# Datathon Semester 2 Preparation Workshop
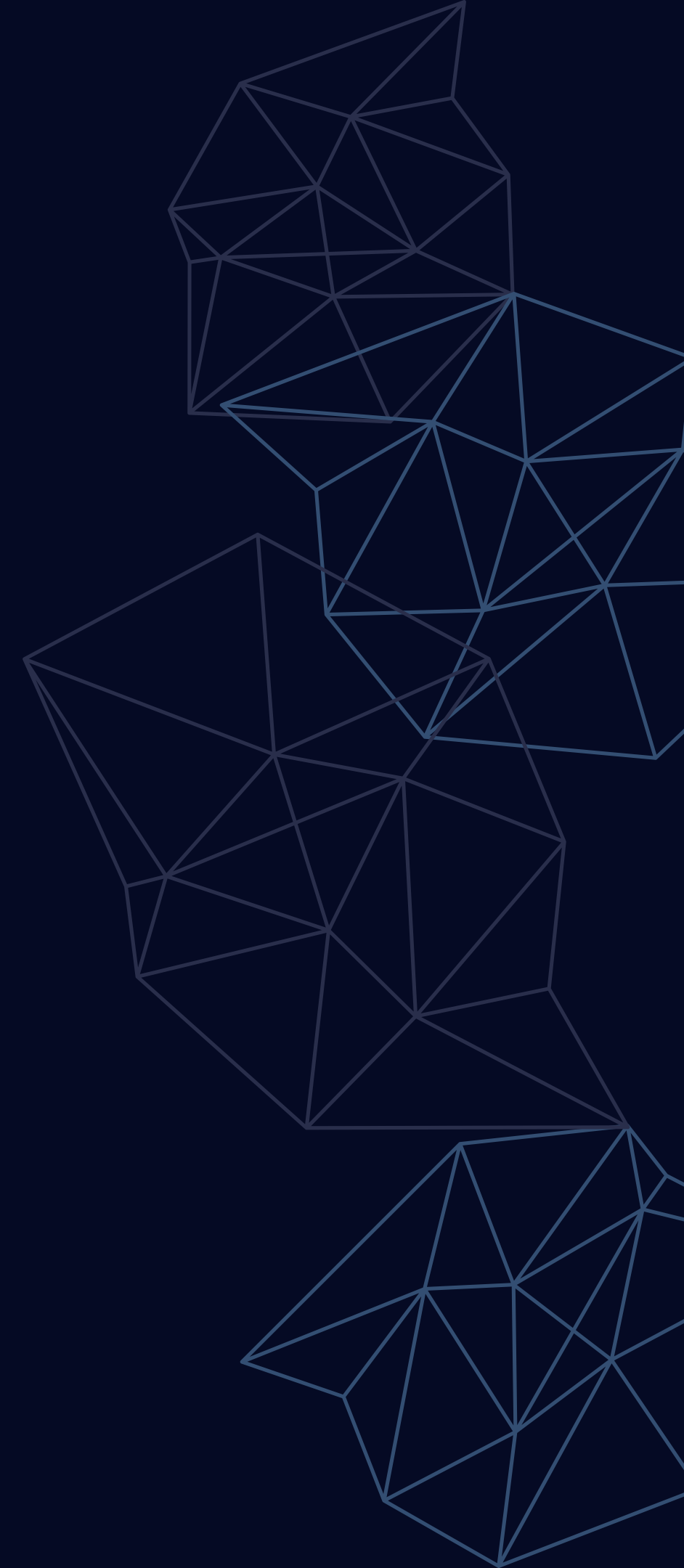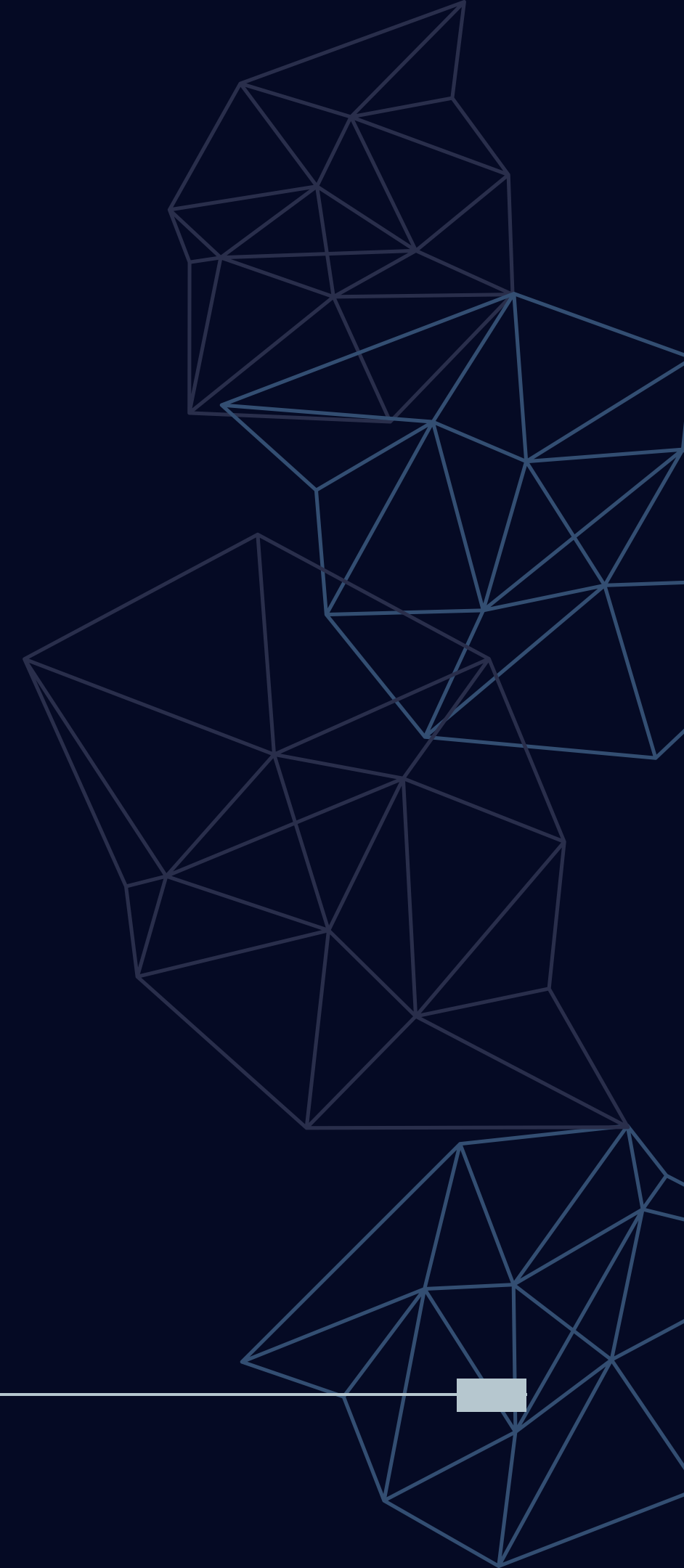
{mdss}

# Today's Agenda

1. Introduction to the MDSS Datathon

2. Overview of Data Cleaning/Wrangling

3. Data Visualization

4. Q&A

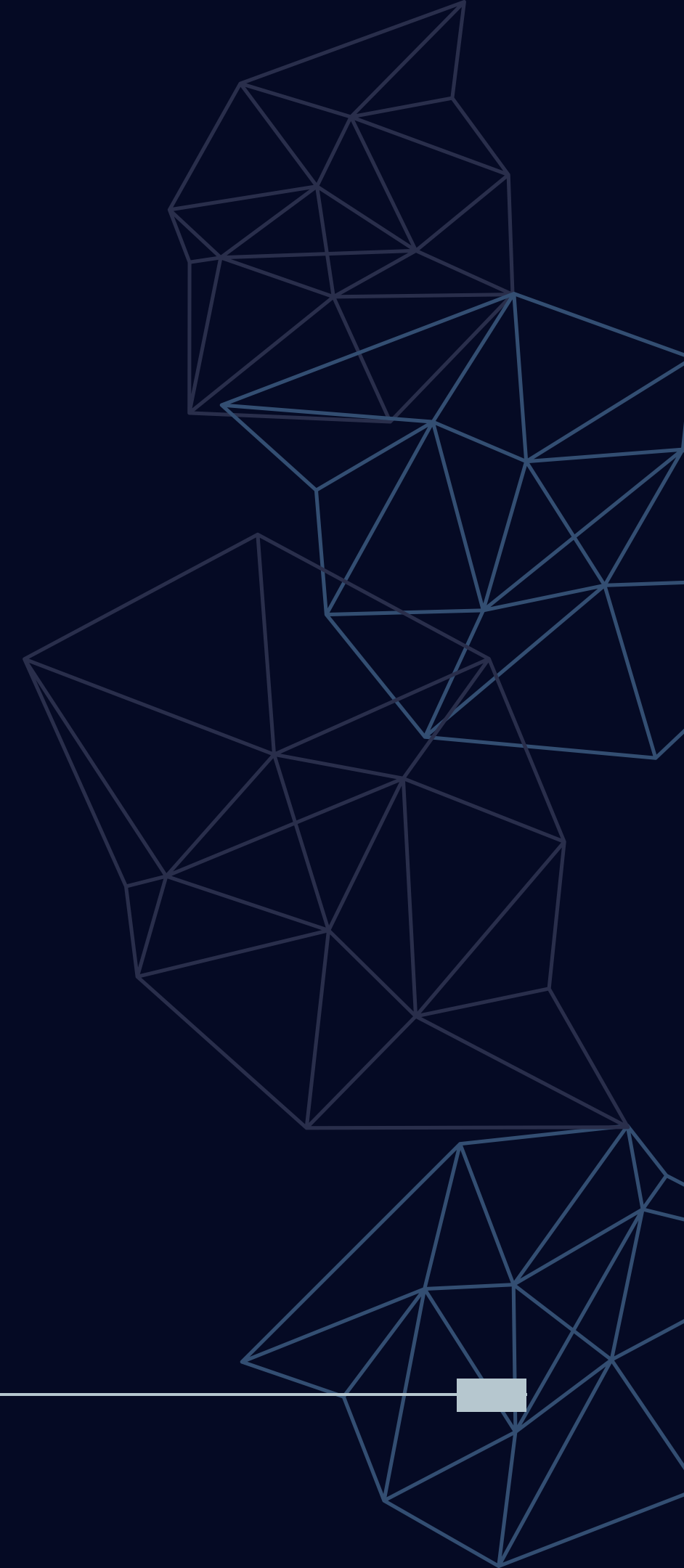{mdss} | MDSS
DATATHON
WORKSHOP

# What is a Datathon?

Collaborate, innovate, and demonstrate data science skills to analyse and visualise data in a insightful and creative way.

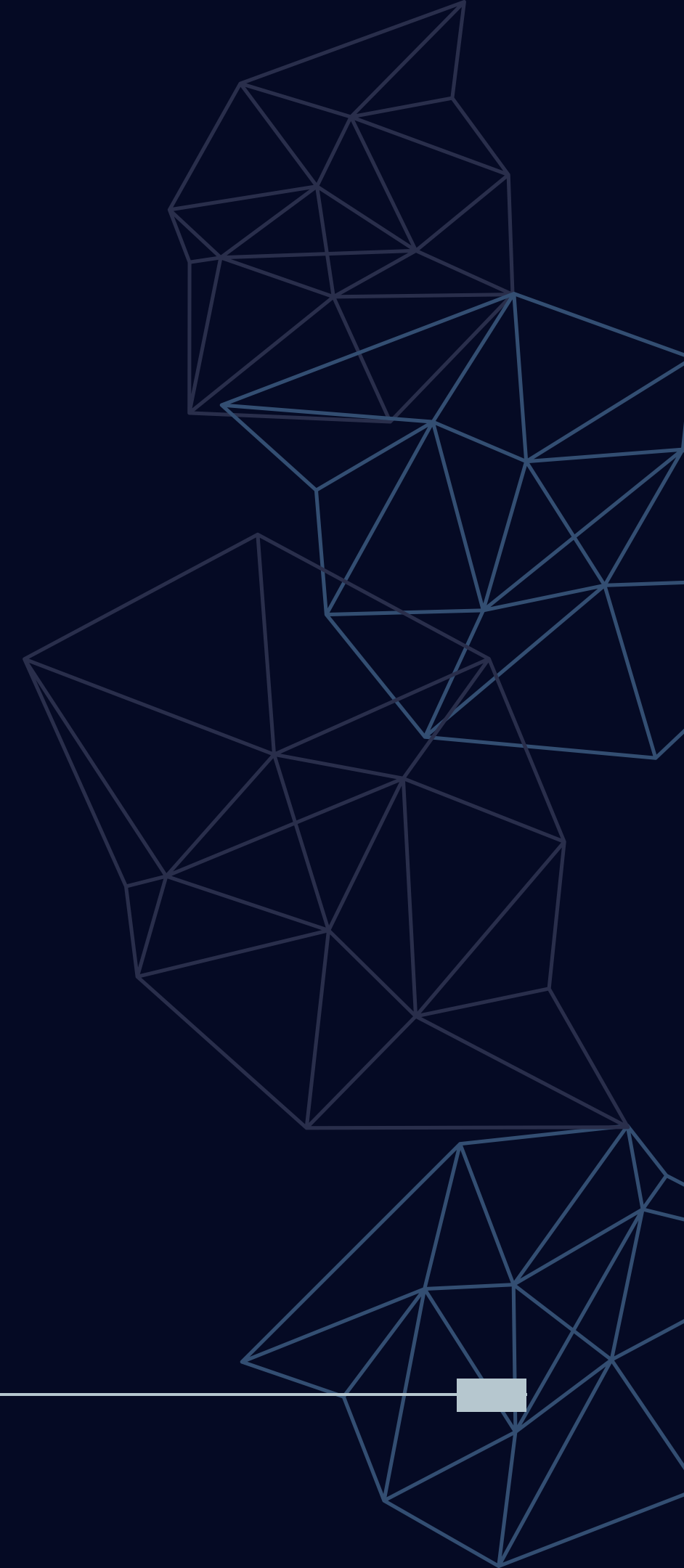This event is for both seasoned data experts and beginners

{mdss}

MDSS
DATATHON
WORKSHOP

# Data Cleaning/ Wrangling

# Data Cleaning Basics

Data cleaning refers to the process of identifying and correcting errors, inconsistencies, and inaccuracies in datasets.

1. Handling duplicates

2. Handling missing values

3. Handling outliers

4. Handling inconsistent data

# Data Wrangling

It is the process of transforming and reshaping raw data into a format suitable for analysis.

1. Data acquisition
2. Data cleaning
3. Data integration
4. Data reshaping
5. Data reduction
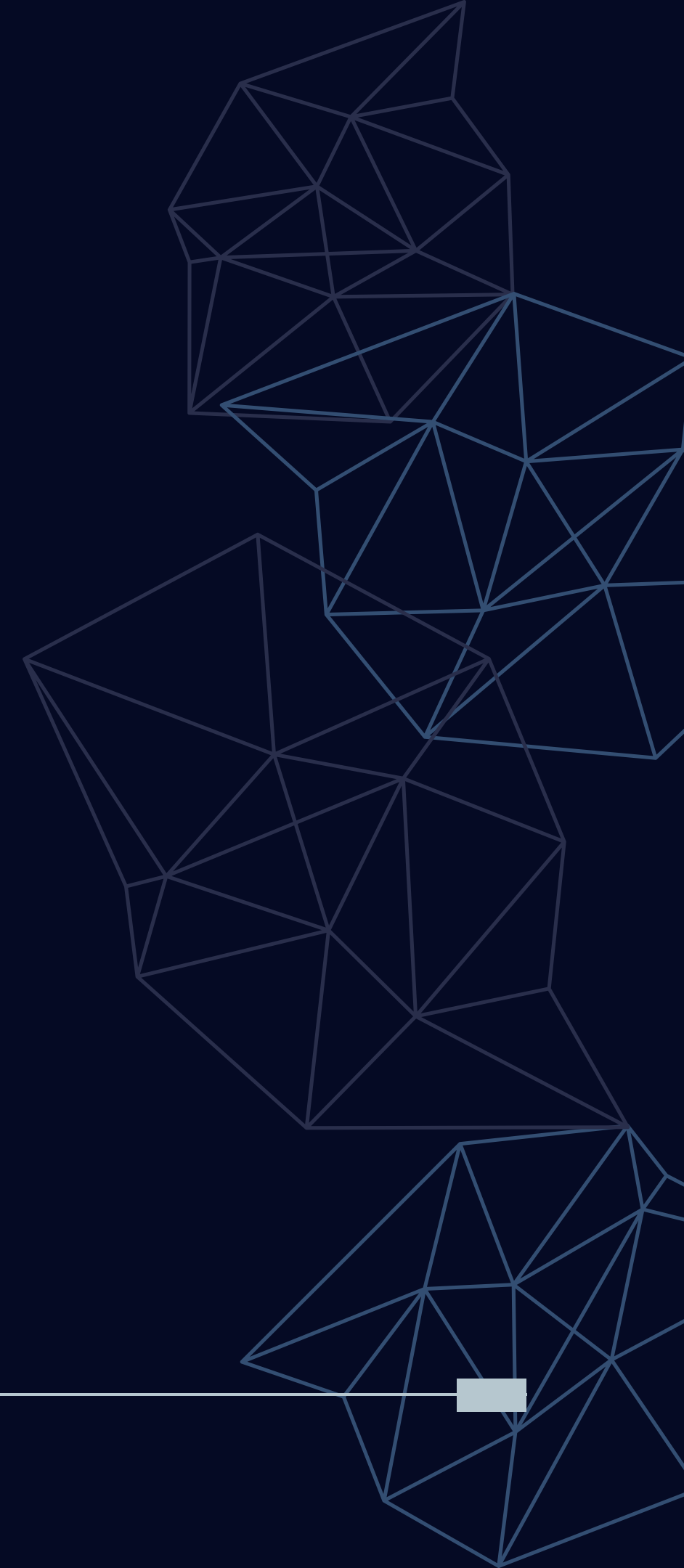6. Data validation

Repeat as many times necessary

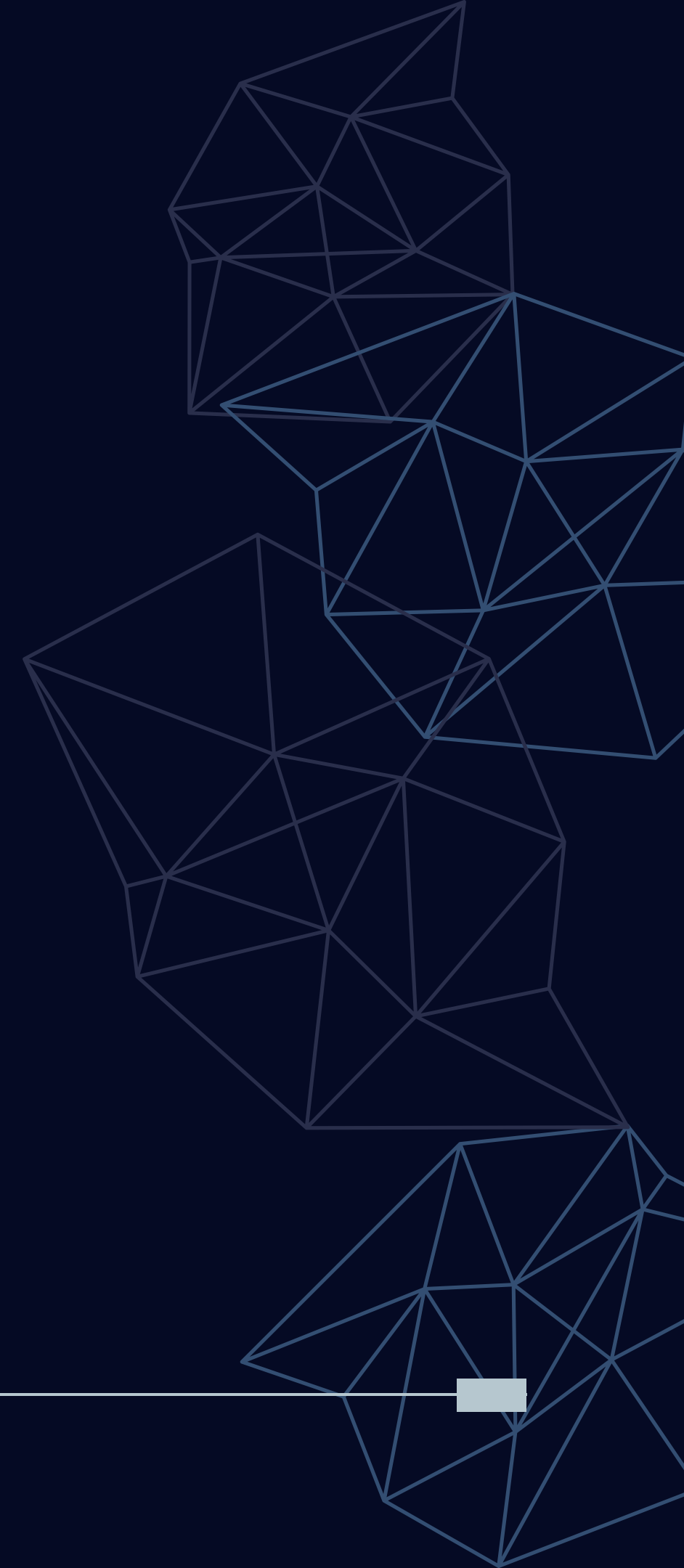{mdss}

MDSS
DATATHON
WORKSHOP

Notebook Demonstration

{mdss}

MDSS
DATATHON
WORKSHOP

Data Exploration

MDSS
DATATHON
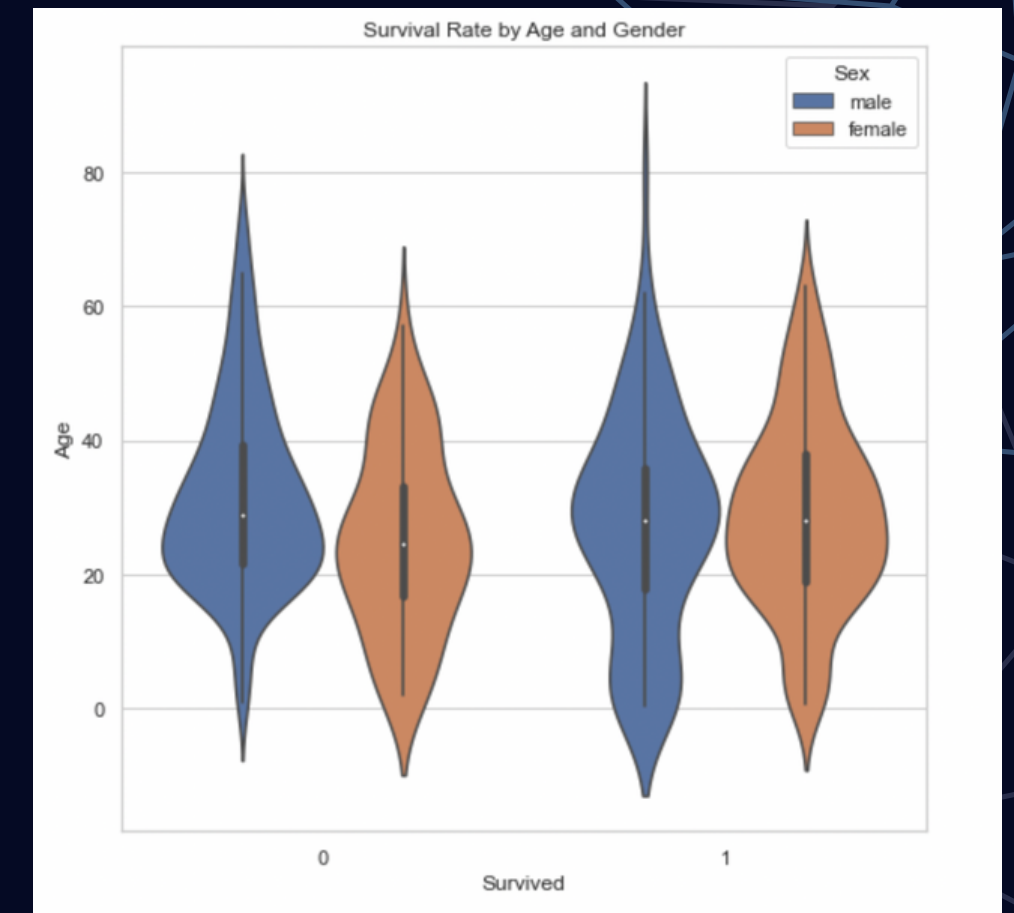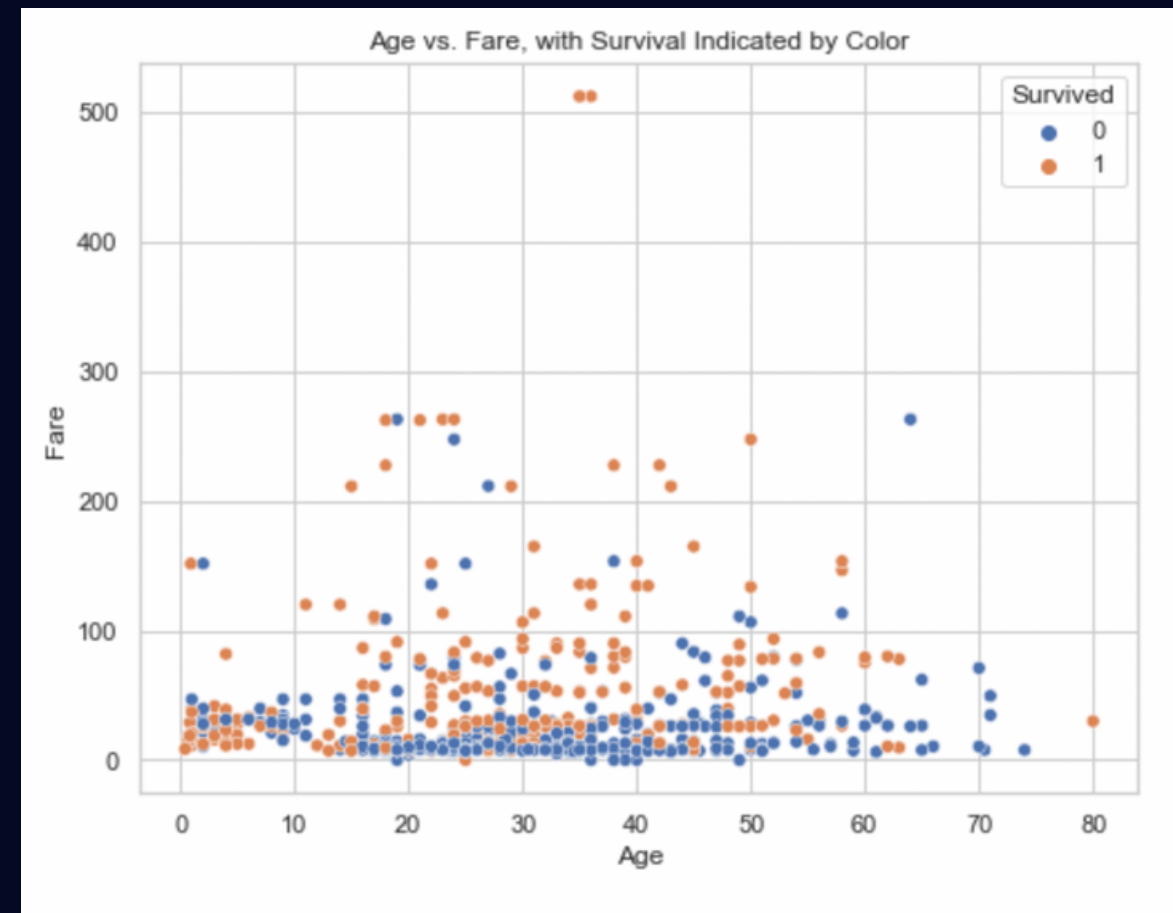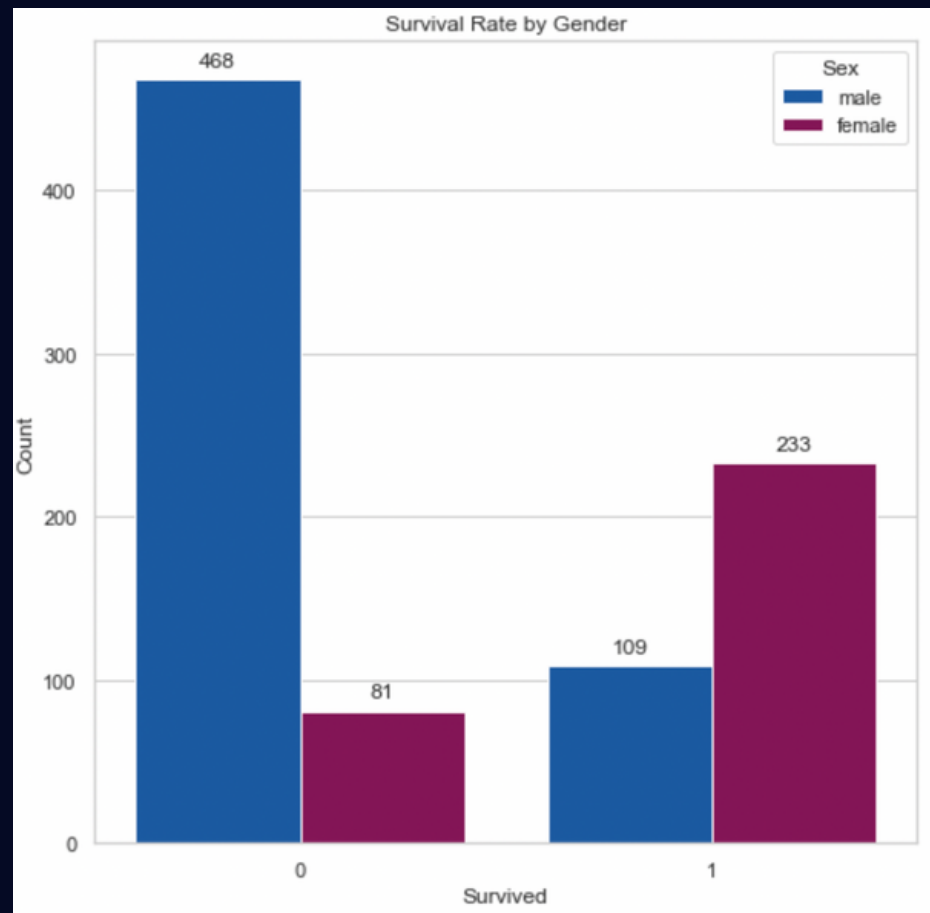WORKSHOP

# Exploration

## Common techniques:

- Summary statistics - basic statistics, such as mean, median, and standard deviation.
- Data visualization - creating graphs and charts, such as histograms, scatterplots, and box plots, to visually explore the data and identify patterns and relationships.
- Data transformation - converting or manipulating the data, such as scaling or normalizing variables, to better understand the relationships between variables.
- Dimensionality reduction - involves reducing the number of variables in the data, such as using principal component analysis (PCA), to simplify the analysis and identify the most important variables.

"the process of examining and understanding the data before performing any formal analysis."

{mdss}

MDSS
DATATHON
WORKSHOP

# Example Graphs



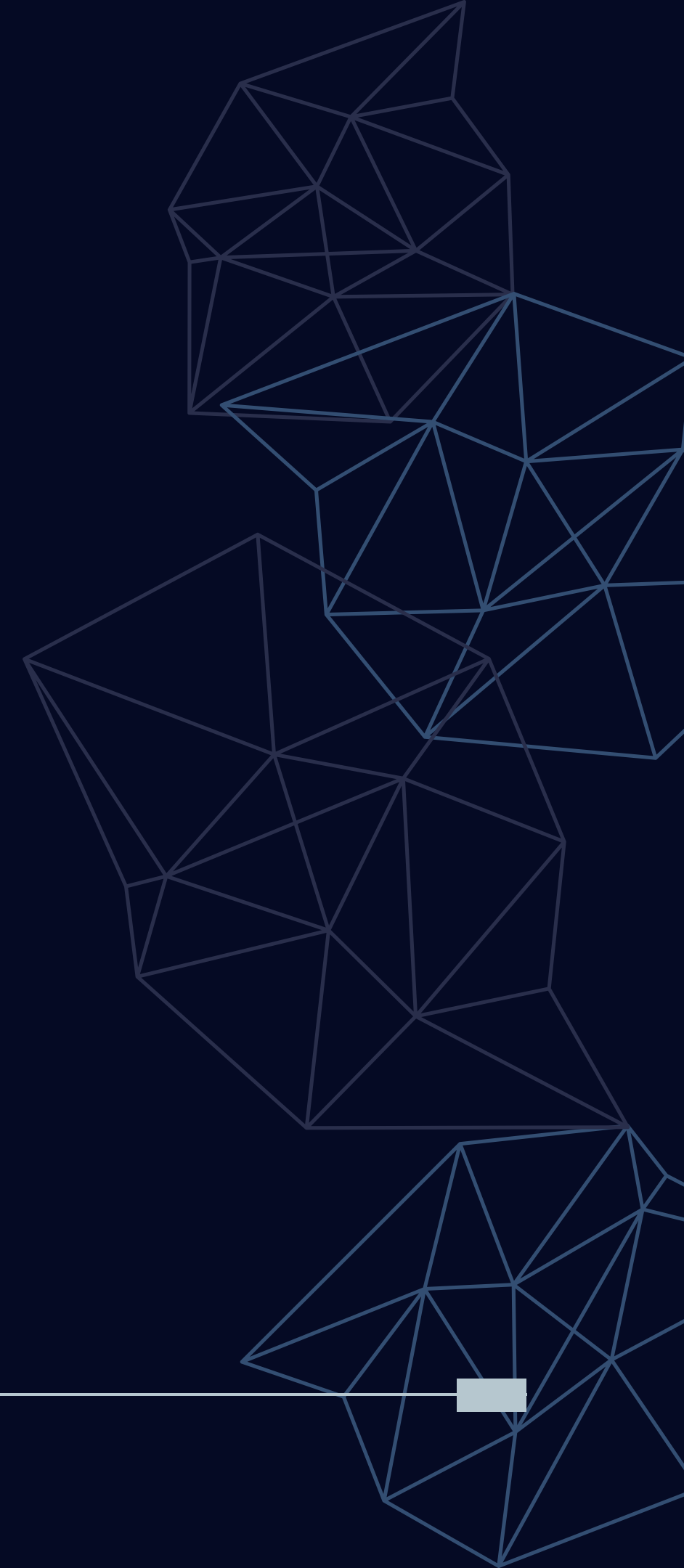These visualisations are created using the famous Titanic ML competition dataset - which we will explore shortly!

{mdss}

MDSS
DATATHON
WORKSHOP

# Data Set Used

https://www.kaggle.com/c/titanic

# Q&A