

DATA MINING

ASSESSMENT 01

ABHIRUPA MITRA
17BCE0437

QUESTION:

Consider the Sales Transaction Dataset given below in this link in CSV format.

Perform the following tasks:-

- a) Read and clean the data first
- b) Try removing all the rows from the dataset that contain missing values (0) for the product transaction per week. How many of the dataset are left? Find out what percent of the dataset is missing? Print the values.
- c) Use Min-Max Normalization to normalize the transaction data for all the weeks
- d) Use the appropriate plot to display the highest sold product (ID) in the first 3 weeks.

```
In [43]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 header = ['Product_code', 'w0', 'w1', 'w2', 'w3', 'w4', 'w5', 'w6', 'w7', 'w8', 'w9', 'w10', 'w11', 'w12', 'w13', 'w14', 'w15', 'w16', 'w17', 'w18', 'w19', 'w20', 'w21', 'w22', 'w23', 'w24', 'w25', 'w26', 'w27', 'w28', 'w29', 'w30', 'w31', 'w32', 'w33', 'w34', 'w35', 'w36', 'w37', 'w38', 'w39', 'w40', 'w41', 'w42', 'w43', 'w44', 'w45', 'w46', 'w47', 'w48', 'w49', 'w50', 'w51', 'min', 'max']
5 df = pd.read_csv("Sales_Transactions_Dataset_Weekly.csv", header=None, names=header)
```

```
In [44]: 1 print(df.head(10))
2 row=df.shape[0]
```

	Product_code	w0	w1	w2	w3	w4	w5	w6	w7	w8	...	w44	w45	w46	w47	\
0	Product_Code	W0	W1	W2	W3	W4	W5	W6	W7	W8	...	W44	W45	W46	W47	
1	P1	11	12	10	8	13	12	14	21	6	...	8	10	12	3	
2	P2	7	6	3	2	7	1	6	3	3	...	5	1	1	4	
3	P3	7	11	8	9	10	8	7	13	12	...	5	5	7	8	
4	P4	12	8	13	5	9	6	9	13	13	...	3	4	6	8	
5	P5	8	5	13	11	6	7	9	14	9	...	7	12	6	6	
6	P6	3	3	2	7	6	3	8	6	6	...	4	3	6	5	
7	P7	4	8	3	7	8	7	2	3	10	...	3	6	2	6	
8	P8	8	6	10	9	6	8	7	5	10	...	4	8	8	6	
9	P9	14	9	10	7	11	15	12	7	13	...	13	3	7	7	

	w48	w49	w50	w51	min	max
0	W48	W49	W50	W51	MIN	MAX
1	7	6	5	10	3	21
2	5	1	6	0	0	10
3	14	8	8	7	3	14
4	14	8	7	8	2	19
5	5	11	8	9	3	18
6	3	3	10	6	0	11
7	2	4	2	1	0	10
8	7	4	9	9	3	15
9	10	12	7	13	3	18

[10 rows x 55 columns]

```
In [45]: 1 df.describe()
2 # df.dtypes
```

Out[45]:

	Product_code	w0	w1	w2	w3	w4	w5	w6	w7	w8	...	w44	w45	w46	w47	w48	w49	w50	w51	min	max	
count		812	812	812	812	812	812	812	812	812	...	812	812	812	812	812	812	812	812	812	812	
unique		812	51	50	55	54	55	54	55	54	...	47	47	51	48	49	46	45	43	26	63	
top	P436	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	2	
freq		1	256	259	258	260	247	250	248	262	251	...	205	190	184	204	176	178	167	154	445	100

4 rows x 55 columns

```
In [46]: 1 df = df.iloc[1:]
2 df[['w0', 'w1', 'w2', 'w3', 'w4', 'w5', 'w6', 'w7', 'w8', 'w9', 'w10', 'w11', 'w12', 'w13', 'w14', 'w15', 'w16', 'w17', 'w18', 'w19', 'w20']
3 print(df)
```

	Product_code	w0	w1	w2	w3	w4	w5	w6	w7	w8	...	w44	w45	w46	w47	\
1	P1	11	12	10	8	13	12	14	21	6	...	8	10	12	3	
2	P2	7	6	3	2	7	1	6	3	3	...	5	1	1	4	
3	P3	7	11	8	9	10	8	7	13	12	...	5	5	7	8	
4	P4	12	8	13	5	9	6	9	13	13	...	3	4	6	8	
5	P5	8	5	13	11	6	7	9	14	9	...	7	12	6	6	
6	P6	3	3	2	7	6	3	8	6	6	...	4	3	6	5	
7	P7	4	8	3	7	8	7	2	3	10	...	3	6	2	6	
8	P8	8	6	10	9	6	8	7	5	10	...	4	8	8	6	
9	P9	14	9	10	7	11	15	12	7	13	...	13	3	7	7	
10	P10	22	19	19	29	20	16	26	20	24	...	11	24	13	16	
11	P11	15	7	15	14	17	7	10	16	11	...	12	11	13	8	
12	P12	3	4	1	6	4	3	7	3	5	...	4	5	1	3	
13	P13	12	10	9	6	10	11	18	8	10	...	5	8	7	19	
14	P14	14	12	9	11	13	12	8	12	13	...	3	9	15	8	
15	P15	19	45	47	42	29	44	43	36	25	...	33	30	37	30	
16	P16	30	27	27	43	29	32	49	41	49	...	33	39	42	45	
17	P17	49	40	40	28	40	47	44	45	39	...	36	35	43	28	
18	P18	40	38	39	38	39	33	28	44	36	...	29	41	35	22	
19	P19	26	21	45	26	21	20	20	24	42	...	24	22	26	20	

```
In [50]: 1 rowN=df.shape[0]
2 # rowN
3 print("Total number of rows:", rowN)
4 print("Percentage of dataset missing:", (rowN/812)*100,"%")
```

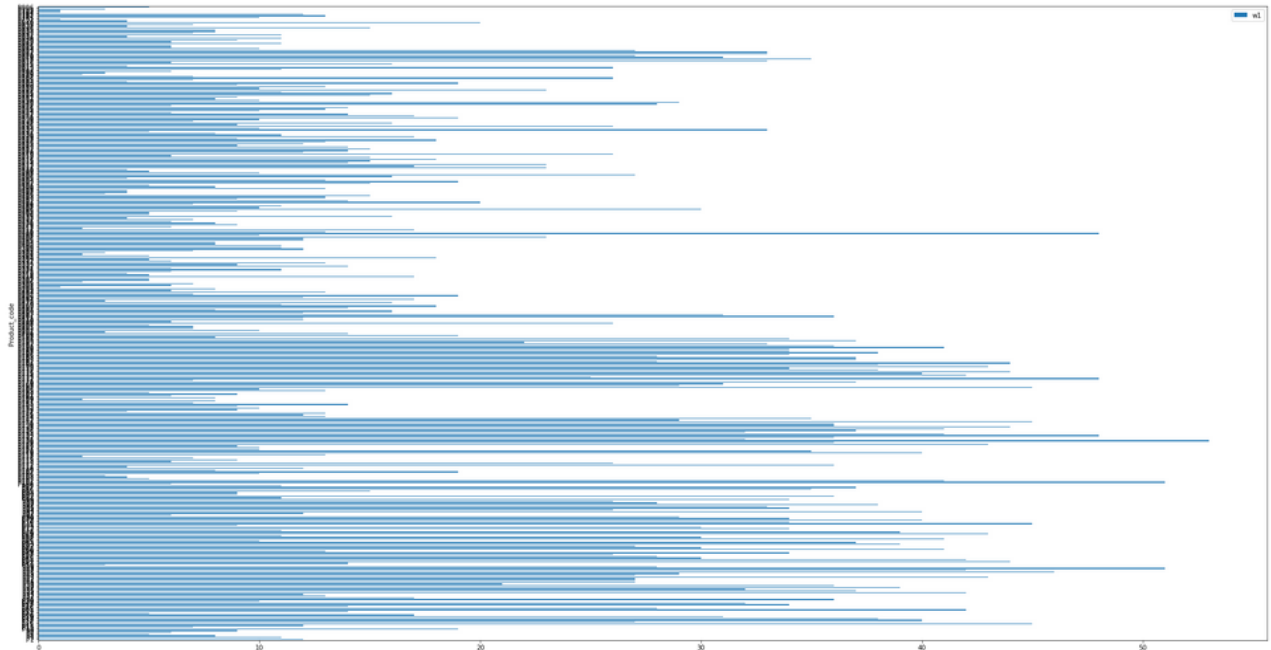
```

In [11]: 1 import matplotlib.pyplot as plt
2 # fig, axes = plt.subplots(3, 4, figsize=(14, 8))
3 df.plot(x="Product_code", y=["w1"], kind="barh",figsize=(34, 18))
4 # df.plot.scatter(x="w0", y="Product_code");
5 m0=result[1].argmax()
6 max0=result.iloc[[m0]]['Product_code']
7 print("Product Code of item sold the most", max0)
8

```

/home/abhirupa/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:5: FutureWarning: 'argmax' is deprecated, use 'idxmax' instead. The behavior of 'argmax' will be corrected to return the positional maximum in the future. Use 'series.values.argmax' to get the position of the maximum now.

Product Code of item sold the most 115 P130
Name: Product_code, dtype: object



```

In [12]: 1 import matplotlib.pyplot as plt
2 # fig, axes = plt.subplots(3, 4, figsize=(14, 8))
3 df.plot(x="Product_code", y=["w2"], kind="barh",figsize=(34, 18))
4 # df.plot.scatter(x="w0", y="Product_code");
5 m0=result[2].argmax()
6 max0=result.iloc[[m0]]['Product_code']
7 print("Product Code of item sold the most", max0)
8

```

/home/abhirupa/anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:5: FutureWarning: 'argmax' is deprecated, use 'idxmax' instead. The behavior of 'argmax' will be corrected to return the positional maximum in the future. Use 'series.values.argmax' to get the position of the maximum now.

Product Code of item sold the most 83 P92
Name: Product_code, dtype: object

