# COMMUNITY CONTRIBUTION
## TITLE- TUTORIAL ON DIVE (OPEN-SOURCE TOOL)
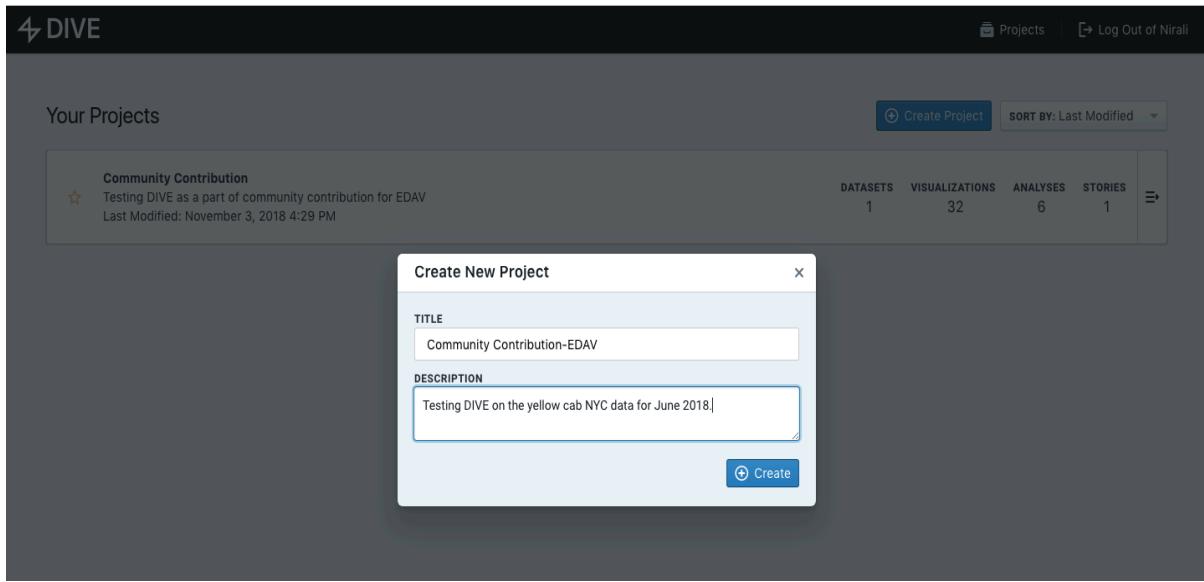## AUTHOR: NIRALI SHAH, UNI: NSS2173

DIVE
- DIVE is an open-source tool for recommendation-driven data exploration and visualization developed by Kevin Hu and César Hidalgo at the MIT Media Lab.
- It is a web-based data exploration system that allows users to visualize and explore their data without going through the hassle of writing any code.
- DIVE provides a single platform to perform:
    1. Semantic Data Ingestion
    2. Recommendation-based data visualization and analysis
    3. Dynamic story sharing
- Dive is an easy-to-learn data exploration tool.
- Dive is a mixed-initiative system combining recommender systems with point-and-click manual specification to support state-of-the-art data model inference, visualization, statistical analysis, and storytelling capabilities. [1]
- Tasks performed in order[2]:
    1. Intelligent Data Ingestion
    2. Semi-automated visualization recommendation
    3. Point-and-Click statistical analysis
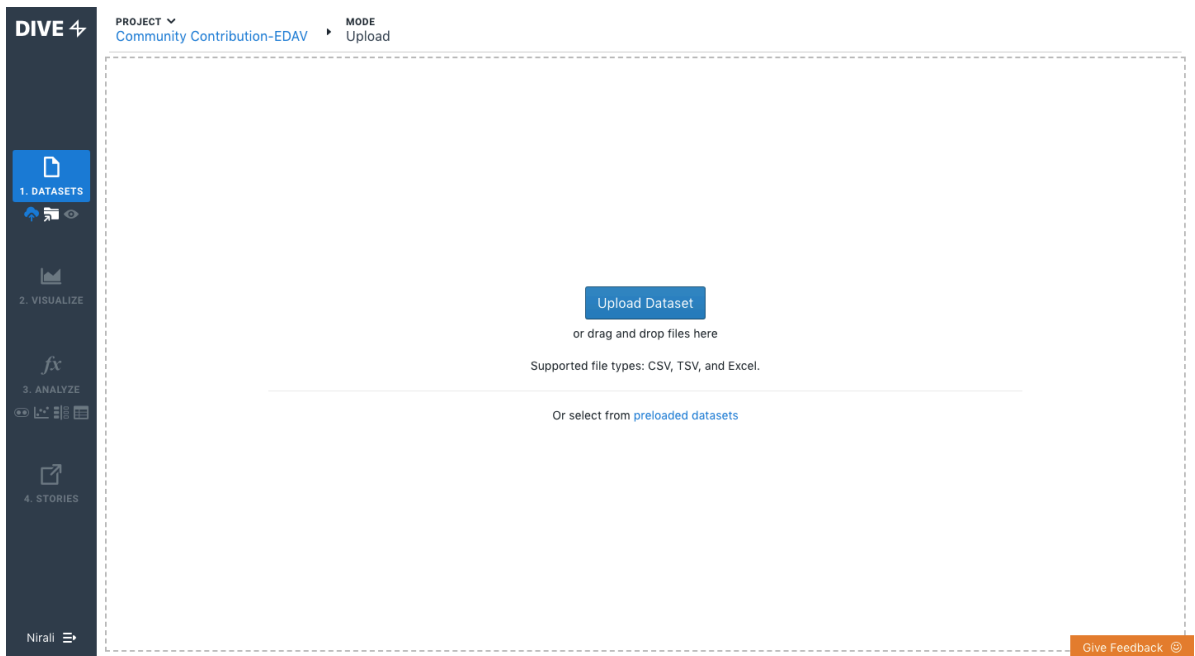    4. WYSIWYG (What You See Is What You Get) visual narratives

Steps to create a new project or work on an exisiting one:
1. To use the tool, visit the site https://dive.media.mit.edu/ .The website also contains a demo video that helps us understand the tool.
2. Create an account. Although the tool can be used without creating an account, it's better to create an account if one wants to switch between different projects while using the tool.
3. Log in to see the projects you've worked on. You can also create a new project.

I will be working on the NYC yellow cab data of June 2018.

# 1. Upload your dataset and inspect it using Inspect Mode:
It supports 3 file types: CSV, TSV and Excel



Once the data is uploaded, it computes dataset field properties and enters the **Inspect** mode.
The Inspect mode shows all the variables in the dataset and their field type.

Additional properties of each numerical or decimal variable:
1. Whether the variable is an ID.
2. Mean value
3.Median Value
4. Standard Deviation
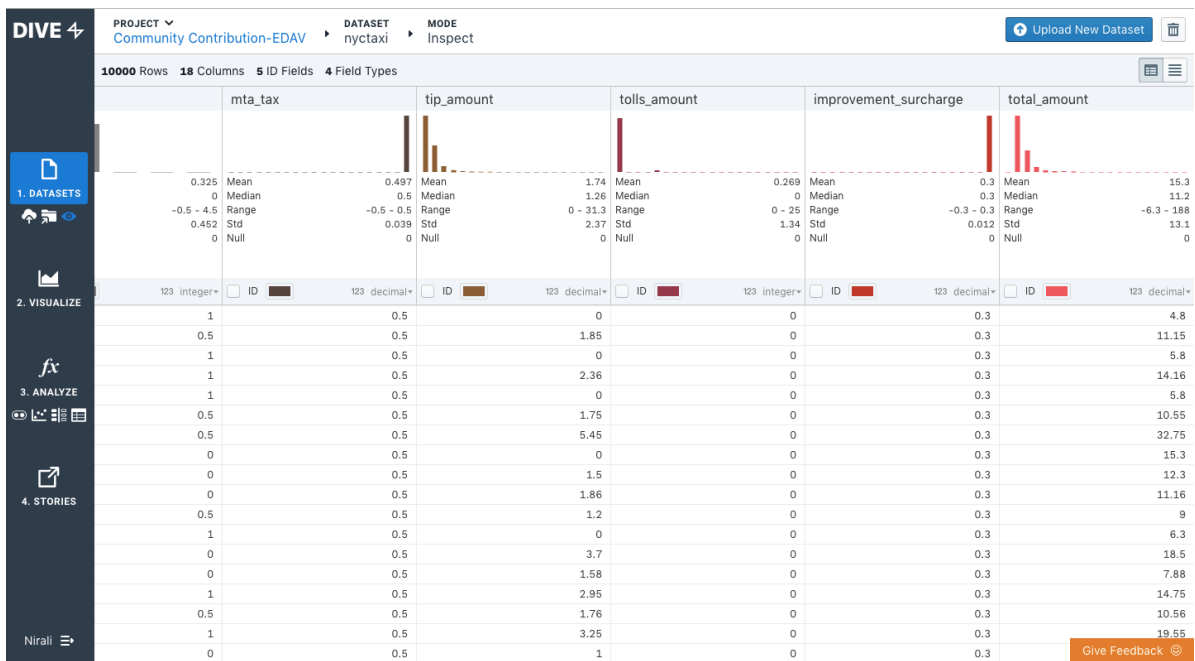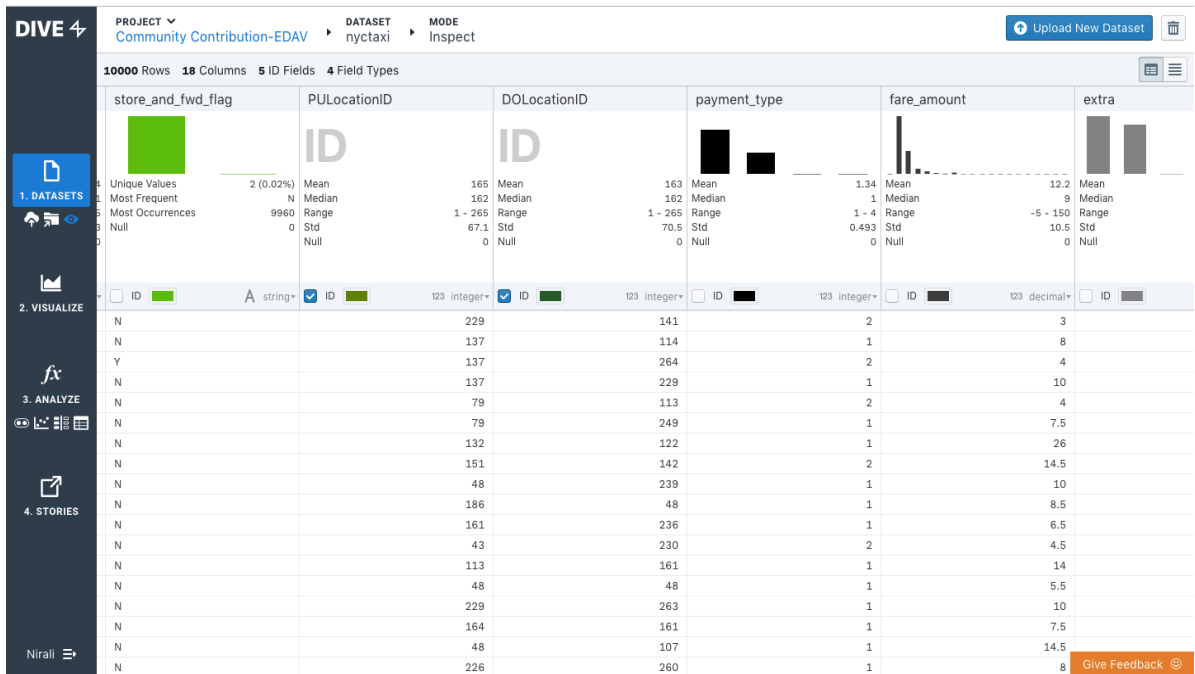5. Number of null values

For a string variable, it displays:
1. Number of unique values
2. Most frequent value
3. Number of occurences of the most frequent value
4. Number of null values

For a Date Time variable, it displays the range of values.
This helps us understand the dataset.
The type of the field can be changed and a field can be marked as a unique identifier if needed.
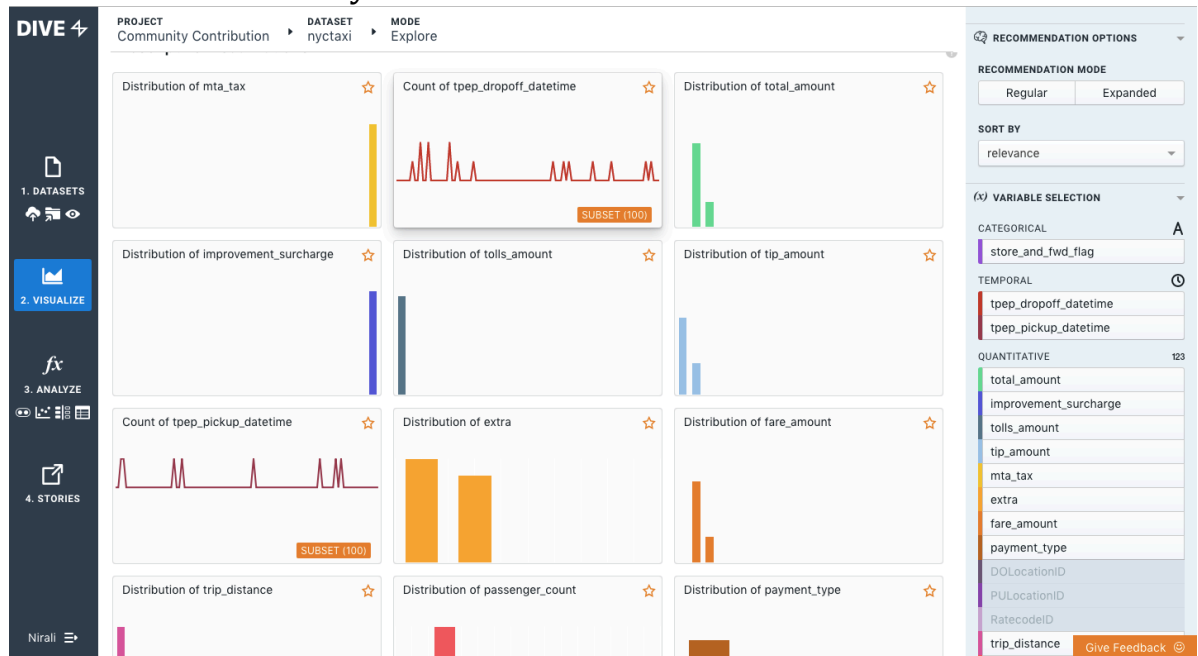
**PROJECT** ✓
Community Contribution-EDAV   **DATASET** nyctaxi   **MODE** Inspect   ⬆ Upload New Dataset  🗑

**10000** Rows  **18** Columns  **5** ID Fields  **4** Field Types

1. DATASETS
2. VISUALIZE
3. ANALYZE
4. STORIES

Nirali

| | store_and_fwd_flag | | PULocationID | | DOLocationID | | payment_type | | fare_amount | | extra |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Unique Values | 2 (0.02%) | Mean | 165 | Mean | 163 | Mean | 1.34 | Mean | 12.2 | Mean | |
| Most Frequent | N | Median | 162 | Median | 162 | Median | 1 | Median | 9 | Median | |
| Most Occurrences | 9960 | Range | 1 - 265 | Range | 1 - 265 | Range | 1 - 4 | Range | -5 - 150 | Range | |
| Null | 0 | Std | 67.1 | Std | 70.5 | Std | 0.493 | Std | 10.5 | Std | |
| | | Null | 0 | Null | 0 | Null | 0 | Null | 0 | Null | |

| ☐ ID | A string▾ | ☑ ID | 123 integer▾ | ☑ ID | 123 integer▾ | ☐ ID | 123 integer▾ | ☐ ID | 123 decimal▾ | ☐ ID | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| N | | 229 | | 141 | | 2 | | 3 | | | |
| N | | 137 | | 114 | | 1 | | 8 | | | |
| Y | | 137 | | 264 | | 2 | | 4 | | | |
| N | | 137 | | 229 | | 1 | | 10 | | | |
| N | | 79 | | 113 | | 2 | | 4 | | | |
| N | | 79 | | 249 | | 1 | | 7.5 | | | |
| N | | 132 | | 122 | | 1 | | 26 | | | |
| N | | 151 | | 142 | | 2 | | 14.5 | | | |
| N | | 48 | | 239 | | 1 | | 10 | | | |
| N | | 186 | | 48 | | 1 | | 8.5 | | | |
| N | | 161 | | 236 | | 1 | | 6.5 | | | |
| N | | 43 | | 230 | | 2 | | 4.5 | | | |
| N | | 113 | | 161 | | 1 | | 14 | | | |
| N | | 48 | | 48 | | 1 | | 5.5 | | | |
| N | | 229 | | 263 | | 1 | | 10 | | | |
| N | | 164 | | 161 | | 1 | | 7.5 | | | |
| N | | 48 | | 107 | | 1 | | 14.5 | | | |
| N | | 226 | | 260 | | 1 | | 8 | | | |

Give Feedback ☺

---

**PROJECT** ✓
Community Contribution-EDAV   **DATASET** nyctaxi   **MODE** Inspect   ⬆ Upload New Dataset  🗑

**10000** Rows  **18** Columns  **5** ID Fields  **4** Field Types

1. DATASETS
2. VISUALIZE
3. ANALYZE
4. STORIES

Nirali

| | | mta_tax | | tip_amount | | tolls_amount | | improvement_surcharge | | total_amount |
|---|---|---|---|---|---|---|---|---|---|---|
| | 0.325 | Mean | 0.497 | Mean | 1.74 | Mean | 0.269 | Mean | 0.3 | Mean | 15.3 |
| | 0 | Median | 0.5 | Median | 1.26 | Median | 0 | Median | 0.3 | Median | 11.2 |
| | -0.5 - 4.5 | Range | -0.5 - 0.5 | Range | 0 - 31.3 | Range | 0 - 25 | Range | -0.3 - 0.3 | Range | -6.3 - 188 |
| | 0.452 | Std | 0.039 | Std | 2.37 | Std | 1.34 | Std | 0.012 | Std | 13.1 |
| | 0 | Null | 0 | Null | 0 | Null | 0 | Null | 0 | Null | 0 |

| 123 integer▾ | ☐ ID | 123 decimal▾ | ☐ ID | 123 decimal▾ | ☐ ID | 123 integer▾ | ☐ ID | 123 decimal▾ | ☐ ID | 123 decimal▾ |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | 0.5 | | 0 | | 0 | | 0.3 | | 4.8 |
| 0.5 | | 0.5 | | 1.85 | | 0 | | 0.3 | | 11.15 |
| 1 | | 0.5 | | 0 | | 0 | | 0.3 | | 5.8 |
| 1 | | 0.5 | | 2.36 | | 0 | | 0.3 | | 14.16 |
| 1 | | 0.5 | | 0 | | 0 | | 0.3 | | 5.8 |
| 0.5 | | 0.5 | | 1.75 | | 0 | | 0.3 | | 10.55 |
| 0.5 | | 0.5 | | 5.45 | | 0 | | 0.3 | | 32.75 |
| 0 | | 0.5 | | 0 | | 0 | | 0.3 | | 15.3 |
| 0 | | 0.5 | | 1.5 | | 0 | | 0.3 | | 12.3 |
| 0 | | 0.5 | | 1.86 | | 0 | | 0.3 | | 11.16 |
| 0.5 | | 0.5 | | 1.2 | | 0 | | 0.3 | | 9 |
| 1 | | 0.5 | | 0 | | 0 | | 0.3 | | 6.3 |
| 0 | | 0.5 | | 3.7 | | 0 | | 0.3 | | 18.5 |
| 0 | | 0.5 | | 1.58 | | 0 | | 0.3 | | 7.88 |
| 1 | | 0.5 | | 2.95 | | 0 | | 0.3 | | 14.75 |
| 0.5 | | 0.5 | | 1.76 | | 0 | | 0.3 | | 10.56 |
| 1 | | 0.5 | | 3.25 | | 0 | | 0.3 | | 19.55 |
| 0 | | 0.5 | | 1 | | 0 | | 0.3 | | |

Give Feedback ☺

## 2. Descriptive Visualizations:

This section showcases univariate summary visualizations. To add a visualization to a story, the graph must be saved.

## Visualizations of every variable



We are interested in the following graphs:
1. trip_distance
2. total_amount
3. tip_amount

The following display options are available:
1. Binning configuration:
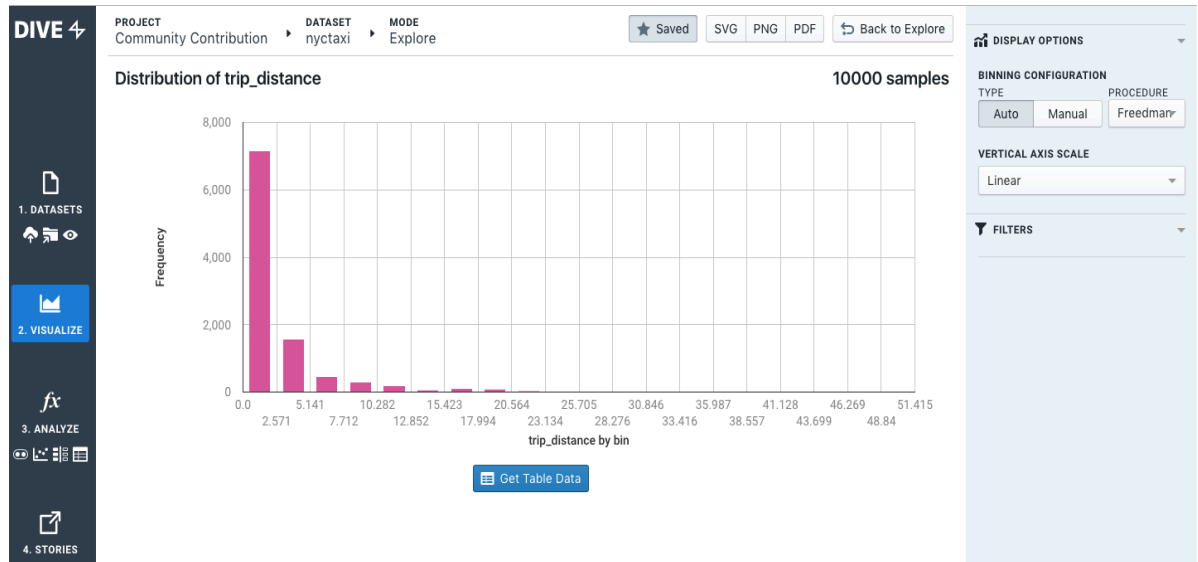   Type- Use automated number of bins or manually set the number of bins
Procedure- Select a procedure from a given list of procedures to create a visualization.
2. Vertical Axis Scale- it can be set to either Linear or Logarithmic.
3. Filter- Filters the data using multiple conditions. These conditions are based on values of different fields. Multiple conditions can be combined using OR or AND.

# trip_distance

## Default distribution of trip_distance



## Distribution of trip_distance with number of bins=21 and a logarithmic vertical scale

# tip_amount

## Default distribution of tip_amount



## Distribution of tip_amount with bins=22 and trip_distance>0 and total_amount>0

# total_amount

## Default distribution of total_amount



## Distribution of total_amount with bins=23 and trip_distance>0 and tip_amount ≥ 1

# 3. Analyze mode:

This mode gives 4 options:

a) Aggregation
b) Correlation
c) Comparison
d) Regression


## a) Aggregation

This mode helps aggregate the data based on the variables selected.
On selecting two variables, the values of the variables are divided into different bins and a count for combinations of bins of the two variables is displayed.
This aggregation can be performed on any other variable of the dataset.
Example: aggregating total_amount and tip_amount on trip_distance.
One can choose to display either the sum of trip_distance or the mean of trip_distance for every combination of bins of total_amount and tip_amount.
The binning of the variables chosen can also be changed and filters can be applied just as explained previously.


Aggregating total_amount and tip_amount on Count



| TIP_AMOUNT | -6.300-3.396 | 3.396-13.091 | 13.091-22.787 | 22.787-32.482 | 32.482-42.178 | 42.178-51.873 | 51.873-61.569 | 61.569-71.264 | 71.2 |
|---|---|---|---|---|---|---|---|---|---|
| 0.000-1.563 | 0 | 1971 | 195 | 15 | 6 | 0 | 1 | 0 | |
| 1.563-3.127 | 0 | 1561 | 1115 | 78 | 17 | 7 | 2 | 0 | |
| 3.127-4.691 | 0 | 6 | 462 | 201 | 14 | 5 | 1 | 2 | |
| 4.691-6.254 | 0 | 7 | 21 | 143 | 73 | 13 | 6 | 14 | |
| 6.254-7.817 | 0 | 0 | 2 | 19 | 76 | 39 | 0 | 3 | |
| 7.817-9.381 | 0 | 0 | 1 | 1 | 13 | 49 | 28 | 6 | |
| 9.381-10.944 | 0 | 0 | 4 | 0 | 4 | 22 | 22 | 21 | |
| 10.944-12.508 | 0 | 0 | 0 | 0 | 0 | 3 | 11 | 55 | |
| 12.508-14.072 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 7 | |
| 14.072-15.635 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 2 | |
| 15.635-17.198 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 17.198-18.762 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 18.762-20.326 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | |
| 20.326-21.889 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 21.889-23.453 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | |
| 23.453-25.016 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 25.016-26.579 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 26.579-28.143 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | |
| 28.143-29.706 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 29.706-31.273 | 27 | 2636 | 586 | 150 | 79 | 30 | 67 | 12 | |
| Column Totals | 27 | 6181 | 2387 | 608 | 283 | 168 | 143 | 122 | |

# Aggregating total_amount and tip_amount by mean of trip_distance
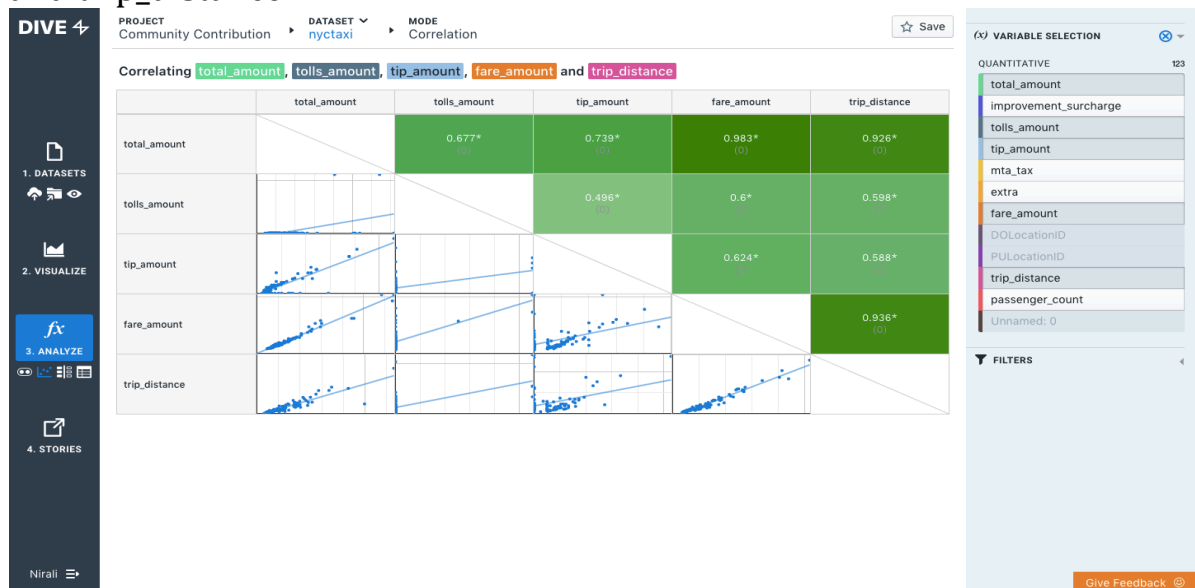


## b) Correlation

This mode helps find the relationship between the variables in a dataset. It displays a correlation matrix that displays the correlation between the variables to indicate if the variables are positively correlated, negatively correlated or not correlated.

We just need to select the variables in the right pane and a correlation matrix will be generated. Filters can also be applied.

## Correlation between total_amount, tip_amount, fare_amount and trip_distance

# Correlation between total_amount, tip_amount, tolls_amount, fare_amount and trip_distance



We see a positive correlation of 0.739 between total_amount and tip_amount.

## c) Comparison

This mode performs a one-way analysis of variance (ANOVA) comparing the values of the dependent variable for different groups of the independent variable. It tells us if the groups of the independent variable are distinct with a p-value. It also performs Post-hoc pairwise comparison between 2 groups of the dependent variable using the Tukey HSD test if those groups are distinct.

## Comparing tip_amount by passenger_count

## d) Regression

Here, we need to select the dependent variable. By clicking on 'Recommend Model', it returns a set of models with regression coefficients and constants along with the contribution of the independent variables to the value of the dependent variable. We can also manually select a list of independent variables of our choice along with the degree (x or $x^2$) to test their contribution to the dependent variable. It also displays $\bar{R}^2$, which is the amount of variance explained by independent variables.

## tip_amount in terms of total_amount



## tip_amount in terms of fare_amount

tip_amount in terms of total_amount and fare_amount
This is the recommended linear regression model with regression type as
'Forward Selection on $R^2$ '





We see that total_amount contributes more as compared to fare_amount.

## 4. Stories
In this mode, all the saved visualizations appear in the right pane.
To include a particular visualization in the story, we just need to click on
it. The visualization appears in the story along with some space alloted for
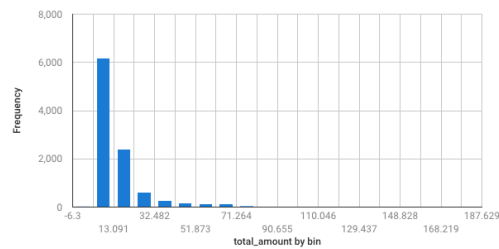text, which can be used to explain the visualization.
In this way, an entire document can be created with all the necessary
visualizations and this document can be shared.
Once you click the 'Share' option, it creates a web page for your document.
This allows dynamic sharing of your work.
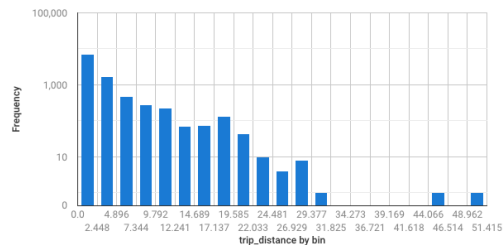
Document created.

## References:

[1] Kevin Hu, Diana Orghian, and César Hidalgo. 2018. DIVE: A Mixed-Initiative System Supporting Integrated Data Exploration Workflows. In ACM SIGMOD Workshop on Human-In-the-Loop Data Analytics (HILDA), June 10, 2018, Houston, TX, USA. ACM, New York, NY, USA, Article 4, 7 pages.
https://doi.org/10.1145/3209900.3209910

[2] DIVE Tool videos
https://dive.media.mit.edu/

[3] DIVE demo video
https://www.youtube.com/watch?v=J3FceN2lYdA&t=17s