# Laboratory #3.  Probability Models

Statistical analysis rests on frequency distributions—either empirical (from data) or theoretical (from a probability model). Empirical distributions display all of the information in a sample or in set of observations.  We'll be using empirical distributions to calculate probabilities in Lab 4. In today's lab we'll be using probability models, such as the normal and binomial distributions, to calculate probabilities and likelihoods.

We'll begin by calculating the probability of the possible outcomes of a learning experiment, using the binomial probability density function (pdf). We'll then calculate cumulative probabilities with the cumulative version of the binomial probability model--the cumulative distribution function, cdf. After that we'll use the binomial probability model to obtain likelihoods and likelihood ratios, which we calculate once we have collected data.
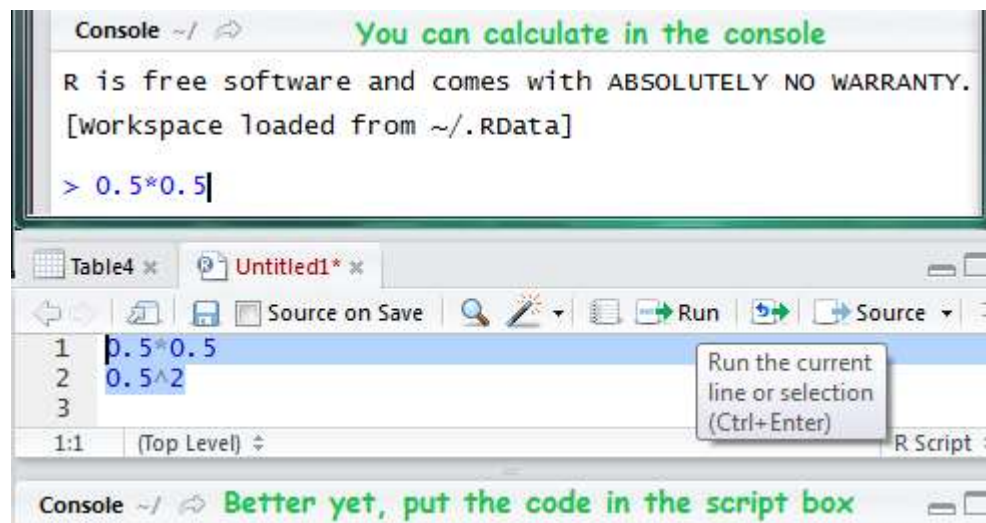
Next, we'll learn how to calculate probabilities for the most commonly used probability models for hypothesis testing: $t$, $F$, and $\chi^2$ (pronounced khai-squared).  We'll also run the calculations "backwards" with inverse cdf functions, which are used for confidence intervals.

In this lab we'll focus on the mechanics of calculation. The lab can be completed in a spreadsheet, in any statistical package with a pull-down menu, or a package based on coding (R). A spreadsheet is the easiest way to complete this lab, especially if you already know how to use a function in a spreadsheet. These functions are a convenient way to calculate probabilities from $t$, $F$, and $\chi^2$ distributions. Spreadsheet functions are not the easiest way to do most of the calculations in the remaining labs in this course, so in future labs you will be using either a menu-based or a code-based package.  Spreadsheet calculations and R-code will be shown in today's lab.  If want to learn how to make probability calculations in a menu-based package such as SPSS, the course instructor can help you with that.

To begin, we'll use your spreadsheet or statistical package as a simple calculator.
For a coin toss, the probability of Heads on the first toss is            $p =$ _____
Now calculate the probability of Heads on 2 successive tosses.     $p^2 =$ _____



Calculation on the R console.

Calculation on the R console gets messy.  It is better to use the R script box, where you can modify the code.

1. Now calculate the probability of 4 Heads on 4 tosses.  _____

   The probability you calculated was for independent events.

   Next, we'll calculate a probability for dependent events: the probability of drawing two successive diamonds from a fair deck of 52 cards  with four suits,
   13 clubs, 13 diamonds, 13 hearts, and 13 spades
   For the first card, what is the probability of a diamond?         _____
   How many diamonds remain, after the draw?                        _____
   What is the probability of another diamond?                      _____
          This is called a conditional probability.
   Calculate the probability of 2 successive draws of a diamond.
          (product of the two successive events).                   _____

2. Now that we have our probability calculator working we'll use the binomial probability model to plan an experiment.
      In 1959 James V. McConnell created  The *Worm Runner's Digest* (W.R.D.) to report his experiments with memory transfer in planarian worms (see  Hartry AL *et al* 1964 Planaria: Memory transfer through cannibalism re-examined *Science* 146: 274-275).

   What is the probability that an untrained planarian worm will correctly guess whether to turn left or right to obtain food in a T-maze on 6 successive trials?  *I.e.*, what is the probability of 6 correct guesses by chance alone?
                   $n = 6$ trials
                   $X = 6$ successes in 6 trials
                   $p = 50\%$   That is, we assume an equal probability of turning left or right.

      The probability of 6 successes in 6 trials is         $Pr\{X=6\} = $ _____

   We can use the binomial a probability model to make the same calculation.

   We are interested in $Pr\{X= x \mid p\}$ the probability (relative frequency) of $X$ successes in $n$ trials given a fixed probability of success of  $p = 0.5$ (1:1 odds) on each trial.
   The symbol $Pr\{X= x \mid p\}\}$ is read: "the probability that successes $X$ will take on a certain value $x$, given a known parameter $p$." We can display $Pr\{X= x \mid p\}$ as a graph of probabilities plotted against $x$, or as a table of numbers that show probabilities $Pr\{X= x\}$ for each value $x$.
   The algebraic expression for a run of $x$ successes in $n$ trials  is:     $f(x) = p^n$
   In the example above          $f(6) = (0.5)^6 = 0.0156$
   Now calculate $f(6)$ using the binomial probability model.         _____

   ---

   Excel Spreadsheet.  Select a cell
        Insert Function  (*fx* icon below the menu)
           Statistical functions
              Binom.Dist
                 number_s = 6
                 trials = 6
                 probability_s  = 0.5
                 cumulative = FALSE (to get the probility, pdf)

   *Calculate binomial probability in excel.*

   ```
   dbinom(6, 6, 0.5)
   ```

   *Calculate binomial probability in R*

Now calculate the probability of a worm making 5 correct choices in 6 trials, assuming a success rate of $p = 50\%$ on each trial.         _____

7 correct choices in 7 trials, assuming a success rate of $p = 80\%$ trial.      _____

3. Having calculated a single probability, we move to calculating and displaying the entire binomial probability distribution. If we carry out the planaria experiment repeatedly, and no learning occurs ($p$ remains fixed at 50%), what is the probability distribution of outcomes $Pr(X = 0$ successes, 1 success, 2 successes, 3 successes, *etc.*) in $n = 6$ trials?

To compute this, we use the functional expression for the binomial distribution, which is often written as:

$$f(x) = \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x}$$

In this expression the number of successes $x$ varies from 0 in $n$ trials to $n$ in $n$ trials
$n!$ means the factorial of the number $n$.   $6! = 6 \cdot 5 \cdot 4 \cdot 3 \cdot 2 \cdot 1$
$0! = 1$       because (any number)$^0 = 1$

To remind ourselves that $p$ is a parameter, we will use a greek letter.

$$f(x) = \frac{n!}{x!\,(n-x)!} \pi^x (1-\pi)^{n-x}$$

From this, we calculate the probability of 6 successes in 6 trials.

$$f(6) = \frac{6!}{6!(6-6)!}\, 0.5^6 (1-0.5)^{6-6} = \Big[ \underline{\hspace{5cm}} \Big]$$

To shorten the calculation you can apply algebra to simplify the expression in the space above
Now calculate the probability.

$$f(6) = \underline{\hspace{5cm}} = [\,\underline{\hspace{2cm}}\,]$$

Write out the expression for $f(0)$ then do the calculation ($x = 0$ instead of $x = 6$)

$$f(0) = \underline{\hspace{7cm}}$$         write expression

$$f(0) = \underline{\hspace{2cm}}$$                         show result

We could continue in this fashion, for $f(1), f(2)$, *etc.* This is laborious so we will calculate a column of probabilities $f(x)$ from a column of possible outcomes $x$. The procedure is nearly the same in spreadsheet and in statistical packages that use spreadsheet input. Because the procedure is so similar among packages, it is first shown as pseudocode–a list of procedures to be carried out in any package. The pseudocode is then translated into the specific procedures of a package or spreadsheet that we are using.

Pseudocode to calculate a probability density function $f(x)$

> Select a column and name it $x$.
> Place the values $x = 0,1,...6$ into this column.
> Select an adjacent column and name it $f(x)$.
> Select the first cell in column $f(x)$.
> Apply the binomial function to calculate $f(0)$ from the cell $x = 0$.
> Apply the function to the rest of the column $f(x)$.

*Calculate a probability density function f(x)*

Once you have looked at the pseudocode, carry out the procedure in any package you like.

> Spreadsheet.  Select top cell in column labelled $f(x)$
>     Insert Function
>         Statistical functions
>             Binom.Dist
>                 number_s = adjacent cell with value of x = 0
>                 trials = 6
>                 probability_s = 0.5
>                 cumulative = false (we want the probability, pdf)
> Select cell with $f(0)$, copy, then paste in rest of column $f(x)$.

*Calculate a probability density function f(x) in excel*

```
#DataDef - Create a column called x, with values 0 to 6
 x <- c(0:6)
#Execution - Calculate Pr(X=x) for each value of X
 dbinom(X, 6, 0.5)
 f.x <- dbinom(x,6,0.5)
 cbind(x,f.x)
```

*Calculate a probability density function f(x) in R*

Now display the result as a graph.

> Select both columns   $x$ and $f(x)$
>     Insert
>         Scatterplot
> Print the plot on paper and then draw in the bars by hand.
> Or find a barchart routine on the web.

*Produce a scatterplot of the function f(x) in excel*
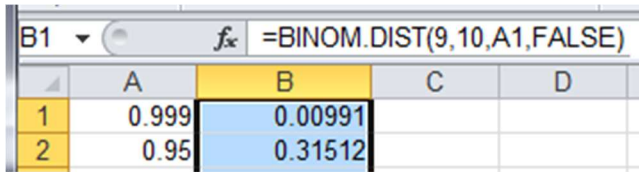
```
barplot(f.x, names.arg=x)
```

*Produce a barchart in R*

Tape or paste your plot here:

4. In games of chance such as roulette, dice, and cards we know the probability parameter $\pi$. In science we often do not know the parameter. So we estimate of the parameter $\hat{\pi}$ from the data. To do this we use our probability model to calculate the most likely value of the parameter $\pi$, given the data. Instead of holding $\pi$ fixed, we let it vary over its possible range, which is 0 to 1. We calculate $L(\hat{\pi}|data)$, the likelihood of the estimate $\hat{\pi}$ given the data, over this range. Note the contrast between this expression compared to our previous expression $Pr(outcome|\pi)$. As you can see we have flipped things around—instead of calculating the probability given a fixed parameter we calculate the likelihood of a parameter, given the data.

To gain a better sense of how this works let's do calculations. For the simple T-maze, set up a column of proportions (Success/Trial) ranging from 0.5 (no learning) up to 0.999 (highly trained) as shown below. Use the binomial distribution to calculate likelihoods $L(\hat{\pi}|data)$ for the results of 2 experiments:  data = 9 successes in 10 trials and data = 7 successes in 10 trials.

| | Experiment 1 | Experiment 2 | Data |
|---|---|---|---|
| | 9 | 7 | = Successes |
| | 10 | 10 | = Trials |
| Success/Trial | $L(\hat{\pi}|data)$ | $L(\hat{\pi}|data)$ | |
| 0.999 (highly trained) | 0.00991 | | |
| 0.95 | 0.315 | | |
| 0.90 | | | |
| 0.85 | | | |
| 0.80 | | | |
| 0.75 | | | |
| 0.70 | | | |
| 0.65 | | | |
| 0.60 | | | |
| 0.55 | | | |
| 0.50 (no learning) | | | |

| B1 | | $f_x$ | =BINOM.DIST(9,10,A1,FALSE) |
|---|---|---|---|
| | A | B | C | D |
| 1 | 0.999 | 0.00991 | | |
| 2 | 0.95 | 0.31512 | | |

Note the similarity between the excel command and the dbinom command.

```
#DataDef    Create a column called pr, with values 0.99 to 0.5
pr <- c(0.999, 0.95,0.9, 0.85, 0.8,0.75,0.7,0.65,0.6,0.55,0.5
)
#Execution Calculate L(pr|data) for each value of pr
dbinom(9, 10, pr)
L.pr <- dbinom(9,10,pr)
```

*Calculate the likelihood of each value in R*

Print your table and tape it in the blank area above and to the right (or fill in columns by hand).

Have a look at the barplot of likelihood across the parameter values from 0.5 to 0.999.

For the first experiment, what is the maximum likelihood?  Max $L(\hat{\pi}|data)$ = _____

What is the degree of learning (Success/Trial) for which likelihood is maximum?  $\hat{\pi}$ = _____

Calculate the likelihood ratio $LR = Max\ L(\hat{\pi}|data) / L(\pi=0.5|\ data)$ = _____

For the second experimental result (7 successes in 10 trials),
          what is the maximum likelihood?        $Max\ \mathcal{L}(\hat{\pi}|data)$ = _____

What is the degree of learning (Success/Trial) for which likelihood is maximum?  $\hat{\pi}$ = _____

Calculate the likelihood ratio $LR = Max\ \mathcal{L}(\hat{\pi}|data)\ /\ \mathcal{L}(\pi{=}0.5|\ data)$ = _____
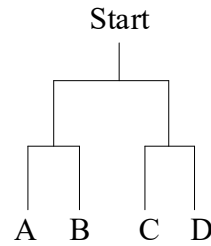
For Experiment2 make a statement about the evidence relative to the following criteria
          $LR < 20$ - little or no evidence
          $LR \geq 20$ – some evidence

5. Next, we move to cumulative frequency distributions, which are used in hypothesis testing. We'll begin with a new maze and a new probability parameter, which we hold constant.

Start

Can a planarian worm learn to avoid endpoint D?  In a maze with 2 tiers (see right) a randomly moving worm will arrive at each of the four endpoints (A, B, C, or D) with equal probability.  What proportion of the time will a randomly moving worm arrive at endpoint D?

A   B    C   D

          $\pi$ = _____

What is the probability that an untrained worm will arrives at D only once in 6 trials.
(*i.e.,* what is the probability of doing this by chance alone ?)
Place your guess (this will be marked as present, not right or wrong)        $Pr\{x \leq 1\}$ = _____

To calculate the answer, we need to accumulate probabilities.
We need to add the probability of no arrivals and the probability of arriving once at D.

What is the chance of never arriving at D in 6 trials?                    $f(0)$ = _____
What is the chance of arriving at D exactly once in 6 trials ?            $f(1)$ = _____

Now add the probabilities to find out the chance
of arriving one or fewer times at endpoint D.                $Pr\{x \leq 1\}\ =\ F(1)$ = _____

Adding up the probabilities in the tail of the distribution gets tedious.
Our spreadsheet or statistical package will do this for us by using the cumulative distribution function $F(x)$.

The functional expression is   $F(x) = \sum f(x)$
In the example above          $F(1) = f(0) + f(1)$

Here is pseudocode for calculating the cumulative distribution function $F(x)$

Pseudocode (applicable to almost any package).
   Select a column and name it $x$.
   Place the values x = 0,1,...6 into this column (vector).
   Select an adjacent column and name it $F(x)$.
   Select the first cell in column $F(x)$.
   Apply the binomial function to calculate $F(0)$ from the cell x = 0.
   Apply the function to the rest of the column $F(x)$.

*Calculate a cumulative distribution function F(x)*

5. (continued).  Now use the pseudocode to make the calculation in your package.

> Excel Spreadsheet.  Select a top cell in column labelled $F(x)$.
>     Insert Function (*fx* icon)
>         Statistical functions
>             Binom.Dist
>                 number_s = adjacent cell with value of $x = 0$
>                 trials = 6
>                 probability = 0.25
>                 cumulative = true (we want the cdf)
> Select cell with $F(0)$, copy, then paste in rest of column $F(x)$.

*Calculate a cumulative distribution function F(x) in excel*

```
F.x <- pbinom(x,6,0.25)
cbind(F.x, x)
```

*Calculate F(x) in R*

Make sure your calculations using the cumulative distribution function match the calculations by summing the probability density function.

Have a look at the barchart of the binomial cdf.
Tape, paste, or fill in by hand this table - - >
showing the pdf and cdf.

| $x$ | pdf $f(x)$ | cdf $F(x)$ |
|---|---|---|
| 0 | | |
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |

6.  In the previous example we looked at $Pr\{X \leq x\}$ the probability of $X$ or less successes in $n$ trials. In statistical hypothesis testing we are interested in $Pr\{X > x\}$ the probability of exceeding $x$.  To obtain this probability we calculate the cumulative probability $Pr\{X \leq x\}$, then subtract it from 1, to obtain $Pr\{X > x\}$.

Can a planarian worm learn (when rewarded with food)  to arrive at endpoint A in our 2-tier T-maze?  To evaluate this we calculate improbable events--the probability of an untrained worm arriving at endpoint A by chance 5 or more times in 6 trials.

We can compute this by accumulating probabilities in the upper tail, using the pdf.

The functional expression is    $Pr\{X > x\} = 1 - F(x) = 1 - \sum f(x)$
In the example above          $Pr\{X > 4\} = 1 - F(4) = 1 - \sum f(x) = f(5) + f(6)$

$$f(5) = \text{_____}$$
$$f(6) = \text{_____}$$
$$Pr\{X > 4\} = \text{_____}$$

We can compute this more efficiently by using the cumulative distribution function.

The functional expression is    $Pr\{X > x\} = 1 - F(x) = 1 - \sum f(x)$
In the example above          $Pr\{X > 4\} = 1 - F(4) = 1 - \text{_____} = \text{_____}$

7. Based on these calculations, we decide that we need a more challenging maze to identify learning with fewer trials. In the space to the left and below, draw a 3 tier maze with 8 possible endpoints. Label them A through H. Success is defined as arrival at point H.

   What is the expected arrive rate at H by chance on each run through the maze?  _____

   What is the chance of 4 correct arrivals at H in 4 trials?   $Pr\{X = 4\} =$ _____

   3 or more correct endpoints in 4 trials?   $Pr\{X > 2\} =$ _____

   Compare the two designs: 2-tier maze and 6 trials *versus* 3 tier maze and 4 trials.
   At what value of $x$ does $Pr\{X > x\}$ fall below 5% in each maze?    2- tier maze.  $x =$ _____
                                                                           3- tier maze.  $x =$ _____
   Which design is more efficient (fewer trials to reach 5% criterion?    _____

8. Next we calculate probabilities for the $\chi^2$ distribution, which is used for hypothesis testing. As an example, we'll use the coin-tossing experiment that resulted in 7 successes in 10 trials. This result ( $LR = 0.2668 / 0.1172$) was 2.3 times more likely than 5 success in 10 trials. Is this too large to be due to chance? To evaluate this we use the $G$ statistic = 2 $ln(LR)$.
   Calculate $G_{Experiment2} =$ _____

   To evaluate this $G$ statistic we use the $\chi^2$ distribution with $df = 1$ degree of freedom.
   We want the probability in the upper tail $Pr\{ \chi^2 (1) > G\}$.
   This is the probability of obtaining a value of $\chi^2 (1) = G$ or more extreme.
   Some packages report the upper tail $Pr\{ \chi^2 (df) > G\}$, which is what we want.
   Most packages report the cumulative distribution from left to right $F(x) = Pr\{ \chi^2 (df) \leq G\}$.
   For these, we plug in the value of $G$ (as above) with $df = 1$ to obtain $Pr\{ \chi^2 (1) \leq G\}$.
   We then subtract to obtain the upper tail: $Pr\{ \chi^2 (df) > G\}. = 1 - Pr\{ \chi^2 (df) \leq G\}$.

   Here is the pseudocode for calculating a p-value for a Chisquare statistic.

| |
|---|
| Select a column and name it *Gstat*.<br>Place the value $G$ in the first cell.<br>Select an adjacent column and name it $Pr(G)$.<br>Select the first cell in column $Pr(G)$.<br>Apply the Chisquare function to calculate $Pr(G)$ from $G$<br>Compute the upper tail from the cdf when necessary. |

*Calculate a p-value from a Chisquare distribution*

   Now use the pseudocode to carry out the calculations in the package you are using.
   As a trial run, we'll use a value of $G = 3.84$.

| |
|---|
| Excel menu<br>   Function<br>      Statistical functions<br>         Chidist<br>            X = cell 1 in adjacent column (3.84)<br>            Deg_freedom = 1 |

*Calculate a p-value from a Chisquare distribution in excel.*

```
pchisq(3.84,df=1)
1-pchisq(3.84,1)
```

*Calculate a p-value from a Chisquare distribution in R*

8. (continued).    The probability of obtaining a $G$-statistic of 3.84 or more, by chance is _____

(not 0.95)

Now calculate the probability of obtaining the following Gstatistic.

$$G = 5.67 \quad df = 3 \quad p = \text{_____}$$

(not 0.8712)

$$G = 10.23 \quad df = 5 \quad p = \text{_____}$$

For the second coin tossing experiment $G_{Experiment2} = 2\ ln(\text{LR}) = 2ln(2.3) \quad df = 1 \quad p = \text{_____}$

9. Graphs are useful in visualizing concepts and the flow of computations.
   Let's compare the probability density function pdf and the cumulative distribution function cdf
   of the chisquare distribution having 1 degree of freedom.

| chisq=x | pdf=f(x) | cdf = F(x) | 1-cdf |
|---------|----------|------------|-------|
| 0.5 | 0.4394 | | 0.4795 |
| 1 | 0.2420 | | 0.3173 |
| 2 | 0.1038 | | 0.1573 |
| 4 | 0.0270 | | 0.0455 |
| 8 | 0.0026 | | 0.0047 |

| Sketch pdf $= f(x)$ | In the box to the left, make a sketch of the pdf versus $x$, from the data above. |
|---|---|
| | Then in the box below the pdf sketch the cdf from the pdf. This will be marked as present or absent, not on whether it is correct. |
| $F(x) = $ cdf, from box above | $F(x) = $ cdf, from computer output |

Now calculate the cdf (fill in the empty column in the table above). Graph the cdf versus $x$, using the same $x$ axis for the pdf and the cdf. Draw the cdf from calculattions in the box to the right at the bottom of this page.

9. Comment on your initial sketch of the cdf, compared to the graph from calculations.

.

|  |
|--|
|  |

10. The *t*-distribution is used to calculate probabilities when comparing two means, two slopes, *etc*. Use your package to compute *p*-values for a *t*-distribution. Use "search" in your package to find the *t*-distribution cdf and how to calculate it.

Two tailed probability     $t_{obs}$ = -2.179  df = 12   $p$ = _____

One-tailed probability   $t_{obs}$ = 1.96  df = 300   $p$ = _____

```
lefttail <- pt(-2.179, df=12)
righttail <- 1-pt(2.179, df=12)
Twotail<-lefttail + righttail
Twotail
```

| C2 | | $f_x$ | =1-T.DIST(A2,B2,TRUE) | | |
|---|---|---|---|---|---|
|  | A | B | C | D | E | F |
| 1 | -2.179 | 12 | 0.0250 | | T.DIST(A1,B1,TRUE) |
| 2 | 2.179 | 12 | 0.0250 | | |
| 3 | | | | | |
| 4 | 2.179 | | 0.0500 | | T.DIST.2T(A5,12) |
| 5 | | | | | |

11. The *F*-distribution is used to calculate probabilities when comparing several means, several slopes, *etc*. Find the *F*-distribution cdf in your package to compute *p*-values for

$F_{obs}$ = 4.56  numerator $df$ = 8, denominator $df$ = 23    $p$ = _____    (not 0.998)

$F_{obs}$ = 2.28   $df$ numerator = 8, $df$ denominator = 23    $p$ = _____    (not 0.9416)

$F_{obs}$ = 1.23   $df$ numerator = 8, $df$ denominator = 23    $p$ = _____    (not 0.6738)

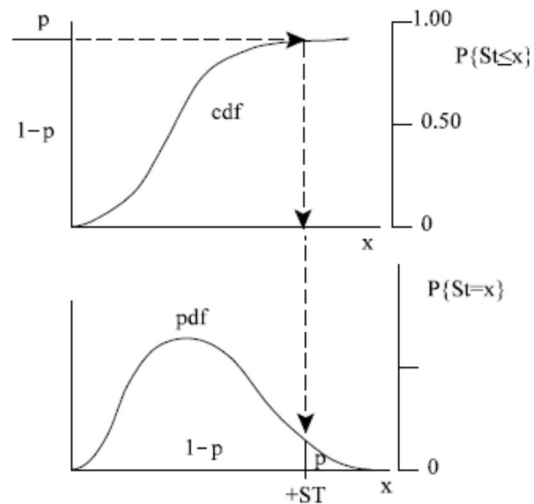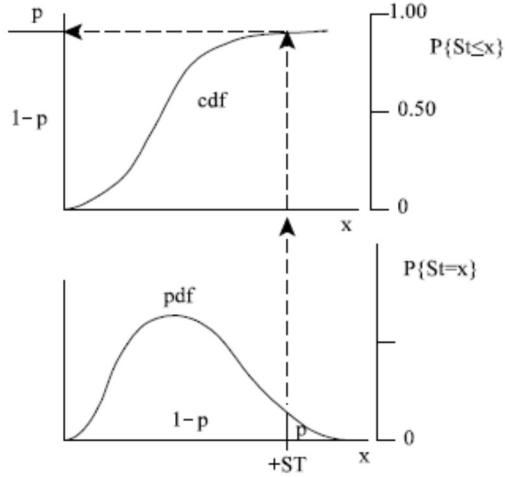| |
|--|
| **HINT for R Users:** Try pf instead of pt |

12. The Z statistic is defined as $(x - \bar{x})/stdev(x)$ where $\bar{x}$ is the mean value of x, and $stdev(x)$ is the standard deviation.   Find the normal distribution cdf in your package and calculate the two-tailed probability for :        $Z$ = 1.96  mu = 0  sigma = 1   $p$ = _____

| |
|--|
| **HINT for R Users:** Try pnorm instead of pf |

13. Critical values.  With modern statistical software we can compute the exact probability of any statistic we like.  Before the days of computers, this was a time consuming chore. To relieve the chore, tables of critical values were constructed and published.  You may have seen these. These tables work backward from the critical *p*-value (1%, 5%, *etc*.) to the critical value of the statistic (F, t, $\chi^2$, *etc*.).  The result is an antique style of statistical practice that compares a statistic to a critical value, rather than reporting the exact p-value.

    To calculate from a critical p-value (*e.g.* 5%) to a statistic, you need to know how your package behaves.  Does it return the probability in the right tail $Pr\{X > St\}$ for a statistic *St*? Excel does this.  If it does, use the critical *p*-value to obtain the critical value of the statistic. Does your package return the cumulative probability $Pr\{X < St\}$ for the statistic St? Most packages do this.  If so, then you will need to use $1–p$ rather than $p$ to obtain the critical value of the statistic.

To the left is a diagram showing the flow of computations for computing a *p*-value from a cumulative distribution function, cdf.



On the right is a diagram showing the reverse flow of computations, going from a *p*-value to a statistic.

Find the inverse form of the $\chi^2$ cdf for your package.

Names will differ among packages. Write the name of the $\chi^2$ cdf in your package_____

Now use the $\chi^2$ (df = 1) distribution to calculate the critical value of the $\chi^2$ statistic, using either p = 5% or 1–p = 95%, as appropriate for your package.

> **HINT for R Users:** Try qchisq

Did it return the critical value of $X^2 = 3.84$ ?          _____

If not, find the appropriate probability function so that the package returns a value of 3.84 corresponding to p = 5%.

14. An extremely conservative morphologist wants to work with Type I error set at $\alpha = 0.0001$, but does not have tables that supply critical values at this extreme probability level. Use the inverse CDF to obtain critical values of $\chi^2$ on a single degree of freedom test at $\alpha = 0.0001$.

Use your statistical package to compute the critical $\chi^2$ value.

For $\alpha = 0.0001$, the critical $\chi^2$ value for *df* = 1 is _____

Now make your own statistical table, for some unusual alpha levels of 0.002 and 0.0002, place critical values of *t* into the following t-table.

Table 3.1. Critical t-values

| d.f. | alpha | |
|---|---|---|
| | .002 | .0002 |
| 1 | 159.156 | 1590 |
| 5 | _____ | _____ |
| 10 | _____ | _____ |

These tables started becoming unnecessary in the 1970s, when calculators with statistical functions became available. Today we can report exact p-values with an application on our phones, or with a spreadsheet. Nevertheless, the tables of critical values are still with us. And unfortunately the practice of reporting critical probability levels (0.05, 0.01, 0.005, 0.001) rather than reporting the exact *p*-value is still with us.

**Write-up for this lab**.
       Fill in blank spaces as requested, on all the pages of this lab and submit to Brightspace.
       NOTE: Adobe Reader can be used to annotate the fill in the blanks on handouts.
       Adobe Reader is free to download at https://get.adobe.com/reader/

Learning goals for this lab.

Mechanics
    Use of spreadsheet functions or statistical package as a calculator  Ex. 1
    Use of probability models to
        calculate probabilities.  Ex 2,3
        calculate likelihoods  Ex 4
    Use of spreadsheet or package to display and print tables and graphs.  Ex 3,4,5

Concepts
    Use of probability models to design experiments Ex 2,3,5,6,7
    The difference between likelihood and probability.Ex 4
    Likelihood ratios as a measure of evidence, and how to calculate them. Ex 4,8
    Cumulative probabilities (the p-value) from cumulative distributions (cdfs) Ex 5-11
    One and two-tailed probabilities Ex 5-10
    Calculating probabilities from log likelihood ratios Ex 8-11
    Critical values and how to calculate them (inverse cdf) Ex 13-14