*MERC 2025 - Session 4*

# EXPLORATORY DATA ANALYSIS

- Introduction to Exploratory Data Analysis.

- EDA – Case Study

# Introduction to Exploratory Data Analysis (EDA)

## QC and Imputation
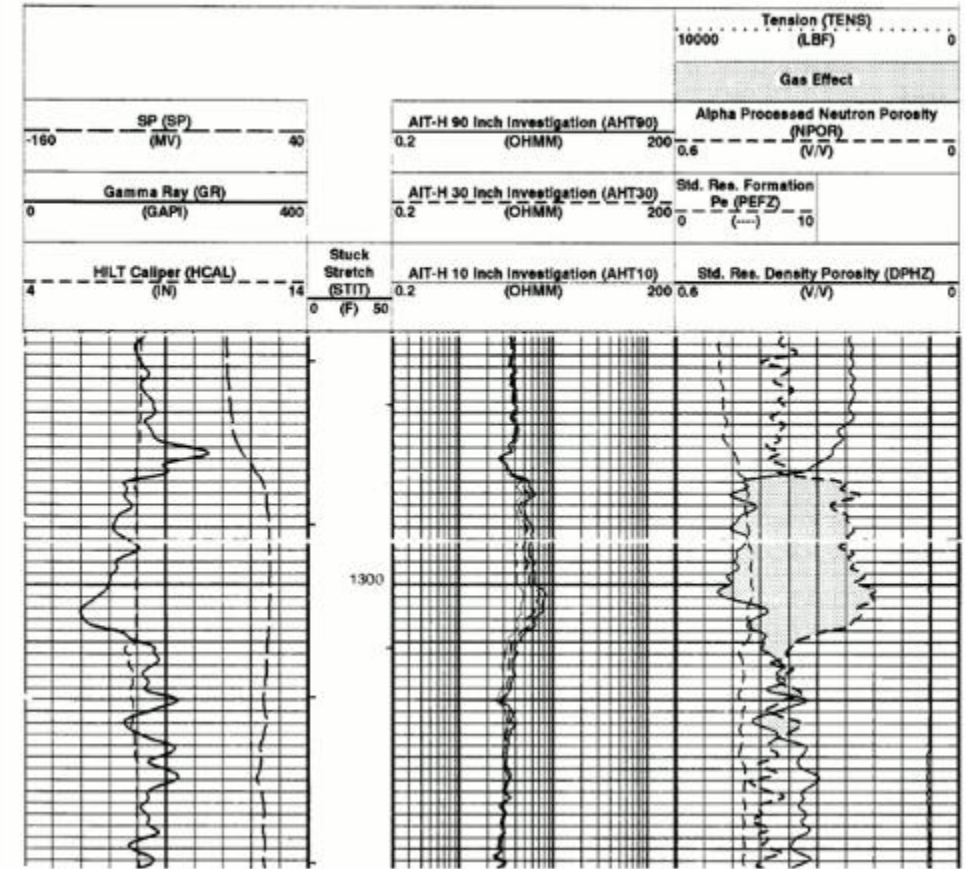- Visual Inspection
- Borehole Conditions
- Fill missing data

## Basic Stats
- Distributions
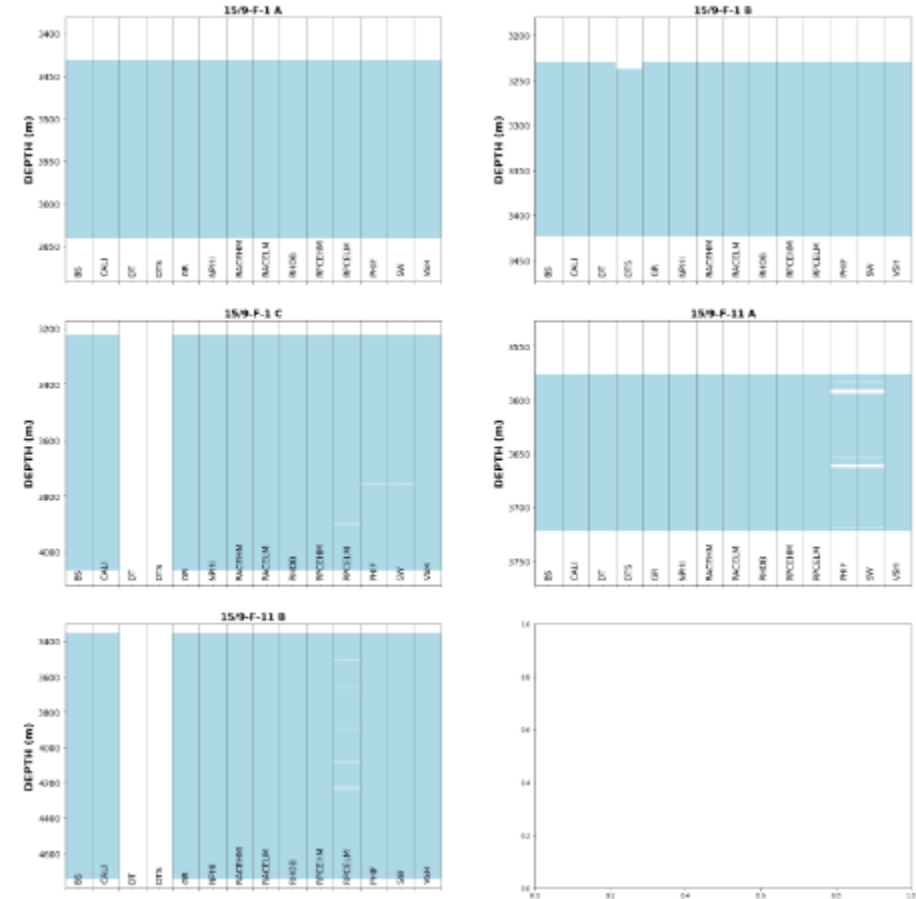- Categorial Statistics
- Cross-plots
- Box and Violin Plots

## Machine Learning Techniqus
- PCA
- t-SNE
- Cluster Analysis
- Self-Organizing Maps (SOMs)

- Manual review of the data.
  - Examining for large spikes or unexplained changes in the data.
  - Are the values expected for the environment?

- Baseline Check
  - Are there shifts or changes to the log that are unrelated to geology (borehole or probe related?)

- Composite View
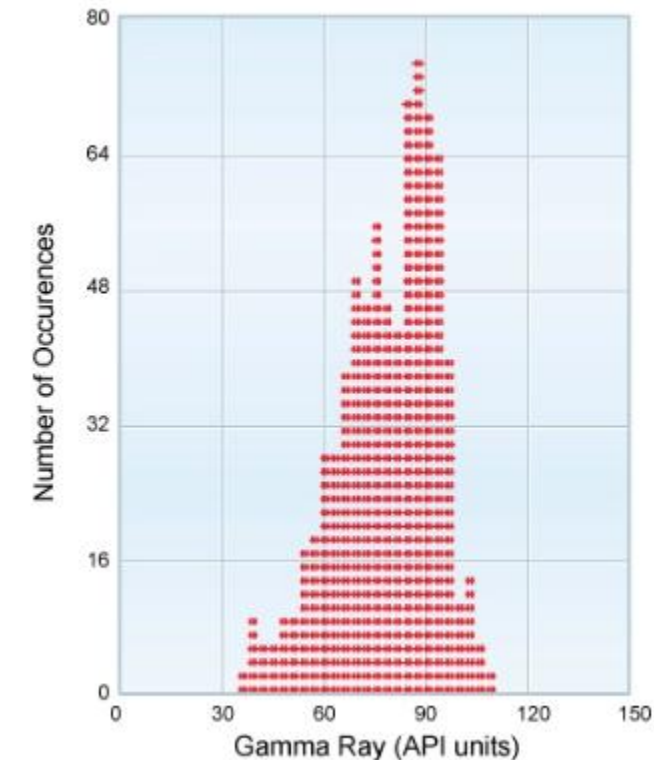  - Viewing all the logs together to get a holistic view of each hole.

- In preparation for advanced methods, ensuring there are no gaps in the data is critical.
  - Interpolation
    - Small gaps where obvious trends are continuous.
  - Imputation
    - Using correlated logs to help fill larger gaps.
  - Manual Correction
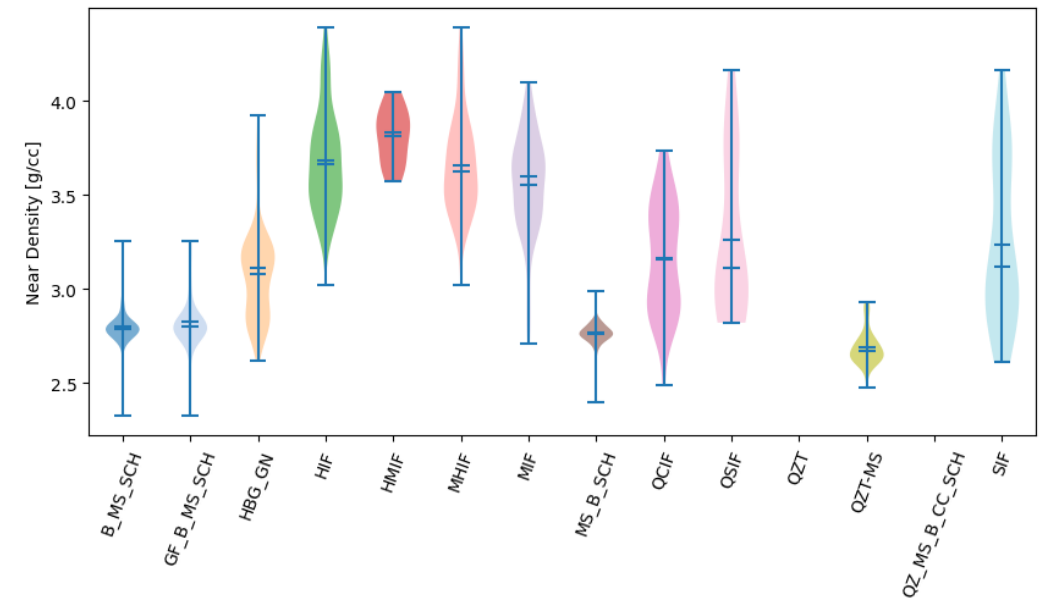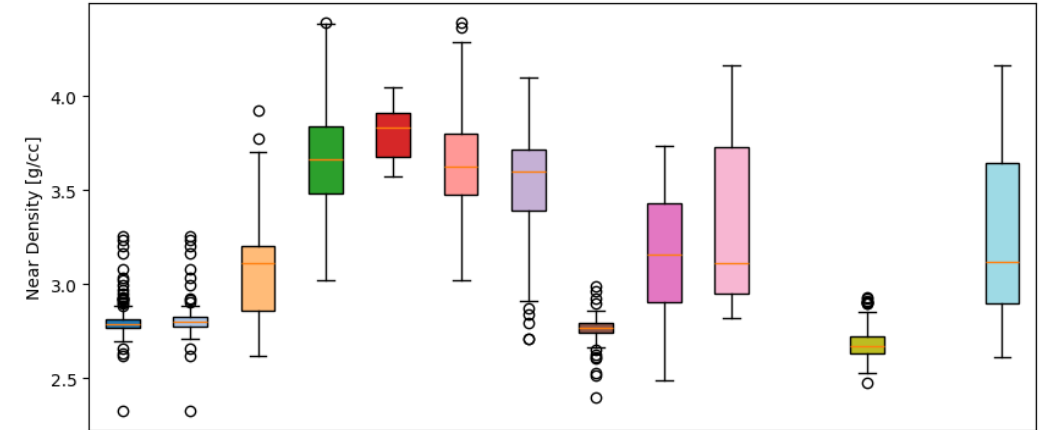    - Manually filling in areas based on observed trends.

## Histograms and Statistics

Analyzing the shape of data distribution helps in normalization and lithology identification.

- **Modality**
  - Bi-Modal or multi-modal distributions often indicate the presence of multiple geologic units or formations.

- **Skewness**
  - Some parameters (mainly magnetic and electrical) will follow a log-normal distribution and require additional transformation.

- **Outliers**
  - Values outside the 10th and 90th percentiles should be examined with care as they may indicate errors in the instruments or small (but significant) observations.

# EDA – Statistics by Geological Unit

- Segment borehole data based on core-logged 'from-to' zones.

- Using visualization such as box or violin plots shows the distribution for each unit.

- Establishes a **characteristic response** for the formation and can be then used in modelling.

## Are there any redundant parameters?

- This is an important consideration before using advanced machine learning techniques as it can introduce bias.

- Correlation Matrix
  - Heat-map showing the Pearson correlation coefficient between parameters.
    - Drop highly co-linear features in improve performance.



Heat map of Pearson's correlation coefficients for well 15/9-F-11

- Linear Dimension Reduction
  - Borehole physical properties are generally multi-parameter (greater than 4) and understanding their relationship can be difficult.
  - PCA reduces the data set into the 2 or 3 components that can be compared with each other.
  - Assumes that all the parameters are **linear**.

- Like PCA but uses a technique called **Manifold Learning** to reduce multi-dimensional data into 2 or 3 components.

- It's a stochastic method, such that each time it's performed it will be different.

- Is a non-linear technique and can often identify trends that are not apparent compared to PCA.

- An alternative interpretation that defines new **domains** based on the natural divisions in the borehole logs.

- Different methods are available such as:
  - K-Means
  - DBSCAN
  - HDSCAN



Well B (Training)

# Physical Property Statistical Analysis Case Study

- 11 holes of physical property data acquired by third party with the following parameters:
  - Magnetic Susceptibility
  - Induction Conductivity
  - Induced Polarization
  - Full waveform sonic
  - Spectral Gamma
- Goal of the project was to examine the data quality, examine the relationships between the parameters and geology, and perform a cluster analysis.

# Data Quality Check – Magnetic Susceptibility

# Data Preparation: Level Chargeability Data

# Data Preparation: Remove Outliers

- Chargeability Capped at 20.
- P Wave Velocity Capped at (3000 & 7300)
- Natural Gamma capped at 450.

# Data Preparation: Log Transform Induction

# Data Preparation: Compositional Scaling of Spectral Gamma

# Dimensional Reduction using TSNE



TSNE Results - HOLEID

# Clustering Tuning Parameters



Epsilon Value



Minimum Samples

# Cluster Analysis: 8 Clusters in the model



TSNE - Clusters results. eps = 4.3; Number Of Clusters = 8

Legend:
- m1_1
- m1_2
- m1_3
- m1_4
- m1_5
- m1_6
- null
- m1_7
- m1_8

# Parallel Coordinate Plot



Model 1: Parallel Coordinate Plot

# Parallel Coordinate Plots of 8 Clusters

# Parallel Coordinate Plots of 8 Clusters



Cluster 2 is one of the smaller groups.
- Highest K response

# Parallel Coordinate Plots of 8 Clusters



Cluster 3 – Low U response.

# Parallel Coordinate Plots of 8 Clusters



Cluster 4 – Largest Group
- From the t-SNE component plot, there appears to be a smaller group that is different.
- Highest Gamma response.

# Parallel Coordinate Plots of 8 Clusters

Cluster 5 – Lowest Th and P-Wave response.

# Parallel Coordinate Plots of 8 Clusters

Cluster 6 – High Th and Low K

# Parallel Coordinate Plots of 8 Clusters

Cluster 7 – Low Gamma, high conductivity and chargeability,

Parallel Coordinate Plots of 8 Clusters

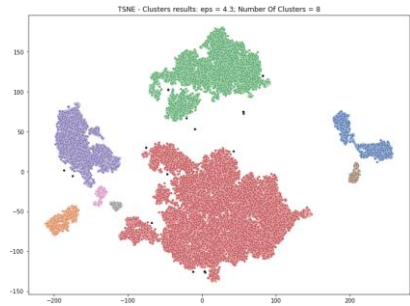Cluster 8 – Similar to cluster 7, but higher P-Wave and Th response.
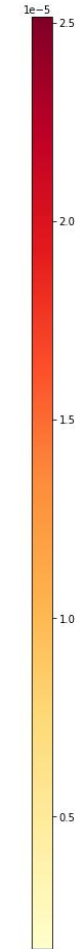
# PRP Parameter Heat Map

# Lithology Heat Map

# Lithology Heat Map

# HOLEID Heat Map

# Remarks

- PRP domains are an alternative interpretation based on the downhole logs.

- Will often match a mix of lithology, alteration, and/or mineralization.

- Same methodology can be applied to other datasets (geochemistry).

- It's an iterative process.
  - Clusters can identify anomalies in the data.
    - How can we explain these anomalies?
      - Geology Related?
      - Borehole Conditions?
      - Probe Related?

# Exercise
# Python Analysis and Visualization
## Click Here to Open