

# Lab Package 1.1: Doing Data Science in ArcGIS Pro

Monday, August 16, 2021 7:51 PM

In this package, we will practice some ArcGIS Pro tools for a simple data science project.

The research project's goal is to analyze the 311 calls in the City of Dallas and identify where, when and what services were requests during the data period (October 1, 2020 to present). Where are the census tracts with high and low calls? What kinds of calls? Any temporal pattern? Open data are now the common policy for many cities, including Dallas. The vast readily available data drive advances in data science.

Go to the Dallas open data site at <https://www.dallasopendata.com/browse> and find the 311 data (below). (yes, I could have just include the direct link to 311 data, but I think that searching for the data will give you an opportunity to learn what data are available at this site.

### 311 Service Requests October 1, 2020 to Present

Services

Updated August 17, 2021  
Views 1,297

There are over 400 service requests types that are reported in the 311 system that affect the quality of life of our citizens, neighborhoods, and communities. The most popular service requests include but are not limited to animal services requests, high weeds, junk motor vehicles, and a number of other code compliance-related issues. Requests regarding streets and mobility (such as street and pot hole repair) are also common as well as requests to address environmental issues such as water conservation and air quality concerns.

This dataset represents all Service Requests created for the fiscal year time period of October 1, 2020 to present.

Less

Tags 311, animal, building, bulk, code concern, and 29 more

API Docs

The data are updated daily, so the data change every day. For this exercise, create a folder in OneDrive first (with a folder name starting with A, such as ASDSPkg1) and please download the data [here](#) to the OneDrive folder

Need to identify inconsistent encoding in the data:

161	161	Test CRM
162	162	TEST CRM
163	163	test sr
164	164	Test SR

Change all values to capitalize the first letter of the first word

Service Request Type

.denominator()

Insert Values

\*

/

+

-

=

Service\_Request\_Type =

Additional spaces

Service\_Request\_Type!.capitalize()

21	21	Bike share - trn
22	22	Boarding home facilities - ccs
23	23	Boarding home facilities - ccs

Service Request Type

Insert Values

\*

/

+

-

=

Service\_Request\_Type =

Missing code

Service\_Request\_Type!.split()

67	67	Homeless encampment	1
68	68	Homeless encampment - ohs	4040

32	32	Code concern	3
33	33	Code concern - ccs	92788

Calculate Service Request Type Simple

1. When you are on the data download page, check out the meta data.

How many columns are in the data? \_\_\_\_\_.

Which column records what has been done to the request? Note the column name here \_\_\_\_\_.

What is the data type for Lat\_Long Location? \_\_\_\_\_

Download the data in CSV file.

## Data Engineering

- Pay attention to the detail in every step
- The lat\_long location is in the format of (latitude, longitude). The comma will cause problems when loading the data into ArcGIS Pro because ArcGIS Pro (and many software tools) because we need to have latitude and longitude in separate columns as x and y to plot locations. The easiest way to do this is to open the file with Notepad, and use "find and replace" to find "( and )" and replace them with nothing (leave the replace blank). Then you also need to find Lat\_Long Location in the end of the first line, and change it to Latitude, Longitude.
- Add the data to ArcGIS Pro and display the locations with latitudes and longitudes, save the data to a feature class for easy access and analysis in ArcGIS Pro, and remove the standalone table csv file to save RAM space.
- Inspect the data, identify issues, and correct the data: remove data with no lat/long and correct data with coding errors.

2. How many service request types are there? \_\_\_\_\_

Which service type has the highest frequency of service calls? \_\_\_\_\_

Is Bike share a significant problem in Dallas? Yes, if more than 500 service calls, or otherwise no \_\_\_\_\_  
(note that the answers will change as the data updated daily.)

- Add a new field column, Service\_Request\_Type\_Simple, with simplified service type by removing "- xxx"  
Inspect the result statistics table, find potential duplicates, and correct the values (hint: street lighting).

3. How many simple service request types are there? \_\_\_\_\_

How many simple service request types are related to water? \_\_\_\_\_

Among the top 10 most calls, how many types are related to sanitation? \_\_\_\_\_

- The questions ask to map the calls in census tracts. However, more people in a census tract may lead to more service calls. Therefore, it is important to consider population when making the comparison of service calls across census tracts. Search census demonstration in the Living Atlas, and add Census Demonstration Products v5 April 2021 > Tracts v5 to the map, and display the tracts in the City of Dallas. Name the extent feature class TractIDDallasPopulation.

- e) The questions ask to map the calls in census tracts. However, more people in a census tract may lead to more service calls. Therefore, it is important to consider population when making the comparison of service calls across census tracts. Search census demonstration in the Living Atlas, and add Census Demonstration Products v5 April 2021 > Tracts v5 to the map, and clip out the tracts in the City of Dallas. Name the output feature class Tractv5DallasPopulation
- f) Cognitively, we need to limit types to 10 or less. More categories can be difficult to perceive any patterns. Look at the top 10 most frequent service request types. Add a new field column, Top\_Service\_Request\_Type, in the statistics table and re-encode all the service types to the following. Use "Top Service Request Types.xlsx" [here](#) reclassify simple service request types to top service request types.
- Code concern or violation
  - Sanitation issues
  - Water or wastewater issues
  - Traffic issues
  - Parking issues
  - Animal issues
  - Noise
  - Others
- g) After you added the Top Service Request Types in the statistics table, the table is now a "look-up" table. The table can be joined with the Service Calls, so that every record on the Service Calls has the Top Service Request Type. After joining the tables, we need to make the result permanent. Export it to a feature class (e.g., TopCategoryServices), and use the new feature class for the rest of the analysis.
- h) There are too many attributes in the feature class. Remove attributes that we are not going to use in the project can speed up analysis. Use delete field to keep only the following attributes: Created Date, Top\_Service\_Request\_Type
- i) Recall our research questions:
- Where are the census tracts with high and low calls? What kinds of calls? Any temporal patterns?

First, we will convert the field column "Created Date" from Text data type to a Time field, using Data Engineering > Format > Convert Time Field.

**A little visualization here** to look for temporal patterns.

Use the time clock and calendar heat chart to show temporal cycles

4. Which month in the data showed the highest number of service calls? \_\_\_\_\_  
 Which day in June had the lowest number of service calls? \_\_\_\_\_  
 Which month had the least numbers of calls throughout the entire month? \_\_\_\_\_

Back to Data Engineering

Now we prepare the data for analysis based on census tracts. Add the census tracts with census 2020 population. Download the Dallas City boundary shapefile from Dallas Open Data at <https://gis.dallascityhall.com/shapefileDownload.aspx> to clip the census tracts in the city (e.g. DallasTracts). Use SF1 Total Population for the analysis, so only keep geographic identification code and SF1 population in DallasTracts to speed up the analysis. (Note that there are suites of demographic data that can help answer many other questions beyond the project).

For spatial distribution of service calls, think about how to arrange data in a table for the analysis.

	Variable1	Variable2	...	...	...
Unit of analysis	value	value	...	...	...
Unit of analysis	value	value	...	...	...
...	...	...	...	...	...
...	...	...	...	...	...

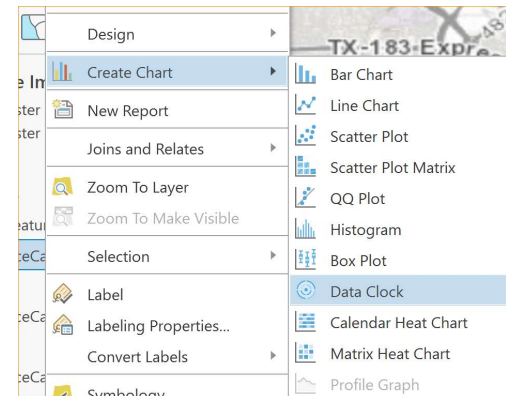
5. Challenge (required for graduate students;10 bonus points for undergraduate students who come up with a working solution): Create a table below for analysis.

33	33	Code concern - ccs	92788
----	----	--------------------	-------

Calculate Service\_Request\_Type\_Simple

Insert Values

Service\_Request\_Type\_Simple =
|Service\_Request\_Type|.split(' -')[0]



Add Data			
Portal > Living Atlas > Search Results for 'census demonstration' >			
Organize > New Item >	census demonstration		
	Title	Type	Date
Project	Census Demonstration Products v5 April 2021	Feature Layer (Hos)	8/5/2021 6:07:07 PM
Databases	Census Demonstration Products v2 May 2020	Feature Layer (Hos)	8/5/2021 6:12:29 PM
Folders	Census Demonstration Products v1 October 2019	Feature Layer (Hos)	8/5/2021 6:11:43 PM
Portal	Census Demonstration Products v3 September 2020	Feature Layer (Hos)	8/5/2021 6:10:32 PM
My Content	Census Demonstration Products v4 November 2020	Feature Layer (Hos)	8/5/2021 6:09:53 PM
My Favorites			

The table presents census tracts with service calls. The Geographic Identification Code is the census tract identifier in the census TIGER file.

Note that if the first Geographic Identification Code is <Null>. Look at the Service Calls table. There are many entries with <Null> in address, latitude, and longitude. These calls are without locations. However, there are some calls with addresses and latitudes/longitudes but without GEOID. Plot these points on the map to see why and find solutions to locate these points to corresponding census tracts.

Raise your hand or Send me a text in the class TEAM channel if you find the solution. Enter the sequence of tools in your solution (e.g., clip to the city of Dallas boundary, sum the total population over tracts, ...) You can only get the credits in class.

Field:		Add	Calculate	Selection:		Select By Attributes	Zoom To	Switch	Clear	Delete	
	OBJECTID_1 *	Geographic...	Animal	Code	Noise	Others	Parking	Sanitation	Traffic	Water	
1	1	48085031649	6	3	<Null>	6	<Null>	<Null>	<Null>	<Null>	
2	2	48085031704	113	153	2	261	23	271	83	31	
3	3	48085031706	19	46	<Null>	56	6	258	11	7	
4	4	48085031708	65	156	5	235	27	395	103	23	
5	5	48085031709	44	92	1	157	37	237	37	15	

Next, use find and replace to replace <Null> with 0, delete the total population (already in the DallasTracts), save the table, and Join the table to DallasTracts, save to a new feature class, named TractCallPop

#### Spatial Analysis and Machine Learning

1. Do an optimized hot spot analysis on the service calls. We must project the data to a Cartesian coordinate system first before the hot spot analysis because the hot spot analysis uses Euclidean distance in the calculation. Use Data Management Tools > Project to reproject to NAD\_1983\_StatePlane\_Texas\_North\_Central\_FIPS\_4202\_Feet (or use the City Limit layer as the project reference). Also project TractCallPop.

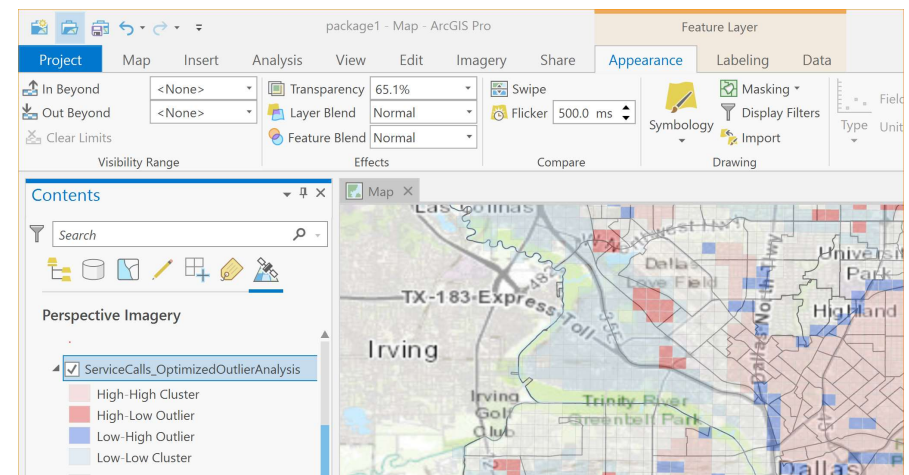
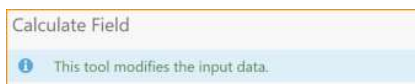
Perform Optimized Hot Spot Analysis, Optimized Outlier Analysis, and Density-based clustering with DBSCAN and 2000 minimum features per cluster. Set layers transparent one at a time, so that we can geographically contextualize the hot spots, cold spots and outliers.

6. What is the use of fishnet grid in Optimized Hot Spot Analysis and Optimized Outlier Analysis? \_\_\_\_\_ (a, b, or c) (a: aggregate service calls, b: assign colors for hot or cold spots, c: rasterization for spatial autocorrelation)  
Is Fair Park in downtown Dallas a hot spot or cold spot of service calls? \_\_\_\_\_  
Optimized Hot Spots are statistically based, but Density-based clustering is a learning-based method. Based on the analysis, is the following statement true or false? Both methods catch the major hotspots but the Optimized Hot Spot method identified much more confined smaller hot spots than DBSCAN. \_\_\_\_\_ (T or F)

- k) To make comparisons among service calls at census tracts, we need to consider population in the tracts, assuming that more people are like to have more calls for services.

Before continuing, check the TractCallPop attribute table. Notice that there are many tracts without any service calls. Delete these tracts first. The best way to do this is to sort the top category fields, select, delete, and save.

Therefore, create a table that includes the following: census tract, geographical identification code, population proportion in the tract against all population in the city, and proportion of each top-service calls in the city. Better to (1) add and set up all new fields make sure that the data type is float, and save and (2) keep both calculate field and chart properties windows open and replace the attribute names for computing through the top categories.



**Calculate Field**

This tool modifies the input data.

Input Table  
TractCallPop

Field Name (Existing or New)  
PopPercent

Expression Type  
Python 3

Expression

Fields Helpers

OBJECTID\_1  
Shape  
Geographic Identification Co  
Total Population (SF1)  
Animal  
Code  
Noise

Insert Values  
PopPercent =  
!sf\_totalpop!/1652237

**Chart Properties - TractCallPop**

Distribution of Total Population (SF1)

Data Axes Guides Format General

Variable  
Number  
Total Population (SF1)  
With transformation  
None  
☐ Show Normal distribution

Bins 20

Statistics

	Dataset
<input checked="" type="checkbox"/> Mean	4,325.2
<input type="checkbox"/> Median	4,076
<input type="checkbox"/> Std. Dev.	2,007.6
Rows	382
Count	382
Nulls	0
Min	10
Max	17,151
Sum	1,652,237
Skewness	1.22
Kurtosis	7.3

Data Labels  
☐ Label bins

I) So that we can map percentage of a top service-call type normalized by percentage of population in a census tract.

**Symbology - TractCallPop**

Primary symbology

Graduated Colors

Field  
SanitationPercent

Normalization  
PopPercent

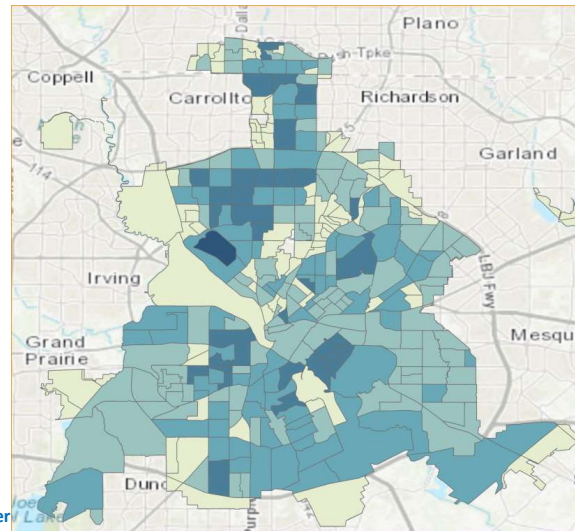
Method  
Natural Breaks (Jenks)

Classes  
5

Color scheme

Classes Histogram Scales

Symbol	Upper value	Label
Lightest Blue	≤ 0.533046	0.000 - 0.5330
Light Blue	≤ 1.377879	0.5331 - 1.378
Medium Blue	≤ 2.563234	1.379 - 2.563
Dark Blue	≤ 6.765117	2.564 - 6.765
Darkest Blue	≤ 82.872686	6.766 - 82.87



If the calls are in proportion to population proportion in census tracts, what proportion of calls we would expect for a census tract with 5% of total population in the City of Dallas? \_\_\_\_\_

Where were the highest calls for sanitation services in the City of Dallas? Adjust the transparency to see the geographic context. Don't give Geographic Identification Identifier, but name the geographic feature at that location \_\_\_\_\_

Check out all the service calls, which service call appears relatively confined to central and NW Dallas \_\_\_\_\_

Note that the default classes are based on natural breaks for each service call type. You can adjust the number of classes to see how they affect the map.

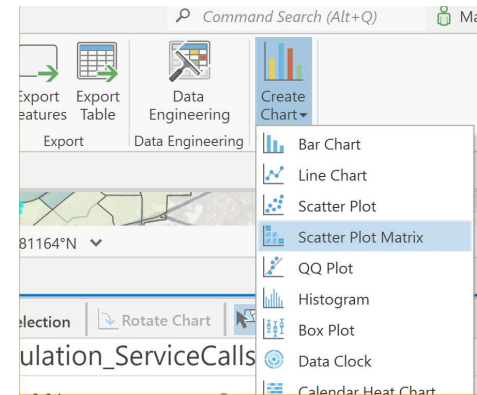
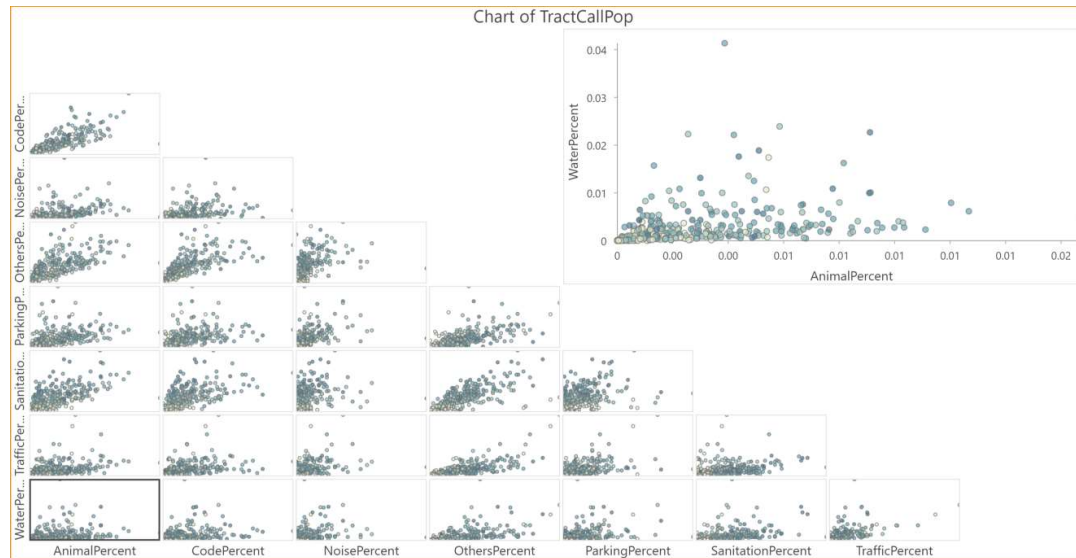
Dallas\_Population\_ServiceCalls - Chart of Tractv5\_Dallas\_Population\_S

Command Search (Alt+Q) May

Export Export Data Create

Where were the highest calls for sanitation services in the City of Dallas? Adjust the transparency to see the geographic context. Don't give Geographic Identification Identifier, but name the geographic feature at that location \_\_\_\_\_  
 Check out all the service calls, which service call appears relatively confined to central and NW Dallas \_\_\_\_\_

Note that the default classes are based on natural breaks in each service request type, so a darker shade in one may be equivalent to a light shade in another service type. Create Scatter Plot Matrix to see how these service calls relate.

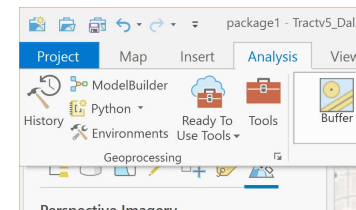


8. From the scatter plot matrix, which pairs of service request calls are more linearly correlated? a) Water and Animal, b) Code and Traffic, c) Animal and Code? \_\_\_\_\_(a, b, or c)  
 The scatter plot matrix suggests that noise, compared to other types of service calls, are relatively infrequent in most census tracts. \_\_\_\_\_(T or F)  
 The scatter plot matrix suggest that the frequency of sanitation calls varies widely across census tracts, but relatively few census tracts have many traffic calls. \_\_\_\_\_(T or F)

Now we want to find out what type of issue people called for service requests most in each census tract. We need to build a table for the analysis. The table should be something like:

Geographic Identification Code for Tracts	Type of Service calls has the max value
48113014132	Animal
.....	....

- m) While we can do this in ArcGIS Pro and Arcpy, but it is too cumbersome and not elegant. Instead, we will do the following
- export the attribute table to an excel file
  - Open the excel file and remove all fields but save the following: id, Animal, Code, Noise, Others, Parking, Sanitation, Traffic, and Water
  - Open the python notebook in ArcGIS Pro
  - Type the following codes in the notebook



```

Edit View Insert Cell Help
+ [Icons] Run Code
Replace the path with yours

In [ ]: import pandas as pd
    
```



```

Edit View Insert Cell Help
+ % Copy Paste Undo Redo Run Code
Replace the path with yours

In [ ]: import pandas as pd

In [ ]: df = pd.read_excel(r'D:\SDSprojects\package1_1\TractCallPop_TableToExcel.xlsx')

In [ ]: df.columns

In [ ]: df.iloc[:, 2:]

In [ ]: df['MaxType'] = df.iloc[:, 2:].idxmax(axis=1)
Replace the path with yours

In [ ]: df['MaxType']

In [ ]: df.to_csv(r'D:\SDSprojects\package1_1\MaxType.csv')

```

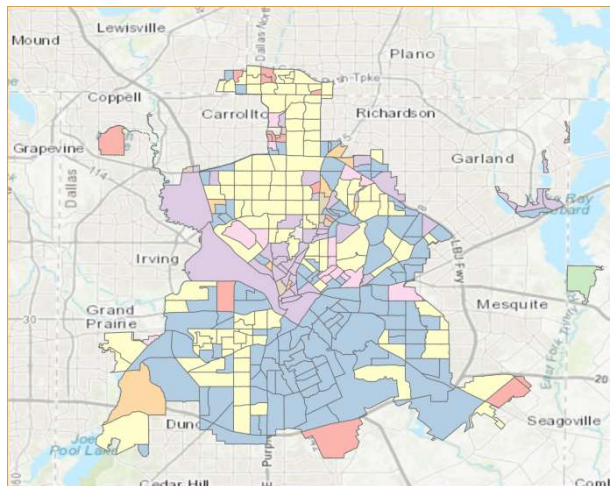
Type equation here.

- Open the csv file and delete all fields but id and MaxType. Open the output csv file in ArcGIS Pro, We will join it with the Tractv5 \_Dallas\_Population\_ServiceCalls feature class. However, we need to first export the csv file (df.csv) to a geodatabase table (DF). Now join the DF table to the Tract feature class. Click on "Validate Join" to see if any issues.

9. Why did the join fail? a) Cannot join fields with different names, b) Cannot join fields with different data types, c) Cannot join a feature class with a table. \_\_\_\_\_ (a, b, or c)

How to fix the error? a) Change "id" in DF to "Geographic Identification Code", b) Change "id" to the same data type as "Geographic Identification Code", c) Add a new field "geoid" with the same data type as "Geographic Identification Code" and copy the values from id. \_\_\_\_\_ (a, b, or c)

- After successfully join the tables, map the MaxType. Your map should look something like below:



10. What was the most frequent types of service requests from most of south Dallas? \_\_\_\_\_  
What was the most frequent types of service requests from most of north Dallas? \_\_\_\_\_
11. The map shows that census tracts with parking issues as the most frequent types of service requests are relatively much smaller (except for one) than other census tracts. What may be the potential explanation?
12. Does the spatial pattern surprise you? Why or why not? What spatial questions does the map prompt you?