# EM384: Analytical Methods for Engineering Management

Lesson 6: Data Exploration and Analysis I

26 January 2023

# Table of contents

# Lesson Objectives

# Lesson 6 Objectives

- Define a database and describe its components.
- Use the five-step process for data exploration to evaluate a database.
- Generate conclusions from examining relationships among variables.
- Create a 'Pivot Table' in Excel from a dataset

# Datasets

# A Database/Dataset

- A database contains a dataset. A dataset is a table of information.

- In this class, we generally use two-dimensional datasets.
  - Each row is a record in the dataset.
  - Each column is a field for the records.

## A Database/Dataset

- Be skeptical of data, and ask:
    - How are fields defined?
    - What types of data are represented? Nominal, Ordinal, binary, etc.
    - What units are the data in?

- 5 Steps to Using Data:
    - Understand the Data
    - Organize and Subset
    - Examine Individual Variables
    - Calculate Summary Measures for the Individual Variables
    - Examine Relationships among the Variables

## The Five Steps to Using Data

The five Steps to Using Data:

- Understand the Data
- Organize and Subset
- Examine Individual Variables
- Calculate Summary Measures for the Individual Variables
- Examine Relationships among the Variables

# Organize and Subset

Two essential tools: Sort and Filter.

On the Home ribbon in the Editing group and the Data Ribbon in the Sort and Filter group

Home ▸ Editing ▸ Sort & Filter ▸ Custom Sort opens the Sort window

(Sort by more than one criterion using Add Level)
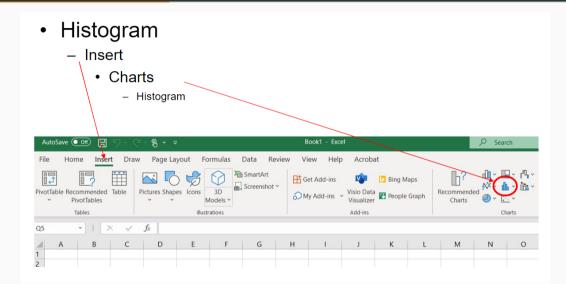
Home ▸ Editing ▸ Sort & Filter ▸ Filter

Filtering allows us to probe a large database and extract what interests us.

Warning!: Using 'average' (for example) on a filtered list will include the records that have been filtered out! Use SUBTOTAL

## Examine Individual Variables

- For numerical variables, we typically want to know the range of records from lowest to highest, and areas where most outcomes lie.

- A common way to summarize a set of numerical values is the histogram, although Excel provides eight choices.

- Excel provides numerous functions useful for investigating individual variables.

- Some can summarize the values of numerical variables; others can be used to identify or count specific variables, both numerical and categorical.

# Histogram
- Insert
  - Charts
    - Histogram

# Calculate Summary Measures for Individual Variables

The most common summary measure of a numerical value is average or mean.

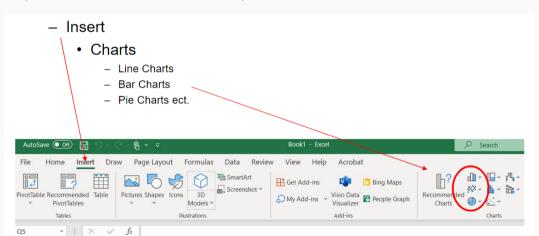Calculate using the AVERAGE function in Excel, for example:

- AVERAGE(C2:C2918) = 28.97
  *Crtl-shift-down, etc!

Other useful summary measures are median, minimum, maximum, and standard deviation.

In many cases relationships among variables are more important in analysis than the properties of one variable.
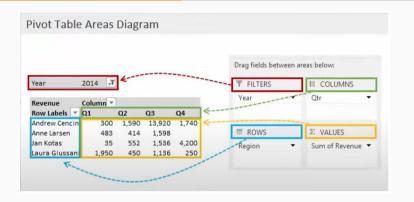
Graphical methods can track relationships.

# Pivot Tables

A pivot table is a table that summarizes data from another table and is made by applying an operation such as sorting, averaging, or summing to data in the first table, typically including grouping of the data.

…A tool to easily and quickly organize and present your data.

Pivot Table Areas Diagram

- Know the four options: "Filters, Rows, Columns, Values"
- You can add the same 'field' to multiple places
- Best technique: Trial and Error between the four options!

## Pivot Tables - Refreshing Data

Important note about Pivot Tables: If the data values change, the table does NOT update automatically. You need to refresh the table. There are 2 ways to do this:
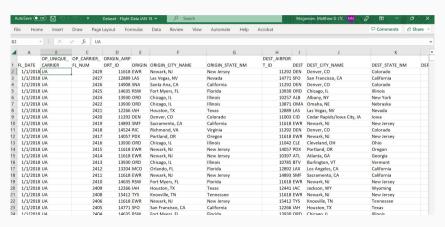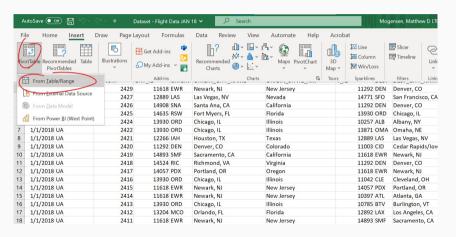
- From the Ribbon, under Pivot Table Tools



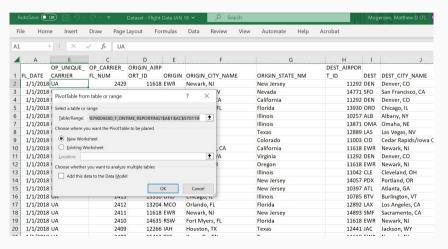- Right Click on a call in the pivot table and choose Refresh

# Pivot Table Practice

Download and open the "Dataset - Flight Data JAN18" file from on Teams. Take a moment to examine the fields and records. How many are there of each?
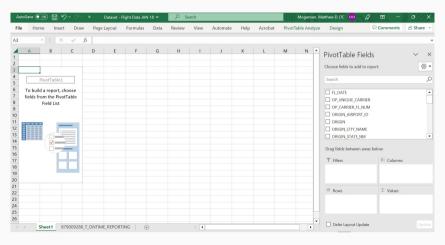
# Pivot Table Practice

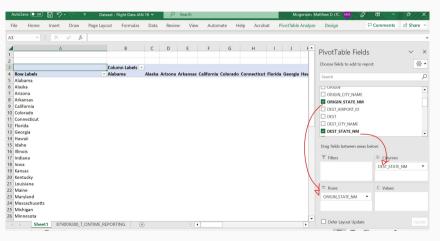Select **insert** then **pivot table**, then From Table/Range.

# Pivot Table Practice

The Table/Range may autofill. Otherwise select the range for your dataset, then select **New Worksheet**, and press **OK**.
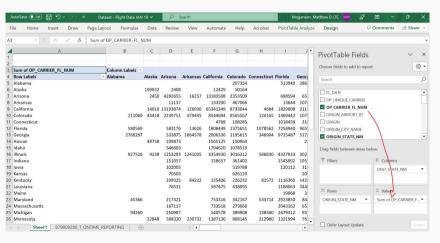
# Pivot Table Practice

A new sheet should be created with the pivot table fields on the right side. If the pivot table fields are not there, click on your pivot table to the left to open the fields.
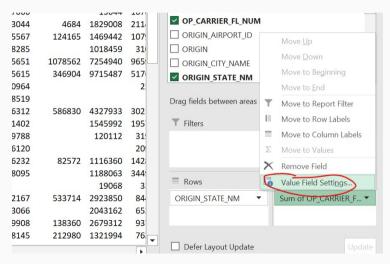
# Pivot Table Practice

Let's attempt to create a pivot table that will show us the number of flights between states. Drag the **ORIGIN_STATE_NM** field to **Rows** and the **DEST_STATE_NM** to **Columns**.
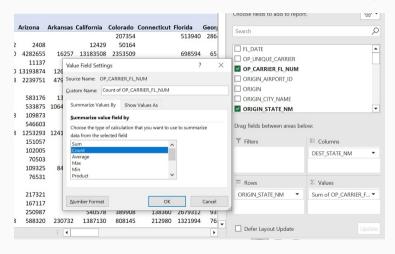
# Pivot Table Practice

Now drag the **OP CARRIER FLIGHT NUM** field to **Values**. The default is sum. Right now this pivot table doesn't tell us much because it's summing the flight numbers.
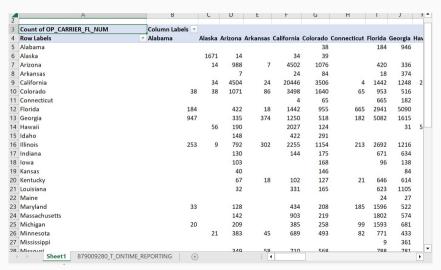
# Pivot Table Practice

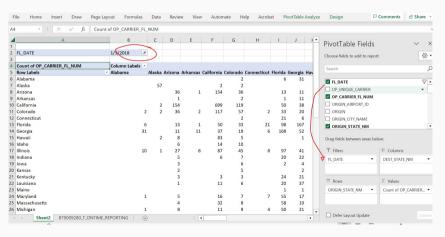Click on **OP CARRIER FLIGHT NUM** in the **Values** field.

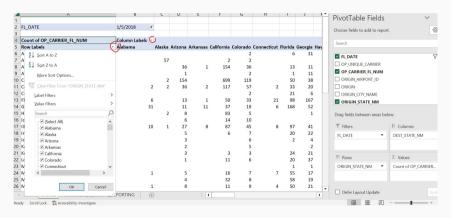Select **Count**. Now it only counts the number of values instead of summing.

We have now completed our pivot table. What if we want to filter by date?
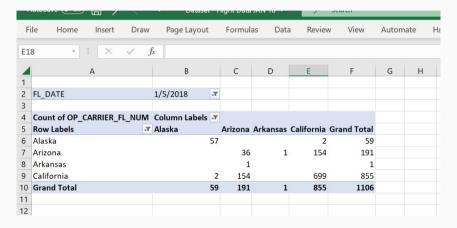
# Pivot Table Practice

Let's display only the flights for 5 January 2018. Drag **FL_DATE** down to **Filter**. The click on the filter box on the left and select only the date we want.

What if we want to only show flights between certain states? We can click the drop down buttons for rows and columns on the pivot table itself and only select the first four states.

# Pivot Table Practice

You should end up with the following pivot table

# Practical Exercise

# Conclusion

## Next Class

#### Homework:

- Watch tutorial video on Python libraries, data structures and control structures (Will be available on Teams).

#### Next Lesson:

- Understand the use of Python libraries.
- Understand basic Python data structures (integer, float, boolean, string, and lists/Numpy arrays), including assignment of values and referencing.
- Understand basic Python control structures (IF,ELSE,FOR).