# AI-AUGMENTED WORKFLOW FOR DATA ANALYSTS – PART 1

Workflow Breakdown & Cursor-based Automation Examples

*Prepared by - Seung-Mi Jeon, Hrishikesh Bhatt*

# PROJECT OBJECTIVE

## Goal

1. Decompose the end-to-end workflow of a data analyst

2. Identify repetitive tasks vs. human judgment points

3. Explore how AI tools like Cursor can boost productivity

## Key Deliverables

1. Workflow framework

2. Pain points and decision points (human vs. automation)

3. Example prompts & execution screenshots

(ex: automated code for missing value handling, outlier removal, natural language → SQL transformation)

# WORKFLOW OVERVIEW

Data Loading

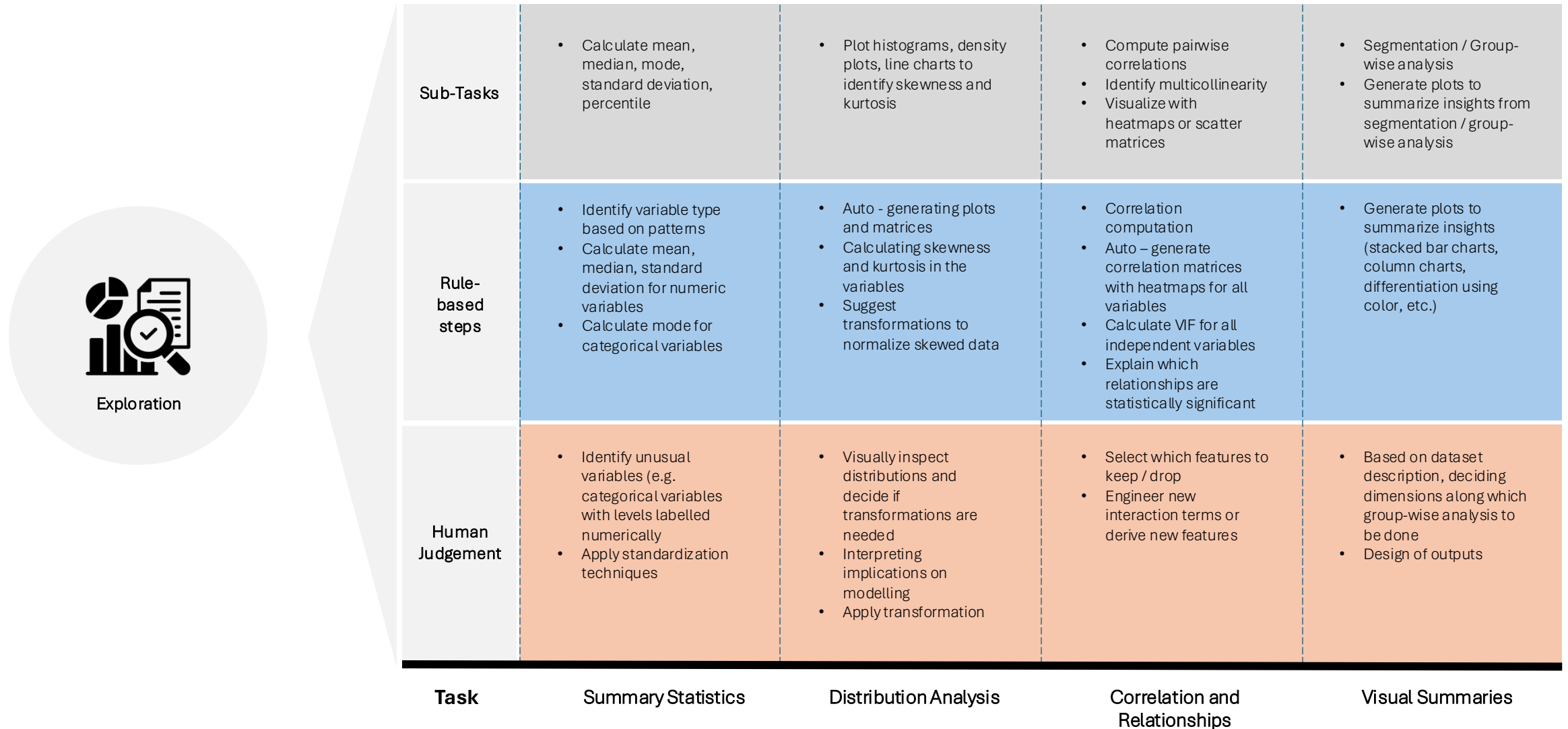Inspection
& Cleaning

Exploration

Modeling

Decision
Support

# DETAILED WORKFLOW BREAKDOWN

| Workflow | Data Loading | Inspection & cleaning | Exploration | Modeling | Decision Support |
|---|---|---|---|---|---|
| **Tasks** | • Structured Data – Querying SQL and NoSQL databases<br>• Structured data – loading csv, xml, json files<br>• Unstructured data – open source and collect from the web through APIs | • Verify dataset size and structure<br>• Duplicate value removal<br>• Handle missing values<br>• Outlier detection and removal<br>• Standardization of formats | • Basic statistical Analysis(Distribution, Unique value analysis, Correlations, VIF)<br>• Feature Engineering (Select features, Create new, Apply transformations) | • Hypothesis Testing<br>• ANOVA<br>• Supervised Learning<br>• Unsupervised Learning<br>• Optimization<br>• A/B Testing, Causal Inference<br>• Hyperparameter Tuning and Model Selection<br>• Interpretability | • Interpreting the data and analytical outputs<br>• Contextualizing insights for business case<br>• Actionable strategies based on insights<br>• Presenting insights to stakeholders |
| **Rule-based steps** | • Syntax and logic for writing SQL queries<br>• Loading data into a DataFrame | • Verifying size and structure<br>• Removing duplicate rows<br>• Standardization of formats (datetime, numbers, etc.) | • EDA - visualize distributions of all key variables<br>• Examine unique values in categorical or ordinal data<br>• Generate heatmaps for all variables<br>• Calculate VIF for all independent variables | • Calculate test statistic and p-value<br>• Initialize and run model<br>• K-elbow visualizer to select number of clusters<br>• Find the optimal points | • Reviewing statistical findings, modelling results – identifying significant trends / anomalies / risks<br>• Generate dashboards |
| **Human Judgement** | • Validate Query outputs<br>• Query optimization given database schema | • Deciding rule(s) for outlier detection and removal given distributions<br>• Deciding rules for handling missing values given counts and distributions | • Specify whether unique value analysis is required<br>• Select which features to keep / drop<br>• Engineer new interaction terms or derive new features<br>• Apply transformations | • Formulate the hypothesis and confidence level<br>• Define independent and dependent variables<br>• Select appropriate model<br>• Specify objective and constraints<br>• Define hyperparameter range<br>• Select the formula of error calculation and criteria<br>• Define the criteria | • Contextualize insights for business case<br>• Present insights to stakeholders - decide format (report, graphs, design, etc.)<br>• Translate insights into plain-language and next steps for business stakeholders |

# METHODOLOGY

1. Under each major Data Analysis Workflow:

   - Develop detailed use cases and process flow charts

   - Outline rule – based and Human Judgement – based steps

2. For each rule – based and Human Judgement – based step:

   - Write prompts for AI augmentation (automate steps or give analyst options where necessary)

   - Create detailed prompt frames for different data and model types

3. Test prompts using Cursor

# EXAMPLE : EXPANDED PROCESS FOR DATA EXPLORATION

Exploration

| | Summary Statistics | Distribution Analysis | Correlation and Relationships | Visual Summaries |
|---|---|---|---|---|
| **Sub-Tasks** | • Calculate mean, median, mode, standard deviation, percentile | • Plot histograms, density plots, line charts to identify skewness and kurtosis | • Compute pairwise correlations<br>• Identify multicollinearity<br>• Visualize with heatmaps or scatter matrices | • Segmentation / Group-wise analysis<br>• Generate plots to summarize insights from segmentation / group-wise analysis |
| **Rule-based steps** | • Identify variable type based on patterns<br>• Calculate mean, median, standard deviation for numeric variables<br>• Calculate mode for categorical variables | • Auto - generating plots and matrices<br>• Calculating skewness and kurtosis in the variables<br>• Suggest transformations to normalize skewed data | • Correlation computation<br>• Auto – generate correlation matrices with heatmaps for all variables<br>• Calculate VIF for all independent variables<br>• Explain which relationships are statistically significant | • Generate plots to summarize insights (stacked bar charts, column charts, differentiation using color, etc.) |
| **Human Judgement** | • Identify unusual variables (e.g. categorical variables with levels labelled numerically<br>• Apply standardization techniques | • Visually inspect distributions and decide if transformations are needed<br>• Interpreting implications on modelling<br>• Apply transformation | • Select which features to keep / drop<br>• Engineer new interaction terms or derive new features | • Based on dataset description, deciding dimensions along which group-wise analysis to be done<br>• Design of outputs |
| **Task** | Summary Statistics | Distribution Analysis | Correlation and Relationships | Visual Summaries |

# CONSTRUCT PROMPT FOR EDA

You are a highly experienced data analyst conducting exploratory data analysis. The dataset is 'Customer_cleaned.csv', which has already been completely preprocessed.

...

Each file you generate should reflect the user's decision. Ensure that each generated file is clearly named and output files from each step are located on the folder named 'eda_step()_output'. Please include descriptive comments above each code block explaining what it does.

...

There are 4 steps for EDA. At each stage, always ask for human input in the chat before proceeding with tasks that require judgment, and generate separate EDA files based on the user's decisions. To use a prompt-chaining method, I will give you prompts for those 4 steps one by one. Wait for the user's natural language response in the chat, and then proceed to generate the corresponding code in a separate EDA analysis file.

Step 1: **Summary Statistics:**

 - Calculate mean, median, mode, standard deviation, and percentiles for numerical variables.

 - *`Rule-based:`* Automatically identify variable types based on patterns and compute these metrics for numeric variables; for categorical variables, compute the mode.

 - *`Human judgement:`* Identify any unusual variables (e.g., categorical variables with numerically labeled levels) and decide if standardization is needed.

## Split EDA Process with 4 steps and distinguish rule-based and human judgement ones with backtick!

# FURTHER STEPS

1. Under each major Data Analysis Workflow:

   - Develop detailed use cases and process flow charts for AI – enabled deployment

   - Outline rule – based and Human Judgement – based steps

2. For each rule – based and Human Judgement – based step:

   - Write prompts for AI augmentation (automate steps or give analyst options where necessary)

   - Create detailed prompt frames for different data and model types

3. Develop metrics to test performance

4. Test prompts using Cursor