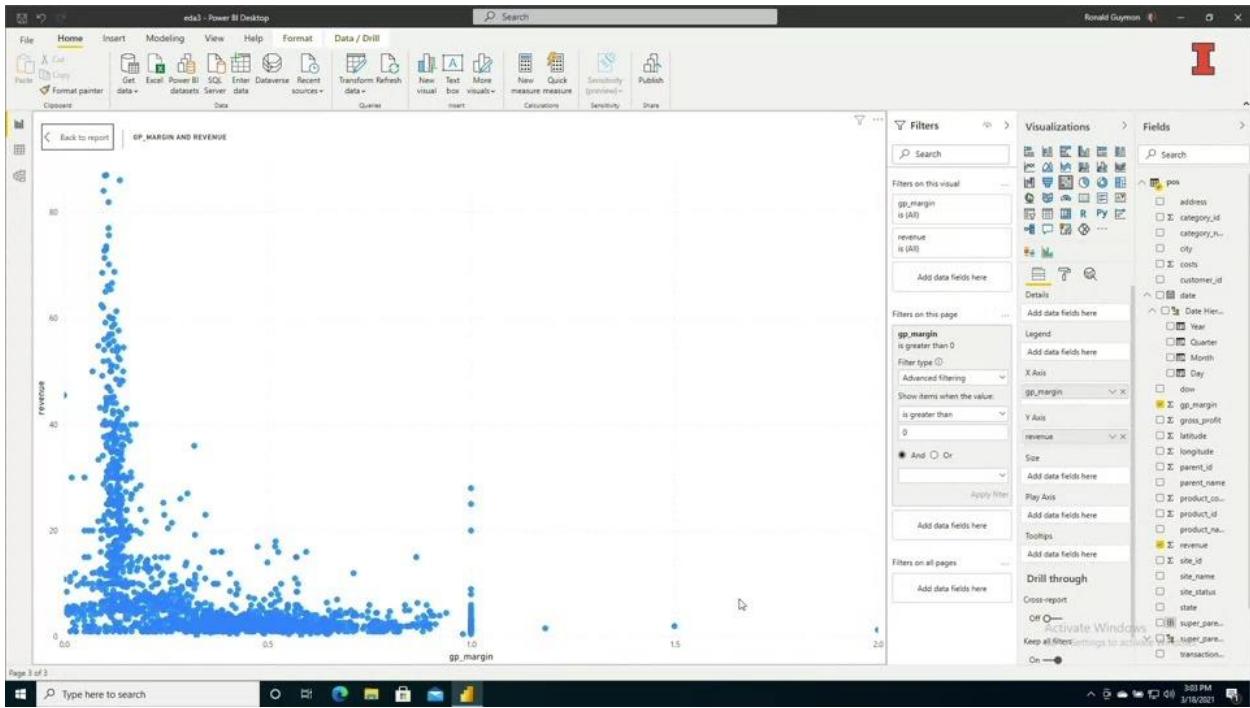
So another thing that you can do is just go over to the date here and select the down arrow and say, hey I don't want the date hierarchy, I just want the date. And if I do that, now you can see that we've got an individual data point for each day in this data set. All right, so this is two columns of data, date and revenue. I can add a third column of data by adding in year. So if I go to the date hierarchy and drag year to the legend. Now I can see all the daily sales and it's colored by date. So that's allowing me to investigate three different columns of data at once and this may not be as useful here. All depends on what you're trying to do, again if I want to maybe see year over year monthly sales, that may be more useful. I could go back over to the date and go to date hierarchy and I can drill down and say I don't want it a year or a quarter, but I wanted that month and now I've got a separate line for each year and a dot for each month. So pretty awesome how we can create a multivariate line plot in Power BI, so quickly and intuitively.



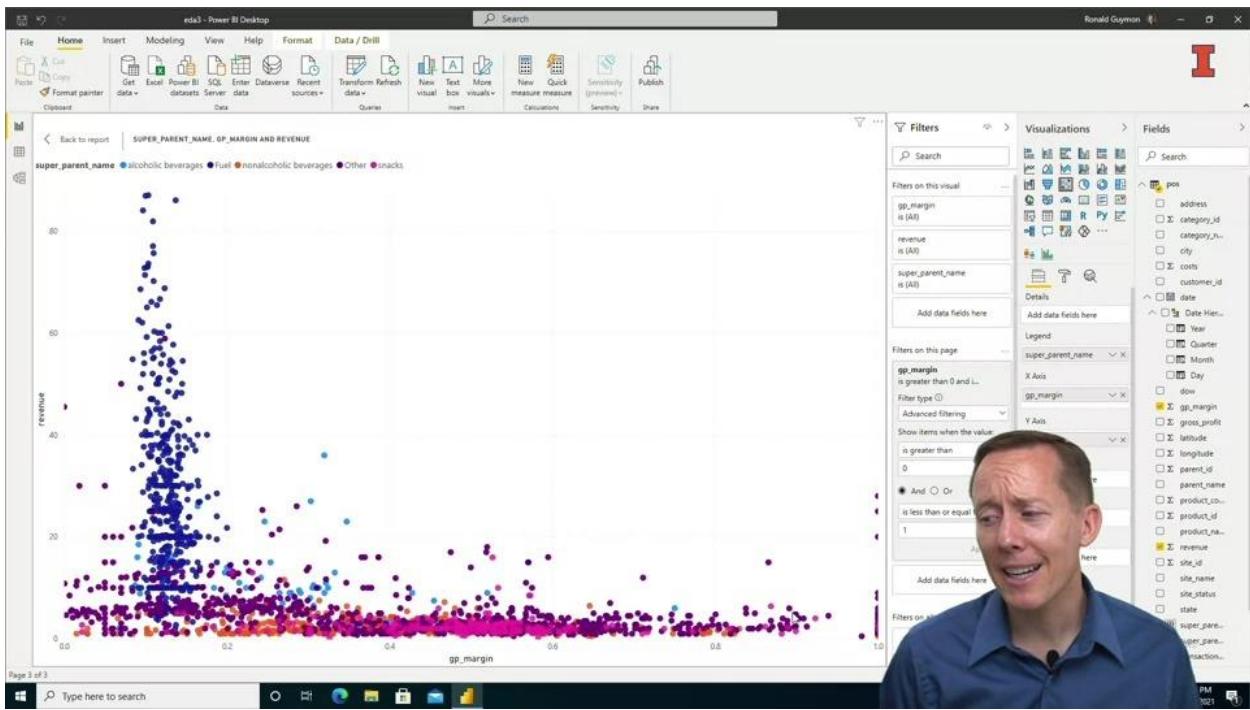
All right, let's go back to the report here and create a scatter plot. So I will go ahead and click on the scatter chart icon and then I will focus on this. And let's look at the relationship between our gross profit margin and revenue, and the rationale for this is we might want to investigate if units or line items that we sell that have a high gross profit margin also have a high revenue. That would be ideal, right? If our highest gross profit margin, if we have 99% gross profit margin also brings in a ton of revenue. So let's go ahead and investigate that.



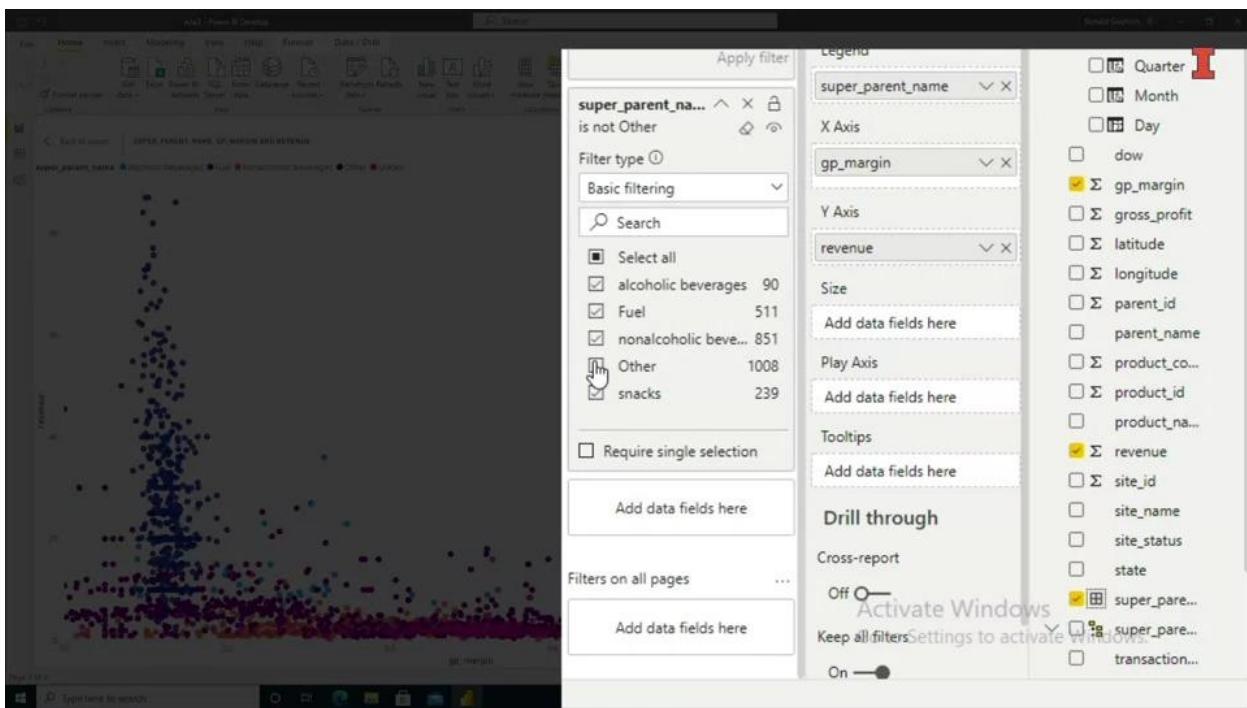
So we will move gp margin to the x axis and on the y axis we will put revenue. Now at this point there's not a whole lot of information here that's because there's some aggregation going on. So we need to click on the down arrows and let's say don't summarize and it's saying it's not playing well, well why is that the case? So if we look at gross profit margin, gp margin, we can see that it is a text data type right now. So we need to go in and change that. So we'll go to the home tab, go to transform data. This brings up the power query editor and we'll scroll over to gross profit margin and we'll click on this data type icon and change it to a decimal number. All right, now we can see that there are 93% valid entries and 7% that are empty, perfect. So this is what we want to go ahead and close and apply this. All right, now we've got a scatter plot that is working.



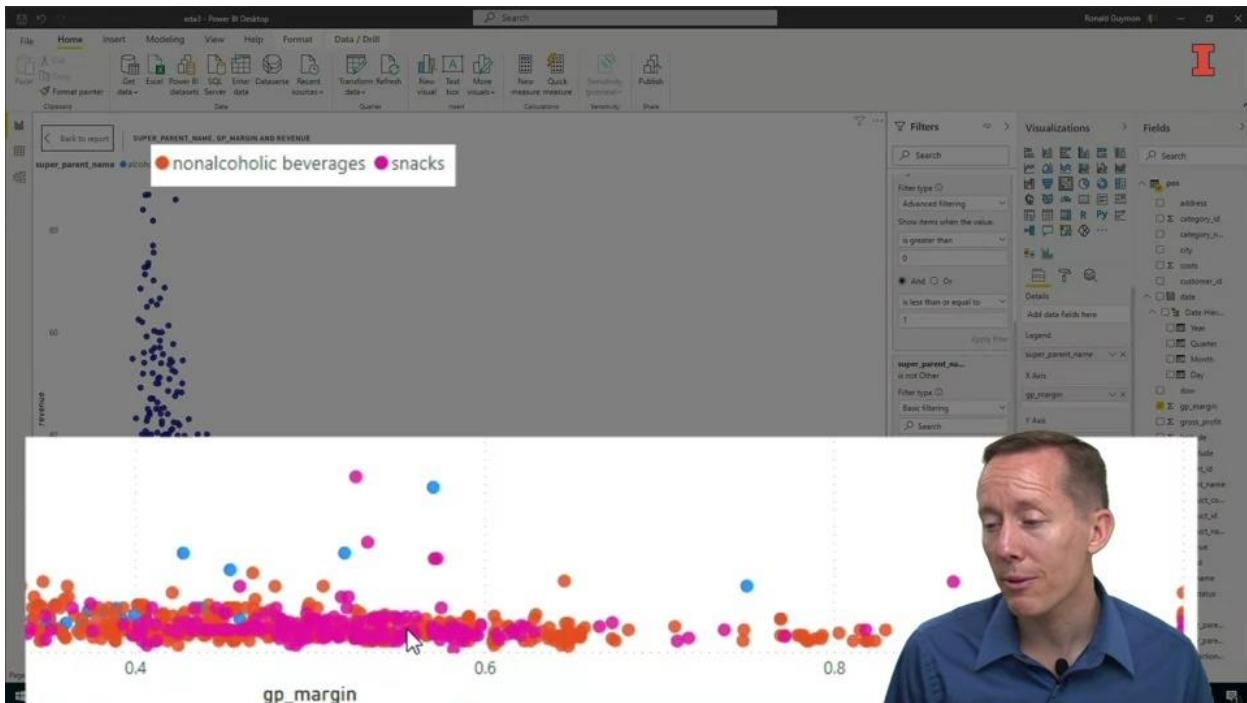
And let's go ahead and make sure that gross profit margin is not being aggregated and it's not so perfect. But this isn't real informative and so let's try to explore what's going on here. As we look at this, we can see that on the x axis there are definitely some outliers on the low end, so negative 400 up to zero. And we may want to remove those outliers, I suspect that most of the observations will be between zero and one. And so we could go back into the power of query editor and remove outliers that way but there's often a better way to do that in Power BI and that's by using filters or slicers. Let's go ahead and just use the filter here. And let's drag gross profit margin onto the filters on this page, and notice that we can use filters with numeric data as well. So we want to say we are using gross profit margin and we want to say it's not less than but greater than zero. And let's go ahead and apply that filter. All right, so this is looking much better. We're not being overwhelmed by the extreme negative outliers, but we do have some positive outliers that we may not want to include either.



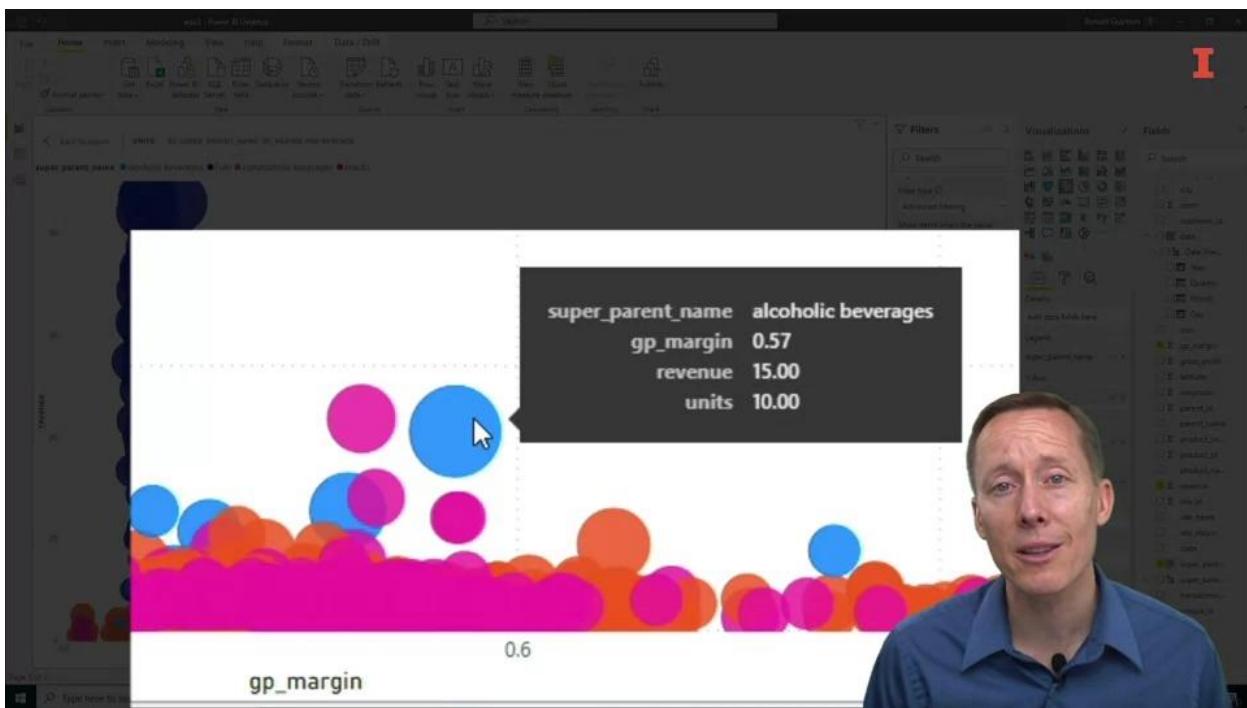
So we can go ahead and go back to our filter and say we also want to filter observations such that gross profit margin is less than 1 and actually might say less than or equal to 1, then apply the filter. Now this lets us see a lot more of the nuance in the data. And at this point we're looking at two different columns of data, gross profit margin and revenue. And the next thing we might want to ask ourselves is does this relationship, which it's not real strong. It looks like we've got a bunch of observations that are kind of low profit margin, but high revenue. And then a bunch of observations that are high profit margin, low revenue. Well, we might want to ask ourselves, is that due to a particular type of product? So let's go ahead and color these data points using the super parent name. So drag super parent into the legend box. All right, now this is pretty cool, we can see that all these points here that are blue are the fuel points and so those are the ones that have the low gross profit margin but high revenue typically. Then we've got some over here that are this dark purple and these are the other category here. So this, if you may recall, this is a bunch of different types of products that are in that other category. Now this may be kind of overwhelming, we're not interested in seeing the other in there.



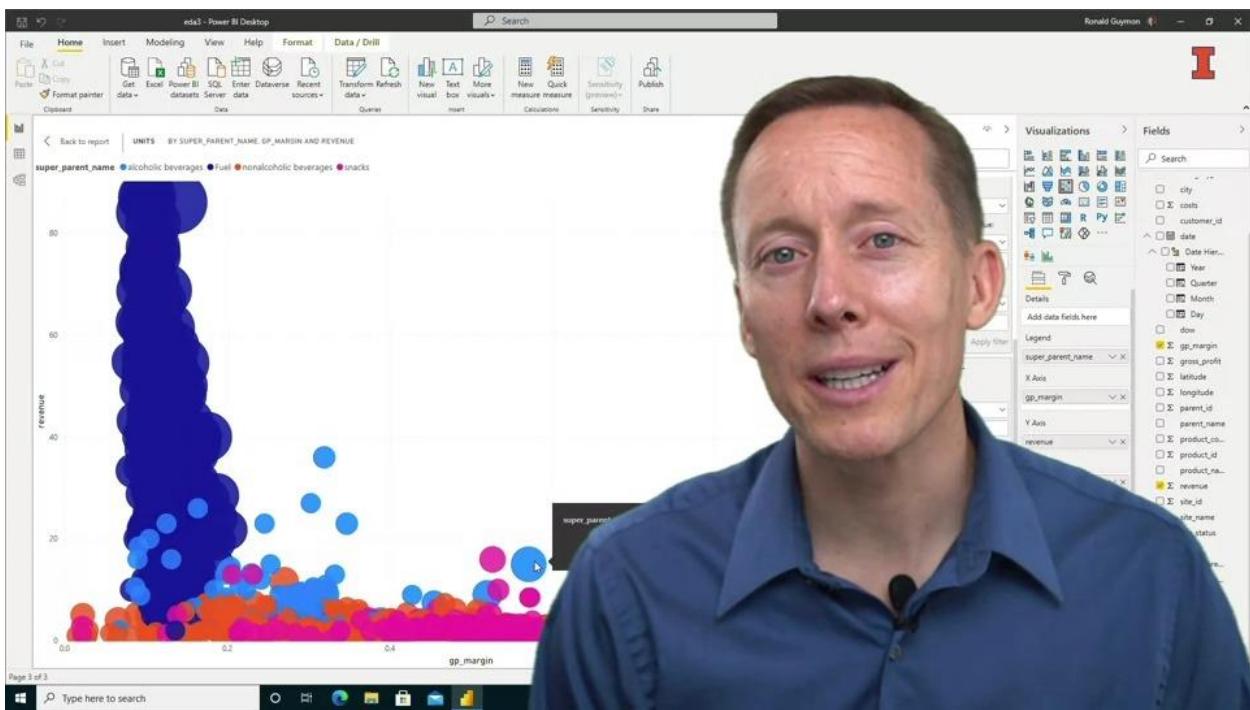
So we could filter out these other observations by also dragging super parent name into this filter And now we have a new filter box and we'll start with selecting all and then we will just uncheck the other.



So now we're not being overwhelmed by all the other observations and it's a little easier to see the gross profit margin and revenue relationship between our four main categories. And we can see that alcoholic beverages, these light blue points often have a higher revenue and a higher gross profit margin, whereas the orange non alcoholic beverages and pink snacks often have a really high gross profit margin, but pretty low revenue.



Now you might ask yourselves as you think about the data. Well, maybe some of these observations have high revenue because there are a lot of units that people are buying at once. Well, we can actually investigate that and include that in here as a fourth variable. So let's go ahead and drag units into this size box here. All right, now that's kind of overwhelming, it doesn't show a whole lot here. That's because remember fuel, the units are the gallons purchase and that's often in the teens or 20's, whereas these snacks are often in the 1's or 2's. But we can see here that this observation for an alcoholic beverage has a really high gross profit margin and high revenue is because they purchase ten units at once.

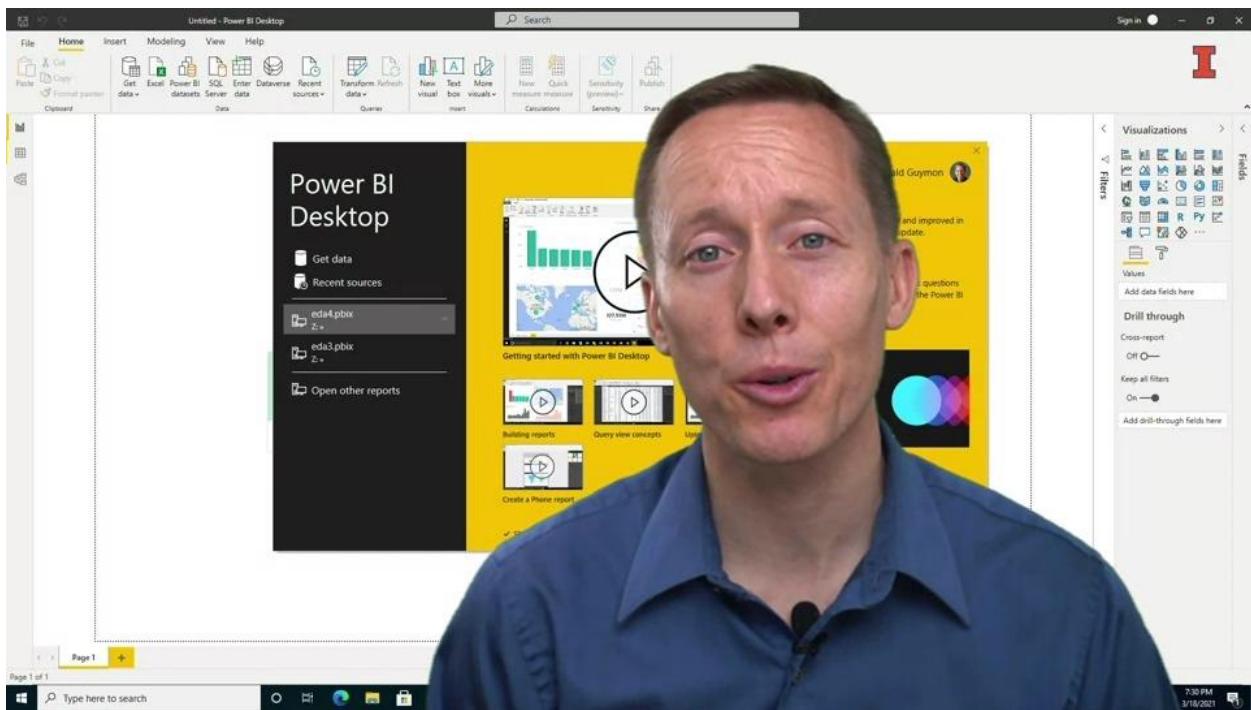


So that helps kind of explain things as well. So real quickly, you can see how we can create a multivariate scatter plot in a multivariate line chart and it's easy to include three or four different columns of data at once. Now, word of warning is to be careful in using too many multivariate plots on a report like this.

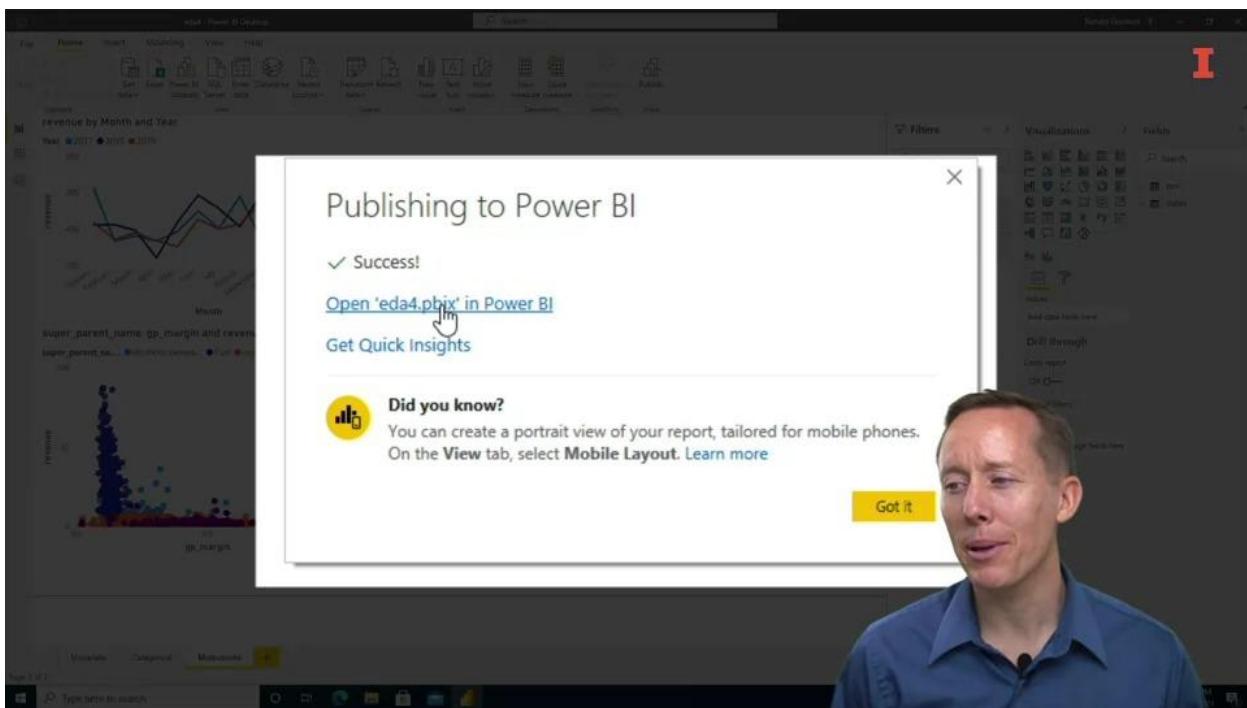


When you do that and you have all these plots that have lots of different variables on there it can get overwhelming very quickly. So I would recommend to maybe just put one multivariate plot on a tab in a report, maybe two, but just consider how much information processing will need to occur by the people interpreting this. And if your purpose is to just explore the data, maybe it does make sense to put all these plots on there. But on the other hand, you could replace one multivariate plot with several bivariate and univariate plots and then filter through and look at the relationships that way.

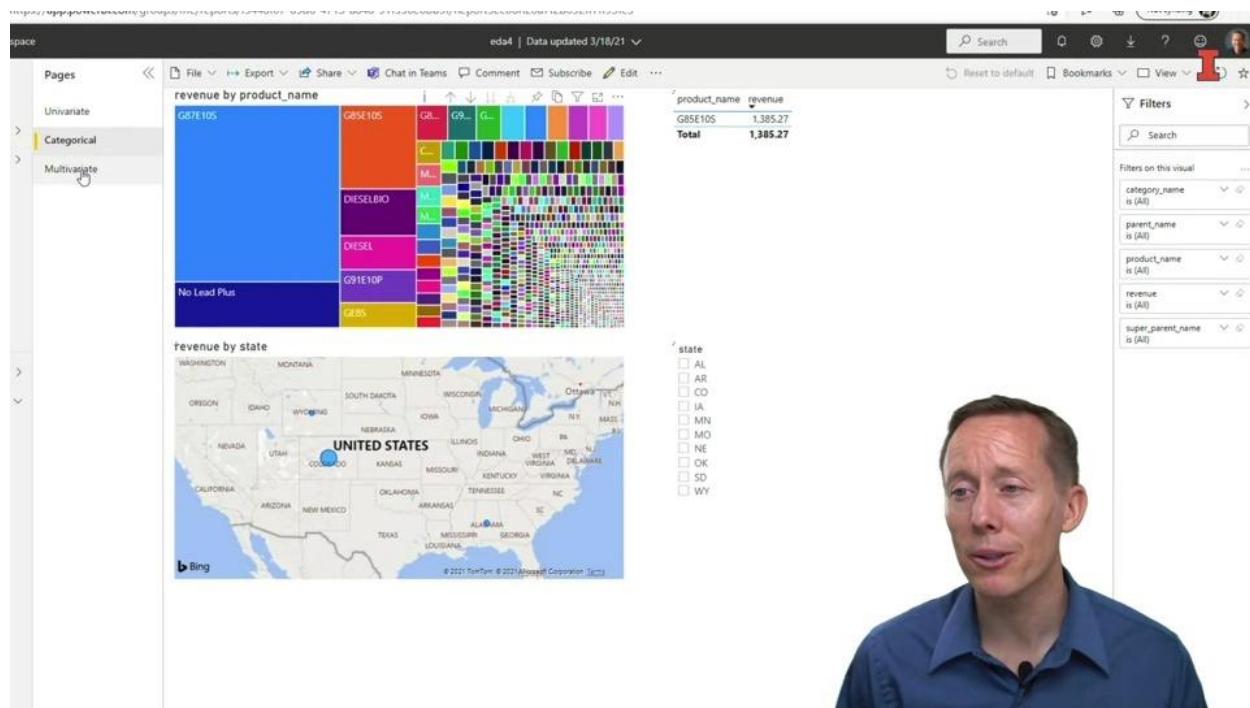
Lesson 2-3.5 EDA 5: Publishing a Power BI Report



In this video, we want to show you how you can really take advantage of the Microsoft environment by publishing your Power BI report online so that everyone in your organization can access it.



First of all, we'll go ahead and open Power BI and open a report that we've been working on. All right, now that we have this report open, we just make sure that we go to the home tab and then over to the public icon, we'll save any changes that had not been saved, and then if you have not logged in you'll be asked to log in, make sure you're logged in to your organization account, and then you can choose if you want to publish this to just your workspace to different workspaces. I'll just publish it to my workspace and click select. All right, it has published it, let's go see what that means. We will open this in Power BI online. All right.



Here's my report in my personal workspace online, and if you explore it a little bit, you'll see that you can access the different tabs of this report, you will notice that our visuals are a Power BI Pro feature, so that's one thing that you have to pay for. That's a limitation, but it's still pretty awesome that for free, I guess your organization needs an account. You can publish these reports and you could share them to an area where others could access them, and they could apply filters, they could click on parts of a treemap or a map to apply filters that way.



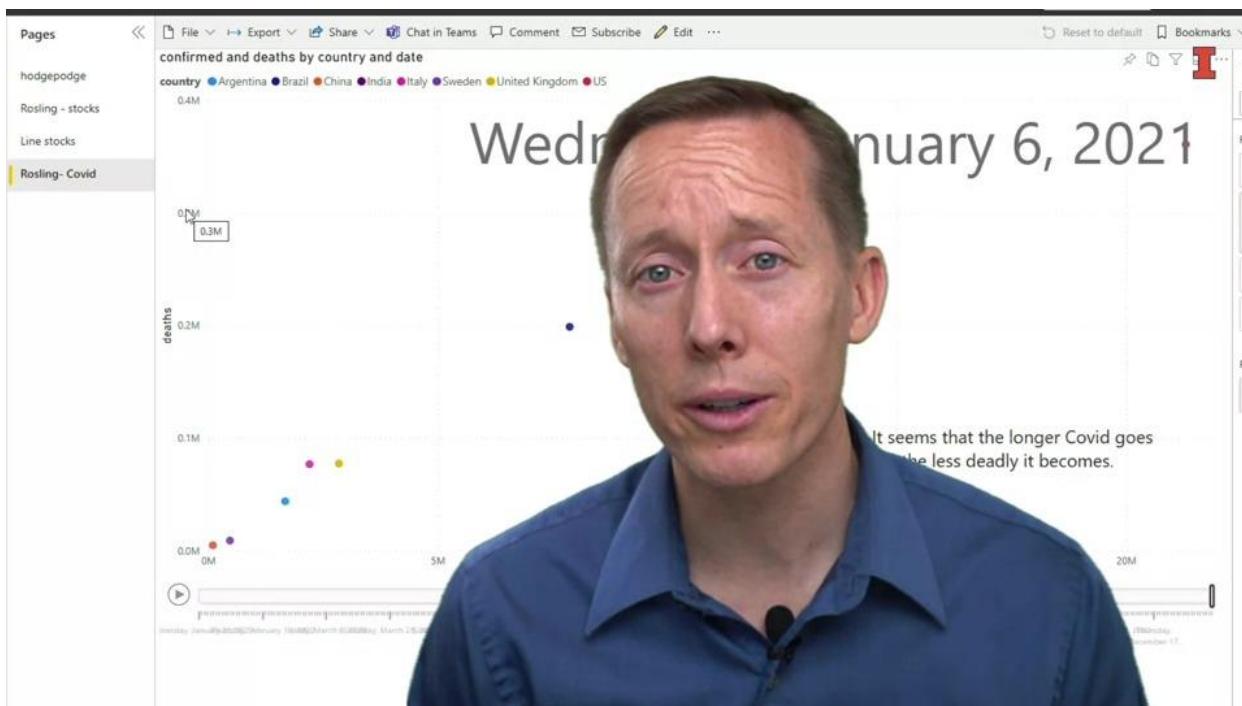
They can't design and edit the layout of the report, but they can apply filters and explore it. This is a really powerful feature, I think, of Power BI, this is just an exploratory data analysis report. But even then, if you're working with lots of people, what a great way to allow them to access your data. Now, let me show you something else about this environment here.

The screenshot shows the Power BI workspace interface. At the top, there's a header with a profile picture, the workspace name "My workspace", and a "New" button. Below the header is a navigation bar with "All", "Content", and "Datasets + dataflows" tabs. The "Content" tab is selected, displaying a list of items:

Name	Type	Owner	Refreshed	Next refresh	Endorsement	Sensitivity
covidDashboard	Report	Ronald Guymon	3/18/21, 3:19:42 PM	—	—	—
covidDashboard	Dataset	Ronald Guymon	3/18/21, 3:19:42 PM	N/A	—	—
eda4	Report	Ronald Guymon	3/18/21, 7:31:20 PM	—	—	—
eda4	Dataset	Ronald Guymon	3/18/21, 7:31:20 PM	N/A	—	—
edaTeca_6	Report	Ronald Guymon	3/3/21, 10:38:27 PM	—	—	—
edaTeca_6	Dataset	Ronald Guymon	3/3/21, 10:38:27 PM	N/A	—	—
tdAlabama	Dataset	Ronald Guymon	1/6/21, 11:55:10 AM	N/A	—	—
tdAlabamaDashboard1	Dashboard	Ronald Guymon	—	—	—	—

A video overlay of a man in a blue shirt is positioned over the right side of the screen.

If I go to the far left sidebar, click on this home tab, I can see favorite's and frequent, I can see data stories from Power BI community, there's some other tutorials. I'll just click on my workspace. In here, I can see other reports that I have created like this COVID dashboard, as well as the data that goes along with it.



If I want to view another dashboard, I can click on that, and now I can access this other report here and explore it. This is enhanced Rosling type visualization where you can look at the evolution of COVID over time. This is an example of posting an exploratory data analysis report online, but certainly, you could use this to post a finished overall report to highlight findings in your analysis.

Module 2 Review

Module 2 Conclusion



In the summer of 1812, the French leader Napoleon created an army of 680,000 soldiers and began an invasion of Russia. At that point, it was the largest assembly of soldiers ever known. Napoleon marched with his troops to Moscow, Russia, and expected the Russian army to surrender. However, they never did.

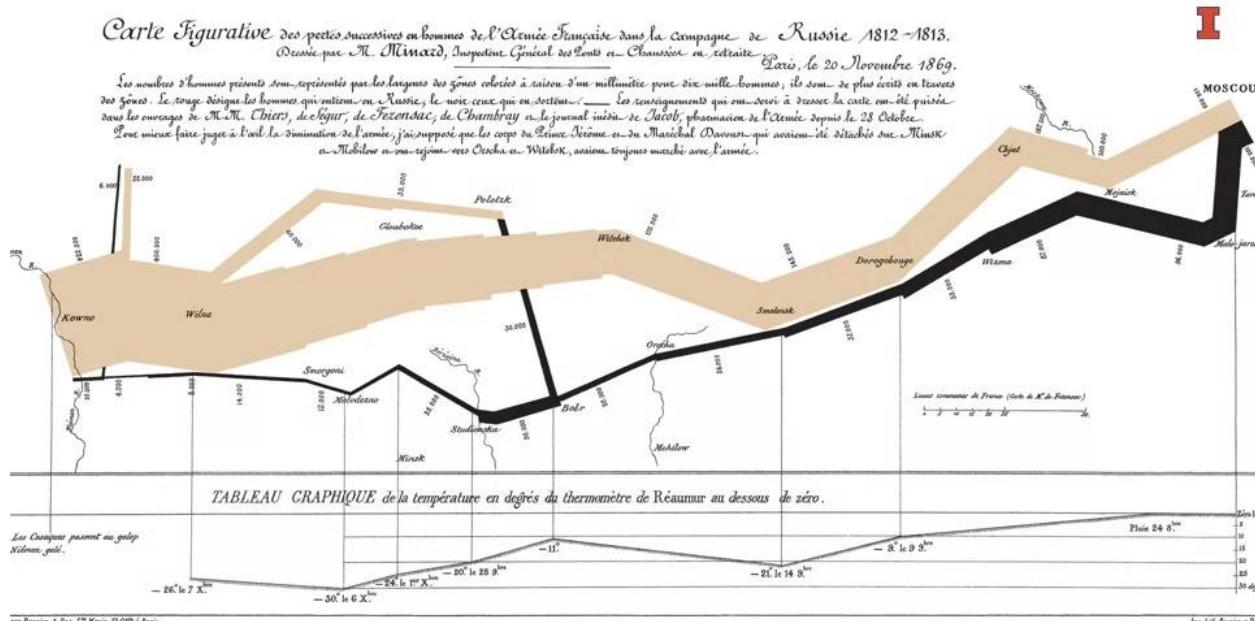


I

Scorched earth tactics in which the armies would burn everything and leave nothing for the opposing army to subsist upon left Napoleon's army lacking supplies. The distance and terrain made it hard for the French to deliver adequate provisions to its army. Eventually, Napoleon retreated and headed back to France.



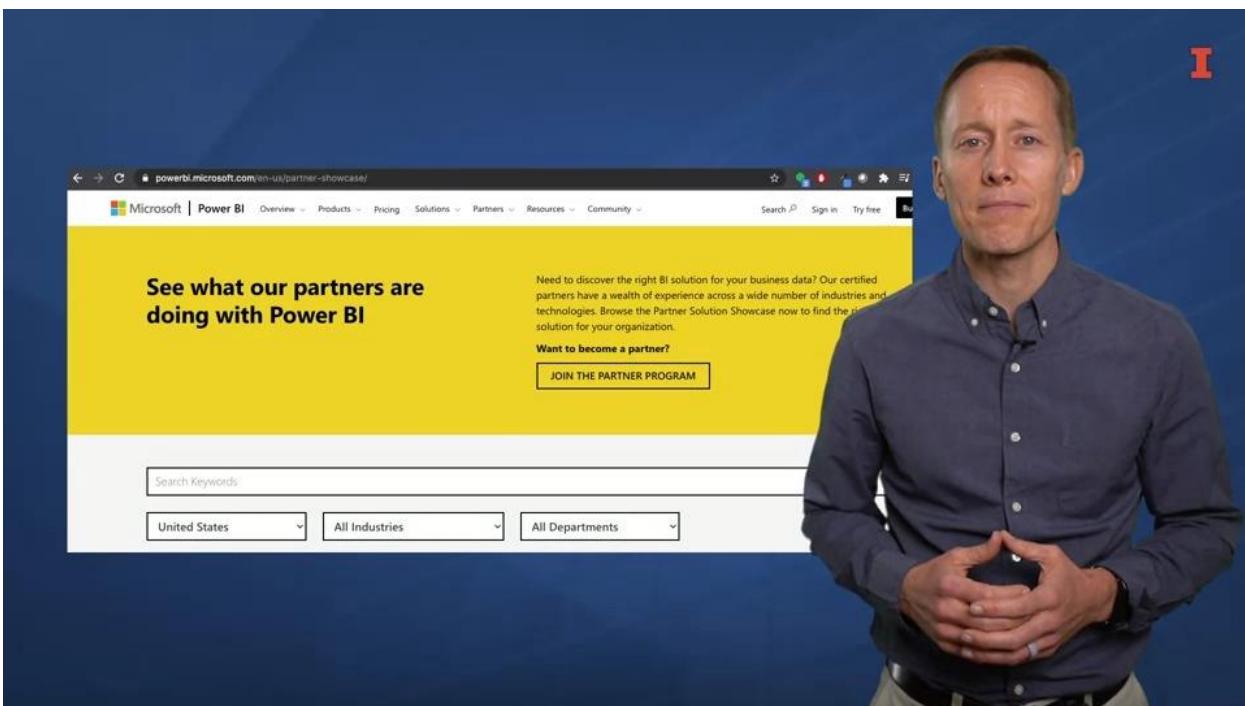
By that time, the temperature had started to drop. As a result of the cold, inadequate provisions, and repeated attacks from Russian peasants, Napoleon returned to Paris with only 27,000 soldiers. I tell you that story because one of the greatest data visualizations of all time is attributed to Charles Minard's depiction of that ill-fated campaign.



Minard used basic elements to succinctly communicate at least five dimensions on the chart. One, bars communicate the location of the paths upon which the troops traveled. Two bars also communicate the distance that the troops traveled. Three, color communicates the direction that the troops traveled. Four, bar width communicates the number of troops in the army, and five, a separate line chart on the bottom is used to communicate the temperature during the French army's return from Moscow.



This is a great example of using data visualizations to tell a story. It's more impressive when you consider that Charles Minard was an engineer who died in 1870. Can you imagine how long it must have taken him to create that chart by hand? Can you imagine how happy he would have been if he had a tool like Power BI? How do you think he would have used the filtering and zooming tools to improve his visual story? I hope at this point you have a better idea of how power BI fits in with the business analytic workflow. It's a powerful tool for assembling and exploring data, as well as for communicating results with others. There are many wonderful modern examples of effective data visualizations that can be created with Power BI.



Microsoft has a partner showcase website where you can explore examples from many different industries and departments. I encourage you to consider these examples inspiration for how to use data visualization tools that you've learned about not only to explore data but also to effectively communicate a story. As you can see by now in these lessons, we've only scratched the surface of what is possible with Power BI.