

















Variance and Standard Deviation

The variance provides a scalar representing the spread of the data set. In a data set X={x₁,x₂,...x_N} an unbiased estimate s_u² for the variance can be calculated as

$$s_u^2 = \frac{\sum_{i=1}^N (x_i - m)^2}{N - 1}$$

• N-1 is often called the number of degrees of freedom of the data set

• The standard deviation s is the square root of the variance

$$s = \sqrt{s^2} = \sqrt{\frac{\sum_{i=1}^{N} (x_i - m)^2}{N - 1}}$$

- · In the case of a sample set, s is often referred to as standard error
- Unlike the variance, the standard deviation is measured in the same units the original data was measured in.

9

<image><image><image><image><image><image><image><image><section-header><image><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header>







































- The goal of hypothesis testing is to establish the viability of a hypothesis about a parameter of the population (very often the mean)
- Recall our example: Is the Suprouter delay on the average higher than Riporouter's?
 - Riporouter characteristics are: μ_R =170ms σ_R =1.2ms
 - Suprouter was measured: {167,174,168,180,173,182,160,170,174,165}
- Define hypothesis (also called *alternative hypothesis*): $H_A: \mu_S > \mu_R$
- Set up **null hypothesis** (i.e. the "opposite" of H) H₀: $\mu_{s} = \mu_{R}$
- Compute the percentile and thus the likelihood H₀
 - If H₀ has more than a small likelihood then the data does not significantly support H_A (since the data could also represent H₀)
 - What are small likelihoods? Usually thresholds (*significance levels*) of 5% or smaller are used.





· If we are not sure we should choose a two-tailed test (which is more stringent)



























- We have looked at setting a statistical relationship between the sample average and the mean of the underlying random variable.
- We could do this as we knew the distribution of the average (based on CLT for z or based on normals adding to normal for t).
- If we knew the sampling distribution of other statistics we could do the same for them as well...
- Variance (arguably the second most important value of a sample set) comes to mind...

















