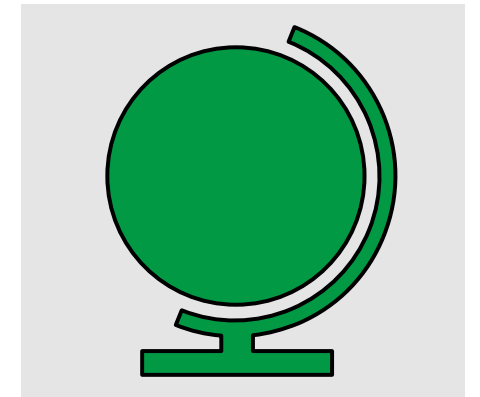


The Wide World of Distributions

BERDC Special Topics Talk 7



DaCCoTA
DAKOTA CANCER COLLABORATIVE
ON TRANSLATIONAL ACTIVITY

Dr. Mark Williamson
Biostatistics, Epidemiology,
and Research Design Core

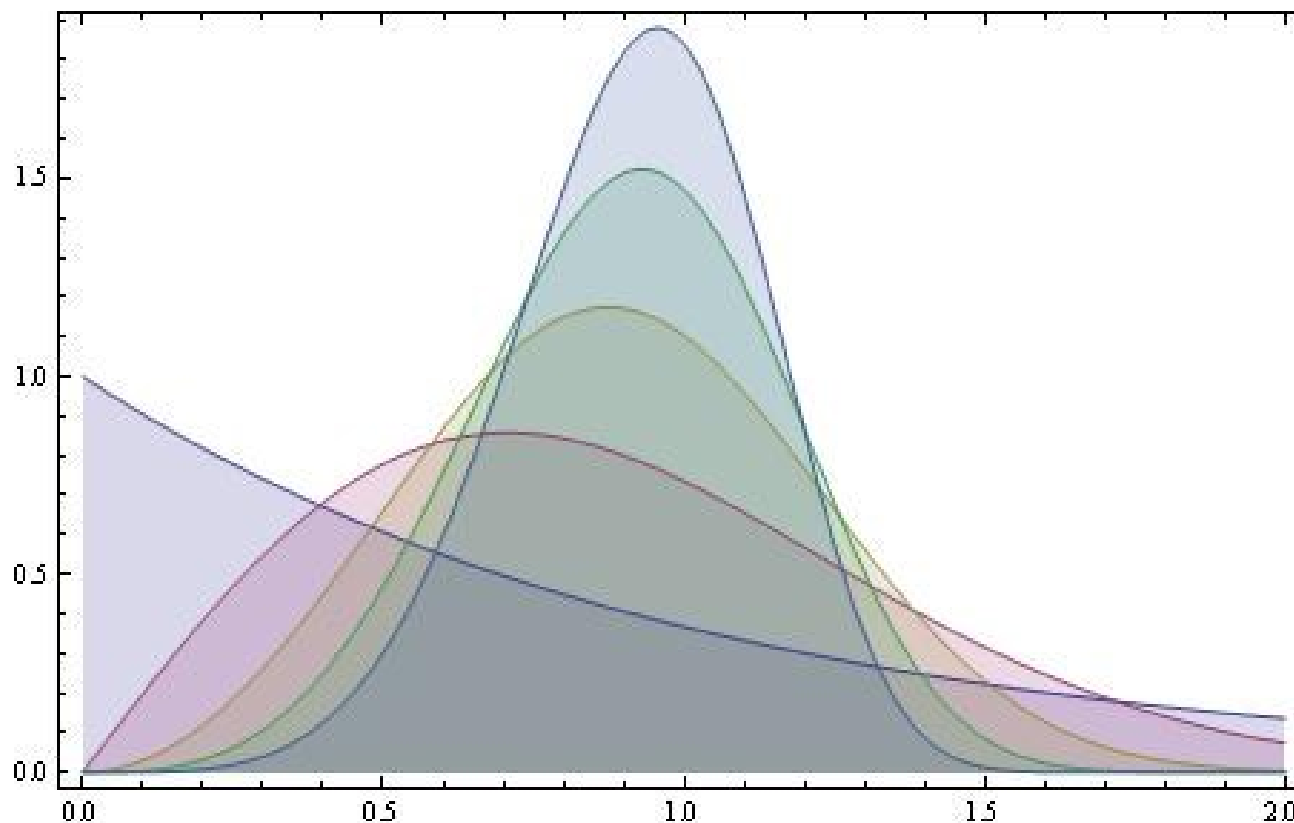
Opening

Goal: learn about statistical distributions

- Better statistical understanding
- Better statistical modeling
- Better statistical insights

Before Moving On:

Pre-test: https://und.qualtrics.com/jfe/form/SV_1QRRQKFLG086iuW

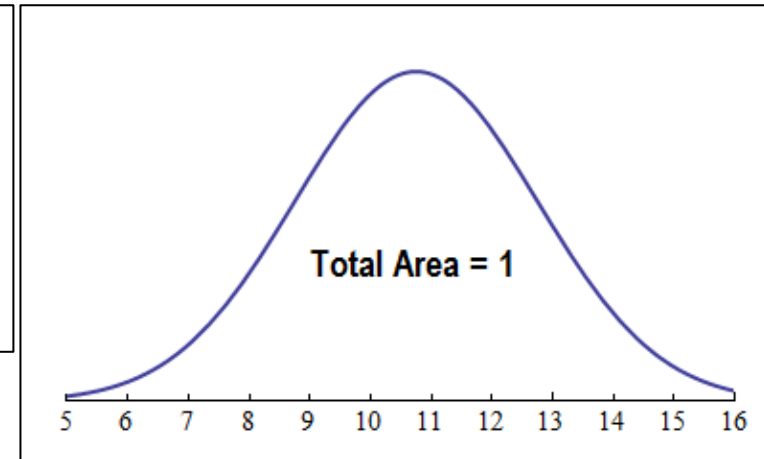
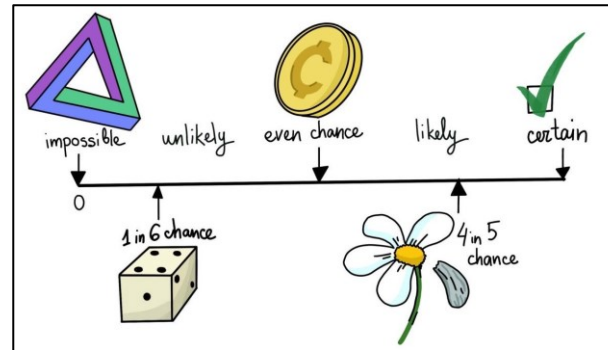


Definition

“Function that shows the possible values for a variable and how often they occur”

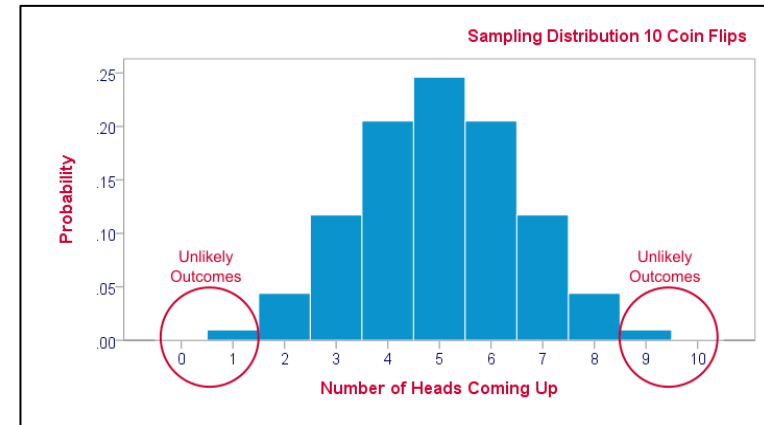
Properties

- probability that x can take a specific value is $p(x)$
- $p(x)$ is non-negative for all real x
- sum of $p(x)$ over all x equals 1



From the variable's perspective

- variables are sampled from random distributions that occur naturally
- each random distribution tells a different story about the processes that produce that variable
- each distribution has a mean, variance, and probability density function



“Mathematical construct to describe how variables are generated”

“Mathematical description of how data conceivably can be produced”

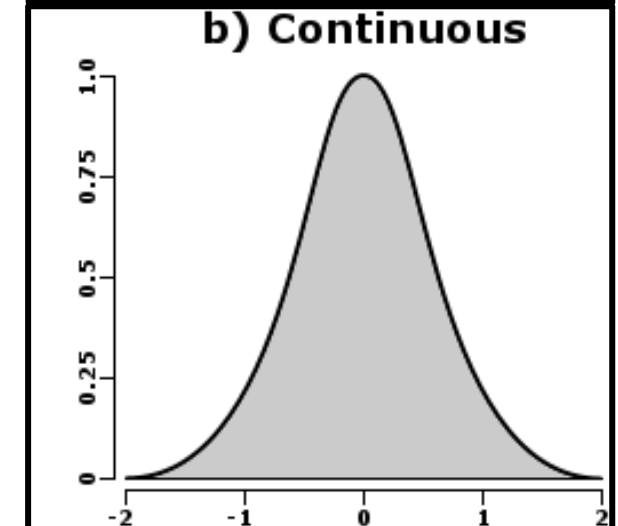
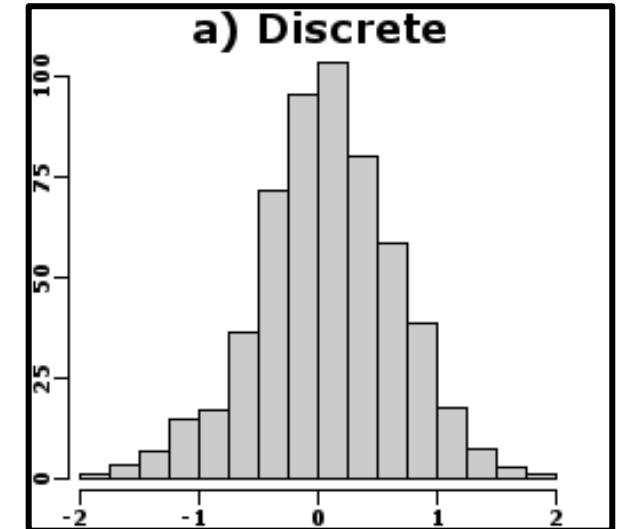
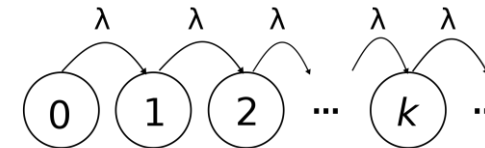
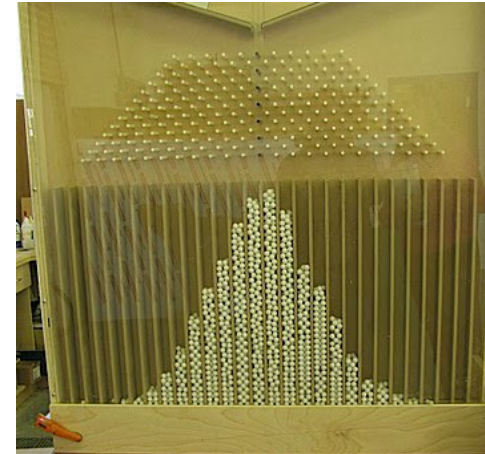
Variable origins

Random generation processes

- Gaussian – random dispersal from central point; random diffusion from mean
- Bernoulli Trial – discrete event with only two outcomes (success/failure) with constant probability of success: $P(\text{success}) = p$.
- Poisson point process – events that occur individually in continuous time (or space)

Broad classification

- Discrete - takes on only integer or 'count' variables
- Continuous – takes on any numerical value within range



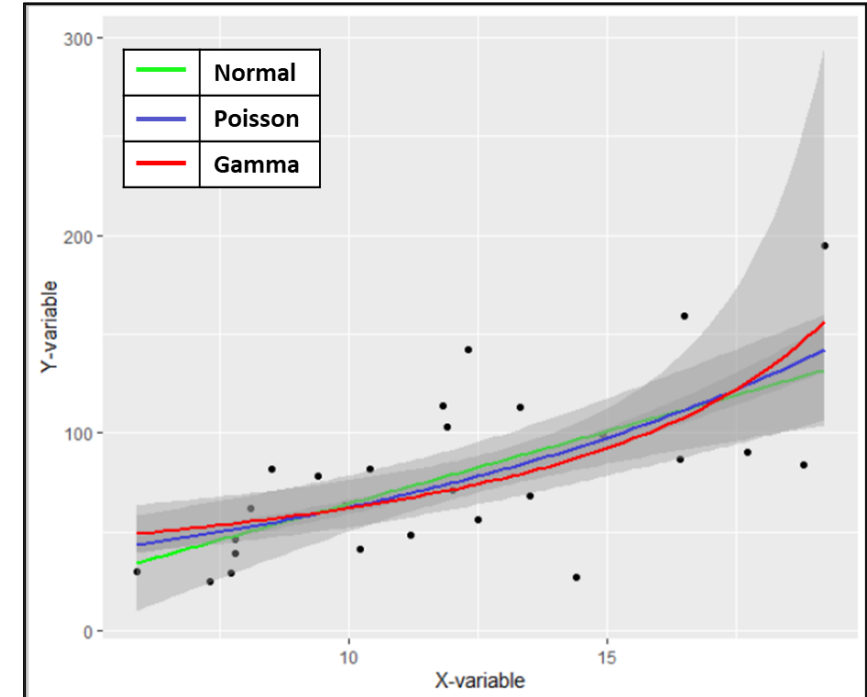
Properties

Name	Notation	Range	Mean(x)	Var(x)	Explanation
Normal	$N(\mu, \sigma^2)$	$-\infty < x < \infty$	μ	σ^2	x=dispersal from a central point, or a diffusion through a Gaussian filter, with variance independent of mean.
Log-Normal	$Lognormal(\mu, \sigma^2)$	$x > 0$	$\exp(\mu + \sigma^2/2)$	$[\exp(\sigma^2) - 1] * \exp(2\mu + \sigma^2)$	x =probability distribution whose logarithm is normally distributed.
Exponential	$exp(\beta)$	$x > 0$	β	β^2	x =time between events that occur at a rate of $\lambda = 1/\beta$.
Gamma	$Gamma(k, \theta)$	$x > 0$	$k\theta$	$k\theta^2$	x =time it takes for k events to occur within a rate of $\lambda = 1/\theta$, or the sum of k exponential events.
Beta	$Beta(a, b)$	$0 < x < 1$	$\frac{a}{a + b}$	$\frac{ab}{(a + b)^2 + (a + b + 1)}$	x =distribution of probabilities based on a successes and b failures, where both a and $b > 1$.
Binomial	$Bin(n, p)$	$x = 0, 1, 2, \dots$	np	$np(1-p)$	x =number of positive events out of n trials each with a probability of success p .
Geometric	$G(p)$	$x = 1, 2, 3, \dots$	$\frac{1}{p}$	$\frac{1-p}{p^2}$	x =number of trials, with probability of success p , that are needed to obtain one success.
Negative Binomial	$NB(n, p)$	$x = 0, 1, 2, \dots$	$\frac{k(1-p)}{p}$	$\frac{k(1-p)}{p^2}$	x =number of failures before k successes occur in sequential independent trials, all with the same probability of success, p .
Poisson	$Poisson(\lambda)$	$x = 0, 1, 2, \dots$	λ	λ	x =count of items in a standardized unit of effort that occur at rate λ .

Statistical tests

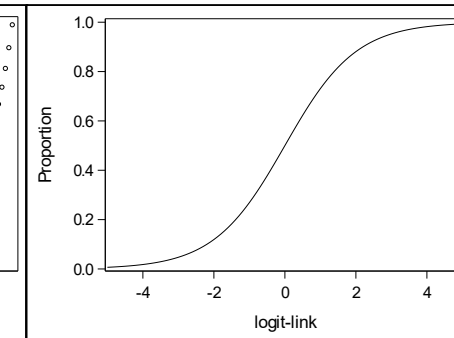
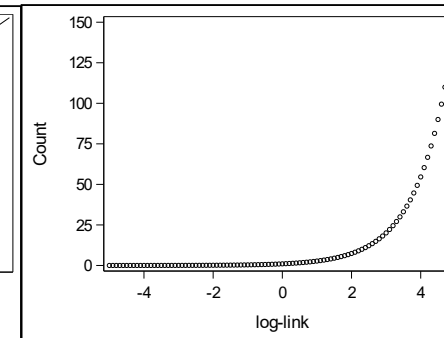
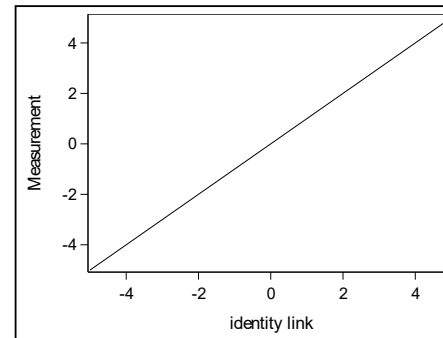
What sort of tests use distributions?

- Parametric tests assume a normal (Gaussian) distribution
 - ANOVA, t-test, linear regression, etc.
- Generalized linear models can use various distributions
 - ‘generalized’ refers to models that don’t assume a normal distribution
 - Common Examples are Poisson and Logistic regression



How to plot results?

- Each distribution has a standard link function
- Common ones are ‘identity’, ‘log’ and ‘logit’

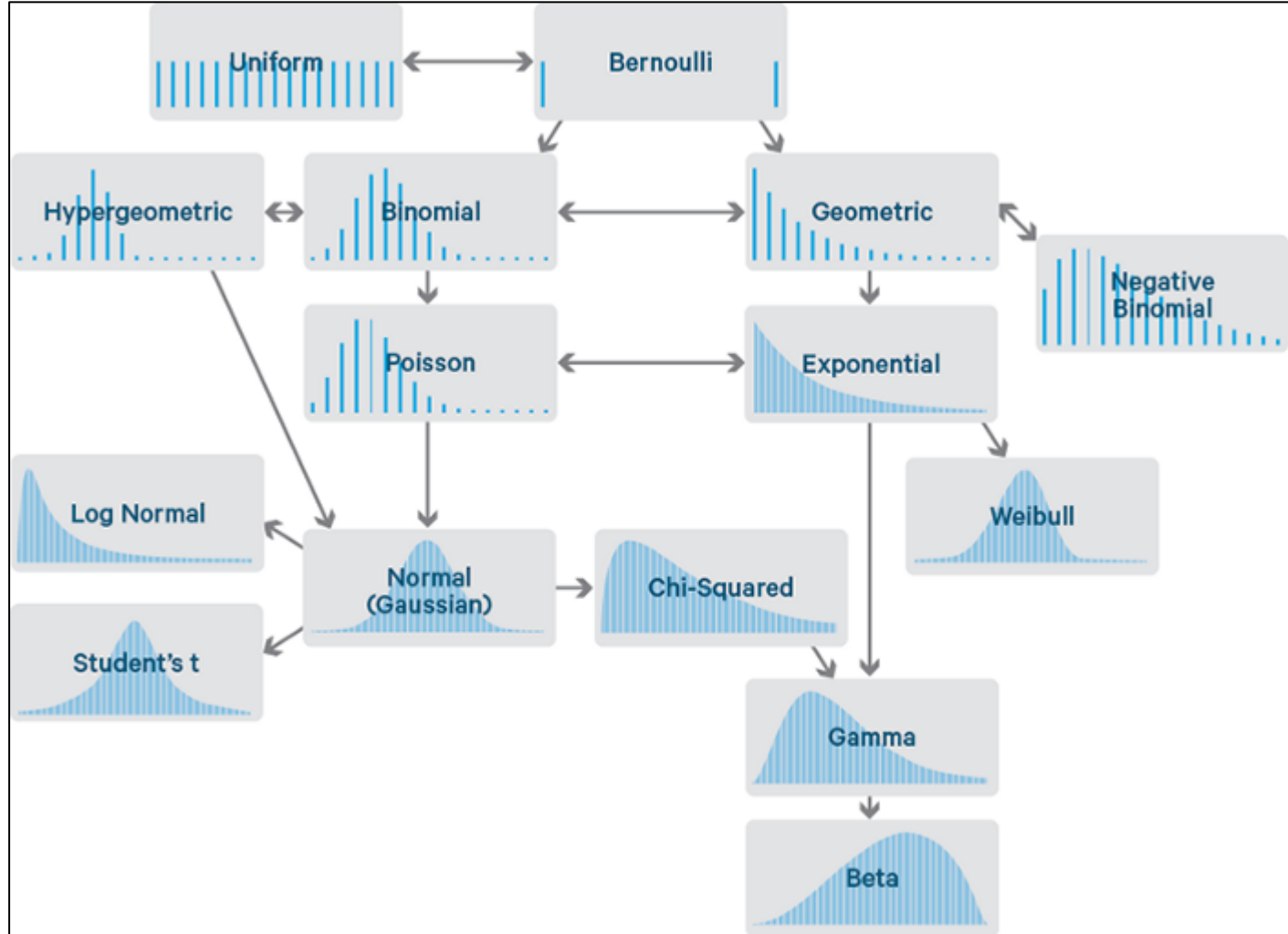


Taxonomy

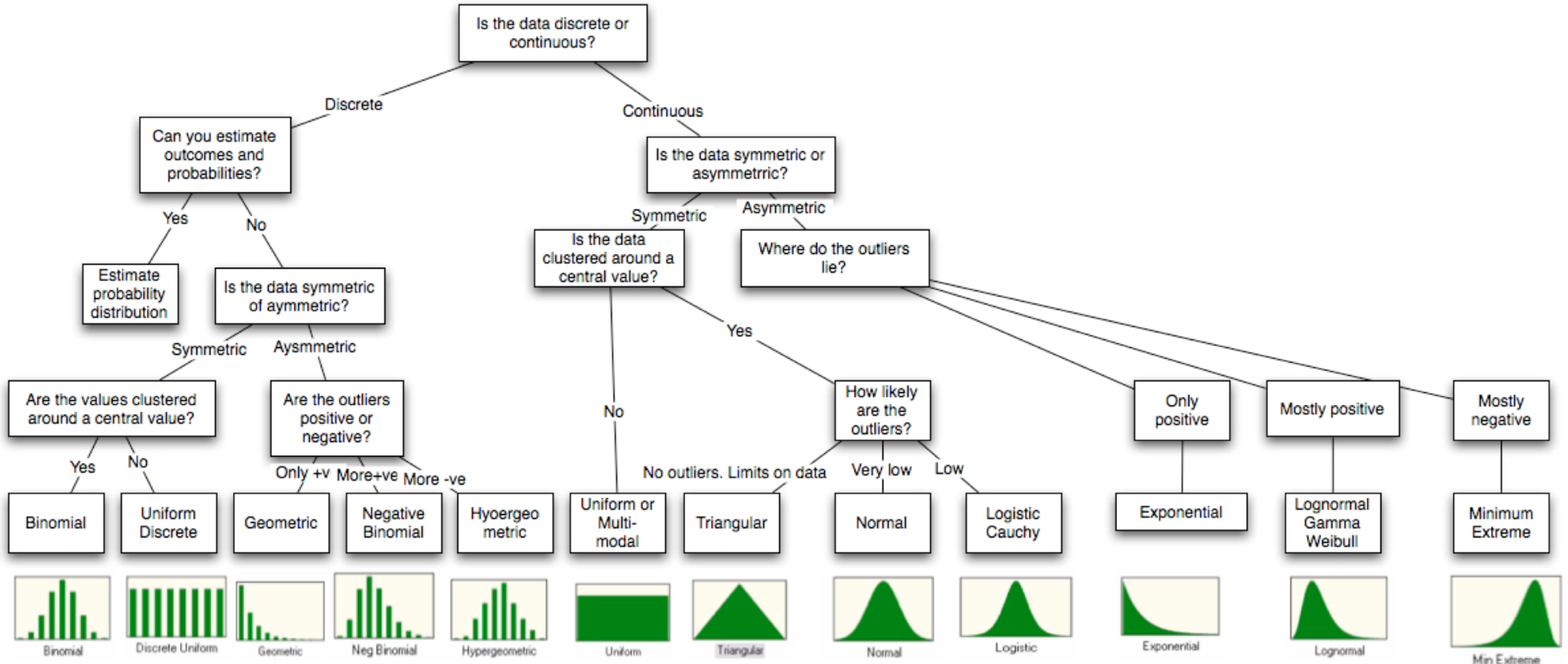
Taxonomy of Random Variables

Below is a taxonomy of common distributions, classified primarily by the topology of their **support**.

- ▶ finite
 - ▶ Dirac
 - ▶ Bernoulli
- ▶ countable
 - ▶ binomial
 - ▶ geometric
 - ▶ Poisson
- ▶ interval
 - ▶ uniform
 - ▶ beta
- ▶ half-line
 - ▶ exponential
 - ▶ gamma
- ▶ unbounded
 - ▶ normal
 - ▶ Cauchy
 - ▶ (Lévy) stable
- ▶ *transforms*
 - ▶ (generalized) Pareto
 - ▶ inverse gamma
 - ▶ lognormal
- ▶ *mixtures*
 - ▶ (Gosset) Student t
 - ▶ negative-binomial
- ▶ *non-parametric*
 - ▶ empirical



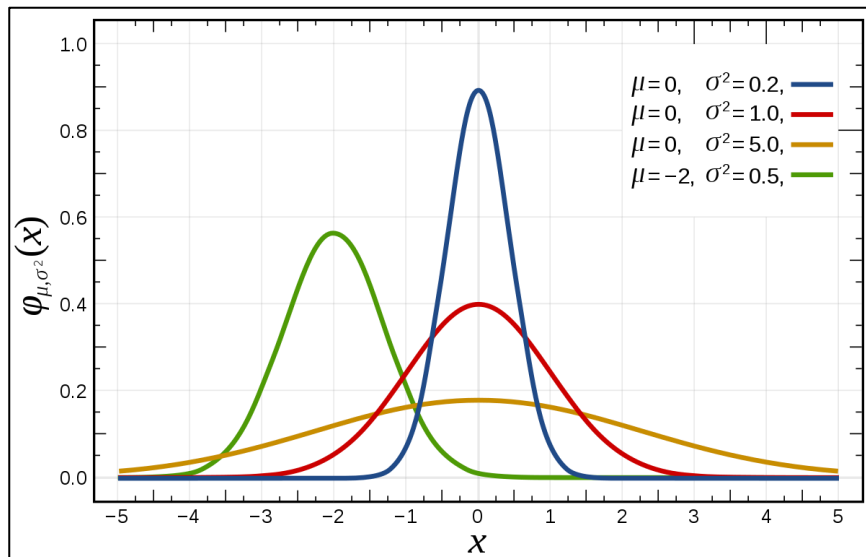
Taxonomy cont.



Normal & Log Normal

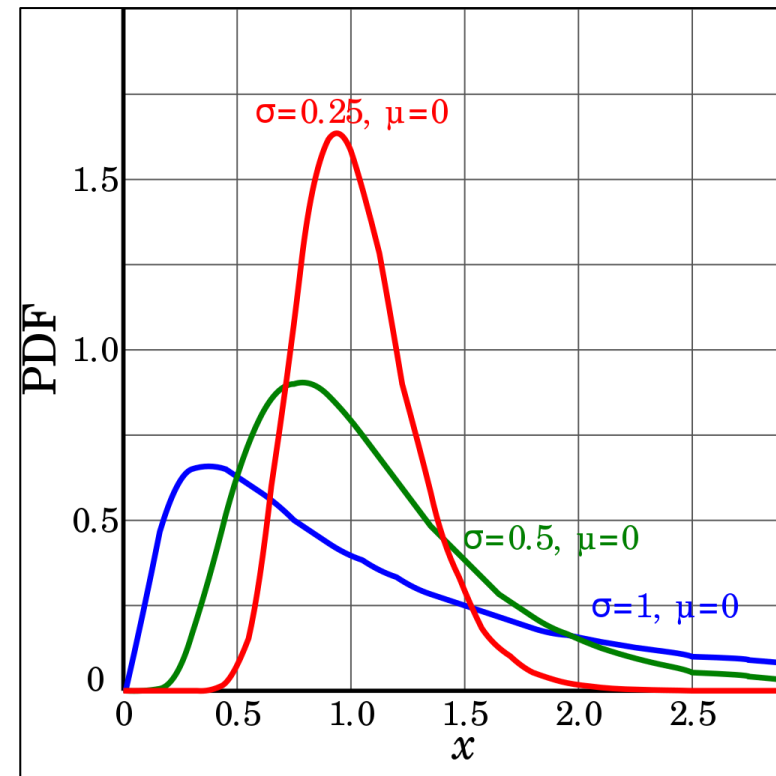
Normal

- Distribution models dispersion from central point
- Continuous, unbounded, and with a variance independent of its mean



Log Normal

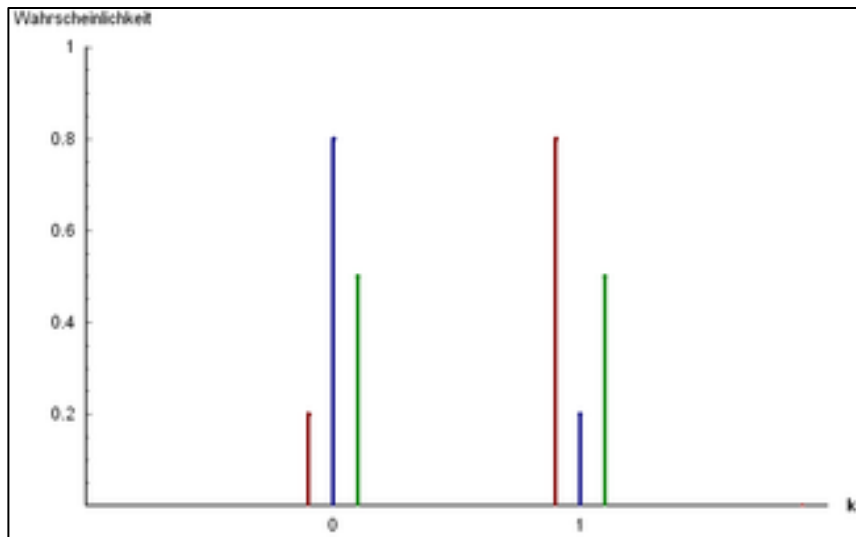
- Distribution whose log-transformation is normal



Binary & Binomial

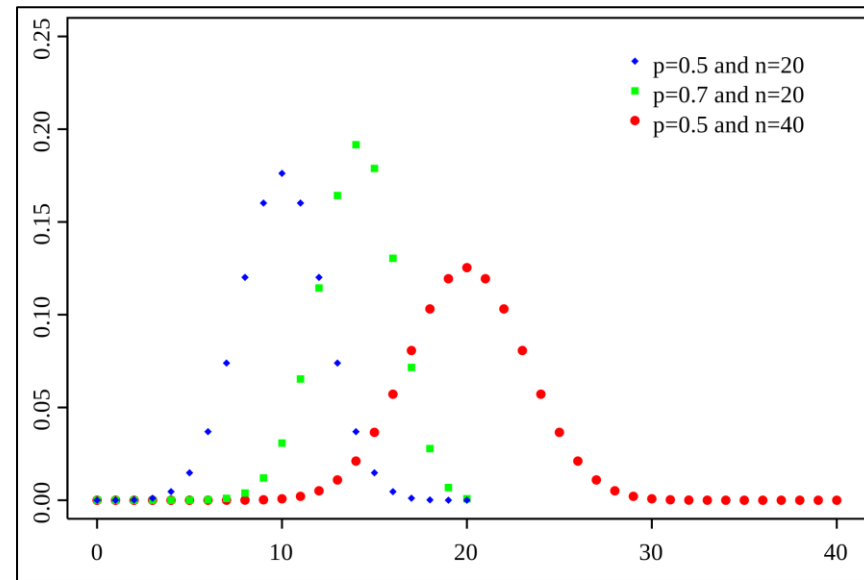
Binary

- Distribution that models binary (two) outcomes
- Aka 'Bernoulli trials'
- Example: coin flips



Binomial

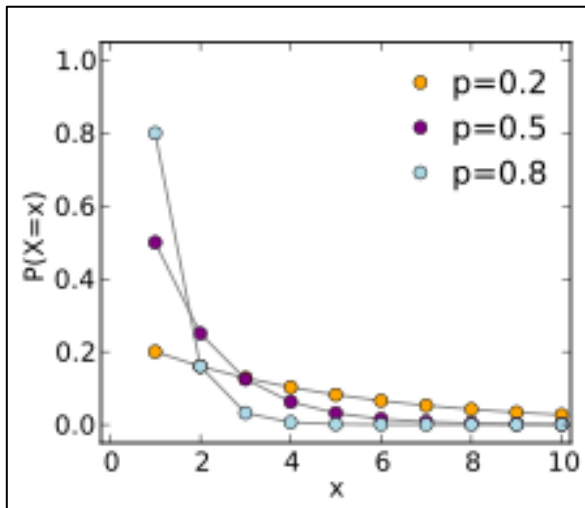
- Distribution that models the number of success from Bernoulli trials
- Example: coin flips



Geometric, Poisson, & Negative Binomial

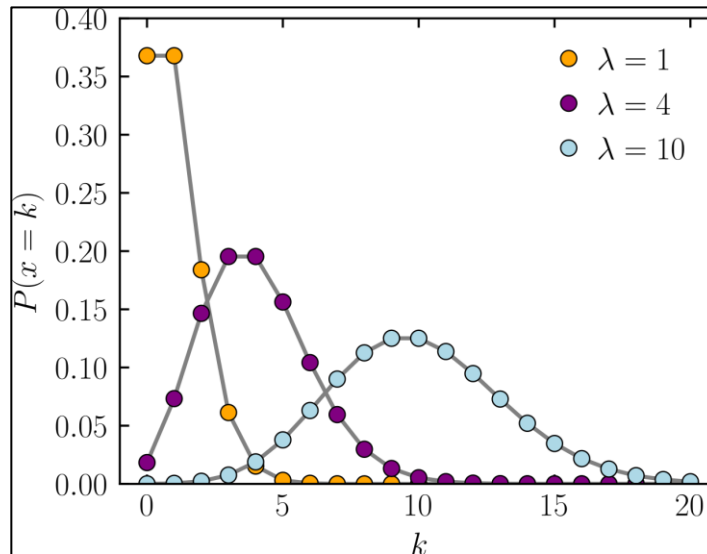
Geometric

- Distribution that models number of Bernoulli trials needed for one success
- Start at 0 or 1
- Example: coin flips until heads



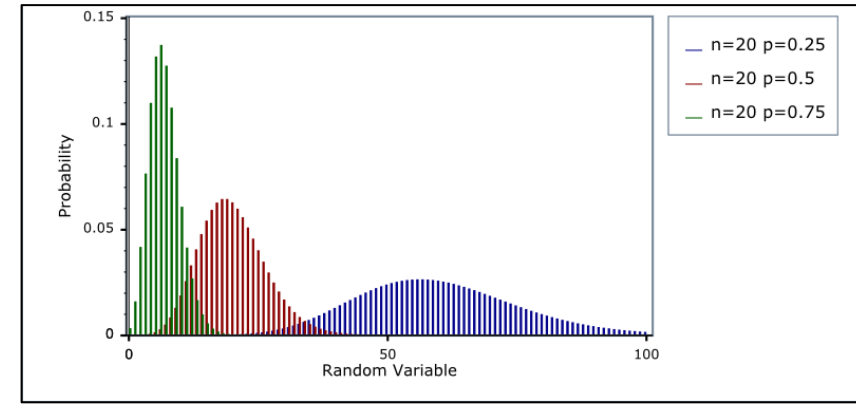
Poisson

- Distribution that models the number of counts occurring at a certain rate
- Example: number of patrons entering restaurant an hour, or number of trees per acre of woods



Negative Binomial

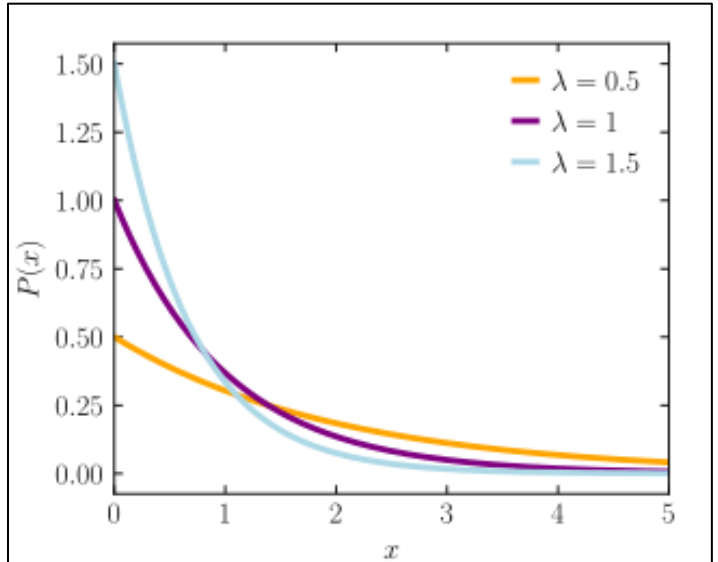
- Distribution that models the number of failures needed for a certain number of successes
- Can also be used for 'count' data as a sort of over-dispersed correction for Poisson
- Example: flip coins until 3 heads, count # of tails



Exponential, Gamma, & Beta

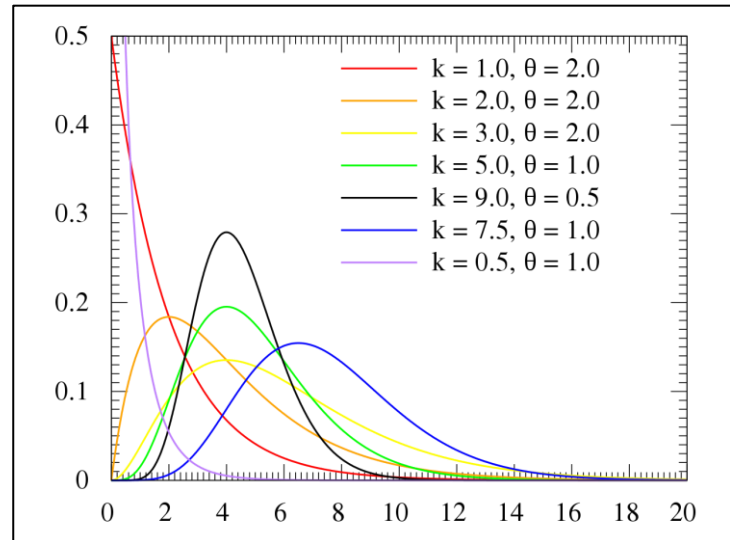
Exponential

- Distribution that models the inter-arrival time between Poisson process events
- Example: time between arrive of patrons into a restaurant



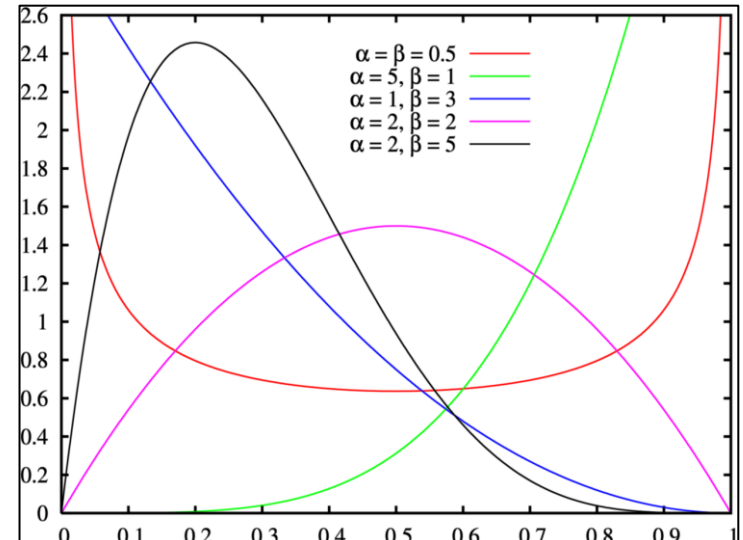
Gamma

- Distribution that models the total arrival time for a number of Poisson process events
- Example: time it takes for 5 patrons to enter a restaurant

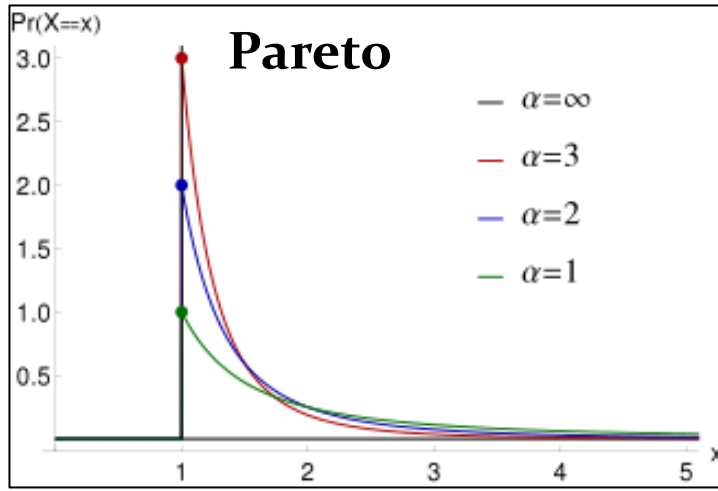
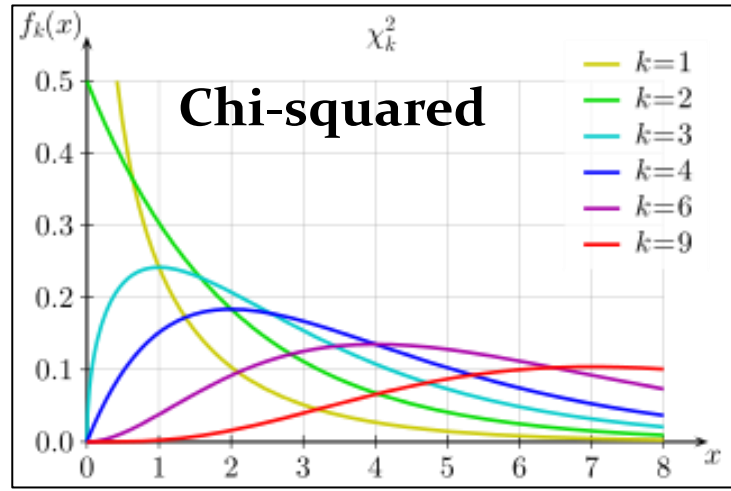
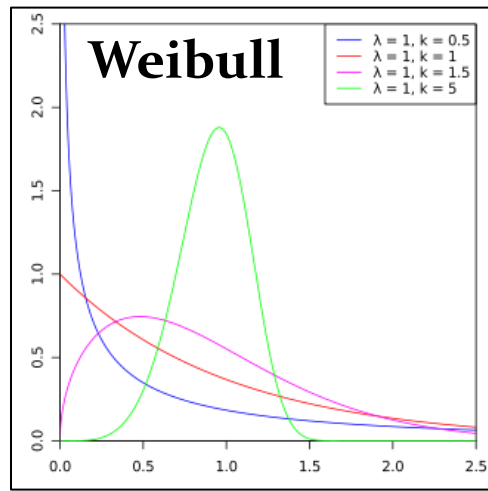
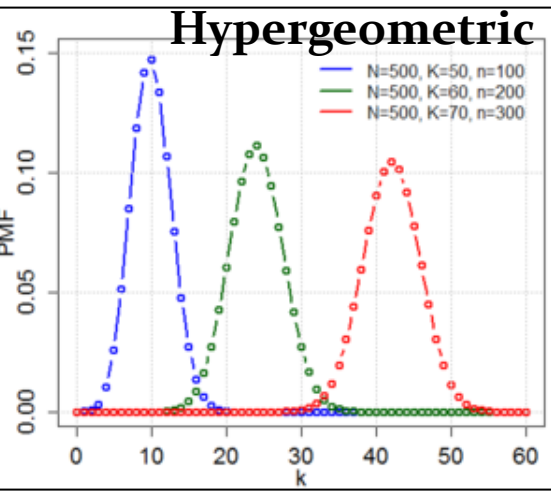
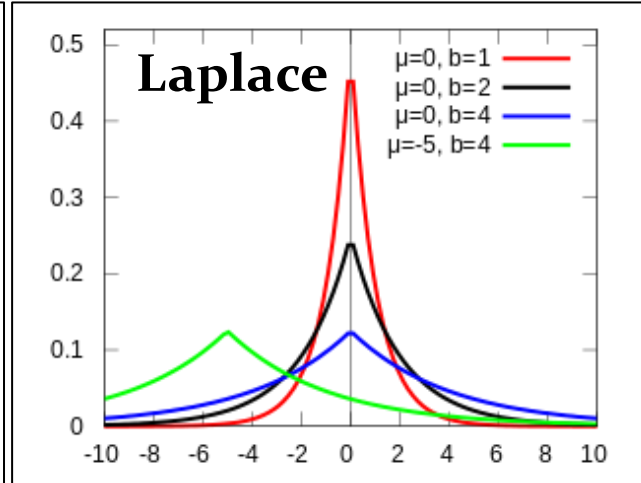
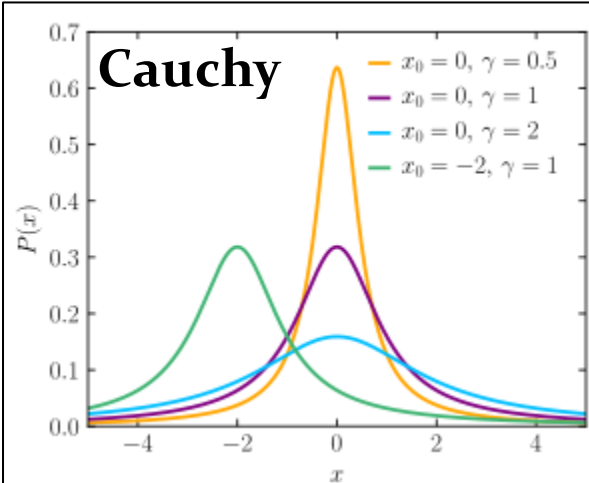
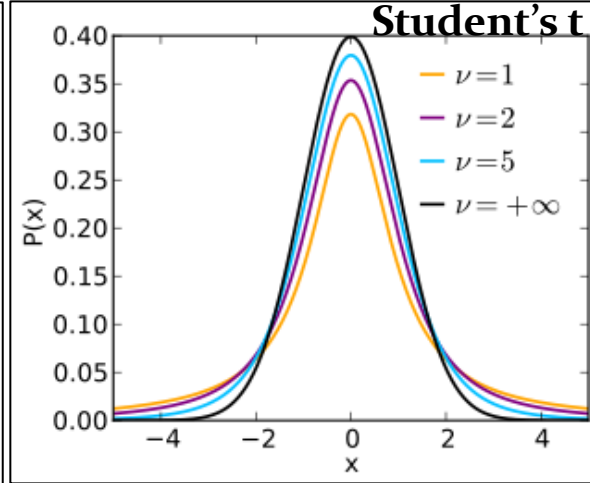
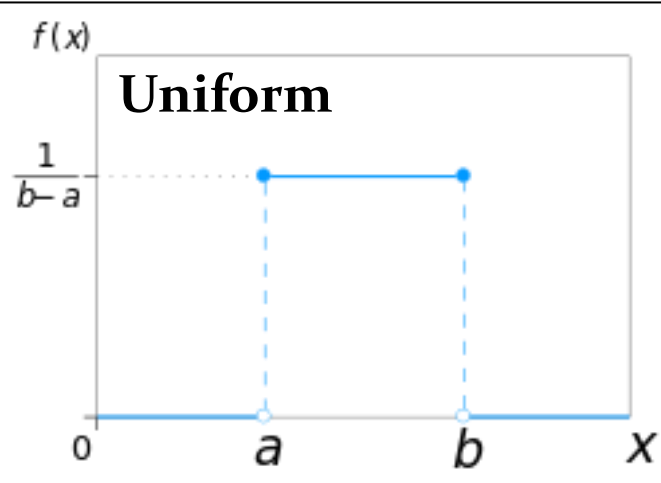


Beta

- Distribution that models a proportion of success between 0 and 1
- Example: proportion of hits out of at-bats



Other distributions



SAS examples

DIST=	Distribution	Default Link Function	Numeric Value
BETA	beta	logit	12
BINARY	binary	logit	4
BINOMIAL BIN B	binomial	logit	3
EXPONENTIAL EXPO	exponential	log	9
GAMMA GAM	gamma	log	5
GAUSSIAN G NORMAL N	normal	identity	1
GEOMETRIC GEOM	geometric	log	8
INVGauss IGAUSSIAN IG	inverse Gaussian	inverse squared	6
		(power(-2))	
LOGNORMAL LOGN	lognormal	identity	11
MULTINOMIAL MULTI MULT	multinomial	cumulative logit	NA
NEGBINOMIAL NEGBIN NB	negative binomial	log	7
POISSON POI P	Poisson	log	2
TCENTRAL TDIST T	t	identity	10
BYOBS(variable)	multivariate	varied	NA

```

Distributions Examples.sas
CODE LOG RESULTS
*Create and view distributions;
DATA rand;
  do i=1 to 10000;
    * generate random values for 10000 observations;
    trt="trt";
    mynormal=rand("NORMAL", 85, 40);
    mylognormal=rand("LOGNORMAL", 1, 0.1);
    mybeta=rand("BETA", 1,99);
    mybinary=rand("BINOM", 0.95, 1);
    mybinomial=rand("BINOM", 0.95, 2000);
    myexponential=rand("EXPONENTIAL", 4.2);
    mygamma=rand("GAMMA", 2, 2);
    mygeometric=rand("GEOMETRIC", 0.10);
    mynegbinomial=rand("NEGBINOMIAL", 0.97, 1000);
    mypoisson=rand("POISSON", 10);
  output;
  end;
PROC PRINT data=rand(obs=20);
PROC UNIVARIATE data=rand;
  var mynormal mylognormal mybeta mybinary mybinomial myexponential mygamma mygeometric
  mynegbinomial mypoisson;
  histogram;
  *can change around the parameters;
  *-----;
*Examples of distributions in action;
*Datasets;
DATA Fish; set sashelp.Fish;
PROC PRINT data=Fish(obs=25);
DATA Baseball; set sashelp.Baseball;
PROC PRINT data=Baseball(obs=25);
*Normal;
DATA INMVAR;
/home/markwilliamson20/my_courses/markwilliamson0/MW_2021_Work/Module Examples/Distributions Examples.sas
Messages User: markwilliamson20
  
```

Code available at: https://med.und.edu/daccota/files/docs/berdc_docs/distributions_sas_code.txt

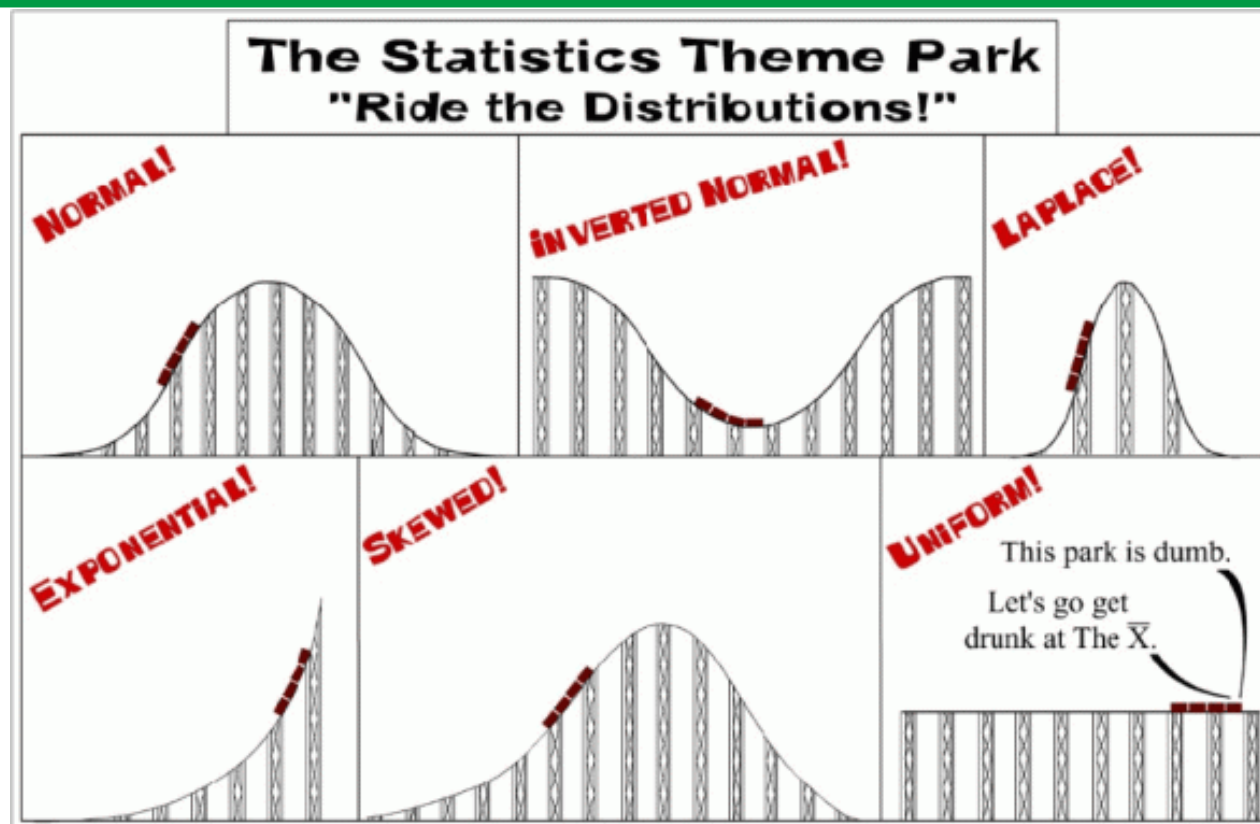
R examples

Family	Default Link Function
binomial	(link = "logit")
gaussian	(link = "identity")
Gamma	(link = "inverse")
inverse.gaussian	(link = "1/mu^2")
poisson	(link = "log")
quasi	(link = "identity", variance = "constant")
quasibinomial	(link = "logit")
quasipoisson	(link = "logit")

```

RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function
Addins
Source on Save
Run Source
In selection Match case Whole word Regex Wrap
1 #Create and view distributions
2
3 mynormal <- rnorm(n=10000, mean = 85, sd = 40)
4 mylognormal <- rlnorm(n=10000, meanlog = 1, sdlog = 0.1)
5 mybeta <- rbeta(n=10000, shape1=1, shape2=99)
6 mybinary <- rbinom(n=10000, size=1, prob=0.95)
7 mybinomial <- rbinom(n=10000, size=2000, prob=0.95)
8 myexponential <- rexp(n=10000, rate=4.2)
9 mygamma <- rgamma(n=10000, shape=2, rate=2) #rgamma(n, shape, rate = 1, scale = 1/rate)
10 mygeometric <- rgeom(n=10000, prob=0.10)
11 mynegbinomial <- rnbinom(n=10000, size=1000, prob=0.97)
12 mypoisson <- rpois(n=10000, lambda=10)
13
14 hist(mynormal)
15 hist(mylognormal)
16 hist(mybeta)
17 hist(mybinary)
18 hist(mybinomial)
19 hist(myexponential)
20 hist(mygamma)
21 hist(mygeometric)
22 hist(mynegbinomial)
23 hist(mypoison)
24
25 #can change around the parameters
  
```

Wrap-up



Please take the post-test and survey:

Post-test: https://und.qualtrics.com/jfe/form/SV_5swKhmpU4tO7hau

Survey: https://und.qualtrics.com/jfe/form/SV_5o370oBxAOTtMr4

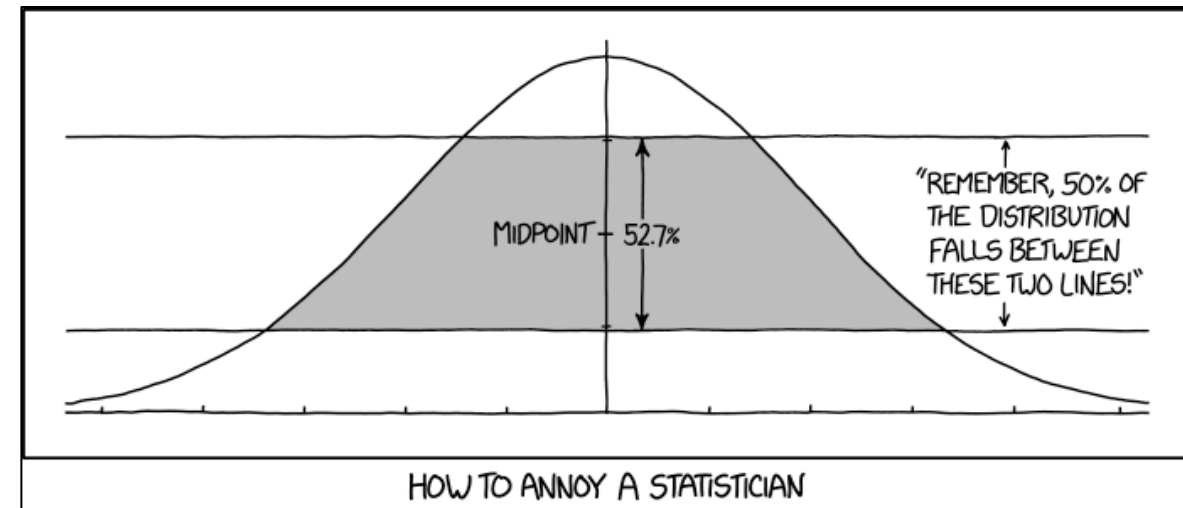
References

Images:

- <https://www.wolfram.com/mathematica/newin6/content/SymbolicStatisticalComputing/BuiltinStatisticalDistributions.html>
- <https://pierpaolo28.github.io/blog/blog19/>
- <https://courses.lumenlearning.com/wmopen-concepts-statistics/chapter/continuous-probability-distribution-2-of-2/>
- <https://www.spss-tutorials.com/sampling-distribution-what-is-it/>
- <https://medium.com/analytics-vidhya/probability-distributions-444e7babf2e1>
- <https://www.kdnuggets.com/2020/02/probability-distributions-data-science.html>
- <https://www-users.cse.umn.edu/~dodso013/fm503/0910/lectures/fall2.pdf>
- <https://tinyheero.github.io/2016/03/17/prob-distr.html>
- https://en.wikipedia.org/wiki/List_of_probability_distributions
- <https://loonylabs.files.wordpress.com/2019/10/distribution-park-joke.gif?w=590>

Materials:

- <https://365datascience.com/tutorials/statistics-tutorials/distribution-in-statistics/>
- <https://www.itl.nist.gov/div898/handbook/eda/section3/eda361.htm>
- <https://www.statisticshowto.com/negative-binomial-experiment/>
- <https://www-users.cse.umn.edu/~dodso013/fm503/0910/lectures/fall2.pdf>
- <https://medium.com/@srowen/common-probability-distributions-347e6b945ce4>
- https://www.causascientia.org/math_stat/Dists/Compendium.pdf



<https://xkcd.com/2118/>

Acknowledgements



- The DaCCoTA is supported by the National Institute of General Medical Sciences of the National Institutes of Health under Award Number U54GM128729.
- For the labs that use the Biostatistics, Epidemiology, and Research Design Core in any way, including this Module, please acknowledge us for publications. ***"Research reported in this publication was supported by DaCCoTA (the National Institute of General Medical Sciences of the National Institutes of Health under Award Number U54GM128729)"***.

DaCCoTA
DAKOTA CANCER COLLABORATIVE
ON TRANSLATIONAL ACTIVITY