

Assignment 1: Data Preprocessing and Visualization

Rana Muhammad Sahil Khan

July 2024

Contents

1	Problem 1	2
1.1	2
1.2	3
1.3	4
1.4	4
1.5	4
2	Problem 2	4
2.1	4
2.2	6
3	Problem 3:	7
3.1	7
3.2	7
3.3	9
3.4	14
3.5	14
3.6	19
3.7	23
3.8	28
3.9	29
3.10	30
3.11	35
3.12	42
3.13	47
3.14	48
3.15	49
3.16	49
3.17	50

1 Problem 1

1.1

Minimum: 2.2

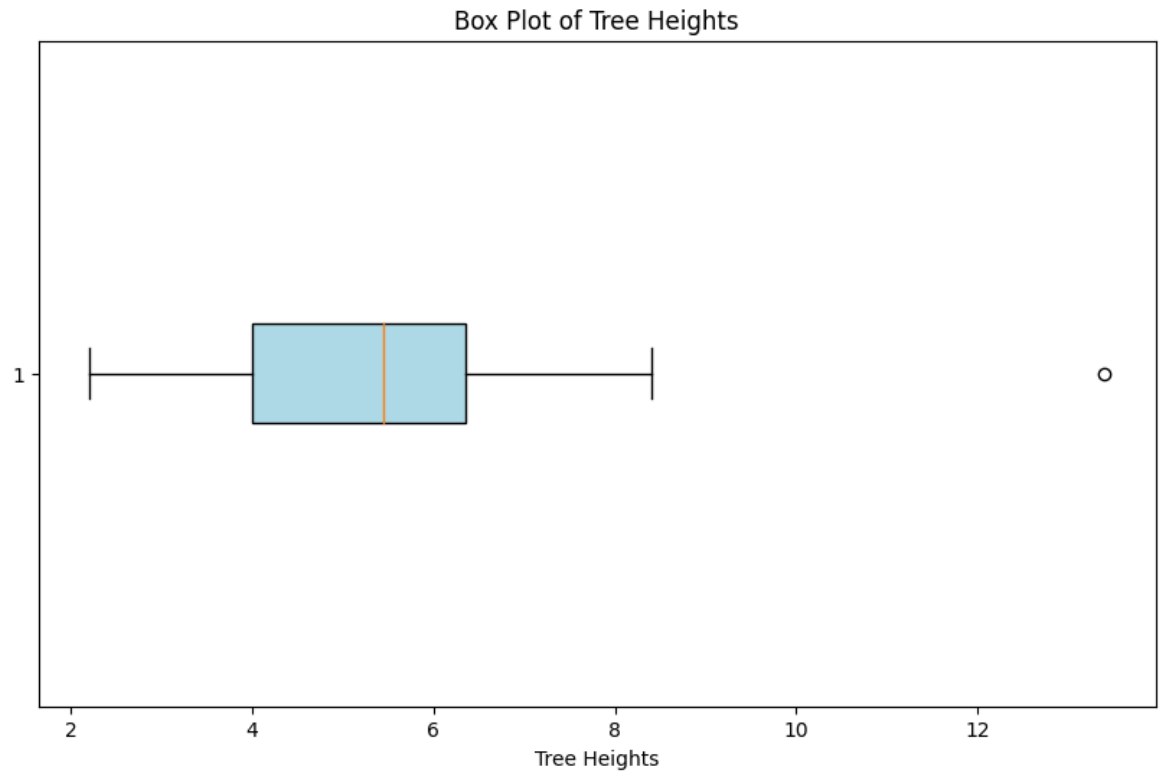
Q1 (25th percentile): 4.0

Median (50th percentile): 5.45

Q3 (75th percentile): 6.35

Maximum: 13.4

1.2



1.3

ordered data: 2.2, 2.5, 2.5, 2.7, 4.0, 4.0, 4.2, 4.3, 4.8, 5.4, 5.5, 5.5, 5.5, 5.8, 6.3, 6.5, 7.9, 7.9, 8.4, 13.4

Equal Frequency Partitioning:

bin 1: [2.2, 2.5, 2.5, 2.7, 4.0]

bin 2: [4.0, 4.2, 4.3, 4.8, 5.4]

bin 3: [5.5, 5.5, 5.5, 5.8, 6.3]

bin 4: [6.5, 7.9, 7.9, 8.4, 13.4]

Equal Width Partitioning:

bin 1: [2.2, 2.5, 2.5, 2.7, 4.0, 4.0, 4.2, 4.3, 4.8]

bin 2: [5.4, 5.5, 5.5, 5.5, 5.8, 6.3, 6.5]

bin 3: [7.9, 7.9, 8.4]

bin 4: [13.4]

1.4

The outliers can be inferred from the box plot, any value lying outside $3Q + 1.5 \times (IQR)$ is considered an outlier which for this specific data set is just one value: 13.4.

1.5

The missing values in the following data can be dealt with in two ways

4.3, 2.5, 2.5, ?, ?, 5.4, 6.3, ?, 4.0, 4.0, ?, ?, 13.4, 8.4, ?, 5.5, 5.5, 5.5, 2.7, 6.5

Since the missing values are missing completely at random, the following methods can be used to deal with them:

Imputation: The missing values can be replaced(imputed) with the mean(average), median(middle value) or mode(most repeated value) of the remaining data and then the completed dataset can be used for data analysis but this could lead to some loss of variation.

Deletion: The missing values can be deleted and not considered for data analysis and since they are missing completely at random, the resulting data analysis will be unbiased.

2 Problem 2

2.1

Answers:

Mean: 58.8

Median: 60.0

Mode: ModeResult(mode=45.0, count=2)

5% Trimmed Mean: 59.7

10% Trimmed Mean: 60.5

(Python) Code for 2.1

```

import numpy as np
from scipy import stats

exam_data = [56.7, 45, 65, 77, 75, 78, 56, 73, 75, 24, 0, 10, 100,
95, 45, 52, 65, 66, 58, 58.5, 59, 61]

mean = round(np.mean(exam_data), 1)
median = np.median(exam_data)
mode = stats.mode(exam_data)
trimmedfive = round(stats.trim_mean(exam_data, 0.05), 1)
trimmedten = round(stats.trim_mean(exam_data, 0.1), 1)

print(f"Mean: {mean}")
print(f"Median: {median}")
print(f"Mode: {mode}")
print(f"5% Trimmed Mean: {trimmedfive}")
print(f"10% Trimmed Mean: {trimmedten}")

```

Measure of centrality:

I think that median is the better measure of centrality because it is the least effected by outliers in the data.

Robust:

I think that median is more robust than others because it does not take into account the outliers and is least effected by outliers in all scenarios.

2.2**Values:**

Mean: 2.7

Median: 2.5

Mode: 1.2

5% Trimmed Mean: 2.6

10% Trimmed Mean: 2.6

(Python) Code for 2.2

```
import numpy as np
from scipy import stats

random_numbers = np.random.uniform(low=1, high=5, size=20)

mean = round(np.mean(random_numbers), 1)
median = round(np.median(random_numbers), 1)
mode_result = stats.mode(random_numbers)
mode_value = mode_result.mode
mode= round(mode_value, 1)

trimmedfive = round(stats.trim_mean(random_numbers, 0.05), 1)
trimmedten = round(stats.trim_mean(random_numbers, 0.1), 1)

print(f"Mean: {mean}")
print(f"Median: {median}")
print(f"Mode: {mode}")
print(f"5% Trimmed Mean: {trimmedfive}")
print(f"10% Trimmed Mean: {trimmedten}")
```

Comments:

The mean and median should be close to 3 as this is a uniformed set of numbers; the mode of this dataset is not much meaningful due to the random nature of it. As expected, the trimmed means are very close to the mean because of the absence of much outliers.

3 Problem 3:

3.1

Numerical Attributes: age, fnlwgt, education-num, capital-gain, capital-loss, hours-per-week

Nominal Attributes: race, sex, native-country

Categorical Attributes: workclass, education, marital-status, occupation, relationship, race, sex, native country

Unique values per attribute:

age: 73
workclass: 9
fnlwgt: 21648
education: 16
education-num: 16
marital-status: 7
occupation: 15
relationship: 6
race: 5
sex: 2
capital-gain: 119
capital-loss: 92
hours-per-week: 94
native-country: 42
income: 2

3.2

Age:

Minimum: 17
Q1 (25th percentile): 28.0
Median (50th percentile): 37.0
Q3 (75th percentile): 48.0
Maximum: 90
Mode: 36
Mean: 38.6

fnlwgt(Final Weight):

Minimum: 12285
Q1 (25th percentile): 117827.0
Median (50th percentile): 178356.0
Q3 (75th percentile): 237051.0
Maximum: 1484705
Mode: 123011

Mean: 189778.4

education-num:

Minimum: 1
Q1 (25th percentile): 9.0
Median (50th percentile): 10.0
Q3 (75th percentile): 12.0
Maximum: 16
Mode: 9
Mean: 10.1

capital-gain:

Minimum: 0
Q1 (25th percentile): 0.0
Median (50th percentile): 0.0
Q3 (75th percentile): 0.0
Maximum: 99999
Mode: 0
Mean: 1077.6

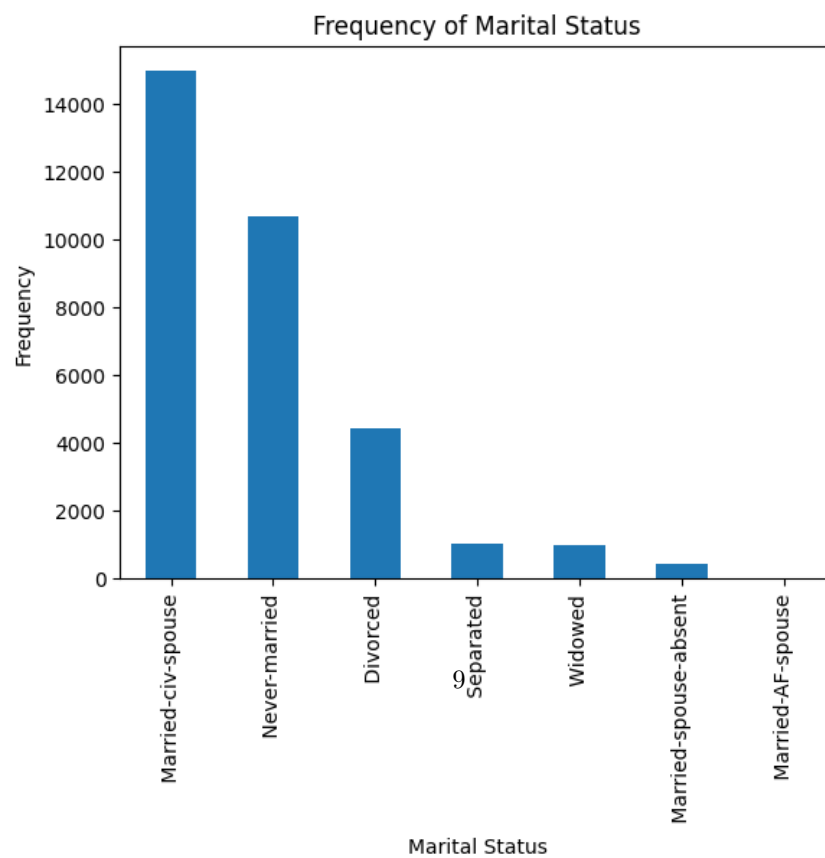
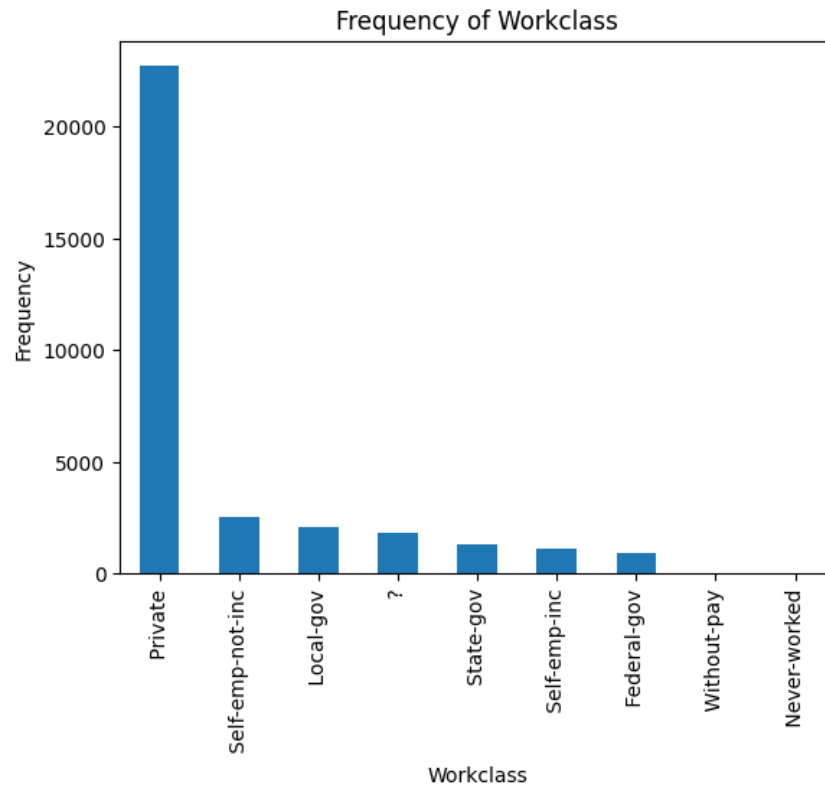
capital-loss:

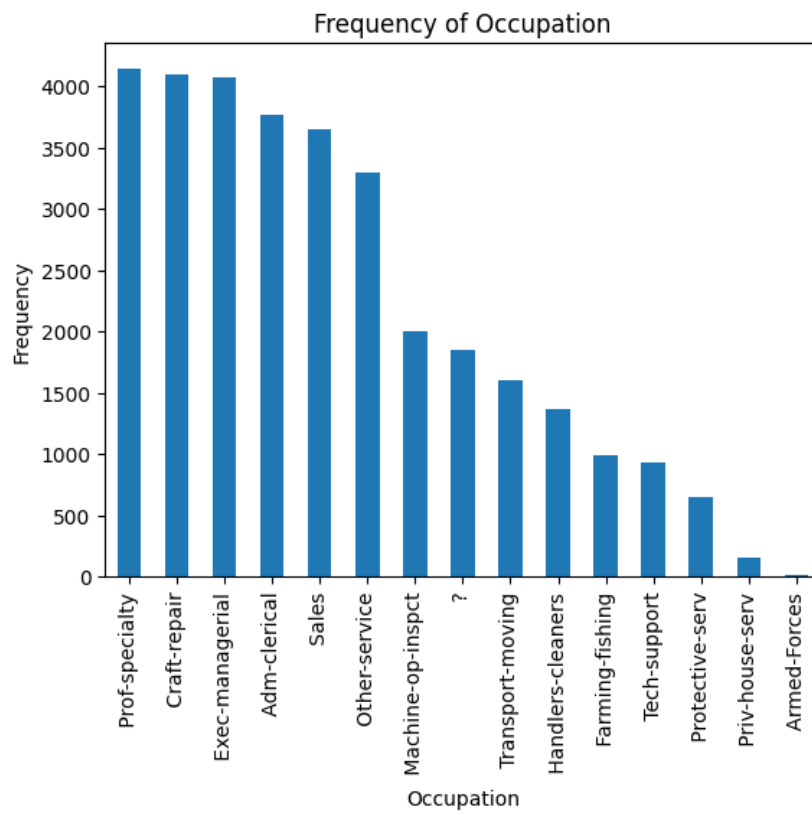
Minimum: 0
Q1 (25th percentile): 0.0
Median (50th percentile): 0.0
Q3 (75th percentile): 0.0
Maximum: 4356
Mode: 0
Mean: 87.3

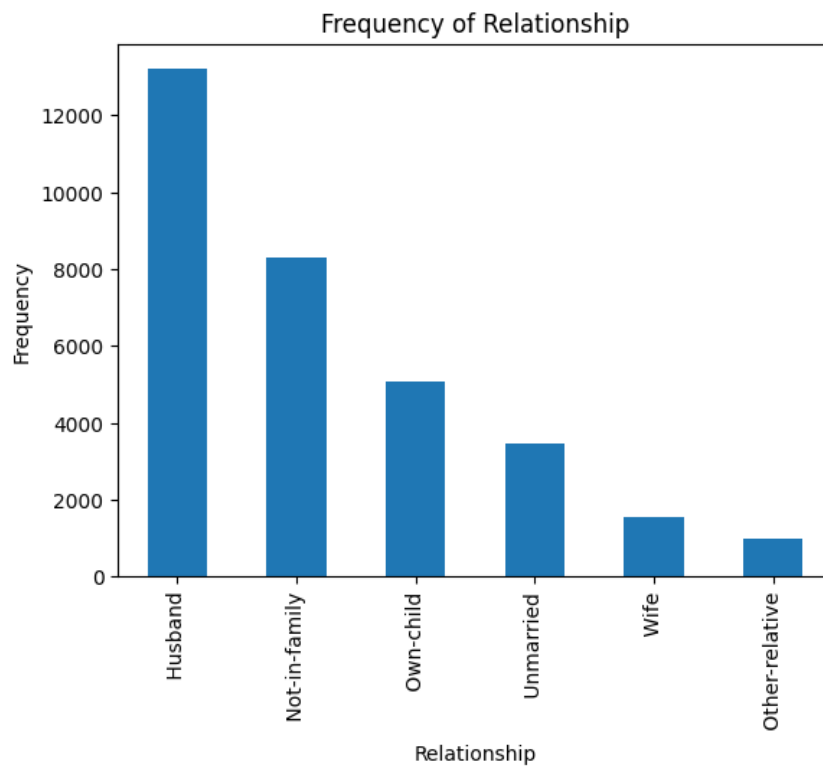
hours-per-week:

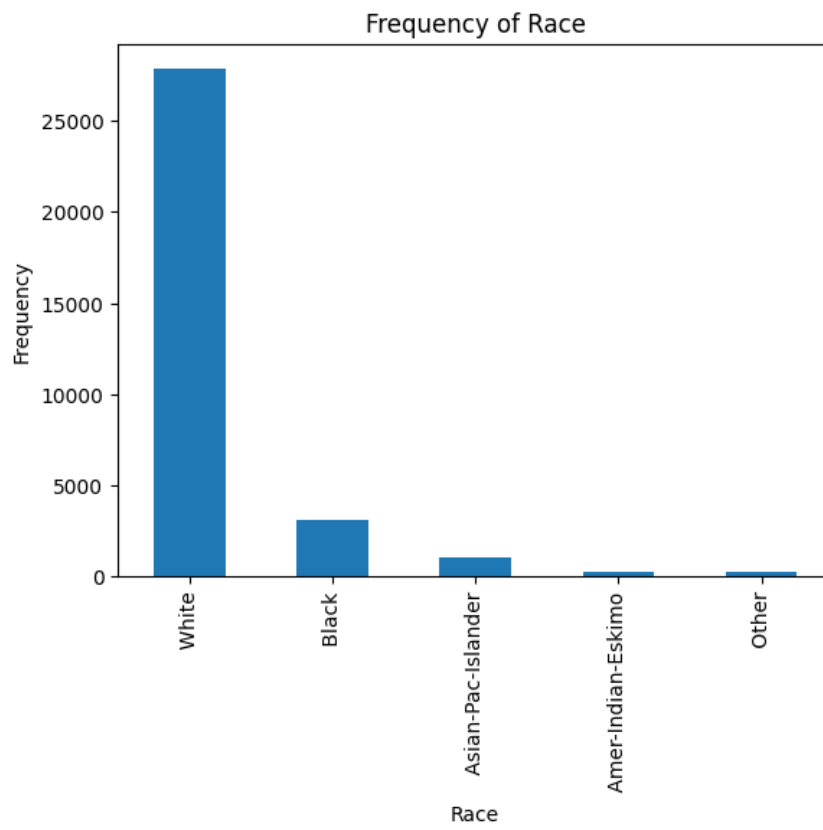
Minimum: 1
Q1 (25th percentile): 40.0
Median (50th percentile): 40.0
Q3 (75th percentile): 45.0
Maximum: 99
Mode: 40
Mean: 40.4

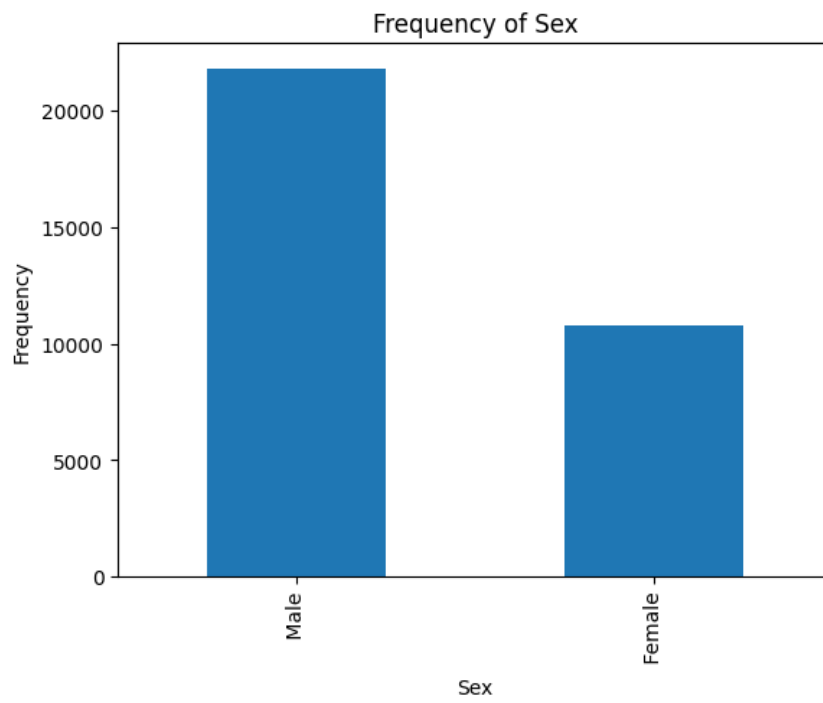
3.3

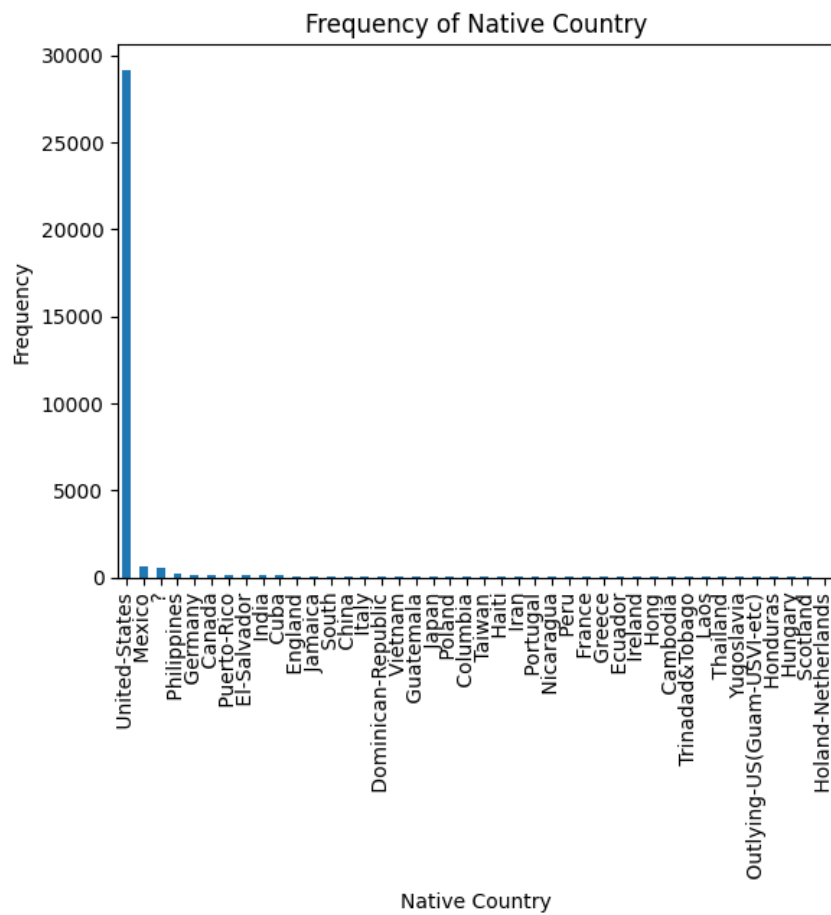












3.4

Unable to find missing data.

3.5

Mean Values:

age: 38.6

fnlwgt: 189778.4

education-num: 10.1

capital-gain: 1077.6

capital-loss: 87.3

hours-per-week: 40.4

Standard Deviation Values:

age: 13.6

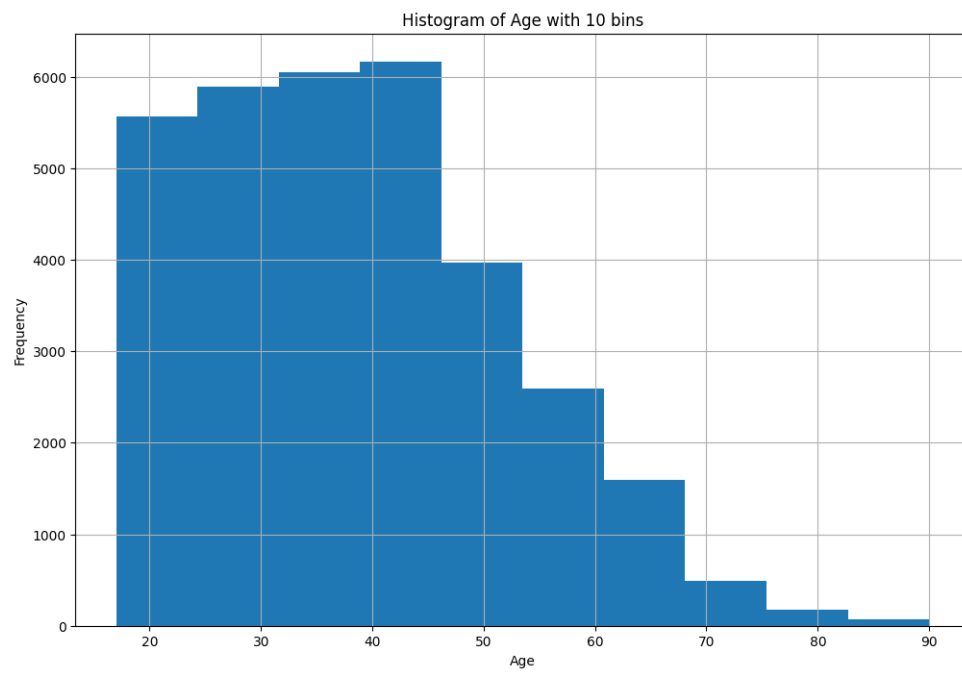
fnlwgt: 105550.0

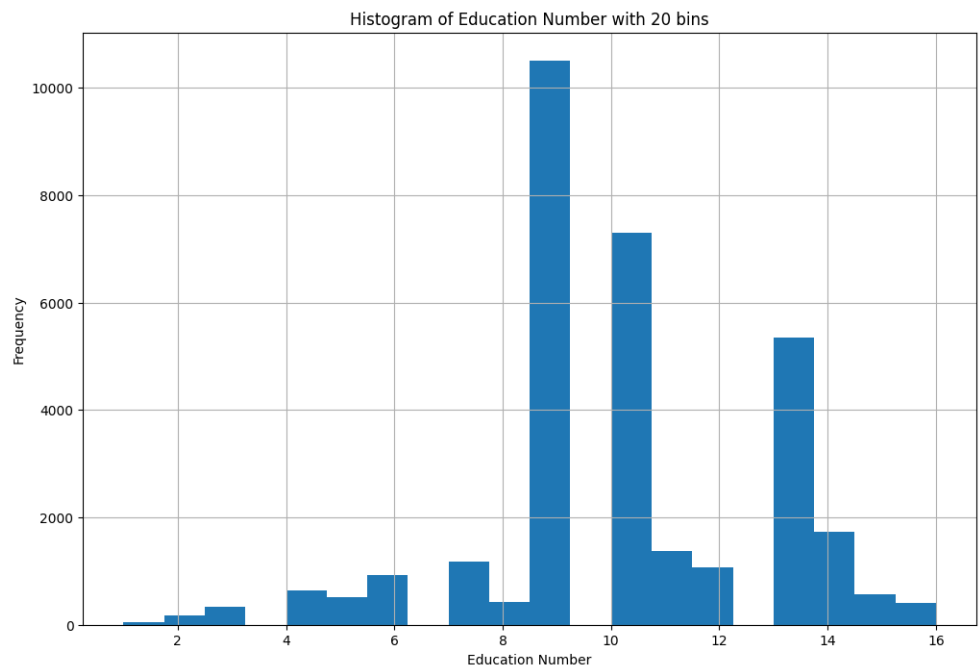
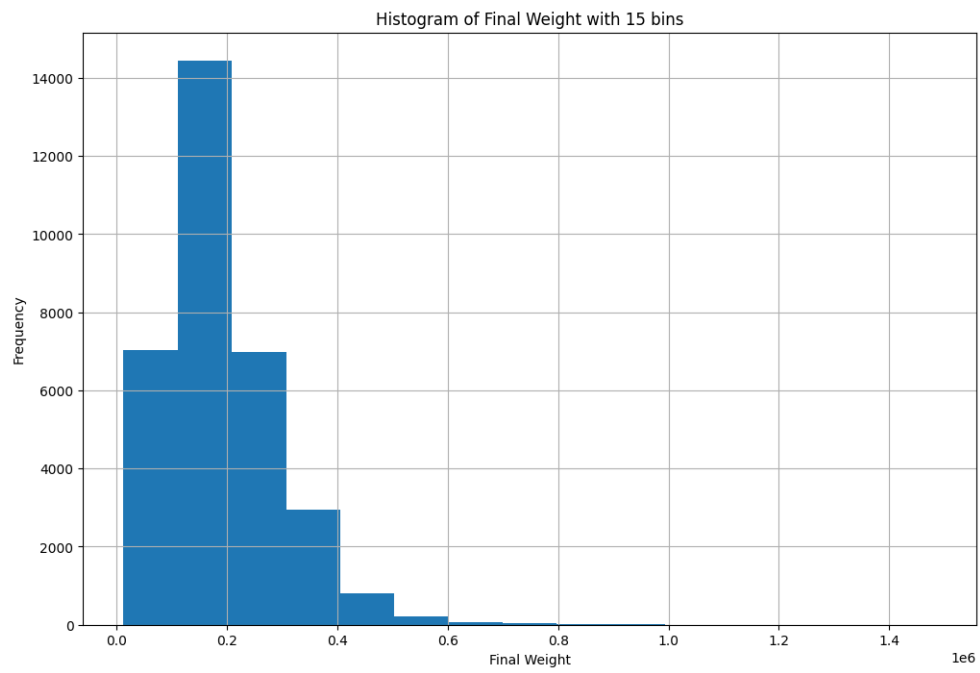
education-num: 2.6

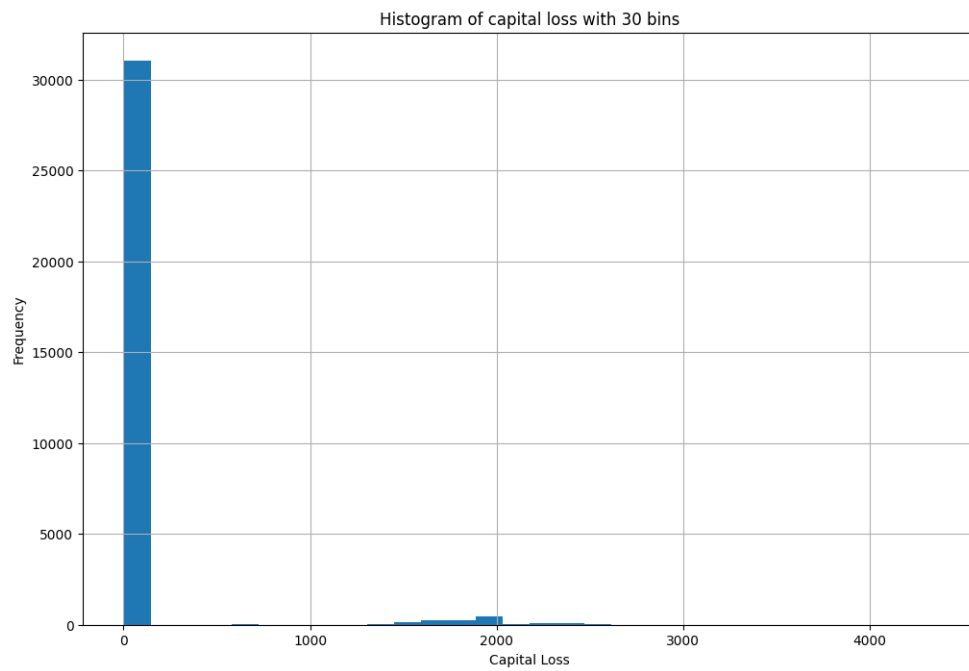
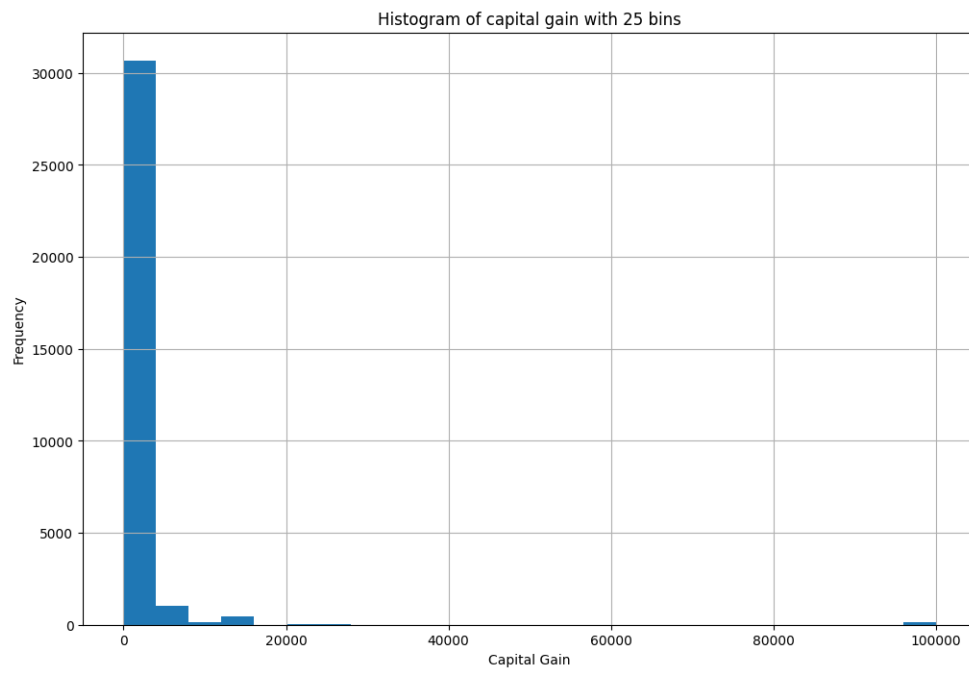
capital-gain: 7385.3

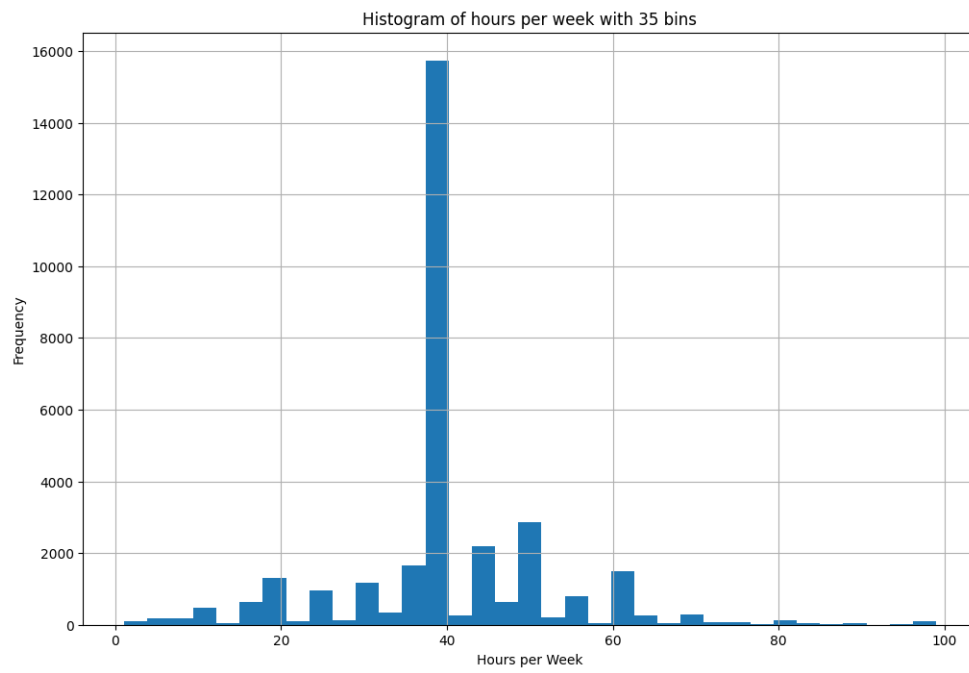
capital-loss: 403.0

hours-per-week: 12.3

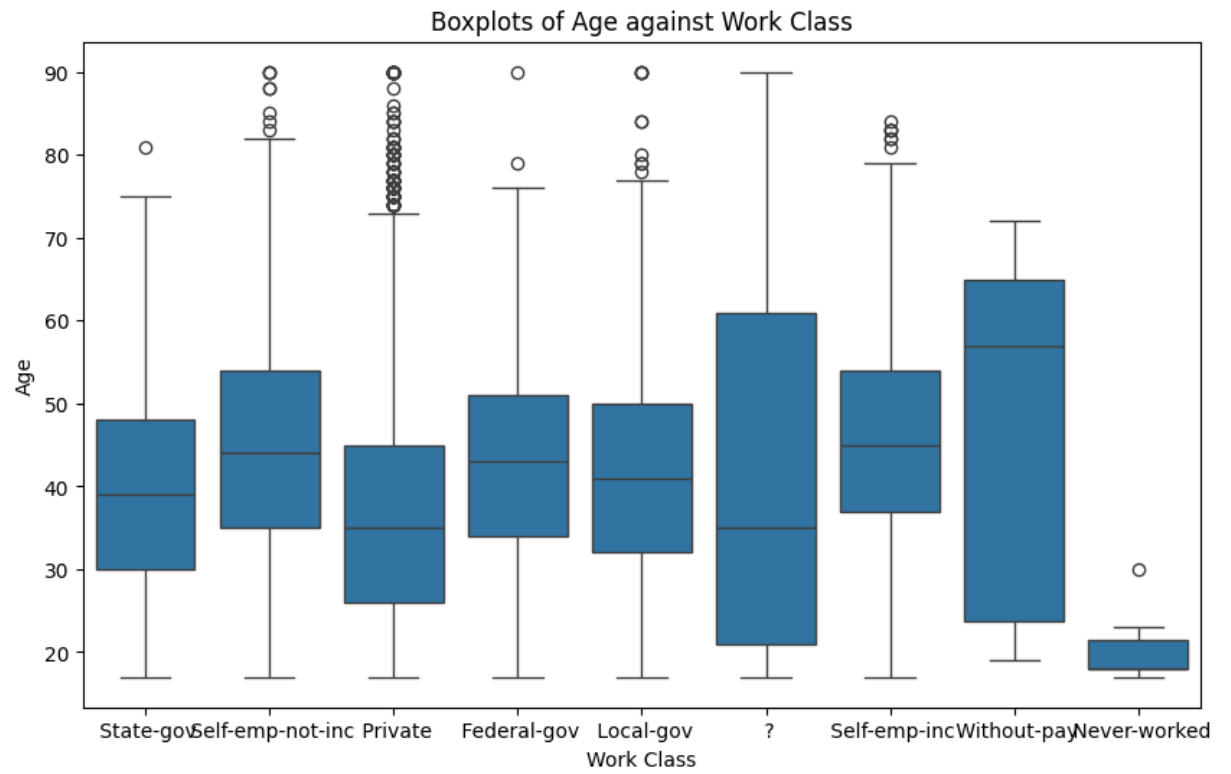


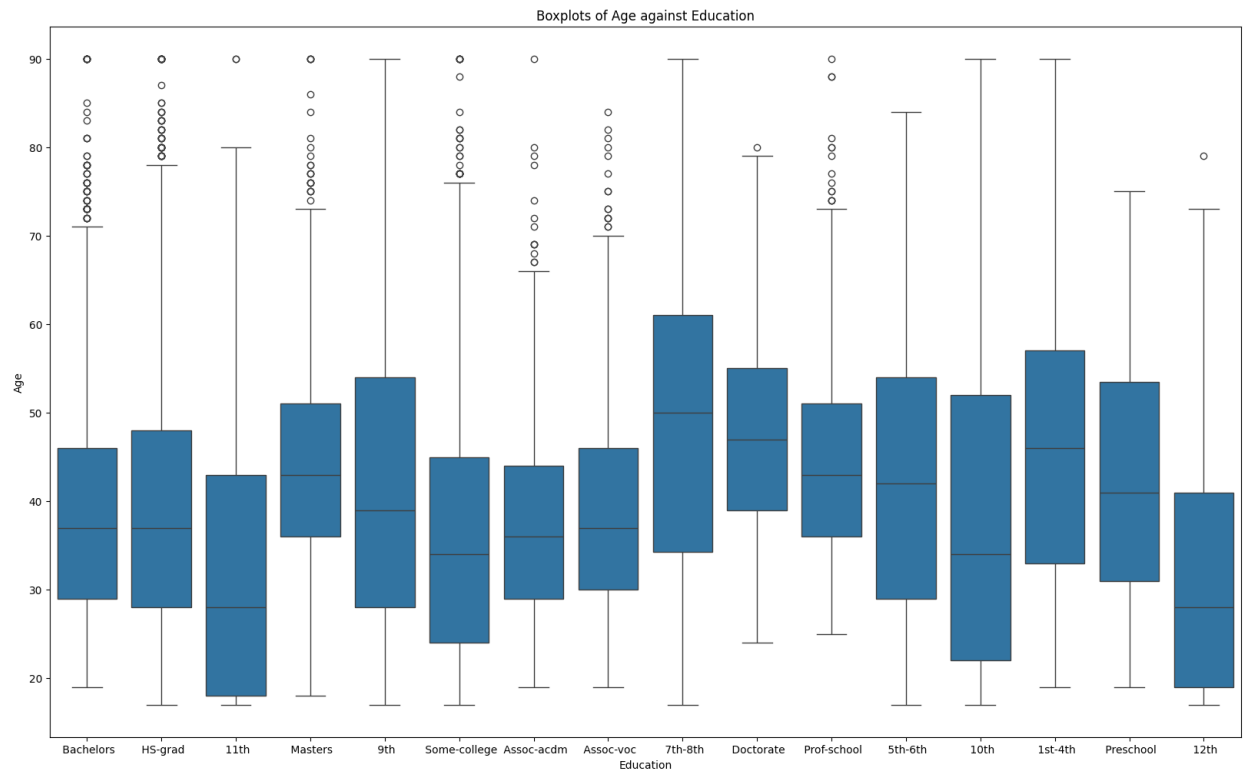


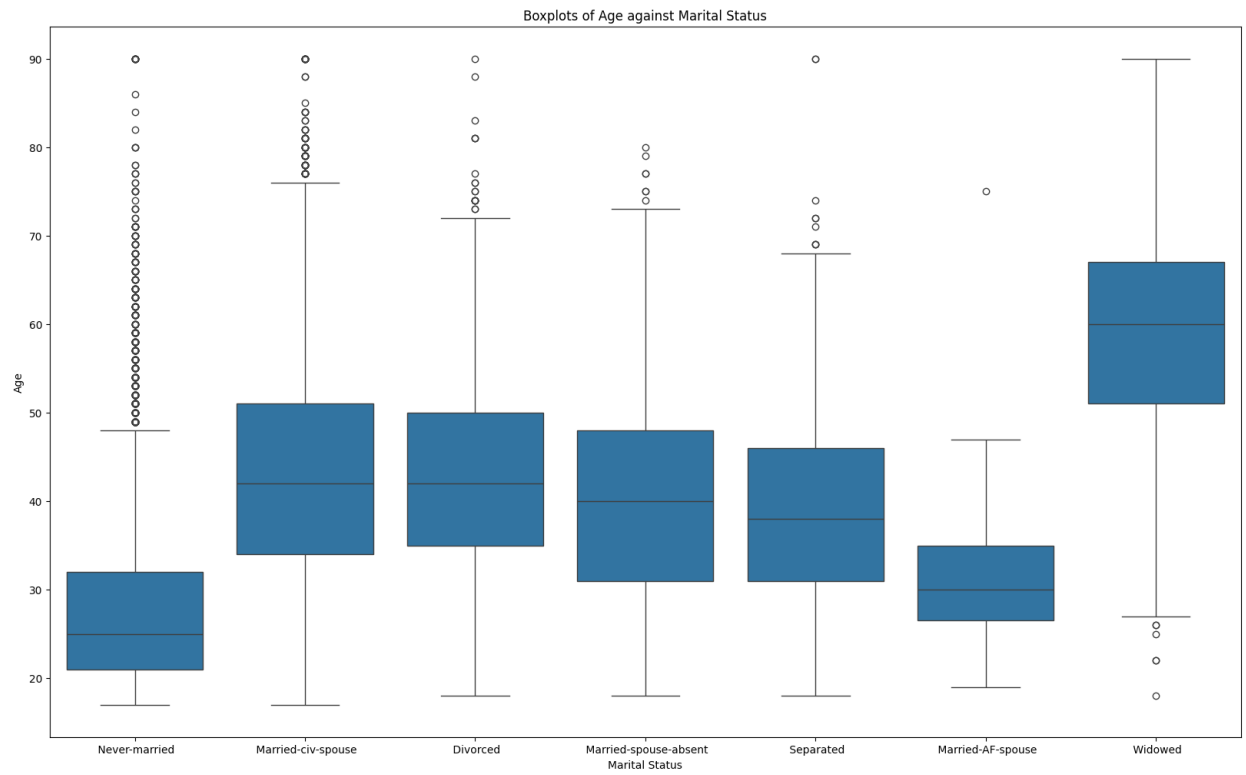


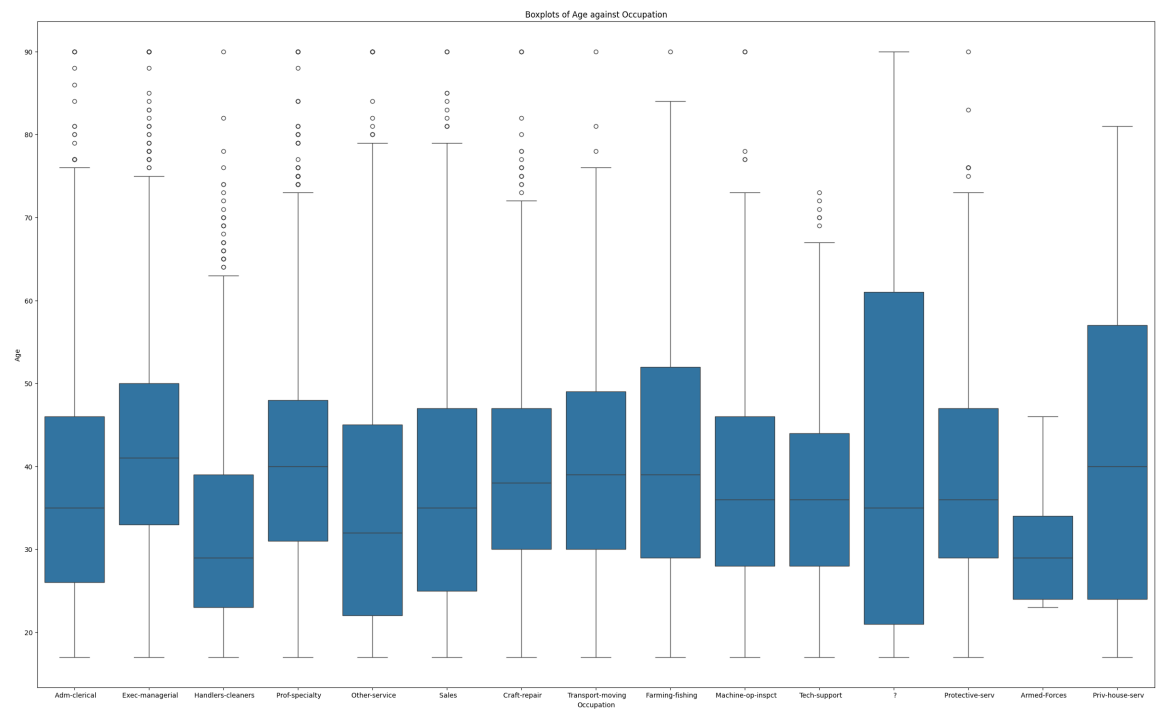


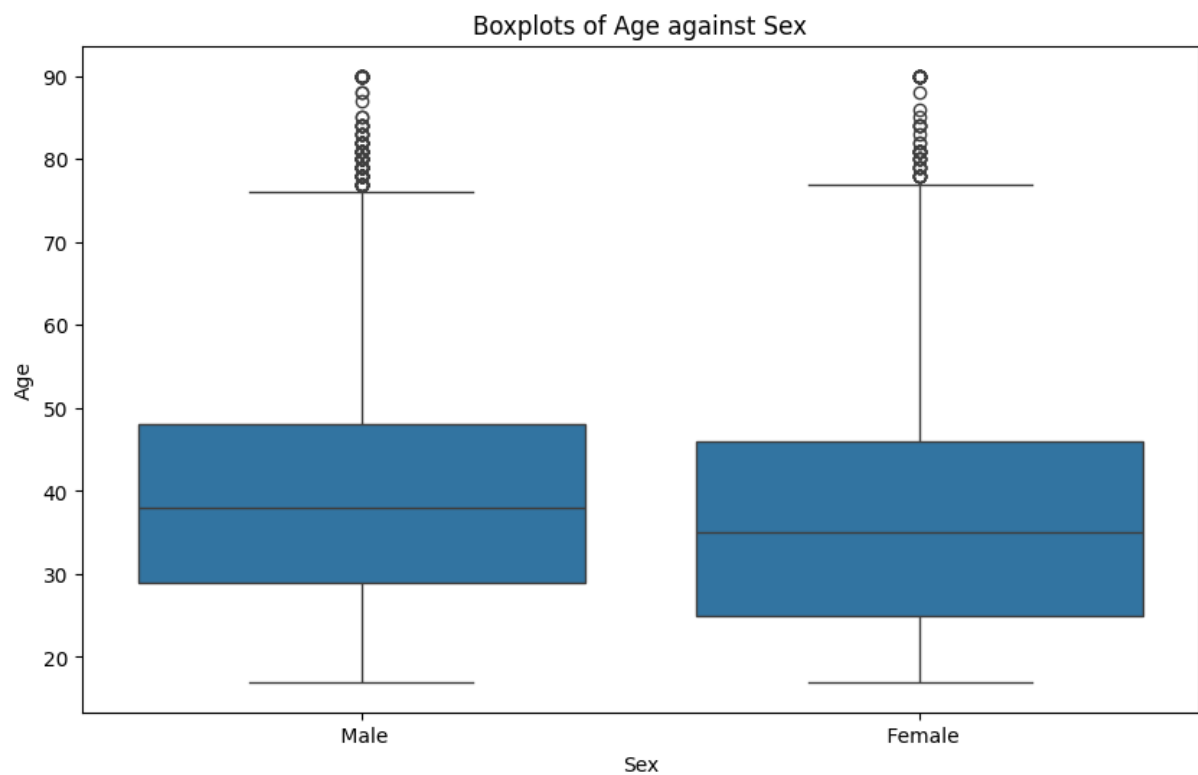
3.6











3.7

Age and Work Class

workclass	?	Federal-gov	Local-gov	Never-worked	Private
age					
17	64	1	14	1	300
18	92	3	10	3	413
19	113	5	11	0	540
20	115	9	12	1	581
21	89	3	9	0	577
...
86	0	0	0	0	1
87	1	0	0	0	0
88	0	0	0	0	1
90	7	1	4	0	28
Total	1836	960	2093	7	22696

workclass	Self-emp-inc	Self-emp-not-inc	State-gov	Without-pay	Total
age					
17	7	6	2	0	395
18	6	14	9	0	550
19	3	17	21	2	712
20	3	11	21	0	753
21	1	10	30	1	720
...
86	0	0	0	0	1
87	0	0	0	0	1
88	0	2	0	0	3
90	0	3	0	0	43
Total	1116	2541	1298	14	32561

Age and Education

education	10th	11th	12th	1st-4th	5th-6th	7th-8th	9th
age							
17	138	180	37	0	1	3	25
18	38	166	60	0	4	5	8
19	16	45	29	3	2	9	6
20	15	27	14	1	5	3	12
21	17	16	8	1	4	7	10
...
86	0	0	0	0	0	0	0
87	0	0	0	0	0	0	0
88	0	0	0	0	0	0	0
90	1	2	0	1	0	3	1
Total	933	1175	433	168	333	646	514

education	Assoc-acdm	Assoc-voc	Bachelors	Doctorate	HS-grad
age					
17	0	0	0	0	8
18	0	0	0	0	152
19	1	3	2	0	266
20	13	10	1	0	237
21	21	29	7	0	225
...
86	0	0	0	0	0
87	0	0	0	0	1
88	0	0	0	0	0
90	1	0	9	0	14
Total	1067	1382	5355	413	10501

education	Masters	Preschool	Prof-school	Some-college	Total
age					
17	0	0	0	3	395
18	1	0	0	116	550
19	0	1	0	329	712
20	1	1	0	413	753
21	1	2	0	372	720
...
86	1	0	0	0	1
87	0	0	0	0	1
88	0	0	2	1	3
90	4	0	1	6	43
Total	1723	51	576	7291	32561

Age and Marital Status

marital-status	Divorced	Married-AF-spouse	Married-civ-spouse
age			
17	0	0	2
18	1	0	7
19	6	2	15
20	7	0	38
21	4	0	53
...
86	0	0	0
87	0	0	0
88	1	0	2
90	1	0	20
Total	4443	23	14976

marital-status	Married-spouse-absent	Never-married	Separated	Widowed
age				
17	0	393	0	0
18	1	539	1	1
19	3	681	5	0
20	5	695	8	0
21	6	651	6	0
...
86	0	1	0	0
87	0	0	0	1
88	0	0	0	0
90	0	14	2	6
Total	418	10683	1025	993

marital-status	Total
age	
17	395
18	550
19	712
20	753
21	720
...	...
86	1
87	1
88	3
90	43
Total	32561

Age and Occupation

occupation	?	Adm-clerical	Armed-Forces	Craft-repair
age				
17	65	23	0	14
18	95	55	0	17
19	113	102	0	40
20	116	117	0	35
21	89	121	0	59
...
86	0	1	0	0
87	1	0	0	0
88	0	1	0	0
90	7	4	0	3
Total	1843	3770	9	4099

occupation	Exec-managerial	Farming-fishing	Handlers-cleaners
age			
17	1	9	40
18	6	14	50
19	12	24	65
20	15	23	81
21	18	25	51
...
86	0	0	0
87	0	0	0
88	1	0	0
90	8	1	1
Total	4066	994	1370

occupation age	Machine-op-inspct	Other-service	Priv-house-serv
17	2	129	8
18	17	152	4
19	30	166	3
20	41	139	3
21	51	142	4
...
86	0	0	0
87	0	0	0
88	0	0	0
90	3	6	0
Total	2002	3295	149

occupation age	Prof-specialty	Protective-serv	Sales	Tech-support
17	10	3	87	1
18	10	5	115	2
19	18	3	112	8
20	28	9	108	14
21	30	7	93	16
...
86	0	0	0	0
87	0	0	0	0
88	1	0	0	0
90	5	1	3	0
Total	4140	649	3650	928

occupation age	Transport-moving	Total
17	3	395
18	8	550
19	16	712
20	24	753
21	14	720
...
86	0	1
87	0	1
88	0	3
90	1	43
Total	1597	32561

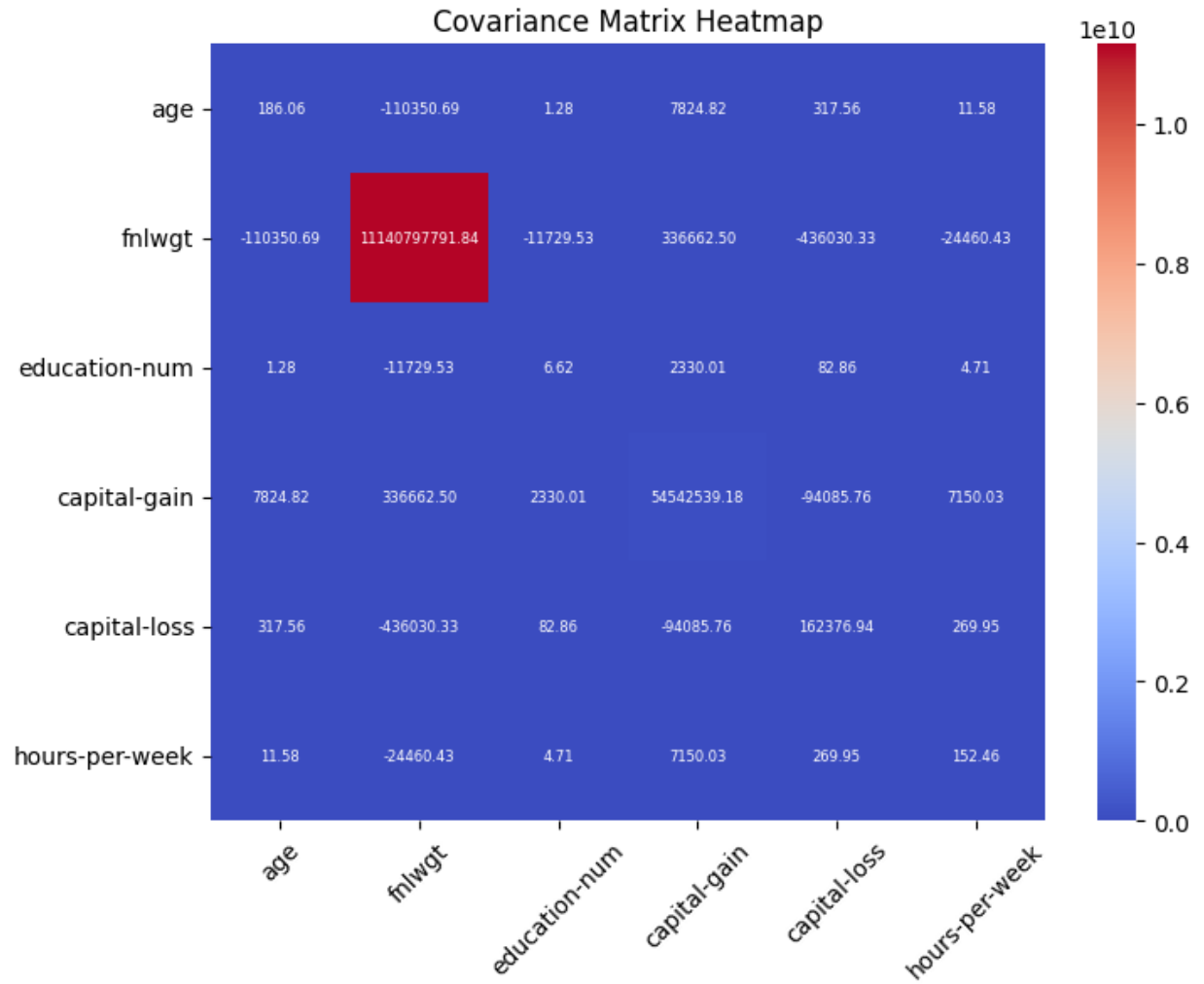
Age and Sex

sex	Female	Male	Total
age			
17	186	209	395
18	268	282	550
19	356	356	712
20	363	390	753
21	329	391	720
...
86	1	0	1
87	0	1	1
88	1	2	3
90	14	29	43
Total	10771	21790	32561

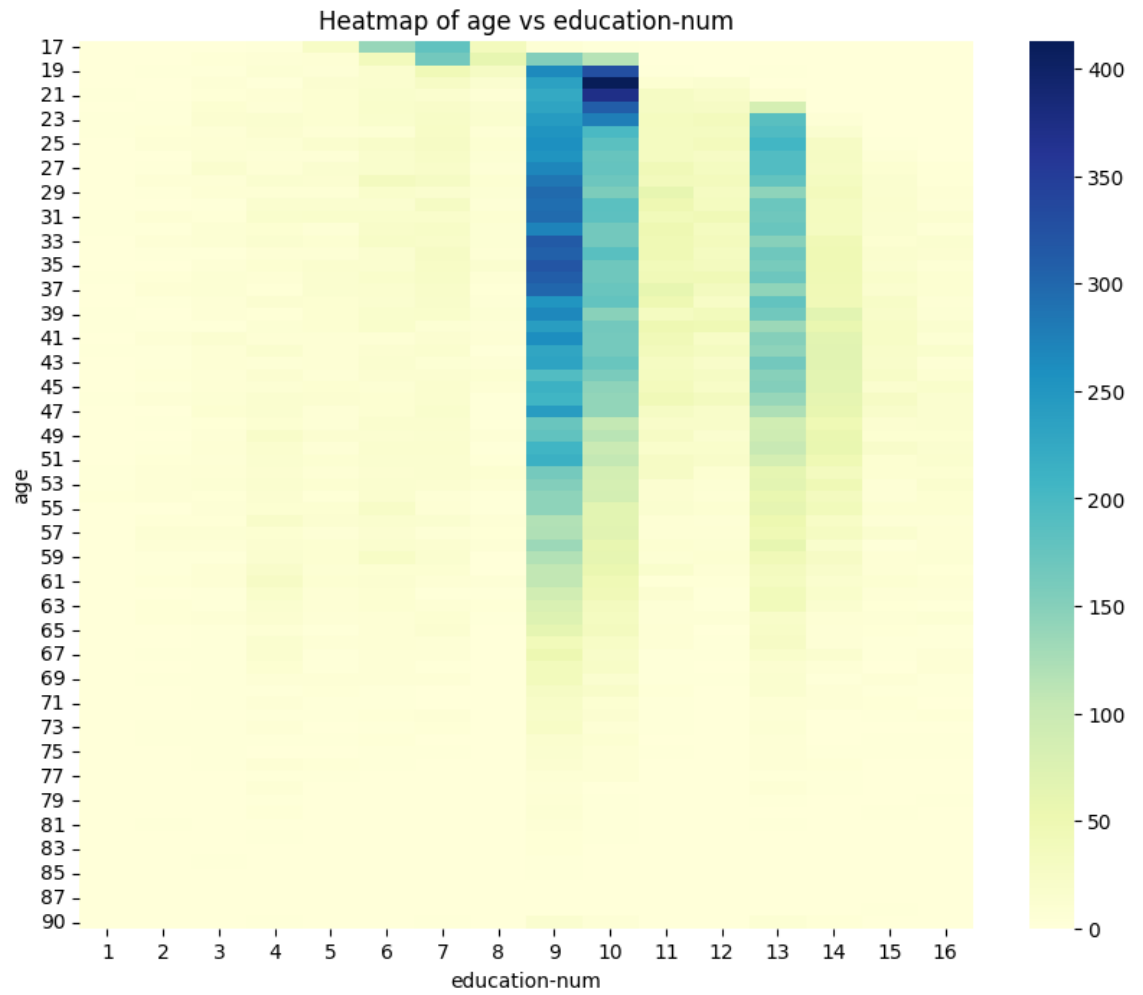
3.8

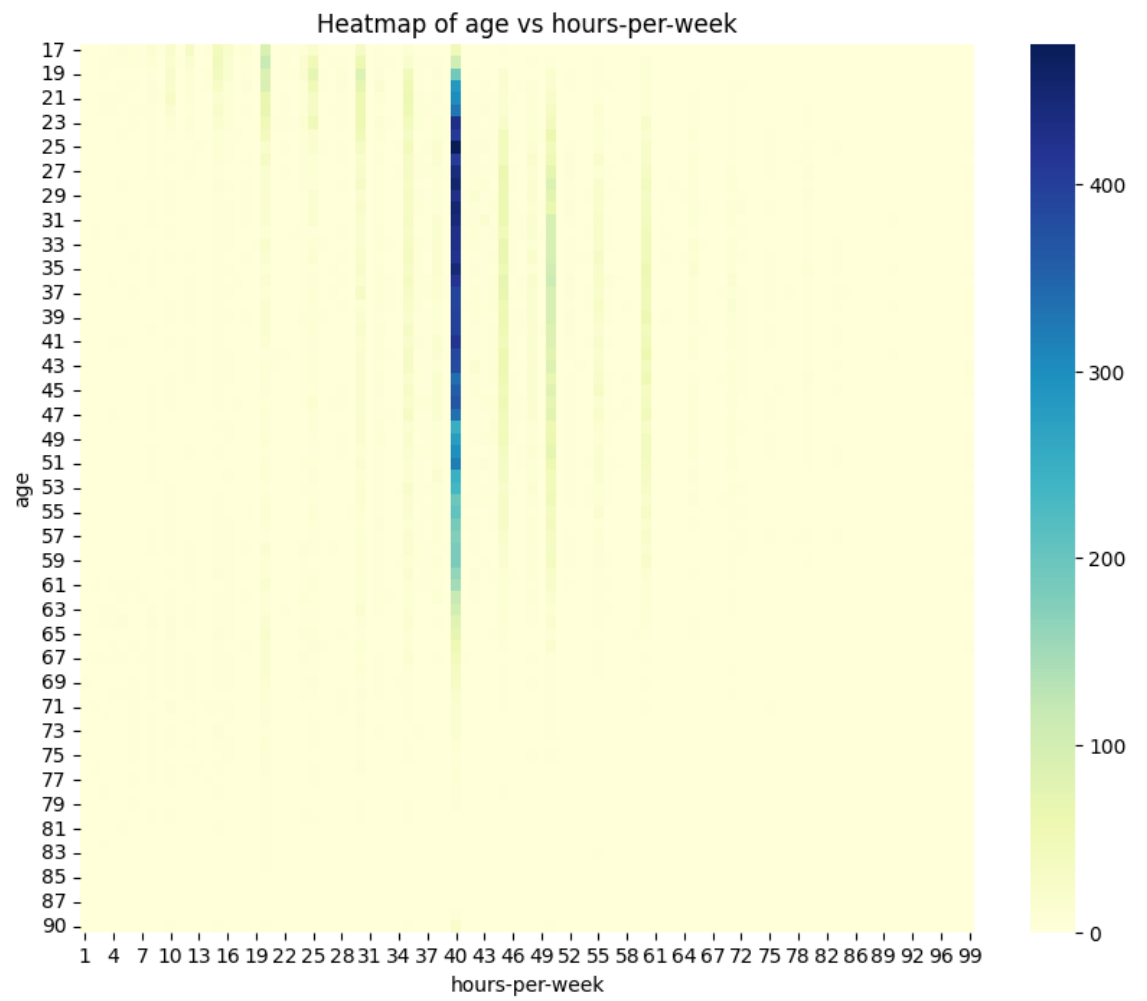
	age	fnlwgt	education-num	capital-gain	capital-loss	hours-per-week
age	186.06	-110350.69	1.28	7824.82	317.56	11.58
fnlwgt	-110350.69	11140797791.84	-11729.53	336662.50	-436030.33	-24460.43
education-num	1.28	-11729.53	6.62	2330.01	82.86	4.71
capital-gain	7824.82	336662.50	2330.01	54542539.18	-94085.76	7150.03
capital-loss	317.56	-436030.33	82.86	-94085.76	162376.94	269.95
hours-per-week	11.58	-24460.43	4.71	7150.03	269.95	152.46

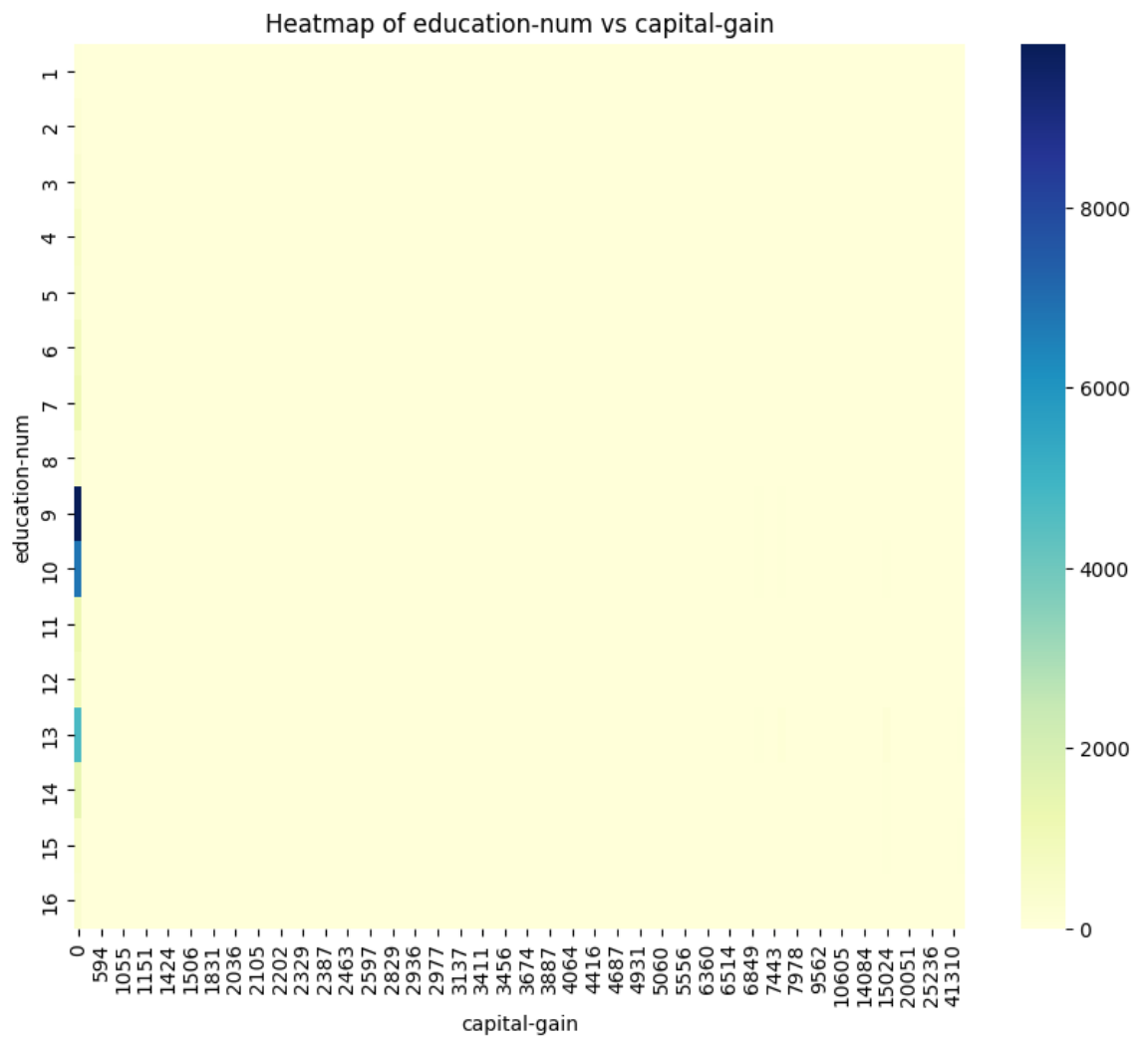
3.9

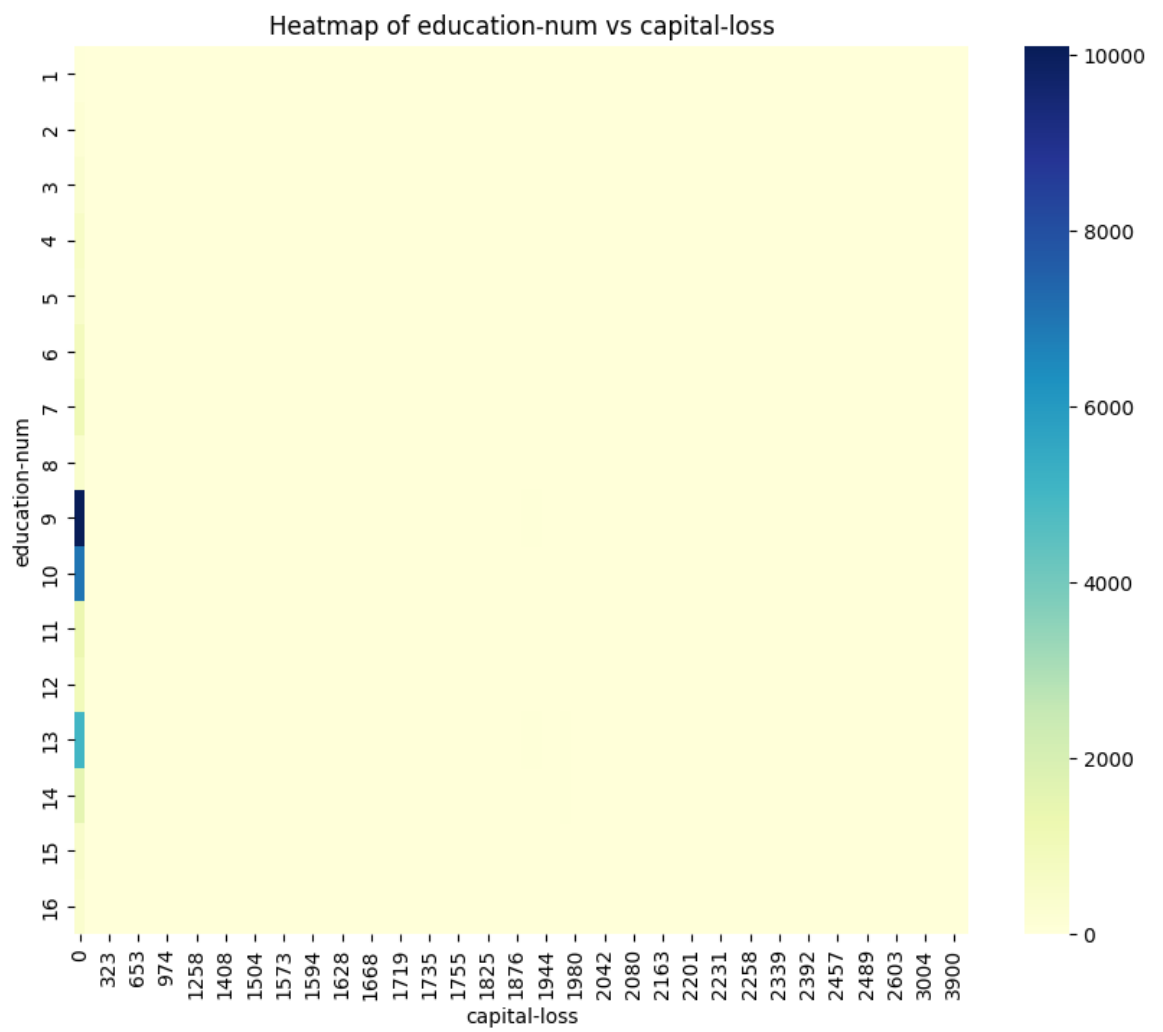


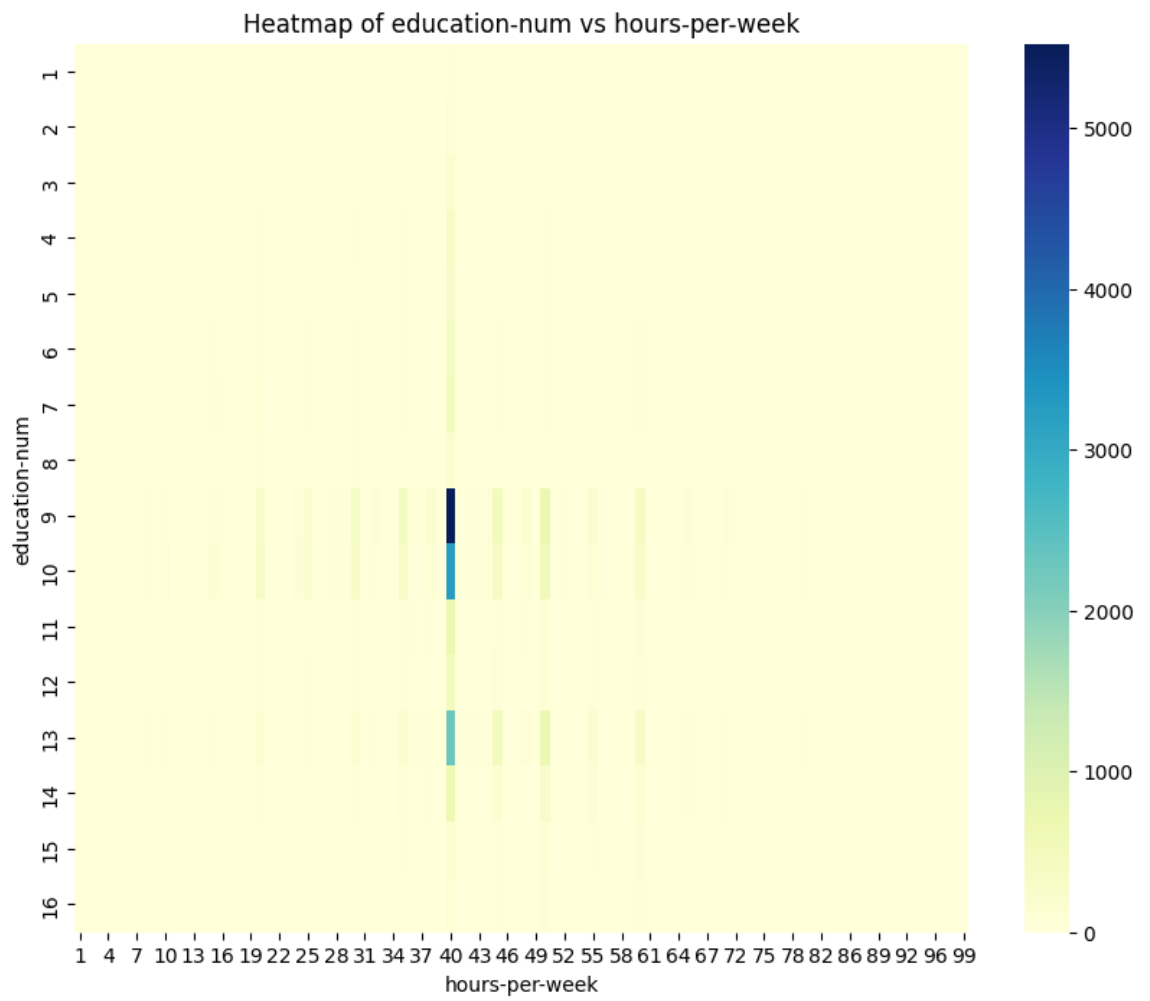
3.10





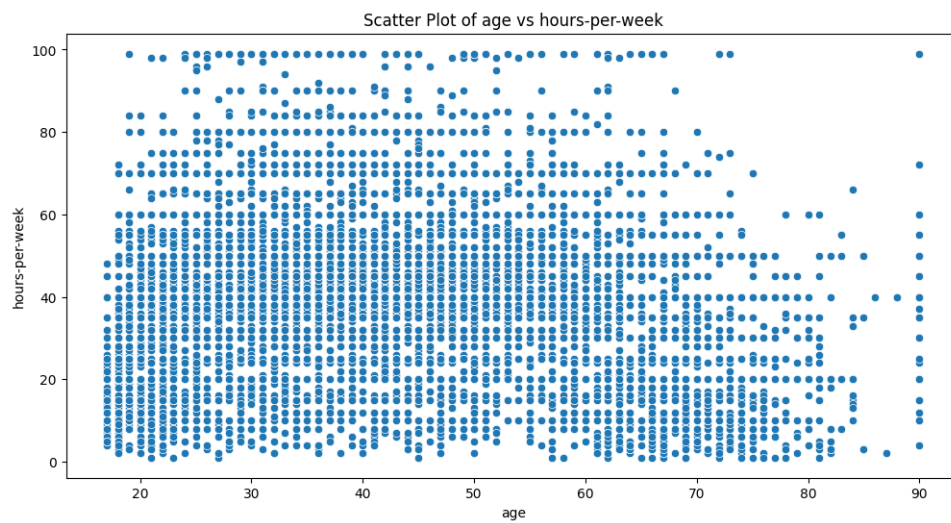
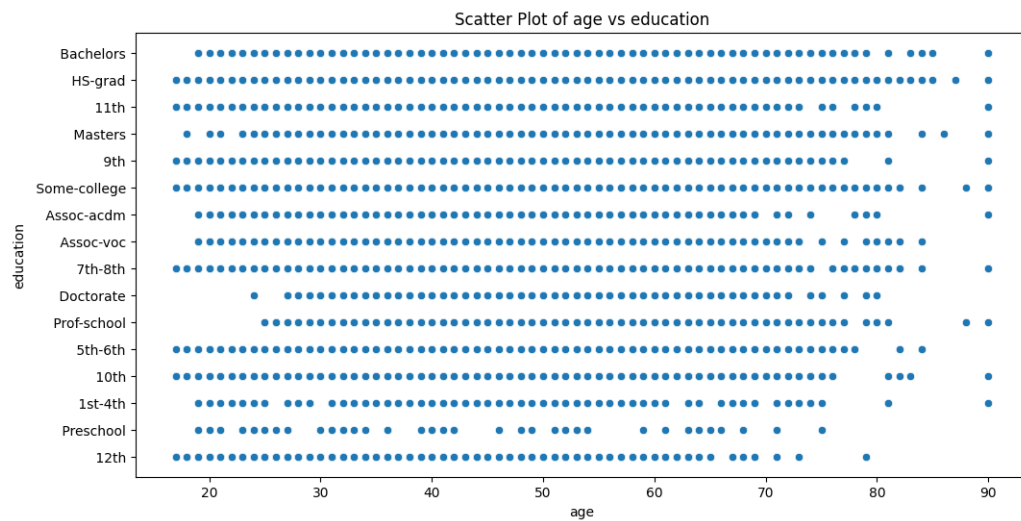


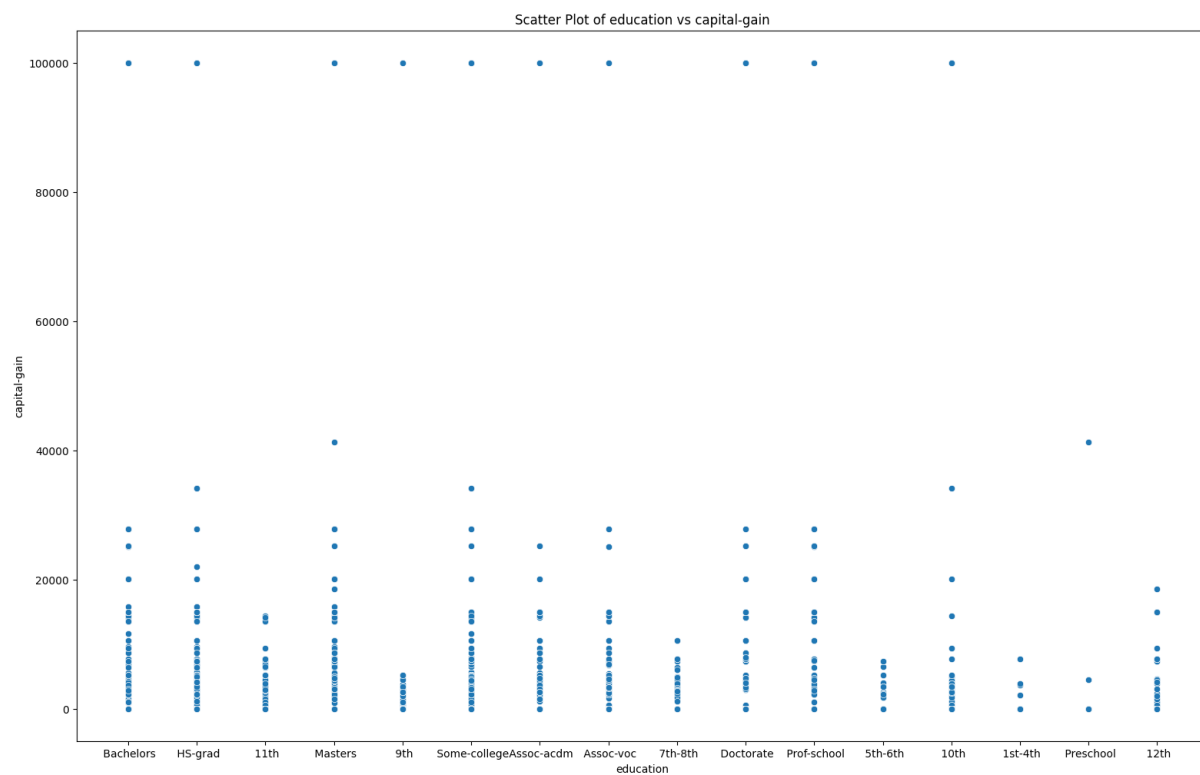


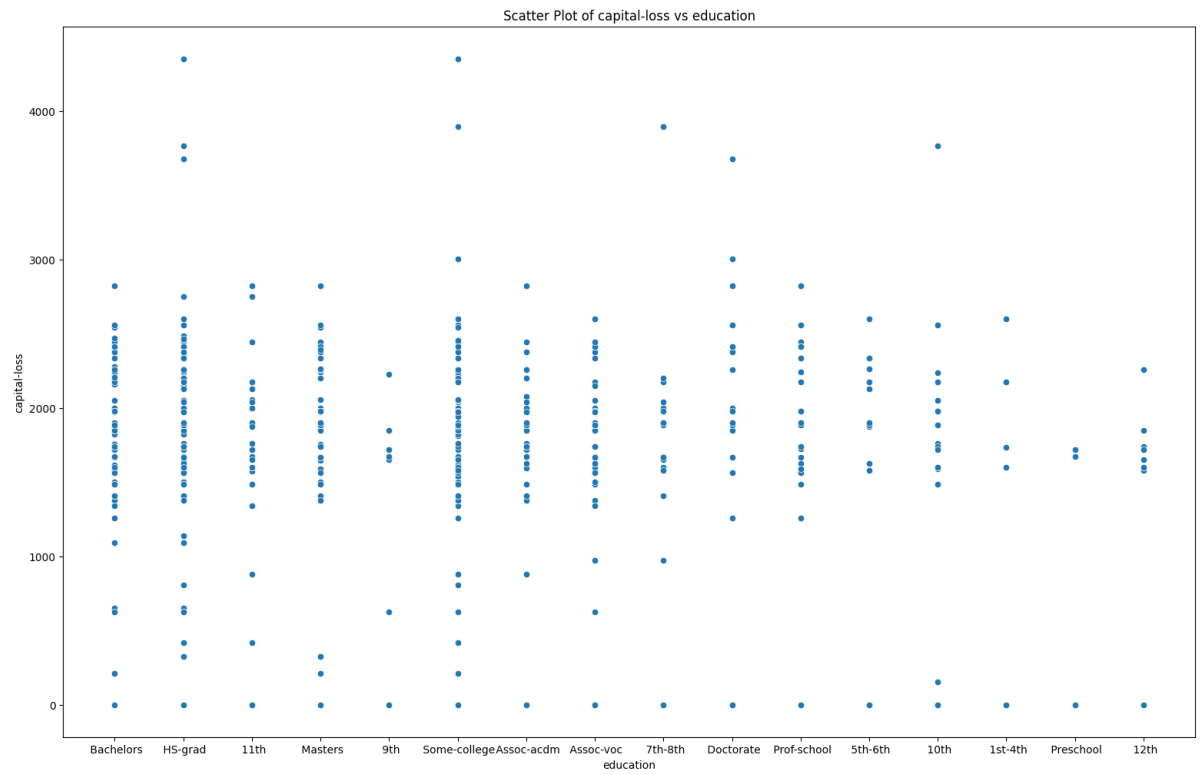


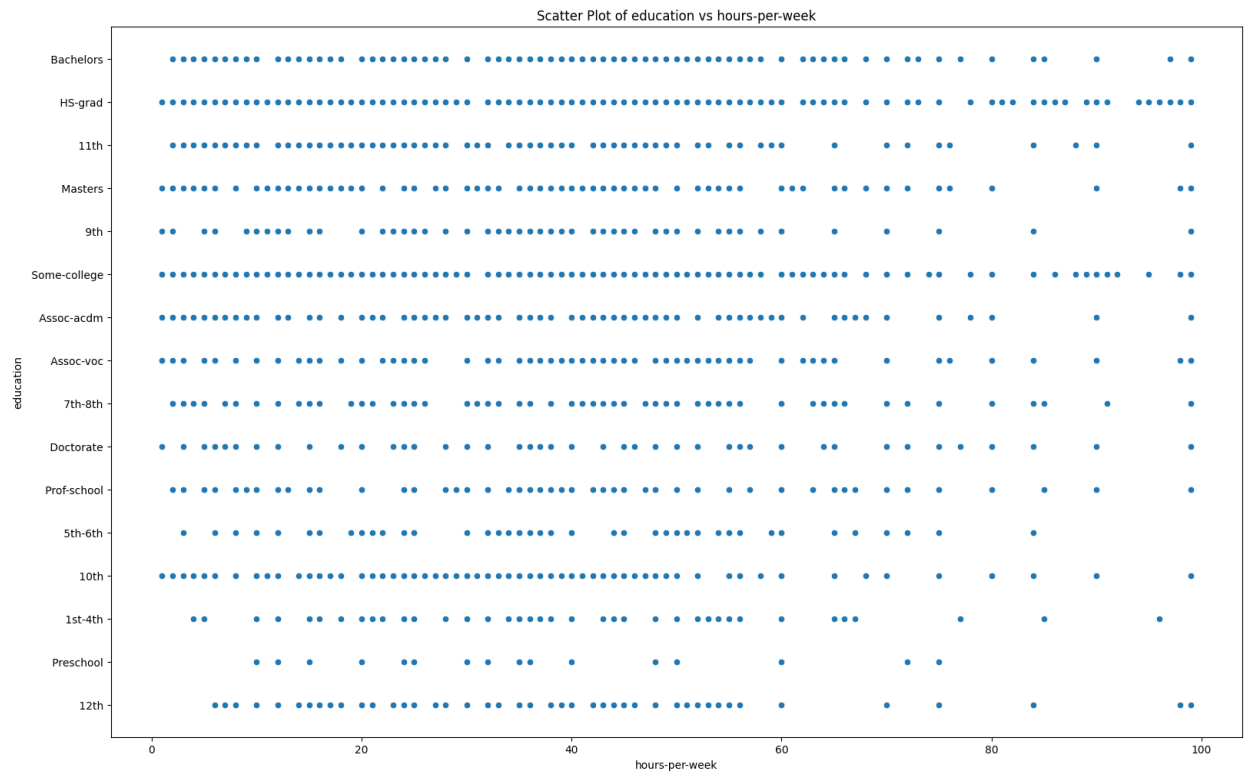
3.11

Scatter Plots

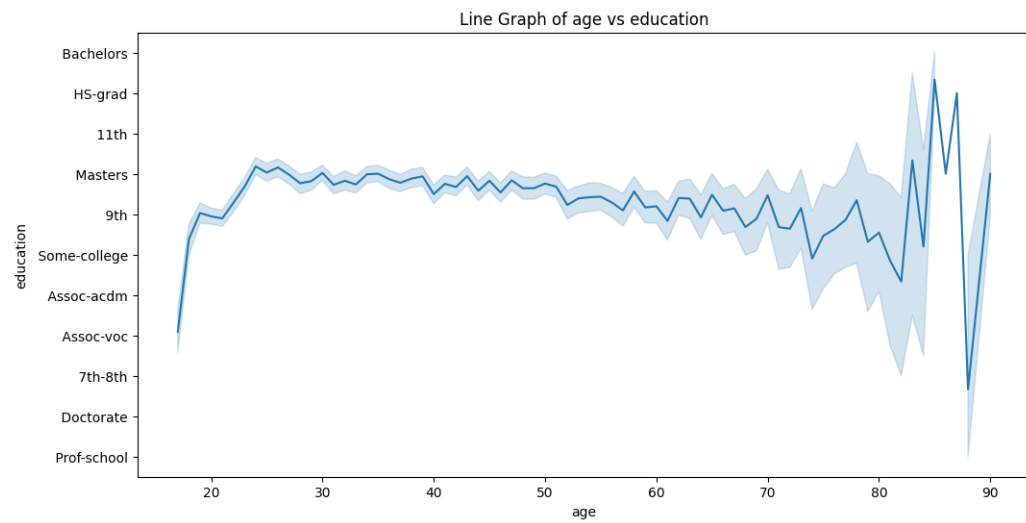


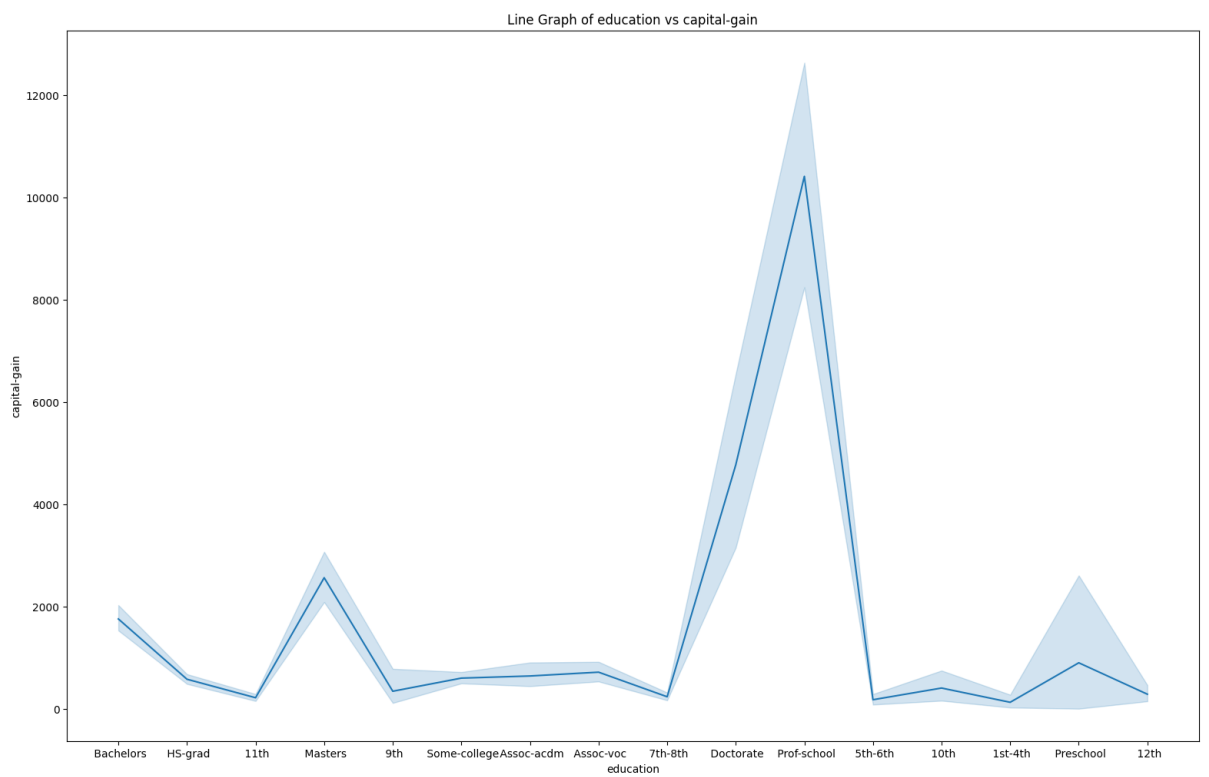
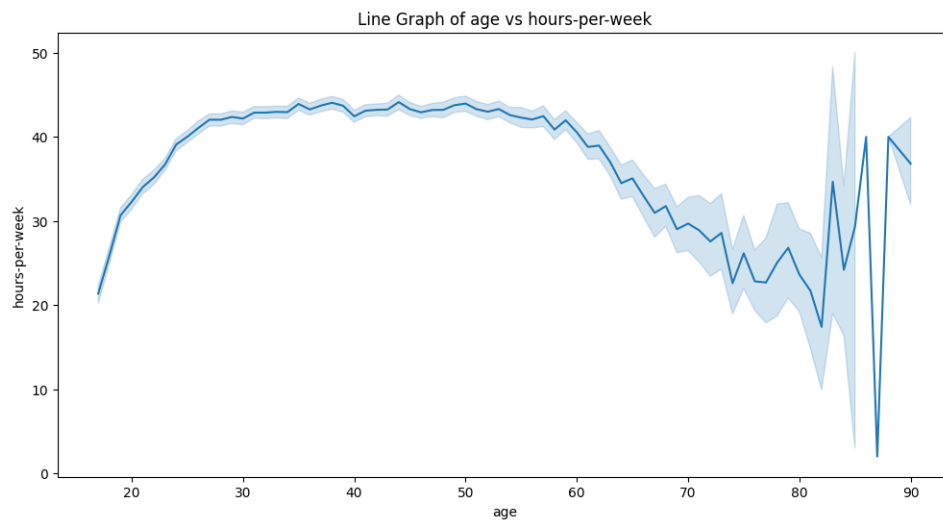


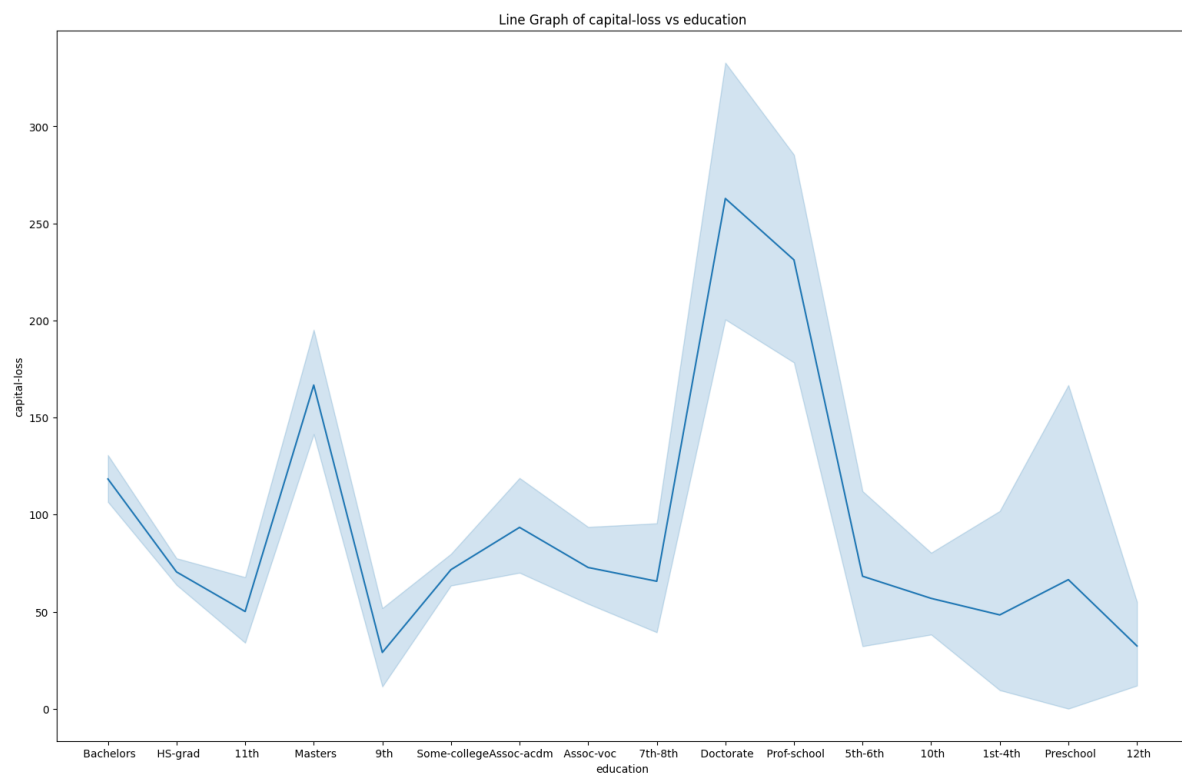


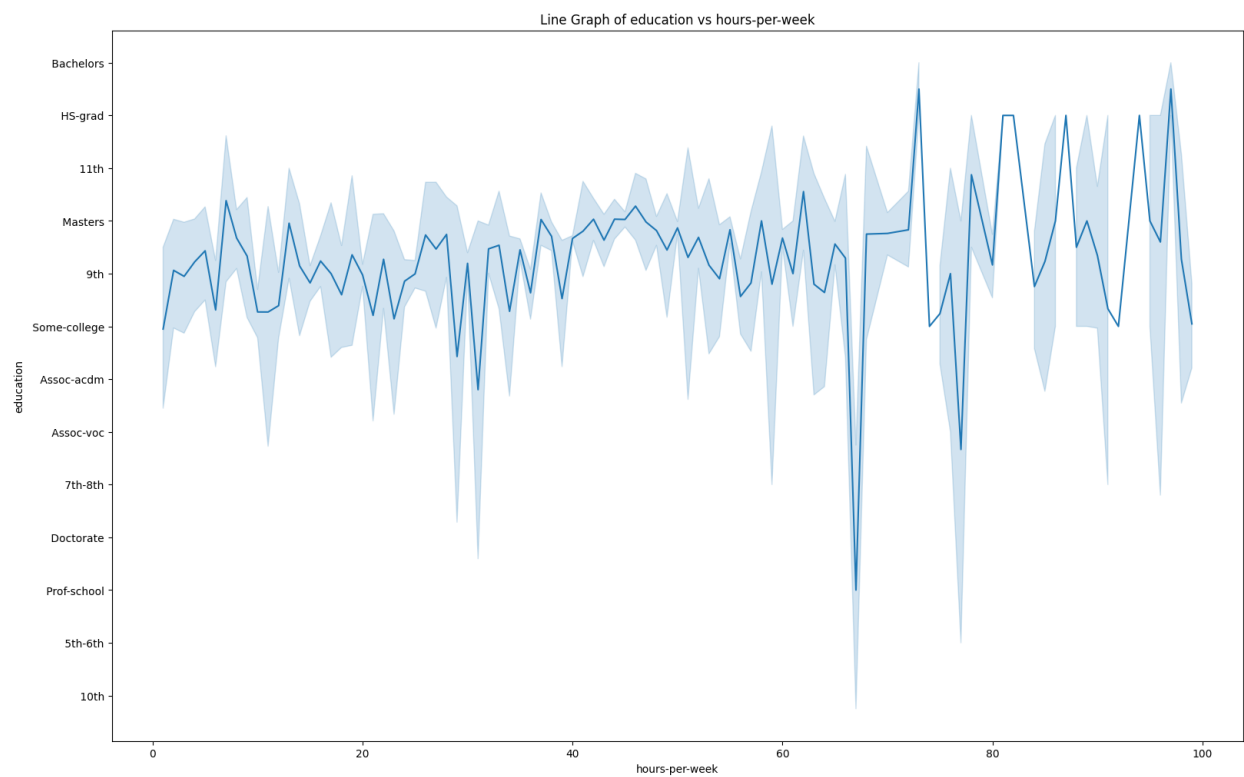


Line Graphs



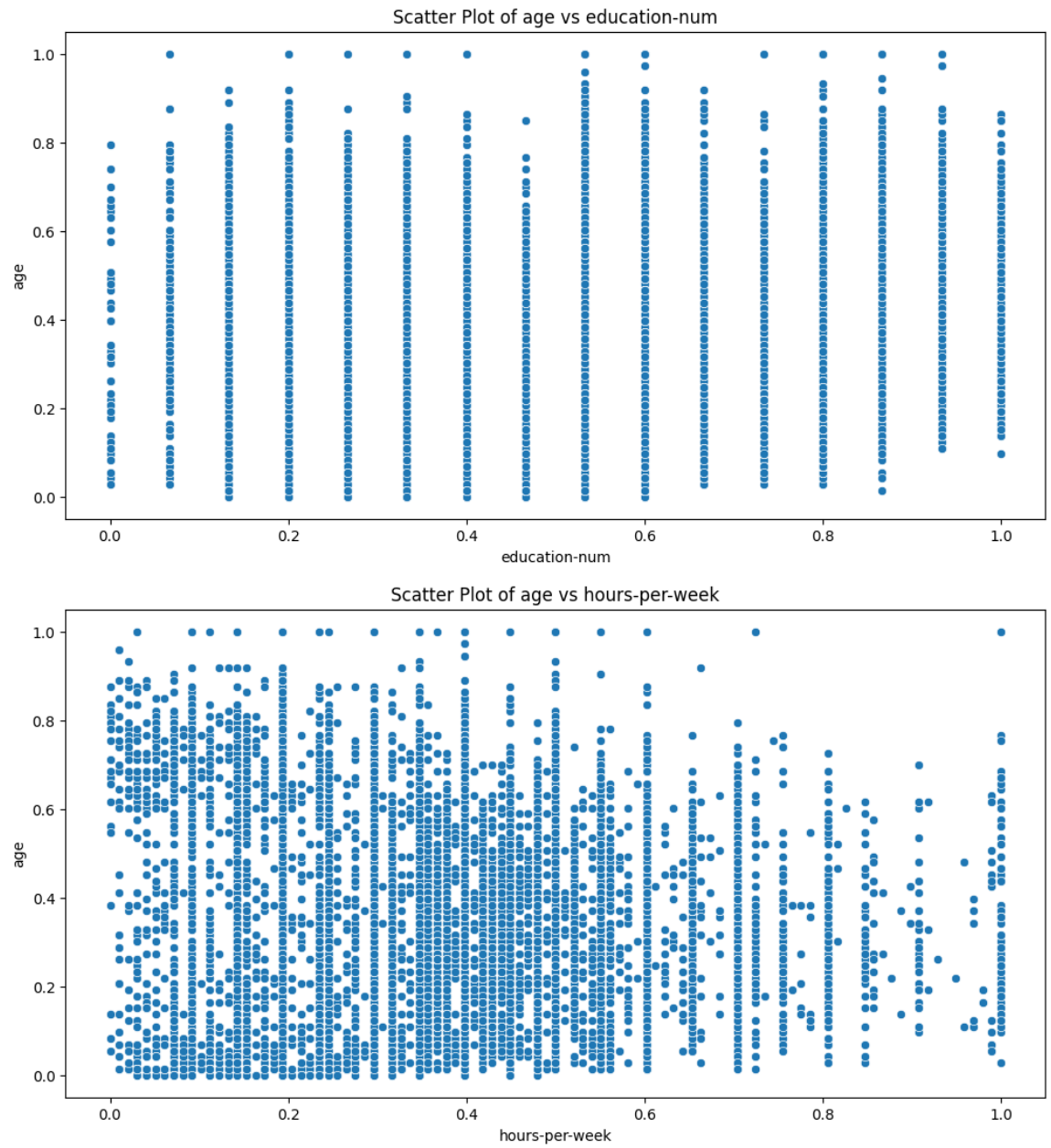


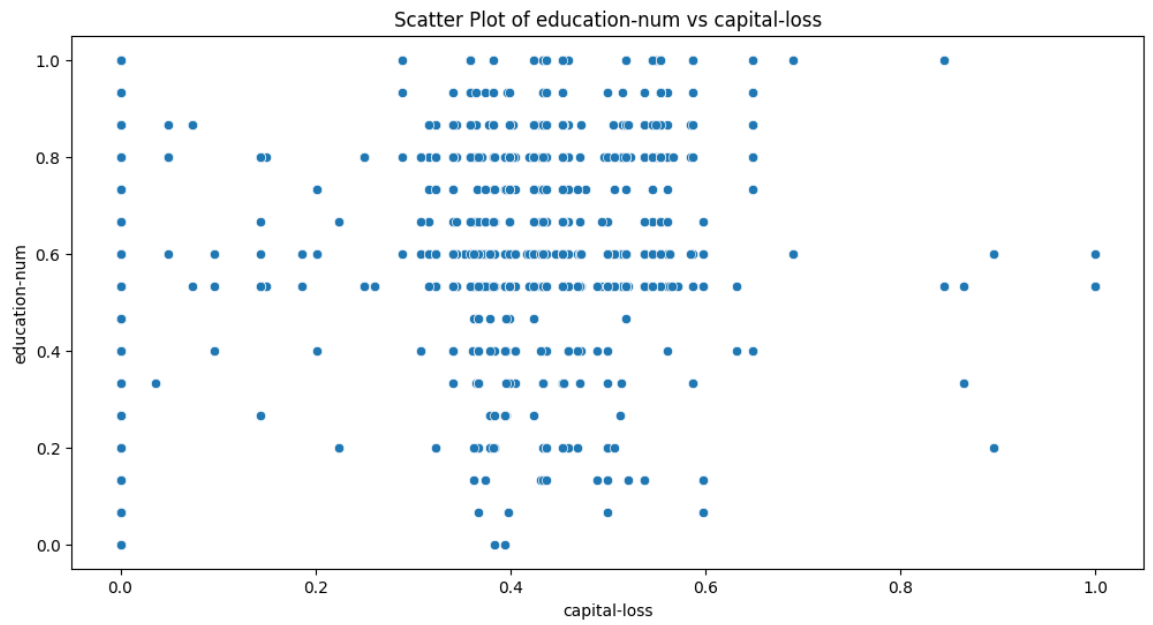
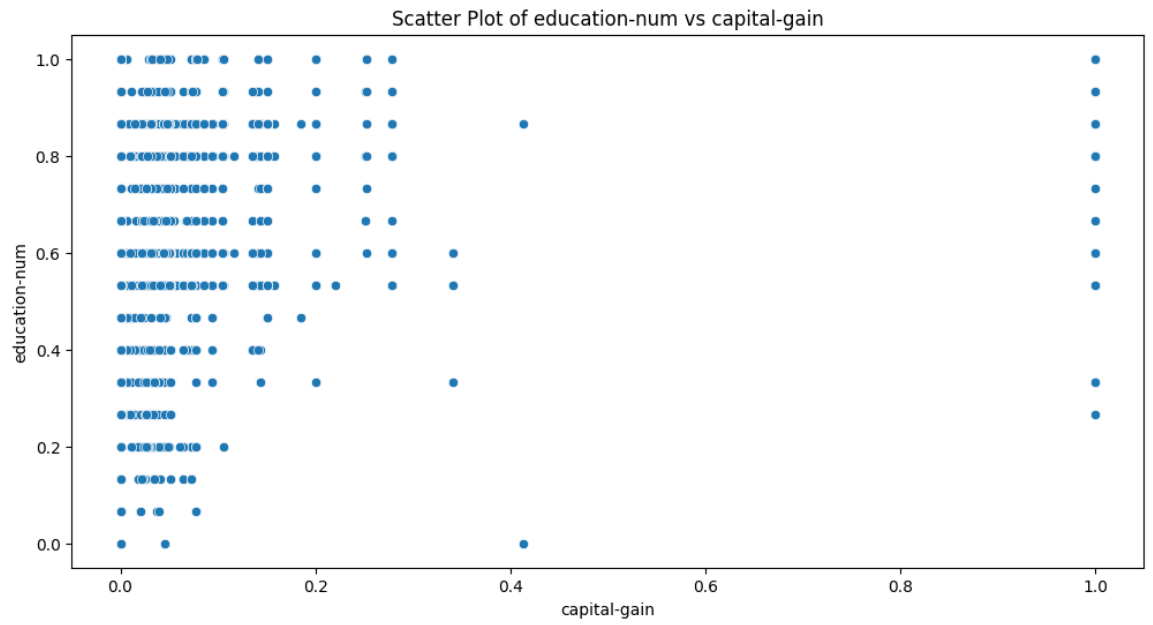


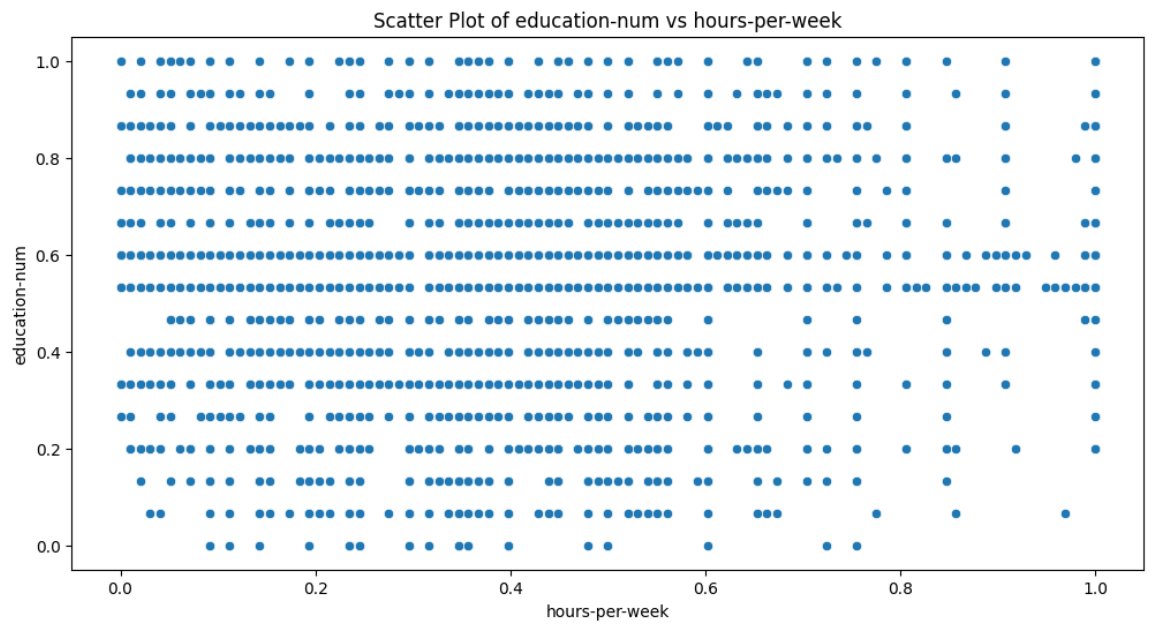


3.12

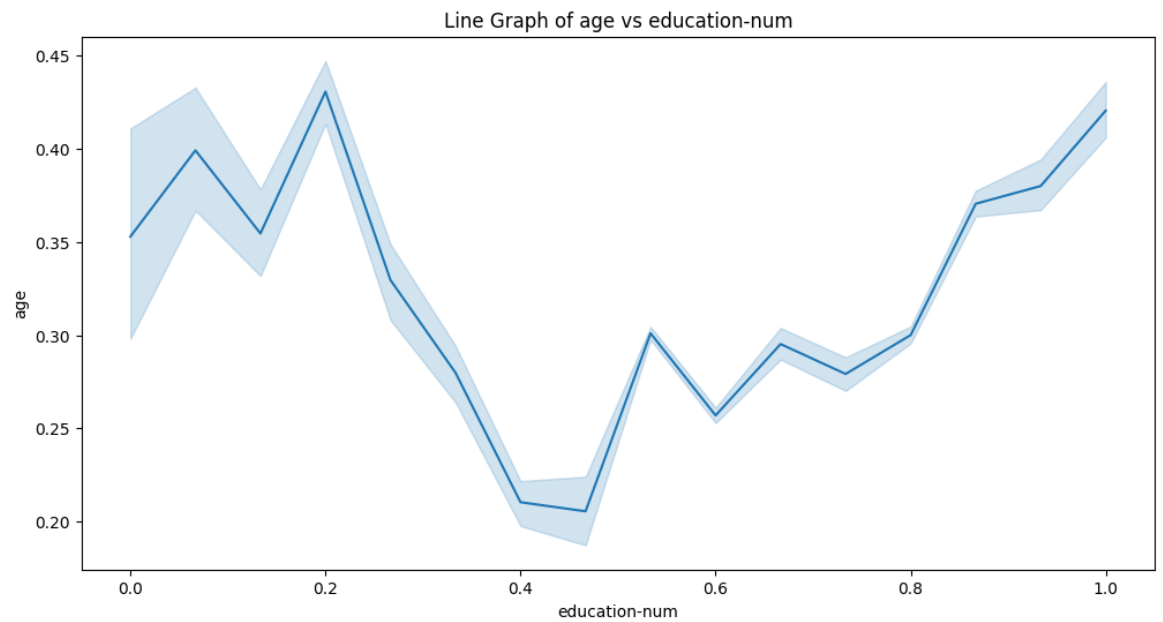
Min-Max Scaled Scatter Plots

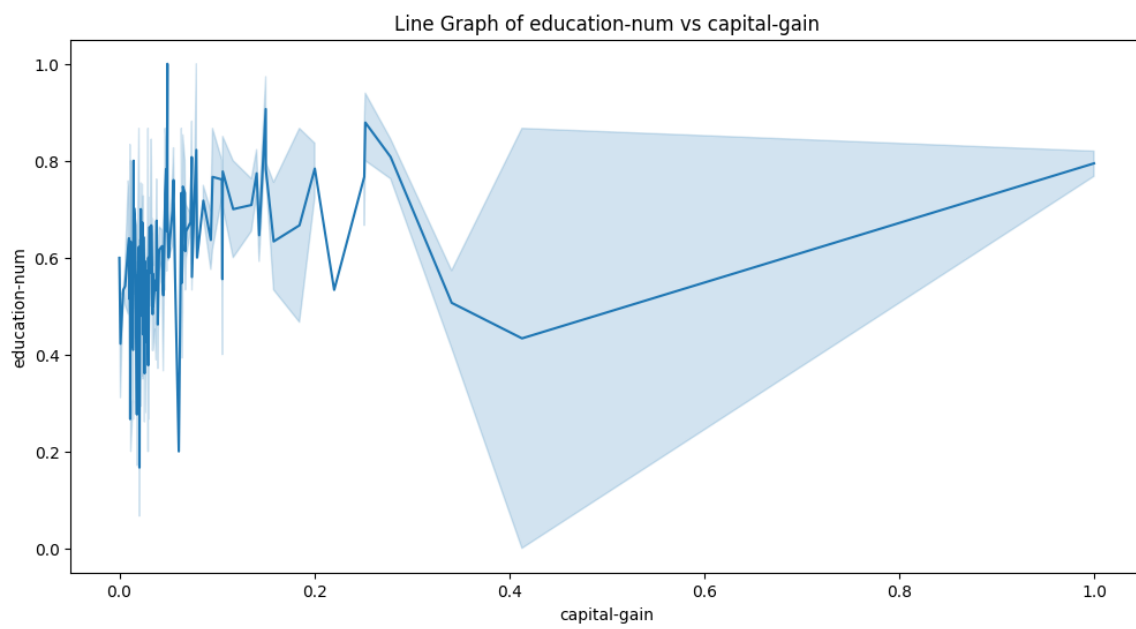
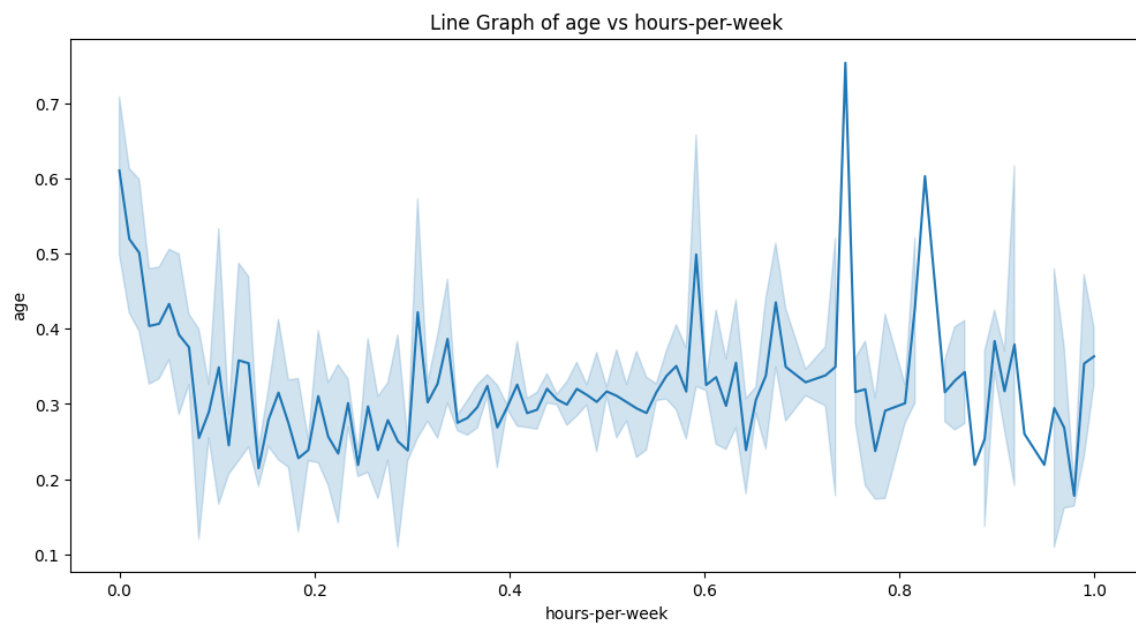


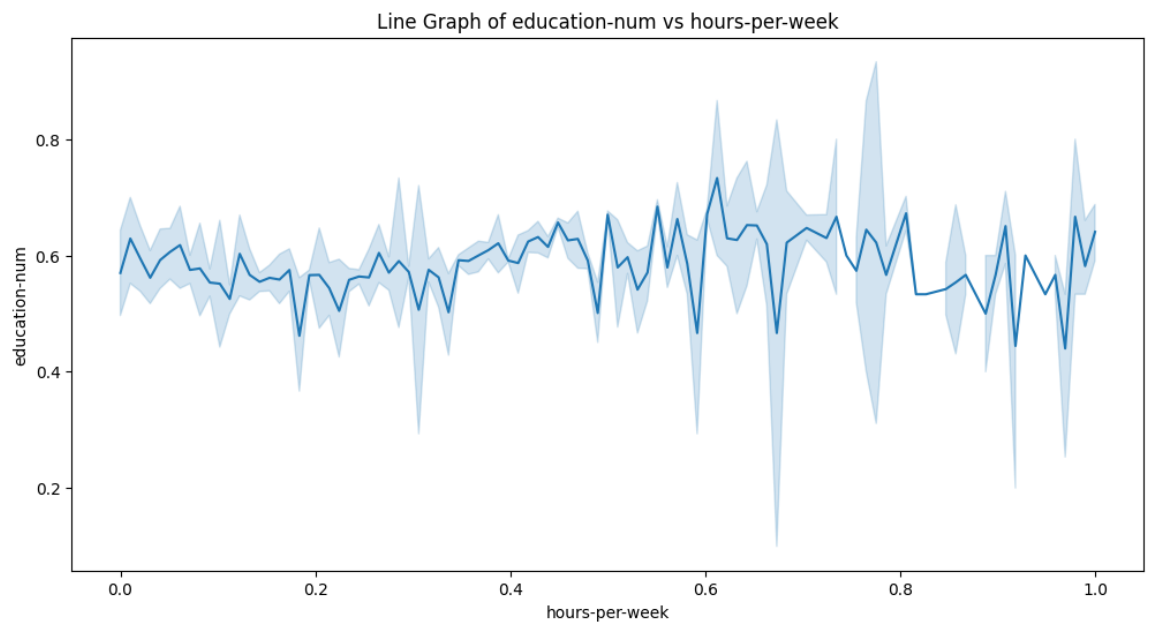
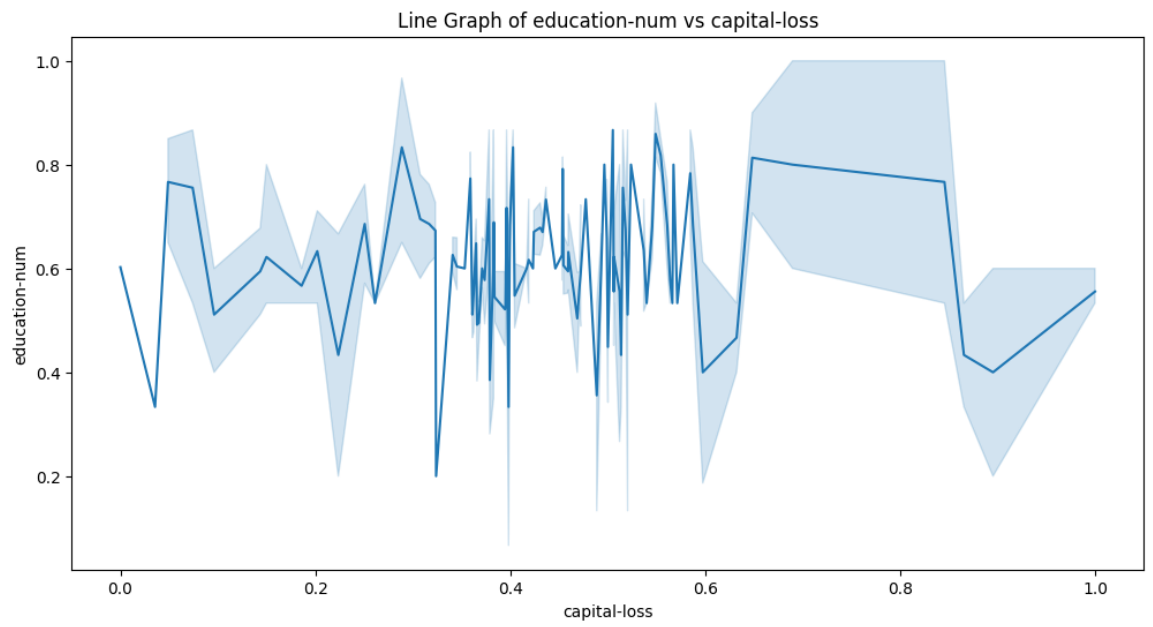




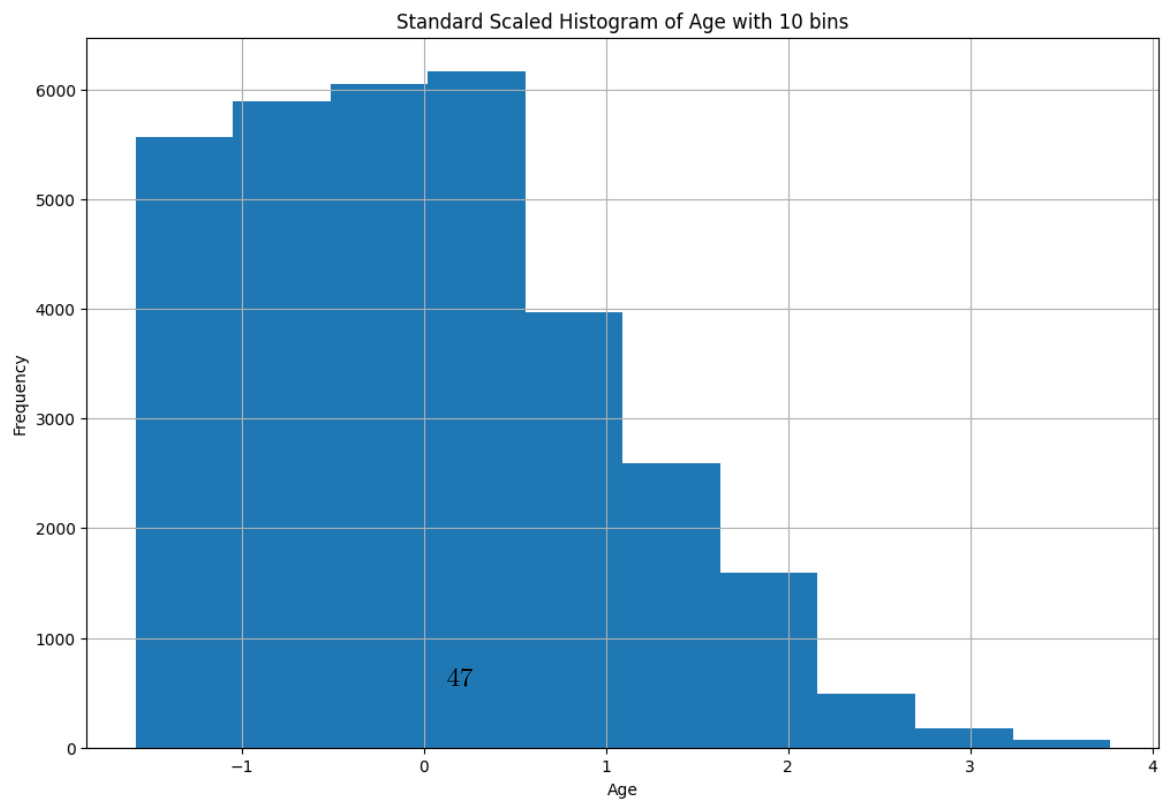
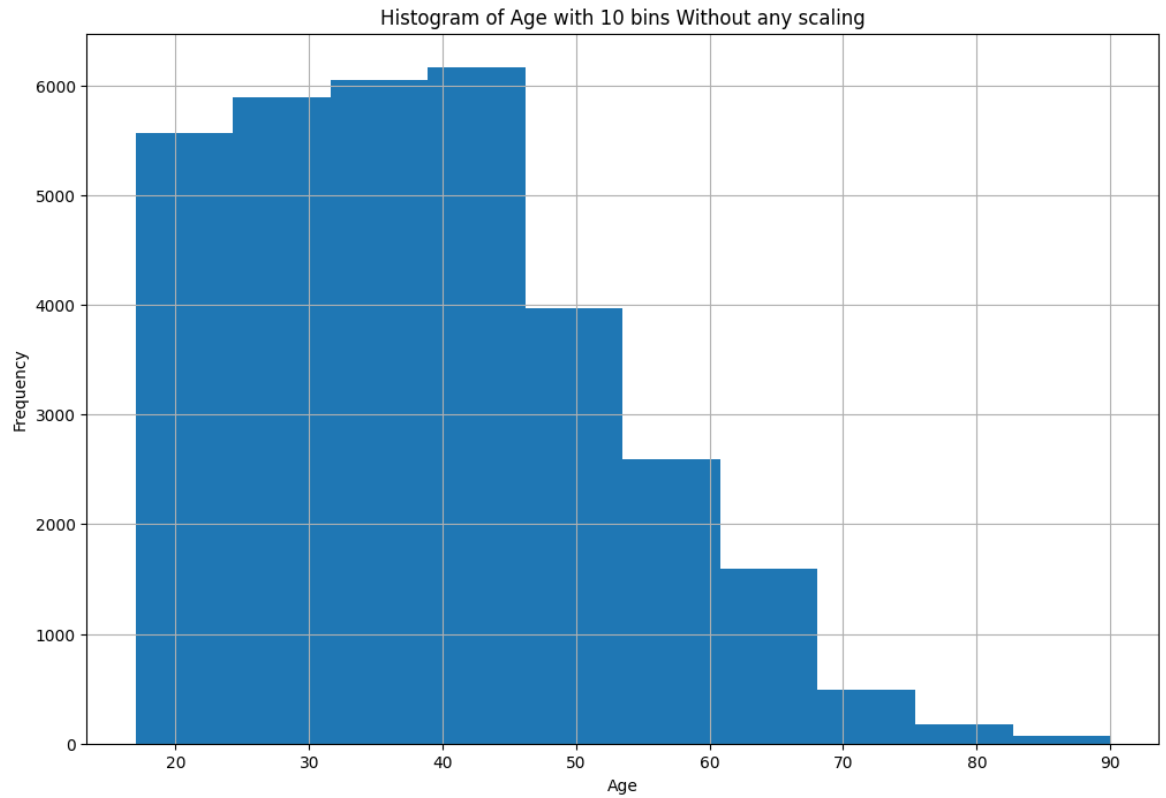
Min-Max Scaled Line Graphs

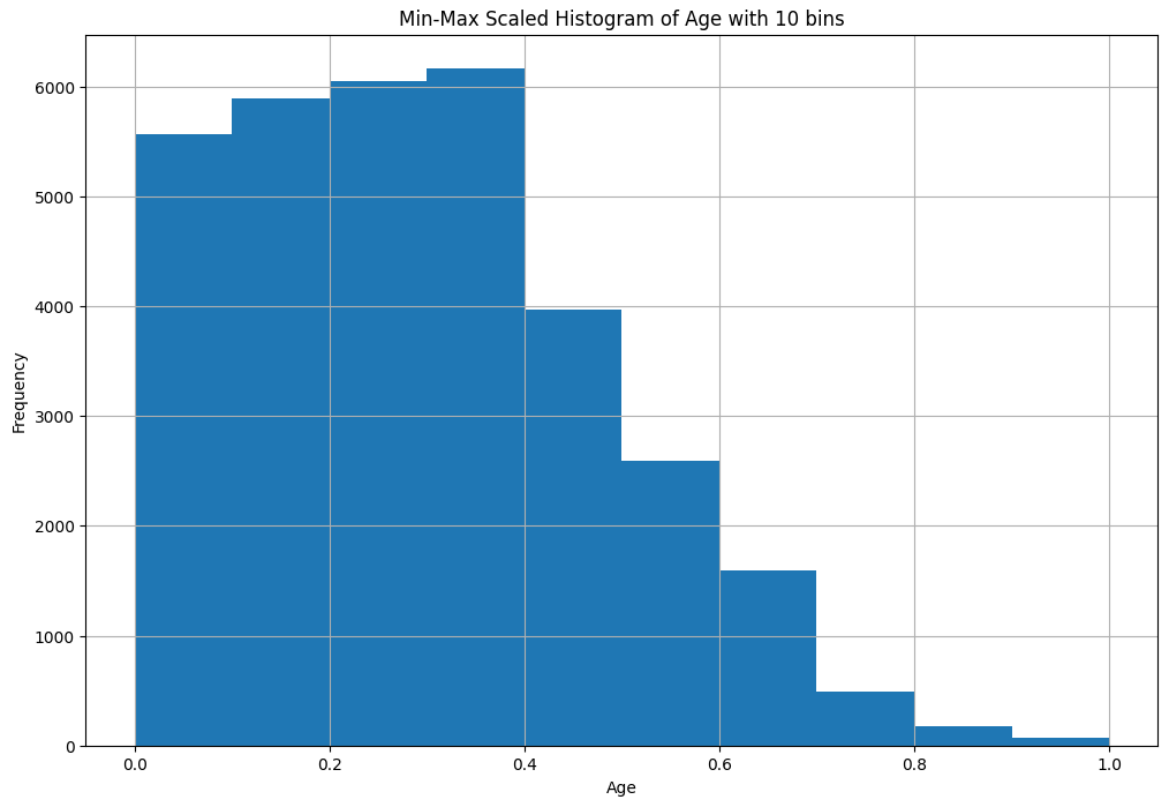






3.13





The above histogram are of the 'age' numerical attribute of the Adult dataset. Through observation it can be concluded that the normal, min-max scaled and standard scaled histograms of the same numerical attribute are of the same shape.

3.14

Before normalization: The heat map shows the raw L2 distances between numeric attributes. Attributes with larger numerical ranges (e.g., 'capital-gain', 'capital-loss') tend to have larger distances compared to attributes with smaller ranges (e.g., 'age', 'education-num'). This can skew the interpretation of distances due to the differing scales.

After Min-Max normalization: After Min-Max scaling, attributes are scaled to a fixed range. The heat map shows distances normalized within the $[0, 1]$ range. It preserves the shape of the distances but compresses them into a uniform scale.

After Z-score normalization: After standard scaling, the heat map shows distances that are more balanced across all attributes. It emphasizes rel-

ative distances based on their distribution rather than their absolute magnitude.

3.15

This database was organised in .name and .data filetypes. The .data filetype contained the data without attribute names while the .name file contained the names of the attributes and some additional information. Data was arranged in columns and rows, the names provided sufficient description for ease of understanding for the reader.

In some use cases the data retrieval from the file took more time than usual i.e when making heatmaps for the numerical attributes but in general the speed of data retrieval was quick.

The data contained 32561 rows and 15 columns which described attributes like age, education level, relationship status etc of adults from varying demographics.

3.16

```
from google.colab import drive
import pandas as pd

drive.mount('/content/drive')
path= '/content/drive/MyDrive/adult/adult.data'

names = ['age', 'workclass', 'fnlwgt', 'education', 'education-num',
         'marital-status', 'occupation', 'relationship', 'race', 'sex',
         'capital-gain', 'capital-loss', 'hours-per-week', 'native-country',
         'income']

df = pd.read_csv(path, header=None, names=names)

Q1 = df['age'].quantile(0.25)
Q3 = df['age'].quantile(0.75)

IQR = Q3 - Q1

lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

outliers = df[(df['age'] < lower_bound) | (df['age'] > upper_bound)]

print("Outliers detected using IQR method:")
print(outliers)
```

Using the above python code, the outliers were detected to be:
Age: 143

Final Weight(fnlwgt): 992
Education Number(education-num): 1198
Capital Gain: 2712
Capital Loss: 1519
Hours Per Week: 9008

3.17

Age:

Minimum: 17
Q1 (25th percentile): 28.0
Median (50th percentile): 37.0
Q3 (75th percentile): 47.0
Maximum: 78
Mode: 36
Mean: 38.4

Final Weight:

Minimum: 12285
Q1 (25th percentile): 116508.0
Median (50th percentile): 175935.0
Q3 (75th percentile): 228570.0
Maximum: 415847
Mode: 123011
Mean: 179631.3

Education Number:

Minimum: 5
Q1 (25th percentile): 9.0
Median (50th percentile): 10.0
Q3 (75th percentile): 13.0
Maximum: 16
Mode: 9
Mean: 10.3

Capital Gain:

Minimum: 0
Q1 (25th percentile): 0.0
Median (50th percentile): 0.0
Q3 (75th percentile): 0.0
Maximum: 0
Mode: 0
Mean: 0.0

Capital Loss:

Minimum: 0
Q1 (25th percentile): 0.0
Median (50th percentile): 0.0
Q3 (75th percentile): 0.0
Maximum: 0
Mode: 0
Mean: 0.0

Hours Per Week:

Minimum: 33
Q1 (25th percentile): 40.0
Median (50th percentile): 40.0
Q3 (75th percentile): 42.0
Maximum: 52
Mode: 40
Mean: 41.6