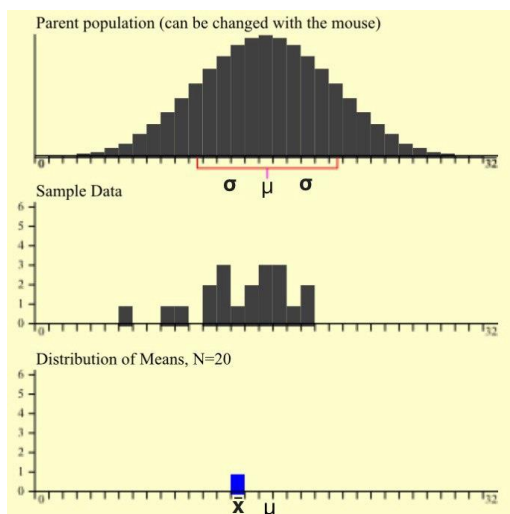Population parameters:
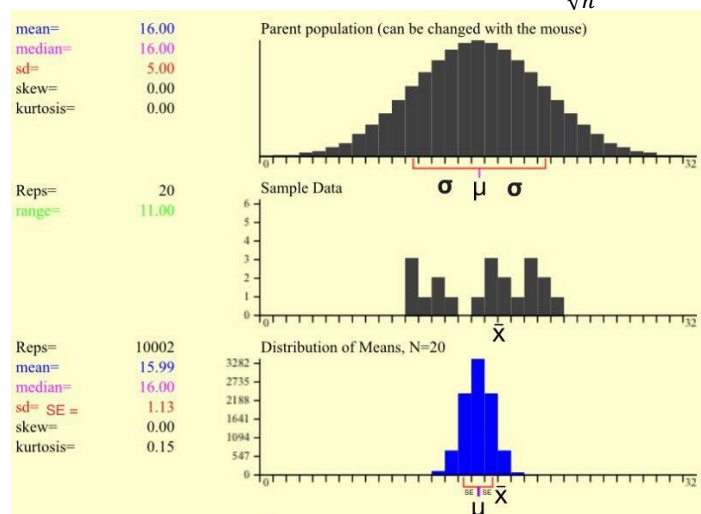- μ is the population mean
- σ is the population standard deviation

According to **CLT for normal distributions** for any sample size n:
- **(1)** The distribution of the sample mean will be normal.
- **(2)** The mean of the distribution of the sample mean will be equal to μ.
- **(3)** The standard deviation (standard error) of the distribution of the sample mean $SE = \frac{\sigma}{\sqrt{n}}$.



| 1 sample of size 20 is run | 10002 samples of size 20 are run and the results follow **(1)**, **(2)**, **(3)** |

In hypothesis testing:

We first assume($H_0$ - null hypothesis) population mean to be equal to some value μ.

Then we get some sample of size **n**, calculate its mean($\bar{x}$) and variance($s^2$) and try to reject it with confidence level c or with significance level a = 1-c.

$z_{score}$ calculates distance between $\bar{x}$(sample mean) and μ(the mean of the distribution of the sample mean, not population mean) taking SE as a unit distance.

$$z_{score} = \frac{\bar{x} - \mu}{SE} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

**But why do we need $z_{score}$, why do we need distance between $\bar{x}$ and μ?**

According to **CLT for normal distributions** the distribution of the sample mean follows **(1)**, **(2)**, **(3)** rules.

(2) The mean of the distribution of the sample mean is equal to μ.

(3) SE(Standard deviation of the distribution of the sample mean) = $\frac{\sigma}{\sqrt{n}}$.

(1) The distribution of the sample mean is normal. **This lets us to know what percentage of the data(sample means) is covered by [μ-n\*SE, μ+n\*SE] interval.**

| $z_{score}$ | Percentage of data covered within $z_{score}$*SE of the μ(mean) | Maximum confidence level with which $H_0$ can be rejected | P value(Minimum significance level with which $H_0$ can be rejected) |
|---|---|---|---|
| 1 | 68% | 68%, 0.68 | 0.32 |
| 2 | 95% | 95%, 0.95 | 0.05 |
| 3 | 99.7% | 99.7%, 0.997 | 0.003 |

It means that if $H_0$ is true and we take a sample, then the probability of that sample mean($\bar{x}$) to be in [μ-SE, μ+SE] interval is 68%, to be in [μ-2\*SE, μ+2\*SE] interval is 95% and so on.

And if x̄ is out of some interval, then we can reject the H₀ with the confidence level that is equal to the probability of x̄ to be in that interval which is equal to the percentage of data(sample means) covered within that interval.

For example, if x̄ is out of [μ-2*SE, μ+2*SE] interval then we can reject the H₀ with the confidence level 95%, because the probability of x̄ to be in that interval is 95%, because 95% of the data(sample means) is covered by that interval.

And to reject the H₀ with the confidence level c%, x̄ has to be out of such interval that covers more than c% of data.

For example, to reject H₀ with confidence level 95%, x̄ has to be out of [μ-2*SE, μ+2*SE], in other words z$_{score}$ has to be greater than 2, because [μ-2*SE, μ+2*SE] covers 95% of the data.

The bigger the confidence level, the bigger the z$_{score}$(distance between x̄ and μ) has to be for us to be able to reject H₀.

When the population variance($\sigma^2$) is unknown, sample variance($s^2$) is used instead and z$_{score}$ "becomes" t$_{score}$.

$$SE = \frac{s}{\sqrt{n}}$$

$$t_{score} = \frac{\bar{x}-\mu}{SE} = \frac{\bar{x}-\mu}{\frac{s}{\sqrt{n}}}$$

There are tables where for each confidence level their corresponding minimum z and t scores are shown to reject the H₀.

### An example problem

At a water-bottling factory, a machine is supposed to put 2 liters of water into the bottles. After an overhaul, management thinks the machine is no longer putting the correct amount of water in. They sample 20 bottles and find an avg of 2.10 L of water with standard deviation of 0.33 L. Test the claim at 0.01 level of significance.

H₀: μ = 2          n(sample size) = 20

Hₐ: μ ≠ 2          x̄(sample mean) = 2.1

c = 0.99           s(sample standard deviation) = 0.9

$$t = \frac{x-\mu}{SE} = \frac{x-\mu}{\frac{s}{\sqrt{n}}} = 0.4969$$

It's obvious that with t<1 we can't reject H₀ with 99% confidence level. Because we know that only 68% of the data is covered within 1 SE of the μ in normal distributions. So if x̄ isn't even out of that interval(x̄ is t*SE distant from μ) then we can't even reject H₀ with 68 % confidence level.

If it's not that obvious, then we can look at t-test table to see how much t at least has to be for us to be able to reject the H₀. T-test table shows that t has to be not less than 2.845 for us to be able to reject H₀ with 99% confidence level.

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 17 | 0.000 | 0.689 | 0.863 | 1.069 | 1.333 | 1.740 | 2.110 | 2.567 | 2.898 | 3.646 |
| 18 | 0.000 | 0.688 | 0.862 | 1.067 | 1.330 | 1.734 | 2.101 | 2.552 | 2.878 | 3.610 |
| 19 | 0.000 | 0.688 | 0.861 | 1.066 | 1.328 | 1.729 | 2.093 | 2.539 | 2.861 | 3.579 |
| 20 | 0.000 | 0.687 | 0.860 | 1.064 | 1.325 | 1.725 | 2.086 | 2.528 | 2.845 | 3.552 |
| 21 | 0.000 | 0.686 | 0.859 | 1.063 | 1.323 | 1.721 | 2.080 | 2.518 | 2.831 | 3.527 |
| 22 | 0.000 | 0.686 | 0.858 | 1.061 | 1.321 | 1.717 | 2.074 | 2.508 | 2.819 | 3.505 |
| 23 | 0.000 | 0.685 | 0.858 | 1.060 | 1.319 | 1.714 | 2.069 | 2.500 | 2.807 | 3.485 |
| 24 | 0.000 | 0.685 | 0.857 | 1.059 | 1.318 | 1.711 | 2.064 | 2.492 | 2.797 | 3.467 |
| 25 | 0.000 | 0.684 | 0.856 | 1.058 | 1.316 | 1.708 | 2.060 | 2.485 | 2.787 | 3.450 |
| 26 | 0.000 | 0.684 | 0.856 | 1.058 | 1.315 | 1.706 | 2.056 | 2.479 | 2.779 | 3.435 |
| 27 | 0.000 | 0.684 | 0.855 | 1.057 | 1.314 | 1.703 | 2.052 | 2.473 | 2.771 | 3.421 |
| 28 | 0.000 | 0.683 | 0.855 | 1.056 | 1.313 | 1.701 | 2.048 | 2.467 | 2.763 | 3.408 |
| 29 | 0.000 | 0.683 | 0.854 | 1.055 | 1.311 | 1.699 | 2.045 | 2.462 | 2.756 | 3.396 |
| 30 | 0.000 | 0.683 | 0.854 | 1.055 | 1.310 | 1.697 | 2.042 | 2.457 | 2.750 | 3.385 |
| 40 | 0.000 | 0.681 | 0.851 | 1.050 | 1.303 | 1.684 | 2.021 | 2.423 | 2.704 | 3.307 |
| 60 | 0.000 | 0.679 | 0.848 | 1.045 | 1.296 | 1.671 | 2.000 | 2.390 | 2.660 | 3.232 |
| 80 | 0.000 | 0.678 | 0.846 | 1.043 | 1.292 | 1.664 | 1.990 | 2.374 | 2.639 | 3.195 |
| 100 | 0.000 | 0.677 | 0.845 | 1.042 | 1.290 | 1.660 | 1.984 | 2.364 | 2.626 | 3.174 |
| 1000 | 0.000 | 0.675 | 0.842 | 1.037 | 1.282 | 1.646 | 1.962 | 2.330 | 2.581 | 3.098 |
| z | 0.000 | 0.674 | 0.842 | 1.036 | 1.282 | 1.645 | 1.960 | 2.326 | 2.576 | 3.090 |
| | 0% | 50% | 60% | 70% | 80% | 90% | 95% | 98% | 99% | 99.8% |
| | | | | | Confidence Level | | | | | |