

Midterm 1

1.1:

Mutually Exclusive: no overlap; $P(A \cap B) = 0$;
Exhaustive: complete coverage

1.2: Sample counting

Ordered: With replacement: n^r ; Without replacement (order important): $nPr = n! / [(n-r)!]$; **Unordered:** $nCr = n! / [(n-r)!r!]$
Distinct permutations of n things, k kinds: $n! / (n_1!n_2!\dots n_k!)$ **1.3 BELOW**

1.4:

$P(A|B) = P(A \cap B) / P(B)$: **conditional** prob that A will occur given that B occurs
 $P(A|B) = N(A \cap B) / N(B)$ equally likely model, propo. of outcomes from B that satisfy criteria for A

$$P(A \cap B) = P(B|A)P(A) = P(A|B)P(B)$$

$$P(B) = P(A \cap B) + P(A' \cap B) = P(B|A)P(A) + P(B|A')P(A')$$

$$P(A' \cap B) = 1 - P(A|B); P(B|A) = P[(A|B)P(B)] / P(A)$$

If E_1, \dots, E_k are M.E. and exhaustive, for any event B and any i: (Bayes')

$$P(E_i|B) = P(B|E_i)P(E_i) / [P(B|E_1)P(E_1) + \dots + P(B|E_k)P(E_k)]$$

independent if: $P(B|A) = P(B)$ or $P(A \cap B \cap C \cap D) = P(A)P(B)P(C)P(D)$

2.1: Random Variables: Sample space S , function $X: S \rightarrow R$ associates a number $X(s)$ to each outcome

pmf: gives the probability for each value of x

2.2 Discrete:

pmf: $f(x) = P(X=x)$; cdf: $F(x) = P(X \leq x) = \sum_{y \leq x} f(y)$
 $\mu = E(X) = \sum_x x f(x)$ = center of distribution mass = mean of X 's probability distribution; To get, $\sum [E(X) * p(x)]$;
(If it's a function $E(g(X))$, multiply that function's value at x by $p(x)$)
 $\sigma^2 = E[(X-\mu)^2] = \sum_x (x-\mu)^2 f(x) = E[X^2] - (E[X])^2 = \sum_x (x^2) f(x) - (\mu)^2$; to get, $\sum [(x-\mu)^2 * p(x)]$; prob multiplied by (distance from mean)²
Prob. dist.: a table with possible x values and their $p(x)$ values

2.3 Linear Transformation

$Y = a + bX$ is the format; $\star E(a + bX) = a + bE(X)$;
 $\star \text{Var}(a + bX) = b^2 \text{Var}(X)$; $\star \text{SD}(a + bX) = |b| \text{SD}(X)$

2.4 Continuous:

pmf becomes **pdf** $f(x) = d/dx [F_x(x)]$;

CDF of any r.v. X : $F_x(x) = P(X \leq x)$

$$P(a < X <= b) = F_x(b) - F_x(a)$$

If we take a linear function $a + bX$ of a c.r.v., formulas \star are the same.

Continuous Random Variable (M2)

pdf $f(x)$: height of the curve at x , probabilities are areas under the curve. Total area = 1. Integrate with an x in a bound logically to get cdf, use to calculate probs on intervals. Prob at one exact point = 0 (area), so must be an interval a to b on cdf.

μ : integrate with extra $*x$ on $[a, b]$: $\int_a^b x f(x) dx$;

σ^2 : integrate with extra $*x^2$ on $[a, b]$, then subtract μ^2 : $E(X^2) - [E(X)]^2$

Properties of pdf $f(x)$:

Single values have no mass, if $a < b$ then:

$$P(a < X < b) = P(a \leq X < b) = P(a < X \leq b) = P(a < X \leq b) = \int_a^b f(x) dx = F(b) - F(a) \leftarrow \text{CDF}$$

2.5 Discrete Joint Distribution:

X and Y are both d.r.v.'s - their joint pmf: $f_{XY}(x, y) = P(X=x, Y=y) = P(X=x \text{ and } Y=y)$;

Range of random vector (X, Y) : $R_{XY} = \{(x, y) : f_{XY}(x, y) \neq 0\}$;

$P(X, Y) \text{ element of } A = \sum_{(x, y) \text{ element of } A} f_{XY}(x, y)$

Marginal pmf of X (same with Y):

$$f_X(x) = P(X=x) = \sum_y f_{XY}(x, y)$$

Independent if $f_{XY}(x, y) = f_X(x)f_Y(y)$ for all x and y

2.6 Covariance and Correlation

X and Y both r.v.'s; expectation of $h(X, Y)$ called $E[h(X, Y)]$ is the ave. expected value w/ repeated trials.

$$E[aX + bY] = aE[X] + bE[Y]; V[X + Y] = V[X] + V[Y] + 2E[(X - \mu_X)(Y - \mu_Y)]$$

$$\text{Covariance between } X \text{ and } Y: \sigma_{XY} = E[(X - \mu_X)(Y - \mu_Y)] = E[XY] - \mu_X \mu_Y$$

Corr. coeff. is between X and Y : $\rho_{XY} = \sigma_{XY} / (\sigma_X \sigma_Y)$; $-1 \leq \rho \leq 1$; if points on (X, Y) are a line then $\rho_{XY} =$ either 1 or -1 (slope +/-); if indie then $\sigma_{XY} = \rho_{XY} = 0$ (reverse not true); if $\neq 0$ then correlated and indie; if $= 0$ then uncorrelated but ? Indie.

3.1 Binomial Dist.: $[X \sim B(n, p)]$

x = number of out of n indie bernie (pass/fail) trials, x has Binomial dist w/ replacement (indie); **Pmf:** $(nC_x)(p^x)(1-p)^{n-x}$ for $x=0, 1, 2, \dots, n$

Cdf: $P(X \leq x) = \sum_{i=0}^x (nC_i)(p^i)(1-p)^{n-i}$; $\mu = np$; $\sigma^2 = np(1-p)$

Bernoulli: Special case with $n=1$, $x=0, 1$;

3.2 Geometric Dist.:

x : #trials needed to get 1st success out of many (indie) bernie trials

Pmf: $(1-p)^{x-1} * p$ for $x=1, 2, 3, \dots$ (p is success and $1-p$ is fail)

$x=1$ is always the most likely, so $P(X > x) = (1-p)^x$ and...

Cdf: $P(X \leq x) = 1 - (1-p)^x$; $\mu = 1/p$; $\sigma^2 = (1-p)/p^2$

3.3 Negative Binomial Dist.:

x : number of indie bern (pass/fail) trials needed to get r^{th} success; $P(\text{success})=p$; $P(\text{fail})=1-p$; x =trial# of the r^{th} success.

Pmf: $({}_{x-1}C_{r-1})(p^r)(1-p)^{x-r}$ for $x=r, r+1, \dots$; $\mu = r/p$; $\sigma^2 = [r(1-p)]/p$

3.4 Poisson Dist.: (# warhorses killed in 30 days)

x : #occurrences of an indie event in a period; rate is constant.

Pmf: $(\lambda^x * e^{-\lambda}) / (x!)$, for $x=0, 1, 2, \dots$; $\mu = \lambda = \sigma^2$;

$\lambda = (\text{events/second}) * (\text{\#seconds in period})$

Poisson Approximation to the Binomial Dist.:

large n and small p ; $\mu = np$;

Pdf $= (nC_x)(p^x)(1-p)^{n-x} \approx e^{-\mu}(\mu^x / (x!))$

good approx if $n \geq 20$ & $p \leq 0.05$.

Poisson Process of λ if: # changes in separate regions are indie; $P(1 \text{ change in int. length } h) \approx \lambda h$; prob of 2 or more changes in a short interval is ≈ 0 . Counting changes in a fixed length interval: X is # changes in interval of length t in a poisson process of rate α . X has a poisson distribution w/ param μ , $\mu = \lambda t$. If a poi proc. exists w/ rate λ , X is the length of interval needed to see a change. The cdf of X : $F(x) = 1 - e^{-\lambda x}$ where $x \geq 0$ (see exponential distribution)

Erlang Distribution: $X \sim \text{Erlang}(r, n)$

x : interval length for r events with rate λ .

cdf $x > 0$: $P(X \leq x) = 1 - P(X > x) = 1 - P(\text{"at most } r-1 \text{ changes in } [0, x]) = 1 - \sum_{k=0}^{r-1} [(e^{-\lambda x} * (\lambda x)^k) / (k!)]$; $\mu = r/\lambda$; $\sigma^2 = r/(\lambda^2)$

| parameter | estimator |
|------------|---|
| μ | $\bar{X} = \sum_{i=1}^n X_i/n$ |
| σ^2 | $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ |
| σ | $S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}$ |
| p | $\hat{p} = \# \text{ of successes} / n$ |

Midterm 2

Uniform Distribution
pdf = $1/(b-a)$ on $[b,a]$; find probs w/ rectangle areas.
Mean: $\mu = E[X] = (a+b)/2$; **Var:** $\sigma^2 = V[X] = [(b-a)^2]/12$
cdf: $(x-a)/(b-a)$ w/ $a <= x <= b$.

Exponential Distribution: $X \sim \exp(\lambda)$
 x : interval length to see change w/ rate λ ;
pdf = $\lambda e^{-\lambda x}$, **cdf** = $1 - e^{-\lambda x}$. $\mu = 1/\lambda$; $\sigma^2 = 1/(\lambda^2)$

Sampling distribution and data descriptions
Characteristics of the pop. are known as variables like X,Y, etc. We take n sample objects. Thus x_1, x_2, \dots, x_n are n indie observations of X. A function of the random sample X_1, X_2, \dots, X_n is a **statistic**. We use the stat $\hat{\theta} = h(X_1, X_2, \dots, X_n)$ to estimate pop param θ . $\hat{\theta}$ is a point estimator of θ , is a rv, and its observed value of the rv $\hat{\theta} = h(x_1, x_2, \dots, x_n)$ is called a point estimate of θ .
Sampling dist of Xbar is the dist of Xbar in all possible samples of a certain size from the population. Sampling dist of sample mean: r.v. Xbar is the mean of n indie observations of a pop with mean μ , sd σ
 Mean of the sample mean's sampling distribution is equal to the population mean ($\mu_{\text{Xbar}} = \mu$). $\sigma_{\text{Xbar}} = \sigma/\text{root}(n)$

Normal Distribution: $X \sim N(\mu, \sigma^2)$
 $P(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-x^2/(2\sigma^2)}$
 Cdf: $\Phi(z) = P(Z <= z) = \int_{-\infty}^z \phi(t) dt = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt$

Standard Normal: $Z \sim N(0,1)$
 $P(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-x^2/(2\sigma^2)}$
 $\phi(z)$ = prob that a standard normal var assumes a value on $[0, z]$.

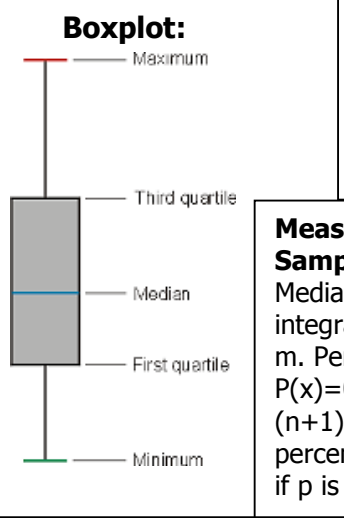
Sampling Distributions: Central Limit Theorem
 Sample mean will be \sim normally distributed for large sample sizes ($n > 30$).
 Xbar: sample mean μ of n indie observations. Sample mean is = population mean. Standev of sampling distribution of Xbar:
 $\sigma_{\text{Xbar}} = \sigma/\text{root}(n)$

Measures of Dispersion/Variability
 The sample Standard Deviation is s and s^2 is the sample variance:

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

 The Sample range, $r = \max(x_i) - \min(x_i) = y_n - y_1$;
 Quartile: $Q_1 = 25^{\text{th}}$ percentile. $\text{IQR} = Q_3 - Q_1$, measures dispersion.

Convert to std. norm: $Z = \frac{X - E[X]}{\sqrt{V(X)}} = \frac{X - \mu}{\sigma}$
 $F_X(x) = P(X \leq x) = P\left(\frac{X - \mu}{\sigma} \leq \frac{x - \mu}{\sigma}\right)$
 $= P\left(Z \leq \frac{x - \mu}{\sigma}\right) = \Phi\left(\frac{x - \mu}{\sigma}\right)$
 $P(a < X \leq b) = F_X(b) - F_X(a)$
 $= \Phi\left(\frac{b - \mu}{\sigma}\right) - \Phi\left(\frac{a - \mu}{\sigma}\right)$



Expectations and Variances of a Linear Function of the Random Sample
 consider rv's X_1, \dots, X_n and constants c_0, \dots, c_n . $Y = c_0 + c_1 X_1 + \dots + c_n X_n$;
 $E[Y] = c_0 + c_1 E[X_1] + \dots + c_n E[X_n]$. Assume all X indie, $\text{Var } \sigma_Y = c_1^2 \sigma_{X_1}^2 + \dots + c_n^2 \sigma_{X_n}^2$. Let X be rv, a & b constant; $E[aX+b] = aE[X]+b$; $V[aX+b] = a^2 \sigma_X^2$.
 Sample mean has $\mu_{\text{Xbar}} = \mu$ mean and $\sigma_{\text{Xbar}}^2 = \sigma^2/n$

Using the table:
 What you get from finding a value: The probability to the left of it.
 To get $P(a < Z < b)$ you do $z(b) - z(a)$; To get $P(Z > a)$, do $P(Z < (1-a))$.
 $P(a < Z < b) = 0.95$, $1 - 0.95 = 0.05$, $0.05/2 = 0.025$, find in table!
 "nth percentile" \rightarrow find 0.n in body of table
Ex: $P(X < 180.0) = P((x - \mu)/\text{sig} < (180 - \mu)/\text{sig}) = P(Z < \text{dec.ans.})$, find $z(\text{dec.ans.}) = \text{answer}$.

Measures of center
Sample mean: avg of the n values, denoted by **xbar**.
 Median: 50th percentile, position is $(n+1)*(1/2)$; to get it integrate pdf (f)x from bottom to m, must=1/2, solve for m. Percentile: find %th percentile of f(x) by finding x at $P(x)=0.\%$. Position of %th percentile is $(n+1)*(k/100) = m(\text{whole}) + p(\text{fractional})$. The kth percentile is y_m if $p=0$, and is equal to $y_m + p(y_{m+1} - y_m)$ if p is fractional.

Independent samples
 2 indie random samples of populations $X_1, \dots, X_{n_1}, Y_1, \dots, Y_{n_2}$ each with mean μ and variance σ^2 . If they're both normal, $N(\mu_1, \sigma_1^2)$ and $N(\mu_2, \sigma_2^2)$, then also
 $\bar{X} - \bar{Y} \sim N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)$
 If we assume n_1 and n_2 are both $>= 30$, we can show that Z has approx. a standard normal dist.:

$$Z = \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}}$$

Estimators and point estimates:
 \hat{p} = sample proportion = estimate of pop. proportion p; sometimes sampling distros are \sim normal; thus their conf intervals are in form Point estimate \pm margin of err. We use estimator, point estimate, sample mean(Xbar), sample variance(S^2), sample standard deviation (S), and sample proportion $\hat{P} = X/n$ (X is # items satisfying attribute among n items).
 Let $\hat{\theta}$ be an estimator for θ . $\hat{\theta}$ is an unbiased estimator if $E[\hat{\theta}] = \theta$. A binomial experiment with $\hat{P} = X/n$ be the sample proportion from a pop with $P(\text{success}) = p$ (population proportion), X is successes among n trials. $E[\hat{P}] = p$; $\sigma^2_{\hat{P}} = [p(1-p)]/n$. Standev of a point estimator is the standard error. serr of the mean and sample proportion are: $\sigma_{\text{Xbar}} = \sigma/\text{root}(n) = s/\text{root}(n)$, $\sigma_{\hat{P}} = \text{root}([p(1-p) / n] = \text{root}[(\hat{p}(1-\hat{p})) / n]$

Confidence Interval for the mean
 Conf. int.: within \pm errmargin, represented by a % between (bla,bla). n indie obs, normal dist, mean $\mu = \sigma = \sigma_{\text{Xbar}} = \sigma/\text{root}(n)$.
 Std.Normal Z: $P(-Z_{\alpha/2} < Z < Z_{\alpha/2}) = 1 - \alpha$. Conf int for μ : $\text{Xbar} \pm Z_{\alpha/2} * (\sigma/\text{root}(n))$, where alpha is 1 - (%confidence), ie 95% con \rightarrow alpha is 0.05; the \pm is the margin of error.

| Parameter | Estimator | Estimate |
|------------|-----------|-----------|
| μ | \bar{X} | \bar{x} |
| σ^2 | S^2 | s^2 |

Confidence Interval for the proportion

$X_{bar} \pm [t_{\alpha/2} * (s/\sqrt{n})]$; $n > 40$ is good, 15 to 40 is less so, < 15 is a no go.

If we have a sample of size n from the population, and we denote by Y the total number of individuals (in the sample) who possess the desired characteristic, then a *point estimator* for p is: $\hat{P} = Y/n$.

If n is large, \hat{p} is distributed normally with mean: p and var: $p*(1-p)/n$.

$$Z = \frac{\hat{P} - p}{\sqrt{p(1-p)/n}}$$
 is approximately a standard normal random variable

For large n , $100(1-\alpha)\%$ conf int for p : $\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

Thus, if \hat{p} is an estimate for p , we are at least $100(1-\alpha)\%$ confident that the error $|\hat{p} - p|$ will not be greater than E with a sample size

$$n \geq \frac{1}{4} \left(\frac{z_{\alpha/2}}{E} \right)^2$$

Hypothesis Testing:

Left-sided ($H_1: \mu < \mu_0$), right- ($H_1: \mu > \mu_0$),

two- ($H_1: \mu \neq \mu_0$);

Type 1 Error: Rejecting H_0 when it is actually true;

Type 2 Error: Not rejecting H_0 when it is false.

Type 1: $P(\text{Type I error} | H_0 \text{ is true}) = \alpha = \text{significance level}$; Type 2: β . **Power** = $1 - P(\text{Type 2 error}) = 1 - \beta$

Hypothesis Testing, σ known:

Z = (orange); If H_0 (null hyp.) is true, you should get a random sample from the standard normal.

Rejection regions:

$\alpha = \text{sig. level}$, probability of rejecting null hyp. if H_0 is true); Find rejection region, reject H_0 if the test falls in area.

For $H_1: \mu \neq \mu_0$, $\alpha = 0.05$:

e.g. $H_0: \mu = \mu_0$, $H_1: \mu \neq \mu_0$; reject outside center region (95% conf. interval area, i.e. areas outside of $z(\alpha/2)$, the "critical value"). So if Z value is ≥ 1.96 or ≤ -1.96 , reject H_0 .

For $H_1: \mu < \mu_0$, $\alpha = 0.05$:

Reject null hypothesis if $Z \leq -1.645$

For $H_1: \mu > \mu_0$, $\alpha = 0.05$:

Reject null hypothesis if $Z \geq 1.645$

P-value method:

P-value: probability of getting the observed value or a value w/better evidence against H_0 , if H_0 is true;

i.e. it is a measure of how strong the evidence against the null hypothesis is.

Smaller p-value \rightarrow greater evidence against H_0 .

If $p\text{-val} < \alpha$, the evidence against H_0 is significant at the α level of significance, so **we reject H_0 .**

ex: $P\text{-value} = P(H_1) = (X_{bar} - \mu_0) / (\sigma / \sqrt{n})$, $P[Z \leq z_x] = z_{table}(z_x)$

i.e. **if $H_1: \mu < \mu_0$** , $p\text{val} = P[Z \leq z_x]$ (finding area to left of crit. value)

and **if $H_1: \mu > \mu_0$** , $p\text{val} = 1 - P[Z \leq z_x]$ (finding area to right of crit. value)

Confidence Interval for the mean μ ; σ known:

$Z_{\alpha/2}$: Find α , divide by 2, use logic, find the z of that% in table. (body of table is the %) "A sample of n people gave an average value of gh " means that $gh = X_{bar}$. For a "% conf int, % of the "% conf ints calculated would capture the true value.

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \quad 1 - \alpha = P(-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma / \sqrt{n}} \leq z_{\alpha/2})$$

$$= P(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}})$$

a $100(1 - \alpha)\%$ confidence interval for μ is

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = [\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}]$$

If \bar{X} is the size- n -sample mean from a normal population with known variance σ^2 , the exact $(1 - \alpha)$ C.I. for μ :

$$\bar{X} \pm z_{\alpha/2} \sigma / \sqrt{n};$$

Precision:

$$\bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} - (\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}) = 2 \frac{z_{\alpha/2} \sigma}{\sqrt{n}}$$

Sample Size: If \bar{x} is used as a point estimate for μ , then we are $100(1 - \alpha)\%$ confident that the error $|\bar{x} - \mu|$ will not be greater than E , if the sample size satisfies

(where E is Error estimation) $n \geq \left(\frac{z_{\alpha/2} \sigma}{E} \right)^2$

Confidence Interval for the mean μ ; σ unknown:

Conf interval: $X_{bar} \pm [t_{\alpha/2} * (s/\sqrt{n})]$; where $s = \text{StandErr}(X_{bar}) = s/\sqrt{n}$; t on table with $n-1$ degrees of freedom ν .

if σ is not known, then

where $t_{\alpha/2, n-1}$ is a Student's t -distribution with $n-1$ degrees of freedom. Consequently,

$$T := \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1} \quad P(-t_{\alpha/2, n-1} < \frac{\bar{X} - \mu}{S/\sqrt{n}} < t_{\alpha/2, n-1}) = 1 - \alpha,$$

where $t_{\alpha/2, n-1}$ is the $100(1 - \alpha/2)$ th percentile of the Student distribution with $n-1$ degrees of freedom. Student's t -distribution It is perhaps not

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$$

T has Student's t -dist with $n-1$ degrees of freedom df.

a $100(1 - \alpha)\%$ confidence interval for μ is

$$\bar{x} \pm t_{\alpha/2, n-1} \left(\frac{s}{\sqrt{n}} \right) \quad (1)$$

where t is the value that we read in Table at level $\nu = n - 1$, such that $P(-t \leq T \leq t) = 1 - \alpha$, i.e $P(T \leq t) = 1 - \alpha/2$.

If \bar{X} is the size- n -sample mean from a normal population with unknown variance, the exact $(1 - \alpha)$ C.I. for μ :

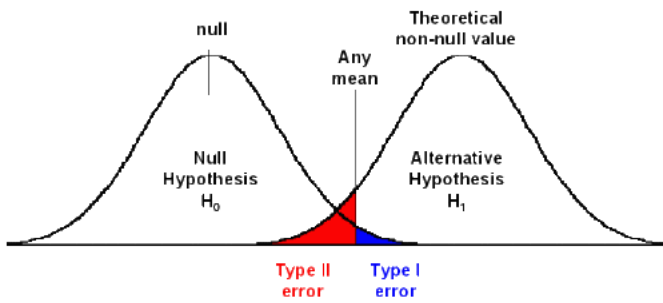
$$\bar{X} \pm t_{\alpha/2, n-1} S / \sqrt{n};$$

If \bar{X} is the size- n -sample mean from a normal or non-normal population with unknown variance and n is 'big', the approximate $(1 - \alpha)$ C.I. is:

$$\bar{X} \pm z_{\alpha/2} S / \sqrt{n};$$

\leftarrow The confidence interval is simply the values of μ_0 for which we do not reject the null hypothesis.

1.3 Subjective: $0 \leq P(E) \leq 1$; Equally likely: $P(E) = N(E)/N(S) = \# \text{good} / \# \text{total}$; Axioms: Positivity: $P(E) \geq 0$; Certainty: $P(S) = 1$; Additivity: add up prob. of mut.excl. events; $P(E') = 1 - P(E)$; $P(A) = P(A \cap B) + P(A \cap B')$
 $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
 $P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C)$



| Parameter name | Population parameter symbol | Sample statistic |
|------------------------|-----------------------------|-----------------------------------|
| Number of cases | N | n |
| Mean | μ (mu) | \bar{x} (Sample mean) |
| Proportion | π (Pi) | P (Sample proportion) |
| Variance | σ^2 (Sigma-square) | s^2 (Sample variance) |
| Standard deviation | σ (Sigma) | s (sample standard deviation) |
| Correlation | ρ (rho) | r (Sample correlation) |
| Regression Coefficient | β (beta) | b (sample regression coefficient) |

Regression:

Predictor: X (independent); Response: Y (dependent). Line of best fit = regression line.

Notation of fitted/estimated line: $E(X)=\hat{y} = \hat{\alpha} + (\hat{\beta}) * x = \text{yintercept} + \text{ystretch} * x$; Fitted value of the i^{th} observation: $\hat{y}_i = \hat{\alpha} + (\hat{\beta}) * x_i$; Difference between fitted value/observed response y_i is called the i^{th} residual observed error: $e_i = y_i - \hat{y}_i = y_i - (\hat{\alpha} + (\hat{\beta}) * x_i) = \text{observed} - \text{predicted}$.

Residual sum of squares:
$$SSE = \sum_{i=1}^n [y_i - (\hat{\alpha} + \hat{\beta}x_i)]^2$$
; Total sum of squares:
$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$
;
$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$$

Regression sum of squares:
$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$
; Error sum of squares:
$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$
;
$$\hat{\beta} = \frac{s_{xy}}{s_x^2}$$

Point $(x, \hat{\alpha} + \hat{\beta}x)$ is always on the regression line. total deviation = explained deviation + unexplained deviation

r: correlation between y and x, r^2 : coeff. of determination, r between +- 1 with exactly positive / and negative \ slope respectively being $\sim = +-1$. (r_{xy} close to zero means no correlation)

$$R^2 = r^2 = \frac{SSR}{SST} = \frac{s_{xy}^2}{(s_x^2 s_y^2)}; s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2; s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2; s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}); r = \frac{s_{xy}}{s_x s_y}$$

standerr($\hat{\beta}$) = $\sqrt{\hat{\sigma}^2 / s_x^2}$; $\hat{\sigma}^2 = SSE / (n-2)$

Other distributions:

Hypergeometric: without replacement -> not independent.

If you are sampling <5% of pop, you can approximate this using binomial.

Example calculation:

12D, 24R, 8I. Pick 6 out of the total (42). What is P(3D 2R 1I)?

$$[(12C3) * (24C2) * (8C1)] / (44C6) = 0.0688$$

Multinomial: x_i = number of occurrences of outcome i out of n indie trials, k m.e. outcomes

For each trial, these k outcomes occur with probabilities p_1, \dots, p_k , $\text{SUM}(p_1, \dots, p_k) = 1$

Pmf: $P(X_1=x_1, \dots, X_k=x_k) = [n! / (x_1! * \dots * x_k!)] (p_1^{x_1} * \dots * p_k^{x_k})$ for $x_i=0, 1, \dots, n$, ($x_1 \neq 1$)

$E(X_i) = n * p_i$; $\text{Var}(X_i) = n * p_i * (1-p_i)$; $\text{SUM}(x_i) = n$