

In [2]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

In [3]:

```
df=pd.read_csv('https://d2beiqkhq929f0.cloudfront.net/public_assets/assets/000/001/125/original/aerofit_treadmill.csv?163999274')
```

In [4]:

```
df
```

Out[4]:

	Product	Age	Gender	Education	MaritalStatus	Usage	Fitness	Income	Miles
0	KP281	18	Male	14	Single	3	4	29562	112
1	KP281	19	Male	15	Single	2	3	31836	75
2	KP281	19	Female	14	Partnered	4	3	30699	66
3	KP281	19	Male	12	Single	3	3	32973	85
4	KP281	20	Male	13	Partnered	4	2	35247	47
...	...	...	...	...	...	...	...	...	...
175	KP781	40	Male	21	Single	6	5	83416	200
176	KP781	42	Male	18	Single	5	4	89641	200
177	KP781	45	Male	16	Single	5	5	90886	160
178	KP781	47	Male	18	Partnered	4	5	104581	120
179	KP781	48	Male	18	Partnered	4	5	95508	180

180 rows × 9 columns

In [7]:

```
df.shape
```

Out[7]:

(180, 9)

In [8]:

```
df.describe(include='all')
```

Out[8]:

	Product	Age	Gender	Education	MaritalStatus	Usage	Fitness	Income	Miles
count	180	180.000000	180	180.000000	180	180.000000	180.000000	180.000000	180.000000
unique	3	NaN	2	NaN	2	NaN	NaN	NaN	NaN
top	KP281	NaN	Male	NaN	Partnered	NaN	NaN	NaN	NaN
freq	80	NaN	104	NaN	107	NaN	NaN	NaN	NaN
mean	NaN	28.788889	NaN	15.572222	NaN	3.455556	3.311111	53719.577778	103.194444
std	NaN	6.943498	NaN	1.617055	NaN	1.084797	0.958869	16506.684226	51.863605
min	NaN	18.000000	NaN	12.000000	NaN	2.000000	1.000000	29562.000000	21.000000
25%	NaN	24.000000	NaN	14.000000	NaN	3.000000	3.000000	44058.750000	66.000000
50%	NaN	26.000000	NaN	16.000000	NaN	3.000000	3.000000	50596.500000	94.000000
75%	NaN	33.000000	NaN	16.000000	NaN	4.000000	4.000000	58668.000000	114.750000
max	NaN	50.000000	NaN	21.000000	NaN	7.000000	5.000000	104581.000000	360.000000

In [9]:

```
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 180 entries, 0 to 179
Data columns (total 9 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   Product         180 non-null    object
 1   Age             180 non-null    int64
 2   Gender          180 non-null    object
 3   Education       180 non-null    int64
 4   MaritalStatus   180 non-null    object
 5   Usage           180 non-null    int64
 6   Fitness         180 non-null    int64
 7   Income          180 non-null    int64
 8   Miles           180 non-null    int64
dtypes: int64(6), object(3)
memory usage: 12.8+ KB
```

In [10]:

```
df_male=df.loc[df['Gender']=='Male']
df_male
```

Out[10]:

	Product	Age	Gender	Education	MaritalStatus	Usage	Fitness	Income	Miles
0	KP281	18	Male	14	Single	3	4	29562	112
1	KP281	19	Male	15	Single	2	3	31836	75
3	KP281	19	Male	12	Single	3	3	32973	85
4	KP281	20	Male	13	Partnered	4	2	35247	47
7	KP281	21	Male	13	Single	3	3	32973	85
...	...	...	...	...	...	...	...	...	...
175	KP781	40	Male	21	Single	6	5	83416	200
176	KP781	42	Male	18	Single	5	4	89641	200
177	KP781	45	Male	16	Single	5	5	90886	160
178	KP781	47	Male	18	Partnered	4	5	104581	120
179	KP781	48	Male	18	Partnered	4	5	95508	180

104 rows × 9 columns

In [11]:

```
df_female=df.loc[df['Gender']=='Female']
```

In [12]:

```
df_female
```

Out[12]:

	Product	Age	Gender	Education	MaritalStatus	Usage	Fitness	Income	Miles
2	KP281	19	Female	14	Partnered	4	3	30699	66
5	KP281	20	Female	14	Partnered	3	3	32973	66
6	KP281	21	Female	14	Partnered	3	3	35247	75
9	KP281	21	Female	15	Partnered	2	3	37521	85
11	KP281	22	Female	14	Partnered	3	2	35247	66
...	...	...	...	...	...	...	...	...	...
152	KP781	25	Female	18	Partnered	5	5	61006	200
157	KP781	26	Female	21	Single	4	3	69721	100
162	KP781	28	Female	18	Partnered	6	5	92131	180
167	KP781	30	Female	16	Partnered	6	5	90886	280
171	KP781	33	Female	18	Partnered	4	5	95866	200

76 rows × 9 columns

In [13]:

```
df_married_single=df.loc[df['MaritalStatus']=='Single']
df_married_single
```

Out[13]:

	Product	Age	Gender	Education	MaritalStatus	Usage	Fitness	Income	Miles
0	KP281	18	Male	14	Single	3	4	29562	112
1	KP281	19	Male	15	Single	2	3	31836	75
3	KP281	19	Male	12	Single	3	3	32973	85
7	KP281	21	Male	13	Single	3	3	32973	85
8	KP281	21	Male	15	Single	5	4	35247	141
...	...	...	...	...	...	...	...	...	...
165	KP781	29	Male	18	Single	5	5	52290	180
172	KP781	34	Male	16	Single	5	5	92131	150
175	KP781	40	Male	21	Single	6	5	83416	200
176	KP781	42	Male	18	Single	5	4	89641	200
177	KP781	45	Male	16	Single	5	5	90886	160

73 rows × 9 columns

In [14]:

```
df_married_partner=df.loc[df['MaritalStatus']=='Partnered']
df_married_partner
```

Out[14]:

	Product	Age	Gender	Education	MaritalStatus	Usage	Fitness	Income	Miles
2	KP281	19	Female	14	Partnered	4	3	30699	66
4	KP281	20	Male	13	Partnered	4	2	35247	47
5	KP281	20	Female	14	Partnered	3	3	32973	66
6	KP281	21	Female	14	Partnered	3	3	35247	75
9	KP281	21	Female	15	Partnered	2	3	37521	85
...	...	...	...	...	...	...	...	...	...
171	KP781	33	Female	18	Partnered	4	5	95866	200
173	KP781	35	Male	16	Partnered	4	5	92131	360
174	KP781	38	Male	18	Partnered	5	5	104581	150
178	KP781	47	Male	18	Partnered	4	5	104581	120
179	KP781	48	Male	18	Partnered	4	5	95508	180

107 rows × 9 columns

In [15]:

```
df['Product'].value_counts()
```

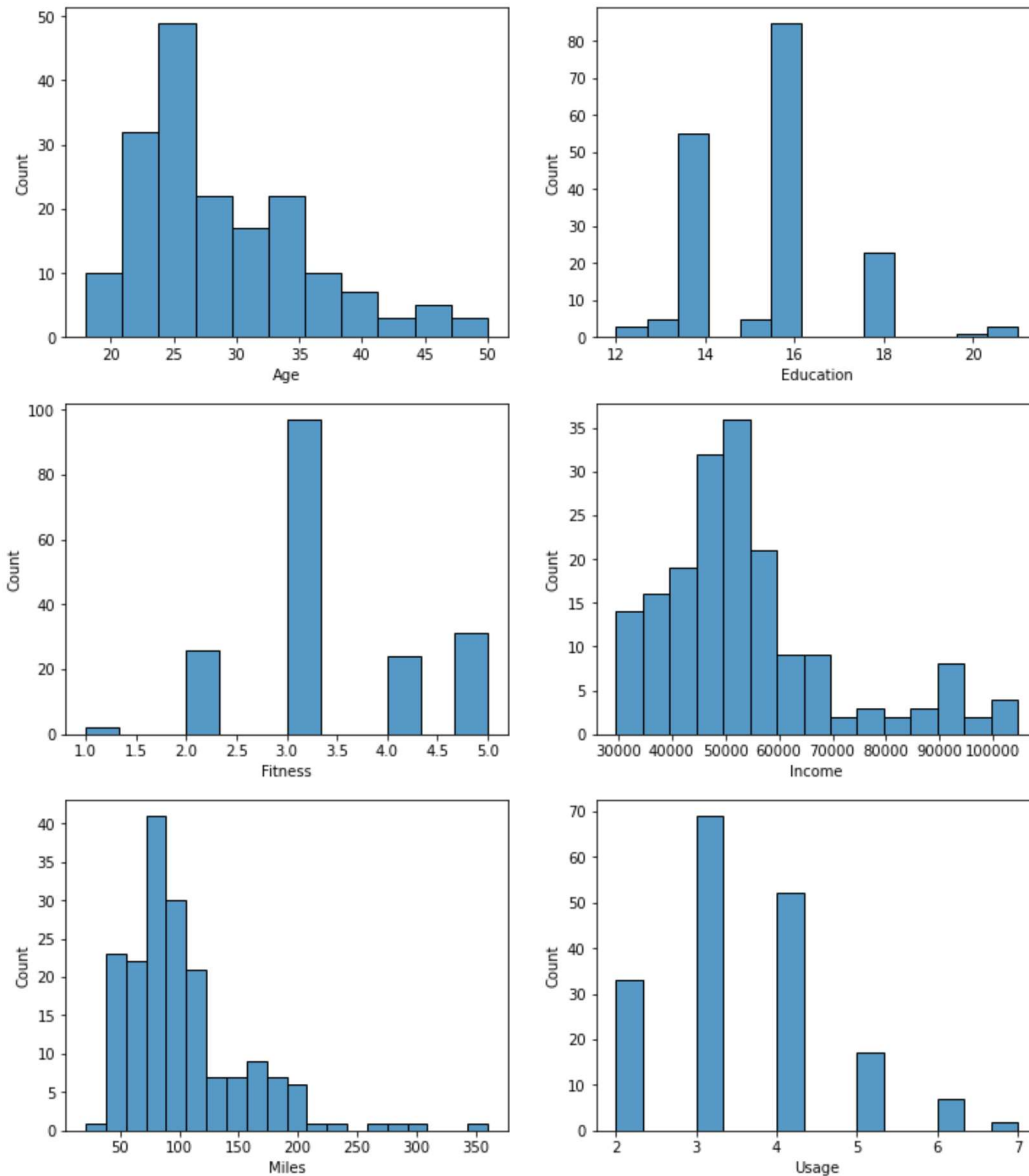
Out[15]:

KP281 80  
KP481 60  
KP781 40  
Name: Product, dtype: int64

# Univariate Analysis using histplot

In [28]:

```
fig, axis = plt.subplots(nrows=3, ncols=2, figsize=(12, 10))
fig.subplots_adjust(top=1.2)
sns.histplot(data=df, x="Age", ax=axis[0,0])
sns.histplot(data=df, x="Education", ax=axis[0,1])
sns.histplot(data=df, x="Fitness", ax=axis[1,0])
sns.histplot(data=df, x="Income", ax=axis[1,1])
sns.histplot(data=df, x="Miles", ax=axis[2,0])
sns.histplot(data=df, x="Usage", ax=axis[2,1])
plt.show()
```



## Outliers detection using boxplot

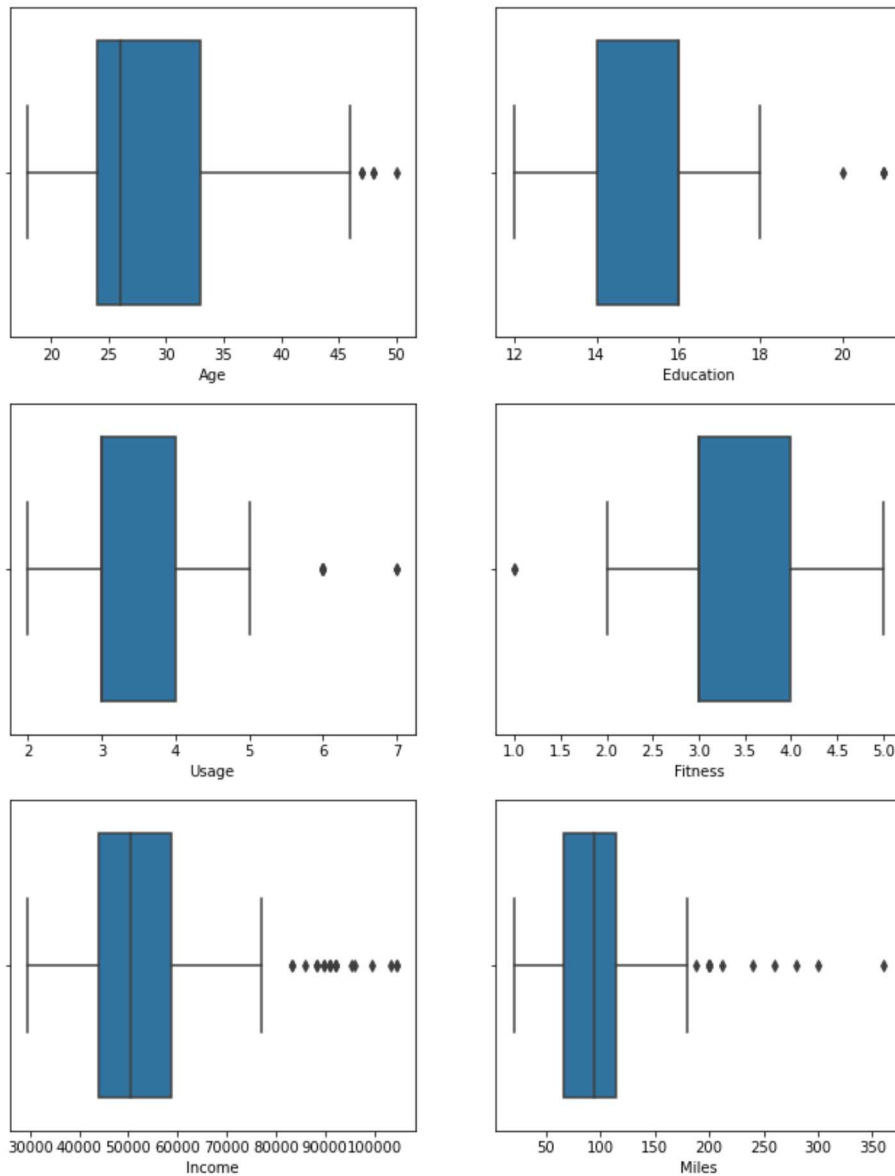
In [30]:

```
fig, axis = plt.subplots(nrows=3, ncols=2, figsize=(11, 10))
fig.subplots_adjust(top=1.2)
```

```
sns.boxplot(data=df, x="Age", orient='h', ax=axis[0,0])
sns.boxplot(data=df, x="Education", orient='vh', ax=axis[0,1])
sns.boxplot(data=df, x="Usage", orient='h', ax=axis[1,0])
sns.boxplot(data=df, x="Fitness", orient='h', ax=axis[1,1])
sns.boxplot(data=df, x="Income", orient='h', ax=axis[2,0])
sns.boxplot(data=df, x="Miles", orient='h', ax=axis[2,1])
plt.show()
```

C:\Anaconda3\lib\site-packages\seaborn\\_core.py:1319: UserWarning: Vertical orientation ignored with only `x` specified.

```
warnings.warn(single_var_warning.format("Vertical", "x"))
```

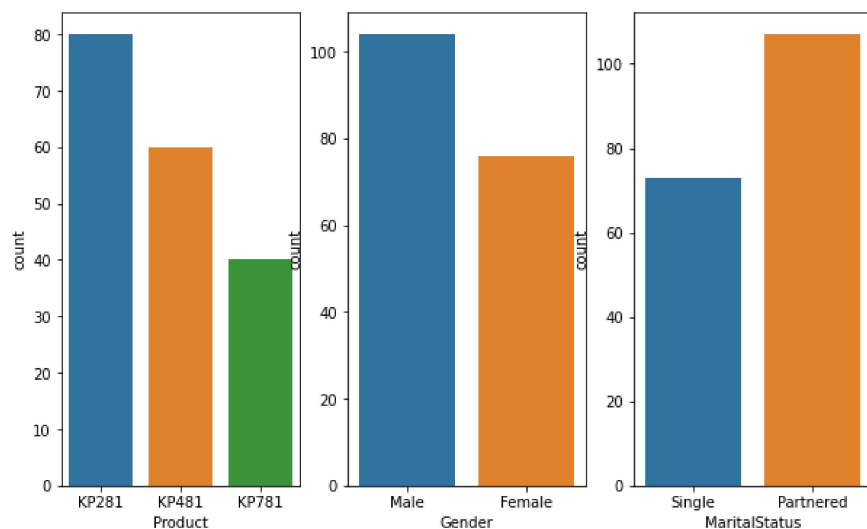


1. Income and Miles are having more outliers
2. Age, Education, Usage and Fitness are having less outliers

## CountPlots

In [35]:

```
fig, axs = plt.subplots(nrows=1, ncols=3, figsize=(10, 6))
sns.countplot(data=df, x='Product', ax=axs[0])
sns.countplot(data=df, x='Gender', ax=axs[1])
sns.countplot(data=df, x='MaritalStatus', ax=axs[2])
plt.show()
```

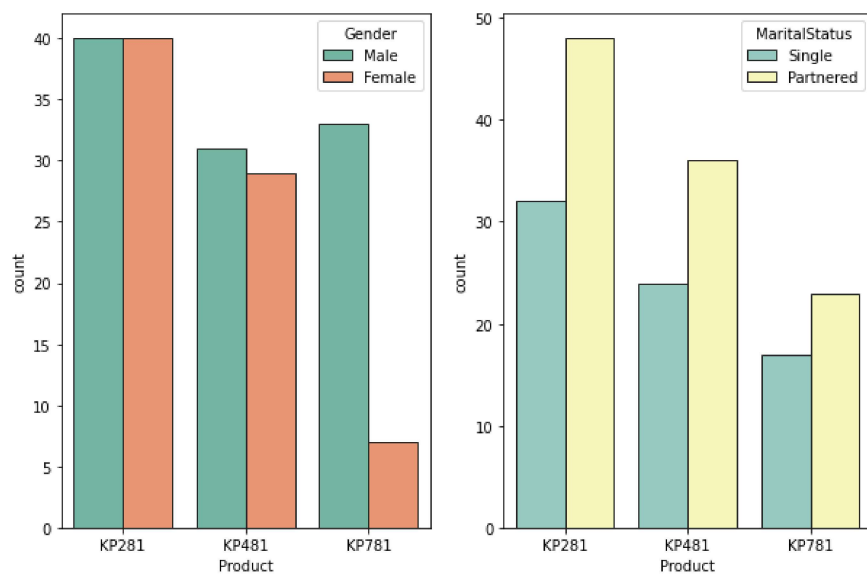


1. KP281 is highest product which is used
2. More males than females
3. Partneres is more than single

## Bivariate Analysis

In [37]:

```
fig, axs = plt.subplots(nrows=1, ncols=2, figsize=(10, 6.5))
sns.countplot(data=df, x='Product', hue='Gender', edgecolor="0.15", palette='Set2', ax=axs[0])
sns.countplot(data=df, x='Product', hue='MaritalStatus', edgecolor="0.15", palette='Set3', ax=axs[1])
plt.show()
```

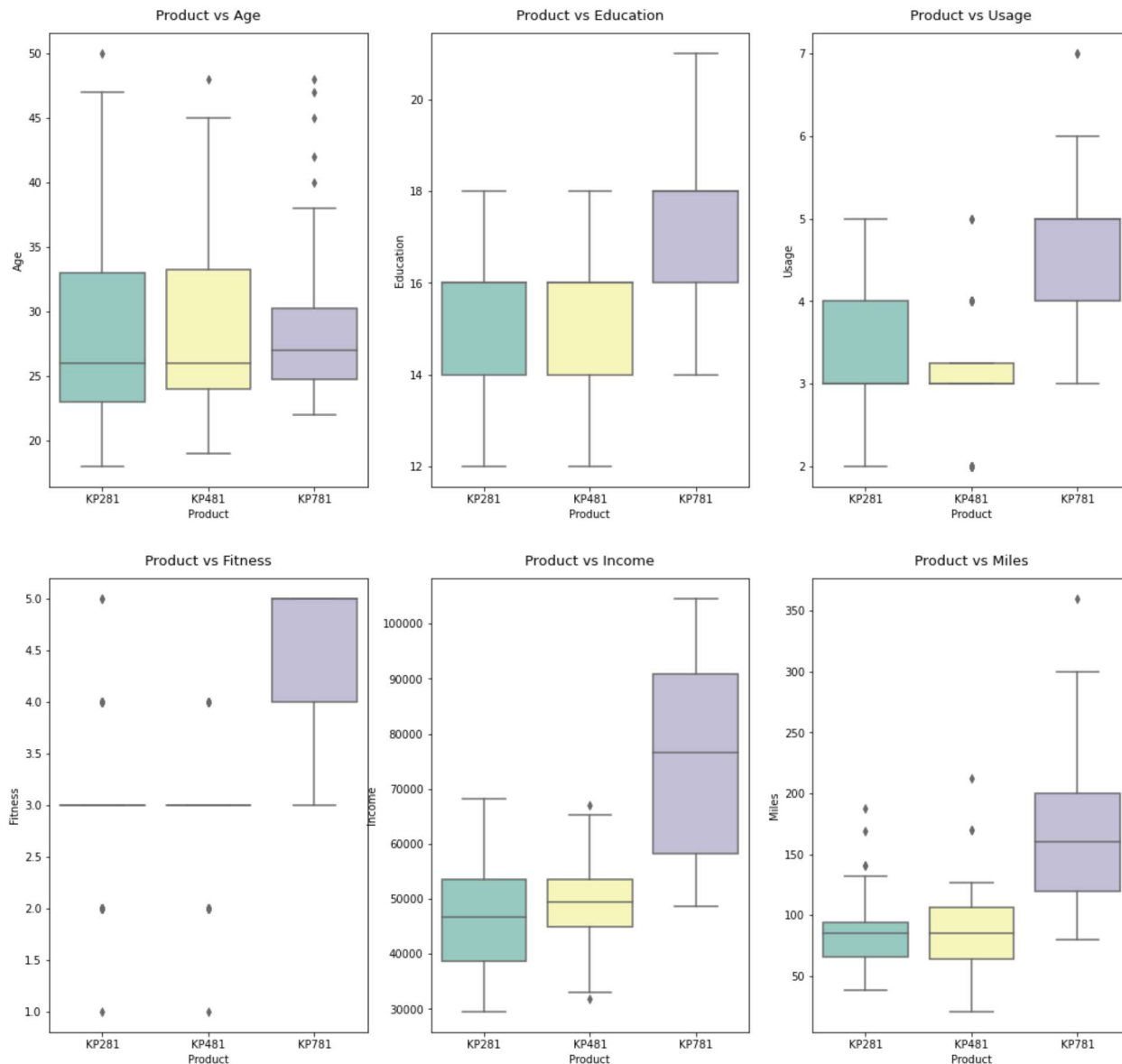


In [39]:

```

attrs = ['Age', 'Education', 'Usage', 'Fitness', 'Income', 'Miles']
fig, axs = plt.subplots(nrows=2, ncols=3, figsize=(18, 12))
fig.subplots_adjust(top=1.2)
count = 0
for i in range(2):
    for j in range(3):
        sns.boxplot(data=df, x='Product', y=attrs[count], ax=axs[i,j], palette='Set3')
        axs[i,j].set_title(f"Product vs {attrs[count]}", pad=12, fontsize=13)
        count += 1

```



## Marginal probability

In [40]:

```
df['Product'].value_counts(normalize=True)
```

Out[40]:

```

KP281    0.444444
KP481    0.333333
KP781    0.222222
Name: Product, dtype: float64

```

## Conditional probability

Probability of each product given gender

In [5]:

```
df1 = df[['Product', 'Gender', 'MaritalStatus']].melt()
df1.groupby(['variable', 'value'])['value'].count() / len(df)
```

Out[5]:

		value
variable	value	
Gender	Female	0.422222
	Male	0.577778
MaritalStatus	Partnered	0.594444
	Single	0.405556
Product	KP281	0.444444
	KP481	0.333333
	KP781	0.222222

In [7]:

```
def p_prod_given_gender(gender, print_marginal=False):
    if gender!= "Female" and gender != "Male":
        return "Invalid gender value."

    df1 = pd.crosstab(index=df['Gender'], columns=df['Product'])
    p_781 = df1['KP781'][gender] / df1.loc[gender].sum()
    p_481 = df1['KP481'][gender] / df1.loc[gender].sum()
    p_281 = df1['KP281'][gender] / df1.loc[gender].sum()

    if print_marginal:
        print(f"P(Male): {df1.loc['Male'].sum()/len(df):.2f}")
        print(f"P(Female): {df1.loc['Female'].sum()/len(df):.2f}\n")

        print(f"P(KP781/{gender}): {p_781:.2f}")
        print(f"P(KP481/{gender}): {p_481:.2f}")
        print(f"P(KP281/{gender}): {p_281:.2f}\n")
    p_prod_given_gender('Male', True)
    p_prod_given_gender('Female')
```

P(Male): 0.58  
P(Female): 0.42

P(KP781/Male): 0.32  
P(KP481/Male): 0.30  
P(KP281/Male): 0.38

P(KP781/Female): 0.09  
P(KP481/Female): 0.38  
P(KP281/Female): 0.53

Probability of each product given MaritalStatus



In [8]:

```
def p_prod_given_mstatus(status, print_marginal=False):
    if status!= "Single" and status!= "Partnered":
        return "Invalid marital status value."

    df1 = pd.crosstab(index=df['MaritalStatus'], columns=[df['Product']])
    p_781 = df1['KP781'][status] / df1.loc[status].sum()
    p_481 = df1['KP481'][status] / df1.loc[status].sum()
    p_281 = df1['KP281'][status] / df1.loc[status].sum()

    if print_marginal:
        print(f"P(Single): {df1.loc['Single'].sum()/len(df):.2f}")
        print(f"P(Partnered): {df1.loc['Partnered'].sum()/len(df):.2f}\n")

    print(f"P(KP781/{status}): {p_781:.2f}")
    print(f"P(KP481/{status}): {p_481:.2f}")
    print(f"P(KP281/{status}): {p_281:.2f}\n")
p_prod_given_mstatus('Single', True)
p_prod_given_mstatus('Partnered')
```

P(Single): 0.41  
P(Partnered): 0.59

P(KP781/Single): 0.23  
P(KP481/Single): 0.33  
P(KP281/Single): 0.44

P(KP781/Partnered): 0.21  
P(KP481/Partnered): 0.34  
P(KP281/Partnered): 0.45

In [ ]:

le with Fitness Level 3 or less are likely to purchase KP281 and KP481. And with Fitness Level 5 are likely to purchase KP781. le who use the treadmill more are more likely to purchase KP781. As, buying the treadmill is directly proportional to it's usage