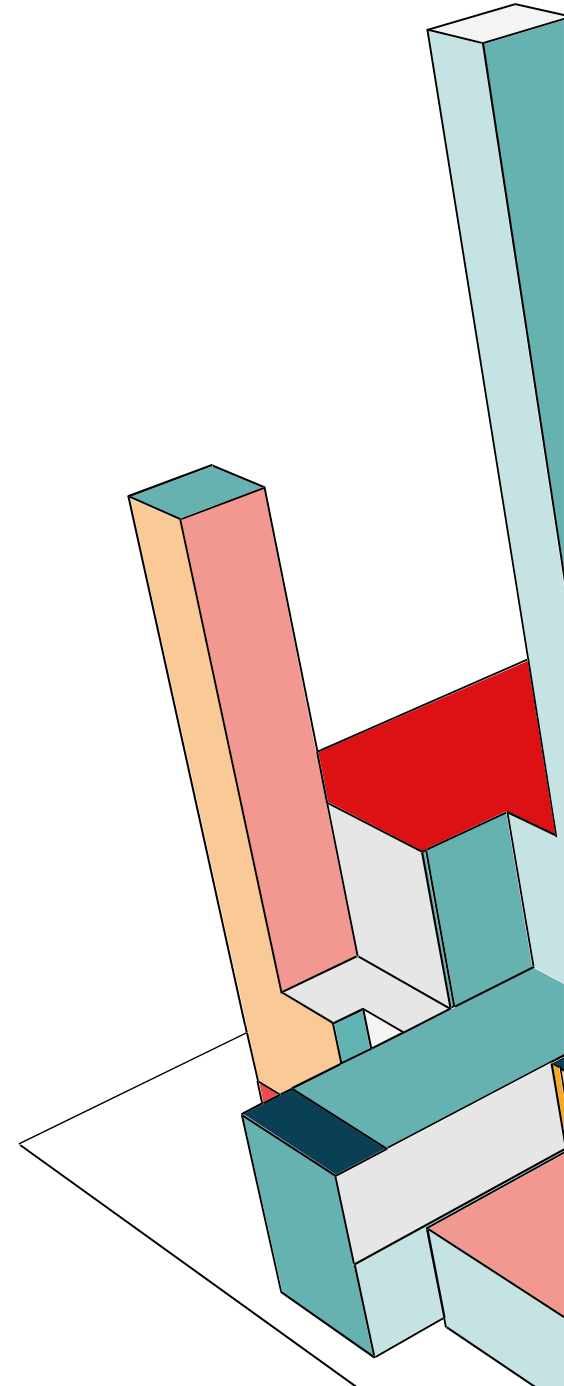


STATISTICAL ANALYSIS AND VISUALIZATION WITH PYTHON

BLOCH CARLO

AGENDA

- Introduction
- Data Generation
- Descriptive statistics
- Data Visualization
- Correlation analysis
- Inferential statistics
- Final considerations



INTRODUCTION

- This project aims to perform statistical analysis and data visualization using Python, generating a synthetic dataset, performing descriptive and inferential statistics, and visualizing the results using Matplotlib and Seaborn.
- Deliverables are a presentation summarizing your findings, visualizations, and interpretations addressing the technical questions provided.





DATA GENERATION

Use Numpy to create a dataset with 1000 samples

The dataset should include the following columns: Age, Height, Weight, Gender, and Income.

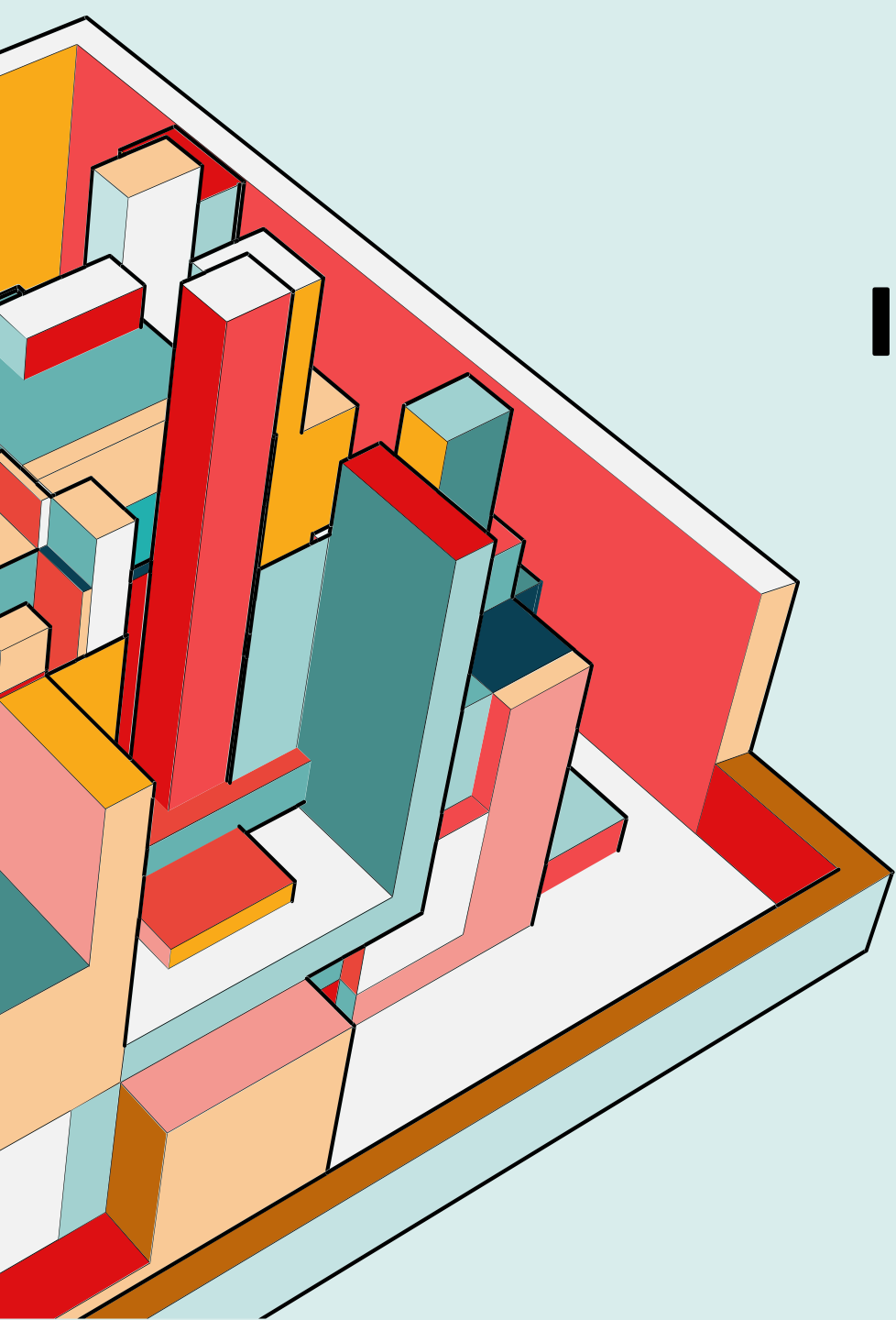
Age: Normally distributed with a mean of 35 and a standard deviation of 10

Height: Normally distributed with a mean of 170 cm and a standard deviation of 15 cm

Weight: Normally distributed with a mean of 70 kg and a standard deviation of 10 kg

Gender: Randomly assigned with 50% probability for 'Male' and 'Female'

Income: Normally distributed with a mean of 50,000 and a standard deviation of 15,000



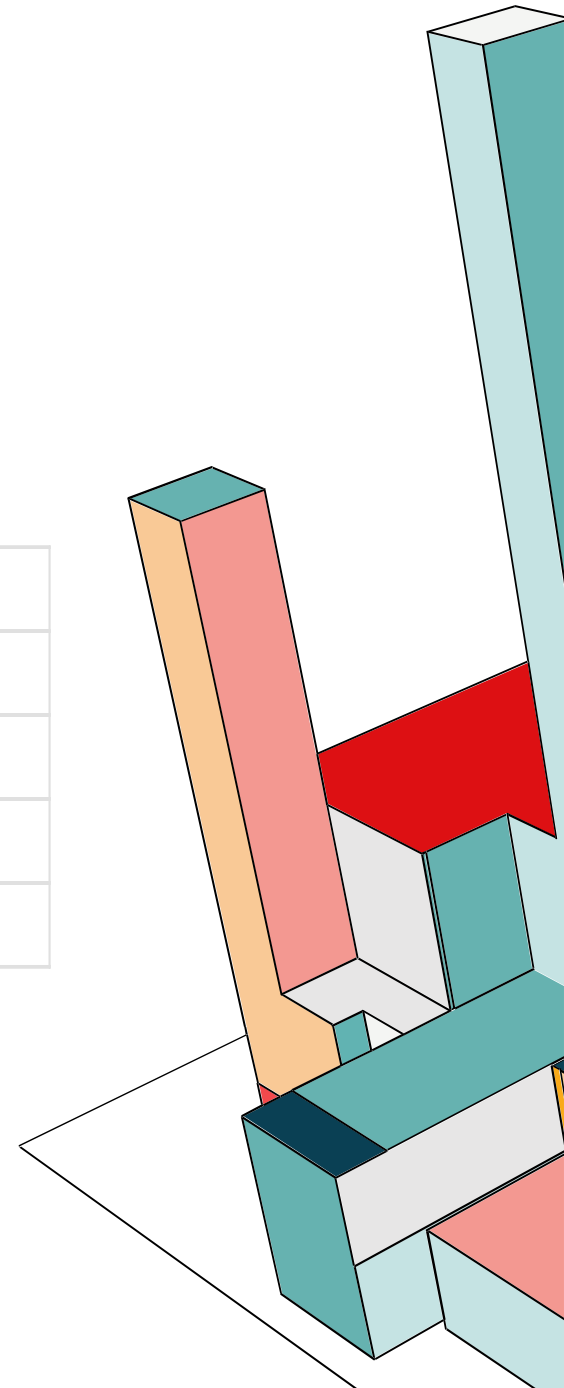
INFERENCE STATISTICS

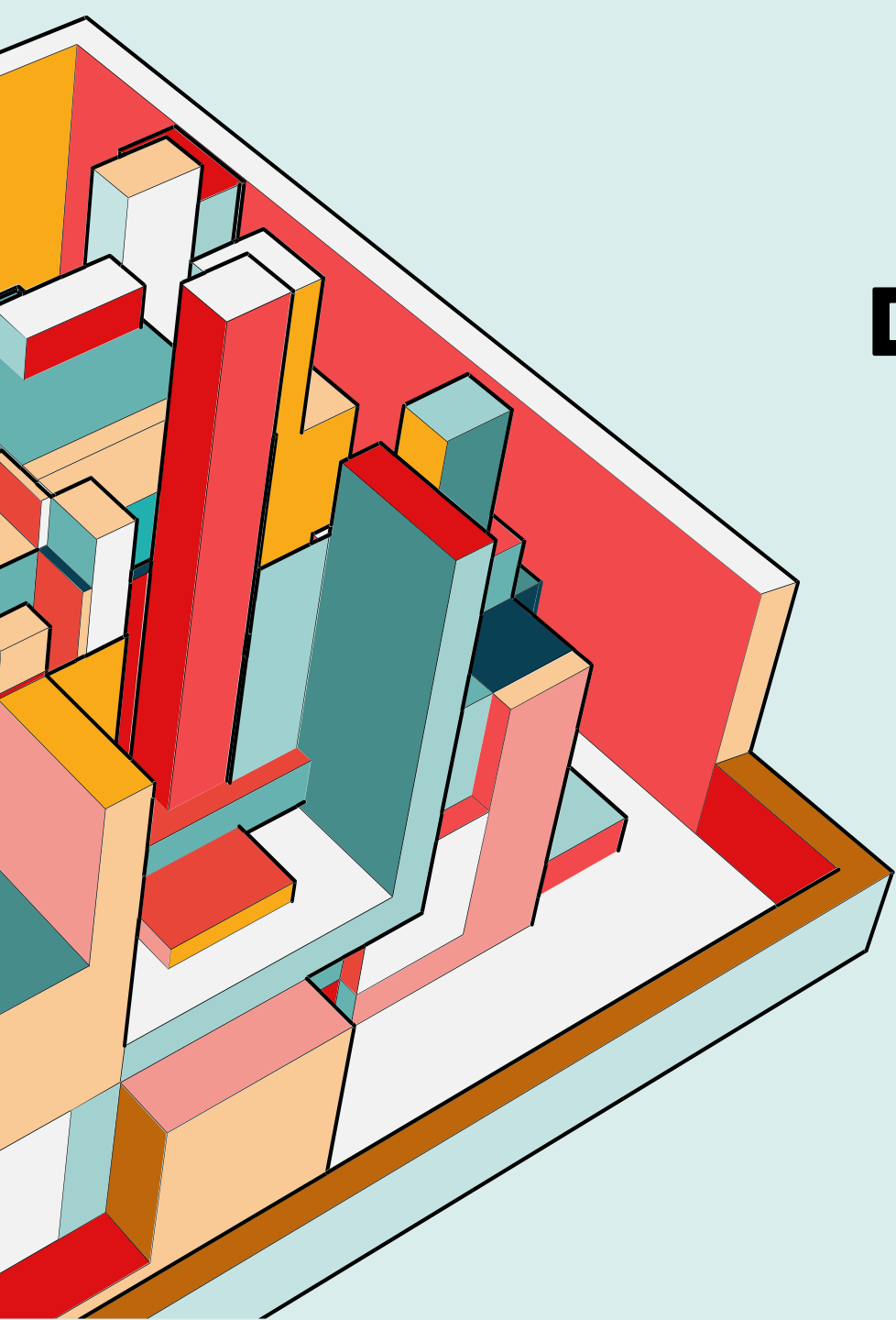
DESCRIPTIVE STATISTICS

Basic descriptive statistics

- Calculate mean, median, standard deviation, and variance for Age, Height, Weight, and Income
- Calculate the mode for Gender.

for column	AGE	Mean=34.46,	Median=34.5,	Standard deviation=9.8,	Variance =96.0	
for column	HEIGHT	Mean=1.7,	Median=1.69,	Standard deviation=0.14,	Variance =0.02	
for column	WEIGHT	Mean=70.17,	Median=69.72,	Standard deviation=10.19,	Variance =103.91	
for column	INCOME	Mean=49583.06,	Median=49065.31,	Standard deviation=15240.22,	Variance =232264164.8	
for column	GENDER	Mode=Male				

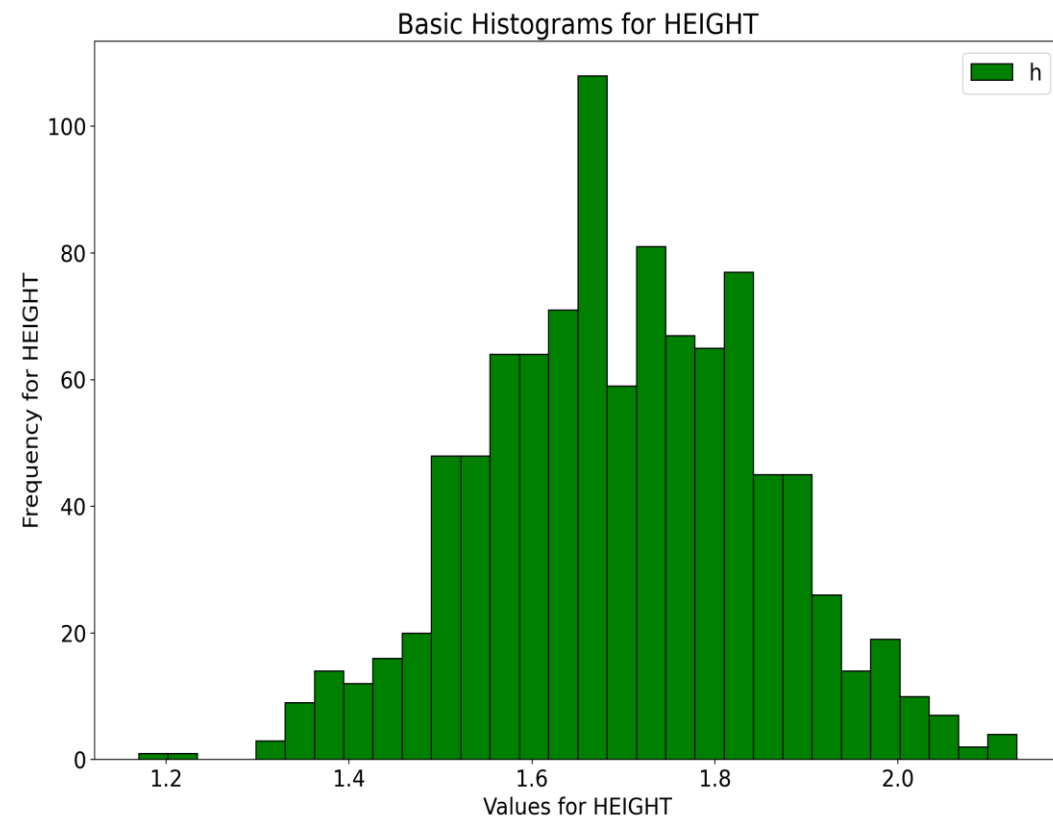
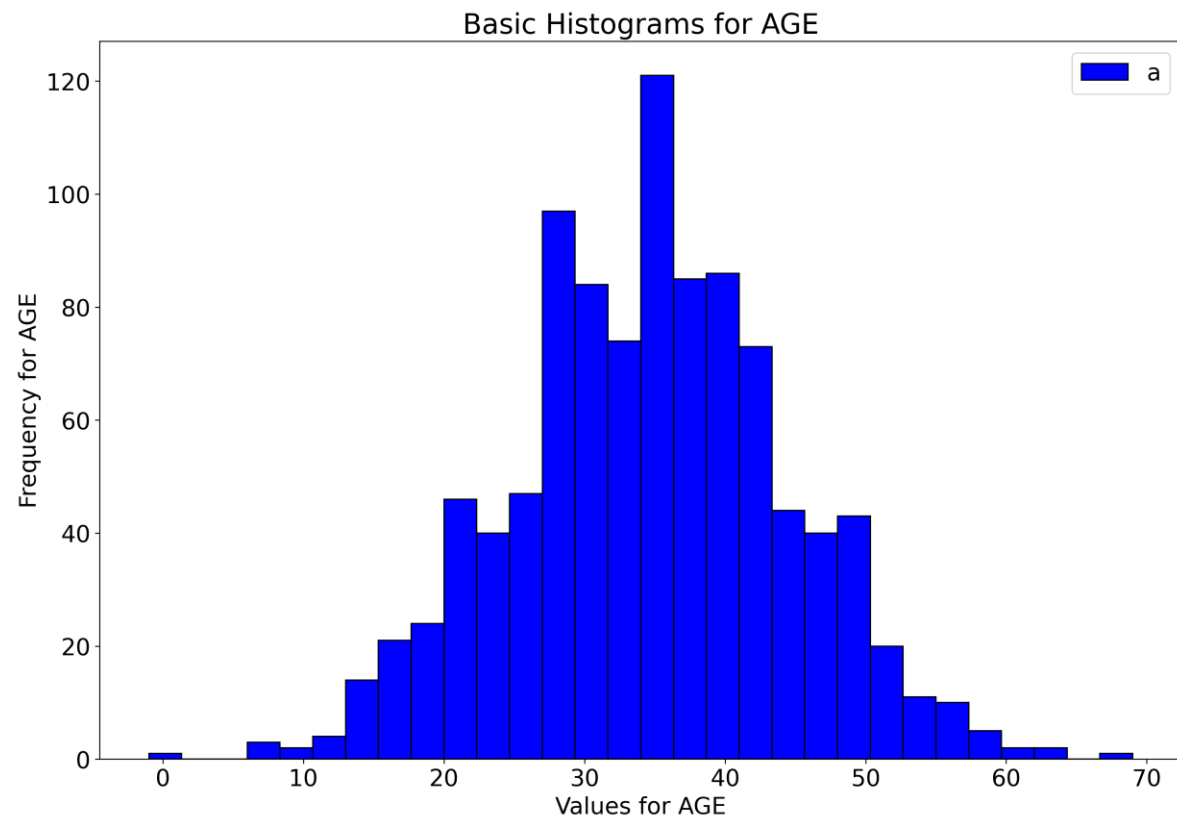




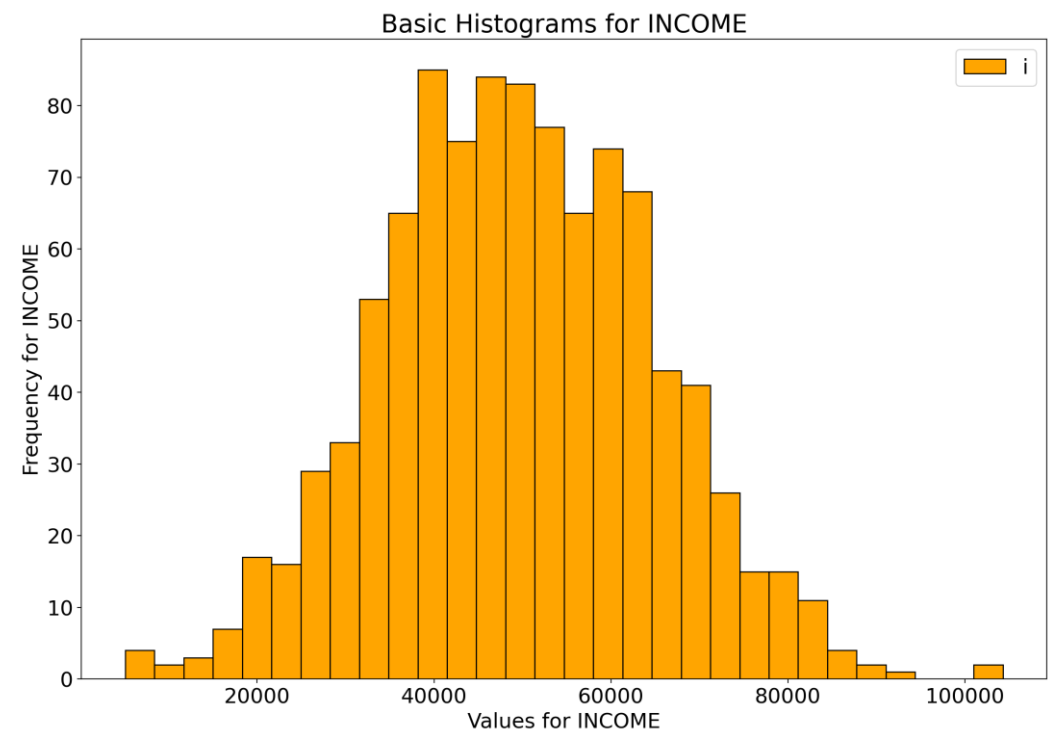
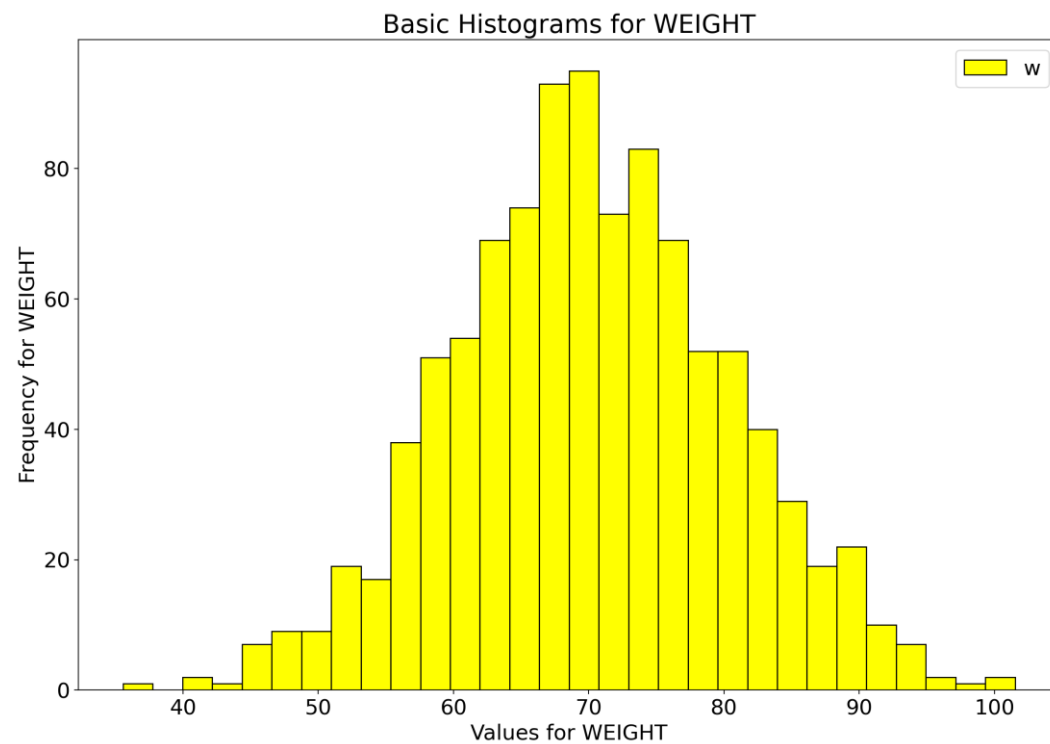
DATA VISUALIZATION

Results following

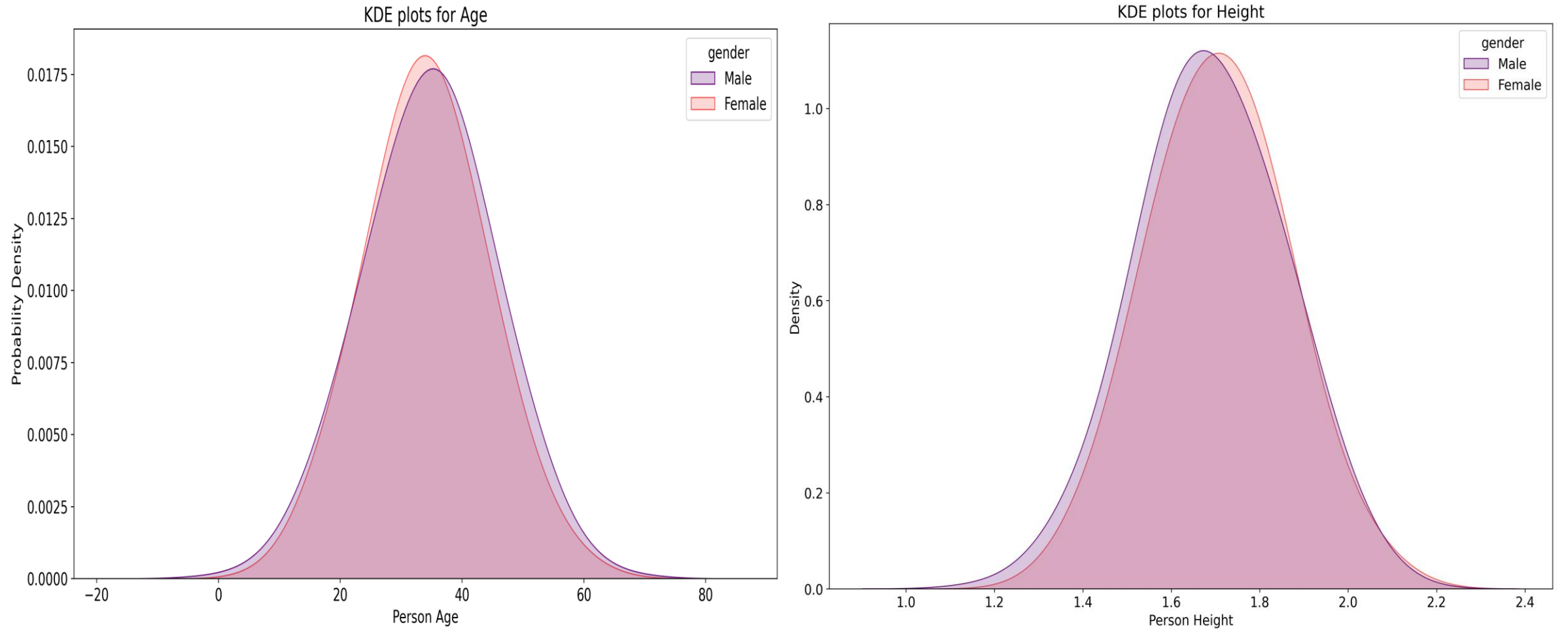
AGE, HEIGHT HISTOGRAMS



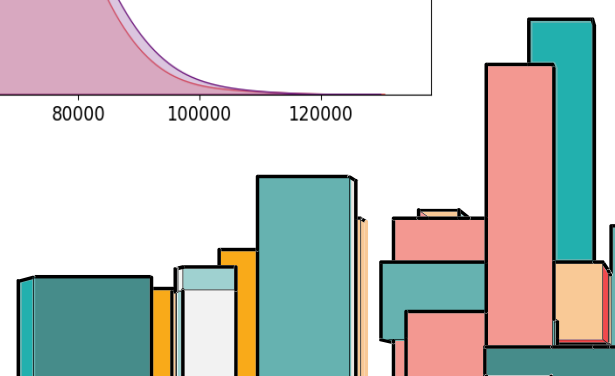
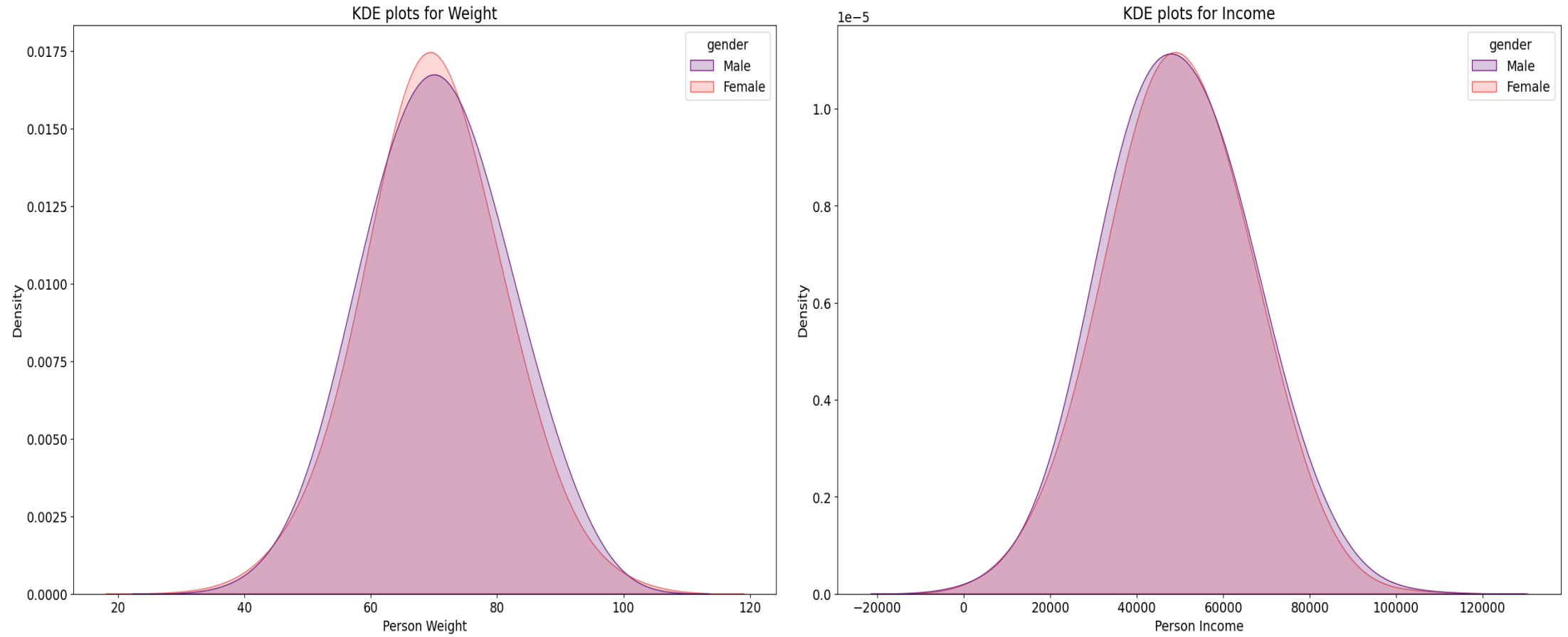
WEIGHT, INCOME HISTOGRAMS



KDE PLOTS

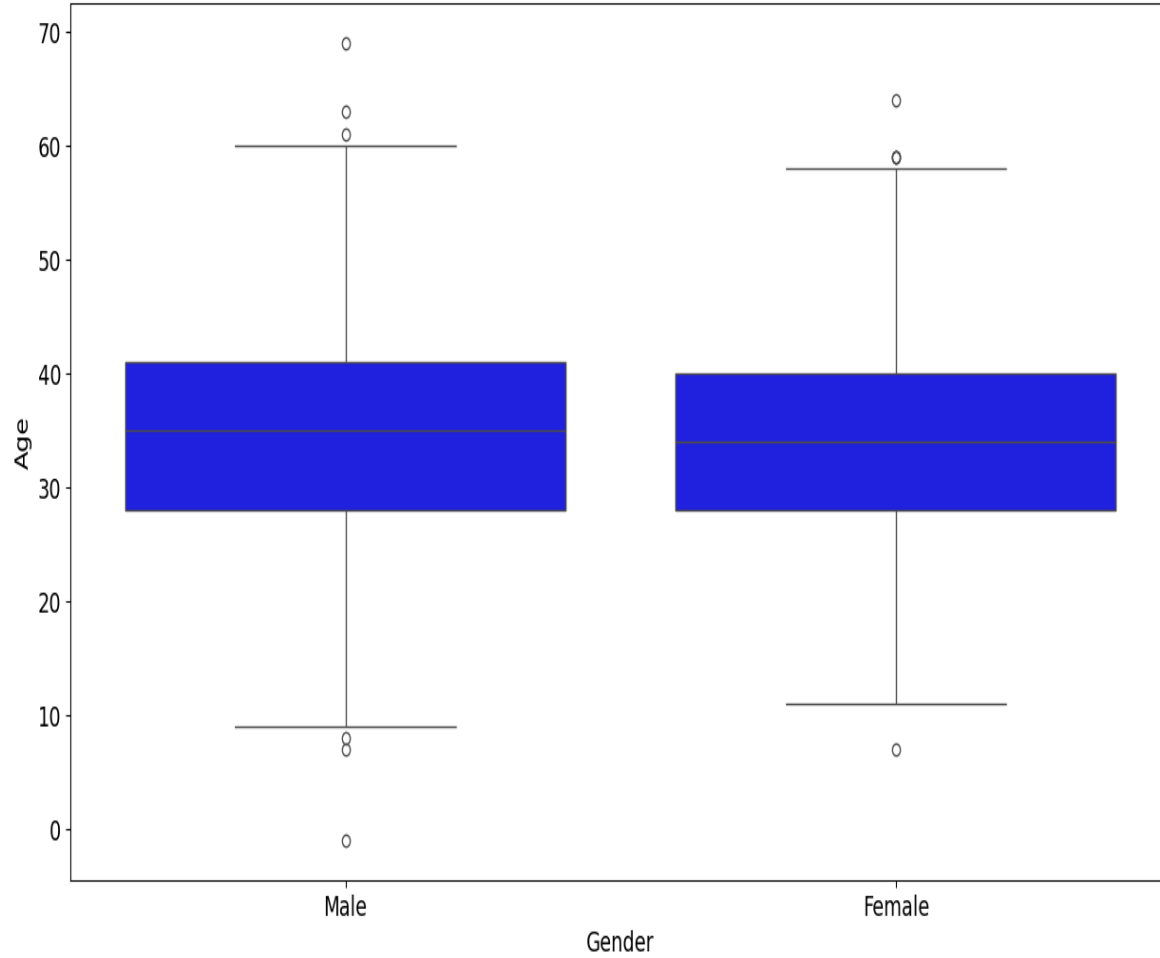


KDE PLOTS

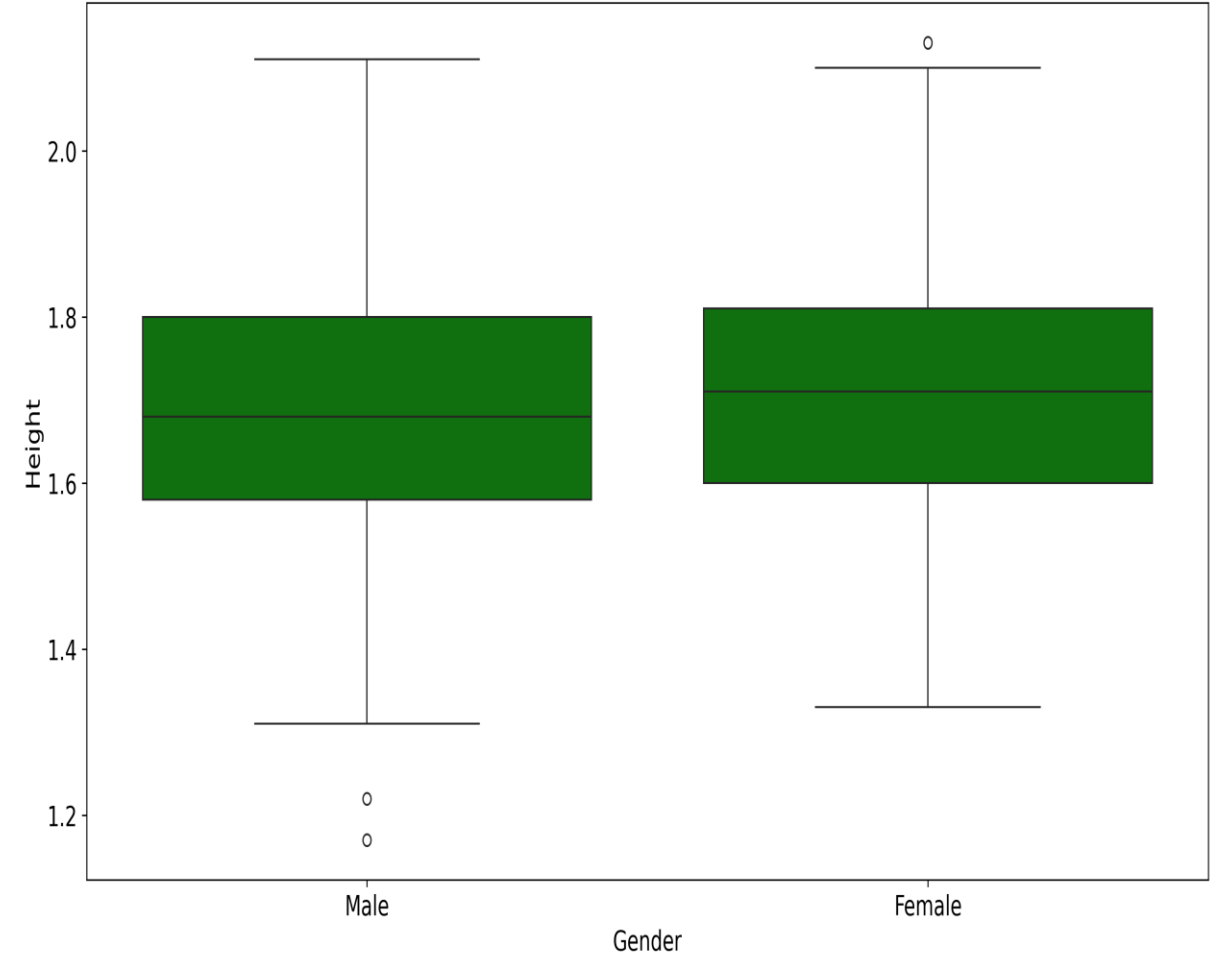


BAR PLOTS – HIGHLIGHTING OUTLIERS

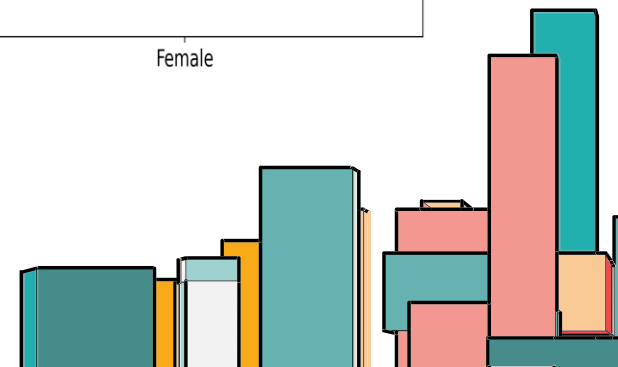
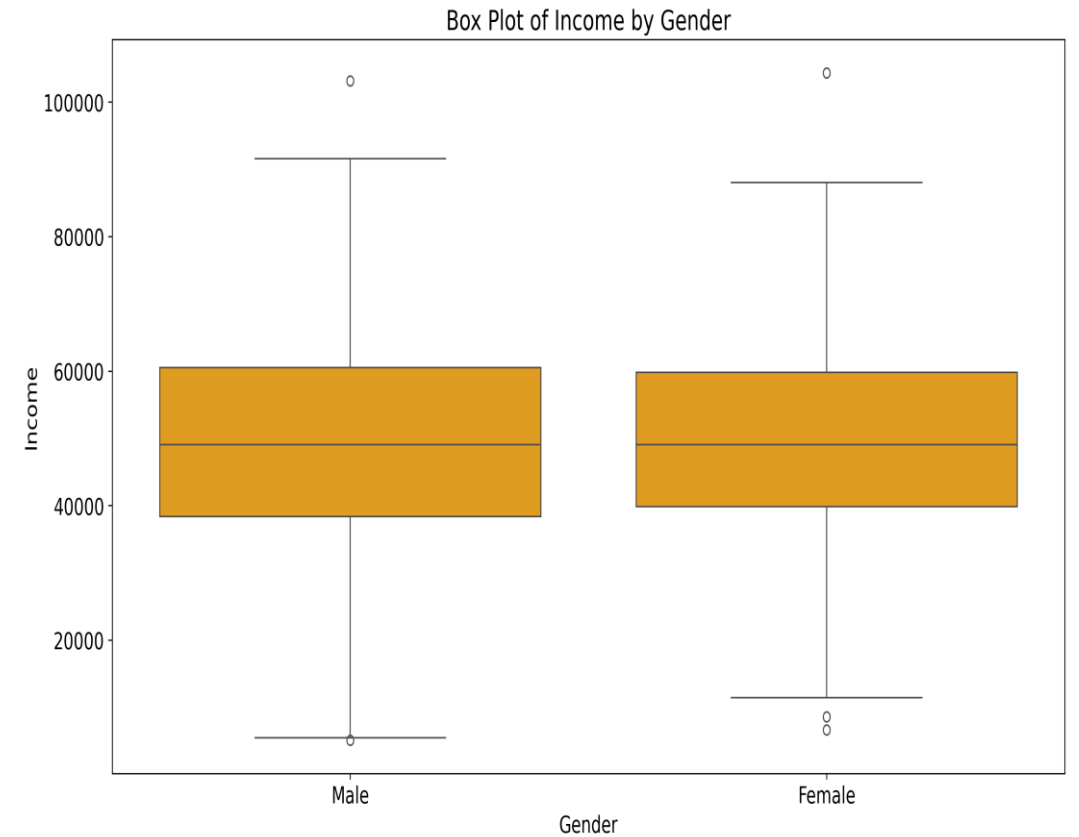
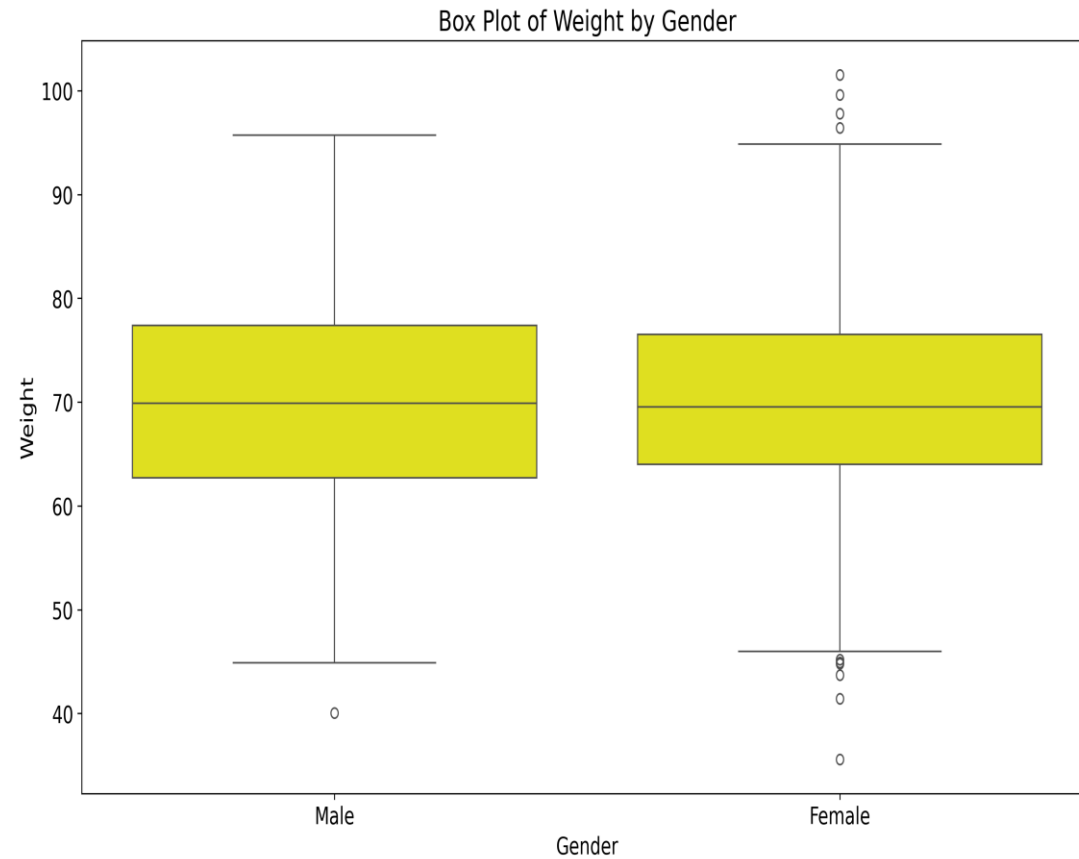
Box Plot of Age by Gender

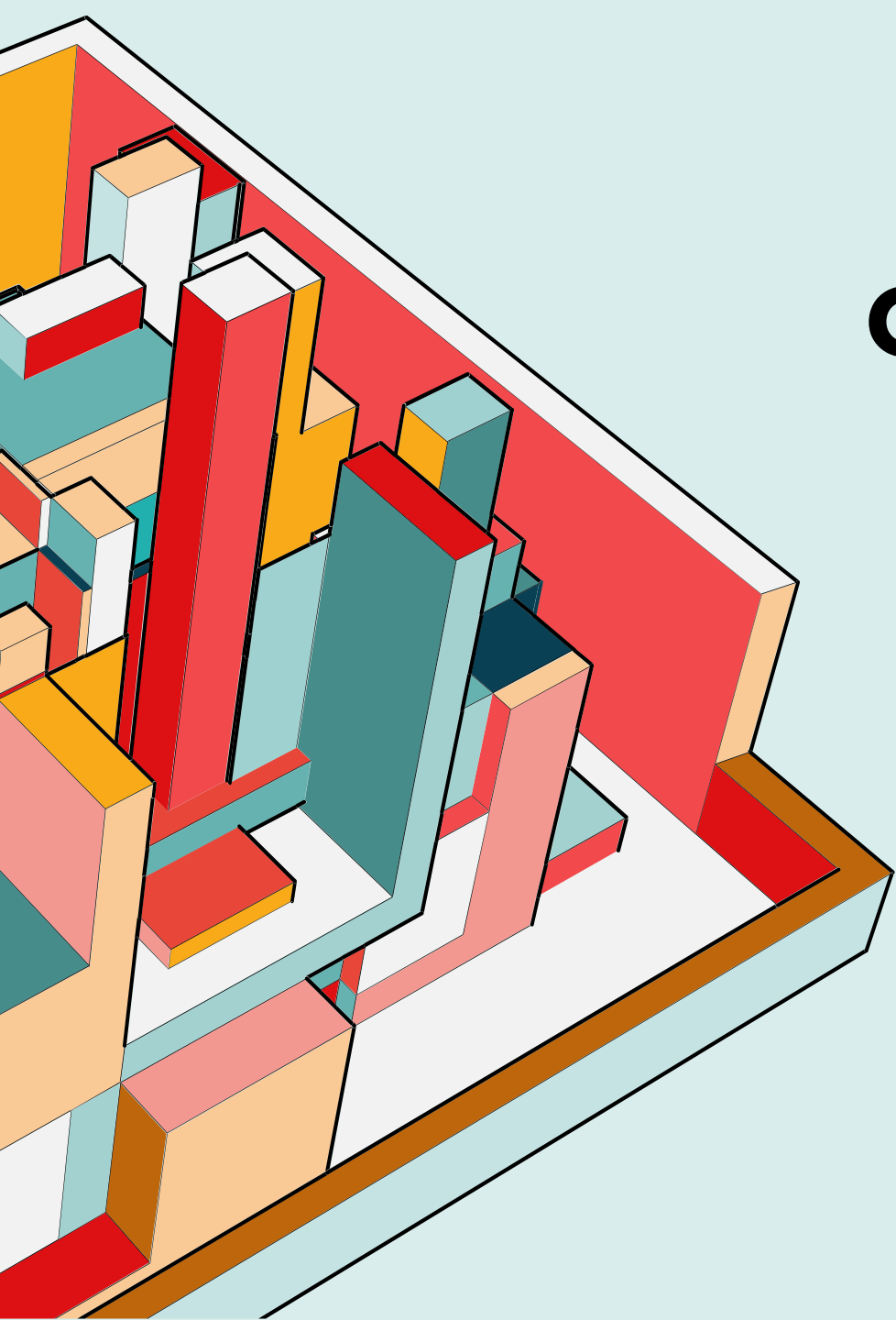


Box Plot of Height by Gender



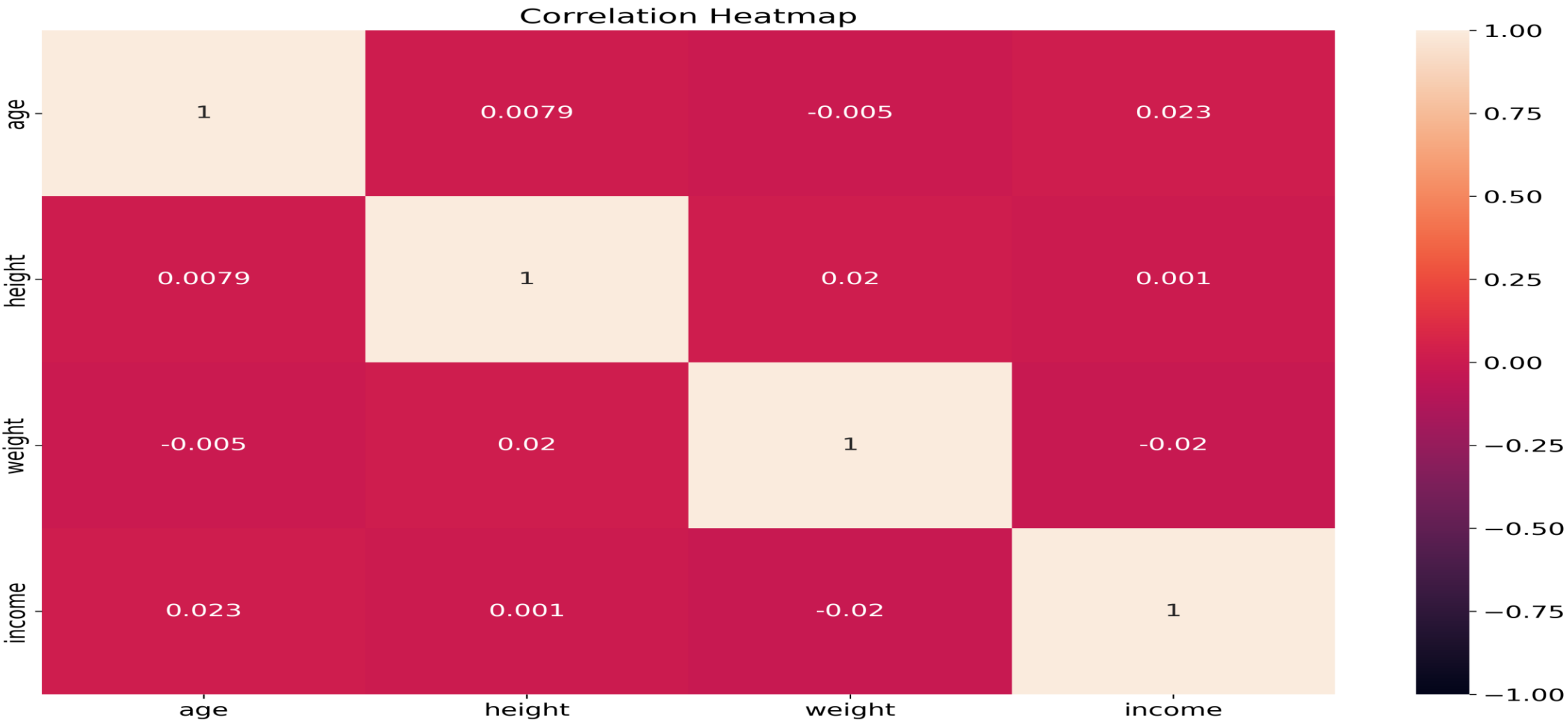
BAR PLOTS – HIGHLIGHTING OUTLIERS

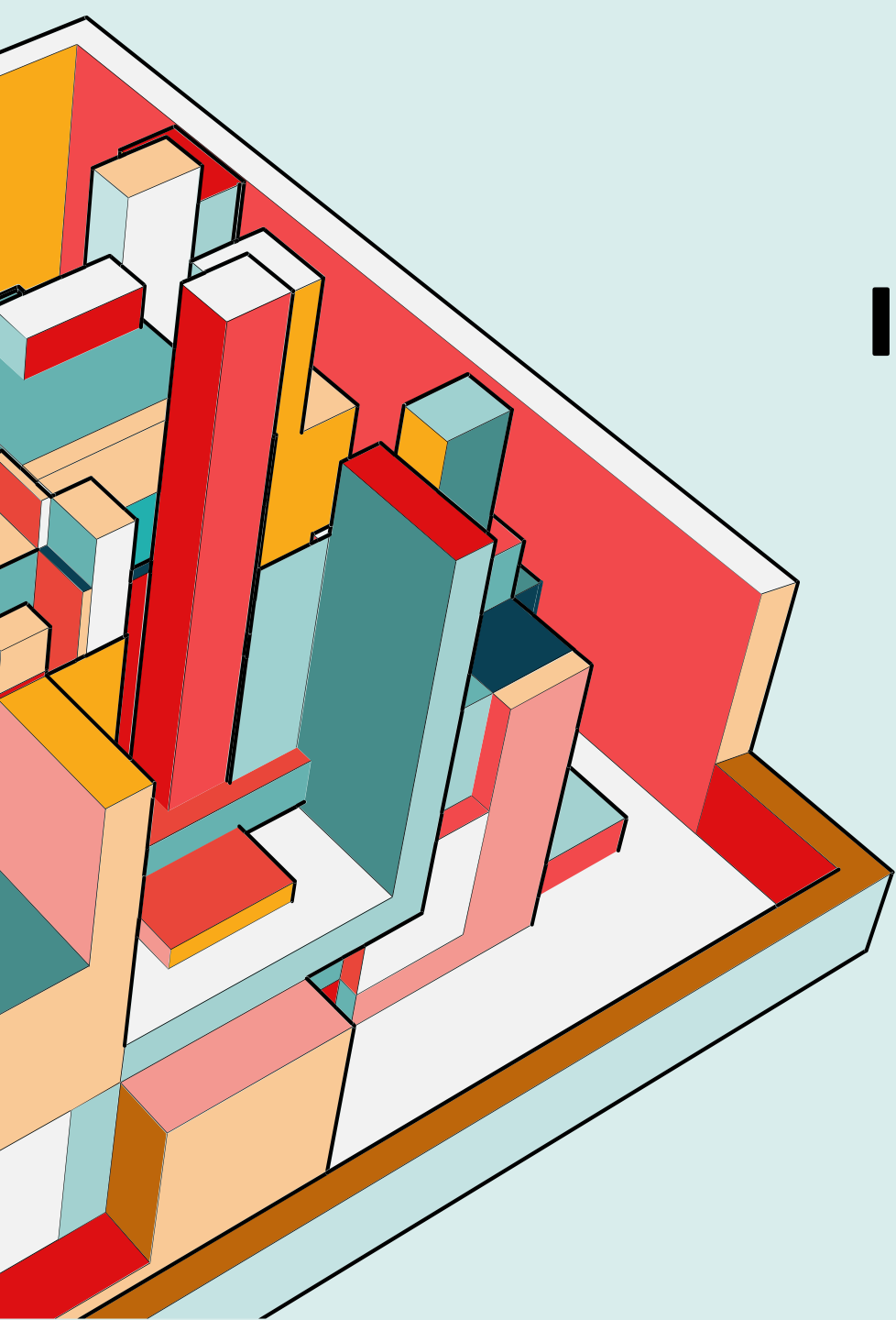




CORRELATION ANALYSIS

PEARSON-CORRELATION MATRIX



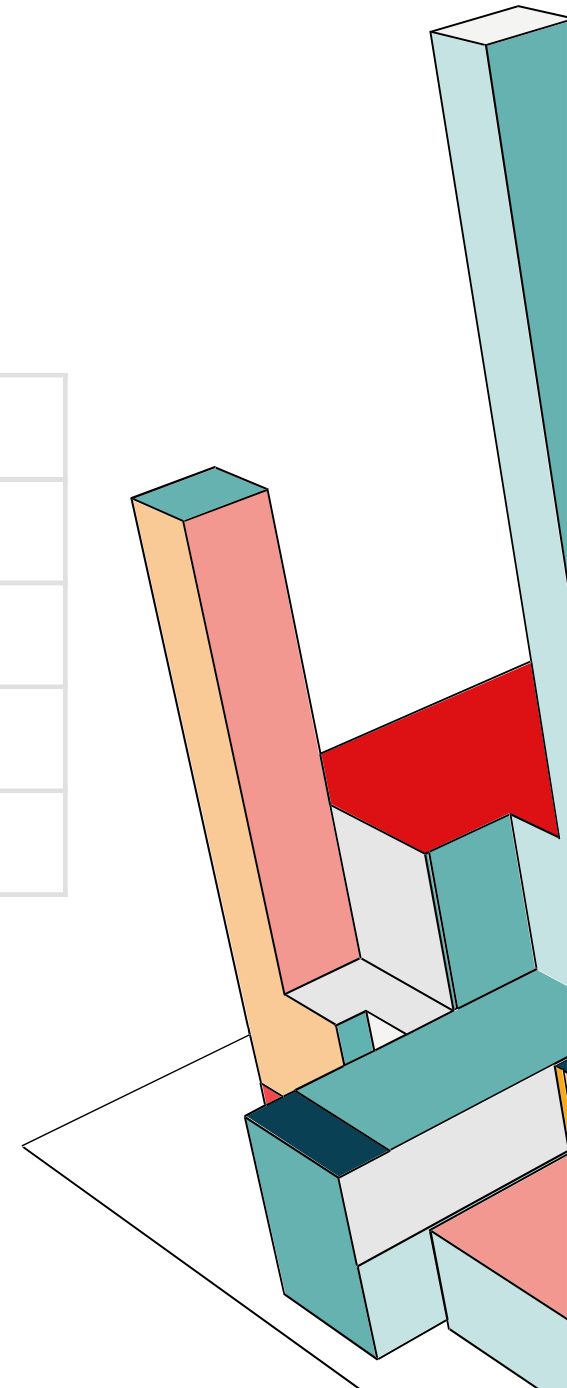


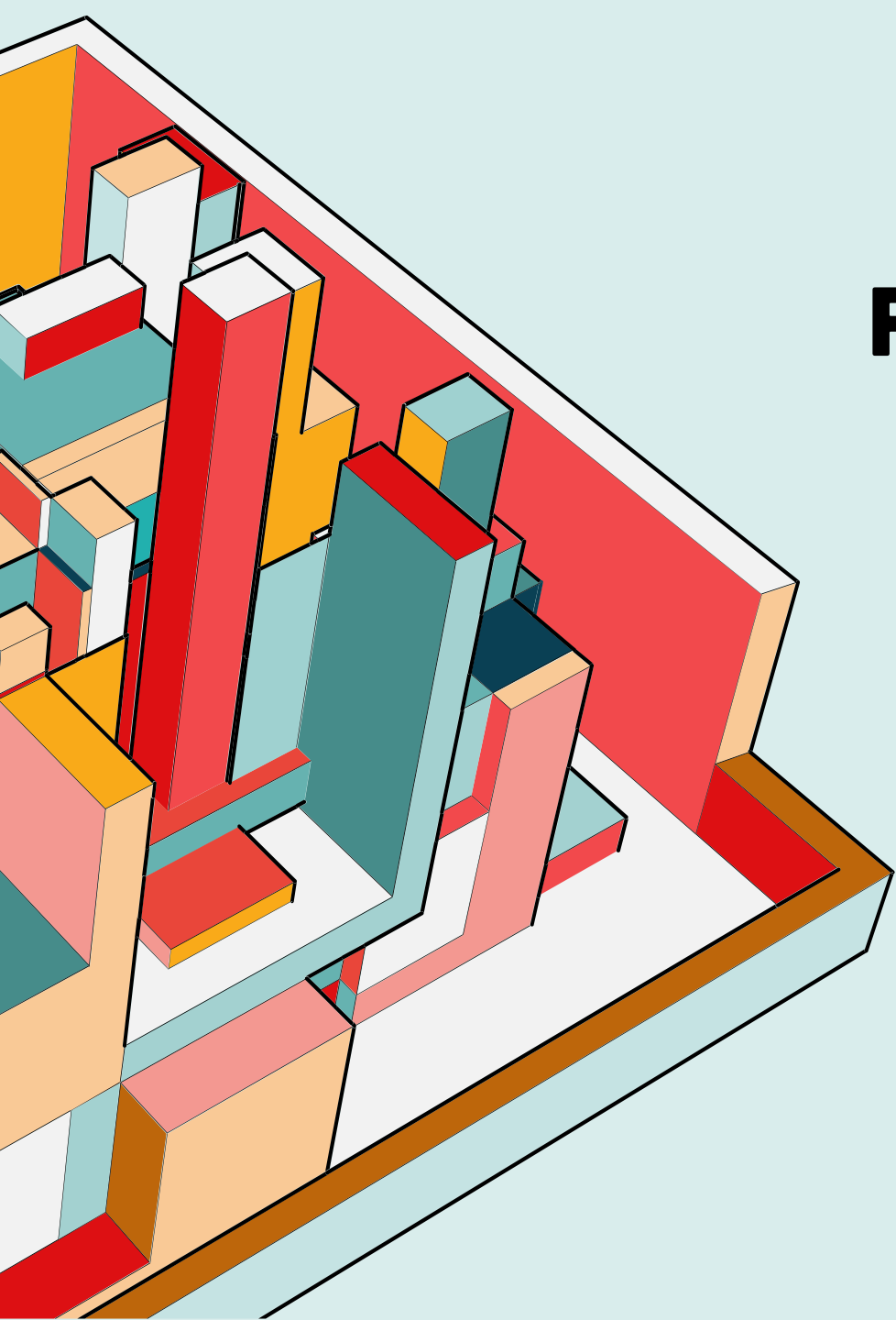
INFERENCE STATISTICS

INFERENCEAL STATISTICS

T-test to see if there is a significant difference in Income between Male and Female

income differences between Female vs Male group)	
income: t-statistics = 0.0141, p_value = 0.9888	
income differences between Male vs Female group)	
income: t-statistics = -0.0141, p_value = 0.9888	





FINAL CONSIDERATIONS

FINAL CONSIDERATIONS

Final questions	Reason	Objective	Test done?
Normal distributions for the variables reason	Having Dataset normalized	Using Numpy to generate random datasets	YES
Descriptive statistics insights	Verify Datasets Requirements	Verify if computed values, for mean, variance, ecc. are correct	YES
Boxplots interpretation, about outliers	N. of outliers very small	Result Expected within gaussain distributions	YES
Correlation Matrix	correlations between data columns	No correlation found	YES
T-test results	The p-value shows the likelihood of your data occurring under the null hypothesis	P-value almost equal to 1, i.e. no income differenceses	YES

THANK YOU

Carlo Bloch

