

```
In [1]: #P2). Repeat b), c), d), e) in Q2, by using Python on the HML.csv data set.
```

```
import Python libraries
import numpy as np
import scipy as sp
import pandas as pd
from IPython.display import display, HTML
```

```
In [2]: df=pd.read_csv('HML.csv')
```

```
In [3]: df.head()
```

```
Out[3]:
```

	Cost ID	Gender	Income	Age	Rentals	Avg Per Visit	Incidentals	Genre
0	1	M	45000	25	27	2.5	Yes	Action
1	2	F	54000	33	12	3.4	No	Drama
2	3	F	32000	20	42	1.6	No	Comedy
3	4	F	59000	70	16	4.2	Yes	Drama
4	5	M	37000	35	25	3.2	Yes	Action

```
In [4]: #Reduce min-max normalization to transform the values of the Income attribute onto the range [1-5].
```

```
# Normalization method: Min-max normalization [1,5]
df['Income_Minmax']=((df['Income']-df['Income'].min())/((df['Income'].max()-df['Income'].min()))+1)*1
display(HTML(df.head(50).to_html()))
```

	Cost ID	Gender	Income	Age	Rentals	Avg Per Visit	Incidentals	Genre	Income_Minmax
0	1	M	45000	25	27	2.5	Yes	Action	3.000000
1	2	F	54000	33	12	3.4	No	Drama	3.400001
2	3	F	32000	20	42	1.6	No	Comedy	1.450001
3	4	F	59000	70	16	4.2	Yes	Drama	3.636364
4	5	M	37000	35	25	3.2	Yes	Action	2.636364
5	6	M	18000	20	33	1.7	No	Action	1.727277
6	7	F	29000	45	19	3.8	No	Drama	2.272727
7	8	M	74000	25	31	2.4	Yes	Action	4.318182
8	9	M	38000	21	18	2.1	No	Comedy	2.618181
9	10	F	65000	40	21	3.3	No	Drama	3.909091
10	11	F	41000	22	48	2.3	Yes	Drama	2.818182
11	12	F	26000	22	29	2.9	Yes	Action	2.136364
12	13	M	83000	46	14	3.6	No	Comedy	4.727273
13	14	M	45000	36	24	2.7	No	Drama	3.000000
14	15	M	68000	30	36	2.7	Yes	Comedy	4.045455
15	16	M	17000	19	26	2.2	Yes	Action	1.727273
16	17	M	36000	35	28	3.5	Yes	Drama	2.590909
17	18	F	6000	15	39	1.8	Yes	Action	1.227273
18	19	F	24000	25	41	3.1	No	Comedy	2.045455
19	20	M	12000	15	23	2.2	Yes	Action	1.500000
20	21	F	47000	52	11	3.1	No	Drama	3.090909
21	22	M	25000	33	16	2.9	Yes	Drama	2.090909
22	23	F	2000	15	30	2.5	No	Comedy	1.045455
23	24	F	79000	35	22	3.8	Yes	Drama	4.545455
24	25	M	1000	18	25	1.4	Yes	Comedy	1.000000
25	26	F	56000	35	40	2.6	Yes	Action	3.500000
26	27	F	62000	47	32	3.6	No	Drama	3.727277
27	28	M	57000	52	22	4.1	No	Comedy	3.545455
28	29	F	15000	18	37	2.1	Yes	Action	1.636364
29	30	M	41000	25	17	1.4	Yes	Action	2.818182
30	31	F	49000	56	15	3.2	No	Comedy	3.181818
31	32	M	47000	30	21	3.1	Yes	Drama	3.090909
32	33	M	23000	25	28	2.7	No	Action	2.000000
33	34	F	29000	32	19	2.9	Yes	Action	2.272727
34	35	M	70000	29	43	4.6	Yes	Action	4.318182
35	36	F	29000	21	34	2.3	No	Comedy	2.272727
36	37	M	89000	46	12	1.2	No	Comedy	5.000000
37	38	M	41000	38	20	3.3	Yes	Drama	2.818182
38	39	F	69000	35	19	3.9	No	Comedy	4.045455
39	40	M	17000	19	30	1.8	No	Action	1.727273
40	41	F	50000	33	17	1.4	No	Drama	3.227273
41	42	M	32000	25	26	2.2	Yes	Action	2.400001
42	43	F	49000	28	48	3.3	Yes	Drama	3.181818
43	44	M	35000	24	24	1.7	No	Drama	2.545455
44	45	M	56000	38	30	3.5	Yes	Drama	3.500000
45	46	F	57000	43	9	1.1	No	Drama	3.545455
46	47	F	69000	35	22	2.8	Yes	Drama	4.090909
47	48	F	52000	47	14	1.6	No	Drama	3.318182
48	49	M	31000	25	42	3.4	Yes	Action	2.363636
49	50	M	24000	20	33	4.7	No	Action	2.045455

```
In [5]: #Reduce z-score normalization to standardize the values of the Age attribute.
```

```
from scipy.stats import zscore
df[['Age_zscore']] = zscore(df[['Age']])
display(HTML(df.head(50).to_html()))
```

	Cost ID	Gender	Income	Age	Rentals	Avg Per Visit	Incidentals	Genre	Income_Minmax	Age_z-score
0	1	M	45000	25	27	2.5	Yes	Action	3.000000	-0.552004
1	2	F	54000	33	12	3.4	No	Drama	3.400001	0.121216
2	3	F	32000	20	42	1.6	No	Comedy	2.400001	-0.973982
3	4	F	59000	70	16	4.2	Yes	Drama	3.636364	3.236782
4	5	M	37000	35	25	3.2	Yes	Action	2.636364	0.289670
5	6	M	18000	20	33	1.7	No	Action	1.727277	-0.873082
6	7	F	29000	45	19	3.8	No	Drama	2.272727	1.131345
7	8	M	74000	25	31	2.4	Yes	Action	4.318182	-0.552004
8	9	M	38000	21	18	2.1	No	Comedy	2.618181	-0.888914
9	10	F	65000	40	21	3.3	No	Drama	3.909091	0.710458
10	11	F	41000	22	48	2.3	Yes	Drama	2.818182	-0.804737
11	12	F	26000	22	29	2.9	Yes	Action	2.136364	-0.804737
12	13	M	83000	46	14	3.6	No	Comedy	4.727273	1.215523
13	14	M	45000	36	24	2.7	No	Drama	3.000000	0.373748
14	15	M	68000	30	36	2.7	Yes	Comedy	4.045455	-0.131317
15	16	M	17000	19	26	2.2	Yes	Action	1.727273	-1.057269
16	17	M	36000	35	28	3.5	Yes	Drama	2.590909	0.289670
17	18	F	6000	15	39	1.8	Yes	Action	1.227273	-1.309801
18	19	F	24000	25	41	3.1	No	Comedy	2.045455	-0.552004
19	20	M	12000	15	23	2.2	Yes	Action	1.500000	-1.309801
20	21	F	47000	52	11	3.1	No	Drama	3.090909	1.705987
21	22	M	25000	33	16	2.9	Yes	Drama	2.090909	0.121216
22	23	F	2000	15	30	2.5	No	Comedy	1.045455	-1.393979
23	24	F	79000	35	22	3.8	Yes	Drama	4.545455	0.289670
24	25	M	1000	18	25	1.4	Yes	Comedy	1.000000	-1.393981
25	26	F	56000	35	40	2.6	Yes	Action	3.500000	0.289670
26	27	F	62000	47	32	3.6	No	Drama	3.727277	1.299700
27	28	M	57000	52	22	4.1	No	Comedy	3.545455	1.705987
28	29	F	15000	18	37	2.1	Yes	Action	1.636364	-1.141446
29	30	M	41000	25	17	1.4	Yes	Action	2.818182	-0.552004
30	31	F	49000	56	15	3.2	No	Comedy	3.181818	2.057297
31	32	M	47000	30	21	3.1	Yes	Drama	3.090909	-0.131317
32	33	M	23000	25	28	2.7	No	Action	2.000000	-0.552004
33	34	F	29000	32	19	2.9	Yes	Action	2.272727	0.037038
34	35	M	74000	29	43	4.6	Yes	Action	4.318182	-0.215484
35	36	F	29000	21	34	2.3	No	Comedy	2.272727	-0.888914
36	37	M	89000	46	12	1.2	No	Comedy	5.000000	1.215523
37	38	M	41000	38	20	3.3	Yes	Drama	2.818182	0.542103
38	39	F	69000	35	19	3.9	No	Comedy	4.045455	0.289670
39	40	M	17000	19	30	1.8	No	Action	1.727273	-1.057269
40	41	F	50000	33	17	1.4	No	Drama	3.227273	0.121216
41	42	M	32000	25	26	2.2	Yes	Action	2.400001	-0.552004
42	43	F	49000	28	48	3.3	Yes	Drama	3.181818	-0.299672
43	44	M	35000	24	24	1.7	No	Drama	2.545455	-0.636382
44	45	M	56000	38	30	3.5	Yes	Drama	3.500000	0.542103
45	46	F	57000	43	9	1.1	No	Drama	3.545455	0.962990
46	47	F	69000	35	22	2.8	Yes	Drama	4.090909	0.289670
47	48	F	52000	47	14	1.6	No	Drama	3.318182	1.299700
48	49	M	31000	25	42	3.4	Yes	Action	2.363636	-0.552004
49	50	M	24000	20	33	4.7	No	Action	2.045455	-0.973982

```
In [6]: # Define the bins and labels
bins = [1, 25, 45, float('inf')]
labels = ['Young', 'MidAge', 'Old']
# Creating a new column called 'Age_Category' for discretized values
df['Age_discretization'] = pd.cut(df['Age'], bins=bins, labels=labels, right=False)
display(HTML(df.head(50).to_html()))
```

	Cost ID	Gender	Income	Age	Rentals	Avg Per Visit	Incidentals	Genre	Income_Minmax	Age_z-score	Age_discretization
0	1	M	45000	25	27	2.5	Yes	Action	3.000000	-0.552004	MidAge
1	2	F	54000	33	12	3.4	No	Drama	3.400001	0.121216	MidAge
2	3	F	32000	20	42	1.6	No	Comedy	2.400001	-0.973982	Young
3	4	F	59000	70	16	4.2	Yes	Drama	3.636364	3.236782	Old
4	5	M	37000	35	25	3.2	Yes	Action	2.636364	0.289670	MidAge
5	6	M	18000	20	33	1.7	No	Action	1.727277	-0.873082	Young
6	7	F	29000	45	19	3.8	No	Drama	2.272727	1.131345	Old
7	8	M	74000	25	31	2.4	Yes	Action	4.318182	-0.552004	MidAge
8	9	M	38000	21	18	2.1	No	Comedy	2.618181	-0.888914	MidAge
9	10	F	65000	40	21	3.3	No	Drama	3.909091	0.710458	MidAge
10	11	F	41000	22	48	2.3	Yes	Drama	2.818182	-0.804737	MidAge
11	12	F	26000	22	29	2.9	Yes	Action	2.136364	-0.804737	MidAge
12	13	M	83000	46	14	3.6	No	Comedy	4.727273	1.215523	Old
13	14	M	45000	36	24	2.7	No	Drama	3.000000	0.373748	MidAge
14	15	M	68000	30	36	2.7	Yes	Comedy	4.045455	-0.131317	MidAge
15	16	M	17000	19	26	2.2	Yes	Action	1.727273	-1.057269	Young
16	17	M	36000	35	28	3.5	Yes	Drama	2.590909	0.289670	MidAge
17	18	F	6000	15	39	1.8	Yes	Action	1.227273	-1.309801	Young
18	19	F	24000	25	41	3.1	No	Comedy	2.045455	-0.552004	Young
19	20	M	12000	15	23	2.2	Yes	Action	1.500000	-1.309801	Old
20	21	F	47000	52	11	3.1	No	Drama	3.090909	1.705987	Old
21	22	M	25000	33	16	2.9	Yes	Drama	2.090909	0.121216	MidAge
22	23	F	2000	15	30	2.5	No	Comedy	1.045455	-1.393979	Young
23	24	F	79000	35	22	3.8	Yes	Drama	4.545455	0.289670	MidAge
24	25	M	1000	18	25	1.4	Yes	Comedy	1.000000	-1.393981	Young
25	26	F	56000	35</							