

1. Importing Package

In [1]:

```
import pandas as pd
import seaborn as sns
import numpy as np
from matplotlib import pyplot as plt
%matplotlib inline
```

2. Loading the Dataset

In [2]:

```
df = pd.read_csv('/content/Churn_Modelling.csv')
df
```

Out[2]:

[illegible]

	RowN umbe r	Custo merI d	Surn ame	Credi tScor e	Geog raph y	Ge nd er	A g e	Te nur e	Bala nce	NumOf Product s	HasC rCar d	IsActive Membe r	Estimat edSalar y	Exi te d
9 9 9 5	9996	1560 6229	Obiji aku	771	Franc e	Ma le	3 9	5	0.00	2	1	0	96270.6 4	0
9 9 9 6	9997	1556 9892	John ston e	516	Franc e	Ma le	3 5	10	5736 9.61	1	1	1	101699. 77	0
9 9 9 7	9998	1558 4532	Liu	709	Franc e	Fe mal e	3 6	7	0.00	1	0	1	42085.5 8	1
9 9 9 8	9999	1568 2355	Sabb atini	772	Ger man y	Ma le	4 2	3	7507 5.31	2	1	0	92888.5 2	1
9 9 9 9	10000	1562 8319	Walk er	792	Franc e	Fe mal e	2 8	4	1301 42.7 9	1	1	0	38190.7 8	0

10000 rows × 14 columns

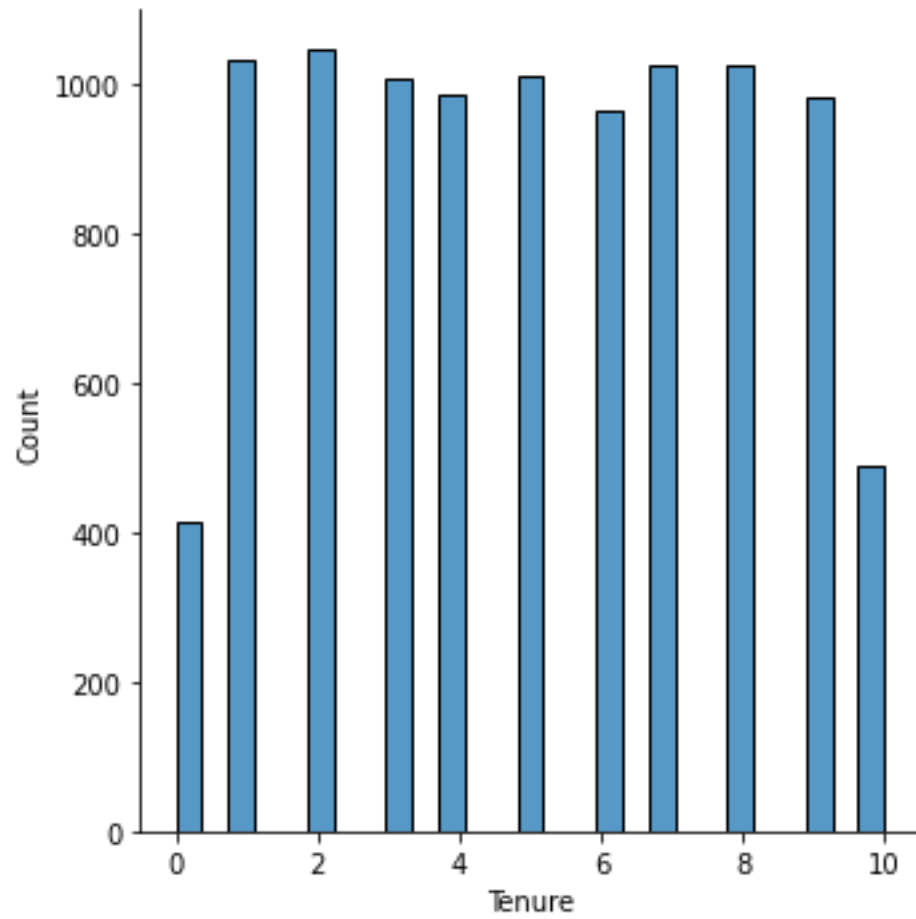
3. Visualizations

3.1 Univariate Analysis

In [3]:

```
sns.displot(df.Tenure)
```

Out[3]:

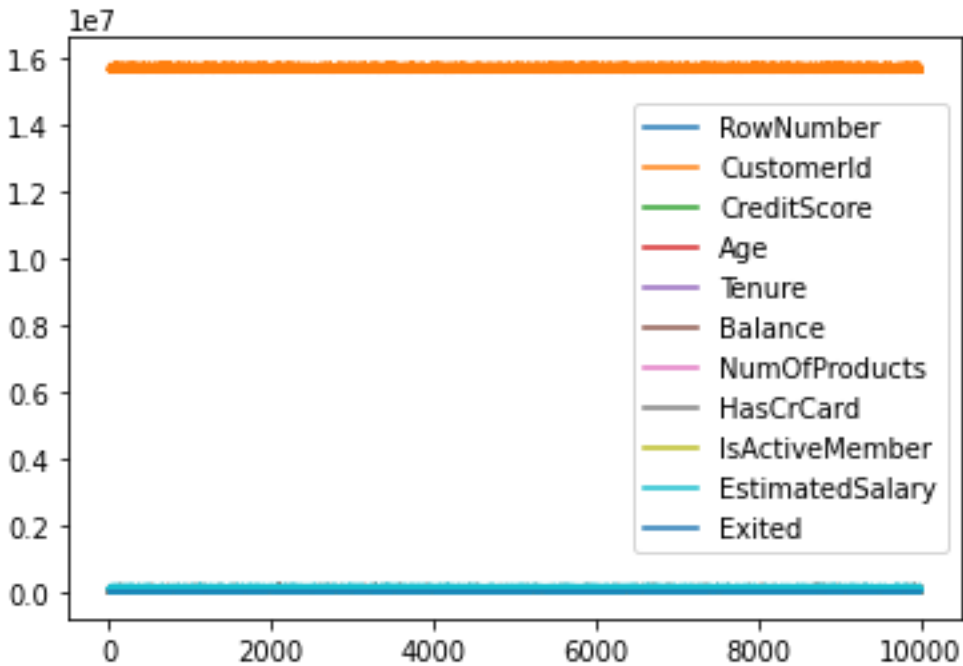


3.2 Bi-Variate Analysis

In [5]:

```
df.plot.line()
```

Out[5]:



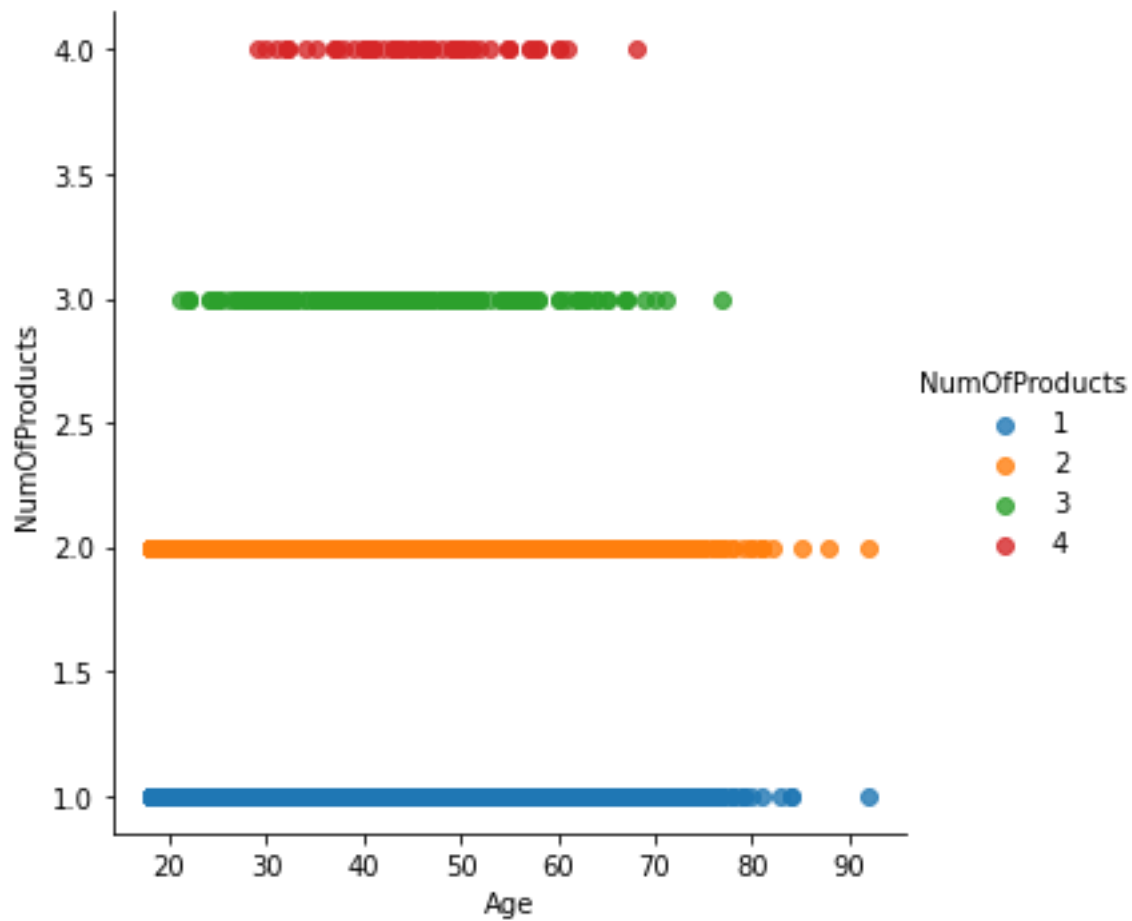
3.3 Multi-Variate Analysis

In [6]:

```
sns.lmplot("Age", "NumOfProducts", df, hue="NumOfProducts", fit_reg=False);
```

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variables as keyword args: x, y, data. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.

FutureWarning



4. Perform descriptive statistics on the dataset

In [7]:

```
df.describe()
```

Out[7]:

	RowN umber	Custo merId	CreditS core	Age	Tenure	Balance	NumOf Product s	HasCr Card	IsActive Member	Estimat edSalar y	Exited
co un t	10000 .0000 0	1.0000 00e+04	10000. 000000	10000. 000000	10000. 000000	10000.0 00000	10000.0 00000	10000. 00000	10000.0 00000	10000.0 00000	10000. 000000
m ea n	5000. 50000	1.5690 94e+07	650.52 8800	38.921 800	5.0128 00	76485.8 89288	1.53020 0	0.7055 0	0.51510 0	100090. 239881	0.2037 00

	RowN umber	Custo merId	CreditS core	Age	Tenure	Balance	NumOf Product s	HasCr Card	IsActive Member	Estimat edSalar y	Exited
std	2886.	7.1936	96.653	10.487	2.8921	62397.4	0.58165	0.4558	0.49979	57510.4	0.4027
	89568	19e+04	299	806	74	05202	4	4	7	92818	69
min	1.000	1.5565	350.00	18.000	0.0000	0.00000	1.00000	0.0000	0.00000	11.5800	0.0000
	00	70e+07	0000	000	00	0	0	0	0	00	00
25%	2500.	1.5628	584.00	32.000	3.0000	0.00000	1.00000	0.0000	0.00000	51002.1	0.0000
	75000	53e+07	0000	000	00	0	0	0	0	10000	00
50%	5000.	1.5690	652.00	37.000	5.0000	97198.5	1.00000	1.0000	1.00000	100193.	0.0000
	50000	74e+07	0000	000	00	40000	0	0	0	915000	00
75%	7500.	1.5753	718.00	44.000	7.0000	127644.	2.00000	1.0000	1.00000	149388.	0.0000
	25000	23e+07	0000	000	00	240000	0	0	0	247500	00
max	10000	1.5815	850.00	92.000	10.000	250898.	4.00000	1.0000	1.00000	199992.	1.0000
	.0000 0	69e+07	0000	000	000	090000	0	0	0	480000	00

5. Handle the Missing values

In [11]:

```
data = pd.read_csv("Churn_Modelling.csv")
pd.isnull (data["Age"])
```

Out[11]:

```
0      False
1      False
2      False
3      False
4      False
...
9995   False
9996   False
9997   False
9998   False
9999   False
Name: Age, Length: 10000, dtype: bool
```


Row Number	Customer Id	Surname	Credit Score	Geography	Tenure	Balance	NumOfProducts	Housing	IsActiveMember	.	Gender	Gender	Gender	Gender	Gender	Gender	Gender	Gender	Gender	Gender
											r_78	r_79	r_80	r_81	r_82	r_83	r_84	r_85	r_88	r_92
2	3	15619304	Onio	502	France	8	159660.80	3	1	0	.	0	0	0	0	0	0	0	0	0
3	4	15701354	Boni	699	France	1	0.00	2	0	0	.	0	0	0	0	0	0	0	0	0
4	5	15737888	Mitchell	850	Spain	2	125510.82	1	1	1	.	0	0	0	0	0	0	0	0	0

5 rows × 84 columns

8. Split the data into dependent and independent variables.

8.1 Split the data into dependent variables.

In [14]:

```
dependent_var = df.iloc[:, -1].values
print(dependent_var)
[1 0 1 ... 1 1 0]
```

8.2 Split the data into independent variables.

In [15]:

```
independent_var = df.iloc[:, :-2].values
print(independent_var)
[[1 15634602 'Hargrave' ... 1 1 1]
 [2 15647311 'Hill' ... 1 0 1]
 [3 15619304 'Onio' ... 3 1 0]
 ...
 [9998 15584532 'Liu' ... 1 0 1]
 [9999 15682355 'Sabbatini' ... 2 1 0]]
```



```
[10000 15628319 'Walker' ... 1 1 0]]
```

9. Scale the independent variables

In [17]:

```
import pandas as pd
from sklearn.preprocessing import MinMaxScaler
Scaler = MinMaxScaler()
df[["RowNumber"]] = Scaler.fit_transform(df[["RowNumber"]])
print(df)
```

	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	\
0	0.0000	15634602	Hargrave	619	France	Female	42	
1	0.0001	15647311	Hill	608	Spain	Female	41	
2	0.0002	15619304	Onio	502	France	Female	42	
3	0.0003	15701354	Boni	699	France	Female	39	
4	0.0004	15737888	Mitchell	850	Spain	Female	43	
...	
9995	0.9996	15606229	Obijiaku	771	France	Male	39	
9996	0.9997	15569892	Johnstone	516	France	Male	35	
9997	0.9998	15584532	Liu	709	France	Female	36	
9998	0.9999	15682355	Sabbatini	772	Germany	Male	42	
9999	1.0000	15628319	Walker	792	France	Female	28	

	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMember	\
0	2	0.00	1	1	1	
1	1	83807.86	1	0	1	
2	8	159660.80	3	1	0	
3	1	0.00	2	0	0	
4	2	125510.82	1	1	1	
...	
9995	5	0.00	2	1	0	
9996	10	57369.61	1	1	1	
9997	7	0.00	1	0	1	
9998	3	75075.31	2	1	0	
9999	4	130142.79	1	1	0	

	EstimatedSalary	Exited
0	101348.88	1
1	112542.58	0
2	113931.57	1
3	93826.63	0
4	79084.10	0
...
9995	96270.64	0
9996	101699.77	0
9997	42085.58	1
9998	92888.52	1
9999	38190.78	0

```
[10000 rows x 14 columns]
```

10. Split the data into training and testing

In [18]:

```
from sklearn.model_selection import train_test_split
train_size = 0.8
x = df.drop(columns = ['Tenure']).copy()
y = df['Tenure']
x_train, x_rem, y_train, y_rem = train_test_split(x,y, train_size=0.8)
test_size = 0.5
x_valid, x_test, y_valid, y_test = train_test_split(x_rem, y_rem, test_size=
0.5)
print(x_train.shape), print(y_train.shape)
print(x_valid.shape), print(y_valid.shape)
print(x_test.shape), print(y_test.shape)
(8000, 13)
(8000,)
(1000, 13)
(1000,)
(1000, 13)
(1000,)
```

Out[18]:

```
(None, None)
```