# Harold's Descriptive Statistics
## Cheat Sheet
21 October 2024
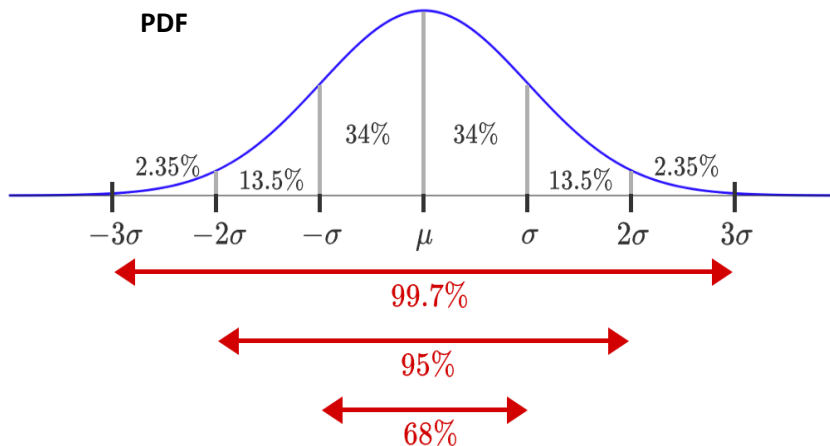
## Descriptive

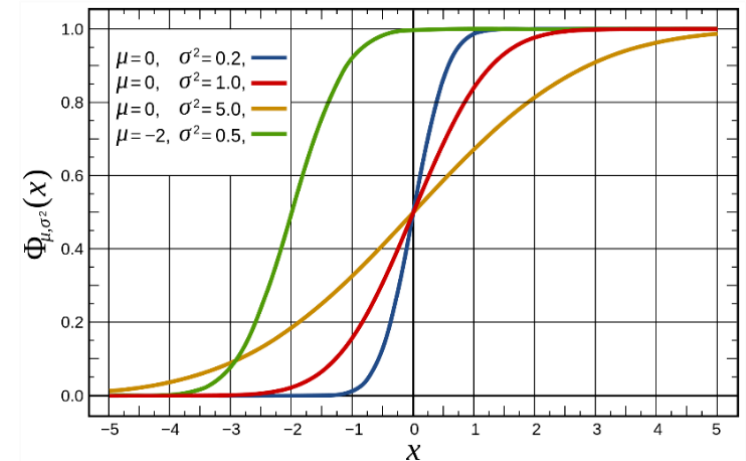| Description | Population | Sample | Used For |
|---|---|---|---|
| **Data** | Parameters | Statistics | Describing and predicting. |
| **Random Variable** | $X, Y$ | $x, y$ | The random value from the evaluated population. |
| **Size** | $N$ | $n$ | Number of observations in the population / sample. |

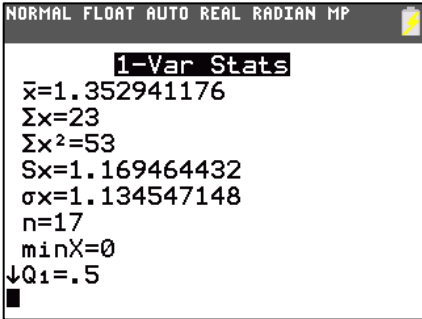| **Measures of Center** | | (Measure of central tendency) | Indicates which value is typical for the data set. |
|---|---|---|---|
| **Mean** | $$\mu = \frac{1}{N}\sum_{i=1}^{N} x_i\, f$$ $f = 1\ if\ samples\ are\ unordered$ | $$\overline{x} = \frac{1}{n}\sum_{i=1}^{n} x_i\, f$$ $$n = \sum f$$ | Measure of center for unordered and frequency distributions. Average, arithmetic mean. Used when same probabilities for each X. Answers "*Where is the center of the data located?*" |
| **Weighted Mean** | $$\mu = \frac{\sum a_i\, x_i}{\sum a_i}$$ | $$\overline{x} = \frac{\sum a_i\, x_i}{\sum a_i}$$ | Some values are counted more than once. $a_i$ = positive integer or percentage. |
| **Median** | $$Md = \frac{n+1}{2}\ if\ n\ is\ odd$$ | $$Md = \frac{n}{2} + 1\ if\ n\ is\ even$$ | The middle element in a <u>sorted</u> dataset. More useful when data are skewed with outliers. |
| **Mode** | $Mo = \max(f)$ | Appropriate for categorical data. | The most frequently-occuring value in a dataset. |
| **Mid-Range** | $$MidRange = \frac{max. + min.}{2}$$ | Not often used, easy to compute. | Highly sensitive to unusual values. |
| **Python** | ```import pandas as pd``` ```data = pd.read_csv('file.csv')``` ```print(data.mean())``` ```print(data[['Header1']].median())``` ```print(data[['Header1', 'Header2']].mode())``` ```mid_range = (data.min() + data.max()) / 2.0``` | | |

| Description | Population | Sample | Used For |
|---|---|---|---|
| **Measures of Dispersion** | | (Measure of dispersion, variability, or spread of the distribution) | Reflect the variability of the data (e.g. how different the values are from each other. |
| **Variance** | $\sigma^2 = \dfrac{1}{N}\sum_{i=1}^{N}(x_i - \mu)^2 f$ $\sigma^2 = \dfrac{1}{N}\left(\sum_{i=1}^{N} f\,x_i^2 - N\mu^2\right)$ | $s^2 = \dfrac{1}{n-1}\sum_{i=1}^{n}(x_i - \bar{x})^2 f$ $s^2 = \dfrac{1}{n-1}\left(\sum_{i=1}^{n} f\,x_i^2 - n\bar{x}^2\right)$ | The average of the sum of the square differences. Not often used. See standard deviation. Special case of covariance when the two variables are identical. |
| **Covariance** | $\sigma(X,Y) = \dfrac{1}{N}\sum_{i=1}^{N}(x - \mu_x)(y - \mu_y)$ $\sigma(X,Y) = \dfrac{1}{N}\sum_{i=1}^{N} x_i y_i - \mu_x \mu_y$ | $g = \dfrac{1}{n-1}\sum(x - \bar{x})(y - \bar{y})$ $\sigma(x,y) = \dfrac{1}{n-1}\left(\sum_{i=1}^{n} x_i y_i - n\bar{x}\bar{y}\right)$ | A measure of how much two random variables change together. Measure of "linear dependence". If X and Y are independent, then their covariance is zero (0). |
| **Standard Deviation** | $\sigma = \sqrt{\sigma^2} = \sqrt{\dfrac{\sum(x_i - \mu)^2}{N}}$ $\sigma = \sqrt{\dfrac{\sum x_i^2}{N} - \mu^2}$ | $s_x = \sqrt{\dfrac{\sum(x_i - \bar{x})^2}{n-1}}$ $s = \sqrt{\dfrac{\sum x_i^2 - n\bar{x}^2}{n-1}}$ | Measure of variation; average distance from the mean. Same units as mean. Answers "*How spread out is the data?*" |
| **Mean Absolute Deviation** | $MAD = \dfrac{1}{N}\sum|x_i - \mu|$ | $MAD = \dfrac{1}{n}\sum|x_i - \bar{x}|$ | Uses the absolute value instead of the square root of a sum of squares to avoid negative distances. |
| **Pooled Standard Deviation** | $\sigma_p = \sqrt{\dfrac{N_1\sigma_1^2 + N_2\sigma_2^2}{N_1 + N_2}}$ | $s_p = \sqrt{\dfrac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1-1)+(n_2-1)}}$ | Inferences for two population means. |
| **Interquartile Range (IQR)** | $IQR = Q3 - Q1$ |  | Less sensitive to extreme values. |
| **Range** | $Range = max. - min.$ | Not often used, easy to compute. | Highly sensitive to unusual values. |
| **Python** | ```<br>import pandas as pd<br>data = pd.read_csv('file.csv')<br>print(data.var())<br>print(data.cov())<br>print(data.std())<br>``` | ```<br>print(data.mad())<br>def IQR(data):     # (import numpy as np)<br>    Q3 = np.quantile(data, 0.75)<br>    Q1 = np.quantile(data, 0.25)<br>    IQR = Q3 - Q1<br>range = data.max() - data.min()<br>``` | |

| Description | Population | Sample | Used For |
|---|---|---|---|
| **Measures of Relative Standing** | | (Measures of relative position) | Indicates how a particular value compares to the others in the same data set. |
| Percentile | Data divided onto 100 equal parts by rank. | | Important in normal distributions. |
| Quartile | Data divided onto 4 equal parts by rank. | | Used to compute IQR. |
| **Z-Score / Standard Score / Normal Score** | $x = \mu + z\,\sigma$ <br><br> $z = \dfrac{x - \mu}{\sigma}$ | $x = \bar{x} + z\,s$ <br><br> $z = \dfrac{x - \bar{x}}{s}$ | The $z$ variable measures how many standard deviations the value is away from the mean. Rule of Thumb: Outlier if $\lvert z \rvert > 2$. |
| **Calculator (TI-84)** | [2nd][VARS][2] normalcdf(-1E99, z) | | |
| **Python** | ```import scipy.stats as st

mean, sd, z = 0, 1, 1.5
print(st.norm.cdf(z, mean, sd))        # P(z <= 1.5)
print(st.norm.sf(z, mean, sd))         # P(z >= 1.5)

mean, sd, x = 55, 7.5, 62
print(st.norm.cdf(x, mean, sd))        # P(x <= 62)
print(st.norm.sf(x, mean, sd))         # P(x >= 62)``` | | ```0.9331927987311419
0.06680720126885780
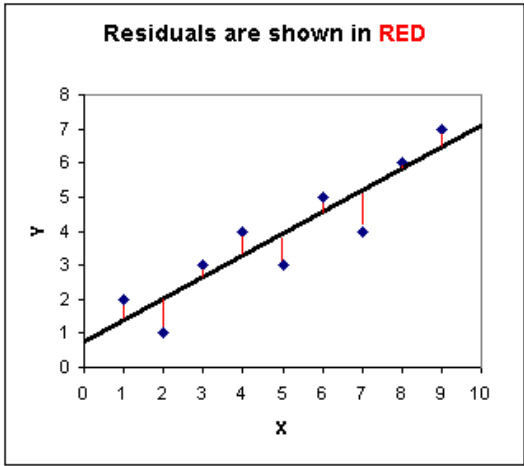
0.8246760551477705
0.1753239448522295``` |



PDF



CDF

| Example | Data | Method | Results |
|---|---|---|---|
| **Example** | | | |
| **Data** | *Unordered Data: 1, 0, 1, 4, 1, 2, 0, 3, 0, 2, 1, 1, 2, 0, 1, 1, 3* | | Relative Frequency: $\quad p(x) = f/n$ |

**Manually**

*Ordered Data:*

| $x$ | $f$ |
|---|---|
| 0 | 4 |
| 1 | 7 |
| 2 | 3 |
| 3 | 2 |
| 4 | 1 |

| $x$ | $f$ | $f/n$ | $x - \bar{x}$ | $(x - \bar{x})^2$ | $(x - \bar{x})^2 f$ |
|---|---|---|---|---|---|
| 0 | 4 | 0.35 | -1.35 | 1.83 | 7.32 |
| 1 | 7 | 0.41 | -0.35 | 0.12 | 0.87 |
| 2 | 3 | 0.18 | 0.65 | 0.42 | 1.26 |
| 3 | 2 | 0.12 | 1.65 | 2.71 | 5.43 |
| 4 | 1 | 0.06 | 2.65 | 7.01 | 7.01 |

$$n = \sum f = 4 + 7 + 3 + 2 + 1 = 17$$

$$\bar{x} = \frac{1}{n}\sum x_i\, f = \frac{(0 \cdot 4) + \cdots + (4 \cdot 1)}{17} = \frac{23}{17} \approx 1.35$$

$$\sigma^2 = \frac{1}{n}\sum (x - \bar{x})^2 f = \frac{7.32 + \cdots + 7.01}{17} \approx 1.21$$

$$\sigma = \sqrt{\sigma^2} \approx 1.13$$

**Calculator (TI-84)**

1. [STAT] [1] selects the list edit screen
2. Move the cursor up to L1
3. [CLEAR] [ENTER] erases L1
4. Repeat for L2
5. Enter $x$ data in L1 and $f$ data in L2
6. [STAT] → [1] to select 1-Var Stats
7. [2nd] [1] [ENTER] for L1
8. [2nd] [2] [ENTER] for L2
9. Calculate [ENTER]

**Output:**



```
NORMAL FLOAT AUTO REAL RADIAN MP
            1-Var Stats
x̄=1.352941176
Σx=23
Σx²=53
Sx=1.169464432
σx=1.134547148
n=17
minX=0
↓Q₁=.5
```

**Python**

```python
import pandas as pd
df = pd.DataFrame(
[1,0,1,4,1,2,0,3,0,2,1,1,2,0,1,1,3])
print(df.describe())
print()
print(df.std(ddof=1))    # Sample SD
print(df.std(ddof=0))    # Population SD
========================================================
from scipy.stats import rv_discrete
x = [0,1,2,3,4,5,6]                  # Outcomes
p = [0.1,0.2,0.3,0.1,0.1,0.0,0.2] # Prob of outcomes
discrete_var = rv_discrete(values=(x,p)) # Links x2p
print(discrete_var.mean())
print(discrete_var.std())
```

**Output:**

```
count   17.000000
mean     1.352941
std      1.169464
min      0.000000
25%      1.000000
50%      1.000000
75%      2.000000
max      4.000000


std1     1.169464
std0     1.134547
```

## Regression and Correlation

| Description | Formula | Used For |
|---|---|---|
| Response Variable | $Y$ | Output |
| Covariate / Predictor Variable | $X$ | Input |
| Least-Squares Regression Line | $\hat{y} = b_0 + b_1 x$ | $b_1$ is the slope<br>$b_0$ is the y-intercept<br>$(\bar{x}, \bar{y})$ is always a point on the line |
| Regression Coefficient (Slope) | $b_1 = \dfrac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x - \bar{x})^2}$<br><br>$b_1 = r\dfrac{s_y}{s_x}$ | $b_1$ is the slope |
| Regression Slope Intercept | $b_0 = \bar{y} - b_1 \bar{x}$ | $b_0$ is the y-intercept |
| Linear Correlation Coefficient (Sample) | $r = \dfrac{1}{n-1}\sum\left(\dfrac{x - \bar{x}}{s_x}\right)\left(\dfrac{y - \bar{y}}{s_y}\right)$<br><br>$r = \dfrac{g}{s_x s_y}$ | Strength and direction of linear relationship between x and y.<br><br>$r = \pm 1$ Perfect correlation<br>$r = +0.9$ Positive linear relationship<br>$r = -0.9$ Negative linear relationship<br>$r = \sim 0$ No relationship<br>$r \geq 0.8$ Strong correlation<br>$r \leq 0.5$ Weak correlation<br><br>Correlation DOES NOT imply causation. |
| Residual | $\hat{e}_i = y_i - \hat{y}$<br>$\hat{e}_i = y_i - (b_0 + b_1 x)$<br>$\sum e_i = \sum(y_i - \hat{y}_i) = 0$ | Residual = Observed – Predicted |
| Standard Error of Regression Slope | $s_{b_1} = \dfrac{\sqrt{\dfrac{\sum e_i^2}{n-2}}}{\sqrt{\sum(x_i - \bar{x})^2}}$<br><br>$s_{b_1} = \dfrac{\sqrt{\dfrac{\sum(y_i - \hat{y}_i)^2}{n-2}}}{\sqrt{\sum(x_i - \bar{x})^2}}$ |  |
| Coefficient of Determination | $r^2$ | How well the line fits the data.<br><br>Represents the percent of the data that is the closest to the line of best fit. Determines how certain we can be in making predictions. |

## Proportions

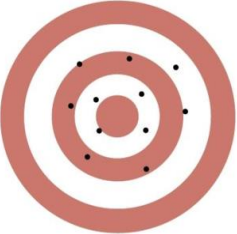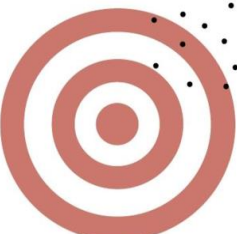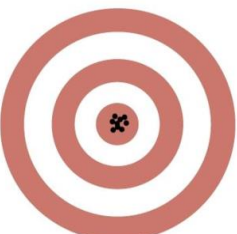| Description | Population | Sample | Used For |
|---|---|---|---|
| **Proportion** | $P = p = \dfrac{x}{N}$ | $\hat{p} = \dfrac{x}{n}$ | Probability of **success**. The proportion of elements that has a particular attribute (x). |
| | $q = 1 - p$ <br> $Q = 1 - P$ | $\hat{q} = 1 - \hat{p}$ | Probability of **failure**. The proportion of elements in the population that does not have a specified attribute. |
| **Variance of Population (Sample Proportion)** | $\sigma^2 = \dfrac{pq}{N}$ <br><br> $\sigma^2 = \dfrac{p(1-p)}{N}$ | $s_p^2 = \dfrac{\hat{p}\hat{q}}{n-1}$ <br><br> $s_p^2 = \dfrac{\hat{p}(1-\hat{p})}{n-1}$ | Considered an unbiased estimate of the true population or sample variance. |
| **Pooled Proportion** | *NA* | $\hat{p}_p = \dfrac{x_1 + x_2}{n_1 + n_2}$ <br><br> $\hat{p}_p = \dfrac{\hat{p}_1 n_1 + \hat{p}_2 n_2}{n_1 + n_2}$ | $x = \hat{p}n =$ frequency, or number of members in the sample that have the specified attribute. |

# Discrete Random Variables

| Description | Formula | Used For |
|---|---|---|
| **Random Variable** | $X$ | A rule that assigns a number to every **outcome** in the sample space, S.<br>$e.g., X(a,b) = a + b = r$<br>Derived from a probability experiment with different probabilities for each X.<br>**Used in discrete or finite PDFs.** |
| **Event** | $X = r$<br>$X(s) = r$ | An event assigns a value to the random variable X with probability:<br>$P(X = r)$ |
| **Expected Value of *X*** | $E[X] = \bar{x} \text{ or } \mu_x$<br>Each event:<br>$E[X] = \sum P(X) \cdot X$<br>$E[X] = \sum_{s \in S} X(s) \cdot P(s)$<br>Groups of like events:<br>$E[X] = \sum_{i=1}^{N} p_i(x) \cdot x_i$<br>$\boldsymbol{E[X]} = \sum_{r \in X(S)} \boldsymbol{r \cdot P(X = r)}$ | E(X) is the same as the mean or average. X takes some countable number of specific values. Discrete.<br>Expectation of a random variable.<br>*P(s)* = probability of outcome *s* from *S*. |
| **Linearity of Expectations** | $E[X + Y] = E[X] + E[Y]$<br>$E[X + Y + Z] = E[X] + E[Y] + E[Z]$<br><br>$E[cX] = cE[X]$ | When carefully applied, linearity of expectations can greatly simplify calculating expectations.<br><br>Does not require that the random variables be independent. |
| **Variance of *X*** | $\boldsymbol{V(X) = \sigma_x^2 = \sum p_i(x) \cdot (x_i - \mu_x)^2}$<br>$\sigma_x^2 = \sum P(X) \cdot (X - E[X])^2$<br>$\sigma_x^2 = \sum X^2 \cdot P(X) - E[X]^2$<br><br>$\sigma_x^2 = E[X^2] - E[X]^2$ | Calculate variances with proportions or expected values. |
| **Standard Deviation of *X*** | $SD(X) = \sqrt{V(X)}$<br><br>$\sigma_x = \sqrt{\sigma_x^2}$ | Calculate standard deviations with proportions. |
| **Sum of Probabilities** | $\sum_{i=1}^{N} p_i(x) = 1$ | If same probability, then $p_i(x) = \frac{1}{N}$. |

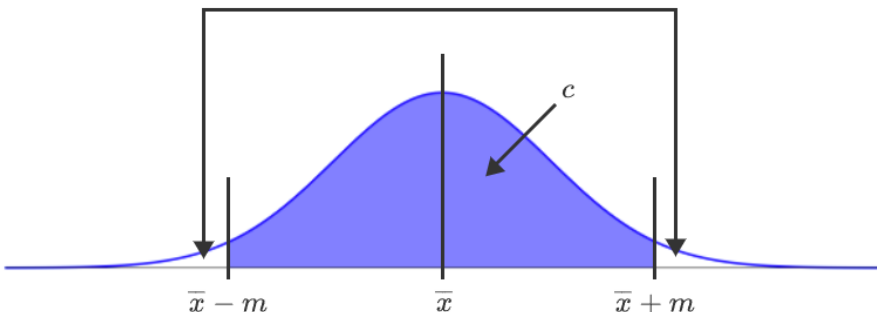NOTE: See also "Discrete Definitions" on Harold's_Stats_Distributions_Cheat_Sheet.

## Sampling Distribution Statistical Inference

| Description | Mean | Standard Deviation |
|---|---|---|
| **Sampling Distribution** | Is the probability distribution of a statistic; a statistic of a statistic. | |
| **Central Limit Theorem (CLT)** | $PDF(\bar{x}) \approx \mathcal{N}\left(0, \dfrac{\sigma^2}{n}\right)$ | As the sample size drawn from the population with distribution **X** becomes larger, the sampling distribution of the means $\overline{X}$ approaches that of a normal distribution $\mathcal{N}\left(0, \dfrac{\sigma^2}{n}\right)$. |
| **Sample Mean** | $\mu_{\bar{x}} = \mu$ | Sampling with replacement: $$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$ Sampling without replacement: $$\sigma_{\bar{x}} = \sqrt{\frac{N-n}{N-1}} \cdot \frac{\sigma}{\sqrt{n}}$$ (2x accuracy needs 4x n) |
| z-Score | $z = \dfrac{\bar{x} - \mu_{\bar{x}}}{\sigma_{\bar{x}}}$ | $z = \dfrac{\bar{x} - \mu}{\sigma / \sqrt{n}}$ |
| Sample Mean Rule of Thumb | Use if $n \geq 30$ **or** if the population distribution is normal | |
| 10% Condition | $n \leq \dfrac{N}{10}$. Sample size must be at most 10% of the population size. | |
| **Sample Proportion** | $\mu = p$ | $\sigma_p = \sqrt{\dfrac{p(1-p)}{n}}$ |
| z-Score | $z = \dfrac{\hat{p} - \mu}{\sigma_p}$ | $z = \dfrac{\hat{p} - p}{\sqrt{\dfrac{p(1-p)}{n}}}$ |
| Sample Proportion Rule of Thumb | Large Counts Condition: Use if $np \geq 5$ **and** $n(1-p) \geq 5$ Use if $np \geq 10$ **and** $n(1-p) \geq 10$ | 10 Percent Condition: Use if $N \geq 10n$ |
| **Difference of Sample Means** | $E(\bar{x}_1 - \bar{x}_2) = \mu_{\bar{x}_1} - \mu_{\bar{x}_2}$ | $\sigma_{\bar{x}_1 - \bar{x}_2} = \sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}$ |
| Special case when $\sigma_1 = \sigma_2$ | | $\sigma_{\bar{x}_1 - \bar{x}_2} = \sigma\sqrt{\dfrac{1}{n_1} + \dfrac{1}{n_2}}$ |
| **Difference of Sample Proportions** | $\Delta\hat{p} = \hat{p}_1 - \hat{p}_2$ | $\sigma = \sqrt{\dfrac{p_1(1-p_1)}{n_1} + \dfrac{p_2(1-p_2)}{n_2}}$ |
| Special case when $p_1 = p_2$ | | $\sigma = \sqrt{p(1-p)}\sqrt{\dfrac{1}{n_1} + \dfrac{1}{n_2}}$ |

| | | |
|---|---|---|
| **Bias** | Caused by non-random samples.<br><br>Selection Bias: Under coverage, Nonresponse, Voluntary response<br>Response Bias: Acquiescence, Extreme, Social desirability | High bias, low variability (a)    Low bias, high variability (b) |
| **Variability** | Caused by too small of a sample.<br>$n < 30$<br><br>Sampling Methods:<br>Simple random, systematic, stratified, cluster, convenience | High bias, high variability (c)    The ideal: low bias, low variability (d) |

# Confidence Intervals for One Population Mean / Proportion (σ is Known)

| Description | Formula |
|---|---|
| **Critical Value (z\*)** | Usually set ahead of time, unless using p-values to determine. <br><br> Set to a threshold value of 0.05 (5%) or 0.01 (1%), but always ≤ 0.10 (10%). <table><tr><td>**Confidence Level**</td><td>**Critical Value**</td></tr><tr><td>c = 0.90</td><td>z\* = 1.645</td></tr><tr><td>c = 0.95</td><td>z\* = 1.960</td></tr><tr><td>c = 0.99</td><td>z\* = 2.576</td></tr></table> |
| **p-value** | Probability of obtaining a sample "more extreme" than the ones observed in your data, assuming $H_0$ is true. |
| **Sample Size** <br> (for estimating μ) | $$n = \left(\frac{z^*\sigma}{SE}\right)^2 = \left(\frac{z^*}{SE}\right)^2 p(1-p)$$ <br> The size of the sample needed to guarantee a confidence interval with a specified margin of error. Rounded up to the nearest whole number. |
| **Margin of Error / Standard Error (SE)** <br> (for the estimate of μ) | $$SE(\bar{x}) = m = z^*\frac{\sigma}{\sqrt{n}} = z^*\sqrt{\frac{p(1-p)}{n}}$$ <br> The estimate $\bar{x}$ differs from the actual value by at most SE. <br> Use p = 0.50 for worst case if no previous estimate is known. <br><br> SE with replacement: $$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \sqrt{\frac{p(1-p)}{n}}$$     SE without replacement (with correction factor): $$\sigma_{\bar{x}} = \sqrt{\frac{N-n}{N-1}} \cdot \frac{\sigma}{\sqrt{n}}$$ |
| **Confidence Interval for μ** <br> **(z interval)** <br> (σ known, normal population or large sample) | $z\ interval = statistic \pm (critical\ value) \bullet (SD\ of\ statistic)$ <br> $z\ interval = \bar{x} \pm SE(\bar{x})$ <br> $\bar{x} \pm m = [\bar{x} - m,\ \bar{x} + m]$ <br><br> $$z\ interval = \bar{x} \pm z^*\frac{\sigma}{\sqrt{n}} = \bar{x} \pm z^*\sqrt{\frac{p(1-p)}{n}}$$  <br> $$\frac{\alpha}{2} = \frac{1-c}{2}$$ <br> $z^* = z\ score\ for\ probabilities\ of\ ^{\alpha}/_2\ (two-tailed)$ |
| **Standardized Test Statistic** <br> (of the variable $\bar{x}$ from the CLT) | $$z = \frac{statistic - parameter}{SD\ of\ statistic}$$ <br> $$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$ |

# Confidence Intervals for One Population Mean / Proportion (σ is Unknown)

| Description | Formula |
|---|---|

| Description | Formula |
|---|---|
| **Critical Value (t\*)** | Usually set ahead of time, unless using p-values to determine. *df = n-1*.<br><br>Set to a threshold value of 0.05 (5%) or 0.01 (1%), but always ≤ 0.10 (10%). |

| df | α = 0.10 | α = 0.05 | α = 0.01 |
|---|---|---|---|
| 5 | 2.015 | 2.571 | 4.032 |
| 10 | 1.812 | 2.225 | 3.169 |
| 15 | 1.753 | 2.131 | 2.947 |
| 24 | 1.711 | 2.064 | 2.797 |
| 32 | 1.309 | 1.694 | 2.449 |

| Description | Formula |
|---|---|
| **p-value** | Probability of obtaining a sample "more extreme" than the ones observed in your data, assuming $H_0$ is true. |
| **Sample Size** (for estimating μ) | Preliminary estimate of n:<br><br>$$n^* = \left(\frac{z^* s}{SE}\right)^2$$<br><br>Actual sample size, n:<br><br>$$n = \left(\frac{t^* s}{SE}\right)^2$$<br><br>The size of the sample needed to guarantee a confidence interval with a specified margin of error. Rounded up to the nearest whole number. |
| **Margin of Error / Standard Error (SE)** (for the estimate of μ) | $$SE(\bar{x}) = m = t^* \frac{s}{\sqrt{n}}$$<br>The estimate $\bar{x}$ differs from the actual value by at most SE.<br><br>SE with replacement: $$s_{\bar{x}} = \frac{s}{\sqrt{n}}$$     SE without replacement (with correction factor): $$s_{\bar{x}} = \sqrt{\frac{N-n}{N-1}} \cdot \frac{s}{\sqrt{n}}$$ |
| **Confidence Interval for μ (t interval)** (σ unknown, t distribution or small sample) | $$t\ interval = statistic \pm (critical\ value) \bullet (SD\ of\ statistic)$$<br><br>$$t\ interval = \bar{x} \pm SE(\bar{x})$$<br>$$\bar{x} \pm m = [\bar{x} - m,\ \bar{x} + m]$$<br>$$\boldsymbol{t\ interval = \bar{x} \pm t^* \frac{s}{\sqrt{n}}}$$<br>$\alpha$<br><br>$c$<br><br>$\bar{x} - m$     $\bar{x}$     $\bar{x} + m$ |
| **Standardized Test Statistic** (of the variable $\bar{x}$ from the CLT) | $$t = \frac{statistic - parameter}{SD\ of\ statistic}$$<br><br>$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$ |

# Confidence Intervals for the Difference Between Two Population Means / Proportions (σ is Known)

| Description | Formula |
|---|---|
| **Critical Value (z\*)** | Usually set ahead of time, unless using p-values to determine.<br><br>Set to a threshold value of 0.05 (5%) or 0.01 (1%), but always ≤ 0.10 (10%).<br><br><table><tr><td>**Confidence Level**</td><td>**Critical Value**</td></tr><tr><td>c = 0.90</td><td>z\* = 1.645</td></tr><tr><td>c = 0.95</td><td>z\* = 1.960</td></tr><tr><td>c = 0.99</td><td>z\* = 2.576</td></tr></table> |
| **p-value** | TI-84: DISTR 2: normalcdf(z_test, 99999999) = p |
| **Margin of Error / Standard Error (SE)** (for the estimate of μ) | $$E(\bar{x}_1 - \bar{x}_2) = \mu_{\bar{x}_1} - \mu_{\bar{x}_2}$$ $$SE(\bar{x}_1 - \bar{x}_2) = \sqrt{SE_1^2 + SE_2^2} = m$$ $$= \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} = \sqrt{\hat{p}(1-\hat{p})}\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$ $\hat{p}$ = Overall probability of success when the two samples are combined.<br>The estimate $\bar{x}_1 - \bar{x}_2$ differs from the actual value by at most SE.<br>Use p = 0.50 for worst case if no previous estimate is known. |
| **Confidence Interval for μ (z interval)** (σ known, normal population or large sample) | $$z\ interval = statistic \pm (critical\ value) \bullet (SD\ of\ statistic)$$ $$z\ interval = (\bar{x}_1 - \bar{x}_2) \pm SE(\bar{x}_1 - \bar{x}_2)$$ $$(\bar{x}_1 - \bar{x}_2) \pm m = [(\bar{x}_1 - \bar{x}_2) - m,\ (\bar{x}_1 - \bar{x}_2) + m]$$ $$\boldsymbol{z\ interval = (\bar{x}_1 - \bar{x}_2) \pm z^* \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$ $$\boldsymbol{z\ interval = (\bar{x}_1 - \bar{x}_2) \pm z^* \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}}$$ $$\frac{\alpha}{2} = \frac{1-c}{2}$$ $$z^* = z\ score\ for\ probabilities\ of\ ^{\alpha}/_2\ (two-tailed)$$ |
| **Standardized Test Statistic** (of the variable $\bar{x}$ from the CLT) | $$z = \frac{observed\ difference - hypothesided\ difference}{SD\ for\ the\ difference}$$ $$z = \frac{(\bar{x}_1 - \bar{x}_2) - 0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$ |
| **Python** | ```
from statsmodels.stats.weightstats import ztest
sample1 = [21, 28, 40, 55, 58, 60]
sample2 = [13, 29, 50, 55, 71, 90]
print(ztest(x1 = sample1, x2 = sample2, value = 0))
```<br>(-0.58017, 0.56179)<br>z-score = -0.5802<br>p-value = 0.5618<br>(two-tailed) |

## Confidence Intervals for the Difference Between Two Population Means / Proportions (σ is Unknown)

| Description | Formula | | | | |
|---|---|---|---|---|---|
| **Critical Value (t\*)** | Usually set ahead of time, unless using p-values to determine. *df = n-1*.<br><br>Set to a threshold value of 0.05 (5%) or 0.01 (1%), but always ≤ 0.10 (10%). | *df* | **α = 0.10** | **α = 0.05** | **α = 0.01** |
| | | **5** | 2.015 | 2.571 | 4.032 |
| | | **10** | 1.812 | 2.225 | 3.169 |
| | | **15** | 1.753 | 2.131 | 2.947 |
| | | **24** | 1.711 | 2.064 | 2.797 |
| | | **32** | 1.309 | 1.694 | 2.449 |

| Description | Formula |
|---|---|
| **p-value** | TI-84: DISTR 6: tcdf(t_test, 99999999) = p |
| **Margin of Error / Standard Error (SE)**<br>(for the estimate of μ) | $$SE(\bar{x}_1 - \bar{x}_2) = m = \frac{s_d}{\sqrt{n}}$$<br>The estimate $\bar{x}_1 - \bar{x}_2$ differs from the actual value by at most SE. |
| **Confidence Interval for μ (t interval)**<br>(σ unknown, t distribution or small sample) | $$t\ interval = statistic \pm (critical\ value) \cdot (SD\ of\ statistic)$$<br><br>$$\boldsymbol{t\ interval = (\bar{x}_1 - \bar{x}_2) \pm \frac{s_d}{\sqrt{n}}}$$ |
| **Standardized Test Statistic**<br>(of the variable $\bar{x}$ from the CLT) | $$t = \frac{mean\ difference\ between\ samples - parameter}{sample\ SD\ of\ the\ differences\ /\ \sqrt{n}}$$<br>Paired t-test:<br>$$t = \frac{\bar{d} - \mu_d}{s_d / \sqrt{n}}$$<br>$$df = n - 1$$<br><br>Unpaired t-test:<br>$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$<br>$$df = n_1 + n_2 - 2$$<br><br>$$t = \frac{\hat{p}_1 - \hat{p}_2}{SE}$$ |

| | | | |
|---|---|---|---|
| **Python** | **Paired** | ```import scipy.stats as st```<br>```import pandas as pd```<br>```df = pd.read_csv('ExamScores.csv')```<br>```st.ttest_rel(df['Exam1'],df['Exam2'])``` | Ttest_relResult (statistic = 1.4179, pvalue = 0.16254) |
| | **Unpaired** | ```import scipy.stats as st```<br>```import pandas as pd```<br>```df = pd.read_csv('Machine.csv')```<br>```st.ttest_ind(df['Old'],df['New'], equal_var=False))``` | Ttest_indResult (statistic = 3.3972, pvalue = 0.00324) |
| | | ```from statsmodels.stats.proportion import proportions_ztest```<br>```counts = [95, 125]```<br>```n = [5000, 5000]```<br>```print(proportions_ztest(counts, n))``` | (statistic = -2.04522, pvalue = 0.04083) |

**Sources:**

- [SNHU MAT-353](#) - Applied Statistics for STEM, zyBooks.

- [SNHU MAT-353](#) - Applied Statistics for STEM, zyBooks.