**Discrete Distribution Families**
**Bernoulli($p$)**
Binary trial (success/failure).
PMF: $P(X = x) = p^x(1-p)^{1-x}$, $x \in \{0,1\}$
Mean: $p$, Var: $p(1-p)$
R: d/p/q/rbinom(n=1, p)
**Binomial($n,p$)**
successes in $n$ Bernoulli trials.
PMF: $P(X = x) = \binom{n}{x}p^x(1-p)^{n-x}$
Mean: $np$, Var: $np(1-p)$
R: dbinom, pbinom, qbinom, rbinom
**Geometric($p$)**
Failures before 1st success.
PMF: $P(X = x) = (1-p)^x p$, $x \geq 0$
Mean: $(1-p)/p$, Var: $(1-p)/p^2$
R: dgeom, pgeom, qgeom, rgeom
**NegBin($k,p$)**
Failures before $k$-th success.
PMF: $P(X = x) = \binom{x+k-1}{x}p^k(1-p)^x$
Mean: $k(1-p)/p$, Var: $k(1-p)/p^2$
R: dnbinom, pnbinom, qnbinom, rnbinom
**Poisson($\lambda$)**
Counts in interval, rate $\lambda$.
PMF: $P(X = x) = \frac{\lambda^x e^{-\lambda}}{x!}$
Mean: $\lambda$, Var: $\lambda$
R: dpois, ppois, qpois, rpois
**Continuous Distribution Families**
**Uniform($a,b$)**
All values equally likely on $[a,b]$.
PDF: $f(x) = \frac{1}{b-a}$, $a \leq x \leq b$
Mean: $(a+b)/2$, Var: $(b-a)^2/12$
Skewness: 0 (symmetric)
R: dunif, punif, qunif, runif

**Normal($\mu,\sigma^2$)**
Bell-shaped, $-\infty < \mu < \infty, \sigma^2 > 0$
PDF: $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}}e^{-(x-\mu)^2/(2\sigma^2)}$
Mean: $\mu$, Var: $\sigma^2$
Skewness: 0 (symmetric)
R: dnorm, pnorm, qnorm, rnorm;

**Lognormal($\mu,\sigma^2$)**
If $\ln X \sim N(\mu,\sigma^2)$.
PDF: $f(x) = \frac{1}{x\sigma\sqrt{2\pi}}e^{-(\ln x-\mu)^2/(2\sigma^2)}$
Mean: $e^{\mu+\sigma^2/2}$, Var: $(e^{\sigma^2}-1)e^{2\mu+\sigma^2}$
Skewness: $(e^{\sigma^2}+2)\sqrt{e^{\sigma^2}-1}$

R: dlnorm, plnorm, qlnorm, rlnorm;

**Exponential($\lambda$)**
Time between Poisson events.
Positive RV, wait time, memoryless
($\lambda$) is average rate
($\beta$) is mean wait time
PDF: $f(x) = \lambda e^{-\lambda x}, x \geq 0$
Mean: $1/\lambda$, Var: $1/\lambda^2$
Skewness: 2
R: dexp, pexp, qexp, rexp;

**Beta($\alpha,\beta$)**
On $[0,1]$, Bayesian priors.
Uniform distribution special case
The Gamma function ($\Gamma$)is a generalization
of the factorial function to non-integer numbers.
PDF: $f(x) = \frac{1}{B(\alpha,\beta)}x^{\alpha-1}(1-x)^{\beta-1}$
Mean: $\alpha/(\alpha+\beta)$, Var: $\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$
Skewness: $\frac{2(\beta-\alpha)\sqrt{\alpha+\beta+1}}{(\alpha+\beta+2)\sqrt{\alpha\beta}}$
R: dbeta, pbeta, qbeta, rbeta;

**Weibull($k,\lambda$)**
Lifetimes, survival, reliability.
Exponential family (when k = 1), longer you wait
Survival Analysis
PDF: $f(x) = \frac{k}{\lambda}(x/\lambda)^{k-1}e^{-(x/\lambda)^k}$
Mean: $\lambda\Gamma(1+1/k)$, Var: $\lambda^2[\Gamma(1+2/k)-\Gamma(1+1/k)^2]$
Skewness: $\frac{\Gamma(1+3/k)\lambda^3-3\mu\sigma^2-\mu^3}{\sigma^3}$
R: dweibull, pweibull, qweibull, rweibull;

**Gamma($\alpha,\theta$)**
Waiting time for $\alpha$ events.
non-negative numbers
$\alpha$ is shape parameter
$\theta$ is scale parameter
PDF: $f(x) = \frac{1}{\Gamma(\alpha)\theta^\alpha}x^{\alpha-1}e^{-x/\theta}$
Mean: $\alpha\theta$, Var: $\alpha\theta^2$
Skewness: $2/\sqrt{\alpha}$
R: dgamma, pgamma, qgamma, rgamma;

**R: Computing Expected Value**
**Discrete RV**
```
X <- 0:2
p <- c(0.2,0.5,0.3)
EV <- sum(X*p)
```

**Continuous RV (numerical integration)**
```
f <- function(x) 2*x
EV <- integrate(function(x) x*f(x), 0,1)$value
```

**From sample (Monte Carlo)**
```
samples <- rnorm(10000, mean=5, sd=2)
mean(samples)
```
**Define the PDF function**
```
f <- function(x)   2*x
prob <- integrate(f, lower = 0.5, upper =
0.75)$value # Integrate over [0.5, 0.75]
```
**Conditional Distributions**
Let $X$ and $Y$ be two random variables.

$$f_{X|Y}(x|y) = \begin{cases} \frac{P(X=x,Y=y)}{P(Y=y)}, & \text{discrete case} \\ \frac{f_{X,Y}(x,y)}{f_Y(y)}, & \text{continuous case, } f_Y(y) > 0 \end{cases}$$

Conditional dist is just a segment of marginal
dist, then re-normalized to have an area
under the curve equal to 1        If X and Y are
independent, in continouse case,

$$f_{Y|X}(y) = f_Y(y)$$

This means conditional PDF of Y and X is
marginal PDF of Y.
**Random Sample**
It is independent and identically distributed
(iid).  Each pair of observations are
independent, and each observation comes from
the same distribution.
**MLE**
Great way to find estimators.  Applied on
multi and univariate.  Relies on random sample
of n observations.  Mean, is a estimator for
univariate, multivariate, linear regression
**Steps for MLE**
1.Nature of variable (discrete or contin)
2.estimate the parameters of a theoretical
distribution (eg $\lambda$ in a Poisson distribution)
3.Choose distribution:  Normal, Exponential,
Poisson, Binomial, etc.)        4.Play with the
parameters for that family of distributions
to find the one that would be most likely
given our data and choose the corresponding
parametric estimates)        5.To obtain these
estimates, we use the likelihood function of
our observed random sample.)