

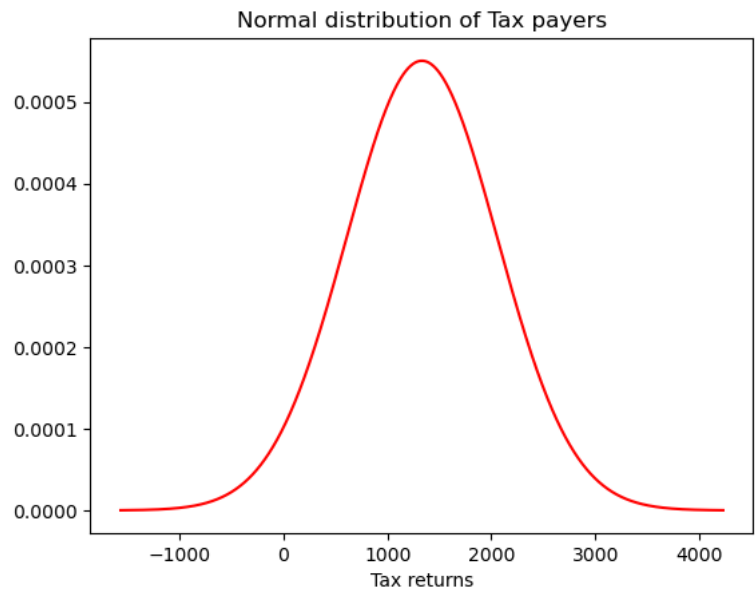
Importing necessary modules

```
In [49]: 1 from scipy.stats import zscore,norm,kurtosis,skew
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import pandas as pd
5 import statistics as st
6 import seaborn as sns
7 import warnings
8 warnings.filterwarnings('ignore')
```

Tax payer problems

```
In [88]: 1 mean=1332
2 stdev=725
3 x1=2000
4 x2=0
5 x3=100
6 x4=700
7 zscore4=(x4-mean)/stdev
8 zscore3=(x3-mean)/stdev
9 zscore2=(x2-mean)/stdev
10 zscore1=(x1-mean)/stdev
11 p1=norm.cdf(zscore1)
12 p2=norm.cdf(zscore2)
13 p3=norm.cdf(zscore3)
14 p4=norm.cdf(zscore4)
15
16 #Task1:
17 print("Proportion of tax payers who pay above 2000$",round(1-p1,2))
18
19 #Task2:
20 print("Proportion of tax payers who don't recieve returns",round(p2,2))
21
22 #Task3:
23 print("Proportion of tax payers who recieve returns between 100 to 700$ ",round(p4-p3,2))
24
25 #Representing in norm graph
26 lower_n= mean-4*stdev
27 upper_n= mean+4*stdev
28
29 norm_n= np.arange(lower_n,upper_n)
30 plt.plot(norm_n, norm.pdf(norm_n, mean , stdev), color='red')
31 plt.title("Normal distribution of Tax payers")
32 plt.xlabel("Tax returns")
33 plt.show()
```

Proportion of tax payers who pay above 2000\$ 0.18
Proportion of tax payers who don't recieve returns 0.03
Proportion of tax payers who recieve returns between 100 to 700\$ 0.15



High end video games

```
In [33]: 1 games=pd.read_csv(r"K:\Desktop\NIIT\tables\DS1_C5_S5_Computers_Data_Challenge.csv")
2 games1=games
3 games1
```

Out[33]:

	index	price	speed	hd	ram	screen	cd	multi	premium	ads	trend	
	0	1	1499	25	80	4	14	no	no	yes	94	1
	1	2	1795	33	85	2	14	no	no	yes	94	1
	2	3	1595	25	170	4	15	no	no	yes	94	1
	3	4	1849	25	170	8	14	no	no	no	94	1
	4	5	3295	33	340	16	14	no	no	yes	94	1

6254	6255	1690	100	528	8	15	no	no	yes	39	35	
6255	6256	2223	66	850	16	15	yes	yes	yes	39	35	
6256	6257	2654	100	1200	24	15	yes	no	yes	39	35	
6257	6258	2195	100	850	16	15	yes	no	yes	39	35	
6258	6259	2490	100	850	16	17	yes	no	yes	39	35	

6259 rows × 11 columns

Task1: To sample data with device screen size with price less than 4000\$

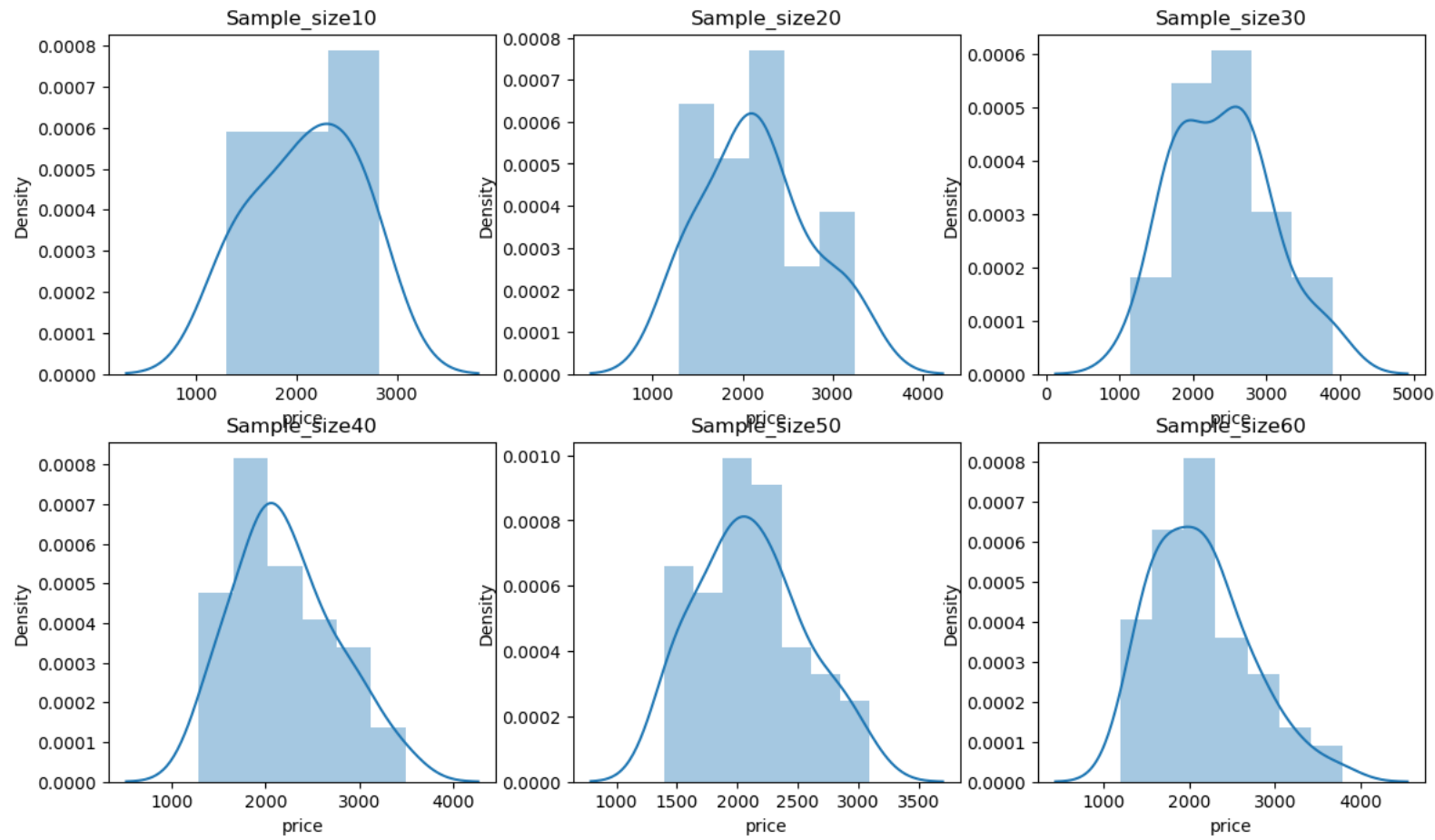
```
In [56]: 1 game_samp=games.price[games.price<4000].sample(50,ignore_index=True)
2 games=games[games.price<4000]
3 measures_game=pd.DataFrame([ [game_samp.mean(),games.price.mean()], [game_samp.median(),games.price.median()], [st.mode(game_samp),st.mode(games.price)], [game_samp.std(),games.price.std()]
4                               [kurtosis(game_samp),kurtosis(games.price)], [skew(game_samp),skew(games.price)]],index=["Mean", "Median", "Mode", "Standard dev", "Kurtosis", "Skewness"],columns=
```

Task 2 : To calculate measures of central tendency and show central limit thoerem with population data

Rule 1

```
In [71]: 1 print("As per central limit theorem's first rules we need to plot 6 different sample sizes to compare normal distribution")
2 fig,ax=plt.subplots(2,3,figsize=(14,8))
3 samps=[10,20,30,40,50,60]
4 j=0
5 i=0
6 for samp in samps:
7     sns.distplot(games.price.sample(samp,ignore_index=True,replace=True),ax=ax[j,i])
8     ax[j,i].set_title("Sample_size"+str(samp))
9     i+=1
10    if(i==3):
11        i=0
12        j=1
13
14
15
```

As per central limit theorem's first rules we need to plot 6 different sample sizes to compare normal distribution



Rule 2

```
In [57]: 1 print("As per central limit theorem 2 we can see the sample and population data have very close mean")
2 measures_game
```

As per central limit theorem 2 we can see the sample and population data have very close mean

Out[57]:

	Sample	Population
Mean	2276.320000	2208.855515
Median	2067.000000	2144.000000
Mode	1499.000000	1999.000000
Standard dev	710.293772	560.501267
Kurtosis	0.036968	-0.126395
Skewness	0.973392	0.518092

Rule 3

```
In [74]: 1 means=[]
2 for i in range(1,21):
3     means.append(games.price.sample(50,ignore_index=True,replace=True).mean())
4 samp_means_std=np.std(means)
5 pop_std=np.std(games.price)
6
7 print("The sample means std = {0} and population std by N sqrt = {1}".format(round(samp_means_std,2),round(pop_std/(50)**0.5),2))
```

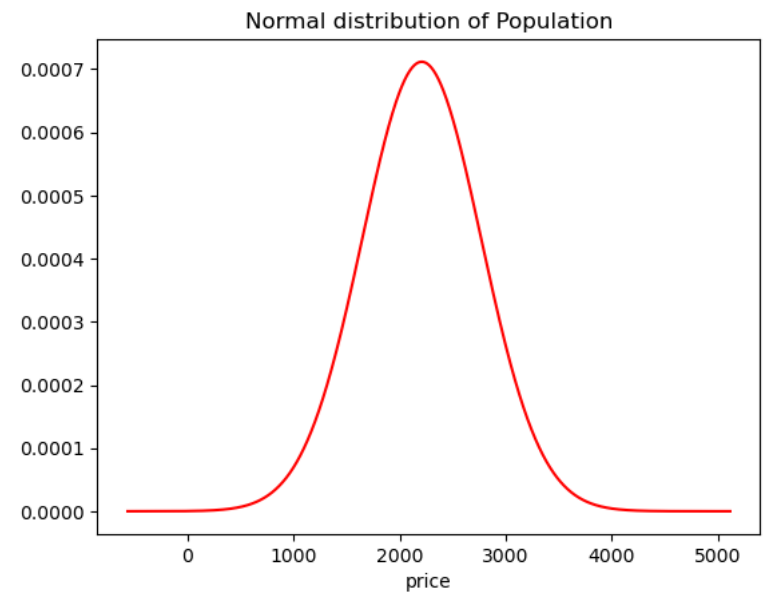
The sample means std = 77.59 and population std by N sqrt = 79

```
In [75]: """When comparing the sample means standard deviation with population's standard deviatio by sqrt of sample size we can see they are very close to each other
Thus the dataset follows centra limit theorem"""
```

Out[75]: "When comparing the sample means standard deviation with population's standard deviatio by sqrt of sample size we can see they are very close to each other\nThus the dataset follows centra limit theorem"

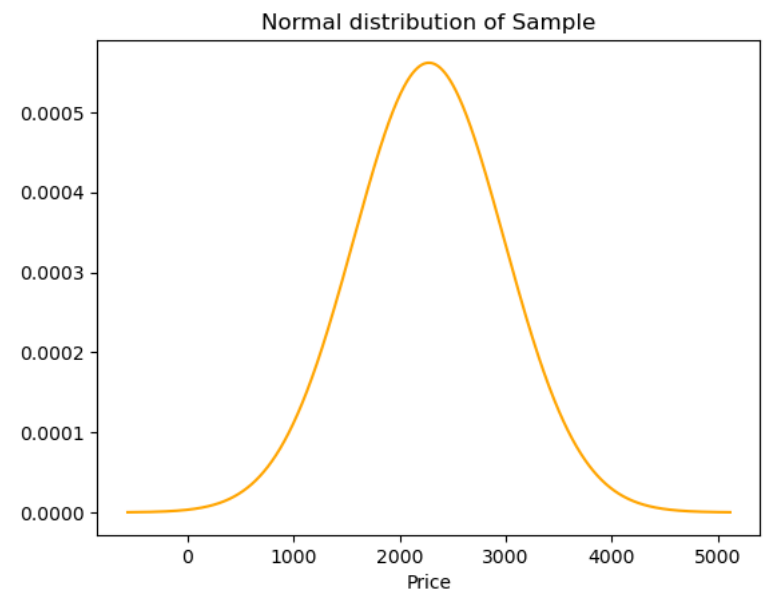
Task3 : To draw normal distribution plot for population

```
In [87]: 1 lower_p= games.price.mean()-4*games.price.std()
2 upper_p= games.price.mean()+4*games.price.std()
3
4 norm_p= np.arange(lower_s,upper_s)
5 plt.plot(norm_p, norm.pdf(norm_p,games.price.mean(),games.price.std()), color='red')
6 plt.title("Normal distribution of Population")
7 plt.xlabel("price")
8 plt.show()
```



Task4: To draw normal distribution plot for sample

```
In [84]: 1 lower_s= game_samp.mean()-4*game_samp.std()
2 upper_s= game_samp.mean()+4*game_samp.std()
3
4 norm_s= np.arange(lower_s,upper_s)
5 plt.plot(norm_s, norm.pdf(norm_s,game_samp.mean(),game_samp.std()), color='orange')
6 plt.title("Normal distribution of Sample")
7 plt.xlabel("Price")
8 plt.show()
```



Task5: To calculate zscores

In [107]:

```
1 price_z=zscore(game_samp)
2 probs=[]
3 for i in price_z:
4     probs.append(norm.cdf(i))
5
6 price_score=pd.DataFrame()
7 price_score["Price"]=game_samp
8 price_score["Zscore"]=price_z
9 price_score["Probability"]=probs
10 price_score
```

Out[107]:

	Price	Zscore	Probability
0	3090	1.157184	0.876401
1	2685	0.581209	0.719450
2	1854	-0.600607	0.274051
3	1879	-0.565053	0.286019
4	1468	-1.149562	0.125162
5	1708	-0.808243	0.209475
6	3334	1.504192	0.933734
7	1839	-0.621940	0.266991
8	2425	0.211447	0.583731
9	1499	-1.105475	0.134477
10	3795	2.159808	0.984606
11	1679	-0.849486	0.197806
12	2575	0.424771	0.664498
13	2445	0.239890	0.594792
14	3799	2.165497	0.984825
15	2495	0.310998	0.622099
16	1799	-0.678826	0.248624
17	1644	-0.899261	0.184257
18	3595	1.875376	0.969629
19	2395	0.168782	0.567016
20	2444	0.238468	0.594241
21	2190	-0.122761	0.451148
22	1740	-0.762734	0.222811
23	2404	0.181582	0.572044
24	1395	-1.253379	0.105034
25	1873	-0.573586	0.283124
26	2299	0.032255	0.512866
27	1595	-0.968947	0.166286
28	1395	-1.253379	0.105034
29	1989	-0.408615	0.341411
30	2145	-0.186758	0.425925
31	2495	0.310998	0.622099
32	3995	2.444241	0.992742
33	1790	-0.691626	0.244586
34	2644	0.522900	0.699478
35	2925	0.922528	0.821873
36	3904	2.314824	0.989689
37	1899	-0.536610	0.295769
38	2390	0.161671	0.564218
39	1894	-0.543721	0.293317
40	1595	-0.968947	0.166286
41	1799	-0.678826	0.248624
42	1499	-1.105475	0.134477
43	2540	0.374996	0.646168
44	1699	-0.821042	0.205811
45	2399	0.174471	0.569252
46	1740	-0.762734	0.222811
47	3440	1.654941	0.951032
48	1728	-0.779800	0.217754
49	1970	-0.435637	0.331550

Task9: To find the probability of getting less than 1301

In [113]:

```
1 X=(1301-game_samp.mean())/game_samp.std()
2
3 print("The probability of getting less than 1301 is ",1-norm.cdf(X))
```

The probability of getting less than 1301 is 0.9151427875687974

Task10: To find the probability of occurences 2000 and 2900

In [114]:

```
1 X1=(2000-game_samp.mean())/game_samp.std()
2 X2=(2900-game_samp.mean())/game_samp.std()
3 print("The probability of getting between 2000 and 2900 is ",norm.cdf(X2)-norm.cdf(X1))
```

The probability of getting between 2000 and 2900 is 0.46141432485830275