

Forward School

Program Code: J620-002-4:2020

Program Name: FRONT-END SOFTWARE DEVELOPMENT

Title : Exe11 - Normal Distribution Exercise

Name: Chong Mun Chen

IC Number: 960327-07-5097

Date : 10/7/2023

Introduction : Practising more on normal distribution of data.

Conclusion : Getting more familiar with the normal distribution and plotting graphs with the normalized data.

Normal Distribution

The normal distribution is defined by the following probability density function, where μ is the population mean and σ^2 is the variance.

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$$

If a random variable X follows the normal distribution, then we write:

$$X \sim N(\mu, \sigma^2)$$

In particular, the normal distribution with $\mu = 0$ and $\sigma = 1$ is called the standard normal distribution, and is denoted as $N(0,1)$. It can be graphed as follows.

The normal distribution is important because of the **Central Limit Theorem**, which states that the population of all possible samples of size n from a population with mean μ and variance σ^2 approaches a normal distribution with mean μ and $\frac{\sigma^2}{n}$ when n approaches infinity.

Read and understand more about **Central Limit Theorem (CLT)** [here](https://statisticsbyjim.com/basics/central-limit-theorem/)
(<https://statisticsbyjim.com/basics/central-limit-theorem/>)

Question 1

Suppose widge weights produced at MS Widge Works have weights that are normally distributed with mean 17.46 grams and variance 375.67 grams. What is the probability that a randomly chosen widge weighs more than 19 grams?

```
In [6]:  ▶ from scipy.stats import norm
import math

# 1 - norm.cdf(19, 17.46, math.sqrt(375.67))
norm.sf(19, 17.46, math.sqrt(375.67))
```

Out[6]: 0.46833563578991133

Question 2

Suppose IQ scores are normally distributed with mean 100 and standard deviation 15. What is the 95th percentile of the distribution of IQ scores?

```
In [7]:  ▶ norm.ppf(0.95, 100, 15)
```

Out[7]: 124.67280440427209

Question 3

Suppose wages are normally distributed with a mean of 1900 and a standard deviation of 150.

1. What percentage of people have wages less than 1800?
2. What percentage of people have wages greater than 2100?
3. What percentage of people have wages between 1800 and 2100?
4. What wages separate the top 10% from the others?
5. What wages separate the lower 25% from the others?

```
In [9]:  ▶ # 1
norm.cdf(1800, 1900, 150)
```

Out[9]: 0.2524925375469229

```
In [11]: ▶ # 2
norm.sf(2100, 1900, 150)
```

Out[11]: 0.09121121972586788

```
In [12]: # 3
norm.cdf(2100, 1900, 150) - norm.cdf(1800, 1900, 150)
```

Out[12]: 0.6562962427272092

```
In [8]: # 4
norm.ppf(0.9, 1900, 150)
```

Out[8]: 2092.23273483169

```
In [15]: # 5
norm.ppf(0.25, 1900, 150)
```

Out[15]: 1798.8265374705877

Question 4

Based on the Ages of Death during the Spanish Flu, 1918.

Demonstration of the central limit theorem, using the distribution of sample mean age at death in samples from a highly non-normal distribution: the frequency distribution of age at death in Switzerland in 1918 during the Spanish flu epidemic.

```
In [2]: ## Question 4

import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline

path="http://whitlockschluter.zoology.ubc.ca/wp-content/data/chapter10/chap
flu = pd.read_csv(path)
flu.head()
```

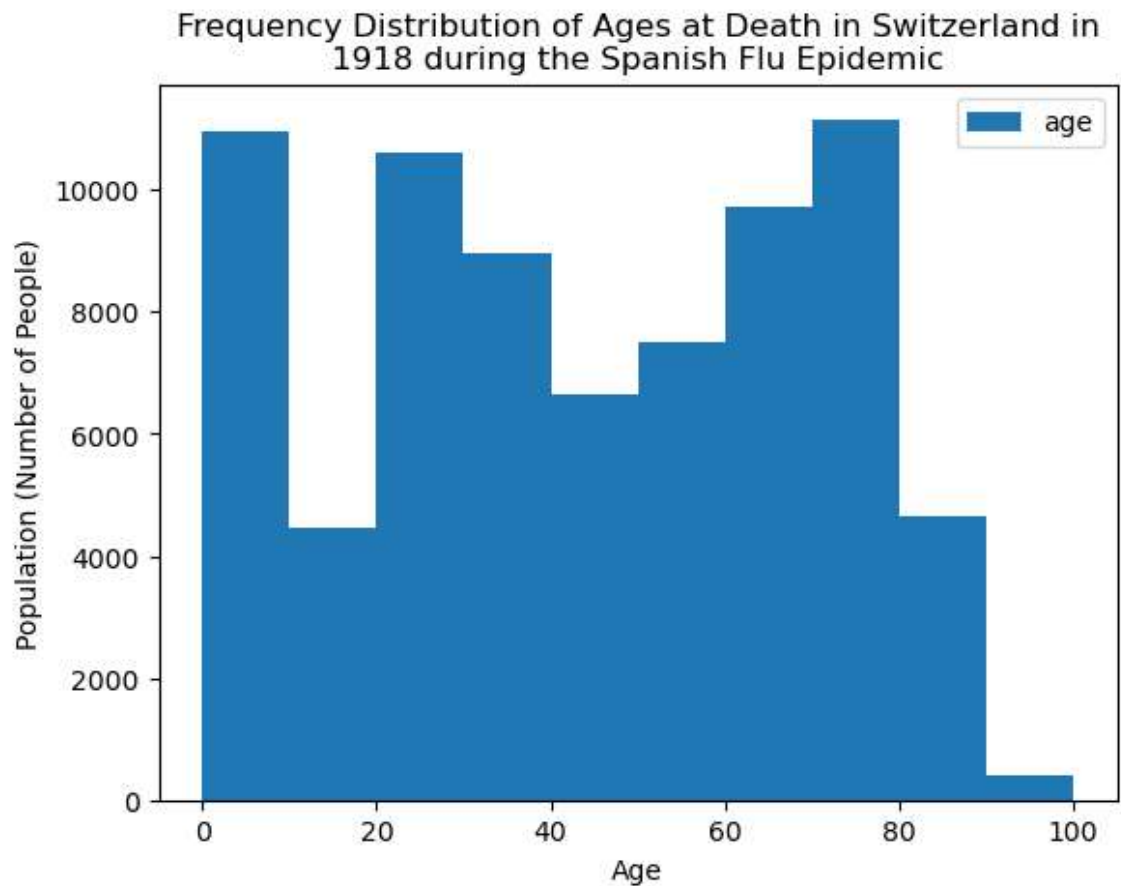
Out[2]:

	age
0	0
1	0
2	0
3	0
4	0

Histogram showing the frequency distribution of ages at death in Switzerland in 1918 during the Spanish flu epidemic.

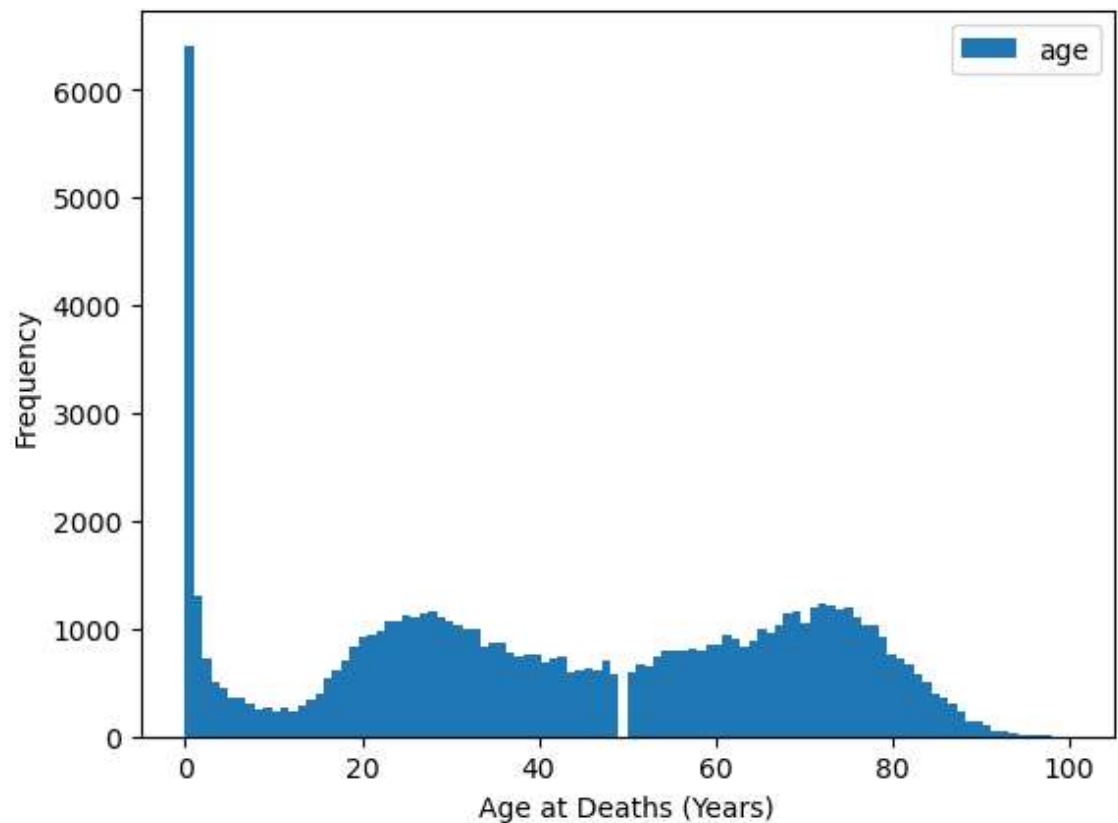
```
In [3]: import textwrap

flu.plot.hist()
title_text = 'Frequency Distribution of Ages at Death in Switzerland in 1918 during the Spanish Flu Epidemic'
plt.title(textwrap.fill(title_text, width=60))
plt.xlabel('Age')
plt.ylabel('Population (Number of People)')
plt.show()
```



Histogram with better binning (0,102,2) and axis labels

```
In [4]: flu.plot.hist(bins=102)
plt.xlabel('Age at Deaths (Years)')
plt.show()
```



Demonstrate the **central limit theorem**. Treat the age at death measurements from Switzerland in 1918 as the population. Take a large number of random samples, each of size n , from the population of age at death measurements and plot the sample means.

Note: your results won't be the identical to the one shown below, because 10,000 random samples is not large enough for extreme accuracy. Change the n below to another number and rerun to see the effects of sample size on the shape of the distribution of sample means.

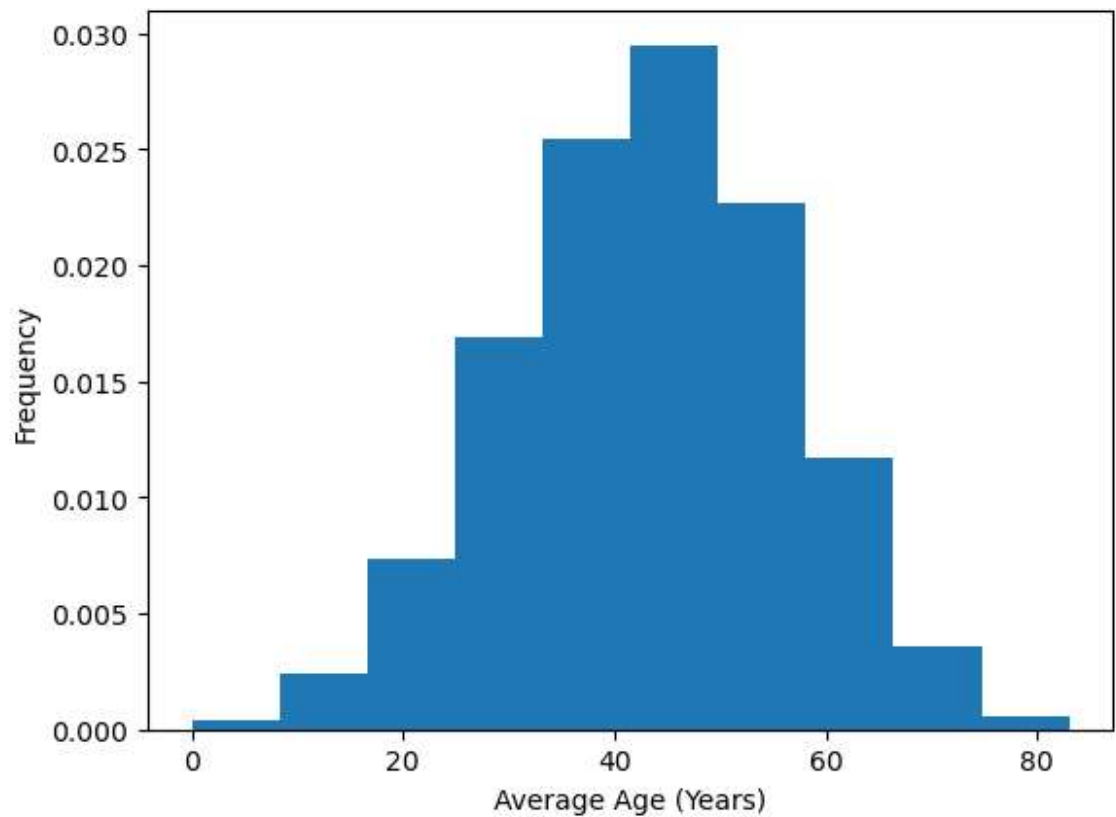
Write a loop to sample 10000 times from 'Age'. Each time, collect 4 samples. Store the average age in a new variable, `age1`. Plot the histogram for `age1`.

```
In [9]: ▶ import random

age1 = []
for i in range(10000):
    sample_age = flu.sample(n=4)
    age1.append(sample_age['age'].mean())

    # Alternate method
#     random_numbers = random.choices(flu['age'].values, k=4)
#     age1.append(statistics.mean(random_numbers))

plt.hist(x=age1, density=True)
plt.xlabel('Average Age (Years)')
plt.ylabel('Frequency')
plt.show()
```



Histogram of the sample means with more options

```
In [10]: ▶ import random

age1 = []
for i in range(10000):
    sample_age = flu.sample(n=4)
    age1.append(sample_age['age'].mean())

plt.hist(x=age1, density=True, bins = 50, edgecolor = 'black', color = 'red')
plt.xlabel('Average Age (Years)')
plt.ylabel('Frequency')
plt.show()
```

