

In [38]:

```
1 from scipy.stats import binom,poisson,expon,mode
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5 import warnings
6 import math
7 warnings.filterwarnings("ignore")
8
9 ### source: tinyurl.com/4zs92xtk
```

Assuming that X represents the data (population), if X has a distribution with average μ and standard deviation σ , and if X is approximately normally distributed or if the sample size n is large.

Then the sample mean \bar{X} is normally distributed with the average μ and standard error σ/\sqrt{n}

The above Distribution is only valid if,

X is approximately normal or sample size n is large, and,

the data (population) standard deviation σ is known.

If X is normal, then \bar{X} is also normally distributed regardless of the sample size n . Central Limit Theorem tells us that even if X is not normal, if the sample size is large enough (usually greater than 30), then \bar{X} 's distribution is approximately normal (Sharpe, De Veaux, Velleman and Wright, 2020, pp. 318–320). If \bar{X} is normal, we can easily standardize and convert it to the standard normal distribution Z .

If the population standard deviation σ is not known, we cannot assume that the sample mean \bar{X} is normally distributed. If certain conditions are satisfied (explained below), then we can transform \bar{X} to another random variable t such that,

$$t = \frac{\bar{X} - \mu}{s/\sqrt{n}}$$

The random variable t is said to follow the t -distribution with $n-1$ degrees of freedom, where n is the sample size. The t -distribution is bell-shaped and symmetric (just like the normal distribution) but has fatter tails compared to the normal distribution. This means values further away from the mean have a higher likelihood of occurring compared to that in the normal distribution.

The conditions to use the t -distribution for the random variable t are as follows (Sharpe et al., 2020, pp. 415–420):

If X is normally distributed, even for small sample sizes ($n < 15$), the t -distribution can be used.

If the sample size is between 15 and 40, the t -distribution can be used as long as X is unimodal and reasonably symmetric. For sample sizes greater than 40, the t -distribution can be used unless X 's distribution is heavily skewed.

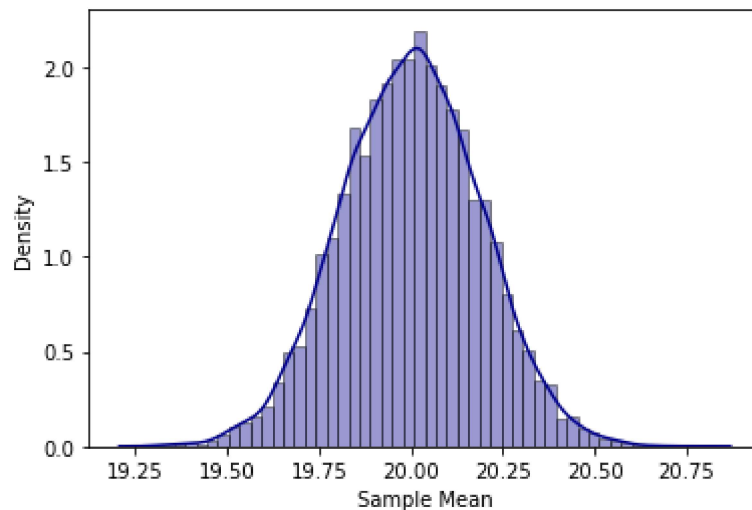
Let's draw a sample of size $n=250$ from the normal distribution. Here we are assuming that our data is normally distributed and has parameters $\mu = 20$ and $\sigma = 3$. Collecting one sample from this population

```
In [6]: 1 mu, sigma = 20, 3
        2 sample = np.random.normal(mu, sigma, 250)
        3 x_bar = sample.mean()
        4 x_bar
```

Out[6]: 20.042507169222105

```
In [7]: 1 ## But if I ran this code 10,000 times and recorded the values of x and plot
        2
        3 import seaborn as sns
        4 samples = np.zeros(10000)
        5 mu, sigma = 20, 3
        6 for s in range(10000):
        7     sample = np.random.normal(mu, sigma, 250) #sample of size 250
        8     x_bar = sample.mean() #calculating sample mean
        9     samples[s] = x_bar
       10 sns.distplot(samples, hist = True, color = 'darkblue',
       11                 hist_kws={'edgecolor': 'black'}).set(xlabel = 'Sample Mean',
       12                                                         ylabel = 'Density')
```

Out[7]: [Text(0.5, 0, 'Sample Mean'), Text(0, 0.5, 'Density')]



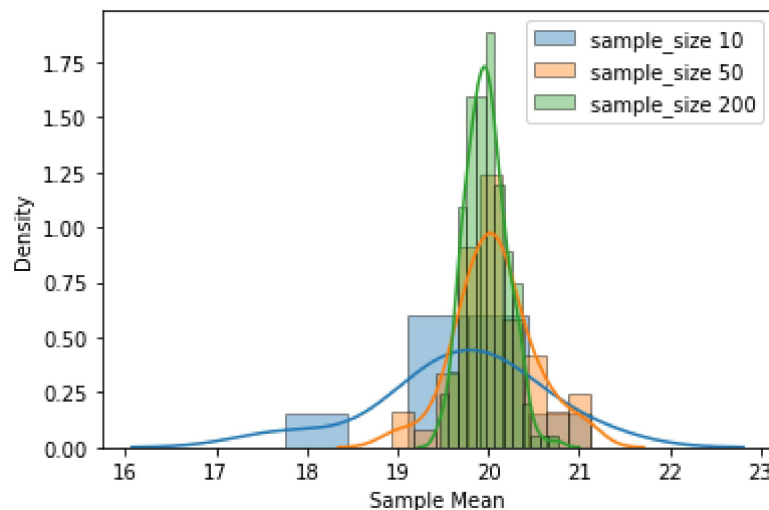
As you can see, the distribution is approximately symmetric and bell-shaped (just like the normal distribution) with an average of approximately 20 and a standard error that is approximately equal to $3/\sqrt{250} = 0.19$.

Sampling from the same population with different sample sizes will result in different measures of spread in the outcome distribution. As we expect, increasing the sample size will reduce the standard error and therefore, the distribution will be narrower around its average. Note that the

distribution of \bar{X} is normal even for extremely small sample sizes. This is because X is normally distributed.

```
In [17]: 1 mu, sigma = 20, 3
2
3 for sample_size in [10,50,200]:
4     samples = np.zeros(sample_size)
5     for s in range(sample_size):
6         sample = np.random.normal(mu, sigma, sample_size) #sample of size 2
7         x_bar = sample.mean() #calculating sample mean
8         samples[s] = x_bar
9     sns.distplot(samples, hist = True,
10                  hist_kws={'edgecolor':'black'},label= 'sample_size ' + str(s
11                                                                    ylabel = 'Density')
12 plt.legend()
```

Out[17]: <matplotlib.legend.Legend at 0x21d2287ca00>

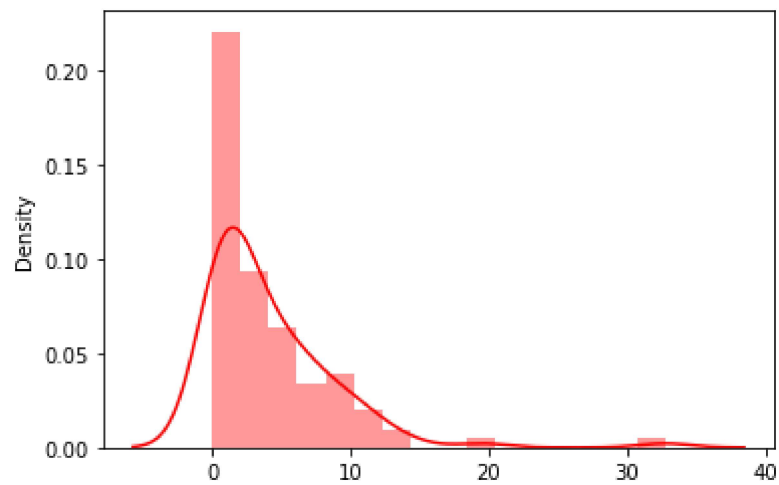


What if the population (data) is not normal?

No worries! Even if your data is not normally distributed, if the sample size is large enough, the distribution of \bar{X} can still be approximated using the normal distribution (according to Central Limit Theorem). The following figure shows the distribution of \bar{X} when X is heavily skewed to the left. As you can see, \bar{X} 's distribution tends to mimic the distribution of X for small sample sizes. However, as sample size grows the distribution of \bar{X} becomes more symmetric and bell-shaped. As mentioned above, if sample size is large (usually larger than 30), \bar{X} 's distribution is approximately normal regardless of what the distribution of X is.

```
In [18]: 1 ### Lets sample from exponential distribution and understand how the mean of  
2 expon_data = expon.rvs(scale=4,size= 100)  
3 sns.distplot(expon_data ,hist=True,kde=True,color="red")
```

Out[18]: <AxesSubplot:ylabel='Density'>

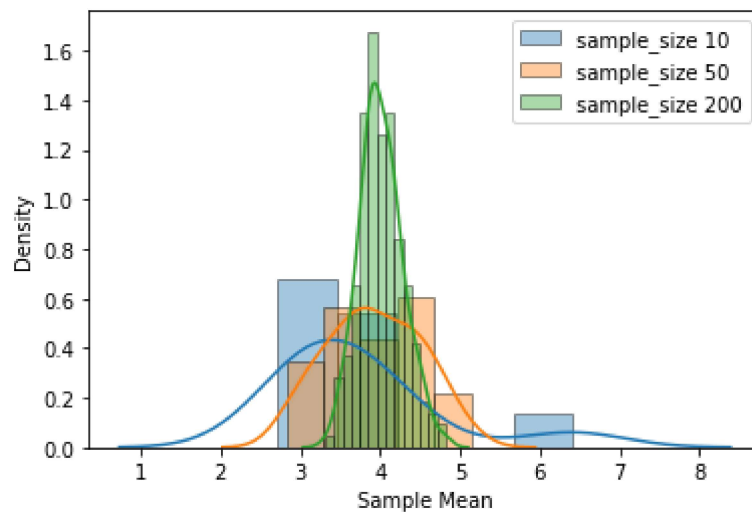


```

In [19]: 1 mu, sigma = 20, 3
          2
          3 for sample_size in [10,50,200]:
          4     samples = np.zeros(sample_size)
          5     for s in range(sample_size):
          6         sample = expon.rvs(scale=4,size= sample_size) #sample of size 250
          7         x_bar = sample.mean() #calculating sample mean
          8         samples[s] = x_bar
          9     sns.distplot(samples, hist = True,
10                        hist_kws={'edgecolor':'black'},label= 'sample_size ' + str(s)
11                        ylabel = 'Density')
12 plt.legend()

```

Out[19]: <matplotlib.legend.Legend at 0x21d2192e280>



As we can see if we increase the sample size the distribution comes down to a normal distribution with mean 4

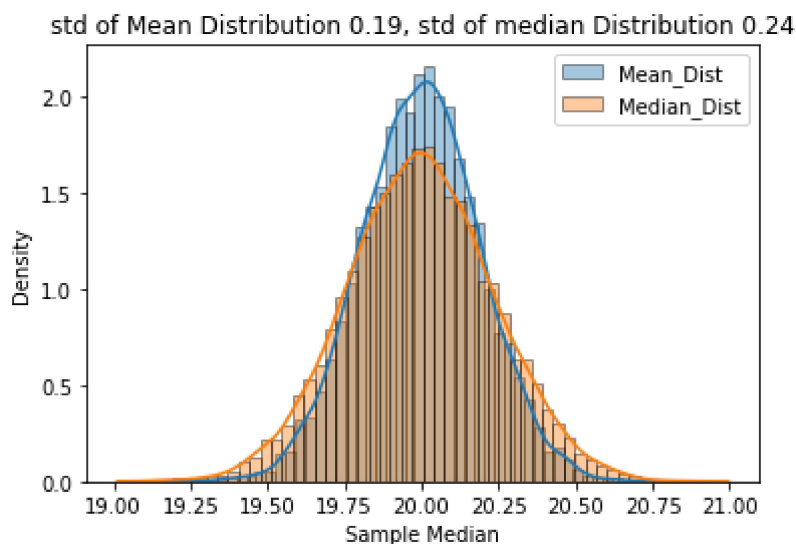
What happens when we compare the median sample distribution with mean sample distribution

```

In [33]: 1  ## But if I ran this code 10,000 times and recorded the values of x and plot
2
3  import seaborn as sns
4  samples_mean = np.zeros(10000)
5  samples_median = np.zeros(10000)
6  mu, sigma = 20, 3
7  for s in range(10000):
8      sample = np.random.normal(mu, sigma, 250)  #sample of size 250
9      x_bar = sample.mean()  #calculating sample mean
10     samples_mean[s] = x_bar
11     samples_median[s] = np.median(sample)
12     sns.distplot(samples_mean, hist = True, label = 'Mean_Dist',
13                 hist_kws={'edgecolor':'black'}).set(xlabel = 'Sample Mean',
14                                                     ylabel = 'Density')
15     sns.distplot(samples_median, hist = True, label = 'Median_Dist',
16                 hist_kws={'edgecolor':'black'}).set(xlabel = 'Sample Median',
17                                                     ylabel = 'Density')
18     plt.title("std of Mean Distribution {}, std of median Distribution {}".format
19             samples_mean, samples_median)
20     plt.legend()

```

Out[33]: <matplotlib.legend.Legend at 0x21d2874fe20>



As we can see mean is same for Distribution of sample Mean and Median, but standard deviation is more for Median by .24/.19 around 1.26

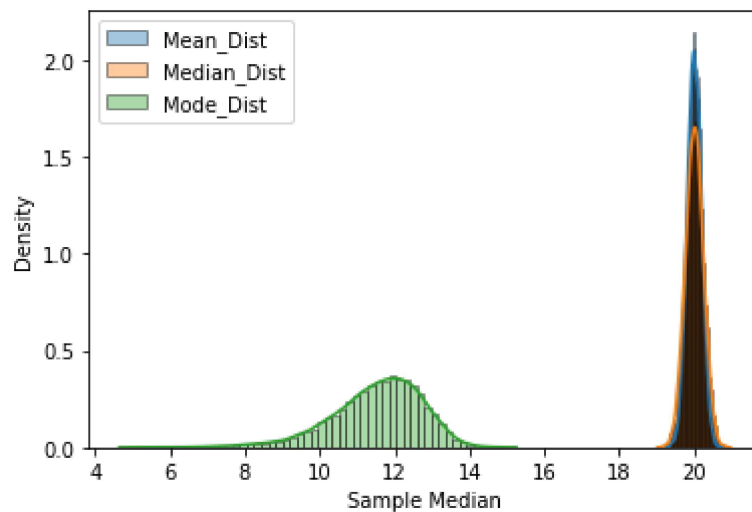
Lets check what happens with the mode

```

In [49]: 1 samples_mean = np.zeros(10000)
          2 samples_median = np.zeros(10000)
          3 samples_mode = np.zeros(10000)
          4 mu, sigma = 20, 3
          5 for s in range(10000):
          6     sample = np.random.normal(mu, sigma, 250) #sample of size 250
          7     x_bar = sample.mean() #calculating sample mean
          8     samples_mean[s] = x_bar
          9     samples_median[s] = np.median(sample)
         10     samples_mode[s] = mode(sample).mode[0]
         11 sns.distplot(samples_mean, hist = True, label = 'Mean_Dist',
         12                 hist_kws={'edgecolor':'black'}).set(xlabel = 'Sample Mean',
         13                 ylabel = 'Density')
         14 sns.distplot(samples_median, hist = True, label = 'Median_Dist',
         15                 hist_kws={'edgecolor':'black'}).set(xlabel = 'Sample Median',
         16                 ylabel = 'Density')
         17 sns.distplot(samples_mode, hist = True, label = 'Mode_Dist',
         18                 hist_kws={'edgecolor':'black'}).set(xlabel = 'Sample Median',
         19                 ylabel = 'Density')
         20 plt.legend()

```

Out[49]: <matplotlib.legend.Legend at 0x21d28f9e8e0>



No conclusive statement form MODE

