

Hypothesis testing:

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug, subjecting each to neurological stimulus, and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats' response times is 1.05 seconds with a sample standard deviation of 0.5 seconds. Do you think that the drug has an effect on response time?

H_0 : drug has no effect

H_1 : drug has an effect

H_0 : Drug has no effect $\Rightarrow \mu = 1.2 \text{ s}$ (even w/ drug)
 H_1 : Drug has an effect $\Rightarrow \mu \neq 1.2 \text{ s}$ when the drug is given

Let's assume the H_0 (null hypothesis) is true:

- The mean of our sample distribution is the same as the mean of the population distribution(1.2 s)
- The standard deviation of our sampling distribution should be equal to $1/(\text{square root of population distribution})$
- We don't know the standard deviation of the population, so we estimate it as the standard deviation of the sample, which is 0.5s (notice that the sd of the sample is not the sd of the sample distribution)

- $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \approx \frac{s}{\sqrt{100}} = \frac{0.5}{10}$

So if the sample(drug has no effect) is come from the population, we need to consider what is the probability for the sample to get a mean of 1.05, or what is the standard deviation for the population

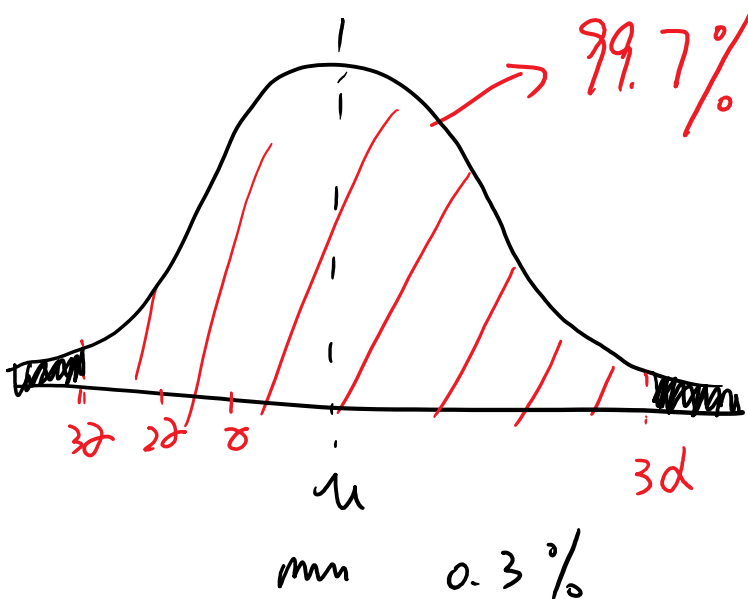
with a mean of 1.2 s to get a sample mean of 1.05s

Z score:

$$Z = (1.2 - 1.05) / 0.05 = 3$$

So the sample mean is 3 sd away from the population mean

The probability of within 3 sd is 99.7% (two-tail)



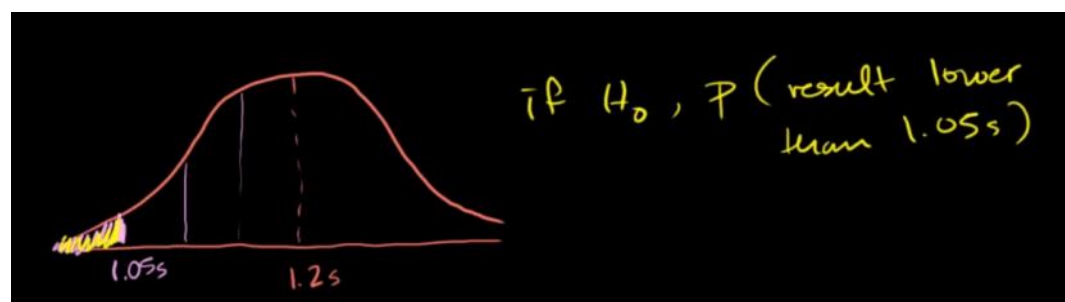
We are going to reject the null hypothesis (two tail test)

One tail test:

H0: drug has no effect

H1: drug lowers response time

H_0 : Drug has no effect: $\mu_{\text{drug}} = 1.2\text{s}$
 H_1 : Drug lowers response: $\mu_{\text{drug}} < 1.2\text{s}$



Z-statistics & t-statistics:

$$Z = \frac{\bar{X} - \mu}{\sigma_{\bar{X}}}$$

sample mean population mean

sample distribution SD

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

sample size

SD of population

SS
S

SD of sample

Here is the problem:

We can only assume the sd of the population is approximately the sd of the sample, when the sample size is relatively big.

And only in this situation, can we apply z score test

In most of the case, we need our sample size to be greater than 30.

If the sample size is small (<30), we use **t test**, look up the t table.

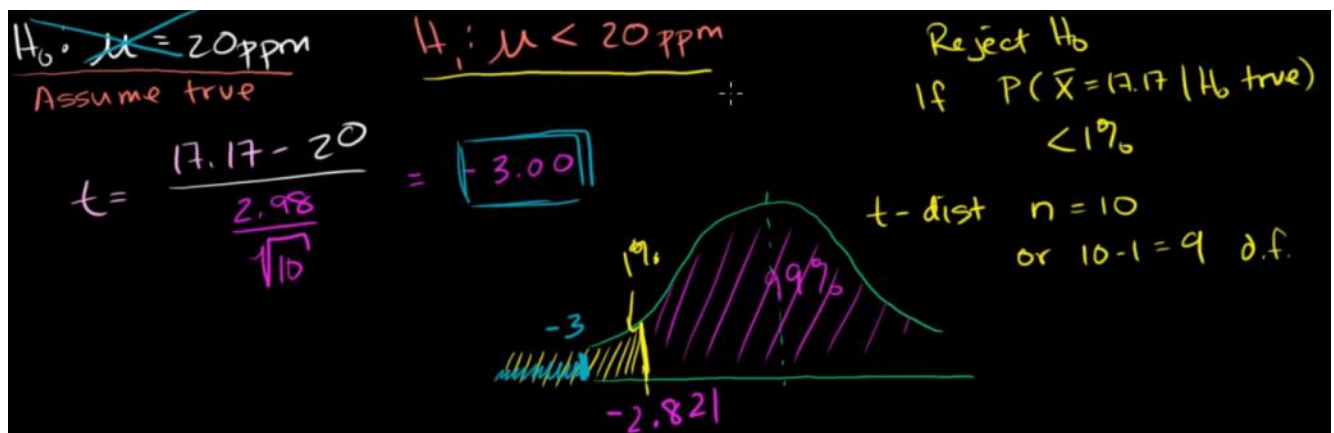
Type 1 error:

Rejecting H_0 even though it is true

The mean emission of all engines of a new design needs to be below 20 ppm if the design is to meet new emission requirements. Ten engines are manufactured for testing purposes, and the emission level of each is determined. The emission data is:

15.6 16.2 22.5 20.5 16.4 19.4 16.6 17.9 12.7 13.9 $\bar{X} = 17.17$ $S = 2.98$

Does the data supply sufficient evidence to conclude that this type of engine meets the new standard? Assume we are willing to risk a Type I error with probability = 0.01



And vice visa, you can define the confidence interval, and calculate the new requirement of engines

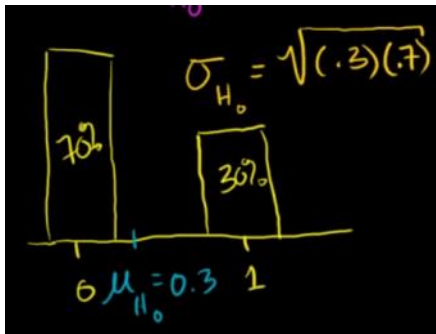
We want to test the hypothesis that more than 30% of U.S. households have Internet access (with a significance level of 5%). We collect a sample of 150 households and find that 57 have access.

$H_0: p \leq 30\%$

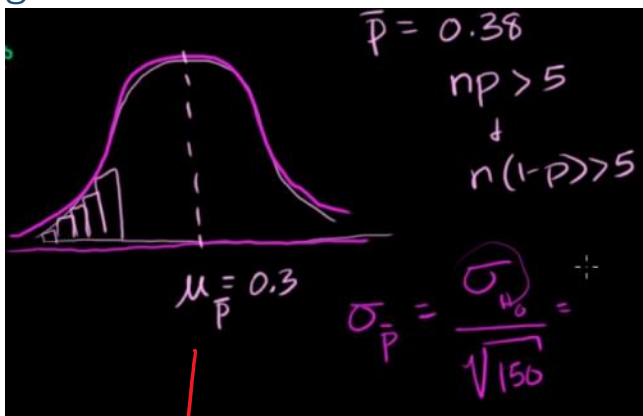
$H_1: p > 30\%$

Assuming H_0 true

This is a Bernoulli distribution, the SD = square root of $(p \cdot (1-p))$



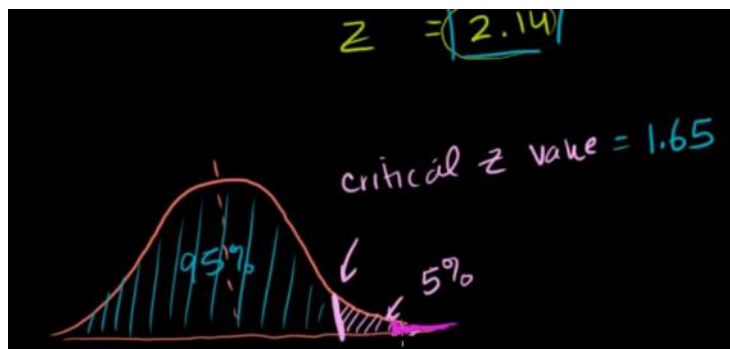
In Bernoulli distribution, you can assume the sample distribution is a normal distribution when $np > 5$ and $n(1-p)$ are greater than 5



sample distribution graph
population mean

sample distribution SD

$$Z = \frac{\bar{P} - \mu_{\bar{P}}}{\sigma_{\bar{P}}} = \frac{0.38 - 0.3}{0.037} = \frac{0.08}{0.037}$$
$$Z = 2.14$$



Variance of differences of random variables:

$$\text{Var}(X) = E[(X - \mu_x)^2] = \sigma_x^2$$

$$\text{Var}(Y) = E[(Y - \mu_y)^2] = \sigma_y^2$$

$$Z = X + Y$$

$$A = X - Y$$

$$\text{Var}(Z) = \text{Var}(X) + \text{Var}(Y)$$

$$\text{Var}(A) = \text{Var}(X) + \text{Var}(Y)$$

Proving:

<https://en.wikipedia.org/wiki/Variance>

$$\text{Var}(X) = \mathbb{E}[(X - \mu)^2].$$

This definition encompasses random variables that are generated by processes that are [discrete](#), [continuous](#), [neither](#), or mixed. The variance can also be thought of as the [covariance](#) of a random variable with itself:

$$\text{Var}(X) = \text{Cov}(X, X).$$

The variance is also equivalent to the second [cumulant](#) of a probability distribution that generates X . The variance is typically designated as $\text{Var}(X)$, σ_X^2 , or simply σ^2 (pronounced "[sigma](#) squared"). The expression for the variance can be expanded:

$$\begin{aligned}\text{Var}(X) &= \mathbb{E}[(X - \mathbb{E}[X])^2] \\ &= \mathbb{E}[X^2 - 2X\mathbb{E}[X] + \mathbb{E}[X]^2] \\ &= \mathbb{E}[X^2] - 2\mathbb{E}[X]\mathbb{E}[X] + \mathbb{E}[X]^2 \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2\end{aligned}$$

<https://stats.stackexchange.com/questions/31177/does-the-variance-of-a-sum-equal-the-sum-of-the-variances>

To see this let X_1, \dots, X_n be random variables (with finite variances). Then,

$$\text{var} \left(\sum_{i=1}^n X_i \right) = E \left(\left[\sum_{i=1}^n X_i \right]^2 \right) - \left[E \left(\sum_{i=1}^n X_i \right) \right]^2$$

Now note that $(\sum_{i=1}^n a_i)^2 = \sum_{i=1}^n \sum_{j=1}^n a_i a_j$, which is clear if you think about what you're doing when you calculate $(a_1 + \dots + a_n) \cdot (a_1 + \dots + a_n)$ by hand. Therefore,

$$E \left(\left[\sum_{i=1}^n X_i \right]^2 \right) = E \left(\sum_{i=1}^n \sum_{j=1}^n X_i X_j \right) = \sum_{i=1}^n \sum_{j=1}^n E(X_i X_j)$$

similarly,

$$\left[E \left(\sum_{i=1}^n X_i \right) \right]^2 = \left[\sum_{i=1}^n E(X_i) \right]^2 = \sum_{i=1}^n \sum_{j=1}^n E(X_i) E(X_j)$$

so

$$\begin{aligned}\text{var} \left(\sum_{i=1}^n X_i \right) &= \sum_{i=1}^n \sum_{j=1}^n (E(X_i X_j) - E(X_i) E(X_j)) = \sum_{i=1}^n \\ &\quad \sum_{j=1}^n \text{cov}(X_i, X_j)\end{aligned}$$

To prove the initial statement, it suffices to show that

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y).$$

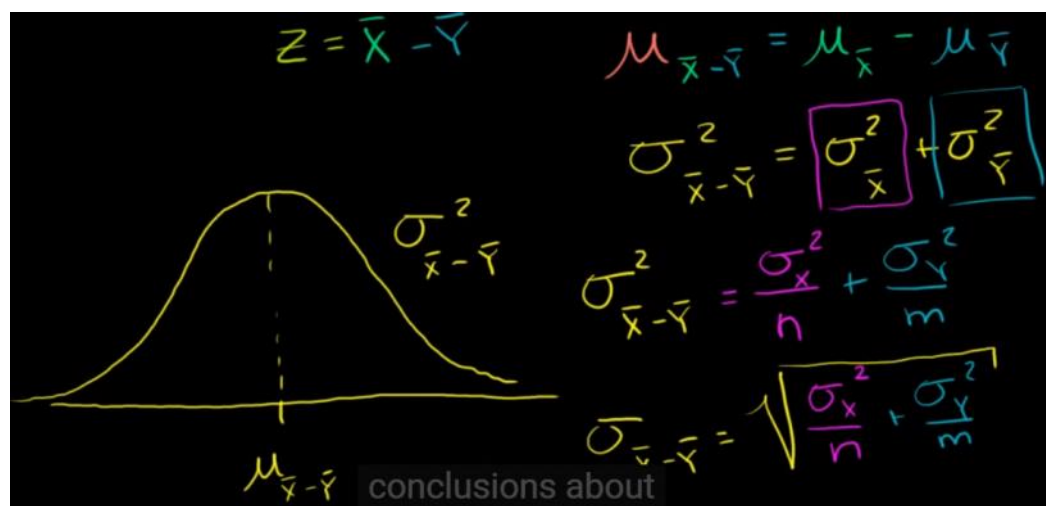
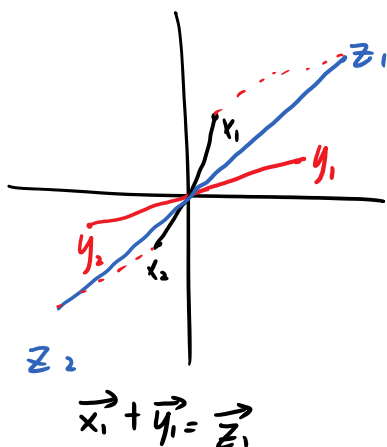
The general result then follows by induction. Starting with the definition,

$$\begin{aligned}\text{Var}(X + Y) &= \text{E}[(X + Y)^2] - (\text{E}[X + Y])^2 \\ &= \text{E}[X^2 + 2XY + Y^2] - (\text{E}[X] + \text{E}[Y])^2.\end{aligned}$$

Using the linearity of the expectation operator and the assumption of independence (or uncorrelatedness) of X and Y , this further simplifies as follows:

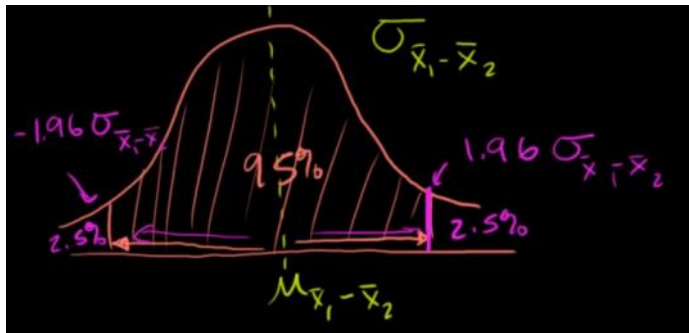
$$\begin{aligned}\text{Var}(X + Y) &= \text{E}[X^2] + 2\text{E}[XY] + \text{E}[Y^2] - (\text{E}[X]^2 + 2\text{E}[X]\text{E}[Y] + \text{E}[Y]^2) \\ &= \text{E}[X^2] + \text{E}[Y^2] - \text{E}[X]^2 - \text{E}[Y]^2 \\ &= \text{Var}(X) + \text{Var}(Y).\end{aligned}$$

The square root of variance is the standard deviation



The standard deviation looks like a distance formula...

We're trying to test whether a new, low-fat diet actually helps obese people lose weight. 100 randomly assigned obese people are assigned to group 1 and put on the low fat diet. Another 100 randomly assigned obese people are assigned to group 2 and put on a diet of approximately the same amount of food, but not as low in fat. After 4 months, the mean weight loss was 9.31 lbs. for group 1 ($s=4.67$) and 7.40 lbs. ($s=4.04$) for group 2.



$$\begin{aligned}\sigma_{\bar{x}_1 - \bar{x}_2}^2 &= \sigma_{\bar{x}_1}^2 + \sigma_{\bar{x}_2}^2 \\ &= \frac{\sigma_{x_1}^2}{100} + \frac{\sigma_{x_2}^2}{100} \\ \sigma_{\bar{x}_1 - \bar{x}_2} &= \sqrt{\frac{s_1^2}{100} + \frac{s_2^2}{100}} = 0.617\end{aligned}$$

$$1.91 \pm 1.21$$

$$0.7 - 3.12$$

We are 95% confident the differences between the two groups (normal diet and low fat diet) is between 0.7 and 3.21. Which indicates that the low fat diet does do something on weight diet, because the lower limit is also positive.

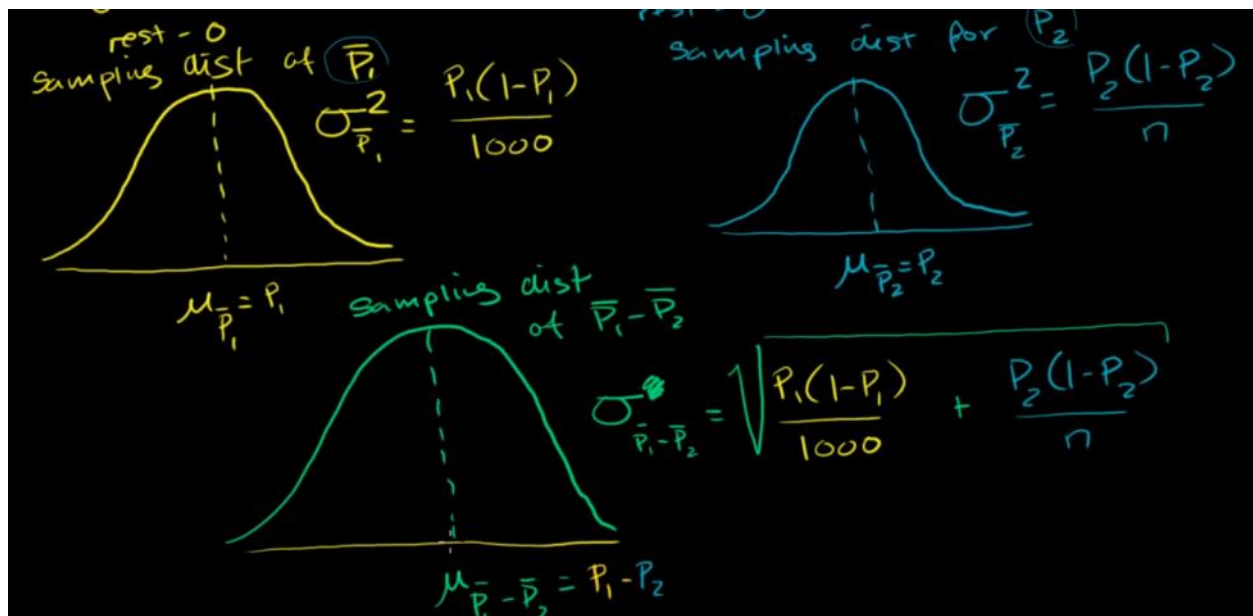
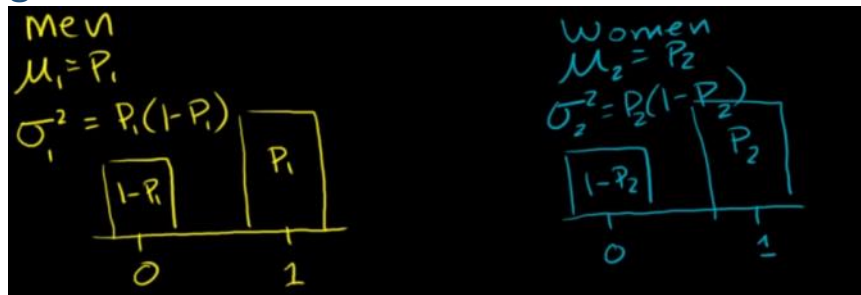
We can also do this using hypothesis test

$$H_0: \mu_1 - \mu_2 = 0$$

$$H_1: \mu_1 - \mu_2 \neq 0$$

Example :

Look for whether there exists difference on the voting rates of different gender



The 95% confidence interval for $p_1 - p_2$ is 0.008 - 0.094

This problem can also be solved by hypothesis test